

1 **The hidden cost of receiving favors:**

2 **A theory of indebtedness**

3
4 *Xiaoxue Gao^{1,2}, Eshin Jolly³, Hongbo Yu⁴, Huiying Liu⁵,*
5 *Xiaolin Zhou^{1,2,6,7,8*}, Luke J. Chang^{3*}*

6
7 ¹ Shanghai Key Laboratory of Mental Health and Psychological Crisis Intervention,
8 School of Psychology and Cognitive Science, East China Normal University,
9 Shanghai, China, 200062

10 ² School of Psychological and Cognitive Sciences, Peking University,
11 Beijing 100871, China

12 ³ Department of Psychological and Brain Sciences, Dartmouth College,
13 Hanover, NH 03755, USA

14 ⁴ Department of Psychological and Brain Sciences, University of California Santa
15 Barbara, Santa Barbara, CA 93106-9660, USA

16 ⁵ Mental Health Education Center, Zhengzhou University,
17 Zhengzhou 450001, Henan, China

18 ⁶ School of Business and Management, Shanghai International Studies University,
19 Shanghai 200083, China

20 ⁷ Beijing Key Laboratory of Behavior and Mental Health, Peking University,
21 Beijing 100871, China

22 ⁸ PKU-IDG/McGovern Institute for Brain Research, Peking University,
23 Beijing 100871, China

24
25 *Correspondence to:

26 Xiaolin Zhou (xz104@pku.edu.cn) and Luke J. Chang (luke.j.chang@dartmouth.edu)

27 **Abstract**

28

29 Receiving help or a favor from another person can sometimes have a hidden cost. In
30 this study, we explore these hidden costs by developing and validating a theoretical
31 model of indebtedness across three studies that combine a large-scale online
32 questionnaire, interpersonal games, computational modeling, and neuroimaging. Our
33 model captures how individuals perceive the altruistic and strategic intentions of the
34 benefactor. These perceptions produce distinct feelings of guilt and obligation that
35 together comprise indebtedness and motivate reciprocity. Perceived altruistic
36 intentions convey care and concern and are associated with activity in the insula,
37 dorsolateral prefrontal cortex and ventromedial prefrontal cortex, while perceived
38 strategic intentions convey expectations of future reciprocity and are associated with
39 activation in the temporal parietal junction and dorsomedial prefrontal cortex. We
40 further develop a neural utility model of indebtedness using multivariate patterns of
41 brain activity that captures the tradeoff between these feelings and reliably predicts
42 reciprocity behavior.

43

44

45 ***Key words:*** indebtedness; guilt; obligation; reciprocity; intention; gratitude

46 **Introduction**

47 Giving gifts and exchanging favors are ubiquitous behaviors that provide a concrete
48 expression of a relationship between individuals or groups^{1,2}. Altruistic favors
49 convey concern for a partner's well-being and signal a communal relationship such as
50 a friendship, romance, or familial tie³⁻⁵. These altruistic favors are widely known to
51 foster the beneficiary's positive feeling of gratitude, which can motivate reciprocity
52 behaviors that reinforce the communal relationship⁶⁻⁹. Yet in daily life, favors and
53 gifts can also be strategic and imply an expectation of reciprocal exchanges,
54 particularly in more transactive relationships^{2,4,5,10-12}. Accepting these favors can
55 have a hidden cost, in which the beneficiary may feel indebted to the favor-doer and
56 motivated to reciprocate the favor at some future point in time¹³⁻²¹. These types of
57 behaviors are widespread and can be found in most domains of social interaction. For
58 example, a physician may preferentially prescribe medications from a pharmaceutical
59 company that treated them to an expensive meal^{22,23}, or a politician might vote
60 favorably on policies that benefit an organization, which provided generous campaign
61 contributions²⁴. However, very little is known about the psychological and neural
62 mechanisms underlying this hidden cost of *indebtedness* and how it ultimately
63 impacts the beneficiary.

64

65 Immediately upon receipt of an unsolicited gift or favor, the beneficiary is likely to
66 engage in a mentalizing process to infer the benefactor's intentions²⁵⁻²⁷. Does this
67 person care about me? Or do they expect something in return? According to appraisal
68 theory²⁸⁻³³, these types of cognitive evaluations can evoke different types of feelings,
69 which will ultimately impact how the beneficiary responds. Psychological Game
70 Theory (PGT)³⁴⁻³⁶ provides tools for modeling these higher order beliefs about
71 intentions, expectations, and fairness in the context of reciprocity decisions^{26,27,37,38}.
72 Actions that are inferred to be motivated by altruistic intentions are more likely to be
73 rewarded, while those thought to be motivated by strategic or self-interested

74 intentions are more likely to be punished ^{26,27,37,38}. These intention inferences can
75 produce different emotions in the beneficiary ³⁹. For example, if the benefactor's
76 actions are perceived to be altruistic, the beneficiary may feel gratitude for receiving
77 help, but this could also be accompanied by the feeling of guilt for personally
78 burdening the benefactor ⁴⁰⁻⁴³. Both feelings motivate reciprocity out of concern for
79 the benefactor, which we refer to as "communal concern" throughout the paper ^{44,45}.
80 In contrast, if the benefactors' intentions are perceived to be strategic or even
81 duplicitous, then the beneficiary is more likely to feel a sense of obligation ^{13,14,21,46,47}.
82 Obligation can also motivate the beneficiary to reciprocate ^{13,14,21,46,47}, but unlike
83 communal concern, it arises from external pressures, such as social expectations and
84 reputational costs ^{48,49} and has been linked to feelings of pressure, burden, anxiety,
85 and resentment ⁴⁹⁻⁵¹. In everyday life, inferences about a benefactor's intentions are
86 often mixed and we propose that indebtedness is a superordinate emotion that
87 includes feelings of guilt for burdening the benefactor ⁴⁰⁻⁴³ and also social obligation
88 to repay the favor ^{13,14,21,46,47}.

89

90 In this study, we propose a conceptual model of indebtedness to capture how a
91 beneficiary's appraisals and emotions lead to reciprocity behaviors (Fig. 1).
92 Specifically, we propose that there are two components of indebtedness - guilt and the
93 sense of obligation, which are derived from appraisals about the benefactor's
94 intentions that can differentially impact the beneficiary's reciprocity behaviors. The
95 guilt component of indebtedness, along with gratitude, arises from appraisals of the
96 benefactor's altruistic intentions (i.e., perceived care from the help) and reflects
97 communal concern. In contrast, the obligation component of indebtedness results
98 from appraisals of the benefactor's strategic intentions (e.g., second-order belief of the
99 benefactor's expectation for repayment). Building on previous models of
100 other-regarding preferences ^{37,38,52}, we develop a computational model of the utility

101 associated with reciprocal behaviors as reflecting the trade-off between these different
102 feelings (Eq. 1).

103

$$104 \quad U(D_B) = \theta_B * \pi_B + (1 - \theta_B) * (\phi_B * U_{Communal} + (1 - \phi_B) * U_{Obligation}) \quad \mathbf{Eq.1}$$

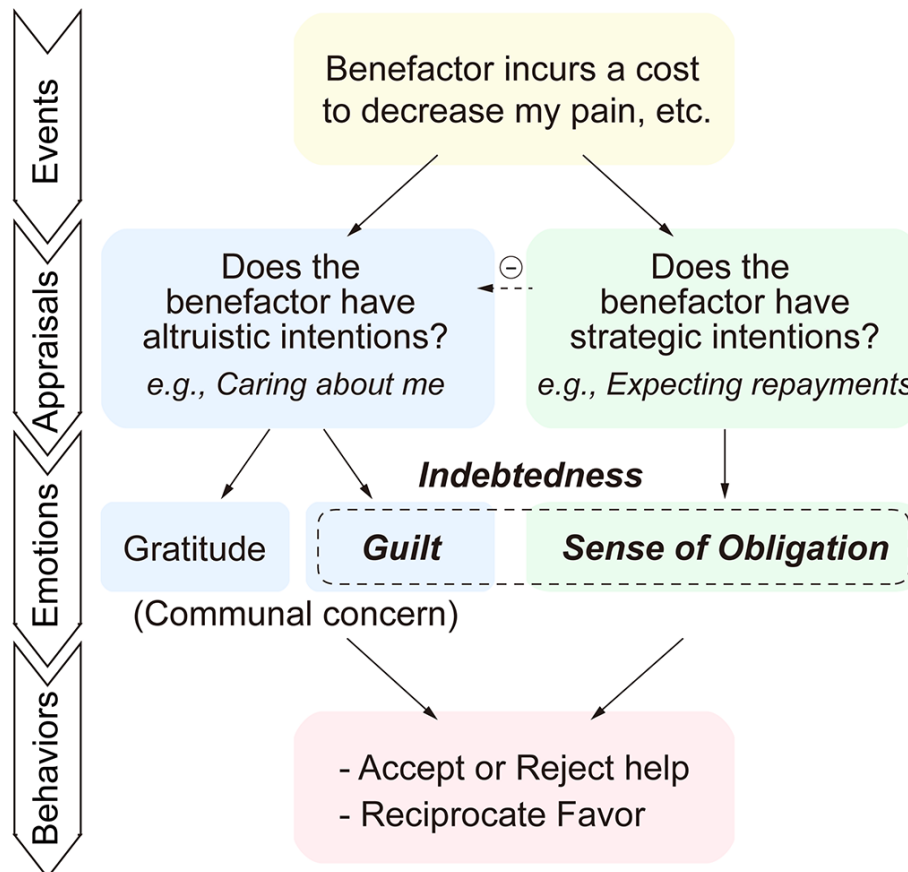
105

106 The central idea of this model is that upon receiving a favor D_A from a benefactor A ,
107 the beneficiary B chooses an action D_B that maximizes his/her overall utility U . This
108 utility is comprised of a mixture of values arising from self-interest π weighted by a
109 greed parameter θ , and feelings of communal concern $U_{Communal}$ and obligation
110 $U_{Obligation}$, which are weighted by the parameter ϕ . Larger ϕ values reflect the
111 beneficiary's higher sensitivity to feelings of communal concern relative to obligation.
112 $U_{Communal}$ reflects a linear combination of both gratitude and guilt components, but we
113 focus on guilt in the present article (see *Computational Modeling in Materials and*
114 *Methods*).

115

116 We validate our conceptual and computational models of indebtedness across a series
117 of studies. In Study 1, we explore lay intuitions of indebtedness using a large-scale
118 online questionnaire to test the hypothesis that indebtedness is a mixed feeling
119 comprised of both guilt and obligation. In Study 2, we evaluate how different
120 components of indebtedness are generated and influence behaviors in an interpersonal
121 game, in which benefactors (co-players) choose to spend some amount of their initial
122 endowments to reduce the amount of pain experienced by the participants. We test the
123 hypothesis that guilt and obligation arise from appraisals of the benefactor's intentions,
124 and specifically that appraisals of altruistic intentions produce guilt while appraisals
125 of strategic intentions lead to obligation. We then evaluate how well our
126 computational model (Eq. 1) captures these appraisal/feeling components and can
127 predict participants' decisions to reciprocate help in the interactive game. In Study 3,
128 we test the hypothesis that the two components of indebtedness are associated with

129 unique brain representations using functional magnetic resonance imaging (fMRI).
130 We create a neural utility model of indebtedness by applying our computational
131 model directly to multivariate brain patterns to demonstrate that neural signals reflect
132 the tradeoff between these feelings and can be used to predict participants'
133 trial-to-trial reciprocity behavior.
134

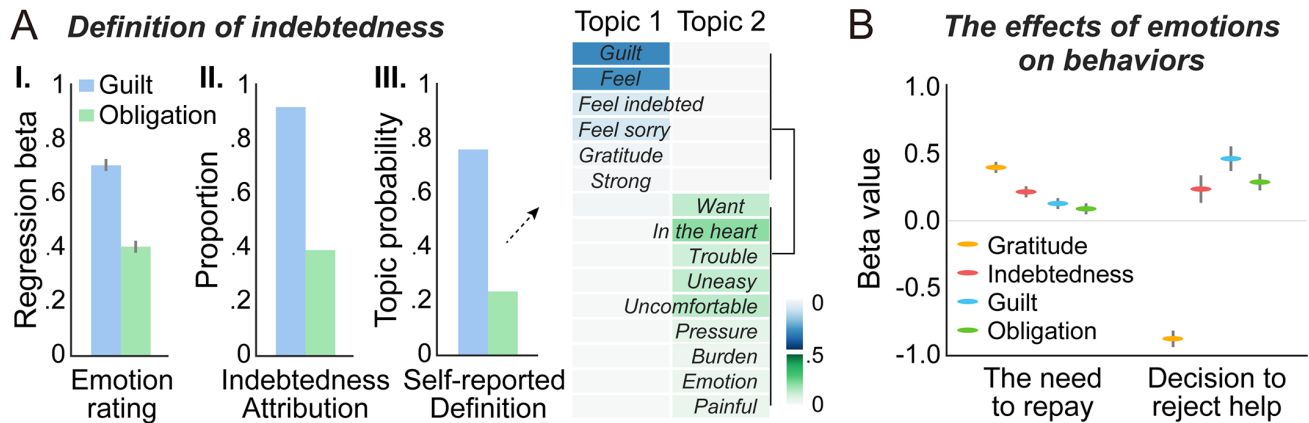


135 **Fig. 1 Conceptual model of indebtedness.** We propose that there are two components of indebtedness,
136 guilt and the sense of obligation, which are derived from appraisals about the benefactor's altruistic
137 and strategic intentions and can differentially impact the beneficiary's reciprocity behaviors. The
138 higher the perception of the benefactor's strategic intention, the lower the perception of the benefactor's
139 altruistic intention. The guilt component of indebtedness, along with gratitude, arises from appraisals of
140 the benefactor's altruistic intentions (i.e., perceived care from the help) and reflects communal concern.
141 In contrast, the obligation component of indebtedness results from appraisals of the benefactor's
142 strategic intentions (e.g., second-order belief of the benefactor's expectation for repayment). The
143 beneficiary makes trade-offs between communal and obligation feelings to determine the reciprocal
144 behaviors to favors (e.g., accept or reject the help and reciprocity after receiving help).

145 **Results**

146 *Indebtedness is a mixed feeling comprised of guilt and obligation*

147 In Study 1, we used an online questionnaire to characterize the subjective experience
148 of indebtedness in Chinese participants. First, participants (N = 1,619) described
149 specific events, in which they either accepted or rejected help from another individual
150 and rated their subjective experiences of these events. A regression analysis revealed
151 that both self-reported guilt and obligation ratings independently and significantly
152 contributed to increased indebtedness ratings ($\beta_{\text{guilt}} = 0.70 \pm 0.02, t = 40.08, p < 0.001$;
153 $\beta_{\text{obligation}} = 0.40 \pm 0.02, t = 2.31, p = 0.021$; Fig. 2A-I; Table S1). Second, participants
154 were asked to attribute sources of indebtedness in their daily lives. While 91.9%
155 participants stated that their feelings of indebtedness arose from feeling guilt for
156 burdening the benefactor, 39.2% participants reported feeling obligation based on the
157 perceived ulterior motives of the benefactor (Fig. 2A-II, Fig. S1A). Third, participants
158 were asked to describe their own personal definitions of indebtedness. We applied
159 Latent Dirichlet Allocation (LDA) based topic modeling⁵³ to the emotion-related
160 words extracted from the 100 words with the highest weight/frequency in the
161 definitions of indebtedness based on annotations from an independent sample of raters
162 (N = 80). We demonstrate that indebtedness is comprised of two latent topics (Fig. S1,
163 B-C). Topic 1 accounted for 77% of the variance of emotional words, including
164 communal-concern-related words such as "guilt," "feel," "feel sorry," "feel indebted,"
165 and "gratitude". In contrast, Topic 2 accounted for 23% of the emotional word
166 variance, including words pertaining to burden and negative bodily states, such as
167 "uncomfortable," "uneasy," "trouble," "pressure," and "burden" (Fig. 2A-III).



168 **Fig. 2 Subjective experiences of indebtedness in Study 1.** (A) Contributions of guilt and obligation to
 169 indebtedness in Study 1 in (I) the emotion ratings in the daily event recalling, (II) attribution of guilt and
 170 obligation as source of indebtedness, and (III) topic modeling of the emotional words in self-reported
 171 definition of indebtedness. The background color underlying each word represents the probability of
 172 this word in the current topic. (B) The influence of emotions on the self-reported need to reciprocate
 173 after receiving help and decisions to reject help. Error bars represent ± 1 SE.

