

1 **Behavioral and neural evidence that robots are implicitly perceived as a threat**

2 Zhengde Wei^{1, +}, Ying Chen^{2, +}, Jiecheng Ren¹, Yi Piao¹, Pengyu Zhang¹, Qian Zhao¹, Rujing Zha¹,

3 Bensheng Qiu³, Daren Zhang¹, Yanchao Bi⁴, Shihui Han⁵, Chunbo Li^{6*}, Xiaochu Zhang^{1, 2, 3, 7*}

4 1 School of Life Sciences, Division of Life Science and Medicine, University of Science & Technology

5 of China, Hefei, Anhui 230027, China

6 2 School of Humanities & Social Science, University of Science & Technology of China, Hefei, Anhui

7 230026, China

8 3 Centers for Biomedical Engineering, School of Information Science and Technology, University of

9 Science & Technology of China, Hefei, Anhui 230027, China

10 4 State Key Laboratory of Cognitive Neuroscience and Learning, Beijing Normal University, Beijing

11 100875, China

12 5 School of Psychological and Cognitive Sciences, Peking University, Beijing, China

13 6 Shanghai Key Laboratory of Psychotic Disorders, Shanghai Mental Health Center, Shanghai Jiao

14 Tong University School of Medicine, Shanghai 200030, China

15 7 Academy of Psychology and Behavior, Tianjin Normal University, Tianjin, 300387, China

16

17 + These authors contributed equally to this work.

18 * Correspondence: licb@smhc.org.cn; zxcustc@ustc.edu.cn

19

20 **Abstract**

21 A deeper understanding of the human side of human-robot interaction determined by studying the
22 human brain when we perceive robots should help solve the biggest challenges of successful social
23 encounters with robots. However, current social neuroscience studies mainly focus on explicit
24 perception of robots, and implicit perception of robots is rather unexplored. Here, our behavioral
25 analysis indicated that despite self-reported positive attitudes, participants had negative implicit
26 attitudes toward humanoid robots. Our neuroimaging analysis indicated that subthreshold presentation
27 of humanoid robot vs. human images led to significant left amygdala activation that was associated
28 with negative implicit attitudes toward robots. After successfully weakening the negative attitudes, the
29 left amygdala response to subthreshold presentation of humanoid robot images decreased, and the
30 decrease in the left amygdala response was positively associated with the decrease in negative attitudes.
31 Our results reveal that the processing of information about humanoid robots displays automaticity with
32 regard to the recruitment of amygdala activation. Our findings that people may implicitly perceive
33 humanoid robots as a threat may guide more appropriate interaction with social robots.

34 **Key words:** robots; threat; implicit perception; amygdala

35 **Significance statement**

36 Social interactions with robots are one of the biggest challenges in robotics, which necessitates a
37 deeper understanding of how people perceive robots. Our results reveal automaticity for processing
38 information about humanoid robots similar to that previously evident for threats. Given the effort
39 currently being put into the development of robots for daily assistance, studying implicit perception of
40 robots could be a step toward building smooth human-robot social relationships.

41 **Introduction**

42 Given that robots are becoming increasingly present in the everyday environment and that social
43 interactions with robots are one of the biggest challenges in robotics ¹, a deeper understanding of how
44 people perceive robots is necessary ². This is especially relevant given the effort currently being put
45 into the development of robots for daily assistance ^{3,4}. When a person perceives a social stimulus,
46 information about that stimulus is immediately and spontaneously activated, including attitudes ⁵ and
47 social stereotypes ⁶, and can consequently influence people's behavior ⁷. Similarly, perceptions toward
48 robots might automatically influence how we interact with robots. Studying the human brain with
49 neuroimaging techniques when we perceive robots will shed light on a deeper understanding of the
50 human side of human-robot interaction ⁸. However, current social neuroscience studies mainly focus on
51 explicit perception of robots ^{8,9}, and the implicit perception of robots is rather unexplored.

52 Researchers have proposed that our ideas about robots are mainly informed by science fiction
53 media ^{10,11}, and robots in science fiction are generally described as a threat ¹⁰. Recent statements from
54 some influential industry leaders have strengthened these fears ¹². These ideas that robots tend to be
55 dangerous entities that will threaten the survival of humanity might be rooted in our long-term
56 socialization experiences, leading to the possibility that people might perceive robots as a threat.
57 Empirical evidence has illustrated that people's explicit perception toward robots is inconsistent,
58 ranging from enthusiasm to fear and anxiety ^{11,13-18}, possibly because people's explicit perceptions or
59 attitudes are subject to recall bias and social desirability bias ¹⁹. Notably, implicit perception is
60 considered to be a highly stable evaluative representation stemming from long-term socialization
61 experiences ²⁰. The most important feature is that implicit perception is irrespective of whether a person
62 consciously considers the stimulus as good or bad ²¹. Decades of work have identified implicit

63 perception as having a crucial influence on thoughts and actions²². Here, we aimed to provide evidence
64 from implicit behavioral and neural measures to support the possibility that people might implicitly
65 perceive robots as a threat.

66 To accurately assess people's implicit perception of robots, it is important to implement implicit
67 multimethod approaches. The implicit association test (IAT) is the most frequently used implicit
68 approach²³ and is based on the response given by a participant when performing an association task²⁴.
69 IATs are commonly used to explore implicit attitudes toward sensitive issues, such as racism²⁵.
70 Furthermore, neuroimaging approaches might be more sensitive, as they reveal qualitative differences
71 in processing without the need for distinct behavioral measures of implicit processing²². The
72 preparedness model^{26, 27} hypothesizes that the neural basis of automatic threat processing is the
73 amygdala. The amygdala has long been known to play a key role in responding to emotionally relevant
74 stimuli, activating in response to images containing threatening or highly arousing features^{28, 29}. When
75 automatic and controlled evaluations of threats differ, more positive controlled processing can
76 moderate more negative automatic processing³⁰, which could account for the absence of significant
77 activation in the amygdala in response to suprathreshold presentations of threats^{31, 32}. Notably,
78 although controlled processing can eliminate amygdala activation caused by consciously presented
79 threats, the amygdala still showed significant activation caused by threats that were presented
80 unconsciously³⁰⁻³². We hypothesized that if humanoid robots are implicitly perceived as a threat, then
81 participants would be expected to show negative implicit attitudes toward robots and an amygdala
82 response to subthreshold presentation of robot images.

83 How people implicitly perceive humanoid robots is an intriguing question. People's perception of
84 robots relies on the robots' behaviors, abilities and appearance³³. More human characteristics are

85 consistently perceived in robots that look and act like humans³⁴. Here, we tested whether people
86 implicitly perceive humanoid robots as a threat. We pursued this topic by first measuring participants'
87 explicit and implicit attitudes toward humanoid robots. We then employed functional magnetic
88 resonance imaging (fMRI) and intervention to test the amygdala response to humanoid robots and its
89 causal relationship to implicit attitudes.

90 **Material and Methods**

91 **Participants.** One hundred fifty-two participants (age: 21.75±2.58 years; 61 females) took part in the
92 original implicit association test (IAT), 22 participants (age: 22.77±1.97 years; 14 females) performed
93 the Black IAT and the mixed IAT, 30 participants (age: 22.23±1.89 years; 20 females) performed the
94 positive adjective IAT and the available humanoid IAT, and 38 participants (age: 21.77±1.67 years; 23
95 females) performed the humanoid weapon IAT and pet-robot weapon IAT. Ninety-nine participants
96 took part in the intervention experiment. Among them, 69 participants were allocated to the positive
97 intervention group, and 30 participants were allocated to the neutral intervention group. Ten
98 participants in the positive intervention group were excluded because of excessive head movement
99 during scanning (TRs with motion over 0.3 mm were censored, and participants who had more than 15%
100 of their TRs censored were removed from further analysis); 36 participants discontinued the positive
101 intervention after completing preintervention tasks due to personal reasons; and 3 participants
102 discontinued the neutral intervention after completing the preintervention tasks because of withdrawal.
103 Finally, data from 23 participants (age: 20.39±2.06 years; 15 females) in the positive intervention group
104 and 27 participants (age: 22.13±1.49 years; 18 females) in the neutral group were analyzed. All
105 participants were right-handed, with normal or corrected-to-normal vision, without alcohol or drug
106 dependence disorders and without prior head injury. All participants and their immediate relatives had
107 no psychiatric disorders or history of psychiatric disorders. All participants were nonpsychology
108 undergraduates, and they were unaware of the purpose and contents of the study before they
109 participated in the experiment.

110 The sample size was determined using a medium effect size (Cohen's $d = 0.5$)³⁵. Power analyses
111 using G*power suggested a sample size of 34 for 0.8 power. This sample size was implemented for all

112 tasks except for the evaluating conditioning task. We performed power analyses using G*power with a
113 priori effect size (Cohen's $d = 0.55$) and 0.8 power to determine the sample size for a one-tailed
114 paired-sample t test, resulting in a sample size of 22 for the evaluating conditioning task. The study was
115 approved by the Research Ethics Committee of the University of Science and Technology of China
116 (NO. 2020-N(H)-099), and written informed consent was obtained from all participants consistent with
117 the Declaration of Helsinki. The methods were carried out in accordance with the approved guidelines.

118 **Robot-related questionnaires and scales**

119 Participants were asked to rate their perceptions toward humanoid robots on a 5-point questionnaire (1-
120 completely disagree, 5-completely agree; Table 1). Items 4, 5 and 6 were reversed items. The sum score
121 of these six items indicates participants' opinion toward humanoid robots. We also obtained
122 participants' explicit attitudes toward robots using the Negative Attitudes Toward Robots Scale (NARS)
123 ³⁶.

124 We determined participants' prior experience with robots by asking about the frequency of daily
125 interactions with robots and which of 20 fictional films (*i.e.*, *Wall-E*, *The Terminator*, *Alita: Battle*
126 *Angel*, *I, Robot*, *AI: Artificial Intelligence*, *Chappie*, *Bicentennial Man*, *RoboCop*, *Short Circuit*, *Ex*
127 *Machina*, *Edward Scissorhands*, *Blade Runner*, *My Robot Girlfriend*, *Autómata*, *Real Steel*, *2001: A*
128 *Space Odyssey*, *Moon*, *Star Wars*, *The Surrogates*, *Forbidden Planet*) portraying robots they have seen.

