

Cost-Efficient Interventions for Promoting Fairness in the Ultimatum Game

Theodor Cimpanu¹, Cedric Perret², and The Anh Han^{1,*}

¹ School of Computing, Engineering and Digital Technologies, Teesside University, Middlesbrough, UK TS1 3BA

² College of Life and Environmental Sciences, University of Exeter, Exeter, UK EX4 4PY

* Corresponding author: The Anh Han (T.Han@tees.ac.uk)

1 **Abstract**

2 Institutions and investors face the constant challenge of making accurate decisions and predic-
3 tions regarding how best they should distribute their endowments. The problem of achieving
4 an optimal outcome at a minimal cost has been extensively studied and resolved using several
5 heuristics. However, these works usually failed to address how an external party can target
6 different types of fair behaviour or do not take into account how limited information can shape
7 this complex interplay. Here, we consider the Ultimatum game in a spatial setting and propose
8 a hierarchy of interference mechanisms based on the amount of information available to an ex-
9 ternal decision-maker and desired standards of fairness. Our analysis reveals that monitoring
10 the population at a macroscopic level requires more strict information gathering in order to ob-
11 tain an optimal outcome and that local observations can mediate this requirement. Moreover,
12 we identify the conditions which must be met for an individual to be eligible for investment
13 in order to avoid unnecessary spending. We further explore the effects of varying mutation or
14 behavioural exploration rates on the choice of investment strategy and total accumulated costs
15 to the investor. Overall, our analysis provides new insights about efficient heuristics for cost-
16 efficient promotion of fairness in societies. Finally, we discuss the differences between our
17 findings and previous work done on cooperation dilemmas and present our suggestions for pro-
18 moting fairness as an external decision-maker.

19

20 **Keywords:** Ultimatum game, interference, cost efficiency, decision making, evolutionary game
21 theory, networks.

1 Introduction

The problem of how collective behaviour, such as cooperation, coordination, safety compliance and fairness among self-interested individuals, emerges in evolving, dynamical systems has fascinated researchers from many disciplines, ranging from Evolutionary Biology, Economics, Physics, Social Sciences and Computer Science (Airiau et al., 2014, Han, 2013, Han et al., 2020, Liu and Chen, 2020, Maynard-Smith, 1982, Nowak, 2006, Perc et al., 2017, Perc and Szolnoki, 2010, Sigmund, 2010, Tuyls and Parsons, 2007, West et al., 2007). Several mechanisms that are responsible for promoting the emergence of cooperation have been proposed, including direct and indirect reciprocity (Okada, 2020, Sigmund, 2010, Trivers, 1971), kin selection (Hamilton, 1964) and network reciprocity (Nowak, 2012, Ohtsuki et al., 2006, West et al., 2007). In these works, the evolution of desired collective behaviour is typically shaped by the combined actions of individuals within the systems.

On the other hand, external interference, where the advocating of certain desired collective behaviour is carried out by an external decision maker, who does not belong to the system, studies how this can be done in a cost-effective way (Chen et al., 2015, Cimpeanu et al., 2019, Han et al., 2018, Han and Tran-Thanh, 2018, Wang et al., 2019). These works aim to identify a broad class of interference strategies, or heuristics, that exploit available information such as global statistics (population behavioural composition), as well as local information such as local behaviour profile and diversity, and graph structures, for budget saving. This line of research is useful to provide insights into the design of self-organised and distributed Multi-Agent Systems (MAS), in order to ensure agents achieve a desired collective state. For instance, one might consider a hybrid system consisting of humans and intelligent machines, in which it is important to ensure a cooperative and trustful relationship amongst each other (Andras et al., 2018, Paiva et al., 2018, Santos et al., 2019). Another example is how international agencies such as the European Union and United Nations might advocate certain preferred political behaviours or resolve international conflicts, given a limited budget (e.g. in terms of cost and military resources) (Marton-Lafevre et al., 2007, Smidt, 2020).

The literature on external interference in evolving, dynamical systems (or populations) has so far focused on cooperation dilemmas, namely the Prisoner's Dilemma (PD) (Cimpeanu et al., 2019, Han et al., 2018, Han and Tran-Thanh, 2018) and the Public Goods Game (PGG) (Chen and Perc, 2014, Chen et al., 2015, Sasaki et al., 2012, Wang et al., 2019). In these games, the interactions are symmetric and the players' roles are equivalent. However, many real-world and

54 MAS interactions are asymmetric, where players may have different baseline characteristics
55 and/or play different roles in the interactions (McAvoy and Hauert, 2015, Ogbo et al., 2021,
56 Tuyls et al., 2018). Examples include conflict resolution (Selten, 1978, Smidt, 2020), technol-
57 ogy adoption by firms (Ogbo et al., 2021), and multiparty resource allocations (Chevaleyre et al.,
58 2005, Lerat et al., 2013), where participants might have different roles (e.g. proposers/dictators
59 vs responders) or bargaining power in the decision making process. In this asymmetric setting,
60 the external decision maker might need to take into account the difference among players' un-
61 derlying characteristics, such as their roles in the interactions, in order to optimise the cost and
62 the level of desired behaviour. In particular, we might ask, is it enough to target a subset of the
63 roles to already achieve a sufficiently good outcome, since collecting information about all the
64 roles might be (very) costly and time consuming?

65 This paper contributes to advancing the state-of-the-art by studying cost-efficient external
66 interference in a spatial Ultimatum Game (UG), a popular bargaining game for investigating
67 fair decision making in many disciplines, such as economics and AI/MAS research (Cimpeanu
68 et al., 2021, de Jong and Tuyls, 2011, De Jong et al., 2008, de Melo et al., 2018, Fehr and
69 Schmidt, 1999, Güth et al., 1982, Rauwolf and Bryson, 2018, Santos et al., 2019, Teixeira
70 et al., 2021). While there is concern over which quantitative definition of fairness to use in
71 several disciplines, such as algorithmic fairness in Artificial Intelligence and Machine Learn-
72 ing (Hutchinson and Mitchell, 2019, Mehrabi et al., 2021), the abstract framework of the UG
73 presents an uncontroversial, unique mathematical criterion for measuring fairness. Fair indi-
74 viduals in this setting are those who choose to donate the larger proportion of their individual
75 endowment. Moreover, the target variable and the model score are unrelated, and we test for
76 each possible set of targets using this unique measure of fairness. Similar to previous works
77 (Chen et al., 2015, Cimpeanu et al., 2019, Duong and Han, 2021a, Han et al., 2018, Han and
78 Tran-Thanh, 2018, Wang et al., 2019), we measure cost-efficiency of an interference strategy
79 by its total cost required over time. Furthermore, we will discuss explicitly the requirements for
80 each strategy to be implementable, such as the availability of information about the population
81 and local neighbourhoods, and budget constraints in each time step. In a standard UG, players
82 have two different roles, proposer and receiver (or responder), with different bargaining pow-
83 ers (See Methods in Section 3 for a detailed description of the game). We consider the spatial
84 version of the game (Page et al., 2000) where players are distributed on a network in order to
85 examine how to exploit the roles' asymmetry in both global and local interference strategies
86 (Han et al., 2018).

87 In general, a cost-efficient interference problem consists of solving a bi-objective optimi-
88 sation problem (Han and Tran-Thanh, 2018, Wang et al., 2019), maximising the overall level
89 of desired behaviours (in the long run) while ensuring the total cost spent being within bud-
90 get and/or minimal. The key challenge is that, for evolving dynamical systems such as those
91 in the above-mentioned examples, the system dynamics are shaped by various stochastic and
92 random effects, such as those resulting from behavioural updates and mutation (behavioural ex-
93 ploration) (Rand et al., 2013, Traulsen et al., 2009). With behavioural updates, such as through
94 social learning or reproduction (Nowak, 2006, Sigmund, 2010), undesired behaviours might
95 resurface over time whenever interference was not sufficiently strong in the past. Through mu-
96 tation, these behaviours might do so even when they were extinct. Hence, the external decision
97 maker needs to take into consideration that they will have to repeatedly interfere in the system,
98 in order to sustain the desired behaviour over time. Note however that, for simplicity, previous
99 works have either omitted mutation (Han and Tran-Thanh, 2018, Wang et al., 2019), or assumed
100 that it is infinitely small (for analytical treatment) (Duong and Han, 2021a, Han and Tran-Thanh,
101 2018). Mutation (behavioural exploration), where agents can freely experiment with new be-
102 haviours, is usually non-negligible in real populations and has been shown to play an important
103 role in enabling cooperation in the context of social dilemmas (Antal et al., 2009, Duong and
104 Han, 2019, 2021b, Han et al., 2012, Rand et al., 2013, Traulsen et al., 2009). Thus, the present
105 work will also advance the state-of-the-art in this respect, where we will closely examine how
106 different regimes of mutation, or agents' propensity for behavioural exploration, influence the
107 manner in which external interference should be carried out. Indeed, our results show that when
108 mutation is sufficiently high, in line with those observed in behavioural experiments, optimal
109 interference strategies can be significantly different.

110 In the next section we describe the models and methods in detail, before presenting the
111 results and a final discussion. We also include with this paper a Supplementary Information
112 (SI) that includes additional results to support the robustness of the paper findings described in
113 the main text.

114 2 Models and Methods

115 2.1 Ultimatum Game (UG)

116 Agents' interaction is modelled using the one-shot Ultimatum Game (UG) (Nowak et al., 2000,
 117 Page et al., 2000). In the UG, two players are offered a chance to win a certain sum of money,
 118 normalised to 1, which they must divide between each other. One player is elected proposer,
 119 and suggests how to split the sum, while the other, the receiver (responder) can accept or reject
 120 the deal. If the deal is rejected, neither player receives any part of the initial sum. As in (Nowak
 121 et al., 2000, Page et al., 2000), we assume that a player is equally likely to perform in one of
 122 the roles (proposer or receiver). A player's strategy is defined by a pair of parameters, p and q .
 123 When acting as proposer, the player offers the amount p , whereas in a receiver's role, the player
 124 rejects any offer smaller than q .

125 As we focus in this paper on the effect of having multiple roles on interference decision
 126 making, we consider a baseline UG model where proposers have two possible strategic offers,
 127 a low (L, with $p = l$) and a high one (fair) (H, with $p = h$), where $l < h \in [0, 1]$. On the other
 128 hand, receivers have two options, a low threshold (L, with $q = l$) and a high threshold (H, with
 129 $q = h$). Thus, overall, there are four possible strategies HH, HL, LH and LL (i.e. HL would
 130 denote proposing high and accepting any offers, etc.). Fairness is measured by calculating
 131 what percentage of the population is representative for either the HH or HL strategies (i.e., fair
 132 proposers), and this allows us to have a clear comparison with previous works—in terms of the
 133 level of population fairness achieved—that have studied the evolution of fairness in the UG, see
 134 e.g. (Nowak et al., 2000, Page et al., 2000, Rand et al., 2013). Unlike our work, they did not
 135 study the cost-efficiency of interference strategies for enhancing fairness.

136 Given evidence from several behavioural experiments (Güth et al., 1982, Rand et al., 2013),
 137 in which people (almost) never offered more than half of the sum in UG, we assume $h \leq 0.5$.
 138 Particularly, we set $h = 0.5$ and $l = 0.1$, as shown in (Page et al., 2000). In this scenario, the
 139 strategy LL is evolutionarily stable. We also confirm this result in our simulations, as shown
 140 in Figure S1 in Supplementary Information (SI) and we note that this result is true for several
 141 mutation rates.

142 The payoff for the four strategies HH, HL, LH and LL reads (for row player):

143

	HH	HL	LH	LL
HH	$\frac{1}{2}$	$\frac{1}{2}$	$\frac{1-h}{2}$	$\frac{1-h}{2}$
HL	$\frac{1}{2}$	$\frac{1}{2}$	$\frac{1-h+l}{2}$	$\frac{1-h+l}{2}$
LH	$\frac{h}{2}$	$\frac{1+h-l}{2}$	0	$\frac{1-l}{2}$
LL	$\frac{h}{2}$	$\frac{1+h-l}{2}$	$\frac{l}{2}$	$\frac{1}{2}$

144 2.2 Population structure and dynamics

145 We consider a population of agents or individuals on a square lattice of size $Z = L \times L$ with
146 periodic boundary conditions— a widely adopted population structure in population dynamics
147 and evolutionary games (Szabó and Fath, 2007). We focus our analysis on the efficiency of
148 various interference strategies in spatial settings, adopting an agent-based model directly com-
149 parable with the setup of recent lab experiments on cooperation (Rand et al., 2014). We set
150 $L = 100$ for all our experiments, resulting in a population size $Z = 10^4$. For the baseline
151 results on well-mixed populations (complete graph), we chose a population size $Z = 100$.

