

1

# When norm change hurts

2

Charles Efferson<sup>1,\*</sup>, Sönke Ehret<sup>1</sup>, Lukas von Flüe<sup>1</sup>, and Sonja Vogt<sup>1,\*</sup>

3

<sup>1</sup>Faculty of Business and Economics, University of Lausanne, Switzerland

4

\*Address correspondence to *charles.efferson@unil.ch* and *sonja.vogt@unil.ch*.

5

**Word count:** c. 9500 (excluding abstract, references, and captions)

6

**Key words:** social tipping, behaviour change, coordination, conformity, social norms

## 7 **Abstract**

8 Applied cultural evolution includes any effort to mobilise social learning and cultural evolution  
9 to promote behaviour change. Social tipping is one version of this idea based on conformity  
10 and coordination. Conformity and coordination can reinforce a harmful social norm, but  
11 they can also accelerate change from a harmful norm to a beneficial alternative. Perhaps  
12 unfortunately, the link between the size of an intervention and social tipping is complex in  
13 heterogeneous populations. A small intervention targeted at one segment of society can in-  
14 duce tipping better than a large intervention targeted at a different segment. We develop  
15 and examine two models showing that the link between social tipping and social welfare is  
16 also complex in heterogeneous populations. An intervention strategy that creates persistent  
17 miscoordination, exactly the opposite of tipping, can lead to higher social welfare than an-  
18 other strategy that leads to tipping. We show that the potential benefits of miscoordination  
19 often hinge specifically on the preferences of people most resistant to behaviour change. Al-  
20 together, ordinary forms of heterogeneity complicate applied cultural evolution considerably.  
21 Heterogeneity weakens both the link between the size of a social planner's intervention and  
22 behaviour change and the link between behaviour change and the well-being of society.

## 23 **1 Introduction**

24 Applied cultural evolution is, to shun euphemism, an attempt to engineer culture. A social  
25 planner wants people to behave differently, and she intervenes in society in pursuit of this  
26 objective. The interesting twist is that, once people exposed to the intervention start to  
27 change behaviour, endogenous cultural evolutionary processes can take effect. Some people  
28 change behaviour because they have direct experience with the intervention. Some people  
29 change behaviour because they observe others doing so. If the social planner knows how this  
30 second process works, she can implement her intervention in a way that maximises the sum of  
31 both the direct effect and the associated indirect cultural evolutionary effect. In particular,  
32 the indirect effect might far outstrip the direct effect, in which case cultural evolutionary  
33 dynamics dramatically amplify the intervention's consequences. This idea is the essence of  
34 applied cultural evolution as an attempt to engineer culture [1].

35 The indirect cultural evolutionary effect occurs because we influence, teach, and learn  
36 from each other, and we do not do so randomly [2]. We pay attention to some people and  
37 ignore others [3, 4]. Sometimes we follow the majority, and sometimes we do not [5, 6]. Some  
38 people provide examples of how to behave, and some people provide examples of how not

39 to behave [7–9]. Some behaviours we simply like, and others we do not [10]. Whatever the  
40 details, the necessary result is some kind of cultural evolution at the aggregate level [1, 11, 12].  
41 If we want to know what kind of cultural evolution, the details are crucial [1, 11, 13, 14].  
42 If we discriminate when we learn from each other in one way, cultural evolution unfolds  
43 accordingly. If we discriminate in some other way, we can expect cultural evolution to unfold  
44 quite differently.

45     Engineering culture sounds unpleasant, even imperialistic, and sometimes it is [15–17].  
46 However, because most people live in a society, any attempt to modify anyone’s behaviour  
47 comes with the potential to induce a secondary cultural evolutionary effect. Every policy  
48 change, persuasion campaign, marketing push, therapy session, and passing advice for a  
49 friend is an intervention that can affect the individuals directly exposed. Because others are  
50 watching, it may also activate subsequent cultural evolution. Thus, the crucial question is  
51 not a matter of whether we want applied cultural evolution. We have it, and we will keep it.  
52 Rather, the crucial question is, do we have the insight and wherewithal to manage cultural  
53 evolutionary processes for the benefit of society? This paper develops a model focused on a  
54 particular version of this second question. We examine the link between behaviour change and  
55 social welfare. We show that in heterogeneous populations, even though everyone faces clear  
56 incentives to behave like others, behaving like others is not always best. As a consequence, an  
57 intervention that triggers large-scale norm change can actually be worse than an intervention  
58 that generates chronic disagreement.

59     We assume that everyone somehow wants to behave like others because of some mix of  
60 conformity and coordination incentives. Conformity and coordination incentives can create  
61 multiple equilibria. Social norms, by which we mean a shared understanding of how people  
62 should behave and how people do behave, help people collectively pick a specific equilibrium.  
63 In the simplest case with two behaviours, one locally stable steady state has everyone choos-  
64 ing one behaviour, and another locally stable steady state has everyone choosing the other  
65 behaviour. The population has converged on one of these equilibria, but the two states may  
66 not be equally good for society. One can be relatively harmful and the other relatively benefi-  
67 cial. Because both are locally stable, the population can get stuck in the harmful equilibrium.  
68 Happily, however, the same conformity and coordination incentives that trap the population  
69 in the harmful equilibrium can create the potential for a rapid transition to the beneficial  
70 equilibrium. A sufficiently large shock, a social planner’s intervention for example, can dis-  
71 lodge the population from the harmful equilibrium and tip it into the basin of attraction  
72 for the beneficial alternative. Once this happens, cultural evolutionary forces finish the job.

73 Conformity and coordination incentives ensure that the population completes the transition  
74 to the a new socially beneficial norm without further inputs from the social planner.

75 This is the basic model of cultural change based on social tipping [9, 18, 19]. It is an  
76 extremely influential model of how social planners can recruit cultural evolutionary processes  
77 to promote behaviour change. The idea has appeared, in one form or another, across a range  
78 of policy-relevant domains [1] related to gender-based violence [16, 20–25] and other forms of  
79 gender bias [26–28], natural resource use [29], health [30, 31], species conservation [32], and  
80 climate change [33–37].

81 If everyone is the same, the social planner’s task is relatively straightforward. She needs  
82 to know how big the initial shock must be. Put differently, she needs to know what proportion  
83 of people must change for conformity and coordination to switch from reinforcing the status  
84 quo norm to reinforcing the social planner’s preferred alternative. The trouble is that people  
85 are usually not all the same [38], and ordinary forms of heterogeneity introduce a number of  
86 challenges and complexities [1, 9, 14, 39–42]. Bare minimum, the social planner must ask both  
87 how big her intervention should be and which segment of the population to target with the  
88 intervention. Interestingly, the best answer to this second question can vary, but targeting  
89 the individuals most amenable to change is often the worst strategy if the objective is to  
90 maximise behaviour change [1, 14]. Here we ask a related but different question. Namely,  
91 if the social planner is considering two different intervention targets, which one maximises  
92 social welfare? Surprisingly, this question can lead to very different conclusions. Targeting  
93 the most amenable segment of the population can actually limit behaviour change, with  
94 frequent miscoordination the outcome, but it can lead to the greatest social welfare. As we  
95 will see, this paradoxical outcome readily occurs in situations, perhaps typical, where some  
96 people want society to transition to a new norm, but others do not.

## 97 **2 Model and results**

98 Assume an infinitely large population of individuals,  $i \in I$ , where  $I$  is some uncountably  
99 infinite indexing set. Everyone is playing a game with two possible choices, SQ and Alt.  
100 Individuals pair off randomly to play the game in a periodic fashion with random rematching  
101 every period. For reasons explained below, we think of SQ as the “status quo”, namely the  
102 behaviour everyone chooses before intervention. The social planner does not like everyone  
103 playing SQ. She would prefer that everyone switch to choosing Alt, and thus at some point  
104 the social planner implements an intervention that promotes Alt as an “alternative” to the

105 status quo. Importantly, we might normally imagine that the social planner wants people  
 106 to choose Alt because Alt is somehow socially beneficial. While a natural interpretation, we  
 107 do not assume or insist on this idea. The social planner might instead simply have her own  
 108 preferences that differ from those of the people. This possibility has proven important in  
 109 discussions about certain cultural traditions like female genital cutting and early marriage  
 110 [15–17, 21]. One view of programmes promoting the abandonment of cutting, for example,  
 111 is that they help the people in a cutting society help themselves by shifting them towards a  
 112 non-cutting equilibrium that respects human rights and improves outcomes for all involved.  
 113 Another view is that such programmes are a form of cultural imperialism, with Europeans  
 114 and their proxies once again imposing European values on the rest of the world [1, 14]. We  
 115 will not venture a resolution of this kind of dispute. We do, however, respect the validity of  
 116 such disputes by allowing the possibility that the social planner’s desired outcome may not  
 117 be socially beneficial. This possibility, in fact, lies at the centre of our claim that a social  
 118 planner’s intervention may induce norm change that hurts.

119 In any case, before intervention, the game is a strict coordination game for everyone, and  
 120 individuals have heterogeneous preferences (Table 1, Pre-intervention (All)). The game is a  
 121 strict coordination game for everyone because we assume that, for each  $i$ ,  $a + x_i > a$ , and  
 122  $b + x_i < d$ , where  $b < d$ . Intuitively, everyone faces incentives to match the choices of their  
 123 partners. This shared interest in coordinating, however, mixes with heterogeneous preferences  
 124 in the sense that each individual has her own idiosyncratic  $x_i$  value. Across individuals,  $x_i$   
 125 values are somehow distributed on the interval  $(0, d - b)$  according to the density function  $f$   
 126 and its associated cumulative probability function  $F$ .

127 Because everyone faces incentives to coordinate, a focal player’s beliefs about what her  
 128 next partner will play are important, and we can summarise an individual’s preferences as  
 129 an indifference point defined in terms of beliefs. Let  $\tilde{q}_i$  be  $i$ ’s belief that her next randomly  
 130 selected partner will play Alt. The expected payoff (Table 1, Pre-intervention (All)) from  
 131 choosing SQ,  $E[\Pi_i(\text{SQ})]$ , and the expected payoff from choosing Alt,  $E[\Pi_i(\text{Alt})]$ , are the  
 132 following,

$$\begin{aligned}
 E[\Pi_i(\text{SQ})] &= (1 - \tilde{q}_i)(a + x_i) + \tilde{q}_i(b + x_i) \\
 E[\Pi_i(\text{Alt})] &= (1 - \tilde{q}_i)(a) + \tilde{q}_i(d).
 \end{aligned}
 \tag{1}$$

133 The individual is indifferent between the two choice options if  $\tilde{q}_i = x_i/(d-b)$ . If  $\tilde{q}_i > x_i/(d-b)$ ,  
 134 the individual prefers to choose Alt. If  $\tilde{q}_i < x_i/(d-b)$ , she prefers SQ.

135 As a kind of reference model, imagine that, for each  $i$ , the belief in the current period

136 is simply the actual distribution of choices from the previous period, and each  $i$  chooses the  
 137 behaviour with the highest expected payoff given this belief. Imagine further that individuals  
 138 choose Alt when indifferent. This model is sometimes called the “threshold” model [14,  
 139 39, 43]. The model supports at least two interpretations. First, individuals form beliefs  
 140 myopically, namely by simply extrapolating from the recent past, and they choose the best  
 141 option given these beliefs [44]. Second, all individuals are conformists in the sense specified  
 142 by Boyd and Richerson [11], but conformist social learning mixes with content biases that are  
 143 heterogeneous across individuals [1]. We motivate our analysis here with the coordination  
 144 game interpretation, but this is not essential. The two models are isomorphic [1], and our  
 145 analysis would remain the same were we to rely on the conformity interpretation. Regardless  
 146 of interpretation, the model is extraordinarily good at predicting behaviour in experimental  
 147 studies with coordination games [19, 44]. Thus, although we do not limit our analyses by  
 148 assuming that social dynamics unfold according to this model, we sometimes highlight the  
 149 steady states of the model as a point of reference.

150 Notice that, if  $x_i < (d - b)/2$ , the set of beliefs for which  $i$  chooses Alt is larger than the  
 151 set of beliefs for which  $i$  chooses SQ. This is one way of saying that  $i$  prefers Alt over SQ.  
 152 More technically, we will say that, for such an individual, coordinating on Alt risk dominates  
 153 coordinating on SQ [45]. If  $x_i > (d - b)/2$ , the opposite holds, and coordinating on SQ risk  
 154 dominates coordinating on Alt.

155 Empirical research has shown that, without special countervailing mechanisms in place,  
 156 risk dominance exerts an extraordinary pull on cultural evolutionary dynamics [46–52]. What  
 157 would this mean in a heterogeneous population? As others have argued [27], we assume the  
 158 population is most likely to converge on the equilibrium that a majority of individuals view  
 159 as risk-dominant. By extension, the situation of interest for our analysis is one in which  
 160  $F((d - b)/2) < 0.5$ . Specifically,  $F((d - b)/2) < 0.5 \Rightarrow 1 - F((d - b)/2) > 0.5$ , where the  
 161 second condition means that before intervention a majority of individuals view coordinating  
 162 on SQ as risk-dominant. In this case, we expect the population to converge on SQ. If the  
 163 distribution of preferences had been otherwise, the population would have probably converged  
 164 on Alt, and the social planner would have had no need to intervene in the first place [1, 27].  
 165 If most individuals view coordinating on SQ as risk-dominant, the distribution of  $x_i$  values  
 166 should be left-skewed because left skew ensures that  $F((d - b)/2) < 0.5$ .

167 Crucially, although we argue that left-skewed  $x_i$  distributions represent the situations of  
 168 interest, this claim is silent about the welfare consequences of coordinating on SQ versus  
 169 coordinating on Alt. Specifically, we have made no claims so far about the relation between

170 the  $a + x_i$  and  $d$ , where the former are the payoffs players get when coordinating on SQ  
171 and the latter the payoff they get when coordinating on Alt. Many possibilities exist. At  
172 one extreme, for each  $i$ ,  $a + x_i > d$ , which would mean that coordinating on SQ is socially  
173 beneficial in the precise sense that coordinating on SQ is strictly better for everyone than  
174 coordinating on Alt. A social planner who intervenes in this situation is simply promoting  
175 her own agenda, as discussed above, to the detriment of the people. At the other extreme,  
176 for each  $i$ ,  $a + x_i < d$ , which would mean that coordinating on SQ is socially harmful because  
177 coordinating on SQ is strictly worse for everyone than coordinating on Alt. Both of these  
178 extremes are consistent with saying that coordinating on SQ is risk-dominant for a majority  
179 of individuals. Risk dominance depends on the relation between  $x_i$  and  $d - b$ , not the relation  
180 between  $a + x_i$  and  $d$ . Situations between the two extremes are also possible, and such  
181 situations figure prominently in our analyses below.

