

1 **Title:**

2 **Decoding of emotion expression in the face, body and voice reveals sensory modality**  
3 **specific representations.**

4

5

6 **Authors:** Maarten Vaessen<sup>1</sup>, Kiki Van der Heijden<sup>1</sup>, Beatrice de Gelder<sup>1,2</sup>

7

8 **Affiliations:**

9 <sup>1</sup> Maastricht University, Maastricht (NL)

10 <sup>2</sup> Department of Computer Science, University College London, London (UK)

11

12

13 **ABSTRACT**

14 A central issue in affective science is whether the brain represents the emotional expressions of  
15 faces, bodies and voices as abstract categories in which auditory and visual information converge  
16 in higher order conceptual and amodal representations. This study explores an alternative theory  
17 based on the hypothesis that under naturalistic conditions where affective signals are acted upon  
18 that rather than reflected upon, major emotion signals (face, body, voice) have sensory specific  
19 brain representations. During fMRI recordings, participants were presented naturalistic dynamic  
20 stimuli of emotions expressed in videos of either the face or the whole body, or voice fragments.  
21 To focus on automatic emotion processing and bypass explicit emotion cognition relying on  
22 conceptual processes, participants performed an unrelated target detection task presented in a  
23 different modality than the stimulus. Using multivariate analysis to assess neural activity patterns  
24 in response to emotion expressions in the different stimuli types, we show a distributed brain  
25 organization of affective signals in which distinct emotion signals are closely tied to the sensory  
26 origin. Our findings are consistent with the notion that under ecological conditions the various  
27 sensory emotion expressions have different functional roles, even when from an abstract  
28 conceptual vantage point they all exemplify the same emotion category.

29

## 30 INTRODUCTION

31 In the course of daily social interactions, emotion signals from the face, voice and body  
32 are recognized effortlessly and responded to spontaneously when rapid adaptive actions are  
33 required. The specifics of the subjective experience in the natural environment determine which  
34 affective signal dominates and triggers the adaptive behavior. Rarely are the face, the whole  
35 body and the voice equally salient. That is, the actual conditions under which we react to an  
36 angry face may be different from those of hearing an angry voice or viewing whole body  
37 movements. For example, we see faces from close by and therefore personal familiarity may play  
38 a role in how we react to the angry face. This is less so for the voice or the whole body, both of  
39 which already prompt reactions when seen or heard from a distance while information about  
40 personal identity is not yet available or needed for action preparation. Thus, the angry body  
41 expression viewed from a distance and the angry face expression seen from closeby may each  
42 trigger a different reaction as behavior needs to be adapted to the concrete context. Therefore, a  
43 representation of affective meaning that is sensitive to the spatiotemporal parameters may seem  
44 desirable rather than an abstract system of higher order concepts as traditionally envisaged by  
45 emotion theorists (Ekman P and D Cordaro 2011); but see (Lindquist KA et al. 2012).

46 Research on the brain correlates of emotion has favored the traditional notion of abstract  
47 neural representations of basic emotions and this has also been the dominant rationale for  
48 multisensory research. Studies comparing how not just the face but also the voice and the whole  
49 body convey emotions have followed this overall basic emotion perspective and asked where in  
50 the brain affective information from different sensory systems converges (Peelen MV et al. 2010;  
51 Klasen M et al. 2011). Such representations were found in high-level brain areas known for their  
52 role in categorization of mental states (Peelen MV *et al.* 2010). Specifically, medial prefrontal  
53 cortex (MPFC) and superior temporal cortex (STS) represent emotions perceived in the face,  
54 voice or body at a modality-independent level. Furthermore, these supramodal or abstract  
55 emotion representations presumably play an important role in multisensory integration by  
56 driving and sustaining convergence of the sensory inputs towards the amodal emotion  
57 representation (Gerdes AB et al. 2014).

58 The present study explores a different perspective, complementary or compatible with  
59 that of high level abstract emotion representations yet motivated by the natural variability of

60 daily emotion perception conditions that do not routinely involve conscious use of verbal labels.  
61 Emotion perception in naturalistic conditions is often driven by a specific context relative to a  
62 behavioral goal, as opposed to the conditions in the lab, where tasks often use explicit evaluation  
63 of emotional expressions. Under natural conditions, each sensory modality may have its own  
64 functionality such that e.g. fear is more effectively conveyed by the face, anger by the body and  
65 happiness by the voice. If so, brain responses would be characterized by specific emotion-  
66 modality combinations.

67 Additionally, the notion that supramodal representations of basic emotions are the pillars  
68 of emotion processing in the brain and allow for smooth convergence between the different  
69 sensory modalities is not fully supported by the literature. First, since the original proposal by  
70 Ekman (Ekman P 1992) and the constructivist alternative argued by Russell (Russell JA 2003)  
71 and most recently Feldman Barrett (Lewis M et al. 2010; Barrett LF 2017), the notion of a set of  
72 basic emotions with discrete brain correlates continues to generate controversy (Kragel PA and  
73 KS LaBar 2016; Saarimaki H et al. 2016). Second, detailed meta-analyses of crossmodal and  
74 multisensory studies, whether they are reviewing the findings about each separate modality or  
75 the results of crossmodal studies (Dricu M and S Fruhholz 2016; Schirmer A and R Adolphs  
76 2017), provide a mixed picture. Furthermore, these meta-analyses also show that a number of  
77 methodological obstacles stand in the way of valid comparisons across studies. That is, taking  
78 into account the role of task (incidental perception, passive perception, and explicit evaluation of  
79 emotional expression) and the use of appropriate control stimuli reduces the number of studies  
80 that can validly be compared. Third, findings from studies that pay attention to individual  
81 differences and to clinical aspects reveal individual differences in sensory salience and  
82 dominance in clinical populations, for example in autism and schizophrenia. For example, (Karle  
83 KN et al. 2018) report an alteration in the balance of cerebral voice and face processing systems  
84 in the form of an attenuated face-vs-voice bias in emotionally competent individuals. This is  
85 reflected in cortical activity differences as well as in higher voice-sensitivity in the left  
86 amygdala. Finally, even granting the existence of abstract supramodal representations – probably  
87 in higher cognitive brain regions - it is unclear how they relate to early stages of affective  
88 processing where the voice, the face and the body information are processed by different sensory  
89 systems comprising distinct cortical and subcortical structures.

