

# The evolution of moral rules in a model of indirect reciprocity with private assessment

Cedric Perret<sup>a,\*</sup>, Marcus Krellner<sup>a</sup>, and The Anh Han<sup>a</sup>

<sup>a</sup>Teesside University, Southfield Rd, Middlesbrough TS1 3BX,  
United Kingdom

\*Corresponding author

## Abstract

Moral rules allow humans to cooperate by indirect reciprocity. Yet, it is not clear which moral rules best implement indirect reciprocity and are favoured by natural selection. Previous studies either considered only public assessment, where individuals are deemed good or bad by all others, or compared a subset of possible strategies. Here we fill this gap by identifying which rules are evolutionary stable strategies (ESS) among all possible moral rules while considering private assessment. We develop an analytical model describing the frequency of long-term cooperation, determining when a strategy can be invaded by another. We show that there are numerous ESSs in absence of errors, which however cease to exist when errors are present. We identify the underlying properties of cooperative ESSs. Overall, this paper provides a first exhaustive evolutionary invasion analysis of moral rules considering private assessment. Moreover, this model is extendable to incorporate higher-order rules and other processes.

## 1 Introduction

2 Morality states which action can be considered good, which action is deemed  
3 to be rewarded and which action should be punished. Moral rules are pervasive  
4 in human societies. They can be observed in a range of examples from the  
5 behaviours of eight month-old infants [1] to the moral norms of societies [2].

6 The pervasiveness of these rules could be explained by their capacity to create  
7 cooperation by indirect reciprocity [3]. Indirect reciprocity describes a form of  
8 reciprocity where the action of an individual is reciprocated by a third party in  
9 future interactions. In its simplest form, cooperators get rewarded by receiving  
10 future cooperation, and defectors get punished by future defection [4]. Indirect  
11 reciprocity can be beneficial as it is one of the few mechanisms that can create  
12 cooperation [3, 5, 6] even when the interactions are not repeated or between  
13 related individuals.

14 Indirect reciprocity explains well the emergence of moral rules but it is not  
15 clear which moral rules best implement indirect reciprocity, and thus, which  
16 moral rules should be observed in the real world. The number of possible rules  
17 can quickly become staggering. When individuals judge the action of another,  
18 they can take into account not only the action observed, but also the reputation  
19 of the individuals involved in the interaction. For instance, helping someone is  
20 generally seen as a positive action, but helping a criminal can be deemed bad.  
21 Do individuals use only a few rules among all the possible rules, or do a wide  
22 variety of rules coexist? Are there features common to all these rules and if yes,  
23 what are they? For instance, ones could expect that successful rules are simple  
24 ones as observed in direct reciprocity [7, 8], while others argued that rules could  
25 be so complex that it drove the evolution of larger brains [9].

26 Tackling this problem, previous works have compared the evolutionary suc-  
27 cess of a large number of rules. Their results show that only few strategies  
28 stand out in term of evolutionary success and the frequency of cooperation  
29 they enforce [4, 10, 11]. These previous works have been a major contribution  
30 but its conclusions are limited. First, they did not consider the evolution of  
31 different assessment rules, i.e. how an individual is judged. Assessment rules  
32 were fixed in a group by moral norms, and all individuals within a group judge  
33 someone else actions in the same way. Although individuals within a group can  
34 share moral rules because they conform to common norms, evidence suggests  
35 that moral rules are also strongly determined by individual characteristics and  
36 thus, can differ between individuals. For instance, infants [12, 13] and toddlers,  
37 which were almost not exposed to social norms, already exhibit signs of indirect  
38 reciprocity [14, 15]. Second, these previous models consider that the opinions  
39 and assessments are public. This means that individuals are considered either  
40 exclusively good or exclusively bad by all others at a given time. Yet, in the  
41 real world, individuals can disagree in their judgements because they have dif-  
42 ferent moral rules or because they get contradictory information. For instance,

43 hunter-gatherers exhibit reciprocity and moral [16] , but often disagree on who  
44 exhibit these moral values [17].

45 The limits of the assumption of public assessment are well acknowledged but  
46 models considering private assessment have been limited by analytical complex-  
47 ity. Indeed, disagreement between individual opinions can itself result in future  
48 interactions being judged differently by the actor and an observer, fuelling more  
49 disagreement. As a result, previous work that considered private assessment or  
50 individual assessment rules limited their analysis to a small set of strategies,  
51 usually the ones that have been shown successful in models with public assess-  
52 ment [18-21]. Recently, [22] developed a model to explore the success of a large  
53 number of assessment rules against strategies that always cooperate or defect  
54 [23]. Yet, an exhaustive study which confronts all possible rules with each other  
55 is still missing.

56 In this paper, we aim to fill this gap by identifying the evolutionary stable  
57 strategies among all possible moral rules. The contributions of this paper are  
58 two-fold. First, we provide the first exhaustive evolutionary invasion analysis of  
59 moral rules considering private assessment. We show that few moral rules stand  
60 out and we identify the common features of these rules. Second, we provide a  
61 model which describes the dynamics of opinions and provide the frequency of  
62 cooperation of an individual given its strategy with private assessment. This  
63 model can be extended to incorporate higher order rules and other processes,  
64 e.g. communication [24].

## 65 **Model description**

66 The model describes a well-mixed and infinitely large population of individuals  
67 that play a one-shot dyadic donation game. In this game, a randomly chosen  
68 individual called 'donor' decides whether to cooperate with another randomly  
69 chosen individual called 'recipient'. If the donor cooperates, **it** pays a cost  $c$  to  
70 provide a benefit  $b$  to the receiver, while if **it** defects nothing happens. This  
71 game is a social dilemma as we consider  $b > c$ , because all would benefit if all  
72 individuals donated, but individuals are not willing to pay the cost.