Initially each agent is designated as one of the four strategies (i.e. HH, HL, LH, HH), with equal probability. At each time step or generation, each agent plays the UG with its (four) immediate neighbours. In the well-mixed baseline, each agent plays the UG with every other agent in the population. The score for each agent is the sum of the payoffs in these encounters. At the end of each generation an agent A with score f_A chooses to copy the strategy of a randomly selected neighbour agent B with score f_B with a probability given by the Fermi function (i.e. *stochastic update*) (Traulsen et al., 2006):

$$(1 + e^{(f_A - f_B)/K})^{-1},$$

152 where K denotes the amplitude of noise in the imitation process (Szabó and Fath, 2007). Vary-
153 ing K allows us to capture a wide range of update rules and levels of stochasticity, including
154 those used by humans, as measured in lab experiments (Rand et al., 2013, Zisis et al., 2015).
155 In line with previous works and lab experiments (Rand et al., 2013, Szabó and Fath, 2007,
156 Zisis et al., 2015), we set $K = 0.1$ in our simulations. With a given probability μ , this pro-
157 cess is replaced instead by a randomly occurring mutation. A mutation is equivalent to be-
158 havioural exploration, where the individual makes a stochastic decision switch to one of the

159 four available strategies. For each interference strategy, we study four different mutation rates,
 160 for $\mu \in \{10^{-4}, 10^{-3}, 10^{-2}, 2 * 10^{-1}\}$, as well as $\mu = 0$ for the deterministic update. We will
 161 explicitly state the values of mutation rates in all figures' captions. Note that we sometimes
 162 only include a certain subset of these results in the main text for clarity. For a comprehensive
 163 set of results, see SI.

164 Although our analysis below will focus on the stochastic update rule (in order to examine
 165 how stochasticity affects interference, as discussed above), we will also provide results for *de-*
 166 *terministic update* to have a clear comparison with previous works (see e.g. (Han et al., 2018)).
 167 For the deterministic update, an agent's strategy is always changed to that of its highest scoring
 168 neighbour (Nowak and May, 1992, Szabó and Fath, 2007). This is a way of approximating the
 169 stochastic update rule where the stochastic effect is infinitely small, i.e. $K \rightarrow 0$.

170 We simulate this evolutionary process until a stationary state or a cyclic pattern is reached.
 171 Similarly to (Nowak and May, 1992), all the simulations in this work (described in next sec-
 172 tions) converge quickly to such a state. For the sake of a clear and fair comparison, all simula-
 173 tions are run for 500 generations. Moreover, for each simulation, the results are averaged over
 174 the final 25 generations, in order to account for the fluctuations characteristic of these stable
 175 states. Furthermore, to improve accuracy, for each set of parameter values, the final results are
 176 obtained from averaging 30 independent realisations. When shown in figures, the error bars rep-
 177 resent the standard error of the mean between replicates. In the case of the well-mixed baseline
 178 results, which we do not use for direct comparison of cost-efficiency, we run the simulations
 179 for 100 generations and average the results over the final 10 generations. This accounts for the
 180 faster convergence time characteristic of these networks, and the lower population size.

181 Note that in the special case of deterministic update (where we also do not consider mu-
 182 tations), simulations can stop early when the proportion of fair proposers reaches 100%. We
 183 note that when maximum fairness is not reached, investment can still be ongoing beyond 500
 184 generations and thus, the total cost of interference is dependent on the chosen stopping point.
 185 However, our results show that the average investment at the 500 generation mark is never more
 186 than 0.2% of the average total investment, for all types of interference. Thus, this arbitrary
 187 number has a limited effect and should not affect these results qualitatively.

188 2.3 Cost-Efficient Interference in Networks

189 We aim to study how one can efficiently interfere in a structured population to achieve high
 190 levels of fairness while minimising the cost of interference. Naturally, the level of fairness is
 191 measured by the fraction of fair offers in the population (Rand et al., 2013), which is the total
 192 of HH and HL frequencies. An investment decision consists of a cost $\theta > 0$ to the external
 193 decision-making agent/investor, this value θ is added as surplus to the payoff of each suitable
 194 candidate. In order to determine cost-efficiency, we vary the individual investment amount θ for
 195 each proposed interference strategy, and we measure the total accumulated costs to the external
 196 investor. Thus, the most efficient interference schemes will be the ones with the lowest relative
 197 total cost.

198 We examine and compare different approaches of interference to induce fairness, based on
 199 ensuring fairness for either role or both, leading to different desirable behaviours to be targeted

- 200 (i) ensure all proposals are fair, thus investing in HH and HL (**Target: HH, HL**);
- 201 (ii) ensure only fair offers are accepted, thus investing in HH and LH (**Target: HH, LH**);
- 202 (iii) ensure both (i) and (ii), i.e. investing in HH only (**Target: HH**).

203 Moreover, in line with previous works on network interference (Chen et al., 2015, Cim-
 204 peanu et al., 2019, Han and Tran-Thanh, 2018), we will compare global interference strategies
 205 where investments are triggered based on network-wide information, and the local ones where
 206 investments are based on local neighbourhood information.

207 In the *population-based* approach, a decision to invest in desirable behaviours is based on
 208 the current composition of the population. We denote x_f the fraction of individuals in the
 209 population with a desirable behaviour, given a targeting approach at hand, i.e. (i), (ii) or (iii)
 210 as defined above. Namely, investment is made if x_f is less or at most equal to a threshold p_f
 211 (i.e. $x_f \leq p_f$), for $0 \leq p_f \leq 1$. They do not invest otherwise (i.e. $x_f > p_f$). The value
 212 p_f describes how rare the desirable behaviours should be to trigger external support. In the
 213 *neighbourhood-based* approach, a decision to invest is based on the fraction x_f calculated at
 214 local level. Investment happens if the proportion of neighbours of a focal individual with the
 215 desirable behaviours is less or at most equal to a threshold n_f (i.e. $x_f \leq n_f$), for $0 \leq n_f \leq 1$;
 216 otherwise, no investment is made.

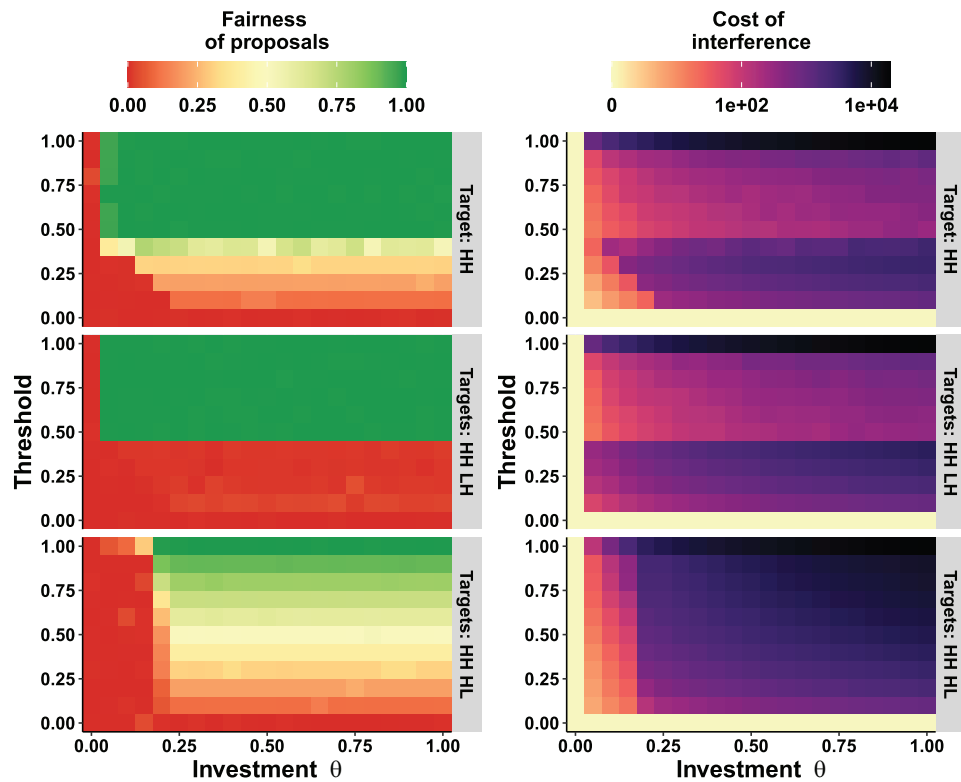


Figure 1. Baseline results after interference in a well-mixed scenario. Average fairness (left) and average cost of interference (right) as a function of the individual endowment θ and the threshold p_f (population-based, well-mixed network, $\mu = 0.01$, stochastic update). Each row represents a different targeting scheme. The cost of interference is shown on a logarithmic scale.

217 **3 Results**

218 When choosing to invest in a population of individuals in an effort to ensure some form of
 219 desirable outcome, an external decision maker must first consider several factors before any
 220 decision is made. Among these, we consider and aim to resolve the questions regarding what
 221 sort of behaviour they should invest in, how large the individual endowment must be, but also
 222 what an investor can do when information about the population or the environment is incom-
 223 plete, or even unknown. As such, we consider that the simplest form of information gathering
 224 evaluates the overall population (in the form of some metrics measuring fairness on average),
 225 as opposed to fine-grained observations on individual neighbourhoods. Likewise, we consider
 226 that ensuring all proposals are fair (i.e. investing in HH or HL) is less demanding on an external
 227 decision-maker than ensuring that only fair offers are accepted (i.e. investing in HH and LH),
 228 which is, in turn, a simpler endeavour than for both the former and latter to be strictly enforced
 229 (choosing to invest in HH only). In this way, we can conceptualise a hierarchy of investment
 230 strategies, in terms of complexity, some of which may simply be impossible for an investor to
 231 follow, merely due to lack of information, funding, or a combination of the two.

232 We consider that there exists a *minimal level of fairness* which the external decision maker is
 233 aiming to enforce in regards to the population's behaviour (Han and Tran-Thanh, 2018), and we
 234 study the least expensive investment strategies for differing preferences of such an acceptable
 235 fairness.

236 **3.1 Population-based results**

237 Firstly, we explore the simplest class of investment strategies, using a macroscopic metric of
 238 the population, measuring average fairness in the whole system (population).

239 **3.1.1 Baseline, Well-Mixed Populations**

240 As the foundation of this analysis, we first introduce a baseline analysis of this interference on
 241 well-mixed (complete graph) populations, in Figure 1. We notice an increase in fairness for all
 242 three different targets, if the threshold for investment is sufficiently high, but there are marked
 243 differences in the cost of interference. Specifically, targeting both fair responses and proposals
 244 (HH), as well as only fair responses (HH LH), reduce the accumulated costs of interference for
 245 the external investor for a broader range of parameters than targeting only fair proposals (HH

246 HL). Furthermore, the threshold for investment is the deciding factor for ensuring high levels of
 247 fairness for all cases. This suggests that if certain levels of fair behaviour are maintained, then
 248 the population will converge to fairness without requiring further investment.

249 Based on the amount of information available to the external decision maker, we confirm
 250 that more information gathering leads to a more flexible investment approach. Respectively,
 251 the strictest approach (targeting HH only) leads to the highest levels of fairness with lowest
 252 accumulated costs, followed by ensuring fair responses, and, lastly, promoting fair proposals.
 253 Targeting both roles or only fair responses produce almost indistinguishable results if the chosen
 254 threshold is sufficiently high ($p_f \gtrsim 40\%$), whereas only targeting fair proposers is very costly
 255 regardless of minimal fairness requirements. These results show that fair responders drive the
 256 dynamics of the system in the well-mixed scenario, and they should be targeted correspondingly
 257 by an external decision maker.

258 3.1.2 Structured Populations

259 We now consider that the population is structured and that individuals interact only with their
 260 neighbours. Figure 2 shows the results for different population-based interference scheme and
 261 clearly demonstrates the difference between the three targets for investment. We would like to
 262 point out the higher levels of fairness obtained using the HH targeting scheme, especially for
 263 a lower threshold p_f . We also notice an increase in the threshold for investment p_f in order to
 264 achieve similar levels of fairness. When it comes to the accumulated cost of interference, we
 265 see that HH is the most cost-effective solution, due to the previously perceived lower threshold
 266 required to maintain fairness.

267

268 Figure 3 further exemplifies the finding that targeting HH is the optimal scheme for population-
 269 based interference. Each row (portraying the different targeting schemes), drifts further away
 270 from the cost-optimal bottom left. As the threshold increases, so does the total cost, so the
 271 regions of high fairness for a lower threshold observed in Figure 2 coincide with the maximal
 272 savings (while still achieving desired levels of fairness).

273 Table 1 shows the most cost-efficient schemes for ensuring specific standards of fairness
 274 when only a population-based approach is possible, under differing rates of mutation (μ). We
 275 observe a definitive bias towards the most complex investment scheme (i.e. targeting HH play-
 276 ers), which reiterates our previous observation. We note that, in order to maintain a desired level

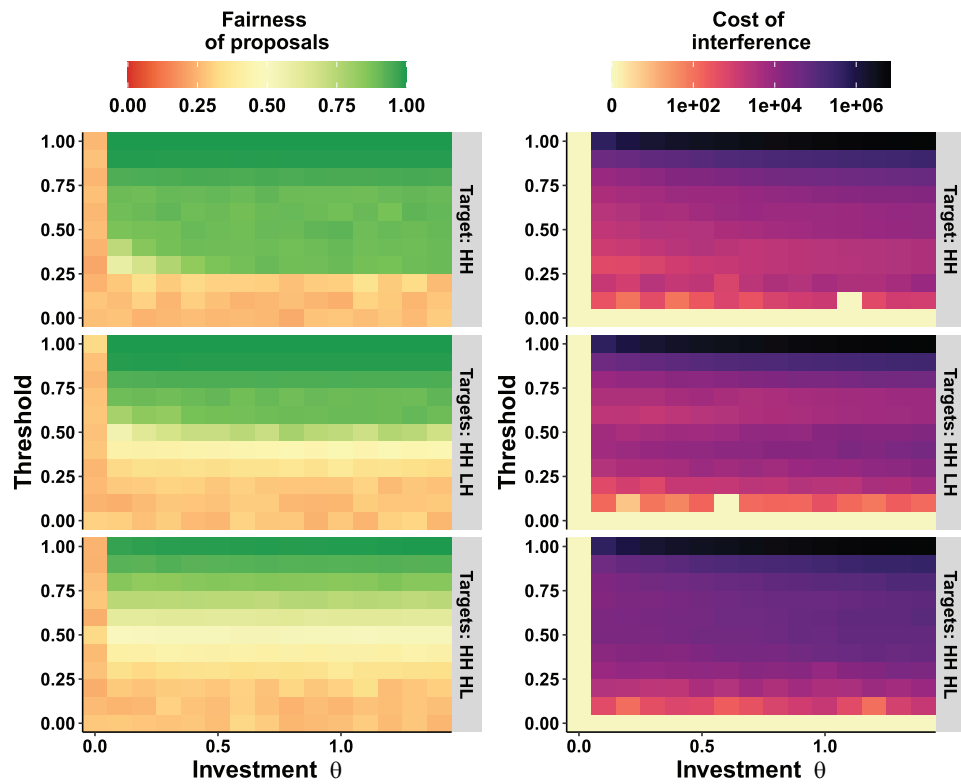


Figure 2. Average fairness (left) and average cost of interference (right) as a function of the individual endowment θ and the threshold p_f (population-based, $\mu = 0.01$, stochastic update). Each row represents a different targeting scheme. The cost of interference is shown on a logarithmic scale.