90  
91 Here we used naturalistic dynamic stimuli to investigate whether the brain represents  
92 different sensory emotion expressions as modality specific or modality-invariant, using fMRI  
93 with dynamic stimuli expressing affect with either the body, the face or the voice. Importantly,  
94 we investigate these processes independently of the explicit evaluation of emotion to study  
95 whether the brain still differentiates between emotional expressions even when they are not in  
96 the direct focus of attention. We perform multivariate pattern analysis to identify cortical regions  
97 containing representations of emotion independently of the explicit evaluation of emotion.

98 For the sake of clarity, we contrast the implications of the centralist view and the  
99 distributed modality view for the neural representation of emotions. Following the first, there  
100 should be localized representations (here multi voxel patterns) for specific emotions that are  
101 independent of modality. These region(s) would show the following behavior: (1) respond to all  
102 stimuli above some threshold (i.e. be activated by all emotions/modalities); (2) have no  
103 discernable activation pattern between modalities, that is, modality cannot be decoded from the  
104 voxel activation patterns; (3) exhibit distinct voxel activation patterns for each emotion (emotion  
105 is decodable from the regions). This would provide strong evidence for a cross-modal or  
106 modality independent emotion representation and correspond to the classical notion of basic  
107 emotions in the literature. The alternative, distributed theory postulates that emotions are  
108 represented in the brain in a modality-emotion specific way. For example, some specific  
109 emotion-modality combination like a fearful body will elicit activation patterns that are different  
110 from fear from the voice or the face. Likewise, hearing a happy voice (laughing) can elicit  
111 different brain responses from seeing a happy expression from the face or body. In terms of  
112 testing this hypothesis, we would expect to find brain regions where emotion can be decoded but  
113 only within a specific modality, not across modalities. That pattern would provide clear evidence  
114 that the decoding will be driven by specific responses to emotion-modality combination.

115

## 116 **METHODS**

### 117 *Participants*

118 Thirteen healthy participants (mean age = 25.3; age range = 21-30; two males) took part in the  
119 study. Participants reported no neurological or hearing disorders. Ethical approval was provided  
120 by the Ethical Committee of the Faculty of Psychology and Neuroscience at Maastricht  
121 University. Written consent was obtained from all participants. The experiment was carried out  
122 in accordance with the Declaration of Helsinki. Participants either received credit points or were  
123 reimbursed with monetary reward after their participation in the scan session.

#### 124 *Stimuli*

125 Stimuli consisted of color video and audio clips of four male actors expressing three different  
126 emotional reactions to specific events (e.g. fear in a car accident or happiness at a party). Images  
127 were shown to the actors during recordings with the goal of triggering spontaneous and natural  
128 reactions of anger, fear, happiness and an additional neutral reaction. A full description of the  
129 recording procedure, the validation and the video selection is given in (Kret ME, S Pichon, J  
130 Grezes, et al. 2011). In total there were 16 video clips of facial expressions, 16 video clips of  
131 body expressions, and 32 audio clips of vocal expressions, half of which were recorded in  
132 combination with the facial expressions and half of which were recorded in combination with the  
133 body expressions (i.e. two audio clips per emotional expression per actor). All actors were  
134 dressed in black and filmed against a green background under controlled lighting conditions.

135 Video clips were computer-edited using Ulead, After Effects, and Lightworks (EditShare). For  
136 the body stimuli, faces of actors were blurred with a Gaussian mask such that only the  
137 information of the body was available. The validity of the emotional expressions in the video  
138 clips was measured with a separate emotion recognition experiment (emotion recognition  
139 accuracy > 80%). For more information regarding the recording and validation of these stimuli,  
140 see (Kret ME, S Pichon, J Grezes, *et al.* 2011; Kret ME, S Pichon, J Grèzes, et al. 2011).

#### 141 *Experimental design and behavioral task*

142 In a slow-event related design, participants viewed series of 1 second video clips on a projector  
143 screen or listened to series of 1 second audio clips through MR-compatible ear buds  
144 (Sensimetrics S14) equipped with sound attenuating circumaural earbuds (attenuation > 29 dB).  
145 The experiment consisted of 12 runs divided over 2 scan sessions. Six runs consisted of blocks of  
146 face and voice stimuli, followed by six runs consisting of body and voice stimuli. Blocks were

147 either *auditory* (consisting of 18 audio clips) or *visual* (consisting of 18 video clips). These 18  
148 trials within each block comprised 16 regular trials, and two catch trials requiring a response.  
149 Catch trials were included to determine that attention was diverted from explicit recognition or  
150 evaluation of the emotional expression by focusing attention on the other modality. That is,  
151 during visual blocks, participants were instructed to detect auditory catch trials, and during  
152 auditory blocks, participants were instructed to detect visual catch trials. For the auditory catch  
153 trial task, a frequency modulated tone was presented and participants had to respond whether the  
154 direction of frequency modulation was up or down. For the visual distractor task, participants  
155 indicated whether the fixation cross turned lighter or darker during the trial. A separate localizer  
156 session was also performed where participants passively viewed stimuli of faces, bodies, houses,  
157 tools and words in blocks; see (Zhan M et al. 2018) for details.

