# Coevolution of Third-Party Punishment and Punishment Reputation Dependency

**Taisho Ozaki\*, Yasuo Ihara**

University of Tokyo, Hongo 7-3-1, Bunkyoku, Tokyo 113-0033, Japan

*Author for correspondence (pqLcy-cool@g.ecc.u-tokyo.ac.jp)

An obstacle in explaining the evolution of cooperation through third-party punishment is the second-order free rider problem, that is, if punishing others is costly, it is adaptive for any individual not to punish. One hypothesis to resolve the issue states that third-party punishment is advantageous for the punisher, because it is evaluated as socially good and thereby enhances future cooperation from the observers. However, it is thus far unclear whether cooperation is truly promoted by third-party punishment in the presence of a social norm favoring punishing behavior, and under what circumstances such social norm is established. In this paper, we present two mathematical models to explore the reputation-enhancing effect of third-party punishment and its role in the evolution of cooperation. Our results suggest that third-party punishment can indeed facilitate the evolution of cooperation when individuals base their evaluation of others on their punishing behavior, and such punishment reputation dependency can coevolve with third-party punishment during the process where large-scale cooperation is established within a population.

## Keywords

## Introduction

One type of social behavior that characterizes humans is third-party punishment, or punishment for transgressions imposed by someone other than the victim of those transgressions (Fehr & Fischbacher, 2004; Riedl et al., 2012). Third-party punishment seems to be universal across cultures and the tendency to punish wrongdoing may emerge in early infancy (Kanakogi et al., 2022), suggesting a possible innate basis. However, the evolution of third-party punishment is a puzzle, because, so long as punishment on others carries a cost, it should be advantageous for anyone to avoid it (i.e., the second-order free rider problem). As a possible solution, it has been pointed out that third-party punishment may bring about a reputational benefit by signaling the punisher's

willingness to retaliate and/or cooperative intent (Raihani & Bshary, 2015). In particular, one hypothesis that is currently understudied is that punishers may gain a reputation as trustworthy individuals, thereby eliciting future cooperation (Barclay, 2006; Kurzban et al., 2007; Nelissen, 2008). So far, however, evolutionary origins of such punishment reputation dependency (PRD), or individuals' tendency to value those who impose third-party punishment highly, remains unclear.

Meanwhile, the role of reputation based on individuals' cooperative behavior has been intensively studied in the context of indirect reciprocity (Alexander, 1987; Sugden, 1986; Trivers, 1971). The idea is that if a good reputation leads to increased cooperation from others, individuals are better off paying the cost of cooperation in order to maintain positive evaluation. A simplest social norm to achieve indirect reciprocity is Image Scoring (Nowak & Sigmund, 1998), under which an actor's reputation rises if he cooperates with any others and falls if he does not. In this paper, we incorporate third-party punishment into the analytical framework of indirect reciprocity, Image Scoring in particular, to explore the coevolution of third-party punishment and the PRD. The goals of the study are to examine whether third-party punishment can facilitate cooperation given the presence of the PRD, and under what circumstances the PRD can be favored by natural selection.

## Methods

### Model 1

We consider an infinitely large population of individuals, who are characterized by one of three strategies: Image scoring (IS; Nowak & Sigmund, 1998), Image scoring with punishment (ISp), and Always defect (ALLD). Each individual also has either Good or Bad reputation, which is known by all individuals in the population. A generation consists of a repeated game with an indefinite number of rounds, in each of which the population breaks up randomly into groups of three individuals. Within such group, one individual is chosen at random to be the donor, another to be the recipient, and the other to be the punisher. The donor may help the recipient by paying a cost, $c$, in order to provide a benefit, $b$, to the recipient ($c > 0$, $b > 0$). As a donor, IS and ISp help the recipient if and only if the recipient is Good, except that they fail to help a Good recipient with error rate $e$ ($0 \leq e < 1$). ALLD never benefits the recipient. In case the donor refrains from helping, the punisher may engage in a costly punishment. ISp punishes a non-helping donor if and only if the recipient is Good, which imposes costs $\beta$ on the donor and $\alpha$ on the punisher ($\beta > 0$, $\alpha > 0$). Neither IS nor ALLD punishes a non-helping donor. Let $x_N$, $x_P$, and $y$ denote the frequencies of IS, ISp, and ALLD, respectively, where $x = x_N + x_P$. We also use $\psi$ to represent the frequency of Good individuals. The expected payoffs per round are shown in Table 1.