129 **Original implicit association test.** We assessed participants' automatic attitude using a standard
130 implicit attitude measure paradigm — the IAT task ^{37, 38}. In the original IAT (Fig. 3a), stimuli consisted
131 of attribute words and concept images. The word stimuli were six adjectives, including three
132 adjectives—threatening, strikes and lethal—that indicate “threatening” and three adjectives—neutral,
133 normal and average—that indicate “nonthreatening” ³⁹. The image stimuli were 20 images of fictional
134 humanoid robots and 20 images of Caucasian individuals (Caucasian individuals were used because the
135 fictional humanoid robot forms were Caucasian; we matched the action of each pair of images, and the

136 gender ratio of fictional humanoid robots and humans was the same; the stimuli are available at
137 <https://rec.ustc.edu.cn/share/1f97a050-7745-11ec-bbe8-61ef14e3189d>). To avoid confounding of
138 character stereotypes, stimuli from the movies were excluded. These images were downloaded from the
139 open source on the internet. Adobe Photoshop CS5 was used for all images to uniformly process the
140 pixels and colors of the images so that the materials were consistent and the pixel size was 300×300.
141 Before the formal experiment, participants were presented with an image set including all 20 images of
142 fictional humanoid robots and 20 images of Caucasian individuals in a pilot experiment. For each
143 image, participants provided a rating of affect valence (1 = extremely negative; 5 = extremely positive),
144 arousal (1 = extremely calm; 5 = extremely aroused) and novelty (1 = not novel at all, 9 = extremely
145 novel). The results showed that the affect valence and arousal for fictional humanoid robot images were
146 neutral (valence: mean = 3.05, SD = 0.69; arousal: mean = 2.94, SD = 0.90). The affect valence and
147 arousal for Caucasian images were neutral as well (valence: mean = 3.67, SD = 0.66; arousal: mean =
148 3.20, SD = 1.03). There was a significant difference in novelty between fictional humanoid robot
149 images and Caucasian images (fictional humanoid robot: mean = 5.64, SD = 0.362; Caucasian: mean =
150 3.95, SD = 0.31; $t_{24} = 16.99$, $p < 0.001$).

151 The original IAT involved a series of five discrimination blocks.

152 Block 1 was a concept discrimination task. In this block, the concept labels (humanoid or human)
153 were displayed on the top left and right sides of the screen, and participants were told to categorize the
154 concept images (humanoid robot images or human images) by pressing a corresponding button. Each
155 trial began with the presentation of a fixation cross for 1 s followed by a target image remaining on the
156 screen until the participant responded or for a maximum of 3 s. There were 40 trials in this block, and
157 each concept image was presented 1 time. The order of stimulus presentation was pseudorandomized

158 across participants in all blocks. The label positions (left or right) of the concepts (humanoid or human)
159 were counterbalanced between participants.

160 Block 2 was an attribute discrimination task. The participants were required to categorize the
161 attribute words according to the attribute labels (threat or nonthreat), and the experimental procedure
162 was consistent with the concept discrimination task. There were 18 trials in this block, and each
163 attribute word was repeated 3 times.

164 Block 3 was a compatible combined task (humanoid + threat; human + nonthreat). In this block,
165 humanoid robot images and threatening words shared the same response, and human images and
166 nonthreatening words shared the same response. There were 76 trials in this block; each concept image
167 was presented 1 time, and each attribute word was repeated 6 times.

168 Block 4 was a reversed attribute discrimination task. Participants performed the same task as in
169 block 2, but the positions of the attribute labels (threatening words or nonthreatening words) were
170 reversed. There were 18 trials in this block, and each attribute word was repeated 3 times.

171 Block 5 was an incompatible combined task. Participants performed the same task as in block 3,
172 but the position of the attribute labels (threat or nonthreat) were swapped (humanoid + nonthreat;
173 human + threat). There were 76 trials in this block; each concept image was presented 1 time, and each
174 attribute word was repeated 6 times.

175 Brief instructions were presented on the screen at the beginning of each block until the participant
176 responded.

177 Before analyzing the data from the IAT, the following data reduction procedure was applied^{37,38}:

178 1) eliminate data from subjects for whom more than 10% of trials have a latency shorter than 300 ms; 2)

179 eliminate data from subjects for whom more than 20% of trials are incorrect; 3) compute the mean
180 latencies for only correct trials for blocks 3 and 5; 4) compute the pooled standard deviation for all
181 trials in blocks 3 and 5; 5) replace each incorrect trial's latency in blocks 3 and 5 with the block mean
182 computed in step 3 plus 600 ms; 6) compute the difference between the mean latencies of blocks 3 and
183 5 using paired sample t test; 7) divide the difference computed in step 6 by the pooled standard
184 deviation computed in step 4 to calculate the IAT effect size - the IAT score.

185 **Black IAT.** The procedure of the Black IAT was the same as that of the original IAT except the 20
186 images of Caucasian individuals were replaced with 20 images of Black individuals.

187 **Mixed IAT.** The procedure of the mixed IAT was the same as that of the original IAT except the 20
188 images of Caucasian individuals were replaced with 20 images of Black, Caucasian and Asian
189 individuals.

190 **Available humanoid IAT.** The procedure of the available humanoid IAT was the same as that of the
191 original IAT except the 20 images of fictional humanoid robots were replaced with 20 images of the
192 currently available humanoid robots (*e.g.*, Pepper, Nao, and iCub).

193 **Positive adjective IAT.** The procedure of the positive adjective IAT was the same as that of the
194 original IAT except the threat-related adjectives were replaced with positive adjectives (*i.e.*, salient,
195 arousing, and affective).

196 **The humanoid weapon IAT.** The procedure of the humanoid weapon IAT was the same that of as the
197 original IAT except the 20 images of Caucasian individuals were replaced with 20 images of weapons
198 and the concept label of "human" was replaced with "weapon".

199 **Pet-robot weapon IAT.** The procedure of the pet-robot weapon IAT was the same as that of the

200 humanoid weapon IAT except the 20 images of fictional humanoid robots were replaced with 20
201 images of pet robots.

202 **Procedures in the intervention experiment.** In this experiment, participants performed the original
203 IAT before and after the intervention. In the positive intervention group, the positive evaluative
204 conditioning task was used to intervene with implicit attitudes, whereas the neutral evaluative
205 conditioning task was used in the neutral intervention group. In the positive intervention group,
206 participants completed a backward masking task with fMRI scanning and a forced-choice detection
207 task before intervention and then completed the backward masking task with fMRI scanning again after
208 intervention (Fig. 2c). In the neutral intervention group, participants completed the positive adjective
209 IAT and the available humanoid IAT before intervention. The implicit association test and evaluating
210 conditioning task were presented by E-Prime 2.0, and the backward masking task and forced-choice
211 detection task were presented by MATLAB (R2015a).

212 **Materials in the fMRI experiment.** The materials of the backward masking task and forced-choice
213 detection task consisted of 20 images of fictional humanoid robots and 20 images of Caucasian
214 individuals, the same images used in the original IAT. The materials of the evaluative conditioning task
215 consisted of conditioned stimuli (CSs), unconditioned stimuli (USs), target stimuli and neutral fillers
216 (Table 2). The CSs consisted of 20 images of fictional humanoid robots and 20 images of Caucasian
217 individuals. Approximately half of the USs, target stimuli and neutral fillers were images, and half were
218 words. We selected the US images, target images and neutral filler images from the Chinese Affective
219 Image System (CAPS), including 3 target images, 5 positive images (USs), 5 neutral images (USs), and
220 16 neutral filler images. Adobe Photoshop CS5 was used for all images to uniformly process the pixels
221 and colors of the images so that the materials were consistent and the pixel size was 300×300. The

222 valence, arousal and dominance scores for these 29 images are displayed in Table 3 (on a scale from 1
223 = low valence/arousal/dominance to 9 = high valence/arousal/dominance). The word stimuli—3 target
224 adjectives, 5 positive adjectives (USs) ⁴⁰, 5 neutral adjectives (USs) ⁴¹ and 17 neutral filler nouns
225 ⁴¹—were selected from the pilot study. We translated these words into Chinese characters by using
226 Collins COBUILD Advanced Learner’s English-Chinese Dictionary ⁴². The stimuli for this task are
227 shown in Table 4.

228 **Backward masking task.** Participants completed the backward masking task in the MRI scanner (Fig.
229 2a-b) ⁴³. For the subthreshold presentation, the target image was presented for 17 ms followed by a
230 mask for 183 ms and a fixation cross for 1800 ms. For the suprathreshold presentation, the target image
231 was presented for 200 ms followed by a fixation cross for 1800 ms. These images were arranged in a
232 block design consisting of 10 images (either humanoid robot or human) in a computer-generated
233 pseudorandom order. To avoid any possible effects on subthreshold presentation, the six suprathreshold
234 blocks (three humanoid blocks and three human blocks) were presented after the six subthreshold
235 blocks (three humanoid blocks and three human blocks). Except for the first and last baseline blocks (a
236 fixation cross displayed on the screen) lasting for 10 s, each target block was separated by a 20 s
237 baseline block.

238 **Forced-choice detection task.** To confirm that participants were aware of the stimulus during
239 suprathreshold presentation but not during subthreshold presentation in the backward masking task, a
240 forced-choice detection task was used. The forced-choice detection task consisted of 80 trials; the first
241 half were subthreshold presentations, and the second half were suprathreshold presentations. The
242 stimuli were presented similarly to those in the backward masking task. The difference was that a 2000
243 ms forced-choice phase followed the target stimuli. Participants were informed that the target stimulus

244 could have humanoid or human images and were told to recognize the content of each image. Data
245 from two participants in this task were excluded because of program crashes. We compared the
246 accuracy, response rate, and reaction time between suprathreshold and subthreshold presentations
247 separately using a paired sample t test. We also compared the accuracy for both conditions to chance
248 level (50%) using a one-sample t test.