73 Individuals hold private opinions on each other individual except themselves.  
74 Opinions are either 1 or 0, i.e. good or bad. Individuals use these opinions to  
75 apply their strategies. A strategy consists of a set of action rules,  $A$ , and two

76 sets of assessment rules,  $C$  and  $D$ ,

$$A = \begin{pmatrix} a_1 \\ a_0 \end{pmatrix}, \quad C = \begin{pmatrix} c_{11} \\ c_{10} \\ c_{01} \\ c_{00} \end{pmatrix}, \quad D = \begin{pmatrix} d_{11} \\ d_{10} \\ d_{01} \\ d_{00} \end{pmatrix}, \quad (1)$$

77 where  $a_i, c_{ij}, d_{ij} \in \{0, 1\} \forall i, j \in \{0, 1\}$ . The action rules determine how the  
 78 individual will behave when it is chosen as a donor and meets a recipient it  
 79 thinks to be good ( $a_1$ ) or bad ( $a_0$ ). For example,  $a_1 = 1$  means to cooperate  
 80 with someone good and  $a_0 = 0$  to defect with someone bad. The assessment  
 81 rules determine how the individual updates its opinions when it observes an  
 82 interaction between two other players. The action of the donor and the current  
 83 opinions of the observing individual (toward the two observed players) are con-  
 84 sidered. For example, the rule  $c_{10} = 1$  means that a good donor cooperating  
 85 with a bad recipient is regarded as good afterwards. **Assessment rules are di-**  
 86 **vided into two here for simplicity ( $C$  applies to the case where donor cooperates**  
 87 **and  $D$  applies to the case where donor defects).** The opinion about the recipi-  
 88 ent is not updated. Errors may occur during assessment or while implementing  
 89 an action. Following literature [25], we consider (i) execution errors, at a rate  
 90  $\mu_e$ , where an individual does the opposite of what it intended (i.e. determined  
 91 by her strategy) and (ii) assessment errors, at a rate  $\mu_a$ , where an individual  
 92 assigns the opposite opinion of what her assessment rules would suggest.

93 A large number of strategies are possible, and each strategy can differ in its  
 94 evolutionary success. We want to compute the evolutionary success of different  
 95 strategies and see if particular strategies stand out. Our method consists in  
 96 deriving the long-run average proportion of good opinions others have on an  
 97 individual, i.e. the individual's h-score. We use this h-score to calculate the  
 98 frequency of cooperation from and towards this individual which determinate  
 99 its fitness. We apply this method to compare the h-score and the fitness of  
 100 individuals in a population with a resident strategy and a single mutant strategy  
 101 to perform an ESS analysis.

## 102 **Results**

### 103 **Monomorphic population**

104 We confirm that the analytical model correctly approximates the h-score and  
105 the probability of cooperation at equilibrium, by comparing the predictions of  
106 the analytical model with agent-based simulations for any possible strategies in  
107 a monomorphic population. Some pairs of strategies are equivalent (as formally  
108 defined in the mirror image section of the extended method). They simply assess  
109 and act upon opinions in an opposite way (what one would call good is called  
110 bad by the other). After keeping one instance of each pair, we are left with 258  
111 strategies including 256 discriminator strategies, unconditional cooperator and  
112 unconditional defector. The simulations implement the aforementioned model  
113 with a population of 100 individuals and one observer per interaction. The  
114 results of the simulations are taken after  $4 \times 10^5$  interactions, and averaged over  
115  $10^5$  interactions and 30 independent replicates. We run analysis considering  
116 that (i) the error rate is negligible and (ii) the error rate is not negligible. In  
117 the former, we consider that the error rate is equal to 0 in the analytical model  
118 and we keep a very low error rate in simulation (namely,  $10^{-4}$ ). In the latter,  
119 we do not vary independently the execution and assessment error rates because  
120 we are interested in testing the robustness of the conclusion obtained from the  
121 model without errors, rather than describing the effect of a particular type of  
122 error.

123 The results show that the analytical model well approximates the h-score  
124 and the probability of cooperation at equilibrium. The mean difference between  
125 the h-score predicted and simulated is 0.014 in absence of error and 0.005 in  
126 presence of error. The mean difference between the frequency of cooperation  
127 predicted and simulated is 0.01 in absence of error and 0.0009 in presence of  
128 error. The detailed results are provided in **SI**. The results of the simulations  
129 are illustrated in supplementary Figure S1 and summarised in supplementary  
130 Figure S2.

### 131 **Evolutionary invasion analysis**

132 We now use the analytical model to conduct an evolutionary invasion analysis (in  
133 short, ESS analysis). As common assumptions in ESS analysis [26], we assume  
134 that i) mutations are rare and thus, there is at most one mutant strategy  $m$  at  
135 a time in a population of individuals with resident strategy  $r$ , ii) the mutant's

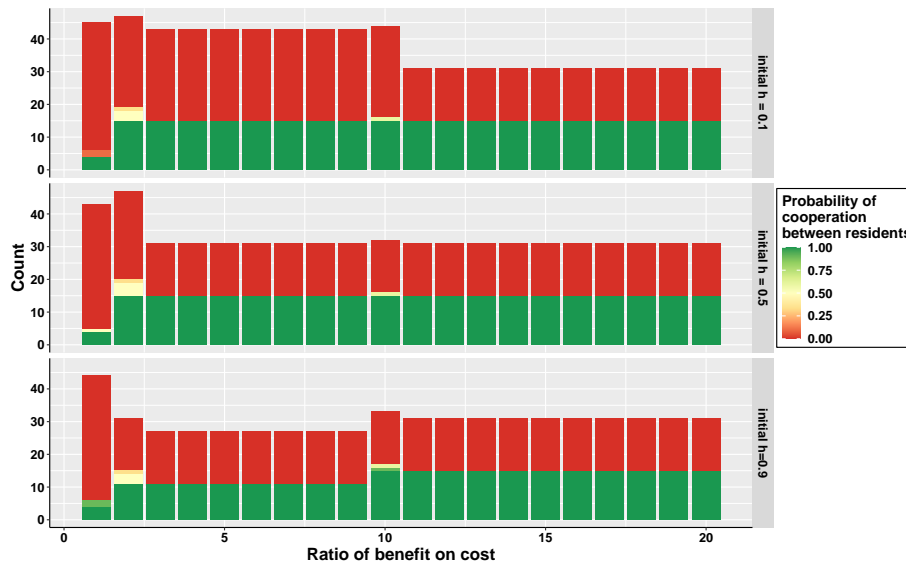


Figure 1: Number of ESS as a function of the benefit to cost ratio,  $b/c$ . The colour represents the probability of cooperation between residents. The results are presented for different initial h-scores  $h(t_0)$  (0.1, 0.5 and 0.9, in top, middle and bottom rows, respectively).

136 effect is negligible on the dynamics and iii) **population size is large enough so that**  
 137 **stochasticity in selection is negligible.** To know if a strategy can be invaded or  
 138 not by another, we compute the difference of absolute fitness between a mutant  
 139 strategy in a population of resident strategy. If the fitness of the mutant is lower  
 140 or strictly equal, the mutant disappears and the resident resists invasion. If a  
 141 strategy resists invasion from all other possible strategies, it is an ESS. **Unlike**  
 142 **the previous section, fitness and h-score are now directly computed rather than**  
 143 **simulated.** We consider that any differences in fitness less than  $10^{-4}$  are equal  
 144 to 0 because of the floating point error.

#### 145 **In absence of errors**

146 First, we consider that the errors are negligible, that is  $\mu_a = 0$  and  $\mu_e = 0$ .  
 147 Figure 1 shows that there are multiple evolutionary stable strategies (ESS),  
 148 that is strategies that can not be invaded by others. Supplementary Figure S3  
 149 shows that some strategies are not ESS for all studied initial h-score  $h(t_0)$ . We  
 150 focus on strategies that are ESS for all three initial h-scores.

