1    **A House Divided: Cooperation, Polarization, and the Power of Reputation**

2    Sebastian Wild[1], Phillip Keldenich[2], Jann Spiess[3], Maximilian Schlund[4], Jano Costard[5], Jonas

3    Radbruch[6], Paul Stursberg[7], Sándor P. Fekete[2*]

4    [1] Department of Computer Science, University of Liverpool, Liverpool L69 3BX, UK.

5    [2] Department of Computer Science, TU Braunschweig, 38106 Braunschweig, Germany.

6    [3] Graduate School of Business, Stanford University, Stanford, CA 94305, USA.

7    [4] School of Computation, Information and Technology, Technical University Munich, 85748

8    Garching, Germany.

9    [5] Federal Agency for Disruptive Innovation - SPRIND, 04103 Leipzig, Germany.

10   [6] School of Business and Economics, Humboldt University Berlin, 10099 Berlin, Germany.

11   [7] Department of Mathematics, Technical University Munich, 85748 München, Germany.

12   [*] Corresponding author: s.fekete@tu-bs.de

13

14   **Altruistic cooperation has enabled humans to thrive[1]. However, the interaction of**

15   **sentient individuals faces the dilemma of limiting the downsides of personally beneficial,**

16   **but globally detrimental selfish behavior without causing even more damage through**

17   **escalating conflicts. The evolution of cooperation has been studied in non-zero sum**

18   **games, with the *Prisoner's Dilemma,* "the E. coli of social psychology"[2], providing a**

19   **fundamental test case. Typically[3-12], interactions between individuals may (i) occur**

20   **repeatedly, (ii) involve groups of individuals, (iii) be subject to evolutionary**

mechanisms, often based on the study of equilibria for homogeneous settings.[13] However, a better understanding of the *non-equilibrium dynamics* of cooperation *in structured environments* is crucial for further progress. Here we consider an inhomogeneous, spatial, dynamic setting, in which evolution occurs not necessarily at an equilibrium. We demonstrate how minimal, publicly observable information on previous behavior can be exploited to outperform alternatives, achieving evolutionary performance similar to clandestine, membership-based strategies. We also show how *polarization* (with a cooperating population disintegrating into competing factions) and *tribalism* (with cooperation solely based on group membership instead of behavior) can arise, how these phenomena can be overcome with two additional mechanisms, and how cooperation can erode. Our results demonstrate how cooperation, reputation, polarization and tribalism are intricately linked, even in a simple mathematical model in which they arise in absence of complex psychological mechanisms. This provides a fundamental explanation for how robust cooperation may break down when faced with eroding universality of globally recognized values and of local, direct reciprocity; it may also help to prevent behavior-based reputation systems from giving way to emergent polarization and, ultimately, purely membership-based tribalism. We also anticipate that our methods will be of critical importance for the design and implementation of artificial structures based on the interaction of many independent, self-interested virtual agents.

Classic work[2-13] on the evolution of cooperation analyzes group scenarios with large populations of interacting agents (subject to the aspects (i)–(iii) stated above), introducing evolutionary mechanisms based on superior payoff (in settings such as the Prisoner's Dilemma) for groups of cooperating individuals following a joint strategy. Typically, well-mixed populations are considered, in which all pairs of individuals can interact, so that the

46    evolutionary outcomes are global equilibria. Among the considered mechanisms[3-12,14] are

47    direct and indirect reciprocity, often making use of an evaluation of observable behavior, as

48    well as additional, hidden information.

49    In contrast, we consider the setting in which interaction (iv) takes place in spatially structured

50    environments. This reflects the fact that evolutionary successful, cooperative groups are often

51    first established locally[15], making it natural to consider populations that interact *spatially*[3,15-

52    21]; see Fig. 1b for a basic model of such spatial interaction. This does not only occur in cell

53    biology, but even in human populations: the impact on partisan sorting of humans was

54    recently highlighted[22]. Other recent work on polarization has considered cognitive aspects[23],

55    and algorithmic complexity[24]. The interplay of reputation and polarization has also been

56    considered[25], but only based on indirect reputation evaluation, which differs from our more

57    general approach; moreover, the resulting form of polarization is different from what we

58    consider here, and inherently unstable.

59    At the higher level, the competition between different subpopulations with the same

60    respective strategy is based on the success of locally competing individuals, according to

61    their respective payoffs; as a consequence, global success is not necessarily based on

62    centrally coordinated global welfare maximization, but through the dynamic, distributed

63    process of local optimization. This makes it natural to extend the notion of evolutionary

64    success by gauging it with the quantitative parameter of *invasion speed*; refer to Fig. 1 c-f and

65    the Methods section for details of the ensuing spatial dynamics. This goes beyond the

66    traditional notion of an *evolutionary stable strategy* (ESS), which requires only stability

67    against a small fraction of mutant strategies[23] in well-mixed populations without spatial

68    considerations. Instead, we consider how such a stable population can *evolve* in the first place

69    by invading and defeating an existing, well-established population in a spatial setting with

localized neighborhoods. Previous work[16,17,21] on the spatial Prisoner's Dilemma shows how unconditional cooperators are able to invade a population of defectors and maintain spatially coherent clusters for mildly adversarial environments with only weak benefit of defecting; however, they quickly die out in more hostile settings.

These limitations of cooperative strategies can be addressed by enhancing them with principles such as indirect reciprocity through the use of natural, publicly available information and algorithmic mechanisms. This involves evaluating observed behavior between individual interactions to gauge trustworthiness, leading to a reputation that is assigned to players[7-10,12,14-15,27-35]. The basic idea is that in large populations, it is possible to observe and learn from the behavior of an individual towards many others, even if the number of interactions between the same two individuals is limited: "Indirect reciprocity describes the interaction between a donor and a recipient. The donor can either cooperate or defect. The basic idea of indirect reciprocity is that cooperation increases one's own reputation, while defection reduces it. The fundamental question is whether natural selection can lead to strategies that base their decision to cooperate (at least to some extent) on the reputation of the recipient." (Nowak[3], supporting online material.) This establishes a setting in which "Each player has an image score, s, which is known to every other player." (Nowak and Sigmund[7]); "An individual's score is known by all group members, for instance because all interactions are publicly observed" (Leimar and Hammerstein[10]). More technically, Nowak and Sigmund[9] noted, "This review of theoretical and empirical studies of indirect reciprocity stresses the importance of monitoring not only partners in continuing interactions but also all individuals within the social network. Indirect reciprocity requires information storage and transfer as well as strategic thinking and has a pivotal role in the evolution of collaboration and communication."

94   Technically, reputation is captured by a mathematical function that uses a spectrum of

95   information on a player (in particular, observed past actions) as input to compute a decision

96   on cooperation or defection when interacting with that player. This expresses how a player

97   can map an opponent's (potentially extensive) sequence of past decisions to an eventual

98   binary decision of trustworthiness, i.e., an action from {cooperate, defect} in a new

99   interaction. To achieve evolutionary success, a considerable variety of functions have been

100  proposed. In some settings, simple reputation systems may suffice[7,8], but often they are not

101  successful under all circumstances[10,30]. More advanced reputation systems are often able to

102  overcome shortcomings of simpler ones[14,30-32], frequently at the expense of using a larger

103  amount of interaction data. Another option is to use additional, hidden information, such as

104  membership in a clandestine organization: the strategy MAFIA is characterized by a secret bit

105  that is only visible to other members. Both aspects encounter limitations: keeping track of

106  vast amounts of interaction data quickly becomes prohibitively costly, and mechanisms that

107  are based on covert coordination or group membership may be undesirable for other reasons,

108  e.g., in the context of organized crime or racism. Further details are discussed in the Methods

109  section.

110  We have developed a simple yet powerful mechanism that uses only a minimal amount of

111  *publicly* visible information. Our strategy GANDHI assigns a reputation value of good or bad

112  (corresponding to worthy or unworthy of cooperation, i.e., the actions {cooperate, defect}

113  in the next interaction) to each individual, and conducts updates only based on two bits of

114  information, corresponding to two past interactions with others: An individual is considered

115  good if both (i) its last interaction with another good individual was cooperation, and (ii) its

116  last interaction with a bad individual was non-cooperation. The key idea behind this strategy

117  is to efficiently promote both desired cooperation with trustworthy individuals and punish

118 undesired support of defectors; the name alludes to a well-known quote by Mahatma

119 Gandhi[36]: "Non-cooperation with evil is as much a duty as is cooperation with good."

120 We have demonstrated the power of this simple strategy in a systematic comparison with a

121 spectrum of other methods for indirect reciprocity that have been proposed in the literature.