249 **Positive evaluative conditioning task.** The evaluative conditioning task is a classic paradigm of
250 changing implicit attitudes by pairing CSs and USs⁴⁴. In the positive evaluative conditioning task, the
251 participants' implicit negative attitude toward humanoid robots may be reduced by pairing the
252 humanoid robot images (CSs) with the positive stimuli (USs) (Fig. 4b). As a control condition, the
253 human images (CSs) were paired with neutral stimuli (USs). Participants were unaware of the repeated
254 conditioned stimulus–unconditioned stimulus (CS-US)⁴⁵.

255 This task included 6 blocks of 61 trials each. For each block, the specific stimuli were arranged,
256 and all stimuli were presented in pseudorandom order. All stimuli appeared for 1.5 s each. During the
257 experiment, the participants were instructed to view a stream of images or words and respond as soon
258 as possible whenever a prespecified target image or words appeared. To ensure that the participants
259 carefully viewed the US-CS pairs during the entire experiment, the participants were told that the
260 accuracy rate had to be at least 95% after the end of the experiment or they would have to restart the
261 task. Before the formal evaluative conditioning task, a training evaluative conditioning task was
262 performed.

263 **Neutral evaluative conditioning task.** The procedure of the neutral evaluative conditioning
264 task was the same as that for the positive evaluating conditioning task except that the humanoid robot

265 and human images were both paired with the neutral stimuli.

266 **MRI acquisition.** Gradient echo-planar imaging data were acquired using a 3.0 T GE Discovery
267 MR750 with a circularly polarized head coil at the Information Science Center of the University of
268 Science and Technology of China. A T2*-weighted echo-planar imaging sequence (FOV = 240 mm, TE
269 = 30 ms, TR = 2000 ms, flip angle = 85°, matrix = 64 × 64) with 33 axial slices (no gaps, voxel size:
270 3.75 × 3.75 × 3.7 mm³) covering the whole brain was used to acquire the functional MR images.
271 High-resolution T1-weighted spin-echo imaging data were also acquired for anatomical overlays, and
272 three-dimensional gradient-echo imaging data were acquired for stereotaxic transformations after
273 functional scanning. Before entering the scanner, participants were instructed to keep their heads still
274 during all scans. Participants were placed in a light head restraint within the scanner to limit head
275 movement. Visual stimuli were projected on a screen and viewed through a mirror attached to the head
276 coil. During the backward masking task, 4 functional scan runs occurred, each lasting 4 min. There was
277 an interval of approximately 1 min between every two runs.

278 **fMRI processing.** Functional data were analyzed using the Analysis of Functional NeuroImages
279 (version: AFNI_21.01.01) software. Time series were realigned to the second volume. The realigned
280 images were normalized to the Talairach coordinate. Raw data were corrected for temporal shifts
281 between slices and for motion (TRs with motion over 0.3 mm were censored, and participants who had
282 more than 15% of their TRs censored were removed from further analysis), spatially smoothed with a
283 Gaussian kernel (full width at half maximum = 8 mm), and temporally normalized (for each voxel, the
284 signal of each volume was divided by the temporally averaged signal). High-pass temporal filtering
285 (using a filter width of 128 s) was also applied to the data.

286 To elucidate neural responses that correlated with humanoid robot images and human images
287 under subthreshold and suprathreshold conditions, a general linear model (GLM) was used. Regressors
288 of interest were subthreshold humanoid blocks, subthreshold human blocks, suprathreshold humanoid
289 blocks and suprathreshold human blocks. The regressor of no interest was the fixation block. These
290 regressors were convolved with a hemodynamic response function (HRF) and simultaneously regressed
291 against the blood oxygenation level-dependent (BOLD) signal in each voxel. The regressors were not
292 orthogonalized, and there was no significant collinearity among the regressors. Six regressors for head
293 motion were also included. Individual contrast images were analyzed for each regressor of interest's
294 responses using one-sample t tests to generate statistical maps.

295 **Region of interest (ROI) analysis.** We conducted ROI analyses on the bilateral amygdala. The ROIs
296 for the bilateral amygdala were identified from the AAL atlas ⁴⁶. We determined parameter estimates
297 for each participant from the local average in a mask back-projected from the ROIs. The differential
298 neural responses between the humanoid robot and human conditions in the ROIs were analyzed.
299 Correlation analysis was performed between neural responses and implicit attitudes.

300 **Results**

301 **Negative implicit attitudes toward humanoid robots.** In the attitudes toward robots questions,
302 participants showed a positive explicit attitude toward humanoid robots ($t_{65} = 8.84$, $p < 0.001$, Cohen's
303 $d = 1.10$). In the original IAT, participants had significant IAT scores (baseline = 0.2; mean = 0.46; SD
304 = 0.48; $t_{151} = 6.79$, $p < 0.0001$, Cohen's $d = 0.54$; Fig. 3b). Participants responded faster to
305 combinations of "humanoid + threatening" and "human + nonthreatening" than to combinations of
306 "humanoid + nonthreatening" and "human + threatening" ($t_{151} = -9.64$, $p < 0.0001$, Cohen's $d = 0.79$;
307 Fig. 3b). Previous studies have indicated that participants showed implicit negative attitudes toward
308 Black individuals relative to Caucasian individuals^{30, 47}. We wondered whether the automatic negative
309 associations to humanoid robots compared with Caucasian individuals would still survive when
310 compared with a biased race (*i.e.*, Black individuals) and could be generalized to common races. Thus,
311 we also recruited another group of participants to perform the Black IAT and the mixed IAT. Consistent
312 results were found. Participants responded faster in the compatible task than in the incompatible task in
313 the Black IAT (RT: $t_{21} = -6.40$, $p < 0.001$, Cohen's $d = 1.36$; mean IAT score = 0.61; Fig. 3c) and mixed
314 IAT (RT: $t_{21} = -7.33$, $p < 0.001$, Cohen's $d = 1.57$; mean IAT score = 0.60; Fig. 3c). Our IAT results
315 indicated that, despite the self-reported findings from the questionnaires, participants have negative
316 implicit attitudes toward humanoid robots.

317 In the positive adjective IAT, participants had no significant IAT scores (baseline = 0.2; mean =
318 0.17; SD = 0.37; $t_{29} = -0.45$, $p = 0.65$). Participants did not respond faster to combinations of
319 "humanoid + positive" and "human + neutral" than to combinations of "humanoid + neutral" and
320 "human + positive" ($t_{29} = -1.30$, $p = 0.21$). These results indicated that humanoid robots were not
321 implicitly more associated with positive words. In the available humanoid IAT, participants had

322 significant IAT scores (baseline = 0.2; mean = 0.45; SD = 0.62; $t_{29} = 2.17$, $p = 0.038$). Participants
323 responded faster to combinations of “humanoid + threatening” and “human + nonthreatening” than to
324 combinations of “humanoid + nonthreatening” and “human + threatening” ($t_{29} = -3.62$, $p = 0.001$,
325 Cohen’s $d = 0.66$). These results indicated that participants also have negative implicit attitudes toward
326 currently available humanoid robots. The more fictional films the participants had seen, the less
327 explicitly negative attitudes toward robots they had ($r = -0.439$, $p = 0.015$). However, no significant
328 relationships between participants’ prior experience with robots and implicit attitudes toward robots
329 (neither fictional humanoid robots nor currently available humanoid robots) were observed (all $p >$
330 0.60). Our results indicated that people’s prior experience with robots might influence their explicit
331 perception of robots rather than their implicit perception of robots.

332 Humanoid robots may be considered more competitive with humans than pet robots are, resulting
333 in more negative implicit attitudes toward humanoid robots than toward pet robots. To test the possible
334 influence of the robot’s appearance on implicit attitudes, we focused on the implicit attitude differences
335 between the humanoid and pet robots by using a well-known threatening stimulus (*i.e.*, weapons⁴⁸) as
336 a baseline. Our results indicated that participants displayed larger IAT scores in the humanoid weapon
337 IAT than in the pet-robot weapon IAT ($t_{37} = 3.07$, $p < 0.01$, Cohen’s $d = 0.50$; Fig. 4a). These results
338 indicated that participants have a more negative implicit attitude toward humanoid robots than toward
339 pet robots.

340 **Greater left amygdala activity was induced by humanoid robot images than by human**
341 **images under subthreshold presentation.** As shown in Fig. 4c, the mean response rate in the
342 forced-choice detection task was more than 90% under both presentations and was higher under
343 suprathreshold presentation ($t_{58} = 4.14$, $p < 0.001$, Cohen’s $d = 0.54$) than under subthreshold

344 presentation. The response time under suprathreshold presentation was significantly shorter than that
345 under subthreshold presentation ($t_{58} = -10.91$, $p < 0.001$, Cohen's $d = 1.42$). The accuracy under
346 suprathreshold presentation was significantly higher than that under subthreshold presentation ($t_{58} =$
347 16.68 , $p < 0.001$, Cohen's $d = 2.17$). Importantly, the accuracy under subthreshold presentation did not
348 differ from random chance ($t_{58} = -0.31$, $p = 0.76$), whereas the accuracy under suprathreshold
349 presentation was higher than random chance ($t_{58} = 19.13$, $p < 0.001$, Cohen's $d = 2.49$). These findings
350 indicate that participants are aware of the stimuli under suprathreshold presentation but are unaware of
351 the stimuli under subthreshold presentation.