174

175 Next, we examined how participants' emotion ratings were related to their
 176 self-reported behavioral responses to the help (Fig. 2B). Participants described events
 177 in which they chose to accept help and reported their experienced emotions. We
 178 found that indebtedness ($\beta = 0.20 \pm 0.04$, $t = 5.60$, $p < 0.001$), guilt ($\beta = 0.12 \pm 0.04$, $t =$
 179 2.98 , $p = 0.002$), obligation ($\beta = 0.09 \pm 0.04$, $t = 2.27$, $p = 0.023$), and gratitude ($\beta =$
 180 0.38 ± 0.04 , $t = 9.86$, $p < 0.001$) independently explained participants' reported need to
 181 repay after receiving help. Participants also described events, in which they chose to
 182 reject help and reported their anticipated counterfactual emotions had they instead
 183 accepted the benefactor's help⁵⁴. Decisions to reject help were negatively associated
 184 with gratitude ($\beta = -0.87 \pm 0.06$, $t = -13.65$, $p < 0.001$), but positively associated with
 185 indebtedness ($\beta = 0.23 \pm 0.10$, $t = 2.40$, $p = 0.017$), guilt ($\beta = 0.46 \pm 0.09$, $t = 5.06$, $p <$
 186 0.001), and obligation ($\beta = 0.28 \pm 0.06$, $t = 4.70$, $p < 0.001$). These results indicate that
 187 the two components of indebtedness (i.e., guilt and obligation) along with gratitude
 188 influence the behavioral responses to others' favors.

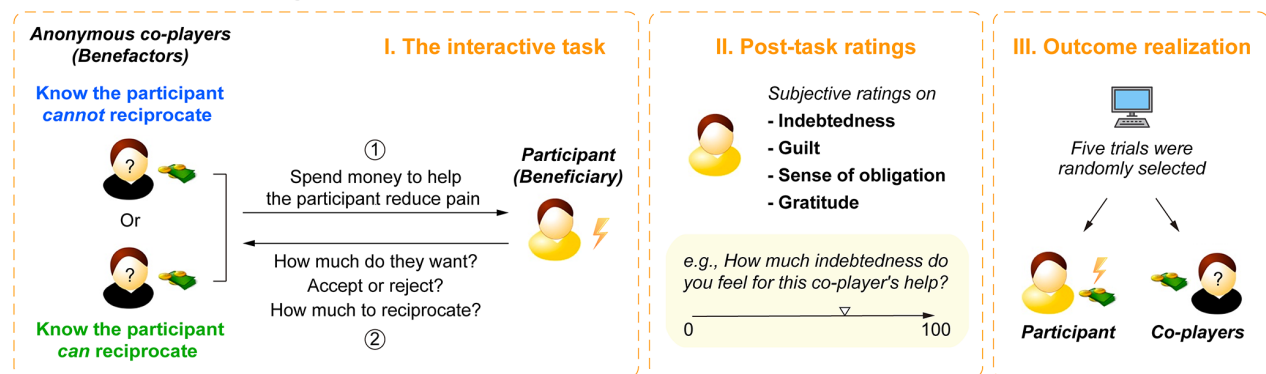
189

190 ***Benefactor's intentions lead to diverging components of indebtedness.***

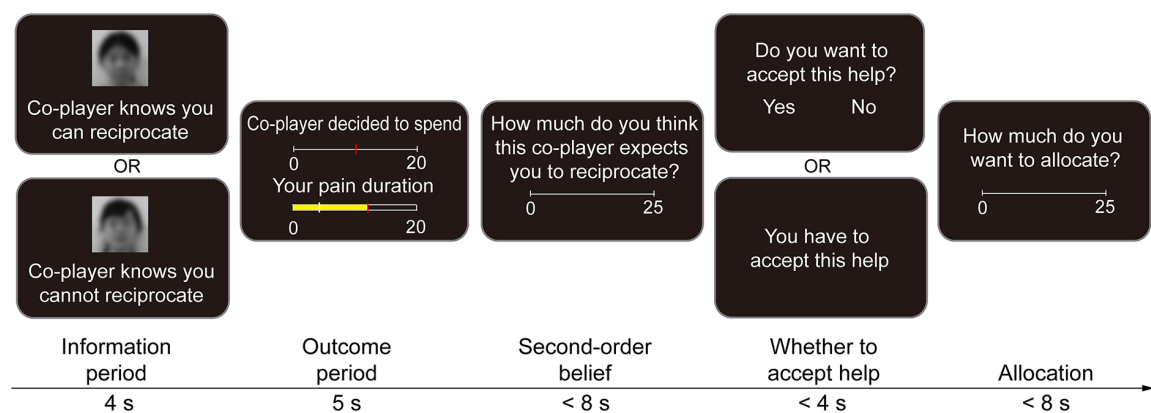
191 Next, we tested the predictions of the conceptual model regarding how different
192 components of indebtedness are generated and influence behaviors using a
193 laboratory-based task involving interactions between participants (Fig. 3). In Study 2a
194 (N = 51), participants were randomly paired with a different anonymous same-sex
195 co-player (benefactor) in each trial and were instructed that they would receive 20
196 seconds of pain stimulation in the form of a burst of medium intensity electrical
197 shocks. The participant was instructed that each benefactor was: (a) informed of the
198 participant's situation, (b) endowed with 20 yuan (~ \$3.1 USD), and (c) could spend
199 any amount of this endowment to help the participant reduce the duration of pain (i.e.,
200 benefactor's cost). The more the benefactor spent, the shorter the duration of the
201 participant's pain experience. After seeing how much money the benefactor chose to
202 spend, the participant reported how much they believed this benefactor expected them
203 to reciprocate for their help (i.e., second-order belief of the benefactor's expectation
204 for repayment). In half of the trials, the participant had to accept the benefactor's help;
205 in the other half, the participant could freely decide whether to accept or reject the
206 benefactor's help. Finally, at the end of each trial, the participant decided how much
207 of their own 25 yuan endowment (~ \$3.8 USD) he/she wanted to allocate to the
208 benefactor as reciprocity for their help. We manipulated the participant's beliefs about
209 the benefactor's intentions by providing additional information regarding the
210 benefactor's expectations of reciprocation. Each participant was instructed that before
211 making decisions, some benefactors knew that the participant would be endowed with
212 25 yuan and could decide whether to allocate some endowments to them as
213 reciprocity (i.e., ***Repayment possible condition***), whereas the other benefactors were
214 informed that the participant had no chance to reciprocate after receiving help (i.e.,
215 ***Repayment impossible condition***). In fact, participants could reciprocate in both
216 conditions during the task. After the task, participants recalled how much they
217 believed the benefactor cared for them, as well as their feelings of indebtedness,

218 obligation, guilt, and gratitude based on the help they received for each trial. At the
 219 end of the experiment, five trials of the interactive task were randomly selected to be
 220 realized and participants received the average number of shocks and money based on
 221 their decisions. Unbeknownst to participants, benefactors' decisions were
 222 pre-determined by a computer program (Table S2). We additionally manipulated the
 223 exchange rate between the benefactor's cost and the participant's benefit (i.e., the help
 224 efficiency) in Study 2b (N = 57) (see *Procedures of Study 2* in *Materials and Methods*,
 225 and Table S2). However, we did not observe any significant interaction effect
 226 between efficiency and any of other experimental variables (Table S3), and thus we
 227 combined the datasets of Studies 2a and 2b when reporting results in the main text for
 228 brevity.
 229

A Procedures for Study 2



B Detailed procedure for the interactive task



230 **Fig. 3 Experimental procedures for Study 2.** (A) General procedures. In the interactive task (I), the
 231 participant was instructed that anonymous co-players (benefactors) made single-shot decisions to help
 232 reduce the duration of the participant's pain, and the participant, in turn, decided whether to accept

233 *help and how much money to return to the benefactor. After the interactive task, all of the decisions in*
234 *the first session were displayed again in a random order. The participant was asked to recall and rate*
235 *their feelings of indebtedness, guilt, obligation, and gratitude when they received the help of the*
236 *benefactor (II. Post-task ratings). At the end of the experiment, five trials in the interactive task were*
237 *randomly selected to be realized to determine the participant's final amount of pain and payoff, and the*
238 *selected benefactor's final payoffs (III. Outcome realization). (B) Detailed procedure for the interactive*
239 *task. In each round, the benefactor, decided how much of their endowment to spend (i.e., benefactor's*
240 *cost) to reduce the participant's pain duration. The more the benefactor spent, the more the duration of*
241 *the participant's pain decreased. Participants indicated how much they thought the benefactor expected*
242 *them to reciprocate (i.e., second-order belief of the benefactor's expectation for repayment). In half of*
243 *the trials, participants could decide whether to accept the help; in the remaining trials, participants had*
244 *to accept help and could reciprocate by allocating monetary points to the benefactor. Unbeknownst to*
245 *participants, benefactors' decisions (i.e., benefactor's cost) were pre-determined by the computer*
246 *program (Table S2). We manipulated the perception of the benefactor's intentions by providing extra*
247 *information about whether the benefactor knew the participant could (i.e., Repayment possible*
248 *condition), or could not (i.e., Repayment impossible condition) reciprocate after receiving help.*

249

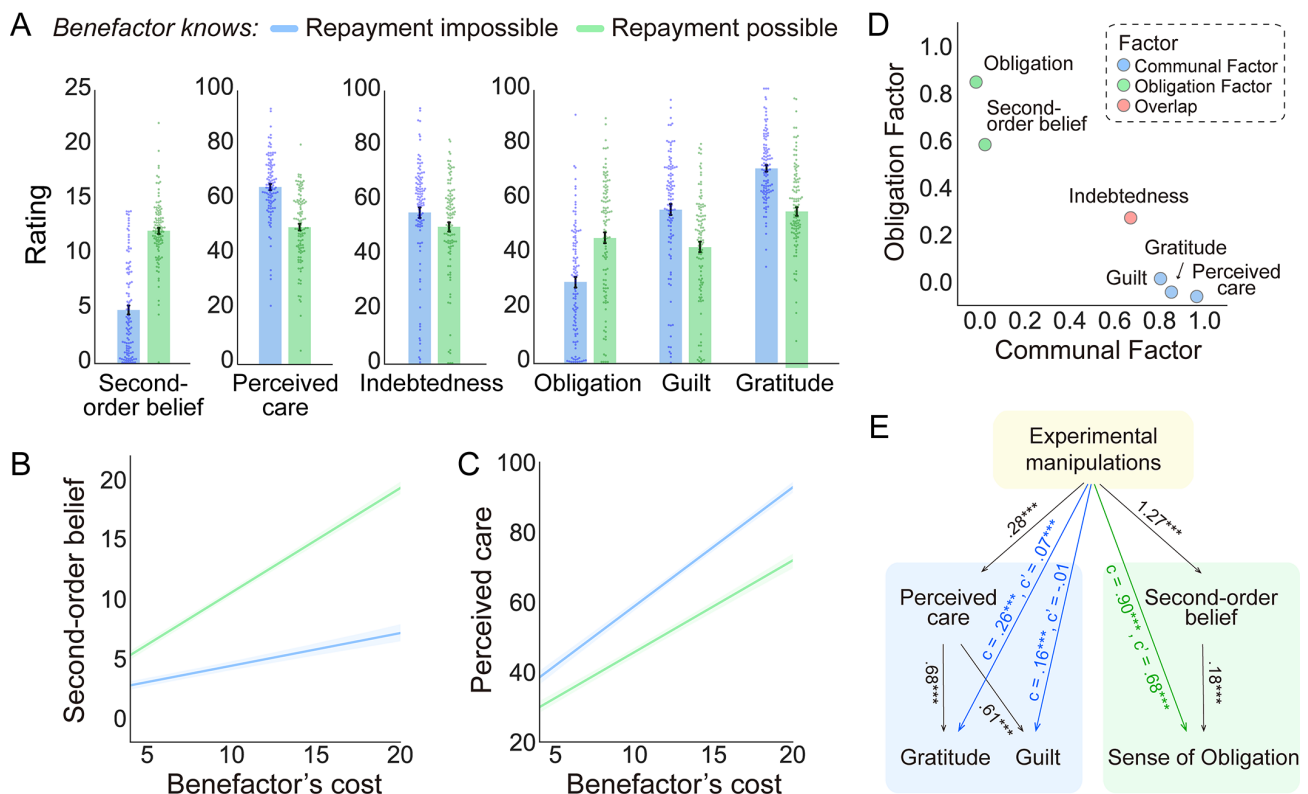
250 Our experimental manipulation successfully impacted participants' appraisals of the
251 benefactors' hidden intentions behind their help. Participants reported increased
252 second-order beliefs of the benefactor's expectations for repayment ($\beta = 0.53 \pm 0.03$, t
253 $= 15.71$, $p < 0.001$) and decreased perceived care ($\beta = -0.31 \pm 0.02$, $t = -13.89$, $p <$
254 0.001) (Fig. 4A, Table S3) when the participant believed the benefactor knew they
255 could reciprocate (Repayment possible) compared to when they could not reciprocate
256 (Repayment impossible). Both of these effects were amplified as the benefactor spent
257 more money to reduce the participant's duration of pain (Fig. 4, B-C; second-order
258 belief: $\beta = 0.22 \pm 0.02$, $t = 13.13$, $p < 0.001$; perceived care: $\beta = -0.08 \pm 0.01$, $t = -6.64$, p
259 < 0.001). In addition, perceived care was negatively associated with second-order
260 beliefs ($\beta = -0.44 \pm 0.04$, $t = -11.29$, $p < 0.001$) controlling for the effects of
261 experimental variables (i.e., extra information about benefactor's intention, cost, and
262 efficiency).

263

264 The belief manipulation not only impacted the participants' appraisals, but also their
265 feelings. Our conceptual model predicts that participants will feel indebted to
266 benefactors who spent money to reduce their pain, but for different reasons depending

267 on the perceived intentions of the benefactors. Consistent with this prediction,
 268 participants reported feeling indebted in both conditions, but slightly more in the
 269 Repayment impossible compared to the Repayment possible condition (Fig. 4A, Fig.
 270 S2A, $\beta = 0.09 \pm 0.03$, $t = 2.98$, $p = 0.003$). Moreover, participants reported feeling
 271 greater obligation (Fig. 4A, Fig. S2B, $\beta = 0.30 \pm 0.03$, $t = 9.28$, $p < 0.001$), but less
 272 guilt ($\beta = -0.25 \pm 0.02$, $t = -10.30$, $p < 0.001$), and gratitude ($\beta = -0.27 \pm 0.02$, $t = -13.18$,
 273 $p < 0.001$) in the Repayment possible condition relative to the Repayment impossible
 274 condition (Fig. 4A, Fig. S2, C-D). Similar to the appraisal results, these effects were
 275 magnified as the benefactor's cost increased (Fig. S2, B-D; obligation: $\beta = 0.11 \pm 0.01$,
 276 $t = 8.85$, $p < 0.001$; guilt: $\beta = -0.05 \pm 0.01$, $t = -4.28$, $p < 0.001$; gratitude: $\beta =$
 277 -0.06 ± 0.01 , $t = -4.20$, $p < 0.001$).

278



279 **Fig. 4 Appraisals and emotional responses to benefactor's help with different intentions.** (A)
 280 Participant's appraisals (i.e., second-order belief of how much the benefactor expected for repayment
 281 and perceived care) and emotion ratings (indebtedness, obligation, gratitude and guilt) in Repayment
 282 impossible and Repayment possible conditions. Each dot represents the average rating in the
 283 corresponding condition for each participant. (B and C) Participant's second-order beliefs of how

284 *much the benefactor expected for repayment and perceived care plotted as functions of extra*
285 *information about benefactor's intention and benefactor's cost. (D) Factor analysis showed that*
286 *participants' appraisals and emotions could be explained by two independent factors, which appeared*
287 *to reflect two distinct subjective experiences. The Communal Factor reflects participants' perception*
288 *that the benefactor cared about their welfare and resulted in emotions of gratitude and guilt, while the*
289 *Obligation Factor reflects participants' second-order beliefs about the benefactor's expectation for*
290 *repayment and the sense of obligation. (E) Simplified schematic representation of mediation analysis.*
291 *See full model in Fig. S3C. Results showed that second-order beliefs and perceived care appraisals*
292 *differentially mediated the effects of the experimental manipulations on emotional responses.*
293 *Second-order belief mediated the effects of the experimental manipulations on the sense of obligation,*
294 *while perceived care mediated the effects of experimental manipulations on gratitude and guilt. Error*
295 *bars represent ± 1 SE.*

296

297 We conducted two separate types of multivariate analyses to characterize the
298 relationships between appraisals and emotions. First, exploratory factor analysis (EFA)
299 on the subjective appraisals and emotion ratings in Study 2 revealed that 66% of the
300 variance in ratings could be explained by two factors (Fig. 4D, and Fig. S2E; Fig. S3,
301 A-B). The Communal Factor reflected participants' perception that the benefactor
302 cared about their welfare and resulted in emotions of guilt and gratitude, while the
303 Obligation Factor reflected participants' second-order beliefs about the benefactor's
304 expectation for repayment and the sense of obligation. Interestingly, indebtedness
305 moderately loaded on both factors. Second, a path analysis revealed that,
306 second-order beliefs and perceived care appraisals differentially mediated the effects
307 of the experimental manipulations on emotional responses (total indirect effect =
308 0.59 ± 0.04 , $Z = 14.49$, $p < 0.001$; Fig. 4E and Fig. S3C). Second-order beliefs
309 mediated the effects of the experimental manipulations on obligation (Indirect effect =
310 0.22 ± 0.03 , $Z = 7.18$, $p < 0.001$), while perceived care mediated the effects of the
311 experimental manipulations on guilt (Indirect effect = 0.17 ± 0.01 , $Z = 13.23$, $p < 0.001$)
312 and gratitude (Indirect effect = 0.19 ± 0.01 , $Z = 13.72$, $p < 0.001$). Together, these
313 results provide further support for the predictions of our conceptual model that
314 indebtedness is comprised of two distinct feelings. The guilt component of
315 indebtedness arises from the belief that the benefactor acts from altruistic intentions

316 (i.e., perceived care from the help), while the obligation component of indebtedness
317 arises when the benefactor's intentions are perceived to be strategic (e.g., expecting
318 repayment).

