1    **Neural mechanisms of modulations of empathy and altruism by beliefs of others' pain**
2
3    **Taoyu Wu, Shihui Han**
4
5    **School of Psychological and Cognitive Sciences,**
6    **PKU-IDG/McGovern Institute for Brain Research,**
7    **Beijing Key Laboratory of Behavior and Mental Health,**
8    **Peking University, Beijing, China**
9
10
11
12
13   **Address correspondence to:**
14   **Shihui Han Ph.D. Prof.**
15   **School of Psychological and Cognitive Sciences**
16   **Peking University**
17   **52 Haidian Road, Beijing 100080, China**
18   **Email: shan@pku.edu.cn**
19   **Phone: (86) 10-6275-9138**
20   **Fax: (86) 10-6276-1081**
21
22

**Abstract**

Perceived cues signaling others' pain induce empathy which in turn motivates altruistic behavior toward those who appear suffering. This perception-emotion-behavior reactivity is the core of human altruism but does not always occur in real life situations. Here, by integrating behavioral and multimodal neuroimaging measures, we investigate neural mechanisms underlying modulations of empathy and altruistic behavior by beliefs of others' pain. We show evidence that lack of beliefs of others' pain reduces subjective estimation of others' painful feelings and decreases monetary donations to those who show pain expressions. Moreover, lack of beliefs of others' pain attenuates neural responses to their pain expressions within 200 ms after face onset and modulates neural responses to others' pain in the insular, post-central, and frontal cortices. Our findings suggest that beliefs of others' pain provide a cognitive basis of human empathy and altruism and unravel the intermediate neural mechanisms.

## Introduction

Aesop's fable 'The boy who cried wolf' tells a story that villagers run or do not run to help a shepherd boy who cries wolf depending on whether or not they believe that the boy's crying indicates his actual emotion and need. This story illustrates an important character of human altruistic behavior, that is, perceived cues signaling others' suffering drives us to do them a favor only when we believe that their suffering is true. Although this character of human altruism was documented over 2000 years ago in Aesop's fable and is widely observed in current human societies, its psychological and neural underpinnings have not been fully understood. The present study investigated how beliefs of others' pain (BOP) modulate human altruistic behavior independently of perceived cues signaling others' suffering and whether the modulation effect, if any, is mediated by changes in empathy for others' pain and relevant brain underpinnings.

Empathy refers to understanding and sharing of others' emotional states (Decety and Jackson, 2004) and has been proposed to provide a key motivation for altruistic behavior in both humans and animals (Batson et al., 2015; De Waal, 2008; Decety et al., 2016). Empathy can be induced by perceived cues signaling others' pain that activate neural responses in brain regions underlying sensorimotor resonance (e.g., the sensorimotor cortex), affective sharing (e.g., the anterior insula (AI) and anterior cingulate cortex (ACC)), and mental state inference/perspective taking (e.g., the medial prefrontal cortex (mPFC) and temporoparietal junction (TPJ)) (Singer et al., 2004; Jackson et al., 2005; Avenanti et al., 2005; Saarela et al., 2007; Fan and Han, 2008;

59   Shamay-Tsoory et al., 2009; Han et al., 2009; Sheng and Han, 2012; Fan et al., 2011;

60   Lamm et al., 2011; Zhou and Han, 2021). Neural responses to others' pain in the

61   empathy network and functional connectivity between its key hubs can predict motives

62   for subsequent altruistic actions (e.g., Hein et al., 2010; 2016; Mathur et al., 2010; Luo

63   et al., 2015). These brain imaging findings revealed neural mechanisms underlying the

64   perception-emotion-behavior reactivity (e.g., perceived pain-empathy-help) that occurs

65   often in everyday lives (Eisenberg et al., 2010; Hoffman, 2008; Penner et al., 2005).

66   However, empathic neural responses are influenced by multiple factors such as

67   perceptual features depicting others' pain (Gu and Han, 2007; Li and Han, 2019),

68   observers' perspectives and attention (Gu and Han, 2007; Li and Han, 2010; Jaunizux et

69   al., 2019), and perceived social relationships between observers and empathy targets

70   (Xu et al., 2009; Avenanti et al., 2010; Hein et al., 2010; Mathur et al., 2010; Sheng and

71   Han, 2012; Azevedo et al., 2013; Sheng et al., 2014; 2016; Han, 2018; Zhou and Han,

72   2021). What remains unclear is whether and how BOP modulates empathic brain

73   activity through which to further influence altruistic behavior. To address these issues is

74   crucial for understanding variations of empathy and altruism during complicated social

75   interactions as that illustrated in the Aesop's fable.

76       Beliefs refer to mental representations of something that is not immediately present

77   to the scenes but allows people to think beyond what is here and now (Fuentes, 2019).

78   Beliefs reflect organism's endorsement of a particular state of affairs as actual (McKay

79   and Dennett, 2009). Beliefs that best approximate reality enable the believers to act

effectively and maximize their survival (Fodor, 1985; Millikan, 1995). Previous

research has shown that beliefs affect multiple mental processes such as visual

awareness (Sterzer et al., 2008) and processing of emotions (Petrovic et al., 2005)

including experiences of pain (Wager et al., 2004; Colloca and Benedetti, 2005). The

function of beliefs is also manifested in increasing efficiency of neural processes

involved in decision making and goal setting (Garces and Finkel 2019; Régner et al.,

2019). Potential effects of beliefs on empathic neural responses were tested by

presenting participants with photographs showing pain inflicted by needle injections

into a hand that was believed to be or not to be anesthetized (Lamm et al., 2007).

Functional magnetic resonance imaging (fMRI) of brain activity suggested modulations

of insular responses to perceived pain by beliefs of anesthetization. However, the results

cannot be interpreted exclusively by BOP because the stimuli (i.e., needles) used to

induce beliefs of numbed and non-numbed hands were different. An ideal paradigm for

testing modulations of empathy by BOP independently of perceived cues signaling

others' pain should compare brain activities in response to identical stimuli under

different beliefs and enable researchers to test how BOP influences altruistic behavior.

In six behavioral, electroencephalography (EEG), and fMRI experiments, the

current study tested the hypothesis that BOP affects empathy and altruistic behavior by

modulating brain activity in response to others' pain. Specifically, we predicted that lack

of BOP may result in inhibition of altruistic behavior by decreasing empathy and its

underlying brain activity. Our behavioral, EEG, and fMRI experiments were designed

101   based on the common beliefs that patients show pain expressions to manifest their

102   actual feelings of pain whereas pain expressions performed by actors/actresses do not

103   indicate their actual emotional states. To examine BOP effects on empathy, we

104   experimentally manipulated BOP by asking participants to learn and remember different

105   identities (i.e., patient or actor/actress) of a set of neutral faces during a learning

106   procedure. Thereafter, we measured self-reports of others' pain and own unpleasantness

107   from the participants when they viewed learned faces with pain or neutral expressions.

108   During EEG/fMRI recording the participants were asked to discriminate patient or

109   actor/actress identities of faces with pain or neutral expressions. We compared

110   self-reports of others' feelings and brain activities related to pain (vs. neutral)

111   expressions of patients' faces with those related to actors/actresses' faces. If perception

112   of patients' pain expressions implicitly activates BOP whereas perception of

113   actors/actresses' pain expressions does not activate BOP, we expected that lack of BOP

114   (i.e., to compare actors/actresses vs. patients) would reduce self-report of empathy,

115   empathic brain activity, and altruistic behavior. We further predicted that BOP effects

116   on altruistic behavior might be mediated by decreased empathy and empathic brain

117   activity due to lack of BOP.

118       Similar to previous research (Jackson et al., 2005; Fan and Han, 2008; Hein et al.,

119   2010; Mathur et al., 2010; Sheng and Han, 2012), we adopted both subjective and

120   objective estimations of empathy for others' pain. Subjective estimation of empathy for

121   pain depends on collection of self-reports of others' painful feelings and ones' own

122  unpleasantness when viewing others' suffering (e.g., Bieri et al., 1990; Jackson et al.,

123  2005; Lamm et al., 2007; Fan and Han, 2008; Sheng and Han, 2012). Objective

124  estimation of empathy for pain relies on recording of brain activities, using fMRI or

125  EEG, that differentially respond to painful versus non-painful stimuli applied to others

126  (e.g., Singer et al., 2004; Jackson et al., 2005; Gu and Han, 2007; Fan and Han, 2008;

127  Hein et al., 2010) or to others' faces with pain versus neutral expressions (Botvinick et

128  al., 2005; Saarela et al., 2007; Han et al., 2009; Sheng and Han, 2012). Brain responses

129  to perceived non-painful stimuli applied to others or neutral expressions were also

130  collected to control empathy-unrelated perceptual or motor processes. fMRI studies

131  revealed greater activations in the ACC, AI, and sensorimotor cortices in response to

132  painful compared to non-painful stimuli applied to others (e.g., Singer et al., 2004;

133  Jackson et al., 2005; Gu and Han, 2007; Hein et al., 2010, see Lamm et al., 2011; Fan et

134  al., 2011, for review). EEG studies showed that event-related potentials (ERPs) in

135  response to perceived painful stimulations applied to others' body parts elicited neural

136  responses that differentiated between painful and neutral stimuli over the frontal region

137  as early as 140 ms after stimulus onset (Fan and Han, 2008; see Coll, 2018 for review).

138  Moreover, the mean ERP amplitudes at 140–180 ms predicted self report of others' pain

139  and ones' own unpleasantness (Fan and Han, 2008).

140      Particularly related to the current work are neuroimaging findings that compared

141  brain responses to pain versus neutral expressions. fMRI studies found that viewing

142  video clips (Botvinick et al., 2005) or pictures (Sheng et al., 2014) showing faces with

143     pain versus neutral expressions or viewing photos of faces of patients who were

144     suffering from provoked pain versus chronic pain (Saarela et al., 2007) induced

145     activations in the ACC, AI, and inferior parietal cortex. Moreover, the cortical areas

146     activated by facial expressions of pain were also engaged by the first-hand experience

147     of pain evoked by thermal stimulation (Botvinick et al., 2005). Moreover, the strengths

148     of AI activations during observation of others' pain were correlated with subjective

149     feelings of others' pain (Saarela et al., 2007). ERP studies found that neural responses to

150     pain expressions occurred as early as 130 ms after face onset over the frontal/central

151     regions as indexed by the increased amplitude of a positive component at 128–188 ms

152     (P2) in response to pain compared neutral expressions (Sheng and Han, 2012; Sheng et

153     al., 2013; 2016; Han et al., 2016; Li and Han, 2019). In addition, the P2 amplitudes in

154     response to others' pain expressions positively predicted subjective feelings of own

155     unpleasantness induced by others' pain and self-reports of one's own empathy traits

156     (Sheng and Han, 2012). In addition, source estimation of the P2 component in response

157     to others' pain expressions suggested a possible origin in the ACC. Taken together,

158     these brain imaging findings suggest effective subjective and objective measures of

159     empathy (i.e., understanding and sharing of others' pain) that are suitable for

160     investigation of neural mechanisms underlying modulations of empathy and altruism by

161     BOP.

162        In Experiment 1 we randomly assigned patient or actor/actress identities to faces to

163     test how experimentally manipulated BOP associated with face identities caused

164  changes in empathy (i.e., subjective evaluation of others' pain) and altruistic behavior

165  (i.e., monetary donations). We predicted that lack of BOP related to actors/actresses (vs.

166  patients) would result in reduced empathy and altruistic behavior. In Experiment 2,

167  based on the common belief that an effective medical treatment reduces a patient's pain,

168  we tested whether decreasing BOP due to knowledge of effective medical treatments of

169  patients also reduced empathy and altruistic behavior.

170       In Experiments 3 and 4 we investigated whether BOP modulates empathic brain

171  activity by recording EEG signals in response to pain or neutral expressions of faces

172  with patient or actor/actress identities. Brain activities related empathy were quantified

173  by comparing neural responses to pain versus neutral expressions to exclude neural

174  processes of facial structures, social attributes (e.g., gender), and other

175  empathy-unrelated information. Given previous findings that the P2 amplitude

176  increased to pain compared to neutral expressions and was associated with self-report of

177  sharing of others' pain (Sheng and Han, 2012; Sheng et al., 2013; 2016; Han et al., 2016;

178  Li and Han, 2019), we focused on how the P2 amplitude in response to pain (vs. neutral)

179  expressions was modulated by facial identities (i.e., patient or actor/actress) that link to

180  different beliefs (i.e., patients' pain expressions manifest their actual feelings whereas

181  actors/actresses' pain expressions do not). Our ERP results showed evidence that

182  actor/actress compared to patient identities of faces decreased the empathic neural

183  responses (i.e., P2 amplitudes in response to pain (vs. neutral) expressions) within 200

184  ms post-stimulus. In Experiment 5 we further revealed behavioral and EEG evidence

185   that neural responses to pain expressions of faces mediate BOP effects on empathy and

186   monetary donations.

187        In Experiment 6 we employed fMRI to examine brain regions in which blood

188   oxygen level dependent (BOLD) signals are modulated by BOP. We examined BOLD

189   responses to faces that had either patient or actor/actress identities, received

190   painful/non-painful stimulations, and showed pain or neutral expressions. fMRI results

191   allowed us to test whether empathic neural responses in the cognitive (i.e., the dorsal

192   mPFC and TPJ, Völlm et al., 2006; Schnell et al., 2011; also see Lamm et al., 2011; Fan

193   et al., 2011; Shamay-Tsoory, 2011), sensorimotor/affective (i.e., the ACC, insula, and

194   sensorimotor cortex, Jackson et al., 2006; Singer et al., 2004; Avenanti et al., 2005), or

195   both nodes of the empathic neural network would be modulated by BOP that was

196   manipulated by assigning different identities (i.e., patient or actor/actress) to empathy

197   targets. In addition, we examined whether neural responses in the empathic network

198   would be able to predict variations of subjective feelings of others' pain due to lack of

199   BOP.

200   .    Together, our behavioral and brain imaging results showed consistent evidence that

201   lack of BOP or decreasing BOP resulted in reduced empathy and altruistic behavior.

202   Our findings suggest that BOP may provide a cognitive basis for human empathy and

203   altruism and uncover intermediate brain mechanisms by which BOP influences empathy

204   and altruistic behavior.

205   **Results**

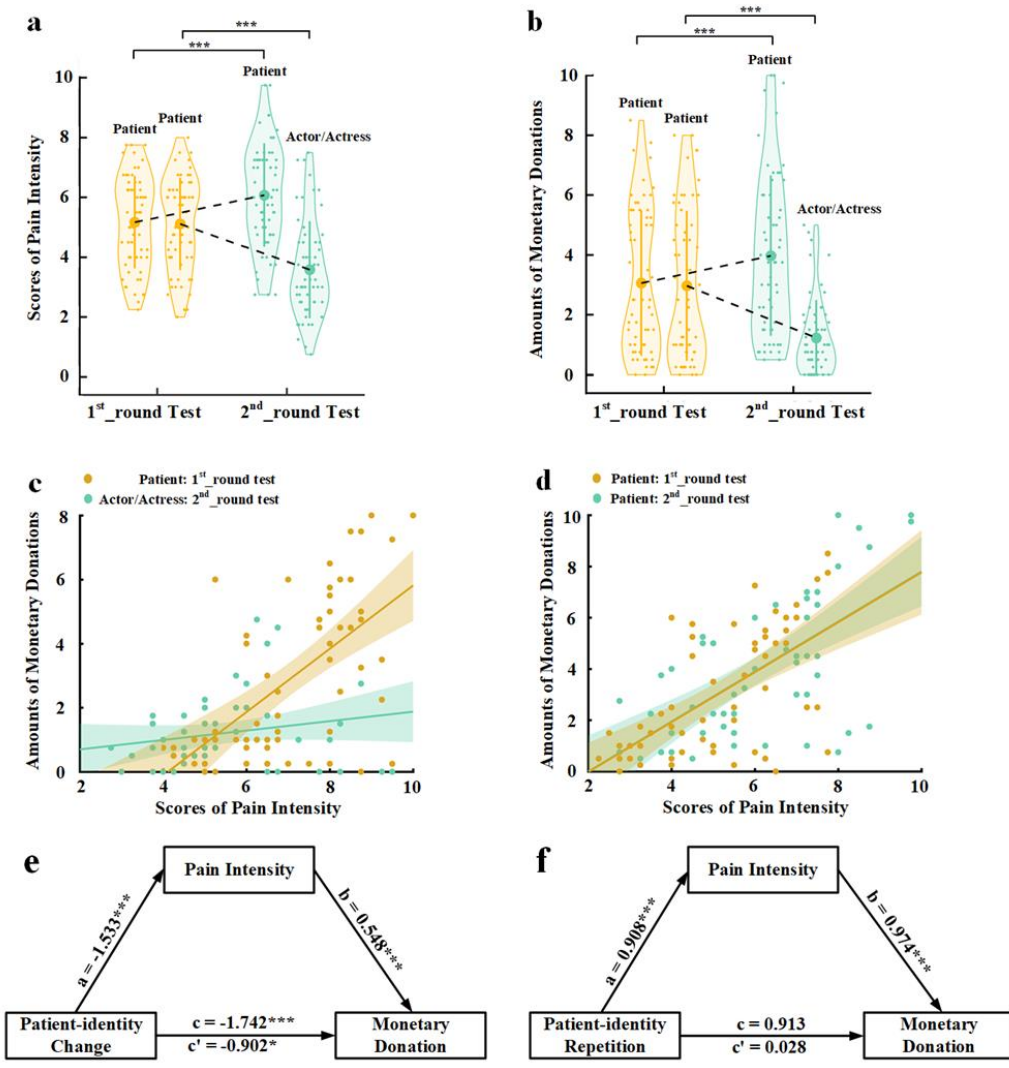**Experiment 1: Lack of BOP reduces subjective estimation of empathy and**

**altruistic behavior**

In Experiment 1 we tested the predictions that lack of BOP decreases empathy and

altruistic behavior by experimentally manipulating individuals' BOP. We presented

participants (N = 60) with photos of faces of 16 models (half males) with pain

expressions (see Methods for details). The participants were informed that these photos

were taken from patients who suffered from a disease. In the $1^{st}$_round test the

participants were shown with each photo and asked to report perceived pain intensity of

each patient by rating on a Likert-type scale (0 = not painful at all; 10 = extremely

painful). This rating task was adopted from previous research (Bieri et al., 1990;

Jackson et al., 2005; Lamm et al., 2007; Fan and Han, 2008; Sheng and Han, 2012) to

assess the participants' understanding of others' pain feeling $—$ a key component of

empathy. Thereafter, the participants were invited to donate money to the patient in the

photo by selecting an amount from an extra bonus payment for their participation (0 to

10 points, 1 point = ¥0.2) as a measure of altruistic behavior. The participants were

informed that the amount of one of their donation decisions would be selected randomly

and endowed to a charity organization to help those who suffered from the same disease.

After the $1^{st}$_round test the participants were asked to perform a 5-minute

calculation task to clean their memory of performances during the $1^{st}$_round test. The

participants were then informed that this experiment actually tested their ability to

recognize facial expressions and the photos were actually taken from 8 patients and 8

227    actors/actresses. We expected that identity changes from patients to actors/actresses

228    would decrease BOP because patients' pain expressions reflect their actual emotional

229    states whereas pain expressions performed by actors/actresses do not indicate an actual

230    painful state. The participants were then asked to perform the 2nd_round test in which

231    each photo was presented again with patient or actor/actress identity indicated by a

232    word (i.e., patient, actor, or actress) below the photo. The participants had to perform

233    the same pain intensity rating and donation tasks as those in the 1st_round test. The

234    participants were told that an amount of money would be finally selected from their

235    2nd_round donation decisions and presented to the same charity organization after the

236    study.

237        The mean rating scores of pain intensity and amounts of monetary donations were

238    subject to repeated-measures analyses of variance (ANOVAs) of Test Phase (1st_round

239    vs. 2nd_round test) $\times$ Identity Change (patient-identity change (patient to actor/actress)

240    vs. patient-identity repetition (patient to patient)) as independent within-subjects

241    variables. As expected, the results revealed that patient-identity change or

242    patient-identity repetition produced opposite effects on both perceived pain intensity

243    and amounts of monetary donations, as indicated by significant interactions of Test

244    Phase $\times$ Identity Change ($F(1,59) = 123.476$ and $60.638$, $ps < 0.001$, $\eta_p^2 = 0.677$ and

245    $0.507$, $90\%$ CI $= (0.555, 0.747)$ and $(0.351, 0.611)$, Fig. 1a and 1b). Specifically,

246    patient-identity change (i.e., from patients to actors/actresses) significantly reduced

247    perceived pain intensity and amounts of monetary donations in the 2nd_round (vs.

