

Serveur Académique Lausannois SERVAL 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: Examining punishment at different explanatory levels.

Authors: dos Santos M., Wedekind C.

Journal: Behavioral and Brain Sciences,

Year: 2012

Volume: 35(1)

Pages: 23-24

DOI: [10.1017/S0140525X1100121X](https://doi.org/10.1017/S0140525X1100121X)

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 Examining punishment at different explanatory levels

2
3 Miguel dos Santos & Claus Wedekind

4
5 Department of Ecology and Evolution, Biophore, University of Lausanne, 1015 Lausanne,
6 Switzerland.

10 12 13 ABSTRACT

14 Experimental studies on punishment have sometimes been overinterpreted not only for
15 the reasons Guala lists but also because of a frequent conflation of proximate and
16 ultimate explanatory levels that Guala's review perpetuates. Moreover, for future
17 analyses we may need a clearer classification of different kinds of punishment.

18
19 When explaining behavioral decisions, it is important to distinguish between different
20 explanatory levels, especially between proximate (mechanistic) and ultimate
21 (evolutionary) explanations (Tinbergen, 1963). Proximate explanations of a given
22 behavior deal with questions about its ontogeny (e.g. how does the behavior change
23 with age and experience) or about its causation, i.e. the physiological, molecular, and
24 cognitive mechanisms underlying the behavior and the stimuli that elicit it. Ultimate
25 explanations either deal with questions about the phylogeny of the behavior (e.g. how
26 does it compare with similar behaviors in related species) or its adaptive value (e.g.
27 what is its impact on the individual's survival and life-time reproductive success).

28 The concept of weak reciprocity, as defined in Guala (2011), is an attempt to
29 explain the adaptive value of cooperation and punishment because it concentrates on
30 the fitness benefits one could get from cooperating, defecting, or punishing (Trivers,
31 1971; Alexander, 1974). This concept is restricted to one explanatory level only. In
32 contrast, strong reciprocity mixes different explanatory levels: it uses proximate
33 arguments to explain ultimate problems (Fehr & Gächter, 2002; Fehr & Fischbacher,
34 2003, 2004; Fehr & Rockenbach, 2003; Gintis et al., 2003; Bowles & Gintis, 2004). Strong
35 reciprocity is, for example, called a "... predisposition to reward others for cooperative,
36 norm-abiding behaviours" and "... a propensity to impose sanctions on others for norm
37 violations" (Fehr & Fischbacher, 2003, p. 785). Such a definition clearly relates to the
38 causal mechanisms of cooperation and punishment. But the concept is then frequently
39 used as to answer ultimate (evolutionary) questions, for example in Bowles & Gintis
40 (2004, p.17): "... cooperation is maintained because many humans have a predisposition
41 to punish those who violate group-beneficial norms". Such a mixing up of different
42 explanatory levels can, from an evolutionary point of view, easily lead to
43 overinterpretations of proximate patterns (Hagen & Hammerstein, 2006; Sigmund,
44 2007; West et al., 2007, in press; Rankin et al., 2009). For example, punishment that can
45 be observed in anonymous one-shot interactions seems truly altruistic and was
46 interpreted as such in Fehr & Gächter (2002). However, until very recently, humans
47 lived in groups where anonymous one-shot interactions were probably very rare, i.e.
48 such interactions are most probably not the context in which human punishment has

49 evolved. If studied within a more natural social context, human punishment may
50 ultimately be self-interested.

51 As discussed in Guala (2011), explaining punishment from an evolutionary point
52 of view requires determining the costs and benefits of punishment. In line with weak
53 reciprocity models, recent studies have shown that punishment can lead to long-term
54 net benefits and hence be evolutionarily stable when punitive actions contribute to a
55 punishment reputation (Hilbe & Sigmund, 2010; dos Santos et al., 2011). Under such
56 conditions, the immediate costs of punishment can be outweighed by the benefits a
57 punisher receives later because of his/her punishment reputation. Experimental studies
58 that ignore the possible effects of a punishment reputation can therefore easily produce
59 artifacts (Hagen & Hammerstein, 2006).

