1   **Faces and voices in the brain: a modality-general person-identity representation in**

2   **superior temporal sulcus**

3   Abbreviated title: Modality-general identity representation in STS

4

5   Maria Tsantani[1], Nikolaus Kriegeskorte[2], Carolyn McGettigan[3], Lúcia Garrido[1]

6   [1] Division of Psychology, Department of Life Sciences, Brunel University London

7   [2] Zuckerman Mind Brain Behavior Institute, Columbia University

8   [3] Department of Psychology, Royal Holloway, University of London

9

10   Corresponding authors:

11   Maria Tsantani, Division of Psychology, Department of Life Sciences, Brunel University

12   London, Uxbridge UB8 3PH, United Kingdom, Telephone: (+44) 01895268155, E-mail:

13   maria.tsantani@gmail.com

14

15   Lúcia Garrido, Division of Psychology, Department of Life Sciences, Brunel University

16   London, Uxbridge UB8 3PH, United Kingdom, Telephone: (+44) 01895265555, E-mail:

17   garridolucia@gmail.com

18

19   Number of pages: 46

20   Number of figures: 5; Number of tables: 3

21   Number of words: Abstract: 244; Introduction: 650; Discussion: 1467

22

27 **Abstract**

28 Face-selective and voice-selective brain regions have been shown to represent face-identity and

29 voice-identity, respectively. Here we investigated whether there are modality-general person-

30 identity representations in the brain that can be driven by either a face or a voice, and that

31 invariantly represent naturalistically varying face and voice tokens of the same identity.

32 According to two distinct models, such representations could exist either in multimodal brain

33 regions (Campanella and Belin, 2007) or in face-selective brain regions via direct coupling

34 between face- and voice-selective regions (von Kriegstein et al., 2005). To test the predictions

35 of these two models, we used fMRI to measure brain activity patterns elicited by the faces and

36 voices of familiar people in multimodal, face-selective and voice-selective brain regions. We

37 used representational similarity analysis (RSA) to compare the representational geometries of

38 face- and voice-elicited person-identities, and to investigate the degree to which pattern

39 discriminants for pairs of identities generalise from one modality to the other. We found no

40 matching geometries for faces and voices in any brain regions. However, we showed

41 crossmodal generalisation of the pattern discriminants in the multimodal right posterior

42 superior temporal sulcus (rpSTS), suggesting a modality-general person-identity representation

43 in this region. Importantly, the rpSTS showed invariant representations of face- and voice-

44 identities, in that discriminants were trained and tested on independent face videos (different

45 viewpoint, lighting, background) and voice recordings (different vocalizations). Our findings

46 support the Multimodal Processing Model, which proposes that face and voice information is

47 integrated in multimodal brain regions.

48

49

50

51

52 **Significance statement**

53 It is possible to identify a familiar person either by looking at their face or by listening to their

54 voice. Using fMRI and representational similarity analysis (RSA) we show that the right

55 posterior superior sulcus (rpSTS), a multimodal brain region that responds to both faces and

56 voices, contains representations that can distinguish between familiar people independently of

57 whether we are looking at their face or listening to their voice. Crucially, these representations

58 generalised across different particular face videos and voice recordings. Our findings suggest

59 that identity information from visual and auditory processing systems is combined and

60 integrated in the multimodal rpSTS region.

61

62

63

64

65

66

67

68

69

70

71

72

73

74

75

76

## Introduction

Looking at a familiar person's face or listening to their voice automatically grants us access to a wealth of information regarding the person's identity, such as their name, our relationship to them, and memories of previous encounters. Knowledge about how the brain processes faces and voices separately has advanced significantly over the past twenty years: functional magnetic resonance imaging (fMRI) revealed cortical regions that are face-selective (Kanwisher et al., 1997) and regions that are voice-selective (Belin et al., 2000). Recent advances using multivariate classification methods have further shown that some of these regions are important for identification. In particular, face-selective regions in the posterior occipitotemporal lobe, anterior temporal lobe, and posterior superior temporal sulcus (pSTS) can discriminate different facial identities (Kriegeskorte et al., 2007; Nestor et al., 2011; Goesaert and Op de Beeck, 2013; Verosky et al., 2013; Axelrod and Yovel, 2015; Collins et al., 2016; Visconti Di Oleggio Castello et al., 2017). Crucially, a number of studies also found representations in these regions that generalised across different images of the same person (Anzelotti et al., 2014; Anzelotti and Caramazza, 2016; Guntupalli et al., 2017), i.e., were able to "tell people together" (Jenkins et al., 2011; Burton, 2013). Similarly for voices, Formisano et al. (2008) found voice-identity representations in the right STS and Heschl's gyrus that could both discriminate between speakers and generalise across different vowel sounds spoken by the same voice.

Despite these advances, we still have a limited understanding of how the brain combines and integrates face and voice information. Two major models have been put forward. The *Multimodal Processing Model* proposes that there are multimodal systems that process information about people and receive input from both face- and voice-responsive regions (Ellis et al., 1997; Campanella and Belin, 2007). Patient (Ellis et al., 1989; Gainotti, 2011) and fMRI

4

102 studies (Shah et al., 2001; Joassin et al., 2011; Watson et al., 2014a) suggest the anterior

103 temporal lobe, the posterior cingulate cortex, the STS, and the hippocampus as candidate

104 multimodal regions. In contrast, the *Coupling of Face and Voice Processing Model* proposes

105 that the direct coupling between face- and voice-responsive brain regions is crucial for the

106 integration of person-identity information (von Kriegstein et al., 2005). In particular, fMRI

107 studies have shown that voice recognition of familiar (or recently learned) people is associated

108 with increased activation in face-responsive regions of the fusiform gyrus (von Kriegstein et

109 al., 2005, 2006, 2008; von Kriegstein and Giraud, 2006).

110

111 In this study, we tested the predictions from these two models by investigating whether there

112 are modality-general person-identity representations in multimodal regions (Multimodal

113 Model) and/or in face- and voice-selective regions (Coupling Model). We used representational

114 similarity analysis — RSA (Kriegeskorte et al., 2008a, 2008b) to compare the representational

115 geometries of face- and voice-elicited person-identities, and to investigate the degree to which

116 pattern discriminants for pairs of identities generalise from one modality to the other. We

117 predict that, if a region shows a modality-general person-identity representation, the

118 representational geometry of face and voice identities will match, and/or pattern discriminants

119 will generalise across faces and voices. Two recent studies found some support for the

120 Multimodal Model by showing that multimodal regions in the STS and inferior frontal gyrus

121 (Hasan et al., 2016; Anzelotti and Caramazza, 2017) could discriminate between the activation

122 patterns of two face-identities based on voice information (and vice-versa). However, these

123 studies did not show that the regions that could decode identities across modalities could also

124 decode identities within each modality, which is a crucial feature of modality-general person-

125 identity representations. Furthermore, these studies used very few identities and tokens per

126 identity. In our study, we included multiple, naturalistically varying face videos and voice

5

127  recordings of 12 different identities. Thus, we were able to sample the variability of visual and

128  auditory appearance that we are exposed to in everyday life, and to better capture processes of

129  person identification, which are distinct from image or sound recognition (Burton, 2013).

130

## Materials and Methods

131

132  **Overview of study**

133  In this study, we measured fMRI activation patterns in response to the faces and voices of 12

134  famous individuals. It was important to use highly familiar individuals because we needed to

135  guarantee that participants were well acquainted with the faces and voices of those individuals.

136  We thus only recruited participants for the full study if they demonstrated that they were

137  familiar with the majority of the famous individuals in an online Recognition Task. The full

138  study consisted of two MRI scanning sessions and one behavioural session, with each session

139  taking approximately 90 minutes. All three sessions took place on separate days. Before

140  entering the scanner at the start of the first MRI session, participants repeated the Recognition

141  Task in the presence of the experimenter and also completed a Familiarity Task in which they

142  rated all face and voice stimuli on perceived familiarity.

143

144  In each MRI session participants completed three functional runs (main experimental runs) in

145  which they viewed the faces and listened to the voices of the famous people in an event-related

146  design. In addition, participants underwent two structural scans (one in each session) and

147  functional localisers for face-selective, voice-selective, and multimodal regions of interest

148  (ROIs). Across both sessions participants completed at least one run (in most cases two) of (1)

149  the temporal voice area (TVA) localiser (Belin et al., 2000), (2) a face localiser, (3) a

150  multimodal (face-voice) localiser, and (4) a voice localiser. Finally, participants completed a

151  behavioural testing session. In this behavioural session they rated the famous faces and voices

6

152  on various social and perceptual dimensions; however, these results are not included here.

153

154  To investigate the existence of modality-general person-identity representations in each of our

155  ROIs we used RSA (Kriegeskorte et al., 2008a, 2008b; Kriegeskorte and Kievit, 2013) to

156  compare the representational geometry of face-identities with the representational geometry of

157  voice-identities (Analysis A), and to investigate the  degree to which pattern discriminants for

158  each pair of identities generalise from one modality to the other (Analysis B). Analysis A

159  focused on the representational geometry of all of identities, i.e. the entire structure of pairwise

160  distances between the activity patterns elicited by these identities in each modality, and

161  compared geometries across modalities. Analysis B focused on the discriminability of pairs of

162  identities, and used a linear discriminant computed in one modality to test discriminability of

163  the same pair of identities in the other modality (in a similar way to traditional pattern

164  classification methods). These two analyses complement each other and allowed us to test

165  different predictions regarding the nature of modality-general person-identity representations.

166

167  For Analysis A (RSA comparing representational geometries), we predicted that brain regions

168  with modality-general person-identity representations would show matching representational

169  geometries for face-identities and voice-identities. This analysis is constrained by two

170  assumptions. The first assumption is that there is sufficient variability in the representational

171  distances between different identities within-modality, i.e. different degrees of similarity

172  between identities. If all identities are equally distinct from each other, we do not expect to find

173  correlations between geometries across the two modalities. The second assumption is that

174  modality-general information dominates over any modality-specific information that may be

175  present in the same voxels. Specifically, it is possible that the voxels comprising the pattern

176  estimates contain both unimodal and multimodal neurons (Quiroga et al., 2009). In this case,

177  the influence of modality-specific information on the representational distances between all

178  identities could override the influence of modality-general information on the representational

179  geometry, and could result in non-matching representational geometries across modalities.

180

181  We thus also conducted Analysis B (RSA investigating identity discriminability), and we

182  predicted that brain regions with modality-general person-identity representations would be

183  able to discriminate between pairs of identities in one modality based on their representational

184  distance in the other modality. This analysis focuses on one pair of identities at a time, and thus

185  is not affected by the degree of variability in the representational distances between all

186  identities. In addition, this analysis is focused on pattern discriminants that generalise across

187  modalities, and therefore we believe that it is more sensitive to detect modality-general person-

188  identity representations even in the presence of modality-specific information.

189

190  **Participants**

191  Participants were recruited at Royal Holloway, University of London and Brunel University

192  London to take part in a behavioural and fMRI experiment. All participants were required to be

193  native English speakers aged between 18 and 30, and to have been resident in the UK for a

194  minimum of 10 years. These requirements were set to increase the likelihood of participants

195  being familiar with the famous people whose faces and voices were presented in the

196  experiment. In addition, participants completed an online Recognition Task (see below) as part

197  of the screening procedure for the study and were only invited if they were able to recognise at

198  least 75% of our set of famous people from both their face and their voice.