**Table 1.** Expected payoffs per round in Model 1.

| Strategy | Reputation | Payoff |
|---|---|---|
| IS | Bad | $-\frac{1}{3}[(1-e)c+e\beta x_P]\psi$ |
| IS | Good | $-\frac{1}{3}[(1-e)c+e\beta x_P]\psi+\frac{1}{3}(1-e)bx$ |
| ISp | Bad | $-\frac{1}{3}[(1-e)c+e\beta x_P]\psi-\frac{1}{3}\alpha(ex+y)\psi$ |
| ISp | Good | $-\frac{1}{3}[(1-e)c+e\beta x_P]\psi+\frac{1}{3}(1-e)bx-\frac{1}{3}\alpha(ex+y)\psi$ |
| ALLD | Bad | $-\frac{1}{3}\beta x_P\psi$ |
| ALLD | Good | $-\frac{1}{3}\beta x_P\psi+\frac{1}{3}(1-e)bx$ |

In each round, the donor's reputation is updated to be Good if it helps the recipient, and Bad if it does not. In addition, we assume that the punisher's reputation is updated to be Good if it punishes a donor who does not help a Good recipient. The recipient's reputation is not altered. Let $x_{N0}$, $x_{N1}$, $x_{P0}$, $x_{P1}$, $y_0$, and $y_1$ denote the frequencies of Bad IS, Good IS, Bad ISp, Good ISp, Bad ALLD, and Good ALLD, respectively, in a given round, where $x_N = x_{N0} + x_{N1}$, $x_P = x_{P0} + x_{P1}$, and $y = y_0 + y_1$ are constant within a generation. The frequencies of Good individuals in the next round are given by

$$x'_{N1}=\frac{2}{3}x_{N1}+\frac{1}{3}x_N(1-e)\psi, \tag{1a}$$

$$x'_{P1}=\frac{2}{3}x_{P1}+\frac{1}{3}x_P(1-e)\psi+\frac{1}{3}x_{P0}(ex+y)\psi, \tag{1b}$$

$$y'_1=\frac{2}{3}y_1. \tag{1c}$$

Given that one round of game is played, another round is played with probability $w$, and the game ends there with probability $1 - w$ ($0 < w < 1$), where the first round in each generation is always played. Using Table 1 and (1), we numerically obtain the expected total payoffs of IS, ISp, and ALLD, denoted by $D_N$, $D_P$, and $D_e$, respectively (Electronic Supplementary Material Appendix A). Assuming that individuals produce offspring, who inherit the parent's strategy, with rates proportional to the expected total payoffs, we describe the evolutionary dynamics across generations by the following replicator equations:

$$\frac{dx_N}{dt}=(D_N-\overline{D})x_N, \tag{2a}$$

$$\frac{dx_P}{dt}=(D_P-\overline{D})x_P, \tag{2b}$$

$$\frac{dy}{dt}=(D_e-\overline{D})y, \tag{2c}$$

where $\overline{D}=D_N x_N+D_P x_P+D_e y$.

### Model 2
Here we perform agent-based simulations involving a population of $n$ individuals. In each time step, a triad is chosen randomly to play the roles as donor, recipient, and punisher. As in Model 1, the donor may help the recipient by paying cost $c$ to give benefit $b$, and in case the donor refrains from helping, the punisher may perform a costly

punishment by paying cost $\alpha$ to incur cost $\beta$ on the non-helping donor. In Model 2, an individual's reputation is captured by two variables: help reputation, $R$, and punishment reputation, $Rp$ ($0 \le R, Rp \le 100$). Donors gain one unit of help reputation when they help a recipient, and lose one unit when they do not. Punishers gain $u_p$ units of punishment reputation by punishing a non-helping donor ($u_p \ge 0$), while one's punishment reputation never decreases.

Individuals also vary in terms of three genetic traits. First, they vary in whether they have the punishment reputation dependency (PRD). Individuals with the PRD are concerned with both help and punishment reputations of other individuals, while those without the PRD give care only to the help reputation. Second, they carry their own helping threshold, $s$ ($0 \le s \le 101$). Donors without the PRD consider a recipient Good and thus provide help if and only if $R \ge s$, where $R$ is the recipient's help reputation and $s$ is the donor's help threshold, except that helping fails with error rate $e$ as in Model 1. Donors with the PRD, on the other hand, consider a recipient Good and hence give help if and only if $\min(R + Rp, 100) \ge s$, also being subject to error rate $e$. In words, donors with the PRD regard whichever is smaller of $R + Rp$ and 100 as the recipient's reputation. Note that a donor with $s > 100$ never provides help. Third, they are either punishing or non-punishing. Being in the role of punisher, punishing individuals punish a donor who does not help a recipient whom they consider Good, while non-punishing individuals never punish anyone.

