1  **Journal:**

2

3

4

5

6

7

8

9

10  **Title: fMRI repetition suppression reveals no sensitivity to trait judgments from faces**

11  **along the ventral visual stream or in the theory-of-mind network**

12

13

14

15  **Emily E. Butler, Rob Ward, Paul E. Downing & Richard Ramsey**

16

17

18  Wales Institute for Cognitive Neuroscience, School of Psychology, Bangor University,

19  Bangor, Gwynedd, Wales, LL57 2AS, United Kingdom

20  Corresponding author: r.ramsey@bangor.ac.uk

21

23

24 **Abstract**

25 The human face cues a wealth of social information, but the neural mechanisms that underpin

26 social attributions from faces are not well known. In the current fMRI experiment, we used

27 repetition suppression to test the hypothesis that populations of neurons in face perception and

28 theory-of-mind neural networks would show sensitivity to faces that cue distinct trait

29 judgments. Although faces were accurately discriminated based on associated traits, our

30 results showed no evidence that face or theory-of-mind networks showed repetition

31 suppression for face traits. Thus, we do not provide evidence for population coding models of

32 face perception that include sensitivity to high and low trait features. Due to aspects of the

33 experimental design, which bolstered statistical power and sensitivity, we have reasonable

34 confidence that we could detect effects of a moderate size, should they exist. The null

35 findings reported here, therefore, add value to models of neural organisation in social

36 perception by showing instances where effects are absent or small. To test the generalisability

37 of our findings, future work should test different types of trait judgment and different types of

38 facial stimuli, in order to further probe the neurobiological bases of impression formation

39 based on facial appearance.

40

**Introduction**

41

42    Faces signal information that guide social interactions (Emery, 2000). Although complex

43    social signals such as emotional states, trait characteristics, and attentional focus are readily

44    perceived from faces (Jack & Schyns, 2017; Todorov et al., 2015), the neural mechanisms

45    that process social dimensions of face perception remain unclear. Here, in a functional

46    magnetic resonance imaging (fMRI) experiment, we use repetition suppression to investigate

47    the neural representation of how trait inferences are arrived at during social perception.

48         The majority of neuroscience research on face perception has focused on detection and

49    recognition of identity and emotion. This research has identified face-selective patches of

50    cortex that respond more to viewing faces than other categories of objects such as houses and

51    cars (Duchaine & Yovel, 2015; Haxby et al., 2000; Kanwisher et al., 1997). Key regions in

52    the face perception network include the fusiform face area (FFA; Kanwisher et al., 1997),

53    occipital face area (OFA; Gauthier et al., 2000) and posterior superior temporal sulcus (pSTS;

54    Allison et al., 2000; Pitcher et al., 2011). These three nodes along the ventral visual stream are

55    suggested to perform core visual analyses of facial features, but also interact with extended

56    circuits in anterior cortex for more elaborate representations of identity and emotional valence

57    (Duchaine & Yovel, 2015; Haxby et al., 2000; Kanwisher, 2010).

58         Face recognition is important for initiating social interactions, but faces cue much

59    more than the mere presence of a social agent. Indeed, impressions of others are partly formed

60    on the basis of stable, non-emotional aspects of facial appearance (Todorov et al., 2015;

61    Zebrowitz, 2011). As such, there is interplay between the perception of facial features and the

62    formation of character impressions (Jack & Schyns, 2017). Models of social impressions from

63    faces have been developed that include dimensions of valence/trustworthiness, dominance

64    and attractiveness (Todorov et al., 2008; Sutherland et al., 2013; Wang et al., 2016). However,

65    there is currently little known regarding the neural bases of such impression formation. For

66    example, faces that cue social evaluations of trustworthiness and attractiveness have been

67    associated with responses in the amygdala and ventral striatum, which have been thought to

68    index the reward value and typicality of faces (Bzdok et al., 2011; Mende-Siedlecki et al.,

69    2013; Said et al., 2010; 2011; Todorov et al., 2013). Additionally, behavioural research has

70    shown that personality characteristics such as extraversion are accurately perceived from

71    static facial features (Borkenau & Liebler, 1992; Borkenau et al., 2009; Kramer & Ward,

72    2010; Penton-Voak et al., 2006). However, beyond brain circuits associated with reward, little

73    is currently known regarding the neural architecture supporting personality judgments that are

74    cued during face perception.

75         Research investigating trait judgments has primarily focused on reading statements,

76    rather than faces (Uleman et al., 2008). For example, reading trait-diagnostic statements, such

77    as "she gave money to charity", engages the theory-of-mind (ToM) network more than trait-

78    neutral statements such as "she sharpened her pencil" (Heleven et al., 2017; Heleven & Van

79    Overwalle, 2016; Ma et al., 2014; Mitchell et al., 2002; 2005; Van Overwalle et al., 2016).

80    The ToM network is engaged when attributing mental states such as beliefs, desires and

81    attitudes to others, as well as judging character and is thought to be central to understanding

82    social cognition (van Overwalle, 2009; Frith & Frith, 1999). The ToM network is largely

83    distinct from the face perception network with key nodes covering temporoparietal junction,

84    medial prefrontal cortex, temporal poles and precuneus (van Overwalle, 2009; Frith & Frith,

85    1999; Saxe & Kanwisher, 2003). However, the potential role of the ToM network in forming

86    impressions based on facial appearance has not been studied in depth. As such, the cognitive

87    and neural systems that identify perceptual features and link them to trait judgments are not

88    well known (Over & Cook, 2018). The current study, therefore, investigates the hypothesis

89    that impression formation from faces relies on a distributed neural architecture that spans the

90    face perception and ToM neural networks.

91    In the current fMRI study, we addressed the extent to which face perception and ToM

92    networks contribute to forming impressions based on facial appearance. The experiment used

93    a repetition suppression (RS) design (Grill-Spector, Henson, & Martin, 2006; Barron et al.,

94    2016). RS designs measure a reduced BOLD response following a repeated stimulus feature

95    and a release from suppression following a novel stimulus feature. Compared to conventional

96    subtraction designs, which can show if a brain region shows magnitude differences between

97    conditions, RS studies hold the potential to study neural processes at the level of neural

98    populations within a given brain region. A brain region that shows RS, therefore, can allow

99    inferences about the organisation of underlying neural populations (Barron et al., 2016;

100   Figure 1). We created face stimuli that cued high and low trait judgments and showed these

101   stimuli to participants in a sequence that created novel and repeated events. To identify

102   functional regions of interest, we used established face and ToM localiser tasks and to bolster

103   statistical power we used an analysis pipeline that has been demonstrated to exhibit high

104   functional resolution and sensitivity (Julian et al., 2012; Nieto-Castanon & Fedorenko, 2012).

105   If the face and ToM networks are engaged in forming impressions based on facial features in

106   the manner that we predict, we would expect to observe repetition suppression for face traits

107   in both networks.