199

200  Thirty-one healthy adult participants were recruited who matched all the above criteria. One

201  participant was excluded from the study after the first MRI session due to excessive head

8

202    movement in the scanner (more than 3 mm in any direction within one run). The final sample

203    consisted of 30 participants (eight males) with mean age of 21.2 years (SD=2.37, range=19-

204    27). All reported normal or corrected-to-normal vision and normal hearing, provided written

205    informed consent and were reimbursed for their participation. The study was approved by the

206    Ethics Committee of Brunel University London.

207

208    **Recognition Task**

209    Participants completed a face and voice Recognition Task to determine whether they could

210    recognise at least 75% of the famous people (i.e. at least 9 out of 12) from both the face and the

211    voice. Face stimuli consisted of single photographs of each of the 12 famous people that were

212    obtained from the Internet through Google Image searches. Photographs included the top part

213    of the body and were front-facing. Voice stimuli consisted of single sound-clips for each of the

214    12 famous people and were obtained from YouTube videos. Sound-clips were approximately

215    8-seconds long and were root-mean-square (RMS) normalized using Praat (version 5.3.80;

216    Boersma and Weenink, 2014; www.praat.org). None of these face or voice stimuli were

217    presented in the main experiment.

218

219    Stimuli were presented using Qualtrics (Qualtrics, Provo, UT). For each stimulus participants

220    had to identify the person shown in the picture or the person speaking (by providing their name

221    or other uniquely identifying biographical information). In the online task participants typed

222    their responses below each stimulus, and in the lab task responses were made verbally.

223

224    **Stimuli for Familiarity Task and Main Experimental Runs**

225    Six silent, non-speaking video clips of moving faces, and six sound clips of voices for each of

226    the 12 famous people (six female, six male) were obtained from videos on YouTube (in total,

9

227    72 stimuli per modality). These people had been identified in our pilot studies as having highly

228    recognisable faces and voices within samples of native English speakers between the ages of

229    18-30 who have been resident in the UK for a minimum of 10 years. This list of famous people

230    included actors, pop stars, politicians, comedians, and TV personalities: Alan Carr, Beyonce

231    Knowles, Daniel Radcliffe, Emma Watson, Arnold Schwarzenegger, Barack Obama, Sharon

232    Osbourne, Kylie Minogue, Graham Norton, Cheryl Cole, Barbara Windsor, and Jonathan Ross.

233

234    The face stimuli were selected so that the background did not provide any cues to the identity

235    of the person. Other than the absence of speech, there were no constraints on the type of face

236    movement. Examples of face movements included nodding, smiling, and rotating the head.

237    However, all stimuli were selected to be primarily front-facing. Face stimuli were edited using

238    Final Cut Pro X (Apple, Inc.) so that they were three seconds long and centred on the bridge of

239    the nose. Six video-clips of the face of the same person were obtained from different original

240    videos set in a different background.

241

242    Voice stimuli were edited using Audacity® 2.0.5 recording and editing software

243    (RRID:SCR_007198) so that they contained three seconds of speech after removing long

244    periods of silence. Voice stimuli were converted to mono with a sampling rate of 44100, low-

245    pass filtered at 10KHz, and RMS normalised using Praat. Six sound clips of the voice of the

246    same person were obtained from different original videos. All of the voice stimuli had a

247    different verbal content and were non-overlapping. The stimuli were selected so that the

248    speakers' identity could not be determined based on the verbal content, conforming to the

249    standards set by Van Lancker et al. (1985) and Schweinberger et al. (1997).

250

251    **Familiarity Task**

10

252   Before entering the scanner, participants rated all stimuli that would be presented in the main

253   experimental runs on perceived familiarity. Participants were presented with the face stimuli

254   first, followed by the voice stimuli, in separate blocks. Stimuli were presented using the

255   Psychophysics Toolbox (version 3; RRID:SCR_002881; Brainard, 1997; Pelli, 1997) running

256   in Matlab (version R2013b; MathWorks; RRID:SCR_001622). Face stimuli were presented in

257   the centre of the screen. Participants listened to the voice stimuli through headphones

258   (Sennheiser HD 202). Participants rated each stimulus on scale from 1 (very unfamiliar) to 7

259   (very familiar). Each block took approximately 5 minutes to complete.

260

261   **MRI data acquisition**

262   Participants were scanned using a 3.0 Tesla Tim Trio MRI scanner (Siemens, Erlangen) with a

263   32-channel head coil at the Combined Universities Brain Imaging Centre (CUBIC) at Royal

264   Holloway, University of London. In each of the two scanning sessions, a whole-brain T1-

265   weighted anatomical scan was acquired using magnetization-prepared rapid acquisition

266   gradient echo (MPRAGE) [1.0 x 1.0 in-plane resolution; slice thickness, 1.0mm; 176 axial

267   interleaved slices; PAT, Factor 2; PAT mode, GRAPPA (GeneRalized Autocalibrating

268   Partially Parallel Acquisitions); repetition time (TR), 1900ms; echo time (TE), 3.03ms; flip

269   angle, 11°; matrix, 256x256; field of view (FOV), 256mm].

270

271   For all functional runs T2*-weighted whole-brain functional scans were acquired using echo-

272   planar imaging (EPI) [3.0 x 3.0 in-plane resolution; slice thickness, 3.0mm; PAT, Factor 2;

273   PAT mode, GRAPPA (GeneRalized Autocalibrating Partially Parallel Acquisitions); 34

274   sequential (descending) slices; repetition time (TR), 2000ms; echo time (TE), 30ms; flip angle,

275   78°; matrix, 64x64; field of view (FOV), 192mm]. For the majority of participants, slices

276   covered all parts of the brain except for the most dorsal part of parietal cortex. In each

11

277    experimental run we obtained 293 brain volumes, in the TVA localiser we obtained 251 brain

278    volumes, and in each run of the face, voice, and multimodal localiser runs we obtained 227

279    brain volumes.

280

281    **fMRI data pre-processing**

282    Data were pre-processed using Statistical Parametric Mapping (SPM12; Wellcome Department

283    of Imaging Science, London, UK; RRID:SCR_007037; http://www.fil.ion.ucl.ac.uk/spm)

284    operating in Matlab. Pre-processing was performed separately for each scanning session. All

285    runs within each session (main experiment or localizer runs) were pre-processed together. The

286    first three EPI images in each run (dummy scans) were discarded to allow for T1-equilibration

287    effects. Images were slice-time corrected based on the middle slice in each volume and then

288    realigned to correct for head movement based on the first image. The structural image in native

289    space was then coregistered with the realigned mean functional image and segmented into grey

290    matter, white matter, and cerebrospinal fluid. No smoothing was performed on the images from

291    the experimental runs. Functional images from the localiser runs were smoothed with a 4-mm

292    Gaussian kernel (full width at half maximum).

293

294    After separate pre-processing of the images in each session, images from the second scanning

295    session were realigned to the structural image from the first session. Specifically, the structural

296    image from session two was coregistered to the structural image from session one, and the

297    transformation was then applied to all functional images from session two. As a result, all

298    functional images were in the same space.

299

300    **Functional localiser runs**

301    **TVA localiser.** We used the TVA localiser developed by Belin et al. (2000) which contains

12

302    vocal and non-vocal auditory stimuli. Stimuli were presented in 40 blocks of 8 seconds each.

303    Vocal stimuli were presented in 20 blocks and included speech and non-speech vocalisations

304    obtained from 47 speakers (Pernet et al., 2015).  Non-vocal stimuli were presented in 20 blocks

305    and consisted of industrial sounds, environmental sounds, and animal vocalisations. Within

306    each block stimuli were presented in a random order that was fixed across participants.

307    Participants were instructed to close their eyes and focus on the sounds. The TVA localiser was

308    presented directly after the main experimental runs. The duration of a single run was

309    approximately 10 minutes.

310

311    **Face, Voice, and Multimodal localisers.** We created new face, multimodal, and voice

312    localiser runs that shared the same experimental design and presented stimuli from comparable

313    categories (people and objects/scenes). Importantly, we used videos and not static images of

314    faces. Dynamic face stimuli have been shown to be more effective that static face stimuli for

315    localising face-selective regions (Fox et al., 2009; Pitcher et al., 2011). Stimuli used for the

316    face localiser were silent, non-speaking video clips of famous and non-famous (French

317    celebrities unknown to our participants) moving faces, and silent video clips of moving large

318    objects and natural or manmade visual scenes (such as videos of airplanes, trains, traffic,

319    rainforests, waves on a beach) obtained from videos on YouTube. For the multimodal localiser

320    the stimuli were audio-visual and included videos clips of the faces of famous and non-famous

321    people speaking, and video clips of moving large objects and natural or manmade scenes (same

322    categories as above). For the voice localiser we presented voice clips of famous and non-

323    famous people, and sound clips of manmade or natural environmental sounds (same categories

324    as used in the other two types of localisers), with no video.

325

326    Videos (640 x 360 pixels) were presented at the centre of the screen. The screen resolution was

13

327 1024 x 768 pixels, and from a distance of 85 cm, the videos subtended 20.83 x 12.27 degrees

328 of visual angle. Audio stimuli were presented via MR-compatible earbuds (S14; Sensimetrics

329 Corp.), which participants used for each entire scanning session. Each stimulus lasted 8

330 seconds and each run presented 48 stimuli. Stimuli were presented in pairs (24 pairs) showing

331 the same person (such as two videos of Brad Pitt) or the same category of objects or scenes

332 (such as two videos of trains). Eight pairs showed stimuli from famous people, eight pairs

333 showed stimuli from non-famous people, and eight pairs showed object/scene stimuli.

334 Participants were encouraged to always fixate at the centre of the screen. Participants

335 performed a one-back task in which they had to detect the exact same stimulus repetition

336 within each pair, which occurred in approximately 15% of the trials. A 16-second period of

337 fixation was presented at the end of each run and twice in the middle of each run (every 16

338 trials).

339

340 The order of the face, voice, and multimodal localisers was counterbalanced across

341 participants. For participants who completed two runs of each localiser, different identities

342 were presented on each run. The duration of each localiser run was approximately 8 minutes.

343

344 **General linear models.** To identify face-selective (face localiser), voice-selective (voice

345 localiser and TVA localiser), and people-selective (multimodal localiser) brain regions, we

346 computed mass univariate time-series models for each participant. Regressors modelled the

347 blood-oxygenation-level-dependent (BOLD) response following the onset of the stimuli and

348 were convolved with a canonical hemodynamic response function (HRF). We also used a high-

349 pass filter cutoff of 128 seconds, and autoregressive AR(1) model to account for serial

350 correlations. For the face, voice, and multimodal localisers there were three experimental

351 regressors: (1) famous faces/voices/people, (2) non-famous faces/voices/people, and (3) objects

14

352    and scenes. For the TVA localiser there were two experimental regressors: (1) voices and (2)

353    non-voices. For all localisers six head motion parameters computed during realignment were

354    included as covariates. Selectivity was defined with a *t*-test contrasting the responses to

355    faces/voices/people (famous and non-famous) *versus* responses to the control stimuli.