### 158 ***Data acquisition***

159 We measured blood-oxygen level-dependent (BOLD) signals with a 3 Tesla Siemens Trio whole  
160 body MRI scanner at the Scannexus MRI scanning facilities at Maastricht University  
161 (Scannexus, Maastricht). Functional images of the whole brain were obtained using T2\*-  
162 weighted 2D echo-planar imaging (EPI) sequences [number of slices per volume = 50, 2 mm in-  
163 plane isotropic resolution, repetition time (TR) = 3000 ms, echo time (TE) = 30 ms, flip angle  
164 (FA) = 90°, field of view (FoV) = 800 x 800 mm<sup>2</sup>, matrix size = 100 x 100, multi-band  
165 acceleration factor = 2, number of volumes per run = 160, total scan time per run = 8 min]. A  
166 three-dimensional (3D) T1-weighted (MPRAGE) imaging sequence was used to acquire high  
167 resolution structural images for each of the participants [1-mm isotropic resolution, TR = 2250  
168 ms, TE = 2.21 ms, FA = 9°, matrix size = 256 x 256, total scan time = 7 min approx.]. The  
169 functional localizer scan also used a T2\*-weighted 2D EPI sequence [number of slices per  
170 volume = 64, 2 mm in-plane isotropic resolution, TR = 2000 ms, TE = 30 ms, FA = 77, FoV =  
171 800 x 800 mm<sup>2</sup>, matrix size = 100 x 100, multi-band acceleration factor = 2, number of volumes  
172 per run = 432, total scan time per run = 14 min approx.].

### 173 ***Analysis***

174 *Pre-processing:* Data were preprocessed and analyzed with BrainVoyager QX (Brain  
175 Innovation, Maastricht, Netherlands) and custom Matlab code (Mathworks, USA) (Hausfeld L et

176 al. 2012; Hausfeld L et al. 2014). Preprocessing of functional data consisted of 3D motion  
177 correction (trilinear/sync interpolation using the first volume of the first run as reference),  
178 temporal high pass filtering (thresholded at five cycles per run), and slice time correction. We  
179 co-registered functional images to the anatomical T1-weighted image obtained during the first  
180 scan session and transformed anatomical and functional data to the default Talairach template.

181 *Univariate analysis:* We estimated a random-effects General Linear Model (RFX GLM) with a  
182 predictor for each stimulus condition of interest (12 conditions in total): four emotion conditions  
183 times three modalities (face, body, voice). We note that modality is used here in a broader sense  
184 than just the physical nature of the stimuli (sound versus visual). Additionally, we included  
185 predictors for the trials indicating the start of a new block and the catch trials. Predictors were  
186 created by convolving stimulus predictors with the canonical hemodynamic function. Finally, we  
187 included six motion parameters resulting from the motion correction as predictors of no interest.  
188 For this analysis, data was spatially smoothed with a 6mm full-width half-maximum (FWMH)  
189 Gaussian kernel. To assess where in the brain the two different experimental factors had an  
190 influence, an ANOVA was run with either modality or emotion as a factor. Additionally the  
191 ANOVA with the emotion factor was run for each modality separately. As effect sizes were  
192 generally low, the final group statistical maps were liberally thresholded at  $p < 0.001$  uncorrected.  
193 For visualization purposes, the group volume maps were mapped to the cortical surface. As this  
194 operation involves resampling the data (during which the original statistical values get lost),  
195 surface maps are displayed with discrete label values instead of continuous statistical values. As  
196 such, we do not include a colorbar in the surface map figures.

197 *Multivariate analysis:* We first estimated beta parameters for each stimulus trial with custom  
198 MATLAB code by fitting an HRF function with a GLM to each trial in the time series. These  
199 beta values were then used as input for a searchlight multivariate pattern analysis (MVPA) with a  
200 Gaussian Naïve Bayes classifier (Ontivero-Ortega M et al. 2017). The searchlight was a sphere  
201 with a radius of 5 voxels. The Gaussian Naïve Bayes classifier is an inherently multi-class  
202 probabilistic classifier that performs similar to the much-used Support Vector Machine classifier  
203 in most scenarios, but is computationally more efficient. The classifier was trained to decode (1)  
204 stimulus modality (visual or auditory) (2) stimulus type (i.e. body, face or voice); (3) stimulus  
205 emotion (e.g. fear in all stimulus types vs angry in all stimulus types); (4) within-modality

206 emotion (e.g. body angry vs. body fear) (5) cross-modal emotion (e.g. classify emotion by  
207 training on body stimuli and testing on the voice stimuli from the body session). Classification  
208 accuracy was computed by averaging the decoding accuracy of all folds of a leave one-run out  
209 cross-validation procedure. We tested the significance of the observed decoding accuracies at the  
210 group level with a one-sample t-test against chance (50% for modality, 33% for stimulus types,  
211 25% for emotion). We also tested the emotion effect with an additional analysis where the  
212 neutral condition was excluded. As in the univariate analyses, these maps were transformed to  
213 the cortical surface for display purposes. The cross-modal decoding revealed several subcortical  
214 structures that are not visualized well on the cortical surface and therefor are displayed on a  
215 volume map ( $p < 0.01$ , uncorrected).