108

109   **Method**

110
111   *Participants*

112   Twenty-eight participants completed the experiment (14 female; $M_{age}$=23.96, SD=5.52). All

113   participants received a monetary reimbursement (£15), had normal or corrected-to-normal

114   vision and gave informed consent according to the local ethics guidelines.

115

116   *Stimuli and experimental tasks*

117      *Stimuli*. Face stimuli were initially selected from a face database created at Bangor

118      University. The Bangor face database comprises photographs of participants with an

119      emotionally neutral expression and self-report measures of various personality and subclinical

120      traits (Kramer and Ward, 2010; Jones et al., 2012; Scott et al., 2013). Individual images were

121      extracted from the database and transformed along four personality or health dimensions

122      (Extraversion, Agreeableness, Neuroticism and Physical Health). These dimensions were

123      chosen because prior work had shown that these dimensions were readily identifiable in

124      composite stimuli, which average faces across multiple identities (Kramer & Ward, 2010).

125      All face transformations were performed in JPsychomorph (Tiddeman et al., 2001).

126      Face stimuli were produced by transforming an original face image from the database towards

127      an average template of a high trait face (High Trait) or towards an average template of a low

128      trait face (Low Trait). Template faces were produced by creating a composite of the 15

129      individuals with the highest or lowest ratings along each of the four dimensions. For example,

130      for physical health, an average composite of the 15 most physically healthy individuals in the

131      database was created as well as an average composite of the 15 least physically healthy

132      individuals. This process was repeated for all four dimensions (Extraversion, Agreeableness,

133      Neuroticism and Physical Health). To avoid skin colour or make-up influencing the

134      construction of composite images, only individuals that were white and not wearing make-up

135      were included. Also, to simplify the design space, we only used images of female individuals.

136      Individual images were selected that were between those included in the high and low

137      composites and also met the above inclusion criteria (i.e., white females who were not

138      wearing make-up). Additionally, the individual in each image had provided consent that their

139      individual face could be shown in later studies. The number of individuals fitting these

140      criteria per trait were: Extraversion = 54, agreeableness = 53, neuroticism = 56, physical

141      health = 54, which made a total of 217 IDs. Note that these were not unique IDs and most

142     were used across traits. An individual face image was then transformed in two ways: towards

143     the high trait composite image by 100% and towards the low trait composite image by 100%

144     (Figure 1). A 100% transform retains the identity cues of the original image whilst shifting

145     the appearance by 100% of the shape, colour, and texture difference between the high and the

146     low composite images. This produced two transformed images per original stimulus (High

147     trait, Low trait), which made 434 images in total.

148         We transformed stimuli in this manner to exaggerate the distinctive facial features

149     associated with particular trait characteristics, whilst maintaining a variety of facial identities

150     by using individual faces rather than composite images. We did not use composite face

151     images, as this would reduce the variety of identities presented during the scanning task,

152     which may lead participants to disengage. Indeed, we wanted to maintain interest in the

153     stimuli and thus encourage a 'fresh' social judgment on every trial and increasing variety of

154     idiosyncratic facial features and identities seemed a concrete way of doing so.

155         *Pilot task*. To assess the extent to which these stimuli would cue distinct trait

156     judgments, we ran a pilot behavioural experiment (see Supplementary Method). The pilot

157     experiment demonstrated that judgements of Low and High Extraversion, Neuroticism and

158     Physical health were perceived distinctly and as anticipated based on prior research (Kramer

159     & Ward, 2010; Supplementary Figure 2). However, there was no difference in the perception

160     of high and low agreeableness (Supplementary Figure 2). Prior work on agreeableness

161     averaged multiple facial identities to create one composite image (Kramer & Ward, 2010). In

162     the current study, we used individual faces that had been transformed towards High or Low

163     trait features. Therefore, after the pilot study, it was unclear if the lack of distinct behavioural

164     judgments based on agreeableness was due to the method of stimulus construction. We

165     decided to leave the agreeableness stimuli in for the scanning experiment in order to see if the

7

166   same pattern of results persisted in new participants and, if so, if there were neural effects in

167   the absence of distinct behavioural judgments.

168        *Main task*. The main task used an event-related design with two types of face stimuli

169   presented (High trait and Low trait faces). The design of the main task is illustrated in Figure

170   1. Each run comprised 17 blocks of 9 trials. On every trial participants were shown a face and

171   asked to make a social judgement. At the start of each block, participants were shown a

172   written statement and a ratings scale for 4 seconds (1=Strongly disagree, 2= Disagree, 3 =

173   Agree, 4=Strongly agree). The task for participants was to rate how well the person matched

174   the statement. Each trial lasted 3s and participants were instructed to make a judgment based

175   on their initial reaction or "gut instinct". The scale was always the same, but was included

176   with the statement before each block as a reminder. Participants responded on a button box

177   within the scanner by pressing the corresponding key. Between blocks a white cross was

178   presented on a black screen for a randomly selected duration of 2, 3 or 4 seconds.

179        Each block contained High and Low versions of stimuli from one category (e.g.,

180   Physical Health) and each trial showed a different person. However, participants were not

181   shown high and low versions of the same person in the same category. Instead, participants

182   were shown either a high or a low version of an individual to avoid confusion with seeing the

183   same person transformed to opposite ends of a single dimension. Statements for each block

184   related to the category of stimuli presented in that block. For example, in a physical health

185   block, participants made judgments based on statements concerning physical health. Four

186   statements per category were taken for Extraversion, Agreeableness and Neuroticism from the

187   corresponding scales of the mini-IPIP (Donnellan et al., 2006). An example of an

188   Extraversion statement is "Is the life of the party". For physical health judgements, items were

189   used from the Short-Form 12-Item Health Survey, which assesses physical health (Ware,

190   Kosinski, & Keller, 1996). An example physical health statement is ''Finds it easy to climb

8

191   the stairs". The first block in a run was randomly selected as a starter block. Subsequently,

192   four blocks of each category were presented in a pseudorandom order such that each block

193   followed each other equally often.

194      Each block began with a starter trial, which was randomly selected from that category.

195   The next 8 trials were sequenced to achieve an even number of novel and repeated trials with

196   novel and repeated trials following each other equally often. Each trial was defined in

197   reference to the preceding trial. For example, a High trait trial that was preceded by a High

198   trait trial would be defined as a repeated trial, whereas a High trait trial that was preceded by a

199   Low trait trial would be defined as a novel trial. This design produced the two conditions of

200   interest, which were modelled as separate regressors in the general linear model:

201   Novel_FaceTrait and Repeated_FaceTrait. The starter trial was included as an additional

202   regressor of no interest since the trial was not preceded by any trial and therefore it was not

203   comparable to the other trials. Each trial was modelled from the onset of the first image for a

204   nominal zero second duration. Across a block there were four trials per condition and across a

205   run there were 68 trials per condition. Each participant completed two runs of the main task,

206   which made 136 trials per condition over the entire experiment. In addition, before entering

207   the scanner, participants completed two practice blocks of the main task.