182 At some point, the social planner rolls out an intervention to promote behaviour change.  
183 She targets (T) some proportion,  $\phi \in (0, 1)$ , of the population and incentivises these people  
184 to switch from SQ to Alt (Table 1, Post-intervention (T)). The intervention is unequivocally  
185 effective in the sense that  $h > g$ , and thus post-intervention all targeted individuals always  
186 choose Alt. This is a strong assumption. It implies that all targeted individuals, regardless  
187 of their initial preferences, effectively acquire new preferences because of their experience  
188 with the intervention. One natural interpretation is that the intervention in question is  
189 an extremely effective persuasion campaign that instils new preferences based on strong  
190 personal values. Imagine an individual who becomes extremely well-informed about climate  
191 change. Green choices like riding a bike and avoiding beef become intrinsically valuable  
192 to this individual ( $h$ ), and brown choices like driving an SUV and eating steaks become  
193 intrinsically painful ( $g$ ). These intrinsic values dominate decision making in the sense that  
194 they are far more important than whether or not the individual manages to coordinate with  
195 others. Later, we relax this assumption with simulation models that assume individuals are  
196 less likely to respond to the intervention in this way as they become more resistant to change  
197 pre-intervention [14, 24].

198 Before that, however, we consider the model in which the intervention leads any targeted  
199 (T) individual, regardless of her pre-intervention  $x_i$ , to change behaviour. The question is,  
200 what do the non-targeted individuals (NT) do? This question lies at the centre of applied  
201 cultural evolution in general and norm change based on social tipping specifically. If endoge-  
202 nous cultural evolutionary processes lead to behaviour change among people having no direct  
203 experience with the intervention, the social planner has activated these processes, whether

204 intentionally or not, to amplify the direct effect of her intervention.

205 To analyse long-run welfare, we consider two alternative intervention strategies. Although  
206 many more possibilities exist, we analyse the two extremes because they bracket the range  
207 of possibilities and intuitively capture the trade-offs the social planner faces [14]. At one  
208 extreme, the social planner targets the segment of the population most amenable to change.  
209 Specifically,  $\exists x_A^1 \in (0, d - b)$  such that  $F(x_A^1) = \phi$ , and the social planner targets everyone  
210 with an  $x_i \leq x_A^1$ . At the other extreme, the social planner targets the segment of the  
211 population most resistant to change. In this case,  $\exists x_R^1 \in (0, d - b)$  such that  $F(x_R^1) = 1 - \phi$ ,  
212 and the social planner targets everyone with an  $x_i > x_R^1$ .

213 After intervention, assume the population stabilises in the long-run on some proportion  
214 choosing Alt. Under an amenable target, we denote this proportion as  $\hat{q}_A$ . Under a resistant  
215 target, we use  $\hat{q}_R$ . Because targeted individuals always choose Alt post-intervention, then  
216  $\hat{q}_A, \hat{q}_R \geq \phi$ . Because of one additional assumption,  $\hat{q}_R \geq \hat{q}_A$  must also hold. Specifically,  
217 of the non-targeted individuals who change from SQ to Alt, we assume they do so in order  
218 from those most amenable towards Alt to those most resistant. Put differently, non-targeted  
219 individuals who change from SQ to Alt do so in order from those with the smallest  $x_i$  values  
220 to those with the largest  $x_i$  values. This assumption is consistent with the threshold model  
221 [14, 39], but more intuitively it simply means that people who are relatively favourable  
222 towards Alt choose Alt at least as early as those who are relatively unfavourable towards Alt.  
223 With this assumption in place, as long as targeted individuals respond the same regardless  
224 of their pre-existing preferences, one can show that  $\hat{q}_R \geq \hat{q}_A$  must hold [14].

225 The intuition is the following. Under an amenable target, by targeting the most amenable  
226 segment of the population, the social planner chooses the easiest possible task for the interven-  
227 tion. This leaves the hardest possible task for subsequent endogenous cultural evolutionary  
228 processes because the non-targeted individuals necessarily comprise a proportion  $1 - \phi$  of  
229 the population as resistant to change as possible. Under a resistant target, in contrast, the  
230 social planner takes the hardest possible task for the intervention, but this does not matter  
231 because we are assuming the intervention is equally effective regardless of the target. A resis-  
232 tant target also leaves the easiest possible task for endogenous cultural evolution because the  
233 non-targeted individuals make up a proportion  $1 - \phi$  as amenable to change as possible. For  
234 this reason, the long-run proportion choosing Alt must be at least as large under a resistant  
235 target as under an amenable target [14].

236 The result of the intervention is a partition of the population (Fig. 1) into either two or  
237 three categories of player. First, targeted players choose Alt, which we designate with (Alt,T),



238 and this category always exists. Second, some or all non-targeted players may choose Alt, a  
 239 category denoted as (Alt,NT). Finally, some or all non-targeted players may stick with SQ,  
 240 a category denoted as (SQ,NT). Our focus is the welfare effects associated with the long-run  
 241 partition of the population. Crucially, the partition and its welfare consequences depend on  
 242 the social planner's intervention strategy.

243 Under an amenable target, targeted individuals constitute a proportion  $\phi$  of the popula-  
 244 tion from the left tail of the  $x_i$  distribution. Specifically,  $F(x_A^1) = \phi$ , and the targeted subset  
 245 thus consists of all individuals with  $x_i$  values in the interval  $(0, x_A^1]$ . These individuals form  
 246 the (Alt,T) category. If additional individuals choose Alt, then  $\hat{q}_A > \phi$ . These individuals  
 247 form the (Alt,NT) category, which make up a proportion  $\hat{q}_a - \phi$  of the population. Because  
 248 non-targeted individuals choose Alt in order from those most amenable to most resistant,  
 249 the most amenable individuals in this category come from somewhere in the middle of the  $x_i$   
 250 distribution. More precisely,  $\exists x_A^2 > x_A^1$  such that  $F(x_A^2) = \hat{q}_A$ . The subset of non-targeted  
 251 individuals who choose Alt thus consists of all individuals with  $x_i$  values in the interval  
 252  $(x_A^1, x_A^2]$  if  $x_A^2 < d - b$  and  $(x_A^1, x_A^2)$  if  $x_A^2 = d - b$ . If  $\hat{q}_A = \phi$ , this category of player does  
 253 not exist. Finally, if  $\hat{q}_A < 1$ , some non-targeted individuals persist in choosing SQ. These  
 254 individuals necessarily come from the right tail of the  $x_i$  distribution. Specifically, individuals  
 255 with  $x_i$  values in the interval  $(x_A^2, d - b)$  form this (SQ,NT) category. If  $\hat{q}_A = 1$ , this category  
 256 does not exist.

257 Under a resistant target, we start at the right tail of the  $x_i$  distribution. Given an  $x_R^1$   
 258 such that  $F(x_R^1) = 1 - \phi$ , the (Alt,T) individuals are those individuals with  $x_i$  values in  
 259 the interval  $(x_R^1, d - b)$ . If  $\hat{q}_R > \phi$ , the (Alt,NT) individuals come from the left tail of the  
 260 distribution. Specifically,  $\exists x_R^2 > 0$  such that  $F(x_R^2) = \hat{q}_R - \phi$ . These (Alt,NT) individuals  
 261 consist of everyone with  $x_i$  values in the interval  $(0, x_R^2]$ . Finally, if  $\hat{q}_R < 1$ , some individuals  
 262 stick with SQ in the long run, and these individuals have  $x_i$  values in the interval  $(x_R^2, x_R^1]$ ,  
 263 which is somewhere in the middle of the  $x_i$  distribution. Fig. 1 shows an example of how an  
 264 amenable versus a resistant target might induce two distinct partitions.

265 To analyse the welfare consequences of our two intervention strategies, we imagine ran-  
 266 domly sampling an individual in a period and calculating this individual's expected payoff.  
 267 Because we simply focus on the expected payoff of a randomly selected individual, we are  
 268 adopting a utilitarian view of social welfare. Intuitively, we are assuming that first and fore-  
 269 most the social planner has an obligation to maximise the aggregate payoffs in society. If  
 270 she has other objectives, like reducing inequality for example, she can redistribute ex post to  
 271 accomplish these objectives. We relax the focus on payoffs when we develop a more elaborate

272 model below.

273 To assist with the logic behind calculating an individual's expected payoff from one period  
 274 of play, Table 2 shows all the ways in which two players can pair off to play, the probabilities of  
 275 the various pairs, and the payoffs generated for each player in a pair. For convenience, define  
 276  $V_A = (x_A^2, d - b)$ , which is the set of  $x_i$  values for (SQ,NT) individuals under an amenable  
 277 target. Analogously, let  $V_R = (x_R^2, x_R^1]$  be the set of  $x_i$  values for (SQ,NT) individuals under  
 278 a resistant target.

279 Under an amenable target, the expected payoff,  $E_A[\Pi_i]$ , takes the form,

$$\begin{aligned}
 E_A[\Pi_i] = & \frac{\phi^2(2h)}{2} + \frac{\phi(\hat{q}_A - \phi)(h + d)}{2} + \frac{\phi(1 - \hat{q}_A)(h + b + E[X_i | x_i \in V_A])}{2} \\
 & + \frac{(\hat{q}_A - \phi)\phi(d + h)}{2} + \frac{(\hat{q}_A - \phi)^2(2d)}{2} + \frac{(\hat{q}_A - \phi)(1 - \hat{q}_A)(a + b + E[X_i | x_i \in V_A])}{2} \\
 & + \frac{(1 - \hat{q}_A)\phi(b + E[X_i | x_i \in V_A] + h)}{2} + \frac{(1 - \hat{q}_A)(\hat{q}_A - \phi)(b + E[X_i | x_i \in V_A] + a)}{2} \\
 & + \frac{(1 - \hat{q}_A)^2(2a + 2E[X_i | x_i \in V_A])}{2}.
 \end{aligned} \tag{2}$$

280 This expression simplifies to

$$\begin{aligned}
 E_A[\Pi_i] = & \phi h + a(1 - \phi)(1 - \hat{q}_A) + b\hat{q}_A(1 - \hat{q}_A) \\
 & + d\hat{q}_A(\hat{q}_A - \phi) + (1 - \hat{q}_A)E[X_i | x_i \in V_A].
 \end{aligned} \tag{3}$$

281 The expected payoff under a resistant target,  $E_R[\Pi_i]$ , is analogous,

$$\begin{aligned}
 E_R[\Pi_i] = & \phi h + a(1 - \phi)(1 - \hat{q}_R) + b\hat{q}_R(1 - \hat{q}_R) \\
 & + d\hat{q}_R(\hat{q}_R - \phi) + (1 - \hat{q}_R)E[X_i | x_i \in V_R].
 \end{aligned} \tag{4}$$

282  $E_A[\Pi_i]$  and  $E_R[\Pi_i]$  look similar, but recall that  $\hat{q}_A$  and  $\hat{q}_R$  can be different. More subtly, the  
 283 terms  $E[X_i | x_i \in V_A]$  and  $E[X_i | x_i \in V_R]$  are conditional expectations over different parts of  
 284 the  $x_i$  distribution. Specifically, with  $\mu$  as the Lebesgue measure,

$$\begin{aligned}
 E[X_i | x_i \in V_A] &= \frac{\int_{V_A} x_i f(x_i) d\mu(x_i)}{\int_{V_A} f(x_i) d\mu(x_i)} \\
 E[X_i | x_i \in V_R] &= \frac{\int_{V_R} x_i f(x_i) d\mu(x_i)}{\int_{V_R} f(x_i) d\mu(x_i)}.
 \end{aligned} \tag{5}$$

285  $E[X_i | x_i \in V_A]$  is the expected  $x_i$  value of (SQ,NT) individuals under an amenable target.

286 If this category exists ( $\hat{q}_A < 1$ ), these individuals will come from the right tail of the  $x_i$

287 distribution and thus be relatively resistant to Alt. For this reason,  $E[X_i | x_i \in V_A]$  will tend  
 288 to be large. In contrast,  $E[X_i | x_i \in V_R]$  is the expected  $x_i$  value of (SQ,NT) individuals under  
 289 a resistant target. If this category exists ( $\hat{q}_R < 1$ ), these individuals will come from the middle  
 290 of the  $x_i$  distribution, and so they will only be moderately resistant to Alt. Consequently,  
 291  $E[X_i | x_i \in V_R]$  will tend take intermediate values.

292 To see which intervention strategy yields the highest expected payoff, subtract one from  
 293 the other,

$$\begin{aligned}
 E_A[\Pi] - E_R[\Pi] &= (1 - \hat{q}_A)E[X_i | x_i \in V_A] - (1 - \hat{q}_R)E[X_i | x_i \in V_R] \\
 &\quad - (\hat{q}_R - \hat{q}_A) \{b + (\hat{q}_A + \hat{q}_R)(d - b) - (1 - \phi)a - \phi d\}.
 \end{aligned}
 \tag{6}$$

294 Note that  $h$  and  $g$  disappear. This happens because we assume the intervention is equally ef-  
 295 fective when targeting resistant versus amenable individuals. The welfare differences between  
 296 the two interventions strategies thus depend exclusively on effects related to non-targeted indi-  
 297 viduals. To gain some intuition about what equation (6) represents, assume the two different  
 298 intervention strategies produce outcomes as different from each other as possible. Specifically,  
 299 let  $\hat{q}_A = \phi$  and  $\hat{q}_R = 1$ . This means that using the intervention to target amenable people  
 300 does not induce any behaviour change via endogenous cultural evolution; the only effect is  
 301 the direct effect of the intervention. In contrast, targeting people resistant to change induces  
 302 the maximum possible change; everyone eventually chooses Alt. Moreover, the rate of mis-  
 303 coordination is relatively high under the amenable target at  $2\phi(1 - \phi)$ , but miscoordination  
 304 never occurs under the resistant target. For this special case,  $E_A[\Pi] - E_R[\Pi] > 0$  if and only  
 305 if the following holds,

$$(1 - \phi)(a + E[X_i | x_i \in V_A]) + \phi(b + E[X_i | x_i \in V_A]) > d.
 \tag{7}$$

306 Condition (7) compares the expected payoffs of non-targeted individuals given the two inter-  
 307 vention strategies. Under a resistant target, all non-targeted players choose Alt. They always  
 308 coordinate, and they always get a payoff of  $d$ . Under an amenable target, all non-targeted  
 309 players choose SQ. A randomly selected player of this type is paired with another non-targeted  
 310 player with probability  $1 - \phi$ . They coordinate, and the expected payoff of the focal player  
 311 is  $a + E[X_i | x_i \in V_A]$ . The focal player is paired with a targeted player with probability  $\phi$ .  
 312 They miscoordinate, and the expected payoff of the focal player is  $b + E[X_i | x_i \in V_A]$ .

313 Because the  $x_i$  are distributed on  $(0, d - b)$ ,  $b + E[X_i | x_i \in V_A] < b + d - b = d$  must be  
 314 true. However,  $a + E[X_i | x_i \in V_A] > d$  is certainly possible, and condition (7) is also possible.

