

1 Acting Without Considering Personal Costs
2 Signals Trustworthiness in Helpers but Not
3 Punishers

4 Nicole C. Engeler^{1*}, Nichola J. Raihani^{1,2}

5 ¹ Experimental Psychology, University College London, London, UK

6 ² School of Psychology, University of Auckland, Auckland, New Zealand

7

8

9 * Corresponding author: Nicole C. Engeler (engeler.nicole@gmail.com).

10 **Abstract**

11 **Third-party punishment and helping can signal trustworthiness, but the interpretation of**
12 **deliberation may vary: uncalculated help signals trustworthiness, but this may not hold for**
13 **punishment. Using online experiments, we measured how deliberation over personal costs and**
14 **impacts to targets affected trustworthiness of helpers and punishers. We expected that personal**
15 **cost-checking punishers and helpers would be trusted less. Conversely, impact deliberation was**
16 **expected to increase perceived trustworthiness of punishers but not helpers. Replicating previous**
17 **work, we found that refraining from checking the personal cost of helping signals trustworthiness,**
18 **although evidence for observers trusting uncalculating over calculating helpers is mixed. This did**
19 **not extend to punishment: only uncalculating *non*-punishers were more trustworthy than cost-**
20 **checking non-punishers. Impact deliberation results were mixed: deliberation affected the trust**
21 **and trustworthiness of non-helpers more than helpers and no conclusive results were found for**
22 **punishment. These results show that deliberation differentially affects assessments of those who**
23 **help or punish others.**

24 Introduction

25 Prosocial behaviours, such as helping and cooperating, can benefit others but often come at a
26 personal cost to the actor¹⁻³. Punishment, which involves an actor paying a cost to impose a cost on a
27 social partner^{4,5}, can encourage and maintain prosocial behaviours by deterring selfish actions⁶⁻¹³.
28 Although punishing anti-social behaviours can increase group-level cooperation, it also imposes a
29 cost on the punisher by requiring effort and time, and puts the punisher at risk of retaliation^{6,14,15}. To
30 understand why people invest in punitive acts, we must explain how punishment might ultimately
31 lead to downstream benefits for the punisher.

32 This question is particularly pertinent when it comes to third-party punishment, where a punisher
33 intervenes to punish a cheat even though they were not personally harmed by the cheat's behaviour
34 and may not interact with the target of punishment again in the future. Third-party punishment can
35 still provide reputation benefits to the punisher¹⁶⁻²¹, either by signalling their formidability (which
36 may deter their current social partners or bystanders from transgressing in the future^{18,22}) or by
37 signalling their cooperative intent (which may result in others being more likely to cooperate with
38 them²³ or choose them as partners for cooperative interactions^{16,24-27}). Third-party punishment can
39 therefore act as a signal that communicates an otherwise unobservable intent to act
40 prosocially^{21,23,26-30}. Accordingly, in some settings, individuals invest more in third-party punishment
41 when they are observed^{20,24} and are evaluated in a preferential manner by others for doing so²⁶.

42 Nevertheless, punishment is, by definition, a harmful act, which complicates inferences about the
43 punisher's intentions^{5,21,24}. Punishment could stem from antisocial, competitive, or spiteful
44 motivations rather than from a desire to cooperate, promote fairness, or uphold social norms⁵.
45 Indeed, compared to those who compensate victims, third-party punishers typically have higher
46 scores for antisocial personality traits such as Machiavellianism, narcissism, and psychopathy³¹.
47 Third-party punishment is therefore a more ambiguous signal of trustworthiness and cooperative
48 intent than helping or compensating a victim^{21,24,26,31,32}; and is most likely to signal cooperative intent
49 in scenarios where the punisher cannot compensate the victim, and where self-serving motives are
50 less likely²¹ (e.g. when the punisher does not increase their own payoffs relative to those of the
51 target when they punish^{5,21}).

52 The potential for helpful acts to signal cooperative intent also depends, to some extent, on context³³.
53 One such context concerns whether the helpful act was calculated or not: uncalculated help is a
54 stronger signal of the helper's cooperative disposition than calculated help. One recent study
55 operationalised uncalculated help by measuring whether individuals looked at the personal cost to
56 themselves before helping, and by recording how long it took individuals to make their helping
57 decision once the cost of helping was revealed²⁶. Helpers who check the cost are ostensibly weighing
58 up the costs and benefits of their actions, suggesting that cooperating is a more strategic and
59 calculated move. Response time is also informative about a person's underlying commitment to
60 cooperation, as slower decisions indicate greater decision conflict³⁴. Observers use information about
61 another's decision time to infer levels of conflict experienced and to make predictions about whether
62 a person is making a calculated or uncalculated decision to cooperate^{26,35-37}, and whether to trust
63 them²⁶. This previous work found that uncalculated cooperation was a more reliable signal of
64 trustworthiness than calculated cooperation²⁶. Moreover, people were apparently aware of the
65 signalling value of uncalculated cooperation and were less likely to check the costs of cooperating
66 and made cooperative decisions more quickly when observed²⁶.

67 Although both third-party punishment and helping are prosocial acts, decision conflict for these
68 behaviours may not be interpreted in the same way, and it is therefore unclear whether findings on
69 uncalculated help²⁶ would be expected to translate directly to the punishment setting. We address
70 this issue here. Decision conflict over whether to help another is likely to stem from self-interested
71 considerations of whether to pay a personal cost. Decision conflict over punishment, by contrast,
72 could also stem from concerns about inflicting harm on the target. This perspective yields nuanced

73 predictions about the two different measures of decision conflict above. Checking the personal cost
74 of administering punishment is likely to indicate a self-interested concern about personal costs. As
75 with helping, such calculated decisions may be perceived negatively by observers. However, because
76 punishment involves imposing costs on another individual, taking longer to decide whether to punish
77 another person might not be viewed negatively and could even be viewed positively. Perhaps
78 carefully thinking about and balancing both the prosocial aspect as well as the negative
79 consequences for the punished is the 'right' thing to do when deciding whether to engage in third-
80 party punishment.

81 To better define the conditions under which punishers are viewed positively and to differentiate
82 between punishment and helping as signals of trustworthiness, we conducted two studies. Study 1
83 aimed to replicate Jordan et al.'s²⁶ research on uncalculated helping and extend it to punishment,
84 asking how deliberation over personal costs affects trustworthiness perceptions in both cases. Study
85 2 extended this by asking how deliberation over the impacts on targets affects perceptions of
86 trustworthiness for both punishment and helping. With this approach we hoped to understand
87 whether and why punishers are evaluated differently to helpers, and to show how deliberation
88 differentially signals trustworthiness in decisions to punish or help others. Overall, we expected to
89 show that deliberating about the *personal cost* of one's actions signals untrustworthiness for both
90 punishing and helping, whereas deliberating about the *impact* of those actions on others constitutes
91 the differentiating factor between helping and punishing behaviours. Specifically, we expected that
92 those who deliberate about the impact of punishment are viewed relatively positively, whereas
93 those who deliberate about the impact of helping are viewed relatively negatively. See Table 1 for
94 descriptions of all preregistered hypotheses.

95 Across two studies, comprising five experiments (Table 2), we investigated whether and when
96 uncalculated punishment and help are used as signals of trustworthiness. As in Jordan et al.²⁶, all
97 experiments had two stages: a first stage, where Player A could pay a cost to help a victim / punish a
98 cheat; and a second stage, where Player B decided whether and how much to trust Player A. Any
99 money entrusted by Player B was tripled by the experimenter and Player A then decided how much
100 to return to Player B, yielding a measure of trustworthiness. We included two conditions: Player B
101 was either be able to make their trusting decision based on (i) Player A's decision and decision
102 process in the first stage (via cost-checking / decision time), or (ii) solely on Player A's help /
103 punishment decision, with the decision *process* remaining concealed.

104 Four of our five experiments operationalised deliberation through cost- or impact-checking
105 behaviours. Due to financial constraints, we included only one decision time study, specifically
106 focussing on personal costs of punishment. This context holds particular significance, as we expected
107 to observe significant differences between decision time and cost checking. Longer decision times
108 may be attributed to concern for the target, in addition to self-interested considerations. Because
109 decision conflict over punishment could also stem from concerns about inflicting harm on the target,
110 this experiment forms a bridge between Studies 1 (personal cost deliberation) and 2 (target impact
111 deliberation).

112 Study 1 investigated whether and how personal cost deliberation signals trustworthiness. Calculated
113 behaviours were operationalised as decisions to check the personal cost of helping / punishing (1.1 &
114 1.2a), or long decision times after the cost of punishing is revealed (1.2b). Study 1 therefore
115 comprises three separate experiments (Table 2), investigating how deliberating over decisions to
116 help (1.1) or to punish (1.2a & 1.2b) signal trustworthiness.

117 Study 2 was similar to Study 1 but investigated whether deliberation over impacts on targets of help
118 / punishment signals trustworthiness. In Study 2, the personal cost of help / punishment was
119 therefore known. Calculated behaviours were instead defined as decisions to check the impact of
120 help / punishment ('impact checking', 2.1 & 2.2).

121

122 When deliberation was operationalised as checking the personal cost of helping / punishing
123 (Experiments 1.1 and 1.2), we predicted that uncalculated decisions signal trustworthiness:
124 individuals who do not check the cost of helping / punishing would be entrusted with more money
125 than those who do check the cost. We also expected participants to be sensitive to these potential
126 reputation benefits and to be less likely to check the cost of helping / punishing when they were
127 observed than when their decision process was hidden from Player B. Further, we expected that
128 individuals who do not check the cost of helping / punishing would be more trustworthy than those
129 who do check the cost. Lastly, we anticipated uncalculated help to be a stronger signal of
130 trustworthiness than uncalculated punishment.

131 When personal cost deliberation was operationalised as decision speed (Experiment 1.2b), our
132 predictions were more nuanced. Here, the cost of punishing was revealed to Player A right before
133 they made their decision. As the sole new information provided to participants prior to measuring
134 their decision-making time related to the personal cost of punishing, we predicted that here too
135 uncalculated punishment would be used as a signal of trustworthiness. Specifically, we expected that
136 participants would exhibit faster decision-making when their decision time was revealed to others, in
137 contrast to when the decision process was concealed and unable to confer any reputation-related
138 advantages, as was observed with helping in Jordan et al²⁶.

139 Nevertheless, we also envisaged some differences in how quick decisions to punish might be
140 perceived by others, which would be driven by the different motives attributed to punishers and
141 helpers. Jordan et al.²⁶ found that fast decisions to help others were perceived positively, likely
142 because fast decisions were associated with less decision conflict. However, because punishment
143 inflicts harm on the target, fast decisions to punish could be evaluated differently^{5,21}. One possibility
144 is that observers interpret a fast decision to punish as the punisher being moral and interested in
145 restoring fairness regardless of the cost to the self. Here, uncalculated punishment would be
146 approved of, and observers would infer that fast punishers were more trustworthy than slow
147 punishers. Another possibility, however, is that observers may approve of more considered decisions
148 to punish others, if they infer that decision conflict stems from concern about the harm caused to the
149 target. Thus, slow punishers may not be evaluated as negatively and, consequently, we expected
150 decision speed to be a weaker signal of trustworthiness than cost-checking decisions for punishment.

151 Despite this ambiguity, we still expected fast punishers to be trusted more than slow punishers in the
152 context of personal cost deliberation. This is because observers were informed that the only new
153 information Player A received right before making their punishing decision was the cost of
154 punishment to themselves. Observers should infer, therefore, that deliberation stems only from the
155 consideration of this personal cost and not the impact to the target. In addition, to disambiguate
156 personal costs from harm aversion, we set the minimum potential cost of helping or punishing to be
157 £0.00 in Study 1. Therefore, we anticipated that observers would send more of their endowment to
158 third-party punishers who made their decision quickly, using uncalculated third-party punishment as
159 a signal of trustworthiness. Similarly, we believed that being slower in the decision to punish would
160 reflect the punisher's conflict about whether paying the cost would be beneficial to themselves,
161 rather than an additional consideration of whether harming the violator is the "right thing to do".
162 Thus, we expected uncalculated (fast) punishers to return more money than calculated (slow)
163 punishers.

164

165 Study 2 was designed to address some of the open questions raised by Study 1 – specifically, whether
166 deliberating over the impact on targets is perceived more positively for punishment than for helping.
167 For both helping and punishment, we expected participants to be sensitive to the reputational
168 consequences of their behaviour, albeit in different ways. If they want to be evaluated positively by
169 an observer, helpers should be less likely to check how much helping will impact targets, whereas
170 punishers should be more likely to check how much their actions will harm another. In other words,
171 unlike Study 1, calculated punishment now served as a signal of trustworthiness: punishers who

172 deliberate about the impact to the target should be entrusted with more money by observers and
173 should be more trustworthy, compared to punishers who do not deliberate in this way. For helpers,
174 as in Study 1, we expected uncalculated decisions to signal trustworthiness.

175

176 We also had several secondary predictions pertaining to decisions *not* to help or punish that can
177 further clarify when and why deliberating over social actions carries reputation consequences.
178 Specifically, we were interested in whether non-punishers are evaluated differently to non-helpers –
179 and to what extent deliberation moderates these perceptions.

180 To understand how non-helpers and non-punishers are perceived, we must consider the potential
181 motives driving decisions not to help or punish others. One primary reason individuals may refrain
182 from helping or punishing in this task is due to personal costs. Additionally, non-helpers may not be
183 especially motivated to help others because they are antisocial or inequity averse (i.e., they do not
184 want someone else to receive more than them). For punishment, individuals may also refrain
185 because they are averse to harming others. To disambiguate personal costs from harm aversion, we
186 set the minimum cost of helping or punishing to be £0.00 in Study 1 (the minimum impact of helping
187 or punishing was set to £0.01, so that some impact of investing in help or punishment was
188 guaranteed). This feature allowed us to make nuanced predictions about how non-helpers and non-
189 punishers would be evaluated.

190 A decision not to help is likely to stem primarily from self-interest and, consequently, unhelpful
191 individuals are generally not trusted by others^{24,26}. Those who decide not to help without checking
192 the cost to themselves (Experiment 1.1) or the impact to the target (Experiment 2.1) might be
193 evaluated especially negatively as it indicates an unwillingness to help, even if helping might impose
194 no personal cost (Study 1) and regardless of the potential benefit to the target (Study 2). Conversely,
195 deliberation indicates that the individual at least considered helping before deciding not to. Thus, in
196 general, uncalculated decisions not to help should be evaluated negatively by observers in both
197 studies, and uncalculating non-helpers should be less trustworthy than calculating non-helpers.

198 The same is not true for punishment: refraining from punishing others could stem from self-interest
199 or from harm aversion. The possibility for non-punishment to stem from harm-aversion may help to
200 explain why non-punishers can sometimes be trusted as much as punishers²⁴. In Study 1, we
201 expected that non-punishers who do not consider the personal cost might be perceived as harm
202 averse – individuals who would not punish even if it were free to do so. Here, non-punishers who do
203 not check the cost of punishing might be perceived as (and should actually be) relatively trustworthy.
204 Conversely, non-punishers that check the personal cost before refraining from punishment should be
205 seen as less trustworthy, because the inference is that these decisions were driven by self-interest
206 (i.e. the personal cost being too large) rather than harm aversion.

207 In Study 2, uncalculated non-punishers (those who did not check the impact of punishment on
208 targets) might either be completely harm averse, or might be unwilling to pay the personal cost
209 associated with punishing. Calculated non-punishers (those who *did* check the impact of punishing)
210 by contrast, are those who might be willing to incur the personal cost of punishing but who wanted
211 to know what impact this would have on the target before doing so. For uncalculated non-
212 punishment to be perceived as more trustworthy than calculated non-punishment, participants
213 would need to believe the target deserved no punishment. However, this would also imply tacit
214 acceptance of the behaviour exhibited by Player 2 (returning nothing after their partner entrusted
215 them with their entire endowment). Since both Players A and Players B knew they would
216 subsequently be playing a Trust Game together, attitudes towards the target and what is considered

217 acceptable behaviour in a Trust Game are relevant. As such, we expected calculated non-punishers
218 to be perceived as and actually be more trustworthy than uncalculated non-punishers.

219 Given that the motives for actions are somewhat more transparent than those for non-actions, and
220 that we expected incurring a cost to punish the cheat or help the cheated to be seen as more
221 prosocial than doing nothing, we anticipated that deliberation (both of personal cost and target
222 impact) would have a more substantial impact on trust and trustworthiness when individuals chose
223 to punish or help than when they did not. Finally, we again expected non-punishers' decision speed
224 to have a weaker effect on trust and trustworthiness than cost-checking decisions.

225 We must note, however, that predictions regarding the trustworthiness of punishers/helpers and
226 non-punishers/non-helpers are considered exploratory, as we did not know whether they would
227 achieve 95% power.

228 See Table 1 for more detailed descriptions of all hypotheses, and Figure 1 for a visualisation thereof.

229 **Methods**

230 ***Ethics information***

231 The research complies with all relevant ethical regulations. The study was approved by the UCL Ethics
232 Board (Project ID: ICN-NH-PWB-7-1-23A). Informed consent was obtained from all participants.
233 Although 'Player 3' / 'The Sender' (Player A) and 'The Receiver' (Player B) really did exist and
234 participants' decisions really did influence their own payoff and that of fellow participants, 'Player 1'
235 (the cheater in punishing contexts or the recipient in helping contexts) and 'Player 2' (the violator)
236 did not actually exist. Therefore, participants were fully debriefed after the study, and only those
237 who previously indicated they are willing to take part in studies involving deception were invited to
238 the study. Participants were compensated at an hourly rate of £9.

239 ***Design***

240 We conducted experiments of highly similar designs to investigate different operationalisations of
241 each calculating behaviour (i.e., (i) checking the cost or impact of punishment or helping, and (ii)
242 punishing cost decision time). Each experiment recruited separate sets of both Players A and Players
243 B, and had two conditions: decision process hidden or decision process observable, with a between-
244 subjects design. The studies were built in Qualtrics (www.qualtrics.com) and consisted of two-stage,
245 incentivized, anonymous economic games (see Figure 2 for a visualisation of the study design).
246 Players A made decisions during both games, whilst Players B only made decisions during the second
247 game. Prior to making any decisions, Players A and Players B read the instructions for both games to
248 ensure that they could make informed decisions and comprehension could be assessed.

249 The following describes the procedure for Study 1 in a punishing context. We outline any differences
250 for the helping context below.