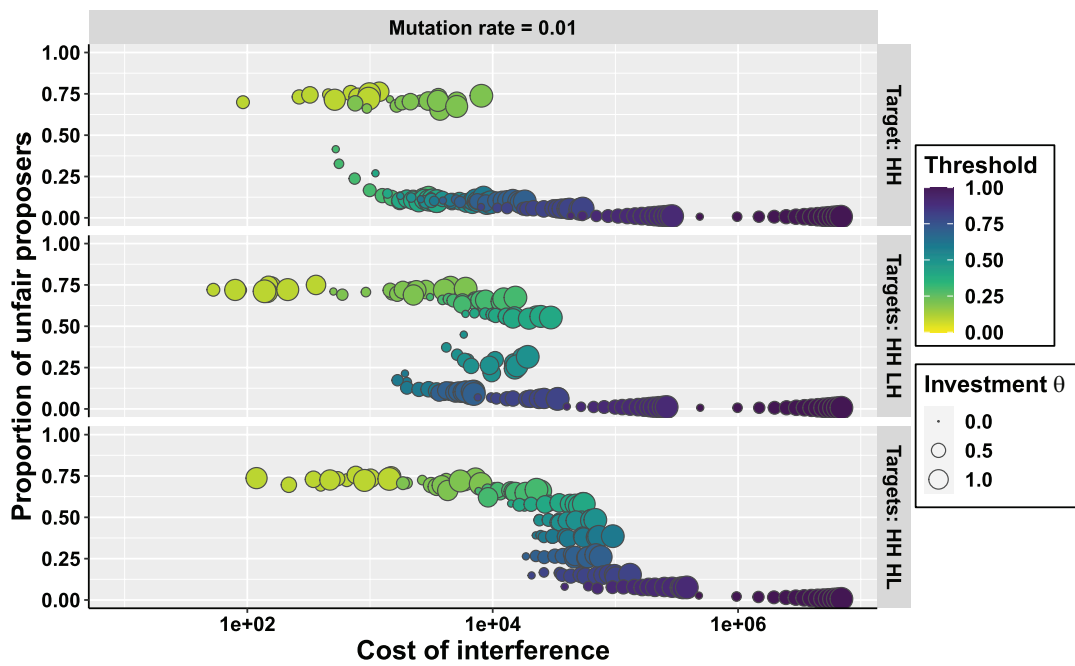


Figure 3. Proportion of unfair proposers as a function of average cost of interference for different targeting scheme (population-based, $\mu = 0.01$, stochastic update). The size and colour of the circles correspond to investment amount and threshold of investment, respectively. We note that the most desirable outcomes are closest to the origin.

Table 1. Most cost-efficient scheme to reach a minimum fairness of proposals for different mutation rates (population-based, stochastic update). There exist no schemes which satisfy the higher minimum fairness requirements in the case of very high mutation rate, written as ‘–’ in the table.

Mutation rate	Minimum fairness	Target	Threshold	θ	Cost (mean \pm 1.96 se)
10^{-4}	75%	HH	0.3	0.1	530 ± 5
10^{-4}	90%	HH	0.3	0.1	530 ± 5
10^{-4}	99%	HH	0.3	0.4	999 ± 7.6
10^{-2}	75%	HH	0.3	0.3	750 ± 5.4
10^{-2}	90%	HH	0.3	0.7	1747 ± 11.2
10^{-2}	99%	HH	1	0.1	487514 ± 93.6
0.2	75%	HH	0.6	0.2	358089 ± 650
0.2	90%	–	–	–	–
0.2	99%	–	–	–	–

277 of fairness, an external decision maker must increase the threshold at which they resume their
 278 investment, but also the individual endowment (θ). It becomes increasingly difficult to maintain
 279 standards of fairness when the population is exposed to high degrees of behavioural exploration
 280 and this naturally attracts an increase in total cost for the investor. We report similar figures for
 281 other values of μ in Figures S2, S3, S4, in SI.

282 Moreover, we observe an increase in fairness for all schemes of interference, across most
 283 values of individual endowment θ , which bodes well when the external decision maker pos-
 284 sesses limited knowledge. If reducing cost is not the main objective, fairness can be maintained
 285 using any targeting scheme (i.e. any relevant observations made about the population), by in-
 286 creasing the minimum threshold p_f .

287 When the external decision maker is limited to the macroscopic metrics associated with
 288 population-based interference, interference is characterised by its strictness. To elaborate, in-
 289 formation gathering should be the main goal for the investor, as ensuring that proposals and
 290 responses are simultaneously fair (i.e. targeting HH) is the optimal outcome. In this way, the
 291 minimum threshold can be kept low, reducing the accumulated cost. These findings are robust
 292 when compared to well-mixed populations, although it is easier for an investor to maintain fair-
 293 ness in the case of structured populations, when targeting fair proposers is the only option for
 294 investment.

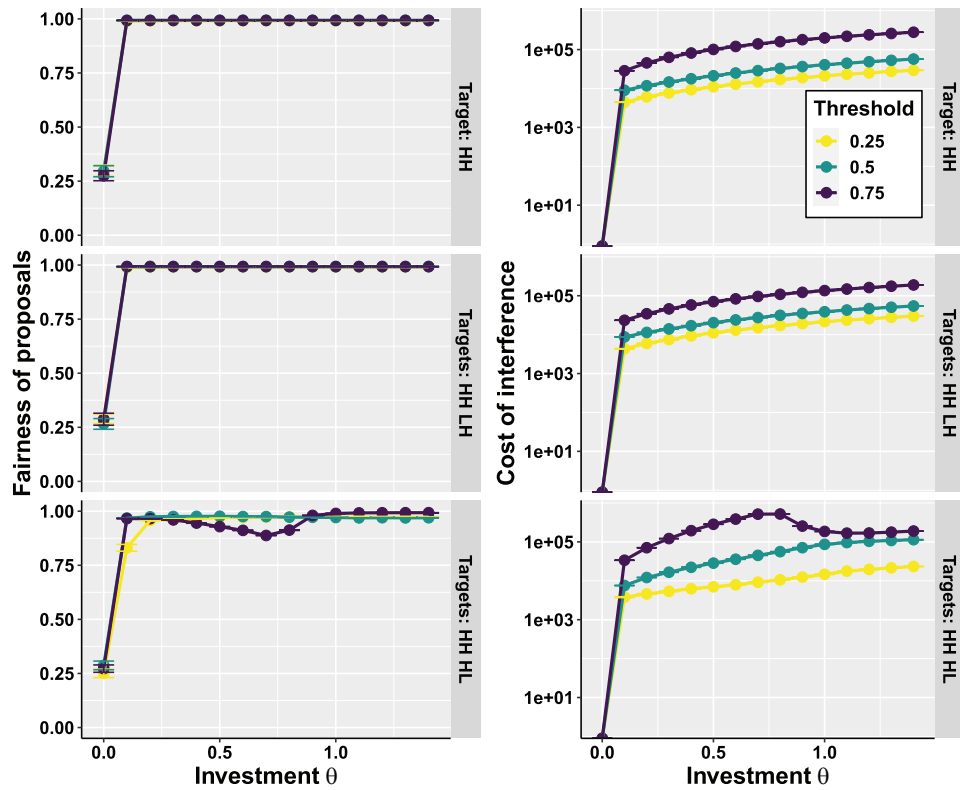


Figure 4. Average fairness (left) and average cost of interference (right) as a function of θ and threshold n_f (neighbourhood-based, $\mu = 0.01$, stochastic update). Each row represents a different targeting scheme. The cost of interference is on a logarithmic scale.

295 **3.2 Neighbourhood-based interference**

296 Previous works on the PD have shown that the greatest gains in cooperation (while maintaining
297 a minimal investment cost) require very detailed observations of individual neighbourhoods,
298 coupled with overly strict investment schemes (Cimpeanu et al., 2019, Han et al., 2018, Han
299 and Tran-Thanh, 2018). In order to decipher whether or not these findings hold for the spatial
300 ultimatum game, we study the outcome when an investor can perceive fairness at the local level.

301 Figure 4 reports the relationship between gains in fairness and increases in cost for an ex-
302 ternal investor, with diverse targets for receiving investments. We observe that fairness is more
303 easily achieved than in population based interference, with only a very low investment required
304 to sustain a majority of fair proposals. Further investment increases the cost of interference, but
305 only slightly. If different thresholds result in fairness, Figure 4 shows that a threshold of 25% is
306 the most cost-efficient. Similarly to population-based interference, the external decision maker
307 should invest only when a large proportion of unfair individuals are present to limit the cost of
308 investment. Finally, there are no significant differences between targeting schemes.

309 Similarly to our findings using a population-based approach, we observe that the more pro-
310 hibitive option, HH, is also the most cost-effective. On the other hand, high fairness can be
311 achieved in all three cases for the same values of endowment. Ensuring that all proposals are
312 fair (thus investing in HH and HL), can lead to an increase in cost of interference, and a de-
313 crease in fairness gains (relative to the other two interference strategies). While all investment
314 schemes evidently succeed in promoting the evolution of fairness, only ensuring the equitable
315 proposals is not as reliable as encouraging discerning responses to offers or both. We note that
316 this effect can only be seen when the threshold for investment is very high (i.e. an investor only
317 invests in neighbourhoods with three or more fair proposers). As discussed earlier, investing in
318 neighbourhoods with at most one fair agent and not investing otherwise, solves this dilemma.

319 Markedly, it is not effective to invest in neighbourhoods with a high percentage of fair
320 proposals. These results point to a key observation, that it is more important to invest in fair
321 proposers when there are few of them in a specific neighbourhood. In this sense, the lonely
322 fair individuals require aid in otherwise competitive, unjust entourages. This result can further
323 be seen in Figure 4. By being very selective with which neighbourhoods the external investor
324 chooses to invest in (i.e. only choosing very fair neighbourhoods), they inadvertently produce
325 a much higher final cost to their own selves. An external decision-maker would then unwit-
326 tingly keep investing in fair proposals ad infinitum because fairness is eventually reached in
327 the ultimatum game, even when individual endowment is relatively low. It is clear, therefore,

328 that to reduce potential costs, only players in unfair groups should be eligible for investment.
 329 Therefore, the defining characteristic of neighbourhood-based interference is the low threshold
 330 for investment (25%).

Table 2. Most cost-efficient scheme to reach a minimum fairness of proposals for different mutation rates (neighbourhood-based, stochastic update). There exist no schemes which satisfy the higher minimum fairness requirements in the case of very high mutation rates, written as ‘-’ in the table.

Mutation rate	Minimum fairness	Target	Threshold	θ	Cost (mean \pm 1.96 se)
10^{-4}	75%	HH	0.25	0.1	1395 ± 36.9
10^{-4}	90%	HH	0.25	0.1	1395 ± 36.9
10^{-4}	99%	HH	0.25	0.1	1395 ± 36.9
10^{-2}	75%	HH HL	0.25	0.1	3794 ± 200.1
10^{-2}	90%	HH LH	0.25	0.1	4352 ± 56.2
10^{-2}	99%	HH LH	0.25	0.2	5957 ± 60.7
0.2	75%	HH	0.25	0.4	150777 ± 121.5
0.2	90%	-	-	-	-
0.2	99%	-	-	-	-

331

332 By varying minimal fairness requirements and rates of mutation, we can gain further insight
 333 into which investment strategies are the most robust and cost-effective. Table 2 highlights some
 334 surprising findings. We see that neighbourhood based interference can result in a higher total
 335 cost than the optimal population-based interference schemes (see Table 1). Previous work has
 336 shown that more specific and restrictive intervention schemes are more effective in the PD
 337 (Cimpeanu et al., 2019, Han et al., 2018), but by being able to target different roles in the
 338 Ultimatum game, these differences can be mitigated. Furthermore, mutation rate serves as an
 339 equaliser between the investment targets, and we observe that less specific schemes (HH & HL
 340 and HH & LH) are the most cost-efficient options. We note that the differences between results
 341 are small enough that different runs could yield any outcome in the case of high or intermediate
 342 mutation rates. The lack of significant variability among the distinct targeting schemes contrasts
 343 strongly with the findings on the PD (Cimpeanu et al., 2019, Han and Tran-Thanh, 2018). We
 344 report similar figures for other values of μ in Figures S6, S7, S8, in SI.