356

357    **ROI definition.** We used probabilistic maps from previous studies to define regional masks in

358    which we predicted that our regions of interest (ROIs) would be located.  We then defined

359    ROIs by extracting all selective voxels within those regional masks for each participant. This

360    approach is similar to the one implemented by  Julian et al. (2012)  and avoids experimenter

361    biases in ROI definition.

362

363    Probabilistic maps were thresholded to only show voxels that were present in 20% of the

364    participants and binarised to create regional masks. We used a probabilistic map of the TVAs

365    created    by    Pernet    et    al.    (2015)    and    obtained    from    neurovault

366    (http://neurovault.org/images/106/) to create separate masks for the right and left TVA (rTVA,

367    lTVA). For all other regional masks, we used probabilistic maps that were obtained from a

368    previous study conducted in the lab (unpublished data). In this previous study we tested 22

369    participants using the same face and voice localisers as the current study (we did not use the

370    multimodal localiser in this previous study). We defined face-selective and voice-selective *t*-

371    test images for each participant, thresholded each image at *p*<.05 (uncorrected), binarised the

372    resulting image, and summed all images across participants to create face-selective and voice-

373    selective probabilistic maps. In cases where there was some overlap between the masks for

374    different regions we manually defined the borders of these masks using anatomical landmarks.

375

376    Regional masks of face-selective regions were created for the right fusiform face area (rFFA),

15

377   the right occipital face area (rOFA), and the right posterior superior temporal sulcus (rpSTS).

378   Regional masks of voice-selective regions were created for the right and the left superior

379   temporal sulcus and gyrus (rSTS/STG, lSTS/STG). Regional masks of multimodal regions

380   were created based on joint face-selective and voice-selective probabilistic maps. These masks

381   were created for a number of regions that showed *both* face-selective and voice-selective

382   responses in most participants: precuneus/posterior cingulate, orbitofrontal cortex (OFC),

383   frontal pole (FP), and right and left temporal pole with anterior inferior temporal cortex (rTP-

384   aIT, lTP-aIT) — we considered the TP and aIT together as the peaks were difficult to separate

385   in most participants. We did not create a mask of the multimodal STS using this method due to

386   the voice-selective STS region being much larger than the face-selective STS region. However,

387   there was large overlap between our mask of the face-selective rpSTS and our masks of the

388   rSTS/STG and rTVA, suggesting that this face-selective rpSTS region also responds to voices.

389

390   All of the regional masks (in MNI space) were registered and resliced to each participant's

391   native space using FSL (version 5.0.9; RRID:SCR_002823; Jenkinson et al., 2012). These

392   masks were then used to extract ROIs from the *t*-test maps obtained from the contrasts of

393   interest from the face, voice, TVA, and multimodal localisers from the current study. All

394   voxels that fell within the boundaries of the mask and that were significantly activated at

395   $p<.001$ (uncorrected) were included in the subject-specific ROI. If there were fewer than 30

396   voxels at $p<.001$ the threshold was lowered to $p<.01$ or $p<.05$. If we could not define 30

397   selective voxels even at $p<.05$, the ROI for that participant was not included in the analyses.

398   We required that all ROIs be present in at least 20 participants (out of 30).

399

400   **Main experimental runs: Experimental design and statistical analysis**

401   **Design and procedure.** Face and voice stimuli were presented using the Psychophysics

16

402     Toolbox via a computer interface inside the scanner. Face and voice clips of all 12 identities

403     were intermixed within each run. A fixation point was always present and participants were

404     asked to fixate. The videos were 640 x 360 pixels and, from a viewing distance of 85cm,

405     videos subtended 20.83 x 12.27 degrees of visual angle. The six face videos and the six voice

406     recordings for each of the 12 identities were evenly distributed among the three runs so that

407     each run contained two different videos of the face and two different recordings of the voice of

408     each identity. Each individual stimulus was presented twice within each run. Therefore, in each

409     run there were 96 experimental trails (48 face trials, 48 voice trials) in total.

410

411     Participants performed an anomaly detection task that involved pressing a button when they

412     saw or heard a novel famous person that was not part of the set of the 12 famous people that

413     they had been familiarised with prior to entering the scanner. Therefore, each run also

414     contained 12 task trials presenting six famous faces and six famous voices that were not part of

415     the set of famous people that the participants had been familiarised with.

416

417     Stimuli were presented in a pseudorandom order that ensured that within each modality each

418     identity could not be preceded or succeeded by one of the other identities more than once, and

419     that each stimulus could not be succeeded by a repetition of the exact same stimulus. Face and

420     voice clips were presented for three seconds with a SOA of four seconds. Thirty-six null

421     fixation trials were added to each run (~25% of the total number of trials). Thus, each run

422     contained 144 trials in total and lasted approximately 10 minutes.

423

424     The presentation order of the three runs was counterbalanced across participants. The same

425     three runs with the same face videos and voice recordings that were presented in scanning

426     session one were also presented in session two. However, the three runs were presented in

17

427   different orders in both sessions (counterbalanced across participants) and stimuli within each

428   run were presented in a new pseudorandom sequence. As an exception, the stimuli for the task

429   trials were different in the two sessions in order to maintain their novelty.

430

431   **General linear models.** We computed mass univariate time-series models for each participant.

432   Models were defined separately for each scanning session and each experimental run (six runs

433   in total). Regressors modelled the BOLD response following the onset of the stimuli and were

434   convolved with a canonical hemodynamic response function (HRF). We also used a high-pass

435   filter cutoff of 128 seconds and autoregressive AR(1) model to account for serial correlations.

436   The 12 different identities in each modality were entered as separate regressors in the model

437   (i.e. 24 regressors). Each of these regressors included the two different face videos and voice

438   recordings of each identity that were presented in the run, as well as the two repetitions of each

439   stimulus. Task trials and six head motion parameters computed during realignment were

440   included as regressors of no interest.

441

442   As part of the crossvalidation procedure used in the RSA analyses described below, separate

443   models were estimated for each partition of each crossvalidation fold, thus resulting in

444   parameter estimates and residual time courses for every possible independent partition. For

445   partitions with two runs, data was concatenated before estimating the model.  In the analyses

446   described below we used the beta estimates computed at each voxel of each ROI for each of

447   the 24 experimental conditions (12 face-identities and 12 voice-identities).

448

449   **Mean response to faces and voices in ROIs.** We conducted an analysis to characterise the

450   responses to faces and voices in each ROI, and to confirm that each ROI showed the expected

451   responsivity to faces and voices. For this analysis, we calculated the mean (across all voxels in

18

452   each ROI, and across all runs) of the parameter estimates for the 12 face-identities and the

453   mean of the parameter estimates for the 12 voice-identities. For each ROI we tested whether

454   the mean for faces and mean for voices were significantly different from zero (across

455   participants) using one-sample $t$-tests. P values were corrected for 24 comparisons (2 tests x 12

456   ROIs) controlling the false discovery rate (FDR), with $q<.05$. We also compared the mean for

457   faces with the mean for voices in each ROI using paired $t$-tests. P values were corrected for

458   multiple comparisons (12 comparisons) using FDR with $q<.05$.

459

460   **Analysis A: RSA comparing representational geometries**. For this analysis we computed

461   representational dissimilarity matrices (RDMs) for faces and voices (12x12 matrices)

462   separately for each participant, each scanning session and each ROI. We then computed the

463   correlations between pairs of these RDMs (for an example, see Figure 4). These analyses were

464   performed using in-house Matlab code and the RSA toolbox (Nili et al., 2014). To compute the

465   RDMs we used the linear discriminant contrast (LDC), a crossvalidated distance measure (Nili

466   et al., 2014; Walther et al., 2016). For each ROI, each modality (i.e. faces and voices

467   separately), and each scanning session, we calculated the LDC between the pattern estimates

468   (beta estimates across all voxels within an ROI) elicited by the different identities. The

469   resulting 12x12 matrices were symmetric around a diagonal of zeros (Figure 4). Each cell in

470   the RDMs showed the discriminability of the pattern estimates corresponding to a pair of

471   identities in the chosen modality and ROI.

472

473   RDMs were computed using leave-one-run-out crossvalidation across the three runs in each

474   session (each run presented the same identities with different stimuli). In each of three

475   crossvalidation folds, the pattern estimates for each identity were computed with data from two

476   runs (partition one) and separately from the pattern estimates from the remaining run (partition

19

477    two). The pattern estimates from each pair of identities from partition one were used to obtain a

478    linear discriminant, which was then applied to differentiate the activity patterns of the same

479    identity pairs in partition two (Nili et al., 2014; Walther et al., 2016). We applied multivariate

480    noise normalisation by computing a noise variance-covariance matrix based on the residual

481    time courses obtained from the model that was estimated with data from partition one. More

482    specifically, to compute the LDC for each pair of identities we first multiplied the contrast

483    between the patterns of a pair of identities in partition one (the discriminant weights) by the

484    inverse of the noise variance-covariance matrix (after regularisation using the optimal

485    shrinkage method: Ledoit and Wolf, 2004), and transformed the resulting weights to unit

486    length. We then computed the dot product between the resulting vector and the vector with the

487    contrast between the patterns of the same pair of identities from partition two (Carlin and

488    Kriegeskorte, 2017), which resulted in a single value showing the discriminability of those

489    identities.  Finally, the resulting RDMs with LDC values from each crossvalidation fold were

490    averaged to create one RDM per scanning session. This procedure resulted in four RDMs per

491    participant per ROI: faces session 1, voices session 1, faces session 2, and voices session 2

492    (Figure 4).  Crossvalidating across runs with different videos of the face and recordings of the

493    voice of each identity ensured that the resulting RDMs represented face- and voice-*identity*,

494    rather than specific face videos and voice recordings.

495

496    In order to compare the representational geometries of the face- and voice-identities, the RDMs

497    for each participant were compared across the two scanning sessions using Pearson's

498    correlation coefficient (Figure 4). We also compared the representational geometries of face

499    and voice-identities within modality across two scanning sessions in order to investigate the

500    stability of the representational geometries across the two scanning sessions. For the

501    *crossmodal comparisons* we compared the face and voice RDMs from session one with the

20

502     RDMs of the *other* modality in session two (i.e. faces session 1 vs. voices session 2, and voices

503     session 1 vs. faces session 2). For the *unimodal comparisons* we compared the face and voice

504     RDMs from session one with RDMs of the *same* modality in session two (i.e. faces session 1

505     vs. faces session 2 and voices session 1 vs. voices session 2). At the group level for each ROI

506     we compared the single-subject correlations for each of the four comparisons (two crossmodal,

507     two unimodal) against zero using one-sample one-tailed Wilcoxon signed-rank tests (because

508     correlations are not normally distributed). P values were corrected for multiple comparisons

509     (48 comparisons: 4 tests x 12 ROIs) controlling for FDR with $q < .05$.

510

511     **Analysis B: RSA investigating identity discriminability.** For this analysis we computed

512     crossmodal RDMs separately for each participant, each scanning session and each ROI. We

513     used the activity patterns of identity pairs in one modality to create a linear discriminant and

514     then applied the discriminant to the activity patterns of the same identity pairs in the other

515     modality. With this exception, the crossvalidation procedure was identical to the procedure for

516     creating face and voice RDMs for the previous analysis. Two crossmodal RDMs for each ROI

517     were computed using this method: one by applying a linear discriminant based on face data to

518     voice data, and one by applying a linear discriminant based on voice data to face data. The

519     LDC provides a continuous measure of discriminability for each pair of stimuli (Nili et al.,

520     2014; Walther et al., 2016; Carlin and Kriegeskorte, 2017). Importantly, under the null

521     hypothesis the LDC is symmetrically distributed around zero, and thus unbiased. By

522     calculating the mean LDC value across all cells in an RDM for a certain ROI we can determine

523     the overall ability of that ROI to discriminate between identities. Mean LDC values for all

524     participants can then be subjected to random-effects inference comparing against zero.

