

**Serveur Académique Lausannois SERVAL [serval.unil.ch](http://serval.unil.ch)**

## **Author Manuscript**

**Faculty of Biology and Medicine Publication**

**This paper has been peer-reviewed but does not include the final publisher proof-corrections or journal pagination.**

Published in final edited form as:

**Title:** Human cooperation based on punishment reputation.

**Authors:** dos Santos M., Rankin D.J., Wedekind C.

**Journal:** Evolution

**Year:** 2013

**Volume:** 67(8)

**Pages:** 2446-2450

**DOI:** [10.1111/evo.12108](https://doi.org/10.1111/evo.12108)

In the absence of a copyright statement, users should assume that standard copyright protection applies, unless the article contains an explicit statement to the contrary. In case of doubt, contact the journal publisher to verify the copyright status of an article.

# 1 Human cooperation based on punishment reputation

2

3 Miguel dos Santos<sup>1</sup>, Daniel J. Rankin<sup>2,3</sup> & Claus Wedekind<sup>1\*</sup>

4

5 <sup>1</sup> Department of Ecology and Evolution, Biophore, University of Lausanne, 1015 Lausanne,  
6 Switzerland.

7 <sup>2</sup> Department of Biochemistry, University of Zurich, Winterthurerstrasse 190, 8057 Zurich,  
8 Switzerland.

9 <sup>3</sup> Swiss Institute of Bioinformatics, Quartier Sorge Bâtiment Génopode, 1015 Lausanne,  
10 Switzerland.

11

12 \* Tel: +41 21 6924250; Fax: +41 21 6924265; E-mail: claus.wedekind@unil.ch

13

14 Running title: Punishment and reputation

15

20

## 21 **ABSTRACT**

22 The threat of punishment usually promotes cooperation. However, punishing itself is costly,  
23 rare in non-human animals, and humans who punish often finish with low payoffs in  
24 economic experiments. The evolution of punishment has therefore been unclear. Recent  
25 theoretical developments suggest that punishment has evolved in the context of reputation  
26 games. We tested this idea in a simple helping game with observers and with punishment and  
27 punishment reputation (experimentally controlling for other possible reputational effects). We  
28 show that punishers fully compensate their costs as they receive help more often. The more  
29 likely defection is punished within a group, the higher the level of within-group cooperation.  
30 These beneficial effects perish if the punishment reputation is removed. We conclude that  
31 reputation is key to the evolution of punishment.

32

33 **Keywords:** experimental game theory, punishment, indirect reciprocity

34 Data to be archived in Dryad

35

## 36 **1. INTRODUCTION**

37 Punishment of non-cooperators often promotes cooperation in humans (Yamagishi 1986; Fehr  
38 and Gächter 2000, 2002; Rockenbach and Milinski 2006; Sigmund 2007) and other animals  
39 (Bshary and Grutter 2005; Raihani et al. 2012). When punished, non-cooperators are more  
40 likely to cooperate not only with the punisher but also with other social partners. Although  
41 usually conferring benefits on the punisher's social group on the long run (Yamagishi 1986;  
42 Fehr and Gächter 2000; Gächter et al. 2008; Wu et al. 2009), punishing defectors is costly to  
43 the punisher. Thus, punishment has often been perceived as truly altruistic (Fehr and Gächter  
44 2002; Barclay 2006) and its evolutionary significance as puzzling (Dreber et al. 2008; Rankin  
45 et al. 2009; Tao et al. 2009; Wu et al. 2009; dos Santos and Wedekind 2012).