251 **Stage 1**

252 Stage 1 is a standard Trust Game that was observed by both Player A and Player B. The outcome of a
253 trust game played between 'Player 1' and 'Player 2' (who do not actually exist) were presented to
254 Player A and B. Participants were told that Player 1 started with a £0.10 endowment and could
255 choose how much of this endowment to send to Player 2, who started the trust game with no
256 money. Participants were told that the amount entrusted was tripled by the experimenter and that
257 Player 2 could then choose how much to return to Player 1. Participants were told that Player 1 and
258 Player 2 already made their decisions: Player 1 sent their entire endowment (£0.10) to Player 2, who
259 returned nothing, now leaving Player 1 with no money and Player 2 with £0.30.

260 Player A started with a £0.10 endowment and could then choose whether to use some of their
261 endowment to 'punish' Player 2. Participants did not know exactly how much it would cost to punish
262 Player 2, except that it would be somewhere between £0.00 and £0.10. Punishing always removed
263 £0.15 from Player 2's bonus. We chose this amount as it leaves Player 2 with a bonus of £0.15, which
264 would have been the 'fair' amount for Player 2 to receive in the trust game. Moreover, Player 2 still
265 finishes Stage 1 with more than Player A, thereby reducing potential that punishment is perceived as
266 stemming from Player A wanting to have more than Player 2.

267 Players A additionally decided whether to make their punishing decision in a calculating or an
268 uncalculating manner. The way this was operationalised varied between participants across
269 experiments (Table 2).

270 **Cost-checking:** Player A chose whether to check the cost of punishing Player 2 before making
271 their punishing decision. The cost of punishment was always £0.05.

272 **Decision time:** Player A was told that the cost of punishing is £0.05, after which they
273 immediately made their punishing decision on the same page. The time spent on this page was
274 recorded to assess their decision time.

275 **Stage 2**

276 In Stage 2, Players A and B participated in a trust game with each other. The trust games 'rules' were
277 the same as in Stage 1. Player B started with a bonus of £0.10 and had to choose how much of their
278 endowment (£0.00 - £0.10) to send to Player A. Any amount sent was tripled. Player A then chose
279 what percentage of the endowment to return to Player B.

280 In the 'process observable' condition, Players B could condition their sending decision both on (i)
281 Player A's Game 1 decision process (whether Player A chose to reveal the cost of punishing in
282 Experiment 1.2a, or Player A's fast/slow decision time in Experiment 1.2b) and (ii) Player A's Game 1
283 decision (whether Player A punished Player 2). In the 'process hidden' condition, Players B could only
284 condition their sending decision on Player A's Game 1 decision (whether they chose to punish Player
285 2). We employed the strategy method for both players: Players B decided how much to send to a
286 Player A who engaged in all possible combinations of punishing decisions and (depending on the
287 condition) processes, without knowing what Player A did. Similarly, Players A decided what
288 percentage of the amount they received from Player B to return, without knowing how much Player
289 B sent. Participants were told that the choice that matched the decision the other player made would
290 determine their payoff.

291 All participants were asked several comprehension questions, primarily to assess their understanding
292 of the incentive structure of both games. In addition to Player A decision time, we recorded the time
293 spent on the three pages with comprehension questions in Experiment 1.2b. These recordings serve
294 as a control for general comprehension and reading speed.

295 After completing the data collection, we randomly matched pairs of Players A and Players B who
296 participated in the same experiment and condition. The decisions participants made during the study
297 then determined their bonus payments.

298 **Differences across studies and contexts**

299 The procedure across punishing and helping contexts was identical, except that instead of choosing
300 whether to punish Player 2, Players A decided whether to use some of their endowment to help
301 Player 1.

302 Whereas Study 1 explored how participants respond to an unknown personal cost of helping or
303 punishing, Study 2 explored how participants respond to the unknown impact of helping or punishing
304 another. In Study 2 participants were told that the personal cost of punishment (or helping) is £0.05
305 but they did not know exactly how much punishing would remove from Player 2 (or helping would
306 benefit Player 1), except that it would be somewhere between £0.01 and £0.30. The maximum
307 impact of £0.30 in Study 2 is equivalent to the maximum personal cost of £0.10 in Study 1: £0.10 was
308 the entire endowment of the helper/punisher in Study 1, whereas £0.30 was the entire endowment
309 of the target (in the punishment condition) in Study 2. Participants were informed that the minimum
310 potential impact of helping/punishing a target is £0.01 because £0.00 would indicate no punishment
311 or no help. When the impact was revealed, punishing still removed £0.15 from the target, and
312 helping still delivered £0.15 to the target, just as in Study 1.

313 The procedure for Stage 2 only differed in that the Stage 1 procedure influenced what decisions and
314 decision processes Players B could condition their sending decision on (i.e., whether it was a
315 punishing or helping decisions and whether it centred around personal cost checking or impact
316 checking).

317 The procedure and instructions as seen by participants can be viewed in the supplementary
318 information under "Supplementary Methods".

319

320 ***Sampling plan***

321 ***Power Analysis***

322 Our power calculation was conducted in R³⁸ using the package ‘pwr’³⁹ with the ‘pwr.f2.test’ function.
323 We used a power of 0.95 with a 0.05 significance level and a one numerator degree of freedom (u ,
324 the number of coefficients in the model without the intercept). While estimating the required sample
325 size, we referred to Jordan et al.'s supplementary materials for effect sizes, but specific effect sizes
326 were not explicitly mentioned. We acknowledge that the available coefficients in their
327 supplementary materials vary considerably, but generally produced small to medium effect sizes. As
328 their study closely matches our experimental design, procedure, and research questions, we used an
329 effect size of $f^2 = 0.02$ in our power analysis. According to Cohen's guidelines⁴⁰, $f^2 \geq 0.02$ represents a
330 small effect size⁴¹. Because our main interests focussed on third-party punishment rather than
331 helping, we expected to find similar or smaller effect sizes. Nevertheless, we must acknowledge that
332 our choice of $f^2 = 0.02$ might be considered a heuristic approximation rather than a precise
333 estimation based on a formal inspection of Jordan et al.'s results. Based on the model with the
334 highest number of predictors (as $n = v + p$, with p being the number of predictors including the
335 intercept and v the degrees of freedom for the denominator), a sample size of 653 would be needed.
336 However, as each of the models upon which this calculation was based involves either a Player A or a
337 Player B participant, taking part in one of five experiments, in either the process hidden or in the
338 process observable condition, a sample size of 13,060 (i.e., 1,306 Player A - Player B pairs per
339 experiment) was needed. As we could not predict how many participants would decide to punish/
340 help or not punish/ help, it was not possible to ascertain before data collection whether 95% power
341 would be achieved for all analyses. Results that did not meet the power requirements are therefore
342 interpreted as suggestive, pending confirmation in future research.

343 **Participants**

344 We ran our experiments on Prolific (<https://prolific.co/>). Participants were invited to take part if they
345 previously indicated on Prolific that they (i) are aged 18 or above, (ii) are from the UK, so that the
346 currency specifications are familiar, (iii) are fluent in English, (iv) have the maximum approval rate of
347 100, and (v) selected “Yes, I would be comfortable to take part in such a study” to the question
348 “Would you be happy to take part in a study where you are intentionally given inaccurate
349 information about other participants and the study? You would be debriefed after the study”. To
350 avoid participants taking part in more than one experiment, we launched the experiments in
351 sequence, and allowed only new participants to take part. As preregistered, we lowered the approval
352 rate to 97, in one-unit increments, as we did not reach enough participants with the maximum
353 approval rate of 100.

354 **Data Exclusion**

355 We used the “force response” feature in Qualtrics to ensure that we did not receive incomplete
356 responses. As in Jordan et al.²⁶, responses by participants who failed more than one attention check
357 were still included in the analyses. However, we re-ran the same analyses excluding those who failed
358 more than one attention check, and reported this version when it lead to significant differences in
359 results. Any duplicate responses were removed.

360

361 **Analysis Plan**

362 Our analyses were conducted in R³⁸, and all analytical decisions for our hypotheses were
363 independent of each other. For hypotheses in which we predicted cost-checking or impact-checking
364 decisions (binary variables), we ran a logistic regression. For all other hypotheses we used linear
365 regressions, as they predict decision speed, sending decisions, or returning decisions (continuous
366 variables). Decision speed was a continuous variable when returning decisions were predicted, but a
367 dummy variable (median split of relatively fast or slow) was used when predicting sending decisions.
368 For ease of interpretation the measure for endowment sent was transformed from an absolute value
369 (pence) sent to the percentage of endowment sent. The return measure was not transformed, as
370 Players A already indicated what percentage, rather than what absolute value, they wished to return.

371 In instances where Players B made sending decisions based on Player A's decision process, analyses
372 were restricted to the process observable condition, as Players B did not know Player A decision
373 processes in the process hidden condition. Due to Players B making multiple sending decisions in
374 each of these analyses (based on the possible decisions made by Player A during the first stage), each
375 sending decision was treated as an observation and robust standard errors were clustered on
376 participant ID to account for the non-independence of repeated observations from the same
377 participant (i.e., two observations per participant for decision process hidden, and four observations
378 for decision process observable conditions). To accomplish this, we utilized the `lmtest` package⁴²,
379 employing the functions `coefest()` with the argument `vcov = vcovCL` to specify the use of the
380 sandwich estimator, and `coefci()`. As data was collected on Player A's decision process in both
381 conditions (even when Player B could not observe it), data from both the observable and hidden
382 condition were used for analyses of returning decisions. As variance in decision time could also be
383 caused by general comprehension ability or reading speed, rather than solely the time taken to reach
384 a punishing decision, we included a control for general comprehension and reading speed when
385 Player A decision time is an independent variable. General comprehension and reading speed was
386 operationalised as the natural log transformed sum of time spent on the three comprehension
387 question pages. All reported coefficients are unstandardised. For more detail of individual analysis
388 methods for our hypotheses, see Table 1.

389 **Preregistered Exploratory Analyses**

390 As we could not know how many participants would decide to punish/ help or not punish/ help, it
391 was impossible to ascertain before data collection whether 95% power would be achieved for all
392 analyses in which Player A returning decisions are predicted. Hypotheses H14.1 to H20.2 are
393 therefore considered exploratory analyses. Results that do not meet the power requirements are
394 interpreted as suggestive, pending confirmation in future research. This applied to hypotheses
395 H14.2a, H14.2b, H15.1, H19.2 and H20.1, as they did not meet power requirements.

396 **Bayesian Analyses**

397 In addition to our frequentist analyses, we conducted equivalent Bayesian analyses to assess the
398 evidence for each hypothesis compared to the null hypothesis. We used the `BayesFactor` package⁴³
399 with the `lmBF` function for linear regressions with return decisions as the response variable. For the
400 logistic regressions and linear regressions with sending decisions as the response variable, we used
401 the `brm()`, `bridge_sampler()` and `bayes_factor()` functions from the `brms` package⁴⁴, which is better
402 suited to handle the repeated observations from Player Bs. For analyses with repeat observations,
403 we fit mixed models in `brm()` with ID as random effect.

404 We constructed effect priors that are zero-centred t-distribution priors with 4 degrees of freedom.
405 The prior width was designed such that only one-third of the prior mass on each side of zero is larger
406 than the desired effect (i.e., the relevant coefficient observed in Jordan et al.²⁶). Specifically, for any
407 desired effect, one-third of the prior mass on each side of zero was more extreme than the absolute
408 value of the desired effect. The total prior mass smaller than the desired effect was calculated as $0.5 + 0.5 * 2/3 = 0.83333$ (e.g. assuming an effect of 5.6 would lead to an effect prior with scale of 5.09).
409 To achieve this, we used the below R code to calculate the scale of the prior width for a given desired
410 effect (e.g. 5.6):

```
412 desiredEffect <- 5.6  
413 myt <- function(x) { abs(extraDistr::qlst(0.5 + 0.5 * 2/3 , df = 4, mu =  
414 0, sigma = x) - desiredEffect) }  
415 calc_scale <- optimize(myt, interval = c(0, 20))  
416 prior_width_scale <- calc_scale$minimum
```

417 The defined function "myt" calculates the absolute difference between the desired effect size (in this
418 example 5.6) and the quantile of the t-distribution with 4 degrees of freedom and zero mean,
419 corresponding to the prior mass of $0.5 + 0.5 * 2/3$. The "optimize" function in R then finds the value of

420 the scale parameter that minimizes the absolute difference between the observed effect size and the
421 quantile of the t-distribution. The resulting prior_width_scale value is what we used as the width of
422 our prior distributions.

423 We chose this specification because previous research from Jordan et al.²⁶ indicates effect sizes may
424 generally be small, and as we investigated punishment as well as helping, effect sizes in punishing
425 contexts may be smaller still. However, we still allowed for the possibility that we could sometimes
426 find larger effects.

427 All other priors were weakly informative, using a zero-centred t-prior with 4 degrees of freedom and
428 a scale of 10. We chose these weakly informative priors to allow for some flexibility in the effect size
429 estimates while still constraining them to reasonable values. The choice of 4 degrees of freedom and
430 a scale of 10 reflects our prior belief that the effect size was unlikely to be very large, but may
431 occasionally have been larger than expected.

432 To ensure that our prior on the effect size was appropriate, we set rscalefixed = 0.5 when using the
433 BayesFactor package, as this is the smallest recommended prior on the effect size⁴⁵.

434 For hypotheses investigating the effect of deliberation on trust and trustworthiness, based on
435 whether it is measured through cost checking or decision time, or takes place in the context of
436 helping or punishing, we could not directly rely on equivalent coefficients from previous research to
437 set priors as we did above. However, as we expected small effects, we set rscalefixed = 0.5 here as
438 well, and to be conservative used the smallest interaction effect found in Jordan et al. to calculate
439 the prior scale for analyses in which sending decisions are the response variable.

440 We also conducted sensitivity analyses for each Bayes factor test by conducting two additional
441 analyses: one with a prior scale of 0.5 times the original value and one with a prior scale of 1.5 times
442 the original value. We report the results of these sensitivity analyses in our supplementary materials,
443 unless they changed the direction of the Bayes factor, in which case they are reported in the main
444 text.

445 In evaluating the strength of evidence for or against the alternative hypothesis compared to the null
446 hypothesis, we used common decision heuristics^{46,47} and considered Bayes factors of 3 as weak
447 evidence in favour of the alternative hypothesis, and Bayes factors of one-third as weak evidence in
448 favour of the null hypothesis over the alternative hypothesis. Bayes factors of 10 or more were
449 considered substantial evidence for the alternative hypothesis. Conversely, Bayes factors of one-
450 tenth or less were considered substantial evidence for the null hypothesis. In cases where the Bayes
451 factor fell between these thresholds, we concluded that the data provided no strong evidence for
452 either the alternative or the null hypothesis and that more data were needed to draw a conclusive
453 inference.

454

455 **Protocol Registration**

456 The Stage 1 protocol for this Registered Report was accepted in principle on 13th November 2023.
457 The protocol, as accepted by the journal, can be found at
458 <https://doi.org/10.6084/m9.figshare.24559462.v1>.

459

460 **Deviations from Stage 1 protocol**

461 Due to some participants starting but not finishing the experiment, some condition cells were
462 unbalanced. We therefore recruited an additional 13 participants (five each in Experiments 1 and 4,
463 and three in Experiment 2) to ensure that each condition reached the preregistered number of
464 participants, bringing the total sample size to 13073 rather than 13060. Originally, we planned to run

465 Bayesian analyses for H9.2b with lmbf(). However, as the function currently does not allow for
466 models containing both continuous and categorical predictors, those models were fit with brm() as
467 specified in the Bayesian Analysis section instead. Lastly, to maintain consistency with our registered
468 analyses, we incorporated additional analyses centring around non-action in the comparison
469 between help and punishment, ensuring comprehensive coverage and completeness across all
470 hypotheses.

471

472 **Data availability**

473 All study data and materials, as well as the laboratory log are available on OSF under this link:
474 <https://osf.io/y2hgu/>.

475

476 **Code availability**

477 The analysis code is available on OSF under this link: <https://osf.io/y2hgu/> (project DOI:
478 10.17605/OSF.IO/Y2HGU).

479 Results

480 Median completion time for experiments ranged between six and seven minutes. Demographics
481 were similar across experiments. Participants in Experiment 1 were aged between 18-80 years ($M =$
482 39.45 , $SD = 12.51$) with 1381 women and 1207 men (17 identified as 'other' and 7 preferred not to
483 say). In Experiment 2, participants were aged between 18-79 years ($M = 39.16$, $SD = 12.36$) with 1519
484 women and 1071 men (18 identified as 'other' and 7 preferred not to say). In Experiment 3,
485 participants were predominantly women (1598 women, 990 men, 22 'other', and 2 preferred not to
486 say) and were aged between 18 and 91 years ($M = 39.7$, $SD = 12.65$). 63% of participants in
487 Experiment 4 were women (1645 women, 949 men, 21 'other', 2 preferred not to say) with a mean
488 age of 37 years ($SD = 11.79$; range: 18-80 years). In Experiment 5, the average age was 38 years ($SD =$
489 12.45 ; range: 18-78 years), and 61% of participants were women (1581 women, 1002 men, 20
490 'other', and 9 preferred not to say). Data largely conformed to the assumptions of the tests used, but
491 decision time was heavily skewed. As pre-registered, punishing decision time was therefore natural-
492 log transformed.

493 Preregistered Primary Hypotheses: Deliberation over help and punishment as a signal of 494 trustworthiness

495 Our main prediction was that people would strategically adjust their decision-making process to gain
496 reputational benefits. When considering the personal cost of helping (Exp. 1) and punishment (Exp. 2
497 & Exp. 3), as well as when considering the impact of helping (Exp. 4), we expected Players A to be
498 more likely to make uncalculating decisions (by deciding quickly or without checking the personal
499 cost or target impact) when their decision process was observed by others (and could therefore
500 potentially confer reputation benefits), compared to when their decision process was hidden. We
501 predicted the opposite when participants considered the impact of punishment on a target (Exp. 5):
502 as punishment is a harmful act, we expected Players A to make calculating decisions (by checking the
503 impact that punishment would have on the target) when their decision process was observed. See
504 Figure 3 for a visualisation of results and Table 2 for the nomenclature of experiments.

505 In line with predictions, participants were around half as likely (odds ratio (OR) = 0.64, 95%
506 confidence interval (CI) = [0.50, 0.82]) to check the personal cost of helping when their decision
507 process was observable than when it was hidden (Exp. 1, H1.1). Specifically, 78% of Players A checked
508 the cost of helping in the decision process hidden condition, whereas only 69% did so in the decision
509 process observable condition ($b = -0.44$, 95% CI = [-0.69, -0.19], $p < 0.001$, BF = 71.29).