122 To this end, we used extensive multi-parameter computer simulations of a standard model

123 from evolutionary game theory (Fig. 1 a-b), complemented with mathematical analysis of

124 Markov chain approximations. This model isolates the features of a social dilemma in which

125 individuals have no immediate incentive to cooperate; additional mechanisms, in particular,

126 public reputation, can help cooperative strategies gain foothold even when the temptation of

127 defecting is very high. We have followed previous work[15-21,29-30,32-35] on the spatial Prisoner's

128 Dilemma, which considered a setting in which reputation-based strategies had to attempt

129 invading a population of unconditional defectors (ALLD, which never cooperate) or a

130 population of unconditional cooperators (ALLC, which always cooperate); see Fig. 1. In the

131 context of this spatial version of the classic Prisoner's Dilemma, GANDHI is never weaker

132 than other reputation-based strategies, and outperforms all of them in terms of the memory

133 required for this success. Furthermore, the performance of GANDHI is comparable to MAFIA,

134 which has access to hidden information; these conclusions are validated by both a detailed

135 mathematical analysis and extensive simulations. Our findings allow us to combine and

136 extend results from previous spatial and indirect reciprocity approaches to scenarios in which

137 cooperation is more costly; our results also contribute to explaining how reputation may have

138 evolved by showing that even a small group of individuals can dominate a large population

139 even in rather adversarial scenarios, establishing reputation as a global mechanism through

140 the course of gradual evolution.

141    While this demonstrates the power of GANDHI in competition with other strategies, its simple

142    mechanism has one significant downside, which can lead to polarization of an otherwise

143    cooperating population: it is *antisymmetric*, i.e., it remains consistent when good and bad

144    reputation are *swapped*. As shown in Fig. 2 (and discussed in more detail in the SI), this can

145    lead to fragmenting the successful GANDHI population into two competing factions ("red" and

146    "blue") that both follow GANDHI, despite implementing the exact same set of rules: while

147    members of RED GANDHI (RG) consider other RG individuals as good, but members of BLUE

148    GANDHI (BG) as bad, members of BG consider other BG individuals as good, but members

149    of RG as bad. Such a split can be induced by an inhomogeneous initialization

150    (corresponding, e.g., to optimistic or pessimistic individuals), but may be triggered even by a

151    seemingly innocuous disturbance, such as a single error in observation. Once this

152    fragmentation happens, further observations only strengthen the respective assignments, as

153    BG will cooperate with BG, but not with RG, and vice versa, confirming the respective (but

154    antisymmetric) labeling as good and bad. Even an ongoing competition between the two

155    factions (with individuals being "turned" when overpowered by their neighbors of the other

156    faction) only leads to stronger polarization in which the spatial separation between factions

157    increases, resembling the process of coarsening of spin glasses from physics[37]. This very

158    gradual reduction in separation length (corresponding to the occurrence of non-cooperation

159    between the factions) still manages to slightly improve global welfare (corresponding to

160    overall average score) based on local competition, and thus still outperforms other strategies.

161    However, neither party can decisively defeat the other; moreover, the coarsening process

162    itself proceeds extremely slowly, when compared to the relatively swift evolutionary of

163    success of GANDHI against other strategies. As it turns out, this becomes completely

164    analogous to a contest between two purely membership-based strategies, in which members

165    of RED MAFIA (RM) only cooperate with other members of RM, while members of BLUE

166  MAFIA (BM) only cooperate with each other. In effect, *polarization* (in which two factions

167  emerge that start to fight each other, based on observable behavior) becomes

168  indistinguishable from *tribalism*, in which cooperation and non-cooperation are not based on

169  behavior, but on group membership alone.

170  As a consequence, we studied additional mechanisms to deal with polarization. The first is a

171  *global* mechanism based on universally recognized authorities, in which two entities

172  ("virtue" and "evil") are uniformly considered as good and bad; players interact with these

173  authorities at random occasions (with a probability of $h$), allowing their reputation to be

174  updated to good even in the perception of the other faction. While this does overcome

175  polarization for sufficiently large values of $h$, the critical threshold (around $h = 0.735$)

176  appears too high for an effective mechanism by itself.

177  A second enhancement is a *local* mechanism based on direct reciprocity, in which a direct

178  neighbor is considered good as long as their last direct interaction was to cooperate.  In the

179  long run, this can lead to more wide-spread cooperation, but it does not overcome the

180  polarization of reputation, leaving the whole population vulnerable to a collapse of

181  cooperation when local reciprocity is weakened.

182  However, using both of the global and the local mechanisms[1] in combination with GANDHI

183  (resulting in GANDHI++) is able to counter polarization: for very small values of $h$, the

184  combination of direct, local reciprocity with global calibration by universally recognized

185  authorities is able to overcome even settings with artificially enhanced polarization (Fig. 3).

186  As a consequence, GANDHI++ is also quite robust against perturbations, making it a very

187  powerful strategy that uses only minimal information and mechanisms.

---

[1] This may be mapped to the classical "Love God and love your neighbor.", Matthew 22:36-40 and Mark 12:30-31.
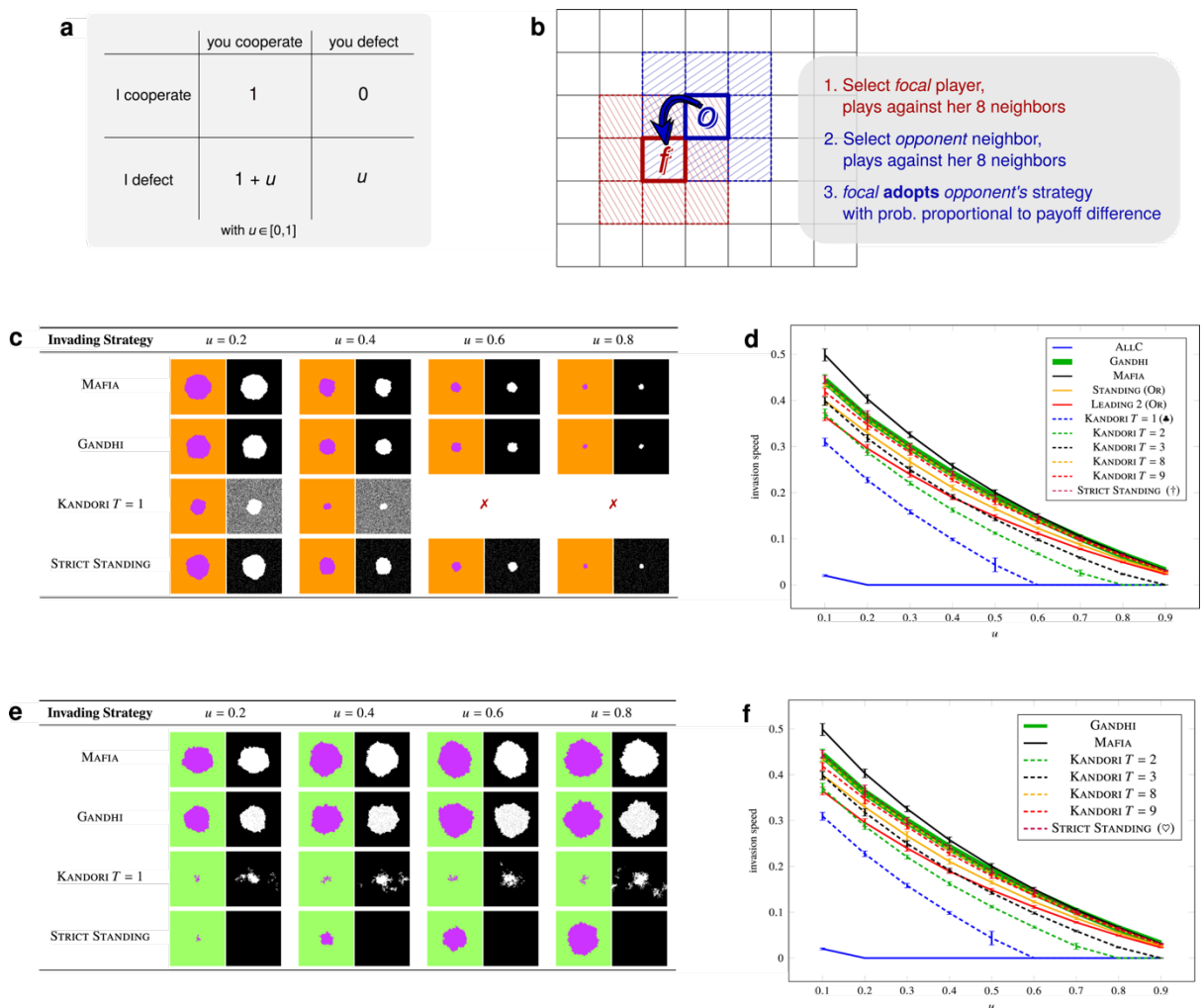
188     While this sequence of insights is rather encouraging for the development of cooperative

189     mechanisms, a further twist and caveat arises from considering a direct competition between

190     GANDHI++, GANDHI and MAFIA: While GANDHI++ is able to both thrive in basic adversarial

191     settings (such as swiftly defeating populations of ALLD) and also to deal with polarization—

192     thereby achieving universal cooperation—it is vulnerable to populations without the

193     stabilizing effects of globally recognized institutions and local reciprocity, leading to an

194     erosion of cooperation: As we show in Fig. 4, a population of GANDHI++ can slowly but

195     surely be defeated by an opposing group of GANDHI. Furthermore, GANDHI in turns falls prey

196     to MAFIA, due to the slightly slower update mechanism when taking over other players.