46 Reputation within a social group, for example by using an individual's "image  
47 scoring" or "standing", has been proposed as a potentially important mechanism to explain  
48 cooperation in humans (Nowak and Sigmund 1998; Wedekind and Milinski 2000; Milinski et

49 al. 2002; Wedekind and Braithwaite 2002; Nowak and Sigmund 2005). Several findings  
 50 suggest that reputation could also be important for the evolution of punishment. For example,  
 51 humans playing an ultimatum game reject lower offers when they know that others will learn  
 52 about their acceptance threshold (Fehr and Fischbacher 2003). They also seem to be more  
 53 likely to punish as the number of observers increases (Kurzban et al. 2007), and punishers in  
 54 public goods games are perceived as more "group focused" and receive more monetary  
 55 benefits in a consecutive trust game (Barclay 2006). However, previous experimental studies  
 56 on human punishment have either not allowed for reputational effects, or the kind of  
 57 reputation that could built up was not clearly defined because reputational effects of punitive  
 58 actions were potentially confounded with reputational effects of cooperative actions (Fehr and  
 59 Gächter 2000; Rockenbach and Milinski 2006; Dreber et al. 2008; Gächter et al. 2008; Rand  
 60 et al. 2009; Ule et al. 2009; Wu et al. 2009).

61 Theory predicts that individuals benefit from taking the likelihood of being punished  
 62 into account (Brandt et al. 2003; Gardner and West 2004; Hilbe and Sigmund 2010).  
 63 Punishing defection can then, on the long run, be advantageous to the punisher and hence  
 64 evolve if it builds up a punishment reputation (Hilbe and Sigmund 2010; dos Santos et al.  
 65 2011; Hilbe and Traulsen 2012). This leads to two key predictions: (i) humans use reputation  
 66 to discriminate between non-punishers and punishers and are more cooperative to the latter,  
 67 and (ii) the immediate costs of punishment are compensated over time by the additional  
 68 cooperation the punisher receives from punished and observers. We tested these predictions in  
 69 a helping game with observers and with the option of punishment. Experimentally controlling  
 70 for potential reputational effects of cooperation and defection allowed us to specifically test  
 71 the significance of a potential punishment reputation.

72

## 73 2. METHODS

74 A total of 163 students played in groups of seven to nine (after written informed consent was  
 75 obtained). Players in isolated cubicles could push buttons inside a box that was connected to a  
 76 switchboard by a tangle of cables (Wedekind and Braithwaite 2002). The experimenter read  
 77 the game instructions (supplementary material) and distributed player IDs in a procedure that  
 78 ensured full anonymity (Wedekind and Braithwaite 2002). Each player received an initial  
 79 amount of 20 CHF that could be used in a game. Final gains were paid out as in Wedekind  
 80 and Braithwaite (2002).

81 One player was put in the Donor role, another in the Receiver role. The Donor  
 82 (indicated to the player via a small bulbs inside the box) could either refuse to donate or  
 83 donate to the Receiver by paying 1 CHF for the other to receive 2 CHF (we donated the  
 84 difference). The Receiver could then decide whether or not to "make the Donor lose money"  
 85 by accepting a cost of 1 CHF for the Donor to lose 2 CHF. Then a next pair of players was  
 86 chosen. Each player played once per round as Donor and twice per two rounds as Receiver  
 87 (e.g. once in rounds 1 and 2 each, or twice in round 2 only). The total number of rounds (16)  
 88 was not communicated to the players.

89 During the first 8 rounds, the Donor's decision (but not his/her ID in order to control  
 90 for potential confounding reputational effects (Wedekind and Milinski 2000; Wedekind and  
 91 Braithwaite 2002)) was displayed on a projector screen to all players. We also displayed the  
 92 Receivers' ID and his/her punishment score, i.e. an arrow wandering on a scale from -5 or +5,  
 93 starting at 0, changing +1 for every punished and -1 for every non-punished defection. From  
 94 round 9 on, we either removed the punishment option but continued to display the last  
 95 punishment reputation the players had earned ("NOPUN/REP", n=36, 4 groups), left the  
 96 option to punish but removed the display of reputation ("PUN/NOREP", n=36, 4 groups), or  
 97 changed nothing ("PUN/REP", n=91, 11 groups; the larger number of groups under these  
 98 control conditions allowed us to further analyze the within-group correlations between net

99 gain and punishment reputation after 16 rounds of undisturbed interactions). The  
 100 “NOPUN/REP” treatment allowed us to test whether donors would either reward punishers  
 101 for incurring the cost of disciplining non-cooperators or just stop discriminating between  
 102 punishers and non-punishers in the absence of the threat of punishment. The “PUN/NOREP”  
 103 treatment allowed us to assess the effect of punishment alone on the cooperation frequency.  
 104 Players in the later two treatments were only informed at the end of round 8 about the change  
 105 of rules.