525     Therefore, we predicted that crossmodal RDMs for regions with modality-general person-

526     identity representations would show mean LDC values that are significantly greater than zero.

527

528     In addition to investigating identity discrimination *across modalities* using crossmodal RDMs,

529     we also investigated the ability of each ROI to discriminate between identities *within modality,*

530     using the face and voice RDMs that were created in the previous analysis. We predicted that

531     face or voice RDMs for regions that represent face or voice identity, respectively, would show

532     mean LDC values that are significantly greater than zero.

533

534     For this analysis the corresponding RDMs (e.g. faces session 1 and faces session 2) for each

535     scanning session were averaged across the two sessions, and then the mean LDC across the

536     vectorised matrix was calculated. Thus, for each participant and each ROI we obtained four

537     mean LDC values representing (1) face discriminability, (2) voice discriminability, (3a)

538     crossmodal discriminability - face discriminant generalised to voices, and (3b) crossmodal

539     discriminability - voice discriminant generalised to faces. For each ROI and each type of

540     discriminability we entered participants' LDC values into a one-sample one-tailed *t*-test

541     comparing them against zero. P values were corrected for all comparisons (48 comparisons: 4

542     tests x 12 ROIs) controlling for FDR with *q*<.05.

543

544     **Code and data accessibility.** All the code and data for the above analyses will be made

545     available after publication.

546

547     **Exploratory whole-brain searchlight analyses.** Despite including a broad range of

548     functionally defined ROIs, it is possible that modality-general person-identity representations

549     may exist in brain regions not included in our ROIs. Specifically, these representations may

550     exist in brain regions that are not face-selective or voice-selective. Therefore, we used an

551     exploratory whole-brain searchlight analysis to identify potential brain regions with person-

552    identity representations using the same methods as in our main ROI analyses. We note that we

553    focused solely on modality-general person-identity representations in this exploratory analysis,

554    as that was the main aim of this study.

555

556    For each participant we created 6mm radius spheres centred on each voxel within a grey-matter

557    mask of their brain (obtained from the segmentation procedure) using the RSA toolbox (Nili et

558    al., 2014) in Matlab. A 6 mm radius resulted in a searchlight sphere of 33 voxels, which

559    matched our requirement for minimum ROI size of 30 voxels in the main analyses. For the

560    analysis comparing representational geometries we computed a face and a voice RDM in each

561    searchlight sphere, averaging the RDMs from both scanning sessions, and then calculated the

562    Pearson correlation between them. Correlations were Fisher z-transformed. The output of this

563    analysis was a whole-brain map of Fisher-transformed correlation coefficients for each

564    participant. For the second analysis investigating identity discriminability we computed a

565    single crossmodal RDM in each searchlight sphere by averaging the crossmodal face-voice

566    RDM with the crossmodal voice-face RDM, and then calculating the mean LDC across the

567    resulting matrix in vector form. The output for each participant was a whole-brain map of mean

568    LDC values.

569

570    The whole-brain searchlight maps from each analysis were normalised to MNI space using the

571    normalisation parameters generated during the segmentation procedure and spatially smoothed

572    with 9-mm Gaussian kernel (full width at half maximum) to correct for errors in intersubject

573    alignment. For group-level analysis, all searchlight maps were entered into a one-sample *t*-test

574    to determine whether the correlation coefficient/mean LDC value was significantly greater than

575    zero at each voxel. We used the randomise tool (Winkler et al., 2014) in FSL for inference on

576    the resulting statistical maps (5,000 sign-flips). Clusters were identified with threshold-

23

577 free cluster enhancement (TFCE), and p-values were corrected for multiple comparisons (FWE

578 < 0.05).

579

580 ## **Results**

581 **Familiarity ratings**

582 Familiarity ratings of both faces and voices were high (Faces: $M = 6.28$, $SD = 0.5$; Voices: $M =$

583 6.2, $SD = 0.49$). Average familiarity of each identity's face and voice are shown in Table 1.

584

585 **Table 1.** Familiarity ratings of the face and voice of each identity averaged across participants showing the mean
586 ($M$) rating of the face and voice of each identity across all face videos and all voice recordings of that identity, the
587 standard deviation ($SD$) of participants' ratings of each identity, and the range of mean ratings for the six face
588 tokens and six voice tokens for each identity. The rating scale ranged from 1 (very unfamiliar) to 7 (very familiar).

| | | AC | AS | BO | DR | GN | JR | BK | BW | CC | EW | KM | SO |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Faces | M | 6.48 | 5.72 | 6.94 | 6.76 | 6.27 | 6.42 | 6.36 | 5.65 | 6.49 | 6.73 | 5.25 | 6.32 |
| | SD | 0.75 | 1.19 | 0.2 | 0.59 | 0.85 | 0.59 | 0.94 | 1.45 | 0.79 | 0.45 | 1.48 | 0.92 |
| | Token range | 6.37-6.60 | 4.83-6.13 | 6.90-7 | 6.67-6.87 | 5.87-6.5 | 6.17-6.57 | 6.07-6.57 | 5.33-5.9 | 6.37-6.60 | 6.47-6.9 | 4.6-5.57 | 6.17-6.47 |
| Voices | M | 6.59 | 5.66 | 6.73 | 6.69 | 6.37 | 6.54 | 6.23 | 5.54 | 6.63 | 6.07 | 5.3 | 6.02 |
| | SD | 0.54 | 1.48 | 0.63 | 0.57 | 0.77 | 0.71 | 1.04 | 1.74 | 0.74 | 0.94 | 1.45 | 1.03 |
| | Token range | 6.37-6.83 | 5.43-5.87 | 6.7-6.8 | 6.57-6.8 | 6.07-6.67 | 6.3-6.7 | 6.07-6.37 | 5.47-5.67 | 6.4-6.77 | 4.6-6.53 | 5-5.53 | 5.5-6.43 |

589

590

591 **ROI definition**

592 Using functional localisers we defined face-selective ROIs (rFFA, rOFA, rpSTS), voice-

593 selective ROIs (rSTS/STG, rTVA, lSTS/STG, lTVA), and multimodal ROIs (OFC, FP, rTP-

594 aIT, lTP-aIT, Prec./P.Cing. [including the retrosplenial cortex]) in each participant. We were

595 able to localise these ROIs with at least 30 voxels in all 30 participants, except for the face-

596 selective rFFA (28 participants) and rOFA (29 participants), the Prec./P.Cing. (26

597 participants), and the OFC (21 participants). We note that the voice-selective ROIs in the right

598 hemisphere (rTVA, rSTS/STG) overlap with each other and with the face-selective rpSTS and

24

599    the multimodal rTP-aIT ROIs. In addition, the voice-selective ROIs in the left hemisphere

600    (lTVA, lSTS/STG) overlap with each other and with the multimodal lTP_aIT ROI. For

601    visualisation purposes only, probabilistic maps of all ROIs were created by normalising the

602    single subject ROIs to MNI space and summing them. Figure 1 shows those maps thresholded

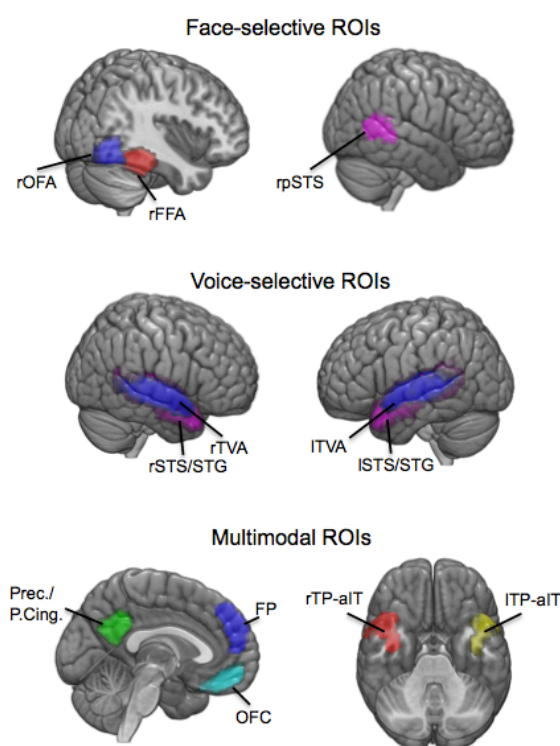603    to display all voxels that were present in at least 20% of the participants.

604



**Figure 1: Face-selective, voice-selective, and multimodal ROIs**. Location of ROIs that resulted from the face, voice, and multimodal localisers in MNI space.

r = right, l = left, FFA = fusiform face area, OFA = occipital face area, pSTS = posterior superior temporal sulcus, STS/STG = superior temporal sulcus/superior temporal gyrus, TVA = temporal voice area, OFC = orbitofrontal cortex, FP = frontal pole, TP = temporal pole, aIT = anterior inferior temporal cortex, Prec = precuneus, P.Cing. = posterior cingulate.

**Mean response to faces and voices in ROIs**

615    In order to confirm that each ROI showed the expected responsiveness to faces and voices, we

616    computed the regional mean of the parameter estimates for faces and for voices across

617    participants for each ROI and modality (Figure 2). As expected, mean beta values for faces

618    were high and significantly greater than zero in all three face-selective ROIs (all one-sample *t*-

25

619    tests with *p*<.0001). Mean beta values for voices were also significantly greater than zero in

620    the rFFA (*p*<.0001) and rpSTS (*p*<.0001), but not in the rOFA. The rFFA and the rOFA

621    showed significantly greater responses to faces compared with voices (both paired-samples *t*-

622    tests with *p*<.0001). In contrast, the rpSTS showed significantly greater responses to voices

623    compared with faces (*p*=.0002) despite being defined using our face localiser. This is most

624    likely due to the large overlap between this ROI and the voice-selective rSTS/STG and rTVA

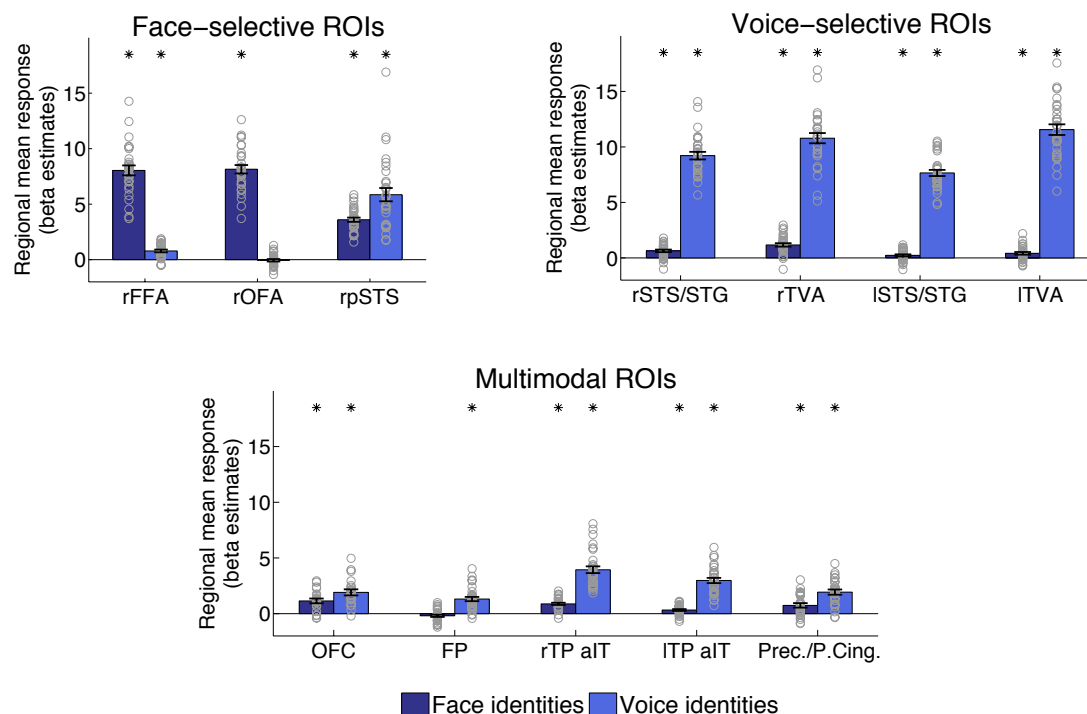625    ROIs. This finding demonstrates that the rpSTS also showed substantial responses to voices.