352 In fMRI analysis, we first tested whether subthreshold presentation of images of humanoid robots
353 leads to a greater amygdala (Fig. 5a) response compared to images of humans. Activation in response
354 to humanoid robot images was significantly stronger than activation in response to human images in
355 the anatomically defined left amygdala under subthreshold presentation (left amygdala: $t_{58} = 2.61$, $p =$
356 0.012 , Cohen's $d = 0.34$; right amygdala: $t_{58} = 0.26$, $p = 0.80$; Fig. 5b); no such difference was observed
357 under suprathreshold presentation (left amygdala: $t_{58} = -0.76$, $p = 0.37$; right amygdala: $t_{58} = 0.26$, $p =$
358 0.80). The activation differences between humanoid robot and human images under subthreshold
359 presentation in the left amygdala were significantly stronger than activation differences under
360 suprathreshold presentation (left amygdala: $t_{58} = 2.23$, $p = 0.03$, Cohen's $d = 0.30$; right amygdala: $t_{58} =$
361 1.32 , $p = 0.19$; Fig. 5b). Importantly, a greater IAT score was associated with greater left amygdala
362 activity under the subthreshold condition ($r = 0.46$, $p < 0.001$; Fig. 5c). Our results indicated that
363 greater left amygdala activity induced by humanoid robot images compared to that induced by human
364 images under subthreshold presentation is associated with negative implicit attitudes toward humanoid
365 robots.

366 We found that there was a significant difference in novelty between humanoid robot images and
367 human images. It has been reported that novelty contributes to amygdala activation⁴⁹. We controlled
368 for the novelty of images by adding this as a covariate in our fMRI general linear model analyses, and
369 this did not alter the results related to amygdala activation (Fig. 6a-b).

370 **The left amygdala response to subthreshold presentation of humanoid robot images changes**

371 **after successfully weakening negative attitude.** We conducted a two-factor (group factor: positive
372 intervention group, neutral intervention group; time factor: pretest, posttest) repeated-measures
373 ANOVA on IAT scores. A significant group×time interaction effect was found ($F_{(1,48)} = 4.99$, $p = 0.038$,
374 $\eta^2 = 0.094$, Fig. 4d). *Post hoc* analysis revealed that the posttest IAT scores were significantly smaller
375 than the pretest scores ($t_{22} = -2.16$, $p = 0.042$, Cohen's $d = 0.45$; Fig. 5d) in the positive intervention
376 group but not in the neutral intervention group ($t_{26} = -0.88$, $p = 0.39$), indicating that participants'
377 negative implicit attitudes toward humanoid robots had been successfully weakened by the positive
378 evaluative conditioning task. We tested whether the left amygdala response to subthreshold
379 presentation of humanoid robot images changed after successfully weakening negative implicit
380 attitudes. We conducted a two-factor (time factor: before modulation, after modulation; presentation
381 factor: subthreshold presentation, suprathreshold presentation) repeated-measures ANOVA on
382 activation differences between humanoid robot and human images in the left amygdala. A significant
383 time×presentation interaction effect was found ($F_{(1,22)} = 8.51$, $p = 0.008$), but no significant time or
384 presentation main effects were found (all $p > 0.29$). *Post hoc* analysis revealed that there was a
385 marginally significant decrease in left amygdala activation between humanoid robot and human images
386 under the subthreshold presentation between the posttest and pretest scans ($t_{22} = -2.02$, $p = 0.056$,
387 Cohen's $d = 0.38$; Fig. 5e). Correlation analysis revealed that there was a significant correlation

388 between IAT score changes and activation value changes in the left amygdala under subthreshold
389 presentation ($r = 0.58$, $p = 0.004$; Fig. 5f). These results demonstrated a causal relationship between
390 implicit attitudes toward humanoid robots and the left amygdala response to subthreshold presentation
391 of humanoid robot images. Similar to the previous section, after controlling for the novelty of images
392 by adding it as a covariate in our fMRI general linear model analyses, the results did not change (Fig.
393 6c-d).

394

395 **Discussion**

396 Despite self-reported positive attitudes in the questionnaires, participants had negative implicit
397 attitudes toward humanoid robots. The left amygdala response to subthreshold presentation of
398 humanoid robots, which was positively associated with implicit attitudes toward humanoid robots,
399 indicates the automatic and quick detection of humanoid robots. After successfully weakening negative
400 attitudes, the decrease in IAT scores was positively associated with the decrease in activation in the left
401 amygdala under subthreshold presentation, demonstrating a causal relationship between implicit
402 attitudes toward humanoid robots and the left amygdala response to subthreshold presentation of
403 humanoid robot images.

404 Our results provide evidence that humanoid robots may be implicitly perceived as a threat. People
405 detect and respond rapidly to threatening stimuli ^{50,51}. Nonhuman primates also respond more rapidly
406 to threatening stimuli (at least some types) than to neutral stimuli ^{52,53}. The amygdala, a subcortical
407 structure in the anterior temporal lobe, is located in an evolutionarily old part of the brain and is shared
408 by other mammals. It is assumed to be the neural basis of the hardwired “fear module” that allows us to
409 automatically and quickly detect threatening stimuli ²⁷. Studies have documented that the amygdala
410 responds selectively to threats, sometimes irrespective of the affective valence while using implicit
411 measures, such as animate entities ⁵⁴⁻⁵⁶ and depictions of humans ⁵⁷⁻⁵⁹. Our results showing no
412 amygdala activity in response to suprathreshold presentation of humanoid robot images seemingly
413 support that people did not perceive humanoid robots as a threat. However, when automatic and
414 controlled evaluations of threat differ, more positive controlled processing can moderate more negative
415 automatic processing ³⁰. With the positive explicit attitudes toward humanoid robots found in the
416 present study, the absence of amygdala response to suprathreshold presentation of humanoid robot

417 images is understandable. Interestingly, although controlled processing can eliminate amygdala activity
418 caused by consciously presented threats, the amygdala still shows greater responses to threats that are
419 presented unconsciously³⁰⁻³². Several studies suggest that the left amygdala might be specifically
420 involved in the processing of facial stimuli^{60, 61}. It has also been reported that the left amygdala shows
421 less habituation to fearful stimuli than the right amygdala, which might make it more likely to capture
422 the blood oxygen level-dependent changes in this area^{62, 63}. However, the lateralization of amygdala
423 activation is still controversial^{64, 65}. Our present study found that greater left amygdala activity was
424 induced by humanoid robot images than by human images under subthreshold presentation. After
425 successfully weakening negative attitudes, the decrease in IAT scores was positively associated with
426 the decrease in activation in the left amygdala under subthreshold presentation. Our results indicate a
427 causal relationship between implicit attitudes toward humanoid robots and the left amygdala response
428 to subthreshold presentation of humanoid robot images. Note that this result did not change after
429 controlling for the novelty of the humanoid robot and human images. These results potentially reflect
430 the automaticity with rapid recruitment of the amygdala for humanoid robot-related stimuli processing.

431 Except for threatening stimuli, many studies have indicated that the amygdala plays an important
432 role in emotional processing, especially negative emotions, including disgust, sadness and pain^{66, 67}.
433 The amygdala is also activated by stimuli involving social information^{68, 69}. Thus, some might argue
434 that the amygdala response to subthreshold presentation of humanoid robots in this study could be
435 explained by factors other than threat or fear. In this study, we used humanoid robots with neutral faces
436 and humans with neutral faces as control stimuli, which might eliminate the confounding factors of
437 social information, such as facial expressions and emotions. Although there was a significant difference
438 in novelty between humanoid robot and human images, the amygdala-related results were not altered

439 after controlling for novelty. Furthermore, salience is also a possibility, and consequently, humanoid
440 robots should be implicitly more associated with other salient adjectives (e.g., positive words).
441 However, our results of the positive adjective IAT were not significant, indicating that humanoid robots
442 were not implicitly more associated with salient adjectives. Of note, the significantly positive
443 correlation between the decrease in the left amygdala response and the decrease in negative attitudes
444 indicates that there is a modulation effect of the positive evaluating conditioning task on amygdala
445 activity, although the reduced saliency of stimuli might also contribute to the decrease in amygdala
446 activity.

447 People perceive robots based upon context, cues, and cultural assumptions³³. Our ideas about
448 robots are mainly informed by science fiction media¹⁰, and robots in science fiction are generally
449 described as a super species with greater intelligence than that of humans, attempting to eliminate
450 humanity¹⁰. The concept that fully autonomous robots are dangerous competitive “living” entities that
451 will threaten the survival of humanity gradually becomes a deep impression. Two species with the
452 closest living habits have the strongest competitiveness between them^{70, 71}, such as Neanderthal
453 extinction by competition with anatomically modern humans⁷²⁻⁷⁴. Combined with our results, robots
454 with human-like faces and intelligence are plausibly implicitly perceived as a new species or even race.
455 Consequently, the fear of humanoid robots is likely similar to threat stimuli related to survival in
456 evolutionary history. However, whether people perceive humanoid robots as an evolutionary threat is
457 equivocal and is worth future research to uncover.

458 A growing body of evidence has indicated that the physical appearance of a robot has a strong
459 impact on people’s perception^{9, 75, 76}. However, our results suggest that a robot’s appearance is likely to
460 influence explicit perception rather than implicit perception (at least the implicit attitudes toward

461 robots). Fictional humanoid robots with highly anthropomorphic appearances and the currently
462 available humanoid robots with less anthropomorphic appearances were used in our experiment. Of
463 note, participants showed negative implicit attitudes toward both kinds of humanoid robots. Consistent
464 with our results, previous studies have shown negative implicit attitudes toward robots by using robot
465 silhouettes^{77, 78}. Taken together, the physical appearance of a robot plausibly has a weak impact on
466 people's implicit perception.

467 Our demonstration that people implicitly perceive humanoid robots as a threat might contribute to
468 some negative biases for robots. Consistent with previous studies^{77, 78}, we found that people have
469 negative implicit attitudes toward robots. A study suggests that the early top-down process of empathy
470 is weaker for humanoid robots than for humans⁷⁹. A systematic review⁸⁰ of anxiety and acceptance
471 toward social robots reveals that people only slightly accept social robots and feel slightly anxious
472 about them in general. It might be that the negative perception toward robots immediately and
473 spontaneously activates negative attitudes and social stereotypes toward robots and consequently
474 causes biased behaviors. The negative perception toward robots would be detrimental to the
475 development of robots and successful social encounters with robots. Of note, using an evaluative
476 conditioning task seems to weaken this negative perception toward humanoid robots. More research
477 should address the issue of where negative perception toward humanoid robots comes from and how to
478 modulate this negative perception.