315 The outcome depends on both  $a$  and the distribution of  $x_i$  values among the non-targeted  
316 individuals given an amenable target. In this sense,  $a + E[X_i | x_i \in V_A]$  is a measure of the  
317 alignment or misalignment between the social planner's objectives and the preferences of a  
318 key segment of the population, namely the individuals the social planner does not target  
319 when she chooses an amenable target.

320 At one extreme, alignment is high, and condition (7) does not hold regardless of the  $x_i$   
321 distribution. Specifically, note that  $\sup V_A = d - b$ . Substituting this value for  $E[X_i | x_i \in V_A]$   
322 and rearranging reveals that condition (7) does not hold if  $a \leq b$ . In this special case, the  
323 population fixed on SQ before intervention really is stuck in a harmful equilibrium. Behaviour  
324 change is unambiguously good in the sense that a resistant target, which maximises behaviour  
325 change, produces larger expected payoffs than an amenable target, which minimises behaviour  
326 change, whatever form preference heterogeneity takes.

327 At the other extreme, misalignment is high, and condition (7) holds under extreme condi-  
328 tions. Specifically, when  $\hat{q}_A = \phi$ , note that  $\inf V_A = x_A^1$ , where  $F(x_A^1) = \phi$ . Substituting this  
329 value and rearranging shows that condition (7) holds if and only if  $a > (d - \phi b - x_A^1)/(1 - \phi)$ .  
330 Note that this condition involves  $x_A^1$ , which depends on the  $x_i$  distribution given  $\phi$ . For this  
331 reason, we do not isolate a situation in which condition (7) holds regardless of the entire  $x_i$   
332 distribution. Instead, in this special case, we focus on an extreme situation in which con-  
333 dition (7) holds even if the distribution of  $x_i$  values among non-targeted players minimises  
334  $a + E[X_i | x_i \in V_A]$ . In this special case, behaviour change is unambiguously harmful for all  
335 non-targeted individuals given an amenable target. Consequently, an amenable target is best  
336 precisely because non-targeted individuals do not change behaviour.

337 Between these two extremes, the distribution of preferences among non-targeted individ-  
338 uals, given an amenable target, is important. If  $a$  is sufficiently large, and if the distribution  
339 of  $x_i$  values is sufficiently left-skewed,  $a + E[X_i | x_i \in V_A]$  can be large enough to ensure  
340 that condition (7) is satisfied. In this case, even if some people clearly benefit from changing  
341 behaviours from SQ to Alt, others do not. Indeed, some individuals do best by sticking  
342 with SQ and tolerating frequent miscoordination. These individuals are exactly the people  
343 who are not targeted under an amenable target, they are exactly the people who do not  
344 change behaviour under an amenable target, and they are exactly the people who ensure  
345 that an amenable target with frequent miscoordination is better for society than a resistant  
346 target with no miscoordination. In this situation, the social planner is at odds with the  
347 most resistant segment of the population under her influence. She leaves these people out of  
348 her intervention given an amenable target, and they maintain their pre-existing preferences

349 as a result. These preferences favour SQ to such an extent that coordinating on SQ with  
 350 other non-targeted individuals more than compensates for miscoordinating when paired with  
 351 targeted individuals. A resistant target could generate genuine norm change, with everyone  
 352 coordinating on Alt, but it would do real harm relative to an amenable target with chronic  
 353 disagreement and persistent miscoordination.

354 More broadly, we can make further progress if we choose specific parameter values and use  
 355 a graphical approach to the unrestricted condition (Eq. 6). Accordingly, Figs. 2 – 4 summarise  
 356 which intervention strategy produces the highest expected payoff (Eq. 6) under two different  
 357  $x_i$  distributions, two different values of  $a$ , and three different values of  $\phi$ . Because these  
 358 figures work with the unrestricted condition (Eq. 6), they show the relative welfare effects of  
 359 the alternative intervention strategies for any possible outcome subject to  $\phi \leq \hat{q}_A \leq \hat{q}_R \leq 1$ .  
 360 In this way, we consider a wide range of steady states, and thus we do not limit attention to a  
 361 particular dynamical process. To create these figures, we set  $b = 0$  and  $d = 1$ . We then vary  
 362  $a \in \{0.25, 0.75\}$ ,  $\phi \in \{0.25, 0.5, 0.75\}$ , and the skewness of the  $x_i$  distribution. The figures  
 363 reveal that, when the social planner and the people have partially misaligned preferences,  
 364 an amenable target often yields a society with higher payoffs than a resistant target. More  
 365 interestingly, this result frequently obtains even though the amenable target leads to less  
 366 behaviour change and more miscoordination than the resistant target. This happens under  
 367 an increasingly broad range of conditions as  $a$  increases and as the skew of the  $x_i$  distribution  
 368 increases. In other words, it happens as the disconnect between the social planner and the  
 369 resistant segment of society becomes more pronounced.

370 To see this result, note that for all figures the red region expands as we move from left  
 371 (a,c) to right (b,d) and from top (a,b) to bottom (c,d). Red represents a steady state in  
 372 which an amenable target produces greater welfare than a resistant target. The move from  
 373 left to right means the left skew of the  $x_i$  distribution increases, and the move from top to  
 374 bottom means  $a$  increases. The expansion of the red region is especially important along the  
 375 top boundaries of each panel because this represents outcomes for which a resistant target  
 376 leads to full-fledged norm change (i.e.  $\hat{q}_R \approx 1$ ), but an amenable target does not (i.e.  $\hat{q}_A < 1$ ).

377 Thus far, we have compared the payoffs under two alternative intervention strategies. We  
 378 have shown that an intervention strategy that produces a complete shift in the population  
 379 from one equilibrium to another, with little or no miscoordination the result, can actually be  
 380 worse in terms of social welfare than an alternative strategy that leads to no norm at all after  
 381 intervention. What if, in contrast, we hold the intervention strategy constant? Specifically,  
 382 what if we hold the size (i.e.  $\phi$ ) and target (i.e. amenable or resistant) constant and allow

383 the steady state to vary? Intuitively, one might imagine that with such comparisons, where  
384 all else is equal, more coordination would be better than less coordination. Intuition is once  
385 again, however, potentially misleading. Fig. 5 shows example results.

386 The results reveal that, when at least some individuals prefer to coordinate on SQ as  
387 opposed to Alt, expected payoffs can actually increase with miscoordination all else equal.  
388 An amenable target is especially prone to this result because an amenable target leaves  
389 the most resistant individuals to persist in choosing SQ. These individuals tolerate frequent  
390 miscoordination because they get especially high payoffs when they manage to coordinate on  
391 SQ, and this effect can drive up the expected payoff for the entire population. As  $a$  increases  
392 (Fig. 5c, d), and as the skew of the  $x_i$  distribution increases (Fig. 5b, d), expected payoffs  
393 can rise with miscoordination rates for the same basic reason. Of particular note, amenable  
394 versus resistant targets do not necessarily lead to the same patterns of variation in expected  
395 payoffs. Expected payoffs can increase as miscoordination rises under an amenable target but  
396 decrease as miscoordination rises under a resistant target (Fig. 5c,  $\phi \in \{0.5, 0.75\}$ ; Fig. 5d,  
397  $\phi = 0.75$ ). As a result, even if we assume the two intervention strategies generate the same  
398 behaviour change (i.e.  $\hat{q}_A = \hat{q}_R$ ), miscoordination may be good for social welfare under an  
399 amenable target but bad for social welfare under a resistant target. Equally challenging for  
400 the social planner, the effects of miscoordination are not even reliably monotonic given an  
401 intervention strategy. The solid red line of Fig. 5b provides an example. As we move from  
402  $\hat{q}_A = 1$  to  $\hat{q}_A = 0.5$ , expected payoffs first decrease and then increase. This means that, given  
403 an amenable target with  $\phi = 0.25$ , increasing miscoordination is first bad for society, but  
404 then it becomes good. This subtlety can be relevant when the social planner has committed  
405 to an intervention strategy, but she is uncertain about the final outcome. Because of the  
406 uncertainty, she will not be able to say ex ante if social welfare increases or decreases as  
407 coordination rises.

### 408 **3 A stylised beef-eating illustration**

409 To clarify the intuition, imagine a population of beef eaters and a social planner who would like  
410 everyone to switch to plant-based alternatives for reasons related to both public health and  
411 climate change. Even though everyone is a regular beef eater before intervention, say because  
412 beef-eating is the local culture, people's tastes naturally vary. Some people have tastes that  
413 favour coordinating on fruits, vegetables, and beans over yet another meal with blood on the  
414 plate and a brick in the stomach. We can call these people the "berry lovers". Other people

415 would prefer to coordinate on stew, cheese, and Waygu instead of fruits, vegetables, and  
416 beans. We can call these people the “steak lovers”. Between these two groups we have the  
417 “omnivores”, a group of people who enjoy sharing a steak dinner with friends just as much as  
418 they enjoy sharing a meal of salad and strawberries. The berry lovers constitute a minority  
419 of the population. Frequencies rise as we move into the omnivores. Steak lovers form the  
420 majority, which is exactly why the population converged on beef eating before intervention.  
421 In terms of the model above, this stylised population could have  $d = 1$ ,  $b = 0$ , and  $a = 0.5$ ,  
422 with a distribution of  $x_i$  values that covers the full support,  $(0, d - b)$ , but is left-skewed. The  
423 berry lovers have  $x_i$  values noticeably less than 0.5, the steak lovers have  $x_i$  values noticeably  
424 greater than 0.5, and the omnivores have  $x_i$  values close to 0.5.

425 Now consider a social planner who implements an intervention of size  $\phi = 0.5$  in either the  
426 amenable tail of the preference distribution or the resistant tail. Imagine further that, under  
427 an amenable target, only the targeted individuals switch to a plant-based diet ( $\hat{q}_A = \phi$ ), but  
428 the entire population switches under the resistant target ( $\hat{q}_R = 1$ ). How can the former lead  
429 to larger expected payoffs than the latter? Because everyone responds to the intervention in  
430 the same way when targeted, the answer depends on what happens among the individuals  
431 who are not targeted.

432 Because steak lovers are common, the social planner targeting the amenable tail has to  
433 target the berry lovers, the omnivores, and maybe even a few half-hearted steak lovers to accu-  
434 mulate 50% of the population for her intervention. The remaining non-targeted 50% consists  
435 of serious steak lovers who simply continue to choose steak. Because this half of the popula-  
436 tion has extreme preferences, they can tolerate some miscoordination as long as they get to  
437 have their preferred food. Sometimes they have a steak at the table with someone enjoying  
438 plant-based alternatives, and perhaps their berry-loving companions chastise them along the  
439 way, but at least they get to eat steak. In particular, their payoffs are only slightly less than  
440 the coordination payoffs they would have received if they had joined their tablemates and  
441 chosen a plant-based meal. That said, when they do coordinate on steak with another serious  
442 steak lover, payoffs are especially high all around. These high payoffs follow from the extreme  
443 preferences of the non-targeted group, given an amenable target, and the high payoffs more  
444 than compensate for the small miscoordination costs that sometimes occur. For exactly this  
445 reason, non-targeted individuals persist in choosing steak instead of switching to plant-based  
446 alternatives. Their decision to do so has important welfare consequences, and in particular  
447 their extreme preferences ensure that choosing steak and occasionally miscoordinating may  
448 generate much higher expected payoffs than switching to plant-based alternatives.

449 In contrast, when the social planner targets the resistant tail, she targets the 50% of the  
450 population composed of the serious steak lovers. The remaining non-targeted 50% consists of  
451 the berry lovers, the omnivores, and few half-hearted steak lovers. When they all switch to  
452 coordinating on plant-based alternatives, the berry lovers experience a gain, but they are a  
453 small group. The omnivores are, by definition, approximately indifferent over the equilibria,  
454 and so they switch from one norm to another norm that is more or less just as good. The  
455 half-hearted steak lovers experience a loss, but they are also a small part of the non-targeted  
456 group. All in all, many non-targeted individuals have moderate preferences. Coordinating is  
457 important; coordinating on a specific behaviour is not. The non-targeted change behaviour  
458 for this reason, but their decision to do so has only moderate welfare effects precisely because  
459 many of them have moderate preferences.

460 As a crucial caveat, the example outlined here is deliberately vague about externalities.  
461 One interpretation is that externalities are not present; the payoff matrix for each individual  
462 captures the full suite of consequences associated with the choices the individual and her  
463 partner can make. If true, the analysis above holds without complication. An intervention  
464 among the resistant may produce a complete shift to a plant-based norm, but an interven-  
465 tion in the amenable tail would have produced chronic miscoordination with higher average  
466 payoffs. If externalities are present, however, the social planner may be justified in choos-  
467 ing an intervention strategy that generates complete norm change even if she knows this  
468 strategy will produce lower perceived payoffs than other intervention strategies. The social  
469 planner would be justified with such an approach, for example, if beef is underpriced, which  
470 is almost certainly true, because the price does not account for all climatic effects associated  
471 with raising cattle instead of plants. As another way to think about this, if the price of beef  
472 was correct, the value of  $a$  would be lower than it is, and the entire distribution of  $a + x_i$   
473 values would shift downward. Most people then would actually perceive a complete shift to  
474 a plant-based norm as beneficial, and this would be true even for many people who really  
475 enjoy steak.

476 Efforts to shift norms related to cultural traditions like female genital cutting and early  
477 marriage [1, 15, 16, 21] also raise critical questions about the extent to which externalities are  
478 present. A viewpoint emphasising cultural relativism might argue that families have exactly  
479 the preferences they should have within the context of cultural traditions that value female  
480 genital cutting or early marriage. An outsider may not understand these preferences, but  
481 this is no reason to discount their legitimacy. In this case, norm change can actually hurt  
482 relative to alternative social planning strategies. As always, details related to the distribution



483 of  $a + x_i$  values and the set of social planning strategies under consideration are crucial. In  
484 stark contrast, a viewpoint emphasising universal human rights might argue that cutting and  
485 early marriage perpetuate cultural systems that devalue girls and women, with a host of  
486 attendant social costs. By this account, externalities are ubiquitous, and the social planner  
487 should promote norm change even if she knows families will perceive themselves as worse off  
488 than they would under different intervention strategies.

## 489 4 Heterogeneous response to intervention

490 The analytical model above (§ 2) makes the strong assumption that the intervention is equally  
491 effective regardless of how amenable or resistant a targeted individual is before intervention.  
492 In practice, however, interventions designed to change behaviour often have heterogeneous  
493 effects [53]. To account for this possibility, we developed an agent-based simulation (Sup-  
494plementary Information) that allows heterogeneous responses to the intervention. These  
495 simulations also allow us to show transient dynamics, the inequalities an intervention pro-  
496 duces, and any stochastic effects that might occur because the population is finite. First, we  
497 explain the generic structure of the model. Then we explain the parameter space we used  
498 and key state variables we recorded when simulating.