After $m$ games are played ($m \ge 1$), a new generation of $n$ individuals are produced, whose parents are chosen at random with probabilities proportional to their lifetime total payoffs. The offspring inherit the parents' genetic traits and have $R = Rp = 0$ at birth. Mutation switches an individual's punishing/non-punishing trait and the presence/absence of the PRD with the rates $\mu_p$ and $\mu_d$, respectively, while it modifies the helping threshold with the rate $\mu_s$ ($0 < \mu_p, \mu_d, \mu_s < 1$). When a mutation hits the helping threshold, a new value for $s$ is drawn from the integers satisfying $0 \le s \le 101$ with equal probabilities.

Below, we consider the population to be in the "cooperative state" when the donor helps the recipient in at least $0.9(1 - e)$ of the $m$ games in a given generation, and in the "non-cooperative state" when the donor helps the recipient in less than $0.1(1-e)$ of the games.

## Results

### Model 1

When $e = 0$, numerical analysis suggests that the evolutionary dynamics converge to either of two equilibrium states: an equilibrium point where all individuals are ALLD, or a neutrally stable line of equilibria, consisting of all points at which ALLD is absent (Figure 1). In the absence of ISp, selection favors IS over ALLD if and only if the frequency of IS exceeds a threshold, given by $x_N^* = (3-2w)c / (wb)$. Thus, a population of IS is stable against invasion by ALLD if
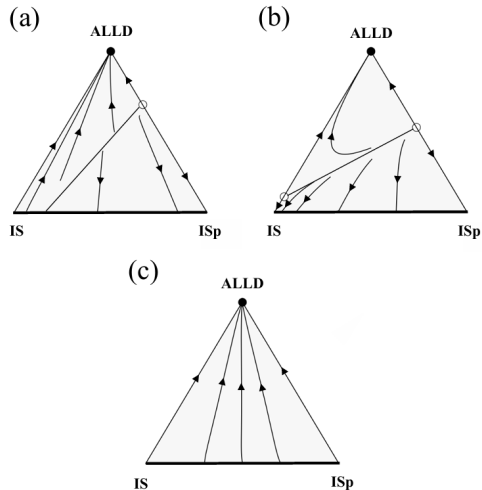
$$\frac{w}{3-2w} > \frac{c}{b}. \tag{3}$$

In the absence of IS, ISp is selectively favored over ALLD if and only if the frequency of ISp is above a threshold, given approximately by $x_P^* \approx (3-2w)(c+\alpha)/[wb+(\alpha+\beta)(3-2w)]$. Hence, a population of ISp is stable against invasion by ALLD if

$$\frac{w}{3-2w} > \frac{c-\beta}{b}. \tag{4}$$

When $e > 0$, it is suggested numerically that one of three equilibrium points is reached eventually, corresponding to the fixation of IS, ISp, or ALLD. The evolutionary dynamics are qualitatively the same as in the case of $e = 0$, except that the neutrally stable line of equilibria no longer exists (Figure A3 in ESM Appendix A).

If the threshold frequency of ISp for it to be selectively favored over ALLD is smaller than that of IS (i.e., $x_P^* < x_N^*$), one can conclude that the third-party punishment facilitates cooperation through indirect reciprocity. Whether or not the error in helping is taken into consideration, this is the case if and only if

$$\frac{w}{3-2w} < \frac{c}{b}\left(1+\frac{\beta}{\alpha}\right). \tag{5}$$

See ESM Appendix A for further details.

As (3) and (4) indicate, both IS and ISp tend to resist invasion by ALLD when the cost-to-benefit ratio, $c/b$, is smaller and/or the probability of playing another round of game, $w$, is larger (see (A8) and (A26) in ESM Appendix A for the case of $e > 0$). When IS receives a Bad reputation, it resumes a Good reputation only after obtaining an opportunity to help a Good recipient, whereas a Bad ISp can recover the reputation also by punishing a non-helping donor. Thus, intuition suggests that there should be a parameter region in which ISp is superior to IS in the competition against ALLD. Indeed, (A26) and (5) jointly specify the condition under which ISp outcompetes ALLD in a broader range of population composition than IS does. Furthermore, (A19) shows that ISp outcompetes IS in direct competition when the cost of punishing, $\alpha$, is small relative to $b$.

