1 **Neural signals in amygdala predict implicit prejudice toward an**
2 **ethnic outgroup**

3

4 **Abbreviated Title:** Decoding ethnic prejudice

5

6 Keise Izuma[1],*, Ryuta Aoki[2], Kazuhisa Shibata[3] & Kiyoshi Nakahara[2]

7

8 [1] Department of Psychology, University of York, Heslington, York, YO10 5DD, UK.

9 [2] Research Center for Brain Communication, [4]School of Information, Kochi

10 University of Technology, Kami, Kochi 782-8502, Japan.

11 [3] Graduate School of Informatics, Nagoya University, Furo-cho, Chikusa-ku, Nagoya

12 464-8601, Japan

13

14 * Correspondence should be addressed to Keise Izuma,

15 Current address

16 Department of Psychology, University of Southampton,

17 University Road, Southampton, SO17 1BJ, UK.

18 Email: K.Izuma@soton.ac.uk

19

22

23 **Abstract**

24 Racial and ethnic prejudice is one of the most pressing problems in modern societies.

25 Although previous social neuroscience research has suggested the amygdala as a key

26 structure in racial prejudice, it still remains elusive whether the amygdala activity

27 reflects negative attitudes toward an outgroup or other unrelated processes. The

28 present study aims to rigorously test the role of the amygdala in negative prejudice

29 toward an outgroup. Seventy Japanese individuals passively viewed images related to

30 an ethnic outgroup (South Korea) inside a functional magnetic resonance imaging

31 scanner. Using Multi-Voxel Pattern Analysis (MVPA), we found that Japanese

32 individuals' level of implicit (but not explicit) evaluations of South Korea could be

33 predicted from neural signals in the left amygdala. Our result further suggested that

34 the medial and lateral parts of amygdala play different roles in implicit evaluations. In

35 contrast to the MVPA findings, conventional univariate analyses failed to find any

36 reliable relationship between brain activation and both implicit and explicit

37 evaluations. Our findings provide evidence for the amygdala's role in representing an

38 implicit form of prejudice and highlight the utility of the multivariate approach to

39 reveal neural signatures of this complex social phenomenon.

40

41 **Keywords:** fMRI, MVPA, IAT, implicit attitude

42

43 **Highlights**

44 • Amygdala is considered as a candidate region involved in negative prejudice

45 • However, past neuroimaging studies on prejudice have generated mixed findings

46 • We investigated the neural signatures of ethnic prejudice using MVPA

47 • Our results demonstrated the link between left amygdala and implicit evaluations

48

49    **Introduction**

50    Historically, prejudice has been of primary interest to social scientists for a long time.

51    Despite the long history, research on intergroup relations remains of high importance

52    as intergroup relations continue to be fluid and volatile in modern societies. In the

53    past two decades, social neuroscientists have investigated neural mechanisms of

54    prejudice by using neuroimaging methods (Amodio, 2014; Chekroud et al., 2014;

55    Kubota et al., 2012). However, this endeavor has turned out to be challenging due to

56    the high complexity of the phenomenon.

57         Because of its well known role in fear learning (Fendt and Fanselow, 1999;

58    Pape and Pare, 2010), amygdala is considered as a primary candidate as a neural

59    substrate underlying negative prejudice (Phelps et al., 2000). In fact, the amygdala has

60    been most frequently reported in previous neuroimaging research on racial prejudice

61    (Amodio, 2014; Chekroud et al., 2014; Kubota et al., 2012). However, past findings

62    are not necessarily consistent, and the amygdala's involvement in negative prejudice

63    toward an outgroup still remains unclear. For example, while some studies found

64    increased amygdala activation in response to racial outgroup faces compared to

65    ingroup faces (Cunningham et al., 2004; Hart et al., 2000; McCutcheon et al., 2018;

66    Wheeler and Fiske, 2005), a number of other studies failed to find such amygdala

67    activations (Brosch et al., 2013; Cassidy and Krendl, 2016; Gilbert et al., 2012; Golby

68    et al., 2001; Li et al., 2016; Mattan et al., 2018; Phelps et al., 2000; Richeson et al.,

69    2003; Stanley et al., 2012; Terbeck et al., 2015). Furthermore, even among those

70    studies that found amygdala activations, its functional interpretations differ. Whereas

71    some studies have provided evidence that amygdala activation in response to

72    outgroup faces reflects negative emotional reaction to an outgroup (i.e., culturally

73    learned negative association) (Lieberman et al., 2005; Phelps et al., 2000; Telzer et al.,

74    2013b), other studies demonstrated that amygdala activity simply reflects novelty of

75    outgroup faces (Cloutier et al., 2014; Hart et al., 2000; Telzer et al., 2013a). Still other

76    studies showed that amygdala activity is modulated by skin-tone (Ronquillo et al.,

77    2007), gaze direction (Richeson et al., 2008) and status (Mattan et al., 2018).

78    Furthermore, it is known that amygdala activity is sensitive to facial features such as

79    subtle differences in pupil size (Demos et al., 2008), trustworthiness (Said et al.,

80    2009; Winston et al., 2002) and general valence evaluation of faces (Todorov and

81    Engell, 2008). These findings suggest that the amygdala activations in response to

82    outgroup faces found in the previous studies (Cunningham et al., 2004; Hart et al.,

83    2000; McCutcheon et al., 2018; Wheeler and Fiske, 2005) might not be directly

84    related to prejudice toward an outgroup.

85       There exist three studies (Brosch et al., 2013; Cunningham et al., 2004; Phelps

86    et al., 2000) that reported a correlation between amygdala activity and individual

87    differences in implicit attitudes toward an outgroup as measured by implicit measures

88    of attitudes such as an Implicit Association Test (IAT) (Greenwald et al., 1998).

89    However, these across-participant correlations were based on a considerably small

90    sample size (n = 12-13) (Brosch et al., 2013; Cunningham et al., 2004; Phelps et al.,

91    2000), and six other studies with larger sample sizes (n = 15-44) failed to find such a

92    link (Cassidy and Krendl, 2016; Cassidy et al., 2016; Gilbert et al., 2012; Li et al.,

93    2016; Richeson et al., 2003; Terbeck et al., 2015), suggesting that these findings need

94    to be interpreted with extreme caution (see Power Analysis below for more

95    discussion).

96       Thus, although several previous studies have reported amygdala activation in

97    response to outgroup faces, evidence for the amygdala's involvement in prejudice is

98    still weak. Furthermore, evidence for the involvement of other brain regions (such as

99    the anterior insula, striatum and fusiform face area [FFA]) is, if anything, even more

100   mixed (Amodio, 2014; Kubota et al., 2012). Given the complex and multifaceted

101   nature of prejudice, the field requires a more powerful approach to unveil its neural

102   signatures, especially the role of the amygdala in racial and ethnic prejudice.

103       In the present study, rather than simply contrasting brain activations in response

104   to outgroup vs. ingroup faces where observed activations can be explained by a

105   variety of different factors as discussed above, we directly tested if negative

106   evaluations of an outgroup as measured by the IAT are related to neural activations,

107   especially in the amygdala. The IAT measures implicit attitude as the strength of

108   associations between representations of groups and valenced semantic concepts

109   (Greenwald et al., 1998). Thus, if the amygdala reflects an *automatic* negative

110   evaluation of an outgroup, which has been acquired via direct or observational fear

111   learning processes (i.e., repeated associations with a social group and negatively-

112   valenced information), its activities should be associated with IAT scores.

113       To test the role of the amygdala in negative prejudice toward an outgroup, we

114   employed a machine learning method (Multi-Voxel Pattern Analysis; MVPA (Haynes,

115   2015; Norman et al., 2006)). MVPA makes it possible to detect a wider variety of

116   signals and has been proven to be effective for detecting differences in cognitive or

117   affective states that cannot be probed by conventional univariate analysis (e.g., (Izuma

118   et al., 2017; Jimura and Poldrack, 2012; Sapountzis et al., 2010) and thus is more

119   suitable for investigating complex neural representations of prejudice. Seventy

120   Japanese university students passively viewed Japan- and South Korea-related images

121   (Figure 1) inside an fMRI scanner. After the scanning, they completed explicit and

122   implicit measures of attitudes toward each of Japan and South Korea. Using the

123   MVPA, we tested if neural signals especially in the amygdala can predict individuals'

124    level of implicit and explicit negative evaluations of South Korea.

125          The present study focuses on the intergroup relation between Japan and South

126    Korea. Despite that prejudice is a world-wide problem (Landis and Albert, 2012;

127    Noor and Montiel, 2009), the vast majority of past neuroimaging studies on prejudice

128    focused on the intergroup relation between White vs. Black Americans, and prejudice

129    in other intergroup contexts has been rarely investigated (see (Bruneau and Saxe,

130    2010; McCutcheon et al., 2018) for notable exceptions). Fiske (2017) recently argued

131    that stereotypes of race, ethnicity and region are more variable across different

132    cultures compared to stereotypes of gender and age, and thus it is important to

133    formally test whether past social neuroscience findings only apply to the White vs.

134    Black Americans intergroup context or can be generalizable to other contexts. The

135    relationship between the two neighboring countries of Japan and South Korea has

136    deteriorated especially in recent years due to several disputes over political, historical

137    and territorial issues. For example, the recent analysis of comments on a Japanese

138    online news site showed that among 1,000 randomly-selected comments about South

139    Korea, 87.6% expressed negative attitudes toward South Korea, while only 0.7%

140    expressed positive attitudes (11.7% neutral) (Cho, 2017). BBC World Service polls

141    taken between 2010-2014 have also shown that Japanese viewed South Korea more

142    and more negatively over the past years (BBC_World_Service, 2010, 2014). Japanese

143    participants' negative implicit and explicit attitudes toward South Korea were also

144    confirmed by our behavioral data (see below). Thus, the current relationship between

145    the two countries is an ideal intergroup context to explore the neural signatures of

146    prejudice outside of the White vs. Black American context.

147

148    **Materials and Methods**

149 **Participants:** Seventy-one right-handed Japanese university students aged 18-22

150 years with no psychiatric history participated the study. One participant was excluded

151 from the analyses due to excessive error rate during the IAT ($> 25\%$). The final

152 sample consists of $n = 70$ (27 females, mean age = 18.9 years, SD = 1.11).

153 Participants were recruited from a subject pool of the Kochi University of Technology.

154 All participants gave written informed consent for participation, and ethics approval

155 for the study was granted by the Kochi University of Technology Ethics Board.

156

157 **Power Analysis:** The three previous studies (Brosch et al., 2013; Cunningham et al.,

158 2004; Phelps et al., 2000) that reported significant correlation between amygdala

159 activities and race IAT scores reported correlations of $r = 0.58$ (Phelps et al., 2000) $r$

160 $= 0.71$ (Cunningham et al., 2004) and $r = 0.62$ (Brosch et al., 2013). However, due to

161 a small sample size ($n = 12$, 13 and 13 respectively) together with likely large

162 measurement error of implicit measures of attitudes (e.g., Nosek et al. (2007) reported

163 that the median test-retest reliability of IAT is 0.56), these correlations are highly

164 likely to be inflated (Loken and Gelman, 2017; Yarkoni, 2009). Thus, we estimated

165 the effect size of $r = 0.30$ (a correlation between actual IAT scores and predicted

166 scores based on neural activation patterns in the amygdala; see "fMRI Data Analysis

167 (MVPA)" below for more details) for the present study (i.e., about half of the smallest

168 of the three correlations). With this effect size, a sample size of $n = 68$ should achieve

169 statistical power of 80% ($\beta = 0.2$) with $\alpha = 0.05$ (one-tailed). To account for potential

170 data loss (e.g., due to excessive head motion inside an fMRI scanner), we aimed to

171 recruit 70 participants (and actually recruited 71 participants).