216 To evaluate the relative contribution of either stimulus type or of emotion in terms of  
217 information content that can be decoded, we performed an analysis where all stimuli for a  
218 specific combination of two emotions (e.g. anger and happy from the same session) were  
219 extracted and a decoder was trained to either classify the two different emotions or the two  
220 different stimulus types (body/voice or face/voice). The resulting accuracy maps were first  
221 thresholded such that regions where both the stimulus type and emotion decoder accuracy were  
222 below chance at the group level were excluded. Next we contrasted the two accuracy maps at the  
223 group level i.e. accuracy for emotion  $>$  accuracy for stimulus type and displayed these as a  
224 volume map at  $p < 0.001$  uncorrected.

225 *ROI analysis:* Lastly, to gain more insight into details of the responses in some of the regions  
226 identified by the previously mentioned analyses, as well as regions known to be important for  
227 emotion or multi-modal integration (Peelen MV *et al.* 2010), we extracted beta values from these  
228 ROIs and made several plots that (1) display the decoding accuracy of the selected voxels for  
229 stimulus type and emotion, (2) display the mean beta values for each of the 12 conditions, (3)  
230 display the multivoxel representational dissimilarity matrix (Kriegeskorte N *et al.* 2008) and (4)  
231 warp the multivoxel patterns to a 2-dimensional space with multidimensional scaling and display  
232 the relative distance between the conditions with graphical indicators for the stimulus type (icon)  
233 and emotion (line color)

## 234 **RESULTS**

235 ***Univariate analysis***

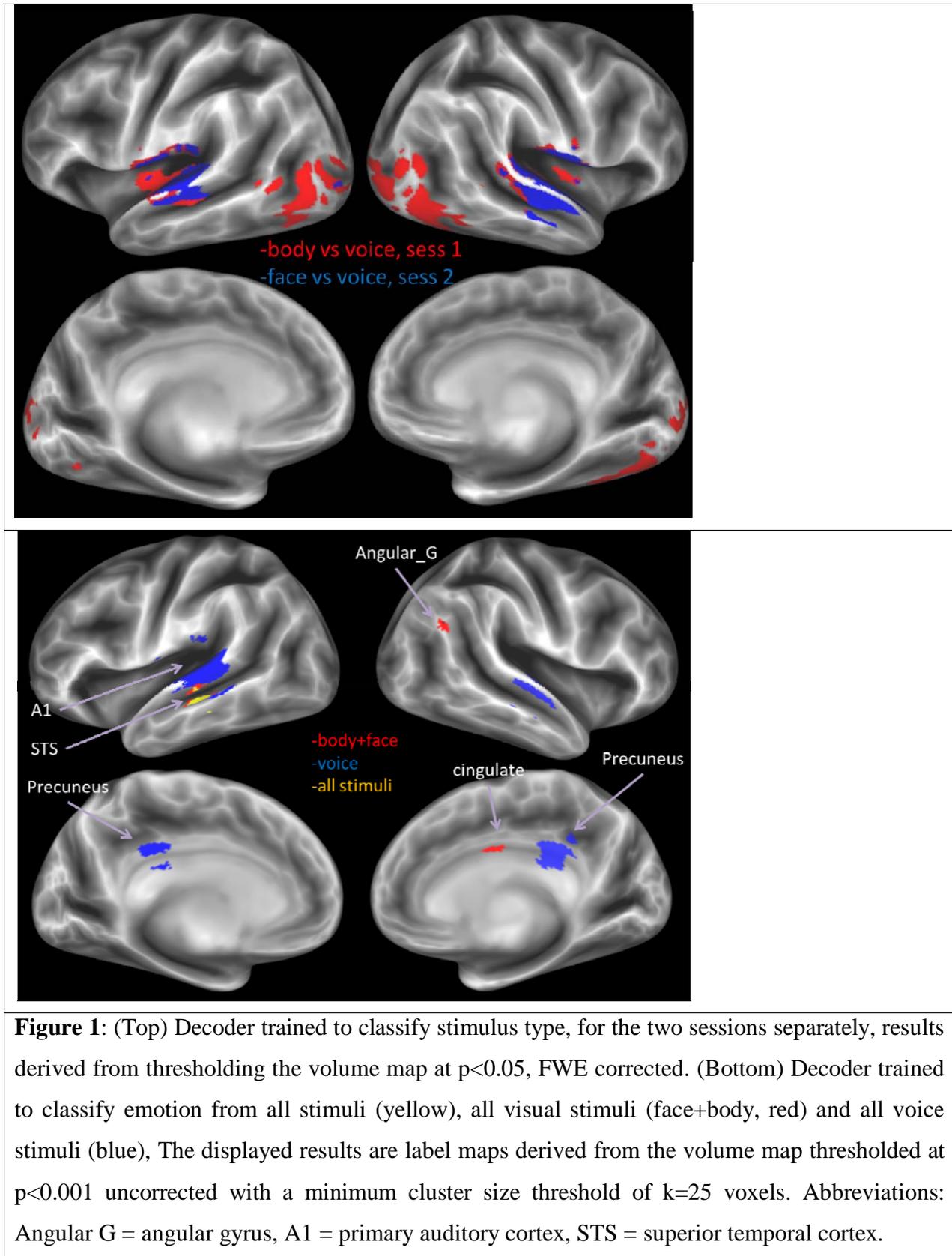
236 We ran an ANOVA with factors ‘stimulus type’ and ‘emotion’ on the beta values estimated with  
237 an RFX GLM on the entire data set, as well as separate ANOVAs with factor ‘emotion’ on the  
238 data of each stimulus type (i.e. faces, bodies, and voices). As expected, the F-map for stimulus  
239 type (see Fig. S1) revealed significant clusters with differential mean activation across stimulus  
240 types in primary and higher-order auditory and visual regions, as well as in motor, pre-motor and  
241 dorsal/superior parietal cortex. The univariate analyses did not reveal cortical regions showing a  
242 significant effect for emotion. Although the F-maps of these emotion analyses showed some  
243 small clusters at lenient thresholds ( $p < 0.001$ , uncorrected), none of these survived correction for  
244 multiple comparisons. See Supplementary results for details.