208      *Face localiser*. To identify face-selective brain regions, we used an established face

209   localiser (Pitcher et al., 2011). Five categories of stimuli were shown to participants (faces,

210   bodies, scenes, objects, scrambled objects), with one category per block. Each block lasted

211   18s and showed six 3s movie clips from that category. A total of two blocks were shown in

212   each functional run. At the start, middle and end of each functional run, there was a rest

213   condition for 18s. In the rest condition, a series of six uniform colour fields were presented

214   for 3s each. The order of blocks was reversed from the first to the second bock (e.g., fixation,

215   faces, objects, scenes, bodies, scrambled objects, fixation, scrambled objects, bodies, scenes,

216    objects, faces, fixation). Throughout all blocks, participants were instructed to watch the

217    movies but were not given an explicit task.

218         *Theory-of-mind localiser*. To localise brain regions associated with ToM, we used an

219    established ToM-localiser (Dodell-Feder et al., 2011; http://saxelab.mit.edu/superloc.php).

220    Participants read 10 short false belief stories, in which the belief characters have about the

221    state of the world is false. Participants also read 10 false photograph stories, where a

222    photograph, map, or sign has out-dated or misleading information. After reading each story,

223    participants had to answer whether the subsequently presented statement is true or false. Each

224    run started with a 12 second rest period, after which the stories (10 seconds) and questions (4

225    seconds) were presented for 14 seconds combined. Each story was separated by a 12 second

226    rest period. The order of items and conditions was identical for each subject. In the first run,

227    stimuli 1 – 5 from each condition were presented, and the remaining stimuli were presented

228    during the second block.

229         *Procedure*. Participants completed two runs of the main task. Two additional

230    functional runs were also completed as part of another experiment – one run included a

231    version of an imitation inhibition task (Brass et al., 2000) and one run included a version of a

232    flanker task (Eriksen & Eriksen, 1974). These runs occurred before each run of the main task

233    in order to add variety and offset boredom. Subsequently, participants then completed one run

234    of the face localiser and two runs of the ToM-localiser. The ToM-localiser was always

235    presented after participants had completed the main task, to ensure that participants were not

236    primed towards making trait inferences during the main task. All participants completed an

237    anatomical scan.

238

239    *Data acquisition*

240    The experiment was conducted on a 3 Tesla scanner (Philips Achieva), equipped with

241    a 32-channel SENSE-head coil. Stimuli were displayed on a MR safe BOLD screen

242    (Cambridge Research Systems: http://www.crsltd.com/) behind the scanner, which

243    participants viewed via a mirror mounted on the head-coil. T2*-weighted functional images

244    were acquired using a gradient echo echo-planar imaging (EPI) sequence with the following

245    parameters: acquisition time (TR) = 2000 ms; echo time (TE) = 30ms; flip angle = 90°;

246    number of axial slices = 35; slice thickness = 4mm; slice gap = 0.8mm; field of view = 230 x

247    230 x 167mm$^3$. After the functional runs were completed, a high-resolution T1-weighted

248    structural image was acquired for each participant (voxel size = 1 mm$^3$, TE = 3.8 ms, flip

249    angle = 8°, FoV = 288 × 232 × 175 mm$^3$). Four dummy scans (4 * 2000 ms) were routinely

250    acquired at the start of each functional run and were excluded from analysis. 291 volumes per

251    functional run were collected, except for participant 1 where 288 and 289 volumes were

252    collected in block 1 and 2 respectively.

253

254    *Behavioural data analysis*

255    During scanning, faces were rated on four dimensions in a similar manner to the pilot

256    experiment. The four dimensions included Extraversion, Agreeableness, Neuroticism and

257    Physical Health and the ratings scale ranged from 1 to 4 (1 = Strongly disagree, 2 = Disagree,

258    3 = Agree, 4 = Strongly agree). Ratings on each of these dimensions were compared between

259    high and low transformed stimuli. We expected high transformed stimuli to be rated in a

260    manner that is more consistent with descriptions of the trait category. For instance, based on

261    prior work (Kramer & Ward, 2010), as well as our behavioural pilot data, we would expect

262    stimuli transformed towards high physical health to be rated in a manner consistent with

263    higher physical heath. To compare high and low transformed stimuli, we computed difference

264    scores between high and low stimulus categories as well as interval estimates using 95%

11

265    confidence intervals (Cumming, 2013). We also computed a paired-samples t-test and a

266    standardised effect size for each difference score (Cohen's $d_z$; Cohen, 1992; Lakens, 2013).

267

268    *fMRI data preprocesing and analysis*

269         *Preprocessing*. Head motion was examined for each participant on each task, with an

270    exclusion criteria if displacement across either task exceeded 3 millimetres. We report for

271    each task how many runs or participants were removed for each experiment. fMRI data were

272    analysed with Statistical Parametric Mapping software (SPM8; Wellcome Trust Department

273    of Cognitive Neurology, London, UK: www.fil.ion.ucl.ac.uk/spm/). Data were realigned,

274    unwarped, corrected for slice timing, and normalised to the MNI template with a resolution of

275    $3mm^3$. Images were then spatially smoothed (5mm).

276         *Analysis*. We used spm_ss to perform our primary analyses (Julian et al., 2012; Nieto-

277    Castanon & Fedorenko, 2012; http://www.nitrc.org/projects/spm_ss). Spm_ss enables a

278    subject-specific approach to fMRI data analysis. Like other ROI approaches, functional

279    regions of interest (fROI) are defined and tested in separate data to ensure that the analyses

280    are not circular (Kriegeskorte et al., 2009). The advantage of spm_ss is that it uses an

281    algorithm (or functional parcels from prior datasets) to define fROIs in a group-constrained

282    and subject-specific manner (GSS). This means that the approach benefits from showing

283    group consistency across participants, without requiring complete voxel-level overlap across

284    participants. As such, the approach integrates single-subject specificity within individuals

285    with group-constrained consistency across individuals.

286         We used GSS to define fROIs using separate localiser data. fROIs were first defined

287    using Face and ToM network localisers before we tested how these fROIs responded in our

288    main task contrasts of interest (RS FaceTraits). To do so, the following steps were taken. 1)

289    Using localiser data, we computed activation maps in individuals, thresholded these images (p

12

290   < 0.001, uncorrected) and overlaid them on top of one another. The resultant overlay map

291   contains information on the percentage of individuals that show an above threshold response.

292   2) The overlay map was then divided into regions by an image parcellation algorithm. 3) The

293   resulting regions are then investigated in terms of the proportion of subjects that show some

294   suprathreshold voxels. 4) Regions that overlap in a substantial number of participants (>50%)

295   are then interrogated using independent data (i.e., data from the main task). Statistical tests

296   across participants were performed on percent signal change values extracted from the fROIs

297   according to contrasts of interest.

298         *Main task contrasts*. For the fMRI data analysis of the main task, we computed our

299   primary contrast of interest: RS Face Traits (Novel_FaceTrait > Repeated_FaceTrait).