319

320 ***Behavioral responses to help are influenced by the benefactor's intentions***

321 Next, we examined participant's behaviors in response to receiving help from a
322 benefactor. Specifically, we were interested in whether participants would reciprocate
323 the favor by sending some of their own money back to the benefactor and also
324 whether they might outright reject the benefactor's help given the opportunity. We
325 found that participants reciprocated more money as a function of the amount of help
326 received from the benefactor, $\beta = 0.63 \pm 0.02$, $t = 25.60$, $p < 0.001$. This effect was
327 slightly enhanced in the Repayment impossible condition relative to the Repayment
328 possible condition, $\beta = 0.03 \pm 0.01$, $t = 2.99$, $p = 0.003$ (Fig. 5A). A logistic regression
329 revealed that when given the chance, participants were more likely to reject help in
330 the Repayment possible condition when they reported more obligation (rejection rate
331 = 0.37 ± 0.10), compared to the Repayment impossible condition (rejection rate =
332 0.30 ± 0.03), $\beta = 0.27 \pm 0.08$, $z = 3.64$, $p < 0.001$ (Fig. 5B).

333

334 ***Computational model captures feelings underlying responses to receiving favors***

335 We performed a more rigorous test of our conceptual model by constructing a
336 computational model of the proposed psychological processes (Eq. 1). This model
337 predicts a beneficiary's reciprocity behavior based on: (a) the benefactor's helping
338 behavior (i.e., benefactor's cost), (b) the belief manipulation (repayment
339 possible/impossible), and (c) a set of free parameters (i.e., θ , ϕ and κ) by simulating
340 appraisals of the benefactor's intentions and the associated feelings of communal
341 concern and obligation. The model then selects the behavior that maximizes the
342 beneficiary's expected utility considering the amount of money they will keep and
343 feelings of communal concern and obligation.

344

345 More specifically, for each trial, we modeled participant's second-order belief E_B'' of
346 how much they believed the benefactor expected them to reciprocate based on how
347 much the benefactor decided to spend to help D_A and whether the benefactor knew
348 repayment was possible (Eq. 2). In the Repayment impossible condition, participants
349 knew that the benefactor did not expect them to reciprocate, so we set E_B'' to zero.
350 However, in the Repayment possible condition, the benefactor knew that the
351 participant had money that they could spend to repay the favor. In this condition, we
352 modeled the E_B'' as proportional to the amount of money the benefactor spent to help
353 the participant.

354

$$E_B'' = \begin{cases} 0 & \text{Repayment impossible condition} \\ D_A & \text{Repayment possible condition} \end{cases} \quad \text{Eq. 2}$$

357 The participant's perceived care ω_B in each trial was defined as a function of the
358 benefactor's cost D_A and second-order belief E_B'' (Eq. 3). Specifically, we assumed
359 that the perceived care from help increased as a linear function of how much the
360 benefactor spent D_A from his/her endowment γ_A . However, this effect was reduced by
361 the second-order belief of the benefactor's expectation for repayment E_B'' . Here, the
362 parameter κ ranges from $[0, 1]$ and represents the degree to which the perceived
363 strategic intention E_B'' reduces the perceived altruistic intention ω_B . This creates a
364 nonlinear relationship between ω_B and E_B'' such that the relationship is negative when
365 κ is high, positive when κ is low, and uncorrelated in the current dataset with $\kappa =$
366 0.32 ± 0.01 , $\beta = -0.03 \pm 0.03$, $t = -1.23$, $p = 0.222$ (Fig. S4).

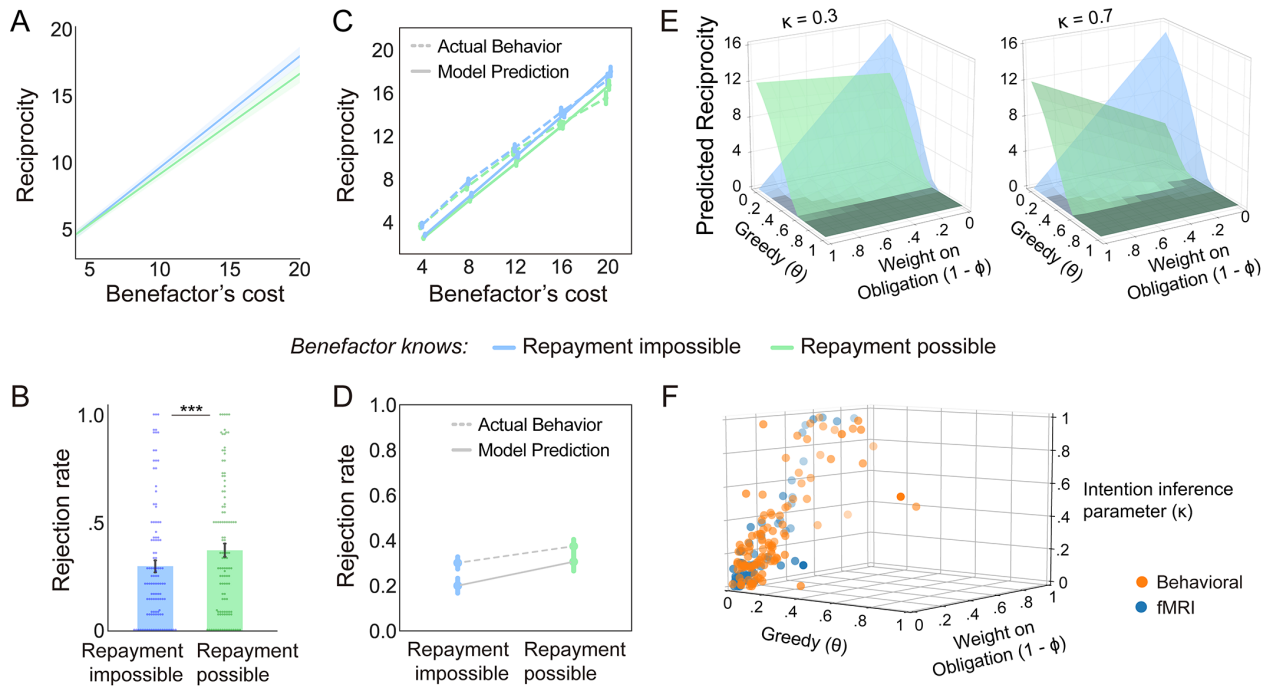
367

$$\omega_B = \frac{D_A - \kappa_B * E_B''}{\gamma_A} \quad \text{Eq. 3}$$

369 To validate our computational model, we tested whether it accurately captured each
370 proposed component process of our conceptual model and successfully predicted
371 participant's behavior. First, we found that each term of our model was able to

372 accurately capture trial-to-trial variations in self-reported appraisals of second-order
373 belief of the benefactor's expectation for repayment ($\beta = 0.68 \pm 0.03$, $t = 21.48$, $p <$
374 0.001 ; Fig. S5, A-B) and perceived care ($\beta = 0.72 \pm 0.03$, $t = 26.76$, $p < 0.001$; Fig. S5,
375 C-D). Moreover, the average value of the model term for perceived care was
376 correlated with the average self-reported perceived care across participants ($r = 0.27$,
377 $p = 0.004$), indicating that κ successfully captured individual differences in perceived
378 care. Our model assumes that appraisals produce their associated feelings, so the
379 perceived care ω_B and second-order belief E_B appraisals should serve as
380 representations of communal and obligation feelings. Supporting our predictions, the
381 perceived care model terms significantly predicted guilt ratings ($\beta = 0.47 \pm 0.03$, $t =$
382 17.21 , $p < 0.001$) as well as the Communal Factor scores obtained from EFA in Fig.
383 4D ($\beta = 0.81 \pm 0.03$, $t = 25.81$, $p < 0.001$), while the second-order belief model terms
384 significantly predicted obligation ratings ($\beta = 0.38 \pm 0.03$, $t = 12.67$, $p < 0.001$) and
385 the Obligation Factor scores ($\beta = 0.64 \pm 0.06$, $t = 15.97$, $p < 0.001$). Second, we found
386 that our indebtedness model was able to successfully capture the patterns of
387 participants' reciprocity behavior after receiving help ($r^2 = 0.81$, $p < 0.001$; Fig. 5C)
388 and significantly outperformed other plausible models, such as: (a) models that solely
389 included terms for communal concern or obligation, (b) a model that independently
390 weighted communal concern and obligation with separate parameters, (c) a model that
391 assumes participants reciprocate purely based on the benefactors helping behavior
392 (i.e., tit-for-tat)^{37,38}, and (d) a model that assumes that participants are motivated to
393 minimize inequity in payments^{52,55} (see *Supplementary Materials*, and Table S7).
394 Parameter recovery tests indicated that the parameters of the indebtedness model were
395 highly identifiable (correlation between true and recovered parameters: reciprocity $r =$
396 0.94 ± 0.07 , $p < 0.001$; Table S9). The accept/reject help model was able to accurately
397 capture decisions of whether to accept help (accuracy = 80.37%; Fig. 5D) but did not
398 significantly outperform models that solely included terms for communal concern or
399 obligation (Table S8). This likely stems a slight instability in the parameterization of

400 the model, which is confirmed by the moderate level of identifiability indicated by the
 401 parameter recovery tests ($r = 0.43 \pm 0.40$, $p < 0.001$; and Table S10) (see detailed
 402 explanations in *Computational Modeling in Materials and Methods*). See Tables S11
 403 and S12 for descriptive statistics, and Fig. S6 for distributions of model parameters.
 404



405 **Fig. 5 Computational model of indebtedness.** (A) Participants' reciprocity behavior in each trial
 406 plotted as function of extra information about benefactor's intention and benefactor's cost. (B) Overall
 407 rate of rejecting help in Repayment impossible and Repayment possible conditions, *** $p < 0.001$. Each
 408 dot represents the average rejection rate in the corresponding condition for each participant. (C) The
 409 observed amounts of reciprocity after receiving help and predictions generated by computational model
 410 at each level of the benefactor's cost in Repayment impossible and Repayment possible conditions. (D)
 411 The observed rates of rejecting help and predictions generated by computational model in Repayment
 412 impossible and Repayment possible conditions. (E) Model simulations for predicted reciprocity
 413 behavior in Repayment impossible and Repayment possible conditions at different parameterizations.
 414 (F) Best fitting parameter estimates of the computational model of indebtedness for each participant.
 415 Error bars represent the standard error of means.

416

417 A simulation of the model across varying combinations of the θ , ϕ and κ parameters
 418 revealed diverging predictions of the beneficiaries' response to favors in Repayment
 419 impossible and Repayment possible conditions (Fig. 5E). Not surprisingly, greedier
 420 individuals (higher θ) are less likely to reciprocate others' favors. However,

421 reciprocity changes as a function of the tradeoff between communal and obligation
422 feelings based on ϕ and interacts with the intention inference parameter κ . Increased
423 emphasis on obligation corresponds to increased reciprocity to favors in the
424 Repayment possible condition, but decreased reciprocity in the Repayment impossible
425 condition; this effect is amplified as κ increases. We found that most participants had
426 low θ values (i.e., greed), but showed a wide range of individual differences in κ and ϕ
427 parameters (Fig. 5F). Interestingly, the degree to which the perceived strategic
428 intention reduced the perceived altruistic intention during intention inference κ , was
429 positively associated with the relative weight on obligation ($1 - \phi$) during reciprocity
430 ($r = 0.79, p < 0.001$). This suggests that the participants who cared more about the
431 benefactor's strategic intentions also tended to be motivated by obligation when
432 deciding how much money to reciprocate.

433

434 ***Communal and obligation feelings are associated with distinct neural processes***

435 Next, we explored the neural basis of indebtedness guided by our computational
436 model and behavioral findings. Participants in Study 3 (N = 53) completed the same
437 task as Study 2 while undergoing fMRI scanning, except that they were unable to
438 reject help. First, we successfully replicated all of the behavioral results observed in
439 Study 2 (see Tables S1 and S4, and Figs. S7 and S8). In addition, we found that the
440 two-factor EFA model we estimated using the self-report data in Study 2 generalized
441 well to the independent sample in Study 3 using confirmatory factor analysis (CFA;
442 Fig. S7G), with comparative fit indices exceeding the > 0.9 acceptable threshold (CFI
443 = 0.986, TLI = 0.970) and the root mean square error of approximation and the
444 standardized root mean squared residual were within the reasonable fit range of $<$
445 0.08 (RMSEA = 0.079, SRMR = 0.019)⁵⁶⁻⁵⁸.

446

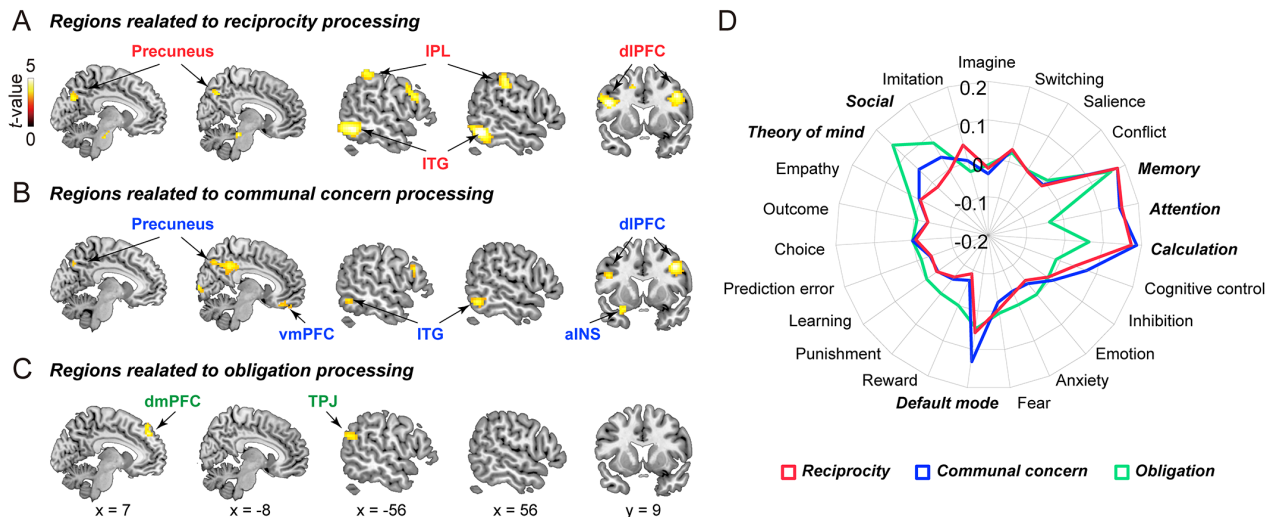
447 Second, we performed univariate analyses to identify brain processes during the
448 outcome period, where participants learned about the benefactor's decision to help.

449 Using a model-based fMRI analytic approach⁵⁹, we fit three separate general linear
450 models (GLMs) to each voxel's timeseries to identify brain regions that tracked
451 different components of the computational model. These included trial-by-trial values
452 for: (1) the amount of reciprocity, (2) communal concern, which depended on the
453 perceived care from the help (ω_B), and (3) obligation, which depended on the
454 second-order belief of the benefactor's expectation for repayment (E_B'') (for details,
455 see *Univariate fMRI Analyses in Materials and Methods*). We found that trial-by-trial
456 reciprocity behavior correlated with activity in bilateral dorsal lateral prefrontal cortex
457 (dlPFC), bilateral inferior parietal lobule (IPL), precuneus, and bilateral inferior
458 temporal gyrus (ITG) (Fig. 6A, Table S13). Trial-by-trial communal feelings tracked
459 with activity in the ventromedial prefrontal cortex (vmPFC), anterior insula,
460 precuneus, bilateral dlPFC, and bilateral ITG (Fig. 6B; Table S13). The processing of
461 obligation was associated with activations in dorsomedial prefrontal cortex (dmPFC)
462 and left temporo-parietal junction (TPJ) (Fig. 6C, Table S13).

463

464 To aid in interpreting these results, we performed meta-analytic decoding⁶⁰ using
465 Neurosynth⁶¹. Reciprocity-related activity was primarily associated with "Attention,"
466 "Calculation," and "Memory" terms. Communal feelings related activity was similar
467 to the reciprocity results, but was additionally associated with "Default mode" term.
468 Obligation activity was highly associated with terms related to "Social," "Theory of
469 mind (ToM)," and "Memory" (Fig. 6D). Together, these neuroimaging results reveal
470 differential neural correlates of feelings of communal concern and obligation and
471 support the role of intention inference in the generation of these feelings. The
472 processing of communal feelings was associated with activity in vmPFC, an area in
473 default mode network that has been linked to gratitude⁶²⁻⁶⁴, positive social value and
474 kind intention,^{65,66} as well as the insula, which has been previously related to guilt
475^{54,67,68}. In contrast, the processing of obligation was associated with activations of

476 theory of mind network, including dmPFC and TPJ, which is commonly observed
 477 when representing other peoples' intentions or strategies^{66,69,70}.