106 The statistical analyses were carried out with R 2.10.1 (R Development Core Team  
 107 2010). We used the lme4 package for linear and logistic mixed-effect models (Bates and  
 108 Sarkar 2007) that it is suitable for unbalanced designs (Baayen et al. 2008). Linear mixed-  
 109 effect models were used to analyze group cooperation and punishment frequency during  
 110 rounds 1 to 8 and 9 to 16, with group as random effect. Generalized linear mixed-effect  
 111 models were used to analyze Donors’ probability of giving, with group and donor as random  
 112 effects.

113

### 114 3. RESULTS

115 By the end of the first 8 rounds, high probabilities of punishing defection lead to high levels  
 116 of cooperation (Fig. 1a), and more cooperative groups reached higher total payoffs than less  
 117 cooperative groups ( $r = 0.74$ ,  $n = 19$ ,  $P < 0.001$ ). Within the controls (“PUN/REP”) the effects  
 118 of reputation did not seem to change over the full 16 rounds (likelihood ratio test (LRT),  
 119 interaction between punishment score and round:  $\chi^2 = 0.12$ ,  $P = 0.73$ ; effect of round:  $\chi^2 =$   
 120  $2.24$ ,  $P = 0.13$ ). The within-group correlation between net gains (i.e. sum of received  
 121 donations - punishment costs) and final punishment reputation ranged from  $r = -0.80$  to  $0.78$   
 122 in the controls. These within-group correlations could be significantly predicted by a  
 123 discriminant score that was the mean difference between the Receivers’ punishment score  
 124 when the Donor gave and did not give (Fig. 1b; the early interactions during a game seemed  
 125 important here: the more individuals punished during the first 4 rounds, the higher their  
 126 discrimination score at the end of the game, linear mixed-effect model,  $t = 2.23$ ,  $P = 0.03$ ).

127 During the experimental stage, i.e. from round 9 to 16, the frequency of cooperation  
 128 did not significantly change in the controls (Fig. 2a; LRT:  $\chi^2 = 2.55$ ,  $P = 0.11$ ) but declined in  
 129 both the PUN/NOREP and NOPUN/REP treatments (Fig. 2a; Table 1). In parallel, the  
 130 probability of punishing defection declined when punishment could no longer affect  
 131 reputation (Fig. 2b; Table 1). In the controls, Donors were more likely to give to Receivers  
 132 with high punishment score than to those with a low punishment score (Fig. 2c, Table 1). In  
 133 the PUN/NOREP treatment where no further punishment reputation could be built up, the  
 134 updated punishment score that would correspond to the Receivers’ actions but was no more  
 135 displayed did also not seem to affect the Donors’ decisions (LRT:  $\chi^2 = 0.25$ ,  $P = 0.61$ ; Fig.  
 136 2c). Correspondingly, the within-group correlation between the players’ account and their  
 137 probability of punishing defection was significantly higher in the controls than in the  
 138 PUN/NOREP treatment (Welch t-test:  $t_{12.95} = 3.14$ ,  $P = 0.007$ ). In the NOPUN/REP treatment,  
 139 where punishment was no longer an option, the reputation that had been built up until round 8  
 140 did not seem to affect the Donors’ decisions either (LRT:  $\chi^2 = 0.09$ ,  $P = 0.76$ ; Fig. 2c ).