245 ***Multivariate analysis with the GNB decoder***

246 We aimed to decode stimulus type by training the classifier separately for the two sessions to  
247 decode visual vs. auditory stimuli: body vs. voice (session 1) and face vs. voice (session2). As  
248 expected, stimulus type could be decoded significantly above chance level (50%,  $p < 0.05$  FWE  
249 corrected) in auditory, visual and fusiform cortex, and large part of the lateral occipital and  
250 temporo-occipital cortex, presumably including the extrastriatal body area (EBA; see Fig. 1 top  
251 panel).

252 Decoding of emotion resulted in qualitatively lower accuracies and smaller clusters compared to  
253 the decoding of modality. Above chance accuracies (33%) for decoding emotion from all  
254 stimulus types together were observed in STS (Fig. 1, bottom panel). Next, we trained and tested  
255 the classifier to decode emotions within a specific stimulus modality, that is, considering either  
256 the combined face and body stimuli (i.e. visual), or the voice stimuli (i.e. auditory). For body and  
257 face this revealed STS, cingulate gyrus and angular gyrus. Emotion could be decoded for voice  
258 stimuli in primary and secondary auditory regions (including the superior temporal gyrus  
259 [STG]), and the precuneus (see Fig 1 bottom panel). See Table 1 for details of these results.

260 We also trained and tested the classifier to decode emotions within each visual category, that is,  
261 considering either the face or the body. This did not reveal any clusters where emotion could be  
262 decoded accurately (suprathreshold at  $p < 0.05$  FWE corrected), although additional results were  
263 obtained at more lenient thresholds (see Fig. S4).



265 **Table 1: Results for the decoding of emotion**

	<i>cluster size</i>	<i>clusgter p(unc)</i>	<i>peak T</i>	<i>peak p(unc)</i>	<i>x</i>	<i>y</i>	<i>z</i>
<b>All stimuli</b>							
STS	38	0.0536	5.8422	0.0000	-3	21	-7
<b>Voice stimuli only</b>							
planum temporale	749	0.001	6.462	0.000	-58	-34	6
posterior cingulate gyrus / precuneus	187	0.058	5.654	0.000	8	-34	30
planum temporale	210	0.046	5.206	0.000	62	-18	2
<b>Face and body stimuli</b>							
STS	59	0.0200	6.8069	0.0000	-52	-26	-12
right caudate	92	0.0052	6.5513	0.0000	16	-12	20
angular gyrus	27	0.0979	6.4339	0.0000	60	-50	34