151 The ESS can be divided in two groups, strategies which cooperate and avoid

Name	$a_1$	$a_0$	$c_{11}$	$c_{10}$	$c_{01}$	$c_{00}$	$d_{11}$	$d_{10}$	$d_{01}$	$d_{00}$	Minimum ratio benefit on cost
C-1	1	0	1	1	1	1	0	1	0	1	1
C-2	1	0	1	1	1	1	0	1	0	0	1
C-3	1	0	1	1	1	1	0	0	0	1	1
C-4	1	0	1	1	1	0	0	1	0	1	1
C-5	1	0	1	1	1	1	0	1	1	1	2
C-6	1	0	1	1	1	1	0	1	1	0	2
C-7	1	0	1	1	1	1	0	0	1	1	2
C-8	1	0	1	1	1	1	0	0	1	0	2
C-9	1	0	1	1	1	0	0	1	1	1	2
C-10	1	0	1	1	1	0	0	1	1	0	2
C-11	1	0	1	1	1	0	0	0	1	1	2
C-12	1	0	1	1	0	1	1	1	0	1	$1/(1-h(t0))$
C-13	1	0	1	1	0	1	1	1	0	0	$1/(1-h(t0))$
C-14	1	0	1	1	0	1	1	0	0	1	$1/(1-h(t0))$
C-15	1	0	1	1	0	0	1	1	0	1	$1/(1-h(t0))$

Figure 2: List of strategies that are cooperators and ESS for any initial h-score and at least one value of the benefit to cost ratio,  $b/c$ . The last column represents the minimum ratio for which the strategy is ESS for any initial h-score.

152 exploitation, and those which defect and efficiently exploit others. There are  
153 strategies that have an intermediary probability of cooperation but they are  
154 only ESS for a specific benefit to cost ratio so we do not discuss them further  
155 here. We present the 15 strategies that are ESS and cooperators in Figure 2,  
156 with the minimum benefit to cost ratio required for a strategy to be an ESS. We  
157 present the 38 strategies that are defectors in supplementary Figure S4, with  
158 the maximum benefit to cost ratio required for a strategy to be ESS. We call  
159 the ESS cooperator and defector strategies respectively C-\* and D-\*. We name  
160 each ESS cooperator, with C1 representing the first ESS cooperator in the table,  
161 C2 the second, and so on.

162 First, we look closely at the ESS cooperators. In term of behaviours, a  
163 distinctive feature of the ESS cooperators is that they fully cooperate with each  
164 other while cooperating less with mutant defectors. By cooperating with each  
165 other, they sustain the highest possible fitness for cooperators. By cooperating  
166 less with mutant defectors, they limit the fitness of the mutant to be less than or  
167 equal to their fitness, providing that the benefit of cooperation is high enough.

168 First, all the ESS cooperators have  $c_{11} = 1$  and  $c_{10} = 1$ . It means that they

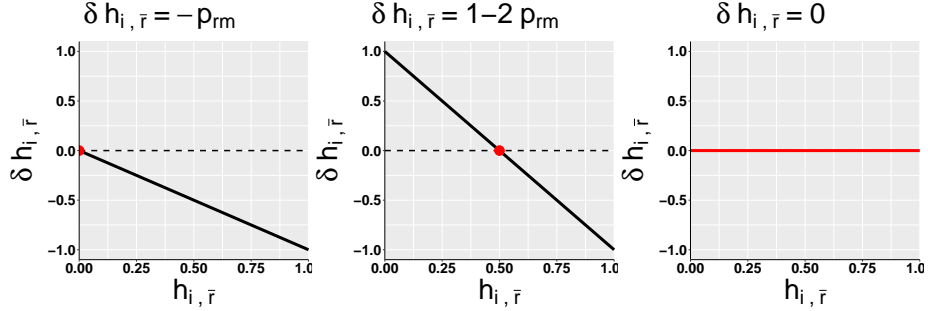


Figure 3: Differential equation of h-score of resident on mutant for the three types of ESS cooperators when mutant is always defect (AllD). From left to right, strategies that are ESS for a ratio of 1, 2 and  $\frac{1}{1-h(t0)}$ .

169 consider that cooperation from an individual seen as good is always rewarded  
 170 by future cooperation. Ultimately, this results in a population of individuals  
 171 which see each other as good and always cooperate with each other. This allows  
 172 the strategies to maintain cooperation once established.

173 Second, the ESS cooperators have either  $d_{11} = 0$  or  $d_{01} = 0$ , or both. This  
 174 means that they will consider individuals defecting towards good individuals as  
 175 (partly or totally) bad. Because the ESS cooperators consider each other as  
 176 good once h-score of 1 is reached, this allows ESS cooperators to defect with  
 177 individuals that defect with them. ESS cooperators differ in their capacity  
 178 to efficiently reciprocate defection. First, there are strategies that have both  
 179  $d_{11} = 0$  and  $d_{01} = 0$ . They are ESS on the whole range of benefit to cost ratios  
 180 (C1-C4). It is because they have an average opinion of 0 on defectors, and thus  
 181 will always defect with them, as shown in the left panel of Figure 3. Second,  
 182 there are strategies that have both  $d_{11} = 0$  and  $d_{01} = 1$ , which are ESS if the  
 183 benefit is at least twice larger than the cost (C5-11). These strategies have half  
 184 of the time good opinion (and cooperating) with mutant defectors, and half  
 185 of the time bad opinion (and defecting) with mutant defectors. Finally, there  
 186 are strategies which have  $d_{11} = 1$  and  $d_{01} = 0$  (C12-15) and for which their  
 187 evolutionary stability depends of the initial h-score. For instance, they are ESS  
 188 for a ratio  $b/c > 10$  if the initial h-score is 0.9 or ESS for a ratio  $b/c > 2$  if the  
 189 initial h-score is 0.5. These strategies have in common that their opinions of  
 190 mutant defecting with good individuals remain roughly the same. For instance,  
 191 one strategy gives 0 to bad individuals defecting and gives 1 to good individuals  
 192 defecting. Thus, the frequency of cooperation received by mutant defectors is



193 approximately the initial h-score  $h(t_0)$  and their fitness  $h(t_0)b$ . Strategies C12-  
194 15 can not be invaded by a defector when their fitness  $b - c > h(t_0)b$ . This  
195 equation can be rearranged, leading to  $b/c > 1/(1 - h(t_0))$ .