248     1$^{st}$_round) test (F(1,59) = 82.664 and 34.542, ps < 0.001, $\eta_p^2$ = 0.584 and 0.369, 90% CI

249     = (0.440, 0.673) and (0.207, 0.495)). By contrast, patient-identity repetition

250     significantly increased both perceived pain intensity and monetary donations in the

251     2$^{nd}$_round (vs. 1$^{st}$_round) test (F(1,59) = 36.060 and 27.457, ps < 0.001, $\eta_p^2$ = 0.379 and

252     0.318, 90% CI = (0.216, 0.503) and (0.159, 0.449)). These results suggest that our

253     manipulations of BOP caused reliable changes in subjective evaluation of others' pain

254     and related monetary donations in opposite directions. Interestingly, to some degree

255     rather than not at all, the participants reported pain and donated to faces with

256     actor/actress identity in the 2$^{nd}$_round test, suggesting that lack of BOP did not fully

257     eliminate empathy and altruistic behavior toward those who showed pain expressions.

**Fig. 1. Behavioral results in Experiment 1.** (a) Mean rating scores of pain intensity in the 1st_ and 2nd_round tests. (b) Mean amounts of monetary donations in the 1st_ and 2nd_round tests. Shown are group means (large dots), standard deviation (bars), measures of each individual participant (small dots), and distribution (violin shape) in (a) and (b). (c) The associations between rating scores of pain intensity and amounts of monetary donations for patients in the 1st_round test and for actors/actresses in the 2nd_round test. (d) The associations between rating scores of pain intensity and amounts of monetary donations for patients in both the 1st_ and 2nd_round tests. (e) Rating scores of pain intensity partially mediate the relationship between patient-identity change and reduced monetary donations. (f) Rating scores of pain intensity mediate the relationship between patient-identity repetition and increased monetary donations. The online version of this article includes the following source data for Figure 1: Figure 1-Source data 1.

To investigate whether perceived pain intensity mediated the relationships between experimentally manipulated BOP and monetary donations, we first conducted Pearson

14

274     correlation analyses of the relationship between empathy and altruism. The results

275     showed that the rating scores of pain intensity of faces whose identities changed from

276     patient in the $1^{st}$\_round test to actor/actress in the $2^{nd}$\_round test significantly predicted

277     the amount of monetary donations in the $1^{st}$\_round but not in the $2^{nd}$\_round test ($r =$

278     0.608 and 0.187, $p < 0.001$ and $p = 0.152$, 95% CI = (0.422, 0.776) and (-0.069, 0.435),

279     all results were FDR-corrected, Fig. 1c). The rating scores of pain intensity also

280     significantly predicted the amount of monetary donations for faces whose patient

281     identities did not change in the $1^{st}$\_round and $2^{nd}$\_round tests ($r = 0.619$ and 0.628, $ps <$

282     0.001, 95% CI = (0.449, 0.776) and (0.417, 0.775), Fig. 1d). We conducted mediation

283     analyses to further test an intermediate role of empathy between BOP and altruistic

284     behavior (see Methods). The first mediation analysis showed that rating scores of pain

285     intensity partially mediated the relationship between patient-identity change and

286     reduced amount of monetary donations (direct effect: $c' = -0.902$, $t(118) = -2.468$, $p =$

287     0.015, 95% CI = (-1.626, -0.178); indirect effect: $a \times b = -0.839$, 95% CI = (-1.455,

288     -0.374), Fig. 1e, see Supplementary file 1 for statistical details). The second mediation

289     analysis showed evidence that the rating scores of pain intensity also mediated the

290     relationship between patient-identity repetition and increased amount of monetary

291     donations (direct effect: $c' = 0.028$, $t(118) = 0.072$, $p = 0.943$, 95% CI = (-0.727, 0.782),

292     indirect effect: $a \times b = 0.885$, 95% CI = (0.314, 1.563), Fig. 1f, see Supplementary file 2

293     for statistical details). These results indicate a key functional role of BOP in altruistic

294    behavior and suggest changes in subjective evaluation of others' pain as an intermediate

295    mechanism underlying the effect of BOP on monetary donations.

296    **Experiment 2: Intrinsic BOP predicts subjective estimation of empathy and**

297    **altruistic behavior**

298         In Experiment 1 BOP was manipulated by randomly assigning patient or

299    actor/actress identities to faces and the results showed that experimentally manipulated

300    BOP changes caused variations of empathy and altruistic behavior. In Experiment 2 we

301    further investigated whether an individual's intrinsic BOP (i.e., various representations

302    of actual emotional states of different faces with pain expressions) can predict empathy

303    and altruistic behavior across different faces. Moreover, as a replication, we tested

304    whether changing the participants' intrinsic BOP causes changes in empathy and

305    altruistic behavior in directions similar to those observed in Experiment 1. In addition,

306    we assessed whether changing intrinsic BOP modulated sharing of others' pain 一

307    another key component of empathy (Bieri et al., 1990; Jackson et al., 2005; Lamm et al.,

308    2007; Fan and Han, 2008; Sheng and Han, 2012). Finally, we tested whether BOP

309    induced emotional sharing mediates the relationship between BOP and altruistic
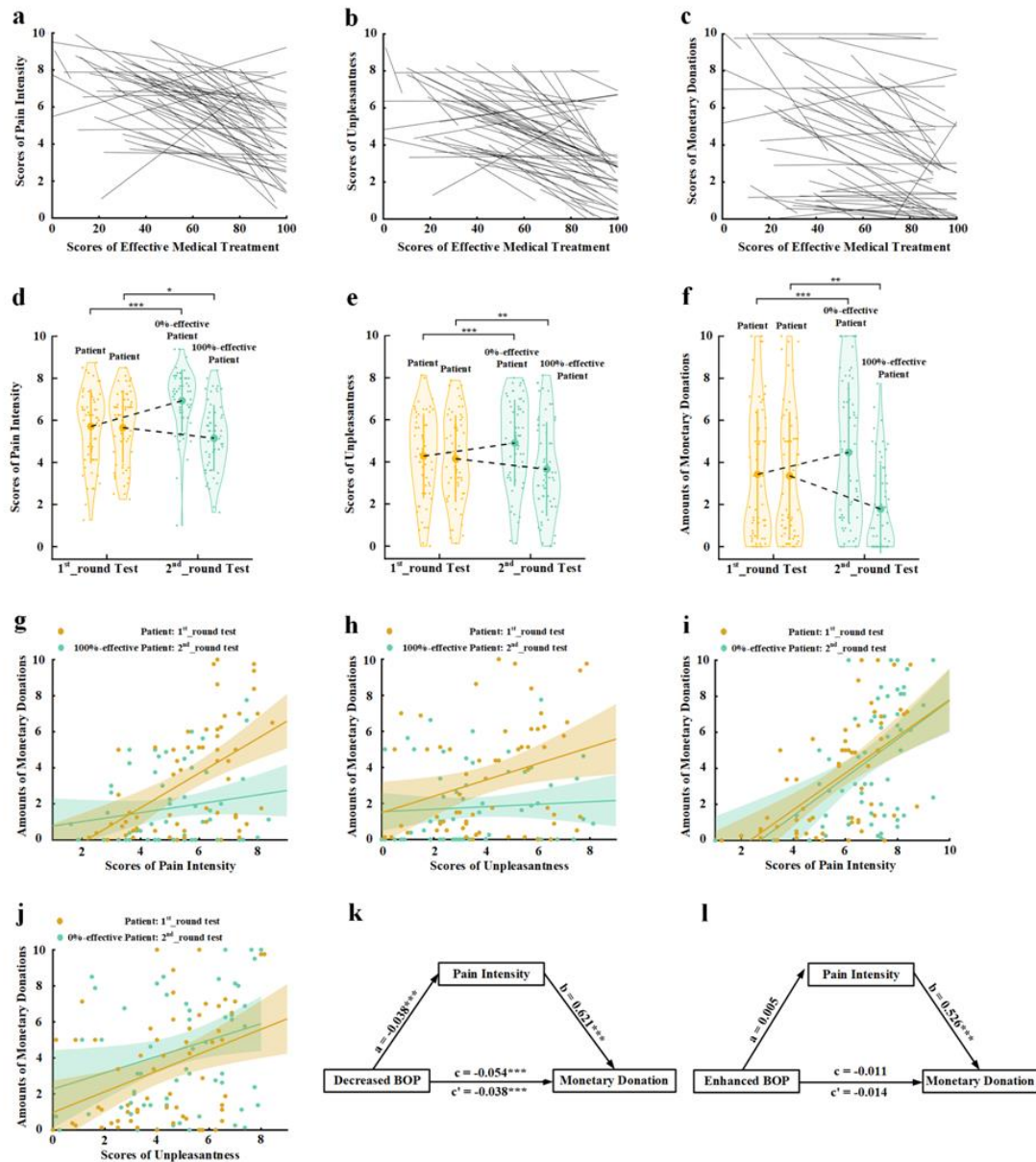
310    behavior.

311         To address these issues, we tested an independent sample (N = 60) using the stimuli

312    and procedure that were the same as those in Experiment 1 except the following. In the

313    1st_round test the participants were informed that they were to be shown with photos

314    with pain expressions taken from patients who suffered from a disease and received a

315    medical treatment. After the presentation of each photo the participants were asked to

316    estimate, based on perceived pain expression of each face, how effective they believed

317    the medical treatment was for each patient by rating on a Likert-type scale (0 = no effect

318    or 0% effective, 100 = fully effective or 100% effective). The rating scores were used to

319    estimate the participants' intrinsic BOP of each face with a higher rating score

320    (indicating more effective treatment) corresponding to a weaker BOP because a more

321    effective medical treatment reduces a patient's pain to a greater degree. In addition to

322    rating pain intensity of each face, the participants were asked to report how unpleasant

323    they were feeling when viewing each photo by rating on a Likert-type scale (0 = not

324    unpleasant at all, 10 = extremely unpleasant). The unpleasantness rating was performed

325    to assess emotional sharing of others' pain. In the $2^{nd}$_round test the participants were

326    told that the medical treatment was actually fully effective for half patients but had no

327    effect for the others. Each photo was then presented again with information that the

328    medical treatment applied to the patient was 100% effective (to decrease the participants'

329    beliefs of the patients' painful states) or 0% effective (to enhance the participants'

330    beliefs of the patients' painful states). Thereafter, the participants were asked to perform

331    the rating tasks and to make monetary donation decisions, similar to those in the

332    $1^{st}$_round test.

333       To assess whether individuals' intrinsic BOP predicted their empathy and altruistic

334    behavior across different target faces, we conducted Pearson correlation analyses of the

335    relationships between intrinsic BOP as indexed by the rating score of treatment

336   effectiveness and empathy rating scores/amounts of monetary donations across the

337   sixteen models in the 1st_round test in each participant. The correlation coefficients

338   were then transformed to Fisher's z values that were further compared with zero.

339   One-sample t-tests revealed that the z values were significantly smaller than zero

340   (correlations between intrinsic BOP and pain intensity/unpleasantness/monetary

341   donation: mean ± s.d. = -0.631 ± 0.531, -0.643 ± 0.524 and -0.469 ± 0.529; $t(59)$ =

342   -9.213, -9.501 and -6.875; $ps < 0.001$; Cohen's $d$ = 1.188, 1.227 and 0.887; 95% CI =

343   (-0.768, -0.494), (-0.778, -0.507), and (-0.606, -0.333), Fig. 2a-c), suggesting that a

344   larger score of treatment effectiveness (i.e., a weaker intrinsic BOP related to a face)

345   predicted weaker empathy and less monetary donations relate to that face. These results

346   provide evidence for associations between intrinsic BOP and empathy/altruism.

**Fig. 2. Behavioral results in Experiment 2.** The relationships between intrinsic BOP
(indexed by the rating score of effective medical treatments) and scores of pain intensity
(a), own unpleasantness (b), and monetary donations (c), respectively, across the sixteen
models in the 1st_round test in each participant. The regression line of each participant
is plotted in (a), (b), and (c). (d-f) Mean rating scores of pain intensity, own
unpleasantness, and monetary donations in the $1^{st}$_ and $2^{nd}$_round tests. (g) The
associations between rating scores of pain intensity and amounts of monetary donations
for patients in the $1^{st}$_round test and for 100%-effective patients in the $2^{nd}$_round tests
across all the participants. (h) The associations between rating scores of own
unpleasantness and amounts of monetary donations for patients in the $1^{st}$_round test and
for-100% effective patients in the $2^{nd}$_round tests across all the participants. (i) The
associations between rating scores of pain intensity and amounts of monetary donations
for patients in the $1^{st}$_round test and for 0%-effective patients in the $2^{nd}$_round tests

361    across all the participants. (j) The associations between rating scores of own

362    unpleasantness and amounts of monetary donations for patients in the 1$^{st}$_round test and

363    for 0%-effective patients in the 2$^{nd}$_round tests across all the participants. (k) Rating

364    scores of pain intensity change partially mediate the relationship between decreased

365    BOP and changes in monetary donations. (l) Rating scores of pain intensity change fail

366    to mediate the relationship between enhanced BOP and changes in monetary donations.

367    Shown are group means (large dots), standard deviation (bars), measures of each

368    individual participant (small dots), and distribution (violin shape) in (d), (e), and (f).

369    The online version of this article includes the following source data for Figure 2: Figure

370    2-Source data 1.

371        Next, we tested whether decreased (or increased) BOP also predicts changes in

372    empathy/altruistic behavior across different target faces for each participant. To do this,

373    we calculated belief changes (decreased BOP: 100%-effective minus the participants'

374    initial estimation; enhanced BOP: the participants' initial estimation minus

375    0%-effective), empathy changes (rating scores in the 2$^{nd}$_round vs. 1$^{st}$_round test), and

376    changes in altruistic behavior (the amount of monetary donation in the 2$^{nd}$_round vs.

377    1$^{st}$_round test) related to each model in each participant. Similarly, we conducted

378    Pearson correlation analyses to examine associations between changes in beliefs and

379    empathy/donation for decreased-BOP patients and enhanced-BOP patients, respectively,

380    in each participant. The correlation coefficients were then transformed to Fisher's z

381    values that were further compared with zero. One-sample t-tests showed that the z

382    values were significantly smaller than zero for decreased-BOP patients (the correlation

383    between changes in belief and pain intensity: z-value (mean ± s.d.) = -0.304 ± 0.370;

384    $t(59) = -6.352$, $p < 0.001$; Cohen's d = 0.822; 95% CI = (-0.400, -0.208); the correlation

385    between changes in belief and unpleasantness: z-value (mean ± s.d.) = -0.277 ± 0.455;

386    $t(59) = -4.706$, $p < 0.001$; Cohen's d = 0.609; 95% CI = (-0.394, -0.159); the correlation

387    between changes in belief and monetary donation: z-value (mean ± s.d.) = -0.236 ±

388    0.410; t(59) = -4.465, p < 0.001; Cohen's d = 0.576; 95% CI = (-0.342, -0.130)). These

389    results suggest that a greater decrease of BOP related to a face predicted greater reduced

390    empathy and less monetary donations. By contrast, one-sample t-tests showed that the z

391    values were significantly larger than zero for enhanced-BOP patients (the correlation

392    between changes in belief and pain intensity: z-value (mean ± s.d.) = 0.286 ± 0.488;

393    t(59) = 4.533, p < 0.001; Cohen's d = 0.586; 95% CI = (0.160, 0.412); the correlation

394    between changes in belief and unpleasantness: z-value (mean ± s.d.) = 0.227 ± 0.470;

395    t(59) = 3.735, p < 0.001; Cohen's d = 0.483; 95% CI = (0.105, 0.348); the correlation

396    between changes in belief and monetary donation: z-value (mean ± s.d.) = 0.162 ± 0.538;

397    t(59) = 2.332, p = 0.023; Cohen's d = 0.301; 95% CI = (0.023, 0.301)). These results

398    suggest that a greater increase of BOP predicted greater increased empathy and more

399    monetary donations across individual empathy targets. These results provide evidence

400    for associations between changes in BOP and empathy/altruism across different faces

401    for each participant.

402         To test whether the results in Experiment 2 replicated those in Experiment 1, we

403    conducted ANOVAs of the mean empathy scores and amounts of monetary donations

404    with Test Phase (1st vs. 2nd_round) and Belief Change (initial self-rated effectiveness to

405    informed 0%-effectiveness vs. initial self-rated effectiveness to informed

406    100%-effectiveness) as independent within-subjects variables. The results showed that

407    decreasing internal BOP (i.e., for 100% effective target faces) resulted in lower

408  subjective evaluation of others' pain and one's own unpleasantness and less monetary

409  donations in the $2^{nd}$_ vs. $1^{st}$_round tests, whereas enhancing BOP (i.e., for 0% effective

410  target faces) produced opposite effects (Fig. 2d-f, see Supplementary file 3 for statistical

411  details). These results replicated those in Experiment 1 and provided further evidence

412  that changing BOP resulted in variations of empathy and altruistic behavior.

413       Pearson correlations analyses of the mean rating scores in the $1^{st}$_round and

414  $2^{nd}$_round tests across the participants showed that, for '100%-effective' patients, the

415  $1^{st}$_round but not the $2^{nd}$_round rating scores of empathy significantly predicted the

416  amount of monetary donations (Pain intensity rating: r = 0.530 and 0.184, p < 0.001 and

417  p = 0.159, 95% CI = (0.334, 0.698) and (-0.057, 0.425), Unpleasantness rating: r =

418  0.307 and 0.074, p = 0.017 and p = 0.576, 95% CI = (0.046, 0.541) and (-0.199, 0.358),

419  Fig. 2g and 2h). For '0%-effective' patients, however, both the $1^{st}$_round and $2^{nd}$_round

420  rating scores of empathy significantly predicted the amount of monetary donations (Pain

421  intensity rating: r = 0.582 and 0.476, ps < 0.001, 95% CI = (0.415, 0.725) and (0.287,

422  0.638); Unpleasantness rating: r = 0.373 and 0.280, p = 0.006 and 0.04, 95% CI =

423  (0.096, 0.590) and (0.011, 0.511), Fig. 2i and 2j).

424       Furthermore, the results of mediation analyses showed that rating scores of pain

425  intensity partially mediated the relationship between decreased BOP (i.e., for

426  '100%-effective' patients) and monetary donations (direct effect: c' = -0.038, t(58) =

427  -3.657, p < 0.001, 95% CI = (-0.059, 0.017); indirect effect: a×b = -0.016, 95% CI =

428  (-0.027, -0.005) , Fig. 2k, see Supplementary file 4 for statistical details). However,

429   rating scores of unpleasantness did not mediate the relationship between decreased BOP

430   and monetary donations (indirect effect: $a \times b$ = -0.002, 95% CI = (-0.009, 0.003)).

431   Neither pain intensity nor unpleasantness ratings mediated the relationship between

432   enhanced BOP (i.e., for '0%-effective' patients) and monetary donations (indirect effect:

433   a*b = 0.003 and -0.002, 95% CI = (-0.009, 0.013) and (-0.007, 0.004) , Fig. 2l, see

434   Supplementary files 5, 6, and 7 for statistical details). These behaviorsl results suggest

435   that decreased BOP influences altruistic decisions possibly via modulations of the

436   cognitive component of empathy (i.e., understanding others' pain) rather than the

437   affective component of empathy (i.e., sharing others' pain).