478 **Fig. 6 Neural processes associated with reciprocity, communal concern and obligation.** (A) Brain
 479 regions responding parametrically to trial-by-trial amounts of reciprocity. (B) Brain regions
 480 responding parametrically to trial-by-trial communal concern, which depended on the perceived care
 481 from the help (ω_B). (C) Brain regions identified in the parametric contrast for obligation (E_B''), the
 482 responses of which monotonically increased in the Repayment possible condition relative to the
 483 Repayment impossible condition. (D) Meta-analytical decoding for the neural correlates of reciprocity,
 484 communal concern and obligation, respectively. All brain maps thresholded using cluster correction
 485 FWE $p < 0.05$ with a cluster-forming threshold of $p < 0.001$ ⁷¹.

486

487 **Neural utility model of indebtedness predicts reciprocity behavior**

488 Having established that our computational model of indebtedness was able to
 489 accurately capture the psychological processes underlying feelings of communal
 490 concern and obligation, we next sought to test whether we could use signals directly
 491 from the brain to construct a utility function and predict reciprocity behavior (Fig.
 492 7A). Using brain activity during the outcome period of the task, we trained two
 493 whole-brain models using principal components regression with 5-fold
 494 cross-validation⁷²⁻⁷⁴ to predict the appraisals associated with communal concern (ω_B)
 495 and obligation (E_B'') separately for each participant. We have previously demonstrated
 496 that this approach is effective in reliably mapping the independent contribution of

497 each voxel in the brain to a psychological state to identify the neural representations
498 of affective states^{73,75,76}. These whole-brain patterns were able to successfully predict
499 the model representations of these feelings for each participant on new trials, though
500 with modest effect sizes (communal concern pattern: average $r = 0.21 \pm 0.03$, $p <$
501 0.001 ; obligation pattern: average $r = 0.10 \pm 0.03$, $p = 0.004$; Fig. 7A). Moreover,
502 these patterns appear to be capturing distinct information as they were not spatially
503 correlated, $r = 0.03$, $p = 0.585$. These results did not simply reflect differences
504 between the Repayment possible and Repayment impossible conditions as the results
505 were still significant after controlling for this experimental manipulation (communal
506 concern: average $r = 0.18 \pm 0.02$, $p < 0.001$; obligation: average $r = 0.04 \pm 0.02$, $p <$
507 0.024). Furthermore, we were unable to successfully discriminate between these two
508 conditions using a whole brain classifier ($accuracy = 55.0 \pm 1.25\%$, permutation $p =$
509 0.746).

510

511 Next, we assessed the degree to which our brain models could account for reciprocity
512 behavior. We used cross-validated neural predictions of communal concern (ω_B) and
513 obligation (E_B) feelings as inputs to our computational model of reciprocity behavior
514 instead of the original terms (Eq. 4):

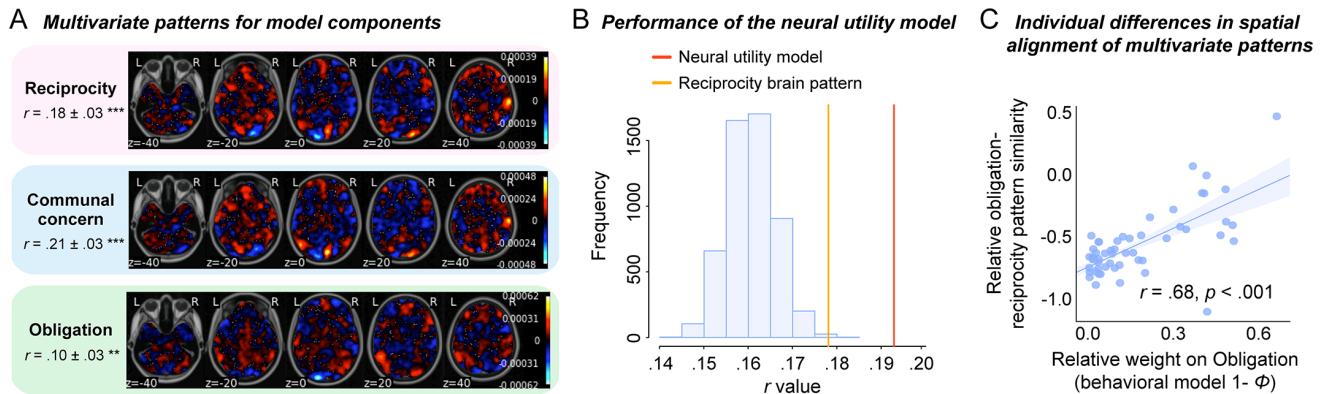
515

$$U(D_B) = \theta_B * \pi_B + (1 - \theta_B) * (\phi_B * \vec{\beta}_{map} \cdot \vec{Communal}_{map} + (1 - \phi_B) * \vec{\beta}_{map} \cdot \vec{Obligation}_{map}) \quad \text{Eq. 4}$$

517

518 where $\vec{\beta}_{map}$ refers to the pattern of brain activity during the Outcome phase of a single
519 trial and $\vec{Communal}_{map}$ and $\vec{Obligation}_{map}$ refer to the multivariate brain models
520 predictive of each participant's communal concern and obligation utilities
521 respectively. We were able to reliably predict reciprocity behavior with our
522 computational model informed only by predictions of communal and obligation
523 feelings derived purely from brain responses (average $r = 0.19 \pm 0.02$, $p < 0.001$, AIC

524 = 317.70 ± 5.00). As a benchmark, this model numerically outperformed a
 525 whole-brain model trained to directly predict reciprocity (average $r = 0.18 \pm 0.03$, $p <$
 526 0.001 , $AIC = 317.54 \pm 5.00$; Fig. 7A), but this difference only approached statistical
 527 significance, $t_{52} = 1.64$, $p = 0.108$.
 528



529 **Fig. 7 Neural utility model of indebtedness.** (A) Unthresholded multivariate patterns used to predict
 530 the amounts of reciprocity, trial-by-trial communal concern (ω_B) and obligation (E_B') separately. (B)
 531 We assessed the importance of the participant-specific model parameters estimated from the neural
 532 utility model (i.e., ϕ) by generating a null distribution of predictions after permuting the estimated ϕ
 533 parameter across participants 5,000 times. The red line indicates the performance of our neural utility
 534 model (r value of prediction), and the yellow line indicates the performance of the whole-brain model
 535 trained to directly predict reciprocity. The subject-specific weightings were important in predicting
 536 behavior as our neural utility model significantly outperformed a null model using parameters
 537 estimated for a different participant. (C) The relationship between the relative weight on obligation ($1 -$
 538 ϕ) derived from behavior and a neurally derived metric of how much obligation vs. communal feelings
 539 influenced reciprocity behavior (Eq. 14).

540

541 We performed several additional validations of the neural utility model to
 542 demonstrate its overall performance. First, we compared the parameter ϕ , which
 543 reflects the tradeoff between guilt and obligation estimated from the neural utility
 544 model and found that it strongly correlated with the same parameter estimated from
 545 the behavioral computational model across participants, $r = 0.88$, $p < 0.001$. Second,
 546 we assessed the individual specificity of ϕ derived from the neural utility model, to
 547 test how uniquely sensitive individuals are to communal concern versus obligation.

548 To do so, we generated a null distribution of predictions after permuting the estimated
549 ϕ parameter across participants 5,000 times. We found that the participant-specific
550 weightings were highly important in predicting behavior as our neural utility model
551 significantly outperformed null models using randomly shuffled ϕ parameters, $p <$
552 0.001 (Fig. 7B). Third, we tested how well our neural-utility model reflected the
553 trade-off between an individual's feelings of communal concern or obligation
554 estimated from the behavioral model. We hypothesized that the relative influence of a
555 particular feeling on behavior should be reflected in the spatial alignment of their
556 corresponding brain patterns⁷⁷. For example, if a participant weights obligation more
557 than communal concern during reciprocity (higher $1 - \phi$ estimated from the
558 behavioral model), then the spatial similarity between their obligation brain pattern
559 and the pattern that directly predicts their reciprocity behavior (reciprocity brain
560 pattern) should be relatively higher compared to the spatial similarity between their
561 communal concern pattern and reciprocity brain pattern (see *Neural Utility Model of*
562 *Indebtedness* in *Materials and Methods*). Our results support this hypothesis.
563 Participants who cared more about obligation relative to communal concern (higher
564 behavioral $1 - \phi$) also exhibited greater spatial alignment between their obligation and
565 reciprocity brain patterns relative to communal concern and reciprocity patterns, $r =$
566 0.68, $p < 0.001$ (Fig. 7C). These results provide evidence at the neural level indicating
567 that individuals appear to trade-off between feelings of communal concern and
568 obligation when deciding how much to reciprocate after receiving help from a
569 benefactor.

570

571 **Discussion**

572 Gift-giving, favor-exchanges, and providing assistance are behavioral expressions of
573 relationships between individuals or groups. While favors from friends and family
574 often engender reciprocity and gratitude, they can also elicit guilt in a beneficiary who
575 may feel that they have burdened a benefactor. Favors in more transactive

576 relationships, however, can evoke a sense of obligation in the beneficiary to repay the
577 favor. In this study, we sought to develop a comprehensive model of indebtedness that
578 outlines how appraisals about the intentions behind a favor are critical to the
579 generation of these distinct feelings, which in turn motivates how willing individuals
580 are to accept or reject help and ultimately reciprocate the favor.

581

582 We provide a systematic validation of this conceptual model of indebtedness across
583 three separate experiments by combining a large-scale online questionnaire,
584 behavioral measurements in an interpersonal game, computational modeling, and
585 neuroimaging. First, we used an open-ended survey to capture lay intuitions about
586 indebtedness based on past experiences. Overall, we find strong support that the
587 feeling of indebtedness resulting from receiving help from others can be further
588 separated into two distinct components – guilt from burdening the favor-doer and
589 obligation to repay the favor. Using topic modeling on lay definitions of indebtedness,
590 we find that guilt and gratitude appear to load on the same topic, while feeling words
591 pertaining to burden and negative body states load on a separate topic. Second, we
592 used a laboratory task designed to elicit indebtedness in the context of an
593 interpersonal interaction and specifically manipulated information intended to shift
594 the benefactor’s perceptions of the beneficiary’s intentions underlying their decisions.
595 Although our manipulation was subtle, we find that it was able to successfully change
596 participants’ appraisals about how much the beneficiary cared about them and their
597 beliefs about how much money the benefactor expected in return. Consistent with
598 appraisal theory²⁸⁻³³, these shifts in appraisals influenced participants’ subjective
599 feelings and ultimately their behaviors. Intentions perceived to be altruistic led to
600 increased guilt and gratitude, while intentions viewed as more strategic increased
601 feelings of obligation. All three feelings were associated with increased monetary
602 reciprocation back to the benefactor after receiving help. However, only the feeling of

603 obligation increased the probability of rejecting help when that option was available
604 to the participant.

605

606 One of the most notable contributions of this work is the development and validation
607 of a computational model of indebtedness, which does not require self-reported
608 ratings of emotions. The majority of empirical research on indebtedness^{21,46,47,78} and
609 other emotions^{79,80} has relied on participants' self-reported feelings in response to
610 explicit questions regarding social emotions, which has significant limitations, such as
611 its dependence on participants' ability to introspect^{81,82}. Formalizing emotions using
612 computational models is critical to advancing theory, characterizing their impact on
613 behavior, and identifying associated neural and physiological substrates^{39,83,84}.
614 However, the application of computational modeling to the study of social emotions is
615 a relatively new enterprise^{39,54,85,86}. Our model demonstrates how emotion appraisal
616 theory²⁸⁻³³ can be integrated with psychological game theory^{36,37} to predict behavior
617³⁹. We model emotions as arising from appraisals about perceived care and beliefs
618 about the beneficiary's expectations, which both ultimately increase the likelihood of
619 the benefactor selecting actions to reciprocate the favor. This model contributes to a
620 growing family of game theoretic models of social emotions such as guilt^{34,54},
621 gratitude⁸⁷, and anger^{88,89}, and can be used to infer feelings in the absence of
622 self-report providing new avenues for investigating other social emotions.

623

624 We provide a rigorous validation of our computational model using behaviors in the
625 interpersonal game, self-reported subjective experiences, and neuroimaging. First, our
626 model performs remarkably well at predicting participants' reciprocity behavior. It
627 also captures our theoretical predictions that participants would be more likely to
628 reject help when they perceived the benefactor to have strategic compared to altruistic
629 intentions. Second, the model's representations of communal concern and obligation
630 accurately captured participant's trial-to-trial self-reported appraisal and feeling

631 ratings. Third, our brain imaging analyses demonstrate that each feeling reflects a
632 distinct psychological process, and that intention inference plays a key role during this
633 process. Consistent with previous work on guilt^{54,67,68,90} and gratitude⁶²⁻⁶⁴, our model
634 representation of communal concern correlated with increased activity in the insula,
635 dlPFC, and default mode network including the vmPFC and precuneus. Obligation, in
636 contrast, captured participants' second order beliefs about expectations of repayment
637 and correlated with increased activation in regions routinely observed in mentalizing
638 including the dmPFC and TPJ^{66,69,70}.

639

640 We provide an even stronger test of our ability to characterize the neural processes
641 associated with indebtedness by deriving a “neural utility” model. Previous work has
642 demonstrated that it is possible to build brain models of preferences that can predict
643 behaviors^{91,92}. Here, we trained multivoxel patterns of brain activity to predict
644 participants' communal and obligation utility. We then used these brain-derived
645 representations of communal concern and obligation to predict how much money
646 participants ultimately reciprocated to the beneficiary. Remarkably, we found that this
647 neural utility model of indebtedness was able to predict individual decisions entirely
648 from brain activity and numerically outperformed (but not significantly) a control
649 model that provided a theoretical upper bound of how well reciprocity behavior can
650 be predicted directly from brain activity. Importantly, the neural utility model was
651 able to accurately capture each participant's preference for communal concern
652 relative to obligation. We observed a significant drop in our ability to predict behavior
653 when we randomly shuffled the weighting parameter across participants. In addition,
654 we find that the more the pattern of brain activity predicting reciprocity behavior
655 resembled brain patterns predictive of communal concern or obligation, the more our
656 behavioral computational model weighted this feeling in predicting behavior,
657 demonstrating that these distinct appraisals/feelings are involved in motivating
658 reciprocity decisions.

659

660 This work provides a substantial advance to our theoretical understanding of social
661 emotions. First, we highlight the complex relationship between gratitude and
662 indebtedness. We propose that feeling cared for by a benefactor, which we call
663 communal concern^{44,45}, is comprised of both guilt and gratitude. Each emotion
664 diverges in valence, with gratitude being positive⁶⁻⁹, and guilt being negative^{40-42,44,54},
665 but both promote reciprocity behavior. When faced with the offer of help, anticipated
666 gratitude should motivate the beneficiary to accept help in order to establish or
667 promote a relationship^{6,7}, whereas anticipated guilt should motivate the beneficiary to
668 reject help out of concern to protect the benefactor from incurring a cost^{44,54,93}.
669 Although we observed some evidence supporting this prediction, our interactive task
670 was not designed to explicitly differentiate guilt from gratitude, which limits the
671 ability of our computational model to capture the specific contributions of guilt and
672 gratitude to communal concern and likely impacted identifiability of the parameters of
673 the model for accepting/rejecting help (see *Computational Modeling in Materials and*
674 *Methods*). Future work might continue to refine the relationship between these two
675 aspects of communal concern both in terms of behaviors in experiments and
676 computations in models^{54,62-64,67,68,90}.

677

678 Second, our model provides a framework to better understand the role of relationships
679 and contexts in generating feelings of indebtedness within a single individual.
680 Different types of relationships (see Clark and Mills's theory of communal and
681 exchange relationships^{4,5}, and Alan Fiske's Relational Models Theory⁹⁴) have been
682 theorized to emphasize different goals and social norms which can impact social
683 emotions^{95,96}. For example, communal relationships prioritize the greater good of the
684 community and are more conducive to altruistic sharing, which can be signaled by
685 altruistic favors³⁻⁵. In contrast, exchange relationships are more transactional in
686 nature^{2,4,5,10-12} and emphasize maintaining equity in the relationship, which can be

687 signaled by strategic favors ⁹⁴. Our model proposes that perceptions of the
688 benefactor's intentions directly impact the feelings experienced by the beneficiary
689 (e.g., guilt & obligation). Although we deliberately attempted to minimize aspects of
690 the relationship between the benefactor and beneficiary by making players
691 anonymous to control for reputational effects, future work might experimentally
692 manipulate these relationships to directly test the hypothesis that relationship types
693 differentially moderate the responses of gratitude and subcomponents of
694 indebtedness.