197     **Discussion**

198     We are confident that our findings will provide useful tools for the field of systems of

199     artificial agents, where cooperation has to be based on explicitly programmed protocols, and

200     the use and availability of a small amount of publicly available information is of crucial

201     importance. This opens up a number of additional mechanisms and aspects beyond the

202     confines of the considered setting with the Prisoner's Dilemma in a spatial setting with fixed

203     neighborhoods of fixed size; in particular, active mechanisms of expanding connectivity and

204     more variable payoffs in other non-zero-sum games (which allow both group support to

205     "frontier" members faced with adversarial individuals, as well as escalation in conflict)

206     promise further relevant insights for theory and practice.

207     While we make no claims in the realm of political or social sciences, it seems inevitable that

208     the simplicity of our reputation-based mechanisms makes them particularly suitable to be

209     studied in these important areas. (After all, even a famous quote such as "*A house divided*

210     *against itself, cannot stand.*" relies on the metaphoric power of gravity.) In particular, it is

211     conceivable that the emergence of increasing tribalism in a society may have some

212  similarities to a transition from GANDHI++ to GANDHI, i.e., the erosion of the polarization-

213  preventing mechanisms of direct reciprocity and universally accepted instances of "virtue"

214  and "evil", which may in turn give way to a transition to the purely membership-based

215  MAFIA. Conversely, successfully overcoming tribalism may hinge on (re-)establishing these
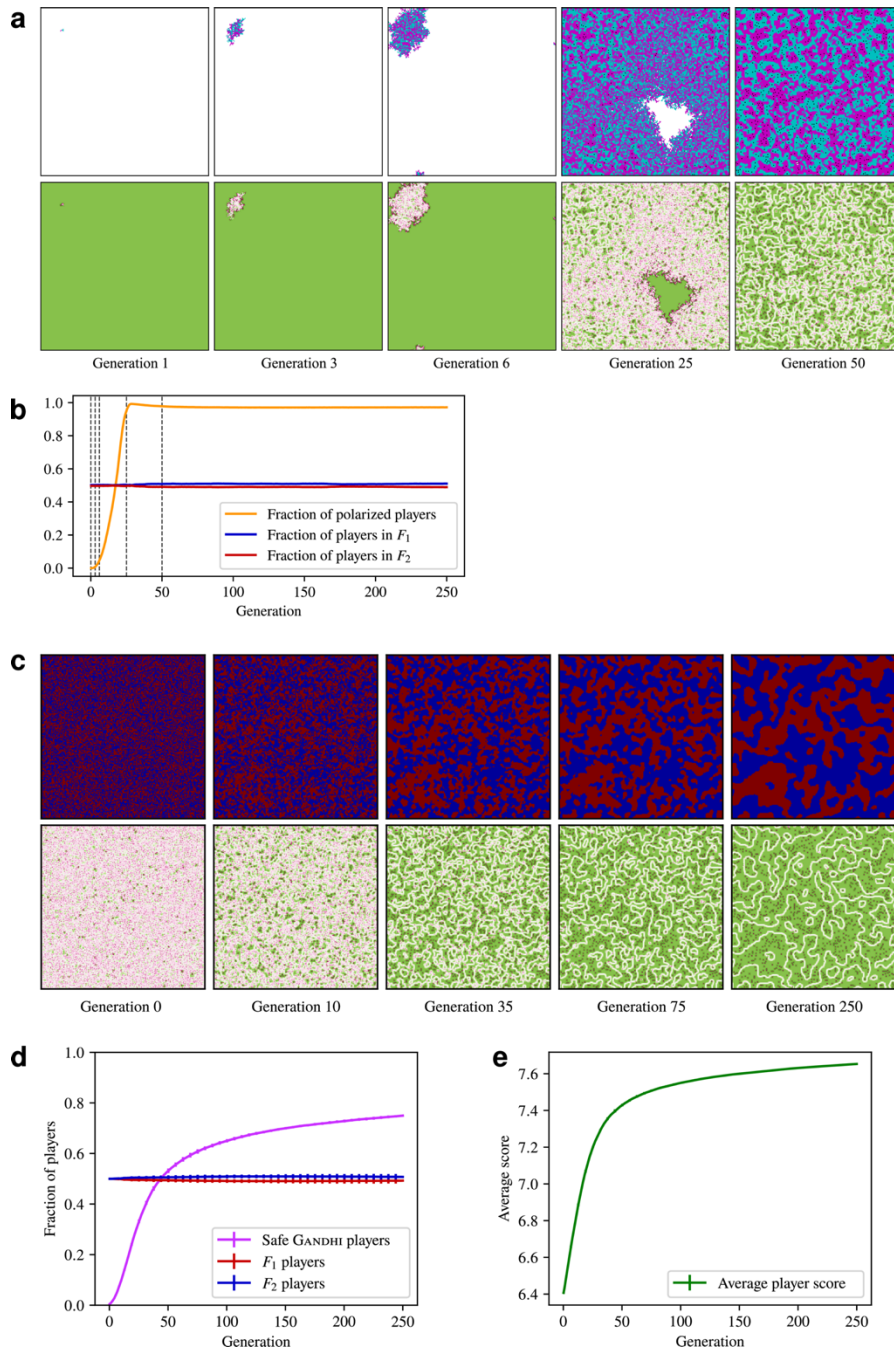
216  global and local mechanisms.

217



219  **Figure 1: Spatial prisoner's dilemma with semi-deterministic replicator rule and public**

220  **reputation, and invasion speed of reputation-based strategies (DISC) in an ALL-DEFECT**

221  **(ALLD) environment.**

222   **a+b**, The underlying model[16] with **a,** the payoff matrix for Prisoner's Dilemma (PD) and **b,**

223  repeated player interaction with their eight neighbors on an N × N square lattice of

224  individuals with periodic boundaries, and adopting more successful strategies in a replicator

225 update process. In addition, players have access to public data based on previous

226 interactions. **c**, Prisoner's Dilemma with ALLD (orange) and DISC invaders (purple) on a 200

227 × 200 grid, and snapshots after 240 generations. Colored tiles: population, B/W tiles:

228 reputation (white = good, black = bad). Rows vary the reputation systems, columns the

229 exploitation benefit $u$. The reputation system KANDORI with $T = 1$ dies out in an ALLD

230 population for $u \geq 0.6$, demonstrating the weakness of single-bit tracking. **d**, Invasion speed

231 of DISC against an ALLD population under all reputation systems and different exploitation

232 conditions; higher speed is stronger. Each data point shows mean and standard deviation of

233 20 independent runs. "STRICT STANDING (†)" represents STANDING and STRICT STANDING

234 reputations, "KANDORI $T = 1$ (♣)" stands for KANDORI with $T = 1$, and LEADING 3, 4 and 5.

235 **e**, **f**, Same as **a, b,** but in an ALLC environment (shown as light green). KANDORI reputation

236 does not die out against cooperators, but fails to convert them effectively, leading to fractal

237 structures in the strategies distribution.

238

239

Figure 2: **Polarization emerges among two symmetric GANDHI factions. a-b**,
Polarization, i.e., players seen as $\mathsf{good}$ by one faction and $\mathsf{bad}$ by the other, spreads from a
single misinterpreted duel (in the top left corner). **a**, simulation after 1, 3, 6, 25, and 50
generations; top tiles: reputation difference, bottom tiles: score. In the reputation-difference
map, players are cyan if considered $\mathsf{good}$ by RED GANDHI (RG) and $\mathsf{bad}$ by BLUE GANDHI
(BG), magenta if considered $\mathsf{bad}$ by RG and $\mathsf{good}$ by BG, and black (resp. white) if
considered $\mathsf{bad}$ (resp. $\mathsf{good}$) by both factions. The score shows the payoff each player

achieved in their last game (greener is better). **b**, Number of polarized players over time. A very small number of players become depolarized; such a player is seen as bad by both factions, because they were the last in a neighborhood to change faction and were hence unable to defect against a bad opponent to regain good reputation with their own faction. **c-e**, Two competing groups of GANDHI, red and blue, over time. **c**, Snapshots of the simulation; top tiles: population, bottom tiles: score. **d**, The number of "safe" players, i.e., players for which all neighbors are in their own faction averaged over n=10 experiments with random initial configurations. This number grows over time through "coarsening" of the boundaries. **e**, Social welfare (average total score that each player gets when playing with their neighbors) over time. This rises in line with safe players, but does not overcome the polarization of the overall population.