300         *Face localiser contrasts*. Each block was modelled from the onset of the first trial for

301   the entire block (18 seconds). A design matrix was fit for each participant with five regressors

302   per block (Faces, Bodies, Scenes, Objects, Scrambled objects). To identify face-selective

303   regions, a Face > All baseline contrast was evaluated in individual participants (Dynamic

304   Faces > Dynamic Scenes + Objects + Scrambled Objects)

305         *ToM localiser contrasts*. A design matrix was fit for each participant with 2

306   regressors, one for each experimental condition (false beliefs and false photographs). The

307   ToM-network was revealed by contrasting false beliefs with false photographs (False Beliefs

308   > False Photographs).

309

310   **Results**

311   *Behavioural data*

312   During scanning, high trait faces were rated more consistent with trait characteristics than low

313   trait faces for extraversion $t(27)=10.88$, $p < 0.001$, $d_z = 2.06$, neuroticism $t(27)=4.50$, $p < 0.001$,

314   $d_z = 0.85$, and physical health $t(27)=3.73$, $p < 0.001$, $d_z = 0.71$ (Figure 2). There was no

13

315    difference between high and low trait faces for judgments of agreeableness t(27)=-0.33, p =

316    0.63, d, = -0.06 (Figure 2). This pattern of results closely replicates our pilot data.

317

318    *fMRI data*

319           The GSS analysis using the face localiser data revealed nine regions where a majority

320    of participants showed a greater response to faces than all other baseline conditions. Three of

321    these regions are of particular interest given our predictions as they represent the core face

322    perception network. These regions include rOFA, rFFA and r STS/STG. None of the three

323    regions of interest showed RS for Face Traits estimated from data from the main task (Figure

324    3A; Table 1). If we widen the search to all nine face responsive regions, we do not find RS for

325    Face Traits in any of the ROIs (Supplementary Table 4).

326           The GSS analysis using the ToM localiser data revealed nine regions where a majority

327    of participants showed a greater response to false belief stories than false photograph stories.

328    Four of these regions are of particular interest given our predictions regarding specific nodes

329    of the ToM network. These regions include rTPJ, mPFC and r anterior STG / temporal pole.

330    None of these regions showed RS for Face Traits estimated from data from the main task

331    (Figure 3B; Table 1). If we widen the search to all nine regions from the ToM localiser, we do

332    not find RS for Face Traits in any of the ROIs (Supplementary Table 4).

333           As judgements of agreeableness showed no behavioural differences in perceptions of

334    trait character or health (Figure 2), we removed agreeableness blocks from the analysis, but

335    the results remained the same in both brain networks of interest.

336           Finally, we completed an exploratory whole-brain analysis, in order to test if regions

337    outside of the Face and ToM networks showed RS for Face Traits. Using SPM8, we

338    calculated Novel > Repeated Face Traits at the single subject level before completing a

339    random effects analysis at the group level using the same contrast. At the group level, no

14

340  significant clusters of activity were found ($p < 0.001$, K=10, $p<0.05$ family wise error

341  corrected). Even at a more liberal threshold ($p < 0.001$, uncorrected for multiple

342  comparisons), no clusters emerged from this contrast.

343       Data from this experiment are freely available, including the behavioural and fROI

344  data (osf.io/7knrp), as well as data from the whole-brain analysis

345  (https://neurovault.org/collections/HDLVMPQU/).

346

**Discussion**

347

348  Here we show that faces readily cued accurate person judgments regarding extraversion,

349  neuroticism and physical health, but the neural networks associated with face perception and

350  ToM showed no sensitivity in terms of repetition suppression to trait judgements. As such, we

351  do not provide evidence that supports population coding models of face perception that

352  include dimensions for high and low trait features along the ventral visual stream and in the

353  ToM network. Due to aspects of the experimental design and analysis pipeline, which

354  bolstered statistical power and sensitivity, we have reasonable confidence that we could detect

355  effects of a moderate size, should they exist. However, it remains possible that these regions

356  are sensitive to other trait dimensions of person perception such as trustworthiness or other

357  types of facial stimuli, such as synthetic stimuli. The null findings reported here, therefore,

358  add value to models of neural organisation by showing instances where effects are absent or

359  small. In addition, by publishing null results, we provide a less biased scientific record, one

360  that future studies can build upon by appropriately powering studies (Open Science

361  Collaboration, 2015; Simmons et al., 2011). Indeed, future work can use these results to guide

362  further interrogation of what is fundamentally an interesting scientific and social question that

363  relates to understanding the neural mechanisms associated with how trait inferences are cued

364  from facial appearance.

15

365

**Understanding the neural basis of impression formation based on facial appearance**

The current experiment provides no evidence that populations of neurons in face perception or ToM networks code for facial features that are associated with distinct trait judgements of extraversion, neuroticism or physical health. Moreover, a whole-brain analysis showed no effects in the amygdala or ventral striatum, which have previously been associated with social evaluations of faces based on valence (Bzdok et al., 2011; Mende-Siedlecki et al., 2013; Said et al., 2010; 2011; Todorov et al., 2013). Observers were able to accurately discriminate faces on the basis of the social trait being displayed for the majority of person dimensions. However, we were unable to uncover the neural substrates for this discrimination. In particular, we could not find evidence for our hypothesis that brain regions representing features and judgements for high traits might be separable from those representing low traits. Rather than distinct populations of neurons in the same neural region coding for high and low trait features and judgments, which a neural response consistent with RS would support (Grill-Spector, Henson, & Martin, 2006; Barron et al., 2016), the results may suggest that face perception and ToM networks have a common neural parameter that codes for the perceptual and judgement space under investigation. If so, the same neural populations would be engaged on all trials, whether novel or repeated. For example, if the same features of the face cue high and low judgements, they would be engaged equally on novel and repeated trials. The ultimate judgment would differ between high and low trait faces, but the underlying neural architecture would be similar. This proposal is speculative, however, and would require further testing and confirmation.

An alternative possibility is that RS may not have been sensitive enough to detect the fine-grained population coding structure that was tested. To bolster statistical power, we included a large number of trials per condition for fMRI research (136), we tested 28

16

390    participants and we used a single-subject analysis pipeline that has been shown to have

391    relatively high sensitivity and functional resolution in multi-subject analyses (Nieto-Castanon

392    & Fedorenko, 2012). Nonetheless, RS may have been smaller than we could detect with

393    reasonable confidence. Future work may consider multi voxel pattern analysis approaches

394    (Kriegeskorte & Kievit, 2013), which have been shown to be more sensitive than RS

395    approaches in the domain of vision (Sapountzis et al., 2010). In addition, future work may

396    consider the relationship between face and ToM networks as prior functional connectivity

397    research has shown that the ToM network functionally couples with nodes of body perception

398    network (Greven et al., 2016; Greven & Ramsey, 2017a, b). The hypothesis that such future

399    connectivity research could pursue is that the representation of trait judgments from faces

400    may span across face perception and ToM networks rather than only within them.