510 Similarly, the odds of checking the personal cost of punishing decreased by around 41% (OR = 0.59,
511 95% CI = [0.47, 0.73]) when the participants' decision process was observable compared to hidden
512 (Exp. 2, H1.2a). Specifically, 67% of Players A checked the cost of punishing in the decision process
513 hidden condition, whereas only 54% did so in the decision process observable condition ($b = -0.53$,
514 95% CI = [-0.76, -0.31], $p < 0.001$, BF = 6972.38).

515 Players A also made significantly faster punishing decisions when their decision process was
516 observable ($M = 2.11$ log-seconds, 95% CI = [2.06, 2.15]) compared to hidden ($M = 2.19$ log-seconds,
517 95% CI = [2.14, 2.23]), with a 7.86% decrease in decision time in the observable condition ($t(1304) = -$
518 2.68 , $p = 0.007$, $b = -0.08$, 95% CI = [-0.14, -0.02], BF = 2.16) (Exp. 3, H1.2b).

519 In Study 2 (checking the impact of helping or punishing on a target), results were less clear-cut. As
520 predicted, participants were significantly more likely to check the impact of helping when their

521 decision process was hidden (83%) compared to observable (78%) ($b = -0.34$, 95% CI = [-0.62, -0.06],
522 $p = 0.016$, BF = 3.86; OR = 0.71, 95% CI = [0.54, 0.94]). Nevertheless, when only those who correctly
523 responded to at least 7 out of 8 comprehension questions were included (hereafter referred to as
524 participants with excellent comprehension), results were in the same direction but no longer
525 significant (81% checked the impact when their decision process was hidden compared to 77% when
526 it was observable; OR = 0.78 [0.51, 1.19]; $b = -0.24$, 95% CI = [-0.66, 0.18], $p = 0.26$, BF = 0.67) (Exp. 4,
527 H5.1).

528 We expected that participants would be more likely to check the impact of punishment when their
529 decision process was observed (Exp. 5, H5.2), but this prediction was not supported by the data (OR =
530 0.84 [0.66, 1.07]; $b = -0.17$, 95% CI = [-0.41, 0.07], $p = .16$, BF = 0.60). Indeed, when only participants
531 with excellent comprehension were included, we found the opposite: 75% of Players A checked the
532 impact of punishing in the decision process hidden condition, whereas only 64% did so in the
533 decision process observable condition (OR = 0.62, 95% CI = [0.42, 0.89]; $b = -0.49$, 95% CI = [-0.86, -
534 0.11], $p = 0.011$, BF = 5.89).

535 **Preregistered Primary Hypotheses: The influence of deliberation over help and** 536 **punishment on perceived trustworthiness**

537 Next, we explored how helping and punishment decisions were interpreted by observers. We
538 expected uncalculated help and punishment in the context of personal cost deliberation, as well as
539 uncalculated help and calculated punishment in the context of target impact deliberation, to confer
540 reputational benefits. Specifically, we expected observers to send helpers and punishers a higher
541 percentage of their endowment in those situations, which we interpret as higher trust. See Figure 4
542 for a visualisation of results.

543 Contrary to predictions, we found no statistically significant difference in the proportion of their
544 endowment that observers sent to helpers who did not check the personal cost of helping ($M =$
545 63.40% , $SD = 34.24$) than to helpers who checked the cost ($M = 60.57\%$, $SD = 33.33$) ($t(1304) = -1.52$,
546 $p = 0.13$, $b = -2.83$, 95% CI = [-6.50, 0.82], BF = 6.00) (Exp. 1, H2.1). Yet, while the preregistered
547 frequentist statistics do not support H2.1 when all participants were included in the analysis, the
548 preregistered Bayesian analysis, with a Bayes Factor > 3 indicates support for H2.1. Importantly,
549 when only participants with excellent comprehension were included in the analysis, we found that
550 observers sent a significantly higher proportion of their endowment to helpers who did not check the
551 personal cost of helping ($M = 68.27\%$, $SD = 33.89$) than to helpers who checked the cost ($M =$
552 62.65% , $SD = 33.73$) ($t(610) = -2.06$, $p = 0.04$, $b = -5.62$, 95% CI = [-10.99, -0.26], BF = 8.97).

553 Our predictions that observers would send more money to punishers who made uncalculating
554 decisions (when considering personal costs) were not supported. If anything, observers entrusted a
555 higher proportion of their endowment to punishers who checked the personal cost of punishment
556 ($M = 51.49\%$, $SD = 36.26$) than to those who did not ($M = 48.76\%$, $SD = 37.62$) (Exp. 2, H2.2a).
557 However, this difference was statistically non-significant ($t(1304) = 1.33$, $p = 0.18$, $b = 2.73$, 95% CI =
558 [-1.28, 6.74], BF = 4.78). When calculating behaviour was operationalised in terms of decision time,
559 observers sent more to relatively slow (more calculating) punishers ($M = 49.17\%$, $SD = 34.75$) than to
560 relatively fast punishers ($M = 47.40\%$, $SD = 36.80$) (H2.2b). Again, this result was not statistically
561 significant ($t(1304) = -0.90$, $p = 0.37$, $b = -1.78$, 95% CI = [-5.66, 2.11], BF = 0.39).

562 We expected that helpers who did not check the impact of helping behaviour would be trusted more
563 by observers. Although observers sent a higher percentage of their endowment to helpers who did

564 not check the impact ($M = 63.12\%$, $SD = 32.39$) than to those who did ($M = 61.53\%$, $SD = 31.68$) (Exp.
565 4, H6.1), this result was statistically non-significant ($t(1304) = -0.90$, $p = 0.37$, $b = -1.60$, 95% CI = [-
566 5.07, 1.88], $BF = 1.53$).

567 Another unsupported prediction was that impact-checking punishers would be trusted more by
568 observers. Although observers did send more of their endowment to punishers who checked the
569 impact of punishing ($M = 48.45\%$, $SD = 34.93$) than to punishers who did not ($M = 45.54\%$, $SD =$
570 35.01) (Exp. 5, H6.2), this difference was also statistically non-significant ($t(1304) = 1.50$, $p = 0.13$, $b =$
571 2.91 , 95% CI = [-0.89, 6.71], $BF = 4.31$). While the preregistered frequentist statistics do not support
572 H6.2, the preregistered Bayesian analysis, with a Bayes Factor > 3 indicates support for H6.2.

573 **Exploratory Preregistered Hypotheses: The influence of deliberation over help and** 574 **punishment on trustworthiness**

575 Next, we asked whether calculated/uncalculated help and punishment decisions reliably signalled
576 trustworthiness (Figure 5). We expected uncalculating helpers in both the personal cost (Exp. 1) and
577 impact checking context (Exp. 4) to be more trustworthy than calculating helpers. Indeed, helpers
578 who did not check the personal cost of helping returned significantly more of the endowment they
579 were sent by observers ($M = 48.74\%$, $SD = 19.76$) than helpers who did check the personal cost ($M =$
580 43.54% , $SD = 19.34$) ($t(1099) = -3.85$, $p < 0.001$, $b = -5.21$, 95% CI = [-7.86, -2.55], $BF = 97.33$) (Exp. 1,
581 H14.1). Similarly, helpers who did not check the impact of helping ($M = 48.23\%$, $SD = 19.04$) returned
582 a higher percentage in the Trust Game than helpers who checked the impact of helping ($M = 45.40\%$,
583 $SD = 18.57$), but this effect was statistically non-significant ($t(1138) = -1.96$, $p = 0.05$, $b = -2.83$, 95% CI
584 = [-5.67, 0.004], $BF = 0.44$) (Exp. 4, H19.1).

585 We expected that punishers who made an uncalculating versus calculating decision in the context of
586 personal cost (Exp. 2 and 3) would be more trustworthy. Conversely, for impact consideration, we
587 predicted that punishers who made calculating decisions would be more trustworthy than punishers
588 who made uncalculating decisions (Exp. 5). Our results did not support these predictions. Although
589 punishers who did not check the personal cost of punishing returned more of the entrusted
590 endowment ($M = 46.11\%$, $SD = 15.48$) than punishers who did check the personal cost ($M = 43.78\%$,
591 $SD = 20.01$), this difference was not statistically significant ($t(506) = -1.04$, $p = 0.30$, $b = -2.33$, 95% CI =
592 [-6.76, 2.09], $BF = 0.17$) (Exp. 2, H14.2a). Conversely, when uncalculating decisions were
593 operationalised as decision time, punishers who made slower (more calculating) decisions returned a
594 slightly higher percentage than those who made faster (uncalculating) punishing decisions ($t(513) =$
595 0.66 , $p = 0.51$, $b = 0.98$, 95% CI = [-1.94, 3.90], $BF = 0.14$) (Exp. 3, H14.2b). This difference was not
596 significant. Punishers who did not check the impact of punishing returned a lower percentage ($M =$
597 38.48% , $SD = 22.53$) than punishers who did check the impact of punishing on the target ($M =$
598 40.03% , $SD = 19.67$), $t(408) = 0.58$, $p = 0.56$, $b = 1.55$, 95% CI = [-3.68, 6.78], $BF = 0.12$ (Exp. 5, H19.2).
599 Although directionally in line with predictions, this difference too was non-significant.

600 It must be noted, that hypotheses H14.2a, H14.2b and H19.2 did not meet power requirements,
601 therefore making their results suggestive, pending confirmation in future research. However, their
602 Bayes Factor values indicate support for the null hypotheses (see Supplementary Table 1 under
603 "Supplementary Notes 1" in the Supplementary Information for sensitivity analyses).

604 **Preregistered Primary Hypotheses: Trust and trustworthiness across experiments**

605 We expected uncalculated decision-making to differentially influence trust and trustworthiness
606 across the experiments (Exp. 1-3) of Study 1 (personal cost). Firstly, for punishment, we predicted
607 that deliberation would have a stronger influence on trust and trustworthiness when calculating
608 behaviour was operationalised as cost-checking (Exp. 2) compared to slow decision time (Exp. 3). In
609 addition, we expected deliberation to have a stronger effect on trust and trustworthiness in the
610 context of helping compared to punishing (Exp. 1 vs Exp. 2).

611 However, the effect of calculated versus uncalculated punishment on trust was not stronger for cost
612 checking than decision time ($t(2608) = 0.33, p = .74, b = 0.95, 95\% \text{ CI} = [-1.99, 3.88], \text{BF} = 0.10$) (H3).
613 The same was true for non-punishment ($t(2608) = -0.19, p = 0.85, b = -0.55, 95\% \text{ CI} = [-3.45, 2.35], \text{BF}$
614 $= 0.75$) (H10). Similarly, the effect of calculated versus uncalculated punishment on trustworthiness
615 was not stronger for cost checking than decision time, $t(1020) = -1.27, p = 0.21, b = -3.44, 95\% \text{ CI} = [-$
616 $8.78, 1.90], \text{BF} = 0.24$ (H16). Again, the same was true for non-punishment: $t(1587) = -1.72, p = 0.09,$
617 $b = -3.80, 95\% \text{ CI} = [-8.13, 0.54], \text{BF} = 0.33$ (H17).

618 Observers trusted helpers significantly more than they trusted punishers ($t(2608) = 7.47, p < 0.001, b$
619 $= 14.64, 95\% \text{ CI} = [10.74, 18.54]$). Moreover, trust was significantly influenced by the interaction
620 between behaviour (helping versus punishing) and decision process (calculating versus uncalculating)
621 ($t(2608) = -2.01, p = 0.04, b = -5.56, 95\% \text{ CI} = [-8.29, -2.83], \text{BF} = 140.10$) (H4). Specifically,
622 uncalculating punishers were trusted the least ($M = 48.76\%, SD = 37.62$), followed by calculating
623 punishers ($M = 51.49\%, SD = 36.26$), calculating helpers ($M = 60.57\%, SD = 33.33$), and uncalculating
624 helpers ($M = 63.40\%, SD = 34.24$). Uncalculating helpers were trusted significantly more than
625 uncalculating punishers ($t(2608) = 7.47, p < 0.001, b = 14.64$) and calculating helpers were trusted
626 significantly more than calculating punishers ($t(2608) = 4.64, p < 0.001, b = 9.08$).

627 This interaction was no longer significant when excluding those who failed more than one
628 comprehension check ($t(1162) = -1.54, p = 0.12, b = -6.63, 95\% \text{ CI} = [-10.66, -2.60], \text{BF} = 20.22$). While
629 the preregistered frequentist statistics no longer support H4 when only participants with excellent
630 comprehension were included, the preregistered Bayesian analysis, with a Bayes Factor > 3 indicates
631 support for H4. Observers still trusted helpers significantly more than punishers ($t(1605) = 4.28, p <$
632 $0.001, b = 13.00, 95\% \text{ CI} = [6.87, 19.13]$). Specifically, observers trusted uncalculating punishers the
633 least ($M = 55.27\%, SD = 40.76$), increasing their levels of trust for calculating punishers ($M = 56.28\%,$
634 $SD = 38.23$), calculating helpers ($M = 62.65\%, SD = 33.73$), and uncalculating helpers ($M = 68.27\%, SD$
635 $= 33.89$).

636 There was also no evidence to suggest that the effect of calculated versus uncalculated decision-
637 making on trustworthiness is stronger in helping compared to punishing contexts ($t(1605) = -1.09, p =$
638 $0.28, b = -2.88, 95\% \text{ CI} = [-8.04, 2.29], \text{BF} = 0.17$) (H18), and there was no significant difference in the
639 trustworthiness of helpers and punishers ($t(1605) = 1.11, p = 0.27, b = 2.63, 95\% \text{ CI} = [-2.01, 7.27]$).

640 **Preregistered Secondary Hypotheses: The influence of deliberation over decisions not to** 641 **help or punish on perceived trustworthiness**

642 Moreover, we had diverging expectations for how uncalculating decisions would be perceived when
643 those decisions result in inaction rather than helping or punishing. We predicted that observers
644 would send more to calculating than uncalculating non-helpers/non-punishers in Experiments 1, 4
645 and 5, but more to uncalculating than calculating non-punishers when personal cost is being
646 considered (Exp. 2 and Exp. 3). However, none of these analyses were statistically significant.

647 Directionally in line with predictions, observers sent more of their endowment to non-helpers who
648 checked the cost of helping ($M = 29.75\%$, $SD = 33.70$) than to non-helpers who did not check the cost
649 of helping ($M = 28.81\%$, $SD = 34.38$) ($t(1304) = 0.50$, $p = 0.61$, $b = 0.95$, 95% CI = [-2.75, 4.64], BF =
650 0.68) (Exp.1, H7.1). Conversely, and again in line with predictions, in Experiment 2 (H7.2a) observers
651 sent directionally less of their endowment to non-punishers who checked the personal cost of
652 punishing ($M = 50.31\%$, $SD = 35.38$) than to non-punishers who did not check the cost of punishing
653 ($M = 51.58\%$, $SD = 37.83$) ($t(1304) = -0.63$, $p = 0.53$, $b = -1.27$, 95% CI = [-5.25, 2.70], BF = 0.80).
654 However, Experiment 3 (H7.2b) found that observers sent more of their endowment to relatively
655 slow (calculating) non-punishers ($M = 49.36\%$, $SD = 37.18$) than to relatively fast (uncalculating) non-
656 punishers ($M = 48.64\%$, $SD = 35.97$) ($t(1304) = 0.36$, $p = 0.72$, $b = 0.72$, 95% CI = [-3.25, 4.69], BF =
657 0.26). In Experiment 4 (H11.1) observers were again in line with predictions and sent more of their
658 endowment to non-helpers who checked the impact of helping ($M = 32.43\%$, $SD = 34.22$) than to
659 non-helpers who did not ($M = 29.71\%$, $SD = 33.17$) ($t(1304) = 1.46$, $p = 0.14$, $b = 2.73$, 95% CI = [-0.93,
660 6.38], BF = 5.05). In Experiment 5 (H11.2) observers sent similar amounts of their endowment to non-
661 punishers who checked the impact of punishing ($M = 52.48\%$, $SD = 34.09$) and to non-punishers who
662 did not ($M = 52.85\%$, $SD = 34.68$) ($t(1304) = -0.19$, $p = 0.85$, $b = -0.37$, 95% CI = [-4.10, 3.37], BF =
663 0.48).

664 **Exploratory Preregistered Hypotheses: The influence of deliberation over decisions not to** 665 **help or punish on trustworthiness**

666 We also had diverging expectations for how uncalculating decisions would be associated with the
667 actual trustworthiness of non-helpers and non-punishers. Specifically, we predicted that calculating
668 non-punishers in the context of impact checking (Exp. 5) and calculating non-helpers in both the
669 context of impact (Exp. 4) and cost checking (Exp. 1) would return more than uncalculating non-
670 helpers/non-punishers. In contrast, we expected uncalculating non-punishers to return more than
671 calculating non-punishers in context of personal cost deliberation (Exp. 2 & Exp. 3). All returning
672 decisions for non-punishers and non-helpers were directionally in line with predictions.

673 In Experiment 1 (H15.1), non-helpers who checked the personal cost of helping returned more of
674 their endowment ($M = 21.29\%$, $SD = 22.95$) than non-helpers who did not check the cost ($M =$
675 15.69% , $SD = 22.88$) ($t(208) = 1.67$, $p = 0.10$, $b = 5.60$, 95% CI = [-1.00, 12.19], BF = 0.56). However,
676 this difference was non-significant. As predicted, in Experiment 2 (H15.2a), non-punishers who did
677 not check the cost of punishing ($M = 42.51\%$, $SD = 21.65$) returned significantly more of their
678 endowment than non-punishers who checked the cost ($M = 38.51\%$, $SD = 22.42$) ($t(799) = -2.57$, $p =$
679 0.01 , $b = -4.0$, 95% CI = [-7.06, -0.94], BF = 1.98). In Experiment 3 (H15.2b) uncalculating (faster) non-
680 punishers again returned more of their endowment than calculating (slower) non-punishers, but this
681 was not significant ($t(787) = -0.29$, $p = 0.77$, $b = -0.43$, 95% CI = [-3.33, 2.47], BF = 0.01). In Experiment
682 4 (H20.1) non-helpers who checked the impact of helping ($M = 25.44\%$, $SD = 23.84$) returned
683 significantly more of their endowment than non-helpers who did not check the impact ($M = 17.24\%$,
684 $SD = 26.56$) ($t(169) = 1.99$, $p = 0.48$, $b = 8.21$, 95% CI = [0.06, 16.35], BF = 1.02). However, the
685 difference (calculating non-helper: 23.40% ($SD = 24.55$), uncalculating non-helper: 16.00% ($SD =$
686 24.11)) was no longer statistically significant when only those with excellent comprehension were
687 included ($t(58) = 1.11$, $p = 0.27$, $b = 7.40$, 95% CI = [-5.98, 20.78], BF = 0.44). Finally, in Experiment 5
688 (H20.2) both non-punishers who checked the impact of punishing ($M = 38.39\%$, $SD = 22.74$) and non-
689 punishers who did not check the impact ($M = 38.26\%$, $SD = 24.68$) returned around 38% of their
690 endowment ($t(822) = 0.08$, $p = 0.94$, $b = 0.13$, 95% CI = [-3.17, 3.43], BF = 0.08).