**a**

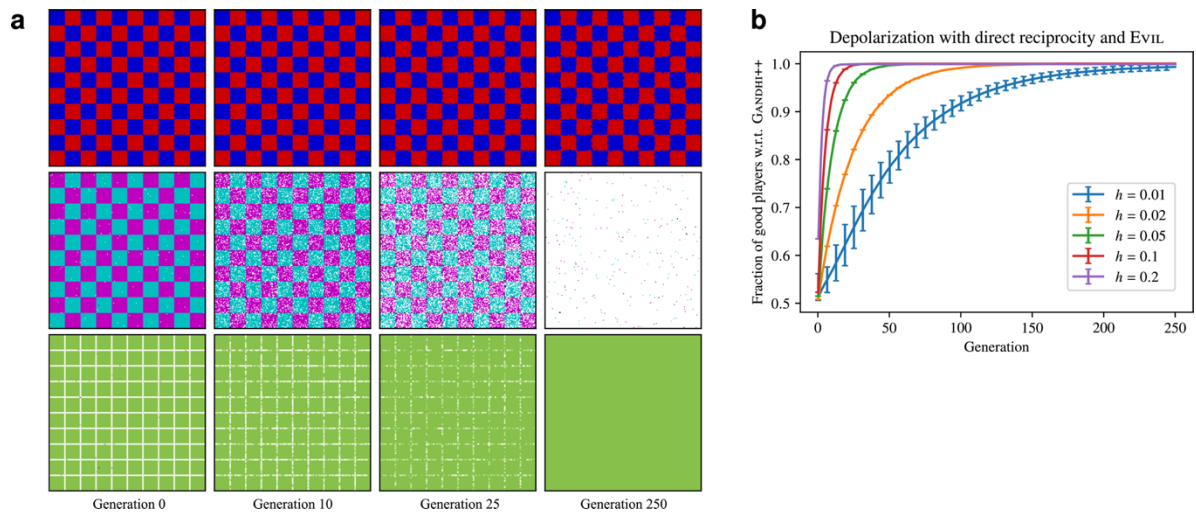Generation 0   Generation 10   Generation 25   Generation 250

**b**

Depolarization with direct reciprocity and EVIL

- $h = 0.01$
- $h = 0.02$
- $h = 0.05$
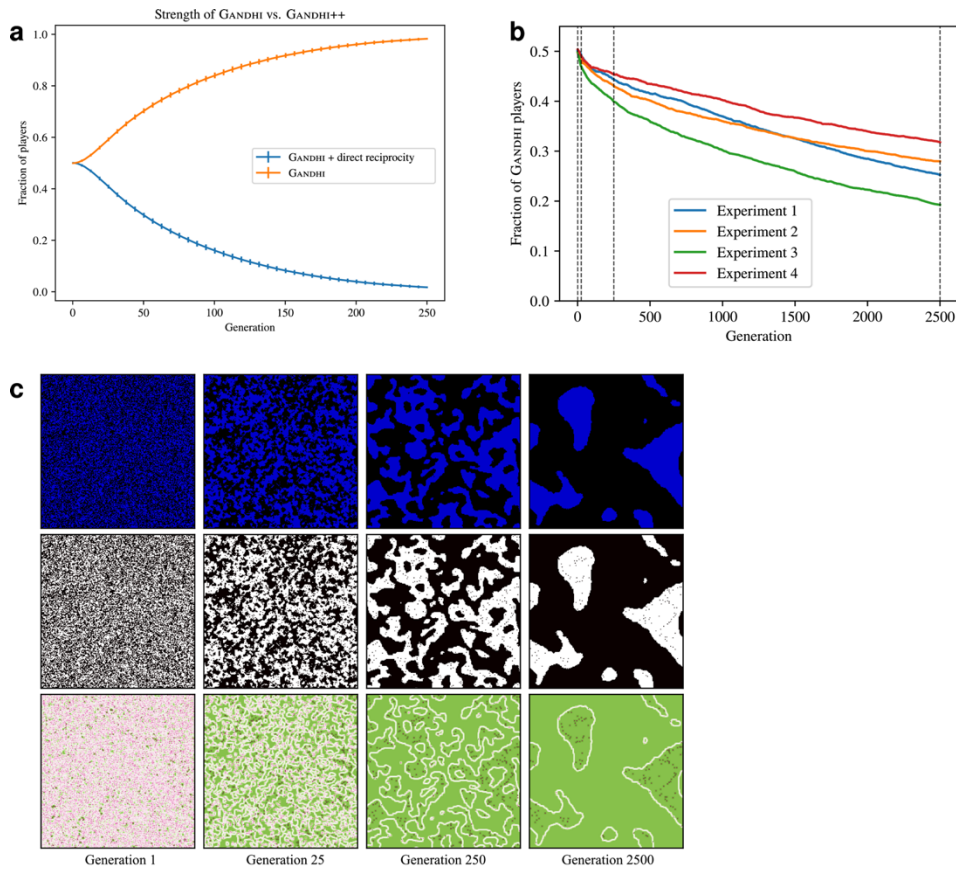- $h = 0.1$
- $h = 0.2$

Fraction of good players w.r.t. GANDHI++

Generation

**Figure 3: Opposing GANDHI++ factions recovering from initially prevalent polarization. a**, Snapshots of a typical simulation on a 200 × 200 grid after 0, 10, 25, and 250 generations. Top tiles: population (red and blue GANDHI factions), middle tiles: reputation difference (colors as in Fig. 3), bottom tiles: scores. Contact probability with virtue and evil is $h = 0.01$. Reciprocity and regular contact with global authorities eventually leads to all players being considered good by both factions and thus to global cooperation. **b**, Number of polarized players over time for different global-authority probabilities $h$. Here, only evil authorities are used, showing that virtue is not necessary in GANDHI++.

278

**Figure 4: Direct competition of GANDHI++, GANDHI and MAFIA. a**, Number of GANDHI and GANDHI++ players over time in the simulation of a direct competition. GANDHI is able to replace GANDHI++ relatively quickly. **b-c**, Direct competition of MAFIA (black) and GANDHI (blue). **b**, The number of GANDHI players over time for four exemplary simulations. **c**, Snapshots from Simulation 1; top tiles: population, middle tiles: GANDHI's reputation, bottom tiles: scores. Similar to a competition between two MAFIA or two GANDHI factions, we observe a coarsening of the strategy distribution. MAFIA eventually overcomes GANDHI, but the process is orders of magnitude slower. Only few GANDHI players on the boundary of the resulting large blocks of GANDHI players are vulnerable.

288

**References and Notes:**

1. Tomasello, M., Melis, A. P., Tennie, C., Wyman, E. & Herrmann, E. Two Key Steps in the Evolution of Human Cooperation: The Interdependence Hypothesis. *Current Anthropology,* **53,** 673-692 (2012).

2. Axelrod, R. Effective choice in the Prisoner's Dilemma. *Journal of Conflict Resolution* **24**, 3–25 (1980).

3. Axelrod, R. & W. D. Hamilton, W. D. The evolution of cooperation. *Science* **211**, 1390–1396 (1981).

4. Nowak, M. A. Five rules for the evolution of cooperation. *Science* **314**, 1560–1563 (2006).

5. Nowak, M. A. Evolving cooperation. *Journal of Theoretical Biology* **299**, 1–8 (2012).

6. Doebeli, M. & Hauert, C. Models of cooperation based on the prisoner's dilemma and the snowdrift game. *Ecology Letters* **8**, 748–766 (2005).

7. Nowak, M. & Sigmund, K. Evolution of indirect reciprocity by image scoring. *Nature* **393**, 573–577 (1998).

8. Nowak, M. A. & Sigmund, K. The dynamics of indirect reciprocity. *Journal of Theoretical Biology*, **194**, 561–574 (1998).

9. Nowak, M. A. & Sigmund, K. Evolution of indirect reciprocity. *Nature* **437**, 1291–1298 (2005).

308     10. Leimar, O. & Hammerstein, P. Evolution of cooperation through indirect reciprocity.
309          *Proceedings of the Royal Society of London. Series B: Biological Sciences* **268**, 745–
310          753 (2001).

311     11. Imhof, L. A., Fudenberg, D. & Nowak, M. A. Evolutionary cycles of cooperation and
312          defection. *PNAS* **102**, 10797-800 (2005).

313     12. Sigmund, K. Moral assessment in indirect reciprocity. *Journal of Theoretical Biology*
314          **299**, 25–30 (2012).

315     13. Nash, J. F. Equilibrium points in n-person games, *PNAS* **36**, 48-49 (1950).

316     14. Ohtsuki, H. & Iwasa, Y. How should we define goodness? — Reputation dynamics in
317          indirect reciprocity. *Journal of Theoretical Biology* **231**, 107–120 (2004).

318     15. Helbing, D., Szolnoki, A. Perc, M. & Szabó, G. Evolutionary establishment of moral
319          and double moral standards through spatial interactions. *PLoS Computational Biology*
320          **6**, e1000758 (2010).

321     16. Fu, F., Nowak, M. A. & Hauert, C. Invasion and expansion of cooperators in lattice
322          populations: Prisoner's dilemma vs. snowdrift games. *Journal of Theoretical Biology*
323          **266**, 358–366 (2010).

324     17. Hauert, C., Spatial effects in social dilemmas. *Journal of Theoretical Biology* **240**,
325          627–636 (2006).

326     18. Szabó, G. & Tőke, C. Evolutionary prisoner's dilemma game on a square lattice.
327          *Physical Review E* **58**, 69–73 (1998).

328    19. Nowak, M. A. & May, R. M., Evolutionary games and spatial chaos. *Nature* **359,**
329         826–829 (1992).

330    20. Hauert, C. & Doebeli, M. Spatial structure often inhibits the evolution of cooperation
331         in the snowdrift game. *Nature* **428**, 643–646 (2004).

332    21. Langer, P., Nowak, M. A. & Hauert, C. Spatial invasion of cooperation. *Journal of*
333         *Theoretical Biology* **250**, 634–641 (2008).