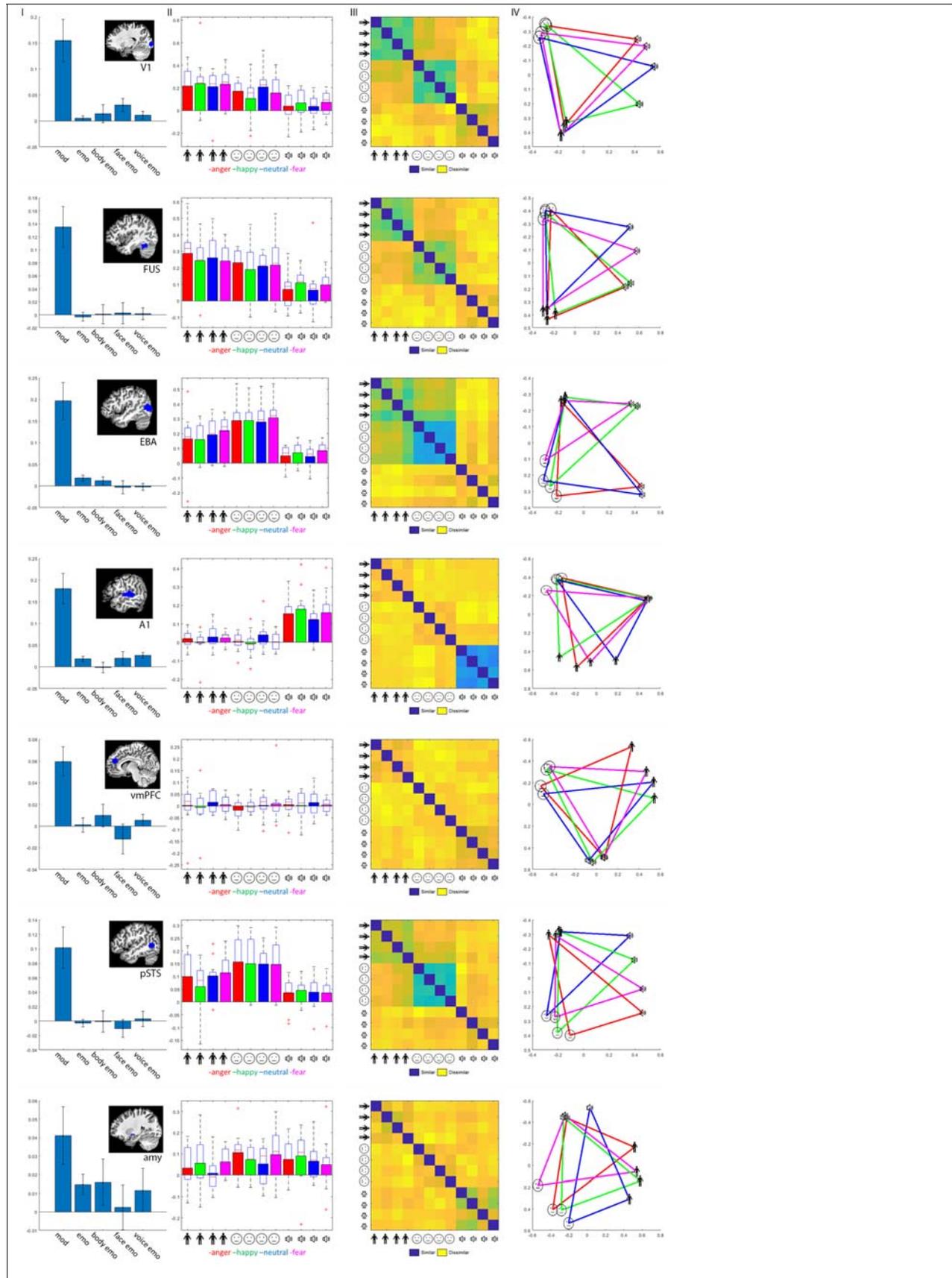
266

267 **ROI Analysis**

268 In addition to the whole-brain searchlight analysis, we performed a more sensitive region-of-  
 269 interest (ROI) analysis. Specifically, we used data of an independent localizer (see *Methods*) to  
 270 identify early auditory cortex, early visual cortex, rEBA, and rFFA. We furthermore included  
 271 two multisensory regions, pSTS and mPFC that were previously identified as regions holding  
 272 supramodal representations of emotion (Peelen MV *et al.* 2010). We used an anatomical  
 273 definition of these areas, defined by a spherical ROI with a radius of 5 voxels centered on the  
 274 reported cluster peak locations. Finally, we included the amygdala given its important role in  
 275 previous studies, most recently in the study by (Whitehead JC and JL Armony 2019) using the  
 276 univariate contrast face fear > face neutral (p<0.01 uncorrected).

277 In all ROIs, stimulus type could be decoded above chance level (one sample t-test against chance  
 278 level, all p<0.0001). Furthermore, when the classifier operated on all stimuli together (that is,  
 279 face, body, and voice), emotion could be successfully decoded in the EBA, auditory cortex, and  
 280 left amygdala (all p<0.02). That is, in line with the results of the searchlight analysis, when the  
 281 classifier operated on the data of each modality in isolation, decoding accuracies for emotion  
 282 were above chance level in early visual cortex for face stimuli (p<0.03), and in auditory cortex  
 283 for voice stimuli (p<0.002). Emotion could not be decoded above chance level in the supramodal

284 regions (mPFC and pSTS). Taken together, these results demonstrate that these regions could not  
285 be identified as supramodal as the responses were not invariant to stimulus type (see Fig. 2).  
286 Qualitatively, the RDM and MDS plots in Figure 2 show that in all tested ROIs there is a strong  
287 effect of stimulus type (blocks on diagonal for the RDM and a large distance between types in  
288 the MDS). Notable, this effect is not always clearly present in just the beta magnitudes (response  
289 levels) and is at least partially caused by the multi-voxel pattern dissimilarities.



**Figure 2:** I) Location of the ROI and graphs indicating decoding accuracy for stimulus type and emotion and emotion for each stimulus type in the ROI. II) Trial-wise beta values for all 12 conditions (3x modality and 4x emotion) averaged over the ROI and group. Error bars indicate SE. III) Representational dissimilarity matrix for the ROI. IV) Multidimensional scaling plot of the group averaged trial-wise beta values. Line colors indicate same emotion, icons display modality. Distance in the plot is related to similarity of the ROI voxel activation patterns. Colors as in the beta plot for emotion.

290

### 291 *Cross-modal decoding*

292 We performed an additional analysis to gain insight into where in the brain supramodal emotion  
293 regions might be found by training a classifier to decode emotion across stimulus modalities.  
294 Being able to predict emotion by training on one modality and testing on another modality would  
295 be a strong indication of supramodal emotion encoding in the brain. Therefore, the cross-modal  
296 classifier was trained (or tested) on either the body or face stimuli and tested (or trained) on the  
297 voice stimuli from the body or face session, respectively. Thus, four classifiers were trained in  
298 total (training on body and testing on voice, training on voice and testing body, training on face  
299 and testing on voice, training on voice and testing on face). In contrast to the successful decoding  
300 of modality and emotion within modality, none of these cross-modal classifiers resulted in  
301 accurate decoding of emotion ( $p < 0.001$ , see Supplementary Fig. S5).