626



627
628
629    **Figure 2: Regional mean responses to faces and voices in ROIs**. Regional mean responses for all face-identities
630    and for all voice-identities in face-selective, voice-selective, and multimodal ROIs (mean beta estimates across all
631    voxels of each ROI, and across all runs). Bars show mean responses across participants, error bars show standard
632    error, and grey circles show individual participants. We tested whether mean responses were significantly greater
633    than zero using one-sample *t*-tests across all 30 participants, and stars show significant results at *p*≤.0209 (FDR
634    corrected for all 24 comparisons). We also tested whether mean beta values for faces were significantly different
635    from mean beta values for voices in each ROI using paired *t*-tests across all participants. In all ROIs mean beta
636    values for faces and voices were significantly different at *p*≤.0011 (FDR corrected for all 12 ROIs).
637

638    It could be that the responses to voices in rpSTS were due to the voices being familiar, and not

639    because of being voices *per se*. To determine whether this region responded to voices more

26

640   generally or just to familiar voices, we investigated the responses in rpSTS to familiar voices,

641   unfamiliar voices, and non-voices during the functional voice localisers. For each participant,

642   we calculated the mean parameter estimates across all voxels of the face-selective rpSTS for

643   each condition of the voice localiser (familiar voices, unfamiliar voices, and auditory scenes)

644   and of the TVA localiser (vocal and non-vocal sounds). For the voice localiser, both the

645   familiar and the unfamiliar voices had significantly higher parameter estimates than the

646   auditory scenes (both $p<.0001$). For the TVA localiser, the rpSTS also showed significantly

647   higher responses to voices than non-voices ($p<.0001$). These results show that the face-

648   selective rpSTS also responds to voices in general and not only familiar voices (for similar

649   results, see Deen et al., 2015), and therefore in the rest of this article we will refer to this rpSTS

650   region as displaying multimodal responses.

651

652   Returning to the analysis of the parameter estimates for faces and voices during the main

653   experimental runs, the mean beta values for voices were significantly greater than zero for all

654   four voice selective ROIs (all $p<.0001$). Mean beta values for faces were also significantly

655   greater than zero for all voice-selective ROIs (all $p≤.0209$), but the parameter estimates were

656   significantly lower than for voices (all $p<.0001$).

657

658   For the multimodal ROIs mean beta values for faces and for voices were significantly greater

659   than zero in all ROIs (all $p≤.0009$) except the frontal pole for faces. This result demonstrates

660   that, although we still included the frontal pole ROI in the main analyses, we cannot be

661   confident about the multimodal responses of this ROI. Also, we note that in all multimodal

662   ROIs (OFC, FP, rTP-aIT, lTP-aIT, Prec./P.Cing.) mean beta values for voices were

663   significantly higher than mean beta values for faces (all $p≤.0011$). We observed this

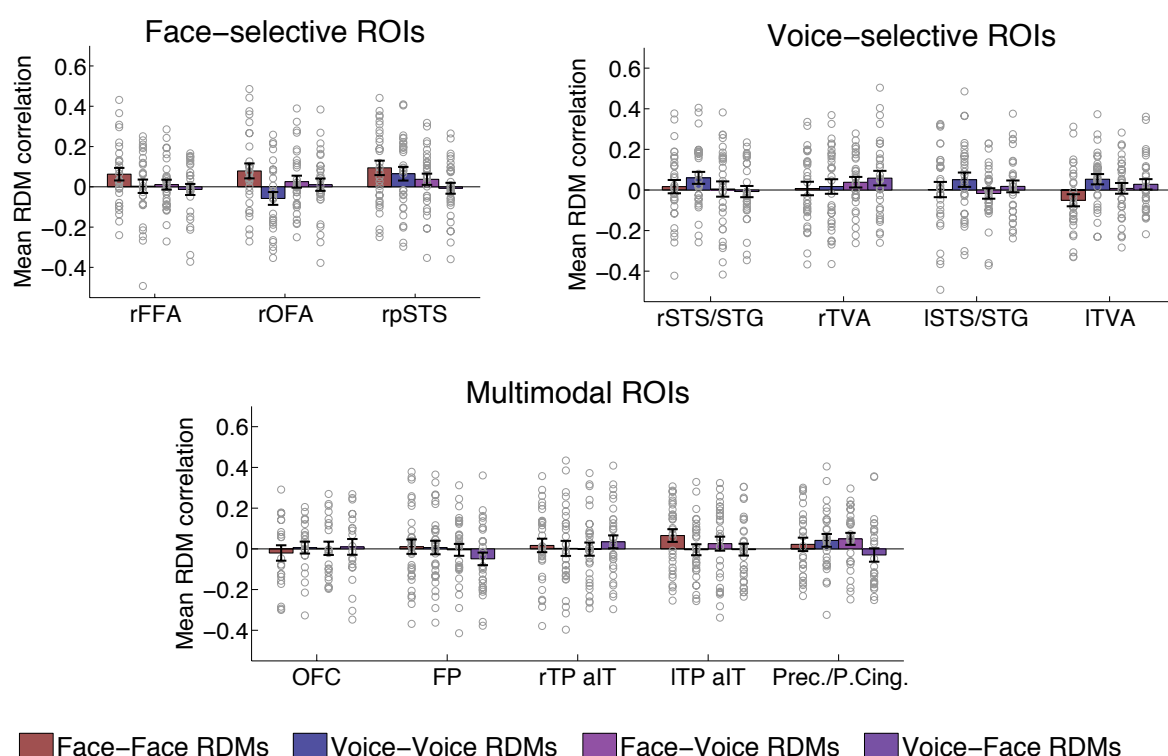664   consistently across all participants.

665

**Analysis A: RSA comparing representational geometries**

Our first main analysis compared the representational geometry of the 12 famous identities across and within modalities in each ROI. We computed face and voice RDMs separately for each session using the LDC and compared the RDMs using Pearson correlation (Figures 3 & 4). We then tested whether these correlations were significantly above zero.

671



**Figure 3: Results of RSA comparing representational geometries**. Comparisons between the representational distance matrices (RDMs) from two scanning sessions using Pearson's correlation coefficient. Bars show mean correlations across participants, error bars show standard error, and grey circles show the correlations of individual participants. Correlations were calculated across scanning sessions and compared face RDMs, voice RDMs, face and voice RDMs, and voice and face RDMs in face-selective, voice-selective, and multimodal ROIs. We tested whether correlations were significantly greater than zero using Wilcoxon signed-rank tests across all 30 participants. No correlations were significant after correction for multiple comparisons at $p \leq .0001$ (FDR corrected for all 48 comparisons). Note that in this figure the rpSTS is classed as a face-selective ROI for consistency purposes only, but in fact it demonstrated multimodal properties.
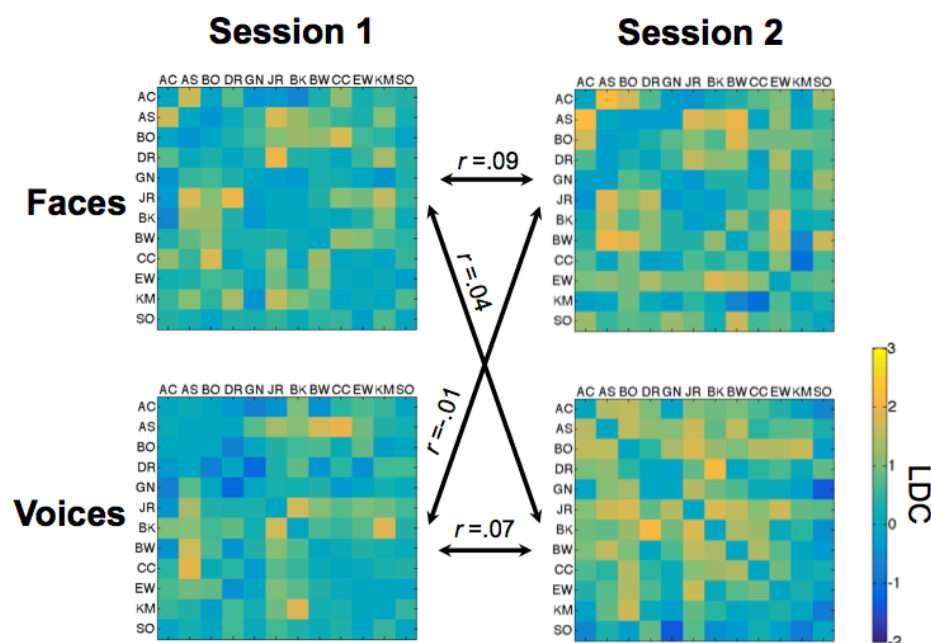
**Figure 4: Representational distance matrix (RDM) comparisons across scanning sessions 1 and 2 in the rpSTS.** Face and voice RDMs for the rpSTS were averaged across all 30 participants for illustration purposes. Each cell shows the discriminability of the brain activity patterns corresponding to a pair of identities (12 identities in total) computed using the linear discriminant contrast (LDC) and crossvalidating across data from three runs. Each matrix is symmetric around a diagonal of zeros. A value of zero or lower indicates no discriminability. For each participant we compared the representational geometry of the face and voice RDMs with the representational geometry in the RDM of the *other* modality (crossmodal comparisons) and in the RDM of the *same* modality (unimodal comparisons) using Pearson's correlation. The figure shows Pearson's correlations for all comparisons averaged across participants.

We predicted that face and voice RDMs would be correlated in ROIs that represent person-identity independently from modality. However, our results showed no significant correlations between face and voice RDMs in face-selective, voice-selective, or multimodal ROIs (Figure 3). It is possible that comparing RDM across different scanning sessions taking place on separate days did not allow us to detect subtle consistencies in the representational geometry for face-identities and voice-identities. To address this concern, we also compared face and voice RDMs within the same scanning session. However, we still found no significant correlations between face and voice RDMs. Therefore, using this method we found no evidence of modality-general person-identity representations in our ROIs.