401

402    **Limitations and constraints on generality**

403    In the current study, we do not show RS for trait inferences based on facial appearance. By

404    contrast, other work using written descriptions of behaviour, which imply trait inferences, have

405    shown that vental medial prefrontal cortex (vmPFC) shows RS for trait implying behaviours

406    (Heleven et al., 2017; Heleven & Van Overwalle, 2016; Ma et al., 2014; Van Overwalle et al.,

407    2016). Indeed, this work shows that vmPFC encodes trait representations for familiar (Heleven

408    & Van Overwalle, 2016) and unfamiliar individuals (Heleven et al., 2017), as well as for distinct

409    traits such as valence and competence (Van Overwalle et al., 2016). Therefore, it is important

410    that we acknowledge relevant constraints on the generality of our findings (Simons et al., 2017).

411    Our data, at least with the stimuli that we used, do not support the view that vmPFC stores a

412    person or trait code, which can be easily accessed or engaged irrespective of the type of input

413    (face or text). It could be that written text is simply a more salient way to engage trait inferences,

414    which could lead to the discrepant results. Alternatively, it be might be that not all sources of

17

415    input (face, text) or all types of person judgment (extraversion, health, trustworthiness) are

416    coded in a similar neural structure. Future work that directly tests interactions between input

417    type and judgments type would be valuable.

418        Of particular interest for future work would be to test judgments from faces that vary

419    on a valence / trustworthiness dimension (Todorov et al., 2008). In the current study, the

420    behavioural data showed that participants' judgments did not distinguish between high and low

421    agreeableness faces, which is the closest dimension to valence / trustworthiness. However,

422    participants were sensitive to other dimensions, such as extraversion, neuroticism and physical

423    health. Importantly, recent models of social judgments from faces have shown that appraising

424    faces has three partly distinct dimensions including valence / trustworthiness, dominance and

425    attractiveness (Sutherland et al., 2013). Since judgments of physical health have been associated

426    with attractiveness (Little et al., 2011), our physical health dimension closely resembles a key

427    dimension in the person perception (attractiveness). Therefore, it may be that health and

428    attractiveness judgments, as well as some other types of traits judgment (extraversion,

429    neuroticism), are not coded in the same way as valence / trustworthiness judgments. Indeed,

430    given the role of valence judgments in guiding approach and avoidance behaviours, it may be

431    that there is a more distinct neural architecture dedicated to perceiving such traits.

432        In the current study, we used morphed images of real human faces. Our findings,

433    therefore, apply most directly to faces that look straightforwardly human. A complementary

434    avenue for future research would be to test models of trait inference from synthetic, computer-

435    generated facial stimuli. The advantage of using computer-generated stimuli would be tighter

436    experimental control, which may boost the ability to detect effects of interest. The obvious

437    disadvantage, however, compared to the current approach of using real photographs, is the

438    artificial limit imposed on ecological validity (Sutherland et al., 2013). Using synthetic images

439    that produce more extreme facial attributes, which differ from the average more, may be

440    important, given research that shows widespread neural responses to faces at high and low ends

441    of continua (Said et al., 2010; 2011). Indeed, even though the majority of trait inferences

442    showed reliable behavioural judgments, it is possible that the similarity between our stimuli

443    reduced the saliency of features that cue trait judgments. Relatedly, we made sure that

444    participants would not see stimuli morphed to different traits in the same block in order to avoid

445    confusion between identities and facial attributes. But, by doing so, this may have made the

446    distinction between high and low exemplars less obvious. An alternative approach would be to

447    show high and low version in the same blocks.

448

449    **Open science and the file drawer problem**

450    Since null results and smaller effect sizes are typically relegated to the file drawer (Rosenthal,

451    1979), the current literature has a publication bias, which prioritises statistically significant

452    results and produces an overestimate of effect sizes. As such, null results from well designed

453    and well powered studies are important if the field is going to move towards a more precise

454    estimate of population effect sizes. Without greater acknowledgement of the value of null

455    results, artificially high estimates of effect sizes will continue to bias models of cognition and

456    brain function, skewing the design of future research and resulting in misallocation of

457    resources (Munafo et al., 2017). Indeed, as outlined above, a null result can make several

458    important contributions to future research (Zwaan et al., 2017). First, replications and

459    extensions can be powered to detect smaller effects or a task can be changed to increase

460    sensitivity. Second, other analysis methods, such as multi-voxel pattern analysis or measures

461    of connectivity (Kriegeskorte & Kievit, 2013; Bullmore & Sporns, 2009), may be prioritised

462    as they may more closely capture the information under investigation. As the data from this

463    study are readily available in online open access repositories, we hope that future research can

464    be guided by this work.

465

466    **Acknowledgements**

467

470

**References**

Barron, H. C., Garvert, M. M., & Behrens, T. E. (2016). Repetition suppression: a means to index neural representations using BOLD?. Phil. Trans. R. Soc. B, 371(1705), 20150355.

Borkenau, P., Brecke, S., Möttig, C., & Paelecke, M. (2009). Extraversion is accurately perceived after a 50-ms exposure to a face. Journal of Research in Personality, 43(4), 703-706.

Borkenau, P., & Liebler, A. (1992). Trait inferences: Sources of validity at zero acquaintance. Journal of personality and social psychology, 62(4), 645.

Brass, M., Bekkering, H., Wohlschläger, A., & Prinz, W. (2000). Compatibility between observed and executed finger movements: comparing symbolic, spatial, and imitative cues. Brain and cognition, 44(2), 124-143.

Bullmore, E., & Sporns, O. (2009). Complex brain networks: graph theoretical analysis of structural and functional systems. *Nature Reviews Neuroscience*, *10*(3), 186-198.

Bzdok, D., Langner, R., Caspers, S., Kurth, F., Habel, U., Zilles, K., ... & Eickhoff, S. B. (2011). ALE meta-analysis on facial judgments of trustworthiness and attractiveness. Brain Structure and Function, 215(3-4), 209-223.

Cohen, J. (1992). A power primer. Psychological bulletin, 112(1), 155.

Cumming, G. (2013). Understanding the new statistics: Effect sizes, confidence intervals, and meta-analysis. Routledge.

Donnellan, M. B., Oswald, F. L., Baird, B. M., & Lucas, R. E. (2006). The mini-IPIP scales: tiny-yet-effective measures of the Big Five factors of personality. Psychological assessment, 18(2), 192.

Dodell-Feder, D., Koster-Hale, J., Bedny, M., & Saxe, R. (2011). fMRI item analysis in a theory of mind task. Neuroimage, 55(2), 705-712.

497    Duchaine, B., & Yovel, G. (2015). A revised neural framework for face processing. Annual

498        Review of Vision Science, 1, 393-416.

499    Emery, N. J. (2000). The eyes have it: the neuroethology, function and evolution of social

500        gaze. Neuroscience & Biobehavioral Reviews, 24(6), 581-604.

501    Eriksen, B. A., & Eriksen, C. W. (1974). Effects of noise letters upon the identification of a

502        target letter in a nonsearch task. Attention, Perception, & Psychophysics, 16(1), 143-

503        149.