172

173 **Stimuli**: Inside an fMRI scanner, participants were presented with 20 pictures each
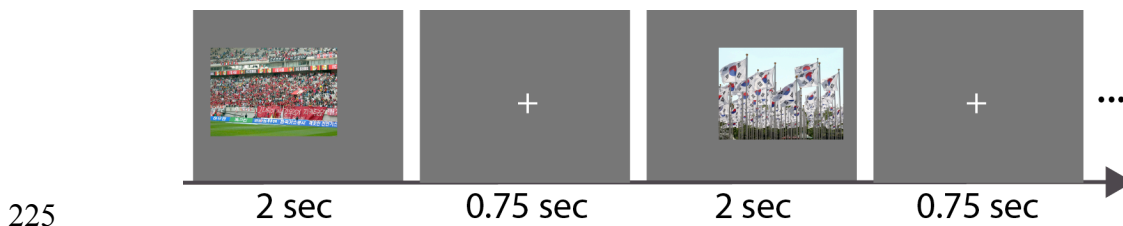
174    for the Japan and South Korea conditions. The images depicted Japanese or Korean

175    people and/or national flags (e.g., South Korea fans chanting and cheering their

176    national team on street, South Korea national football team approaching fans after

177    winning a match, South Korea national flag overlaid onto a map of South Korea, etc.)

178    so that it is evident which country each picture represents (based on uniforms they

179    wear and/or national flags). It should be noted that although there are images that

180    depict sports scenes, each image depicted either Japan or South Korea only, and no

181    picture depicted, for example, a scene where the two countries compete with each

182    other. All pictures were taken from the Internet (i.e., Google image search). Unlike

183    the majority of past studies, we avoided using faces as experimental stimuli to

184    minimize the effect of individual differences in facial trustworthiness/attractiveness

185    judgment on the activity of the amygdala (and other brain regions) (Todorov and

186    Engell, 2008). Furthermore, since ethnic prejudice is based on evaluations of an

187    ethnic group, stimuli symbolizing South Korea as a whole (rather than each

188    unfamiliar Korean individual) could evoke clearer neural representations related to

189    evaluations of South Korea. Note that visual stimuli used in the Japan and South

190    Korea conditions were not equated for lower visual features (e.g., contrast and

191    luminance) as the main purpose of the current study is to test whether activation

192    patterns in the amygdala can explain *individual differences* in the implicit evaluations

193    of the outgroup (all participants viewed the same stimuli).

194        It should be noted that since most of the images we used depicted sports

195    contexts, one might think that it may be possible that the relation between implicit

196    prejudice toward South Korea and neural signals in the amygdala (or any other

197    regions) may be explained by individual differences in a sense of rivalry rather than

198    implicit prejudice. To refute this possibility, we conducted an additional behavioral

199    experiment (n = 49) and confirmed that there is no correlation between implicit

200    prejudice and a sense of rivalry toward South Korea so that any relation found

201    between implicit prejudice and neural signals cannot be attributed to a sense of rivalry

202    felt while viewing South Korea-related images (see the "Additional behavioral results:

203    between implicit prejudice and a sense of rivalry toward South Korea" section below).

204

205    **Procedure and tasks:** During the fMRI scanning, participants viewed 30 blocks

206    comprising 1) Japan blocks, 2) South Korea blocks, and 3) rest (i.e., a fixation cross)

207    blocks (10 blocks each). Each of the Japan and South Korea blocks started with a cue

208    (either "Japan" or "South Korea"; presented for 1 sec) indicating which pictures they

209    were going to see in that trial. After the cue, four different pictures were presented in

210    each block (each picture was presented for 2 sec followed by inter-trial-interval of

211    0.75 sec; see Figure 1). Each block lasted 12 sec. In each trial, an image was

212    presented slightly to the left or right side of the screen, and participants were asked to

213    press one of two keys to indicate whether an image appeared to the left or right side of

214    the screen. Importantly, to prevent participants from suppressing emotions which are

215    automatically evoked by viewing each picture, before the scanning, they were told

216    that the study was interested in investigating how we perceive various social images

217    (they were never explicitly asked to think about how they feel about each group

218    during the scanning). Participants completed one 6-min fMRI run (12 sec × 30

219    blocks). All participants also completed another fMRI run for an independent project

220    (where they were presented with food images [the data will not be reported here]).

221    Since we are interested in how individual differences in ethnic prejudice are related to

222    brain activations, the order of blocks (and the order of trials within each block) was

223    fixed across all participants.

224



225

**Figure 1. Experimental task.** Sample stimuli in the South Korea condition. Within each South Korea block, four South Korea-related images were presented. Participants were asked to indicate whether an image appeared to the left or right side of the screen. Note that due to copyright restrictions, the two South Korea-related images shown here are similar, but not the same as the ones used in the actual experiment.

After fMRI scanning, participants were asked to perform a total of four behavioral tasks; 1) Trust Game, 2) Single-Category Implicit Association Test (SC-IAT) (Karpinski and Steinman, 2006) which measures implicit attitude toward South Korea, 3) SC-IAT which measures implicit attitude toward Japan, and 4) Explicit measure of attitude (semantic differential) toward Japan and South Korea. Since we are interested in whether the amygdala is related to attitudes toward South Korea (not *relative* attitudes toward South Korea compared to Japan), we used the SC-IAT (Karpinski and Steinman, 2006) instead of the conventional IAT (Greenwald et al., 1998). All participants performed the four tasks in this fixed order.

First, in order to measure their trust behavior toward Japan and South Korea in an incentive-compatible manner, we asked each participant to perform a series of single-shot modified Trust Games with 28 Japanese and 28 South Korean male students (Stanley et al., 2011). All participants played a role of a trustor, and in each trial, they were presented with a picture of an East Asian male and a South Korea or Japanese national flag to indicate the nationality of each partner. In reality, all pictures were Japanese individuals (the post-experimental interview confirmed that no

10

249   participant had any doubt on the nationality of each individual presented during the

250   Trust Game). The two face sets for Japan vs. South Korea were matched on

251   trustworthiness based on an independent pilot study [n = 23, 8 females]). In each trial,

252   participants were asked to decide how much of 500 Japanese yen (~$5) to offer to the

253   partner (trustee; between 0 to 500 yen in increments of 100 yen). Before the Trust

254   Games, participants were instructed that each partner would receive quadruple the

255   amount they offered. They were further told that we had conducted a behavioral

256   experiment in several South Korean and Japanese universities, and each of the 56

257   individuals had participated the study and already made the decision to return half or

258   keep all of the money they received (in reality, such an experiment was not

259   conducted). During the Trust Games, participants were not informed of each partner's

260   decision, and they were told that one trial would be selected randomly at the end of

261   the experiment, and the outcome of the randomly-selected trial would be implemented.

262   In reality, partner's decision in the randomly selected trial was determined randomly.

263       Following the Trust Games, each participant was asked to perform an SC-IAT

264   (Karpinski and Steinman, 2006) which measure implicit attitude toward South Korea.

265   The IAT included eight positive (e.g., *Joy, Love, Wonderful*) and eight negative words

266   (e.g., *Agony, Terrible, Nasty*; all words were translated into Japanese). The South

267   Korean category included typical Korean names (e.g., *Han, Kim, Myong*). In addition,

268   each participant performed another SC-IAT that measures implicit attitude toward

269   Japan where typical Japanese names (e.g., *Shima, Nakata, Ono*) were presented

270   instead of the Korean names. All Korean and Japanese names were matched on word

271   length. Each of the two SC-IATs consisted of four blocks; Blocks 1 & 3 (24 trials

272   each) were practice blocks, while Blocks 2 & 4 (72 trials each) were test blocks.

273   During the South Korea SC-IAT, in Blocks 1 & 2, two category labels of Good and

274    South Korea, were presented in the top left corner, while Bad was presented in the top

275    right corner. Participants were instructed to categorize each target word presented in

276    the center of the screen as soon as they could by pressing either "E" or "I" key on a

277    keyboard. An incorrect response was followed by a red "X" presentation in the center

278    of the screen, and it remained on the screen until a participant pressed the other

279    (correct) key. No feedback was presented after a correct response. Reaction times

280    (RTs) for correct responses were used for data analysis. The inter-trial interval was

281    300 ms. In Blocks 3 & 4, a category label Good was presented in the top left corner,

282    while Bad and South Korea were presented in the top right corner. Just like Karpinski

283    and Steinman (2006), during Blocks 1 & 2, South Korea words, good words, and bad

284    words were presented in a 7:7:10 ratio so that 42% of correct responses were on the

285    "I" key and 58% of correct responses were on the "E" key. Similarly, during Blocks 3

286    & 4, South Korea words, good words, and bad words were presented in a 7:10:7 ratio

287    so that 58% of correct responses were on the "I" key and 42% of correct responses

288    were on the "E" key. The Japan SC-IAT was the same as South Korea SC-IAT except

289    that the category South Korea was replaced with Japan, and Japanese names were

290    presented instead of the Korean names as a target word.

291        After completing the two SC-IATs, each participant was asked to complete a

292    semantic differential scale (Greenwald et al., 1998; Karpinski and Steinman, 2006),

293    which measures explicit attitudes toward Japan and South Korea. Participants rated

294    each of South Korea and Japan on six bipolar dimensions using a 7-point scale; *ugly-*

295    *beautiful, bad-good, unpleasant-pleasant, honest-dishonest, awful-nice and*

296    *unfavorable-favorable.* Finally, after completing a demographic questionnaire, they

297    were debriefed, thanked and paid 1,500-2,500 yen for their participation.

298

299    **Behavioral Data Analysis:** For each of the two SC-IATs, a score for each participant

300    was calculated using the D-score algorithm developed by Greenwald et al.

301    (Greenwald et al., 2003). First, after excluding trials whose RT was greater than

302    10,000 ms, we computed the mean RT for each of the four blocks. Second, we

303    computed a pooled standard deviation (SD) for Blocks 1 & 2 and another pooled SD

304    for Blocks 3 & 4. Third, we computed two differences; 1) mean RT in Block 1 - mean

305    RT in Block 3, and 2) mean RT in Block 2 - mean RT in Block 4. Fourth, each

306    difference was divided by its associated pooled SD from step 2 above. Finally, we

307    computed the mean of the two quotients from Step 4, which is a SC-IAT D score for

308    an individual.

309         Semantic differential scores for each participant were computed by averaging

310    the six bipolar scales separately for Japan and South Korea (Japan Cronbach's $\alpha$ =

311    0.80, South Korea Cronbach's $\alpha$ = 0.92). For both IAT and semantic differential

312    scores, positive numbers indicate more positive evaluation of a target group. To

313    compute correlations among all behavioral variables, in addition to SC-IAT and

314    semantic differentials scores for each of the Japan and South Korea conditions, we

315    computed the disparity in these scores between Japan and South Korea by subtracting

316    scores for the South Korea condition from those for the Japan condition. These

317    disparity indices represent *relative* implicit or explicit attitudes toward Japan relative

318    to South Korea. Similarly, we computed the difference between the average amounts

319    of money transferred to Japan vs. South Korea partners during the Trust Games. Note

320    that the average amount of money transferred in the Japan and South Korea

321    conditions were highly correlated with each other ($r(68) = 0.91$, $p < 0.001$), indicating

322    that decisions made during the Trust Game depended more strongly on the individual

323    differences in personality traits such as general trust or risk-taking than ethnic

324   attitudes. Accordingly, the average amount offered in each of the Japan and South

325   Korea conditions were not included in the correlational analyses (note that the average

326   amount offered in each of the Japan and South Korea conditions were not correlated

327   with SC-IAT or semantic differential scores [$-0.20 < rs < 0.04$, all $ps > 0.10$]). There

328   was one outlier in the Trust Game data (more than 3 SD from the mean), and this

329   participant was excluded when analyzing the Trust Game data. Finally, for paired t

330   tests, reported effect sizes are based on Dunlap et al. (Dunlap et al., 1996)

331

332   **fMRI data acquisition:** All fMRI data were acquired using a Siemens 3.0 Tesla

333   Verio scanner with a 32 channel phased array headcoil. For functional imaging,

334   interleaved T2*- weighted gradient-echo echo-planar imaging (EPI) sequences were

335   used to produce 40 contiguous 3.0-mm-thick trans-axial slices covering nearly the

336   entire cerebrum (repetition time [TR] = 2,500 ms; echo time [TE] = 25 ms; flip angle

337   [FA] = 90°; field of view [FOV] = 192 mm; 64 × 64 matrix; voxel dimensions = 3.0 ×

338   3.0 × 3.0 mm). A high-resolution anatomical T1-weighted image (1 mm isotropic

339   resolution) was also acquired for each participant.