We also predicted that there would be correlations between RDMs within the same modality in

29

706      regions that represent only face-identity or only voice-identity. No correlations between face

707      RDMs or between voice RDMs in any ROI were significant after correction for multiple

708      comparisons.

709

710      **Analysis B: RSA investigating identity discriminability**

711      Our second main analysis tested the generalisation of pattern discriminants from one modality

712      to the other. More specifically, we computed crossmodal RDMs and we tested whether linear

713      discriminants computed on pairs of faces could be used to discriminate between pairs of

714      voices, and vice-versa. We also tested whether each ROI could discriminate between pairs of

715      stimuli within the same modality. Mean LDC distances across all cells in crossmodal, face, and

716      voice RDMs were compared against zero.
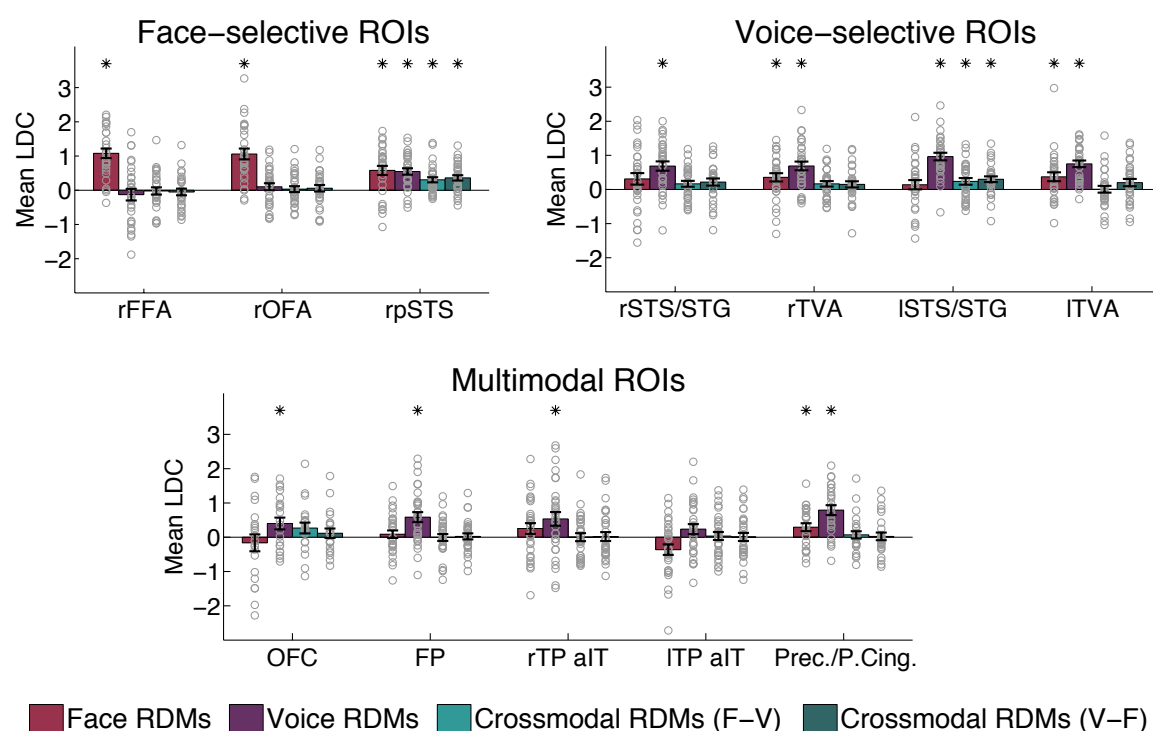
717



718
719
720 **Figure 5: Results of RSA investigating identity discriminability.** Mean LDC between identities in face RDMs,
721 voice RDMs, and crossmodal RDMs in face-selective, voice-selective, and multimodal ROIs. There are two types
722 of crossmodal RDMs: (a) face discriminant applied to voices (F-V), and (b) voice discriminant applied to faces
723 (V-F). Bars show mean LDC values averaged across participants, error bars show standard error, and grey circles
724 show mean LDC values for individual participants. We tested whether the mean LDC values were significantly
725 greater than zero using one-sample $t$-tests across all 30 participants. Stars represent significant tests at $p \leq .0150$

30

726 (FDR corrected for all 48 comparisons). These results show generalisation of the pattern discriminants from one
727 modality to the other in the rpSTS and in the lSTS/STG. In addition, face-selective ROIs discriminate between
728 face-identities, and voice-selective ROIs discriminate between voice-identities. Note that in this figure the rpSTS
729 is classed as a face-selective ROI for consistency purposes only, but in fact it demonstrated multimodal properties.
730

731 **Table 2.** One-sample *t*-test results for mean LDC values in crossmodal RDMs. Stars represent statistical
732 significance at $p \leq .0150$ (FDR corrected for all 48 comparisons in face, voice, and crossmodal RDMs). The rpSTS
733 is classed as a face-selective ROI for consistency purposes only, but in fact it demonstrated multimodal properties.

| | df | t | Crossmodal RDMs (face-voice) Sig. (1-tailed) | d | t | Crossmodal RDMs (voice-face) Sig. (1-tailed) | d |
|---|---|---|---|---|---|---|---|
| Face-selective ROIs | | | | | | | |
| rFFA | 27 | -0.198 | .5779 | 0.04 | -0.529 | .6993 | 0.10 |
| rOFA | 28 | 0.374 | .3557 | 0.07 | 0.624 | .2689 | 0.12 |
| rpSTS | 29 | 4.091 | .0002* | 0.75 | 4.582 | .0001* | 0.84 |
| Voice-selective ROIs | | | | | | | |
| rSTS/STG | 29 | 1.928 | .0319 | 0.35 | 2.093 | .0226 | 0.38 |
| rTVA | 29 | 2.064 | .0240 | 0.38 | 1.662 | .0537 | 0.30 |
| lSTS/STG | 29 | 2.443 | .0104* | 0.45 | 3.543 | .0007* | 0.65 |
| lTVA | 29 | 0.062 | .4755 | 0.01 | 1.891 | .0343 | 0.35 |
| Multimodal ROIs | | | | | | | |
| OFC | 20 | 1.698 | .0525 | 0.37 | 0.841 | .0250 | 0.18 |
| FP | 29 | -0.062 | .5244 | 0.01 | 0.285 | .3888 | 0.05 |
| rTP-aIT | 29 | 0.023 | .4910 | 0.00 | 0.153 | .4398 | 0.03 |
| lTP-aIT | 29 | 0.301 | .3830 | 0.05 | 0.075 | .4703 | 0.01 |
| Prec./P.Cing. | 25 | 0.660 | .2577 | 0.13 | 0.220 | .4138 | 0.04 |

734

735

736 **Table 3.** One-sample *t*-test results for mean LDC values in face and voice RDMs. Stars represent statistical
737 significance at $p \leq .0150$ (FDR corrected for all 48 comparisons in face, voice, and crossmodal RDMs). The rpSTS
738 is classed as a face-selective ROI for consistency purposes only, but in fact it demonstrated multimodal properties.

| | df | t | Face RDMs Sig. (1-tailed) | d | t | Voice RDMs Sig. (1-tailed) | d |
|---|---|---|---|---|---|---|---|
| Face-selective ROIs | | | | | | | |
| rFFA | 27 | 7.764 | .0001* | 1.47 | -0.753 | .7711 | 0.14 |
| rOFA | 28 | 6.707 | .0001* | 1.25 | 0.995 | .1641 | 0.18 |
| rpSTS | 29 | 4.378 | .0001* | 0.80 | 5.871 | .0001* | 1.07 |
| Voice-selective ROIs | | | | | | | |
| rSTS/STG | 29 | 1.850 | .0373 | 0.34 | 5.025 | .0001* | 0.92 |
| rTVA | 29 | 2.945 | .0031* | 0.54 | 5.447 | .0001* | 0.99 |
| lSTS/STG | 29 | 1.019 | .1583 | 0.19 | 8.667 | .0001* | 1.58 |
| lTVA | 29 | 2.846 | .0040* | 0.52 | 7.834 | .0001* | 1.43 |
| Multimodal ROIs | | | | | | | |

| | | | | | | | |
|---|---|---|---|---|---|---|---|
| OFC | 20 | -0.662 | .7424 | 0.14 | 2.337 | .0150* | 0.51 |
| FP | 29 | 0.799 | .2153 | 0.15 | 4.007 | .0002* | 0.73 |
| rTP-aIT | 29 | 1.617 | .0583 | 0.30 | 2.685 | .0059* | 0.49 |
| lTP-aIT | 29 | -2.369 | .9877 | 0.43 | 1.630 | .0570 | 0.30 |
| Prec./P. Cing. | 25 | 2.538 | .0089* | 0.50 | 5.524 | .0001* | 1.08 |

739

740

741  We predicted that in brain regions with modality-general person identity representations the

742  mean LDC values for crossmodal RDMs would be significantly greater than zero. Our results

743  showed that mean LDC values in these RDMs were significantly greater than zero in the

744  rpSTS, and in the voice-selective lSTS/STG (Figure 5; Table 2). These results show that the

745  rpSTS could discriminate pairs of face-identities based on pattern discriminants computed from

746  pairs of voice-identities (and vice-versa), and therefore appears to form modality-independent

747  person-identity representations.

748

749  We note that while the mean LDC values for crossmodal RDMs in the lSTS/STG were

750  significant, the mean LDC value for face RDMs was not. While this result suggests that this

751  region was able to discriminate identities based on crossmodal information, it is unlikely that a

752  modality-general representation could exist without face-identity discrimination. Therefore,

753  this result should be interpreted with caution. It is possible that in addition to the rpSTS, the

754  lpSTS also contains a modality-general person-identity representation and it could be driving

755  the positive result in the lSTS/STG. However, we were not able to test this because we could

756  not localise the lpSTS in our participants using our face localiser.

757

758  We also predicted that mean LDC values for face RDMs and voice RDMs would be

759  significantly greater than zero in ROIs that represent face-identity and voice-identity,

760  respectively. We found that mean LDC values in face RDMs were significantly greater than

761  zero in all ROIs originally defined as face-selective (rFFA, rOFA, rpSTS), in the TVAs, and in

32

762   the multimodal Prec./P. Cing. (Figure 5; Table 3). These results show that all these regions

763   could discriminate between face-identities. A follow up analysis in which all overlapping

764   rpSTS voxels were removed from the rTVA showed that the significant result for faces in

765   rTVA was driven by the rpSTS. Mean LDC values in voice RDMs were significantly greater

766   than zero in all voice-selective ROIs (TVAs, STS/STG), in the rpSTS (originally defined as

767   face-selective), and in the multimodal OFC, FP, rTP-aIT and Prec./P. Cing. (Figure 5; Table 3).