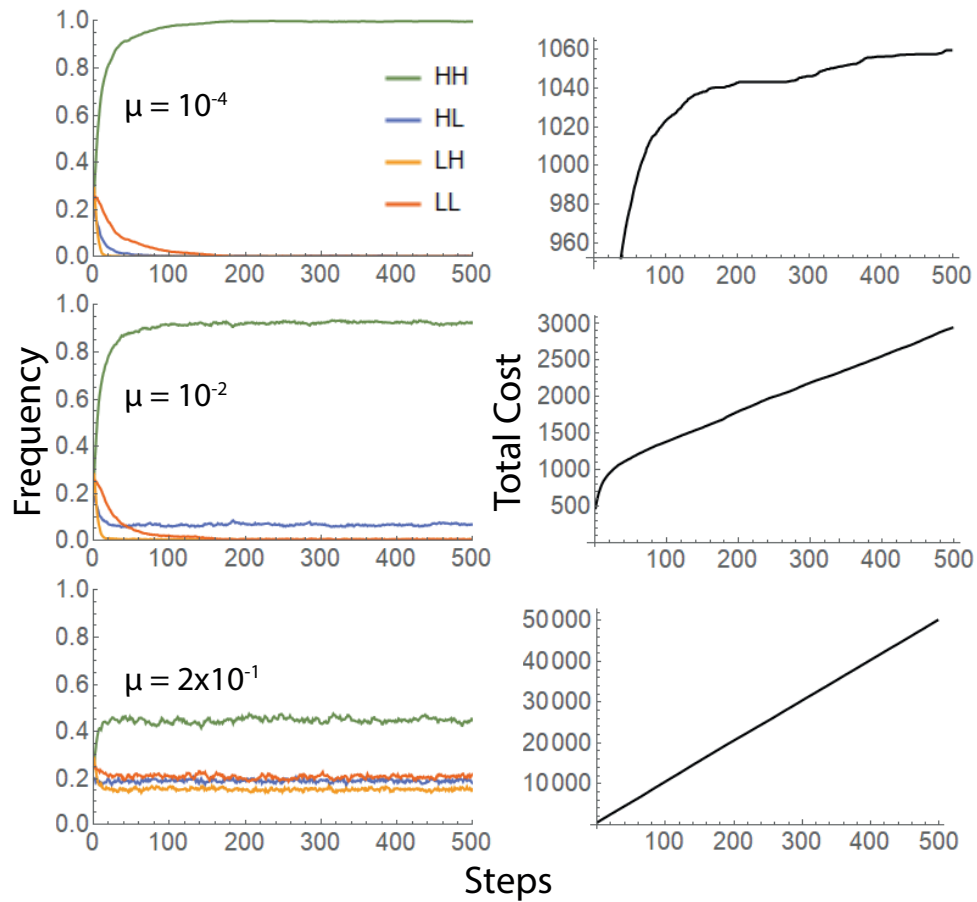


Figure 5. Typical runs showing the evolution of fairness and the associated total cost of interference for various mutation rates (top row $\mu = 10^{-4}$, middle row $\mu = 10^{-2}$, bottom row $\mu = 2 * 10^{-1}$; neighbourhood-based, stochastic update). Parameters: $n_f = 0.25$, $\theta = 0.1$, $Target = HH$. The choice of parameter values was motivated by selecting the optimal solutions in Table 2.

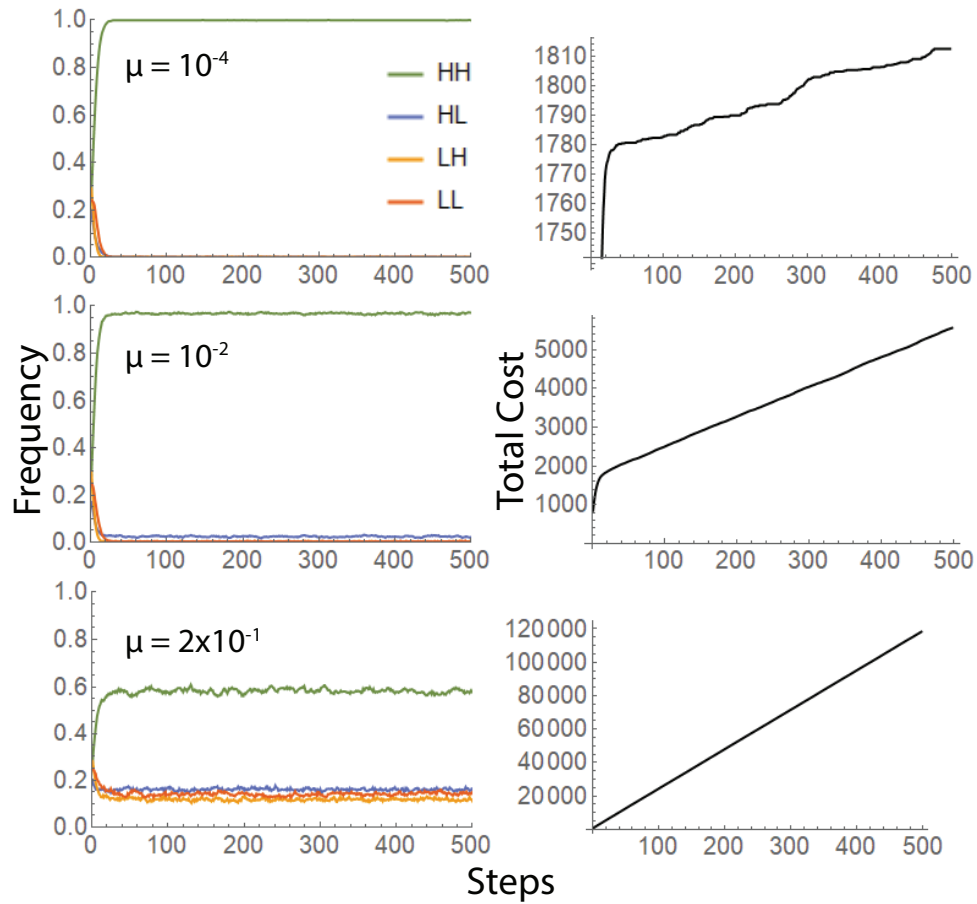


Figure 6. Typical runs showing the evolution of fairness and the associated total cost of interference for various mutation rates (top row $\mu = 10^{-4}$, middle row $\mu = 10^{-2}$, bottom row $\mu = 2 * 10^{-1}$; neighbourhood-based, stochastic update). Parameters: $n_f = 0.5$, $\theta = 0.1$, $Target = HH$. The choice of parameter values was motivated by selecting the optimal solutions in Table 2.

345 **3.3 Evolution of strategies over time**

346 We make use of the optimal parameter values identified in Tables 1 and 2 to explore the evo-
 347 lution of fairness over time for all the strategies in the population, as well as any associated
 348 accumulated costs. Through this analysis, we clarify some of the dynamics differentiating the
 349 different decisions for investment, as well as the effects of varying mutation rates upon the
 350 outcomes and the options available to investors.

351 The effects of mutation on the optimality of different interference schemes can be seen
 352 in Figure 5. As the mutation rate (μ) increases, the capacity of maintaining a threshold of
 353 fairness decreases (as also seen in Table 2). An external investor must increase their individual
 354 investment amount in order to meet these new demands set out by the increased mutation rates,
 355 and by doing so they can maintain fairness levels to a respectable standard.

356 To better highlight the sharp increases in the cost associated with the non-optimal threshold
 357 (i.e. when it is greater than 25%) for neighbourhood-based interference, we show such typical
 358 runs for varying mutation rates for the 50% threshold in Figure 6. When comparing Figures
 359 5 and 6, we note the relative differences in total accumulated costs attributed to the choice of
 360 the threshold for investment n_f . We also note that increasing rates of behavioural exploration
 361 (mutation) amplifies this discrepancy.

362 We show how less specific interference strategies, which require less information gather-
 363 ing, can be effective in facilitating the evolution of fairness, when local monitoring is possible
 364 (Figure 8). Promoting fair proposals may often not be sufficient for low individual investment
 365 budgets (which are also the optimal solution) — in such cases fairness does not evolve. This
 366 occurs due to the inability of indiscriminate fair proposers to protect themselves against un-
 367 fair proposers. Investing in fair proposers, in this case, artificially protects them against very
 368 competitive selection pressures.

369 Figure 7 showcases how different mutation rates call for different approaches to interfer-
 370 ence. As shown previously, optimal interference strategies vary according to the mutation rate.
 371 We point out the three different cases in which an investor might find themselves in. First, when
 372 few initial rounds of investment are enough for the system to converge and stabilise to a desired
 373 state. Second, an investor might be required to reinvest when the population tends to revert back
 374 to its initial condition. Lastly, constant investment is required to maintain a desired level of fair-
 375 ness, with the total cost skyrocketing accordingly. To some extent, a fair population can better
 376 deal with unfair invaders and this explains the need for a sufficiently high initial investment
 377 when mutation rates increase.

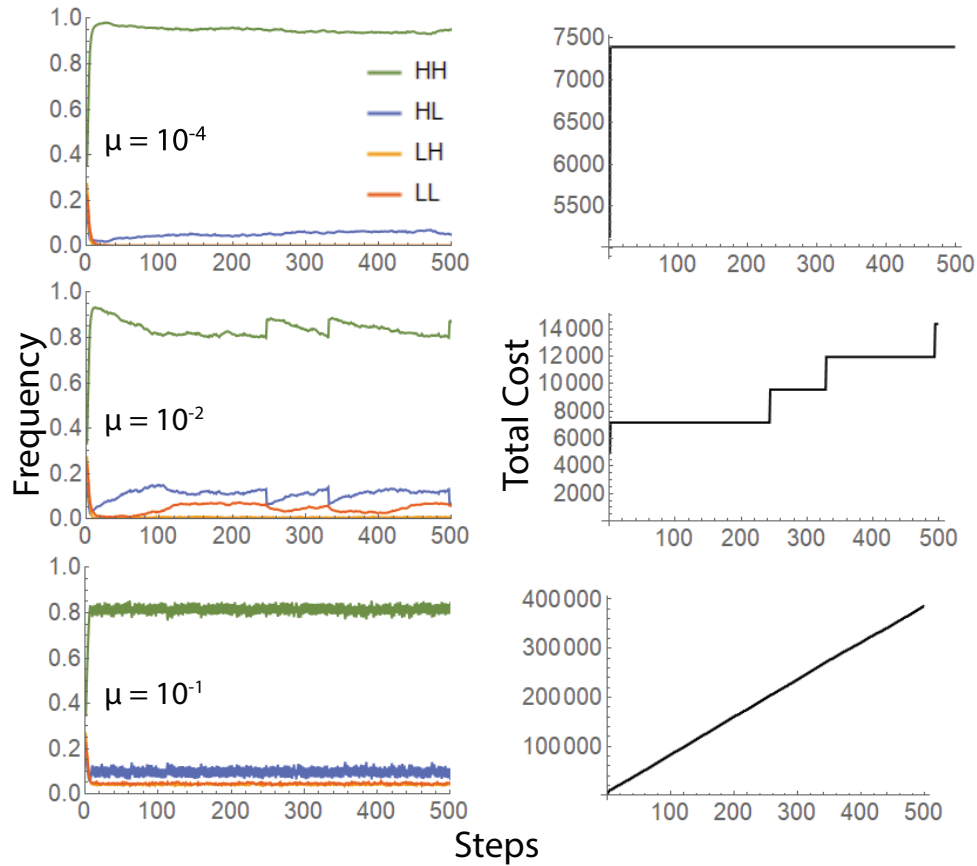


Figure 7. Typical runs showing the evolution of fairness and the associated total cost of interference for various mutation rates (top row $\mu = 10^{-4}$, middle row $\mu = 10^{-2}$, bottom row $\mu = 2 * 10^{-1}$; population-based, stochastic update). Higher mutation rates leads to an increasing need for interference over time. Parameters: $p_f = 0.8$, $\theta = 0.3$, $Target = HH$. The choice of parameter values was motivated by selecting the optimal solutions in Table 1.