479 In previous studies, a particular robot-related experience might have a significant influence on
480 people's explicit perception of robots^{36, 81, 82}. Consistent with these findings, our study found that the
481 more fictional films portraying robots the participants had seen, the less explicitly negative attitudes
482 toward robots they had. However, no significant relationships between participants' prior experience

483 with robots and implicit attitudes toward robots were observed. Our results indicated that people's prior
484 experience with robots might influence their explicit perception of robots rather than their implicit
485 perception of robots. Thus, future research addressing the question of the impact of a priori experience
486 on perception of robots might distinguish between explicit and implicit perception.

487 One limitation was that the humanoid robot stimuli were only in the form of images. Although
488 fictional and currently available humanoid robot images were used in the present study, the ecological
489 effect of researching human-robot interactions is somewhat weak. It is better to investigate how people
490 perceive and interact with robots in a socially dynamic environment. Thus, stimuli in the form of
491 videos and human-robot interaction research in the real world with the help of mobile neuroimaging
492 should be considered in future studies.

493 **Conclusions**

494 In summary, this study demonstrates that humanoid robots are implicitly perceived as a threat.
495 This sheds light on how people perceive and interact with robots that are increasingly entering our
496 social environment. The future of social robots is undeniably exciting, and insights from
497 neuropsychology research will guide the future direction of robot development and bring us closer to
498 interacting with social robots.

499

500 **Acknowledgments**

501 This work was supported by grants from the National Key Basic Research Program
502 (2018YFC0831101), the National Natural Science Foundation of China (71942003, 31771221,
503 61773360, 32100886, and 71874170), the Major Project of Philosophy and Social Science Research,
504 Ministry of Education of China (19JZD010), the CAS-VPST Silk Road Science Fund 2021
505 (GLHZ202128), the Collaborative Innovation Program of Hefei Science Center, CAS
506 (2020HSC-CIP001), and the China Postdoctoral Science Foundation (2016M592051). A portion of the
507 numerical calculations in this study were performed with the supercomputing system at the
508 Supercomputing Centre of USTC.

509 **Author contributions**

510 ZDW, YC, and XCZ conceived and designed the study. ZDW and YC obtained the findings. YC
511 was responsible for the acquisition of data. ZDW and YC analyzed and interpreted the data. PYZ, YP,
512 JCR, QZ, RJZ, BSQ, YCB, SHH, CBL and DRZ provided administrative, technical, or material
513 support. CBL and XCZ supervised the study. ZDW and YC drafted the paper, and all authors
514 contributed to critical revision for intellectual content.

515 **Conflicts of interest**

516 The authors declare no conflicts of interest.

517 **Data and materials availability statement**

518 The data and scripts are available at
519 <https://rec.ustc.edu.cn/share/1f97a050-7745-11ec-bbe8-61ef14e3189d>.

520 **References**

- 521 1. Yang, G.Z., *et al.* The grand challenges of Science Robotics. *Science robotics* **3** (2018).
- 522 2. Bossi, F., *et al.* The human brain reveals resting state activity patterns that are predictive of biases
523 in attitudes toward robots. *Science robotics* **5** (2020).
- 524 3. Scassellati, B., Admoni, H. & Mataric, M. Robots for Use in Autism Research. *Annu Rev Biomed*
525 *Eng* **14**, 275-294 (2012).
- 526 4. Broadbent, E. Interactions With Robots: The Truths We Reveal About Ourselves. *Annual review of*
527 *psychology* **68**, 627-652 (2017).
- 528 5. Fazio, R.H., Jackson, J.R., Dunton, B.C. & Williams, C.J. Variability in automatic activation as an
529 unobtrusive measure of racial attitudes: a bona fide pipeline? *J Pers Soc Psychol* **69**, 1013-1027 (1995).
- 530 6. Hills, P.J., Lewis, M.B. & Honey, R.C. Stereotype priming in face recognition: Interactions between
531 semantic and visual information in face encoding. *Cognition* **108**, 185-200 (2008).
- 532 7. Ferguson, M.J. & Bargh, J.A. How social perception can automatically influence behavior. *Trends*
533 *in cognitive sciences* **8**, 33-39 (2004).
- 534 8. Henschel, A., Hortensius, R. & Cross, E.S. Social Cognition in the Age of Human-Robot Interaction.
535 *Trends in neurosciences* **43**, 373-384 (2020).
- 536 9. Cross, E.S. & Ramsey, R. Mind Meets Machine: Towards a Cognitive Science of Human-Machine
537 Interactions. *Trends in cognitive sciences* **25**, 200-212 (2021).
- 538 10. Cave, S. & Dihal, K. Hopes and fears for intelligent machines in fiction and reality. *Nature Machine*
539 *Intelligence* **1**, 74-78 (2019).
- 540 11. Broadbent, E., *et al.* Attitudes and reactions to a healthcare robot. *Telemedicine journal and*
541 *e-health : the official journal of the American Telemedicine Association* **16**, 608-613 (2010).

- 542 12. Preparing for the Future of Artificial Intelligence (Executive Ofce of the President National Science
543 and Technology Council.
- 544 13. Timo, G. & Markus, A. Are robots becoming unpopular? Changes in attitudes towards
545 autonomous robotic systems in Europe. *Computers in human behavior* **93**, 53-61 (2019).
- 546 14. Baobao, Z. & Dafoe, A. Artificial Intelligence: American Attitudes and Trends. *Oxford, UK: Center
547 for the Governance of AI, Future of Humanity Institute, University of Oxford* (2019).
- 548 15. Liang, Y.H. & Lee, S.A. Fear of Autonomous Robots and Artificial Intelligence: Evidence from
549 National Representative Data with Probability Sampling. *International Journal Of Social Robotics* **9**,
550 379-384 (2017).
- 551 16. Fast, E. & Horvitz, E. Long-Term Trends in the Public Perception of Artificial Intelligence.
552 *Thirty-First Aaai Conference on Artificial Intelligence*, 963-969 (2017).
- 553 17. Smith & Anderson. Automation in everyday life. *Pew Research Center, Washington, DC Retrieved
554 from <http://www.pewinternet.org/2017/10/04/automation-in-everyday-life/> Google Scholar
555 (2017).*
- 556 18. Social, T.O. Public attitudes toward robots. *Special Eurobarometer* **382** (2012).
- 557 19. Rosenman, R., Tennekoon, V. & Hill, L.G. Measuring bias in self-reported data. *International
558 journal of behavioural & healthcare research* **2**, 320-332 (2011).
- 559 20. Petty, R.E., Tormala, Z.L., Brinol, P. & Jarvis, W.B. Implicit ambivalence from attitude change: an
560 exploration of the PAST model. *J Pers Soc Psychol* **90**, 21-41 (2006).
- 561 21. Gawronski, B. & Bodenhausen, G.V. Associative and propositional processes in evaluation: an
562 integrative review of implicit and explicit attitude change. *Psychological bulletin* **132**, 692-731 (2006).
- 563 22. Hannula, D.E., Simons, D.J. & Cohen, N.J. Imaging implicit perception: promise and pitfalls.

- 564 *Nature reviews. Neuroscience* **6**, 247-255 (2005).
- 565 23. Bar-Anan, Y. & Nosek, B.A. A comparative investigation of seven indirect attitude measures.
- 566 *Behavior research methods* **46**, 668-688 (2014).
- 567 24. Schnabel, K., Asendorpf, J.B. & Greenwald, A.G. Assessment of Individual Differences in Implicit
- 568 Cognition A Review of IAT Measures. *Eur J Psychol Assess* **24**, 210-217 (2008).
- 569 25. McConnell, A.R. & Leibold, J.M. Relations among the implicit association test, discriminatory
- 570 behavior, and explicit measures of racial attitudes. *Journal of experimental social psychology* **37**,
- 571 435-442 (2001).
- 572 26. Mineka, S. & Ohman, A. Phobias and preparedness: The selective, automatic, and encapsulated
- 573 nature of fear. *Biological psychiatry* **52**, 927-937 (2002).
- 574 27. Ohman, A. & Mineka, S. Fears, phobias, and preparedness: Toward an evolved module of fear
- 575 and fear learning. *Psychological review* **108**, 483-522 (2001).
- 576 28. Zald, D.H. The human amygdala and the emotional evaluation of sensory stimuli. *Brain research.*
- 577 *Brain research reviews* **41**, 88-123 (2003).
- 578 29. Tamietto, M. & de Gelder, B. Neural bases of the non-conscious perception of emotional signals.
- 579 *Nature Reviews Neuroscience* **11**, 697-709 (2010).
- 580 30. Cunningham, W.A., *et al.* Separable neural components in the processing of black and white faces.
- 581 *Psychological science* **15**, 806-813 (2004).
- 582 31. Williams, L.M., *et al.* Amygdala-prefrontal dissociation of subliminal and supraliminal fear. *Human*
- 583 *brain mapping* **27**, 652-661 (2006).
- 584 32. Felmingham, K., *et al.* Dissociative responses to conscious and non-conscious fear impact
- 585 underlying brain function in post-traumatic stress disorder. *Psychological medicine* **38**, 1771-1780

586 (2008).

587 33. Bigman, Y.E., Waytz, A., Alterovitz, R. & Gray, K. Holding Robots Responsible: The Elements of
588 Machine Morality. *Trends in cognitive sciences* **23**, 365-368 (2019).

589 34. Waytz, A., Heafner, J. & Epley, N. The mind in the machine: Anthropomorphism increases trust in
590 an autonomous vehicle. *Journal of experimental social psychology* **52**, 113-117 (2014).

591 35. Cohen, J. *Statistical power analysis for the behavioral sciences* (Routledge, 2013).

592 36. Riek, L.D., Adams, A., & Robinson, P. Exposure to Cinematic Depictions of Robots and Attitudes
593 Towards Them. (2011).

594 37. Greenwald, A., McGhee, D. & Schwartz, J. Measuring individual differences in implicit cognition:
595 the implicit association test[J]. *JSPS* **74**, 1464-1480 (1998).