334    22. Brown, J. R. & Enos, R. D. The measurement of partisan sorting for 180 million
335         voters. *Nature Human Behaviour* 1–11 (2021).

336    23. Wu, J. S.-T. Wu, Hauert, C.,  Kremen, C., & Zhao, J. A framework on polarization,
337         cognitive inflexibility, and rigid cognitive specialization. *Frontiers in Psychology* 13,
338         Article 776891 (2022).

339    24. Chatterjee, K., Ibsen-Jensen, R. Jecker, I., & Svoboda, J. Complexity of spatial
340         games. *42nd IARCS Annual Conference on Foundations of Software Technology and*
341         *Theoretical Computer Science (FSTTCS),* 11:1--11:14 (2022).

342    25. Gross, J., & De Dreu, C.K.W. The rise and fall of cooperation through reputation and
343         group polarization. *Nature Communications* 10, 776 (2019).

344    26. Maynard Smith, J. & Price, G. R., The logic of animal conflict. *Nature* **246**, 15–18
345         (1973).

346    27. Sugden, R. The economics of rights, co-operation and welfare, B.Blackwell, 1986.

347    28. Kandori, M. Social norms and community enforcement. *The Review of Economic*
348         *Studies* **59**, 63–80 (1992).

29. Nakamaru, M. & Kawata, M. Evolution of rumours that discriminate lying defectors. *Evolutionary Ecology Research* **6**, 261–283 (2004).

30. Panchanathan, K. & Boyd, R. A tale of two defectors: the importance of standing for evolution of indirect reciprocity. *Journal of Theoretical Biology* **224**, 115– 126 (2003).

31. Sigmund, K. & Brandt, H. The logic of reprobation: assessment and action rules for indirect reciprocation. *Journal of Theoretical Biology* **231**, 475–486 (2004).

32. Panchanathan, K. Two wrongs don't make a right: The initial viability of different assessment rules in the evolution of indirect reciprocity. *Journal of Theoretical Biology* **277**, 48–54 (2011).

33. Brandt, H. & Sigmund, K. The good, the bad and the discriminator — errors in direct and indirect reciprocity. *Journal of Theoretical Biology* **239**, 183–194 (2006).

34. Yang, W., Juan W., & Chengyi X. "Evolution of cooperation in the spatial public goods game with the third-order reputation evaluation." Physics Letters A 383.26 (2019): 125826.

35. Quan, J., Qin, Y., Zhou, Y., Wang, X., & Yang, J. B. (2020). How to evaluate one's behavior toward "bad" individuals? Exploring good social norms in promoting cooperation in spatial public goods games. *Journal of Statistical Mechanics: Theory and Experiment*, 2020(9), 093405.

36. Gandhi, M., Statement before Mr. C. N. Broomfield, I. C. S., District and Sessions Judge, Ahmedabad, 18 March, 1922.

37. Matsuda, H. N., Ogita, A., Sasaki, A. & Sato, K. Statistical mechanics of population: The lattice Lotka-Volterra model. *Progress in Theoretical Physics* **88**, 1035-1049 (1992).

# Methods

The overall model in the main text is based on three different aspects: Player actions and interactions lead to payoffs according to a *non-zero-sum game-theoretic setting* (Fig. 1a), the relative success of different strategies leads to their spread in a *spatial setting* (Fig. 1b); choosing cooperation or non-cooperation when interacting with another player can be based on information from previous actions, condensed in *reputation systems*. Here we introduce the main technical aspects; a full account of ensuing data is provided in the Supplementary Information.

**Spatial Replicator Dynamics**

Interaction between players occurs in a setting in which the global population is structured in local environments; this corresponds to a spatial setting with geometric neighborhoods. We make use of the model by Fu et al.[16], who consider an $N \times N$ square lattice of individuals with periodic boundaries, in which two players repeatedly interact with their eight neighbors by playing a symmetric $2 \times 2$ game, as shown in Fig. 1b. To evaluate the evolutionary success of different strategies, we model their spread by using the replicator rule (called semi-deterministic updating by Fu et al.[16]): We randomly choose one focal player out of the $N \times N$ square lattice and an opponent among its eight neighbors. (Using randomly selected duels for potential updates avoids artifacts of synchronization; for settings with stronger parallelization, the expected values for spread can simply be adjusted to filter out the effects of expected waiting times for a duel to occur.) Both play against all their neighbors, resulting in accumulated payoffs $P_f$ and $P_o$ for focal and opponent, respectively. Then the focal player adopts the opponent's strategy with probability

$$f(P_f, P_o) = \begin{cases} \dfrac{P_o - P_f}{8(1 + u)} & P_f < P_o \\ 0 & \text{otherwise.} \end{cases}$$

410

411    Here $8(1+u)$ is the maximal payoff difference in Prisoner's Dilemma. The replicator rule can

412    be seen as a way to apply the classic replicator dynamics for infinite well-mixed

413    populations[38] to finite structured populations; in both cases, the spreading rate is linear in the

414    payoff differences and the payoffs are based on the mean (neighboring) opponent player[39].

415    Finally, focal and opponent's reputation is updated according to the respective reputation

416    system. (The reputation of the other neighbors involved in the duels, i.e., the neighbors of

417    focal and opponent, remains unchanged; modifying this assumption would make reputation-

418    based mechanisms only stronger.) We stress that the reputation update is done irrespective of

419    whether the focal player adopts the opponent's strategy; in particular, her reputation is not

420    newly initialized to some reputation score nor is it copied from the opponent's reputation.

421    The spreading of strategies is a mechanism of learning or imitating behavior; such a strategy

422    change is an internal, hidden event that can only be observed by others through subsequent

423    actions. This reflects a setting in which distinction between individuals and their actions is

424    based on location, not on publicly announced strategies.

425    **Reputation Systems**

426    Keeping track of the trustworthiness of players leads to assigning a reputation to players, i.e.,

427    a function that uses a spectrum of information on a player (in particular, observed previous

428    actions) to result in a decision on cooperation or defection when interacting with that player:

429    Every player follows a strategy, which is a function that takes that player's and her

430    opponent's reputation as arguments and returns an action from {cooperate, defect}. The

431    simplest strategies are the unconditional cooperators (ALLC) and unconditional defectors

432 (ALLD), which do not make use of any reputation; more sophisticated are discriminating

433 (DISC) strategies, which cooperate if the opponent has a good reputation and defect

434 otherwise. The meaning of label good depends on the specific reputation system (such as

435 GANDHI). Because the action space is binary, it suffices to consider only binary reputation

436 values, i.e., players are always either good or bad in the eyes of a DISC strategy. Note that

437 even though the labels good and bad may seem to suggest a moral verdict, our setting does

438 not a priori reward conformal behavior.

439 In our base model, we assume that all players have the same information as their neighbors,

440 modeling a well-connected world with rapid information dissemination and perfectly

441 observable actions; as discussed further down, there may be additional, hidden information.

442 Different discriminating strategies can use different rules to assign reputation labels and may

443 come to a different verdict based on (the same) past behavior. Formally, a reputation system

444 determines a label good or bad for each individual, based on the history of interactions the

445 individual was involved in. It is important to note that we allow for the possibility of

446 including the reputation of former opponents as well, i.e., players have access to higher-order

447 information. For example, the reputation system may rate defection against good or bad

448 players differently. In the base model, we assume that an individual's reputation is globally

449 agreed upon and based on public information. To model the equivalence in the parallel

450 interaction with all neighbors, we update reputation only after all eight neighbor duels of one

451 propagation round have taken place. This also accounts for a delay in the exchange of

452 information between neighbors until more tangible outcomes are visible; more responsive

453 update rules only enhance the advantage of discriminator systems.

454 In previous work, a wide spectrum of reputation functions have been proposed; these include

455 IMAGE SCORING by Nowak and Sigmund[7,8], which tracks the balance of previous cooperate

456 and DEFECT actions, but is unable to distinguish between defecting from cooperative or non-

457 cooperative players, GOOD STANDING by Sugden[24] and Leimar and Hammerstein[10], which

458 performs one-bit updates, making it unable to sanction cooperation with non-cooperative

459 players, KANDORI[25], which tallies a player's score over T rounds and only cooperates when

460 desirable behavior is maintained (requiring $\lceil \log_2 (T + 1) \rceil$ bits and punishing one-time

461 noncompliance through T rounds), and the LEADING EIGHT of Ohtsuki and Iwasa[14], which

462 are based on various 1-bit updates. All these differ from our strategy GANDHI, which only

463 uses two bits, but achieves better performance, as demonstrated in the sequel.

464 **Success of GANDHI**

465 To compare the discriminatory efficacy of different reputation systems, we study the

466 following questions: (1) Can a cluster of individuals who follow a joint reputation-based

467 strategy convince members of other strategies to imitate their discriminating behavior? (2) If

468 so, how does evolutionary success compare quantitatively, i.e., how fast is this invasion?