438   **Experiment 3: Lack of BOP decreased empathic brain activity**
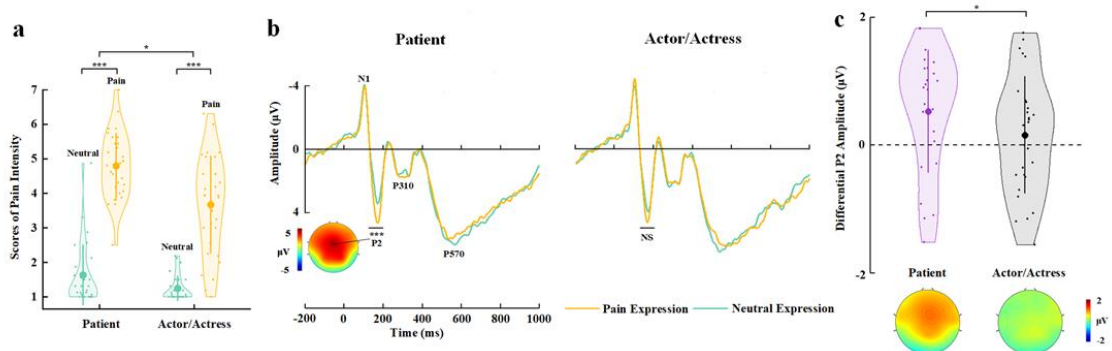
439   Experiments 1 and 2 showed evidence that self-report measures of empathy for pain

440   were affected by BOP. In Experiment 3 we further investigated whether and how

441   changing BOP modulates brain activity in response to perceived cues signaling others'

442   pain as an objective estimation of empathy. If BOP provides a basis of empathy of

443   others' pain, lack of BOP should reduce empathic neural responses to visual stimuli

444   signaling others' pain. We tested this assumption by recording EEG to faces of 16

445   models from an independent sample (N = 30). The participants were first presented with

446   these faces with neutral expressions and were informed that these photos were taken

447   from 8 patients who suffered from a disease and from 8 actors/actresses. The

448   participants were asked to remember patient or actor/actress identity of each neutral face

449   and had to pass a memory test with a 100% recognition accuracy. Thereafter, the

450    participants were informed that they would be presented with photos of these faces with

451    either neutral or pain expressions, and photos of pain expressions were taken from the

452    patients who were suffering from the disease or from the actors/actresses who imitated

453    patients' pain. The participants were asked to make judgments on identity of each face

454    (i.e., patient vs. actor/actress) with a neutral or pain expression by pressing one of two

455    buttons while EEG was recorded. After EEG recording, the participants were asked to

456    rate pain intensity of each face with a pain or neutral expression on a Likert-type scale

457    (0 = not painful at all; 7 = extremely painful) and to what degree they believed in the

458    identity of each face with a pain expression on a 15-point Likert-type scale (-7 =

459    extremely believed as an actor/actress, 0 = not sure, 7 = extremely believed as a patient).

460    Because the same set of stimuli were perceived as patients or actors/actresses across the

461    participants, modulations of brain activity in response to pain expressions only reflected

462    the effects of BOP concomitant with the face identity (i.e., real pain for patients but fake

463    pain for actors/actresses).

464        The participants reported a positive mean belief score corresponding to faces with a

465    patient identity (2.496 $\pm$ 2.51) but a negative mean belief score corresponding to faces

466    with an actors/actresses identity (-2.210 $\pm$ 3.25) (t(29) = 4.932, p < 0.001, Cohen's d

467    = 0.900, 95% CI = (2.755, 6.658)), suggesting successes of our manipulations of face

468    identities. An ANOVA of the mean rating scores of pain intensity with Identity (patient

469    vs. actor/actress) and Expression (pain vs. neutral) as within-subject variables revealed a

470    significant Identity $\times$ Expression interaction (F(1,29) = 4.905, p = 0.035, $\eta_p^2$ = 0.145, 90%

471    CI = (0.006, 0.330), Fig. 3a), suggesting greater subjective feelings of pain intensity for

472    faces with patient compared to actor/actress identity. Moreover, a larger score of belief

473    of patient identities significantly predicted greater subjective feelings of pain intensity

474    related to patients' pain (vs. neutral) expressions (r = 0.384, p = 0.036, 95% CI = (0.074,

475    0.627)), whereas there was no significant association between belief scores and

476    subjective feelings of pain intensity related to actors/actresses' pain (vs. neutral)

477    expressions (r = 0.264, p = 0.159, 95% CI = (-0.162, 0.605)). These results provide

478    further evidence for a link between BOP and empathy for patients' pain.

479



480    **Fig. 3. EEG results of Experiment 3. (a)** Mean rating scores of pain intensity to pain
481    versus neutral expressions of faces with patient or actor/actress identities. (b) ERPs to
482    faces with patient or actor/actress identities at frontal electrodes. The voltage
483    topography shows the scalp distribution of the P2 amplitude with the maximum over the
484    central/frontal region. (c) Mean differential P2 amplitudes to pain versus neutral
485    expressions of faces with patient or actor/actress identities. The voltage topographies
486    illustrate the scalp distribution of the P2 difference waves to pain versus neutral
487    expressions of faces with patient or actor/actress identities, respectively. Shown are
488    group means (large dots), standard deviation (bars), measures of each individual
489    participant (small dots), and distribution (violin shape) in (a) and (c). The online version
490    of this article includes the following source data for Figure 3: Figure 3-Source data 1.

491        The participants responded to face identities with high accuracies during EEG

492    recording (>81% across all conditions, see Supplementary file 8 for details). ERPs to

493    face stimuli in Experiment 3 were characterized by an early negative activity at 95–115

494 ms (N1) and a positive activity at 175–195 ms (P2) at the frontal/central regions, which

495 were followed by two positive activities at 280-340 ms (P310) over the parietal region

496 and 500–700 ms (P570) over the frontal area (Fig. 3b). Previous ERP studies have

497 shown that empathic neural responses to pain expressions are characterized by an

498 increased P2 amplitude and the P2 amplitude to pain (vs. neutral) expressions predicts

499 self-report of affective sharing (Sheng and Han, 2012; Sheng et al., 2016; Luo et al.,

500 2018; Li and Han, 2019). Therefore, our ERP data analyses focused on whether BOP

501 modulates the P2 amplitude to pain (vs. neutral) expressions given the previous ERP

502 findings. ANOVAs of the P2 amplitudes with Identity (patient vs. actor/actress) and

503 Expression (pain vs. neutral) as within-subject variables revealed a significant Identity $\times$

504 Expression interaction ($F(1,29) = 7.490$, $p = 0.010$, $\eta_p^2 = 0.205$, 90% CI = (0.029,

505 0.391), see Supplementary file 9 for statistical details). Simple effect analyses verified

506 significantly greater P2 amplitudes to pain versus neutral expressions of patients' faces

507 ($F(1,29) = 18.059$, $p < 0.001$, $\eta_p^2 = 0.384$, 90% CI = (0.150, 0.546)), whereas the P2

508 amplitude did not differ significantly between pain and neutral expressions of

509 actors/actresses' faces ($F(1,29) = 0.334$, $p = 0.568$, $\eta_p^2 = 0.011$, 90% CI = (0.000, 0.135),

510 Fig. 3b and 3c). We further conducted Bayes factor analyses to examine the null effect

511 of pain expressions on the P2 amplitudes to actors/actresses' faces. The Bayes factor

512 represents the ratio of the likelihood of the data fitting under the alternative hypothesis

513 versus the likelihood of fitting under the null hypothesis. The results showed a Bayes

514 factor of 0.227 which provided further evidence for the null hypothesis. The results

515   indicate that, while the effect of pain (vs. neutral) expression on the P2 amplitudes to

516   patients' faces was similar to our previous findings that the P2 amplitudes increased to

517   pain (vs. neutral) expressions of face without patient identities (Sheng and Han, 2012;

518   Sheng et al., 2016), the P2 amplitude was less sensitive to pain versus neutral

519   expressions of faces with actor/actress identities. This finding indicate that lack of BOP

520   significantly weakens early empathic neural responses to others' pain within 200 ms

521   after stimulus onset.

522   **Experiment 4: BOP is necessary for modulations of empathic brain activity**
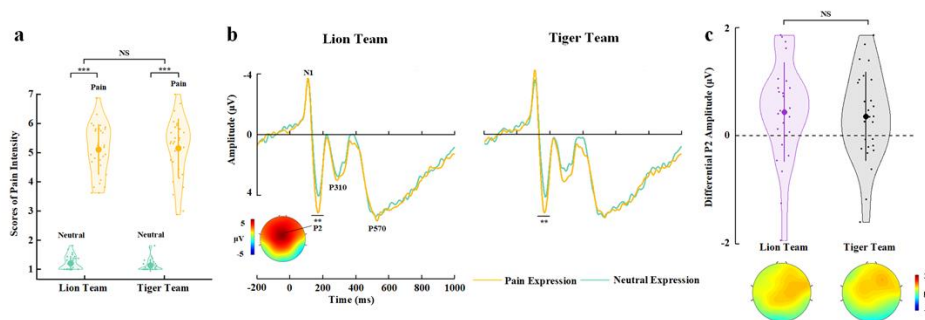
523   The learning and EEG recording procedures in Experiment 3 consisted of multiple

524   processes, including learning, memory and recognition of face identities, assignment to

525   different social groups (e.g., patient or actor groups), etc. The results of Experiment 3

526   left an open question of whether these processes, even without BOP changes induced

527   through these processes, would be sufficient to result in modulations of the P2

528   amplitude in response to pain (vs. neutral) expressions of faces with different identities.

529   In Experiment 4 we addressed this issue using the same learning and identity

530   recognition procedures as those in Experiment 3 except that the participants in

531   Experiment 4 had to learn and recognize identities of faces of two baseball teams and

532   that there is no prior difference in BOP associated with individual faces from the two

533   baseball teams. If the processes involved in the learning and reorganization procedures

534   rather than the difference in BOP were sufficient for modulations of the P2 amplitude in

535   response to pain (vs. neutral) expressions of faces, we would expect similar P2

536    modulations in Experiments 4 and 3. Otherwise, if the difference in BOP produced

537    during the learning procedure was necessary for the modulation of empathic neural

538    responses, we would not expect modulations of the P2 amplitude in response to pain (vs.

539    neutral) expressions in Experiment 4.

540        We clarified these predictions in an independent sample (N = 30) in Experiment 4.

541    We employed the stimuli and procedure that were the same as those in Experiment 3

542    except that, during the learning phase, the participants were informed that the 16 models

543    were from two baseball teams (half from a Tiger team and half from a Lion team) and

544    they suffered from a disease. After the participants had remembered team identity of

545    each neutral face in a procedure similar to that in Experiment 3, they performed identity

546    (i.e., Tiger vs. Lion team) judgments on the faces with neutral or pain expressions

547    during EEG recording. This manipulation built team identities should not influence

548    self-report and EEG estimation of empathy because the Tiger/Lion team identities did

549    not bring any difference in BOP between pain expressions of faces from the two teams.

550        The participants responded to face identities with high accuracies during EEG

551    recording (> 79% across all conditions). Rating scores of pain intensity did not differ

552    significantly between faces from the two teams ($F(1,29) = 1.608$, $p = 0.215$, $\eta_p^2 = 0.053$,

553    90% CI = (0, 0.216), , Bayes factors = 0.261, Fig. 4a, see Supplementary file 10 for

554    details). ANOVAs of the mean P2 amplitudes over the frontal electrodes revealed a

555    significant main effect of facial expression ($F(1,29) = 12.182$, $P = 0.002$, $\eta_p^2 = 0.296$, 90%

556    CI = (0.081, 0.473), Fig. 4b and 4c, see Supplementary file 11 for details), as the P2

557  amplitude was enlarged by pain compared to neutral expressions. However, this effect

558  did not differ significantly between faces from the two teams (F(1,29) = 0.040, P =

559  0.843, $\eta_p^2$ = 0.001, 90% CI = (0, 0.053), Bayes factors = 0.258). The null interaction

560  effect on either self-report of empathy and the P2 amplitudes to pain (vs. neutral)

561  expressions in Experiment 4 was not simply due to an underpowered sample size

562  because the same sample size in Experiment 3 revealed reliable BOP effects on

563  self-report and EEG (i.e., the P2 amplitude) estimation of empathy. Together, the results

564  in Experiments 3 and 4 suggest a key role of BOP, but not other cognitive processes

565  involved in the experimental manipulations, in modulations of neural responses to

566  others' pain.



**Fig. 4. EEG results of Experiment 4. (a)** Mean rating scores of pain intensity to pain
versus neutral expressions of faces with Lion Team or Tiger Team identities. (b) ERPs
to faces with Lion/Tiger team identities at frontal electrodes. The voltage topography
shows the scalp distribution of the P2 amplitude with the maximum over the
central/frontal region. (c) Mean differential P2 amplitudes to pain versus neutral
expressions of faces with Lion/Tiger Team identities. The voltage topographies
illustrate the scalp distribution of the P2 difference waves to pain versus neutral
expressions of faces with the Lion/Tiger Team identities, respectively. Shown are group
means (large dots), standard deviation (bars), measures of each individual participant
(small dots), and distribution (violin shape) in (a) and (c). The online version of this
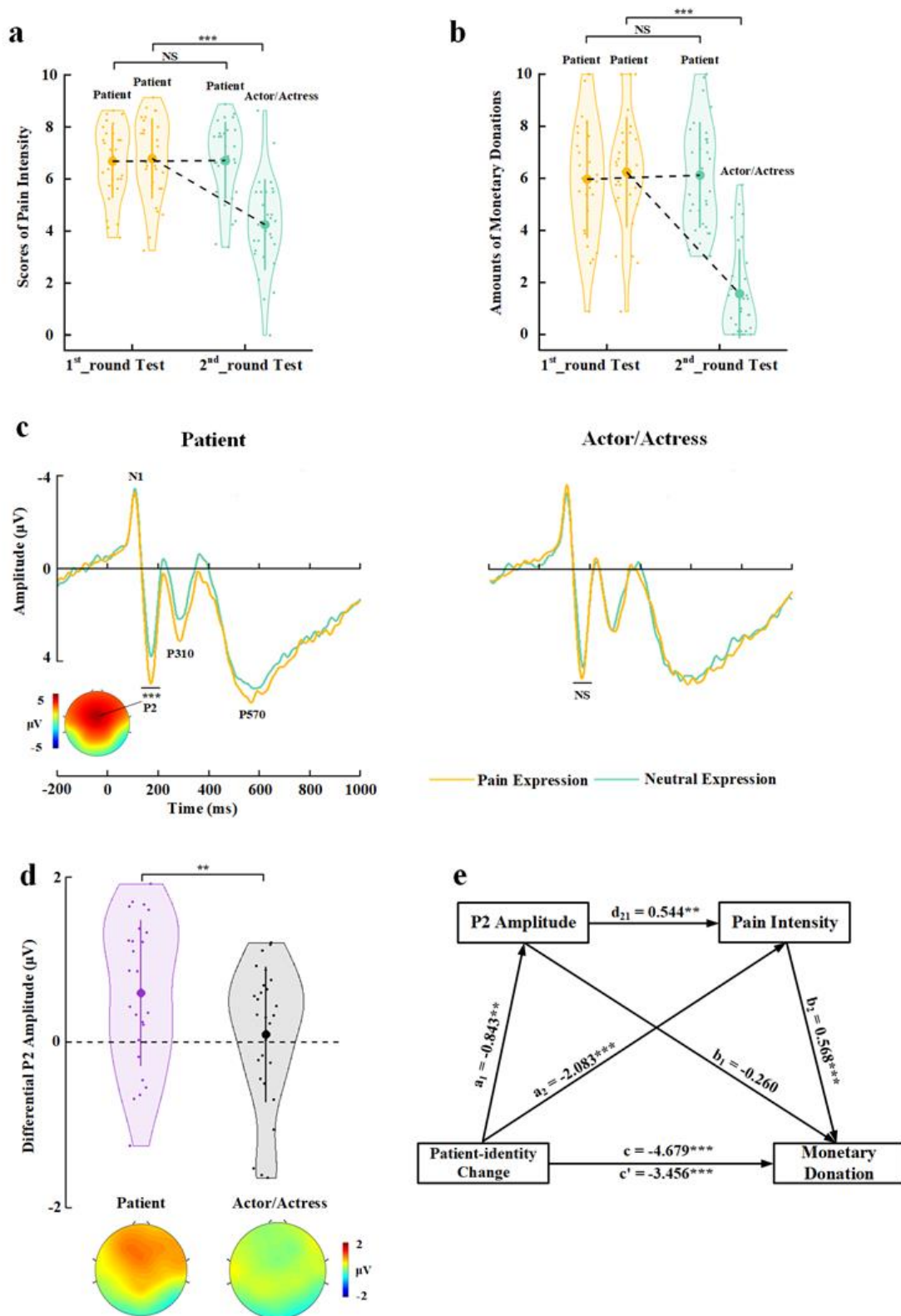article includes the following source data for Figure 4: Figure 4-Source data 1.

**Experiment 5: Empathic brain activity mediates relationships between BOP and**

**empathy/altruistic behavior**

581       Given that Experiments 1 to 4 showed consistent evidence for BOP effects on

582    subjective feelings of others' pain, altruistic behavior, and empathic neural responses, in

583    Experiment 5, we further examined whether BOP-induced changes in empathic brain

584    activity plays a mediator role in the pathway from belief changes to altered subjective

585    feelings of others' pain and altruistic decisions. To this end, we conducted two-session

586    tests of an independent sample (N = 30). In the first session we employed the stimuli

587    and procedure that were identical to those in Experiment 1 to assess BOP effects on

588    empathy and altruistic behavior. In the second session we recorded EEG from the

589    participants using the same stimuli and procedure as those in Experiment 3 to examine

590    BOP effects on empathic neural responses. BOP-induced changes in empathic brain

591    activity, rating scores of pain intensity, and amounts of monetary donations recorded in

592    the two-session tests were then subject to mediation analyses.

593       To assure the participants' beliefs about patient and actor/actress identities of

594    perceived faces, after EEG recording, we asked the participants to complete an implicit

595    association test (IAT) (Greenwald et al., 1998) that measured reaction times to faces

596    with patient and actor/actress identities and words related to patients and

597    actors/actresses (see Methods). The D score was then calculated based on response

598    times (Greenwald et al., 2003) to assess implicit associations between patient and

599    actor/actress faces and the relevant words. One-sample t-test revealed that the D score

600    was significantly larger than zero ($0.929 \pm 0.418$, $t(29) = 12.178$, $p < 0.001$, Cohen's d =

601    2.223, 95% CI = (0.773, 1.085)), suggesting that patient faces were more strongly

602    associated with patient relevant words whereas actor/actress faces were more strongly

603    associated with actor/actress relevant words. The results indicate successful belief

604    manipulations during the two-session tests.

605    　　The behavioral results in the first-session test replicated the findings of Experiment

606    1. In particular, decreasing BOP (i.e., changing patient identity in the $1^{st}$_round test to

607    actor/actress identity in the $2^{nd}$_round test) significantly reduced self-report of others'

608    pain and monetary donations (Test Phase × Identity Change interactions on rating

609    scores of pain intensity and amounts of monetary donations: ($F(1,29) = 59.654$ and

610    129.696, $ps < 0.001$, $\eta_p^2 = 0.673$ and 0.817, 90% CI = (0.479, 0.764) and (0.694, 0.868);

611    Effects of patient-to-actor/actress identity change on rating scores of pain intensity and

612    amounts of monetary donations: $F(1,29) = 58.196$ and 180.022, $ps < 0.001$, $\eta_p^2 = 0.667$

613    and 0.861, 90% CI = (0.472, 0.760) and (0.765, 0.900), Fig. 5a and 5b). However,

614    patient-identity repetition failed to significantly increase rating scores of pain intensity

615    and amounts of monetary donations ($F(1,29) = 0.016$ and 0.209, $p = 0.901$ and 0.651,

616    $\eta_p^2 = 0.001$ and 0.007, 90% CI = (0, 0.022) and (0, 0.119)), possibly due to ceiling

617    effects of our measures in the participants (i.e., larger mean rating scores of pain

618    intensity and mean amounts of monetary donations in the $1^{st}$_round test in Experiment 5

619    than in Experiment 1).

Fig. 5. Behavioral and EEG results of Experiment 5. (a) Mean rating scores of pain intensity in the $1^{st}$_ and $2^{nd}$_round tests. (b) Mean amounts of monetary donations in the $1^{st}$_ and $2^{nd}$_round tests. (c) ERPs to faces with patient or actor/actress identities at frontal electrodes. The voltage topography shows the scalp distribution of the P2 amplitude with the maximum over the central/frontal region. (d) Mean differential P2

626 amplitudes to pain versus neutral expressions of faces with patient or actor/actress
627 identities. The voltage topographies illustrate the scalp distribution of the P2 difference
628 waves to pain versus neutral expressions of faces with patient or actor/actress identities,
629 respectively. (e) Illustration of the serial mediation model of the relationship between
630 decreased BOP and changes in monetary donations. Shown are group means (large
631 dots), standard deviation (bars), measures of each individual participant (small dots),
632 and distribution (violin shape) in (a), (b) and (d). The online version of this article
633 includes the following source data for Figure 5: Figure 5-Source data 1.