596 38. Kumaran, D., Banino, A., Blundell, C., Hassabis, D. & Dayan, P. Computations underlying social
597 hierarchy learning: distinct neural mechanisms for updating and representing self-relevant
598 information. *Neuron* **92**, 1135-1147 (2016).

599 39. Huijding, J. & de Jong, P. Beyond fear and disgust: The role of (automatic) contamination-related
600 associations in spider phobia[J]. *Journal of Behavior Therapy and Experimental Psychiatry* **38**, 200-211
601 (2007).

602 40. Olson, M. & Fazio, R. Reducing automatically activated racial prejudice through implicit
603 evaluative conditioning[J]. *Personality and Social Psychology Bulletin* **32**, 421-433 (2006).

604 41. Scott, G., O'Donnell, P., Leuthold, H. & Sereno, S. Early emotion word processing: Evidence from
605 event-related potentials[J]. *Biol Psychol* **80**, 95-104 (2009).

606 42. Ke, K.e. Collins COBUILD Advanced Learner's English-Chinese Dictionary[M]. in *Beijing: Foreign*
607 *Language Teaching and Research Press & Harper Collins Publishers Ltd* (2011).

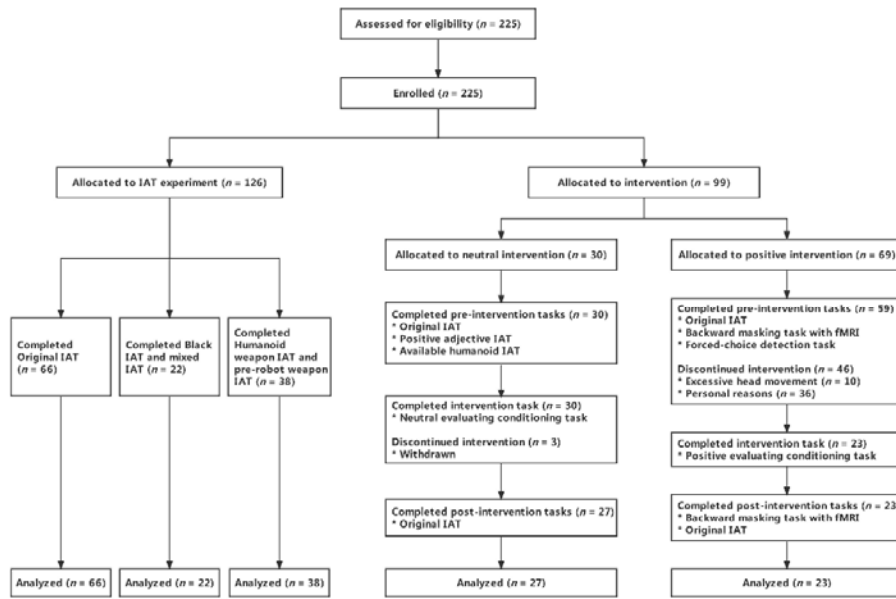
- 608 43. Zhang, X., *et al.* Masked smoking-related images modulate brain activity in smokers. *Hum. Brain*
609 *Mapp.* **30**, 896-907 (2009).
- 610 44. McNulty, J.K., Olson, M.A., Jones, R.E. & Acosta, L.M. Automatic associations between one's
611 partner and one's affect as the proximal mechanism of change in relationship satisfaction: Evidence
612 from evaluative conditioning. *Psychol. Sci.* **28**, 1031-1040 (2017).
- 613 45. Olson, M.A. & Fazio, R.H. Reducing automatically activated racial prejudice through implicit
614 evaluative conditioning. *Personality and Social Psychology Bulletin* **32**, 421-433 (2006).
- 615 46. Rolls, E.T., Huang, C.C., Lin, C.P., Feng, J. & Joliot, M. Automated anatomical labelling atlas 3.
616 *NeuroImage* **206**, 116189 (2020).
- 617 47. Mori, K. Assessment of the Implicit Attitude of Japanese People toward Blacks and Little Black
618 Sambo. *Open Journal of Social Sciences* **06. 1-13. 10.4236** (2018).
- 619 48. Anderson, C.A., Benjamin, A. J., & Bartholow, B. D. Does the gun pull the trigger? Automatic
620 priming effects of weapon pictures and weapon names. *Psychological science* **9**, 308-314 (1998).
- 621 49. Blackford, J.U., Buckholz, J.W., Avery, S.N. & Zald, D.H. A unique role for the human amygdala in
622 novelty detection. *NeuroImage* **50**, 1188-1193 (2010).
- 623 50. Gomes, N., Silva, S., Silva, C.F. & Soares, S.C. Beware the serpent: the advantage of
624 ecologically-relevant stimuli in accessing visual awareness. *Evol Hum Behav* **38**, 227-234 (2017).
- 625 51. Soares, S.C. The Lurking Snake in the Grass: Interference of Snake Stimuli in Visually Taxing
626 Conditions. *Evol Psychol-US* **10**, 187-197 (2012).
- 627 52. Kawai, N. & Koda, H. Japanese monkeys (*Macaca fuscata*) quickly detect snakes but not spiders:
628 Evolutionary origins of fear-relevant animals. *Journal of comparative psychology* **130**, 299-303 (2016).
- 629 53. Shibasaki, M. & Kawai, N. Rapid detection of snakes by Japanese monkeys (*Macaca fuscata*): an

- 630 evolutionarily predisposed visual system. *Journal of comparative psychology* **123**, 131-135 (2009).
- 631 54. Mormann, F., *et al.* A category-specific response to animals in the right human amygdala. *Nature*
632 *neuroscience* **14**, 1247-1249 (2011).
- 633 55. Rutishauser, U., *et al.* Single-Unit Responses Selective for Whole Faces in the Human Amygdala.
634 *Current Biology* **21**, 1654-1660 (2011).
- 635 56. Yang, J.J., Bellgowan, P.S.F. & Martin, A. Threat, domain-specificity and the human amygdala.
636 *Neuropsychologia* **50**, 2566-2572 (2012).
- 637 57. Breiter, H.C., *et al.* Response and habituation of the human amygdala during visual processing of
638 facial expression. *Neuron* **17**, 875-887 (1996).
- 639 58. Wheatley, T., Milleville, S.C. & Martin, A. Understanding animate agents - Distinct roles for the
640 social network and mirror system. *Psychological science* **18**, 469-474 (2007).
- 641 59. Bonda, E., Petrides, M., Ostry, D. & Evans, A. Specific involvement of human parietal systems and
642 the amygdala in the perception of biological motion. *Journal Of Neuroscience* **16**, 3737-3744 (1996).
- 643 60. Morris, J.S., *et al.* A differential neural response in the human amygdala to fearful and happy
644 facial expressions. *Nature* **383**, 812-815 (1996).
- 645 61. Blair, R.J.R., Morris, J.S., Frith, C.D., Perrett, D.I. & Dolan, R.J. Dissociable neural responses to
646 facial expressions of sadness and anger. *Brain : a journal of neurology* **122**, 883-893 (1999).
- 647 62. Wright, C.I., *et al.* Differential prefrontal cortex and amygdala habituation to repeatedly
648 presented emotional stimuli. *Neuroreport* **12**, 379-383 (2001).
- 649 63. Phillips, M., *et al.* Time courses of left and right amygdalar responses to fearful facial expressions.
650 *NeuroImage* **13**, S458-S458 (2001).
- 651 64. Baas, D., Aleman, A. & Kahn, R.S. Lateralization of amygdala activation: a systematic review of

- 652 functional neuroimaging studies. *Brain Res Rev* **45**, 96-103 (2004).
- 653 65. Wager, T.D., Phan, K.L., Liberzon, I. & Taylor, S.F. Valence, gender, and lateralization of functional
654 brain anatomy in emotion: a meta-analysis of findings from neuroimaging. *NeuroImage* **19**, 513-531
655 (2003).
- 656 66. Tamietto, M. & de Gelder, B. Neural bases of the non-conscious perception of emotional signals.
657 *Nature reviews. Neuroscience* **11**, 697-709 (2010).
- 658 67. Fang, Z.Y., Li, H., Chen, G. & Yang, J.J. Unconscious Processing of Negative Animals and Objects:
659 Role of the Amygdala Revealed by fMRI. *Frontiers in human neuroscience* **10** (2016).
- 660 68. Adolphs, R. What does the amygdala contribute to social cognition? *Year In Cognitive*
661 *Neuroscience 2010* **1191**, 42-61 (2010).
- 662 69. Frith, C.D. & Frith, U. Mechanisms of Social Cognition. *Annual Review Of Psychology, Vol 63* **63**,
663 287-313 (2012).
- 664 70. Hardin, G. The competitive exclusion principle. *Science* **131**, 1292-1297 (1960).
- 665 71. Wangersky, P.J. Lotka-Volterra Population Models. *Annu Rev Ecol Syst* **9**, 189-218 (1978).
- 666 72. Klein, R.G. Neanderthals and modern humans - An ecological and evolutionary perspective.
667 *Science* **305**, 45-45 (2004).
- 668 73. Mellars, P. Neanderthals and the modern human colonization of Europe. *Nature* **432**, 461-465
669 (2004).
- 670 74. Zilhao, J. Neandertals and moderns mixed, and it matters. *Evol Anthropol* **15**, 183-195 (2006).
- 671 75. Baraka, K.A.-O., Patrícia & Ribeiro, Tiago. An Extended Framework for Characterizing Social
672 Robots. (2020).
- 673 76. Bartneck, C., Croft, E., & Kulic, D. Measurement Instruments for the Anthropomorphism, Animacy,