### Model 2

Figure 2 shows the numbers of generations in which the population is in the cooperative state among 100,000 simulated generations for four variants of Model 2. The cooperative state was observed least frequently when all individuals were forced to be non-punishing and without the PRD (red), and slightly more frequently when the punishing trait was allowed to evolve under the condition that no individuals have the PRD (dark blue). The number of generations in the cooperative state increased considerably when both punishing and PRD traits were allowed to evolve (light blue), and slightly more so when the punishing trait was allowed to evolve under the condition that all individuals have the PRD (yellow). These results suggest that third-party punishment can promote



**Figure 1.** Sample trajectories of the evolutionary dynamics described by (2).
(a) Inequality (3) is not satisfied and (4) is satisfied. (b) All (3), (4) and (5) are satisfied. (c) Neither (3) nor (4) is satisfied. Parameter values used are $e = 0$, $c = 10$, $w = 0.9$, (a) $b = 11$, $\alpha = 1$, $\beta = 20$, (b) $b = 15$, $\alpha = 10$, $\beta = 20$, (c) $b = 11$, $\alpha = 1$, $\beta = 1$.
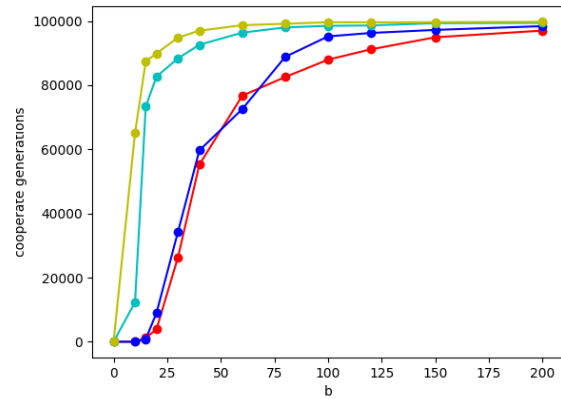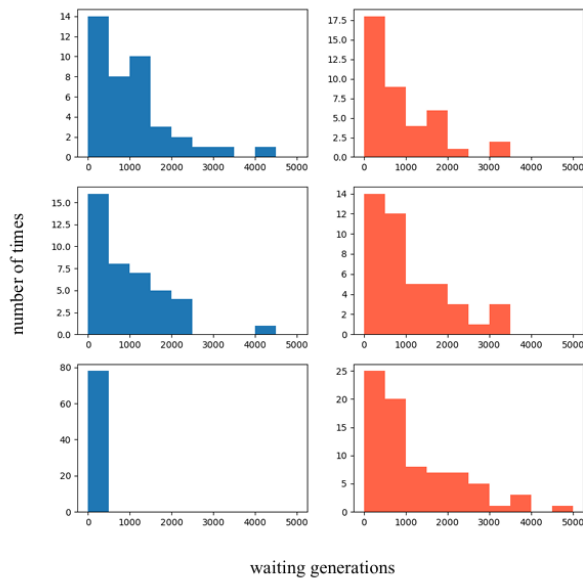


**Figure 2.** The numbers of generations in which the population was in the cooperative state across varying $b$. Different colors represent the results when all individuals were forced to be non-punishing and without the PRD (red), when the punishing trait was allowed to evolve under the condition that no individuals have the PRD (dark blue), when both punishing and PRD traits were allowed evolve (light blue), and when the punishing trait was allowed to evolve under the condition that all individuals have the PRD (yellow). Parameter values used are $n = 100$, $m = 4000$, $c = 10$, $\alpha = 10$, $\beta = 20$, $\mu_s = \mu_p = \mu_d = 0.01$, $u_p = 10$, $e = 0.05$.