634 The participants responded to face identities with high accuracies during EEG

635 recording (> 83% across all conditions). The EEG results replicated those in Experiment

636 3 by showing significantly deceased P2 amplitudes to pain (vs. neutral) expressions of

637 actor/actress compared to patient faces (Identity × Expression interaction: $F(1,29) =$

638 $9.494$, $p = 0.004$, $\eta_p^2 = 0.247$, 90% CI = (0.050, 0.429), Fig. 5c and 5d, see

639 Supplementary file 12 for statistical details). Simple effect analyses verified

640 significantly greater P2 amplitudes to pain vs. neutral expressions for patients' faces

641 $(F(1,29) = 17.409$, $p < 0.001$, $\eta_p^2 = 0.375$, 90% CI = (0.142, 0.539)) but not for faces of

642 actors/actresses $(F(1,29) = 0.270$, $p = 0.607$, $\eta_p^2 = 0.009$, 90% CI = (0, 0.127), Bayes

643 factor = 0.220). These behavioral and EEG results are consistent with those in

644 Experiments 1 and 3 and provide repeated evidence for BOP effects on subjective

645 feelings of others' pain, altruistic behavior, and empathic brain activity in the same

646 sample.

647 Next, we tested a serial mediation model of the relationship between decreased

648 BOP (i.e., identity change from patient to actor/actress) and changes in monetary

649 donations with two mediator variables including empathic neural responses (as indexed

650 by the differential P2 amplitude to pain versus neutral expressions) and changes in

33

651     subjective feelings of others' pain (as indexed by differential rating scores of pain

652     intensity) (see Methods for details). This model includes three paths: (1) the indirect

653     effect of patient-identity change on monetary donation via the P2 amplitude ($a_1 \times b_1$ =

654     0.219, 95% CI = (-0.141, 0.745)); (2) the indirect effect of patient-identity change on

655     monetary donation via pain intensity ($a_2 \times b_2$ = -1.182, 95% CI = (-2.048, -0.510)); (3)

656     the indirect effect of patient-identity change on monetary donation via P2 amplitude $\times$

657     pain intensity ($a_1 \times d_{21} \times b_2$ = -0.261, 95% CI = (-0.584, -0.059), Fig. 5e, see

658     Supplementary file 13 for statistical details). The total indirect effect of patient-identity

659     change on the monetary donation after controlling all indirect effect was c' = -1.223, 95%

660     CI = (-2.145, -0.400), which explained 26.14% variance of total effect of

661     patient-identity change on monetary donation. The effect sizes of the indirect path (2)

662     and (3) were 25.26% and 5.58%, respectively, indicating that subjective feelings of

663     others' pain mediated the association between patient-identity change and reduced

664     monetary donations. Moreover, this mediator role was partially mediated by BOP

665     induced variations of empathic brain activity in response to others' pain expressions.

666     Together, the results of these mediation analyses suggest a pathway from changes in

667     BOP to varied empathic brain activity and changes in subjective report of empathy for

668     other's pain (i.e., the degree of perceived pain in others), which further accounted for

669     BOP-induced changes in monetary donations.

670     **Experiment 6: Neural structures underlying BOP effects on empathy**
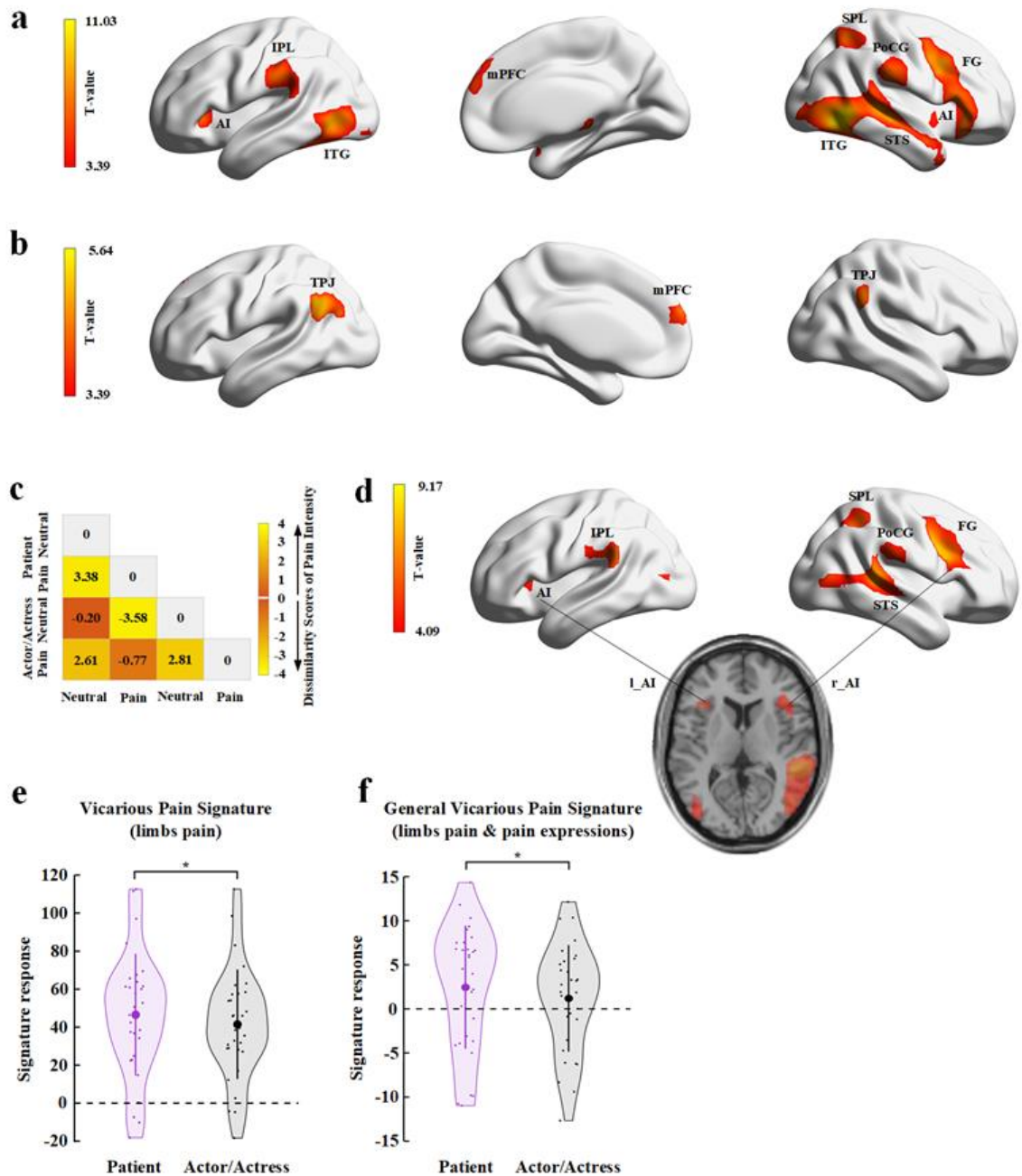
While our EEG results revealed evidence for modulations of empathic neural

responses by BOP, neural structures underlying these modulation effects remain unclear.

In particular, it is unknown whether brain responses underlying cognitive and affective

components of empathy are similarly sensitive to the influence of BOP. Therefore, in

Experiment 6, we used fMRI to record BOLD signals from an independent sample (N =

31) to examine neural architectures in which empathic activities are modulated by BOP.

Similarly, the participants were first shown with photos of neutral faces of 20 models

and had to remember their patient (10 models) or actor/actress (10 models) identities.

After the participants had performed 100% correct in a memory task to recognize the

models' identities, they were scanned using fMRI when viewing video clips of the

models whose faces received painful (needle penetration) stimulation and showed pain

expressions or received non-painful (cotton swab touch) stimulation and showed neutral

expressions, similar to those used in the previous studies (Han et al., 2009; Luo et al.,

2014; Han et al., 2017). Before scanning the participants were informed that these video

clips were recorded from 10 patients who were receiving medical treatment and 10

actors/actresses who practiced to imitate patients' pain expressions. The participants

responded to face identity (patient vs. actor/actress) of each model after viewing each

video clip by pressing one of two buttons with high accuracies (> 80% across all

conditions, see Supplementary file 14 for details).

After fMRI scanning the participants were presented with each video clip again and

had to rate the model's pain intensity and their own unpleasantness. The participants

692    were also asked to rate the degree to which they believed in the models' patient or

693    actor/actress identities in painful video clips on a 15-point Likert-type scale (-7 =

694    extremely believed as an actor/actress, 0 = not sure, 7 = extremely believed as a patient)

695    (see Method, Supplementary file 14 for results). The mean rating scores confirmed

696    significant differences in beliefs of patient and actors/actresses identities ($2.776 \pm 3.20$

697    vs. $-4.890 \pm 1.44$; $t(30) = 10.526$, $p < 0.001$, Cohen's $d = 1.890$, 95% CI = (6.178,

698    9.153)), indicating successful identity manipulations.

699      We first localized empathic neural responses by conducting a whole-brain analysis

700    of BOLD responses to perceived painful versus non-painful stimuli applied to targets

701    (collapsed faces with patient and actor/actress identities). This analysis revealed

702    significant activations in the cognitive, affective, and sensorimotor nodes of the

703    empathy network, including the bilateral anterior insula/inferior frontal cortex (MNI

704    peak coordinates x/y/z = -45/17/-5 and 45/26/-8), bilateral inferior and superior

705    temporal gyri (-48/-70/-2 and 51/-58/-5), mPFC (3/56/25), left inferior parietal lobe

706    (-63/-25/31), right superior parietal lobe (30/-58/55), and right post-central gyrus and

707    posterior insula (58/-25/26, Fig. 6a; all activations were identified using a combined

708    threshold of voxel level $p < 0.001$, uncorrected, and cluster level $p < 0.05$, FWE

709    corrected). These brain activations are similar to those observed in previous research

710    (e.g., Luo et al., 2014). To examine brain activity engaged in representing facial

711    identities independent of perceived painful stimulation and pain expressions, we

712    conducted a whole-brain analysis of the contrast of the stimuli showing non-painful

713    stimulations to patient versus actor/actress. This analysis showed significant activations

714    in the mPFC (-6/59/25) and bilateral TPJ (-54/-58/28 and 57/-67/31, Fig. 6b, all

715    activations were identified using a combined threshold of voxel level p < 0.001,

716    uncorrected, and cluster level p < 0.05, FWE corrected).

**Fig. 6. fMRI results of Experiment 6.** (a) Brain activations in response to perceived
painful (vs. non-painful) stimuli applied to targets (collapsed faces with patient and

720  actor/actress identities). (b) Brain activations in response to non-painful stimuli to
721  patients compared to actors/actresses. (c) Illustration of the behavioral dissimilarity
722  matrix derived from the rating scores of pain intensity across all participants. Each cell
723  in the dissimilarity matrix represents the mean difference in rating scores of pain
724  intensity between each pair of conditions. (d) Brain activations that were correlated with
725  the behavioral dissimilarity matrix revealed in the searchlight RSA. (e) Illustration of
726  the vicarious pain signature (defined by response to perceived noxious stimulation of
727  body limbs) responses to patients' and to actors/actresses' pain. (f) Illustration of the
728  general vicarious signature (defined by response to perceived noxious stimulation of
729  body limbs  and painful facial expressions) responses to patients' and actors/actresses'
730  pain. AI = Anterior Insula; IPL = Inferior Parietal Lobe; ITG = Inferior Temporal Gyrus;
731  mPFC = medial Prefrontal Cortex; SPL = Superior Parietal Lobe; PoCG = Post-Central
732  Gyrus; FG = Frontal Gyrus; STS = Superior Temporal Sulcus; MFC = Middle Frontal
733  Cortex; TPJ = Temporoparietal Junction. The online version of this article includes the
734  following source data for Figure 6: Figure 6-Source data 1.

735      We conducted a whole-brain univariate analysis to examine the interaction effect

736  (patient vs. actor x pain vs. neutral) on brain activities in response to video clips but did

737  not find a significant effect. Therefore, we further conducted multivariate analyses of

738  BOLD signals to assess neural correlates of BOP effects on subjective feeling of others'

739  pain. Specifically, we conducted a representational similarity analysis (RSA) (Nili et al.,

740  2014) of brain activity using a dissimilarity matrix (DM) constructed from scores of

741  pain intensity in different conditions. The RSA sought to find patterns of brain activities

742  in the empathy neural network which can predict the pattern of subjective feeling of

743  others' pain that varied due to BOP. To do this, we first conducted ANOVAs of the

744  mean rating scores and found a significant Identity (patient vs. actor/actress) ×

745  Expression (pain vs. neutral) interaction on the rating scores of pain intensity ($F(1,30) =$

746  $5.370$, $p = 0.027$, $\eta_p^2 = 0.152$, 90% CI = (0.029, 0.391)) but not on the rating scores of

747  unpleasantness ($F(1,30) = 3.945$, $p = 0.056$, $\eta_p^2 = 0.116$, 90% CI = (0, 0.296), see

748  Supplementary file 14 for statistical details). Simple effect analyses showed

749 significantly larger scores of pain intensity for pain expressions of patients (vs.

750 actors/actresses) (F(1,30) = 9.823, p = 0.004, $\eta_p^2$ = 0.247, 90% CI = (0.053, 0.427)),

751 whereas scores of pain intensity did not differ significantly between neutral faces with

752 patient and actor/actress identifies (F(1,30) = 2.829, p = 0.103, $\eta_p^2$ = 0.086, 90% CI = (0,

753 0.260)). The results suggested a clear boundary between subjective feelings of pain

754 intensity in different conditions. Thus we constructed a $4 \times 4$ DM for each participant

755 with each cell in the DM representing the mean difference in rating scores of pain

756 intensity between each pair of conditions, as illustrated in Fig. 6c.

757     Next, we conducted a searchlight RSA to identify brain regions in which the

758 pairwise similarity of neural responses in the 4 conditions (2 Expressions $\times$ 2 Identities)

759 corresponded to the behavioral DM in each participant (see Methods for details). We

760 first conducted a whole-brain searchlight RSA for each participant. The searchlight

761 results of all participants were then subject to a second group-level analysis to examine

762 the voxels in the empathy network, defined based on the results of the whole-brain

763 contrast of painful vs. non-painful stimuli applied to targets, that passed a threshold of

764 voxel level p < 0.05, FWE corrected. The results revealed significant activations in the

765 left anterior insula (MNI peak coordinates x/y/z = -39/20/8) and inferior parietal cortex

766 (-60/-19/29), and the right anterior insula/frontal cortex (36/23/11), superior temporal

767 gyrus (54/-37/11), inferior post-central gyrus (63/-40/26), and superior parietal cortex

768 (39/-49/50) (Fig. 6d).

769    Finally, we estimated BOP effects on neural responses in a vicarious pain signature

770    (VPS) map that was identified to be sensitive to perceived painful stimulations applied

771    to others but not to self-experienced pain (Krishnan et al., 2016). We calculated the VPS

772    pattern responses to video clips showing patient or actor/actress faces that received

773    painful (needle penetration) or non-painful (cotton swab touch) stimulation using both

774    the body-specific VPS map in response to perceived noxious stimulation of body limbs

775    (Krishnan et al., 2016) and the general VPS in response to both perceived noxious

776    stimulation of body limbs and painful facial expressions (Zhou et al., 2020). We tested

777    the hypothesis of decreased VPS responses to actors/actresses' compared to patients'

778    pain (i.e., lack of BOP reduces empathic brain activities) by conducting t-tests of BOLD

779    signals in VPS maps. The results showed that activities in the VPS pattern were

780    significantly decreased in response to video clips showing actors/actresses' compared to

781    patients' pain (Fig. 6e and 6f, body-specific VPS: mean $\pm$ s.d. = 41.487 $\pm$ 28.794 vs.

782    46.548 $\pm$ 32.051, $t(30) = -2.059$, $p_{(one-tailed)} = 0.024$, $BF_{+0} = 2.361$; general VPS: mean $\pm$

783    s.d. = 1.188 $\pm$ 6.058 vs. 2.462 $\pm$ 6.997, $t(30) = -2.447$, $p_{(one-tailed)} = 0.010$, $BF_{+0} = 4.820$).

784    These results provide further evidence for decreased empathic brain activities due to

785    lack of BOP for actors/actresses' pain in the empathic neural network.

786    **Discussion**

787    We conducted six experiments to investigate psychological and neural mechanisms

788    underlying BOP impacts on empathy and altruistic behavior in humans. We

789    manipulated individuals' BOP by randomly assigning patient or actor/actress identities

790    to faces as there was a lack of BOP for actors/actresses' faces but not for patients' faces.

791    We also estimated individuals' intrinsic BOP by asking the participants to estimate

792    effectiveness of medical treatments of patients to trigger BOP as an effective medical

793    treatment reduces a patient's pain. We further measured brain activity using EEG and

794    fMRI to examine BOP effects on empathic neural responses with high temporal and

795    spatial resolutions, respectively. Our behavioral and neuroimaging findings showed

796    evidence for a functional role of BOP in modulations of the

797    perception-emotion-behavior reactivity by illustrating how BOP predicted and affected

798    self-reports of empathy, empathic brain activities, and monetary donations. Our findings

799    suggest that BOP may provide a cognitive basis for empathy and altruistic behavior in

800    humans.

801        Experiments 1 and 2 showed behavioral evidence that manipulated changes in BOP

802    caused subsequent variations of self-report of empathy and altruistic behavior along the

803    directions as predicted. Specifically, decreasing BOP concomitant with changes in face

804    identities (from patient to actor/actress) or changes in effective medical treatments

805    (from suffering due to a disease to recovery due to medical treatment) significantly

806    reduced self-report of both cognitive (perceived intensity of others' pain) and affective

807    (own unpleasantness induced by perceived pain in others) components of empathy.

808    Decreasing BOP also inhibited following altruistic behavior that was quantified by the

809    amount of monetary donations to those who showed pain expressions. By contrast,

810    reassuring patient identities in Experiment 1 or by noting the failure of medical

811    treatment related to target faces in Experiment 2 increased subjective feelings of others'

812    pain and own unpleasantness and prompted more monetary donations to target faces.

813    The increased monetary donations might be due to that repeatedly confirming patient

814    identity or knowing the failure of medical treatment increased the belief of authenticity

815    of targets' pain and thus enhanced cognitive and affective components of empathy.

816    Alternatively, repeatedly confirming patient identity or knowing the failure of medical

817    treatment might activate other emotional responses to target faces such as pity or

818    helplessness, which might also influence altruistic decisions. The increased empathy

819    rating scores and monetary donations might also reflect a contrast effect due to rating

820    patient and actor/actress targets alternately. These possible accounts can be clarified in

821    future work by asking participants to report their emotions and performing rating tasks

822    on patient and actor/actress targets in separate blocks of trials. In consistent with the

823    effects of manipulated BOP on empathy and altruism across the participants, the results

824    of Experiment 2 showed that individuals' intrinsic BOP related to each target face

825    predicted their self-report of empathy and altruistic behavior across different target

826    faces. Moreover, decreased (or increased) intrinsic BOP also predicted changes in

827    empathy/altruistic behavior across different target faces. These converging behavioral

828    findings across different participants and across different target faces provide evidence

829    for causal relationships between BOP and empathy/altruism.

830        Our results showed that self-reports of others' pain intensity and own

831    unpleasantness elicited by perception of others' pain were able to positively predict

832    altruistic behavior across individuals. Previous research using questionnaire measures of

833    empathy ability found that empathy as a trait is positively correlated with the amount of

834    money shared with others in economic games (Edele et al., 2013; Li et al., 2019).

835    Together, these findings are consistent with the proposition that empathy, as either an

836    instant emotional response to others' suffering (e.g., estimated in our study) or a

837    personality trait (e.g., estimated in Edele et al. (2013) and Li et al. (2019)), plays a key

838    role in driving altruistic behavior (Batson, 1987; Batson et al., 2015; Eisenberg et al.,

839    2010; Hoffman, 2008; Penner et al., 2005). Our mediation analyses of the behavioral

840    data in both Experiments 1 and 2 further revealed that the effects of decreased BOP on

841    monetary donations were mediated by self-report of others' pain intensity. These results

842    further suggest empathy as an intermediate mechanism of the BOP effects on altruistic

843    behavior.