340

341   **fMRI Data Pre-processing:** The fMRI data were analyzed using SPM8 (Wellcome

342   Department of Imaging Neuroscience) implemented in MATLAB (MathWorks).

343   Before data processing and statistical analysis, we discarded the first four volumes to

344   allow for T1 equilibration. After correcting for differences in slice timing within each

345   functional image volume, images were realigned using the mean image as a reference.

346   (note that no participant showed excessive head motion [i.e., 3mm] during the

347   scanning). Following realignment, the volumes were normalized to MNI space using

348   a transformation matrix obtained from the normalization of the first EPI image of

349    each individual participant to the EPI template using an affine transformation

350    (resliced to a voxel size of 3.0 × 3.0 × 3.0 mm). The normalized fMRI data were

351    spatially smoothed with an isotropic Gaussian kernel of 8 mm (full-width at half-

352    maximum). We used the smoothed fMRI data for both univariate as well as MVPA

353    analyses following Op de Beeck et al. (Op de Beeck, 2010) who showed that

354    smoothing can improve decoding performance when large-scale activation patterns

355    are assumed (e.g., (Chang et al., 2015). Note that since the voxel size of the amygdala

356    ROIs are small (left amygdala = 54 voxels, right amygdala = 63 voxels), we also

357    reported decoding results from unsmoothed data only for the amygdala ROIs.

358

359    **fMRI Data Analysis (Univariate Analysis):** We first ran a conventional general

360    linear model (GLM) analysis. In the GLM, each of the Japan and South Korea blocks

361    was separately modeled (duration = 12 sec). Six head motion parameters were also

362    included in the model as nuisance regressors. Three contrast images were created for

363    each participant; 1) a contrast image for the South Korea blocks (vs. implicit rest), 2)

364    a contrast image for the Japan blocks (vs. implicit rest) and 3) a contrast for the South

365    Korea vs. Japan blocks. These contrast images were submitted to the second level

366    analysis. The same three contrast images were also used in subsequent MVPA

367    analyses (see below).

368       In the second level analysis, for the South Korea contrast, South Korea SC-IAT

369    scores and South Korea semantic differential scores were entered as covariates to test

370    whether implicit or explicit attitudes were linearly related to activations in the brain.

371    Similarly, for the Japan contrast, Japan SC-IAT scores and Japan semantic differential

372    scores were entered as covariates. Finally, to test whether relative attitudes between

373    Japan and South Korea are linearly related to brain activations, the South Korea vs.

374    Japan contrasts were submitted into the second level analysis, and disparity in SC-

375    IATs scores (Japan SC-IAT minus South Korea SC-IAT), disparity in semantic

376    differential scores (Japan semantic differential minus South Korea semantic

377    differential) and disparity in Trust Game (the average amount offered to Japanese

378    partners minus South Korea partners) were entered as covariates. For the univariate

379    analysis, the statistical threshold was set at $p < 0.001$ voxelwise (uncorrected) and

380    cluster $p < 0.05$ (FWE corrected for multiple comparisons).

381         In addition to the whole-brain analysis, we also conducted the ROI analysis.

382    Previous neuroimaging studies have identified a network of brain regions implicated

383    in racial prejudice (the so-called "prejudice network"), which includes the anterior

384    insula, striatum, ventral medial prefrontal cortex, orbitofrontal cortex as well as the

385    amygdala (Amodio, 2014). Among these brain regions, we focused especially on the

386    amygdala as previous studies have reported the link between amygdala activities and

387    race IAT scores (Brosch et al., 2013; Cunningham et al., 2004; Phelps et al., 2000).

388    The same left and right amygdala masks were used as the MVPA (see below). We

389    extracted beta values from each of the left and right amygdala ROIs and tested

390    whether the activities were correlated with each of the behavioral indices of prejudice.

391

392    **fMRI Data Analysis (MVPA):** In order to decode implicit and explicit

393    attitudes toward South Korea (and Japan) from neural signals, we employed support

394    vector regression (SVR) (Drucker et al. (1997), as implemented in LIBSVM

395    (http://www.csie.ntu.edu.tw/~cjlin/libsvm/), with a linear kernel and the default

396    regularization parameter of c = 1 (default; Note that all other parameters were also set

397    to their default values). We previously used the SVR and successfully decoded

398    individuals' attitudes toward familiar celebrities (Izuma et al., 2017) and implicit self-

399    esteem (implicit attitude toward the self) (Izuma et al., 2018).

400        As stated above, since the amygdala plays a major role in fear learning (e.g., a

401    learned association between an outgroup and negativity) and is one of the brain

402    regions that have been most frequently reported in previous neuroimaging research on

403    racial prejudice (Amodio, 2014; Chekroud et al., 2014; Kubota et al., 2012), we

404    primarily focused on the amygdala in the present study. In addition, we also ran

405    exploratory MVPA using signals from each of 79 anatomical regions (see below).

406    Each ROI was defined by using a WFU pickatlas toolbox for SPM (Maldjian et al.,

407    2003).

408        To test whether neural signals in each of right and left amygdala ROIs can

409    predict implicit and explicit evaluations of South Korea, we used SC-IAT scores and

410    semantic differential scores for South Korea to decode implicit and explicit attitudes,

411    respectively. Contrast images for the South Korea blocks were used to decode implicit

412    and explicit attitudes toward South Korea. Although our primary interest is Japanese

413    individuals' attitudes toward South Korea, we repeated the same analyses for the

414    Japan condition to test whether neural signals in the prejudice network also represent

415    implicit and explicit attitudes toward Japan. Furthermore, we investigated whether we

416    can decode *relative* attitudes between Japan and South Korea based on relative neural

417    signals (i.e., the South Korea vs. Japan contrast images). For this purpose, we used

418    SC-IAT disparity scores, semantic differential disparity scores and Trust Game

419    disparity scores as labels.

420        For each MVPA analysis, we computed decoding performance using the 10-

421    fold balanced cross-validation procedure (Cohen et al., 2010; Izuma et al., 2018); we

422    first divided participants into 10 groups (7 individuals in each group) in a way so that

423    when decoding implicit attitude toward South Korea, these 10 groups had roughly the

424   same means and variances of South Korea SC-IAT scores (or semantic differential

425   scores toward South Korea when decoding explicit evaluations). In each cross-

426   validation, one group was left out, and the SVR was performed using the data from

427   participants in all other groups and then tested on the participants in the left-out group.

428   This procedure was repeated for each group (a total of 10 times), and a Pearson's

429   correlation coefficient between actual SC-IAT scores (or semantic differential scores)

430   and predicted scores was computed.

431        Furthermore, to explore whether brain regions outside of the prejudice network

432   can predict implicit, explicit evaluations, or disparity scores, we repeated the above-

433   mentioned analyses using neural signals from each of a total of 73 regions outside of

434   the prejudice network.

435        Prediction performance in each ROI was evaluated using a permutation test. We

436   created 5,000 randomly shuffled permutations of IAT or semantic differential scores

437   and ran the SVR using the permutated data in each ROI to obtain a distribution of

438   correlations between predicted and actual scores under the null hypothesis. Note that

439   behavioral scores were shuffled within each of the 10 fold groups so that the averages

440   scores in the 10 fold groups were maintained across the permutations. For the regions

441   outside of the amygdala, false discovery rate (FDR) (Benjamini and Hochberg (1995)

442   correction for multiple comparisons was applied  ($q < 0.05$).

443

444   **Results**

445   **Behavioral results:** We confirmed that Japanese participants had negative implicit

446   and explicit attitudes toward South Korea, and their negative attitudes were also

447   reflected in trust behavior as measured by the Trust Game. First, not surprisingly,

448   participants had more positive explicit attitudes toward Japan (mean = 5.34, SD =

449     0.99) compared to South Korea (mean = 3.63, SD = 1.12). Semantic deferential scores

450     were significantly more positive for Japan compared to South Korea ($t(69)$ = 11.14, $p$

451     < 0.001, $d$ = 1.62, paired t-test; Figure 2a). While semantic differential scores for the

452     Japan were significantly higher than the midpoint (i.e., 4) ($t(69)$ = 11.36, $p$ < 0.001,

453     Cohen's $d$ = 1.36, one-sample t-test), scores for South Korea were significantly lower

454     than the midpoint ($t(69)$ = 2.79, $p$ = 0.007, Cohen's $d$ = 0.33, one-sample t-test),

455     indicating that Japanese participants possess positive and negative explicit attitudes

456     toward Japan and South Korea, respectively.

457        Similar results were obtained with the implicit attitude measure. Japan SC-IAT

458     scores were significantly higher than South Korea SC-IAT scores (Japan mean = 0.05,

459     SD = 0.32, South Korea mean = -0.30, SD = 0.31; $t(69)$ = 7.72, $p$ < 0.001, $d$ = 1.12,

460     paired t-test; Figure 2b), indicating more positive implicit attitude toward Japan

461     compared to South Korea. While Japan SC-IAT scores were not significantly

462     different from zero ($t(69)$ = 1.28, $p$ = 0.21, Cohen's $d$ = 0.15, one-sample t-test),

463     South Korea SC-IAT scores were significantly negative ($t(69)$ = 8.13, $p$ < 0.001,

464     Cohen's $d$ = 0.97, one-sample t-test), suggesting the presence of negative implicit

465     attitude (i.e., implicit prejudice) toward South Korea and, on average, neutral implicit

466     attitude toward Japan. It should be noted that since all participants completed the two

467     SC-IATs in the same order (South Korea SC-IAT first), the result might be

468     confounded with the effect of task order (e.g., practice effect) (however, see Study 4

469     of Karpinski and Steinman (2006) for data showing that the order of two SC-IATs has,

470     if anything, a very limited effect).

471        Participants' differential implicit and explicit evaluations of Japan vs. South

472     Korea were also reflected in trust behavior. The data from the Trust Games revealed

473     that participants transferred significantly more money to Japanese partners compared

474    to South Korean partners (Japan mean = 194.4 yen, SD = 114.7, South Korea mean =

475    172.2 yen, SD = 110.0; $t(68) = 5.54$, $p < 0.001$, $d = 0.20$, paired t-test; Figure 2c).
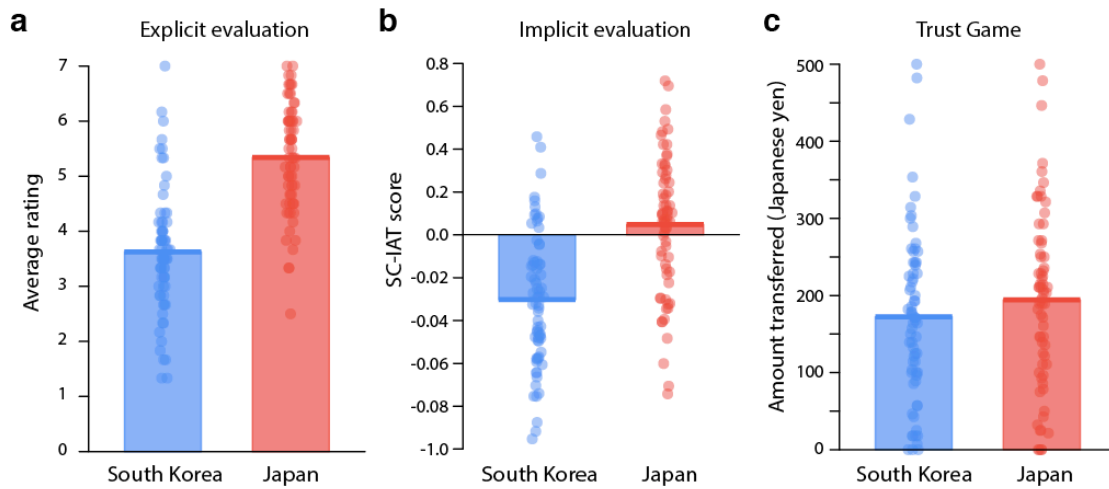


476

477    **Figure 2. Behavioral results**. **a**). Explicit evaluation (average semantic differential
478    scores) of South Korea and Japan. Higher numbers indicate more positive explicit
479    evaluation of each ethnic group. **b**). Implicit evaluation (average SC-IAT scores) of
480    South Korea and Japan. Higher numbers indicate more positive implicit evaluation of
481    each group. **c**). The average amount of money transferred to South Korean and
482    Japanese partners in the Trust Games.
483

484         Implicit and explicit evaluations were uncorrelated with each other, which is

485    largely consistent with previous research (Hofmann et al., 2005). All correlations

486    between implicit (SC-IATs) and explicit evaluations (semantic differential) were non-

487    significant (-0.06 < $rs$ < 0.14, $ps$ > 0.23). Unlike Stanley et al. (2011), our data

488    showed that differential behaviors during the Trust Games with Japan vs. South Korea

489    partners were not correlated with implicit evaluations (disparity in SC-IAT scores)

490    (but see (Oswald et al., 2013) for a meta-analysis demonstrating the low predictive

491    validity of the IAT). Trust Game disparity scores were also not correlated with

492    explicit evaluations (disparity in semantic differential scores). Two SC-IATs were

493    significantly correlated with each other ($r(68) = 0.27$, $p = 0.026$), suggesting that

494    those who possess more positive implicit attitudes toward Japan tend to have more

495    positive implicit attitudes toward South Korea as well. Similarly, two semantic

496    differential scores were significantly correlated with each other ($r(68) = 0.26$, $p =$

497    0.030). All correlations across the behavioral variables are shown in Table 1.