196 To summarise, the rules of ESS cooperators make them efficiently reciprocate  
197 once cooperation is fully established. The first pattern maintains cooperation  
198 while the second makes them defect with mutant defectors. Yet, this does not  
199 assure that they reach cooperation in the first place. For that, we can observe  
200 that ESS cooperators, besides  $c_{11} = 1$  and  $c_{10} = 1$ , judge a number of other  
201 encounters as good. This number and the type of encounters can vary but they  
202 judge enough cases as good so that the h-score of an individual with a resident  
203 strategy seen by other individuals with the same strategy, increases towards 1  
204 (the differential equation is always positive). This ensures that they go towards  
205 full cooperation even in presence of initial disagreement. For instance, the first  
206 three strategies (C1-C3) consider that cooperating is good and at least one other  
207 case as good. Because these strategies cooperate with a probability  $h$ , it ensures  
208 that the differential equation is always positive  $h + p(d_{?0}) - h \geq 0$ . Another  
209 example is C4, which might appear surprising as they consider one cooperation  
210 as bad  $c_{00} = 0$ . However, this case is very rare and it leads to the differential  
211 equation remaining positive.

212 In short, the strategies that are evolutionary stable and cooperators have  
213 rules that (i) establish full cooperation with each other, (ii) sustain full coop-  
214 eration when established, and (iii) reduce the frequency of cooperation with  
215 mutant defectors. Note that the presence of a single of these features in a  
216 strategy does not mean that the strategy will be ESS. Indeed, we looked at the  
217 common patterns among the ESS rather than correlating the rules with the suc-  
218 cess of strategies. Finally, we observe a diversity of rules because first, strategies  
219 can differ in their capacity to defect with defectors, and second, different rules  
220 can lead a population to full cooperation on the long term.

221 We now look at defector strategies that are ESS. Again, there are numerous  
222 strategies fulfilling these conditions but they have similar behaviour. Their  
223 distinctive feature is that they have a lower probability to cooperate with the  
224 mutant, than the mutant have with the resident. The rules have in common  
225 the pattern that  $d_{10} = d_{00} = 0$ . In other words, they always defect with  
226 individuals defecting with individuals they see as bad. This allows them to avoid  
227 cooperation when individuals do not cooperate with them. Again, defectors can  
228 be separated into different types as a function of the maximum ratio of benefit  
229 to cost required for them to be ESS. First, the defectors that are ESS for the

230 whole range of benefit studied (D1 - D16) never cooperate with each other nor  
231 with the mutant. They are behaviourally equivalent to strategies that always  
232 defect. This means that they do not pay any cost and thus no strategies can  
233 have a higher fitness than them. They also all have in common that  $c_{10} = c_{00} =$   
234 0. This means that any strategy interacting with them (that they consider  
235 as bad) will be considered bad, and receive future defection. Second, some  
236 defectors are ESS only for a very limited range of the benefit to cost ratio.  
237 These strategies cooperate with mutant, but at a lesser rate than mutant with  
238 resident. This means that mutant can invade if the benefit they received by  
239 cooperation outweighs the cost of their cooperation. Without going into the  
240 details, these strategies do not have the rules of  $c_{10} = c_{00} = 0$  as previous  
241 strategies, and thus are not perfect defectors.

242 We notice that there are a large number of cases of polymorphism among  
243 these strategies. In the cases where these strategies can be invaded, there are  
244 between 75.5% to 88.4% of cases with polymorphism against 5% to 17% when  
245 looking at any strategies. The reason is that if the benefit that mutant provides  
246 to resident is negligible when mutant are a minority, it is not the case anymore  
247 when they compose most of the population. Thus, defector strategies that are  
248 ESS for only a limited set of the benefit to cost ratios could still be frequent for  
249 other ratios.

### 250 **In presence of errors**

251 When the errors are not negligible, the previously identified ESS are not evo-  
252 lutionary stable anymore except for always defect (AllD). This is because the  
253 errors in assessment lead discriminating cooperators to cooperate less, and dis-  
254 criminating defectors to defect less. For instance, the previously ESS cooper-  
255 ators maintain a lower level of cooperation between each other (and thus are  
256 easier to invade). In addition, the errors create disagreement and can have an  
257 effect on the long term. As a result, C-\* can cooperate more with mutant than  
258 with themselves, even when mutant are strong cooperators e.g. judge good  
259 all except  $d_{00}$ . This is because the cooperation of such strategy breaks down  
260 only in specific cases, which allow them to quickly get their h-score back to  
261 1. AllD is still an ESS even if it sometimes cooperates by mistake, because it  
262 is not affected by assessment errors and thus has still the lowest frequency of  
263 cooperation possible.

264 To gain further insights, we look at the difference of fitness between mutant

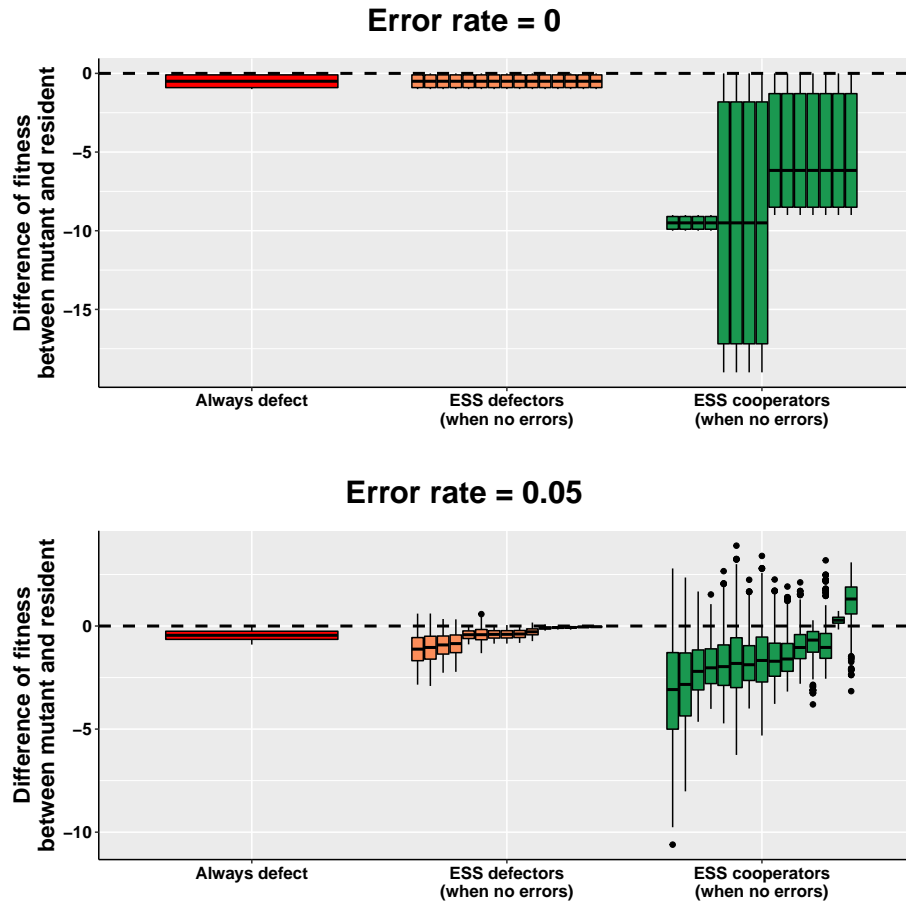


Figure 4: Difference of fitness between mutant and resident  $w_m - w_r$ , for different strategies that are ESS when there are no errors. We differentiate between strategies that were cooperators, defectors and the strategy that always defect (AllD), which is the only ESS in presence of errors. The results are presented for a high benefit to cost ratio ( $b/c = 20$ ) to highlight the difference of fitness. Results for other, smaller benefit to cost ratios, can be found in supplementary figures S5 and S6, showing the same trend.