504    Frith, C. D., & Frith, U. (1999). Interacting minds--a biological basis. Science, 286(5445),

505        1692-1695.

506    Gauthier, I., Tarr, M. J., Moylan, J., Skudlarski, P., Gore, J. C., & Anderson, A. W. (2000).

507        The fusiform "face area" is part of a network that processes faces at the individual

508        level. Journal of cognitive neuroscience, 12(3), 495-504.

509    Greven, I. M., Downing, P. E., & Ramsey, R. (2016). Linking person perception and person

510        knowledge in the human brain. Social cognitive and affective neuroscience, 11(4),

511        641-651.

512    Greven, I. M., & Ramsey, R. (2017a). Person perception involves functional integration

513        between the extrastriate body area and temporal pole. Neuropsychologia, 96, 52-60.

514    Greven, I. M., & Ramsey, R. (2017b). Neural network integration during the perception of in-

515        group and out-group members. Neuropsychologia, *106,* 225-235.

516    Grill-Spector, K., Henson, R., & Martin, A. (2006). Repetition and the brain: neural models of

517        stimulus-specific effects. Trends in cognitive sciences, 10(1), 14-23.

518    Haxby, J. V., Hoffman, E. A., & Gobbini, M. I. (2000). The distributed human neural system

519        for face perception. Trends in cognitive sciences, 4(6), 223-233.

520   Heleven, E., Boukhlal, S., & Van Overwalle, F. (2017). A stranger in my brain: Neural

521        representation for unfamiliar persons using fMRI repetition suppression. *Social*

522        *Neuroscience*, 1-11. doi: 10.1080/17470919.2017.1358663

523   Heleven, E., & Van Overwalle, F. (2016). The person within: memory codes for persons and

524        traits using fMRI repetition suppression. *Social Cognitive and Affective Neuroscience,*

525        *11*(1), 159-171. doi: 10.1093/scan/nsv100

526   Jack, R. E., & Schyns, P. G. (2017). Toward a social psychophysics of face communication.

527        Annual review of psychology, 68, 269-297.

528   Jones, A. L., Kramer, R. S. S., & Ward, R. (2012). Signals of personality and health: The

529        contributions of facial shape, skin texture, and viewing angle. Journal of Experimental

530        Psychology: Human Perception and Performance, 38(6), 1352-1361.

531   Julian, J. B., Fedorenko, E., Webster, J., & Kanwisher, N. (2012). An algorithmic method for

532        functionally defining regions of interest in the ventral visual pathway. Neuroimage,

533        60(4), 2357-2364.

534   Kanwisher, N., McDermott, J., & Chun, M. M. (1997). The fusiform face area: a module in

535        human extrastriate cortex specialized for face perception. Journal of neuroscience,

536        17(11), 4302-4311.

537   Kanwisher, N. (2010). Functional specificity in the human brain: a window into the functional

538        architecture of the mind. Proceedings of the National Academy of Sciences, 107(25),

539        11163-11170.

540   Kramer, R. S., & Ward, R. (2010). Internal facial features are signals of personality and

541        health. The Quarterly Journal of Experimental Psychology, 63(11), 2273-2287.

542   Kriegeskorte, N., & Kievit, R. A. (2013). Representational geometry: integrating cognition,

543        computation, and the brain. Trends in cognitive sciences, 17(8), 401-412.

544   Kriegeskorte, N., Simmons, W. K., Bellgowan, P. S., & Baker, C. I. (2009). Circular analysis

545       in systems neuroscience: the dangers of double dipping. Nature neuroscience, 12(5),

546       535-540.

547   Lakens, D. (2013). Calculating and reporting effect sizes to facilitate cumulative science: a

548       practical primer for t-tests and ANOVAs. Frontiers in psychology, 4.

549   Little, A. C., Jones, B. C., & DeBruine, L. M. (2011). Facial attractiveness: evolutionary

550       based research. *Philosophical Transactions of the Royal Society B: Biological*

551       *Sciences, 366*(1571), 1638-1659. doi: 10.1098/rstb.2010.0404

552   Ma, N., Baetens, K., Vandekerckhove, M., Kestemont, J., Fias, W., & Van Overwalle, F.

553       (2014). Traits are represented in the medial prefrontal cortex: an fMRI adaptation

554       study. *Social Cognitive and Affective Neuroscience, 9*(8), 1185-1192. doi:

555       10.1093/scan/nst098

556   Mende-Siedlecki, P., Verosky, S. C., Turk-Browne, N. B., & Todorov, A. (2013). Robust

557       selectivity for faces in the human amygdala in the absence of expressions. Journal of

558       cognitive neuroscience, 25(12), 2086-2106.

559   Mitchell, J. P., Macrae, C. N., & Banaji, M. R. (2005). Forming impressions of people versus

560       inanimate objects: social-cognitive processing in the medial prefrontal cortex.

561       Neuroimage, 26(1), 251-257.

562   Mitchell, J. P., Heatherton, T. F., & Macrae, C. N. (2002). Distinct neural systems subserve

563       person and object knowledge. Proceedings of the National Academy of Sciences,

564       99(23), 15238-15243.

565   Munafò, M. R., Nosek, B. A., Bishop, D. V., Button, K. S., Chambers, C. D., du Sert, N. P.,

566       ... & Ioannidis, J. P. (2017). A manifesto for reproducible science. Nature Human

567       Behaviour, 1, 0021.

24

568    Nieto-Castañón, A., & Fedorenko, E. (2012). Subject-specific functional localizers increase

569        sensitivity and functional resolution of multi-subject analyses. Neuroimage, 63(3),

570        1646-1669.

571    Open Science Collaboration. (2015). Estimating the reproducibility of psychological science.

572        Science, 349(6251), aac4716.

573    Over, H., & Cook, R. (2018). Where do spontaneous first impressions of faces come from?.

574        *Cognition*, *170*, 190-200.

575    Pitcher, D., Dilks, D. D., Saxe, R. R., Triantafyllou, C., & Kanwisher, N. (2011). Differential

576        selectivity for dynamic versus static information in face-selective cortical regions.

577        Neuroimage, 56(4), 2356-2363.

578    Penton-Voak, I. S., Pound, N., Little, A. C., & Perrett, D. I. (2006). Personality judgments

579        from natural and composite facial images: More evidence for a "kernel of truth" in

580        social perception. Social Cognition, 24(5), 607-640.

581    Rosenthal, R. (1979). The file drawer problem and tolerance for null results. Psychological

582        bulletin, 86(3), 638.

583    Said, C. P., Dotsch, R., & Todorov, A. (2010). The amygdala and FFA track both social and

584        non-social face dimensions. *Neuropsychologia, 48*(12), 3596-3605. doi:

585        https://doi.org/10.1016/j.neuropsychologia.2010.08.009

586    Said, C. P., Haxby, J. V., & Todorov, A. (2011). Brain systems for assessing the affective

587        value of faces. *Philosophical Transactions of the Royal Society B: Biological

588        Sciences, 366*(1571), 1660-1670. doi: 10.1098/rstb.2010.0351

589    Sapountzis, P., Schluppeck, D., Bowtell, R., & Peirce, J. W. (2010). A comparison of fMRI

590        adaptation and multivariate pattern classification analysis in visual cortex.