### 302 *Contrasting accuracies for emotion and modality*

303 There was little overlap in regions identified by the within stimulus type classifier (see Fig. 2 and  
304 S4). Specifically, regions where emotion was decoded successfully from body stimuli did not  
305 converge with regions where emotion was decoded successfully from the face and/or voice.  
306 Therefore, to identify regions that have a purely supramodal representation of emotion, we  
307 contrasted the accuracy map for modality versus emotion directly. Here, finding regions having  
308 higher decoder accuracy for emotion compared to modality would be a strong indication for  
309 supramodal emotion encoding. This analysis revealed in general that, as expected, modality  
310 could be decoded with greater accuracy than emotion in primary auditory and visual regions, as  
311 well as in parietal and pre-motor regions. We could not clearly identify any regions where binary

312 combination of two specific emotions could be decoded with higher accuracy than the modality.  
313 Only for fear and neutral in the face-voice session did we find a region in the medial motor  
314 cortex and for happy and neutral in the body-voice session in the white matter (at  $p < 0.001$   
315 uncorrected, see Fig. S6).

## 316 **DISCUSSION**

317 Our results indicate that the brain regions involved in emotion processing are modality specific.  
318 That is, in the regions where emotion can be decoded, it could only be decoded within modality  
319 but not across modalities, indicating that the decoding is driven by a specific response to a  
320 specific emotion-modality combination. In a departure from the few previous studies using a  
321 partially comparable approach we found evidence for sensory specific rather than abstract  
322 supramodal representations that sustain perception of various affective signals as a function of  
323 the modality (visual or auditory) and the stimulus category (face, voice or body). Because of the  
324 three stimulus categories used, the different emotions studied, the task conditions, and the  
325 converging results from different analysis techniques, our study presents a novel approach to  
326 understanding the specific contribution and the neural basis of emotional signals provided by  
327 different sensory modalities.

328 To understand our findings against the background of the literature, some specific aspects  
329 of our study must be highlighted. We used dynamic realistic face and body stimuli instead of  
330 point light displays or static images. The latter are also known to complicate comparisons with  
331 dynamic auditory stimuli (Campanella S and P Belin 2007). Next, our stimuli do not present  
332 prototypical emotion representations obtained by asking actors to portray emotions but present  
333 spontaneous whole body reactions to images of familiar events. The images we used may  
334 therefore be more spontaneous and trigger more sensorimotor processes in the viewer than posed  
335 expressions. Third, many previous studies used explicit emotion recognition (Lee KH and GJ  
336 Siegle 2012), passive viewing (Winston JS et al. 2003), implicit tasks like gender categorization  
337 (Dricu M and S Fruhholz 2016) or oddball tasks presented in the same modality as the stimulus.  
338 In contrast, our modality specific oddball task is presented in the alternate modality of the  
339 stimulus presentation thereby diverting attention not only from the emotion content but also from  
340 the perceptual modality in which the target stimuli of that block are shown. This task was  
341 intended to approximate the naturalistic experience of emotional signals, where often one is

342 engaged in one activity (visual perception) when another event intrudes (an auditory event). We  
343 discuss separately the findings on the major research questions.

344

### 345 *Univariate analysis*

346 Although our goal was to characterize neural responses with MVPA techniques, for the  
347 sake of comparisons with the literature, we also briefly discuss our univariate results. How do  
348 these results compare to findings and meta-analyses in the literature? As a matter of fact, there  
349 are no previous studies that used comparable materials (four emotion categories, three stimulus  
350 types, two modalities) and an other modality centered task like the present. The studies that did  
351 include bodies used only neutral actions, not whole body emotion expressions (Dricu M and S  
352 Fruhholz 2016) except for one study comparing face and body expression videos by (Kret ME, S  
353 Pichon, J Grezes, *et al.* 2011). Only the study by Peelen *et al.* used faces, bodies, and voices, but  
354 with a very different task as we discuss below (Peelen MV *et al.* 2010).

355 Compared to the literature, the findings of the univariate analysis present  
356 correspondences as well as differences. A previous study (Kret ME, S Pichon, J Grezes, *et al.*  
357 2011) with face and body videos used only neutral, fear and anger expression and a visual  
358 oddball task. They reported that EBA and STS show increased activity to threatening body  
359 expressions and FG responds equally to emotional faces and bodies. For the latter, higher activity  
360 was found in cuneus, fusiform gyrus, EBA, tempo-parietal junction, superior parietal lobe, as  
361 well in as the thalamus while the amygdala was more active for facial than for bodily  
362 expressions, but independently of the facial emotion. Here we replicate that result for faces and  
363 bodies and found highly significant clusters with differential mean activation across stimulus  
364 types in primary and higher-order auditory and visual regions, as well as in motor, pre-motor and  
365 dorsal/superior parietal cortex (Fig. S1).

366

367 Regions sensitive to stimulus category were not only found in primary visual and auditory cortex  
368 as expected but also in motor, pre-motor and dorsal/superior parietal cortex consistent with the  
369 findings in Kret *et al.* 2011. In view of the literature on perception of emotion expressions in  
370 either the face, voice or body it is not surprising that dorsal parietal cortex, pre- motor cortex and  
371 anterior insula (Fig. S2) differentially respond to emotions as different emotions trigger different  
372 adaptive actions (de Gelder B 2006; Grezes J *et al.* 2007; Pichon S *et al.* 2008; Whitehead JC and

373 JL Armony 2019). Interestingly, the interaction effect between modality and emotion also  
374 revealed the retrosplenial cortex (Fig. S3). Retrosplenial cortex receives input from areas known  
375 to play a role in processing salient information (prefrontal cortex, superior temporal sulcus,  
376 precuneus, thalamus, and claustrum (Maddock RJ 1999). The retrosplenial cortex is associated  
377 with navigation and memory functions and may be part of a network that conveys predatory  
378 threat information to the cerebral cortex (de Lima MAX et al. 2019). A similar functionality may  
379 be reflected in this activity here. Activity in this area is also consistent with the recent findings  
380 that the retrosplenial cortex contributed to the accurate classification of fear stimuli. (Caruana F  
381 et al. 2018).