265 and resident for different resident strategies. This difference of fitness gives us  
266 hints on the success of strategies when evolution is stochastic, that is when in-  
267 vasion is a probability based on the difference of fitness. We show the results  
268 for a high benefit to cost ratio in Figure 4 and for other ratios in supplementary  
269 Figures S5 and S6. First, Figure 4 shows that in absence of error, ESS coop-  
270 erators have a higher difference of fitness as the benefit increases. This result  
271 suggests that ESS cooperators could be more prevalent than ESS defectors when  
272 selection is weak and if the benefit of cooperation is sufficiently high. Second,  
273 Figure 4 shows the same (but weaker) trend when errors are not negligible. In  
274 particular, the difference of fitness is higher for C15 than the only ESS with  
275 error: always defect against the majority of mutants. This suggests that if in  
276 presence of errors, C15 is not ESS anymore, it could still be a very frequent  
277 strategy (and more frequent than always defect).

## 278 Discussion

279 Among the large number of possible moral rules, previous work shows that only  
280 a few of them stand out and should be observed in the real world [11]. Yet,  
281 models studying the evolution of moral rules considered either public assessment  
282 or a limited number of strategies and it still lacks of an exhaustive evolutionary  
283 analysis of moral rules with private assessment. In this paper, we fill this gap  
284 by building an analytical model to describe the change in opinions as a function  
285 of time. We used this model to study the invasibility of any strategies by any  
286 other strategies up to **third-order** assessment rules, and identify the evolutionary  
287 stable strategies.

288 Previous results suggested that considering private information breaks down  
289 cooperation and limits the evolution of cooperative moral rules by creating  
290 disagreement [27]. However, our results show that there are evolutionary stable  
291 rules implementing cooperation even when assessment is private. This result is  
292 explained by the fact that some rules are capable ~~to suppress~~ to suppress disagreement on  
293 the long term. Second, our results show that the number of ESS in our study  
294 can even be higher than the number of ESS previously found when considering  
295 public information [11]. This is because multiple rules can end up implementing  
296 the same level of cooperation on the long term. For instance, strategies C-1 and  
297 C-3 differ in their rules about good individuals defecting with bad individuals  
298 ( $d_{10}$ ) but they still end up with full cooperation at equilibrium. Note that these

Table 1: Presentation of the main strategies identified in the literature, i.e. image scoring, standing strategy and the leading eight; and the ESS cooperators identified in our analysis C1-15. \* mark wildcards. Note that all C1-15 have an additional restriction:  $p(o, r, r) > h_{r, \bar{r}}$ , which ensures that the h-score  $h_{r, \bar{r}}$  increases up to 1 (see Equation 4). This means that a maximum of one of the wildcards can be equal to 0 for C1-4 and C12-15, and a maximum of two wildcards can be equal to 0 for C5-11.

	$c_{11}$	$c_{10}$	$c_{01}$	$c_{00}$	$d_{11}$	$d_{10}$	$d_{01}$	$d_{00}$
Image scoring	1	1	1	1	0	0	0	0
Standing strategy	1	1	1	1	0	1	0	1
Leading 1-8	1	*	1	*	0	1	0	*
C1-4	1	1	1	*	0	*	0	*
C5-11	1	1	1	*	0	*	1	*
C12-15	1	1	0	*	1	*	0	*

299 conclusions rely on the assumption that interactions are long enough so that  
300 initial disagreement is negligible. When interactions are shorter, it is likely that  
301 other mechanisms are required for indirect reciprocity to evolve such as empathy  
302 [20] or public institutions [21].

303 Second, we identify the important properties of the rules that are ESS and  
304 cooperators. These strategies consider (i) good cooperators as good, (ii) all or a  
305 part of defectors towards good individuals as bad, and (iii) a varying number of  
306 other cases as good. The last rule allows the strategy to converge toward a full  
307 population of good cooperators, and the first two rules allow them to efficiently  
308 reciprocate once good reputation is established.

309 How do these successful rules compare to rules previously identified in the  
310 literature? To answer this question, we present the main rules in the compar-  
311 ative Table 1. A significant part of previous work has focused on two famous  
312 strategies, image scoring [4], that is cooperate with cooperators and defect with  
313 defectors, and standing [10, 28], that is cooperate with cooperators, cooperate  
314 with defectors towards defectors and defect with defector towards cooperators  
315 (see Table 1). Image scoring was historically one of the first strategies to suc-  
316 cessfully implement indirect reciprocity [4], but later work showed that standing  
317 is more evolutionary successful [10]. Our results concur as we found the stand-  
318 ing strategy to be evolutionary stable for any benefit superior to cost (standing  
319 strategy is C1), while image scoring is not an ESS. Note that image scoring  
320 would be an ESS if initial h-score is exactly 1.

321 In addition, our results show that the important rules of the standing strat-

322 egy are "cooperate with good cooperators" and "defect with defectors against  
323 good", and that the part "cooperate with defectors against bad" can vary. This  
324 distinction provides important new insights into the ongoing debate. Indeed,  
325 some experimental evidence supports the presence of standing strategy in nat-  
326 ural populations [29] while others [30] appear to be against. In laboratory  
327 experiments [30], researchers compared the amount of defection received by an  
328 unconditional defectors and a "justified" defectors, that is individuals that de-  
329 fect with previous defectors. Their results showed that the difference is not  
330 strong enough to support standing strategy. As pointed out [10], these conclu-  
331 sions could be limited as they do not measure the amount of defection received  
332 by individuals that refuse to help previous donors. Our results show rigorously  
333 here that it is this amount of defection towards defectors against good indi-  
334 viduals that matters for the success of the standing strategy rather than the  
335 cooperation toward justified defectors as measured in the experiments. Thus,  
336 our results suggest that further experiments with different measures is required  
337 to reject or accept the prevalence of standing strategy.