141

### 142 4. DISCUSSION

143 Recent theory predicts that punishment has either evolved in another context than cooperation  
 144 (Dreber et al. 2008; Wu et al. 2009) or that reputational effects compensate for the costs of  
 145 punishment (Hilbe and Sigmund 2010; dos Santos et al. 2011). We found the latter to be true.  
 146 High levels of cooperation were maintained when punishment could build up a punishment  
 147 reputation, and, on average, the increased cooperation fully compensated for the costs of  
 148 punishment. In groups with a high degree of discrimination between punishers and non-

149 punishers, the additional cooperation that punishers received even lead to net benefits, i.e. the  
150 costs of punishment were then overcompensated.

151 Punishers in public goods games may often be perceived as trustworthy and group  
152 focused, and may enjoy similar reputational benefits than generous people do in simple  
153 helping games (Barclay 2006). Human punishment has even been called “altruistic”  
154 (Yamagishi 1986; Fehr and Gächter 2000; Fehr and Gächter 2002) because people may  
155 punish defectors even in anonymous one-shot interactions where no benefit could be gained.  
156 However, anonymous one-shot interactions were probably very rare in human history, i.e.  
157 punishment is unlikely to have been evolved in such interactions (dos Santos and Wedekind  
158 2012). In our game where players could built up a reputation in repeated interactions,  
159 punishers seemed feared rather than rewarded for altruistic behavior: Donors stopped  
160 discriminating between punishers and non-punishers when the opportunity to punish had been  
161 removed, i.e. the punishment reputation that had been built up before had no more effect  
162 when the threat of punishment was removed. As a consequence, those who had invested into  
163 their punishment reputation could not get compensated during the second part of our  
164 experiment and finally finished with relatively low payoffs.

165 When humans can observe the others’ actions and can chose with whom to interact,  
166 they sometimes seem to weight cooperation higher than punishment (Rockenbach and  
167 Milinski 2011). Punishment turned out to be important in our experiments, but participants  
168 could not choose their partners and it is possible that the likelihood of punishing defection  
169 depends on whether partner choice is allowed for. Humans also tend to punish more often as  
170 the number of observers increases (Kurzban et al. 2007). This suggests that punishment is a  
171 strategic decision that takes aspects of the social environment into account. We observed that  
172 the outcome of the first interactions within a newly built social group influenced later  
173 dynamics: the participants’ willingness to give to punishers depends on what they experienced  
174 at the beginning of a social interaction.

175

## 176 **ACKNOWLEDGEMENTS**

177 We thank the students for participation, R. Bergmüller, P. Bize, R. Bshary, M. Chapuisat, C.  
178 Clavien, C. El Mouden, A. Gardner, C. Hauert, M. Hochberg, L. Keller, C. Metzger, S.  
179 Nusslé, A. Ross-Gillespie, and S. West for comments or discussion, and the Swiss NSF for  
180 funding.