691 It must be noted that power requirements were not met for hypotheses H15.1 (Exp. 1) and H20.1
692 (Exp. 4), making those results suggestive, pending confirmation in future research.

693 **Preregistered Secondary Hypotheses: The influence of deliberation on perceived and**
694 **actual trustworthiness when decisions result in helping or punishing versus inaction**

695 Lastly, for all experiments we predicted that the effect of uncalculating behaviour on trust and
696 trustworthiness would be larger for action than inaction, meaning that deliberation would more
697 strongly influence sending and returning decisions when Player A decided to help/punish compared
698 to when Player A decided *not* to help/punish.

699 However, for sending decisions this was not the case in Experiment 1 (H8.1; $t(2608) = -1.43$, $p = 0.15$,
700 $b = -3.78$, 95% CI = [-8.99, 1.42], BF = 1.16), Experiment 2 (H8.2a; $t(2608) = 1.39$, $p = 0.17$, $b = 4.0$, 95%
701 CI = [-1.65, 9.64], BF = 1.17), Experiment 3 (H8.2b; $t(2608) = -0.88$, $p = 0.38$, $b = -2.50$, 95% CI = [-5.40,
702 0.40], BF = 0.28), or Experiment 5 (H12.2; $t(2608) = 1.21$, $p = 0.23$, $b = 3.28$, 95% CI = [-2.05, 8.60], BF
703 = 0.79). Yet, when only participants with excellent comprehension were included, there was a
704 significant interaction between deliberation and helping decision in Experiment 4 (H12.1; $t(1060) = -$
705 2.06 , $p = 0.04$, $b = -7.97$, 95% CI = [-15.54, -0.40], BF = 5.89). Specifically, observers entrusted
706 uncalculating non-helpers with only 21.73% ($SD = 29.24$) of their endowment, and calculating non-
707 helpers with 26.77% ($SD = 31.86$) of their endowment. Helpers were sent more than twice as much:
708 calculating helpers were entrusted with 61.20% ($SD = 32.02$) and uncalculating helpers received the
709 most with 64.14% ($SD = 32.72$). Hereby, the differences between uncalculating helpers versus
710 uncalculating non-helpers ($t(1060) = 15.53$, $p < 0.001$, $b = 42.41$) and calculating helpers versus
711 calculating non-helpers ($t(1060) = 12.61$, $p < 0.001$, $b = 34.44$) were statistically significant. Moreover,
712 in Experiment 1 there was a main effect for helping ($t(2608) = 18.43$, $p < 0.001$, $b = 34.59$, 95% CI =
713 [30.87, 38.32]), as observers entrusted more than twice as much to helpers than to non-helpers, and
714 in Experiment 5 observers sent significantly less to punishers than to non-punishers ($t(2608) = -3.81$,
715 $p = 0.0001$, $b = -7.30$, 95% CI = [-11.09, -3.52]).

716 Furthermore, we found no evidence to suggest that deliberation had a larger effect on actual
717 trustworthiness for punishers compared to non-punishers, as the interaction effects were non-
718 significant in Experiment 2 (H9.2a; $t(1305) = 0.58$, $p = 0.56$, $b = 1.67$, 95% CI = [-3.96, 7.29], BF = 0.14),
719 Experiment 3 (H9.2b; $t(1301) = 0.88$, $p = 0.38$, $b = 1.85$, 95% CI = [-2.29, 6.00], BF = 0.3) and
720 Experiment 5 (H13.2; $t(1302) = 0.42$, $p = 0.67$, $b = 1.42$, 95% CI = [-5.17, 8.01], BF = 0.13). For
721 participants with excellent comprehension there were, however, main effects for both punishing
722 ($t(615) = 2.13$, $p = 0.03$, $b = 6.33$) and checking ($t(615) = -2.30$, $p = 0.02$, $b = -4.49$) in Experiment 2,
723 with punishers and uncalculating decision makers returning significantly more than non-punishers
724 and calculating decision makers.

725 In Experiment 1 (H9.1) the effect of uncalculating decision making on trustworthiness was
726 significantly larger when Players A decided to help compared to when Players A decided not to help
727 ($t(1307) = -3.34$, $p < 0.001$, $b = -10.81$, 95% CI = [-17.16, -4.45], BF = 18.29). In line with predictions,
728 non-helpers who did not check the personal cost of helping were the least trustworthy, returning
729 only an average of 15.69% ($SD = 22.88$), whilst cost-checking non-helpers returned 21.29% ($SD =$
730 22.95). Cost-checking helpers were substantially more trustworthy, returning an average of 43.54%
731 ($SD = 19.34$), whilst helpers who did not check the personal cost returned the most, with an average
732 of 48.74% ($SD = 19.34$). Post hoc tests on the estimated marginal means, accounting for multiple
733 comparisons with the multivariate t-test (mvt) adjustment, revealed significant differences between

734 uncalculating helpers and uncalculating non-helpers ($t(1307) = 12.40, p < 0.001, b = 33.05$),
735 calculating helpers and calculating non-helpers ($t(1307) = 12.11, p < 0.001, b = 22.25$) as well as
736 calculating helpers and uncalculating helpers ($t(1307) = -3.74, p < 0.001, b = -5.21$), but not for
737 calculating non-helpers and uncalculating non-helpers ($t(1307) = 1.92, p = 0.18, b = 5.60$).

738 Conversely, and against predictions, in Experiment 4 (H13.1) the effect of uncalculating decision
739 making on trustworthiness was significantly *smaller* when Player A decided to help compared to
740 when Player A decided not to help ($t(1307) = -3.07, p = 0.002, b = 11.04, 95\% \text{ CI} = [-18.10, -3.98], \text{BF} =$
741 8.59). Nevertheless, in line with predictions, non-helpers who did not check the impact of helping
742 were the least trustworthy, returning only an average of 17.24% ($SD = 26.56$) of the endowment
743 observers entrusted them with, whilst calculating non-helpers returned an average of 25.44% ($SD =$
744 23.84), calculating helpers an average of 45.40% ($SD = 18.57$) and uncalculating helpers an average of
745 48.23% ($SD = 19.04$). Hereby, the differences between uncalculating helpers and uncalculating non-
746 helpers ($t(1307) = 10.13, p < 0.001, b = 31.00$), calculating helpers and calculating non-helpers
747 ($t(1307) = 10.54, p < 0.001, b = 19.96$), as well as uncalculating non-helpers and calculating non-
748 helpers ($t(1307) = 2.51, p = 0.04, b = 8.21$) were statistically significant.

749 **Exploratory Unregistered Analyses: The influence of deliberation on trust and** 750 **trustworthiness for non-helpers versus non-punishers**

751 To provide a comprehensive perspective, a final unregistered analysis tested whether the effect of
752 uncalculating behaviour on trust and trustworthiness differs for non-punishers compared to non-
753 helpers. For trust there was no interaction between deliberation and behaviour ($t(2608) = 0.80, p =$
754 $0.42, b = 2.22, 95\% \text{ CI} = [-0.64, 5.08], \text{BF} = 0.98$), nor a significant main effect for deliberation ($t(2608)$
755 $= -0.65, p = 0.52, b = -1.27, 95\% \text{ CI} = [-3.39, 0.85]$). However, observers sent significantly more of
756 their endowment to non-punishers than to non-helpers ($t(2608) = -11.64, p < 0.001, b = -22.77, 95\%$
757 $\text{CI} = [-26.70, -18.85]$).

758 Furthermore, non-punishers returned a significantly higher proportion in the Trust Game than non-
759 helpers did ($t(1007) = -9.43, p < 0.001, b = -26.82, 95\% \text{ CI} = [-32.40, -21.24]$), and non-actors who
760 checked the cost returned significantly less than those who made an uncalculated decision not to
761 help/punish ($t(1007) = -2.54, p = 0.01, b = -4.00, 95\% \text{ CI} = [-7.09, -0.91]$). There was also a significant
762 interaction between deliberation and experiment ($t(1007) = 2.67, p = 0.008, b = 9.60, 95\% \text{ CI} = [2.53,$
763 $16.66], \text{BF} = 3.45$), with uncalculated non-helpers returning the least ($M = 15.69\%, SD = 22.88$),
764 followed by calculated non-helpers ($M = 21.29\%, SD = 22.95$), calculated non-punishers ($M = 38.51\%,$
765 $SD = 22.42$), and uncalculated non-punishers ($M = 42.51\%, SD = 21.65$). However, when only
766 participants with excellent comprehension were included, there no longer was a significant
767 interaction ($t(486) = 1.51, p = 0.13, b = 7.89, 95\% \text{ CI} = [-2.40, 18.17], \text{BF} = 0.46$) although the average
768 percentages returned remained similar (uncalculated non-helpers: ($M = 13.74\%, SD = 25.00$),
769 calculated non-helpers: ($M = 17.14\%, SD = 20.61$), calculated non-punishers: ($M = 36.27\%, SD =$
770 21.50), uncalculated non-punishers: ($M = 40.76\%, SD = 20.06$)).

771 Discussion

772 Previous work²⁶ has shown that helping behaviour that is performed in a reflexive or uncalculating
773 manner can yield reputation benefits, since observers infer that these actions reflect genuinely
774 prosocial motives, rather than stemming from rational calculation of costs and benefits. Accordingly,
775 uncalculated help signals trustworthiness and people are more likely to behave in an uncalculated
776 way when they are observed²⁶. Over five experiments, we replicate this study and extend it by
777 examining whether uncalculated punishment also leads to reputation improvements. In a further
778 extension of previous work, we also ask whether punishers and helpers deliberate over the *impact to*
779 *the target* (rather than the personal cost to themselves) and how such ‘impact deliberation’ is viewed
780 by bystanders. In Study 1 (personal cost deliberation) we expected both uncalculated help and
781 punishment to signal trustworthiness. In Study 2 (target impact deliberation) we expected
782 uncalculated help to signal trustworthiness. Conversely, we expected *calculated* punishment to signal
783 trustworthiness. As punishment inflicts harm on another, we expected that people would observe an
784 implicit moral directive to deliberate over the harm they could inflict on another individual – and that
785 individuals who inflict harm reflexively would be viewed negatively. Replicating previous results²⁶, we
786 found that uncalculated help signals trustworthiness: helpers who did not consider the personal cost
787 of helping were both more trusted and trustworthy than helpers who deliberated over the cost. Our
788 punishment results were more mixed. Although punishers were more likely to perform uncalculated
789 actions when observed, uncalculated punishment was not reliably associated with either perceptions
790 of trustworthiness or with trustworthiness itself. Only uncalculating *non*-punishers were more
791 trustworthy than calculating non-punishers. In contrast to the cost checking context, considering the
792 impact of helping had a larger impact on the trust and trustworthiness of non-helpers than helpers.
793 Lastly, we found no conclusive evidence to suggest that checking the impact of punishing influences
794 perceived or actual trustworthiness.

795 In Experiment 1, we replicated Jordan et al.’s²⁶, finding that uncalculating helpers were perceived as
796 significantly more trustworthy than calculating helpers. Uncalculated helping provides a reliable
797 signal of trustworthiness as it indicates that people are not considering the personal costs of helping
798 and that helping stems from other-regarding rather than strategic motives. As in Jordan et al.²⁶,
799 people were sensitive to these reputation benefits and were less likely to check the personal cost of
800 helping when their decision process was observed than when it was hidden (H1.1). Finally, as in
801 Jordan et al.²⁶, these reputation benefits were restricted to those who helped: deliberation had no
802 effect on trust (H7.1) or trustworthiness (H15.1) when participants decided *not* to help.

803 We similarly expected uncalculated punishers to be perceived as, and to actually be, more
804 trustworthy than those who deliberated over the personal cost of punishing (Exp. 2 & 3). We also
805 expected people to be sensitive to these reputation consequences and to be less likely to check the
806 personal cost (or to decide more quickly) when observed. These predictions were only partially
807 supported. Participants were half as likely to check the cost of punishing when their decision process
808 was observed (H1.2a) and were also significantly faster in their decision-making (H1.2b). In contrast
809 to predictions, observers directionally trusted *calculating* punishers more than uncalculating
810 punishers (while the Bayes Factor value for H2.2a indicated support for this effect, frequentist
811 statistics - which were preregistered as the primary decision criterion - did not support H2.2a or
812 H2.2b). Trustworthiness results were also mixed. Uncalculating punishers were directionally more
813 trustworthy than calculating punishers in Experiment 2 (H14.2a), but directionally less trustworthy in

814 Experiment 3 (H14.2b). Note that Bayes Factor values indicate support for the null hypotheses,
815 although power requirements were not met for H14.2a and H14.2b.

816 Whilst we expected both uncalculated help and uncalculated punishment to signal trustworthiness,
817 we had diverging predictions around decisions *not* to act. Decision conflict over whether to
818 help/punish could stem from self-interested considerations of whether to pay a cost. But unlike
819 helping, punishment decision conflict could also stem from concerns about inflicting harm on the
820 target. As participants initially believed punishing could potentially be free, we expected
821 uncalculating non-punishers to be perceived as harm averse. Conversely, calculating decisions not to
822 punish would indicate a selfish decision (the personal cost of punishing being too high). Support for
823 these predictions was mixed. As expected, uncalculating non-punishers were more trustworthy than
824 calculating non-punishers (though effects were only significant for Exp. 2, H15.2a and not Exp. 3,
825 H15.2b). Perceived trustworthiness was not reliably affected, as observers directionally trusted
826 uncalculating non-punishers more in Experiment 2 (H7.2a; the Bayes Factor value indicates null
827 findings are inconclusive) but directionally trusted calculating non-punishers more in Experiment 3
828 (H7.2b; the Bayes Factor value indicates support for the null hypothesis).

829 Uncalculated punishment does not therefore seem to be perceived as a signal of trustworthiness –
830 and uncalculated punishers were not more trustworthy. As deliberative decisions are often
831 considered to be wise⁴⁸⁻⁵⁰, uncalculated punishment might conceivably signal reduced competence⁵¹,
832 which could have affected perceived trustworthiness. While possible, this explanation is unlikely as
833 the same ought to have been true for the helping context in Experiment 1. Alternatively, it is possible
834 that the signalling effect of uncalculating punishment was too small to have been captured by the
835 present work. However, several of the Bayes Factor values for null results in Experiments 2 and 3
836 were less than 0.33, supporting the absence of an effect as opposed to a need for more data.
837 Moreover, we frequently found the directional opposite of our predictions, especially when
838 deliberation was operationalised as decision time.

839 Helping may enhance reputation more than punishment because, even though third-party
840 punishment is often viewed as a morally justified form of harm, people may still be unsure about
841 those who engage in it^{21,23,31}. Observers may therefore be unsure whether to trust punishers over
842 non-punishers in the first place. Prior research has found that non-punishers can sometimes be
843 trusted as much as punishers^{24,30}, and occasionally third-party punishers are even trusted less than
844 non-punishers^{32,52-55}. We found no significant difference in the perceived trustworthiness of
845 punishers and non-punishers. Nevertheless, trustworthiness did vary. When restricting our sample to
846 participants with excellent comprehension, punishers returned significantly more than non-
847 punishers.

848 Perhaps punishment needs to be seen as the ‘right thing to do’ for the decision process to matter as
849 a signal. This can be difficult, as punishment – unlike helping – is morally bad when undeserved, and
850 there are also questions around legitimacy in that a fellow participant in an economic game may not
851 be seen as an appropriate person to intervene^{54,56}. Further, it has been argued that defection in
852 economic games can be considered ‘fair game’, making the punishment of it less justified^{31,57}. This is
853 additionally important because the appropriateness of non-action decreases for more serious
854 infractions⁵⁸. Furthermore, even third-party punishers can be perceived as spiteful or competitive
855 rather than prosocial, particularly when punishment is excessive³¹. However, punishment in this
856 study is unlikely to be seen as excessive: punishing still leaves the defector with £0.15, the amount
857 that they would have received had they themselves acted fairly. It is also unlikely that punishers are

858 perceived as being competitive (aiming to increase their own payoffs relative to the defector's^{5,59-60})
859 because, whether participants choose to punish or not, they always end with only one-third of the
860 amount the defector receives (when participants punish, they finish the game with £0.05, whilst the
861 defector finishes with £0.15, and when participants do not punish, they finish with £0.10, whilst the
862 defector finishes with £0.30).

863 It should be noted, that although helpers were indeed trusted significantly more than punishers (H4),
864 they were not actually more trustworthy than punishers (H18). Our additional unregistered analyses
865 further showed that non-punishers were significantly more trustworthy than non-helpers, which is in
866 line with prior research²⁴. Just as help is a more reliable signal of trustworthiness than punishment,
867 not helping is a more reliable signal of *untrustworthiness*²¹.

868 We expected deliberating about the *impact* of help / punishment to reflect negatively on helpers but
869 positively on punishers, on the basis that considering both the prosocial aspect of third-party
870 punishment and the negative consequences to the defector may be perceived as the 'right' thing to
871 do when making punishing decisions. Helping on the other hand should be perceived as the 'right'
872 thing to do, whether it helps a little or a lot. Helping without checking the impact to targets was
873 therefore expected to signal trustworthiness. Conversely, deciding not to help without this
874 consideration was predicted to result in particularly negative evaluations.

875 As predicted, uncalculating helpers were both more trusted (H6.1) and trustworthy (H19.1) than
876 calculating helpers – but only directionally and with inconclusive Bayes Factors. Further, non-helpers
877 who made calculating decisions were indeed more trustworthy than uncalculating non-helpers, but
878 only significantly so when all participants were included (note: power requirements were not met for
879 H20.1). Non-helpers were also directionally more trusted when they made calculating compared to
880 uncalculating decisions (H11.1; the effect was statistically non-significant (the primary decision
881 criterion) but supported by the Bayes Factor value). Counter to predictions, impact checking had a
882 larger effect on the trust (H12.1) and trustworthiness (H13.1) of non-helpers compared to helpers.
883 We do not find this result especially surprising. Indeed, we also argued that helping, regardless of the
884 impact, is what matters the most for perceptions of (and actual) trustworthiness. It is possible that
885 impact consideration matters more for those who decided *not* to help, as those who consider the
886 impact at least considered helping, whereas those who do not consider the impact and do not help
887 may be perceived as unwilling to help no matter how much another may benefit from it.