338 Our work also follows the exhaustive evolutionary analysis which showed  
339 eight successful strategies (called leading eight [11]). A direct comparison of the  
340 leading eight and the ESS described here is limited because this previous study  
341 focused on the evolutionary success of different action rules, i.e. what is the best  
342 action rule for a given assessment rule, while the model presented here focuses  
343 on assessment rules, i.e. how to judge someone. Yet, it can shed light on main  
344 differences between private and public assessment. First, the C-\* strategies  
345 require that  $c_{10} = 1$ , a rule which is shared by only half of the leading eight.  
346 This rule is crucial with private assessment to avoid cooperators losing their  
347 good standing. Those leading eight strategies which do not share this rule were  
348 shown to suffer greatly by private assessments before [27]. Second, C\*- rules  
349 also require that enough cases are considered good so that the full population  
350 converge towards being good cooperators. On the other side, C-\* rules can vary  
351 in cases where leading eight can not. For instance, the leading eights always  
352 consider defection towards a good individual as bad. This is shared by the most  
353 successful strategies found here. However, the C-\* strategies can also consider  
354 defection in one situation ( $d_{11}$  or  $d_{01}$ ) as good and still be ESS given that the  
355 benefit of cooperation is high enough. This is because such a rule in public  
356 assessment would lead to all individuals cooperating with defectors while this  
357 happens partially with private assessment.

358 Third, we find that the presence of errors breaks down the evolutionary

359 stability of the previously identified strategies. This is because the property of  
360 these rules that allow them to converge toward full cooperation, also makes them  
361 vulnerable to errors. **This result suggests that private assessment rules could**  
362 **not evolve when errors are frequent, and that public assessment for instance**  
363 **supported by an institution could be preferred in this case [21].** However, this  
364 result is mitigated by two points. First, we have considered that any difference  
365 of fitness, however small, leads to invasion or not. This is a classic assumption  
366 of ESS analyses but in the real world, selection can be weaker and stochasticity  
367 can result in non-ESS to be frequent. Important first steps have been made  
368 by a recent paper which considered stochastic evolution of a population mixing  
369 one discriminator strategy, with unconditional cooperators and defectors [22].  
370 An extension could consider a population of different discriminator strategies  
371 in co-presence. Second, the effect of errors could be suppressed by additional  
372 mechanisms. Evidence shows that not only the outcome of an action plays a  
373 role in assessments but also the intention behind this action [31, 32], and thus  
374 errors in actions could have a limited effect. Another example is the role of  
375 communication and conformity which could counterbalanced the effect of errors  
376 and drive the h-score towards a general agreement. Further work integrating  
377 these mechanisms would provide a more realistic model and test if the strategies  
378 identified here could be frequent in presence of errors.

379 Results from models of indirect reciprocity can be confronted to the donor  
380 game conducted in laboratory experiments. For instance, experiments con-  
381 ducted by [33] showed that information about the partners' previous partners'  
382 reputation increases the level of cooperation. This is in agreement with our  
383 results that all the C-\* use second-order information. Second, recent exper-  
384 iments have studied the strategies employed by individuals [34]. They show  
385 that individuals often requested second-order information, and at a higher fre-  
386 quency when their partner has previously defected. We find some similarity in  
387 our results. All the C-\* require to know the past interactions of their partner  
388 to judge its action when the partner defected, while only eight strategies re-  
389 quire this information when the partner cooperated. This goes to 3 against 1  
390 if we considered the most successful strategies C1-4. Last but not least, they  
391 showed a strong variation in behavioural strategies. This is in line with ours  
392 results, which show that diverse behavioural strategies can be employed. How-  
393 ever, these comparisons are limited as our model considers a large group size  
394 and long interactions, which are both assumptions often absent in laboratory  
395 experiments. A more promising path to test our results would be in study in

396 natural populations.

397 We have made a number of assumptions in this model that need to be  
398 discussed. First, we approximated the reputation dynamics by a deterministic  
399 approach. This required two main assumptions, that the size of the population is  
400 infinitely large and that the number of observers is finite. The first assumption  
401 means that the results in this paper are applicable only to cases where the  
402 population is sufficiently large. The second assumption results naturally from  
403 physical limits of the number of observers, e.g. it is likely that an increase of  
404 ten fold of group size does not mean an increase of ten fold of the number of  
405 observers. However, it is important to note that there are possible exceptions, in  
406 particular systems where actions are widely shared e.g. e-commerce or medieval  
407 merchant guilds [35].

408 Second, we have considered that the initial h-scores are the same for all indi-  
409 viduals, including the resident and mutant strategies. In real world, the opinion  
410 of an individual on a newly met individual could be part of the individual strat-  
411 egy (in the same way that tit-for-tat could play cooperate or defect at the first  
412 round). We did not consider this here to keep the number of strategies reason-  
413 able and we focused on the strategies that are ESS for diverse initial h-scores.  
414 Future work could integrate the initial h-score in the strategy and replicate the  
415 evolutionary analysis.

416 To conclude, the contribution of this paper is two fold. First, it provides a  
417 first exhaustive evolutionary analysis of moral rules with private assessment. It  
418 provides more realistic results, as a large number of real-world situations (in-  
419 cluding most of laboratory experiments) includes private assessment. Second, it  
420 provides an analytical model that describes the opinion dynamics when assess-  
421 ment is private and allows further investigation of the issue accurately and at  
422 a faster speed than with simulations, enabling exhaustive analyses. The model  
423 can easily be extended to integrate other mechanisms. A natural progression  
424 of this work is (i) to study strategies up to second-order action rules **where**  
425 **action also depends of actor's reputation** to compare results with the previous  
426 exhaustive evolutionary analysis with public assessment [11], and (ii) integrate  
427 the effect of communication and conformity [36] as it plays a prevalent role in  
428 indirect reciprocity [24] and can be easily integrated in the model [37].



429 **Method**

430 We build a deterministic model that approximates the average fitness of an  
 431 individual of a given strategy. We first consider a *monomorphic* population  
 432 where all individuals have the same resident strategy. We do so to introduce  
 433 the method in its simplest form. We consider that the number of interactions  
 434 is large enough, and thus, the fitness  $w_i^*$  of a focal individual  $i$  is its average  
 435 payoff ( $N$  is the population size):

$$w_i^* = \frac{1}{N-1} \sum_{r=1}^{N-1} (bp^*(c_{i,r}) - cp^*(c_{r,i})). \quad (2)$$

436 The fitness of an individual  $i$  is the benefit  $b$  received when other individuals  
 437 cooperate with the individual  $i$ , discounted by the cost  $c$  paid when the focal  
 438 individual  $i$  cooperates. The probability that an individual  $r$  cooperates with  
 439 individual  $i$  is denoted by  $p(c_{i,r})$ . The superscript  $*$  denotes that the fitness  
 440 and probability of cooperation considered are at equilibrium. This probability  
 441 itself depends on the many opinions that individuals have on each other, which  
 442 is difficult to track analytically. Instead of describing all the opinions, we define  
 443 a *h-score* of an individual  $i$  as the proportion of other individuals with opinion  
 444 1 on  $i$ , or the average reputation of  $i$ .