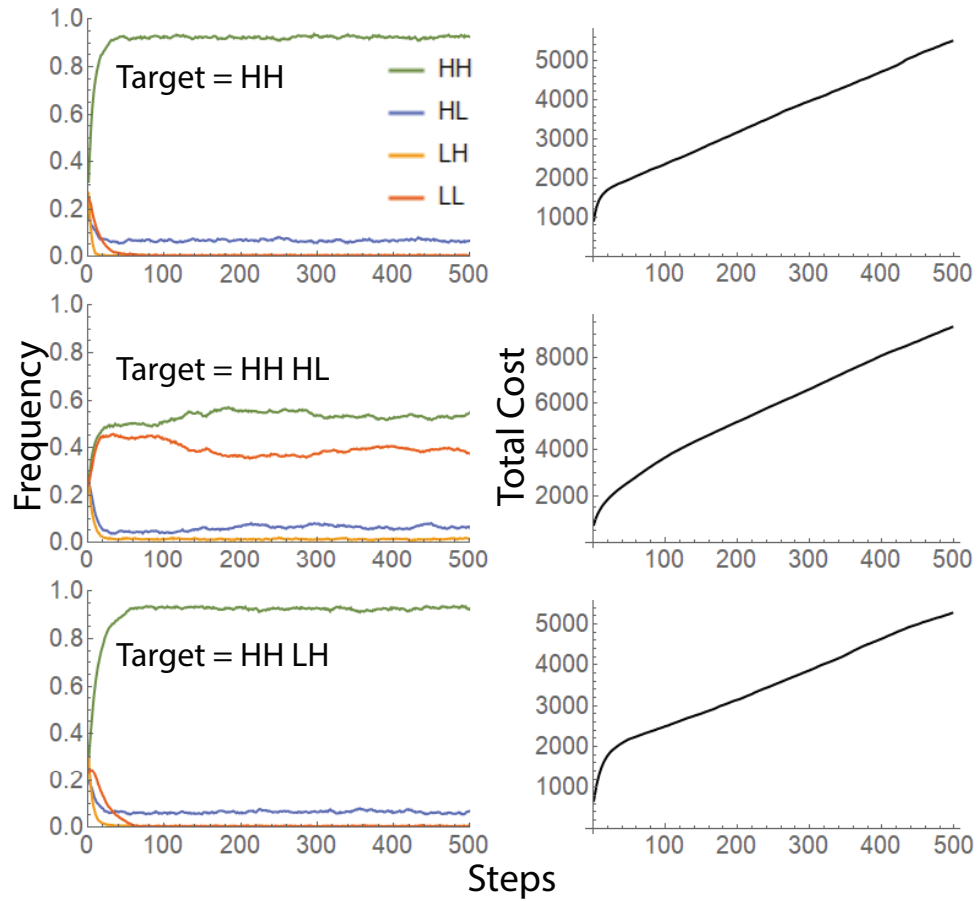


Figure 8. Typical runs showing the evolution of fairness and the associated total cost of interference for various targeting schemes (neighbourhood-based, stochastic update). Parameters: $n_f = 0.25$, $\theta = 0.2$, $\mu = 10^{-2}$. The choice of parameter values was motivated by selecting the optimal solutions in Table 2.

378 Finally, behavioural exploration motivates the manner or strength (in terms of individual
 379 endowment) of any initial efforts to moderate unfair behaviour. Figures 5, 6 and 7 show that the
 380 increase in cost is linear and ever-growing for high mutation-rates and gradually sharper at the
 381 beginning for lower mutation, eventually plateauing when the population is exposed to little or
 382 no behavioural exploration.

Table 3. Most cost-efficient population-based scheme (deterministic update) to reach a minimum fairness of proposals.

Minimum fairness	Target	Threshold	θ	Cost (mean \pm 1.96 se)
75%	HH	0.5	0.5	1251 \pm 10.8
90%	HH	0.6	0.9	2228 \pm 22.6
99%	HH	0.9	1.1	5488 \pm 22.9

383

Table 4. Most cost-efficient neighbourhood-based scheme (deterministic update) to reach a minimum fairness of proposals.

Minimum fairness	Target	Threshold	θ	Cost (mean \pm 1.96 se)
75%	HH	0.25	0.8	2146 \pm 56.3
90%	HH	0.25	0.8	2146 \pm 56.3
99%	HH	0.25	1	2513 \pm 16.4

384

385 3.4 Deterministic update

386 The results and findings reported so far were based on the stochastic update rule. We now
 387 take a step back and consider whether our findings would still hold for the deterministic rule
 388 (see again Methods section). It is not only for the sake of a direct comparison with a previous
 389 analysis reported in (Han et al., 2018), where cost-efficient interference was studied for the
 390 spatial PD in a deterministic setting (with no mutation). It would also allow us to examine if
 391 the findings above would remain robust for the deterministic update, a popular approximation
 392 for rare stochastic effect (or infinite intensity of selection) that is regularly used in the literature

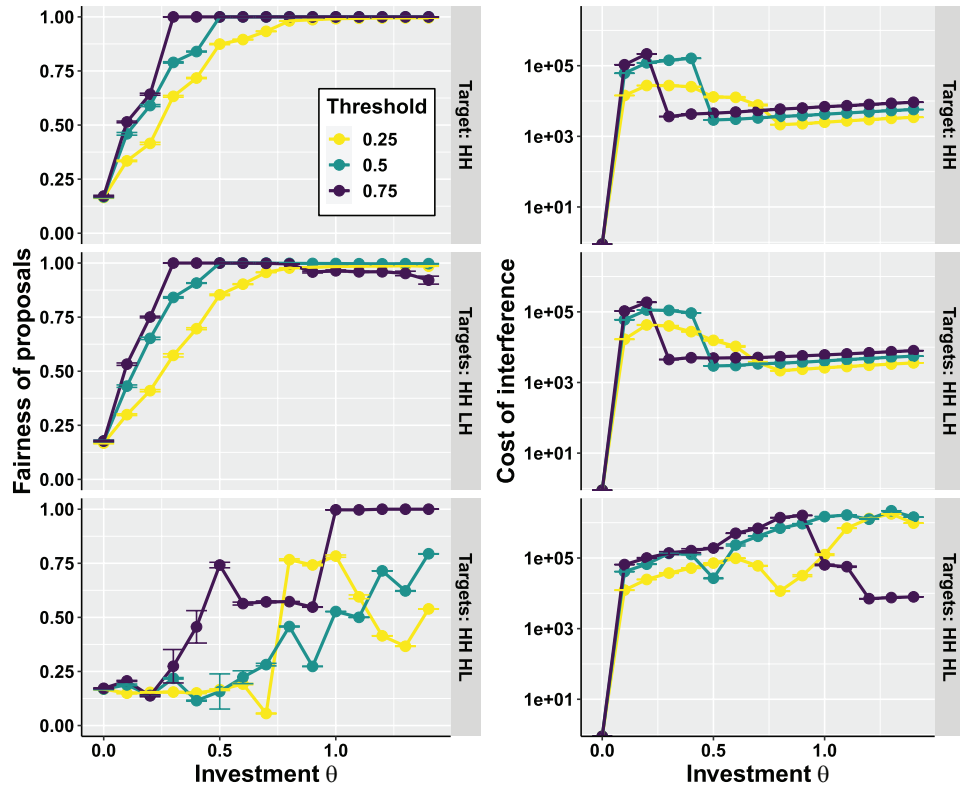


Figure 9. Average fairness (left) and average cost of interference (right) as a function of θ and threshold n_f (neighbourhood-based, deterministic update). Each row represents a different targeting scheme. The cost of interference is on a logarithmic scale for clarity.

393 (Szabó and Fáth, 2007). Tables 3 and 4 report results for the optimal interference strategies,
 394 in population-based (for a full report, see Figure S9, in SI) and neighbourhood-based schemes,
 395 respectively. We observe that for both schemes, targeting HH is always the best option. This
 396 is the same as the stochastic approach for population-based schemes but different from the
 397 neighbourhood-based ones. However, for the latter ones, the optimal threshold $n_f = 0.25$
 398 remains the same as in the case of stochastic update, see also Figure 9. This is in stark contrast
 399 with the PD results where $n_f = 0.75$ was always the optimal choice.

400 4 Discussion

401 In summary, this paper has advanced the state of the art of the literature on external interference
 402 in dynamical systems, or populations of self-interested individuals, in two main respects: i) we
 403 have addressed an asymmetric interaction setting, in the form of the Ultimatum game, where
 404 players have different roles in the interaction. We have shown that it is crucial to consider the
 405 roles' asymmetry to provide cost-efficient investment strategies. This important analysis was
 406 not possible in previous works where symmetric games were studied (Chen and Perc, 2014,
 407 Chen et al., 2015, Cimpanu et al., 2019, Han et al., 2018, Han and Tran-Thanh, 2018, Wang
 408 et al., 2019); ii) we have incorporated realistic levels of mutation or behavioural exploration in
 409 our analysis and have shown that they strongly affect the manner in which interference should be
 410 carried out. Previous works have always omitted mutation or assumed that it is infinitely small,
 411 thereby being unable to address this important issue for real-world populations and applications.

412 We have identified several key features that are required for a cost-effective interference
 413 scheme. On the one hand, population-based schemes are characterised by the need of extensive
 414 information gathering about both roles, as targeting HH always leads to the optimal strategy.
 415 On the other hand, neighbourhood-based schemes are characterised by their flexibility, where
 416 the optimal strategy always entails that investment is only made when there is at most one
 417 player with the desirable behaviour in the neighbourhood (i.e. no investment should be made
 418 when there is a half or larger fraction of such behaviour the neighbourhood). Our findings stand
 419 out in stark contrast with previous works on cooperation dilemmas, where both population and
 420 neighbourhood-based schemes require a highly strict investment approach.

421 The ultimatum game has been widely studied, whether with theoretical models (see Debove
 422 et al. 2016 for a review) or behavioural experiments (see Güth and Kocher 2014 for a review).
 423 The main motivation arises from the gap between theoretical predictions, in which rational indi-

424 viduals keep most of the endowment, and in which responders accept any positive proposition,
425 however small it may be, and experiments, in which individuals propose 40% to 50% of their
426 endowment (and often get punished if they propose less) (Güth et al., 1982). That said, previ-
427 ous works have investigated how fairness can evolve in models of the ultimatum game, wherein
428 several mechanisms promoting the emergence of fairness have been identified. We note that
429 we align our definition of fairness with these previous works, where generous proposers are
430 deemed as fair, regardless of their behaviour when acting in the role of the responder.

431 Among others, Nowak et al. (2000) have studied the evolution of fairness in the Ultimatum
432 game under indirect reciprocity, i.e. when players can observe others' interactions and have in-
433 formation about the reputations of others. We do not rely on reputation building mechanisms, as
434 the role of this mechanism can be limited in large groups, where one-shot interactions between
435 strangers are common. Page et al. (2000), Sinatra et al. (2009) have developed spatial models
436 of the ultimatum game, where interactions happen only between neighbours. It has been shown
437 that a spatial structure can promote the emergence of higher levels of fairness, but an equal split
438 between proposers and responders is yet to be reached. The model developed in the present
439 work has also considered a spatial model because (i) it captures an essential feature of many
440 real-world networks of contacts (Barabasi, 2014), and (ii) it allows us to explore the effects
441 of interference localised in particular neighbourhoods, which has been shown to be more cost
442 efficient (Han et al., 2018). Furthermore, Rand et al. (2013) have shown that even if the popu-
443 lation is well-mixed, fairness can result directly from the effect of randomness due to mutation
444 and stochastic strategy updates. Intuitively, the uncertainty in the responders' choice forces the
445 proposers to offer a high proposal so as to avoid rejection. We show that these stochastic fac-
446 tors also strongly influence the manner in which external interference may be performed while
447 maintaining cost-efficiency. Moreover, the Ultimatum game has also been used to study how
448 fairness can emerge in a hybrid population of human-agent interactions (de Melo et al., 2018,
449 Santos et al., 2019). These works, however, do not consider external interference.

450 The problem of how to externally influence a system of multiple interacting agents to
451 achieve a certain desired behaviour has been of significant interest in mechanism design, net-
452 work theory and control theory literature. For example, how to maximise influence in net-
453 works has been studied in (Bloembergen et al., 2014, Riehl and Cao, 2016, Wilder et al., 2018).
454 Moreover, Endriss et al. (2011) have investigated how to tax games so as incentivise certain
455 behaviours at system equilibrium ; while Wooldridge (2012) has presented potential ways to
456 manipulate games in order to achieve desired behaviours. These works, however, assume that

457 the decision-maker possesses full control of the agents within the systems. With our approach,
458 the decision-makers have little or no direct control on the agents' behaviour, so they can rely
459 only on rewarding schemes and their effects as ways of motivating the evolution of fairness.
460 It is noteworthy that these works do not focus on the cost efficiency problem, whereas cost
461 optimisation is one of our main goals.

462 Although the main focus of our discussion centres around structured populations in the
463 four-neighbour lattice, our results show that similar trends can be observed in well-mixed pop-
464 ulations, as well. When we consider population-based interference that can be enacted both on
465 lattice and well-mixed populations, our observations are robust across network types, and we
466 see that targeting both roles leads to the most cost-effective path towards fairness. In order to
467 determine whether or not spatial heterogeneity plays a key role when multiple roles can be tar-
468 geted by an external investor, future works will need to also consider more complex structures,
469 such as scale-free or multiplex networks.

470 Our future work will examine other asymmetric games with multiple roles, such as the trust
471 and anticipation games, where the bargaining nature is different from the ultimatum game ([Gut,](#)
472 [2009](#), [Han et al., 2021](#), [Rauwolf and Bryson, 2018](#), [Zisis et al., 2015](#)), to see how this bargaining
473 factor might affect the way interference needs to be made. We are also interested in how dif-
474 ferent network structures influence the interference strategies in asymmetric interactions, which
475 has been studied for symmetric games ([Cimpeanu et al., 2019](#)).

476 **Acknowledgements**

477 Some preliminary results from this work was previously published as an extended abstract in
478 AAMAS 2021 conference, pages 1840–1842, see ref. ([Cimpeanu et al., 2021](#)). T.C., C.P. and
479 T.A.H. were supported by Future of Life Institute grant RFP2-154. T.A.H. is also supported by
480 a Leverhulme Research Fellowship (RF-2020-603/9).