60 We also believe that the term “punishment” is currently used too broadly in the
61 literature on cooperation. If “punishment” is the subtraction of resources from free-
62 riders in order to reduce the frequency of further free-riding, there are at least three
63 different kinds of punishment that may need to be distinguished both for ultimate and
64 proximate analyses. Many of these analyses deal with what could be called “simple
65 costly punishment”, i.e. punishers pay a cost to induce a cost on the punished (Fehr &
66 Gächter, 2000; Rockenbach & Milinski, 2006; Dreber et al., 2008; Rand et al., 2009; Wu et
67 al., 2009). Another form of punishment could be called “punishment by taking
68 something away” (e.g. Cephu’s example in Guala, 2011). Here, the punisher takes
69 something from the punished in order to induce a cost to the punished. Regardless of
70 whether the punisher thereby experiences an immediate reduction of the own welfare
71 or not, “punishment by taking something away” and the upper “simple costly
72 punishment” are likely to differ in their cost-benefit ratios (relevant for ultimate
73 analyses) and may involve, for example, different kinds of emotions (relevant for
74 proximate analyses). A third category could be called “punishment by refusal”. The
75 punisher then punishes by refusing to cooperate with the punished in a repeated game
76 like, for example, an iterated Prisoner’s Dilemma (Fudenberg et al., 1994). The examples
77 of ostracism discussed in Guala relate to this kind of punishment. Such defection may
78 typically be a reaction to non-provoked defection and could be called “punishment” if it
79 reduces the income of the punished (i.e. his/her benefits from what would otherwise be
80 cooperative interactions) in order to possibly improve the punisher’s long-term benefits
81 from future cooperative interactions with a refined punished or with others. This third
82 kind of punishment could be immediately costly for the punisher, for example, if it
83 delays the resumption of beneficial mutual cooperation. Such immediate costs would
84 have to be compensated on the long run in order to maintain “punishment by refusal” as
85 an evolutionary successful behavioral strategy. However, a possible alternative function
86 of defection in response to defection may be to simply avoid the losses of anticipated
87 further defection (e.g. avoiding the sucker’s payoff in the Prisoner’s Dilemma). It is
88 probably not useful to call this later form of defection “punishment” if it usually does not
89 ultimately increase the level of cooperation within a group or directly with the defector
90 (from an ultimate point of view), or if it is just a precautionary measure to avoid further
91 losses (from a proximate point of view). Therefore, purely punitive actions may not
92 always be easy to identify. Multidisciplinary approaches that carefully exploit the
93 specific advantages of proximate and ultimate analyses are therefore often necessary to
94 better understand human behavior.

95
96 REFERENCES

97 Alexander, R. D. (1974). The evolution of social behaviour. *Annual Review of Ecological*
98 *Systematics*, 5, 325-383.

99 Bowles, S. & Gintis, H. (2004). The evolution of strong reciprocity: cooperation in
100 heterogeneous populations. *Theoretical Population Biology*, 65, 17-28.

101 dos Santos, M., Rankin, D. J. & Wedekind, C. (2011). The evolution of punishment through
102 reputation. *Proceedings of the Royal Society B-Biological Sciences*, 278, 371-377.

103 Dreber, A., Rand, D. G., Fudenberg, D. & Nowak, M. A. (2008). Winners don't punish.
104 *Nature*, 452, 348-351.

105 Fehr, E. & Fischbacher, U. (2003). The nature of human altruism. *Nature*, 425, 785-791.

106 Fehr, E. & Fischbacher, U. (2004). Third-party punishment and social norms. *Evolution*
107 *and Human Behavior*, 25, 63-87.

108 Fehr, E. & Gächter, S. (2000). Cooperation and punishment in public goods experiments.
109 *American Economic Review*, 90, 980-994.

110 Fehr, E. & Gächter, S. (2002). Altruistic punishment in humans. *Nature*, 415, 137-140.

111 Fehr, E. & Rockenbach, B. (2003). Detrimental effects of sanctions on human altruism.
112 *Nature*, 422, 137-140.

113 Fudenberg, D., Levine, D. & Maskin, E. (1994). The Folk Theorem with Imperfect Public
114 Information. *Econometrica*, 62, 997-1039.

115 Gintis, H., Bowles, S., Boyd, R. & Fehr, E. (2003). Explaining altruistic behavior in humans.
116 *Evolution and Human Behavior*, 24, 153-172.

117 Hagen, E. H. & Hammerstein, P. (2006). Game theory and human evolution: A critique of
118 some recent interpretations of experimental games. *Theoretical Population Biology*,
119 69, 339-348.

120 Hilbe, C. & Sigmund, K. (2010). Incentives and opportunism: from the carrot to the stick.
121 *Proceedings of the Royal Society B-Biological Sciences*, 277, 2427-2433.

122 Rand, D. G., Dreber, A., Ellingsen, T., Fudenberg, D. & Nowak, M. A. (2009). Positive
123 interactions promote public cooperation. *Science*, 325, 1272-1275.

124 Rankin, D. J., dos Santos, M. & Wedekind, C. (2009). The evolutionary significance of
125 costly punishment is still to be demonstrated. *Proceedings of the National Academy*
126 *of Sciences USA*, 106, E135-E135.

127 Rockenbach, B. & Milinski, M. (2006). The efficient interaction of indirect reciprocity and
128 costly punishment. *Nature*, 444, 718-723.

129 Sigmund, K. (2007). Punish or perish? Retaliation and collaboration among humans.
130 *Trends in Ecology & Evolution*, 22, 593-600.

131 Tinbergen, N. (1963). On aims and methods of ethology. *Zeitschrift fur Tierpsychologie*,
132 20, 410-433.

133 Trivers, R. L. (1971). Evolution of reciprocal altruism. *Quarterly Review of Biology*, 46,
134 35-57.

135 West, S. A., El Mouden, C. & Gardner, A. (In press). 16 common misconceptions about the
136 evolution of cooperation in humans. *Evolution and Human Behavior*,

137 West, S. A., Griffin, A. S. & Gardner, A. (2007). Evolutionary explanations for cooperation.
138 *Current Biology*, 17, R661-R672.

139 Wu, J. J., Zhang, B. Y., Zhou, Z. X., He, Q. Q., Zheng, X. D., Cressman, R. & Tao, Y. (2009).
140 Costly punishment does not always increase cooperation. *Proceedings of the*
141 *National Academy of Sciences USA*, 106, 17448-17451.

142

143