499 To begin, a simulation creates a population of  $N = 1000$  agents by drawing an  $x_i$  value for  
500 each agent from a left-skewed beta distribution. In time  $t = 1$ , everyone chooses SQ. Agents  
501 pair off randomly and play coordination games based on their pre-intervention preferences  
502 (Table 1, Pre-intervention (All)). Because everyone plays SQ, everyone coordinates and  
503 receives a payoff of  $a + x_i$ . Between  $t = 1$  and  $t = 2$ , the simulation implements an intervention  
504 of size  $\phi$  in either the amenable (A) or resistant (R) tail of the  $x_i$  distribution. Targeted  
505 individuals respond to the intervention with a probability that is a decreasing function of  
506 their  $x_i$  values. Specifically, a targeted agent  $i$  responds with probability  $s_i = 1 - x_i/(d - b)$ .  
507 If a targeted  $i$  responds to the intervention, she gets a new payoff matrix (Table 1, Post-  
508 intervention (T)), and given this new payoff matrix she changes from SQ to Alt because  
509  $h > g$ . If a targeted  $i$  does not respond, she retains her original payoff matrix. Non-targeted  
510 agents also retain their original payoff matrices. Agents pair off randomly again and play  
511 coordination games. At this point, only targeted individuals who have responded to the  
512 intervention choose Alt. Agents receive payoffs based on their individual payoff matrices,  
513 their choices, and the choices of their partners.

514 Agents update their beliefs myopically by treating the proportion of others choosing Alt

515 in  $t = 1$  as their beliefs about the probability a randomly selected partner will choose Alt  
516 in  $t = 2$ . Specifically, let  $q_t$  be the proportion choosing Alt in  $t$ . If  $i$  chose SQ in  $t = 1$ , she  
517 believes her partner in  $t = 2$  will choose Alt with probability  $q_1 N / (N - 1)$ . If  $i$  chose Alt  
518 in  $t = 1$ , her belief is instead  $(q_1 N - 1) / (N - 1)$ . Agents again pair off randomly to play  
519 coordination games in  $t = 2$ . When deciding how to play, each agent best responds given her  
520 myopically updated belief. Agents receive payoffs. The algorithm, which consists of random  
521 matching and myopic best responding, repeats until  $t = 100$ .

522 Altogether, we ran simulations over the following parameter space.

- 523 1. For beta-distributed  $x_i$  values, we fixed  $\beta = 2$  and allowed  $\alpha$  to vary according to  
524  $\alpha \in \{2.25, 2.5, 2.75, 3, 8\}$ . Left skew increases as  $\alpha$  increases.
- 525 2. We chose intervention sizes based on  $\phi \in \{0.25, 0.5, 0.75, 0.9\}$ .
- 526 3. Interventions targets were relatively amenable (A) or resistant (R).
- 527 4. For targeted agents who respond to the intervention (Table 1, Post-intervention (T)),  
528 we fixed  $g = 0$  and let  $h$  vary based on  $h \in \{2, 3\}$ .
- 529 5. For all agents pre-intervention, non-targeted agents post-intervention, and targeted  
530 agents who do not respond post-intervention (Table 1), we fixed  $b = 0$ ,  $d = 1$ , and let  
531  $a$  vary according to  $a \in \{0.25, 0.75, 0.9\}$ .

532 When all of these parameter values are fully crossed, the result is 240 unique combinations.  
533 For each combination, we simulated 1000 independent populations. These populations dif-  
534 fered in terms of the realised distribution of  $x_i$  values. For each simulation, in a given  $t$ , we  
535 calculated the fraction of agents in the population choosing Alt, the frequency with which  
536 agents miscoordinated, the average payoff in the population, and the Gini coefficient. The  
537 Gini coefficient [54] is a normalised inequality measure that ranges from zero to one. A value  
538 of zero indicates perfect equality, while a value of one indicates maximum inequality with a  
539 single individual holding all the wealth in the population. Given sorted payoff values in  $t$  for  
540 a population of size  $N$ , i.e.  $\pi_{1,t} \leq \pi_{2,t} \leq \dots \pi_{N,t}$ , we calculated the Gini as

$$G_t = \frac{2 \sum_{i=1}^N i \pi_{i,t}}{N \sum_{i=1}^N \pi_{i,t}} - \frac{N+1}{N}. \quad (8)$$

541 Our approach provided four key quantities for characterising any given population in  $t$ . We  
542 then averaged over our 1000 independently simulated populations to get global averages of our  
543 four key quantities in  $t$ . We further estimated 95% confidence intervals using a non-parametric

544 bootstrap procedure based on resampling from the set of independently simulated populations  
545 (Supplementary Information). We provide our code as Supplementary Information, so the  
546 interested reader can repeat the simulation exercise with different parameter values.

547 Figs. 6 – 9 present some of the key results from our simulations. One key result is that, if  
548 the intervention has heterogeneous effects, an amenable target can produce more behaviour  
549 change than a resistant target. This cannot happen when the intervention has homogeneous  
550 effects and agents myopically best respond [14]. However, if the probability of responding  
551 to the intervention declines with resistance to behaviour change, as in our simulations, the  
552 social planner can expect to face an important trade-off [1, 14, 42]. An amenable target  
553 maximises the direct effect of the intervention, but it minimises the secondary indirect effect.  
554 A resistant target minimises the direct effect, but conditional on a direct effect of a given size  
555 it maximises the secondary indirect effect.

556 Because of this trade-off, a resistant target is no longer guaranteed to maximise the sum  
557 of the direct and indirect effects. In particular, an amenable target tends to produce more  
558 overall behaviour change than a resistant target in situations where neither generates much  
559 behaviour change [42], and our results here are consistent with this idea. In Figs. 6 – 8, for  
560 example, the initial  $x_i$  distributions are only moderately skewed, and total behavioural change  
561 ranges from moderate to complete in the sense that the final proportion choosing Alt ranges  
562 from a bit more than 0.4 (Fig. 8a, Amenable target) to 1.0 (Fig. 6a, Resistant target). In all  
563 of these cases, a resistant target generates more behaviour change than an amenable target.  
564 As the skew of the initial  $x_i$  distribution becomes more extreme, however, behaviour change  
565 becomes very limited in general, with the final proportion choosing Alt remaining below 0.2  
566 Given this limitation, however, an amenable target generates more behaviour change than a  
567 resistant target (Fig. 9a).

568 Crucially, one of our key conclusions from the analytical model above (§ 2) continues to  
569 hold in the more complex setting captured by our simulations. Namely, in heterogeneous  
570 populations, an amenable target readily leads to more miscoordination but higher average  
571 payoffs than a resistant target. In particular, all parameter combinations we considered are  
572 consistent with what we might call a non-degenerate form of heterogeneity in the population.  
573 By non-degenerate we mean that before intervention some agents view coordinating on SQ as  
574 risk-dominant, and some agents view coordinating on Alt as risk-dominant. Similarly, some  
575 agents view coordinating on SQ as payoff-dominant, and some agents view coordinating on  
576 Alt as payoff-dominant. Given these characteristics of the parameter space we consider,  
577 Figs. 6–9 represent a typical set of outcomes.

578 First, an amenable target always produces at least as much miscoordination as a resistant  
579 target and often strictly more. Second, even though amenable targets produce more misco-  
580 ordination than resistant targets, amenable targets tend to produce average payoffs that are  
581 at least as high as resistant targets and often strictly higher. We only identified two cases  
582 in which an amenable target produced slightly lower average payoffs than a resistant target.  
583 For these two exceptions,  $h$  is relatively small ( $h = 2$ ),  $a$  is large ( $a = 0.9$ ), the  $x_i$  distribution  
584 is strongly skewed ( $\alpha = 8$ ), and interventions are relatively small ( $\phi \in \{0.25, 0.5\}$ ). In these  
585 two cases, a resistant target creates very little change among the agents targeted, and agents  
586 continue to coordinate on SQ at a high rate after intervention. Miscoordination rates are low,  
587 but behaviour change is also low, exactly the opposite of what the social planner is trying to  
588 accomplish.

589 Finally, in terms of inequality, an amenable target typically generates higher inequality  
590 than a resistant target. This pattern has clear implications for the social planner who wants  
591 to limit inequality. If the social planner has a policy tool to redistribute ex post, then she can  
592 choose the intervention strategy that maximises average payoffs and then redistribute later,  
593 as we assumed for the analytical model (§ 2). If she does not have access to such a policy tool,  
594 however, the social planner must somehow resolve the resulting trade-off between competing  
595 social objectives. That said, sometimes a resistant target produces the same degree of inequal-  
596 ity or even slightly more than an amenable target. We found this pattern in our simulations  
597 for 28 combinations of parameter values. Perhaps more important than the associated details,  
598 however, the differences in Gini coefficients between the two intervention strategies tend to be  
599 small. In our simulations, over all periods and all parameter combinations, the maximum dif-  
600 ference in Gini values between an intervention in the amenable tail versus the resistant tail is  
601 approximately 0.108. This is about the same as the difference between the United States and  
602 Austria in 2018 (<https://data.worldbank.org/indicator/SI.POV.GINI/>). The  
603 mean difference in Gini values averaged over all simulation is far smaller than this at 0.036.  
604 The tendency for both intervention strategies to produce similar degrees of inequality but-  
605 tresses our use of expected payoffs as a way to evaluate social welfare.

## 606 5 Discussion

607 As with any model, we have ignored most of what matters. For example, we have ignored  
608 other types of interventions that incentivise behaviour change differently. The distinguishing  
609 feature of the models above is that payoffs for targeted individuals who respond to the inter-

610 vention do not depend on the choices of others. Nor do payoffs depend on the pre-intervention  
611 preferences of these individuals. In contrast, a social planner could simply subsidise Alt by  
612 paying  $s > 0$  to targeted individuals who choose Alt, whatever “paying” may mean. With  
613 an intervention like this, a targeted individual playing Alt gets  $a + s$  or  $d + s$ . A targeted  
614 individual playing SQ gets  $a + x_i$  or  $b + x_i$ , just like a non-targeted player, and for this rea-  
615 son the pre-existing preferences of targeted players affect payoffs so long as targeted players  
616 choose SQ after intervention with some positive probability.

617 Additionally and perhaps most importantly, we have ignored the possibility that non-  
618 targeted individuals may change their preferences after they change their behaviour. This  
619 possibility would open up an entirely new world of possibilities in terms of social welfare,  
620 but it is certainly feasible. One of us, for example, grew up in a run-of-the-mill U.S. city,  
621 geographically extensive with no meaningful public transport. When he became a teenager, he  
622 really wanted his own vehicle, so he got a job and bought one. The same person has now lived  
623 in Switzerland without a car for many years. He is coordinating on an alternative equilibrium.  
624 If his preference for his own vehicle had persisted, he would now be suffering relative to the  
625 car-based equilibrium of his adolescence and early adulthood. Exactly the opposite has  
626 happened. Experiencing the alternative has dramatically crystallised this person’s view of  
627 how dysfunctional the U.S. equilibrium is, and it has done so in a way that was not possible,  
628 bizarrely, back then while sitting in traffic. We can easily imagine an analogous change in  
629 preferences, for example, when a family abandons female genital cutting [1]. Regardless, we  
630 must be clear about what this mechanism implies. If we believe ex post preference change is  
631 the typical mechanism generating welfare improvements after norm change, we are assuming  
632 that ex ante the social planner usually knows what is best for people; the people themselves  
633 do not.

634 In spite of the mechanisms we have ignored, we hope to have gained some understanding of  
635 how risk dominance, payoff dominance, and behaviour change combine to shape social welfare  
636 in heterogeneous populations. We have argued that the relevant situation is one in which a  
637 majority of individuals view the status quo tradition (SQ) as the risk-dominant equilibrium  
638 ( $x_i > (d - b)/2$ ). If this were not so, the population probably would have never converged on  
639 the status quo tradition. With this backdrop in place, we examine welfare effects due to the  
640 payoff ranking over equilibria. Some situations are straightforward. Imagine that more or less  
641 everyone prefers coordinating on SQ over coordinating on Alt ( $a + x_i > d$ ). In addition to risk  
642 dominating Alt, coordinating on SQ also payoff dominates coordinating on Alt. In this case,  
643 assuming no externalities, people are doing what they should be doing. A well-intentioned

644 and well-informed social planner would know this and focus on other issues in society.

645 Alternatively, imagine that more or less everyone prefers coordinating on Alt over coor-  
646 dinating on SQ ( $a + x_i < d$ ). Coordinating on SQ risk dominates coordinating on Alt, but  
647 coordinating on Alt clearly payoff dominates coordinating on SQ. People really are stuck  
648 in a harmful equilibrium in this case, and individuals cannot afford to deviate as isolated  
649 decision makers. A social planner can do real good in this case by engineering a coordinated  
650 change in behaviour to the alternative equilibrium. Crucially, however, she should probably  
651 not intervene and then just walk away. She should also consider mechanisms [49, 55–57] to  
652 counteract the potential fragility of the Alt equilibrium because of any residual tendency to  
653 treat coordinating on SQ as risk-dominant. Moreover, we do not know how common situa-  
654 tions of this sort actually are. Experimental evidence shows that, although risk dominance  
655 does bias choice dynamics, payoff dominance also matters; populations sometimes find a way  
656 to the beneficial equilibrium [58].

657 In both of the situations immediately above, behaviour change and social welfare relate  
658 in a simple way. Behaviour change is uniformly bad in the first case and good in the sec-  
659 ond. Our analysis focuses on situations between these simple cases, situations in which the  
660 heterogeneity in the population really matters. We suspect that these cases are not just  
661 complicated, but probably also common. In effect, some people view coordinating on Alt  
662 as risk-dominant, and some people view coordinating on SQ as risk-dominant. Similarly,  
663 some people view coordinating on Alt as payoff-dominant, and some people view coordinat-  
664 ing on SQ as payoff-dominant. In scenarios of this sort, the link between behaviour change  
665 and social welfare can be varied and counterintuitive. In particular, we have shown that an  
666 amenable target can produce a lot of miscoordination but relatively high payoffs because it  
667 maximises the chances that the people most resistant to change are exactly the people who do  
668 not change behaviour. These individuals have the strongest preferences for coordinating on  
669 SQ. They can choose SQ after intervention and tolerate the costs of miscoordination because,  
670 when they do coordinate on SQ, they get especially large payoffs. These payoffs can be so  
671 large, in fact, that they boost the payoffs of the entire society. If these individuals could pair  
672 off to play together at rates above chance, essentially a form of homophily, this effect would  
673 be even stronger than it is under the random matching we consider.