844        Our neuroimaging experiments went beyond subjective estimation of the

845    relationships between BOP and empathy/altruism by investigating neural mechanisms

846    underlying BOP effects on empathy for others' pain. It is necessary to conduct objective

847    estimation of empathy to examine BOP effects because self-report measures of empathy

848    can be influenced by social contexts and are unable to unravel brain mechanisms

849    underlying BOP effects on empathy (e.g., Sheng and Han, 2012). Our EEG results in

850    Experiments 3 and 5 repeatedly showed that neural responses to pain (vs. neutral)

851    expressions over the frontal regions within 200 ms after face onset (indexed by the P2

852    amplitude over the frontal/central electrodes) were significantly reduced to faces with

853 actor/actress identities compared to those with patient identities. The results in

854 Experiments 3 and 4 indicate that BOP concomitant with face identity (i.e., patients'

855 pain expressions manifest their actual painful emotional states whereas actors/actresses'

856 pain expressions do not) rather than face identity (e.g., Tiger or Lion team identities)

857 alone resulted in modulations of the P2 amplitudes to pain expressions in the direction

858 as expected. Numerous EEG studies have shown that the frontal P2 component

859 responds with enlarged amplitudes to various facial expressions such as fear, anger,

860 happy (Williams et al. 2006; Luo et al. 2010; Calvo et al. 2013) and pain (Sheng and

861 Han, 2012; Sheng et al., 2013; 2016) expressions compared to neutral faces. These

862 findings uncovered early affective processing by differentiating emotional and neutral

863 expressions. ERPs to others' pain within 200 ms post-stimulus occur regardless of task

864 demands and are associated with spontaneous empathy for pain (Fan and Han, 2008).

865 Our ERP results indicate that BOP may provide a cognitive basis for early spontaneous

866 neural responses to others' suffering reflected in pain expressions. Moreover, the results

867 in Experiment 5 showed that the early spontaneous empathic neural responses in the P2

868 time window mediated the BOP effect on self-report of others' pain intensity, which

869 further mediated the relationship between the P2 empathic responses and the amount of

870 monetary donations. These results highlight both early spontaneous neural responses to

871 others' pain and subjective feelings of others' pain as intermediate mechanisms by

872 which BOP influences altruistic behavior.

873    To identify neural architectures underlying BOP effects on empathy, we recorded

874    BOLD responses, using fMRI, to perceived painful and non-painful stimuli applied to

875    individuals with patient or actor/actress identities in Experiment 6. We showed that the

876    contrast of perceived painful (vs. non-painful) stimulations activated the sensory (i.e.,

877    post-central gyrus), affective (i.e., insula), and cognitive (i.e., mPFC) nodes of the

878    empathy network, similar to the findings of previous studies (Singer et al., 2004;

879    Jackson et al., 2005; Saarela et al., 2007; Shamay-Tsoory et al., 2009; Han et al., 2009;

880    Fan et al., 2011; Lamm et al., 2011; Zhou and Han, 2021; Luo et al., 2014). Viewing

881    non-painful stimulations applied to neutral faces with patient versus actor/actress

882    identities revealed increased activity in the mPFC and bilateral TPJ, suggesting possible

883    neural representation of facial identities in the brain regions. Most importantly, the

884    results of searchlight RSA that was sensitive to both stimuli and subjective feelings

885    evoked by the stimuli revealed significant variations of activities in the insula,

886    post-central gyrus, and lateral frontal cortex in correspondence with the patterns of

887    self-reports of empathy for patients and actors/actresses' pain. In other words, the

888    patterns of the activities in the insula, post-central gyrus, and lateral frontal cortex were

889    able to predict distinct subjective feelings of patients' and actors/actresses' pain.

890    Moreover, the results of our VPS analyses showed consistent evidence for decreased

891    neural activities in the empathy-related neural network due to lack of BOP. These fMRI

892    results together suggest that activities in the brain regions supporting affective sharing

893    (e.g., insula, Shamay-Tsoory et al., 2009; Fan et al., 2011; Lamm et al., 2019), empathic

894 sensorimotor resonance (e.g., post-central gyrus, Avenanti et al., 2005; Zhou and Han,

895 2021), and emotion regulation (e.g., lateral frontal cortex, Ochsner and Gross, 2005;

896 Etkin et al., 2015) may provide intermediate mechanisms underlying variations of

897 subjective feelings of others' pain intensity due to lack of BOP.

898 Numerous studies have shown evidence for modulations of empathy by social

899 contexts. Contextual variables that influence perception of others' pain and empathy

900 include empathy targets' posture (Martel et al., 2008), identifiable pain pathology

901 (Twigg and Byrne, 2015), moral valence (Cui et al., 2016; Nicolardi et al., 2020), etc.

902 Empathizers' prior exposure to pain (Prkachin and Rocha, 2010), socioeconomic status

903 (Varnum et al., 2015), and cultural experiences (Wang et al., 2015; Hampton and

904 Varnum, 2018) also influence empathy and its underlying brain activities. Perceived

905 information about social relationships between observers and empathy targets also

906 modulates empathic neural responses such that, relative to viewing own-race or

907 own-team individuals' pain, viewing other-race or opponent-team individuals' pain

908 decreased empathic neural responses in the affective (e.g., ACC, AI), cognitive (e.g.,

909 mPFC, TPJ), and sensorimotor (e.g., motor cortex) nodes of the empathy network (Xu

910 et al., 2009; Avenanti et al., 2010; Hein et al., 2010; Mathur et al., 2010; Sheng and Han,

911 2012; Sheng et al., 2014; 2016; Han, 2018; Zhou and Han, 2021). The perceived

912 intergroup (racial) relationships between empathizers and empathy targets also

913 influenced altruistic behavior such as medical treatment (Drwecki et al., 2011). These

914 findings uncovered how social information perceived from stimuli and social experience

915    modulate empathic neural responses to others' suffering and subsequent social behavior.

916    The results of our current work complemented the findings of previous studies by

917    uncovering how beliefs, as preexisting internal mental representations of something that

918    is not immediately present to the scenes (Fuentes, 2019), also modulate people's

919    empathy and following altruistic behavior. Specifically, in the current study,

920    participants' beliefs (i.e., pain expressions of patients manifest their actual feelings

921    whereas pain expressions performed by actors/actresses do not) weakened the

922    participants' empathy for others' pain and reduced their monetary donations to those

923    who appeared suffering. BOP effects on empathy and altruistic behavior can be

924    understood as modulations of empathy by preexisting internal information (e.g., beliefs)

925    whereas previous findings revealed modulations of empathy by instantly perceived

926    social information in a specific social context. These findings together help to construct

927    neurocognitive models of empathy that take into consideration of both perceived social

928    information and preexisting internal information and their interactions that lead to

929    modulations of empathy and altruistic behavior during real-life social interactions.

930        It should be noted that our experimental manipulations changed the participants'

931    mind about the models' identities (e.g., patient vs. actor/actress) rather than explicitly

932    asking them to alter their BOP. BOP altered implicitly with target persons' identities due

933    to observers' knowledge about individuals with different identities (e.g., painful stimuli

934    applied to actors/actresses do not really hurt them and they show facial expressions to

935    pretend a specific emotional state). Therefore, the BOP effects on empathy and altruistic

936    behavior identified in our study might take place implicitly. This is different from the

937    placebo effects on first-hand pain experiences that are produced by explicitly perceived

938    verbal, conditioned, and observational cues that induce expectations of effective

939    analgesic treatments (Meissner et al., 2011). Similar explicit manipulations of making

940    individuals believe receiving oxytocin also promotes social trust and preference for

941    close social distances (Yan et al., 2018). Moreover, the placebo treatment relative to a

942    control condition significantly attenuated activations in the ACC, AI, and subcortical

943    structures (e.g., the thalamus) in response to painful electric shocks but increased the

944    prefrontal activity during anticipation of painful stimulations possibly to inhibit activity

945    in pain processing regions (Wager et al., 2004; 2015). The brain regions in which

946    empathic neural responses altered due to BOP (e.g., the lateral frontal cortex) as

947    unraveled in the current study do not overlap with those in which activities are

948    modulated by placebo analgesia (Atlas and Wager, 2014). These results suggest there

949    may be distinct neural underpinnings of BOP effects on empathic brain activity and

950    placebo effects on brain responses to first-hand pain experiences.

951        Do beliefs also provide a cognitive basis for the widely documented ingroup bias in

952    empathy for pain? Previous studies suggest that multiple neurocognitive mechanisms

953    are involved in ingroup bias in empathy for pain such as lack of attention (Sheng and

954    Han, 2012) and early group-based categorization of outgroup faces (Zhou et al., 2020,

955    see Han, 2018 for review). There has been behavioral evidence that white individuals

956    who more strongly endorsed false beliefs about biological differences between blacks

957   and whites (e.g., "black people's skin is thicker than white people's skin") reported

958   lower pain ratings for a black (vs. white) target and suggested less accurate treatment

959   recommendations (Hoffman et al., 2016). These behavioral findings suggest that other

960   beliefs may also provide a basis for modulations of empathy for others' pain and

961   relevant altruistic behavior. The underlying brain mechanisms, however, remain

962   unknown. The paradigms developed in the current study may be considered in future

963   research to examine neural underpinnings of the effects of false beliefs on empathy for

964   pain.

965         Another question arising from the findings of the current study is whether the belief

966   effect is specific to neural underpinnings of empathy for pain or is also evident for

967   neural responses to other facial expressions. To address this issue, we conducted an

968   additional EEG experiment in which we tested (1) whether beliefs of authenticity of

969   others' happiness influence brain responses to perceived happy expressions, and (2)

970   whether lack of beliefs of others' happiness also modulate neural responses to happy

971   expressions in the P2 time window, similar to the BOP effect on ERPs to pain

972   expressions (see Appendix 1 for methods). Similar to the paradigm used in Experiment

973   3, participants in the additional experiment had to first remember face identities

974   (awardees or actors/actresses). Thereafter these faces with happy or neutral faces were

975   presented with contextual information that the awardees showed happy expressions

976   when receiving awards whereas actors/actresses imitated others' happy expressions. The

977   participants also performed identity judgments on the faces while EEG was recorded.

978    Behavioral results in this experiment showed that participants reported less feelings of

979    actors' happiness compared to awardees' happiness. ERP results in this experiment

980    showed that lack of beliefs of authenticity of others' happiness (e.g., actors simulating

981    others' happy expressions vs. awardees smiling when receiving awards) reduced the

982    amplitudes of a long-latency positive component (i.e., P570) over the frontal region in

983    response to happy expressions. However, the face identities did not affect the P2

984    amplitudes in response to happy (vs. neutral) expressions (see Appendix 1 for statistical

985    details). These findings suggest that belief effects are evident for subjective feelings and

986    brain activities in response to happy expressions. However, beliefs of others' pain or

987    happiness affect neural responses to facial expressions in different time windows after

988    face onset. Future research should examine neural mechanisms underlying belief effects

989    on neural responses to other emotions to deep our understanding of general belief

990    effects on neural processes of others' emotional states.

991        Our behavioral and neuroimaging findings have implications for how we

992    understand the general functional role of beliefs in social cognition and interaction.

993    Empathy is supposed to originate from an evolved adaptation to quickly and

994    automatically respond to others' emotional states during parental care that is necessary

995    for offspring survival in humans and other species (De Waal, 2008; Decety, 2011). In

996    most cases of interactions among family members (i.e., between parents and offspring

997    or between siblings) perceived cues signaling pain in a person manifest his/her actual

998    emotional states that urge help from other family members. Such life experiences may

999    set up a default belief that perceived painful stimulation to others and their facial

1000   expressions reflect individuals' actual emotional states. This default belief provides a

1001   fundamental cognitive basis of reflexive and automatic empathy and empathic brain

1002   activity that further generates autonomic and somatic responses, as suggested by the

1003   perception-action model of empathy (Preston and de Waal, 2002). Nevertheless, when

1004   social interactions expand beyond family members to non-kin members and even

1005   strangers, perceived pain expressions or painful stimuli applied to others may not

1006   always manifest others' actual emotional states because perceived painful cues may be

1007   fake in some cases. BOP in such situations may function as cognitive gate-control to

1008   modulate neural responses to perceived pain in others. This is necessary for monitoring

1009   social interactions to determine whether to help or to coordinate with those who appear

1010   suffering. Our findings illustrate how the perception-emotion-behavior reactivity occurs

1011   under the cognitive constraint of BOP to keep empathy and altruistic decision/behavior

1012   for the right target who is really in need of help. In this sense, BOP also provides an

1013   important cognitive basis for survival and social adaption during social interactions.

1014       Some limitations of the current work create future research opportunities. For

1015   example, a recent approach to hierarchical Bayesian models of cognition assumes that

1016   the brain represents information probabilistically and people represent a state or feature

1017   of the world not using a single computed value but a conditional probability density

1018   function (Knill and Pouget, 2004; Friston, 2005; Clark, 2013; Tappin and Gadsby,

1019   2019). Our manipulations of BOP, however, had only two conditions (patient vs.

1020    actor/actress) and thus lack a model of effects of probability-based belief-updating on

1021    empathy and relevant altruistic behavior. Future research should examine how empathy

1022    and relevant altruistic behavior vary as a function of the degree of BOP. Other

1023    interesting research questions arising from our work include how the brain represents

1024    BOP. It has been proposed that different types of beliefs (e.g., empirical beliefs,

1025    conceptual beliefs, relational beliefs) exist in human mind and may have distinct neural

1026    underpinnings (Harris et al., 2009; Seitz and Angel, 2020). To address neural

1027    representations of BOP will allow researchers to further explore and construct neural

1028    models of the interaction between beliefs and empathic brain activity in the key nodes

1029    of the empathy network. Another interesting issue related to our findings is individual

1030    differences in BOP and BOP effects on empathy and altruism. Since specific degrees of

1031    beliefs differ widely across individuals (Ais et al. 2016), it is crucial to examine what

1032    personality/psychopathic traits or biological factors make individuals hold strong or

1033    weak BOP and exhibit large or small BOP effects on empathy and altruistic behavior. It

1034    is also important to clarify what environmental factors modify individuals' default BOP

1035    and consequently change their motivations to help those who appear suffering. To

1036    clarify these issues will advance our understanding of individual and contextual factors

1037    that shape the functional role of BOP in modulations of empathy and altruistic behavior.

1038    Finally, a general issue arising from the current work is whether beliefs affect the

1039    processing of other emotions such as fear, sad, and happy, and, if yes, whether there are

1040    common underlying psychological and neural mechanisms.

**Conclusion**

Our behavioral and neuroimaging findings provide a new cognitive framework for understanding human empathy and altruism. Our findings indicate that lack of BOP or decreasing BOP weakened human empathy and altruistic behavior. Changing BOP affected both subjective feelings of others' emotional states and the underlying brain activity. BOP effects on altruistic behavior were mediated by two serial mediators, i.e., empathic neural responses and subjective feelings of others' pain. Our behavioral and brain imaging findings suggest that BOP provides a cognitive basis of the perception-emotion-behavior reactivity that underlies human altruism. The methods developed in our study open a new avenue for testing functional roles of beliefs as cognitive-gate control of other emotion processing and relevant social behavior.

## Methods

**Participants**

Sixty Chinese students were recruited in Experiment 1 as paid volunteers (29 males, mean age ± s.d. = 21.15 ± 2.31 years). The sample size was estimated using G*Power (Faul et al., 2007) with a middle effect size of 0.25. To test the difference in pain intensity rating scores or monetary donations between the $1^{st}$_ and $2^{nd}$_round tests, we conducted ANOVAs with Test Phase ($1^{st}$ vs. $2^{nd}$_round) and Identity Change (patient to actor/actress vs. patient to patient) as independent within-subjects variables. To detect a significant Test x Identity interaction requires a sample size of 36 with an error probability of 0.05 and a power of 0.95, given the correlation among repeated measures

1062    (0.5) and the nonsphericity correction (1). Sixty Chinese students were recruited in

1063    Experiment 2 as paid volunteers (30 males, $21.55 \pm 2.45$ years). Thirty Chinese students

1064    were recruited in Experiment 3 (all males, $22.23 \pm 2.51$ years) as paid volunteers. The

1065    sample size was determined based on our previous EEG research on empathy for pain

1066    using the same set of stimuli (Sheng and Han, 2012). We recruited only male

1067    participants to exclude potential effects of gender difference in empathic neural

1068    responses. Thirty-one Chinese students were recruited in Experiment 4 as paid

1069    volunteers. One participant was excluded from data analyses due to his lower response

1070    accuracy during EEG recording ($< 50\%$). This left 30 participants (all males, $20.70 \pm$

1071    1.97 years) for behavioral and EEG data analyses. Thirty Chinese students were

1072    recruited in Experiment 5 (all males, $20.60 \pm 1.75$ years). Thirty-two Chinese students

1073    were recruited in Experiment 6 as paid volunteers. One participant was excluded from

1074    data analyses due to excessive head movement during fMRI scanning. There were 31

1075    participants left (all males, $22.23 \pm 2.59$ years) for behavioral and fMRI data analyses.

1076    The sample size in Experiment 6 was determined based on our previous fMRI research

1077    using similar stimuli (Luo et al., 2014). All participants had normal or

1078    corrected-to-normal vision and reported no history of neurological or psychiatric

1079    diagnoses. This study was approved by the local Research Ethics Committee of the

1080    School of Psychological and Cognitive Sciences, Peking University. All participants

1081    provided written informed consent after the experimental procedure had been fully

1082 explained. Participants were reminded of their right to withdraw at any time during the

1083 study.

1084 **Experiment 1: Lack of BOP reduces subjective estimation of empathy and**

1085 **altruistic behavior**

1086 **Stimuli and procedure**

1087 The stimuli were adopted from our previous work (Sheng and Han, 2012), which

1088 consisted of photos of 16 Chinese models (half males) with each model contributing one

1089 photo with pain expression and one with neutral expression.

1090 After reporting demographic information, the participants were informed that they

1091 would be paid with ¥10 as a basic payment for their participation. They would be able

1092 to obtain an extra bonus payment as much as ¥2 depending on their decisions in the

1093 following procedure. In the 1st_round test the participants were informed that they

1094 would be shown photos with pain expressions taken from patients who suffered from a

1095 serious disease. After the presentation of each photo the participants were asked to

1096 evaluate intensity of each patient's pain based on his/her expression by rating on a

1097 Likert-type scale ("How painful do you think this person is feeling?", 0 = not painful at

1098 all; 10 = extremely painful). This rating task was adopted from previous research (Bieri

1099 et al., 1990; Jackson et al., 2005; Lamm et al., 2007; Fan and Han, 2008; Sheng and

1100 Han, 2012) to assess the participants' understanding of others' pain feeling — a key

1101 component of empathy. The instructions of the rating tasks focused on emotional states

1102 of faces and had nothing to do with face identities (i.e., patients or actors/actresses).

1103 Therefore, BOP effects on empathy, if observed, occurred implicitly and automatically.

1104 Immediately after the pain intensity rating, the participants were asked to decide how

1105 much from the extra bonus payment they would like to donate to the patient (0 to 10

1106 points, 1 point = ¥0.2). The participants were informed that the amount of one of their

1107 donation decisions would be selected randomly and endowed to a charity organization

1108 to help those who suffered from the same disease.

1109 After the $1^{st}$_round test the participants were asked to perform a short (5 mins)

1110 calculation task (10 arithmetic calculations, e.g. 25-3×7=?) to clean their memory of the

1111 $1^{st}$_round ratings. Thereafter, the participants were told that the photos were actually

1112 taken from 8 patients and 8 actors/actresses and this experiment actually tested their

1113 ability of recognizing social identities by examination of facial expressions. Faces

1114 assigned with patient or actor/actress identities were counterbalanced across the

1115 participants. The participants were then asked to conduct the $2^{nd}$_round test in which

1116 each photo was presented again with a word below to indicate patient or actor/actress

1117 identity of the face in the photo. The participants had to report again pain intensity of

1118 each face and how much they would like to donate to the person shown in the photo.

1119 The participants were informed that an amount of money would be finally selected

1120 randomly from their $2^{nd}$_round decisions and donated to one of the patients through the

1121 same charity organization. After the experiments had been finished, the total amount of

1122 the participants' donations were subject to a charity organization.

1123    We conducted ANOVAs of rating scores of pain intensity and amounts of monetary

1124    donations with Test Phase ($1^{st}$ vs. $2^{nd}$_round) × Identity Change (patient to actor/actress

1125    vs. patient to patient) as independent within-subjects variables to assess whether and

1126    how beliefs of others' pain (BOP) influenced empathy and altruistic behavior toward

1127    those who suffered. Finally, the participants completed two questionnaires to estimate

1128    individual differences in trait empathy (Davis, 1983) and interpersonal trust (Wright and

1129    Tedeschi, 1975). We analyzed the relationship between our empathy/altruistic measures

1130    and individuals' trait empathy/interpersonal trust but failed to find significant results

1131    and thus were not reported in the main text.