Figure 3 histogram (left column blue, right column orange), y-axis "number of times", x-axis "waiting generations".
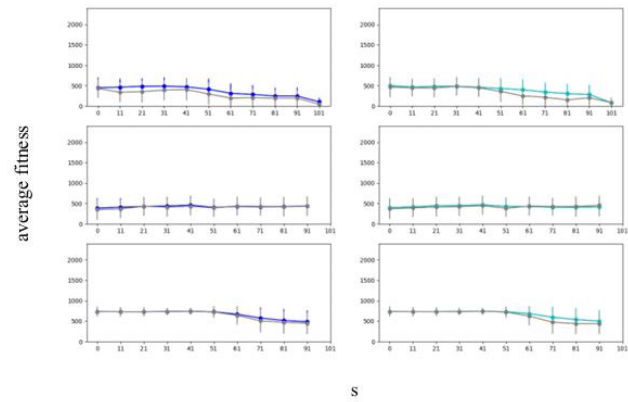
**Figure 3.** The numbers of waiting generations before the formation and collapse of cooperation.
The histogram shows the numbers of generations it took for a population in the non-cooperative state to reach the cooperative state (left) and for a population in the cooperative state to reach the non-cooperative state (right) when all individuals were forced to be non-punishing and without the PRD (top), when the punishing trait was allowed to evolve under the condition that no individuals have the PRD (middle), and when both the punishing and PRD traits were allowed to evolve (bottom). Parameter values used are $n = 100$, $m = 4000$, $b = 30$, $c = 10$, $\alpha = 10$, $\beta = 20$, $\mu_s = \mu_p = \mu_d = 0.01$, $u_p = 10$, $e = 0.05$.



Figure 4 panels, y-axis "average fitness", x-axis "s".

**Figure 4.** Average lifetime total payoffs of individuals with different helping thresholds.
For each class of *s* value, average lifetime total payoffs are shown for the case when the population was moving from the non-cooperating to cooperative state (top), moving from the cooperative to non-cooperative state (middle), and being in the cooperative state (bottom). In the left column, the results are compared between individuals with (blue) and without (gray) the punishing trait. In the right, the comparison is made between individuals with (light blue) and without (gray) the PRD trait. Error bars represent standard deviations. Parameter values used are the same as in Figure 3 except $\alpha = 2$.

cooperation through indirect reciprocity, and this is true even when the PRD is not presupposed to be present, but merely allowed to evolve.

A population that has reached the non-cooperative state typically stays there for a while and eventually moves to the cooperative state, after which it transits back to the non-cooperative state again. Figure 3 shows the distribution of the number of generations it took for a population in the non-cooperative state to reach the cooperative state, and for a population in the cooperative state to reach the non-cooperative state. The figure suggests that the third-party punishment along with the PRD considerably facilitates the transition to the cooperative state, but does not necessarily promote its maintenance.

Figure 4 shows the average lifetime total payoffs in a population when it is moving from the non-cooperative to cooperative state, moving from the cooperative to non-cooperative state, and in the cooperative state. Both punishing and PRD traits are suggested to give advantage to those individuals who have them over those who do not when cooperation is being established, whereas the effect is, if anything, less conspicuous when cooperation is being collapsed or maintained. See ESM Appendix B for further analyses.

## Discussion

Major findings of the study are summarized as follows. First, evolution of third-party punishment is possible in the presence of the punishment reputation dependency (PRD), that is, if punishing behavior against wrongdoers are socially valued (Model 1). Second, the PRD can coevolve with third-party punishment during the process where population-wide cooperation is established (Model 2). These results provide theoretical support for the hypothesis that third-party punishment has played a signaling role in the evolution of cooperation via reputation enhancement, which is consistent with the claim that third-party punishment serves as a costly signal of trustworthiness (Gordon et al., 2014; Jordan et al., 2016; Nelissen, 2008). Our models also highlight an interaction between punishment and indirect reciprocity, two key elements in the theory for the evolution of cooperation, and may play the role as a springboard for future research to deepen our understanding of the relationship between punishment, reputation, and cooperation.

Image scoring, as well as other strategies implementing indirect reciprocity, exclude individuals with bad reputation from the network of cooperation and thereby enjoy the benefit of mutual help. By introducing third-party punishment, selfish behavior is further suppressed through direct sanctioning, given that punishers are not lost from the population. Despite their cost of administering punishment, punishers gain in reputation when sufficiently many individuals have the PRD, where

having the PRD can be advantageous if punishing serves as an honest signal of the punisher's tendency to help. In Model 2, since individuals who punish more frequently tend to have a lower helping threshold and thus tend to help more frequently, an association between punishing and helping behaviors arises, promoting the evolution of the PRD. Hence, our result that third-party punishment is most effective when cooperation is being established may be because the punishment-help association is strongest at that stage, for it harbors a large variation in the helping threshold among individuals.