591        Neuroimage, 49(2), 1632-1640.

592  Saxe, R., & Kanwisher, N. (2003). People thinking about thinking people: the role of the
593      temporo-parietal junction in "theory of mind". Neuroimage, 19(4), 1835-1842.
594  Scott, N. J., Kramer, R. S. S., Jones, A. L., & Ward, R. (2013). Facial cues to depressive
595      symptoms and their associated personality attributions. Psychiatry Research, 30, 47-
596      53.
597  Simmons, J. P., Nelson, L. D., & Simonsohn, U. (2011). False-positive psychology:
598      Undisclosed flexibility in data collection and analysis allows presenting anything as
599      significant. Psychological science, 22(11), 1359-1366.
600  Simons, D. J., Shoda, Y., & Lindsay, D. S. (2017). Constraints on Generality (COG): A
601      Proposed Addition to All Empirical Papers. *Perspectives on Psychological Science,*
602      *12*(6), 1123-1128. doi: 10.1177/1745691617708630
603  Sutherland, C. A. M., Oldmeadow, J. A., Santos, I. M., Towler, J., Michael Burt, D., &
604      Young, A. W. (2013). Social inferences from faces: Ambient images generate a three-
605      dimensional model. *Cognition, 127*(1), 105-118. doi:
606      https://doi.org/10.1016/j.cognition.2012.12.001
607  Todorov, A., Olivola, C. Y., Dotsch, R., & Mende-Siedlecki, P. (2015). Social attributions
608      from faces: Determinants, consequences, accuracy, and functional significance.
609      Annual Review of Psychology, 66.
610  Todorov, A., Mende-Siedlecki, P., & Dotsch, R. (2013). Social judgments from faces. Current
611      opinion in neurobiology, 23(3), 373-380.
612  Todorov, A., Said, C. P., Engell, A. D., & Oosterhof, N. N. (2008). Understanding evaluation
613      of faces on social dimensions. *Trends in Cognitive Sciences, 12*(12), 455-460. doi:
614      http://dx.doi.org/10.1016/j.tics.2008.10.001
615  Uleman, J. S., Adil Saribay, S., & Gonzalez, C. M. (2008). Spontaneous inferences, implicit
616      impressions, and implicit theories. Annu. Rev. Psychol., 59, 329-360.

617  Van Overwalle, F. (2009). Social cognition and the brain: a meta-analysis. Human brain

618      mapping, 30(3), 829-858.

619  Van Overwalle, F., Ma, N., & Baetens, K. (2016). Nice or nerdy? The neural representation of

620      social and competence traits. *Social Neuroscience, 11*(6), 567-578. doi:

621      10.1080/17470919.2015.1120239

622  Wang, H., Hahn, A. C., DeBruine, L. M., & Jones, B. C. (2016). The Motivational Salience of

623      Faces Is Related to Both Their Valence and Dominance. *PLoS ONE, 11*(8), e0161114.

624      doi: 10.1371/journal.pone.0161114

625  Ware Jr, J. E., Kosinski, M., & Keller, S. D. (1996). A 12-Item Short-Form Health Survey:

626      construction of scales and preliminary tests of reliability and validity. Medical care,

627      34(3), 220-233.

628  Zebrowitz, L. A. (2011). Ecological and social approaches to face perception. The Oxford

629      handbook of face perception, 31-50.

630  Zwaan, R. A., Etz, A., Lucas, R. E., & Donnellan, B. (2017, October 20). Making Replication

631      Mainstream. Retrieved from psyarxiv.com/4tg9c

632
633

634 **Table 1.** Main task ROI data.

635

| Region | | | | Novel>Repeated | | |
|---|---|---|---|---|---|---|
| | ROI size (voxels) | Average localiser mask size (voxels) | Inter-subject overlap (%) | Percent signal change (SEM) | t | p(fdr) |
| *Face localiser* | | | | | | |
| Right OFA | 412 | 84 | 93 | -.005 (.19) | -.03 | .82 |
| Right FFA | 223 | 44 | 86 | .067 (.17) | -.41 | .82 |
| Right pSTS | 143 | 24 | 75 | -.127 (.17) | -.73 | .82 |
| | | | | | | |
| *ToM localiser* | | | | | | |
| Right TPJ | 828 | 230 | 96 | -.099 (.19) | -.83 | .80 |
| Right temporal pole | 115 | 22 | 82 | -.054 (.06) | -.88 | .80 |
| Right ant. temp cortex | 225 | 58 | 93 | .028 (.08) | .34 | .80 |
| Anterior mPFC | 50 | 8 | 57 | .017 (.12) | .15 | .80 |

636

637 Abbreviations: ROI = Region of interest; fdr = false discovery rate; OFA = occipital face
638 area; FFA = right fusiform face area; pSTS = posterior superior temporal sulcus; TPJ =
639 temporoparietal junction; mPFC = medial prefrontal cortex; ant. Temp. = anterior temporal.

640

641 Note: 'ROI size' is the total number of voxels in each ROI based on data from a face
642 perception localiser or a theory-of-mind localiser. 'Average localiser mask size' is the number
643 of voxels that overlap in more than 50% of participants within each ROI. Right OFA, for
644 example, consists of a 412 voxel ROI, with 84 voxels showing overlap in 93% of participants.
645 Analyses were performed on the subset of voxels in each ROI that show overlap in a majority
646 of participants (>50%).