674 More broadly, what is the value of using our knowledge of cultural evolution to engineer  
675 beneficial behaviour change? One straightforward answer follows from the idea that beneficial  
676 norms can attenuate social welfare loss when markets fail [59, 60]. In today’s world, every  
677 contract is incomplete, every price is wrong, and externalities are pervasive. Even if we

678 consider a simple transaction like buying milk, the entire global economy is implicated. From  
679 the farmer and his cows to the truck driver who drives the milk to the big city, from the  
680 company making the vinyl for the seat of the farmer’s tractor to the firm extracting the  
681 petroleum for the plastic cap on the milk carton, every production step involves dozens  
682 of contracts with some piece missing. No one knows the correct price of milk, and thus  
683 we have no reason to think that our milk-buying decisions are socially beneficial. Every  
684 day, each of us makes countless decisions that affect countless people in ways we do not  
685 understand. Activating the cultural evolution of conventions and norms that support socially  
686 beneficial choices can attenuate the challenges that follow from pervasive externalities. Our  
687 analysis, however, suggests that we still have much to learn about the best way to do so.  
688 When people differ in ordinary ways, the task of activating cultural evolution for good can  
689 become unexpectedly complex. Recent research has shown that heterogeneity can disrupt any  
690 simple monotonic relationship between the size of an intervention and the degree of behaviour  
691 change that follows [14, 42]. We have shown here that heterogeneity can also disrupt any  
692 simple monotonic relationship between the degree of behaviour change and social welfare.  
693 More surprisingly, heterogeneity can even disrupt any simple relationship between the rate  
694 of miscoordination and social welfare. Sorting through the complexities, however, will be  
695 essential precisely because widespread externalities imply that we need prosocial norms and  
696 related informal institutions to fill the welfare gap externalities leave behind.

## 697 **References**

- 698 [1] Charles Efferson, Sonja Vogt, and Lukas von Flüe. Activating cultural evolution for  
699 good when people differ from each other. In Jeremy Kendal, Rachel Kendal, and Jamie  
700 Tehrani, editors, *Oxford Handbook of Cultural Evolution*, chapter TBD. Oxford Univer-  
701 sity Press, 2023.
- 702 [2] Rachel L Kendal, Neeltje J Boogert, Luke Rendell, Kevin N Laland, Mike Webster, and  
703 Patricia L Jones. Social learning strategies: Bridge-building between fields. *Trends in*  
704 *Cognitive Sciences*, 22(7):651–665, 2018.
- 705 [3] Joseph Henrich and Francisco Gil-White. The evolution of prestige: freely conferred  
706 deference as a mechanism for enhancing the benefits of cultural transmission. *Evolution*  
707 *and Human Behavior*, 22(3):165–196, 2001.
- 708 [4] Maciej Chudek, Sarah Heller, Susan Birch, and Joseph Henrich. Prestige-biased cul-

- 709 tural learning: bystander’s differential attention to potential models influences children’s  
710 learning. *Evolution and Human Behavior*, 33(1):46–56, 2012.
- 711 [5] Charles Efferson, Rafael Lalive, Peter J. Richerson, Richard McElreath, and Mark  
712 Lubell. Conformists and mavericks: the empirics of frequency-dependent cultural trans-  
713 mission. *Evolution and Human Behavior*, 29(1):56–65, 2008.
- 714 [6] Jacob K. Goeree and Leeat Yariv. Conformity in the lab. *Journal of the Economic*  
715 *Science Association*, 1(1):15–28, 2015. ISSN 2199-6776. doi: 10.1007/s40881-015-0001-  
716 7. URL <http://dx.doi.org/10.1007/s40881-015-0001-7>.
- 717 [7] Charles Efferson, Rafael Lalive, Maria Paula Cacault, and Deborah Kistler. The evolu-  
718 tion of facultative conformity based on similarity. *PLoS One*, 11(12):e0168551, 2016.
- 719 [8] Aysha Bellamy, Ryan McKay, Sonja Vogt, and Charles Efferson. What is the extent of  
720 a frequency-dependent social learning strategy space? *Evolutionary Human Sciences*, 4,  
721 2022.
- 722 [9] Sönke Ehret, Sara Constantino, Elke Weber, Charles Efferson, and Sonja Vogt. Group  
723 identities can undermine social tipping after intervention. *Nature Human Behaviour*,  
724 Forthcoming:TBD, 2022. doi: <https://doi.org/10.1038/s41562-022-01440-5>.
- 725 [10] Olivier Morin. *How Traditions Live and Die*. Oxford University Press, 2016.
- 726 [11] Robert Boyd and Peter J. Richerson. *Culture and the Evolutionary Process*. Chicago:  
727 University of Chicago Press, 1985.
- 728 [12] Peter J. Richerson and Richerson Boyd. *Not By Genes Alone: How Culture Transformed*  
729 *the Evolutionary Process*. Chicago: University of Chicago Press, 2005.
- 730 [13] T. J. H. Morgan and K. N. Laland. The biological bases of conformity. *Frontiers in*  
731 *Neuroscience*, 6(87), 2012. doi: 10.3389/fnins.2012.00087.
- 732 [14] Charles Efferson, Sonja Vogt, and Ernst Fehr. The promise and the peril of using social  
733 influence to reverse harmful traditions. *Nature Human Behaviour*, 4:55–68, 2020.
- 734 [15] Ellen Gruenbaum. *The Female Circumcision Controversy: An Anthropological Perspec-*  
735 *tive*. Philadelphia: University of Pennsylvania Press, 2001.
- 736 [16] Karisa Cloward. *When Norms Collide: Local Responses to Activism Against Female*  
737 *Genital Mutilation and Early Marriage*. Oxford University Press, 2016.



- 738 [17] David Lawson and Mhairi Gibson. Evolutionary approaches to population health: In-  
739 sights on polygynous marriage, ‘child marriage’ and female genital mutilation/cutting.  
740 In O. Burger, R. Lee, and R. Sear, editors, *Human Evolutionary Demography*, page  
741 TBD. 2023.
- 742 [18] Damon Centola, Joshua Becker, Devon Brackbill, and Andrea Baronchelli. Experimental  
743 evidence for tipping points in social convention. *Science*, 360(6393):1116–1119, 2018.
- 744 [19] James Andreoni, Nikos Nikiforakis, and Simon Siegenthaler. Predicting social tipping  
745 and norm change in controlled experiments. *Proceedings of the National Academy of*  
746 *Sciences*, 118(16), 2021.
- 747 [20] Gerry Mackie. Ending Footbinding and Infibulation: A Convention Account. *American*  
748 *Sociological Review*, 61:999–1017, 1996.
- 749 [21] Bettina Shell-Duncan and Ylva Hernlund. Female “Circumcision” in Africa: Dimensions  
750 of the Practice and Debates. In Bettina Shell-Duncan and Ylva Hernlund, editors, *Female*  
751 *“Circumcision” in Africa: Culture, Controversy, and Change*, pages 1–40. Boulder, CO:  
752 Lynne Rienner, 2000.
- 753 [22] Marc F. Bellemare, Lindsey Novak, and Tara L. Steinmetz. All in the  
754 family: Explaining the persistence of female genital cutting in West  
755 Africa. *Journal of Development Economics*, 116:252 – 265, 2015. ISSN  
756 0304-3878. doi: <http://dx.doi.org/10.1016/j.jdeveco.2015.06.001>. URL  
757 <http://www.sciencedirect.com/science/article/pii/S0304387815000620>.
- 758 [23] Charles Efferson, Sonja Vogt, Amy Elhadi, Hilal El Fadil Ahmed, and Ernt Fehr. Female  
759 genital cutting is not a social coordination norm. *Science*, 349(6255):1446–1447, 2015.  
760 doi: 10.1126/science.aaa7978.
- 761 [24] Sonja Vogt, Nadia Ahmed Mohammed Zaid, Hilal El Fadil Ahmed, Ernst Fehr, and  
762 Charles Efferson. Changing cultural attitudes towards female genital cutting. *Nature*,  
763 538:506–509, 2016.
- 764 [25] Lindsey Novak. Persistent norms and tipping points: The case of female genital cutting.  
765 *Journal of Economic Behavior & Organization*, 177:433–474, 2020.
- 766 [26] Susan Lee-Rife, Anju Malhotra, Ann Warner, and Allison McGonagle Glinski. What  
767 works to prevent child marriage: A review of the evidence. *Studies in Family Planning*,

- 768 43(4):287–303, 2012. ISSN 1728-4465. doi: 10.1111/j.1728-4465.2012.00327.x. URL  
769 <http://dx.doi.org/10.1111/j.1728-4465.2012.00327.x>.
- 770 [27] Jean-Philippe Platteau, Giulia Camilotti, and Emmanuelle Auriol. Eradicating women-  
771 hurting customs. In Siwan Anderson, Lori Beaman, and Jean-Philippe Platteau, editors,  
772 *Towards Gender Equity in Development*, pages 319–356. Oxford University Press, 2018.
- 773 [28] Matthias Schief, Sonja Vogt, and Charles Efferson. Investigating the structure of son bias  
774 in Armenia with novel measures of individual preferences. *Demography*, 58:1737–1764,  
775 2021.
- 776 [29] Juan Carlos Castilla-Rho, Rodrigo Rojas, Martin S Andersen, Cameron Holley, and  
777 Gregoire Mariethoz. Social tipping points in global groundwater management. *Nature*  
778 *Human Behaviour*, 1(9):640–649, 2017.
- 779 [30] Nicholas A Christakis and James H Fowler. The spread of obesity in a large social  
780 network over 32 years. *New England Journal of Medicine*, 2007(357):370–379, 2007.
- 781 [31] Nicholas A Christakis and James H Fowler. The collective dynamics of smoking in a  
782 large social network. *New England Journal of Medicine*, 358(21):2249–2258, 2008.
- 783 [32] Henry Travers, James Walsh, Sonja Vogt, Tom Clements, and EJ Milner-Gulland. De-  
784 livering behavioural change at scale: What conservation can learn from other fields.  
785 *Biological Conservation*, 257:109092, 2021.
- 786 [33] Karine Nyborg, John M Anderies, Astrid Dannenberg, Therese Lindahl, Caroline Schill,  
787 Maja Schlüter, W Neil Adger, Kenneth J Arrow, Scott Barrett, Stephen Carpenter, et al.  
788 Social norms as solutions. *Science*, 354(6308):42–43, 2016.
- 789 [34] J Doyne Farmer, Cameron Hepburn, Matthew C Ives, T Hale, Thomas Wetzer, Penny  
790 Mealy, Ryan Rafaty, Sugandha Srivastav, and Rupert Way. Sensitive intervention points  
791 in the post-carbon transition. *Science*, 364(6436):132–134, 2019.
- 792 [35] Marwil J Dávila-Fernández and Serena Sordi. Attitudes towards climate policies in a  
793 macrodynamic model of the economy. *Ecological Economics*, 169:106319, 2020.
- 794 [36] Ilona M Otto, Jonathan F Donges, Roger Cremades, Avit Bhowmik, Richard J Hewitt,  
795 Wolfgang Lucht, Johan Rockström, Franziska Allerberger, Mark McCaffrey, Sylvanus SP  
796 Doe, et al. Social tipping dynamics for stabilizing earth’s climate by 2050. *Proceedings*  
797 *of the National Academy of Sciences*, 117(5):2354–2365, 2020.

- 798 [37] Sara M Constantino, Gregg Sparkman, Gordon T Kraft-Todd, Cristina Bicchieri, Da-  
799 mon Centola, Bettina Shell-Duncan, Sonja Vogt, and Elke U Weber. Scaling up change:  
800 A critical review and practical guide to harnessing social norms for climate action. *Psy-*  
801 *chological Science in the Public Interest*, 23(2):50–97, 2022.
- 802 [38] Charles Efferson. Policy to activate cultural change to amplify policy. *Proceedings of the*  
803 *National Academy of Sciences*, 118(23), 2021.
- 804 [39] Mark Granovetter. Threshold models of collective behavior. *American Journal of Soci-*  
805 *ology*, 83(6):1420–1443, 1978.
- 806 [40] H Peyton Young. Innovation diffusion in heterogeneous populations: contagion, social  
807 influence, and social learning. *American Economic Review*, 99(5):1899–1924, 2009.
- 808 [41] Matthew O Jackson and Dunia López-Pintado. Diffusion and contagion in networks  
809 with heterogeneous agents and homophily. *Network Science*, 1(1):49–67, 2013.
- 810 [42] Robin Schimmelpfennig, Sonja Vogt, Sönke Ehret, and Charles Efferson. Promotion  
811 of behavioural change for health in a heterogeneous population. *Bulletin of the World*  
812 *Health Organization*, 99(11):819, 2021.
- 813 [43] Duncan J. Watts and Peter Dodds. Threshold models of social influence. *The Oxford*  
814 *Handbook of Analytical Sociology*, pages 475–497, 2009.
- 815 [44] Michael Mäs and Heinrich H Nax. A behavioral study of “noise” in coordination games.  
816 *Journal of Economic Theory*, 162:195–208, 2016.
- 817 [45] John C. Harsanyi and Reinhard Selten. *A General Theory of Equilibrium Selection in*  
818 *Games*. Cambridge: The MIT Press, 1988.
- 819 [46] John B Van Huyck, Raymond C Battalio, and Richard O Beil. Tacit coordination games,  
820 strategic uncertainty, and coordination failure. *The American Economic Review*, 80(1):  
821 234–248, 1990.
- 822 [47] Marc Knez and Colin Camerer. Creating expectational assets in the laboratory: Coordi-  
823 nation in ‘weakest-link’ games. *Strategic Management Journal*, 15(S1):101–119, 1994.
- 824 [48] Paul G Straub. Risk dominance and coordination failures in static games. *The Quarterly*  
825 *Review of Economics and Finance*, 35(4):339–363, 1995.

- 826 [49] Kenneth Clark, Stephen Kay, and Martin Sefton. When are nash equilibria self-  
827 enforcing? an experimental analysis. *International Journal of Game Theory*, 29:495–515,  
828 2001.
- 829 [50] Roberto A Weber. Managing growth to achieve efficient coordination in large groups.  
830 *American Economic Review*, 96(1):114–126, 2006.
- 831 [51] Charles Efferson, Ryan McKay, and Ernst Fehr. The evolution of distorted beliefs vs. mis-  
832 taken choices under asymmetric error costs. *Evolutionary Human Sciences*, 2:e27, 2020.  
833 doi: 10.1017/ehs.2020.25.
- 834 [52] Jasmina Arifovic, Cars Hommes, Anita Kopányi-Peuker, and Isabelle Salle. Ten isn’t  
835 large! group size and coordination in a large-scale experiment. *American Economic*  
836 *Journal: Microeconomics*, 15(1):580–617, 2023.
- 837 [53] Eva Vivalt. Heterogeneous treatment effects in impact evaluation. *American Economic*  
838 *Review*, 105(5):467–70, 2015.
- 839 [54] Ulrich Schmidt and Philipp C Wichardt. Inequity aversion, welfare measurement and  
840 the Gini index. *Social Choice and Welfare*, 52:585–588, 2019.
- 841 [55] Russell Cooper, Douglas V DeJong, Robert Forsythe, and Thomas W Ross. Commu-  
842 nication in coordination games. *The Quarterly Journal of Economics*, 107(2):739–771,  
843 1992.
- 844 [56] Andreas Blume and Andreas Ortmann. The effects of costless pre-play communication:  
845 Experimental evidence from games with pareto-ranked equilibria. *Journal of Economic*  
846 *Theory*, 132(1):274–290, 2007.
- 847 [57] Catherine C Eckel and Rick K Wilson. Social learning in coordination games: does  
848 status matter? *Experimental Economics*, 10:317–329, 2007.
- 849 [58] Raymond Battalio, Larry Samuelson, and John Van Huyck. Optimization incentives and  
850 coordination failure in laboratory stag hunt games. *Econometrica*, 69(3):749–764, 2001.
- 851 [59] Samuel Bowles. *Microeconomics: Behavior, Institutions, and Evolution*. New York:  
852 Russell Sage, 2004.
- 853 [60] Sergey Gavrillets and Peter J Richerson. Collective action and the evolution of social  
854 norm internalization. *Proceedings of the National Academy of Sciences*, page 201703857,  
855 2017.