695

696 Third, we present new evidence exploring the relationship between indebtedness and
697 guilt. Guilt and indebtedness are interesting emotions in that they are both negatively
698 valenced, yet promote prosocial behaviors. In previous work, we have operationalized
699 guilt as arising from disappointing a relationship partner's expectations ^{39,54,55,97},
700 which is conceptually related to the feeling of obligation in this paper. This feeling
701 results from disappointing a relationship partner or violating a norm of reciprocity and
702 is a motivational sentiment evoked by social expectations reflecting a "sense of
703 should" that is associated with other negative affective responses such as feelings of
704 pressure, burden, anxiety, and even resentment ⁴⁹⁻⁵¹. In other work, we have
705 investigated how guilt can arise from causing unintended harm to a relationship
706 partner ^{68,98}. This is conceptually more similar to how we frame guilt here, which
707 arises from the feeling that one has unnecessarily burdened a relationship partner even
708 though the help was never explicitly requested by the benefactor. We believe that
709 continuing efforts to refine mathematical models of emotions across a range of
710 contexts, will eventually allow the field to move beyond relying on the restrictive and
711 imprecise semantics of linguistic labels to define emotions (e.g., guilt, gratitude,
712 indebtedness, obligation, feeling, motivation, etc.).

713

714 Our study has several potential limitations, which are important to acknowledge. First,
715 although we directly and conceptually replicate our key findings across multiple
716 samples, all of our experiments recruit experimental samples from a Chinese
717 population. It is possible that there are cultural differences in the experience of
718 indebtedness, which may not generalize to other parts of the world. For example,
719 compared with Westerners who commonly express gratitude when receiving
720 benevolent help, Japanese participants (East Asian population) often respond with
721 "Thank you" or "I am sorry", indicating their higher experience of guilt after receiving
722 favors ^{40,41}. Cultural differences may perhaps reflect how the two components of
723 indebtedness are weighted, with guilt being potentially more prominent in East Asian
724 compared to Western populations, reflecting broader cultural differences in
725 collectivism and individualism. Second, our computational model may oversimplify
726 the appraisal and emotion generating processes. Our model operationalizes the
727 appraisals of perceived care and second order belief using information available to
728 each participant in the task (i.e., benefactor's helping behavior and manipulation
729 about the participants' ability to reciprocate). These appraisals are likely
730 context-dependent and our model may not generalize to other experimental contexts
731 without significant modification to how these appraisals are operationalized.
732 Although our model performed well capturing the patterns of participants' reciprocity
733 behaviors in this task, we believe it is important to continue to refine this model in
734 future studies.

735

736 In summary, in this study we develop a comprehensive and systematic model of
737 indebtedness and validate it across three studies combining large-scale online
738 questionnaire, an interpersonal interaction task, and neuroimaging. A key aspect to
739 this work is the emphasis on the role of appraisals about the intentions behind a favor
740 in generating distinct feelings of guilt and obligation, which in turn motivates how
741 willing beneficiaries are to accept or reject help and ultimately reciprocate the favor.

742 Together these findings highlight the psychological and neural mechanisms
743 underlying the hidden costs of receiving favors²²⁻²⁴.

744 **Materials and Methods**

745 ***Study 1 - Online Questionnaire***

746 **Participants.** Participants (1,808 graduate and undergraduate students) were recruited
747 from Zhengzhou University, China to complete an online questionnaire. None of the
748 participants reported any history of psychiatric, neurological, or cognitive disorders.
749 Participants were excluded if they filled in information irrelevant (e.g., this question is
750 boring, or I don't want to answer this question) to the question or experiment in the
751 essay question (189 participants), leaving 1,619 participants (812 females, 18.9 ± 2.0
752 (SD) years). While 98.7% participants reported the events of receiving help, 24.4%
753 participants reported the events of rejecting help within the past one year, which
754 resulted in 1,991 effective daily events. To extract the words related to emotions and
755 feelings in the definition of indebtedness, 80 additional graduate and undergraduate
756 students (45 females, 22.6 ± 2.58 years) were recruited from different universities in
757 Beijing to complete the word classification task. This experiment was carried out in
758 accordance with the Declaration of Helsinki and was approved by the Ethics
759 Committee of the School of Psychological and Cognitive Sciences, Peking University.
760 Informed written consent was obtained from each participant prior to participating.

761

762 **Experimental Procedures.** Participants reported their responses on the Questionnaire
763 Star platform (<https://www.wjx.cn/>) using their mobile phones. The questionnaire
764 consisted of two parts (see Appendix S1 for full questionnaire). Each participant was
765 asked to recall a daily event in which they **received help** (part 1) or **rejected help**
766 (part 2) from others, and to answer the questions regarding their appraisals, emotions,
767 and details of this event. Events were required to be clearly recalled and to have
768 occurred within the past year. Appraisal questions included: "To what extent do you
769 think the benefactor cared about your welfare? (i.e., perceived care)", and "To what
770 extent do you think the benefactor expected you to repay? (i.e., second-order belief)".
771 Emotion ratings included: indebtedness, guilt, obligation, and gratitude. The questions

772 for guilt⁴⁰⁻⁴³ and obligation^{13,14,21,46,47} were designed according to the operational
773 definitions used by previous research. For events in which participants accepted help
774 (Part 1), questions for behaviors included: "To what extent did you think you needed
775 to reciprocate?", "To what extent are you willing to reciprocate a favor to this
776 benefactor?", "To what extent do you want to accept/reject this offer?", and "To what
777 extent are you willing to interact with this benefactor in the future?" Questions were
778 the same for Part 2 (i.e., events in which participants rejected help), except that
779 participants were asked to *imagine* how they would feel or behave if they accepted
780 this help.

781

782 To explore how participants defined indebtedness, participants answered the
783 following two multiple-choice questions about the definition of indebtedness after
784 recalling the event: (1) In the context of helping and receiving help, what is your
785 definition of indebtedness? (2) In daily life, what do you think is/are the source(s) of
786 indebtedness? With four options "Negative feeling for harming the benefactor",
787 "Negative feeling for the pressure to repay caused by other's ulterior intentions",
788 "Both" and "Neither" (see details in Appendices S1 in *Supplementary Material*).

789

790 ***Study 2 - Interactive Task***

791 **Participants.** *For Study 2a (behavioral study)*, 58 graduate and undergraduate
792 Chinese Han students were recruited from Zhengzhou University, China, and 7
793 participants were excluded due to equipment malfunction, leaving 51 participants (33
794 females, 19.9 ± 1.6 years) for data analysis. *For Study 2b (behavioral study)*, 60
795 graduate and undergraduate Chinese Han students were recruited from Zhengzhou
796 University, China, and 3 participants were excluded due to failing to respond in more
797 than 10 trials, leaving 57 participants (45 females, 20.1 ± 1.8 years) for data analyses.
798 None of the participants reported any history of psychiatric, neurological, or cognitive
799 disorders. This experiment was carried out in accordance with the Declaration of

800 Helsinki and was approved by the Ethics Committee of the School of Psychological
801 and Cognitive Sciences, Peking University. Informed written consent was obtained
802 from each participant prior to participating.

803

804 **Experimental Procedure.** In Study 2a and Study 2b, seven participants came to the
805 experiment room together. An intra-epidermal needle electrode was attached to the
806 left wrist of each participant for cutaneous electrical stimulation⁹⁹. The first pain
807 stimulation was set as 8 repeated pulses, each of which was 0.2 mA and lasted for 0.5
808 ms. A 10-ms interval was inserted between pulses. Then we gradually increased the
809 intensity of each single pulse until the participant reported 6 on an 8-point pain scale
810 (1 = not painful, 8 = intolerable). Participants reported that they would only
811 experience the whole pulse train as a single stimulation, rather than as separate shocks.
812 The final intensity of pain stimulation was calibrated to a subjective pain rating of “6”,
813 which was a moderate punishment for the participants.

814

815 Both Study 2a and Study 2b consisted of two sessions. All stimuli were presented
816 using PsychToolBox 3.0.14 (www.psychtoolbox.org) in Matlab 2016a (Mathworks,
817 Natick, MA, USA). Participants were instructed as following:

818 *“In this experiment, you will play an interpersonal game, which is composed of two*
819 *roles: the Decider and the Receiver. The Receiver will be in some trouble and the*
820 *Decider can decide whether to help the Receiver at the cost of his/her own interests.*
821 *Several previous participants have come to our lab during Stage 1 of our study and*
822 *made decisions as the Deciders. Now this experiment belongs to Stage 2 of this*
823 *study. In the two sessions of the experiment, you will perform as the Receiver, facing*
824 *the decisions made by each previous Decider in Stage 1 and make your own*
825 *decisions.”*

826

827 During Session 1 (the main task), each participant played multiple single-shot rounds
828 of this interpersonal game as a Receiver with unique same-sex anonymous Deciders
829 (the co-player) (Fig. 3). The participant was instructed that the co-player in each trial
830 was distinct from the ones in any other trials and only interacted with the participant
831 once during the experiment. In each round, the participant was to receive a 20-second
832 pain stimulation with the intensity of 6. Each co-player was informed of the
833 participant's situation in Stage 1 and was endowed with 20 yuan (~ 3.1 USD). The
834 co-player could decide whether to spend some of their endowment to reduce the
835 duration of the participant's pain – more money resulted in shorter durations of pain.
836 The maximum pain reduction was 16 seconds to ensure that participants felt some
837 amount of pain on each trial.

838

839 Each trial began by informing the participant which Decider from Stage 1 was
840 randomly selected as the co-player for the current trial with a blurred picture of the
841 co-player and their subject id. The co-player's decision on how much they chose to
842 spend to help the participant was presented. Next, the participant indicated how much
843 he/she thought this co-player expected him/her to reciprocate (i.e., second-order belief
844 of the co-player's expectation for repayment; continuous rating scale from 0 to 25
845 using mouse, step of 0.1 yuan). In half of the trials, the participant had to accept the
846 co-player's help; in the other half, the participant could decide whether or not to
847 accept the co-player's help. If the participant accepted the help, the co-player's cost
848 and the participant's pain reduction in this trial would be realized according to the
849 co-player's decision; if the participant did not accept the help, the co-player would
850 spend no money and the duration of participant's pain stimulation would be the full
851 20 seconds. At the end of each trial, the participant was endowed with 25 yuan (~ 3.8
852 USD) and decided how much they wanted to allocate to the co-player as reciprocity in
853 this trial from this endowment (continuous choice from 0 to 25 using mouse, step of
854 0.1 yuan).

855

856 We manipulated the perceived intention of the co-player (i.e., the benefactor) by
857 providing participants with extra information regarding the co-player's expectation of
858 reciprocity (i.e., **extra information about benefactor's intention**) below the
859 co-player's subject id at the beginning of each trial. Each participant was instructed
860 that before making decisions, some co-players were informed that the participant
861 would be endowed with 25 yuan and could decide whether to allocate some
862 endowments to them as reciprocity (i.e., Benefactor knows repayment possible,
863 **Repayment possible condition**). The other co-players were informed that the
864 participant had no chance to reciprocate after receiving help (i.e., Benefactor knows
865 repayment is impossible, **Repayment impossible condition**). In fact, participants
866 could reciprocate in both conditions during the task. The endowment of the co-player
867 (γ_A) was always 20 yuan, and the endowment of the participant (γ_B) in each trial was
868 always 25 yuan. The endowment of the participant was always larger than the
869 endowment of the co-player to make the participant believe that the co-player
870 expected repayments in Repayment possible condition. Unbeknownst to the
871 participant, the co-players' decisions about how much of their endowment to allocate
872 to help reduce the participant's pain (i.e., **Benefactor's cost**) were uniformly sampled
873 from the available choices from an unpublished pilot study on helping behaviors. See
874 Table S2 for details about differences across experiments.

875

876 In Study 2b, to dissociate the effect of the benefactor's cost and participant's benefit
877 (i.e., pain reduction), we manipulated the exchange rate between the co-player's cost
878 and participant's pain reduction (i.e., **Efficiency**, 0.5, 1, and 1.5), whereas Efficiency
879 always 1 in Study 2a. Thus, the participant's pain reduction was calculated by: Pain
880 reduction = co-player's cost / co-player's endowment \times Efficiency \times Maximum pain
881 reduction (16s). For both Study 2a and Study 2b, each condition included one trial for

882 each Benefactor's cost – Efficiency combination. Therefore, there were 48 trials in
883 Study 2a and 56 trials in Study 2b (Table S2).

884

885 During Session 2, all of the decisions in the first session were displayed again in a
886 random order. After being shown the co-player's information and his/her decision, the
887 participant was asked to recall their feelings when they received the help of the
888 co-player. The rating order was counter-balanced across trials. The questions for
889 self-reported ratings on guilt⁴⁰⁻⁴³ and obligation^{13,14,21,46,47} were designed according
890 to the operational definitions built by previous research.

891

- 892 • "How much gratitude do you feel for this co-player's decision?" (Gratitude)
- 893 • "How much indebtedness do you feel for this co-player's decision?"
894 (Indebtedness)
- 895 • "How much do you think this decider cares about you?" (Perceived care)
- 896 • "How much pressure did you feel for the decider's expectation for repayment?"
897 (Obligation)
- 898 • "How much guilt do you feel for this co-player's decision?"(Guilt)

899

900 At the end of the experiment, five trials in Session 1 were randomly selected to be
901 realized. The participant received the average pain stimulation in these five trials. The
902 participant's final payoff was the average amount of endowment the participant left
903 for him/herself across the chosen trials. The participant was instructed that the final
904 payoff of each co-player was the amount of endowment the co-player left plus the
905 amount of endowment the participant allocated to him/her. Participants were informed
906 of this arrangement before the experiment began. After the experiment, participants
907 were further debriefed that the co-players' decisions they were faced with during the
908 experiment were actually pre-selected from participants' decisions in a previous

909 experiment by experimenters, and the co-players' decisions did not necessarily reflect
910 the natural distributions of others' helping behaviors.

911

912 ***Study 3 - fMRI Study***

913 **Participants.** For Study 3, 57 right-handed healthy graduate and undergraduate
914 Chinese Han students from Beijing, China took part in the fMRI scanning. Four
915 participants with excessive head movements (>2mm) were excluded, leaving 53
916 participants (29 females, 20.9 ± 2.3 years) for data analysis. None of the participants
917 reported any history of psychiatric, neurological, or cognitive disorders. This
918 experiment was carried out in accordance with the Declaration of Helsinki and was
919 approved by the Ethics Committee of the School of Psychological and Cognitive
920 Sciences, Peking University. Informed written consent was obtained from each
921 participant prior to participating.

922

923 **Experimental Procedure.** Each participant came to the scanning room individually.
924 The pain-rating procedure and the two sessions of the task in the fMRI study were
925 identical to Study 2a, except that participants always had to accept their co-player's
926 help. Session 1 (the main task) was conducted in the fMRI scanner, while Sessions 2
927 was conducted after participants exited the scanner. The scanning session consisted of
928 three runs (in total 54 trials) and lasted for approximately 39 min. Each run lasted for
929 13 min and consisted of 18 trials (including the 9 levels of the benefactor's cost in
930 Repayment possible condition and Repayment impossible conditions respectively),
931 trial order was pseudorandomized. See Table S2 for additional details about the
932 experimental design.

933

934 Each trial began with a 4 sec Information period, which showed the randomly
935 selected co-player's subject id, blurred picture, and information of whether this
936 co-player knew that the participant could or could not repay. This was followed by the

937 5 sec Outcome period, which included the co-player's decision on how much they
938 spent to help the participant. Participants then had up to 8 sec to report how much
939 he/she thought this co-player expected him/her to reciprocate (i.e., second-order belief
940 of the co-player's expectation for repayment; rating scale from 0 to 25 using left and
941 right buttons to move the cursor, step of 1 yuan). Next, participants had 8 sec to
942 decide how much of their 25 yuan endowment (~ 3.8 USD) to reciprocate to the
943 co-player (from 0 to 25 using left and right buttons to move the cursor, step of 1 yuan).
944 Before and after each period, a fixation cross was presented for a variable interval
945 ranging from 2 to 6 s, which was for the purpose of fMRI signal deconvolution.

946

947 ***Data Analyses in Study 1 (Online Questionnaire)***

948 **Validating Conceptual Model with Emotion Ratings.** We first attempted to validate
949 the conceptual model using the emotional ratings for daily-life events of receiving and
950 rejecting help obtained from online-questionnaire in Study 1. We conducted
951 between-participant linear regressions predicting indebtedness from guilt and
952 obligation ratings. We additionally examined the degree of multicollinearity between
953 guilt and obligation ratings using the variance inflation factor (VIF). The VIF reflects
954 the degree that any regressor can be predicted by a linear combination of the other
955 regressors (VIF = 5 serves as informal cutoff for multicollinearity – lower numbers
956 indicate less collinearity). Results demonstrated an acceptable level of
957 multicollinearity between guilt and obligation ratings (Table S1). To rule out the
958 possibility that these emotion ratings might covary with other related factors in Study
959 1 (e.g., benefactor's cost, the participant's benefit and the social distance between the
960 participant and the benefactor), we estimated a model with these additional variables,
961 which did not appreciably change the results (Table S1).

962

963 **Validating Conceptual Model with Self-Reported Appraisals.** Next, we
964 summarized participants' self-reported sources of their feelings of indebtedness. We

965 calculated the frequency that participants selected each of the four options in the
966 question "In daily life, what do you think is/ are the source(s) of indebtedness?" in
967 Study 1 (Fig, S1A), as well as how often that participants attributed "Negative feeling
968 for harming the benefactor" and "Negative feeling for the pressure to repay caused by
969 other's ulterior intentions" as the sources of indebtedness (i.e., the frequency of
970 choosing each single option plus the frequency of choosing "Both of the above").