469 In answer to these questions, we provide both qualitative and quantitative evidence that

470 GANDHI outperforms other similar strategies.

471 **Qualitative Evidence**

472 Initially, all players in the $N \times N$ ($N = 70$) grid use the same incumbent strategy (either ALLD

473 or ALLC), except for a $5 \times 5$ square cluster of invading DISC individuals in the middle. We

474 explore every possible combination of incumbent strategy and reputation system of the

475 invading DISC players. Moreover, we vary the exploitation surplus parameter $u \in \{0.1, 0.2,$

476 … , 0.9\}. (We have also carried out a large range of additional experiments against mixed

477 populations. These results are not included here, as they do not provide any additional

478 insights.) For each of these setups, the simulation runs until either the invaders die out or the

479 first invader touches the boundary. By the time the boundary is reached, the invasion's final

480 success can be reliably assessed; further progress would be artificially slowed down by

481 boundary effects.

482 Some examples are shown in Fig. 1 c+e, with a full overview listed in Extended Data Fig. 1;

483 in addition to Prisoner's Dilemma (PD), the latter also include analogous results for

484 Snowdrift (SD), a two-player non-zero sum game in which cooperation with a non-

485 cooperating opponent is less detrimental. It can be seen that only a limited number of

486 strategies succeed in defeating both ALLD and ALLC populations: KANDORI (with at least $T$

487 = 8, i.e., higher-order interaction data), MAFIA (which uses hidden information) and

488 GANDHI.

**Quantitative Evidence**

489

490 Similar to the observations of Fu et al.[16], expansion basically proceeds at constant speed in

491 both dimensions. Therefore, the square root of the number of DISC players grows linearly in

492 the number of played duels. We accordingly define invasion speed as the corresponding rate

493 of change, i.e., by how much the square root of the number of DISC players grows on average

494 in one generation. A generation is here defined as $N^2$ simulated duels, which corresponds to

495 one chance per player to reproduce on average. As the snapshots in Fig. 1 c+e show, the

496 region occupied by DISC players is of roughly circular shape. Thus, the invasion speed

497 corresponds to the average growth rate of the radius of this circle.

498 Fig. 1 d+f show the invasion speed of DISC players using various reputation systems. For

499 reputation systems that could never invade, no line is shown. Each point shows the average

500 invasion speed of 20 independent runs of the corresponding simulation. Error bars show one

501 standard deviation around the mean. The narrow error bars show that invasion speed is a

502    robust measure: It is reliably reproduced in independent runs. As invasion speed is a global

503    measure determined from many independent random variables, low variance was to be

504    expected.

505    Again, GANDHI dominates all other strategies, with the exception of MAFIA, which achieves

506    faster update speed through hidden information.

**Mathematical Evidence**

508    Additional mathematical evidence can be obtained by analyzing the behavior of a Markov

509    chain that models the strategy transition of individuals in a mixed population. For the speed

510    $\psi_+$ of MAFIA vs. ALLD, this yields

$$\psi_+ \;=\; \frac{1}{8}\frac{1-u}{1+u}$$

511

512    For the analogous case of GANDHI vs. ALLD, we get a speed of

$$T^{-1} \;=\; \frac{-5u^2 - 3u + 8}{41u^2 + 111u + 72}$$

513

514    which works out to

$$T^{-1} \;\approx\; \frac{1}{9}\frac{1-u}{1+u}.$$

515

516    See the Supplementary Information for details of this analysis. As Extended Data Fig. 2

517    shows, this quantitative correspondence is supported by numerical evidence.

518

519

## Tribalism

The success of MAFIA relies not on a sophisticated strategy, but on strong group coherence, purely based on membership, i.e., *tribalism*. As a consequence, evolutionary success corresponds to the ability of the group to deal with adversarial groups, including other groups that also pursue MAFIA. Remarkably, two different groups of MAFIA cannot overcome each other, but still manage to improve global welfare (corresponding to overall average score) based on local competition. We demonstrate this with a number of experiments; see Extended Data Fig. 3 for an overview. Starting with an initial random distribution of two different group (REDMAFIA and BLUEMAFIA), running the replicator dynamics leads to a process resembling coarsening of spin glasses from physics[35]. More precisely, local competition between the two populations leads to a shortening of the separating boundary, as a weakly connected member of one population will be surrounded by a majority of members of the other; therefore, such an outlier will perform worse than a duel opponent, which is better connected to members of its own group. As a consequence, local majorities will take over their opposing neighbors, leading to smoother, shorter boundaries between the populations, corresponding to improved average score. (Note that this is only the case in the absence of escalation in the interaction with the opposing group.) However, this growing separation and local symmetry also makes it harder to take over neighbors, so that no subpopulation can defeat the other.

## Polarization

GANDHI is not based on membership, so it is more open to cooperating with (and thus benefiting from) neighbors, regardless of their strategy. However, its reputation system is subject to antisymmetry in the following sense. Suppose that there are two factions that both play according to GANDHI, with each faction perceiving its own players as good and the

544     players of the other faction as bad. The players of each faction then consistently cooperate

545     with players of their own faction, but defect against players of the other faction. We call such

546     a population *polarized*. As a consequence, the dynamics play out analogously to two MAFIA

547     factions; see Fig. 2. This implies that there is no inherent mechanism in GANDHI to overcome

548     polarization — once a population is polarized, it remains polarized, and only local boundary

549     minimization (and thus, local improvement of average scores) occurs; refer to Extended Data

550     Fig. 4.

551     There are several possible sources for polarization. Firstly, polarization may stem from

552     differences in initialization: If one ("REDGANDHI") faction $F_1$ "pessimistically" initializes all

553     players to a bad reputation and another ("BLUEGANDHI") faction $F_2$ "optimistically"

554     initializes all players to a good reputation, players in $F_1$ will defect in their first game;

555     similarly, players in $F_2$ will cooperate. Both actions are perceived as bad by the other faction,

556     leading to polarization. Secondly, even a single misperception can lead to a global

557     polarization, exposing the fragility of non-polarized populations in the base model. Suppose

558     we start with two GANDHI factions $F_1$ and $F_2$ sharing the same initialization and then play a

559     duel for which the action of a single player is perceived as cooperate by $F_1$ and defect by $F_2$.

560     This results in a single polarized player who is seen as good by one faction and as bad by the

561     other. Starting from this player, polarization spreads with every duel involving unpolarized

562     and polarized players, until the entire population is polarized; see Fig. 2.

563     **Global Authority**

564     Overcoming polarization in GANDHI requires breaking the antisymmetry between any kind of

565     split into REDGANDHI and BLUEGANDHI. One way to achieve this is by introducing global

566     authorities, virtue and evil, that are unequivocally seen as good resp. bad by any player

567     irrespective of their reputation system. In our simulation, we add these as artificial players

568  that focal and opponent encounter with a probability $h$ after playing the 8 duels with their

569  neighbors. The outcome of the (imaginary) duels with virtue and evil are only used in

570  updating a player's reputation; no payoff results from these encounters. As Extended Data

571  Fig. 5 shows, polarization can be dissolved for sufficiently large values of $h$: If players see a

572  global authority after at least 73% of the duels, polarization vanishes. Below this threshold,

573  some polarized players remain present (Extended Data Fig. 6) and continuously act as seed

574  for new polarization. Extended Data Fig. 5 also demonstrates that the fraction of polarized

575  players remains relatively stable over time. In isolation, global authorities are only an

576  effective cure for polarization if they are nearly omnipresent.

577  **Local Reciprocity**

578  Another mechanism to potentially counter polarization is to complement globally reported

579  information with direct observations, so that sporadic friendly acts among neighbors may be

580  rewarded and perpetuated. We incorporate this in our model in the form of local reciprocity:

581  Each player remembers for her 8 neighbors the last action they played against her, and

582  considers a neighbor $p$ as good whenever $p$ cooperated with her or when $p$'s global

583  reputation is good.

584  Extended Data Fig. 7 shows that this added leniency allows two polarized factions of

585  REDGANDHI and BLUEGANDHI to cooperate with each other: After an initial period, all

586  players manage to establish trust at the local level and hence benefit from the maximal social

587  welfare that cooperation entails. However, in the global reputation, polarization still looms

588  large: about half of the players are still polarized, and almost all other players are now

589  globally seen as bad (for not defecting from evil). Local reciprocity can effectively stop

590  polarization from affecting actions, but it does not cure the underlying divide.