## 856 **Acknowledgements and funding**

857 CE and SV would like to thank the Swiss National Science Foundation (Nr. 100018\_185417)  
858 for financial support.

## 859 **Author contributions**

860 CE, SE, and SV developed the initial idea. CE developed and analysed the analytical model.  
861 LvF developed and analysed the simulation model with input from SE and CE. CE wrote  
862 the paper with feedback from SE, LvF, and SV.

## 863 **Competing interests**

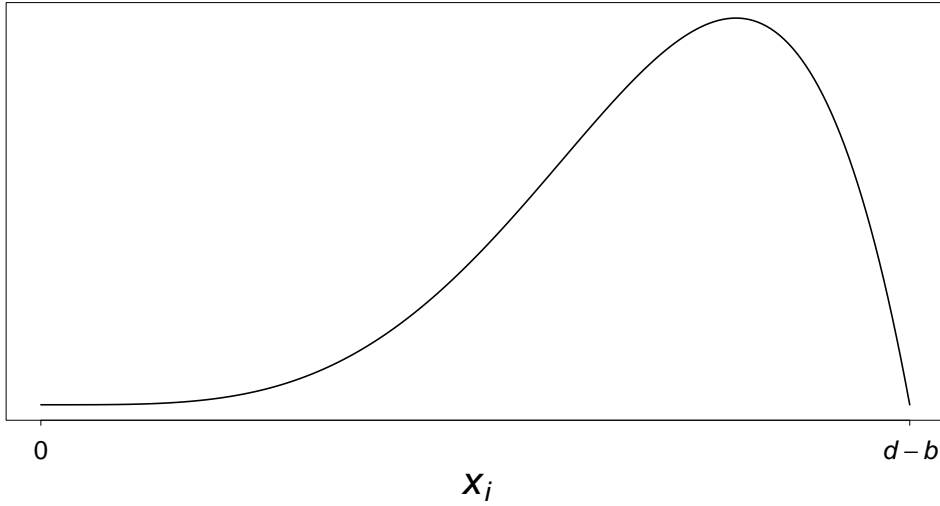
864 We have no competing interests.

Table 1: **Payoff matrices before and after intervention.** The matrices show row player payoffs. Before intervention, everyone plays a coordination game, but individuals also have heterogeneous preferences (Pre-intervention (All)). Specifically, the  $x_i$  are distributed somehow on  $(0, d - b)$ , where  $b < d$ . This ensures that, for each  $i$ ,  $a + x_i > a$ , and  $b + x_i < d$ . The social planner then targets a subset of the population with an intervention and incentivises these individuals to choose Alt (Post-intervention (T)) with a new payoff matrix, where  $h > g$ . Individuals who are not targeted retain their original payoff matrices (Post-intervention (NT)).

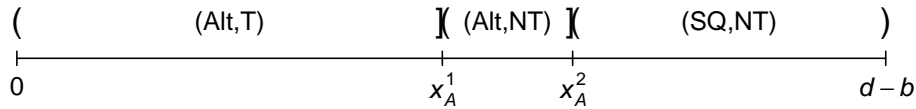
Pre-intervention (All)		Post-intervention (T)		Post-intervention (NT)	
	SQ	Alt		SQ	Alt
SQ	$a + x_i$	$b + x_i$	SQ	$g$	$g$
Alt	$a$	$d$	Alt	$h$	$h$

Table 2: **Player matches.** When the population stabilises after intervention, each individual falls into one of three categories of player. Targeted individuals play Alt (Alt,T). Some or all non-targeted players may also play Alt (Alt,NT). Some or all non-targeted players may play SQ (SQ,NT). The table shows all ways in which two players can pair off to play the game  $(P_1, P_2)$ , the probabilities associated with these pairings  $(P(P_1, P_2))$ , and the payoffs that result  $(\pi_1, \pi_2)$ . The table uses a generic  $\hat{q}_z$ , where  $z \in \{A, R\}$ , because all terms here have the same structure regardless of whether the target is amenable ( $z = A$ ) or resistant ( $z = R$ ). We use a prime to distinguish between the preferences of Player 1 and Player 2 when the  $x_i$  appear.

$P_1$	$P_2$	$P(P_1, P_2)$	$\pi_1$	$\pi_2$
(Alt,T)	(Alt,T)	$\phi^2$	$h$	$h$
(Alt,T)	(Alt,NT)	$\phi(\hat{q}_z - \phi)$	$h$	$d$
(Alt,T)	(SQ,NT)	$\phi(1 - \hat{q}_z)$	$h$	$b + x_{i'}$
(Alt,NT)	(Alt,T)	$(\hat{q}_z - \phi)\phi$	$d$	$h$
(Alt,NT)	(Alt,NT)	$(\hat{q}_z - \phi)^2$	$d$	$d$
(Alt,NT)	(SQ,NT)	$(\hat{q}_z - \phi)(1 - \hat{q}_z)$	$a$	$b + x_{i'}$
(SQ,NT)	(Alt,T)	$(1 - \hat{q}_z)\phi$	$b + x_i$	$h$
(SQ,NT)	(Alt,NT)	$(1 - \hat{q}_z)(\hat{q}_z - \phi)$	$b + x_i$	$a$
(SQ,NT)	(SQ,NT)	$(1 - \hat{q}_z)^2$	$a + x_i$	$a + x_{i'}$



**Amenable target:**



**Resistant target:**



**Figure 1: An example of two different partitions of the population.** The figure shows how an amenable target versus a resistant target might partition the population in different ways given the distribution of  $x_i$  values shown. For an amenable target, the example assumes the social planner targets 10% of the population ( $\phi = 0.1$ ). Under an amenable target, the population stabilises post-intervention on 30% choosing Alt ( $\hat{q}_A = 0.3$ ). Specifically, (Alt,T) are the targeted individuals who change behaviour due to direct experience with the intervention, and (Alt,NT) are the non-targeted individuals who change behaviour due to coordination incentives post-intervention. The remaining 70% of the population are the non-targeted individuals who continue choosing the status quo behaviour, denoted (SQ,NT). For a resistant target, the example also assumes the social planner targets 10%, but in the long-run 50% end up choosing Alt. The partition is completely different from the amenable case because the (Alt,T) individuals come from the right tail of the  $x_i$  distribution, which leaves the left tail for the (Alt,NT) individuals and the middle for the (SQ,NT) individuals. The parentheses and square brackets denote where the intervals associated with the partition are open or closed respectively.



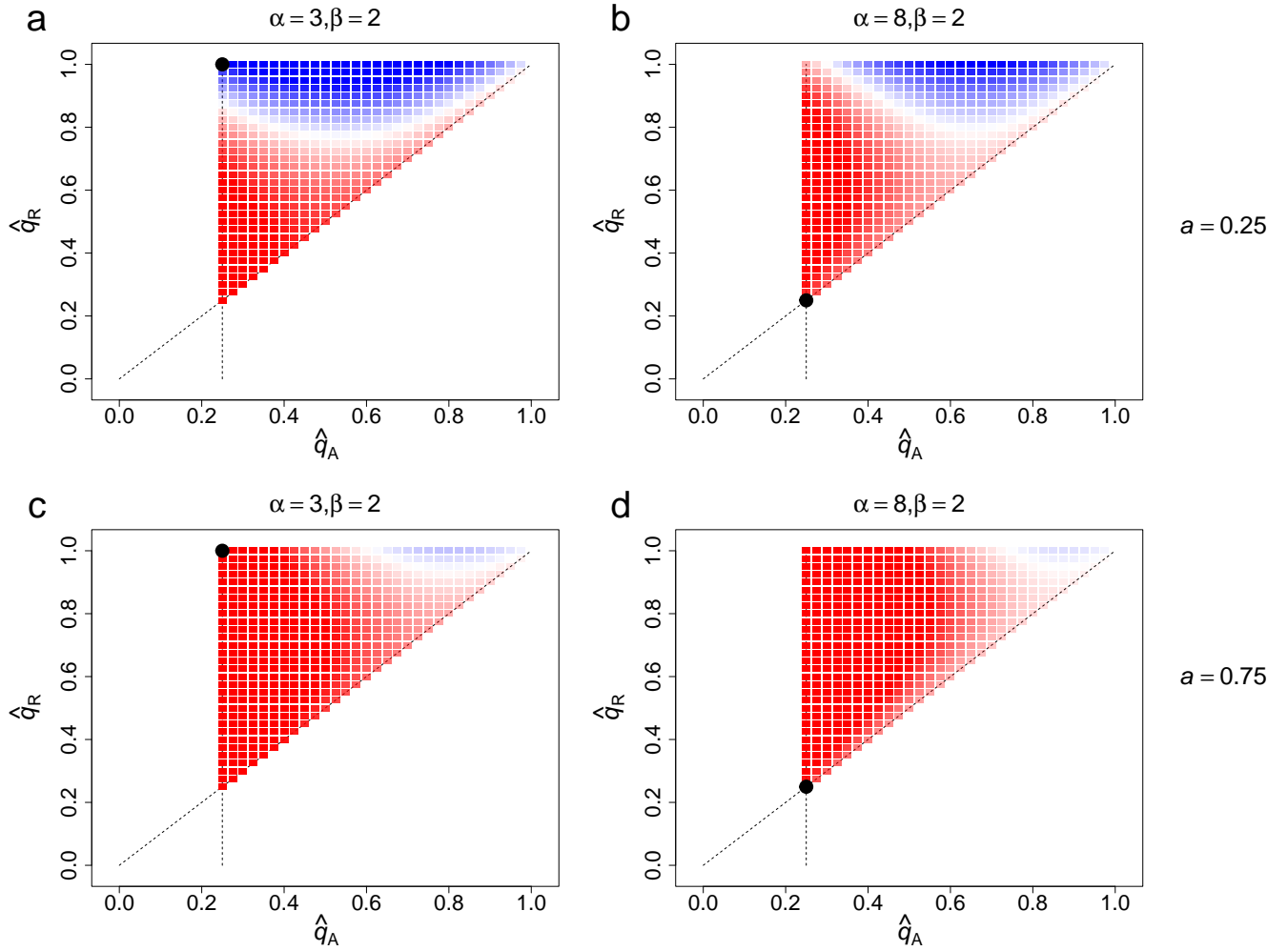


Figure 2: **Comparison of intervention strategies given an intervention of size  $\phi = 0.25$ .** The graphs represent the values of  $E_A[\Pi] - E_R[\Pi]$  (Eq. 6) for all steady states consistent with  $\phi \leq \hat{q}_A \leq \hat{q}_R$ . Red indicates that  $E_A[\Pi] > E_R[\Pi]$ , and thus an amenable target is better in terms of social welfare than a resistant target. Blue indicates the opposite. Colour intensity reflects the magnitude of the differences. The black dot shows the steady state under the threshold model [14, 39]. To generate these graphs,  $b = 0$ , and  $d = 1$ . The parameter  $a$  takes either a moderately low value (**a**, **b**) or a moderately high value (**c**, **d**). The  $x_i$  values are beta distributed with shape parameters  $\alpha$  and  $\beta$  on the interval  $(0, 1)$ . The distribution is either weakly left-skewed (**a**, **c**) or somewhat strongly left-skewed (**b**, **d**).

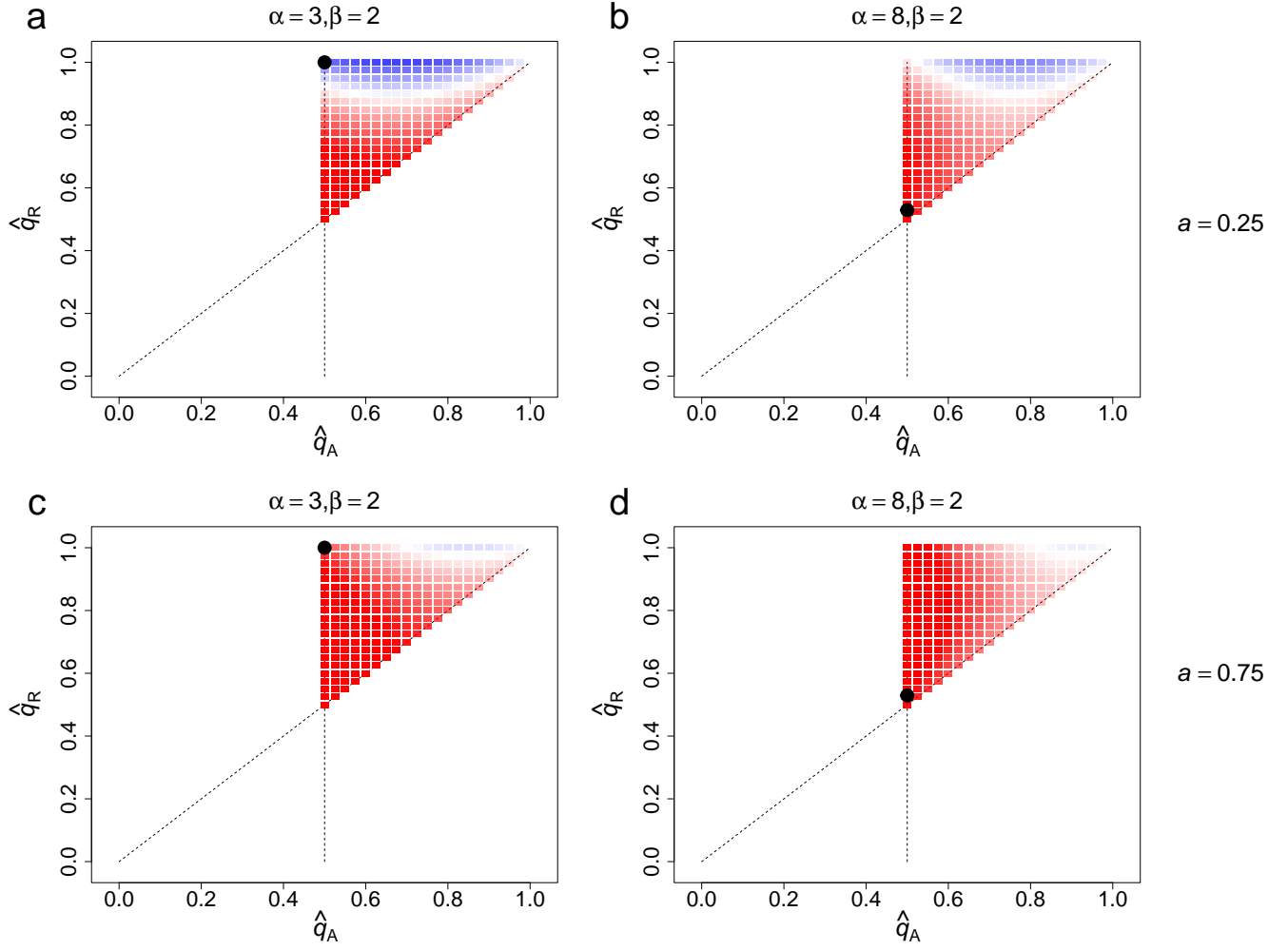


Figure 3: **Comparison of intervention strategies given an intervention of size  $\phi = 0.5$ .** The graphs represent the values of  $E_A[\Pi] - E_R[\Pi]$  (Eq. 6) for all steady states consistent with  $\phi \leq \hat{q}_A \leq \hat{q}_R$ . Red indicates that  $E_A[\Pi] > E_R[\Pi]$ , and thus an amenable target is better in terms of social welfare than a resistant target. Blue indicates the opposite. Colour intensity reflects the magnitude of the differences. The black dot shows the steady state under the threshold model [14, 39]. To generate these graphs,  $b = 0$ , and  $d = 1$ . The parameter  $a$  takes either a moderately low value (**a**, **b**) or a moderately high value (**c**, **d**). The  $x_i$  values are beta distributed with shape parameters  $\alpha$  and  $\beta$  on the interval  $(0, 1)$ . The distribution is either weakly left-skewed (**a**, **c**) or somewhat strongly left-skewed (**b**, **d**).