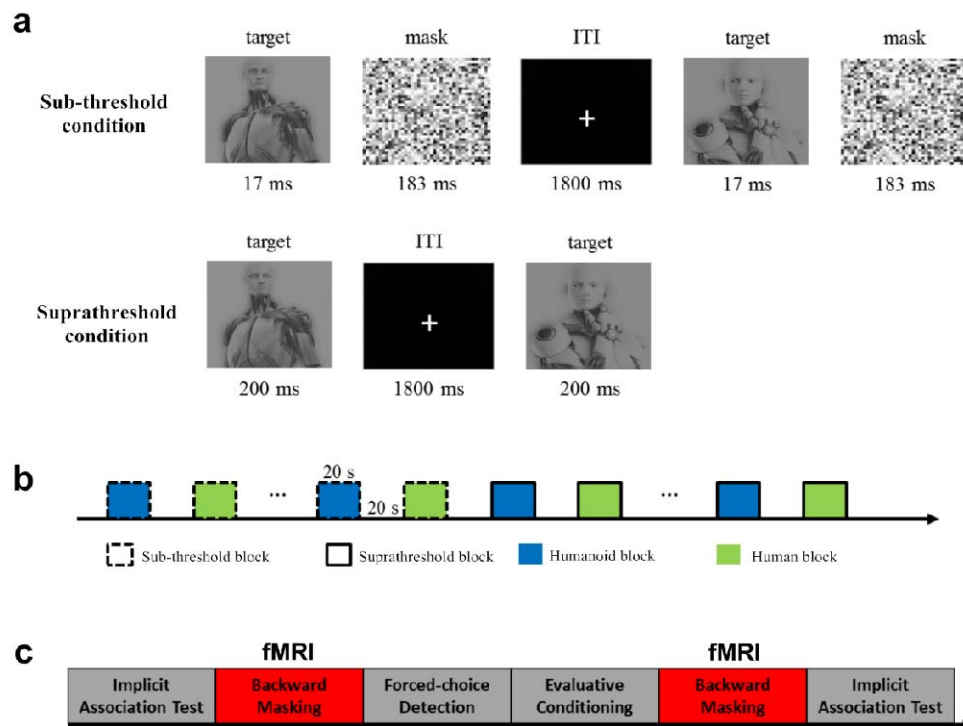
- 674 Likeability, Perceived Intelligence, and Perceived Safety of Robots. *International Journal of Social*
675 *Robotics* **1(1)**, 71-81 (2009).
- 676 77. MacDorman, K.F., Vasudevan, S.K. & Ho, C.-C. Does Japan really have robot mania? Comparing
677 attitudes by implicit and explicit measures. *Ai & Society* **23**, 485-510 (2008).
- 678 78. Graaf, M.M.A.d., Allouch, S.B. & Lutfi, S. What are People's Associations of Domestic Robots?:
679 Comparing Implicit and Explicit Measures. *25th IEEE International Symposium on Robot and Human*
680 *Interactive Communication (RO-MAN) August 26-31, 2016. Columbia University, NY, USA*, 1077-1083
681 (2016).
- 682 79. Suzuki, Y., Galli, L., Ikeda, A., Itakura, S. & Kitazaki, M. Measuring empathy for human and robot
683 hand pain using electroencephalography. *Scientific reports* **5** (2015).
- 684 80. Naneva, S., Gou, M.S., Webb, T.L. & Prescott, T.J. A Systematic Review of Attitudes, Anxiety,
685 Acceptance, and Trust Towards Social Robots. *International Journal Of Social Robotics* **12**, 1179-1201
686 (2020).
- 687 81. Nomura, T., Suzuki, T., Kanda, T. & Kato, K. Measurement of negative attitudes toward robots.
688 *Interact Stud* **7**, 437-454 (2006).
- 689 82. Bartneck, C., Kulic, D., Croft, E. & Zoghbi, S. Measurement Instruments for the
690 Anthropomorphism, Animacy, Likeability, Perceived Intelligence, and Perceived Safety of Robots.
691 *International Journal Of Social Robotics* **1**, 71-81 (2009).
- 692

693 **Figure Legends**



694

695 **Figure 1. Overview diagram of the study flow.**



696

697 **Figure 2. Description of the backward masking task and the procedure for the fMRI experiment.**

698 (a) Time setting of the backward masking task. In the unconscious condition, the target image was

699 presented for 17 ms followed by a mask for 183 ms and a fixation cross for 1800 ms. In the conscious

700 condition, the target image was presented for 200 ms followed by a fixation cross for 1800 ms. (b)

701 Block design of the backward masking task. There were six unconscious blocks (three humanoid

702 blocks and three human blocks) followed by six conscious blocks (three humanoid blocks and three

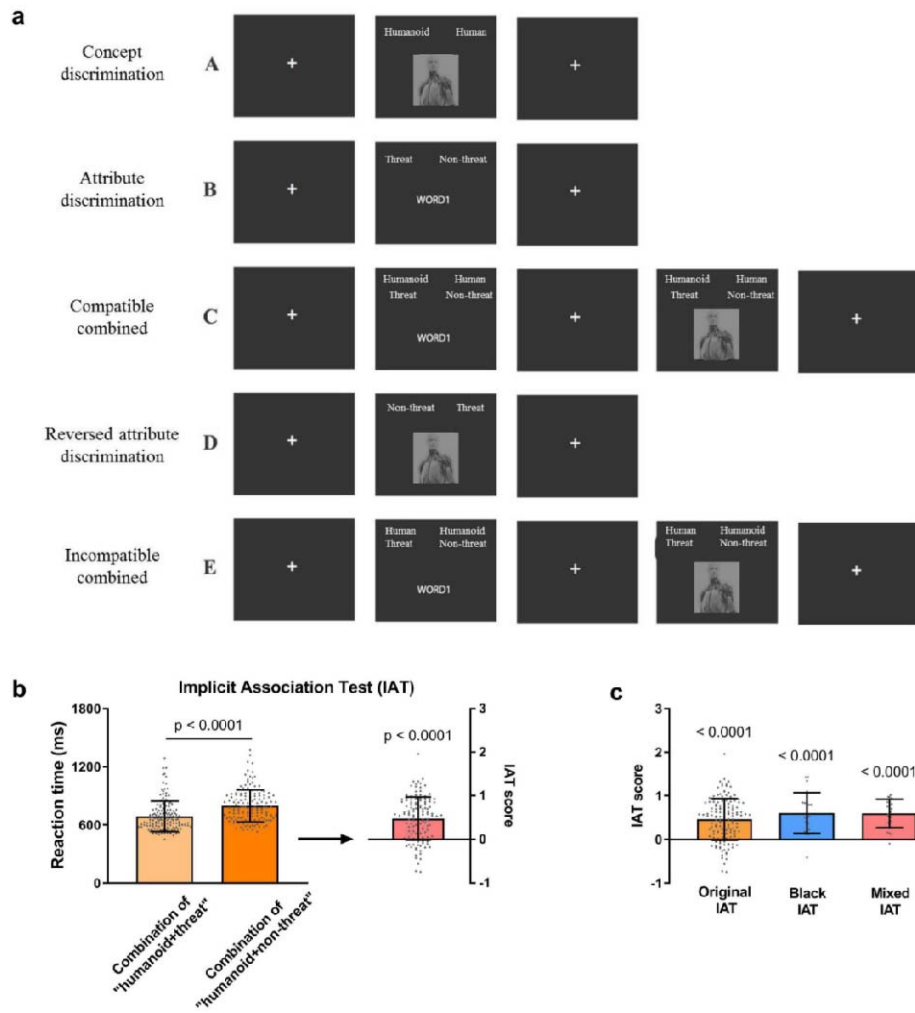
703 human blocks). (c) Procedure for the fMRI experiment. Participants performed the original implicit

704 association test outside the scanner and then completed a backward masking task with fMRI scanning

705 followed by a forced-choice detection task and an evaluating conditioning task in the scanner; then,

706 they performed the backward masking task during a second fMRI scan. Following the end of fMRI

707 scanning, participants completed the original implicit association test outside the scanner again.



708

709 **Figure 3. Implicit association test and its results.** (a) Procedure of the original implicit association

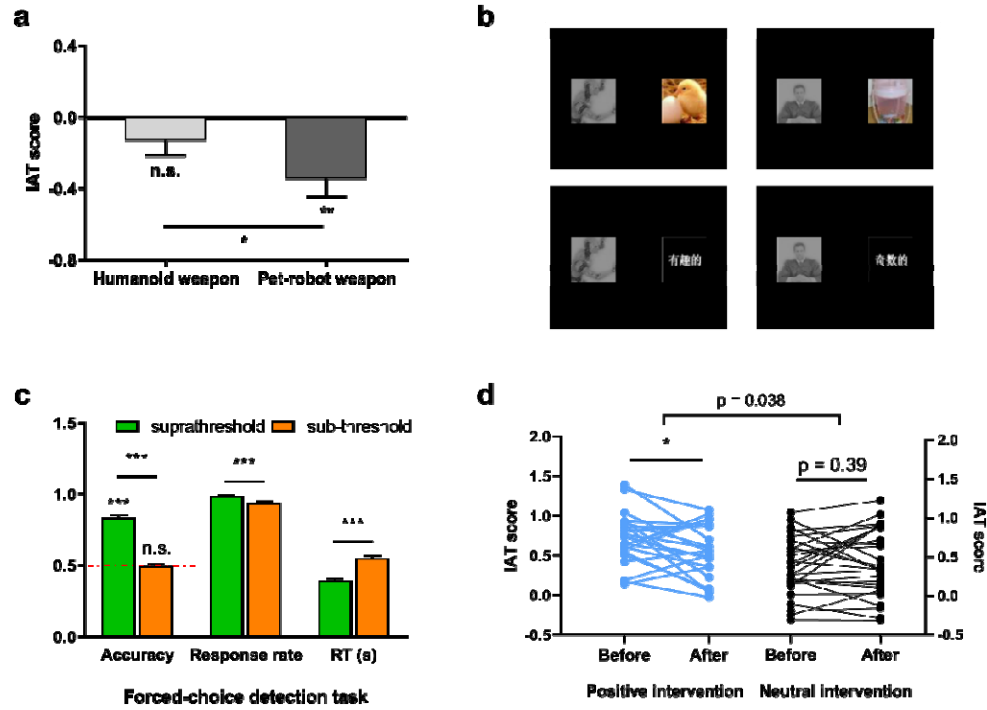
710 test. (b) Negative implicit attitudes toward humanoid robots. In the original implicit association test