481 **References**

- 482 Airiau, S., Sen, S., and Villatoro, D. (2014). Emergence of conventions through social learning.
483 *Autonomous Agents and Multi-Agent Systems*, 28(5):779–804.
- 484 Andras, P., Esterle, L., Guckert, M., Han, T. A., Lewis, P. R., Milanovic, K., Payne, T., Perret,

- 485 C., Pitt, J., Powers, S. T., Urquhart, N., and Wells, S. (2018). Trusting Intelligent Machines:
486 Deepening Trust Within Socio-Technical Systems. *IEEE Technology and Society Magazine*,
487 37(4):76–83.
- 488 Antal, T., Traulsen, A., Ohtsuki, H., Tarnita, C. E., and Nowak, M. A. (2009). Mutation-
489 selection equilibrium in games with multiple strategies. *Journal of theoretical biology*,
490 258(4):614–622.
- 491 Barabasi, A.-L. (2014). *Linked-how Everything is Connected to Everything Else and what it*
492 *Means F*. Perseus Books Group.
- 493 Bloembergen, D., Sahraei, B. R., Bou-Ammar, H., Tuyls, K., and Weiss, G. (2014). Influencing
494 social networks: An optimal control study. In *ECAI*, volume 14, pages 105–110.
- 495 Chen, X. and Perc, M. (2014). Optimal distribution of incentives for public cooperation in
496 heterogeneous interaction environments. *Frontiers in behavioral neuroscience*, 8:248.
- 497 Chen, X., Sasaki, T., Brännström, Å., and Dieckmann, U. (2015). First carrot, then stick: how
498 the adaptive hybridization of incentives promotes cooperation. *Journal of the royal society*
499 *interface*, 12(102):20140935.
- 500 Chevaleyre, Y., Dunne, P. E., Endriss, U., Lang, J., Lemaitre, M., Maudet, N., Padget, J., Phelps,
501 S., Rodríguez-Aguilar, J. A., and Sousa, P. (2005). Issues in multiagent resource allocation.
- 502 Cimpéanu, T., Han, T. A., and Santos, F. C. (2019). Exogenous rewards for promoting cooper-
503 ation in scale-free networks. In *Artificial Life Conference Proceedings*, pages 316–323. MIT
504 Press.
- 505 Cimpéanu, T., Perret, C., and Han, T. A. (2021). Promoting fair proposers, fair responders
506 or both? cost-efficient interference in the spatial ultimatum game. In *In Proc. of the 20th*
507 *International Conference on Autonomous Agents and Multiagent Systems (AAMAS 2021)*,
508 pages 1480–1482.
- 509 de Jong, S. and Tuyls, K. (2011). Human-inspired computational fairness. *Autonomous Agents*
510 *and Multi-Agent Systems*, 22(1):103–126.
- 511 De Jong, S., Uyttendaele, S., and Tuyls, K. (2008). Learning to reach agreement in a continuous
512 ultimatum game. *Journal of Artificial Intelligence Research*, 33:551–574.

- 513 de Melo, C. M., Marsella, S., and Gratch, J. (2018). Social decisions and fairness change when
514 people's interests are represented by autonomous agents. *Autonomous Agents and Multi-*
515 *Agent Systems*, 32(1):163–187.
- 516 Debove, S., Baumard, N., and André, J. B. (2016). Models of the evolution of fairness in the
517 ultimatum game: A review and classification. *Evolution and Human Behavior*, 37(3):245–
518 254.
- 519 Duong, M. H. and Han, T. A. (2019). On equilibrium properties of the replicator–mutator
520 equation in deterministic and random games. *Dynamic Games and Applications*, pages 1–
521 23.
- 522 Duong, M. H. and Han, T. A. (2021a). Cost efficiency of institutional incentives in finite popula-
523 tions. *Proceedings of the Royal Society A: Mathematical, Physical and Engineering Sciences*.
524 (In Press).
- 525 Duong, M. H. and Han, T. A. (2021b). Statistics of the number of equilibria in random social
526 dilemma evolutionary games with mutation. *European Physical Journal B*. (In Press).
- 527 Endriss, U., Kraus, S., Lang, J., and Wooldridge, M. (2011). Incentive engineering for boolean
528 games. *IJCAI '11*, pages 2602–2607.
- 529 Fehr, E. and Schmidt, K. M. (1999). A theory of fairness, competition, and cooperation. *The*
530 *quarterly journal of economics*, 114(3):817–868.
- 531 Gut, A. (2009). *An Intermediate Course in Probability*. Springer Publishing Company, Incor-
532 porated, 2nd edition.
- 533 Güth, W. and Kocher, M. G. (2014). More than thirty years of ultimatum bargaining experi-
534 ments: Motives, variations, and a survey of the recent literature. *Journal of Economic Be-*
535 *havior and Organization*, 108:396–409.
- 536 Güth, W., Schmittberger, R., and Schwarze, B. (1982). An experimental analysis of ultimatum
537 bargaining. *Journal of economic behavior & organization*, 3(4):367–388.
- 538 Hamilton, W. (1964). The genetical evolution of social behaviour. i. *Journal of Theoretical*
539 *Biology*, 7(1):1 – 16.

- 540 Han, T. A. (2013). *Intention Recognition, Commitments and Their Roles in the Evolution of*
541 *Cooperation: From Artificial Intelligence Techniques to Evolutionary Game Theory Models*,
542 volume 9. Springer SAPERE series.
- 543 Han, T. A., Lynch, S., Tran-Thanh, L., and Santos, F. C. (2018). Fostering cooperation in
544 structured populations through local and global interference strategies. In *IJCAI-ECAI'2018*,
545 pages 289–295.
- 546 Han, T. A., Pereira, L. M., and Santos, F. C. (2012). The emergence of commitments and
547 cooperation. In *Proceedings of the 11th International Conference on Autonomous Agents*
548 *and Multiagent Systems (AAMAS'2012)*, pages 559–566. ACM.
- 549 Han, T. A., Pereira, L. M., Santos, F. C., and Lenaerts, T. (2020). To Regulate or Not: A Social
550 Dynamics Analysis of an Idealised AI Race. *Journal of Artificial Intelligence Research*,
551 69:881–921.
- 552 Han, T. A., Perret, C., and Powers, S. T. (2021). When to (or not to) trust intelligent machines:
553 Insights from an evolutionary game theory analysis of trust in repeated games. *Cognitive*
554 *Systems Research*, 68:111–124.
- 555 Han, T. A. and Tran-Thanh, L. (2018). Cost-effective external interference for promoting the
556 evolution of cooperation. *Scientific reports*, 8(1):1–9.
- 557 Hutchinson, B. and Mitchell, M. (2019). 50 years of test (un)fairness: Lessons for machine
558 learning. In *Proceedings of the Conference on Fairness, Accountability, and Transparency*,
559 FAT* '19, page 49–58, New York, NY, USA. Association for Computing Machinery.
- 560 Lerat, J.-S., Han, T. A., and Lenaerts, T. (2013). Evolution of common-pool resources and
561 social welfare in structured populations. In *IJCAI'2013*, pages 2848–2854.
- 562 Liu, L. and Chen, X. (2020). Evolutionary game dynamics in multiagent systems with prosocial
563 and antisocial exclusion strategies. *Knowledge-Based Systems*, 188:104835.
- 564 Marton-Lafevre, J. et al. (2007). *Peace parks: conservation and conflict resolution*. Mit Press.
- 565 Maynard-Smith, J. (1982). *Evolution and the Theory of Games*. Cambridge University Press,
566 Cambridge.

- 567 McAvoy, A. and Hauert, C. (2015). Asymmetric evolutionary games. *PLoS Comput Biol*,
568 11(8):e1004349.
- 569 Mehrabi, N., Morstatter, F., Saxena, N., Lerman, K., and Galstyan, A. (2021). A survey on bias
570 and fairness in machine learning. *ACM Computing Surveys (CSUR)*, 54(6):1–35.
- 571 Nowak, M. A. (2006). *Evolutionary Dynamics: Exploring the Equations of Life*. Harvard
572 University Press, Cambridge, MA.
- 573 Nowak, M. A. (2012). Evolving cooperation. *Journal of Theoretical Biology*, 299:1–8.
- 574 Nowak, M. A. and May, R. M. (1992). Evolutionary games and spatial chaos. *Nature*,
575 359(6398):826–829.
- 576 Nowak, M. A., Page, K. M., and Sigmund, K. (2000). Fairness versus reason in the ultimatum
577 game. *Science*, 289(5485):1773–1775.
- 578 Ogbo, N. B., Elgarig, A., and Han, T. A. (2021). Evolution of coordination in pairwise and
579 multi-player interactions via prior commitments. *Adaptive Behavior (In Press)*. Preprint
580 arXiv:2009.11727.
- 581 Ohtsuki, H., Hauert, C., Lieberman, E., and Nowak, M. A. (2006). A simple rule for the
582 evolution of cooperation on graphs and social networks. *Nature*, 441(7092):502–505.
- 583 Okada, I. (2020). A review of theoretical studies on indirect reciprocity. *Games*, 11(3):27.
- 584 Page, K. M., Nowak, M. A., and Sigmund, K. (2000). The spatial ultimatum game. *Proceedings
585 of the Royal Society of London. Series B: Biological Sciences*, 267(1458):2177–2182.
- 586 Paiva, A., Santos, F. P., and Santos, F. C. (2018). Engineering pro-sociality with autonomous
587 agents. In *Thirty-second AAAI conference on artificial intelligence*.
- 588 Perc, M., Jordan, J. J., Rand, D. G., Wang, Z., Boccaletti, S., and Szolnoki, A. (2017). Statistical
589 physics of human cooperation. *Physics Reports*, 687:1–51.
- 590 Perc, M. and Szolnoki, A. (2010). Coevolutionary games—a mini review. *BioSystems*,
591 99(2):109–125.
- 592 Rand, D. G., Nowak, M. A., Fowler, J. H., and Christakis, N. A. (2014). Static network structure
593 can stabilize human cooperation. *Proc Natl Acad Sci USA*, 111(48):17093–17098.

- 594 Rand, D. G., Tarnita, C. E., Ohtsuki, H., and Nowak, M. A. (2013). Evolution of fairness in
595 the one-shot anonymous ultimatum game. *Proceedings of the National Academy of Sciences*,
596 110(7):2581–2586.
- 597 Rauwolf, P. and Bryson, J. J. (2018). Expectations of fairness and trust co-evolve in environ-
598 ments of partial information. *Dynamic Games and Applications*, 8(4):891–917.
- 599 Riehl, J. R. and Cao, M. (2016). Towards optimal control of evolutionary games on networks.
600 *IEEE Transactions on Automatic Control*, 62(1):458–462.
- 601 Santos, F. P., Pacheco, J. M., Paiva, A., and Santos, F. C. (2019). Evolution of collective
602 fairness in hybrid populations of humans and agents. In *Proceedings of the AAAI Conference*
603 *on Artificial Intelligence*, volume 33, pages 6146–6153.
- 604 Sasaki, T., Brännström, Å., Dieckmann, U., and Sigmund, K. (2012). The take-it-or-leave-it
605 option allows small penalties to overcome social dilemmas. *Proceedings of the National*
606 *Academy of Sciences*, 109(4):1165–1169.
- 607 Selten, R. (1978). A note on evolutionarily stable strategies in asymmetric animal conflicts.
- 608 Sigmund, K. (2010). *The Calculus of Selfishness*. Princeton University Press.
- 609 Sinatra, R., Iranzo, J., Gomez-Gardenes, J., Floria, L. M., Latora, V., and Moreno, Y. (2009).
610 The ultimatum game in complex networks. *Journal of Statistical Mechanics: Theory and*
611 *Experiment*, 2009(09):P09012.
- 612 Smidt, H. M. (2020). United nations peacekeeping locally: enabling conflict resolution, reduc-
613 ing communal violence. *Journal of Conflict Resolution*, 64(2-3):344–372.
- 614 Szabó, G. and Fath, G. (2007). Evolutionary games on graphs. *Physics reports*, 446(4-6):97–
615 216.
- 616 Szabó, G. and Fáth, G. (2007). Evolutionary games on graphs. *Phys Rep*, 97-216(4-6).
- 617 Teixeira, A. S., Santos, F. C., Francisco, A. P., and Santos, F. P. (2021). Eliciting fairness in
618 n-player network games through degree-based role assignment. *Complexity*, 2021.

- 619 Traulsen, A., Hauert, C., De Silva, H., Nowak, M. A., and Sigmund, K. (2009). Explo-
620 ration dynamics in evolutionary games. *Proceedings of the National Academy of Sciences*,
621 106(3):709–712.
- 622 Traulsen, A., Nowak, M. A., and Pacheco, J. M. (2006). Stochastic dynamics of invasion and
623 fixation. *Phys Rev E*, 74(1):011909.
- 624 Trivers, R. L. (1971). The evolution of reciprocal altruism. *Quarterly Review of Biology*, 46:35–
625 57.
- 626 Tuyls, K. and Parsons, S. (2007). What evolutionary game theory tells us about multiagent
627 learning. *Artificial Intelligence*, 171(7):406–416.
- 628 Tuyls, K., Perolat, J., Lanctot, M., Ostrovski, G., Savani, R., Leibo, J. Z., Ord, T., Graepel, T.,
629 and Legg, S. (2018). Symmetric decomposition of asymmetric games. *Scientific Reports*,
630 8(1):1–20.
- 631 Wang, S., Chen, X., and Szolnoki, A. (2019). Exploring optimal institutional incentives
632 for public cooperation. *Communications in Nonlinear Science and Numerical Simulation*,
633 79:104914.
- 634 West, S., Griffin, A., and Gardner, A. (2007). Evolutionary explanations for cooperation. *Cur-*
635 *rent Biology*, 17:R661–R672.
- 636 Wilder, B., Immorlica, N., Rice, E., and Tambe, M. (2018). Maximizing influence in an un-
637 known social network. In *AAAI conference on Artificial Intelligence (AAAI-18)*.
- 638 Wooldridge, M. (2012). Bad equilibria (and what to do about them). *ECAI '12*, pages 6–11.
- 639 Zisis, I., Guida, S. D., Han, T. A., Kirchsteiger, G., and Lenaerts, T. (2015). Generosity mo-
640 tivated by acceptance - evolutionary analysis of an anticipation games. *Scientific reports*,
641 5(18076).