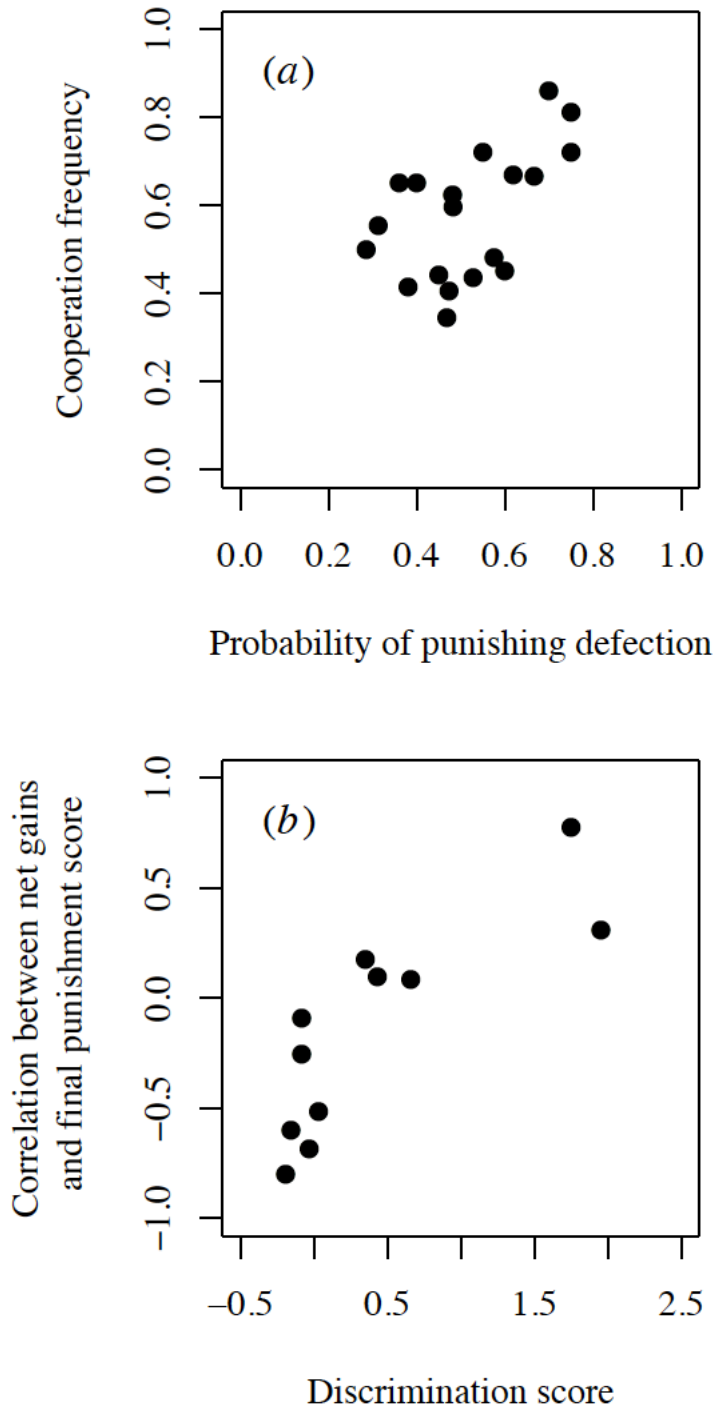
181

## 182 **REFERENCES**

- 183 Baayen, R. H., D. J. Davidson, and D. M. Bates. 2008. Mixed-effects modeling with crossed  
184 random effects for subjects and items. *J. Mem. Lang.* 59:390-412.
- 185 Barclay, P. 2006. Reputational benefits for altruistic punishment. *Evol. Hum. Behav.* 27:325-  
186 344.
- 187 Bates, D. M., and D. Sarkar. 2007. lme4: Linear mixed-effects models using S4 classes, R  
188 package version 0.99875-9.
- 189 Brandt, H., C. Hauert, and K. Sigmund. 2003. Punishment and reputation in spatial public  
190 goods games. *Proc. R. Soc. B* 270:1099-1104.
- 191 Bshary, R., and A. S. Grutter. 2005. Punishment and partner switching cause cooperative  
192 behaviour in a cleaning mutualism. *Biol. Lett.* 1:396-399.
- 193 dos Santos, M., D. J. Rankin, and C. Wedekind. 2011. The evolution of punishment through  
194 reputation. *Proc. R. Soc. B* 278:371-377.
- 195 dos Santos, M., and C. Wedekind. 2012. Examining punishment at different explanatory  
196 levels. *Behav. Brain Sci.* 35.
- 197 Dreber, A., D. G. Rand, D. Fudenberg, and M. A. Nowak. 2008. Winners don't punish. *Nature*  
198 452:348-351.