888 As predicted, considerations of target impact generally produced the opposite pattern in the
889 punishing context. Counter to predictions, participants with excellent comprehension were
890 significantly *less* likely to check the impact of punishing when they were observed than when their
891 decision process was hidden (H5.2). Directionally, punishers who checked the impact to the target
892 were perceived as (H6.2) and directionally were (H19.2) more trustworthy than uncalculating
893 punishers (though effects were statistically non-significant). No conclusive results were found for
894 non-punishers (H11.2 & H20.2). Bayes Factor values of non-significant analyses mostly indicated
895 support for the null hypotheses (except H11.2 and H12.2, which were inconclusive, and H6.2 which
896 indicates support for the alternative hypothesis).

897

898 **Limitations**

899 The findings presented here should be interpreted within the context of certain limitations,
900 particularly regarding ecological validity. The experimental design was highly abstract and therefore
901 may not have fully captured the complexities of real-world decision-making processes related to
902 trust and trustworthiness. Future studies could enhance ecological validity by employing scenarios or
903 tasks more directly applicable to everyday situations, thereby potentially yielding more
904 representative results. Additionally, past research emphasizes the importance of motive attributions
905 in shaping evaluations of helpers and punishers^{21,26,33}, which the present study did not explicitly
906 explore. Understanding the motives observers attribute to actors, as well as eliciting self-reported
907 motives behind actors' calculated versus uncalculated decisions, could provide further insight into
908 the mechanisms underlying third-party punishment and the extent to which it can be interpreted as
909 a prosocial act. This could also help to differentiate between punishment and helping as signals of
910 trustworthiness.

911 Lastly, there may be some concerns inherent to the Trust Game itself. We used this game to measure
912 attitudes towards punishers as previous work has shown that punishment increases trustworthiness
913 whereas results on whether punishers are 'liked' or rewarded for their actions are more
914 mixed^{16,23,28,67}. Nevertheless, it is important to acknowledge that the Trust Game does not fully
915 disambiguate between trust and other underlying mechanisms. While decisions in the Trust Game
916 can be partly attributable to risk attitudes⁶⁴, decision patterns in Trust and Risk games differ⁶⁶.
917 Differences in responses to help and punishment in the present study also indicate that observers did
918 not make decisions based purely on risk preferences (if risk preferences were the key driver behind
919 Trust Game decisions, we would not expect to observe any differences in help / punishment
920 conditions).

921 **References**

- 922 1. Fehr, E., & Fischbacher, U. (2003). The nature of human altruism. *Nature*, 425(6960), 785-
923 791. <https://doi.org/10.1038/nature02043>
- 924 2. Penner, L. A., Dovidio, J. F., Piliavin, J. A., & Schroeder, D. A. (2005). Prosocial behavior:
925 Multilevel perspectives. *Annual Review of Psychology*, 56, 365-392.
926 <https://doi.org/10.1146/annurev.psych.56.091103.070141>
- 927 3. West, S. A., Griffin, A. S., & Gardner, A. (2007). Evolutionary explanations for cooperation.
928 *Current biology*, 17(16), R661-R672. <https://doi.org/10.1016/j.cub.2007.06.004>
- 929 4. Clutton-Brock, T. H., & Parker, G. A. (1995). Punishment in animal societies. *Nature*,
930 373(6511), 209-216. <https://doi.org/10.1038/373209a0>
- 931 5. Raihani, N. J., & Bshary, R. (2019). Punishment: one tool, many uses. *Evolutionary Human*
932 *Sciences*, 1, e12. <https://doi.org/10.1017/ehs.2019.12>
- 933 6. Balafoutas, L., Grechenig, K., & Nikiforakis, N. (2014). Third-party punishment and counter-
934 punishment in one-shot interactions. *Economics letters*, 122(2), 308-310.
935 <https://doi.org/10.1016/j.econlet.2013.11.028>
- 936 7. Balliet, D., Mulder, L. B., & Van Lange, P. A. (2011). Reward, punishment, and cooperation: a
937 meta-analysis. *Psychological bulletin*, 137(4), 594. <https://doi.org/10.1037/a0023489>
- 938 8. Boyd, R., Gintis, H., Bowles, S., & Richerson, P. J. (2003). The evolution of altruistic
939 punishment. *Proceedings of the National Academy of Sciences*, 100(6), 3531-3535.
940 <https://doi.org/10.1073/pnas.0630443100>
- 941 9. Boyd, R., & Richerson, P. J. (1992). Punishment allows the evolution of cooperation (or
942 anything else) in sizable groups. *Ethology and sociobiology*, 13(3), 171-195.
943 [https://doi.org/10.1016/0162-3095\(92\)90032-Y](https://doi.org/10.1016/0162-3095(92)90032-Y)
- 944 10. Charness, G., Cobo-Reyes, R., & Jiménez, N. (2008). An investment game with third-party
945 intervention. *Journal of Economic Behavior & Organization*, 68(1), 18-28.
946 <https://doi.org/10.1016/j.ejpoleco.2017.08.001>
- 947 11. Feinberg, M., Willer, R., & Schultz, M. (2014). Gossip and ostracism promote cooperation in
948 groups. *Psychological science*, 25(3), 656-664. <https://doi.org/10.1177/0956797613510184>
- 949 12. Henrich, J., & Boyd, R. (2001). Why people punish defectors: Weak conformist transmission
950 can stabilize costly enforcement of norms in cooperative dilemmas. *Journal of theoretical*
951 *biology*, 208(1), 79-89. <https://doi.org/10.1006/jtbi.2000.2202>
- 952 13. Mathew, S., & Boyd, R. (2011). Punishment sustains large-scale cooperation in prestate
953 warfare. *Proceedings of the National Academy of Sciences*, 108(28), 11375-11380.
954 <https://doi.org/10.1073/pnas.1105604108>
- 955 14. Dreber, A., Rand, D. G., Fudenberg, D., & Nowak, M. A. (2008). Winners don't
956 punish. *Nature*, 452(7185), 348-351. <https://doi.org/10.1038/nature06723>
- 957 15. Nikiforakis, N. (2008). Punishment and counter-punishment in public good games: Can we
958 really govern ourselves? *Journal of Public Economics*, 92(1), 91-112.
959 <http://doi.org/10.1016/j.jpubeco.2007.04.008>
- 960 16. Barclay, P. (2006). Reputational benefits for altruistic punishment. *Evolution and Human*
961 *Behavior*, 27(5), 325-344. <https://doi.org/10.1016/j.evolhumbehav.2006.01.003>
- 962 17. dos Santos, M. D., Rankin, D. J., & Wedekind, C. (2011). The evolution of punishment through
963 reputation. *Proceedings of the Royal Society B: Biological Sciences*, 278(1704), 371-377.
964 <https://doi.org/10.1098/rspb.2010.1275>
- 965 18. dos Santos, M. D., Rankin, D. J., & Wedekind, C. (2013). Human cooperation based on
966 punishment reputation. *Evolution*, 67(8), 2446-2450. <https://doi.org/10.1111/evo.12108>

- 967 19. dos Santos, M., & Wedekind, C. (2015). Reputation based on punishment rather than
968 generosity allows for evolution of cooperation in sizable groups. *Evolution and Human*
969 *Behavior*, 36(1), 59-64. <https://doi.org/10.1016/j.evolhumbehav.2014.09.001>
- 970 20. Kurzban, R., DeScioli, P., & O'Brien, E. (2007). Audience effects on moralistic punishment.
971 *Evolution and Human behavior*, 28(2), 75-84.
972 <https://doi.org/10.1016/j.evolhumbehav.2006.06.001>
- 973 21. Raihani, N. J., & Bshary, R. (2015a). The reputation of punishers. *Trends in ecology &*
974 *evolution*, 30(2), 98-103. <https://doi.org/10.1016/j.tree.2014.12.003>
- 975 22. Hilbe, C., & Traulsen, A. (2012). Emergence of responsible sanctions without second order
976 free riders, antisocial punishment or spite. *Scientific reports*, 2(1), 1-4.
977 <https://doi.org/10.1038/srep00458>
- 978 23. Raihani, N. J., & Bshary, R. (2015b). Third-party punishers are rewarded—but third-party
979 helpers even more so. *Evolution*, 69(4), 993-1003. <https://doi.org/10.1111/evo.12637>
- 980 24. Batistoni, T., Barclay, P., & Raihani, N. J. (2022). Third-party punishers do not compete to be
981 chosen as partners in an experimental game. *Proceedings of the Royal Society B*, 289(1966),
982 20211773. <https://doi.org/10.1098/rspb.2021.1773>
- 983 25. Dhaliwal, N. A., Skarlicki, D. P., Hoegg, J., & Daniels, M. A. (2020). Consequentialist motives
984 for punishment signal trustworthiness. *Journal of Business Ethics*, 1-16.
985 <https://doi.org/10.1007/s10551-020-04664-5>
- 986 26. Jordan, J. J., Hoffman, M., Nowak, M. A., & Rand, D. G. (2016). Uncalculating cooperation is
987 used to signal trustworthiness. *Proceedings of the National Academy of Sciences*, 113(31),
988 8658-8663. <https://doi.org/10.1073/pnas.1601280113>
- 989 27. Nelissen, R. M. (2008). The price you pay: Cost-dependent reputation effects of altruistic
990 punishment. *Evolution and Human Behavior*, 29(4), 242-248.
991 <https://doi.org/10.1016/j.evolhumbehav.2008.01.001>
- 992 28. Horita, Y. (2010). Punishers may be chosen as providers but not as recipients. *Letters on*
993 *Evolutionary Behavioral Science*, 1(1), 6-9. <https://doi.org/10.5178/lebs.2010.2>
- 994 29. Jordan, J. J., & Rand, D. G. (2017). Third-party punishment as a costly signal of high
995 continuation probabilities in repeated games. *Journal of theoretical biology*, 421, 189-202.
996 <http://doi.org/10.1016/j.jtbi.2017.04.004>
- 997 30. Przepiorka, W., & Liebe, U. (2016). Generosity is a sign of trustworthiness—the punishment
998 of selfishness is not. *Evolution and human behavior*, 37(4), 255-262.
999 <https://doi.org/10.1016/j.evolhumbehav.2015.12.003>
- 1000 31. Dhaliwal, N. A., Patil, I., & Cushman, F. (2021). Reputational and cooperative benefits of
1001 third-party compensation. *Organizational Behavior and Human Decision Processes*, 164, 27-
1002 51. <https://doi.org/10.1016/j.obhdp.2021.01.003>
- 1003 32. Heffner, J., & FeldmanHall, O. (2019). Why we don't always punish: Preferences for non-
1004 punitive responses to moral violations. *Scientific reports*, 9(1), 1-13.
1005 <https://doi.org/10.1038/s41598-019-49680-2>
- 1006 33. Raihani, N. J., & Power, E. A. (2021). No good deed goes unpunished: the social costs of
1007 prosocial behaviour. *Evolutionary Human Sciences*, 3, e40.
1008 <https://doi.org/10.1017/ehs.2021.35>
- 1009 34. Nishi, A., Christakis, N. A., Evans, A. M., O'Malley, A. J., & Rand, D. G. (2016). Social
1010 environment shapes the speed of cooperation. *Scientific reports*, 6(1), 1-10.
1011 <https://doi.org/10.1038/srep29622>
- 1012 35. Critcher, C. R., Inbar, Y., & Pizarro, D. A. (2013). How quick decisions illuminate moral
1013 character. *Social Psychological and Personality Science*, 4(3), 308-315.
1014 <https://doi.org/10.1177/1948550612457688>

- 1015 36. Evans, A. M., & Van De Calseyde, P. P. (2017). The effects of observed decision time on
1016 expectations of extremity and cooperation. *Journal of Experimental Social Psychology*, *68*,
1017 50-59. <http://doi.org/10.1016/j.jesp.2016.05.009>
- 1018 37. Van de Calseyde, P. P., Keren, G., & Zeelenberg, M. (2014). Decision time as information in
1019 judgment and choice. *Organizational Behavior and Human Decision Processes*, *125*(2), 113-
1020 122. <https://doi.org/10.1016/j.obhdp.2014.07.001>
- 1021 38. R Core Team (2020). R: A language and environment for statistical computing. R Foundation
1022 for Statistical Computing, Vienna, Austria. URL <https://www.R-project.org/>.
- 1023 39. Champely, S. (2020). Pwr: basic functions for power analysis (R package version 1.3–
1024 0)[Computer software]. The Comprehensive R Archive Network. Retrieved from
1025 <https://CRAN.R-project.org/package=pwr>
- 1026 40. Cohen, J. E. (1988). *Statistical Power Analysis for the Behavioral Sciences*. Hillsdale, NJ:
1027 Lawrence Erlbaum Associates, Inc.
- 1028 41. Selya, A. S., Rose, J. S., Dierker, L. C., Hedeker, D., & Mermelstein, R. J. (2012). A practical
1029 guide to calculating Cohen's f^2 , a measure of local effect size, from PROC MIXED. *Frontiers in*
1030 *psychology*, *3*, 111. <https://doi.org/10.3389/fpsyg.2012.00111>
- 1031 42. Zeileis, A., & Hothorn, T. (2002). Diagnostic Checking in Regression Relationships. *R News*,
1032 *2*(3), 7–10. Retrieved from <https://CRAN.R-project.org/doc/Rnews/>
- 1033 43. Morey, R., & Rouder, J. (2022). *BayesFactor: Computation of Bayes Factors for Common*
1034 *Designs* (Version 0.9.12-4.4) [Software]. Retrieved from [https://CRAN.R-](https://CRAN.R-project.org/package=BayesFactor)
1035 [project.org/package=BayesFactor](https://CRAN.R-project.org/package=BayesFactor)
- 1036 44. Bürkner, P.-C. (2017). Brms: An R Package for Bayesian Multilevel Models Using Stan. *Journal*
1037 *of Statistical Software*, *80*(1), 1-28. <https://doi.org/10.18637/jss.v080.i01>
- 1038 45. Singmann, H., Kellen, D., Cox, G. E., Chandramouli, S. H., Davis-Stober, C. P., Dunn, J. C., ... &
1039 Shiffrin, R. M. (2023). Statistics in the service of science: Don't let the tail wag the dog.
1040 *Computational Brain & Behavior*, *6*(1), 64-83. <https://doi.org/10.1007/s42113-022-00129-2>
- 1041 46. Jeffreys, H. (1939). *Theory of Probability*. Oxford: Clarendon Press.
- 1042 47. Lee, M. D., & Wagenmakers, E. J. (2014). *Bayesian cognitive modeling: A practical course*.
1043 Cambridge University Press. <https://doi.org/10.1017/CBO9781139087759>
- 1044 48. Bazerman, M. H., & Chugh, D. (2006). Decisions without blinders. *Harvard business*
1045 *review*, *84*(1), 88. PMID: 16447372.
- 1046 49. Grossmann, I., Brienza, J. P., & Bobocel, D. R. (2017). Wise deliberation sustains
1047 cooperation. *Nature Human Behaviour*, *1*(3), 0061. [https://doi.org/10.1038/s41562-017-](https://doi.org/10.1038/s41562-017-0061)
1048 [0061](https://doi.org/10.1038/s41562-017-0061)
- 1049 50. Pinker, S. (2022). *Rationality: What it is, why it seems scarce, why it matters*. Penguin.
- 1050 51. Jordan, J. J., & Kteily, N. S. (2023). How reputation does (and does not) drive people to punish
1051 without looking. *Proceedings of the National Academy of Sciences*, *120*(28), e2302475120.
1052 <https://doi.org/10.1073/pnas.2302475120>
- 1053 52. Eriksson, K., Andersson, P. A., & Strimling, P. (2016). Moderators of the disapproval of peer
1054 punishment. *Group Processes & Intergroup Relations*, *19*(2), 152-168.
1055 <https://doi.org/10.1177/1368430215583519>
- 1056 53. Rai, T. S. (2022). Material benefits crowd out moralistic punishment. *Psychological Science*,
1057 *33*(5), 789-797. <https://doi.org/10.1177/09567976211054786>
- 1058 54. Strimling, P., & Eriksson, K. (2014). Regulating the regulation: Norms about punishment.
1059 *Reward and punishment in social dilemmas*, 52-69.
1060 <https://doi.org/10.1093/acprof:oso/9780199300730.003.0004>
- 1061 55. Sun, B., Jin, L., Yue, G., & Ren, Z. (2022). Is a punisher always trustworthy? In-group
1062 punishment reduces trust. *Current Psychology*, *1-11*. [https://doi.org/10.1007/s12144-022-](https://doi.org/10.1007/s12144-022-03395-2)
1063 [03395-2](https://doi.org/10.1007/s12144-022-03395-2)

- 1064 56. Eriksson, K., Andersson, P. A., & Strimling, P. (2017). When is it appropriate to reprimand a
1065 norm violation? The roles of anger, behavioral consequences, violation severity, and social
1066 distance. *Judgment and decision making*, 12(4), 396-407.
1067 <https://doi.org/10.1017/S1930297500006264>
- 1068 57. Salcedo García, J. C. (2020). Moralistic punishment signaling as a function of proportionality.
1069 [Doctoral Dissertation, Universidad de los Andes] <http://hdl.handle.net/1992/48393>
- 1070 58. Eriksson, K., Strimling, P., Gelfand, M., Wu, J., Abernathy, J., Akotia, C. S., ... & Van Lange, P.
1071 A. (2021). Perceptions of the appropriate response to norm violation in 57 societies. *Nature*
1072 *communications*, 12(1), 1481. <https://doi.org/10.1038/s41467-021-21602-9>
- 1073 59. Bone, J. E., & Raihani, N. J. (2015). Human punishment is motivated by both a desire for
1074 revenge and a desire for equality. *Evolution and Human Behavior*, 36(4), 323-330.
1075 <https://doi.org/10.1016/j.evolhumbehav.2015.02.002>
- 1076 60. Price, M. E., Cosmides, L., & Tooby, J. (2002). Punitive sentiment as an anti-free rider
1077 psychological device. *Evolution and Human Behavior*, 23(3), 203-231.
1078 [https://doi.org/10.1016/S1090-5138\(01\)00093-9](https://doi.org/10.1016/S1090-5138(01)00093-9)
- 1079 61. Tosi, J., & Warmke, B. (2016). Moral grandstanding. *Philosophy & Public Affairs*, 44(3), 197-
1080 217. <https://doi.org/10.1111/papa.12075>
- 1081 62. Haidt, J., & Rose-Stockwell, T. (2019). The dark psychology of social networks. *The Atlantic*, 6-
1082 60. Retrieved from [https://www.theatlantic.com/magazine/archive/2019/12/social-media-](https://www.theatlantic.com/magazine/archive/2019/12/social-media-democracy/600763/)
1083 [democracy/600763/](https://www.theatlantic.com/magazine/archive/2019/12/social-media-democracy/600763/)
- 1084 63. Sunstein, C. R. (2020). Lapidation and apology. *U. Chi. Legal F.*, 295. Retrieved from
1085 <https://chicagounbound.uchicago.edu/uclf/vol2020/iss1/12>
- 1086 64. Chetty, R., Hofmeyr, A., Kincaid, H., & Monroe, B. (2021). The trust game does not (only)
1087 measure trust: The risk-trust confound revisited. *Journal of Behavioral and Experimental*
1088 *Economics*, 90, 101520. <https://doi.org/10.1016/j.socec.2020.101520>
- 1089 65. Brühlhart, M., & Usunier, J. C. (2012). Does the trust game measure trust?. *Economics Letters*,
1090 115(1), 20-23. <https://doi.org/10.1016/j.econlet.2011.11.039>
- 1091 66. Houser, D., Schunk, D., & Winter, J. (2010). Distinguishing trust from risk: An anatomy of the
1092 investment game. *Journal of economic behavior & organization*, 74(1-2), 72-81.
1093 <https://doi.org/10.1016/j.jebo.2010.01.002>
- 1094 67. Kiyonari, T., & Barclay, P. (2008). Cooperation in social dilemmas: Free riding may be
1095 thwarted by second-order reward rather than by punishment. *Journal of Personality and*
1096 *Social Psychology*, 95(4), 826-842. <https://doi.org/10.1037/a0011381>

1097

1098 **Acknowledgements**

1099 The authors received no specific funding for this work.

1100

1101 **Author contributions**

1102 N.R. and N.E. developed the study concept. N.E. designed the study with revisions from N.R. and
1103 conducted data collection and analysis; N.E. drafted the initial manuscript; N.E. and N.R. revised and
1104 reviewed the manuscript and approved the final manuscript for submission.