498      Inside the fMRI scanner, each participant performed a simple button press task

499    (i.e., indicate whether an image appeared to the left or right side of the screen), and

500    their performance was nearly perfect (Japan block = 98.4%, South Korea block =

501    97.9%). There was no significant difference in performance between the Japan vs.

502    South Korea blocks ($t(69) = 1.22$, $p = 0.23$, $d = 0.11$, paired t-test). Participants'

503    reaction times (RTs) were significantly slower in the South Korea blocks (mean RT =

504    670 ms) compared to the Japan blocks (mean RT = 661 ms; $t(69) = 2.06$, $p = 0.043$, $d$

505    $= 0.06$, paired t-test). Importantly, performance and RTs during the button press task

506    inside the fMRI scanner were not correlated with both implicit (SC-IAT) and explicit

507    evaluations (semantic differential) (Japan $-0.09 < rs(68) < 0.16$, $ps > 0.18$; South

508    Korea $-0.08 < rs(68) < 0.15$, $ps > 0.22$).

509

510    **fMRI Results (MVPA):** We first investigated whether neural signals in the amygdala

511    can predict implicit (and explicit) attitudes toward South Korea. The results revealed

512    that neural signals in the left amygdala significantly predicted the level of implicit

513    evaluations of South Korea ($r(68) = 0.31$, $p_{perm} = 0.021$; Figures 3). Since one data

514    point in the predicted scores was identified as an outlier based on the Grubbs' test (p

515    $= 0.028$), we also computed Spearman's rank correlation, and the result remains

516    significant ($rs(52) = 0.28$, $p_{perm} = 0.035$). In contrast, the prediction based on the right

517    amygdala activations was not significant ($p_{perm} = 0.55$; see Table 2). Thus, the result

518    indicates that passive viewing of South Korea-related images automatically evoked

519    similar neural responses in the left (but not right) amygdala across different

520    individuals depending on their level of implicit evaluations. Given the small voxel

521     size of the amygdala ROIs, we also tried the same MVPA analysis on unsmoothed

522     data and found that prediction performance was improved in the left amygdala ($r$(68)

523     = 0.40, $p_{perm}$ = 0.005), but that in the right amygdala remain non-significant ($p_{perm}$ =

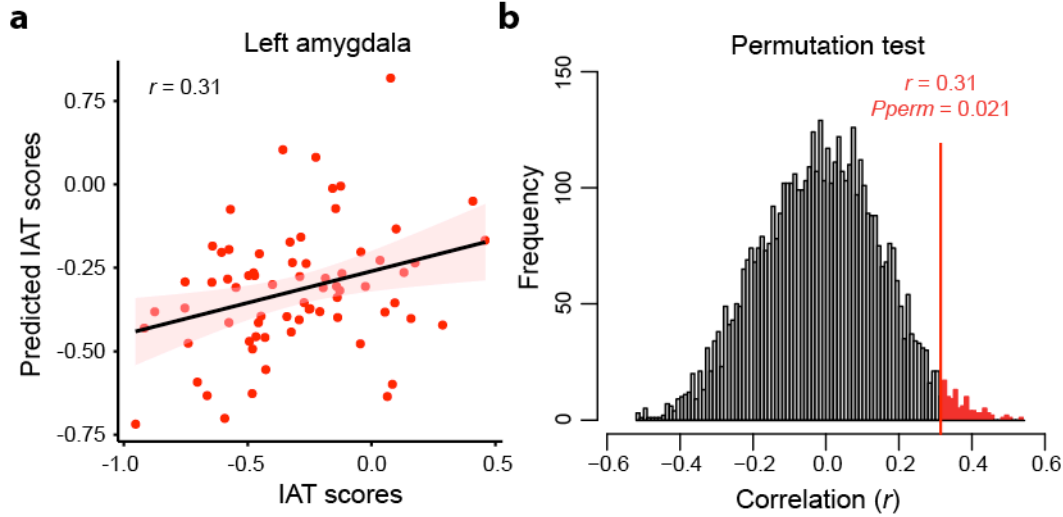524     0.98).

525



526

527     **Figure 3. MVPA results. a**). A scatter plot showing a correlation between
528     participants' SC-IAT scores for South Korea (i.e., implicit evaluations) and predicted
529     scores based on neural activations in the left amygdala ROI ($r$(68) = 0.31). The
530     predicted scores were cross-validated by using a 10-fold cross-validation procedure.
531     **b**). A histogram showing the distributions of correlation coefficients between actual
532     and predicted SC-IAT scores with randomly permutated data (5,000 times). The
533     correlation with actual data was significant at $p_{perm}$ = 0.021.

534

535     Since the possibility that different sub-regions of the amygdala may play

536     different roles in prejudice has been suggested previously (Chekroud et al., 2014), we

537     looked at a distribution of weight values in the left amygdala ROI. It revealed that

538     weight values of voxels in the medial part of the amygdala tend to be positive, while

539     weight values of voxels in the lateral part of the amygdala tend to be negative (see

540     Figure 4a). This is confirmed by a significant Spearman's rank correlation between

541     weight values (converted to 1 = positive weight value or -1 = negative weight value)

542     and x-coordinates ($r_s$(52) = 0.31, $p$ = 0.025). On the other hand, weight values were

543 not correlated with both y- and z-coordinates (both *ps* > 0.49).

544     Although the above analysis suggested that the medial and lateral parts of the

545 left amygdala play different roles in implicit evaluations, the classifier weights need

546 to be interpreted with caution as they do not simply reflect the amplitude of a signal in

547 each voxel (see (Haufe et al., 2014; Haynes, 2015)). Thus, to further investigate how

548 different parts within the left amygdala are related to the individual differences in

549 implicit evaluations, we computed a Pearson correlation between South Korea SC-

550 IAT scores and the amplitude of the signals in the South Korea blocks (i.e., univariate

551 activations) in each of 54 voxels within the left amygdala ROI. 54 correlation

552 coefficients ranged from $r$ = -0.08 to $r$ = 0.30 (average $r$ = 0.09, standard deviation =

553 0.10). Since all of x-, y-, and z-coordinate values were not normally distributed, we

554 computed a Spearman's rank correlation between the correlation coefficients and each

555 of the x-, y-, and z-coordinates. Consistent with the weight map results reported above,

556 the results revealed that there was a significant positive correlation with x-coordinates

557 ($r_s$(52) = 0.77, $p$ < 0.001), indicating that the across-subject correlations between the

558 amplitudes of the South Korea univariate contrast and South Korea SC-IAT scores

559 tended to be more positive in more medial parts (i.e., central nuclei) of the left

560 amygdala (see Figure 4b). Note that since x-, y-, z-coordinates are not completely

561 independent (i.e., the left amygdala ROI is not a complete sphere), we also found a

562 significant positive correlation with y-coordinates ($r_s$(52) = 0.47, $p$ = 0.003), and z-

563 coordinates also showed a significant trend ($r_s$(52) = 0.23, $p$ = 0.09) (see Figure 4).
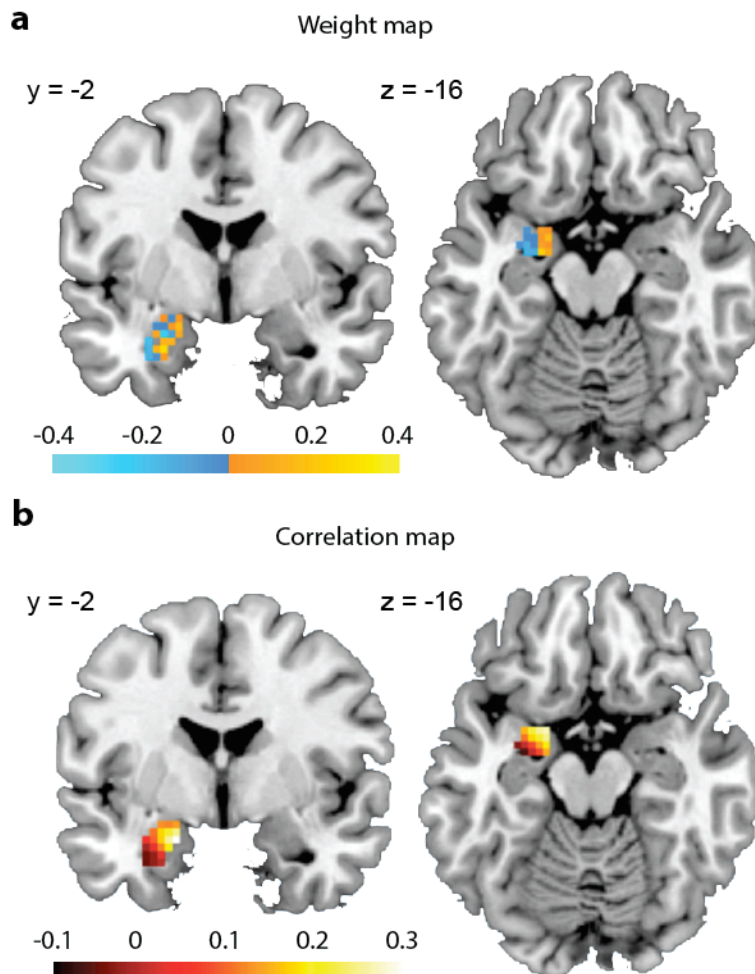
**Figure 4. Weight and correlation maps within the left amygdala ROI. a**). Weight values of the implicit evaluation decoder. **b**). Across-subject correlational values between univariate activations (i.e., the South Korea contrast) and implicit-evaluations.

We repeated the same MVPA analysis with South Korea SC-IAT scores for each region outside of the amygdala ROIs. Out of a total of 79 regions, only 4 regions (right anterior insula, left middle frontal gyrus, right fusiform gyrus, and right lingual gyrus) could predict implicit evaluations of South Korea at $p_{perm} < 0.05$ (uncorrected), but none of them survived the FDR correction (see Table 3). Although it did not survive the corrected threshold, significant prediction by neural signals in the right anterior insula ($r(68) = 0.28$, $p_{perm} = 0.038$) may be notable as the anterior insula is one of the regions previously implicated in racial prejudice (known as the prejudice network) (Amodio, 2014).

579       Interestingly, in contrast to implicit evaluations of South Korea, both left and

580    right amygdala could not predict explicit evaluations (i.e., semantic differential

581    scores) of South Korea. Prediction of explicit evaluations was not significant even

582    when some of the 8 ROIs in the prejudice network were combined (the highest

583    decoding performance among a total of 247 combinations was $r(68) = 0.24$, $p_{perm} =$

584    0.071 [based on neural signals from the left and right caudate nucleus ROIs]).

585       Neural signals in each of the left and right amygdala ROIs were also unrelated

586    to both implicit and explicit evaluations of Japan (Table 2). Neural signals in regions

587    outside of the amygdala were also largely unrelated to explicit evaluations of South

588    Korea as well as both implicit and explicit evaluations of Japan (see Table 3).

589    Furthermore, neural signals both inside and outside of the amygdala ROIs were

590    unrelated to *relative* implicit and explicit evaluations (i.e., SC-IAT disparity scores

591    and semantic differential disparity scores) and Trust Game disparity scores.