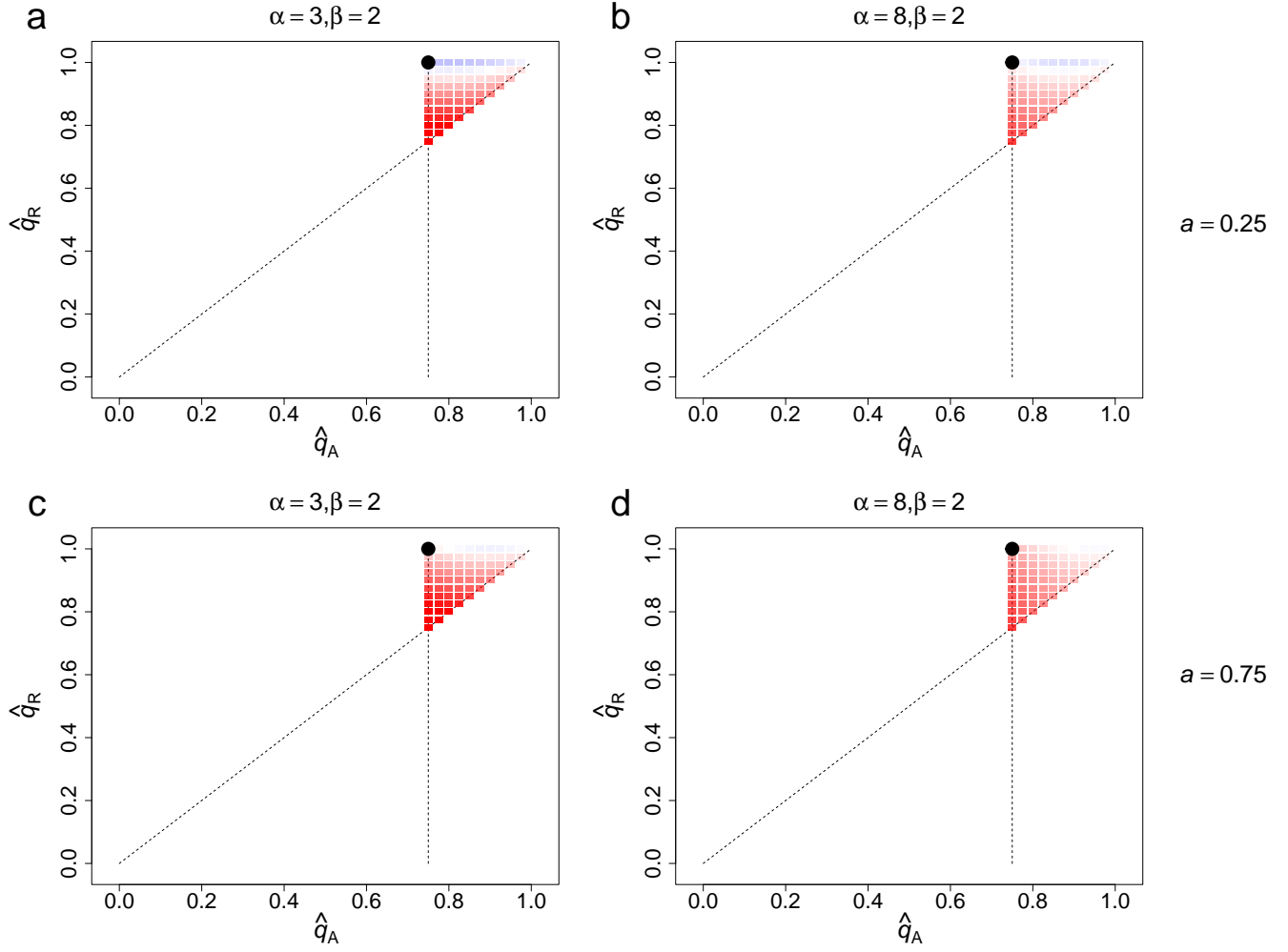


Figure 4: **Comparison of intervention strategies given an intervention of size  $\phi = 0.75$ .** The graphs represent the values of  $E_A[\Pi] - E_R[\Pi]$  (Eq. 6) for all steady states consistent with  $\phi \leq \hat{q}_A \leq \hat{q}_R$ . Red indicates that  $E_A[\Pi] > E_R[\Pi]$ , and thus an amenable target is better in terms of social welfare than a resistant target. Blue indicates the opposite. Colour intensity reflects the magnitude of the differences. The black dot shows steady state values for both intervention strategies under the threshold model [14, 39]. To generate these graphs,  $b = 0$ , and  $d = 1$ . The parameter  $a$  takes either a moderately low value (**a**, **b**) or a moderately high value (**c**, **d**). The  $x_i$  values are beta distributed with shape parameters  $\alpha$  and  $\beta$  on the interval  $(0, 1)$ . The distribution is either weakly left-skewed (**a**, **c**) or somewhat strongly left-skewed (**b**, **d**).

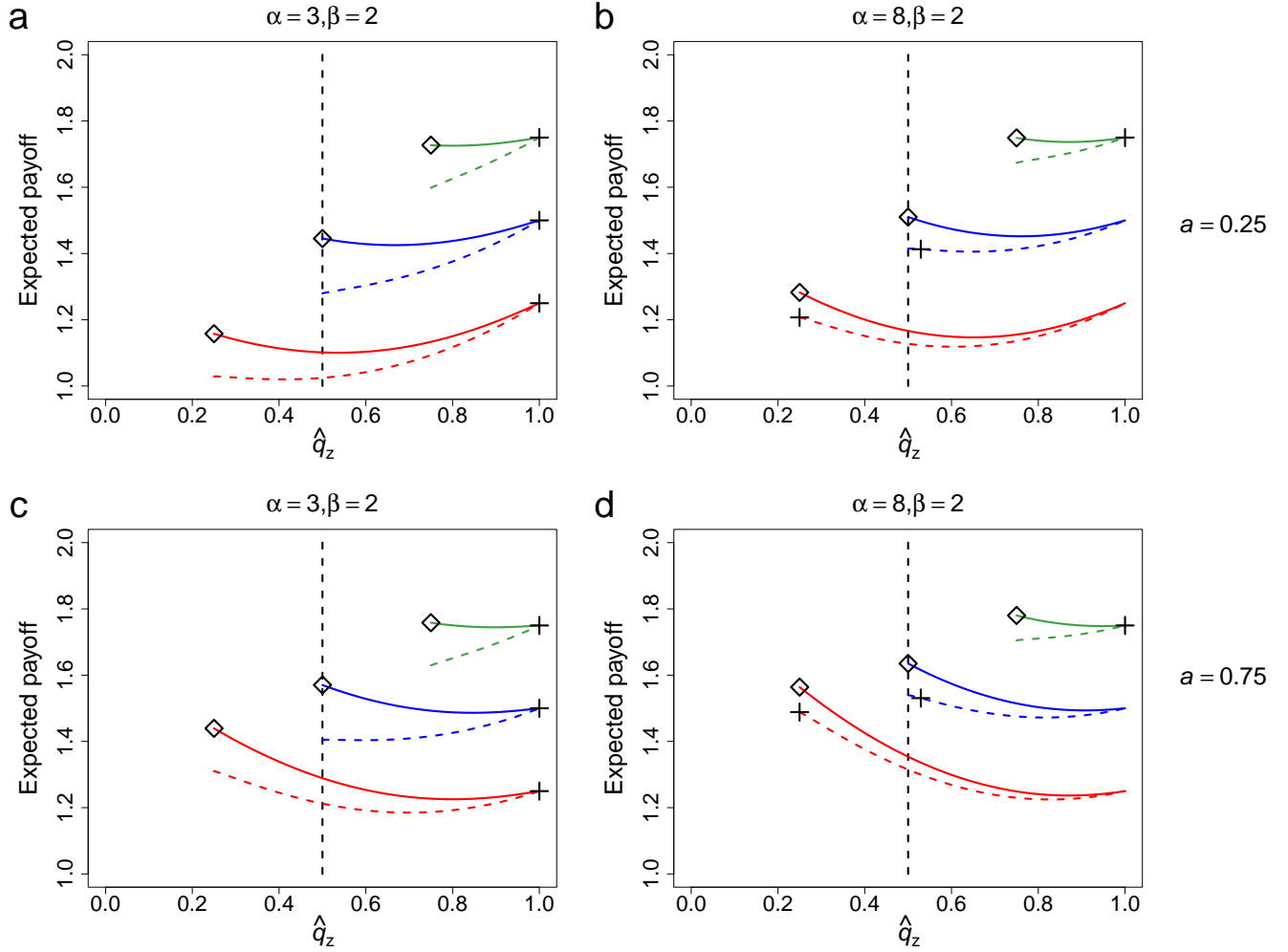


Figure 5: **Social welfare and (mis-)coordination.** The graphs show the expected payoff of a randomly selected individual as a function of the steady state proportion choosing Alt,  $\hat{q}_z$ , where  $z \in \{A, R\}$ . Solid lines show the expected payoffs under an amenable target (Eq. 3) and dashed lines under a resistant target (Eq. 4). Intervention sizes are  $\phi = 0.25$  in red,  $\phi = 0.5$  in blue, and  $\phi = 0.75$  in green. The  $\diamond$  shows the steady state value under the threshold model [14, 39] given an amenable target and the  $+$  given a resistant target. To generate these graphs,  $b = 0$ , and  $d = 1$ . The parameter  $a$  takes either a moderately low value (**a**, **b**) or a moderately high value (**c**, **d**). The  $x_i$  values are beta distributed with shape parameters  $\alpha$  and  $\beta$  on the interval  $(0, 1)$ . The distribution is either weakly left-skewed (**a**, **c**) or somewhat strongly left-skewed (**b**, **d**). Miscoordination reaches the maximum possible rate under random matching when  $\hat{q}_z = 0.5$ . Thus, any time expected payoffs increase as the steady state approaches this value, social welfare is increasing even though miscoordination is also increasing.

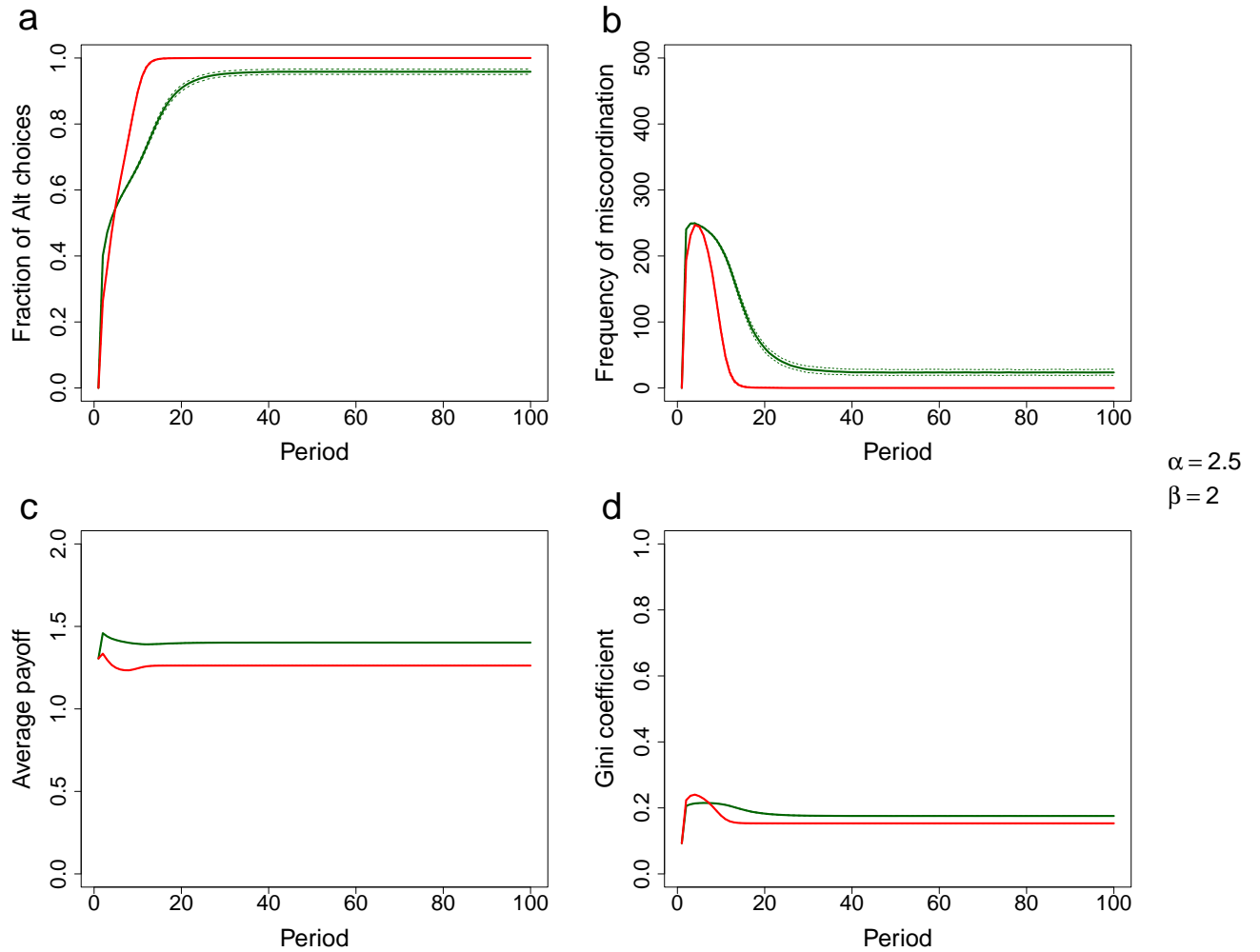


Figure 6: The figure shows a comparison between the dynamics of amenable and resistant target. The solid lines show values averaged over all 1000 simulations. The dashed lines show 95% bootstrapped confidence intervals. Populations with amenable targets are shown in green, and populations with resistant targets are shown in red. Graphs show the fraction of agents choosing Alt (a), the frequency of miscoordination (b), average payoffs (c), and the Gini coefficient (d). The initial conditions at  $t = 1$  are pre-intervention. Thus, a comparison between  $t = 1$  and any  $t > 1$  shows how populations changed, conditional on an amenable or resistant target, as a result of the intervention. Aside from  $\alpha$  and  $\beta$ , parameter values are  $\phi = 0.75$ ,  $a = 0.75$ , and  $h = 2$ .