1105

1106 **Competing interests**

1107 The authors declare no competing interest.

Question	Hypothesis	Sampling plan (e.g. power analysis)	Analysis Plan	Interpretation given to different outcomes
Primary Hypotheses				
Q1. Are uncalculated decisions around the personal cost of helping / punishing used as a signal of trustworthiness?	<p>Helping: H1.1) Participants will be significantly less likely to check the cost of helping in the decision process observable condition than in the decision process hidden condition.</p> <p>Punishing: H1.2a) Participants will be significantly less likely to check the cost of punishment in the decision process observable condition than in the decision process hidden condition.</p> <p>H1.2b) Participants will make significantly faster punishing decisions in the decision process observable condition than in decision the process hidden condition.</p>	<p>Please refer to the Sampling plan in Methods for detail.</p> <p>H1.1 & H1.2a) The sample size for this model will be N = 1306 (653 Players A per condition in Experiment 1/2).</p> <p>H1.2b) The sample size for this model will be N = 1306 (653 Players A per condition in Experiment 3).</p>	<p>H1.1 & H1.2a) We will run a logistic regression with checking decision (0 = did not check the cost, 1 = checked the cost) as a function of decision process observability (0 = process hidden, 1 = process observable).</p> <p>H1.2b) We will run a linear regression, predicting decision time as a function of decision process observability (0 = process hidden, 1 = process observable). If the amount of time spent deciding whether to punish is highly skewed, punishing decision time will be natural log transformed.</p>	<p>H1.1 & H1.2a) A significant negative coefficient for observability (0 = decision process hidden, 1 = decision process observable) will be interpreted as evidence that participants are less likely to check the personal cost of helping/punishing when their decision process is observable compared to hidden (and therefore, that they are more likely to act uncalculatingly when their decision process can be observed). Otherwise, there is no evidence for H1.1/H1.2a.</p> <p>H1.2b) A significant negative coefficient for observability (0 = decision process hidden, 1 = decision process observable) will be interpreted as evidence that participants make faster punishing decisions (i.e., act uncalculatingly) when their decision process is observable compared to hidden. Otherwise, there is no evidence for H1.2b.</p>
Q2. Are uncalculated decisions around	<p>Helping: H2.1) Observers will send significantly more of their</p>	<p>Please refer to the Sampling plan in Methods for detail.</p>	<p>Analyses will be restricted to the observable condition because Players B can only condition their</p>	<p>H2.1 & H2.2a) A significant negative coefficient for cost checking (0 = did not check the cost, 1 = checked the</p>

<p>the personal cost of helping / punishing perceived as a signal of trustworthiness?</p>	<p>endowment to helpers who did not check the personal cost of helping than to helpers who checked the cost.</p> <p>Punishing: H2.2a) Observers will send significantly more of their endowment to punishers who did not check the personal cost of punishing than to punishers who checked the cost.</p> <p>H2.2b) Observers will send significantly more of their endowment to punishers who made relatively fast (vs relatively slow) decisions to punish.</p>	<p>H2.1 & H2.2a) The sample size for this model will be N = 653 (Players B in the observable condition in Experiment 1/2).</p> <p>H2.2b) The sample size for this model will be N = 653 (Players B in the observable condition in Experiment 3).</p>	<p>trust on Player A decision processes in this condition. As Players B make two sending decisions (based on the two possible decisions made by Player A during the first stage), each sending decision will be treated as an observation and robust SEs will be clustered on observer ID to account for the non-independence of repeated observations from the same participant. The endowment sent will be transformed from pence sent to percentage of endowment sent for ease of interpretation.</p> <p>H2.1 & H2.2a) We will run a linear regression predicting the percentage of endowment sent by Players B as a function of helpers/punishers cost-checking decisions (0 = did not check the personal cost, 1 = checked the cost).</p> <p>H2.2b) We will run a linear regression predicting the percentage of endowment Players B sent to punishers in the observable condition as a</p>	<p>cost) will be interpreted as evidence that observers send a higher proportion of their endowment to helpers/punishers who did not check the cost of helping/ punishing compared to helpers/punishers who checked the cost of helping/punishing (i.e., that observers trust uncalculating helpers/punishers more than calculating helpers/punishers). Otherwise, there is no evidence for H2.1/H2.2a.</p> <p>H2.2b) A significant positive coefficient of decision speed (0 = relatively slow, 1 = relatively fast) will be interpreted as evidence that observers send a higher proportion of their endowment to punishers who decided relatively quickly compared to punishers who decided relatively slowly (i.e., that observers trust uncalculating punishers more than calculating punishers). Otherwise, there is no evidence for H2.2b.</p>
---	---	--	--	---

			function of decision time (0 = relatively slow, 1 = relatively fast).	
Q3. Does the operationalisation of uncalculating behaviour differentially influence the perceived trustworthiness of punishers in the context of personal cost?	H3) Observers will send significantly less of their endowment to punishers who check the personal cost of punishing than to punishers who take a long time to decide to punish.	Please refer to the Sampling plan in Methods for detail. H3) The sample size for this model will be N = 1306 (Players B in the observable condition in Experiment 2 and 3).	H3) We will run a linear regression predicting the percentage of endowment sent by Players B as a function of experiment (decision time vs cost checking) and deliberation (uncalculated vs calculated), as well as the interaction between experiment and deliberation.	H3) A significant negative coefficient for the interaction between experiment (0 = decision time, 1 = cost checking) and deliberation (0 = uncalculated, 1 = calculated) will be interpreted as evidence that the effect of calculated vs uncalculated punishment on trust is stronger for cost checking than decision time. Otherwise, there is no evidence for H3.
Q4. Do uncalculated helping and punishing decisions differentially influence perceived trustworthiness in the context of personal cost?	H4) Observers will send significantly less of their endowment to helpers who check the personal cost than to punishers who check the personal cost.	Please refer to the Sampling plan in Methods for detail. H4) The sample size for this model will be N = 1306 (Players B in the observable condition in Experiment 1 and 2).	H4) We will run a linear regression predicting the percentage of endowment sent by Players B as a function of behaviour (punishing vs helping) and deliberation (uncalculated vs calculated), as well as the interaction between behaviour and deliberation.	H4) A significant negative coefficient for the interaction between behaviour (0 = punishing, 1 = helping) and deliberation (0 = uncalculated, 1 = calculated) will be interpreted as evidence that the effect of calculated vs uncalculated decisions on trust is stronger for helping than punishing. Otherwise, there is no evidence for H4.
Q5. Are uncalculated decisions around target impact used as a signal of trustworthiness?	Helping: H5.1) Participants will be significantly less likely to check the impact of helping in the decision process observable condition than in the decision process hidden condition.	Please refer to the Sampling plan in Methods for detail. H5.1 & H5.2) The sample size for this model will be N = 1306 (653 Players A	H5.1 & H5.2) We will run a logistic regression with checking decision (0 = did not check the impact, 1 = checked the impact) as a function of decision process observability (0 = process hidden, 1 = process observable).	H5.1) A significant negative coefficient for observability (0 = decision process hidden, 1 = decision process observable) will be interpreted as evidence that participants are less likely to check the impact of helping when their decision process is observable compared to hidden (and

	<p>Punishing: H5.2) Participants will be significantly more likely to check the impact of punishment in the decision process observable condition than in the decision process hidden condition.</p>	<p>per condition in Experiment 4/5).</p>		<p>therefore, that they are more likely to act uncalculatingly when their decision process can be observed). Otherwise, there is no evidence for H5.1.</p> <p>H5.2) A significant positive coefficient for observability (0 = decision process hidden, 1 = decision process observable) will be interpreted as evidence that participants are more likely to check the impact of punishing when their decision process is observable compared to hidden (and therefore, that they are more likely to act calculatingly when their decision process can be observed). Otherwise, there is no evidence for H5.2.</p>
<p>Q6. Are uncalculated decisions around target impact perceived as a signal of trustworthiness?</p>	<p>Helping: H6.1) Observers will send significantly more of their endowment to helpers who did not check the impact of helping on targets than to helpers who checked the impact.</p> <p>Punishing: H6.2) Observers will send significantly more of their endowment to punishers who checked the impact of punishing on targets than to punishers who did not check the impact.</p>	<p>Please refer to the Sampling plan in Methods for detail.</p> <p>H6.1 & H6.2) The sample size for this model will be N = 653 (Players B in the observable condition in Experiment 4/5).</p>	<p>Analyses will be restricted to the observable condition because Players B can only condition their trust on Player A decision processes in this condition. As Players B make two sending decisions (based on the two possible decisions made by Player A during the first stage), each sending decision will be treated as an observation and robust SEs will be clustered on observer ID to account for the non-independence of repeated observations from the same participant. The endowment sent</p>	<p>H6.1) A significant negative coefficient for impact checking (0 = did not check the impact, 1 = checked the impact) will be interpreted as evidence that observers send a higher proportion of their endowment to helpers who did not check the impact of helping compared to helpers who checked the impact of helping (i.e., that observers trust uncalculating helpers more than calculating helpers). Otherwise, there is no evidence for H6.1.</p> <p>H6.2) A significant positive coefficient for impact checking (0 = did not check the impact, 1 = checked the impact)</p>

			<p>will be transformed from pence sent to percentage of endowment sent for ease of interpretation.</p> <p>H6.1 & H6.2) We will run a linear regression predicting the percentage of endowment sent by Players B as a function of helpers/punishers impact-checking decisions (0 = did not check the impact, 1 = checked the impact).</p>	<p>will be interpreted as evidence that observers send a higher proportion of their endowment to punishers who checked the impact of punishing compared to punishers who did not check the impact of punishing (i.e., that observers trust calculating punishers more than uncalculating punishers). Otherwise, there is no evidence for H6.2.</p>
Secondary Hypotheses				
<p>Q7. Do uncalculated decisions around personal cost affect the perceived trustworthiness of non-helpers / non-punishers?</p>	<p>Helping: H7.1) Observers will send significantly more of their endowment to non-helpers who checked the cost of helping than to non-helpers who did not check the cost of helping.</p> <p>Punishing: H7.2a) Observers will send significantly more of their endowment to non-punishers who did not check the cost of punishing than to non-punishers who checked the cost of punishing.</p> <p>H7.2b) Observers will send significantly more of their</p>	<p>Please refer to the Sampling plan in Methods for detail.</p> <p>H7.1 & H7.2a) The sample size for this model will be N = 653 (Players B in the observable condition in Experiment 1/2).</p> <p>H7.2b) The sample size for this model will be N = 653 (Players B in the observable condition in Experiment 3).</p>	<p>Analysis will be restricted to the observable condition because Players B can only condition their trust on Player A decision processes in this condition. As each observer makes two sending decisions, we will cluster robust SEs on observer ID.</p> <p>H7.1 & H7.2a) We will run a linear regression predicting the percentage of endowment Players B sent to non-helpers/ non-punishers as a function of cost checking decisions.</p> <p>H7.2b) We will run a linear regression predicting the percentage of endowment</p>	<p>H7.1) A significant positive coefficient for cost checking (0 = did not check the cost, 1 = checked the cost) will be interpreted as evidence that observers send more of their endowment to non-helpers who checked the cost of helping compared to non-helpers who did not check the cost of helping (i.e., observers trust calculating non-helpers more than uncalculating non-helpers). Otherwise, there is no evidence for H7.1.</p> <p>H7.2a) A significant negative coefficient for cost checking (0 = did not check the cost, 1 = checked the cost) will be interpreted as evidence that observers send more of their endowment to non-punishers who did</p>

	<p>endowment to relatively fast non-punishers than to relatively slow non-punishers.</p>		<p>Players B sent to non-punishers as a function of decision speed.</p>	<p>not check the cost of punishing compared to non-punishers who checked the cost (i.e., observers trust uncalculating non-punishers more than calculating non-punishers). Otherwise, there is no evidence for H7.2a.</p> <p>H7.2b) A significant positive coefficient for decision speed (0 = relatively slow, 1 = relatively fast) will be interpreted as evidence that observers send more of their endowment to non-punishers who take little time in their decision not to punish compared to those who take a long time to decide not to punish (i.e., observers trust uncalculating non-punishers more than calculating non-punishers). Otherwise, there is no evidence for H7.2b.</p>
<p>Q8. Do uncalculated decisions around personal cost have a stronger effect on the perceived trustworthiness of helpers/ punishers than non-helpers / non-punishers?</p>	<p>Helping: H8.1) Players B will send significantly less of their endowment to helpers who checked the cost of helping than to non-helpers who checked the cost of helping.</p> <p>Punishing: H8.2a) Players B will send significantly less of their endowment to punishers who</p>	<p>Please refer to the Sampling plan in Methods for detail.</p> <p>H8.1 & H8.2a) The sample size for this model will be N = 653 (Players B in the observable condition in Experiment 1/2).</p> <p>H8.2b) The sample size for this model will be N = 653 (Players B in the</p>	<p>H8.1 & H8.2a) We will run a linear regression predicting the percentage of endowment sent by Players B as a function of helping/punishing decision, cost checking decision, and the interaction between the two. Robust SEs will be clustered on observer ID, as Players B will make four sending decisions based on each of the four possible Player A choices.</p>	<p>H8.1 & H8.2a) A significant negative coefficient for the interaction between helping/punishing (0 = did not help/punish, 1 = helped/punished) and cost checking decisions (0 = did not check the cost, 1 = checked the cost) will be interpreted as evidence that the effect of uncalculating behaviour on trust is larger when Players A decided to help/ punish compared to when Players A decided</p>