1132    **Mediation analysis**

1133    We performed mediation analyses to examine whether pain intensity mediates the

1134    pathway from BOP to monetary donation. To do this, we first dummy coded

1135    patient-identity change (i.e., 0 (patient identity in the $1^{st}$_round test) and 1 (actor/actress

1136    in the $2^{nd}$_round test) or patient-identity repetition (i.e., as 0 (patient identity in the

1137    $1^{st}$_round test) and 1 (patient identity in the $2^{nd}$_round test). Then, we estimated four

1138    regression models: 1) whether the independent variable (BOP) significantly accounts

1139    for the dependent variable (monetary donation) when not considering the mediator (e.g.,

1140    Path c); 2) whether the independent variable (BOP) significantly accounts for the

1141    variance of the presumed mediator (pain intensity) (e.g., Path a); 3) whether the

1142    presumed mediator (pain intensity) significantly accounts for the variance of the

1143    dependent variable (monetary donation) when controlling the independent variable

1144 (BOP) (e.g., Path b); 4) whether the independent variable (BOP) significantly accounts

1145 for the variance of the dependent variable (monetary donation) when controlling the

1146 presumed mediator (pain intensity) (e.g., Path c'). To establish the mediation, the path c

1147 is not required to be significant. The only requirement is that the indirect effect a×b is

1148 significant. Given a significant indirect effect, if Path c is not significant, the mediation

1149 is classified as indirect-only mediation which is the strongest full mediation (Kenny et

1150 al., 1998; Zhao et al.., 2019). A bootstrapping method was used to estimate the

1151 mediation effect. Bootstrapping is a nonparametric approach to estimate effect-sizes and

1152 hypotheses of various analyses, including mediation (Shrout and Bolger,2002;

1153 Mackinnon et al., 2004). Rather than imposing questionable distributional assumptions,

1154 a bootstrapping analysis generates an empirical approximation of the sampling

1155 distribution of a statistic by repeated random resampling from the available data, which

1156 is then used to calculate p-values and construct confidence intervals. 5,000 resamples

1157 were taken for our analyses. Moreover, this procedure supplies superior confidence

1158 intervals (CIs) that are bias-corrected and accelerated (Preacher et al., 2007; Preacher

1159 and Hayes, 2008a, 2008b). The analyses were performed using Hayes's PROCESS

1160 macro (Model 4, Hayes, 2017).

1161 **Statistical comparison**

1162     Behavioral data were assumed to have a normal distribution but this was not

1163 formally tested. 95% Confidence intervals (95% CIs) were reported for t-tests of the

1164 mean difference between two conditions and for correlation analyses of correlation

1165 coefficients. 90% CIs were reported for effect sizes ($\eta_p^2$) of ANOVA analyses.

1166 According to Steiger (2004), the general rule of thumb to use CIs to test a statistical

1167 hypothesis (H0) is to use a $100\times(1-\alpha)\%$ / $100\times(1-2\alpha)\%$ CI when testing a two-sided /

1168 one-sided hypothesis at alpha level. We thus reported 90% CIs of $\eta^2$ in ANOVAs

1169 because $\eta^2$ is always positive.

1170 **Experiment 2: Intrinsic BOP predicts subjective estimation of empathy and**

1171 **altruistic behavior**

1172     The face stimuli and the procedure were the same as those in Experiment 1 except

1173 the following. The participants were informed that they were to be shown photos with

1174 pain expressions taken from patients who had suffered from a serious disease and

1175 received medical treatment. After the presentation of each photo the participants were

1176 asked to estimate how effective the medical treatment was for each patient by rating on

1177 a Likert-type scale (0 = no effective or 0% effective, 100 = fully effective or 100%

1178 effective). Besides rating pain intensity of each face in the 1$^{st}$_round test, the

1179 participants were asked to report how unpleasant they were feeling when they viewed

1180 the photo (i.e., own unpleasantness) by rating on a Likert-type scale ("How unpleasant

1181 do you feel when viewing this person?" 0 = not unpleasant at all, 10 = extremely

1182 unpleasant). The unpleasantness rating was performed to evaluate emotional sharing of

1183 others' pain — another key component of empathy (Jackson et al., 2005; Fan and Han,

1184 2008; Sheng and Han, 2012). The order of the two empathy rating tasks was

1185 counterbalanced across the participants. Immediately after the empathy rating tasks, the

1186      participants were asked to decide how much from the extra bonus payment they would

1187      like to donate to the patient (0 to 10 points, 1 point = ¥0.2).

1188      In the $2^{nd}$_round test the participants were told that the medical treatment was

1189      actually effective for only half of the patients. Each photo was then presented again with

1190      information that the medical treatment applied to the patient was 100% effective or 0%

1191      effective. Thereafter, the participants were asked to perform the rating tasks and

1192      monetary donations as those in the $1^{st}$_round test. The participants were told that an

1193      amount of money would be finally selected from their $2^{nd}$_round decisions and donated

1194      to one of the patients.

1195      **Mediation analysis**

1196      This was the same as that in Experiment 1 except that we tested whether changes of

1197      pain intensity mediate the pathway from decreased BOP or enhanced BOP to changes of

1198      monetary donation. To do this, we first calculated belief update (decreased BOP:

1199      100%-effect minus the participants' initial estimation; enhanced BOP: the participants'

1200      initial estimation minus 0%-effect). Then, we estimated four regression models: 1)

1201      whether the independent variable (BOP) significantly accounts for the dependent

1202      variable (changes of monetary donation) when not considering the mediator (e.g., Path

1203      c); 2) whether the independent variable (BOP) significantly accounts for the variance of

1204      the presumed mediator (changes of pain intensity) (e.g., Path a); 3) whether the

1205      presumed mediator (changes of pain intensity) significantly accounts for the variance of

1206      the dependent variable (changes of monetary donation) when controlling the

1207 independent variable (BOP) (e.g., Path b); 4) whether the independent variable (BOP)

1208 significantly accounts for the variance of the dependent variable (changes of monetary

1209 donation) when controlling the presumed mediator (changes of pain intensity) (e.g.,

1210 Path c').

1211 **Experiment 3: Lack of BOP decreased empathic brain activity**

1212 **Stimuli and procedure**

1213     Face stimuli were adopted from our previous work (Sheng and Han, 2012) and used

1214 in Experiments 3, 4 and 5 in this study. The stimuli consisted of 32 faces of 16 Chinese

1215 models (half males) with each model contributed one photo with pain expression and

1216 one with neutral expression. During behavioral tests or EEG recording, each photo was

1217 presented in the center of a gray background on a 21-inch color monitor, subtending a

1218 visual angle of $3.8° \times 4.7°$ (width $\times$ height: $7.94 \times 9.92$ cm) at a viewing distance of 60

1219 cm.

1220     Before EEG recording the participants were asked to perform an identity memory

1221 task in which faces with neutral expressions were presented. Eight faces were marked as

1222 patients and 8 faces as actors/actresses. After viewing photos with marked identity for

1223 15 minutes, the participants performed a discrimination task in which each neutral face

1224 was displayed for 200 ms and the participants had to press the left or right button using

1225 the left or right index finger to indicate identity of each face (i.e., patient or actor/actress)

1226 within two seconds. After their response accuracies reached 100%, the participants were

1227 moved into an acoustically- and electrically-shielded booth for EEG recording.

1228      During EEG recording each trial consisted of a painful or neutral face with a

1229    duration of 200 ms, which was followed by a fixation cross with a duration varying

1230    randomly between 800 and 1400 ms. There were 8 blocks of 64 trials (each of the 32

1231    photographs was presented twice in a random order in each block). The participants

1232    were asked to press the left or right button using the left or right index finger to indicate

1233    the identity of the face (i.e., patient or actor/actress) as fast and accurately as possible.

1234    The relation between responding hand and face identity was counterbalanced across

1235    different blocks of trials.

1236      After EEG recording, the participants were presented with each face again with a

1237    neutral or pain expression and asked to rate how painful the person is feeling (i.e., pain

1238    intensity) by rating on a Likert-type scale (1 = not painful at all; 7 = extremely painful).

1239    To estimate the participants' BOP, they were also asked to answer the question of "To

1240    what extent do you believe the identity of this model (either patient or actor/actress)?"

1241    on a 15-point Likert-type scale (-7 = extremely believed as an actor/actress, 0 = not sure,

1242    7 = extremely believed as a patient).

1243    **EEG data acquisition and analysis**

1244      A NeuroScan system (CURRY 7, Compumedics Neuroscan) was used for EEG

1245    recording and analysis. EEG was continuously recorded from 32 scalp electrodes and

1246    was re-referenced to the average of the left and right mastoid electrodes offline.

1247    Impedances of individual electrodes were kept below 5 kΩ. Eye blinks and vertical eye

1248    movements were monitored using electrodes located above and below the left eye. The

1249 horizontal electro-oculogram was recorded from electrodes placed 1.5-cm lateral to the

1250 left and right external canthi. The EEG signal was digitized at a sampling rate of 1,000

1251 Hz and subjected to an online band-pass filter of 0.01–400 Hz. EEG data were filtered

1252 with a low-pass filter at 30 Hz offline. Artefacts related to eye movement or eye blinks

1253 were removed using the covariance analysis tool implemented in CURRY 7 (Semlitsch

1254 et al., 1986). Only trials with correct responses to face identity were included for data

1255 analyses (see Supplementary file 15 for the numbers of trials included for data analyses

1256 in Experiments 3-5). ERPs in each condition were averaged separately offline with an

1257 epoch beginning 200 ms before stimulus onset and continuing for 1200 ms. Trials The

1258 baseline for all ERP measurements was the mean voltage of a 200-ms prestimulus

1259 interval and the latency was measured relative to the stimulus onset.

1260     Face stimuli in the identity judgment task elicited an early negative activity at

1261 95-115 ms (N1) and a positive activity at 175-195 ms (P2), followed by a positive

1262 activity at 280-340 ms (P310) and a long-latency positivity at 500–700 ms (P570) over

1263 the frontal area. The mean ERP amplitudes were subject to ANOVAs with Identity

1264 (patient vs. actor/actress) and Expression (pain vs. neutral) as within-subject variables.

1265 To avoid potential significant but bogus effects on ERP amplitudes due to multiple

1266 comparisons (Luck and Gaspelin, 2017), the mean values of the amplitudes of the N1,

1267 P2, P310, and P570 components were calculated at frontocentral electrodes (i.e., F3, Fz,

1268 F4, FC3, FCz and FC4).

1269    To further assess the null hypothesis regarding the difference in the P2 amplitude in

1270    response to pain and neutral expressions of actors/actress' faces, we conducted Bayes

1271    factor analyses for repeated-measures ANOVA and paired t-tests. We calculated the

1272    Bayes factor in the program R v.3.5.1 (www.r-project.org) using the function anovaBF

1273    and ttestBF from the package BayesFactor (Morey and Rouder, 2015). We conducted

1274    Bayes factor analyses based on the default priors for ANOVA and paired t-test design

1275    (scale r on an effect size of 0.707). A Bayes factor indicates how much more likely each

1276    alternative model is supported compared with the null hypothesis.

1277    **Experiment 4: BOP is necessary for modulations of empathic brain activity**

1278    **Stimuli and procedure**

1279    These were the same as those in Experiment 3 except the following. Before EEG

1280    recording, the participants were informed that all the 16 faces were patients and they

1281    were from two baseball teams (half from Tiger team and half from Lion team). After the

1282    identity memory task, they performed identity judgments on faces with neutral or pain

1283    expressions by pressing one of two buttons while EEG was recorded.

1284    **EEG data acquisition and analysis**

1285    These were the same as those in Experiment 3.

1286    **Experiment 5: Empathic brain activity mediates relationships between BOP and**

1287    **empathy/altruistic behavior**

1288    **Stimuli and procedure**

1289    The stimuli and behavioral tests were the same as those in Experiment 1 to assess

1290    BOP effects on self-report of perceived pain intensity and altruistic decisions.

1291    Thereafter, the participants went through the EEG session that was the same as that in

1292    Experiment 3 to examine BOP effects on empathic brain activity. These designs

1293    allowed us to test whether BOP induced changes of empathic brain activity plays a

1294    mediator role in the pathway from belief changes to altered subjective feelings of others'

1295    pain and altruistic decisions.

1296    **Behavioral and EEG data recording and analyses**

1297    These were the same as those in Experiments 1 and 3.

1298    **Multiple mediation model analysis**

1299    We constructed a serial mediation model to test the hypothesis that BOP (dummy

1300    coded as 0 for patients and 1 for actors/actresses) effect on monetary donations was

1301    sequentially mediated by two chain mediators, i.e., empathic neural responses and

1302    subjective feelings of others' pain. This model includes three indirect paths: (1) indirect

1303    effect of BOP on monetary donation via empathic neural responses (i.e. P2 amplitude);

1304    (2) indirect effect of BOP on monetary donation via subjective feelings of others' pain

1305    (pain intensity); (3) indirect effect of BOP on monetary donation via P2 amplitude $\times$

1306    pain intensity. To do this, we estimated seven regression models: 1) whether the

1307    independent variable (BOP) significantly accounts for the dependent variable (monetary

1308    donation) when not considering the mediator (e.g., Path c); 2) whether the independent

1309    variable (BOP) significantly accounts for the variance of the presumed mediator (P2

1310    amplitude) (e.g., Path $a_1$); 3) whether the independent variable (BOP) significantly

1311    accounts for the variance of the presumed mediator (pain intensity) (e.g., Path $a_2$); 4)

1312    whether the first independent mediator (P2 amplitude) significantly accounts for the

1313    variance of the second mediator (pain intensity) (e.g., Path $d_{21}$); 5) whether the

1314    presumed mediator (P2 amplitude) significantly accounts for the variance of the

1315    dependent variable (monetary donation) when controlling the independent variable

1316    (BOP) (e.g., Path $b_1$); 6) whether the presumed mediator (pain intensity) significantly

1317    accounts for the variance of the dependent variable (monetary donation) when

1318    controlling the independent variable (BOP) (e.g., Path $b_2$); 7) whether the independent

1319    variable (BOP) significantly accounts for the variance of the dependent variable

1320    (monetary donation) when controlling the presumed the two mediators (e.g., Path c').

1321    To test the significance of the three paths, we separately conducted to examine the

1322    significance of indirect effect ($a_1 \times b_1$) of BOP on monetary donation via the P2

1323    amplitude; indirect effect ($a_2 \times b_2$) of BOP on monetary donation via pain intensity;

1324    indirect effect ($a_1 \times d_{21} \times b_2$) of BOP on monetary donation via P2 amplitude $\times$ pain

1325    intensity. Similarly, the bootstrapping method was used to estimate the effect-size and

1326    test the hypothesis.

1327    **Implicit association test**

1328        To assure our experimental manipulation of patient and actor/actress identities, after

1329    the EEG recording, participants were asked to complete a modified implicit association

1330    test (IAT, Greenwald et al., 1998). The participants were asked to respond to faces with

1331  patient identifies and patient related words (e.g. ache, weak) with one key and to faces

1332  with actor/actress identities and actor/actress related words (e.g. imitation) with another

1333  key in two blocks of trials (60 trials in each block). They were then asked to respond to

1334  faces with patient identities and actor/actress related words with one key and to faces

1335  with actor/actress identities and patient related words with another key in two additional

1336  blocks of trials. A D score was then calculated based on response times according to the

1337  established algorithm (Greenwald et al., 2003). A positive D score significantly larger

1338  than zero would suggest that patient faces were more strongly associated with patient

1339  (vs. actor/actress) relevant words whereas actor/actress faces were more strongly

1340  associated with actor/actress (vs. patient) relevant words.

1341  **Experiment 6: Neural structures underlying BOP effects on empathy**

1342  **Stimuli and procedure**

1343  We adopted 24 video clips from 6 models from our previous work (Luo et al., 2014)

1344  and recorded 56 video clips from 14 Chinese models (half males) in Experiment 6. Each

1345  model contributed four video clips, in which a face with pain expressions receiving

1346  painful stimulation (needle penetration) or with neutral expressions receiving

1347  non-painful stimulation (cotton swab touch) applied to the left or right cheeks. Each

1348  video subtended a visual angle of $21° \times 17°$ (width $\times$ height) at a viewing distance of 80

1349  cm during fMRI scanning.

1350  A photo of each model with a neutral expression was obtained from each video clip.

1351  These photos were then used in the identity memory task, which was the same as that in

1352 Experiment 3. After the identity memory task the participants underwent fMRI scanning.

1353 An event-related design was employed in 6 functional scans. Each scan consisted of 20

1354 video clips (half patients (5 pain and 5 neutral expressions) and half actors/actresses (5

1355 pain and 5 neutral expressions)) that were presented in a random order. Each video clip

1356 lasted for 3 s. There was a 9-s interstimulus interval between two successive video clips

1357 when the participants fixated at a central cross and had to judge the identity (patient or

1358 actor/actress) of each model in the video clip by pressing one of two buttons using the

1359 right index or middle finger. The relation between responding finger and face identity

1360 was counterbalanced across participants.

1361 After fMRI scanning, the participants were presented with each video clip again

1362 outside the scanner. They were asked to rate pain intensity of each model (1 = not

1363 painful at all; 7 = extremely painful) and own unpleasantness (1 = not unpleasant at all,

1364 7 = extremely unpleasant). Finally, we assessed the participants' beliefs of models'

1365 identities by asking them to answer the question of "To what extent do you believe the

1366 identity of this model (either patient or actor/actress)?" on a 15-point Likert-type scale

1367 (-7 = extremely believed to be an actor/actress, 0 = not sure, 7 = extremely believed to

1368 be a patient).

1369 **fMRI data acquisition and analysis**

1370 Imaging data were acquired using a 3.0 T Siemens scanner with a standard head

1371 coil. Head motion was controlled to the maximum extent by using foam padding.

1372 Functional images were acquired by using T2-weighted, gradient-echo, echo-planar

1373    imaging (EPI) sequences sensitive to Siemens scanner contrast ($64\times64\times32$ matrix with

1374    $3.75\times3.75\times5$ mm$^3$ spatial resolution, repetition time = 2000 ms, echo time = 30 ms, flip

1375    angle = 90°, field of view = $24\times24$ cm). Anatomical images were subsequently obtained

1376    using a standard 3D T1-weighted sequence ($256\times256\times144$ matrix with a spatial

1377    resolution of $1\times1\times1.33$ mm3, TR = 2530 ms, TE = 3.37 ms, inversion time (TI) = 1100

1378    ms, FA = 7°).

1379        Functional images were preprocessed using SPM12 software (the Wellcome Trust

1380    Centre for Neuroimaging, London, UK, http://www.fil.ion.ucl.ac.uk/spm). Functional

1381    scans were first corrected for within-scan acquisition time differences between slices

1382    and then realigned to the first volume to correct for inter-scan head motions. This

1383    realigning step provided a record of head motions within each fMRI run. Head

1384    movements were corrected within each run and six movement parameters (translation; x,

1385    y, z and rotation; pitch, roll, yaw) were extracted for further analysis in the statistical

1386    model. The functional images were resampled to $3 \times 3 \times 3$ mm$^3$ voxels, normalized to

1387    the MNI space using the parameters of anatomical normalization and then spatially

1388    smoothed using an isotropic of 8 mm full-width half-maximum (FWHM) Gaussian

1389    kernel.

1390        Whole-brain analyses was conducted to examine brain regions in which activities

1391    increased in response to pain versus neutral stimuli regardless of patient or actor/actress

1392    identities. This contrast pooled video clips of patient and actor/actress models together

1393    to focus on BOLD responses to painful versus neutral stimuli. The general linear model

1394 (GLM) had four regressors including patients receiving pain stimuli, patients receiving

1395 neutral stimuli, actors/actresses receiving pain stimuli, and actors/actresses receiving

1396 neutral stimuli. The GLM also included the realignment parameters to account for any

1397 residual movement-related effect. A box-car function was used to convolve with the

1398 canonical hemodynamic response in each condition. Random-effect analyses were

1399 conducted based on statistical parameter maps from each participant to allow population

1400 inference. The contrast values were compared using whole-brain paired t-tests to

1401 identify activations, which were defined using a threshold of voxel-level $p < 0.001$,

1402 uncorrected, cluster-level $p < 0.05$, FWE corrected. We also conducted a whole-brain

1403 analysis to calculate the contrast of patient versus actor/actress non-painful stimuli to

1404 test whether BOP may motivate inference of patients' mental states independently of

1405 any perceived painful cues.