592

593    **fMRI Results (Univariate Analyses):** The contrast between the South Korea vs.

594    Japan blocks revealed a significant activation only in the visual cortex (right: x = 12,

595    y = -97, z = -2, left: x = -18, y = -106, z = -5; 1,849 voxels). Even when a statistical

596    threshold was lowered to $p < 0.01$, we didn't find any activation in any regions in the

597    prejudice network. We tested whether activity in each of the left and right amygdala

598    ROIs was significantly correlated with SC-IAT disparity scores (i.e., relative implicit

599    prejudice), as reported previously (Brosch et al., 2013; Cunningham et al., 2004;

600    Phelps et al., 2000). The results revealed that amygdala activities in both hemispheres

601    were not correlated with relative implicit prejudice, and if anything, the correlations

602    were in the opposite direction (the more positive their attitudes were toward South

603    Korea relative to Japan, the higher the activity in the amygdala in response to South

604   Korea-related images relative to Japan-related images; left amygdala $r(68) = 0.19$, $p =$

605   0.11, and right amygdala $r(68) = 0.09$, $p = 0.44$). We further conducted the whole-

606   brain analysis and investigated whether SC-IAT disparity scores were significantly

607   correlated with activations in regions other than the amygdala. However, SC-IAT

608   disparity scores were not significantly related to activations in any of the brain regions

609   (neither positively nor negatively). Even when the threshold was lowered to $p < 0.01$,

610   we did not find any significant activation in the amygdala ROIs or any of the regions

611   in the prejudice network.

612         Similarly, activities in each of the left and right amygdala ROIs were not

613   correlated with disparity in semantic differential between Japan vs. South Korea (left

614   amygdala $r(68) = -0.01$, $p = 0.95$, and right amygdala $r(68) = 0.07$, $p = 0.59$). The

615   whole brain analysis showed that the explicit disparity scores were not correlated with

616   activation in any region except for right intraparietal sulcus (IPS; x = 45, y = -28, z =

617   61; 199 voxels). Those who have more positive explicit evaluations of South Korea

618   relative to Japan showed higher IPS activations in response to South Korea images

619   compared to Japan images.

620         We further tested whether activations during each of the Japan and South Korea

621   blocks (vs. implicit rest) were correlated with corresponding SC-IAT scores. South

622   Korea SC-IAT scores were not significantly related to activations in any of the brain

623   regions in response to South Korea-related images (neither positively nor negatively)

624   (see Figure 4b for the same analysis only within the left amygdala ROI). Similarly,

625   Japan SC-IAT scores were not significantly related to activations in any of the brain

626   regions in response to Japan-related images. The same analyses were repeated using

627   semantic differential scores instead of SC-IAT scores, but it revealed no significant

628   correlation between explicit attitudes and brain activations (for both Japan and South

629 Korea). Thus, our univariate analyses failed to show any reliable association between

630 implicit (and explicit) evaluations and neural activities. We also didn't find any

631 significant association between Trust Game behaviors and neural activities.

632

633 **Additional behavioral results: between implicit prejudice and a sense of rivalry**

634 **toward South Korea**

635      As stated above, we selected South Korea- and Japan-related images in a way so

636 that no image depicted a direct competition between Japan and South Korea.

637 Nonetheless, since most of the images we used depicted sports contexts, it may be

638 possible that the relation between implicit prejudice toward South Korea and neural

639 signals in the amygdala (or any other regions) may be explained by a sense of rivalry

640 rather than implicit prejudice toward South Korea. To refute this possibility, we

641 conducted an additional behavioral experiment and tested whether implicit prejudice

642 as measured by the Single-Category Implicit Association Test (SC-IAT) is related to a

643 sense of rivalry toward South Korea.

644      We recruited an additional independent sample of 49 university students (23

645 females, 18-24 years old, mean age = 19.8 years, SD = 1.36) for a behavioral

646 experiment (without neuroimaging). Data from one additional participant was

647 excluded from the analyses because his data was not correctly saved due to a

648 malfunction of the task presentation program. Like the main fMRI experiment, all

649 participants were recruited from a subject pool of the Kochi University of Technology.

650      Participants in this behavioral experiment completed the following three tasks;

651 1) SC-IAT which measures implicit attitude toward South Korea, 2) Explicit measure

652 of attitude (semantic differential) toward South Korea, and 3) the Sport Rivalry Fan

653 Perception Scale (SRFPS) (Havard et al. 2013). The SRFPS consists of four subscales

654    measuring Outgroup Competition against others (Indirect) (OIC; e.g., "I want my

655    favorite team's rival to win all games except when they play my favorite team"),

656    Outgroup Academic Prestige (OAP; e.g., "The academic prestige of my favorite

657    team's rival is poor"), Outgroup Sportsmanship (OS; e.g., "Fans of my favorite team's

658    rival demonstrate poor sportsmanship at games"), and Sense of Satisfaction when the

659    favorite team defeats the rival team in direct competition (SoS; e.g., "I feel a sense of

660    belonging when my favorite team beats my favorite team's rival') (3 items for each of

661    the four subscales). In each item, we replaced "my favorite team" with Japan or the

662    Japanese national team and "my favorite team's rival" with South Korea or the South

663    Korean national team. Furthermore, we removed the OAP subscale (all 3 items) and

664    the following 2 items from the OIC subscale ("I would support my favorite team's

665    rival in a championship game," and "I would support my favorite team's rival in out-

666    of-conference play") because they are specific to an American college sports context.

667    Accordingly, participants answered a total of 7 items, and three following scores were

668    computed for each participants; 1) OIC score (1 item), 2) OS score (average of 3

669    items; Cronbach's $\alpha = 0.81$), and 3) SoS score (average of 3 items; Cronbach's $\alpha =$

670    0.64).

671        Although we found a significant positive correlation between the South Korea

672    SC-IAT scores and the OIC score at a $p < 0.05$ level ($r(47) = 0.30$, $p = 0.017$; one-

673    tailed, no correction for multiple comparison), the same OIC scores were more

674    strongly related to explicit attitudes toward South Korea ($r(47) = 0.52$, $p < 0.001$).

675    The other two subscales were not related to the SC-IAT scores ($rs < 0.03$), while the

676    SoS subscale was related to explicit attitudes toward South Korea ($r(47) = 0.37$, $p =$

677    0.005; but not the OS subscale $r(47) = 0.01$). Thus, our results showed that in general,

678    the more positive explicit attitude a Japanese individual has toward South Korea, the

679    higher the sense of rivalry they have toward South Korea. These results indicate that it

680    is highly unlikely that the link between implicit prejudice toward South Korea (i.e.,

681    the SC-IAT scores) and neural signals in the left amygdala we found (Figure 3) can be

682    explained by the sense of rivalry toward South Korea.

683          Although our additional data indicate that a sense of rivalry is an unlikely

684    explanation, one may still argue that there was an important difference in

685    experimental procedures between the original fMRI study and this behavioral study.

686    While participants in this behavioral study (n = 49) completed only behavioral

687    measures (e.g., the SRFPS scale, SC-IAT, and semantic differential scales),

688    participants in the original fMRI study (n= 70) had been exposed to the pictures

689    depicting South Korea (some of which depicted sports related scenes) inside an fMRI

690    scanner before they completed the behavioral measures. Thus, it might be possible

691    that feelings of competitiveness evoked by viewing these pictures might have affected

692    both their SC-IAT scores as well as the amygdala activation. However, it should be

693    stressed that given the stronger link between a sense of rivalry and explicit

694    evaluations found in the behavioral study, if a sense of rivalry evoked by the pictures

695    were a major factor, we should have found a stronger correlation between the

696    amygdala activations and *explicit* evaluations.

697

698    **Discussion**

699          The present study showed that using MVPA, neural activation patterns in the

700    left amygdala could predict Japanese participants' level of implicit negative

701    evaluations of South Korea. With the much larger sample size (n = 70) than the

702    previous studies (Brosch et al., 2013; Cunningham et al., 2004; Phelps et al., 2000),

703    the present study provides reliable evidence that the left amygdala plays a key role in

704    representing negative implicit evaluations of an ethnic outgroup. Our results also

705    suggest that despite that ethnic prejudice is more variable across different cultures

706    compared to stereotypes of sex and age (Fiske, 2017), the link between the amygdala

707    and implicit negative evaluations of an ethnic or racial outgroup might be

708    generalizable across different intergroup contexts beyond the intergroup relation

709    between White vs. Black Americans. Even more generally, given the well-known role

710    of the amygdala in fear learning (learning of the associations between initially neutral

711    stimuli and aversive events) (Fendt and Fanselow, 1999; Pape and Pare, 2010), it

712    seems likely that the amygdala plays a key role in implicit attitudes toward not only

713    social groups but also non-social objects, although this idea should be formally tested

714    in future research. Implicit measures of attitudes such as IAT are thought to measure

715    the strength of automatically activated evaluative associations (e.g., the association

716    between an outgroup and negatively-valence words) stored in memory (i.e., the brain)

717    (Greenwald et al., 1998). Thus, the present results suggest that the associations

718    between South Korea and negativity possessed by Japanese individuals are stored in

719    the left amygdala. In contrast, explicit evaluations of South Korea were not robustly

720    related to neural signals. Thus, neural activation automatically evoked by the passive

721    viewing of South Korea-related images predicted implicit (automatic) evaluations of

722    South Korea, but not explicit (controlled) evaluations.

723        The result further showed that activations in the medial part of the amygdala

724    contribute positively to the prediction of implicit evaluations, whereas those in the

725    lateral part contribute negatively to the prediction (Figure 4a). Similarly, the mass-

726    univariate correlational analyses revealed that univariate signals were more positively

727    related to implicit evaluations in the more medial and anterior part of the left

728    amygdala (Figure 4b). This medial-lateral distinction is largely consistent with the

729  anatomical organization of the amygdala (the medial part = centromedial nuclei of the

730  amygdala, the lateral part = basolateral nuclei of the amygdala) (Sah et al., 2003). The

731  results may suggest that medial vs. lateral regions of the amygdala play different roles

732  in implicit evaluations, and this is consistent with past research showing the

733  functional distinction between basolateral nuclei and central nuclei of the amygdala

734  (Balleine and Killcross, 2006). As the lateral nuclei of the amygdala play a key role in

735  the formation of memories during fear conditioning (Rodrigues et al., 2004), our

736  results may suggest that neural signals in the centromedial nuclei of the amygdala

737  reflecting other factors such as negative affect or motivational salience of the stimuli

738  (some of which are reflected in the IAT scores) might contribute to the prediction of

739  implicit evaluation.

740       While neural signals in the left amygdala significantly predicted implicit

741  evaluations of South Korea, the right amygdala was not associated with implicit

742  evaluations. The findings from past neuroimaging studies on racial prejudice were

743  inconsistent as to the lateralization of the amygdala responses (see Chekroud et al.

744  (2014) for review). While some studies reported activations in bilateral amygdala in

745  response to black faces vs. white faces (e.g., (Hart et al., 2000); Phelps et al. (2000)),

746  other studies observed activations only in left (e.g., (Wheeler and Fiske, 2005)) or

747  right amygdala (e.g., (Cunningham et al., 2004; McCutcheon et al., 2018)).

748  Furthermore, Phelps et al. (2000) found that activations in the amygdala in both

749  hemispheres were correlated with IAT scores, while activations only in the left

750  amygdala were correlated with the startle eyeblink potentiation bias (a physiological

751  measure of indirect racial bias). In contrast, both Cunningham et al. (2004) and

752  Brosch et al. (2013) found the correlation between amygdala activities and IAT scores

753  only in the right amygdala. Past meta-analyses found that left amygdala activation is

754    more consistently observed than right amygdala activation during the processing of

755    affective stimuli (Baas et al., 2004) and that the left amygdala is more likely to be

756    activated when stimuli contain language whereas the right amygdala is more likely to

757    be activated when stimuli were masked to prevent conscious awareness (Costafreda et

758    al., 2008). Although these explanations could account for some of the past findings

759    (e.g., in Cunningham et al. (2004), faces were briefly presented [30 ms] and masked,

760    and they found activations only in the right amygdala), it is unlikely that differences

761    in the use of language or affective stimuli can explain the lateralization of amygdala

762    activations found in other previous studies as well as the present study. Future studies

763    should systematically manipulate these factors and test whether the left and right

764    amygdala play different roles in racial and ethnic prejudice.