591

**592** **GANDHI++**

**593** The GANDHI++ reputation system consists of simultaneously using global authorities and

**594** local reciprocity in GANDHI. Neither of these two additions in isolation cures polarization in

**595** GANDHI; however, GANDHI++ not only stops polarization from emerging, but can restore

**596** unity in an existing, completely polarized state (Fig. 3a). Note that any positive probability $h$

**597** $> 0$ for encountering global authorities will eventually lead to the eradication of polarization

**598** (Fig. 3b). Incidentally, having regular contact to `virtue` is not needed for this; a global `evil`

**599** authority (i.e., a universally regarded adversary) is sufficient.

**600** **Eroding Cooperation**

**601** While GANDHI++ is able to overcome both basic adversarial settings (such as swiftly

**602** defeating populations of ALLD) and deal with polarization, thereby achieving universal

**603** cooperation, it is also vulnerable to populations without the stabilizing effects of globally

**604** recognized institutions and local reciprocity, leading to an erosion of cooperation: As we

**605** demonstrate in the following, a population of GANDHI++ can be defeated by an opposing

**606** group of GANDHI, which in turns falls prey to Mafia.

**607** We observe that in a direct confrontation, GANDHI++ loses against GANDHI; see Fig. 4a for a

**608** typical outcome. This does not change if we remove the global recognition for `evil` and

**609** `virtue` by setting h = 0, i.e., even though local reciprocity alone does not suffice to overcome

**610** polarization, its presence is already sufficient to lose out against unmodified GANDHI; see

**611** Fig. 4a. This phenomenon can be attributed to the following mechanisms. Suppose that a

**612** GANDHI++ player $p$ has both GANDHI and GANDHI++ neighbors, gets beaten in a duel by

**613** some GANDHI player $q$ in the neighborhood, and changes membership to GANDHI. Player $p$

**614** now considers the GANDHI++ neighbors `bad` and defects against them; each GANDHI++

615     neighbor $r$ will cooperate with $p$ until $r$ itself has been betrayed by $p$. This makes GANDHI++

616     more vulnerable to defectors than GANDHI, which learns not to trust $p$ after a single defection

617     against any GANDHI player. This demonstrates the crucial role of both the existence of

618     universally recognized instances of good and bad, as well as local reciprocity. Therefore,

619     protecting cooperation against polarization hinges on protecting these mechanisms.

620     **From Polarization to Tribalism**

621     While GANDHI exhibits similar power against simple-minded strategies (such as ALLD or

622     ALLC), which are defeated almost as swiftly as by MAFIA, it slowly loses out to MAFIA in a

623     direct confrontation. The speed at which this happens can vary considerably, based on

624     random initialization and duel selection; see Fig. 4b. However, the eventual outcome is

625     inevitable, as long as the update speed for GANDHI (which relies on publicly visible

626     reputation information) is slightly slower than for MAFIA (which only needs to update a

627     hidden bit of information), as observed and analyzed above. A remedy to address this could

628     be to delay the adoption of MAFIA membership by exposed individuals.

629     **Additional References:**

630     38. Taylor, P. D., Jonker, L. B. Evolutionary stable strategies and game dynamics.

631         *Mathematical Biosciences*, 40(1-2):145–156, July 1978.

632     39. Roca, C. P., Cuesta, J. A. & Sánchez, A. Evolutionary game theory: Temporal and

633         spatial effects beyond replicator dynamics. *Physics of Life Reviews*, 6(4):208–49,

634         2009.

635

636

645
646 **Author contributions**

647 All authors conceived the study, performed the analysis, discussed the results and wrote the

648 manuscript.

649
650 **Competing interests**

651 The authors declare no competing interests.
652
653 **Supplementary information**

654 S.I. is available for this paper, and submitted in parallel.
655
656 **Data Availability**

657 All described data is available upon request and will be posted at a public repository.

658
659 **Ethics & Inclusion**

660 The nature of this work does not involve resource-poor settings.
661
662 **Correspondence and requests for materials** should be addressed to S.P.F.
663
664
665
666
667

668

# Extended Data Figures

| Reputation System | PD ALLD | PD ALLC | SD ALLD | SD ALLC |
|---|---|---|---|---|
| none (ALLC) | $[0.1, 0.2]$ | — | $[0.7, 0.8]$ | — |
| IMAGE SCORING | ✗ | ✗ | ✗ | ✗ |
| STRICT STANDING | ✓(†) | ✓(♡) | ✓(†) | ✓(♡) |
| STANDING | ✓(†) | ✗ | ✓(†) | ✗ |
| STANDING (OR) | ✓ | ✗ | ✓ | ✗ |
| LEADING 2 (OR) | ✓ | ✗ | ✓ | ✗ |
| LEADING 3 | $[0.5, 0.6]$(♣) | ✗ | ✓(♣) | ✗ |
| LEADING 4 | $[0.5, 0.6]$(♣) | ✗ | ✓(♣) | ✗ |
| LEADING 5 | $[0.5, 0.6]$(♣) | ✗ | ✓(♣) | ✗ |
| LEADING 8 | ✗ | ✓(♡) | ✗ | ✓(♡) |
| KANDORI $T = 1$ | $[0.5, 0.6]$(♣) | fractal | ✓(♣) | fractal |
| KANDORI $T = 2$ | $[0.7, 0.8]$ | ✓ | ✓ | ✓ |
| KANDORI $T = 3$ | $[0.8, 0.9]$ | ✓ | ✓ | ✓ |
| KANDORI $T = 8$ | ✓ | ✓ | ✓ | ✓ |
| KANDORI $T = 9$ | ✓ | ✓ | ✓ | ✓ |
| GANDHI | ✓ | ✓ | ✓ | ✓ |
| MAFIA | ✓ | ✓ | ✓ | ✓ |

669

670 **Extended Data Figure 1: Qualitative results on discriminatory efficacy.** Each entry

671 shows whether the corresponding reputation system allows DISC to take over the incumbent

672 population in the corresponding setting. Rows with (OR) correspond to scenarios where the

673 OR strategy is used instead of DISC, see Supplementary Information. An entry ✓ means

674 invasion is successful, ✗ means no invasion. An interval $[a, b]$ indicates that invasion

675 depends on the exploitation benefit and the threshold value lies in this interval. The term

676 "fractal" is used when the DISC region forms a fractal-like shape. As only a small fraction of

677 the players joined Disc here, "fractal" counts as ✗. For some settings, several reputation

678 systems become strongly equivalent, i.e., they behave exactly the same in every single step.

679 These equivalence classes are marked by ♥, ♣, and †, respectively.

680

**a** $\frac{8}{9}$ · Mafia speed vs. GANDHI (PD)
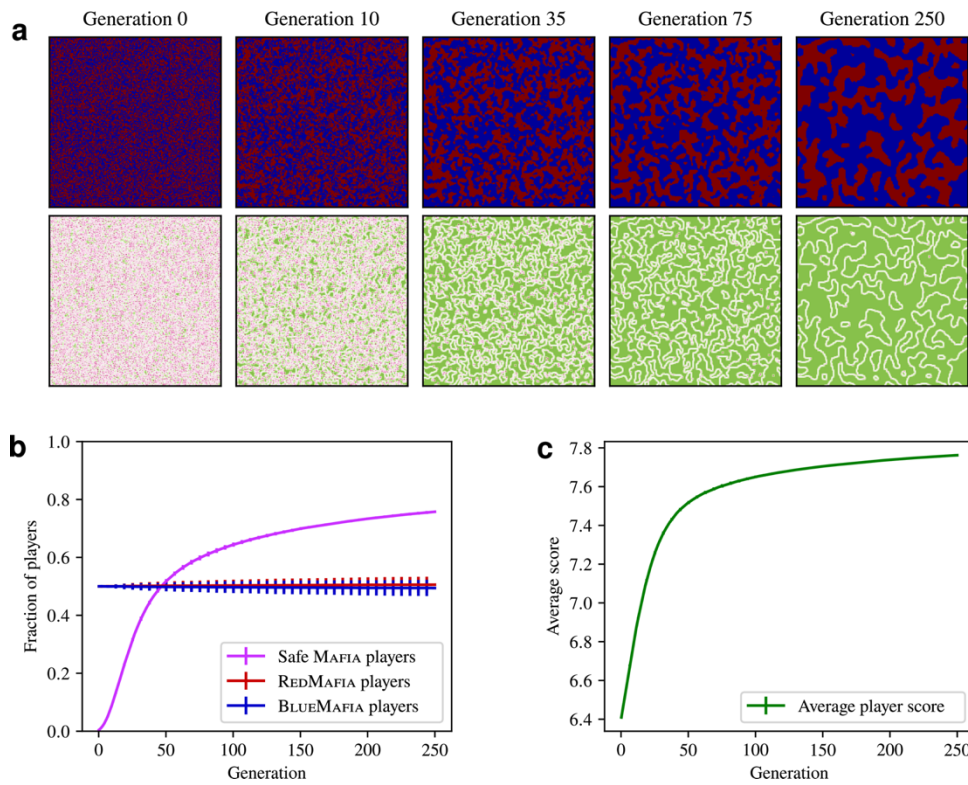
**b** 1D Theory vs. 2D Data (PD)

681

682 **Extended Data Figure 2: Validation of prediction of the one-dimensional Markov model**

683 **on simulation data. a**, The predicted invasion speed of GANDHI (green) from the 1D Markov

684 model for the Prisoner's Dilemma against ALLD (see Supplementary Information) as a

685 function of $u$ against 8/9 times the predicted invasion speed of MAFIA (dashed black). **b**, Plot

686 of the (scaled) predicted invasion speed from the 1D Markov model with the actual invasion

687 speed determined from our simulation (as in Fig. 1c+d) for both GANDHI (green) and MAFIA

688 (dashed black). The dependency in $u$ matches the theoretical prediction extremely well.

689

690

Extended Data Figure 3: Two competing groups of MAFIA over time. a, The distribution of strategies (top tiles, red or blue) and the score each player achieved (bottom tiles, greener is better) in the last round they played in an exemplary experiment after the stated number of generations. An initially fine-grained distribution of players, assigned to a group uniformly at random, coarsens over time. b, The number of players that are "safe", i.e., completely surrounded by players in their own group (average over n=10 experiments), increases over time in this coarsening process. c, The average player score likewise increases over time.