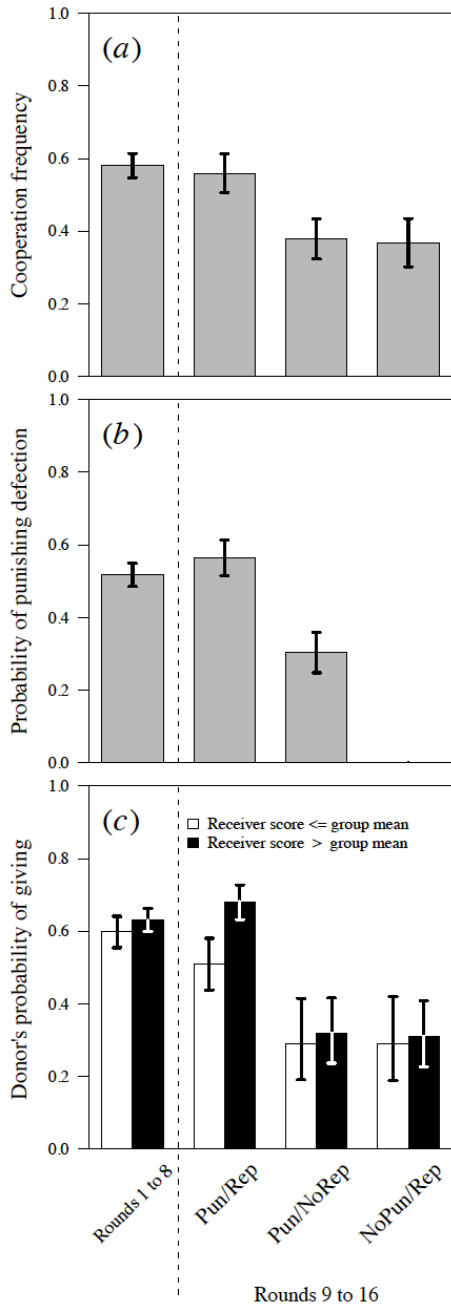
- 199 Fehr, E., and U. Fischbacher. 2003. The nature of human altruism. *Nature* 425:785-791.
- 200 Fehr, E., and S. Gächter. 2000. Cooperation and punishment in public goods experiments.
- 201 *Am. Econ. Rev.* 90:980-994.
- 202 Fehr, E., and S. Gächter. 2002. Altruistic punishment in humans. *Nature* 415:137-140.
- 203 Gächter, S., E. Renner, and M. Sefton. 2008. The long-run benefits of punishment. *Science*
- 204 322:1510-1510.
- 205 Gardner, A., and S. A. West. 2004. Cooperation and punishment, especially in humans. *Am.*
- 206 *Nat.* 164:753-764.
- 207 Hilbe, C., and K. Sigmund. 2010. Incentives and opportunism: from the carrot to the stick.
- 208 *Proc. R. Soc. B* 277:2427-2433.
- 209 Hilbe, C., and A. Traulsen. 2012. Emergence of responsible sanctions without second order
- 210 free riders, antisocial punishment or spite. *Sci Rep* 2.
- 211 Kurzban, R., P. DeScioli, and E. O'Brien. 2007. Audience effects on moralistic punishment.
- 212 *Evol. Hum. Behav.* 28:75-84.
- 213 Milinski, M., D. Semmann, and H. J. Krambeck. 2002. Reputation helps solve the 'tragedy of
- 214 the commons'. *Nature* 415:424-426.
- 215 Nowak, M. A., and K. Sigmund. 1998. Evolution of indirect reciprocity by image scoring.
- 216 *Nature* 393:573-577.
- 217 Nowak, M. A., and K. Sigmund. 2005. Evolution of indirect reciprocity. *Nature* 437:1291-
- 218 1298.
- 219 R Development Core Team. 2010. R: A Language and Environment for Statistical
- 220 Computing. R Foundation for Statistical Computing, Vienna, Austria.
- 221 Raihani, N. J., A. Thornton, and R. Bshary. 2012. Punishment and cooperation in nature.
- 222 *Trends. Ecol. Evol.* 27:288-295.
- 223 Rand, D. G., A. Dreber, T. Ellingsen, D. Fudenberg, and M. A. Nowak. 2009. Positive
- 224 interactions promote public cooperation. *Science* 325:1272-1275.
- 225 Rankin, D. J., M. dos Santos, and C. Wedekind. 2009. The evolutionary significance of costly
- 226 punishment is still to be demonstrated. *Proc. Natl Acad. Sci. USA* 106:E135-E135.
- 227 Rockenbach, B., and M. Milinski. 2006. The efficient interaction of indirect reciprocity and
- 228 costly punishment. *Nature* 444:718-723.
- 229 Rockenbach, B., and M. Milinski. 2011. To qualify as a social partner, humans hide severe
- 230 punishment, although their observed cooperativeness is decisive. *Proc. Natl Acad. Sci.*
- 231 *USA* 108:18307-18312.
- 232 Sigmund, K. 2007. Punish or perish? Retaliation and collaboration among humans. *Trends*
- 233 *Ecol. Evol.* 22:593-600.
- 234 Tao, Y., C. Li, J. J. Wu, and R. Cressman. 2009. Reply to Rankin et al.: The efficiency ratio
- 235 of costly punishment. *Proc. Natl Acad. Sci. USA* 106:E136-E136.
- 236 Ule, A., A. Schram, A. Riedl, and T. N. Cason. 2009. Indirect punishment and generosity
- 237 toward strangers. *Science* 326:1701-1704.
- 238 Wedekind, C., and V. A. Braithwaite. 2002. The long-term benefits of human generosity in
- 239 indirect reciprocity. *Curr. Biol.* 12:1012-1015.
- 240 Wedekind, C., and M. Milinski. 2000. Cooperation through image scoring in humans. *Science*
- 241 288:850-852.
- 242 Wu, J. J., B. Y. Zhang, Z. X. Zhou, Q. Q. He, X. D. Zheng, R. Cressman, and Y. Tao. 2009.
- 243 Costly punishment does not always increase cooperation. *Proc. Natl Acad. Sci. USA*
- 244 106:17448-17451.
- 245 Yamagishi, T. 1986. The provision of a sanctioning system as a public good. *J. Pers. Soc.*
- 246 *Psychol.* 51:110-116.
- 247