Let us discuss some limitations of our study. First, Model 2 suggests that third-party punishment does not facilitate the maintenance of the cooperative state. In fact, punishment seems to accelerate the collapse of the cooperative state by imposing costs for defect due to errors, similar to the collapse of cooperation as detailed in Panchanathan & Boyd (2003). In order to explain long-term maintenance of large-scale cooperation, some additional mechanisms need to be considered. Second, the present study used Image Scoring as the social norm under which individuals' behaviors are evaluated. While some studies have shown that behavior of human subjects in laboratory experiments is consistent with Image Scoring (Milinski et al., 2001), others have found that people behave in a more context-dependent way (Okada et al., 2018; Swakman et al., 2016). Theoretical studies have proposed other social norms within which to evaluate a donor's behavior toward a recipient, which take into account not only the donor's behavior, but also the recipient's reputation (Leimer & Hammerstein, 2001; Ohtsuki & Iwasa, 2004; Panchanathan & Boyd, 2003; Sugden, 1986), or even the donor's reputation (Ohtsuki & Iwasa, 2006). Alternative analyses based on other social norms may yield different results. Finally, our assumption that third-party punishment enhances the punisher's reputation has not been fully supported empirically. While the aforementioned literature suggests that exercisers of third-party punishment tend to be trusted by others, it is thus far not clear if the punisher gets rewarded (Balafoutas et al., 2014; dos Santos et al., 2013; Horita, 2010; Kiyonari & Barclay, 2008). Further empirical research in this regarded is awaited.

As future directions, it seems necessary to explore at least the following three aspects. First, as already mentioned, further analysis using social norms other than Image Scoring is worthwhile. Taking third-party punishment into consideration, there may be more social norms that can sustain cooperation than those currently considered within the framework of indirect reciprocity. Second, as mentioned earlier, our model does not explain the long-term maintenance of cooperation in the presence of errors in behavior. As large-scale cooperation seems stably maintained in many human societies, future work should explore the missing elements necessary to explain reality. Third, while we assumed that an individual's reputation is fully public, this is unrealistic in a large population. Third-party punishment may or may not play a greater role in a larger population, where reputation is less functional. Of potential relevance is the result of a behavioral experiment, which suggests that punishment is more likely to be exercised in larger groups (Marlowe et al., 2008). Further investigation in this line would be interesting.

In sum, we have shown the theoretical possibility of the coevolution of third-party punishment and punishment reputation dependency, and proposed a new framework for considering the relationship between punishment, reputation, and evolution of cooperation.

## Acknowledgments

## Author contribution
TO developed the study concept and design, and TO and YI did the analysis and wrote the manuscript.

## Supplementary material
Electronic supplementary material is available online.

## References
Alexander, R. D. (1987). *The biology of moral systems*. Aldine de Gruyter.

Balafoutas, L., Nikiforakis, N., & Rockenbach, B. (2014). Direct and indirect punishment among strangers in the field. *Proceedings of the National Academy of Sciences of the United States of America*, *111*(45), 15924–15927. https://doi.org/10.1073/pnas.1413170111

Barclay, P. (2006). Reputational benefits for altruistic punishment. *Evolution and Human Behavior*, *27*(5), 325–344. https://doi.org/10.1016/j.evolhumbehav.2006.01.003

dos Santos, M., Rankin, D. J., & Wedekind, C. (2013). Human cooperation based on punishment reputation. *Evolution*, *67*(8), 2446–2450. https://doi.org/10.1111/evo.12108

Fehr, E., & Fischbacher, U. (2004). Third-party punishment and social norms. *Evolution and Human Behavior*, *25*(2), 63–87. https://doi.org/10.1016/S1090-5138(04)00005-4

Gordon, D. S., Madden, J. R., & Lea, S. E. (2014). Both loved and feared: Third party punishers are viewed as formidable and likeable, but these reputational benefits may only be open to dominant individuals. *PLoS ONE*, *9*(10), e110045. https://doi.org/10.1371/journal.pone.0110045

Horita, Y. (2010). Punishers may be chosen as providers but not as recipients. *Letters on Evolutionary Behavioral Science*, *1*(1), 6–9. https://doi.org/10.5178/lebs.2010.2

Jordan, J. J., Hoffman, M., Bloom, P., & Rand, D. G. (2016). Third-party punishment as a costly signal of trustworthiness. *Nature*, *530*(7591), 473–476. https://doi.org/10.1038/nature16981