971

972 **Validating Conceptual Model with Topic Modeling.** We also attempted to validate
973 the conceptual model by applying topic modeling to participant's open-ended
974 responses describing their own definition of indebtedness in Study 1. We used the
975 "Jieba" (<https://github.com/fxsjy/jieba>) package to process the text and excluded
976 Chinese stopwords using the stopwords-json dataset
977 (<https://github.com/6/stopwords-json>). Because Chinese retains its own characters of
978 various structures, we also combined synonyms of the same word as an additional
979 preprocessing step¹⁰⁰. Next, we computed a bag of words for each participant, which
980 entailed counting the frequency that each participant used each word and transformed
981 these frequencies using Term Frequency-Inverse Document Frequency (TF-IDF)
982^{101,102}. This method calculates the importance of a word in the whole corpus based on
983 the frequency of its occurrence in the text and the frequency of its occurrence in the
984 whole corpus. The advantage of this method is that it can filter out some common but
985 irrelevant words, while retaining important words that affect the whole text. Using
986 this method, the 100 words with the highest weight/frequency in the definitions of
987 indebtedness were extracted (Appendices S2). Words beyond these 100 had TF-IDF
988 weights < 0.01 (Fig. S1B), indicating that the words included in the current analysis
989 explained vast majority of variance in the definition of indebtedness. These 100 words
990 were then classified by an independent sample of participants (N = 80) into levels of
991 appraisal, emotion, behavior, person and other (see *Supplementary Materials*). We
992 conducted Latent Dirichlet Allocation (LDA) based topic modeling on the emotional

993 words of indebtedness using collapsed Gibbs sampling implemented in "lda" package
994 (<https://lda.readthedocs.io/en/latest/>)¹⁰³. LDA is a generative probabilistic model for
995 collections of discrete data such as text corpora, which is widely used to discover the
996 topics that are present in a corpus⁵³. It finds latent factors of semantic concepts based
997 on the co-occurrence of words in participant's verbal descriptions without
998 constraining participants' responses using rating scales, which currently dominates
999 emotion research¹⁰⁴. We selected the best number of topics by comparing the models
1000 with topic numbers ranging from 2 to 15 using 5-fold cross validation. Model
1001 goodness of fit was assessed using perplexity¹⁰⁵, which is a commonly used
1002 measurement in information theory to evaluate how well a statistical model describes
1003 a dataset, with lower perplexity denoting a better probabilistic model. We found that
1004 the two-topic solution performed the best (Fig. S1C).

1005

1006 **Validating Conceptual Model with Self-Reported Behaviors.** We next sought to
1007 test the predictions of the conceptual model using the self-reported behaviors from
1008 Study 1. First, we used data from Part 1 of the questionnaire and used linear
1009 regression to predict self-reported need to reciprocate from self-reported feelings of
1010 indebtedness, guilt, obligation and gratitude. Second, we combined the data of the
1011 events associated with receiving (Part 1) and rejecting help (Part 2) and used logistic
1012 regression to classify reject from accept behavior using self-reported counterfactual
1013 ratings of indebtedness, guilt, and obligation, and gratitude.

1014

1015 ***Data analyses in Study 2 (Interactive Task)***

1016 **Validating Conceptual Model with Emotion Ratings.** Similar to Study 1, we tested
1017 whether guilt and obligation contribute to indebtedness using the trial-by-trial
1018 emotional ratings in Study 2. We fit mixed effects regressions using lme4 predicting
1019 indebtedness ratings from guilt and obligation ratings with random intercepts and
1020 slopes for participants and experiments (e.g., 2a, 2b). Hypothesis tests were conducted

1021 using the lmerTest package ¹⁰⁶ in R. We additionally examined the degree of
1022 multicollinearity between guilt and obligation ratings using VIF (Table S1). To rule
1023 out the possibility that these emotion ratings might covary with other related factors
1024 the experimental variables in Studies 2 and 3 (e.g., benefactor's cost, extra
1025 information about the benefactor's intention and efficiency), we fit additional models
1026 controlling for these factors. Results of Study 2 replicated those in Study 1, and did
1027 not change after controlling for these variables (Table S1).

1028

1029 **The Effects of Experimental Conditions on Participants' Appraisal, Emotional**
1030 **and Behavioral Responses.** To test the effects of the benefactor's cost and extra
1031 information about benefactor's intention on beneficiary's appraisals (i.e., second-order
1032 belief and perceived care), emotions (i.e., gratitude, indebtedness, guilt, and
1033 obligation) and behaviors (reciprocity and whether reject help), in Study 2a we
1034 conducted LMM analyses for each dependent variable separately with participant as a
1035 random intercept and slope ¹⁰⁷ (Table S3).

1036

1037 **Relationships between Appraisals and Emotions.** To reveal the relationships
1038 between appraisals (i.e., second-order belief and perceived care) and emotions (i.e.,
1039 indebtedness, guilt, obligation, and gratitude), we estimated the correlations between
1040 these variables at both within-participant and between-participant levels. For
1041 within-participant analysis, for each pair of these six variables, we estimated the
1042 pearson correlation for each participant, transformed the data using a fisher r to z
1043 transformation, and then conducted a one-sample test using z values of all participants
1044 to evaluate whether the two variables were significantly correlated at the group level.
1045 This analysis captured the variability of appraisals and emotions across trials within
1046 participants (Fig. S3A). For between-participant analysis, for each of the six variables,
1047 we computed the average value of the variable across all trials for each participant.

1048 We then estimated the correlations between each pair of variables based on variability
1049 across participants (Fig. S3B).

1050

1051 Given the strong correlations between appraisals and emotions (Fig. S3, A-B, Tables
1052 S5 and S6), we conducted a factor analysis to examine the relationship between
1053 appraisals and emotions ¹⁰⁸. The Kaiser-Meyer-Olkin (KMO) Measure of Sampling
1054 Adequacy ¹⁰⁹ and Bartlett's test of sphericity ¹¹⁰ showed that the current data sets in
1055 Studies 2 and 3 were adequately sampled and met the criteria for factor analysis
1056 (Study 2: KMO value = 0.76, Bartlett's test $\chi^2 = 8801.85$, $df = 15$, $p < 0.001$; Study 3:
1057 KMO value = 0.77, Bartlett's test $\chi^2 = 2970.53$, $df = 15$, $p < 0.001$). All the variables
1058 were centered within participant to exclude the influences of individual differences in
1059 the range of ratings. We first applied exploratory factor analysis (EFA) in Study 2 to
1060 identify the number of common factors and the relationships between appraisals and
1061 emotions. To determine the number of components to retain, the correlation matrix
1062 between the 6 variables was submitted to a parallel analysis using the “psych”
1063 package ¹¹¹ for R. Parallel analysis performs a principal factor decomposition of the
1064 data matrix and compares it to a principal factor decomposition of a randomized data
1065 matrix. This analysis yields components whose eigenvalues (magnitudes) are greater
1066 in the observed data relative to the randomized data. The nScree function was used to
1067 determine the number of factors to retain. The result pointed to a two-factor solution
1068 (Fig. S2E). Factors were then estimated and extracted by combining ML factor
1069 analysis with oblique rotation using the “GPArotation ” package for R ¹⁰⁸. Next, we
1070 conducted confirmatory factor analysis (CFA) using the data of Study 3 to test the
1071 two-factor model built by Study 2 in an independent sample. CFA was conducted
1072 using “lavaan” package ¹¹² for R. Results remained the same after controlling for the
1073 experimental variables.

1074

1075 To test whether the two appraisals mediated the observed effects of experimental
1076 variables on emotional responses, we conducted a multivariate mediation analysis
1077 using structural equation modeling using the ‘lavaan’ package in R¹¹². In this analysis,
1078 experimental variables (extra information about benefactor’s intention, benefactor’s
1079 cost, information-cost interaction, and efficiency) were taken as independent variables,
1080 ratings of second-order belief and perceived care were taken as mediators, and ratings
1081 of guilt, gratitude and the sense of obligation were taken as dependent variables. First,
1082 we built a full model that included all pathways between variables. Then,
1083 non-significant pathways in the full model were excluded from the full model to
1084 improve the fitness of the model. In the final model, experimental variables included
1085 extra information about the benefactor’s intentions, the benefactor’s cost, and their
1086 interaction; efficiency was excluded due to the non-significant effects. Moreover, in
1087 the final model, second-order beliefs mediated the effects of the experimental
1088 variables on obligation, whereas perceived care mediated the effects of experimental
1089 variables on guilt and gratitude. This model performed well (RSMEA = 0.023, SRMR
1090 = 0.004, CFI = 1.000, TLI = 0.997, BIC = 27496.99) and explained participants’
1091 responses better than the full model (RSMEA = 0.046, SRMR = 0.004, CFI = 1.000,
1092 TLI = 0.986, BIC = 27543.52).

1093

1094 **Using Communal and Obligation Factors as Predictors for Behaviors.** To
1095 investigate how participants’ appraisals and emotions influenced their behavioral
1096 responses, we conducted two separate LMMs to predict participants’ reciprocity
1097 behavior, and decisions of whether or not to accept help. Each model included the
1098 scores for Communal and Obligation Factors estimated from the factor analysis as
1099 fixed effects and random intercepts and slopes for participants. See detailed results in
1100 the *Supplementary Material*.

1101

1102 **Computational Modeling.** We built separate models predicting participant's
 1103 reciprocity and rejection behaviors based on the conceptual model of indebtedness
 1104 (see Table S14 for all model object definitions). The utility of each reciprocity
 1105 behavior for player B $U(D_B)$ was modeled using Eq. 1 (Model 1.1), where self-interest
 1106 π_B is the percentage of money kept by player B out of their endowment γ_B .

1107

$$1108 \quad \pi_B = \begin{cases} \frac{\gamma_B - D_B}{\gamma_B} & \text{Reciprocity} \\ \frac{D_A * \mu}{\max(D_A * \mu)} & \text{Accept/Reject help} \end{cases} \quad \text{Eq. 5}$$

1109

1110 Based on our conceptual model (Fig. 1), we define $U_{Communal}$ as a mixture of feelings
 1111 of gratitude $U_{Gratitude}$ and guilt U_{Guilt} , in which the parameter δ_B ranges from [0,1] and
 1112 specifies how much player B cares about gratitude relative to guilt. As the focus of
 1113 this paper is on indebtedness, we set δ_B to zero and leave it to future work to build a
 1114 model of gratitude $U_{Gratitude}$ and explore its relationship with guilt (see also
 1115 *Discussion*). Thus, for this paper $U_{Communal}$ is synonymous with U_{Guilt} .

1116

$$1117 \quad U_{Communal} = \delta_B * U_{Gratitude} + (1 - \delta_B) * U_{Guilt} \quad \text{Eq. 6}$$

1118

1119 We separately modeled the appraisals of second-order beliefs E_B'' of the benefactor's
 1120 expectation for repayment (Eq. 2) and perceived care ω_B (Eq. 3), and used them to
 1121 capture guilt and obligation feelings (Eq. 7 and Eq. 8). We defined the
 1122 appraisal/feelings of U_{Guilt} and $U_{Obligation}$ as:

1123

$$1124 \quad U_{Guilt} = \begin{cases} -\left(\frac{\omega_B * \gamma_B - D_B}{\gamma_B}\right)^2 & \text{Reciprocity} \\ \omega_B & \text{Accept/Reject Help} \end{cases} \quad \text{Eq. 7}$$

1125

$$1126 \quad U_{Obligation} = \begin{cases} -\left(\frac{E_B'' - D_B}{\gamma_B}\right)^2 & \text{Reciprocity} \\ \frac{E_B''}{\gamma_B} & \text{Accept/Reject help} \end{cases} \quad \text{Eq. 8}$$

1127

1128 Participants maximized U_{Guilt} by minimizing the difference between the benefactor's
1129 reciprocity D_B and their perception of how much they believed the benefactor cared
1130 about them ω_B , scaled by the endowment size γ_B . In contrast, participants maximized
1131 $U_{Obligation}$ by minimizing the difference between the amount they reciprocated D_B and
1132 their second-order belief of how much they believed the benefactor expected them to
1133 return (E_B''). We note our mathematical operationalization of obligation here is more
1134 akin to how we have previously modeled guilt from disappointing others in previous
1135 work^{34,39,54,55} (see also *Discussion*).

1136

1137 We modeled the utility U associated with the participants' amounts of reciprocity D_B
1138 after receiving help in Eq. 1 (Model 1.1), where ϕ is a free parameter between 0 and 1,
1139 which captures the trade-off between feelings of communal concern and obligation.
1140 The model selects the participant's decision D_B associated with the highest utility. We
1141 estimated the model parameters for Eq. 1 by minimizing the sum of squared error of
1142 the percentages that the model's behavioral predictions deviate actual behaviors over
1143 all the trials that participants had to accept help using Matlab's `fmincon` routine. More
1144 formally, for each participant we minimized the following objective function:

1145

$$SSE = \sum_{t=1}^n \left(\frac{D_B(t) - \max(U(D_B(t)))}{\gamma_B} * 100 \right)^2 \quad \text{Eq. 9}$$

1146

1147

1148 with t indicating trial number. To avoid ending the fitting procedure at a local
1149 minimum, the model-fitting algorithm was initialized at 1000 random points in
1150 theta-phi-kappa parameter space for each participant.

1151

1152 We created a separate model for decisions to accept or reject help. Self-interest π_B for
1153 accepting help was defined as the percentage of pain reduction from the maximum
1154 amount possible, which depended on how much the benefactor spent to help D_A and
1155 the exchange rate between the benefactor's cost and the participant's benefit μ (see Eq.

1156 5). U_{Guilt} and $U_{Obligation}$ were defined as functions of ω_B and E_B respectively (Eq. 7
1157 and Eq. 8). We model the utility of accepting and rejecting help as:
1158

$$1159 \quad \begin{cases} U(Accept) = \theta_B * \pi_B + (1 - \theta_B) * (\phi_B * U_{Communal} - (1 - |\phi_B|) * U_{Obligation}) \\ U(Reject) = 0 \end{cases}$$

1160 **Eq. 10 (Model 2.1)**

1161 In this model, $U(Reject)$ was set to zero, because the situation would remain and the
1162 participant's emotional responses would not change if the participant did not accept
1163 help. Increased obligation reduces the likelihood of accepting help to avoid being in
1164 the benefactor's debt^{13,14,113}. In contrast, $U_{Communal}$ has a more complex influence on
1165 behavior, with guilt decreasing the likelihood of accepting help to avoid burdening a
1166 benefactor^{34,54}, and gratitude motivating accepting help to build a communal
1167 relationship^{6,7}. However, because $U_{Communal} = U_{Guilt} = \omega_B$ in this formulation,
1168 there is no variability in the design for the model to be able to disentangle the effect of
1169 gratitude from that of guilt. To address this complexity, we constrain ϕ to be within
1170 the interval of [-1, 1], and explicitly divide up the parameter space such that $\phi > 0$
1171 indicates a preference for gratitude and motives the participants to accept the help,
1172 while $\phi < 0$ indicates a preference for guilt and motives the participants to reject the
1173 help.

1174

$$1175 \quad \begin{cases} \phi_B > 0 & \textit{Gratitude} \\ \phi_B < 0 & \textit{Guilt} \end{cases} \quad \text{Eq. 11}$$

1176

1177 Regardless of whether the participant is motivated primarily by guilt or gratitude,
1178 participants can still have a mixture of obligation captured by $1 - |\phi|$, which ranges
1179 from [0,1]. Unfortunately, if participants are equally sensitive to gratitude and guilt, ϕ
1180 will reduce to zero and the weight on obligation increases, which decreases the model
1181 fit and leads to some instability in the parameters (see *Results* and *Discussion*).

1182

1183 We computed the probability of the decision of whether to accept or reject help using
1184 a softmax specification with inverse temperature parameter λ , which ranges from $[0,1]$.
1185 In each trial, the probability of the participant choosing to accept help is given by

1186

$$1187 \quad P(\textit{Accept}) = \frac{e^{U(\textit{Accept})/\lambda}}{e^{U(\textit{Accept})/\lambda} + e^{U(\textit{Reject})/\lambda}} \quad \text{Eq. 12}$$

1188

1189 We then conducted maximum likelihood estimation at the individual level by
1190 minimizing the negative log likelihood of the decision that the participant made ($D_B =$
1191 Accept or Reject) over each trial t with 1000 different starting values:

1192

$$1193 \quad LLE = - \sum_{t=1}^n \log(P(D_B(t))) \quad \text{Eq. 13}$$

1194

1195 Covariance between model terms implies that there might be multiple configurations
1196 of parameters that can produce the same predicted behavior. This means that, in
1197 practice, the more that these constructs covary, the less identifiable our parameters
1198 will become. We conducted parameter recovery analyses to ensure that our model was
1199 robustly identifiable¹¹⁴. To this end, we simulated data for each participant using our
1200 models and the data from each trial of the experiment and compared how well we
1201 were able to recover these parameters by fitting the model to the simulated data. We
1202 refit the model using 1000 random start locations to minimize the possibility of the
1203 algorithm getting stuck in a local minimum. We then assessed the degree to which the
1204 parameters could be recovered by calculating the similarity between all the
1205 parameters estimated from the observed behavioral data and all the parameters
1206 estimated from the simulated data using a Pearson correlation.