647

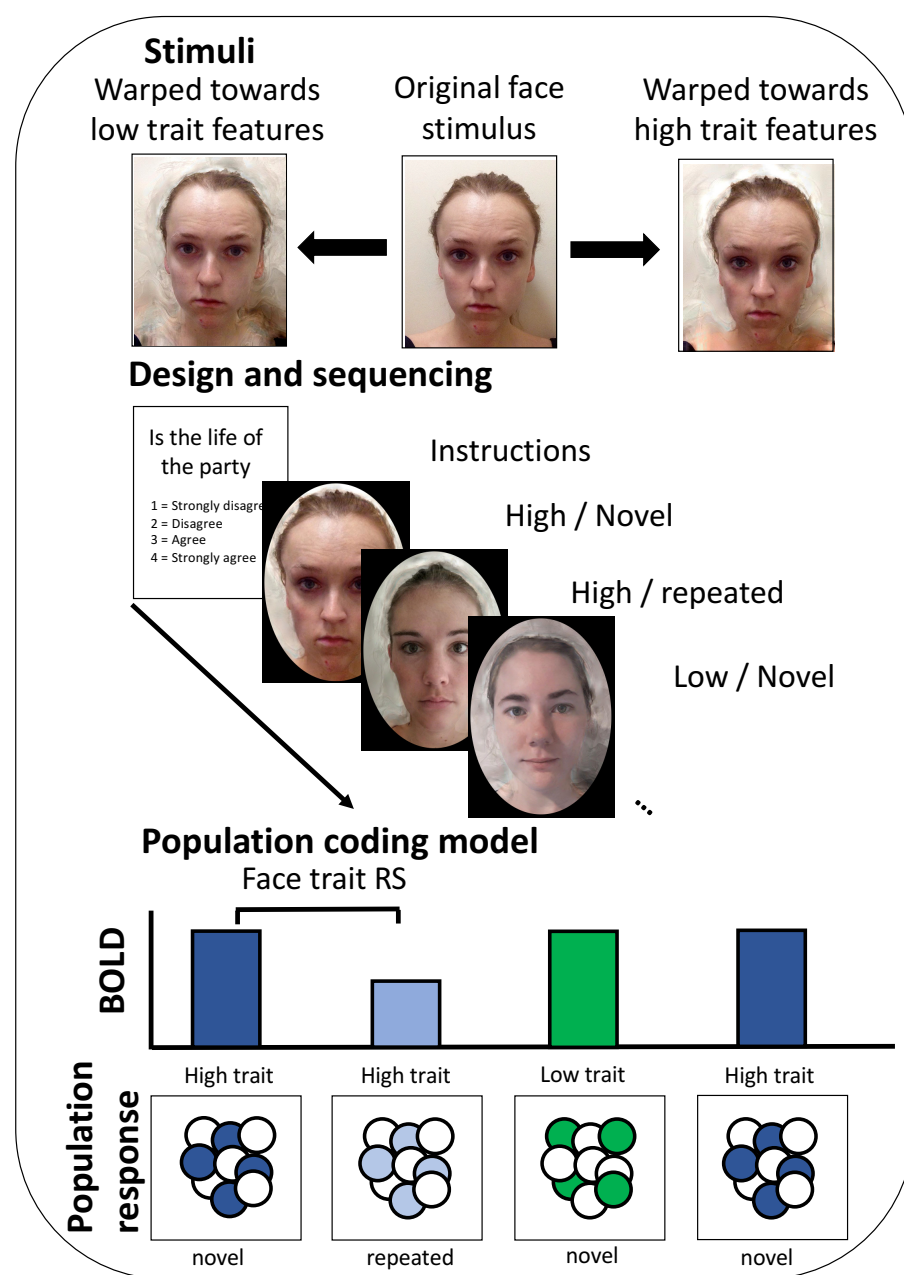648    **Figure 1.** Method and design.



649
650    **Figure 1.** Method and design. A) Individual face images were transformed towards high and
651    low composite templates of trait variables (Extraversion, Agreeableness, Neuroticism,
652    Physical health). The example shown is extraversion. The images used are for illustrative
653    purposes and were not used in the experiment. B) During scanning, each block began with an
654    instruction screen, which provided a statement and a reminder of the ratings scale. On each
655    subsequent trial, participants had to make a judgment based on the face presented. As such, all
656    trials in a mini-block were from the same category (e.g., extraversion), but all trials showed a
657    different individual. C) An illustration of the population coding model of face perception that
658    the repetition suppression design was testing. High and low trait features are presented in blue
659    and green, respectively. Novel and repeated trials are presented in darker and lighter colours,
660    respectively.
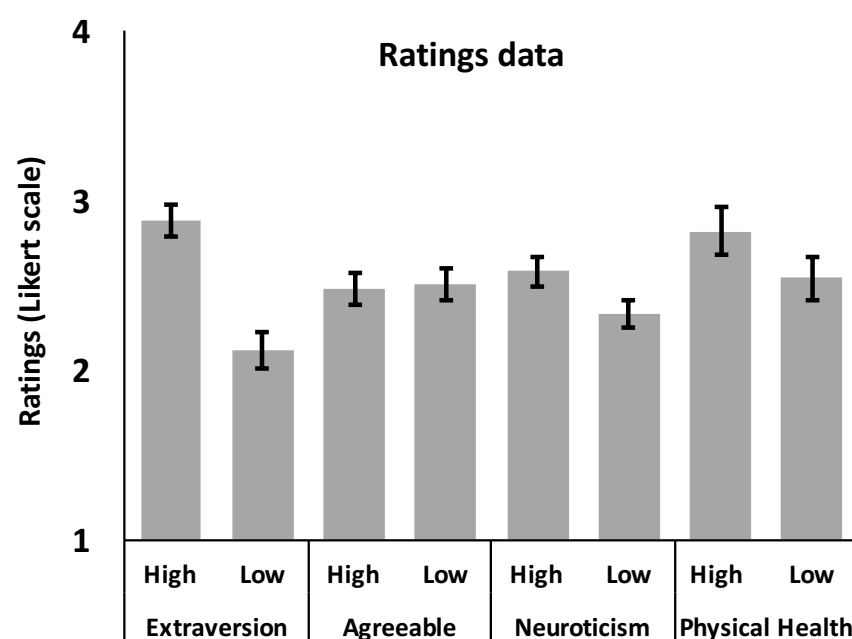
661 **Figure 2.** Ratings data



662

663 **Figure 2.** Mean average face ratings during scanning. Error bars are 95% confidence
664 intervals.

665 **Figure 3.** Percent signal change in our functional regions of interest.
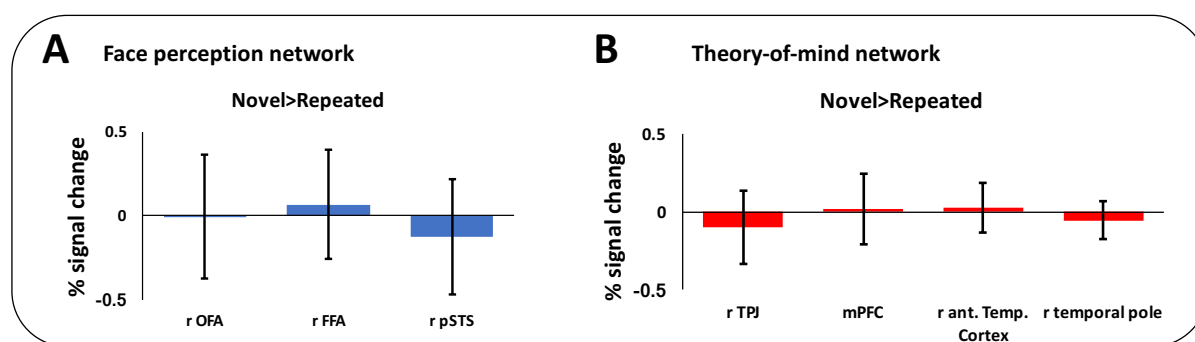


666
667
668 **Figure 3.** Percent signal change for novel compared to repeated trials in the face perception
669 (A) and theory-of-mind network (B). Error bars are standard error of the mean.
670 Abbreviations: r = right; OFA = occipital face area; FFA = right fusiform face area; pSTS =
671 posterior superior temporal sulcus; TPJ = temporoparietal junction; mPFC = medial prefrontal
672 cortex; ant. Temp. = anterior temporal.