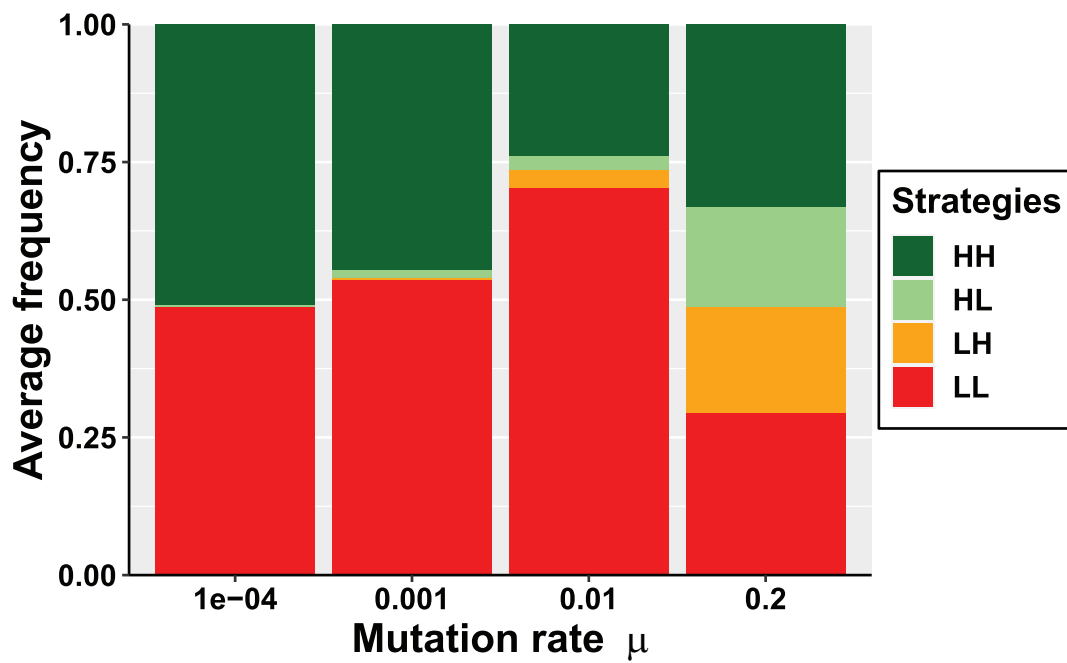
642 **5 Supplementary Information**

Figure S1. Average frequencies of the four strategies HH, HL, LH and LL as a function of mutation rate μ in absence of interference.

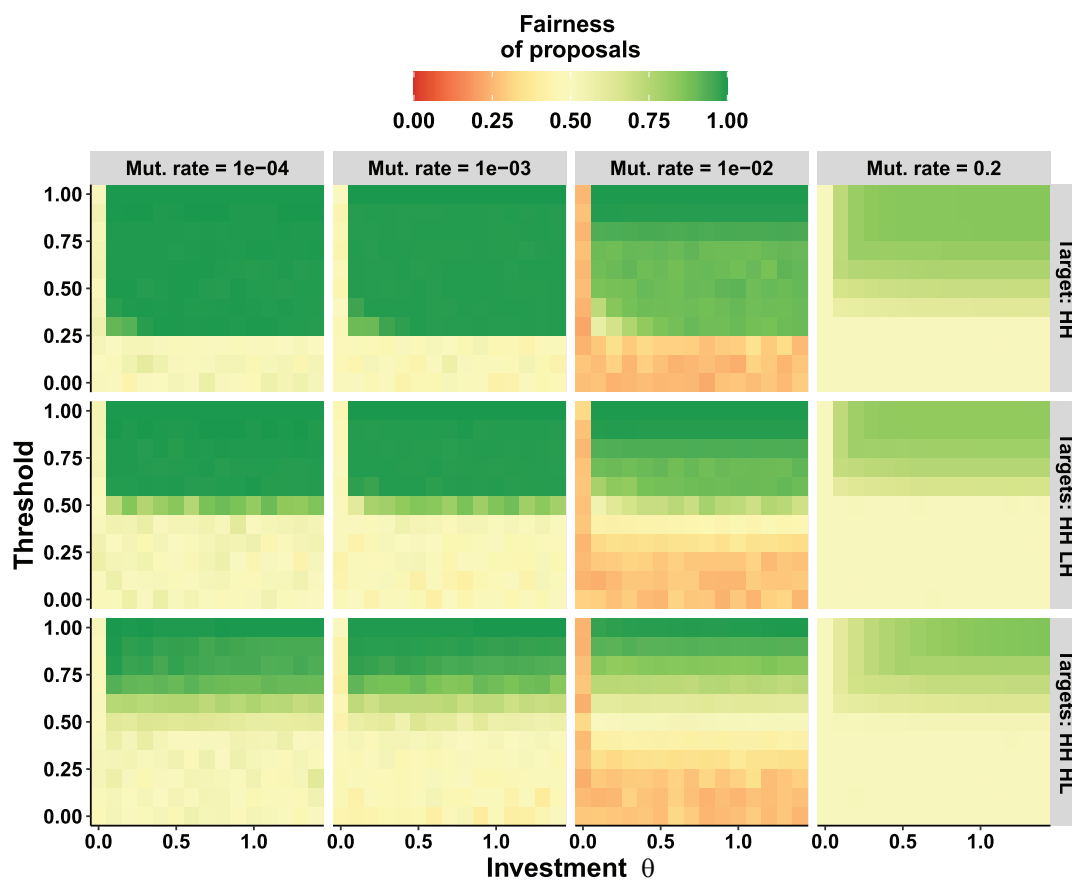


Figure S2. Average fairness as a function of the individual endowment θ , the threshold p_f and the mutation rate μ (population-based, stochastic update). Each row represents a different targeting scheme.

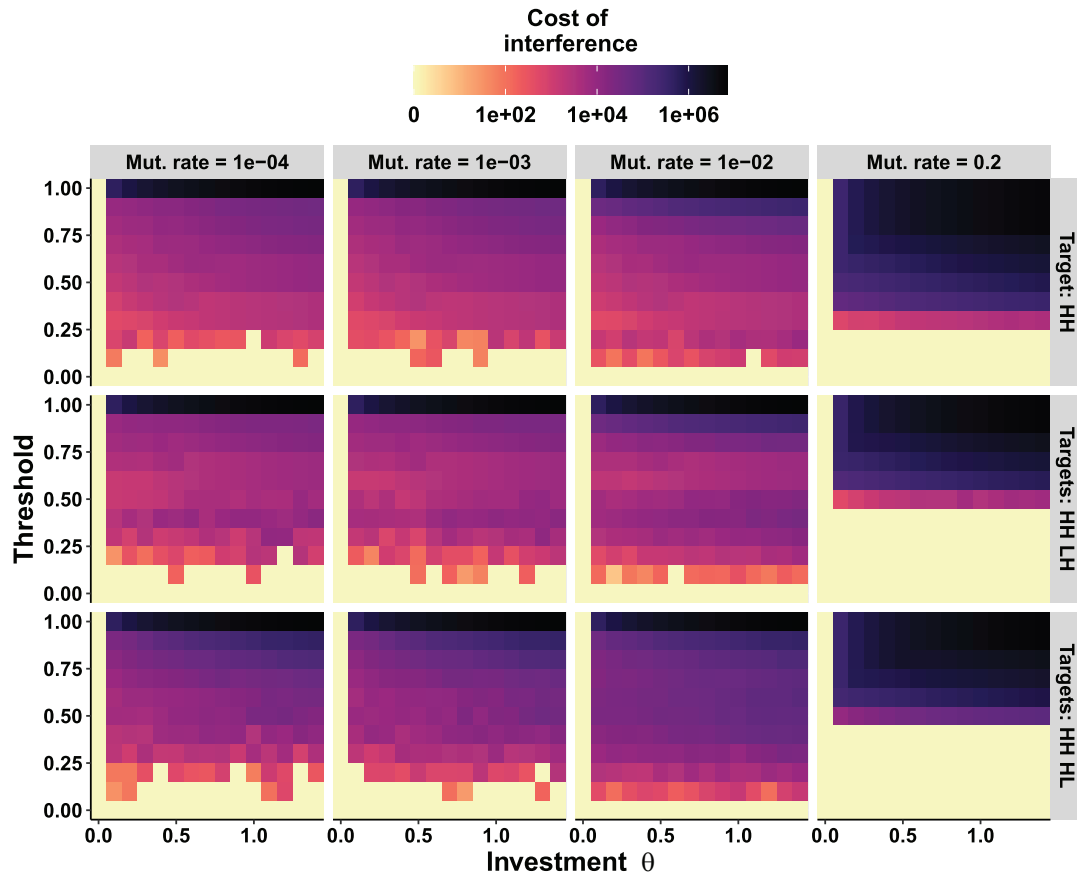


Figure S3. Average cost of interference as a function of the individual endowment θ , the threshold p_f and the mutation rate μ (population-based, stochastic update). Each row represents a different targeting scheme. The cost of interference is shown on a logarithmic scale.

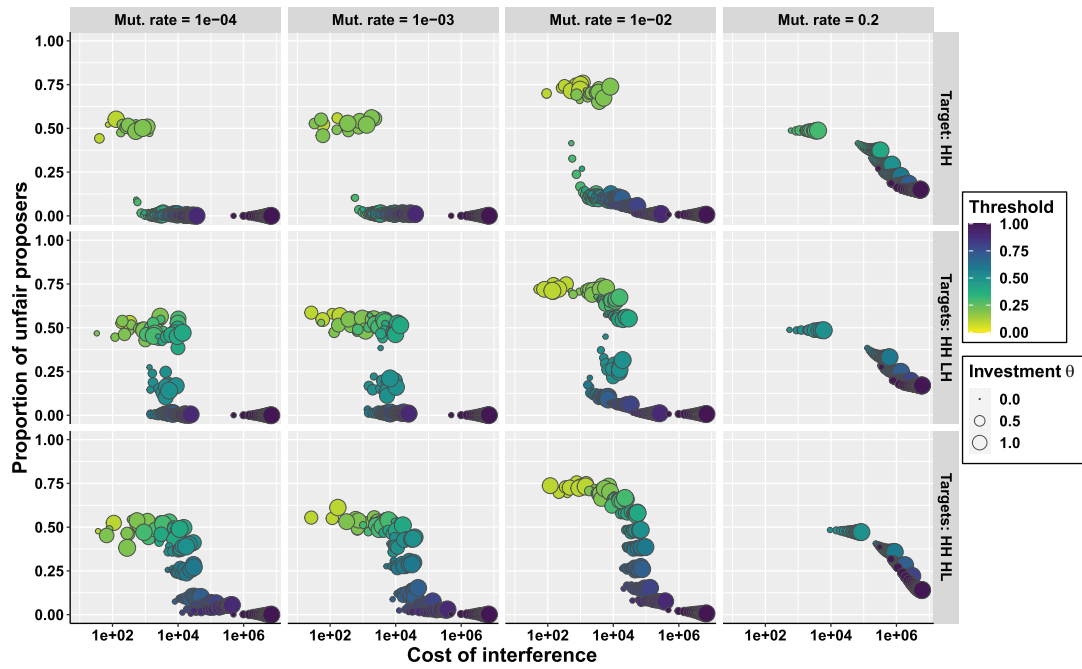


Figure S4. Proportion of unfair proposers as a function of average cost of interference for different targeting scheme and mutation rate μ (population-based, stochastic update). The size and colour of the circles correspond to investment amount and threshold of investment, respectively. We note that the most desirable outcomes are closest to the origin.

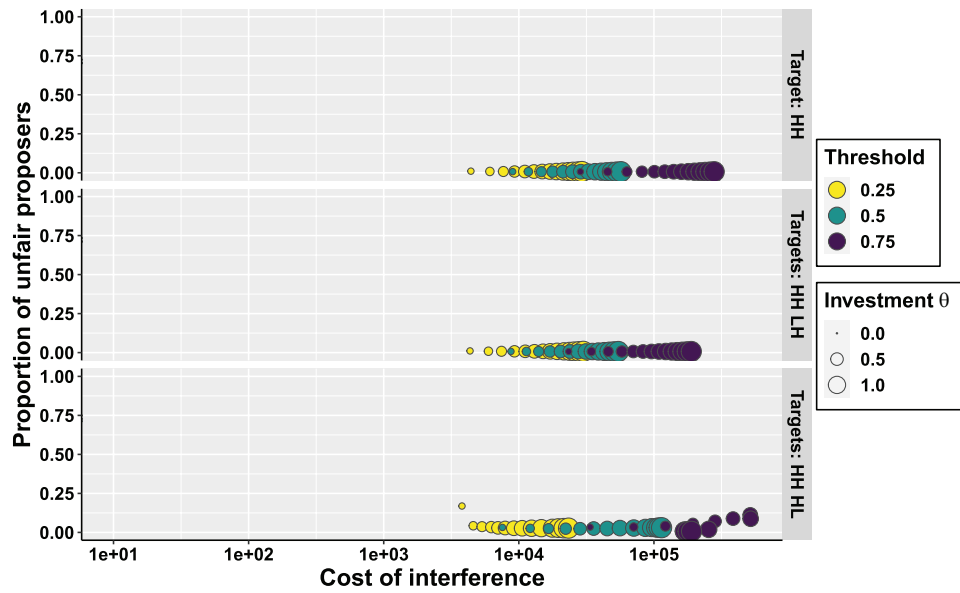


Figure S5. Proportion of unfair proposers as a function of average cost of interference for different targeting scheme (neighbourhood-based, $\mu = 0.01$, stochastic update). The size and colour of the circles correspond to investment amount and threshold of investment, respectively. We note that the most desirable outcomes are closest to the origin.