248 Figure 1. Cooperation and net gains when punishment reputation was unconstrained. (a)  
249 Frequency of cooperation and of punishing defection per group at the end of the first 8 rounds  
250 (Pearson's  $r = 0.54$ ,  $n = 19$  groups,  $P = 0.015$ ). (b) Relationship between the within-group  
251 discriminant score (= the mean difference between the Receivers' punishment score when the  
252 Donor gave and did not give) and the correlation between net gains (i.e. sum of received  
253 donations - punishment costs) and punishment score in the PUN/REP treatment (Spearman's  
254 rank order correlation coefficient =  $0.84$ ,  $n = 11$ ,  $P = 0.002$ ).  
255



256

257 Figure 2. Treatment effects on (a) cooperation frequency, (b) the probability of punishing  
 258 defection, and (c) the Donors' probability of giving to Receivers' with high or low  
 259 punishment scores (corrected for group and round effects). PUN/REP = punishment was  
 260 always possible and punishment reputation was continuously updated and displayed;  
 261 PUN/NOREP = punishment reputation was no longer displayed in Part 2; NOPUN/REP =  
 262 punishment was no longer possible in Part 2, but the reputation that had been built up during  
 263 Part 1 was displayed. See Table 1 for statistics.  
 264



265  
 266  
 267



268 Table 1. The effects of experimental treatment and the punishment score that built up during  
 269 the second part of the experiment (i.e. rounds 9-16) on (a) the frequency of cooperation, (b)  
 270 the probability of punishing defection, and (c) the Donors' probability of giving. Linear (a and  
 271 b) and generalized mixed-effects models (c) were fitted with and without a given effect (or  
 272 interaction) in order to test if the goodness of fit between both models differed in a likelihood  
 273 ratio test.  
 274

	$\chi^2$	<u>d.f.</u>	P
<i>a) Cooperation frequency</i>			
Treatment	6.1	2	0.048
<i>b) Probability of punishing defection</i>			
Treatment	4.9	1	0.027
<i>c) Donor's probability of giving</i>			
Punishment score	8.2	1	0.004
Treatment	6.8	2	0.033
Punishment score x Treatment	2.9	2	0.235

275