445 The h-score is useful because considering that the number of individuals is  
 446 large enough and that the donor, recipient and observers are chosen randomly,  
 447 the h-score also describes the probability that a random individual has an opin-  
 448 ion of 1 on the focal individual  $i$ , that is  $h_{i,\bar{r}} = p(o_{i\bar{r}} = 1)$ . Thus, we can combine  
 449 h-score and the action rules which describe how individuals act upon a given  
 450 encounter, to describe the probability of cooperation:

$$p^*(c_{i,r}) = h_{i,r}a_1 + (1 - h_{i,r})a_0. \quad (3)$$

451 Similarly, using the assessment rules, we can calculate the probability that h-  
 452 score increases or decreases after an interaction  $p(o_{r,r,r})$ , and thus, the dynamics  
 453 of h-score over time. The formula is given in Equation 9 in the detailed method  
 454 in **SI**. Execution and assessment errors modify respectively the probability of  
 455 cooperation or the probability of h-score to increase or decrease after an inter-  
 456 action as described in the equations 12 and 13 of the extended method section  
 457 (**SI**).

458 So far, we derived the change in h-score for an individual with a given strat-

459 egy but we would like to derive the change in h-score for any individuals. Let  
 460 us note that because individuals have the same strategy, the direction of change  
 461 is similar across individuals and their h-score will converge towards the same  
 462 equilibrium points. In addition, we make the assumption that the number of  
 463 observers is small and independent of the population size. We also assume that  
 464 the initial h-score of all individuals are the same. Following these two assump-  
 465 tions, the differences in h-score between individuals due to stochasticity is small  
 466 and negligible on the dynamics. Because the change in h-score is very small af-  
 467 ter an interaction, it can be approximated by the following differential equation  
 468 [38] (see details in SI)

$$\frac{d(h_{r,\bar{r}})}{dt} = p(o_{r,r,r}) - h_{r,\bar{r}}. \quad (4)$$

469 The average h-score at equilibrium can be found by solving the equation  $\frac{d(h_{r,\bar{r}})}{dt} =$   
 470 0. This equation is a polynomial of  $h_{r,\bar{r}}$  of a maximum degree of three (see Equa-  
 471 tions 19 to 22 in **SI**). The stability of equilibrium points is determined by looking  
 472 at the sign of the derivative at the equilibrium points [26].

473 We now extend the analytical model to conduct an evolutionary invasion  
 474 analysis (in short, ESS analysis). To know if a strategy can be invaded or  
 475 not by another, we need to compute the difference of absolute fitness between a  
 476 mutant strategy in a population of resident strategy. If the fitness of the mutant  
 477 is greater than that of the resident, the mutant invades the population and  
 478 becomes resident. If the fitness of the mutant is lower, the mutant disappears  
 479 and the resident resists invasion. When the two values of fitness are equal, the  
 480 resident also resists invasion because in an infinitely large population, a mutant  
 481 strategy can not invade by drift. If two strategies can mutually invade, there  
 482 will be polymorphism.

483 The difference of fitness between a mutant  $w_m$  and a resident  $w_r$  is given as  
 484 follows:

$$\Delta w = w_m - w_r = p^*(c_{m,r})b - p^*(c_{r,m})c - p^*(c_{r,r})(b - c). \quad (5)$$

485 The fitness of the mutant is the benefit received when a resident cooperates with  
 486 the mutant discounted by the cost of the cooperation from mutant to resident.  
 487 There are three different probabilities of cooperation. The probability of coop-  
 488 eration between residents  $p(c_{r,r})$  is calculated as in the case of a monomorphic  
 489 population. The two remaining probabilities of cooperation can be computed as  
 490 previously using h-score and action rules (Equation 15 in **SI**). To find the prob-

491 ability of cooperation at equilibrium, we describe the dynamics of the h-score  
492 as previously

$$\begin{aligned}\frac{d(h_{m,\bar{r}})}{dt} &= p(o_{m,r,r}) - h_{m,\bar{r}}, \\ \frac{d(h_{\bar{r},m})}{dt} &= p(o_{r,r,m}) - h_{\bar{r},m}.\end{aligned}\tag{6}$$

493 The probabilities of h-score to increase after the observation of an interaction,  
494  $p(o)$ , can be described using the h-score and the assessment rules as previously  
495 (Equation 17 in **SI**). This system of two polynomial equations with two un-  
496 knowns are solved numerically. To determinate the stability of the equilibrium  
497 points, we look at the Jacobian matrix at the equilibrium of interest. The equi-  
498 librium is locally stable if the real part of the leading eigenvalue is negative  
499 [26]. Errors are integrated in the same way as in the case of monomorphic  
500 populations.

501 **Acknowledgments.** C.P. and T.A.H. were supported by Future of Life  
502 Institute (grant RFP2-154). T.A.H. is also supported by Leverhulme Research  
503 Fellowship (RF-2020-603/9).

504 **Data availability statement.** The code is available online at  
505 ‘<https://github.com/CedricPerret>’ in the project ‘RepDyn’.

506 **Author contributions.** C.P, M.K and T.A.H designed the research. C.P.  
507 and MK developed the model. C.P. implemented software and carried out the  
508 analysis. C.P, M.K and T.A.H wrote the manuscript

509 **Competing interests.** The authors declare no competing interests.

## 510 References

- 511 1. Hamlin, J. K., Wynn, K., Bloom, P. & Mahajan, N. How infants and  
512 toddlers react to antisocial others. *Proceedings of the National Academy of  
513 Sciences of the United States of America* **108**, 19931–19936 (2011).
- 514 2. Harms, W. & Skyrms, B. in *The Oxford Handbook of Philosophy of Biology*  
515 (Oxford University Press, 2009).
- 516 3. Alexander, R. D. *The biology of moral systems* (ed Routledge) (Taylor &  
517 Francis Inc, Somerset, USA, 1987).
- 518 4. Nowak, M. A. & Sigmund, K. Evolution of indirect reciprocity by image  
519 scoring. *Nature* **393**, 573–577 (1998).