	<p>checked the cost of punishing than to non-punishers who checked the cost of punishing.</p> <p>H8.2b) Players B will send significantly less of their endowment to relatively slow punishers than to relatively slow non-punishers.</p>	<p>observable condition in Experiment 3).</p>	<p>H8.2b) We will run a linear regression predicting the percentage of endowment sent by Players B as a function of punishing decision, decision speed and the interaction between the two. Robust SEs will be clustered on participant ID, accounting for repeated observations (four per participant).</p>	<p>not to help/punish. Otherwise, there is no evidence for H8.1/H8.2a.</p> <p>H8.2b) A significant negative interaction between punishing decision (0 = did not punish, 1 = punished) and decision speed (0 = relatively slow, 1 = relatively fast) will be interpreted as evidence that the effect of uncalculating behaviour on trust is larger when Players A decided to punish compared to when Player As decided not to punish. Otherwise, there is no evidence for H8.2b.</p>
<p>Q9. Do uncalculated decisions around personal cost have a stronger effect on the trustworthiness of helpers/ punishers than non-helpers / non-punishers?</p>	<p>Helping: H9.1) Helpers who checked the cost of helping will return significantly less of their endowment than non-helpers who checked the cost of helping.</p> <p>Punishing: H9.2a) Punishers who checked the cost of punishing will return significantly less of their endowment than non-punishers who checked the cost of punishing.</p> <p>H9.2b) Fast punishers will return significantly less of their endowment than fast non-punishers.</p>	<p>Please refer to the Sampling plan in Methods for detail.</p> <p>H9.1 & H9.2a) The sample size for this model will be N = 1306 (653 Players A per condition in Experiment 1/2).</p> <p>H9.2b) The sample size for this model will be N = 1306 (653 Players A per condition in Experiment 3).</p>	<p>H9.1 & H9.2a) We will run a linear regression predicting the percentage of endowment returned by Players A as a function of helping/punishing decision, cost checking decision, as well as the interaction between the two.</p> <p>H9.2b) We will run a linear regression predicting the percentage of endowment returned by Players A as a function of punishing decision, log-transformed punishing decision time, their interaction, as well as log-transformed general comprehension speed. As the analysis is correlational, we</p>	<p>H9.1 & H9.2a) A significant negative coefficient for the interaction between helping/punishing (0 = did not help/punish, 1 = helped/punished) and cost checking decisions (0 = did not check the cost, 1 = checked the cost) will be interpreted as evidence that the effect of uncalculating decision making on trustworthiness is larger when Players A decide to help/punish compared to when Players A decide not to help/punish. Otherwise, there is no evidence for H9.1/H9.2a.</p> <p>H9.2b) A significant negative interaction between punishing decision (0 = did not punish, 1 = punished) and log-transformed</p>

			wish to avoid concerns that the punishing decision time is reflective of general comprehension and reading speed rather than only of the time taken to consider whether to punish. Therefore, the natural log-transformed time spent reading the comprehension questions (i.e., the sum of time spent on the three comprehension question pages) will be included as a control for comprehension and reading speed.	decision time will be interpreted as evidence that decision time is a stronger predictor of untrustworthiness when Player A punished versus did not punish. Otherwise, there is no evidence for H9.2b.
Q10. Does the operationalisation of uncalculating behaviour differentially influence the perceived trustworthiness of non-punishers in the context of personal cost?	H10) Observers will send significantly less of their endowment to non-punishers who check the personal cost of punishing than non-punishers who take a long time to decide.	Please refer to the Sampling plan in Methods for detail. H10) The sample size for this model will be N = 1306 (Players B in the observable condition in Experiment 2 and Experiment 3).	H10) We will run a linear regression predicting the percentage of endowment sent by Players B as a function of experiment (decision time vs cost checking) and deliberation (uncalculated vs calculated), as well as the interaction between experiment and deliberation.	H10) A significant negative coefficient for the interaction between experiment (0 = decision time, 1 = cost checking) and deliberation (0 = uncalculated non-punishment, 1 = calculated non-punishment) will be interpreted as evidence that the effect of calculated vs uncalculated non-punishment on trust is stronger for cost checking than decision time. Otherwise, there is no evidence for H10.
Q11. Do uncalculated decisions around target impact affect the perceived	Helping: H11.1) Observers will send significantly more of their endowment to non-helpers who checked the impact of helping	Please refer to the Sampling plan in Methods for detail. H11.1 & H11.2) The sample size for this model	Analysis will be restricted to the observable condition because Players B can only condition their trust on Player A decision processes in this condition. As each observer makes two sending	H11.1 & H11.2) A significant positive coefficient for impact checking (0 = did not check the impact, 1 = checked the impact) will be interpreted as evidence that observers send more of their endowment to non-helpers/non-

trustworthiness of non-helpers / non-punishers?	<p>than to non-helpers who did not check the impact.</p> <p>Punishing: H11.2) Observers will send significantly more of their endowment to non-punishers who checked the impact of punishing than to non-punishers who did not check the impact.</p>	will be N = 653 (Players B in the observable condition in Experiment 4/5).	<p>decisions, we will cluster robust SEs on observer ID.</p> <p>H11.1 & H11.2) We will run a linear regression predicting the percentage of endowment Player Bs sent to non-helpers/ non-punishers as a function of impact checking decisions.</p>	punishers who checked the impact of helping/punishing compared to non-helpers/non-punishers who did not check the impact of helping/punishing (i.e., observers trust calculating non-helpers/non-punishers more than uncalculating non-helpers/non-punishers). Otherwise, there is no evidence for H11.1/H11.2.
Q12. Do uncalculated decisions around target impact have a stronger effect on the perceived trustworthiness of helpers/ punishers than non-helpers / non-punishers?	<p>Helping: H12.1) Observers will send significantly less of their endowment to helpers who checked the impact of helping than to non-helpers who checked the impact of helping.</p> <p>Punishing: H12.2) Observers will send significantly less of their endowment to punishers who checked the impact of punishing than to non-punishers who checked the impact of punishing.</p>	<p>Please refer to the Sampling plan in Methods for detail.</p> <p>H12.1 & H12.2) The sample size for this model will be N = 653 (Players B in the observable condition in Experiment 4/5).</p>	<p>H12.1 & H12.2) We will run a linear regression predicting the percentage of endowment sent by Players B as a function of helping/punishing decision, impact checking decision, and the interaction between the two. Robust SEs will be clustered on observer ID, as Players B will make sending decision based on each of the four possible Player A choices.</p>	<p>H12.1 & H12.2) A significant negative coefficient for the interaction between helping/punishing (0 = did not help/punish, 1 = helped/punished) and impact checking decisions (0 = did not check the impact, 1 = checked the impact) will be interpreted as evidence that the effect of uncalculating behaviour on trust is larger when Players A decided to help/punish compared to when Players A decided not to help/punish. Otherwise, there is no evidence for H12.1/H12.2.</p>
Q13. Do uncalculated decisions around target impact have a stronger effect on the actual trustworthiness of helpers/ punishers	<p>Helping: H13.1) Helpers who checked the impact of helping will return significantly less of their endowment than non-helpers who checked the impact of helping.</p>	<p>Please refer to the Sampling plan in Methods for detail.</p> <p>H13.1 & H13.2) The sample size for this model will be N = 1306 (653</p>	<p>H13.1 & H13.2) We will run a linear regression predicting the percentage of endowment returned by Players A as a function of helping/punishing decision, impact checking decision, as well as the interaction between the two.</p>	<p>H13.1 & H13.2) A significant negative coefficient for the interaction between helping/punishing (0 = did not help/punish, 1 = helped/punished) and impact checking decisions (0 = did not check the impact, 1 = checked the impact) will be interpreted as evidence that the effect of</p>

than non-helpers / non-punishers?	Punishing: H13.2) Punishers who checked the impact of punishing will return significantly less of their endowment than non-punishers who checked the impact of punishing.	Players A per condition in Experiment 4/5).		uncalculating decision making on trustworthiness is larger when Players A decide to help/punish compared to when Players A decide not to help/punish. Otherwise, there is no evidence for H13.1/H13.2.
Preregistered Exploratory Hypotheses				
For all exploratory hypotheses, if power requirements are not achieved, the results will be reported as suggestive, pending confirmation in future research.				
Q14. Do uncalculated decisions around the personal cost of helping / punishing predict trustworthiness?	<p>Helping: H14.1) Helpers who did not check the cost of helping will return significantly more of their endowment than helpers who checked the cost of helping.</p> <p>Punishing: H14.2a) Punishers who did not check the cost of punishing will return significantly more of their endowment than punishers who checked the cost of punishing.</p> <p>H14.2b) Punishers who made faster decisions to punish will return significantly more of their endowment than punishers who took a longer time to decide to punish.</p>	<p>Please refer to the Sampling plan in Methods for detail.</p> <p>H14.1 & H14.2a) As we do not know how many Player As will decide to help/punish, the sample size for this model will be up to N = 1306 (653 Player As per condition in Experiment 1/2).</p> <p>H14.2b) As we do not know how many Player As will decide to punish, the sample size for this model will be up to N = 1306 (653 Players A per condition in Experiment 3).</p>	<p>Here, both the observable and the hidden condition will be used, as we collect the data on Player A's decision process, even when Player B cannot observe it.</p> <p>H14.1 & H14.2a) We will run a linear regression predicting the percentage returned as a function of the helpers/punishers cost checking decision.</p> <p>H14.2b) We will run a linear regression predicting the percentage of endowment returned by punishers as a function of log-transformed decision time and log-transformed general comprehension speed.</p>	<p>H14.1 & H14.2a) A significant negative coefficient of cost checking (0 = did not check the cost, 1 = checked the cost) will be interpreted as evidence that helpers/punishers who do not check the personal cost before deciding to help/punish return more of their endowment than helpers/punishers who do check the cost (i.e., that uncalculating helpers/punishers are more trustworthy than calculating helpers/punishers). Otherwise, there is no evidence for H14.1/H14.2a.</p> <p>H14.2b) A significant negative coefficient of decision time will be interpreted as evidence that punishers who take a short time to make their decision to punish return more of their endowment than punishers who are slower in making their decision (i.e., that uncalculating punishers are</p>

				<p>more trustworthy than calculating punishers). Otherwise, there is no evidence for H14.2b.</p> <p>If power requirements are not achieved, the results will be reported as suggestive, pending confirmation in future research.</p>
<p>Q15. Do uncalculated decisions around personal cost predict the actual trustworthiness of non-helpers/ non-punishers?</p>	<p>Helping: H15.1) Non-helpers who checked the cost of helping will return significantly more of their endowment than non-helpers who did not check the cost.</p> <p>Punishing: H15.2a) Non-punishers who did not check the cost of punishing will return significantly more of their endowment than non-punishers who checked the cost.</p> <p>H15.2b) Fast deciding non-punishers will return significantly more of their endowment than slow deciding non-punishers.</p>	<p>Please refer to the Sampling plan in Methods for detail.</p> <p>H15.1 & H15.2a) As we do not know how many Player As will decide not to help/punish, the sample size for this model will be N = 1306 (653 Players A per condition in Experiment 1/2).</p> <p>H15.2b) As we do not know how many Player As will decide not to punish, the sample size for this model will be N = 1306 (653 Players A per condition in Experiment 3).</p>	<p>H15.1 & H15.2a) We will run a linear regression predicting the percentage of endowment returned by non-helpers/non-punishers as a function of cost checking behaviour.</p> <p>H15.2b) We will run a linear regression predicting the percentage of endowment returned by non-punishers as a function of log-transformed punishing decision time, controlling for log-transformed general comprehension speed.</p>	<p>H15.1) A significant positive coefficient for cost checking behaviour (0 = did not check the cost, 1 = checked the cost) will be interpreted as evidence that non-helpers who checked the cost of helping return more of their endowment compared to non-helpers who did not check the cost of helping (i.e., calculating non-helpers are more trustworthy than uncalculating non-helpers). Otherwise, there is no evidence for H15.1.</p> <p>H15.2a) A significant negative coefficient for cost checking behaviour (0 = did not check the cost, 1 = checked the cost) will be interpreted as evidence that non-punishers who did not check the cost of punishing return more of their endowment than non-punishers who checked the cost (i.e., uncalculating non-punishers are more trustworthy than calculating non-punishers). Otherwise, there is no evidence for H15.2a.</p>

				<p>H15.2b) A significant negative coefficient for log-transformed decision time will be interpreted as evidence that fast deciding non-punishers return more of their endowment than slow deciding non-punishers (i.e., uncalculating non-punishers are more trustworthy than calculating non-punishers). Otherwise, there is no evidence for H15.2b.</p> <p>If power requirements are not achieved, the results will be reported as suggestive, pending confirmation in future research.</p>
<p>Q16. Does the operationalisation of uncalculating behaviour differentially influence the actual trustworthiness of punishers in the context of personal cost?</p>	<p>H16) Punishers who check the personal cost of punishing will return significantly less of their endowment than punishers who take a long time to decide to punish.</p>	<p>Please refer to the Sampling plan in Methods for detail.</p> <p>H16) As we do not know how many Players A will decide to punish, the sample size for this model will be up to N = 2612 (653 Players A per condition in Experiments 2 and 3).</p>	<p>H16) We will run a linear regression predicting the percentage of endowment returned by Players A as a function of experiment (decision time vs cost checking) and deliberation (uncalculated vs calculated), as well as the interaction between experiment and deliberation.</p>	<p>H16) A significant negative coefficient for the interaction between experiment (0 = decision time, 1 = cost checking) and deliberation (0 = uncalculated, 1 = calculated) will be interpreted as evidence that the effect of calculated vs uncalculated punishment on trustworthiness is stronger for cost checking than decision time. Otherwise, there is no evidence for H16.</p> <p>If power requirements are not achieved, the results will be reported as suggestive, pending confirmation in future research.</p>

<p>Q17. Does the operationalisation of uncalculating behaviour differentially influence the actual trustworthiness of non-punishers in the context of personal cost?</p>	<p>H17) Non-punishers who check the personal cost of punishing will return significantly less of their endowment than non-punishers who take a long time to decide.</p>	<p>Please refer to the Sampling plan in Methods for detail.</p> <p>H17) As we do not know how many Player As will decide not to punish, the sample size for this model will be up to N = 2612 (653 Players A per condition in Experiments 2 and 32).</p>	<p>H17) We will run a linear regression predicting the percentage of endowment returned by Players A as a function of experiment (decision time vs cost checking) and deliberation (uncalculated vs calculated), as well as the interaction between experiment and deliberation.</p>	<p>H17) A significant negative coefficient for the interaction between experiment (0 = decision time, 1 = cost checking) and deliberation (0 = uncalculated non-punishment, 1 = calculated non-punishment) will be interpreted as evidence that the effect of calculated vs uncalculated non-punishment on trustworthiness is stronger for cost checking than decision time. Otherwise, there is no evidence for H17.</p> <p>If power requirements are not achieved, the results will be reported as suggestive, pending confirmation in future research.</p>
<p>Q18. Do uncalculated helping and punishing decisions differentially influence actual trustworthiness in the context of personal cost?</p>	<p>H18) Helpers who check the personal cost of helping will return significantly less of their endowment than punishers who check the personal cost of punishing.</p>	<p>Please refer to the Sampling plan in Methods for detail.</p> <p>H18) As we do not know how many Players A will decide to help/punish, the sample size for this model will be up to N = 2612 (653 Players A per condition in Experiments 1 and 2).</p>	<p>H18) We will run a linear regression predicting the percentage of endowment returned by Players A as a function of behaviour and deliberation, as well as the interaction between behaviour and deliberation.</p>	<p>H18) A significant negative coefficient for the interaction between behaviour (0 = punishing, 1 = helping) and deliberation (0 = uncalculated, 1 = calculated) will be interpreted as evidence that the effect of calculated vs uncalculated decisions on trustworthiness is stronger for helping than punishing. Otherwise, there is no evidence for H18.</p> <p>If power requirements are not achieved, the results will be reported as suggestive, pending confirmation in future research.</p>

<p>Q19. Do uncalculated decisions around target impact predict the actual trustworthiness of helpers/punishers?</p>	<p>Helping: H19.1) Helpers who did not check the impact of helping will return significantly more of their endowment than helpers who checked the impact of helping.</p> <p>Punishing: H19.2) Punishers who checked the impact of punishing will return significantly more of their endowment than punishers who did not check the impact of punishing.</p>	<p>Please refer to the Sampling plan in Methods for detail.</p> <p>H19.1 & H19.2) As we do not know how many Players A will decide to help/punish, the sample size for this model will be up to N = 1306 (653 Players A per condition in Experiment 4/5).</p>	<p>Here, both the observable and the hidden condition will be used, as we collect the data on Player A's decision process, even when Player B cannot observe it.</p> <p>H19.1 & H19.2) We will run a linear regression predicting the percentage returned as a function of the helpers/punishers impact checking decision.</p>	<p>H19.1) A significant negative coefficient of impact checking (0 = did not check the impact, 1 = checked the impact) will be interpreted as evidence that helpers who do not check the impact on a target before deciding to help return more of their endowment than helpers who do check the impact (i.e., that uncalculating helpers are more trustworthy than calculating helpers). Otherwise, there is no evidence for H19.1.</p> <p>H19.2) A significant positive coefficient of impact checking (0 = did not check the impact, 1 = checked the impact) will be interpreted as evidence that punishers who do check the impact of punishing on targets before deciding to punish return more of their endowment than punishers who do not check the impact (i.e., that calculating punishers are more trustworthy than uncalculating punishers). Otherwise, there is no evidence for H19.2.</p> <p>If power requirements are not achieved, the results will be reported as suggestive, pending confirmation in future research.</p>
---	---	---	--	---

<p>Q20. Do uncalculated decisions around target impact predict the actual trustworthiness of non-helpers/ non-punishers?</p>	<p>Helping: H20.1) Non-helpers who checked the impact of helping will return significantly more of their endowment than non-helpers who did not check the impact.</p> <p>Punishing: H20.2) Non-punishers who checked the impact of punishing will return significantly more of their endowment than non-punishers who did not check the impact.</p>	<p>Please refer to the Sampling plan in Methods for detail.</p> <p>H20.1 & H20.2) As we do not know how many Player As will decide not to punish/help, the sample size for this model will be N = 1306 (653 Players A per condition in Experiment 4/5).</p>	<p>H20.1 & H20.2) We will run a linear regression predicting the percentage of endowment returned by non-helpers/non-punishers as a function of target impact checking behaviour.</p>	<p>H20.1 & H20.2) A significant positive coefficient for impact checking behaviour (0 = did not check the impact, 1 = checked the impact) will be interpreted as evidence that non-helpers/non-punishers who checked the impact of helping/punishing return more of their endowment compared to non-helpers/non-punishers who did not check the impact (i.e., calculating non-helpers/non-punishers are more trustworthy than uncalculating non-helpers/non-punishers). Otherwise, there is no evidence for H20.1/ H20.2.</p> <p>If power requirements are not achieved, the results will be reported as suggestive, pending confirmation in future research.</p>
--	---	---	---	---

1109 **Table 1. Design Tabe**

1110

1111

Study Identification				
Study 1 Personal Cost Deliberation			Study 2 Target Impact Deliberation	
1.1 Help	1.2 Punish		2.1 Help	2.2 Punish
cost checking (E1)	1.2a cost checking (E2)	1.2b decision time (E3)	impact checking (E4)	impact checking (E5)

1112 **Table 2. Nomenclature for studies investigating whether helping and punishment decisions signal**
1113 **trustworthiness.** We recruited 1,306 Player A - Player B pairs for each of the five experiments above
1114 (i.e., each experiment contains 1,306 Player As and 1,306 Player Bs). In each experiment, half of the
1115 players were assigned to the observable decision process condition, while the other half was
1116 assigned to the hidden decision process condition.

1117

1118

1119 **Figure 1. Hypotheses for Study 1 (personal cost deliberation) and Study 2 (target impact**
1120 **deliberation).** In both studies, participants made an uncalculated or a calculated helping/punishing
1121 decision. Lower boxes indicate our expectations regarding whether a calculated or uncalculated
1122 decision is associated with a comparatively higher level of trust and trustworthiness. Green text and
1123 plus signs indicate our expectations of increased trust and trustworthiness, while red text and minus
1124 signs indicate our expectations of decreased trust and trustworthiness. This figure demonstrates our
1125 expectation that uncalculated helping signals trustworthiness in the same way across studies, while
1126 we expected uncalculated punishment to be similar to uncalculated helping when personal cost is
1127 deliberated, but to differ when considering the impact on targets.

1128 **Figure 2. Illustration of our two-stage experimental design investigating uncalculated punishment**
1129 **and helping in Studies 1 and 2, for both checking behaviour and decision speed.** In Game 1 Player A
1130 could pay a cost to punish/help another player. Player A decided (i) whether to make their decision in
1131 a calculated or uncalculated way (operationalised via their cost-checking (Experiments 1.1 & 1.2a) or
1132 impact-checking (Experiments 2.1 & 2.2) decisions, or their decision time (Experiment 1.2b), and (ii)
1133 whether to punish/help. In Game 2, Player B decided how much to send Player A (indicating how
1134 much they trust Player A), and Player A decided how much to return to Player B (indicating how
1135 trustworthy Player A is). There were two conditions in all experiments: in the process observable
1136 condition Player B could base their sending decisions both on Player A's decision process (i.e., their
1137 checking decision or decision time) as well as Player A's punishing/helping decision, whilst in the
1138 process hidden condition Player B could only make their decisions based on Player A's
1139 punishing/helping decision.