711 (IAT), participants responded faster to the combination of “humanoid + threat” than to the combination

712 of “humanoid + nonthreat”, and the computed IAT scores (effect size) were significant. (c) Participants

713 displayed larger IAT scores in the original, Black, and mixed IATs. Plotted data represent the mean \pm

714 SD across participants.



715

716 **Figure 4. Evaluative conditioning task and its results.** (a) More negative implicit attitudes toward

717 humanoid robots compared to animal robots. Participants displayed larger IAT scores in the humanoid

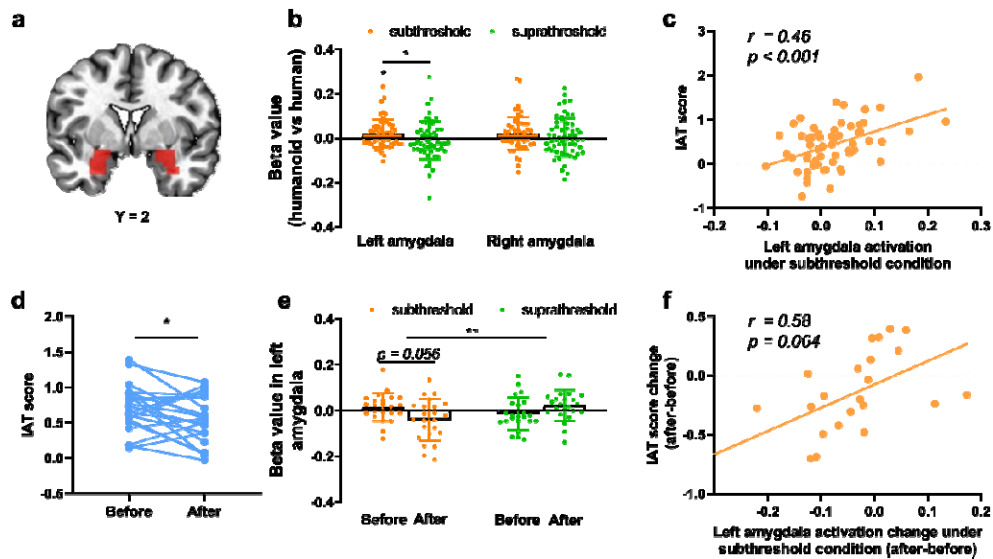
718 weapon IAT than in the pet-robot weapon IAT ($t_{37} = 3.07$, $p < 0.01$, Cohen's $d = 0.50$). (b) Procedure of

719 evaluative conditioning task. (c) Results of the forced-choice detection task. (d) A significant

720 group×time interaction effect was found, indicating that the weakened negative implicit attitudes

721 toward humanoid robots were not due to the practice effect. * $p < 0.05$, ** $p < 0.01$, *** $p < 0.001$.

722 Plotted data represent the mean \pm s.e.m. across participants. IAT = implicit association test.



723

724 **Figure 5. Humanoid robot image-related amygdala activity and amygdala activity changes.** (a)

725 The region of interest for the bilateral amygdala. (b) Although no amygdala activity differences were

726 detected for consciously presented humanoid robot vs. human images, greater left amygdala activity

727 was induced by humanoid robot images than by images of humans under unconscious presentation. (c)

728 A greater IAT score was associated with greater left amygdala activity under unconscious conditions. (d)

729 Significantly smaller IAT scores were found after modulation than those before modulation, indicating

730 that participants' negative implicit attitude toward humanoid robots was successfully weakened. (e)

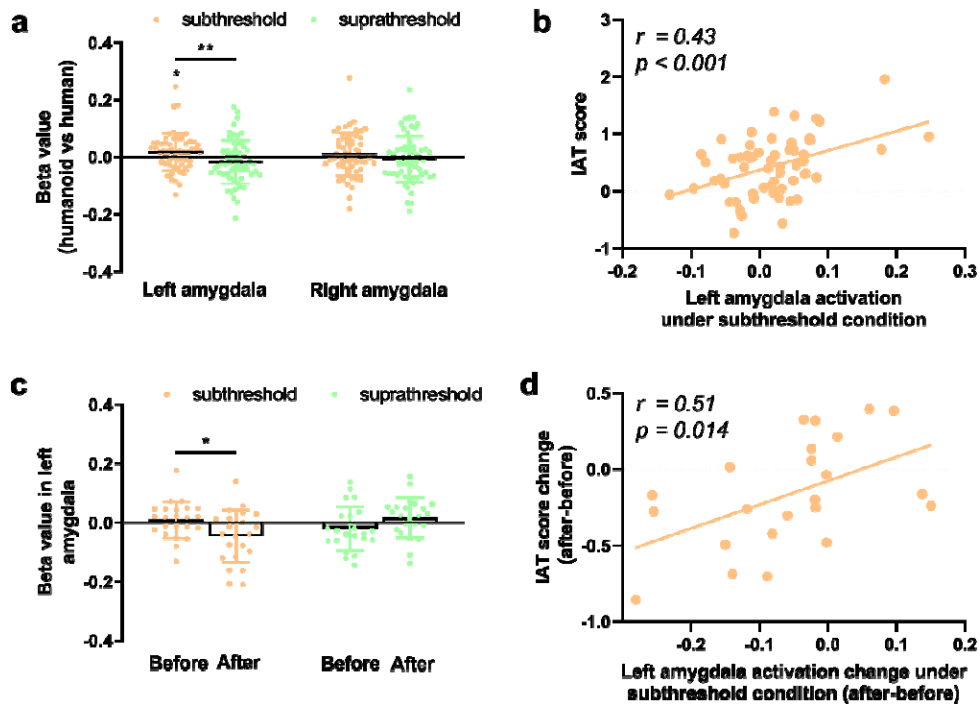
731 The left amygdala activity differences of humanoid robot vs. human images did change under

732 unconscious presentation after successfully weakening the negative implicit attitudes toward humanoid

733 robots. (f) There was a significant correlation between IAT score changes and activation changes in the

734 left amygdala under unconscious presentation. * $p < 0.05$, ** $p < 0.01$, *** $p < 0.001$. For b and e,

735 plotted data represent the mean \pm SD across participants. IAT = implicit association test.



736

737 **Figure 6. Results related to the left amygdala after controlling for novelty.** (a) Greater left

738 amygdala activity was induced by humanoid robot images than by images of humans under

739 subthreshold presentation after controlling for the novelty of the humanoid robot and human images (t_{58}

740 = 2.10, $p = 0.04$, Cohen's $d = 0.27$). (b) A greater IAT score was associated with greater left amygdala

741 activity under the subthreshold condition after controlling for the novelty of the humanoid robot and

742 human images ($r = 0.43$, $p < 0.001$). (c) Controlling for the novelty of the humanoid robot and human

743 images, the left amygdala activity differences of humanoid robot vs. human images did change under

744 subthreshold presentations after successfully weakening the negative implicit attitudes toward

745 humanoid robots ($t_{22} = -2.28$, $p = 0.033$, Cohen's $d = 0.47$). (d) There was a significant correlation

746 between IAT score change and activation change in the left amygdala under subthreshold presentation

747 after controlling for the novelty of the humanoid robot and human images ($r = 0.51$, $p = 0.014$).

748 **Tables**

749 Table 1. Attitudes toward robots questions.

Items	5-point rating (1- completely disagree, 5-completely agree)
1. Humanoid robots could improve work efficiency.	
2. Humanoid robots could do jobs that human can't finish.	
3. Humanoid robots could improve the quality of life for humans.	
4. Humanoid robots would consume a lot of resources.	
5. Humanoid robots could have unexpected dangers.	
6. Humanoid robots would disrupt human life.	

750 Table 2. Stimulus of evaluative conditioning task.

Label	Stimulus
CSs	20 humanoid images, 20 human images
USs	
Positive	
Images	Chicken ^a (14 ^b), Dog 2 (18), Island 1 (29), Apple (77), Cat 8 (781)
Words	Fantastic, Enjoyable, Fabulous, Excellent, Magnificent
Neutral	
Images	Bug 8 (234), Clothes Rack 1 (318), Plastic Cup (386), Butterfly 1 (451), Graph 2 (724)
Words	Odd, Stiff, Cold, Material, Muddy
Target	
Images	Antique 1 (292), Antique 2 (293), Antique 3 (295)
Words	Three meanings of antique
Fillers	
Images	Dock (424), Iron Bridge (406), Wolf 1 (547), Insect 10 (603), Plant 2 (842), Tool (785), Tree 4 (534), Locust (310), Frame (298), Shanghai 2 (387), City 4 (665), Tortoise (601), Hippo (482), Hair Drier (329), River (401), Mountain 7 (732)
Words	Bowl, Wine, Rock, Bench, Glass, Avenue, Boxer, Trunk, Rattle, Spray, Icebox, Ketchup, Radiator, Whistle, Nursery, Pamphlet, Thermometer

751 ^a CAPS image name.

752 ^b CAPS image numbers.

753

754 Table 3. Stimulus score of evaluative conditioning task.

	Valence		Arousal		Dominance	
	Mean	SD	Mean	SD	Mean	SD
USs						
Positive	7.27	0.18	6.25	0.45	7.11	0.42
Neutral	5.06	0.01	4.33	0.43	5.91	0.62
Filler	5.01	0.20	4.51	0.53	5.24	0.59
Target	5.24	0.08	4.10	0.16	5.59	0.01

755 Table 4. Stimulus arrangement in each block of evaluative conditioning task.

Trial	Arrangement
s	
10	US + CS ^a
3	Target stimulus ^b
3	Target stimulus + Neutral filler
10	Neutral filler
10	Neutral filler + Neutral filler
5	Blank screen
20	Blank screen (preceding and following CS-US pairings)

756 ^a US appeared simultaneously with CS for 10 trials in each block.

757 ^b Target stimulus appeared alone for 3 trials in each block.