768

769   It is possible that the discrimination of identities in our ROIs was driven by different-gender

770   identity pairs (female-male). To investigate this possibility, for each ROI and condition that

771   showed mean LDC values significantly greater than zero (Figure 5 & Tables 2,3) and for each

772   participant we compared the mean LDC values for different-gender identity pairs (calculated

773   across 36 pairs) with the mean LDC values for same-gender identity pairs (calculated across 30

774   pairs: female-female & male-male) in each RDM (we used paired $t$-tests, and used FDR

775   correction for all 19 comparisons). Results for the rpSTS showed no significant difference

776   between the discriminability of different-gender and same-gender identity pairs for face, voice,

777   or crossmodal RDMs (all $p>.0533$), demonstrating that person-identity discrimination in this

778   region was not driven by discriminating gender. In contrast, mean LDC values for different-

779   gender identity pairs were significantly higher than mean LDC values for same-gender identity

780   pairs for face RDMs in the rFFA and rOFA (both $p \leq.0010$), and for voice RDMs in the

781   bilateral TVAs and STS/STG (all $p \leq.0005$), suggesting that gender contributed to the

782   discrimination in these regions. However, mean LDC values for same-gender identity pairs

783   were still significantly greater than zero (one-sample $t$-tests) for face RDMs in the rFFA and

784   rOFA (both $p<.0001$) and for voice RDMs in the bilateral TVAs and STS/STG (all $p \leq.0239$),

785   suggesting that identity discrimination in these regions is not solely driven by differences in

786   gender.

787

**Exploratory whole-brain searchlight analyses.**

789 Finally, we conducted additional exploratory searchlight analyses across the whole brain to

790 determine whether there were brain regions with modality-general person-identity

791 representations that are not included in our ROIs. The first searchlight analysis investigated

792 correlations between face and voice RDMs across the whole brain, and we did not find any

793 regions showing such correlations between face and voice representational geometries.

794

795 The second searchlight analysis investigated crossmodal generalization of discriminants for

796 pairs of identities across the whole brain. We found a number of clusters in which the mean

797 LDC in crossmodal RDMs was significantly greater than zero (FWE corrected threshold p ≤

798 .05), and below we report t-values and MNI coordinates for the peak grey matter voxels in each

799 cluster. Anatomical labels for peak voxels are based on the Harvard-Oxford cortical and

800 subcortical structural atlases. The results showed a large cluster (*k*=1927, *p*=.007) with peaks in

801 the right putamen (*t*=4.33, x=21, y=20, z=-1), the left posterior middle temporal gyrus

802 (*t*=4.04,x=-57, y=-19, z=-7), and the right precentral gyrus (*t*=3.89, x=54, y=8, z=32).

803 Significant clusters were also found in the right paracingulate gyrus (k=1340, *p*=.003, *t*=4.34,

804 x=6, y=47, z=23), in the left hippocampus (k=160, *p*=.017, *t*=4.45, x=-24, y=-37, z=2), in the

805 right anterior supramarginal gyrus (k=84, *p*=.006, *t*=6.18, x=48, y=-22, z=38), in the left cuneal

806 cortex (k=48, *p*=.036, *t*=3.99, x=-18, y=-76, z=29), and a cluster (*k*=100, *p*=.039) with peaks in

807 the left temporooccipital middle temporal gyrus (*t*=3.58, x=-48, y=-46, z=5) and inferior lateral

808 occipital cortex (*t*=3.45, x=-48, y=-67, z=8). Finally, we also found a significant cluster in the

809 rpSTS at an uncorrected threshold of p ≤ .005 (k=592, *p*=.001, *t*=4.05, x=48, y=-49, z=11) that

810 overlapped with our rpSTS ROI.

811

34

# Discussion

812

813 We show evidence of a modality-general person-identity representation in a multimodal region

814 of the rpSTS, demonstrating that this region was able to discriminate familiar identities based

815 on modality-general information in faces and voices. More specifically, the rpSTS could

816 discriminate pattern estimates for pairs of face-identities based on linear discriminants

817 computed from pattern estimates for pairs of voice-identities, and vice-versa. A crucial and

818 novel aspect of our study is that we showed that the rpSTS not only discriminates between

819 identities, but also generalises across multiple naturalistically varying face videos and voice

820 recordings of the same identity. By always comparing pattern estimates across independent

821 runs with different face and voice tokens for the same identities, we showed that the face- and

822 voice-elicited person-identity representations in the rpSTS are stimuli-invariant. Invariant

823 identity representations were also found for face-identities in face-selective regions (rFFA and

824 rOFA) and for voice-identities in voice-selective regions (bilateral TVA and STS/STG).

825 Finally, we did not find evidence of matching representational geometries for faces and voices,

826 across or within modalities.

827

828 **A modality-general and invariant person-identity representation in the rpSTS**

829 Our finding of a modality-general person-identity representation in a multimodal region of the

830 rpSTS supports the Multimodal Processing Model, which proposes that face and voice

831 information is integrated in multimodal brain regions (Ellis et al., 1997; Campanella and Belin,

832 2007). In contrast, we did not find support for the prediction from the Coupling of Face and

833 Voice Processing Model (von Kriegstein et al., 2005; von Kriegstein and Giraud, 2006) that

834 there would be a modality-general identity representation in face-selective regions of the

835 fusiform gyrus.

836

837 The rpSTS has previously been associated with crossmodal representations of emotion from

838 faces and voices and shows a preference for people-related stimuli regardless of modality

839 (Watson et al., 2014a, 2014b). Furthermore, multiple studies have demonstrated overlap

840 between face-selective and voice-selective regions in the rpSTS (Wright et al., 2003; Watson et

841 al., 2014a; Deen et al., 2015; Anzellotti and Caramazza, 2017;). It has been proposed that the

842 STS integrates person-specific patterns of movement from faces, voices, and bodies to assist in

843 person-identity recognition (Yovel and O'Toole, 2016). It is possible that the intrinsic

844 relationship between a person's idiosyncratic facial movements and manner of speech

845 contributes to the integration of face- and voice-identity information in the rpSTS.

846

847 Our finding of a modality-general identity representation in the rpSTS is also in agreement

848 with two recent studies showing across-modality classification of pattern estimates for familiar

849 faces and voices in the rpSTS (Anzellotti and Caramazza, 2017) and a more anterior part of the

850 STS (Hasan et al., 2016). In contrast to these previous studies, we showed that the rpSTS also

851 demonstrates face- and voice-elicited representations of person-identity that are invariant to

852 different tokens of the same face and voice. The ability to "tell people together" by identifying

853 different tokens of a face and voice as belonging to the same person is as important as the

854 ability to "tell people apart" (i.e. discriminate between different people) (Burton, 2013;

855 Anzellotti and Caramazza, 2014). Hasan et al. (2016) were unable to investigate invariant

856 representations because they used a single face image and a single voice recording for each

857 identity, which in turn were derived from the same original stimulus, making interpretation of

858 their results difficult (Lavan, 2017). Anzellotti and Caramazza (2017) used two face and voice

859 tokens for each identity but did not train and test their classifier on different tokens, and

860 therefore did not demonstrate representations that were invariant to different tokens of the

861 same face and voice in this study. Finally, in contrast to these previous studies, we used a

36

862 larger set of identities and multiple naturalistically varying tokens in order to better capture the

863 level of robust invariant recognition required in everyday life. While behavioural studies have

864 shown the importance of within-person variability for recognition (Jenkins et al., 2011; Burton,

865 2013; Burton et al., 2016), this is rarely taken into account in neuroimaging experiments, which

866 typically use highly similar or artificial stimuli for the same person.

867

868 **Invariant representations of face-identity and voice-identity**

869 The face-selective rFFA and rOFA were able to discriminate between the faces of different

870 people while also showing invariance to the different videos of each person's face. This finding

871 is in agreement with Anzellotti et al. (2014) and Guntupalli et al. (2017), who showed

872 representations of face-identity in the FFA (and OFA, in Anzelotti et al., 2014) that generalise

873 across different viewpoints of the face. However, in contrast with these studies, which used

874 stimuli with low within-person variability, we show that representations in these regions

875 generalise across highly variable face videos, and can thus discriminate between different face-

876 *identities*, rather than between individual face *images*.

877

878 Voice-selective regions in STS/STG and the TVAs bilaterally could discriminate between

879 different speakers while showing invariance to the different recordings of each voice. These

880 findings are in line Formisano et al. (2008), who showed representations of speaker identity

881 that generalise across utterances of different vowels in the lateral Heschl's gyrus/sulcus and in

882 the right STS. We extend this finding by showing that generalisation across different

883 recordings of the same voice is possible even when using short sentences with variable speech

884 content that were recorded in different settings.

885

886 We also found invariant discrimination of face- and voice-identity in a multimodal region in

37

887    the precuneus/posterior cingulate. This region has been previously associated with the

888    processing of familiar faces and voices (Shah et al., 2001), and has been found to discriminate

889    between different face-identities (Visconti Di Oleggio Castello et al., 2017). Our results

890    suggest that representations of faces and voices may be interspersed in this region, but are not

891    shared across modalities. Finally, we showed invariant representations of voice-identity, but

892    not face-identity, in the frontal pole, a region that has been previously associated with the

893    processing of familiar voices (Nakamura et al., 2001). It should be noted that although we

894    initially localised the frontal pole as a multimodal region, our results showed that it did not

895    respond significantly to faces in the main experimental runs.

896

897    **Representational geometries**

898    We did not find matching representational geometries across faces and voices in rpSTS despite

899    finding crossmodal generalisation of the pattern discriminants. It is possible that all identities

900    were equally distinct from each other within each modality (i.e. the nature of person-identity

901    code in these regions does not result in variable representational distances between identities).

902    In addition, the rpSTS shows both modality-specific and modality-general representations, and

903    it is possible that the former had stronger influence on the representational geometry.

904    Beauchamp et al. (2004) showed that the pSTS contains intermixed visual, auditory, and

905    multisensory patches, and future studies could use higher-resolution neuroimaging methods to

906    probe person-identity representations in this region.

907

908    In all other ROIs, we also did not find any evidence of stable representational geometries for

909    face-identities or voice-identities only. Again, it could be that identities were equally distinct

910    across from each other within each modality, or it could be that experimental conditions would

911    need to be improved to obtain more reliable representational geometries.

912

**Anterior temporal lobe and searchlight results**

We did not find evidence of face-, voice-, or person-identity representations in the anterior temporal lobe. This was surprising given that this region has been previously associated with the processing of person-identity (Ellis et al., 1989; Gainotti, 2011). The fact that our TP-aIT ROIs responded more to voices that to faces suggests that our multimodal region localiser was not optimal for detecting multimodal responses in the anterior temporal lobe. Moreover, our sequences were not tailored to detect fMRI responses in this region (Axelrod and Yovel, 2013), and therefore more research using specialised scanning parameters for the localisation of this region is warranted.

922

It is possible that modality-general representations exist outside face- and/or voice-selective regions, and our exploratory searchlight results revealed person-identity representations in the paracingulate gyrus, right insular cortex, left nucleus accumbens, left anterior postcentral gyrus, and left hippocampus. Quiroga et al. (2005, 2009) found that cells in the hippocampus (and also amygdala and entorhinal cortex) were highly responsive to specific identities, and responded to both the face and name of that person. It will be interesting to further probe the role of the hippocampus (and the other regions found during the searchlight analyses) in person-identity recognition.