1207

1208 We compared the indebtedness model with both communal and obligation feelings
1209 with other plausible models, such as: (a) models that solely included $U_{Communal}$ and
1210 $U_{Obligation}$ terms, (b) a model that independently weighted $U_{Communal}$ and $U_{Obligation}$ with

1211 separate parameters, (c) a model that assumes participants reciprocate purely based on
1212 the benefactors helping behavior (i.e., tit-for-tat)^{37,38}, and (d) a model that assumes
1213 that participants are motivated to minimize inequity in payments^{52,55}. See
1214 *Supplementary Materials* for details.

1215

1216 To validate the model representations of appraisals/feelings, we predicted participants
1217 self-reported appraisals, emotions and the two factors extracted from EFA using the
1218 trial-to-trial model representations using LMMs that included random intercepts and
1219 slopes for each participant.

1220

1221 *Data analyses in Study 3 (fMRI Experiment)*

1222 **fMRI Data Acquisition and Preprocessing.** Images were acquired using a 3T
1223 Prisma Siemens scanner (Siemens AG, Erlangen, Germany) with a 64-channel head
1224 coil at Peking University (Beijing, China). T2-weighted echoplanar images (EPI)
1225 were obtained with blood oxygenation level-dependent (BOLD) contrast. Sixty-two
1226 transverse slices of 2.3 mm thickness that covered the whole brain were acquired
1227 using multiband EPI sequence in an interleaved order (repetition time = 2000 ms,
1228 echo time = 30 ms, field of view = 224×224 mm², flip angle = 90°). The fMRI data
1229 preprocessing and univariate analyses were conducted using Statistical Parametric
1230 Mapping software SPM12 (Wellcome Trust Department of Cognitive Neurology,
1231 London). Images were slice-time corrected, motion corrected, resampled to 3 mm × 3
1232 mm × 3 mm isotropic voxels, and normalized to MNI space using the EPInorm
1233 approach in which functional images are aligned to an EPI template, which is then
1234 nonlinearly warped to stereotactic space¹¹⁵. Images were then spatially smoothed
1235 with an 8 mm FWHM Gaussian filter, and temporally filtered using a high-pass filter
1236 with a cutoff frequency of 1/128 Hz.

1237

1238 **Univariate fMRI Analyses.** We used a model-based fMRI analytic approach ⁵⁹ to
1239 identify brain regions that parametrically tracked different components of the
1240 computational model during the Outcome phase of the task (5s). GLM 1 predicted
1241 brain responses based on the participant's reciprocity behavior D_B . GLM 2 predicted
1242 brain responses based on communal concern, which we modeled as the participant's
1243 appraisal of the co-player's perceived care ω_B . GLM 3 predicted brain responses
1244 based on obligation, which we modeled as a linear contrast of the participant's
1245 second-order belief of the benefactor's expectation for repayment E_B'' . We chose to
1246 use the appraisals rather than the $U_{Communal}$ and the $U_{Obligation}$ terms, as those terms
1247 create costs based on the squared deviation from reciprocity behavior, which results in
1248 a large proportion of trials where the deviations are near zero as a result of
1249 participant's decisions, making them inefficient for parametric analysis to capture
1250 how successfully participants behaved in accordance with their feelings. Instead, ω_B
1251 and E_B'' better captured the inferences that comprised participants' feelings and were
1252 more suitable for testing our hypotheses about brain responses. Regressors of no
1253 interest for GLM1 and GLM 2 included: (a) Outcome phase (onset of the presentation
1254 of the benefactor's decision, 5s), (b) Information period (onset of the presentation of
1255 the benefactor's picture and extra information regarding intention, 4s), (c)
1256 Second-order belief rating period (starting from the time the rating screen presented
1257 and spanning to the time that the participant made choice), (d) Allocation period
1258 (starting from the time the rating screen presented and spanning to the time that the
1259 participant made choice), (e) Missed responses (the missing decision period for
1260 second-order belief or allocation, 8s), and (f) six head motion realignment parameters.
1261 Contrasts were defined as the positive effect of the parametric modulator of interest.
1262
1263 For GLM3, because our computational model's representation of second order beliefs
1264 E_B'' had a very non-normal distribution, we constructed a piecewise linear contrast.
1265 This entailed creating four separate regressors modeling different parts of the function

1266 during the Outcome phase: (1) Repayment impossible, (2) Repayment possible and
1267 low benefactor's cost (i.e., 4, 6, or 8), (3) Repayment possible and medium
1268 benefactor's cost (i.e., 10, 12, or 14), (4) Repayment possible and high benefactor's
1269 cost (i.e., 16, 18, or 20). Subsequently, for each participant, we constructed a contrast
1270 vector of $c = [-6, 1, 2, 3]$. This piecewise linear contrast ensures that brain responses
1271 to the Repayment impossible trials are lower than all of the Repayment possible trials.
1272 We have successfully used this approach in previous work modeling guilt using
1273 similar psychological game theoretic utility models ⁵⁴.

1274

1275 For all GLMs, events in each regressor were convolved with a double gamma
1276 canonical hemodynamic response function. Second-level models were constructed as
1277 one-sample t tests using contrast images from the first-level models. For whole brain
1278 analyses, all results were corrected for multiple comparisons using cluster correction
1279 $p < 0.05$ with a cluster-forming threshold of $p < 0.001$, which attempts to control for
1280 family wise error (FWE) using Gaussian Random Field Theory. This approach
1281 attempts to estimate the number of independent spatial resels or resolution elements in
1282 the data necessary to control for FWE. This calculation requires defining an initial
1283 threshold to determine the Euler Characteristic of the data. It has been demonstrated
1284 that an initial threshold of $p < 0.001$ does a reasonable job of controlling for false
1285 positives at 5% using this approach ⁷¹.

1286

1287 **Meta-analytical Decoding.** To reveal the psychological components associated with
1288 the processing of reciprocity, communal concern and obligation, we conducted
1289 meta-analytic decoding using the Neurosynth Image Decoder ⁶¹
1290 (<http://neurosynth.org>). This allowed us to quantitatively evaluate the spatial
1291 similarity ⁶⁰ between any Nifti-format brain image and selected meta-analytical
1292 images generated by the Neurosynth database. Using this online platform, we
1293 compared the unthresholded contrast maps of reciprocity, communal concern and

1294 obligation against the reverse inference meta-analytical maps for 23 terms generated
1295 from this database, related to basic cognition (i.e., Imagine, Switching, Saliency,
1296 Conflict, Memory, Attention, Cognitive control, Inhibition, Emotion, Anxiety, Fear,
1297 and Default mode) ¹¹⁶, social cognition (Empathy, Theory of mind, Social, and
1298 Imitation) ¹¹⁷ and decision-making (Reward, Punishment, Learning, Prediction error,
1299 Choice, and Outcome) ¹¹⁸.

1300

1301 **Neural Utility Model of Indebtedness.** We constructed a neural utility model of
1302 indebtedness by combining our computational model of indebtedness with
1303 multivariate pattern analysis (MVPA) ¹¹⁹. First, using principal components
1304 regression with 5-fold cross-validation, we trained two separate multivariate
1305 whole-brain models predictive of communal concern (ω_B) and obligation (E_B'') terms
1306 in our behavioral model separately for each participant ⁷²⁻⁷⁴. This analysis was carried
1307 out in Python 3.6.8 using the NLTools package ¹²⁰ version 0.3.14 (<https://nltools.org/>).
1308 This entailed first performing temporal data reduction by estimating single-trial beta
1309 maps of the Outcome period for each participant. Then for each participant, we
1310 separately predicted ω_B and E_B'' from a vectorized representation of the single trial
1311 beta maps. Because these models have considerably more voxel features (~328k) than
1312 trial observations, we performed a principal components analysis to reduce the feature
1313 space and used the principal components to predict the model appraisal
1314 representations (e.g., ω_B and E_B''). We then back-projected the estimated beta
1315 components from the regression back into the full voxel feature space, and then back
1316 to 3-D space. For each whole-brain model, we extracted the cross-validated prediction
1317 accuracy (r value) for each participant, conducted r -to- z transformation, and then
1318 conducted a one-sample sign permutation test to evaluate whether each model was
1319 able to significantly predict the corresponding term.

1320

1321 We used the cross-validated models to generate predictions for each trial for each
1322 participant and then input the brain-predicted communal concern and second-order
1323 beliefs into our neural utility model (Eq 4. in main text). We estimated the θ values
1324 (i.e., weight on greed) and ϕ weighting parameters (i.e., relative trade-off between on
1325 communal concern and obligation) using the same procedure described in the
1326 behavioral computational modeling section by fitting the neural utility model directly
1327 to participant's reciprocity behavior by minimizing the SSE (Eq. 9).

1328

1329 As a benchmark for our neural utility model, we were interested in determining how
1330 well we could predict participant's reciprocity behavior directly from brain activity.
1331 We used the same training procedure described above, but predicted trial-to-trial
1332 reciprocity behavior using principal components regression separately for each
1333 participant. In theory, this should provide a theoretical upper bound of the best we
1334 should be able to predict reciprocity behavior using brain activity. If our neural utility
1335 model is close, then it means that we are able to predict reciprocity behavior using
1336 brain representations of communal concern and obligation as well as the optimal
1337 weighting of brain weights that can predict trial-to-trial reciprocity behavior. To
1338 determine the importance of the participant-specific model parameters, we ran a
1339 permutation test to determine how well we could predict reciprocity behavior for each
1340 participant using parameters from a randomly selected different participant. We ran
1341 5,000 permutations to generate a null distribution of average prediction accuracy after
1342 randomly shuffling the participant weights. The empirical p -value is the proportion of
1343 permutations that exceed our average observed correlation.

1344

1345 Finally, we were interested in evaluating how well we could estimate how much each
1346 participant had a relative preference for communal concern or obligation by
1347 computing the relative spatial alignment of their communal and obligation predictive

1348 spatial maps with their reciprocity predictive spatial map. We operationalized this
1349 relative pattern similarity as:

$$1350 \text{ relative pattern similarity} = \text{corr}(\vec{Obligation}_{map}, \vec{Reciprocity}_{map}) - \text{corr}(\vec{Communal}_{map}, \vec{Reciprocity}_{map})$$

1351 **Eq. 14**

1352 The intuition for this analysis is that if the optimal brain map for predicting a
1353 participant's decision is relatively more similar to their communal concern or
1354 obligation map, then we would expect that the participant cared more about that
1355 particular component of indebtedness during behavioral decision-making. For
1356 example, if a participant weights obligation more than communal concern during
1357 reciprocity (higher $1 - \phi$ estimated from the behavioral model), then the spatial
1358 similarity between their obligation brain pattern and the pattern that directly predicts
1359 their reciprocity behavior (reciprocity brain pattern) should be relatively higher
1360 compared to the spatial similarity between of their communal concern pattern and
1361 reciprocity brain pattern. We tested the correlation between this relative pattern
1362 similarity and the $(1 - \phi)$ parameters estimated by fitting the computational model (Eq.
1363 1) directly to the participants' behaviors.

1364 **Software**

1365 Behavioral data analyses were carried out in RStudio Version 1.1.383¹²¹ and
1366 IPython/Jupyter Notebook (Python 3.6.8)¹²², and was plotted using matplotlib¹²³, and
1367 seaborn 0.9.0 (<https://seaborn.pydata.org/index.html>). The fMRI data preprocessing
1368 and univariate analyses were conducted using Statistical Parametric Mapping
1369 software SPM12 (Wellcome Trust Department of Cognitive Neurology, London).
1370 Unless otherwise noted, all of fMRI multivariate analyses were performed with our
1371 open source Python NLTools package¹²⁰ version 0.3.14 (<https://nltools.org/>).

1372

1373 **Data availability**

1374 Behavioral data from all the three studies is available on github
1375 (https://github.com/xiaoxuepsy/Indebtedness_Gao2021). First and second level maps
1376 from the fMRI study is available on OSF (<https://osf.io/k8rxh/>). Raw imaging data is
1377 available from the corresponding author upon reasonable request.

1378

1379 **Code availability**

1380 The codes used in the current study are available on github
1381 (https://github.com/xiaoxuepsy/Indebtedness_Gao2021).

1382 **Acknowledgements**

1383 We thank Dr. Matthew Rushworth and Dr. Christian C. Ruff for their comments and
1384 suggestions on this article, Ms. Wan Wang, Mr. Shuaiqi Li and Mr. Sensen Song for
1385 their assistances in data collection, Ms. Yunyan Duan's for her advice in topic
1386 modeling, and Ms. Zhewen He for the preparation of the manuscript. This work was
1387 supported by National Natural Science Foundation of China (31900798, 31630034,
1388 71942001), National Basic Research Program of China (973 Program:
1389 2015CB856400), China Postdoctoral Science Foundation (2019M650008), the
1390 National Science Foundation of USA (CAREER 1848370), and the National Institute
1391 of Health (R01MH116026). We also acknowledge support from the Graduate School
1392 of Peking University to fund Dr. Gao's training at Dartmouth College.

Reference

- 1 Sherry Jr, J. F. Gift giving in anthropological perspective. *J. Consum. Res.* **10**, 157-168 (1983).
- 2 Carmichael, H. L. & MacLeod, W. B. Gift giving and the evolution of cooperation. *Int. Econ. Rev.*, 485-509 (1997).
- 3 Nowak, M. A. & Sigmund, K. Evolution of indirect reciprocity. *Nature* **437**, 1291 (2005).
- 4 Clark, M. S. & Mills, J. The difference between communal and exchange relationships: What it is and is not. *Pers. Soc. Psychol. Bull.* **19**, 684-691 (1993).
- 5 Clark, M. S. & Mills, J. R. A theory of communal (and exchange) relationships. in *Handbook of theories of social psychology, Vol. 2* 232-250 (Sage Publications Ltd, 2012).
- 6 Algoe, S. B. Find, remind, and bind: The functions of gratitude in everyday relationships. *Soc. Pers. Psychol. Compass* **6**, 455-469 (2012).
- 7 Algoe, S. B., Haidt, J. & Gable, S. L. Beyond reciprocity: gratitude and relationships in everyday life. *Emotion* **8**, 425 (2008).
- 8 Elfers, J. & Hlava, P. *The Spectrum of Gratitude Experience*. (Springer, 2016).
- 9 McCullough, M. E., Kilpatrick, S. D., Emmons, R. A. & Larson, D. B. Is gratitude a moral affect? *Psychol. Bull.* **127**, 249 (2001).
- 10 Trivers, R. L. The evolution of reciprocal altruism. *Q. Rev. Biol.* **46**, 35-57 (1971).
- 11 Neilson, W. S. The economics of favors. *J. Econ. Behav. Organ.* **39**, 387-397 (1999).
- 12 Akerlof, G. A. Labor contracts as partial gift exchange. *Q. J. Econ.* **97**, 543-569 (1982).
- 13 Greenberg, M. S. A theory of indebtedness. in *Social exchange* 3-26

- (Springer, 1980).
- 14 Greenberg, M. S. & Westcott, D. R. Indebtedness as a mediator of reactions to aid. *New directions in helping* **1**, 85-112 (1983).
- 15 Regan, D. T. Effects of a favor and liking on compliance. *J. Exp. Soc. Psychol.* **7**, 627-639 (1971).
- 16 Kolm, S.-C. *Reciprocity: An economics of social relations*. (Cambridge University Press, 2008).
- 17 Nadler, A. The other side of helping: Seeking and receiving help. in *The Oxford handbook of prosocial behavior*. *Oxford library of psychology*. 307-328 (Oxford University Press, 2015).
- 18 Fisher, J. D., Nadler, A. & Whitcher-Alagna, S. Recipient reactions to aid. *Psychol. Bull.* **91**, 27-54 (1982).
- 19 Fisher, J. *New Directions in Helping: Recipient reactions to aid*. Vol. 1 (Elsevier, 1983).
- 20 Nadler, A., Mayseless, O., Peri, N. & Chemerinski, A. Effects of opportunity to reciprocate and self-esteem on help-seeking behavior. *J. Pers.* **53**, 23-35 (1985).
- 21 Watkins, P. C., Scheer, J., Ovnicek, M. & Kolts, R. The debt of gratitude: Dissociating gratitude and indebtedness. *Cognition Emotion* **20**, 217-241 (2006).
- 22 Bal, A. Doctors and drug companies. *N. Engl. J. Med.* **352**, 733-734 (2005).
- 23 Malmendier, U. & Schmidt, K. You owe me. (National Bureau of Economic Research, 2012).
- 24 Fehr, E. & Gächter, S. Fairness and retaliation: The economics of reciprocity. *J. Econ. Perspect.* **14**, 159-181 (2000).
- 25 Gonzalez, B. & Chang, L. J. Computational models of mentalizing. (2019).
- 26 Falk, A., Fehr, E. & Fischbacher, U. On the nature of fair behavior. *Econ. Inq.* **41**, 20-26 (2003).