765        In contrast to evaluations of South Korea (outgroup), we found that both

766    implicit and explicit evaluations of Japan (ingroup) were not robustly associated with

767    neural signals in any of the brain regions. There are at least two possible

768    interpretations for this dissociation between Japan and South Korea. One idea is that

769    as the amygdala plays a key role in fear learning, neural signals in the left amygdala

770    are related to negative implicit evaluations, but not neutral evaluations. Our

771    behavioral data showed that participants' implicit evaluations of Japan were not

772    clearly negative or positive, while their implicit evaluations of South Korea were

773    largely negative (Figure 2b). Note that it might have been possible to decode

774    individuals' level of implicit evaluations of Japan if their attitudes toward Japan were

775    clearly positive, as the amygdala has been implicated in processing positive as well as

776    negative values of stimuli (Murray, 2007). The other idea is that neural

777    representations evoked by Japan-related images are more complex than those evoked

778    by South Korea-images. While all South Korea-related images may be automatically

779     evaluated in a similar manner across Japanese individuals (i.e., outgroup homogeneity

780     effect), each Japan-related image is likely to evoke a variety of different

781     psychological and emotional reactions in each Japanese individual, which in turn

782     made across-subject decoding of attitudes toward Japan more difficult. Thus, it is

783     important to further investigate the role of the amygdala and other brain regions in

784     negative (and positive) implicit evaluations across a variety of different intergroup

785     contexts in future research.

786        While past social neuroscience studies reported the involvements of the anterior

787     cingulate cortex (ACC) and dorsolateral prefrontal cortex (DLPFC) in racial prejudice

788     (Amodio, 2014; Kubota et al., 2012), and one study (Richeson et al., 2003) reported

789     significant correlations between implicit evaluations and activities in these regions,

790     our results showed that both of these regions were not related to implicit and explicit

791     evaluations. In past studies, ACC and DLPFC activations in response to outgroup

792     faces were often interpreted as reflecting a conflict between automatic affective

793     responses and intention to respond fairly to outgroup faces and cognitive regulation of

794     evoked negative affective responses to the outgroup, respectively. Thus, although

795     speculative, our results might suggest that the level of implicit evaluations people

796     possess is not related to how much conflict they feel and how much they try to

797     suppress their prejudiced responses.

798        Our findings also showed that despite Japanese participants' clear negative

799     implicit as well as explicit evaluations of South Korea (Figure 2), conventional

800     univariate fMRI data analyses failed to find differences in activations in any of the

801     prejudice network between the South Korea vs. Japan blocks. Thus, our results are in

802     line with previous studies that did not find such amygdala activations (Brosch et al.,

803     2013; Cassidy and Krendl, 2016; Gilbert et al., 2012; Golby et al., 2001; Li et al.,

804    2016; Mattan et al., 2018; Phelps et al., 2000; Richeson et al., 2003; Stanley et al.,

805    2012; Terbeck et al., 2015).  However, it should be noted that the experimental design

806    of the present study was optimized for individual difference analyses (i.e., across-

807    subject correlation). For example, we fixed the block order for all participants (so that

808    the order effect, if there is any, should affect all participants in a similar manner). In

809    addition, we did not match lower visual features of the stimuli across the two

810    conditions. These differences might explain the lack of significant activations in the

811    South Korea vs. Japan contrast. Nonetheless, the order effect and the differences in

812    visual features of the stimuli are unlikely to explain the lack of a significant across-

813    subject correlation between univariate activations in the amygdala and implicit

814    evaluations. Contrary to the three small studies (Brosch et al., 2013; Cunningham et

815    al., 2004; Phelps et al., 2000), our univariate analysis revealed no correlation between

816    amygdala activities and implicit evaluations of South Korea, which is in agreement

817    with the six studies with larger sample sizes (Cassidy and Krendl, 2016; Cassidy et al.,

818    2016; Gilbert et al., 2012; Li et al., 2016; Richeson et al., 2003; Terbeck et al., 2015).

819        Thus, in contrast to MVPA, conventional univariate fMRI data analyses may

820    not be sensitive enough to detect the neural signatures of implicit evaluations, and it

821    may be a reason for the inconsistent findings in previous research (see (Amodio,

822    2014; Chekroud et al., 2014; Kubota et al., 2012). Thus, the present study

823    demonstrates the utility of MVPA, which may be able to refine the past findings and

824    provide important insights into the role played by each region within the prejudice

825    network in future research. For example, previous studies demonstrated, using

826    univariate fMRI data analysis, that amygdala activity in response to a racial outgroup

827    is modulated by perceiver's goals (Lieberman et al., 2005; Van Bavel et al., 2008;

828    Wheeler and Fiske, 2005), and one study found no difference in the amygdala activity

829    when participants performed a simple dot detection task while ingroup or outgroup

830    images were presented on the screen (Wheeler and Fiske, 2005). Since MVPA is

831    capable of detecting differences between conditions even when there is no overall

832    difference in the average amplitude of fMRI signals (e.g., (Harrison and Tong, 2009;

833    Kohler et al., 2013), it is interesting to test in future research whether MVPA can

834    decode implicit and/or explicit evaluations of an outgroup regardless of perceiver's

835    goals, which can provide an important insight into the automaticity of stereotyping

836    and prejudice (Bargh, 1999).

837        While the present study provides evidence for the link between the amygdala

838    and implicit prejudice in the intergroup context of Japan vs. South Korea, it is

839    important to investigate the same link in other intergroup contexts including White vs.

840    Black Americans in future research (i.e., using MVPA and with a larger sample size).

841    Although it is conceivable to think that the amygdala plays a major role in implicit

842    prejudice (i.e., association between an outgroup and negativity) in all intergroup

843    contexts given its well-known role in fear learning (Fendt and Fanselow, 1999; Pape

844    and Pare, 2010), other brain regions may play additional roles in a different intergroup

845    context. Given that prejudice is a world-wide problem (Landis and Albert, 2012; Noor

846    and Montiel, 2009), future research should investigate similarities and differences in

847    neural signatures of implicit prejudice across different intergroup contexts to have a

848    comprehensive understanding of its neural mechanisms.

849        It is also important to use various stimuli in future research (e.g., pictures

850    unrelated to sports scenes). Although our additional behavioral study found only a

851    very weak link between a sense of rivalry or competitiveness and implicit evaluations

852    toward the outgroup, there was an important difference in experimental procedures

853    between the original fMRI study and this behavioral study. While participants in this

854    additional behavioral study (n = 49) completed only behavioral measures (e.g., the

855    SRFPS scale, SC-IAT, and semantic differential scales), participants in the original

856    fMRI study (n= 70) had been exposed to the pictures depicting South Korea (some of

857    which depicted sports related scenes) inside an fMRI scanner before they completed

858    the behavioral measures. Thus, it is possible that feelings of competitiveness evoked

859    by viewing these pictures might have affected their SC-IAT scores as well as the

860    amygdala activation. Nonetheless, it should be stressed that given the stronger link

861    between a sense of rivalry and *explicit* evaluations, it is highly unlikely that our main

862    amygdala findings (Figure 3) can be explained by the sense of rivalry.

863        Finally, an important implication of the present finding is that neural signals in

864    the amygdala could be used as an independent neural index of implicit attitudes

865    toward an outgroup. Past social psychological studies have reported that a variety of

866    simple behavioral interventions could reduce implicit attitudes toward an outgroup as

867    measured by IAT (Lai et al., 2014). However, it has been debated whether such

868    interventions actually reduced implicit attitudes or just IAT scores (i.e., implicit

869    attitudes remain unchanged) (Han et al., 2010). An independent neural index has a

870    potential to provide a unique insight into this debate, and thus, the present finding is

871    the important step toward formulating effective interventions to regulate and reduce

872    various prejudice in societies.

873

881

882          **References**

883    Amodio, D.M., 2014. The neuroscience of prejudice and stereotyping. Nature
884    Reviews Neuroscience 15, 670-682.
885    Baas, D., Aleman, A., Kahn, R.S., 2004. Lateralization of amygdala activation: a
886    systematic review of functional neuroimaging studies. Brain Research Reviews 45,
887    96-103.
888    Balleine, B.W., Killcross, S., 2006. Parallel incentive processing: an integrated view
889    of amygdala function. Trends in Neurosciences 29, 272-279.
890    Bargh, J.A., 1999. The cognitive monster: The case against controllability of
891    automatic stereotype effects. In: Chaiken, S., Trope, Y. (Eds.), Dual-process theories
892    in social psychology. Guilford Press, New York, pp. 361-382.
893    BBC_World_Service, 2010. Global views of United States improve while other
894    countries decline.
895    BBC_World_Service, 2014. Negative views of Russia on the rise: Global poll.
896    Benjamini, Y., Hochberg, Y., 1995. Controlling the False Discovery Rate - a Practical
897    and Powerful Approach to Multiple Testing. Journal of the Royal Statistical Society
898    Series B-Methodological 57, 289-300.
899    Brosch, T., Bar-David, E., Phelps, E.A., 2013. Implicit race bias decreases the
900    similarity of neural representations of black and white faces. Psychol Sci 24, 160-166.
901    Bruneau, E.G., Saxe, R., 2010. Attitudes towards the outgroup are predicted by
902    activity in the precuneus in Arabs and Israelis. Neuroimage 52, 1704-1711.
903    Cassidy, B.S., Krendl, A.C., 2016. Dynamic neural mechanisms underlie race
904    disparities in social cognition. Neuroimage 132, 238-246.
905    Cassidy, B.S., Lee, E.J., Krendl, A.C., 2016. Age and executive ability impact the
906    neural correlates of race perception. Soc Cogn Affect Neurosci 11, 1752-1761.
907    Chang, L.J., Gianaros, P.J., Manuck, S.B., Krishnan, A., Wager, T.D., 2015. A
908    sensitive and specific neural signature for picture-induced negative affect. PLoS Biol
909    13, e1002180.
910    Chekroud, A.M., Everett, J.A., Bridge, H., Hewstone, M., 2014. A review of
911    neuroimaging studies of race-related prejudice: does amygdala response reflect
912    threat? Frontiers in Human Neuroscience 8, 179.
913    Cho, K., 2017. Quantitative Text Analysis of "Yahoo! News" : Focusing on
914    Comments on Koreans The Journal of Applied Sociology 59, 113 - 127.
915    Cloutier, J., Li, T., Correll, J., 2014. The impact of childhood experience on amygdala
916    response to perceptually familiar black and white faces. J Cogn Neurosci 26, 1992-
917    2004.
918    Cohen, J.R., Asarnow, R.F., Sabb, F.W., Bilder, R.M., Bookheimer, S.Y., Knowlton,
919    B.J., Poldrack, R.A., 2010. Decoding developmental differences and individual
920    variability in response inhibition through predictive analyses across individuals.
921    Frontiers in Human Neuroscience 4, 47.
922    Costafreda, S.G., Brammer, M.J., David, A.S., Fu, C.H.Y., 2008. Predictors of
923    amygdala activation during the processing of emotional stimuli: A meta-analysis of
924    385 PET and fMRI studies. Brain Research Reviews 58, 57-70.
925    Cunningham, W.A., Johnson, M.K., Raye, C.L., Chris Gatenby, J., Gore, J.C., Banaji,
926    M.R., 2004. Separable neural components in the processing of black and white faces.
927    Psychol Sci 15, 806-813.
928    Demos, K.E., Kelley, W.M., Ryan, S.L., Davis, F.C., Whalen, P.J., 2008. Human
929    amygdala sensitivity to the pupil size of others. Cereb Cortex 18, 2729-2734.
930    Drucker, H., Burges, C.J.C., Kaufman, L., Smola, A., Vapnik, V., 1997. Support
931    vector regression machines. Advances in Neural Information Processing Systems 9 9,