Kanakogi, Y., Miyazaki, M., Takahashi, H., Yamamoto, H., Kobayashi, T., & Hiraki, K. (2022). Third-party punishment by preverbal infants. *Nature Human Behaviour*, *6*(9), 1234–1242. https://doi.org/10.1038/s41562-022-01354-2

Kiyonari, T., & Barclay, P. (2008). Cooperation in social dilemmas: Free riding may be thwarted by second-order reward rather than by punishment. *Journal of Personality and Social Psychology*, *95*(4), 826–842. https://doi.org/10.1037/a0011381

Kurzban, R., DeScioli, P., & O'Brien, E. (2007). Audience effects on moralistic punishment. *Evolution and Human Behavior*, *28*(2), 75–84. https://doi.org/10.1016/j.evolhumbehav.2006.06.001

Leimar, O., & Hammerstein, P. (2001). Evolution of cooperation through indirect reciprocity. *Proceedings of the Royal Society B: Biological Sciences*, *268*(1468), 745–753. https://doi.org/10.1098/rspb.2000.1573

Marlowe, F. W., Berbesque, J. C., Barr, A., Barrett, C., Bolyanatz, A., Cardenas, J. C., Ensminger, J., Gurven, M., Gwako, E., Henrich, J., Henrich, N., Lesorogol, C., McElreath, R., & Tracer, D. (2008). More 'altruistic' punishment in larger societies. *Proceedings of the Royal Society B: Biological Sciences*, *275*(1634), 587–592. https://doi.org/10.1098/rspb.2007.1517

Milinski, M., Semmann, D., Bakker, T. C. M., & Krambeck, H. -J. (2001). Cooperation through indirect reciprocity: Image scoring or standing strategy? *Proceedings of the Royal Society B: Biological Sciences*, *268*(1484), 2495–2501. https://doi.org/10.1098/rspb.2001.1809

Nelissen, R. M. A. (2008). The price you pay: Cost-dependent reputation effects of altruistic punishment. *Evolution and Human Behavior*, *29*(4), 242–248. https://doi.org/10.1016/j.evolhumbehav.2008.01.001

Nowak, M. A., & Sigmund, K. (1998). Evolution of indirect reciprocity by image scoring. *Nature*, *393*(6685), 573–577. https://doi.org/10.1038/31225

Ohtsuki, H., & Iwasa, Y. (2004). How should we define goodness?—Reputation dynamics in indirect reciprocity. *Journal of Theoretical Biology*, *231*(1), 107–120. https://doi.org/10.1016/j.jtbi.2004.06.005

Ohtsuki, H., & Iwasa, Y. (2006). The leading eight: Social norms that can maintain cooperation by indirect reciprocity. *Journal of Theoretical Biology*, *239*(4), 435–444. https://doi.org/10.1016/j.jtbi.2005.08.008

Okada, I., Yamamoto, H., Sato, Y., Uchida, S., & Sasaki, T. (2018). Experimental evidence of selective inattention in reputation-based cooperation. *Scientific Reports*, *8*, Article 14813. https://doi.org/10.1038/s41598-018-33147-x

Panchanathan, K., & Boyd, R. (2003). A tale of two defectors: The importance of standing for evolution of indirect reciprocity. *Journal of Theoretical Biology*, *224*(1), 115–126. https://doi.org/10.1016/s0022-5193(03)00154-1

Raihani, N. J., & Bshary, R. (2015). The reputation of punishers. *Trends in Ecology & Evolution*, *30*(2), 98–103. https://doi.org/10.1016/j.tree.2014.12.003

Riedl, K., Jensen, K., Call, J., & Tomasello, M. (2012). No third-party punishment in chimpanzees. *Proceedings of the National Academy of Sciences of the United States of America*, *109*(37), 14824–14829. https://doi.org/10.1073/pnas.1203179109

Sugden, R. (1986). *The Economics of rights, co-operation and welfare*. Blackwell. https://doi.org/10.1057/9780230536791

Swakman, V., Molleman, L., Ule, A., & Egas, M. (2016). Reputation-based cooperation: Empirical evidence for behavioral strategies. *Evolution and Human Behavior*, *37*(3), 230–235. https://doi.org/10.1016/j.evolhumbehav.2015.12.001

Trivers, R. L. (1971). The evolution of reciprocal altruism. *Quarterly Review of Biology*, *46*(1), 35–57. https://doi.org/10.1086/406755