1406 **Representational similarity analysis**

1407 We conducted a representational similarity analysis (RSA) of brain activity (Nili et

1408 al., 2014) to examine neural correlates to BOP effects on subjective feelings of others'

1409 pain. We constructed a $4 \times 4$ dissimilarity matrix (DM) for each participant with each

1410 cell in the DM represents the mean difference in rating scores of pain intensity between

1411 each pair of conditions. The DM was then used for a whole-brain searchlight RSA to

1412 identify brain regions in which the pairwise similarity of neural responses in the 4

1413 conditions (2 Expressions $\times$ 2 Identities) corresponded to the behavioral DM of

1414 condition dissimilarity in each participant. To do this, functional images were similarly

preprocessed using a GLM but were not smoothed and normalized. We then estimated a

GLM for each participant with Identity (patient vs. actor/actress) and Expression (pain

vs. neutral) as experimental regressors. The estimated beta images corresponding to

each condition were then averaged across runs at each voxel and were used as activity

patterns in the RSA toolbox (Nili et al., 2014). We compared the neural-pattern

similarity (i.e., the neural DM) with the behavioral DM in each voxel of the brain using

the searchlight procedure (Kriegeskorte et al., 2006). The neural DM was constructed

by 1 minus the correlation coefficient between the pattern vectors of each condition pair.

The Spearman rank correlations between the neural DM and behavioral DMs were

computed and assigned to the central voxel of the sphere. As such, the searchlight

procedure produced Spearman p values on each voxel for each participant, which were

then subject to Fisher's z transformation for statistical tests. The resulting z maps were

then normalized to standard space (resampled to 3 x 3 x 3 mm$^3$ voxels), smoothed

(FWHM= 8mm), and entered into a random effect analysis using one-sample t tests

against zero. The searchlight results of all participants were then subject to a second

group-level analysis to examine the voxels in the empathy network, defined based on

the results of the whole-brain contrast of painful versus non-painful stimuli applied to

targets, that passed a threshold of voxel level $p < 0.05$, FWE corrected.

**Neural signature analysis**

We conducted vicarious pain signature (VPS) analyses (Krishnan et al., 2016) to

further assess BOP effects on empathic brain activity. We first calculated contrast

1436    images in the condition of patient-pain (or actor/actress-pain) versus an implicit baseline

1437    (e,g., using a design matrix of [1, 0, 0, 0]) since the test-retest reliability was higher

1438    when examining brain activations to painful stimulation using an implicit baseline than

1439    using a control condition (Han et al., 2021). The VPS map, which was sensitive to

1440    perceived painful stimulations applied to others' body limbs but not to self-experienced

1441    pain (Krishnan et al., 2016), was then converted into the image space using the ImCalc

1442    function of SPM. Thereafter, the VPS map was dot-multiplied with the contrast of

1443    patient-pain versus baseline and the contrast of actor/actress-pain versus baseline,

1444    respectively. These yielded a scalar VPS response value in each condition. The VPS

1445    response values were then subject to a one-tailed t-test to test the hypothesis of

1446    decreased VPS responses related to actor/actress-pain relative to patient-pain. To further

1447    validate the results of VPS analyses, we conducted a similar analysis using the general

1448    vicarious pain signature, which was identified to respond to both perceived noxious

1449    stimulation of body limbs and painful facial expressions (Zhou et al., 2020).

1450    **Supplementary File legends**

1451    • **Source code 1.** Scripts for plotting Figure 1a, 1b, 2d, 2e, 2f, 3a, 4a, 5a, 5b.

1452    • **Source code 2.** Scripts for plotting Figure 3c, 4c, 5d.

1453    • **Source code 3.** Scripts for the whole-brain analysis in Figure 6a and 6b.

1454    • **Source code 4.** Scripts for plotting Figure 6c.

1455    • **Source code 5.** Scripts for plotting Figure 6d.

1456 • **Supplementary File 1.** Statistical results of the mediation analysis (pain intensity

1457 mediated the relationship between decreased BOP and monetary donations) in

1458 Experiment 1.

1459 • **Supplementary File 2.** Supplementary file 2. Statistical results of the mediation

1460 analysis (pain intensity mediated the relationship between enhanced BOP and monetary

1461 donations) in Experiment 1.

1462 • **Supplementary File 3.** Pain intensity, unpleasantness, and monetary donation (mean

1463 ± SD) in Experiment 2.

1464 • **Supplementary File 4.** Statistical results of the mediation analysis (pain intensity

1465 mediated the relationship between decreased BOP and monetary donations) in

1466 Experiment 2.

1467 • **Supplementary File 5.** Statistical results of the mediation analysis (pain intensity

1468 mediated the relationship between enhanced BOP and monetary donations) in

1469 Experiment 2.

1470 • **Supplementary File 6.** Statistical results of the mediation analysis (unpleasantness

1471 mediated the relationship between decreased BOP and monetary donations) in

1472 Experiment 2.

1473 • **Supplementary File 7.** Statistical results of the mediation analysis (unpleasantness

1474 mediated the relationship between enhanced BOP and monetary donations) in

1475 Experiment 2.

1476   • **Supplementary file 8.** Statistical results of reaction times, accuracies, and rating

1477   scores (mean ± SD) in Experiment 3.

1478   • **Supplementary file 9.** Statistical results of mean ERP amplitudes (mean ± SD) in

1479   Experiment 3.

1480   • **Supplementary file 10.** Statistical results of reaction times, accuracies, and rating

1481   scores (mean ± SD) in Experiment 4.

1482   • **Supplementary file 11.** Statistical results of mean ERP amplitudes (mean ± SD) in

1483   Experiment 4.

1484   • **Supplementary file 12.** Statistical results of reaction times, accuracies, and mean

1485   ERP amplitudes (mean ± SD) in Experiment 5.

1486   • **Supplementary file 13.** Results of the serial mediation analysis in Experiment 5.

1487   • **Supplementary file 14.** Statistical results of reaction times, accuracies and rating

1488   scores (mean ± SD) in Experiment 6

1489   • **Supplementary file 15.** Number of ERP trials for analyses (mean ± SD) in

1490   Experiments 3-5.

1491   **Data availability**

1492   All data generated or analyzed for figures of this study are included in the manuscript

1493   and supporting files. Source data files have been provided for Figures 1-6 and Appendix

1494   1 Figure 1.

1495   **Code availability**

1496 Code files used to analyze the data and to generate the figures that support the findings

1497 of this study have been uploaded.

## References

1499 Ais, J., Zylberberg, A., Barttfeld, P. & Sigman, M. Individual consistency in the
1500 accuracy and distribution of confidence judgments. *Cognition* **146**, 377-386 (2016).
1501

1502 Atlas, L. Y. & Wager, T. D. A Meta-analysis of brain mechanisms of placebo analgesia:
1503 consistent findings and unanswered questions. *Handb. Exp. Pharmacol.* **225**, 37-69
1504 (2014).
1505

1506 Avenanti, A., Bueti, D., Galati, G. & Aglioti, S. M. Transcranial magnetic stimulation
1507 highlights the sensorimotor side of empathy for pain. *Nat. Neurosci.* **8**, 955-960 (2005).
1508

1509 Avenanti, A., Sirigu, A. & Aglioti, S. M. Racial bias reduces empathic sensorimotor
1510 resonance with other-race pain. *Curr. Biol.* **20**, 1018-1022 (2010).
1511

1512 Azevedo, R. T., Macaluso, E., Avenanti, A., Santangelo, V., Cazzato, V., & Aglioti, S.
1513 M. Their pain is not our pain: brain and autonomic   correlates of empathic resonance
1514 with the pain of same and different race individuals. *Hum. Brain Mapp.* **34**, 3168-3181
1515 (2013).
1516

1517 Batson, C. D. Prosocial motivation: is it ever truly altruistic? *Adv. Exp. Soc. Psychol.* **20**,
1518 65-122 (1987).
1519

1520 Batson, C. D., Lishner, D. A. & Stocks, E. L. In *The Oxford Handbook of Prosocial*
1521 *Behavior* (eds Schroeder, D. A. & Graziano, W. G.) 259-268 (Oxford Univ. Press,
1522 2015).
1523

1524 Bieri, D., Reeve, R. A., Champion, G. D., Addicoat, L., & Ziegler, J. B. The Faces Pain
1525 Scale for the self-assessment of the severity of pain experienced by children:
1526 development, initial validation, and preliminary investigation for ratio scale properties.
1527 *Pain*, **41**, 139-150 (1990).
1528

1529 Botvinick, M., Jha, A. P., Bylsma, L. M., Fabian, S. A., Solomon, P. E., & Prkachin, K.
1530 M.. Viewing facial expressions of pain engages cortical areas involved in the direct
1531 experience of pain. *NeuroImage*, **25**, 312-319 (2005).
1532

1533 Calvo, M. G., Marrero, H. & Beltrán, D. When does the brain distinguish between
1534 genuine and ambiguous smiles? An ERP study. *Brain Cogn.* **81**, 237-246 (2013).

1535

1536   Clark, A. Whatever next? Predictive brains, situated agents, and the future of cognitive
1537   science. *Behav. Brain Sci.* **36**, 181-204 (2013).

1538

1539   Cohen, J. *Statistical Power Analysis for the Behavioral Sciences* 2nd edn (Lawrence
1540   Erlbaum Associates, 1988).

1541

1542   Coll, M. P. Meta-analysis of ERP investigations of pain empathy underlines
1543   methodological issues in ERP research. *Soc. Cogn. Affect. Neurosci.* **13**, 1003-1017
1544   (2018).

1545

1546   Colloca, L. & Benedetti, F. Placebos and painkillers: is mind as real as matter? *Nat. Rev.*
1547   *Neurosci.* **6**, 545-552 (2005).

1548

1549   Cui, F., Ma, N., & Luo, Y. J. Moral judgment modulates neural responses to the
1550   perception of other's pain: an ERP study. *Sci. Rep.*, **6**, 1-8. (2016).

1551

1552   Davis, M. H. Measuring individual differences in empathy: Evidence for a
1553   multidimensional approach. *J. Pers. Soc. Psychol.* **44**, 113-126 (1983).

1554

1555   Decety, J., & Jackson, P. L. The functional architecture of human empathy. *Behav.*
1556   *Cogn. Neurosci. Rev.* **3**, 71-100 (2004)..

1557

1558   Decety, J. The neuroevolution of empathy. A*nn. NY Acad. Sci.* **1231**, 35-45 (2011).

1559

1560   Decety, J., Bartal, I. B., Uzefovsky, F. & Knafo-Noam, A. Empathy as a driver of
1561   prosocial behaviour: highly conserved neurobehavioural mechanisms across species.
1562   *Philos. Trans. R. Soc. Lond. B Biol. Sci.* **371**, 1686 (2016).

1563

1564   De Waal, F. B. Putting the altruism back into altruism: the evolution of empathy. *Annu.*
1565   *Rev. Psychol.* **59**, 279-300 (2008).

1566

1567   Drwecki, B. B., Moore, C. F., Ward, S. E. & Prkachin, K. M. Reducing racial disparities
1568   in pain treatment: the role of empathy and perspective-taking. *Pain* **152**, 1001-1006
1569   (2011).

1570

1571   Edele, A., Dziobek, I. & Keller, M. Explaining altruistic sharing in the dictator game:
1572   the role of affective empathy, cognitive empathy, and justice sensitivity. *Learn. Indi.*
1573   *Diffe.* **24**, 96-102 (2013).

1574

Eisenberg, N., Eggum, N. D. & Di Giunta, L. Empathy-related responding: associations with prosocial behavior, aggression, and intergroup relations. *Soc. Issues Policy Rev.* **4**, 143-180 (2010).

Etkin, A., Büchel, C. & Gross, J. J. The neural bases of emotion regulation. *Nat. Rev. Neurosci.* **16**, 693-700 (2015).

Fan, Y. & Han, S. Temporal dynamic of neural mechanisms involved in empathy for pain: an event-related brain potential study. *Neuropsychologia* **46**, 160-173 (2008).

Fan, Y., Duncan, N. W., De Greck, M. & Northoff, G. Is there a core neural network in empathy? An fMRI based quantitative meta-analysis. *Neurosci. Biobehav. Rev.* **35**, 903-911 (2011).

Faul, F., Erdfelder, E., Lang, A. G. & Buchner, A. G*Power 3: a flexible statistical power analysis program for the social, behavioral, and biomedical sciences. *Behav. Res. Methods* **39**, 175-191 (2007).

Fodor, J. A. Précis of the modularity of mind. *Behav. Brain Sci.* **8**, 1-42 (1985).

Friston, K. A theory of cortical responses. *Philos. Trans. R. Soc. Lond. B Biol. Sci.* **360**, 815-836 (2005).

Fuentes, A. *Why We Believe: Evolution and the Human Way of Being.* 1st edn (Yale Univ. Press, 2019).

Garcés, M. & Finkel, L. Emotional Theory of Rationality. *Front. Integr. Neurosci.* **13**, 11-35 (2019).

Greenwald, A. G., McGhee, D. E. & Schwartz, J. L. Measuring individual differences in implicit cognition: the implicit association test. *J. Pers. Soc. Psychol.* **74**, 1464-1480 (1998).

Greenwald, A. G., Nosek, B. A. & Banaji, M. R. Understanding and using the implicit association test: I. An improved scoring algorithm. *J. Pers. Soc. Psychol.* **85**, 197-216 (2003).

Gu, X. & Han, S. Attention and reality constraints on the neural processes of empathy for pain. *NeuroImage* **36**, 256-267 (2007).

Hampton, R. S., & Varnum, M. E. The cultural neuroscience of emotion regulation. *Cult. Brain*, **6**, 130-150 (2018).

Han, S. et al. Empathic neural responses to others' pain are modulated by emotional contexts. *Hum. Brain Mapp.* **30**, 3227-3237 (2009).

Han, S. Neurocognitive basis of racial ingroup bias in empathy. *Trends Cogn. Sci.* **22**, 400-421 (2018).

Han, X., Luo, S., & Han, S. Embodied neural responses to others' suffering. *Cogn. Neurosci.* **7**, 114-127. (2016).

Han, X. et al. Empathy for pain motivates actions without altruistic effects: evidence of motor dynamics and brain activity. *Soc. Cogn. Affect. Neuro.* **12**, 893-901 (2017).

Han, X., Ashar, YK., Kragel, P., Petre, B., Schelkun, V., Atlas, LY., Chang, LJ., Jepma, M., Koban, L., Losin, ERA., Roy, M., Woo, CW. & Wager TD. Effect sizes and test-retest reliability of the fMRI-based Neurologic Pain Signature. *bioRxiv: the preprint server for biology, 2021.05.29.445964.*

Harris, S. et al. The Neural Correlates of Religious and Nonreligious Belief. *PLoS One* **4**, e0007272 (2009).

Hayes, A. F. *Introduction to mediation, moderation, and conditional process analysis: A regression-based approach* 2nd edn (The Guilford Press, 2017).

Hein, G. et al. Neural responses to ingroup and outgroup members' suffering predict individual differences in costly helping. *Neuron* **68**, 149-160 (2010).

Hein, G. et al. The brain's functional network architecture reveals human motives. *Science* **351**, 1074-1078 (2016).

Hofman, M. L. In *Handbook of Emotions* (eds Lewis, M., Haviland-Jones, J. M. & Barrett, L. F.) 440-455 (The Guilford Press, 2008).

Hoffman, K. M., Trawalter, S., Axt, J. R., & Oliver, M. N. Racial bias in pain assessment and treatment recommendations, and false beliefs about biological differences between blacks and whites. *Proc. Natl. Acad. Sci. U.S.A.* **113**, 4296-4301 (2016).

Jackson, P. L., Meltzoff, A. N. & Decety, J. How do we perceive the pain of others? A window into the neural processes involved in empathy. *NeuroImage* **24**, 771-779 (2005).

1659  Jauniaux, J., Khatibi, A., Rainville, P. & Jackson, P. L. A meta-analysis of
1660  neuroimaging studies on pain empathy: investigating the role of visual information and
1661  observers' perspective. *Soc. Cogn. Affect. Neurosci.* **14**, 789-813 (2019).
1662
1663  Kenny, D. A., Kashy, D. A. & Bolger, N. In *The Handbook of Social Psychology* (eds
1664  Gilbert, D. T., Fiske, S. T. & Lindzey, G.) 233-265 (Oxford Univ. Press, 1998).
1665
1666  Knill, D. C. & Pouget, A. The Bayesian brain: the role of uncertainty in neural coding
1667  and computation. *Trends Neurosci.* **27**, 712-719 (2004).
1668
1669  Kriegeskorte, N., Goebel, R. & Bandettini, P. Information-based functional brain
1670  mapping. *Proc. Natl. Acad. Sci. U.S.A.* **103**, 3863-3868 (2006).
1671
1672  Krishnan, A., Woo, C. W., Chang, L. J., Ruzic, L., Gu, X., López-Solà, M., Jackson,
1673  P.L., Pujol, J., Fan, J., & Wager, T. D. Somatic and vicarious pain are represented by
1674  dissociable multivariate brain patterns. *eLife*, **5**, e15166 (2016).
1675
1676  Lamm, C., Nusbaum, H. C., Meltzoff, A. N. & Decety, J. What are you feeling? Using
1677  functional magnetic resonance imaging to assess the modulation of sensory and
1678  affective responses during empathy for pain. *PLoS One* **2**, e1292 (2007).
1679
1680  Lamm, C., Decety, J. & Singer, T. Meta-analytic evidence for common and distinct
1681  neural networks associated with directly experienced pain and empathy for pain.
1682  *NeuroImage* **54**, 2492-2502 (2011).
1683
1684  Lamm, C., Rütgen, M. & Wagner, I. C. Imaging empathy and prosocial emotions.
1685  *Neurosci. Lett.* **693**, 49-53 (2019).
1686
1687  Li, W. & Han, S. Perspective taking modulates event-related potentials to perceived
1688  pain. *Neurosci. Lett.* **469**, 328-332 (2010).
1689
1690  Li, W. & Han, S. Behavioral and electctrophysiological evidence for enhanced
1691  sensitivity to subtle variations of pain expressions of same-race than other-race faces.
1692  *Neuropsychologia* **129**, 302-309 (2019).
1693
1694  Li, Z., Yu, J., Yang, X. & Zhu, L. Associations between empathy and altruistic sharing
1695  behavior in Chinese adults. *J. Gen. Psychol.* **146**, 1-16 (2019).
1696
1697  Luck, S. J. & Gaspelin, N. How to get statistically significant effects in any ERP
1698  experiment (and why you shouldn't). *Psychophysiology* **54**, 146-157 (2017).
1699

1700 Luo, S. et al. Reminders of mortality decrease midcingulate activity in response to
1701 others' suffering. *Soc. Cogn. Affect. Neuro.* **9**, 477-486 (2014).

1702

1703 Luo, S. et al. Oxytocin receptor gene and racial ingroup bias in empathy-related brain
1704 activity. *NeuroImage* **110**, 22-31 (2015).

1705

1706 Luo, S., Han, X., Du, N. & Han, S. Physical coldness enhances racial in-group bias in
1707 empathy: Electrophysiological evidence. *Neuropsychologia* **116**, 117-125 (2018).

1708

1709 Luo, W. et al. Three stages of facial expression processing: ERP study with rapid serial
1710 visual presentation. *NeuroImage* **49**, 1857-1867 (2010).

1711

1712 Mackinnon, D. P., Lockwood, C. M. & Williams, J. Confidence limits for the indirect
1713 effect: distribution of the product and resampling methods. *Multivariate Behav. Res.* **39**,
1714 99-128 (2004).

1715

1716 Martel, M. O., Thibault, P., Roy, C., Catchlove, R., & Sullivan, M. J. L. Contextual
1717 determinants of pain judgments. *Pain* **139**, 562-568 (2008).

1718

1719 Mathur, V. A., Harada, T., Lipke, T. & Chiao, J. Y. Neural basis of extraordinary
1720 empathy and altruistic motivation. *NeuroImage* **51**, 1468-1475 (2010).

1721

1722 Mckay, R. T. & Dennett, D. C. The evolution of misbelief. *Behav. Brain Sci.* **32**,
1723 493-510 (2009).

1724

1725 Meissner, K. et al. The placebo effect: advances from different methodological
1726 approaches. *J. Neurosci.* **31**, 16117-16124 (2011).

1727

1728 Millikan, R. G. *White Queen Psychology and Other Essays for Alice.* 1st edn (MIT
1729 Press, 1995).

1730

1731 Morey, R. D. & Rouder, J. N. BayesFactor: Computation of Bayes Factors for common
1732 designs. R package version 0.9.11-11 (2015).

1733

1734 Nicolardi, V., Panasiti, M. S., D'Ippolito, M., Pecimo, G. L., & Aglioti, S. M. Pain
1735 perception during social interactions is modulated by self-related and moral contextual
1736 cues. *Sci. Reports* **10**, 1-12 (2020).