932     155-161.
933     Dunlap, W.P., Cortina, J.M., Vaslow, J.B., Burke, M.J., 1996. Meta-analysis of
934     experiments with matched groups or repeated measures designs. Psychological
935     Methods 1, 170-177.
936     Fendt, M., Fanselow, M.S., 1999. The neuroanatomical and neurochemical basis of
937     conditioned fear. Neuroscience and Biobehavioral Reviews 23, 743-760.
938     Fiske, S.T., 2017. Prejudices in Cultural Contexts: Shared Stereotypes (Gender, Age)
939     Versus Variable Stereotypes (Race, Ethnicity, Religion). Perspectives on
940     Psychological Science 12, 791-799.
941     Gilbert, S.J., Swencionis, J.K., Amodio, D.M., 2012. Evaluative vs. trait
942     representation in intergroup social judgments: distinct roles of anterior temporal lobe
943     and prefrontal cortex. Neuropsychologia 50, 3600-3611.
944     Golby, A.J., Gabrieli, J.D., Chiao, J.Y., Eberhardt, J.L., 2001. Differential responses
945     in the fusiform region to same-race and other-race faces. Nat Neurosci 4, 845-850.
946     Greenwald, A.G., McGhee, D.E., Schwartz, J.L.K., 1998. Measuring individual
947     differences in implicit cognition: The implicit association test. Journal of Personality
948     and Social Psychology 74, 1464-1480.
949     Greenwald, A.G., Nosek, B.A., Banaji, M.R., 2003. Understanding and using the
950     implicit association test: I. An improved scoring algorithm. Journal of Personality and
951     Social Psychology 85, 197-216.
952     Han, H.A., Czellar, S., Olson, M.A., Fazio, R.H., 2010. Malleability of attitudes or
953     malleability of the IAT? Journal of Experimental Social Psychology 46, 286-298.
954     Harrison, S.A., Tong, F., 2009. Decoding reveals the contents of visual working
955     memory in early visual areas. Nature 458, 632-635.
956     Hart, A.J., Whalen, P.J., Shin, L.M., McInerney, S.C., Fischer, H., Rauch, S.L., 2000.
957     Differential response in the human amygdala to racial outgroup vs ingroup face
958     stimuli. Neuroreport 11, 2351-2355.
959     Haufe, S., Meinecke, F., Gorgen, K., Dahne, S., Haynes, J.D., Blankertz, B.,
960     Biessgmann, F., 2014. On the interpretation of weight vectors of linear models in
961     multivariate neuroimaging. Neuroimage 87, 96-110.
962     Haynes, J.D., 2015. A Primer on Pattern-Based Approaches to fMRI: Principles,
963     Pitfalls, and Perspectives. Neuron 87, 257-270.
964     Hofmann, W., Gawronski, B., Gschwendner, T., Le, H., Schmitt, M., 2005. A meta-
965     analysis on the correlation between the implicit association test and explicit self-
966     report measures. Pers Soc Psychol Bull 31, 1369-1385.
967     Izuma, K., Kennedy, K., Fitzjohn, A., Sedikides, C., Shibata, K., 2018. Neural
968     activity in the reward-related brain regions predicts implicit self-esteem: A novel
969     validity test of psychological measures using neuroimaging. J Pers Soc Psychol 114,
970     343-357.
971     Izuma, K., Shibata, K., Matsumoto, K., Adolphs, R., 2017. Neural predictors of
972     evaluative attitudes towards celebrities. Social Cognitive and Affective Neuroscience
973     12, 382-390.
974     Jimura, K., Poldrack, R.A., 2012. Analyses of regional-average activation and
975     multivoxel pattern information tell complementary stories. Neuropsychologia 50,
976     544-552.
977     Karpinski, A., Steinman, R.B., 2006. The single category implicit association test as a
978     measure of implicit social cognition. J Pers Soc Psychol 91, 16-32.
979     Kohler, P.J., Fogelson, S.V., Reavis, E.A., Meng, M., Guntupalli, J.S., Hanke, M.,
980     Halchenko, Y.O., Connolly, A.C., Haxby, J.V., Tse, P.U., 2013. Pattern classification
981     precedes region-average hemodynamic response in early visual cortex. Neuroimage

982    78, 249-260.
983    Kubota, J.T., Banaji, M.R., Phelps, E.A., 2012. The neuroscience of race. Nat
984    Neurosci 15, 940-948.
985    Lai, C.K., Marini, M., Lehr, S.A., Cerruti, C., Shin, J.E.L., Joy-Gaba, J.A., Ho, A.K.,
986    Teachman, B.A., Wojcik, S.P., Koleva, S.P., Frazier, R.S., Heiphetz, L., Chen, E.E.,
987    Turner, R.N., Haidt, J., Kesebir, S., Hawkins, C.B., Schaefer, H.S., Rubichi, S.,
988    Sartori, G., Dial, C.M., Sriram, N., Banaji, M.R., Nosek, B.A., 2014. Reducing
989    Implicit Racial Preferences: I. A Comparative Investigation of 17 Interventions.
990    Journal of Experimental Psychology-General 143, 1765-1785.
991    Landis, D., Albert, R.D. (Eds.), 2012. Handbook of Ethnic Conflict: International
992    Perspectives. Springer.
993    Li, T., Cardenas-Iniguez, C., Correll, J., Cloutier, J., 2016. The impact of motivation
994    on race-based impression formation. Neuroimage 124, 1-7.
995    Lieberman, M.D., Hariri, A., Jarcho, J.M., Eisenberger, N.I., Bookheimer, S.Y., 2005.
996    An fMRI investigation of race-related amygdala activity in African-American and
997    Caucasian-American individuals. Nat Neurosci 8, 720-722.
998    Loken, E., Gelman, A., 2017. Measurement error and the replication crisis. Science
999    355, 584-585.
1000   Maldjian, J.A., Laurienti, P.J., Kraft, R.A., Burdette, J.H., 2003. An automated
1001   method for neuroanatomic and cytoarchitectonic atlas-based interrogation of fMRI
1002   data sets. Neuroimage 19, 1233-1239.
1003   Mattan, B.D., Kubota, J.T., Dang, T.P., Cloutier, J., 2018. External motivation to
1004   avoid prejudice alters neural responses to targets varying in race and status. Soc Cogn
1005   Affect Neurosci 13, 22-31.
1006   McCutcheon, R., Bloomfield, M.A.P., Dahoun, T., Quinlan, M., Terbeck, S., Mehta,
1007   M., Howes, O., 2018. Amygdala reactivity in ethnic minorities and its relationship to
1008   the social environment: an fMRI study. Psychol Med, 1-8.
1009   Murray, E.A., 2007. The amygdala, reward and emotion. Trends Cogn Sci 11, 489-
1010   497.
1011   Noor, N.M., Montiel, C.J. (Eds.), 2009. Peace Psychology in Asia. Springer.
1012   Norman, K.A., Polyn, S.M., Detre, G.J., Haxby, J.V., 2006. Beyond mind-reading:
1013   multi-voxel pattern analysis of fMRI data. Trends Cogn Sci 10, 424-430.
1014   Nosek, B.A., Greenwald, A.G., Banaji, M., 2007. The Implicit Association Test at
1015   Age 7: A Methodological and Conceptual Review. In: Bargh, J.A. (Ed.), Automatic
1016   processes in social thinking and behavior. Psychology Press.
1017   Op de Beeck, H.P., 2010. Against hyperacuity in brain reading: Spatial smoothing
1018   does not hurt multivariate fMRI analyses? Neuroimage 49, 1943-1948.
1019   Oswald, F.L., Mitchell, G., Blanton, H., Jaccard, J., Tetlock, P.E., 2013. Predicting
1020   ethnic and racial discrimination: a meta-analysis of IAT criterion studies. J Pers Soc
1021   Psychol 105, 171-192.
1022   Pape, H.C., Pare, D., 2010. Plastic synaptic networks of the amygdala for the
1023   acquisition, expression, and extinction of conditioned fear. Physiol Rev 90, 419-463.
1024   Phelps, E.A., O'Connor, K.J., Cunningham, W.A., Funayama, E.S., Gatenby, J.C.,
1025   Gore, J.C., Banaji, M.R., 2000. Performance on indirect measures of race evaluation
1026   predicts amygdala activation. J Cogn Neurosci 12, 729-738.
1027   Richeson, J.A., Baird, A.A., Gordon, H.L., Heatherton, T.F., Wyland, C.L., Trawalter,
1028   S., Shelton, J.N., 2003. An fMRI investigation of the impact of interracial contact on
1029   executive function. Nat Neurosci 6, 1323-1328.
1030   Richeson, J.A., Todd, A.R., Trawalter, S., Baird, A.A., 2008. Eye-gaze direction
1031   modulates race-related amygdala activity. Group Processes & Intergroup Relations 11,

1032 233-246.
1033 Rodrigues, S.M., Schafe, G.E., LeDoux, J.E., 2004. Molecular mechanisms
1034 underlying emotional learning and memory in the lateral amygdala. Neuron 44, 75-91.
1035 Ronquillo, J., Denson, T.F., Lickel, B., Lu, Z.L., Nandy, A., Maddox, K.B., 2007. The
1036 effects of skin tone on race-related amygdala activity: an fMRI investigation. Soc
1037 Cogn Affect Neurosci 2, 39-44.
1038 Sah, P., Faber, E.S., Lopez De Armentia, M., Power, J., 2003. The amygdaloid
1039 complex: anatomy and physiology. Physiol Rev 83, 803-834.
1040 Said, C.P., Baron, S.G., Todorov, A., 2009. Nonlinear Amygdala Response to Face
1041 Trustworthiness: Contributions of High and Low Spatial Frequency Information.
1042 Journal of Cognitive Neuroscience 21, 519-528.
1043 Sapountzis, P., Schluppeck, D., Bowtell, R., Peirce, J.W., 2010. A comparison of
1044 fMRI adaptation and multivariate pattern classification analysis in visual cortex.
1045 Neuroimage 49, 1632-1640.
1046 Stanley, D.A., Sokol-Hessner, P., Banaji, M.R., Phelps, E.A., 2011. Implicit race
1047 attitudes predict trustworthiness judgments and economic trust decisions. Proceedings
1048 of the National Academy of Sciences of the United States of America 108, 7710-7715.
1049 Stanley, D.A., Sokol-Hessner, P., Fareri, D.S., Perino, M.T., Delgado, M.R., Banaji,
1050 M.R., Phelps, E.A., 2012. Race and reputation: perceived racial group trustworthiness
1051 influences the neural correlates of trust decisions. Philos Trans R Soc Lond B Biol Sci
1052 367, 744-753.
1053 Telzer, E.H., Flannery, J., Shapiro, M., Humphreys, K.L., Goff, B., Gabard-Durman,
1054 L., Gee, D.D., Tottenham, N., 2013a. Early experience shapes amygdala sensitivity to
1055 race: an international adoption design. J Neurosci 33, 13484-13488.
1056 Telzer, E.H., Humphreys, K.L., Shapiro, M., Tottenham, N., 2013b. Amygdala
1057 sensitivity to race is not present in childhood but emerges over adolescence. J Cogn
1058 Neurosci 25, 234-244.
1059 Terbeck, S., Kahane, G., McTavish, S., McCutcheon, R., Hewstone, M., Savulescu, J.,
1060 Chesterman, L.P., Cowen, P.J., Norbury, R., 2015. beta-Adrenoceptor blockade
1061 modulates fusiform gyrus activity to black versus white faces. Psychopharmacology
1062 (Berl) 232, 2951-2958.
1063 Todorov, A., Engell, A.D., 2008. The role of the amygdala in implicit evaluation of
1064 emotionally neutral faces. Social Cognitive and Affective Neuroscience 3, 303-312.
1065 Van Bavel, J.J., Packer, D.J., Cunningham, W.A., 2008. The neural substrates of in-
1066 group bias: a functional magnetic resonance imaging investigation. Psychol Sci 19,
1067 1131-1139.
1068 Wheeler, M.E., Fiske, S.T., 2005. Controlling racial prejudice: social-cognitive goals
1069 affect amygdala and stereotype activation. Psychol Sci 16, 56-63.
1070 Winston, J.S., Strange, B.A., O'Doherty, J., Dolan, R.J., 2002. Automatic and
1071 intentional brain responses during evaluation of trustworthiness of faces. Nat
1072 Neurosci 5, 277-283.
1073 Yarkoni, T., 2009. Big Correlations in Little Studies: Inflated fMRI Correlations
1074 Reflect Low Statistical Power-Commentary on Vul et al. (2009). Perspectives on
1075 Psychological Science 4, 294-298.
1076