- 27 Sul, S., Guroglu, B., Crone, E. A. & Chang, L. J. Medial prefrontal cortical thinning mediates shifts in other-regarding preferences during adolescence. *Sci. Rep.* **7**, 8510 (2017).
- 28 Ellsworth, P. C. & Scherer, K. R. Appraisal processes in emotion. *Handbook of affective sciences* **572**, V595 (2003).
- 29 Frijda, N. H. The Place of Appraisal in Emotion. *Cognition Emotion* **7**, 357-387 (1993).
- 30 Frijda, N. H., Kuipers, P. & Ter Schure, E. Relations among emotion, appraisal, and emotional action readiness. *J. Pers. Soc. Psychol.* **57**, 212 (1989).
- 31 Lazarus, R. S. & Smith, C. A. Knowledge and appraisal in the cognition—emotion relationship. *Cognition Emotion* **2**, 281-300 (1988).
- 32 Scherer, K. R. Appraisal theory. (1999).
- 33 Smith, C. A. & Ellsworth, P. C. Patterns of cognitive appraisal in emotion. *J. Pers. Soc. Psychol.* **48**, 813 (1985).
- 34 Battigalli, P. & Dufwenberg, M. Dynamic psychological games. *J. Econ. Theory.* **144**, 1-35 (2009).
- 35 Battigalli, P., Corrao, R. & Dufwenberg, M. Incorporating belief-dependent motivation in games. *J. Econ. Behav. Organ.* (2019).
- 36 Geanakoplos, J., Pearce, D. & Stacchetti, E. Psychological games and sequential rationality. *Game. Econ. Behav.* **1**, 60-79 (1989).
- 37 Dufwenberg, M. & Kirchsteiger, G. A theory of sequential reciprocity. *Game. Econ. Behav.* **47**, 268-298 (2004).
- 38 Rabin, M. Incorporating fairness into game theory and economics. *Am. Econ. Rev.*, 1281-1302 (1993).
- 39 Chang, L. J. & Smith, A. Social emotions and psychological games. *Curr. Opin. Behav. Sci.* **5**, 133-140 (2015).
- 40 Benedict, R. Chrysanthemum and the Sword. *Patterns of Japanese Culture*,

- Cleveland, New York (The World Publishing Company) 1946. (1946).
- 41 Kotani, M. Expressing gratitude and indebtedness: Japanese speakers' use of "I'm sorry" in English conversation. *Res. Lang. Soc. Interac.* **35**, 39-72 (2002).
- 42 Naito, T. & Washizu, N. Note on cultural universals and variations of gratitude from an East Asian point of view. *J. Behav. Sci.* **10**, 1-8 (2015).
- 43 Washizu, N. & Naito, T. The emotions *sumanai*, gratitude, and indebtedness, and their relations to interpersonal orientation and psychological well-being among Japanese university students. *International Perspectives in Psychology: Research, Practice, Consultation* **4**, 209 (2015).
- 44 Baumeister, R. F., Stillwell, A. M. & Heatherton, T. F. Guilt: an interpersonal approach. *Psychol. Bull.* **115**, 243-267 (1994).
- 45 Le, B. M., Impett, E. A., Lemay Jr, E. P., Muise, A. & Tskhay, K. O. Communal motivation and well-being in interpersonal relationships: An integrative review and meta-analysis. *Psychol. Bull.* **144**, 1-25 (2018).
- 46 Naito, T. & Sakata, Y. Gratitude, Indebtedness, and Regret on Receiving a Friend's Favor in Japan. *Psychologia* **53**, 179-194 (2010).
- 47 Tsang, J. A. The effects of helper intention on gratitude and indebtedness. *Motiv. Emotion* **30**, 199-205 (2006).
- 48 Rotella, A., Sparks, A. M. & Barclay, P. Feelings of obligation are valuations of signaling-mediated social payoffs. *Behav. Brain Sci.* **43**, e85 (2020).
- 49 Tomasello, M. The Moral Psychology of Obligation. *Behav. Brain Sci.*, 1-33 (2019).
- 50 Beeler-Duden, S., Yucel, M. & Vaish, A. The role of affect in feelings of obligation. *Behav. Brain Sci.* **43**, e60 (2020).
- 51 Theriault, J. E., Young, L. & Barrett, L. F. The sense of should: A biologically-based framework for modeling social pressure. *Phys. Life Rev.* **36**, 100-136 (2021).
- 52 Fehr, E. & Schmidt, K. M. A theory of fairness, competition, and cooperation.

- Q. J. Econ.* **114**, 817-868 (1999).
- 53 Blei, D. M. & Lafferty, J. D. Dynamic topic models. in *Proceedings of the 23rd international conference on Machine learning*. 113-120 (ACM).
- 54 Chang, L. J., Smith, A., Dufwenberg, M. & Sanfey, A. G. Triangulating the neural, psychological, and economic bases of guilt aversion. *Neuron* **70**, 560-572 (2011).
- 55 van Baar, J. M., Chang, L. J. & Sanfey, A. G. The computational and neural substrates of moral strategies in social decision-making. *Nat. Commun.* **10** (2019).
- 56 Browne, M. W. & Cudeck, R. Alternative Ways of Assessing Model Fit. *Sociological Methods & Research* **21**, 230-258 (1992).
- 57 Hu, L. Evaluating model fit. *Structural equation modelling : concepts, issues and applications*, 76-99 (1995).
- 58 West, S. G., Taylor, A. B. & Wu, W. Model fit and model selection in structural equation modeling. in *Handbook of structural equation modeling*. 209-231 (The Guilford Press, 2012).
- 59 O'doherty, J. P., Hampton, A. & Kim, H. Model - based fMRI and its application to reward learning and decision making. *Ann. N. Y. Acad. Sci.* **1104**, 35-53 (2007).
- 60 Chang, L. J., Yarkoni, T., Khaw, M. W. & Sanfey, A. G. Decoding the role of the insula in human cognition: functional parcellation and large-scale reverse inference. *Cereb. Cortex* **23**, 739-749 (2013).
- 61 Yarkoni, T., Poldrack, R. A., Nichols, T. E., Van Essen, D. C. & Wager, T. D. Large-scale automated synthesis of human functional neuroimaging data. *Nat. Meth.* **8**, 665 (2011).
- 62 Fox, G. R., Kaplan, J., Damasio, H. & Damasio, A. Neural correlates of gratitude. *Front. psychol.* **6** (2015).
- 63 Yu, H., Cai, Q., Shen, B., Gao, X. & Zhou, X. Neural substrates and social

- consequences of interpersonal gratitude: Intention matters. *Emotion* **17**, 589-601 (2017).
- 64 Yu, H., Gao, X., Zhou, Y. & Zhou, X. Decomposing gratitude: representation and integration of cognitive antecedents of gratitude in the brain. *J. Neurosci.*, 2944-2917 (2018).
- 65 Cooper, J. C., Kreps, T. A., Wiebe, T., Pirkl, T. & Knutson, B. When giving is good: ventromedial prefrontal cortex activation for others' intentions. *Neuron* **67**, 511-521 (2010).
- 66 Ruff, C. C. & Fehr, E. The neurobiology of rewards and values in social decision making. *Nat. Rev. Neurosci.* **15**, 549 (2014).
- 67 Koban, L., Corradi-Dell'Acqua, C. & Vuilleumier, P. Integration of error agency and representation of others' pain in the anterior insula. *J. Cogn. Neurosci.* **25**, 258-272 (2013).
- 68 Yu, H., Hu, J., Hu, L. & Zhou, X. The voice of conscience: neural bases of interpersonal guilt and compensation. *Soc. Cogn. Affect. Neurosci.* **9**, 1150-1158 (2014).
- 69 Hampton, A. N., Bossaerts, P. & O'Doherty, J. P. Neural correlates of mentalizing-related computations during strategic interactions in humans. *Proc. Natl. Acad. Sci. U. S. A.* **105**, 6741-6746 (2008).
- 70 Van Overwalle, F. & Baetens, K. Understanding others' actions and goals by mirror and mentalizing systems: A meta-analysis. *Neuroimage* **48**, 564-584 (2009).
- 71 Woo, C.-W., Krishnan, A. & Wager, T. D. Cluster-extent based thresholding in fMRI analyses: Pitfalls and recommendations. *Neuroimage* **91**, 412-419 (2014).
- 72 Woo, C. W., Chang, L. J., Lindquist, M. A. & Wager, T. D. Building better biomarkers: brain models in translational neuroimaging. *Nat. Neurosci.* **20**, 365-377 (2017).

- 73 Chang, L. J., Gianaros, P. J., Manuck, S. B., Krishnan, A. & Wager, T. D. A sensitive and specific neural signature for picture-induced negative affect. *PLoS Biol.* **13**, e1002180 (2015).
- 74 Wager, T. D. *et al.* An fMRI-based neurologic signature of physical pain. *N. Engl. J. Med.* **368**, 1388-1397 (2013).
- 75 Chang, L. J. *et al.* Endogenous variation in ventromedial prefrontal cortex state dynamics during naturalistic viewing reflects affective experience. *Sci Adv* **7** (2021).
- 76 Woo, C.-W., Chang, L. J., Lindquist, M. A. & Wager, T. D. Building better biomarkers: brain models in translational neuroimaging. *Nat. Neurosci.* **20**, 365-377 (2017).
- 77 Kriegeskorte, N., Mur, M. & Bandettini, P. A. Representational similarity analysis-connecting the branches of systems neuroscience. *Front. Syst. Neurosci.* **2**, 4 (2008).
- 78 Mathews, M. A. & Green, J. D. Looking at me, appreciating you: Self-focused attention distinguishes between gratitude and indebtedness. *Cognition Emotion* **24**, 710-718 (2010).
- 79 Lench, H. C., Flores, S. A. & Bench, S. W. Discrete emotions predict changes in cognition, judgment, experience, behavior, and physiology: A meta-analysis of experimental emotion elicitation. *Psychol. Bull.* **137**, 834-855 (2011).
- 80 Lindquist, K. A., Siegel, E. H., Quigley, K. S. & Barrett, L. F. The hundred-year emotion war: are emotions natural kinds or psychological constructions? Comment on Lench, Flores, and Bench (2011). *Psychol. Bull.* **139**, 255-263 (2013).
- 81 Larsen, R. J. & Fredrickson, B. L. Measurement issues in emotion research. in *Well-being: The foundations of hedonic psychology.* 40-60 (Russell Sage Foundation, 1999).
- 82 Nisbett, R. E. & Wilson, T. D. Telling more than we can know: Verbal reports

- on mental processes. *Psychol. Rev.* **84**, 231-259 (1977).
- 83 Jolly, E. & Chang, L. J. The Flatland Fallacy: Moving Beyond Low-Dimensional Thinking. *Top. Cogn. Sci.* **11**, 433-454 (2019).
- 84 Chang, L. J. & Jolly, E. Emotions as computational signals of goal error. *The nature of emotion: Fundamental questions*, 343-348 (2018).
- 85 Xiang, T., Lohrenz, T. & Montague, P. R. Computational substrates of norms and their violations during social exchange. *J. Neurosci.* **33**, 1099-1108a (2013).
- 86 Gao, X. *et al.* Distinguishing neural correlates of context-dependent advantageous- and disadvantageous-inequity aversion. *Proc. Natl. Acad. Sci. U. S. A.* **115**, E7680-E7689 (2018).
- 87 Khalmetski, K., Ockenfels, A. & Werner, P. Surprising gifts: Theory and laboratory evidence. *J. Econ. Theory.* **159**, 163-208 (2015).
- 88 Battigalli, P., Dufwenberg, M. & Smith, A. Frustration and Anger in Games. (2015).
- 89 Chang, L. J. & Sanfey, A. G. Great expectations: neural computations underlying the use of social norms in decision-making. *Soc. Cogn. Affect. Neurosci.* **8**, 277-284 (2013).
- 90 Krajbich, I., Adolphs, R., Tranel, D., Denburg, N. L. & Camerer, C. F. Economic games quantify diminished sense of guilt in patients with damage to the prefrontal cortex. *J. Neurosci.* **29**, 2188-2192 (2009).
- 91 Smith, A., Bernheim, B. D., Camerer, C. & Rangel, A. Neural Activity Reveals Preferences Without Choices. *Nber Working Papers* **6**, 1-36 (2014).
- 92 Knutson, B., Rick, S., Wimmer, G. E., Prelec, D. & Loewenstein, G. Neural Predictors of Purchases. *Neuron* **53**, 147-156 (2007).
- 93 Haidt, J. The moral emotions. *Handbook of affective sciences* **11**, 852-870 (2003).
- 94 Fiske, A. P. The four elementary forms of sociality: Framework for a unified

- theory of social relations. *Psychol. Rev.* **99**, 689-723 (1992).
- 95 Rai, T. S. & Fiske, A. P. Moral psychology is relationship regulation: Moral motives for unity, hierarchy, equality, and proportionality. *Psychol. Rev.* **118**, 57-75 (2011).
- 96 Fiske, A. P. Socio-moral emotions motivate action to sustain relationships. *Self and Identity* **1**, 169-175 (2002).
- 97 van Baar, J. M., Klaassen, F. H., Ricci, F., Chang, L. J. & Sanfey, A. G. Stable distribution of reciprocity motives in a population. *Sci. Rep.* **10**, 18164 (2020).
- 98 Yu, H. *et al.* A Generalizable Multivariate Brain Pattern for Interpersonal Guilt. *Cereb. Cortex* (2020).
- 99 Inui, K., Tran, T. D., Hoshiyama, M. & Kakigi, R. Preferential stimulation of Adelta fibers by intra-epidermal needle electrode in humans. *Pain* **96**, 247-252 (2002).
- 100 Liu, Q. A novel Chinese text topic extraction method based on LDA. in *International Conference on Computer Science & Network Technology*. (2016).
- 101 Neto, J. L., Santos, A. D., Kaestner, C. A., Alexandre, N. & Santos, D. Document clustering and text summarization. (2000).
- 102 Salton, G. & Buckley, C. Term-weighting approaches in automatic text retrieval. *Inform. Process. Manag.* **24**, 513-523 (1988).
- 103 Blei, D. M., Ng, A. Y. & Jordan, M. I. Latent dirichlet allocation. *J. Mach. Learn. Res.* **3**, 993-1022 (2003).
- 104 Cowen, A. S. & Keltner, D. Self-report captures 27 distinct categories of emotion bridged by continuous gradients. *Proceedings of the National Academy of Sciences* **114**, E7900 (2017).
- 105 Blei, D. M., Ng, A. Y. & Jordan, M. I. Latent dirichlet allocation. *J. Mach. Learn. Res.* **3**, 993–1022 (2003).
- 106 Bates, D., Mächler, M., Bolker, B. & Walker, S. Fitting linear mixed-effects

- models using lme4. *arXiv preprint arXiv:1406.5823* (2014).
- 107 Blair, R. J. The neurobiology of psychopathic traits in youths. *Nat. Rev. Neurosci.* **14**, 786-799 (2013).
- 108 Fabrigar, L. R., Wegener, D. T., MacCallum, R. C. & Strahan, E. J. Evaluating the use of exploratory factor analysis in psychological research. *Psychol. Methods* **4**, 272-299 (1999).
- 109 Hair, J. F., Black, W. C., Babin, B. J., Anderson, R. E. & Tatham, R. *Multivariate data analysis*. (Uppersaddle River, 2006).
- 110 Tobias, S. & Carlson, J. E. BRIEF REPORT: BARTLETT'S TEST OF SPHERICITY AND CHANCE FINDINGS IN FACTOR ANALYSIS. *Multivar. Behav. Res.* **4**, 375-377 (1969).
- 111 Revelle, W. An overview of the psych package. *Department of Psychology Northwestern University*. Accessed on March 3, 2012 (2011).
- 112 Rosseel, Y. lavaan: An R package for structural equation modeling. *J. Stat. Softw.* **48**, 1-36 (2012).
- 113 Greenberg, M. S. & Shapiro, S. P. Indebtedness: An adverse aspect of asking for and receiving help. *Sociometry*, 290-301 (1971).
- 114 Fareri, D. S., Chang, L. J. & Delgado, M. R. Computational substrates of social value in interpersonal collaboration. *J. Neurosci.* **35**, 8170-8180 (2015).
- 115 Calhoun, V. D. *et al.* The impact of T1 versus EPI spatial normalization templates for fMRI data analyses. *Hum. Brain Mapp.* **38**, 5331-5342 (2017).
- 116 Barrett, L. F. & Satpute, A. B. Large-scale brain networks in affective and social neuroscience: towards an integrative functional architecture of the brain. *Curr. Opin. Neurobiol.* **23**, 361-372 (2013).
- 117 Adolphs, R. The social brain: neural basis of social knowledge. *Annu. Rev. Psychol.* **60**, 693-716 (2009).
- 118 Ruff, C. C. & Fehr, E. The neurobiology of rewards and values in social decision making. *Nat. Rev. Neurosci.* **15**, 549-562 (2014).

- 119 Haynes, J.-D. & Rees, G. Neuroimaging: decoding mental states from brain activity in humans. *Nat. Rev. Neurosci.* **7**, 523 (2006).
- 120 cosanlab/nltools: 0.3.11 v. 0.3.11 (Zenodo, 2018).
- 121 Racine, J. S. RStudio: A Platform - Independent IDE for R and Sweave. *J. Appl. Economet.* **27**, 167-172 (2012).
- 122 Pérez, F. & Granger, B. E. IPython: a system for interactive scientific computing. *Comput. Sci. Eng.* **9** (2007).
- 123 Hunter & John, D. Matplotlib: A 2D Graphics Environment. *Comput. Sci. Eng.* **9**, 90-95 (2007).