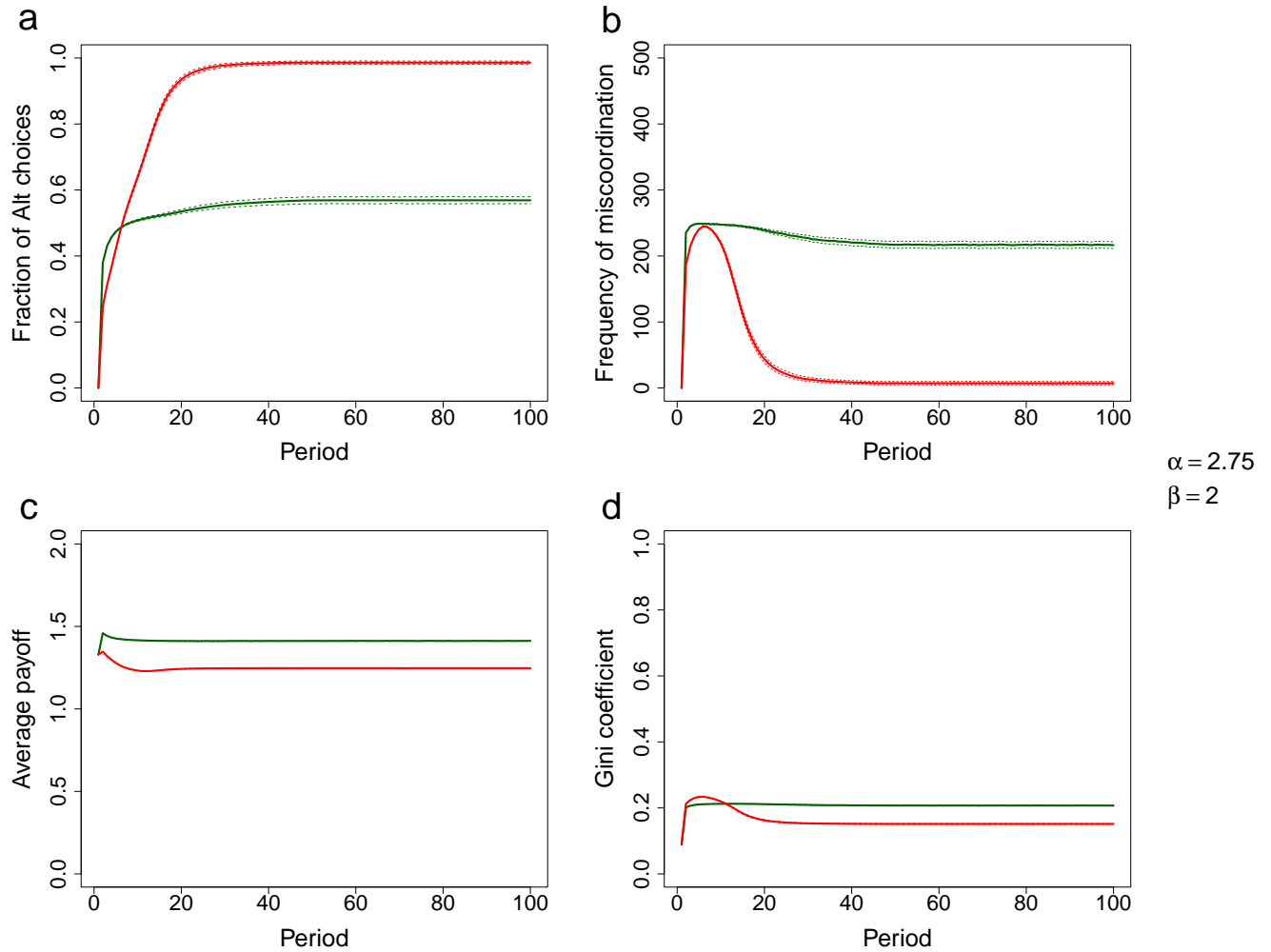


Figure 7: The figure shows a comparison between the dynamics of amenable and resistant target. The solid lines show values averaged over all 1000 simulations. The dashed lines show 95% bootstrapped confidence intervals. Populations with amenable targets are shown in green, and populations with resistant targets are shown in red. Graphs show the fraction of agents choosing Alt (a), the frequency of miscoordination (b), average payoffs (c), and the Gini coefficient (d). The initial conditions at  $t = 1$  are pre-intervention. Thus, a comparison between  $t = 1$  and any  $t > 1$  shows how populations changed, conditional on an amenable or resistant target, as a result of the intervention. Aside from  $\alpha$  and  $\beta$ , parameter values are  $\phi = 0.75$ ,  $a = 0.75$ , and  $h = 2$ .

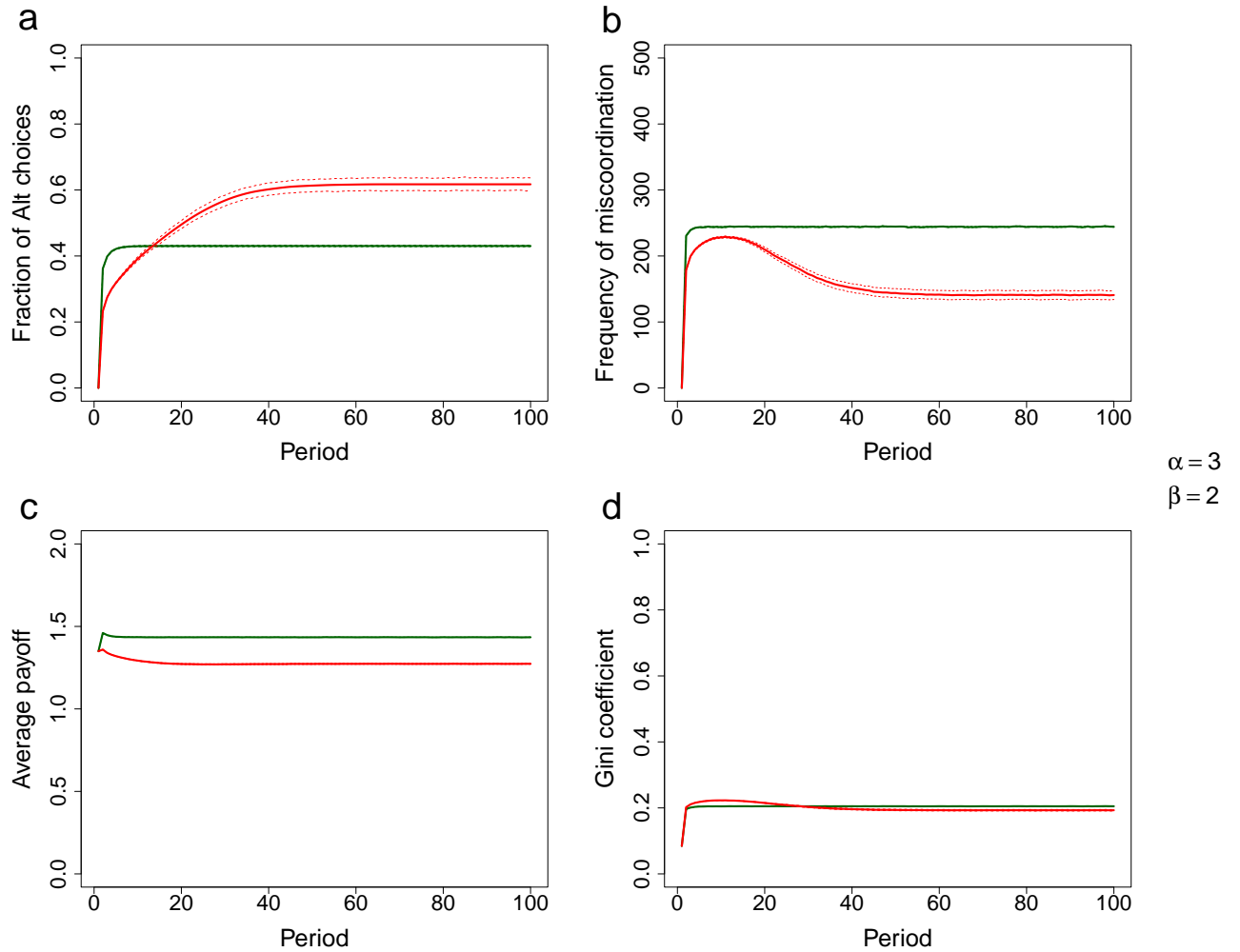


Figure 8: The figure shows a comparison between the dynamics of amenable and resistant target. The solid lines show values averaged over all 1000 simulations. The dashed lines show 95% bootstrapped confidence intervals. Populations with amenable targets are shown in green, and populations with resistant targets are shown in red. Graphs show the fraction of agents choosing Alt (a), the frequency of miscoordination (b), average payoffs (c), and the Gini coefficient (d). The initial conditions at  $t = 1$  are pre-intervention. Thus, a comparison between  $t = 1$  and any  $t > 1$  shows how populations changed, conditional on an amenable or resistant target, as a result of the intervention. Aside from  $\alpha$  and  $\beta$ , parameter values are  $\phi = 0.75$ ,  $a = 0.75$ , and  $h = 2$ .

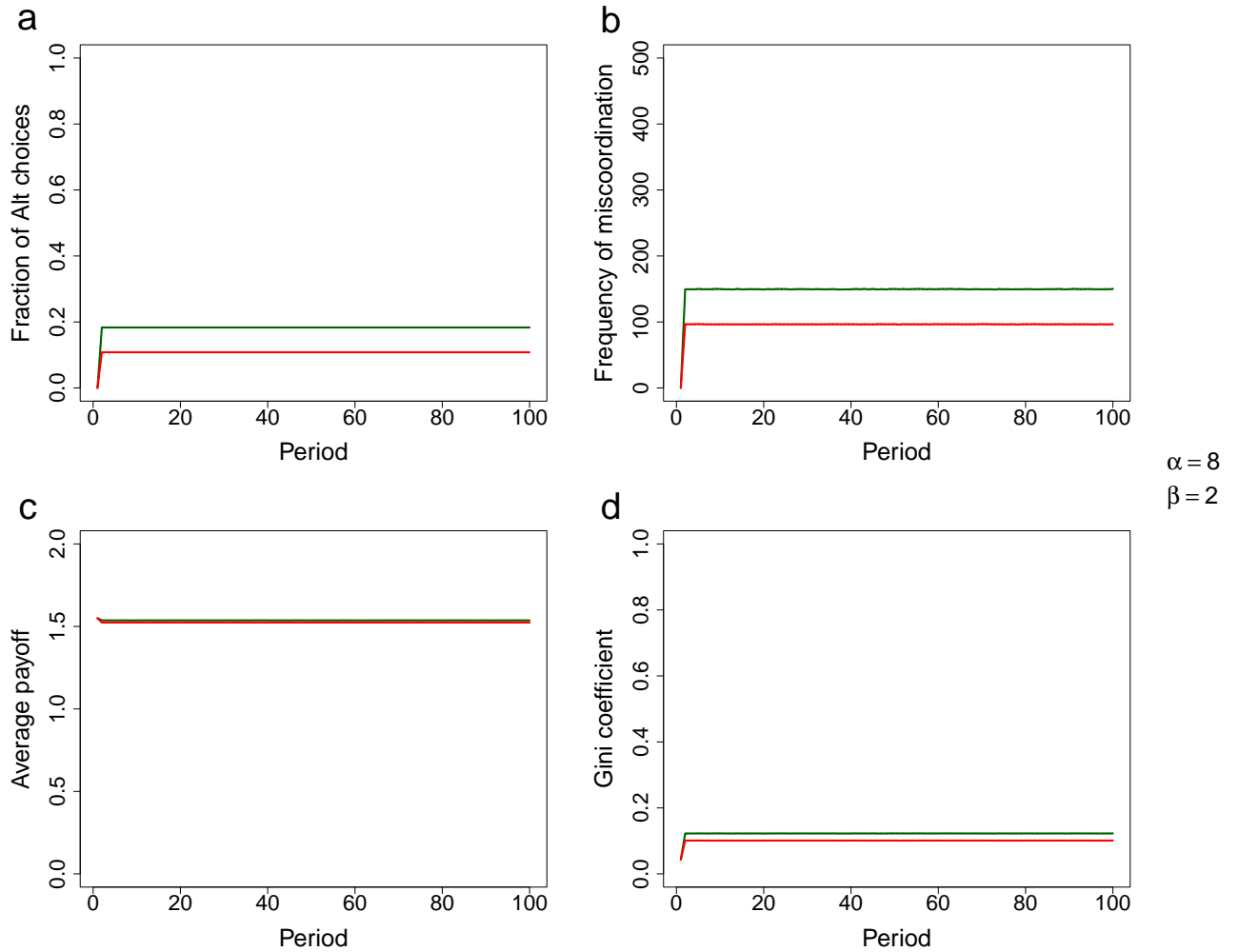


Figure 9: The figure shows a comparison between the dynamics of amenable and resistant target. The solid lines show values averaged over all 1000 simulations. The dashed lines show 95% bootstrapped confidence intervals. Populations with amenable targets are shown in green, and populations with resistant targets are shown in red. Graphs show the fraction of agents choosing Alt (a), the frequency of miscoordination (b), average payoffs (c), and the Gini coefficient (d). The initial conditions at  $t = 1$  are pre-intervention. Thus, a comparison between  $t = 1$  and any  $t > 1$  shows how populations changed, conditional on an amenable or resistant target, as a result of the intervention. Aside from  $\alpha$  and  $\beta$ , parameter values are  $\phi = 0.75$ ,  $a = 0.75$ , and  $h = 2$ .