1077 **Table 1. Correlations across behavioral measures**

| Measure | 1 | 2 | 3 | 4 | 5 | 6 | 7 |
|---|---|---|---|---|---|---|---|

| | | | | | | | |
|---|---|---|---|---|---|---|---|
| 1. Japan SC-IAT | 1 | | | | | | |
| 2. South Korea SC-IAT | 0.27* | 1 | | | | | |
| 3. Disparity in SC-IAT scores (Japan - South Korea) | N/A | N/A | 1 | | | | |
| 4. Japan Semantic differential | 0.12 | 0.12 | 0.00 | 1 | | | |
| 5. South Korea Semantic differential | -0.06 | -0.02 | -0.03 | 0.26* | 1 | | |
| 6. Disparity in semantic differential ratings (Japan - South Korea) | 0.14 | 0.12 | 0.02 | N/A | N/A | 1 | |
| 7. Trust Game score† | 0.02 | -0.07 | 0.07 | 0.03 | -0.08 | 0.10 | 1 |

1078  SC-IAT: Single-Category Implicit Association Test. Note that the Trust Game score
1079  is the difference between the average amount offered to Japanese partners minus
1080  South Korean partners so that higher numbers indicate that they are more trusting
1081  toward Japanese relative to South Korean individuals. N/A: The four correlations are
1082  omitted from the table because they are correlations between two non-independent
1083  variables (and not surprisingly, they are all highly correlated at $p < 0.001$ level). * $p <$
1084  0.05. † One outlier (more than 3 SD from the mean) was excluded when analyzing the
1085  Trust Game data.
1086

1087 **Table 2: Decoding performance in the amygdala**

| ROI name | The number of voxels | Decoding performance (r): | | | |
|---|---|---|---|---|---|
| | | South Korea | | Japan | |
| | | Implicit attitude | Explicit attitude | Implicit attitude | Explicit attitude |
| L amygdala | 54 | 0.31* (0.09, 0.51) | -0.13 (-0.37, 0.11) | -0.03 (-0.27, 0.20) | -0.24 (-0.45, -0.00) |
| R amygdala | 63 | -0.04 (-0.27, 0.20) | 0.11 (-0.12, 0.34) | 0.07 (-0.17, 0.30) | -0.30 (-0.50, -0.07) |

1088 Implicit attitudes were measured by SC-IAT, while explicit attitude was measured by
1089 the semantic differential. Amygdala masks were taken from the Anatomical
1090 Automatic Labeling (AAL) masks implemented in the WFU pickatlas toolbox. Voxel
1091 size = 3 × 3 × 3 mm. * $p_{perm}$ < 0.05 (*p*-value based on permutation test [5,000 times]).
1092 Numbers in parentheses are 95% confidence interval.

**Table 3: Decoding of implicit and explicit attitudes toward each of South Korea and Japan in each of a total of 79 regions outside of the prejudice network.**

| | Name of mask in the WFU pickatlas toolbox | The number of voxels | Decoding performance (r): | | | |
| | | | South Korea | | Japan | |
| | | | Implicit attitude | Explicit attitude | Implicit attitude | Explicit attitude |
|---|---|---|---|---|---|---|
| | *vmPFC* | 364 | -0.46 | 0.02 | 0.08 | -0.35 |
| | *mPFC* | 1683 | -0.17 | -0.25 | 0.01 | 0.03 |
| | Rectus | 314 | -0.20 | 0.16 | 0.10 | 0.04 |
| | Frontal_Mid_L | 1277 | 0.30* | 0.05 | -0.03 | -0.04 |
| | Frontal_Mid_R | 1379 | -0.03 | -0.27 | 0.09 | -0.07 |
| | Frontal_Inf_Oper_L | 274 | 0.09 | 0.19 | -0.15 | 0.11 |
| | Frontal_Inf_Oper_R | 367 | -0.01 | -0.14 | 0.22 | -0.06 |
| | Frontal_Inf_Orb_L | 358 | 0.19 | -0.01 | 0.09 | 0.14 |
| | Frontal_Inf_Orb_R | 351 | -0.05 | -0.26 | -0.04 | 0.12 |
| Frontal lobe | Frontal_Inf_Tri_L | 608 | -0.19 | 0.17 | 0.03 | 0.01 |
| | Frontal_Inf_Tri_R | 475 | -0.20 | 0.19 | -0.07 | 0.02 |
| | Frontal_Sup_L | 930 | 0.08 | 0.07 | -0.21 | 0.21 |
| | Frontal_Sup_R | 1067 | 0.15 | -0.08 | -0.01 | 0.28* |
| | Frontal_Sup_Orb_L | 58 | 0.21 | -0.22 | 0.08 | -0.11 |
| | Frontal_Sup_Orb_R | 51 | 0.00 | 0.20 | 0.02 | -0.14 |
| | Olfactory_L | 42 | 0.16 | 0.01 | 0.05 | -0.29 |
| | Olfactory_R | 49 | 0.08 | 0.13 | 0.21 | 0.14 |
| | Precentral_L | 874 | -0.26 | 0.01 | -0.23 | 0.32* |
| | Precentral_R | 845 | -0.05 | -0.17 | -0.15 | -0.11 |
| | Supp_Motor_Area | 1617 | -0.05 | -0.28 | -0.14 | 0.10 |
| | Angular_L | 329 | 0.06 | 0.09 | 0.02 | 0.11 |
| | Angular_R | 422 | 0.06 | 0.15 | -0.12 | 0.48* |
| | Parietal_Inf_L | 669 | -0.08 | -0.06 | 0.11 | 0.21 |
| | Parietal_Inf_R | 384 | -0.04 | -0.41 | 0.04 | 0.14 |
| | Parietal_Sup_L | 499 | 0.02 | -0.15 | -0.09 | 0.25 |
| | Parietal_Sup_R | 426 | 0.01 | -0.19 | -0.20 | 0.06 |
| Parietal lobe | Postcentral_L | 979 | -0.16 | -0.01 | 0.03 | 0.06 |
| | Postcentral_R | 898 | 0.08 | 0.09 | -0.22 | 0.00 |
| | Rolandic_Oper_L | 271 | 0.00 | 0.09 | -0.22 | 0.41* |
| | Rolandic_Oper_R | 358 | 0.23 | 0.15 | 0.24 | 0.00 |
| | SupraMarginal_L | 317 | -0.13 | 0.14 | 0.26 | 0.15 |
| | SupraMarginal_R | 471 | -0.09 | 0.23 | -0.12 | 0.09 |
| | Paracentral_Lobule | 726 | 0.00 | 0.06 | 0.01 | 0.03 |
| | Precuneus | 2291 | 0.07 | 0.03 | -0.14 | -0.06 |
| Temporal lobe | Fusiform_L | 445 | 0.12 | -0.28 | 0.27 | -0.08 |
| | Fusiform_R | 424 | 0.31* | -0.21 | -0.14 | 0.01 |

|  |  |  |  |  |  |
|---|---|---|---|---|---|
| | Heschl_L | 69 | -0.28 | 0.02 | 0.02 | 0.17 |
| | Heschl_R | 70 | 0.02 | -0.12 | 0.15 | 0.10 |
| | Temporal_Inf_L | 596 | -0.02 | -0.01 | 0.21 | 0.02 |
| | Temporal_Inf_R | 414 | 0.14 | -0.20 | -0.02 | -0.24 |
| | Temporal_Mid_L | 1121 | 0.00 | -0.16 | -0.05 | 0.06 |
| | Temporal_Mid_R | 964 | -0.34 | 0.13 | -0.32 | 0.11 |
| | Temporal_Pole_Mid_ L | 89 | -0.05 | 0.08 | -0.18 | 0.01 |
| | Temporal_Pole_Mid_ R | 138 | -0.09 | 0.10 | 0.13 | -0.15 |
| | Temporal_Pole_Sup_ L | 197 | 0.08 | -0.14 | -0.24 | 0.12 |
| | Temporal_Pole_Sup_ R | 160 | 0.12 | 0.12 | 0.08 | 0.14 |
| | Temporal_Sup_L | 560 | -0.12 | 0.01 | -0.05 | 0.19 |
| | Temporal_Sup_R | 730 | -0.03 | -0.29 | -0.07 | 0.19 |
| Occipital lobe | Calcarine_L | 494 | 0.08 | -0.07 | 0.18 | -0.16 |
| | Calcarine_R | 504 | -0.04 | 0.26 | -0.21 | -0.16 |
| | Cuneus_L | 365 | 0.14 | -0.24 | 0.09 | 0.04 |
| | Cuneus_R | 402 | 0.07 | -0.16 | -0.46 | -0.09 |
| | Lingual_L | 509 | -0.12 | 0.07 | -0.01 | -0.22 |
| | Lingual_R | 453 | 0.30* | -0.11 | -0.23 | -0.05 |
| | Occipital_Inf_L | 211 | 0.15 | -0.21 | -0.07 | -0.23 |
| | Occipital_Inf_R | 204 | -0.09 | 0.05 | 0.01 | -0.08 |
| | Occipital_Mid_L | 916 | 0.15 | -0.11 | -0.07 | 0.13 |
| | Occipital_Mid_R | 498 | 0.06 | -0.33 | 0.17 | 0.10 |
| | Occipital_Sup_L | 372 | -0.02 | -0.08 | 0.10 | 0.18 |
| | Occipital_Sup_R | 359 | 0.25 | -0.05 | -0.28 | 0.15 |
| Insular lobe | *Ant_Insula_L* | 391 | 0.02 | -0.06 | -0.01 | 0.12 |
| | *Ant_Insula_R* | 353 | 0.28* | -0.03 | 0.05 | -0.03 |
| | Post_Insula_L | 210 | -0.47 | -0.02 | -0.17 | 0.26* |
| | Post_Insula_R | 212 | 0.00 | 0.03 | -0.07 | 0.04 |
| Limbic lobe /Subcortical structures | *L Caudate nucleus* | 270 | 0.26 | 0.23 | -0.15 | 0.06 |
| | *R Caudate nucleus* | 283 | 0.08 | 0.13 | -0.03 | 0.15 |
| | Cingulum_Ant | 1151 | -0.15 | 0.08 | 0.16 | -0.05 |
| | Cingulum_Mid | 1540 | -0.01 | -0.22 | 0.07 | -0.05 |
| | Cingulum_Post | 309 | -0.12 | -0.24 | 0.06 | -0.17 |
| | Hippocampus_L | 240 | 0.09 | 0.14 | -0.05 | -0.01 |
| | Hippocampus_R | 251 | 0.02 | 0.10 | 0.08 | -0.05 |
| | ParaHippocampal_L | 190 | -0.12 | -0.11 | -0.04 | -0.02 |
| | ParaHippocampal_R | 248 | 0.14 | -0.19 | 0.13 | 0.10 |
| | Pallidum_L | 67 | -0.04 | -0.05 | -0.16 | -0.04 |
| | Pallidum_R | 64 | 0.07 | -0.14 | -0.02 | 0.06 |

| | | | | | |
|---|---|---|---|---|---|
| Putamen_L | 317 | -0.16 | -0.01 | -0.05 | -0.05 |
| Putamen_R | 321 | -0.01 | 0.02 | -0.15 | -0.04 |
| Thalamus_L | 313 | 0.13 | -0.18 | -0.05 | -0.29 |
| Thalamus_R | 291 | 0.03 | 0.23 | -0.01 | -0.13 |

1095 Implicit attitudes were measured by SC-IAT, while explicit attitude was measured by
1096 the semantic differential. All masks were taken from the Anatomical Automatic
1097 Labeling (AAL) masks implemented in the WFU pickatlas toolbox. For each of the
1098 midline regions (e.g., vmPFC, mPFC, Rectus, Supp_Motor_Area, Paracentral_Lobule,
1099 Precuneus, etc.), mask images in both hemispheres are combined to create a single
1100 mask image. Voxel size = $3 \times 3 \times 3$ mm. * $p < 0.05$ (based on permutation test [5,000
1101 times], uncorrected for multiple comparisons). *Italics* indicates regions included in
1102 the prejudice network reported previously (Amodio, 2014). Note that since the
1103 anterior part of the insula, rather than its posterior part, has been more frequently
1104 implicated in prejudice (Amodio, 2014), we separated each insula mask into anterior
1105 (y coordinate $\geq$ 0) and posterior (y coordinate < 0) parts.