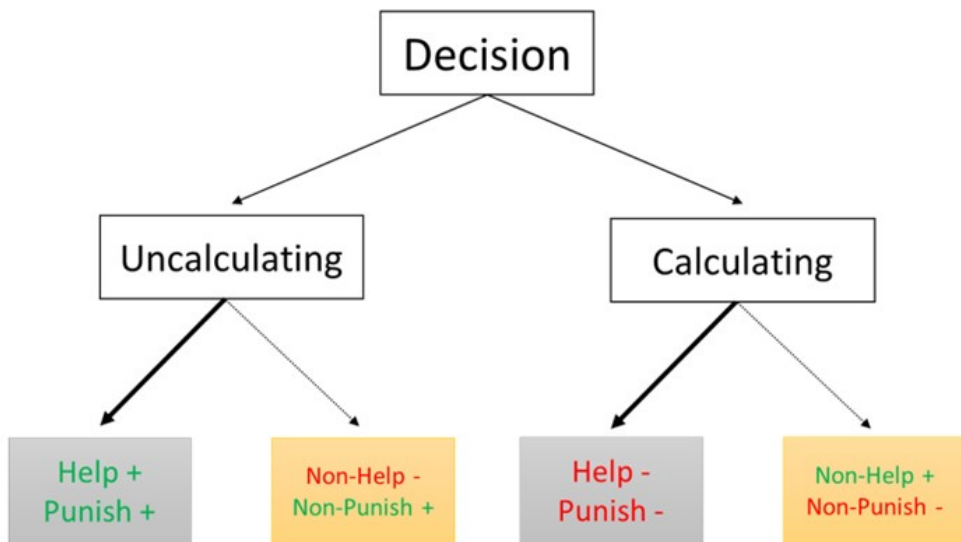
1140 **Figure 3. Decision processes (uncalculating vs calculating) across all five experiments, in the**
1141 **decision process observable and hidden conditions.** Checking (versus not checking) the personal
1142 cost or target impact (Exp. 1, 2, 4 & 5), as well as taking a long (versus a short) time to decide (Exp. 3),
1143 reflect calculated decision making. Error Bars indicate 95% CI. Due to changes in significance levels,
1144 bar charts for Exp. 4 and Exp. 5 only include participants with excellent comprehension (n = 1311 for
1145 Exp. 1, n = 1309 for Exp. 2, n = 1306 for Exp. 3, n = 534 for Exp. 4, n = 535 for Exp. 5). Differences are
1146 significant for all but Exp. 4 (help impact checking).

1147 **Figure 4. Percentage of endowment sent to uncalculating and calculating helpers (Exp. 1 & Exp. 4)**
1148 **and punishers (Exp. 2, Exp. 3 & Exp. 5) by observers.** Checking (versus not checking) the personal
1149 cost or target impact, as well as taking a long (versus a short) time to decide reflect calculated
1150 decision making. The width of the violins indicate the distribution of observations, error bars indicate
1151 95% CI, dots represent the mean. Due to changes in significance levels, Exp. 1 (help cost checking)
1152 only includes participants with excellent comprehension (n = 612 for Exp. 1, n = 1306 for Exp. 2, n =
1153 1306 for Exp. 3, n = 1306 for Exp. 4, n = 1306 for Exp. 5). Differences are significant for Exp. 1 (help
1154 cost checking).

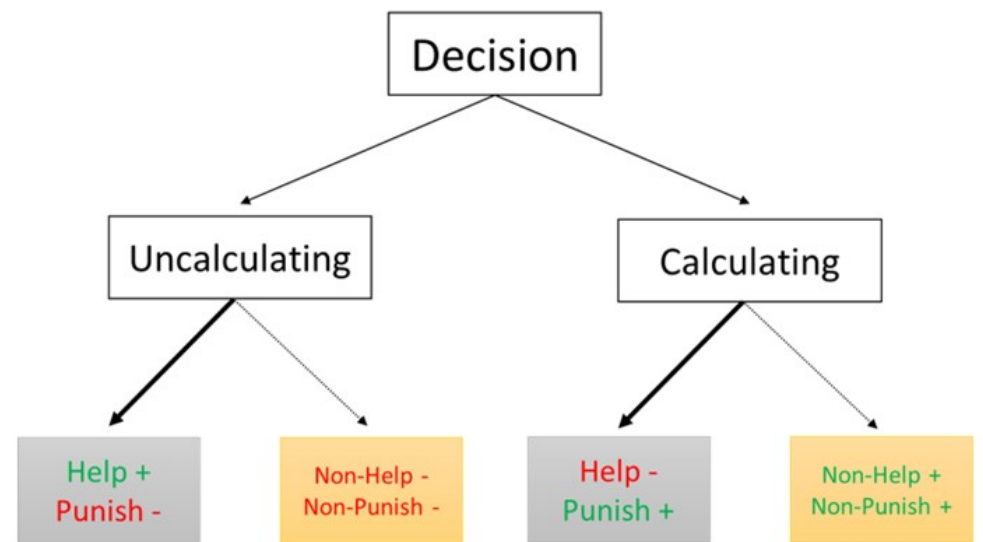
1155 **Figure 5. Percentage of endowment returned to observers by uncalculating and calculating helpers**
1156 **(Exp. 1 & Exp. 4) and punishers (Exp. 2, Exp. 3 & Exp. 5) in the Trust Game.** Checking (versus not
1157 checking) the personal cost or target impact (Exp. 1, 2, 4 & 5; Panel a), as well as taking a longer time
1158 to decide (Exp. 3; Panel b) reflect calculated decision making. In Panel A the width of the violins
1159 indicate the distribution of observations, error bars indicate 95% CI, dots represent the mean. Panel
1160 b shows a scatterplot with regression line. Differences are significant for Exp. 1 (help cost checking).
1161 Participant numbers vary across experiments (n = 1101 for Exp. 1, n = 508 for Exp. 2, n = 515 for Exp.
1162 3, n = 1140 for Exp. 4, n = 410 for Exp. 5).

1163

Study 1: Personal Cost Deliberation

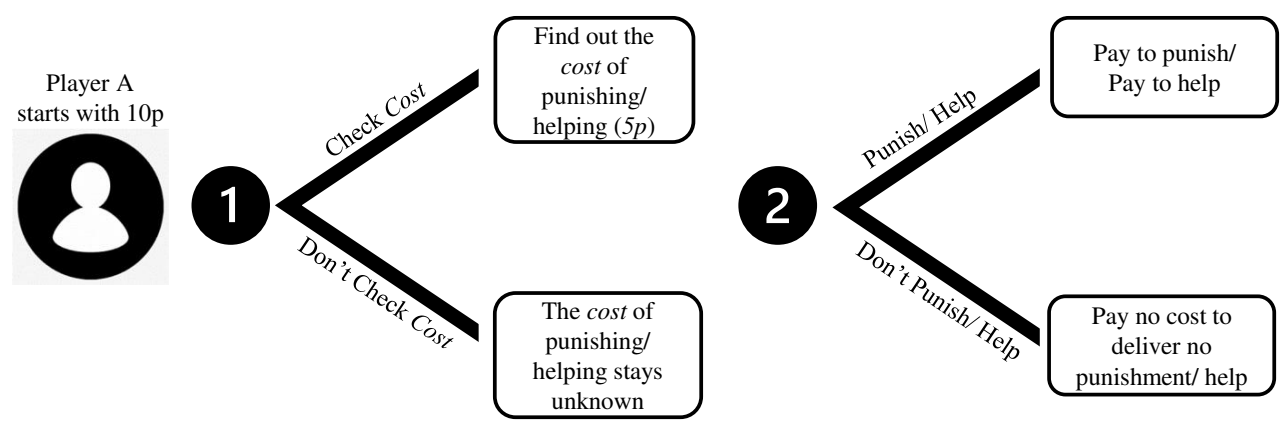


Study 2: Target Impact Deliberation

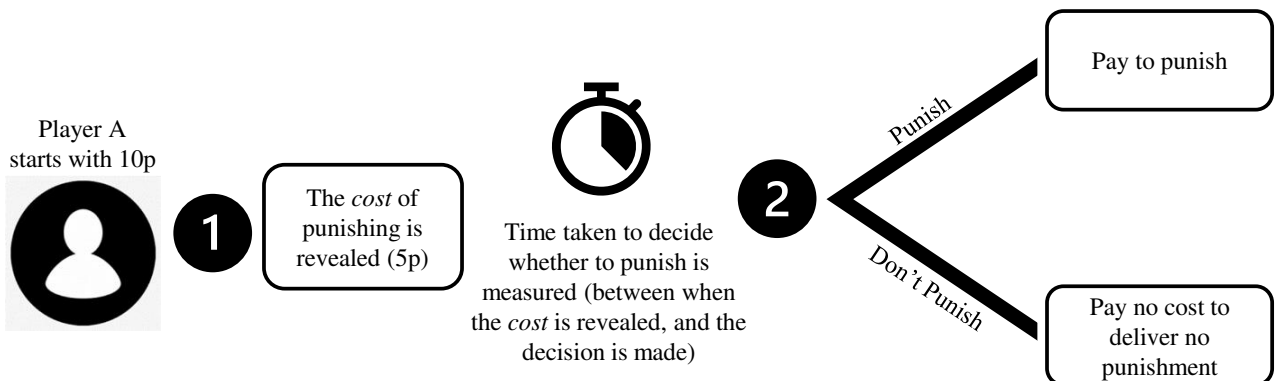


Stage 1: Punishing/ Helping Decision

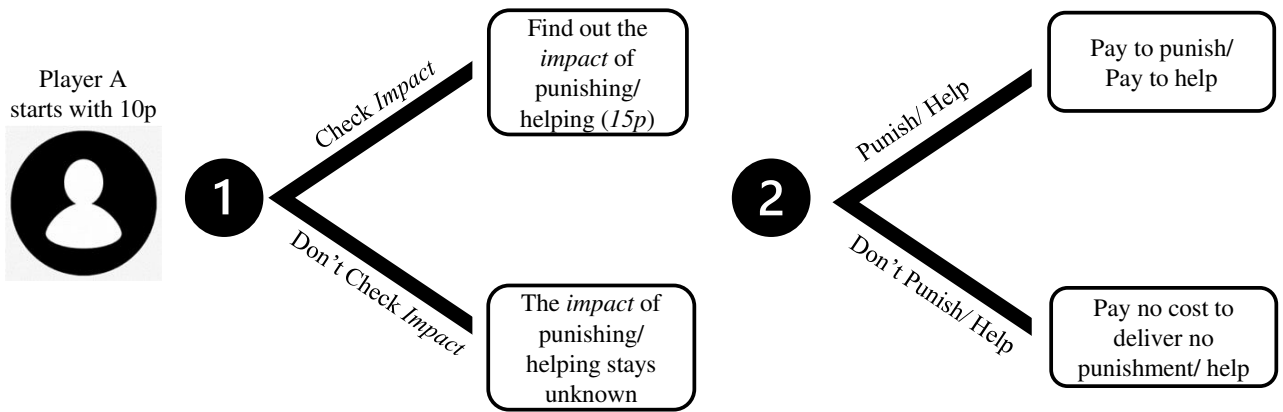
Study 1: Personal Cost Checking



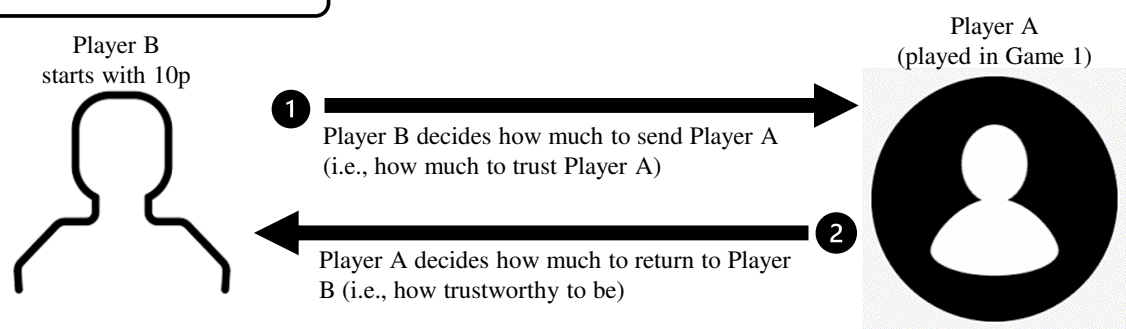
Study 1: Decision Time



Study 2: Target Impact Checking

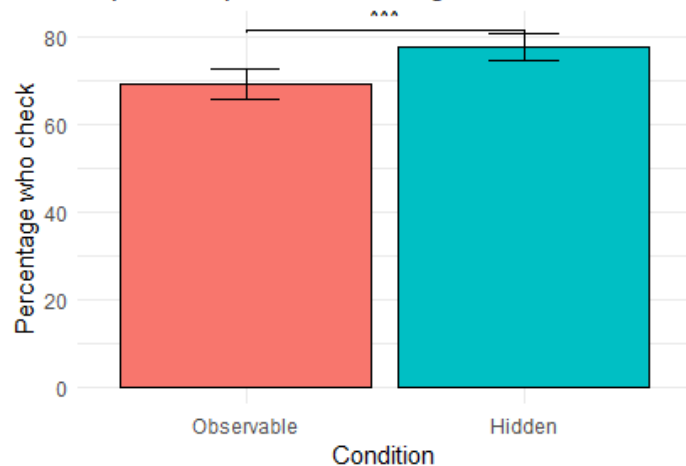
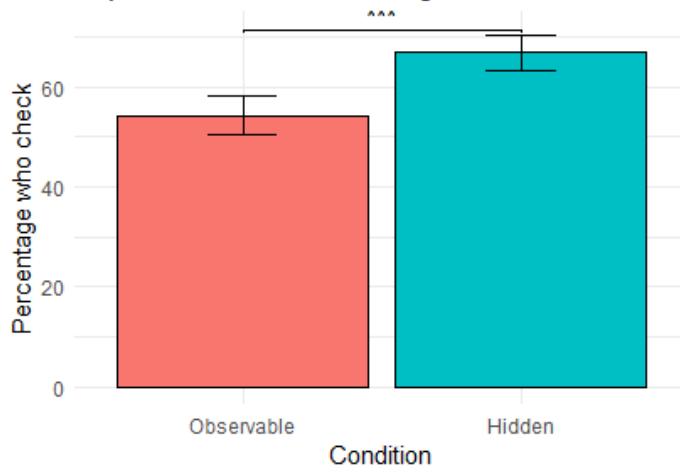
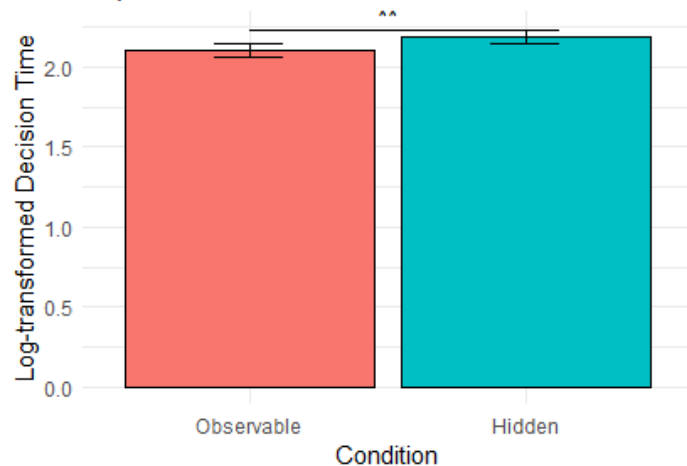
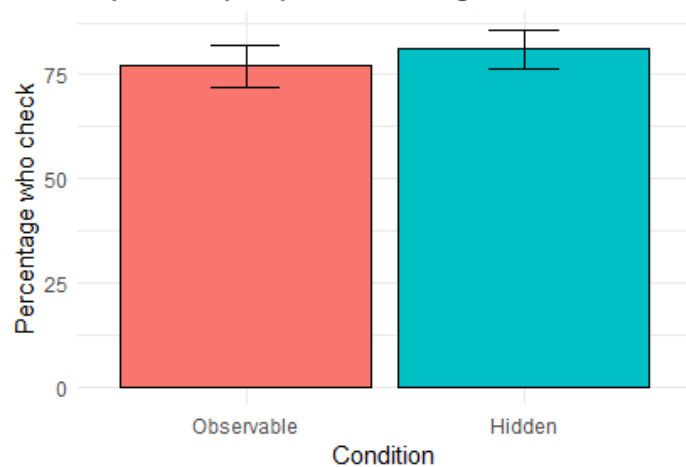
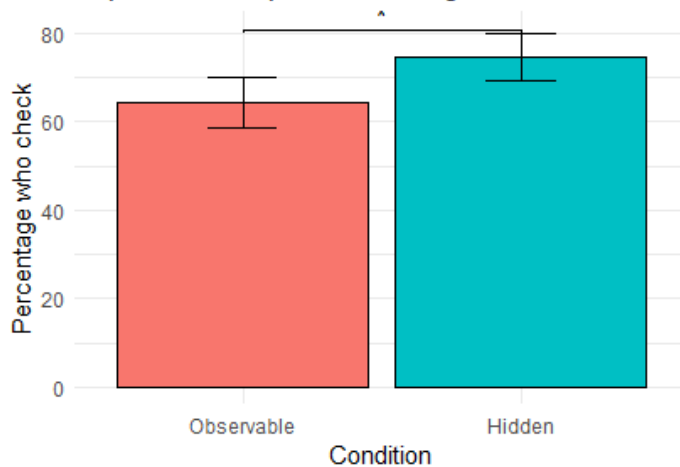


Stage 2: Trust Game



Process Observable Condition: Player B can base their sending decision both on Player A's helping/ punishing decision and decision process in Game 1

Process Hidden Condition: Player B can only base their sending decision on Player A's helping/ punishing decision in Game 1 (not their decision process)

Exp. 1: Help Cost Checking**Exp. 2: Pun Cost Checking****Exp. 3: Pun Cost Decision Time****Exp. 4: Help Impact Checking****Exp. 5: Pun Impact Checking**

Condition █ Observable █ Hidden