931

**Conclusion**

To conclude, we showed a modality-general person-identity representation that generalises across different, naturalistically varying face videos and voice recordings of the same person in a multimodal region of the rpSTS. This supports the Multimodal Processing Model for face and voice integration. We also found evidence of video-invariant face-identity representations

39

937  in face-selective regions (rFFA, rOFA), and sound-invariant voice-identity representations in

938  voice-selective regions (TVA, STS/STG). Future studies could focus on the nature and type of

939  face and voice information that is represented in these different regions, and in how these

940  representations are formed, both through development, and during the process of becoming

941  familiar with someone.

942

943

944

945

946

947

948

949

950

951

952

953

954

955

956

# References

958  Anzellotti S, Caramazza A (2014) The neural mechanisms for the recognition of face identity
959      in humans. Front Psychol 5:1–6.

960  Anzellotti S, Caramazza A (2016) From parts to identity: invariance and sensitivity of face
961      representations to different face halves. Cereb Cortex 26:1900–1909.

962  Anzellotti S, Caramazza A (2017) Multimodal representations of person identity individuated
963      with fMRI. Cortex 89: 85–97.

964  Anzellotti S, Fairhall SL, Caramazza A (2014).Decoding representations of face identity that
965      are tolerant to rotation. Cereb Cortex 24:1988–1995.

966  Axelrod V, Yovel G (2013) The challenge of localizing the anterior temporal face area: A
967      possible solution. Neuroimage 81:371–380.

968  Axelrod V, Yovel G (2015) Successful decoding of famous faces in the fusiform face area.
969      PLoS One 10: e0117126.

970  Beauchamp MS, Argall BD, Bodurka J, Duyn JH, Martin A (2004) Unraveling multisensory
971      integration: patchy organization within human STS multisensory cortex. Nat Neurosci
972      7:1190–1192.

973  Belin P, Zatorre RJ, Lafaille P, Ahad P, Pike B (2000) Voice-selective areas in human auditory
974      cortex. Nature 403:309–312.

975  Brainard DH (1997) The psychophysics toolbox. Spat vis 10:433-436.

976  Burton AM (2013) Why has research in face recognition progressed so slowly? The
977      importance of variability. Q J Exp Psychol 66:1467–1485.

978  Burton AM, Kramer RSS, Ritchie KL, Jenkins R (2016) Identity from variation:
979      Representations of faces derived from multiple instances. Cogn Sci 40:202–223.

980  Campanella S, Belin P (2007) Integrating face and voice in person perception. Trends Cogn Sci
981      11:535–543.

982    Carlin JD, Kriegeskorte N (2017) Adjudicating between face-coding models with individual-

983        face fMRI responses. PLoS Comput Biol 13: e1005604.

984    Collins JA, Koski JE, Olson IR (2016) More than meets the eye: The merging of perceptual

985        and conceptual knowledge in the anterior temporal face area. Front Hum Neurosci 10:1–

986        11.

987    Deen B, Koldewyn K, Kanwisher N, Saxe R (2015) Functional organization of social

988        perception and cognition in the superior temporal sulcus. Cereb Cortex 25:4596–4609.

989    Ellis AW, Young AW, Critchley EMR (1989) Loss of memory for people following temporal

990        lobe damage. Brain 112:1469–1483.

991    Ellis HD, Jones DM, Mosdell N (1997) Intra- and inter-modal repetition priming of familiar

992        faces and voices. Br J Psychol 88:143–156.

993    Formisano E, De Martino F, Bonte M, Goebel R (2008) "Who" is saying "what"? Brain-based

994        decoding of human voice and speech. Science 322:970–973.

995    Fox, CJ, Iaria G, Barton JJS (2009) Defining the face processing network: Optimization of the

996        functional localizer in fMRI. Hum Brain Mapp 30:1637–1651.

997    Gainotti G (2011) What the study of voice recognition in normal subjects and brain-damaged

998        patients tells us about models of familiar people recognition. Neuropsychologia 49:2273–

999        2282.

1000    Goesaert E, Op de Beeck HP (2013) Representations of facial identity information in the

1001        ventral visual stream investigated with multivoxel pattern analyses. J Neurosci 33:8549–

1002        8558.

1003    Guntupalli JS, Wheeler KG, Gobbini MI (2017) Disentangling the representation of identity

1004        from head view along the human face processing pathway. Cereb Cortex 27:46–53.

1005    Hasan BAS, Valdes-Sosa M, Gross J, Belin P (2016) "Hearing faces and seeing voices":

1006        Amodal coding of person identity in the human brain. Sci Rep 6:37494.

1007 Jenkins R, White D, Van Montfort X, Burton AM (2011) Variability in photos of the same

1008    face. Cognition 121:313–323.

1009 Jenkinson M, Beckmann CF, Behrens TE, Woolrich MW, Smith SM (2012)

1010    Fsl. Neuroimage 62:82-790.

1011 Joassin F, Pesenti M, Maurage P, Verreckt E, Bruyer R, Campanella S (2011) Cross-modal

1012    interactions between human faces and voices involved in person recognition. Cortex

1013    47:367–376.

1014 Julian JB, Fedorenko E, Webster J, Kanwisher N (2012) An algorithmic method for

1015    functionally defining regions of interest in the ventral visual pathway. Neuroimage

1016    60:2357–2364.

1017 Kanwisher N, McDermott J, Chun MM (1997) The fusiform face area: a module in human

1018    extrastriate cortex specialized for face perception. J Neurosci 17:4302–4311.

1019 Kriegeskorte N, Formisano E, Sorger B, Goebel R (2007) Individual faces elicit distinct

1020    response patterns in human anterior temporal cortex. Proc Natl Acad Sci U S A,

1021    104:20600–20605.

1022 Kriegeskorte N, Kievit RA (2013) Representational geometry: Integrating cognition,

1023    computation, and the brain. Trends Cogn Sci 17:401–412.

1024 Kriegeskorte N, Mur M, Bandettini PA (2008a) Representational similarity analysis -

1025    connecting the branches of systems neuroscience. Front Syst Neurosci 2:1–28.

1026 Kriegeskorte N, Mur M, Ruff DA, Kiani R, Bodurka J, Esteky H, Tanaka K, Bandettini PA

1027    (2008b) Matching categorical object representations in inferior temporal cortex of man

1028    and monkey. Neuron 60:1126–1141.

1029 Lavan N (2017) Commentary: "Hearing faces and seeing voices": Amodal coding of person

1030    identity in the human brain. Front Neurosci 11:303.

1031 Ledoit O, Wolf M (2004) A well-conditioned estimator for large-dimensional covariance

1032    matrices. J Multivar Anal 88:365-411.

1033    Nakamura K, Kawashima R, Sugiura M, Kato T, Nakamura A, Hatano K, Nagumo S, Kubota

1034        K. Fukuda H, Ito K, Kojima S (2001) Neural substrates for recognition of familiar voices:

1035        A PET study. Neuropsychologia 39:1047–1054.

1036    Nestor A, Plaut DC, Behrmann M (2011) Unraveling the distributed neural code of facial

1037        identity through spatiotemporal pattern analysis. Proc Natl Acad Sci U S A 108:9998–

1038        10003.

1039    Nili H, Wingfield C, Walther A, Su L, Marslen-Wilson W, Kriegeskorte N (2014) A toolbox

1040        for representational similarity analysis. PLoS Comput Biol 10:e1003553.

1041    Pelli DG (1997) The VideoToolbox software for visual psychophysics: Transforming numbers

1042        into movies. Spat vis 10:437-442.

1043    Pernet CR, McAleer P, Latinus M, Gorgolewski KJ, Charest I, Bestelmeyer PEG, Watson RH,

1044        Fleming D, Crabbe F, Valdes-Sosa M, Belin P (2015) The human voice areas: Spatial

1045        organization and inter-individual variability in temporal and extra-temporal cortices.

1046        NeuroImage 119:164–174.

1047    Pitcher D, Dilks DD, Saxe RR, Triantafyllou C, Kanwisher N (2011) Differential selectivity

1048        for dynamic versus static information in face-selective cortical regions. NeuroImage

1049        56:2356–2363.

1050    Quiroga RQ, Reddy L, Kreiman G, Koch C, Fried I (2005) Invariant visual representation by

1051        single neurons in the human brain. Nature 435:1102–1107.

1052    Quiroga RQ, Kraskov A, Koch C, Fried I (2009) Explicit encoding of multimodal percepts by

1053        single neurons in the human brain. Curr Biol 19:1308-1313.

1054    Schweinberger SR, Herholz A, Sommer W (1997) Recognizing famous voices. J Speech Lang

1055        Hear Res 40:453–463.

1056    Shah NJ, Marshall JC, Zafiris O, Schwab A, Zilles K, Markowitsch HJ, Fink GR (2001) The

44

1057     neural correlates of person familiarity. A functional magnetic resonance imaging study

1058     with clinical implications. Brain 124:804–815.

1059    Verosky SC, Todorov A, Turk-Browne NB (2013) Representations of individuals in ventral

1060     temporal cortex defined by faces and biographies. Neuropsychologia 51:2100–2108.

1061    Visconti Di Oleggio Castello M, Halchenko YO, Guntupalli JS, Gors JD, Gobbini MI (2017)

1062     The neural representation of personally familiar and unfamiliar faces in the distributed

1063     system for face perception. Sci Rep 7:1–14.

1064    Van Lancker D, Krieman J, Emmorey K (1985) Familiar voice recognition: Patterns and

1065     parameters. Part I: Recognition of backward voices. J Phon 13:19-38.

1066    von Kriegstein K, Kleinschmidt A, Sterzer P, Giraud, AL (2005) Interaction of face and voice

1067     areas during speaker recognition. J Cogn Neurosci 17:367–376.

1068    von Kriegstein K, Giraud AL (2006) Implicit multisensory associations influence voice

1069     recognition. PLoS Biol 4:e326.

1070    von Kriegstein K, Kleinschmidt A, Giraud AL (2006) Voice recognition and cross-modal

1071     responses to familiar speakers' voices in prosopagnosia. Cereb Cortex 16:1314–1322.

1072    von Kriegstein K, Dogan O, Grüter M, Giraud AL, Kell CA, Grüter T, Kleinschmidt A, Kiebel

1073     SJ (2008) Simulation of talking faces in the human brain improves auditory speech

1074     recognition. Proc Natl Acad Sci U S A 105:6747–6752.

1075    Walther A, Nili H, Ejaz N, Alink A, Kriegeskorte N, Diedrichsen J (2016) Reliability of

1076     dissimilarity measures for multi-voxel pattern analysis. Neuroimage 137:188–200.

1077    Watson R, Latinus M, Charest I, Crabbe F, Belin P (2014a) People-selectivity, audiovisual

1078     integration and heteromodality in the superior temporal sulcus. Cortex 50:125–136.

1079    Watson R, Latinus M, Noguchi T, Garrod O, Crabbe F, Belin P (2014b) Crossmodal adaptation

1080     in right posterior superior temporal sulcus during face-voice emotional integration. J

1081     Neurosci 34:6813–6821.

1082    Winkler AM, Ridgway GR, Webster MA, Smith SM, Nichols TE (2014) Permutation inference

1083        for the general linear model. Neuroimage 92:381-397.

1084    Wright TM, Pelphrey KA, Allison T, McKeown MJ, McCarthy G (2003) Polysensory

1085        interactions along lateral temporal regions evoked by audiovisual speech. Cereb Cortex

1086        13:1034–1043.

1087    Yovel G, O'Toole AJ (2016) Recognizing people in motion. Trends Cogn Sci 20:383–395.