**a**



**b** Reputation in generation 50

700

701 **Extended Data Figure 4: Two competing GANDHI factions over time. a,** The number of

702 polarized players, i.e., players that are seen as good by one faction and bad by the other,

703 over time, as well as the number of players seen as bad by both GANDHI factions $F_1$ and $F_2$.

704 A generation is a number of rounds that corresponds to the total number of players. **b,** The

705 difference in reputation between the two factions. Cyan players are considered good by $F_1$

706 and bad by $F_2$, magenta players are considered bad by $F_1$ and good by $F_2$. Black players are

707 seen as bad by both factions. A very small number of players become depolarized, who are

708 now seen as bad. These players are seen as bad by both factions as a result of being the last

709 player in a neighborhood that changed faction — they are unable to defect against a bad

710 opponent to gain good reputation with their own faction because all their neighbors are

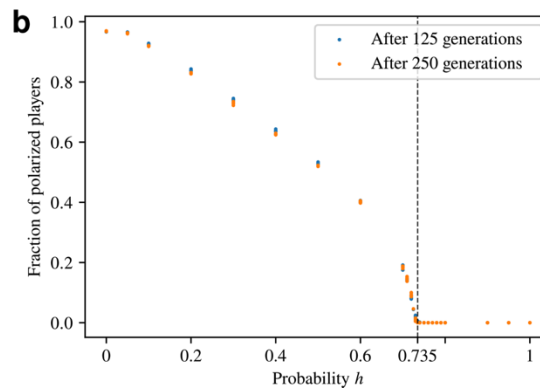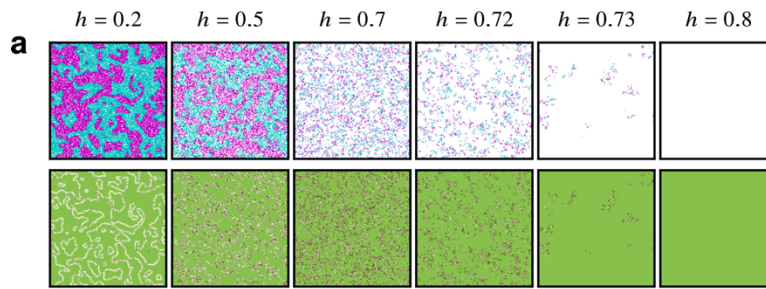711 good. No players are considered good by both factions.

712

713

714

**Extended Data Figure 5: The effect of global authorities on the number of polarized players in two competing GANDHI factions. a, b**, The number of polarized players over time for two competing GANDHI factions and various values of $h$, the probability of `virtue` and `evil` participating in a duel. Each point is the average of n=10 independent simulations, error bars show one standard deviation. If $h$ is not high enough, a part of the population remains polarized. For our grid model, the sufficient probability for completely removing polarization from a fully polarized population seems to be between $h = 0.72$ and $h = 0.74$.

h = 0.2  h = 0.5  h = 0.7  h = 0.72  h = 0.73  h = 0.8

724

**Extended Data Figure 6: The effect of global authorities on two polarized GANDHI factions is stable over time. a**, The reputation difference (top tiles) and the average player score (bottom tiles) after 250 generations for several values of the probability $h$ of encountering virtue and evil in a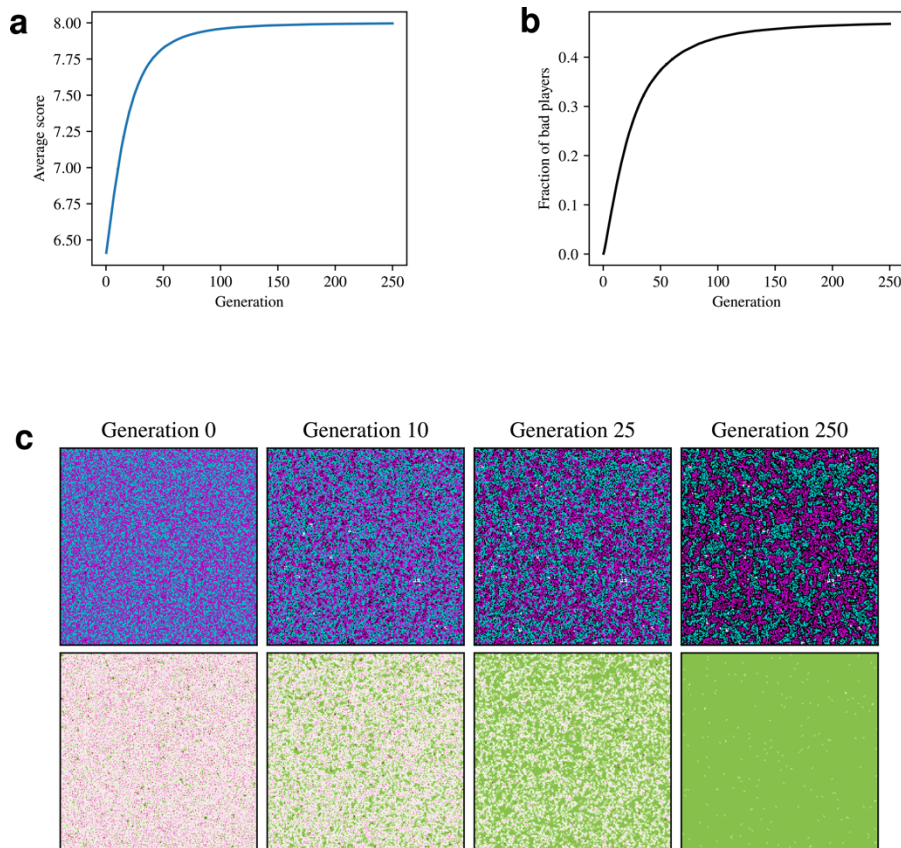 duel of two competing Gandhi factions. **b**, The number of polarized players that remain after 125 and 250 generations for varying values of $h$. Below the depolarization threshold of roughly 0.735, some polarized players remain present and continuously act as seed for new polarization; this fraction remains stable over time.
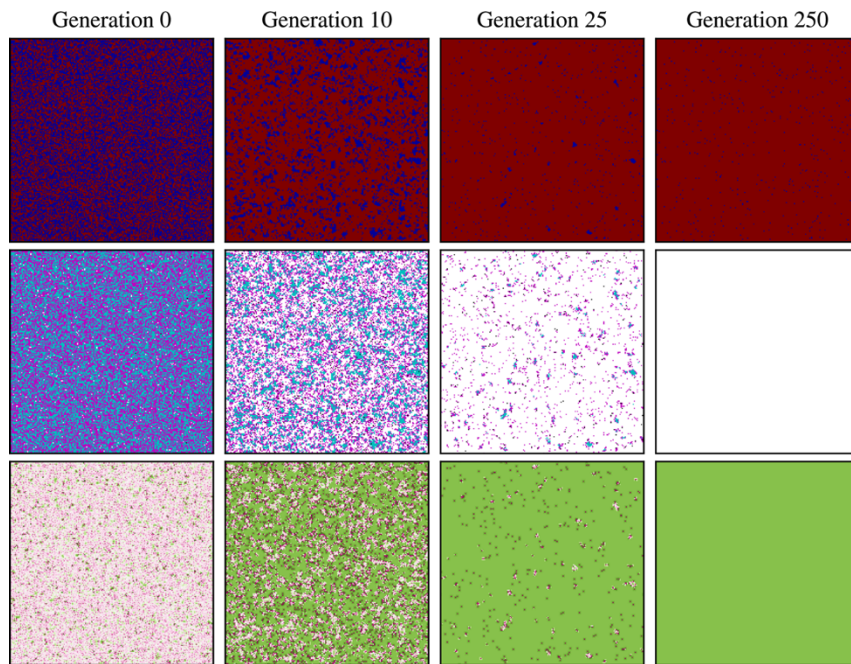
725

726

727

728

729

730

731

732

733

**Extended Data Figure 7: The effect of local reciprocity on two competing factions of GANDHI over time.** Strategies are randomly assigned at the start. Both factions follow the GANDHI strategy, but cooperate with any player that cooperated with them during the last encounter of these two players. **a**, Average score of players over time. **b**, Fraction of players seen as bad by both GANDHI factions over time. **c**, Reputation difference (top tiles) and last score (lower tiles) at different times of the simulation. Frequent strategy changes lead to some players becoming bad — they only cooperate with their neighbors due to direct reciprocity and hence cannot defect against bad players. However, direct reciprocity ensures that eventually, all players cooperate despite the bad reputation, which leads to a high average score.

|Generation 0|Generation 10|Generation 25|Generation 250|

**Extended Data Figure 8: G**ANDHI**++ loses to G**ANDHI **in a direct competition.** The strategy distribution (top tiles), reputation difference (middle tiles) and player score (bottom tiles) in a typical run of GANDHI++ (blue) against GANDHI (red) with probability $h = 0.1$ per duel for contact with the global authorities `Evil` and `Virtue`. The GANDHI++ population quickly collapses and is taken over by GANDHI.