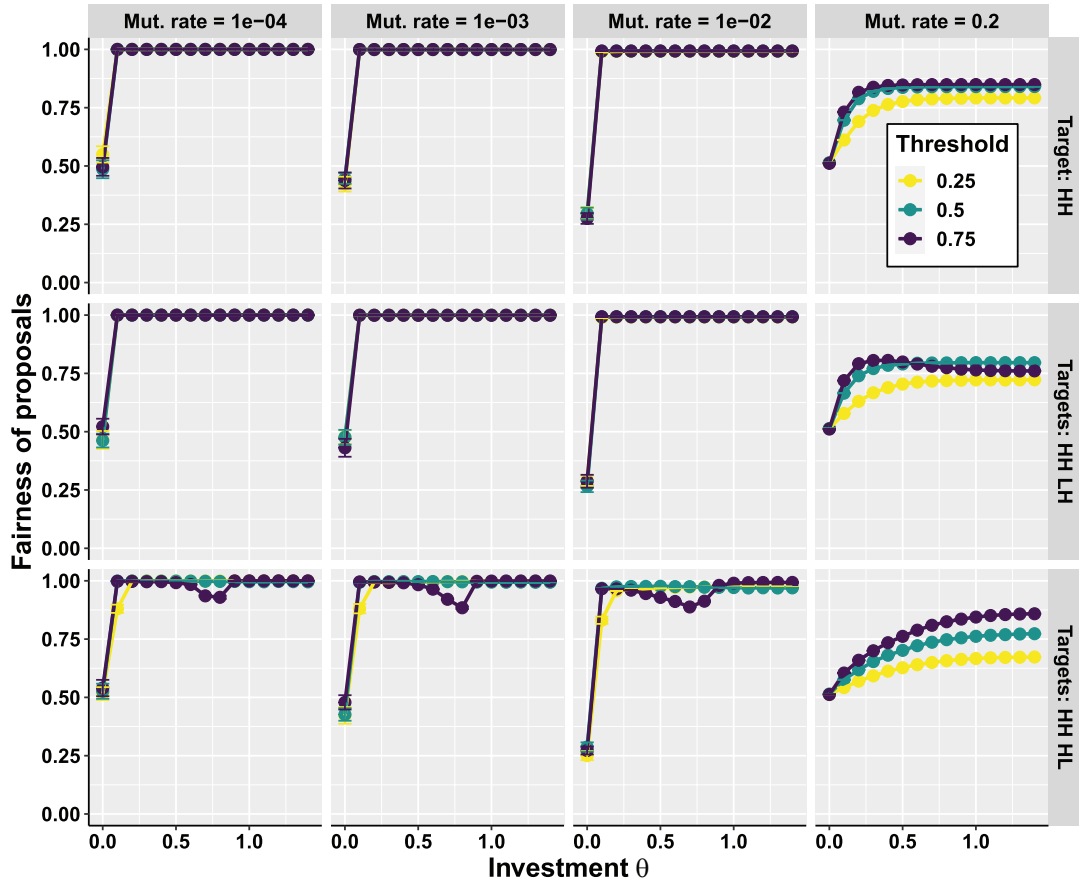


Figure S6. Average fairness measured by the sum of frequencies of HH and HL as a function of the individual endowment θ , the threshold p_f and the mutation rate μ (neighbourhood-based, stochastic update). Each row represents a different targeting scheme.

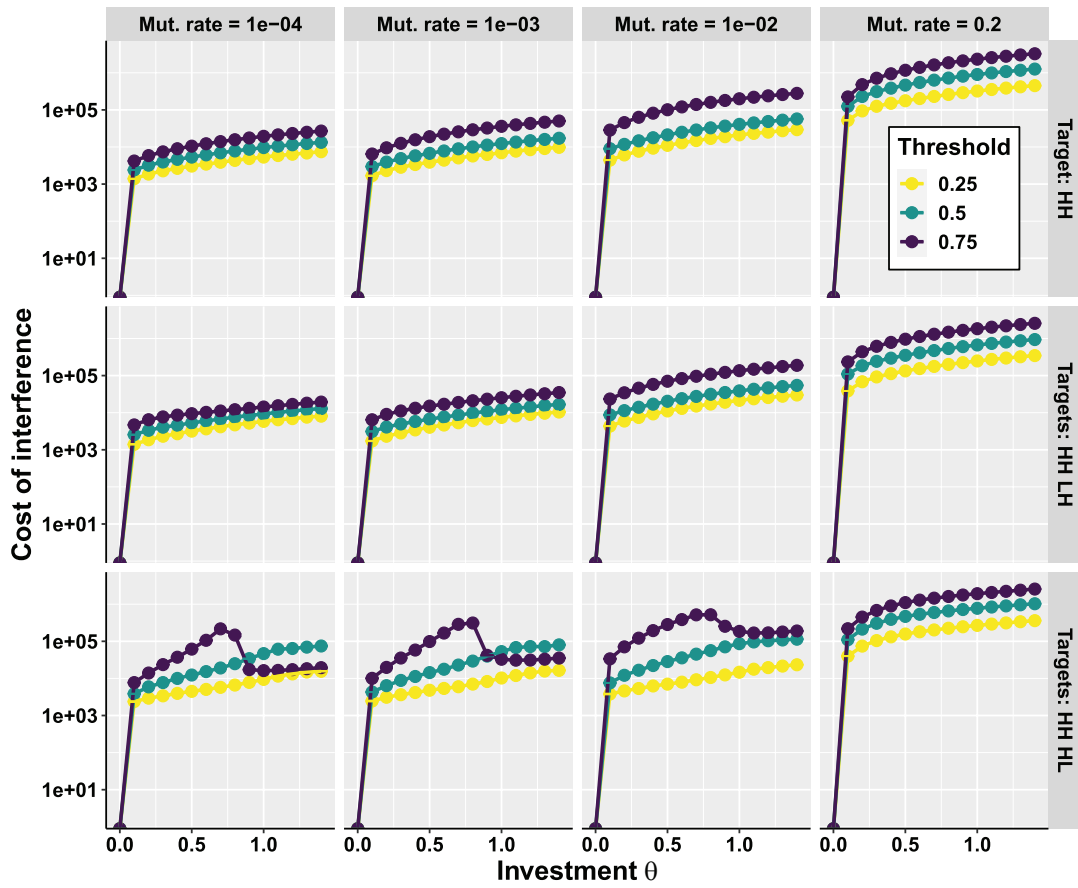


Figure S7. Average cost of interference as a function of the individual endowment θ , the threshold p_f and the mutation rate μ (neighbourhood-based, stochastic update). Each row represents a different targeting scheme. The cost of interference is on a logarithmic scale for clarity.

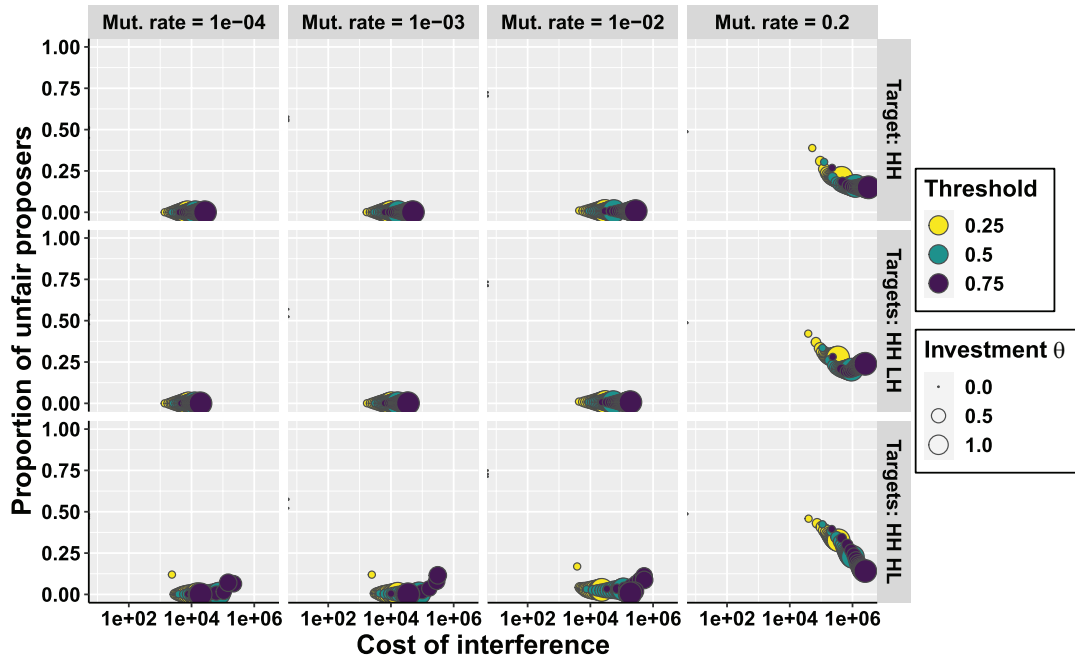


Figure S8. Proportion of unfair proposers as a function of average cost of interference for different targeting scheme and mutation rate μ (neighbourhood-based, stochastic update). The size and colour of the circles correspond to investment amount and threshold of investment, respectively. We note that the most desirable outcomes are closest to the origin.

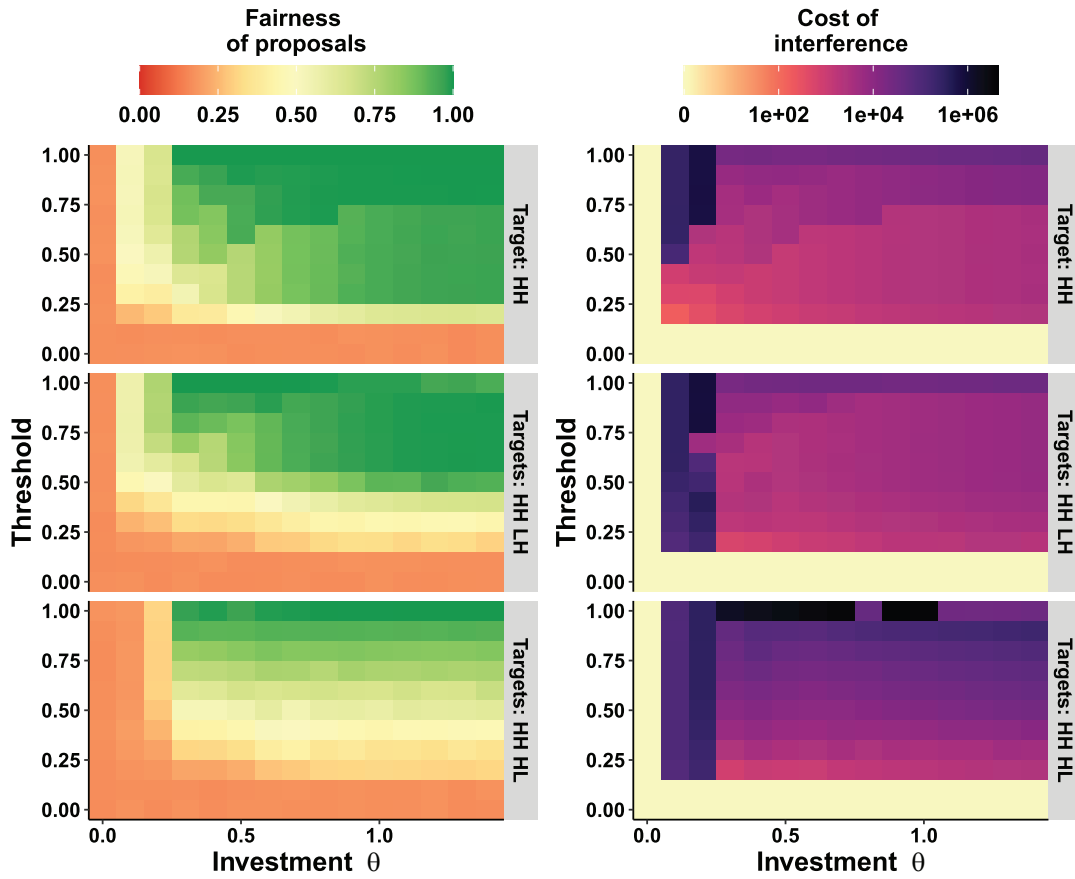


Figure S9. Average fairness (left) and average cost of interference (right) as a function of the individual endowment θ and the threshold p_f (population-based, deterministic update). Each row represents a different targeting scheme. The cost of interference is shown on a logarithmic scale.