1737

1738 Nili, H. et al. A toolbox for representational similarity analysis. *PLoS. Comput. Biol.* **10**,
1739 e1003553 (2014).

1740

1741 Ochsner, K. N. & Gross, J. J. The cognitive control of emotion. *Trends Cogn. Sci.* **9**,
1742 242-249 (2005).
1743

1744 Penner, L. A., Dovidio, J. F., Piliavin, J. A. & Schroeder, D. A. Prosocial behavior:
1745 multilevel perspectives. *Annu. Rev. Psychol.* **56**, 365-392 (2005).
1746

1747 Petrovic, P. et al. Placebo in emotional processing--induced expectations of anxiety
1748 relief activate a generalized modulatory network. *Neuro.* **46**, 957-969 (2005).
1749

1750 Preacher, K. J., Rucker, D. D. & Hayes, A. F. Addressing Moderated Mediation
1751 hypotheses: theory, methods, and prescriptions. *Multivariate Behav. Res.* **42**, 185-227
1752 (2007).
1753

1754 Preacher, K. J. & Hayes, A. F. Asymptotic and resampling strategies for assessing and
1755 comparing indirect effects in multiple mediator models. *Behav. Res. Methods* **40**,
1756 879-891 (2008).
1757

1758 Preacher, K. J. & Hayes, A. F. In *The SAGE Sourcebook of Advanced Data Analysis*
1759 *Methods for Communication Research* (eds Hayes, A. F., Slater, M. D. & Snyder, L. B.)
1760 13-54 (SAGE Publications, 2008).
1761

1762 Preston, S., D. & De Waal, F. B. Empathy: Its ultimate and proximate bases. *Behav.*
1763 *Brain Sci.* **25**, 1-20 (2002).
1764

1765 Prkachin, K. M., & Rocha, E. M. High levels of vicarious exposure bias pain judgments.
1766 *J. Pain* **11**, 904-909 (2010).
1767

1768 Régner, I., Thinus-Blanc, C., Netter, A., Schmader, T., & Huguet, P. Committees with
1769 implicit biases promote fewer women when they do not believe gender bias exists. *Nat.*
1770 *Hum. Behav.* **3**, 1171-1179 (2019).
1771

1772 Saarela, M. V., Hlushchuk, Y., Williams, A. C. D. C., Schürmann, M., Kalso, E., &
1773 Hari, R. The compassionate brain: humans detect intensity of pain from another's face.
1774 *Cereb. Cortex*, **17**, 230-237 (2007).
1775

1776 Schnell, K., Bluschke, S., Konradt, B., & Walter, H. (2011). Functional relations of
1777 empathy and mentalizing: an fMRI study on the neural basis of cognitive empathy.
1778 *NeuroImage* **54**, 1743-1754.
1779

1780 Seitz, R. J. & Angel, H. F. Belief formation – A driving force for brain evolution. *Brain*
1781 *Cogn.* **140**, 105548 (2020).
1782

1783  Semlitsch, H. V., Anderer, P., Schuster, P. & Presslich, O. A solution for reliable and
1784  valid reduction of ocular artifacts, applied to the P300 ERP. *Psychophysiology* **23**,
1785  695-703 (1986).

1787  Shamay-Tsoory, S. G., Aharon-Peretz, J. & Perry, D. Two systems for empathy: a
1788  double dissociation between emotional and cognitive empathy in inferior frontal gyrus
1789  versus ventromedial prefrontal lesions. *Brain* **132**, 617-627 (2009).

1791  Shamay-Tsoory, S. G. The neural bases for empathy. *Neuroscientist* **17**, 18-24 (2011).

1793  Sheng, F. & Han, S. Manipulations of cognitive strategies and intergroup relationships
1794  reduce the racial bias in empathic neural responses. *NeuroImage* **61**, 786-797 (2012).

1796  Sheng, F., Liu, Y., Zhou, B., Zhou, W., & Han, S. Oxytocin modulates the racial bias in
1797  neural responses to others' suffering. *Bio. Psychol.* **92**, 380-386 (2013).

1799  Sheng, F. et al. Task modulations of racial bias in neural responses to others' suffering.
1800  *NeuroImage* 263-270 (2014).

1802  Sheng, F., Han, X. & Han, S. Dissociated Neural Representations of Pain Expressions
1803  of Different Races. *Cereb. Cortex* **26**, 1221-1233 (2016).

1805  Shrout, P. E. & Bolger, N. Mediation in experimental and nonexperimental studies: new
1806  procedures and recommendations. *Psychol. Methods* **7**, 422-445 (2002).

1808  Singer, T. et al. Empathy for pain involves the affective but not sensory components of
1809  pain. *Science* **303**, 1157-1162 (2004).

1811  Steiger, J. H. Beyond the F test: Effect size confidence intervals and tests of close fit in
1812  the analysis of variance and contrast analysis. *Psychol. Methods* **9**, 164-182 (2004).

1814  Sterzer, P., Frith, C. & Petrovic, P. Believing is seeing: expectations alter visual
1815  awareness. *Curr. Biol.* **18**, 697-698 (2008).

1817  Tappin, B. M. & Gadsby, S. Biased belief in the Bayesian brain: A deeper look at the
1818  evidence. *Conscious Cogn.* **68**, 107-114 (2019).

1820  Twigg, O. C., & Byrne, D. G. The influence of contextual variables on judgments about
1821  patients and their pain. Pain Med. **16**, 88-98 (2015).

1823  Varnum, M. E., Blais, C., Hampton, R. S., & Brewer, G. A. Social class affects neural
1824  empathic responses. *Cult. Brain* **3**, 122-130 (2015).

1825

1826 Völlm, B.A., Taylor, A.N., Richardson, P., Corcoran, R., Stirling, J., McKie, S., Deakin,
1827 J.F. & Elliott, R. Neuronal correlates of theory of mind and empathy: a functional
1828 magnetic resonance imaging study in a nonverbal task. *NeuroImage* **29**, 90-98 (2006).
1829

1830 Wager, T. D. et al. Placebo-induced changes in FMRI in the anticipation and experience
1831 of pain. *Science* **303**, 1162-1167 (2004).
1832

1833 Wager, T. D. & Atlas, L. Y. The neuroscience of placebo effects: connecting context,
1834 learning and health. *Nat. Rev. Neurosci.* **16**, 403-418 (2015).
1835

1836 Wang, C., Wu, B., Liu, Y., Wu, X., & Han, S. Challenging emotional prejudice by
1837 changing self-concept: priming independent self-construal reduces racial in-group bias
1838 in neural responses to other's pain. *Soc Cogn. Affect. Neurosci.* **10**, 1195-1201 (2015).
1839

1840 Williams, L. M. et al. The 'when' and 'where' of perceiving signals of threat versus
1841 non-threat. *NeuroImage* **31**, 458-467 (2006).
1842

1843 Wright, T. L. & Tedeschi, R. G. Factor anlaysis of the Interpersonal Trust Scale. *J.*
1844 *Consult. Clin. Psychol.* **43**, 470-477 (1975).
1845

1846 Xu, X., Zuo, X., Wang, X. & Han, S. Do you feel my pain? Racial group membership
1847 modulates empathic neural responses. *J. Neurosci.* **29**, 8525-8529 (2009).
1848

1849 Yan, X., Yong, X., Huang, W. & Ma, Y. Placebo treatment facilitates social trust and
1850 approach behavior. *Proc. Natl. Acad. Sci. U.S.A.* **115**, 5732-5737 (2018).
1851

1852 Zhou, F., Li, J., Zhao, W., Xu, L., Zheng, X., Fu, M., Yao, S., Kendrick, K.M., Wager,
1853 T.D., & Becker, B. Empathic pain evoked by sensory and emotional-communicative
1854 cues share common and process-specific neural representations. *eLife*, **9**, e56929 (2020).
1855

1856 Zhao, X., Lynch JR, J. G. & Chen, Q. Reconsidering Baron and Kenny: Myths and
1857 truths about mediation analysis. *J. Consumer Research* **37**, 197-206 (2010).
1858

1859 Zhou, Y., Gao, T., Zhang, T., Li, W., Wu, T., Han, X., & Han, S. Neural dynamics of
1860 racial categorization predicts racial bias in face recognition and altruism. *Nat. Hum.*
1861 *Behav.* **4**, 69-87. (2020).
1862

1863 Zhou, Y. & Han, S. Neural dynamics of pain expression processing: Alpha-band
1864 synchronization to same-race pain but desynchronization to other-race pain.
1865 *NeuroImage* **224**, 117400 (2021).
1866

1875 **Competing interests**

1876 The authors declare no competing interests.

# Appendix 1

Our results in Experiments 1 to 6 showed consistent evidence for modulations of

both subjective (self-report) and objective (EEG/fMRI) measures of empathy for others'

suffering. An interesting question arising from these findings is whether the belief

effects are specific to neural underpinnings of empathy for pain. We addressed this issue

by examining belief effects on neural responses to other facial expressions in an

additional experiment. Specifically, in this experiment, we sought to test (1) whether

beliefs of authenticity of others' happiness influence brain responses to perceived happy

expressions, and (2) whether beliefs also modulate neural responses to happy

expressions in the P2 time window, similar to the BOP effect on ERPs to pain

expressions. The paradigm used in the additional experiment was the same as that used

in Experiment 3 except the following. We asked an independent sample of participants

to remember identities (awardees or actors/actresses) of neutral faces. Thereafter, EEG

signals to happy and neutral expressions of awardees or actors/actresses were recorded

after informing the participants that photos of happy faces were taken from awardees

who were smiling when receiving awards whereas actors/actresses imitated others'

smiling and showed happy expressions. We predicted that beliefs that actors/actresses'

expressions do not reflect their actual emotional states would decrease brain response to

happy expressions. We tested this prediction by comparing ERPs to happy/neutral faces

with awardee or actor/actress identities.

21       We recorded EEG signals from an independent sample of healthy young adults (N

22   = 30 males, mean age ± s.d. = 22.30 ± 2.73 years). Face stimuli with happy or neutral

23   expressions were adopted from the previous study (Wang and Han, 2021). There were

24   photos of 16 Chinese models (half males) and each model contributed one photo with

25   happy expression and one with neutral expression.

26       The participants were first presented with the faces with neutral expressions and

27   were informed that these photos were taken from 8 awardees who recently obtained

28   awards and from 8 actors/actresses. After the identity memory task, in which the

29   participants were able to correctly recognize all faces with awardee or actor/actress

30   identities, they were asked to perform identity judgments on faces with neutral or happy

31   expressions by pressing one of two buttons while EEG was recorded. After EEG

32   recording, the participants were presented with each happy face again and had to rate

33   how happy the person is feeling (i.e., happiness intensity) by rating on a Likert-type

34   scale (1 = not happy at all; 7 = extremely happy).

35       An ANOVA of the mean rating scores of happiness intensity with Identity (awardee

36   vs. actor/actress) and Expression (happy vs. neutral) as within-subject variables revealed

37   significant main effects of Identity ($F(1,29) = 19.512$, $p < 0.001$, $\eta_p^2 = 0.402$, 90% CI =

38   (0.166, 0.560)) and Expression ($F(1,29) = 422.774$, $p < 0.001$, $\eta_p^2 = 0.936$, 90% CI =

39   (0.889, 0.953)), and a significant Identity × Expression interaction ($F(1,29) = 6.610$, $p =$

40   0.016, $\eta_p^2 = 0.186$, 90% CI = (0.021, 0.372), see Appendix 1-Figure 1a, and Appendix
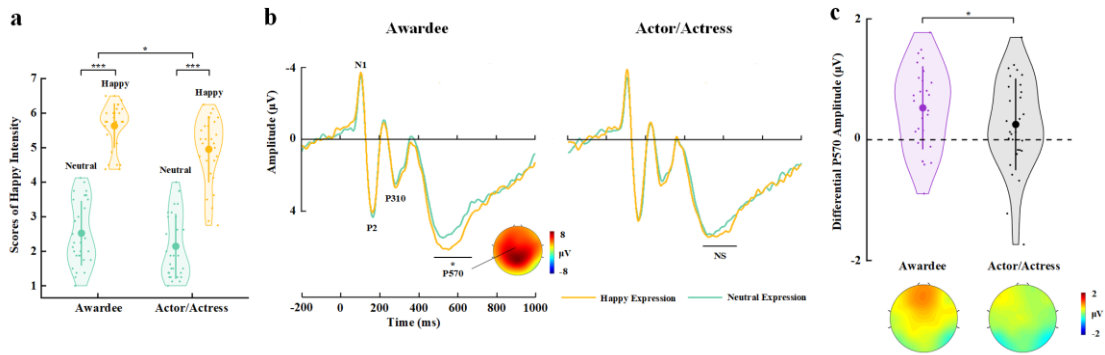
41    1-table 1 for details). The results suggest weaker subjective feelings of happiness

42    intensity for faces with actor/actress identities compared to awardee identities.

43        The participants responded to face identities with high accuracies during EEG

44    recording (>88% across all conditions, see Appendix 1-table 1 for details). Similarly,

45    ERPs to face stimuli in this experiment were characterized by an early negative activity

46    at 90-120 ms (N1) and a positive activity at 175-195 ms (P2) at the frontal/central

47    regions, which were followed by two positive activities at 280-340 ms (P310) over the

48    parietal region and 500–700 ms (P570) over the frontal area (Appendix 1-Figure 1b).

49    ANOVAs of the P2 amplitudes with Identity (awardee vs. actor/actress) and Expression

50    (happy vs. neutral) as within-subject variables did not reveal a significant Identity $\times$

51    Expression interaction ($F(1,29) = 0.441$, $P = 0.512$, $\eta_p^2 = 0.015$, 90% CI = (0, 0.145),

52    Bayes factors = 0.303).

53        Importantly, ANOVAs of the later P570 amplitudes showed a significant Identity $\times$

54    Expression interaction ($F(1,29) = 4.832$, $P = 0.036$, $\eta_p^2 = 0.143$, 90% CI = (0.005,

55    0.328), Appendix 1-Figure 1b and 1c, see Appendix 1-table 1 for statistical details).

56    Simple effect analyses indicated significantly larger P570 amplitudes in response to

57    happy versus neutral expressions of awardees' faces ($F(1,29) = 20.880$, $p < 0.001$, $\eta_p^2 =$

58    0.419, 90% CI = (0.181, 0.573)), but not of actors/actresses' faces ($F(1,29) = 3.375$, $p =$

59    0.076, $\eta_p^2 = 0.104$, 90% CI = (0, 0.285), Bayes factor = 0.858).

**Appendix 1-Figure 1.** EEG results of the additional experiment. (a) Mean rating scores of happy intensity related to happy and neutral expressions of faces with awardee or actor/actress identities. (b) ERPs to faces with awardee or actor/actress identities at frontal electrodes. The voltage topography shows the scalp distribution of the P570 amplitude with the maximum over the central/parietal region. (c) Mean differential P570 amplitudes to happy versus neutral expressions of faces with awardee or actor/actress identities. The voltage topographies illustrate the scalp distribution of the P570 difference waves to happy (vs. neutral) expressions of faces with awardee or actor/actress identities, respectively. Shown are group means (large dots), standard deviation (bars), measures of each individual participant (small dots), and distribution (violin shape) in (a) and (c). The online version of this article includes the following source data for Appendix 1-Figure 1: Appendix 1-Figure 1-Source data 1.

Our behavioral and ERP results in this experiment suggest reduced subjective feelings and brain responses to happy (vs. neutral) expressions of actors/actresses' faces compared to awardees' faces. These results support the prediction that beliefs that actors/actresses' expressions do not reflect their actual emotional states decrease brain response to happy expressions. However, belief effects on brain responses to happy expressions were observed on the P570 amplitudes but not on the P2 amplitudes. This is different from our ERP results in in Experiments 3-5, in which we showed evidence that BOP modulated the P2 amplitudes. These results suggest general belief modulation effects on brain activities involved in processing of facial expressions. In addition, our

84   results suggest that the time window in which beliefs modulate brain responses to facial

85   expressions depends on the nature of facial expressions (e.g., pain or happiness

86   expressions).

87   **Reference**

88   Wang, X., Han, S. Processing of facial expressions of same-race and other-race faces:

89   distinct and shared neural underpinnings. *Soc. Cogn. Affect. Neurosci*. **16**, 576-592

90   (2021).

91

**Appendix 1-table 1**. RTs, accuracies, rating scores, numbers of ERP trials, and ERP amplitudes (mean ± SD) in the additional experiment.

| | Awardee | | Actor/Actress | |
| --- | --- | --- | --- | --- |
| | Neutral | Happy | Neutral | Happy |
| **RT (ms)** | 654±63 | 657±60 | 666±64 | 680±66 |
| **Accuracy (%)** | 92±4.9 | 90±7.5 | 92±5.4 | 88±8.7 |
| **Happy Intensity** | 2.525±0.94 | 5.638±0.64 | 2.146±0.94 | 4.95±0.96 |
| **N1 amplitude (μV)** | -2.267±1.69 | -2.606±1.75 | -2.297±1.43 | -2.620±1.52 |
| **P2 amplitude (μV)** | 2.544±2.64 | 2.375±2.30 | 2.940±2.56 | 2.593±2.56 |
| **P310 amplitude (μV)** | 3.449±3.45 | 3.445±3.30 | 3.492±3.38 | 3.376±3.38 |
| **P570 amplitude (μV)** | 4.677±2.22 | 5.379±2.15 | 4.696±2.16 | 4.950±2.11 |
| **ERP trials** | 114±10 | 110±13 | 113±11 | 108±12 |

| | Statist Value | ANOVA | | | Simple effect (Identity) | |
| --- | --- | --- | --- | --- | --- | --- |
| | | Identity | Expression | Identity*Expression | Awardee | Actor/Actress |
| **RT (ms)** | F | 13.229 | 11.256 | 4.733 | 0.915 | 13.230 |
| | P | 0.001 | 0.002 | 0.038 | 0.347 | 0.001 |
| | $\eta_p^2$ | 0.313 | 0.280 | 0.140 | 0.031 | 0.313 |
| | 90% | (0.094, 0.488) | (0.071, 0.459) | (0.004, 0.326) | (0, 0.180) | (0.094, 0.488) |
| **Accuracy (%)** | F | 0.496 | 40.590 | 0.595 | | |
| | P | 0.487 | <0.001 | 0.447 | | |
| | $\eta_p^2$ | 0.017 | 0.583 | 0.020 | | |
| | 90% | (0, 0.150) | (0.362, 0.698) | (0, 0.158) | | |
| **Happy Intensity** | F | 19.512 | 422.774 | 6.610 | 433.364 | 302.128 |
| | P | <0.001 | <0.001 | 0.016 | <0.001 | <0.001 |
| | $\eta_p^2$ | 0.402 | 0.936 | 0.186 | 0.937 | 0.912 |
| | 90% | (0.166, 0.560) | (0.889, 0.953) | (0.021, 0.372) | (0.892, 0.955) | (0.849, 0.937) |
| **N1 (95-115ms)** | F | 0.031 | 9.890 | 0.005 | | |
| | P | 0.862 | 0.004 | 0.944 | | |
| | $\eta_p^2$ | 0.001 | 0.254 | 0.0002 | | |
| | 90% | (0, 0.041) | (0.055, 0.436) | (0, 0.007) | | |
| **P2 (175-195ms)** | F | 6.476 | 2.822 | 0.441 | | |
| | P | 0.017 | 0.104 | 0.512 | | |
| | $\eta_p^2$ | 0.183 | 0.089 | 0.015 | | |
| | 90% | (0.019, 0.369) | (0, 0.266) | (0, 0.145) | | |
| **P310 (280-340ms)** | F | 0.012 | 0.140 | 0.252 | | |
| | P | 0.913 | 0.711 | 0.619 | | |
| | $\eta_p^2$ | 0.0004 | 0.005 | 0.009 | | |
| | 90% | (0, 0.017) | (0, 0.106) | (0, 0.125) | | |
| **P570 (500-700ms)** | F | 1.948 | 20.752 | 4.832 | 20.880 | 3.375 |
| | P | 0.173 | <0.001 | 0.036 | <0.001 | 0.076 |

| | | | | | |
|---|---|---|---|---|---|
| $\eta_p^2$ | 0.063 | 0.417 | 0.143 | 0.419 | 0.104 |
| **90%** | (0, 0.232) | (0.180, 0.572) | (0.005, 0.328) | (0.181, 0.573) | (0, 0.285) |

94    Note: Effect size is indexed as the partial eta-squared value. The 90% CIs are reported for partial

95    eta-squared value.