- 520 5. Brandt, H., Ohtsuki, H., Iwasa, Y. & Sigmund, K. in *Mathematics for*  
521 *Ecology and Environmental Sciences* 21–49 (Springer Berlin Heidelberg,  
522 2007).
- 523 6. Nowak, M. A. & Sigmund, K. Evolution of indirect reciprocity. *Nature*  
524 **437**, 1291–1298 (2005).
- 525 7. Axelrod, R. & Hamilton, W. D. The Evolution of cooperation. *Evolution*  
526 **211**, 1390–1396 (1981).
- 527 8. Hilbe, C., Martinez-Vaquero, L. A., Chatterjee, K. & Nowak, M. A. Memory-  
528 n strategies of direct reciprocity. *Proceedings of the National Academy of*  
529 *Sciences of the United States of America* **114**, 4715–4720 (2017).
- 530 9. Dunbar, R. I. The social brain hypothesis. *Evolutionary Anthropology: Is-*  
531 *ssues, News, and Reviews* **6**, 178–190 (1998).
- 532 10. Leimar, O. & Hammerstein, P. Evolution of cooperation through indirect  
533 reciprocity. *Proceedings of the Royal Society B: Biological Sciences* **268**,  
534 745–753 (2001).
- 535 11. Ohtsuki, H. & Iwasa, Y. How should we define goodness? - Reputation  
536 dynamics in indirect reciprocity. *Journal of Theoretical Biology* **231**, 107–  
537 120 (2004).
- 538 12. Olson, K. R. & Spelke, E. S. Foundations of cooperation in young children.  
539 *Cognition* **108**, 222–231 (2008).
- 540 13. Kenward, B. & Dahl, M. Preschoolers distribute scarce resources accord-  
541 ing to the moral valence of recipients’ previous actions. *Developmental*  
542 *Psychology* **47**, 1054–1064 (2011).
- 543 14. Meristo, M. & Surian, L. Do infants detect indirect reciprocity? *Cognition*  
544 **129**, 102–113 (2013).
- 545 15. Hamlin, J. K., Wynn, K. & Bloom, P. Social evaluation by preverbal in-  
546 fants. *Nature* **450**, 557–559 (2007).
- 547 16. Smith, K. M., Larroucau, T., Mabulla, I. A. & Apicella, C. L. Hunter-  
548 gatherers maintain assortativity in cooperation despite high levels of resi-  
549 dential change and mixing. *Current Biology* **28**, 3152–3157 (2018).
- 550 17. Smith, K. M. & Apicella, C. L. Hadza hunter-gatherers disagree on per-  
551 ceptions of moral character. *Social Psychological and Personality Science*  
552 **11**, 616–625 (2020).

- 553 18. Okada, I., Sasaki, T. & Nakai, Y. A solution for private assessment in indi-  
554 rect reciprocity using solitary observation. *Journal of Theoretical Biology*  
555 **455**, 7–15 (2018).
- 556 19. Uchida, S. Effect of private information on indirect reciprocity. *Physical*  
557 *review E* **82**, 1–8 (2010).
- 558 20. Radzvilavicius, A. L., Stewart, A. J. & Plotkin, J. B. Evolution of empa-  
559 thetic moral evaluation. *Elife* **8**, 1–17 (2019).
- 560 21. Radzvilavicius, A. L., Kessinger, T. A. & Plotkin, J. B. Adherence to  
561 public institutions that foster cooperation. *Nature Communications* **12**.  
562 <http://dx.doi.org/10.1038/s41467-021-23783-9> (2021).
- 563 22. Okada, I. Two ways to overcome the three social dilemmas of indirect  
564 reciprocity. *Scientific Reports*, 1–9 (2020).
- 565 23. Okada, I. A review of theoretical studies on indirect reciprocity. *Games*,  
566 1–17 (2020).
- 567 24. Dunbar, R. *Grooming, gossip, and the evolution of language* (Harvard Uni-  
568 versity Press, 1998).
- 569 25. Santos, F. C., Santos, F. P. & Pacheco, J. M. Social norm complexity  
570 and past reputations in the evolution of cooperation. *Nature* **555**, 242–245  
571 (2018).
- 572 26. Otto, S. P. & Day, T. *A biologist’s guide to mathematical modeling in*  
573 *ecology and evolution* (Princeton University Press, Princeton, NJ, 2007).
- 574 27. Hilbe, C., Schmid, L., Tkadlec, J., Chatterjee, K. & Nowak, M. A. Indirect  
575 reciprocity with private, noisy, and incomplete information. *Proceedings of*  
576 *the National Academy of Sciences of the United States of America* **115**,  
577 12241–12246 (2018).
- 578 28. Sugden, R. *The economics of rights, co-operation and welfare* 1–243 (Pal-  
579 grave Macmillan, 1986).
- 580 29. Hilbe, C., Šimsa, Š., Chatterjee, K. & Nowak, M. A. Evolution of cooper-  
581 ation in stochastic games. *Nature* **559**, 246–249 (2018).
- 582 30. Milinski, M., Semmann, D., Bakker, T. C. & Krambeck, H. J. Cooperation  
583 through indirect reciprocity: Image scoring or standing strategy? *Proceed-*  
584 *ings of the Royal Society B: Biological Sciences* **268**, 2495–2501 (2001).

- 585 31. Han, T. A. *Intention recognition, commitment and their roles in the evolu-*  
586 *tion of cooperation: From artificial intelligence techniques to evolutionary*  
587 *game theory models* (Springer Publishing Company, Incorporated, 2013).
- 588 32. Barrett, H. C. *et al.* Small-scale societies exhibit fundamental variation  
589 in the role of intentions in moral judgment. *Proceedings of the National*  
590 *Academy of Sciences of the United States of America* **113**, 4688–4693  
591 (2016).
- 592 33. Bolton, G. E., Katok, E. & Ockenfels, A. Cooperation among strangers  
593 with limited information about reputation. *Journal of Public Economics*  
594 **89**, 1457–1468 (2005).
- 595 34. Swakman, V., Molleman, L., Ule, A. & Egas, M. Reputation-based cooper-  
596 ation: Empirical evidence for behavioral strategies. *Evolution and Human*  
597 *Behavior* **37**, 230–235 (2016).
- 598 35. Greif, A. *Institutions and the Path to the Modern Economy: Lessons from*  
599 *Medieval Trade* (Cambridge University Press, 2006).
- 600 36. Cavalli-Sforza, L. L., Feldman, M. W., Chen, K. H. & Dornbusch, S. M.  
601 Theory and observation in cultural transmission. *Science* **218**, 19–27 (1982).
- 602 37. Denton, K. K., Ram, Y., Liberman, U. & Feldman, M. W. Cultural evolu-  
603 tion of conformity and anticonformity. *Proceedings of the National Academy*  
604 *of Sciences of the United States of America* **117**, 13603–13614 (2020).
- 605 38. Bortolussi, L., Hillston, J., Latella, D. & Massink, M. Continuous approxi-  
606 mation of collective system behaviour: A tutorial. *Performance Evaluation*  
607 **70**, 317–349 (2013).