382 To summarize, this univariate analysis including three stimulus modalities and four  
383 emotion categories replicates some main findings about brain areas involved respectively in face,  
384 body and voice expressions while also revealing parietal and motor area activity but does not  
385 provide evidence for overlap in brain activity neither for modality nor emotion category.

### 386 387 *Multivariate analysis*

388 The goal of our multivariate approach was to reveal the areas that contribute most  
389 strongly to an accurate distinction between the modalities and the stimulus emotion. Our results  
390 of the MVPA searchlight show that modality type can be decoded from the sensory cortex and  
391 that emotion can be decoded in STG for voice stimuli and in STS for face and body stimuli. We  
392 found no overlap in brain regions that contribute to the classification of emotion when using  
393 either the face, the body or the voice decoder (Fig. S4). Thus the brain areas that are involved in  
394 discriminating between face, voice or body expressions irrespective of the emotion are different  
395 from each other. Nor did we find evidence for neural activity overlap the other way around when  
396 using a cross-modal emotion decoder and looking for possible brain areas common to the  
397 modalities (Fig. S5). Lastly, we could not clearly identify regions where emotion could be  
398 decoded with higher accuracy than modality (Fig. S6). To put it negatively, we could not clearly  
399 identify supramodal emotion regions, defined by voxel patterns where emotion could be decoded  
400 and that would show very similar voxel patterns for the same emotion in the different modalities.  
401 This clearly indicated that the brain responds to facial, body and vocal emotion expression in a  
402 unique fashion. Thus the overall direction pointed to by our results seems to be that that being  
403 exposed to emotional stimuli (that are not task relevant and while performing a task requiring

404 attention to the other modality than that in which the stimulus is presented) is associated with  
405 brain activity that shows both an emotion specific and a modality specific pattern.

406

### 407 ***ROI analysis***

408 To follow up on the whole-brain analysis we performed a detailed and specific analysis  
409 of a number of ROIs. For the ROIs based on the localizer scans (early visual, auditory areas as  
410 well as FFA and EBA) stimulus type could be decoded and these results are consistent with the  
411 MVPA searchlight analysis. We also defined ROIs based on the literature with the goal of  
412 comparing our results to findings about higher order areas that had been identified as  
413 supramodal areas, mPFC and pSTS. This analysis revealed accurate decoding of stimulus type  
414 but no evidence of supramodal representations.

415 Two previous MVPA studies addressed similar issues investigated in this study using  
416 faces, bodies and voices (Peelen MV *et al.* 2010) or bodies and voices (Whitehead JC and JL  
417 Armony 2019). The first study reported medial prefrontal cortex (MPFC) and posterior superior  
418 temporal cortex as the two areas hosting supramodal emotion representations. These two areas  
419 do not emerge in our MVPA searchlight analysis. To understand this very different result it is  
420 important to remember that in their study participants were encouraged to actively evaluate the  
421 perceived emotional states. The motivation was that explicit judgments would increase activity  
422 in brain regions involved in social cognition and mental state attribution (Peelen MV *et al.*  
423 2010). In contrast, the motivation of the present study was to approximate naturalistic perception  
424 conditions. Our design and task were intended to promote spontaneous non-focused processes of  
425 the target stimuli and did not promote amodal conceptual processing of the emotion content. It is  
426 likely that using an explicit recognition task would have activated higher level representations  
427 e.g. posterior STS, prefrontal cortex and posterior cingulate cortex that would then feedback to  
428 lower level representations and modulate these towards more abstract representations (Schirmer  
429 A and R Adolphs 2017). Note that no amygdala activity was reported in that study. The second  
430 study using passive listening or viewing of still bodies and comparing fear and neutral  
431 expressions also concludes about a distributed network of cortical and subcortical regions  
432 responsive to fear in the two stimulus types they used (Whitehead JC and JL Armony 2019). Of  
433 interest is their finding concerning the amygdalae and fear processing. While in their study this is  
434 found across stimulus type for body and voice, the classification accuracy when restricted to the

435 amygdalae was not significantly above chance. They concluded that fear processing by the  
436 amygdalae heavily relies on contribution of a distributed network of cortical and subcortical  
437 structures.

438 Our findings suggest a different and a novel perspective on the role of the different  
439 sensory systems and the different stimulus categories that convey affective signals in daily life  
440 and fits with the role of emotions as seen from an evolutionary perspective. Our results are  
441 compatible with an ecological and context sensitive approach to brain organization (Cisek P and  
442 JF Kalaska 2010; Mobbs D et al. 2018) rather than with an exclusive focus on high order  
443 representation of emotion categories grounded in concepts and verbal labels. For comparison, a  
444 similar approach not to emotion concepts but to cognitive concepts was argued by Barsalou et al.  
445 (Barsalou LW et al. 2003). This type of distributed organization or emotion representation may  
446 be more akin to what is at stake in the daily experience of affective signals and how they are  
447 flexibly processed for the benefit of ongoing action and interaction in a broader perspective of  
448 emotions as states of action readiness (Frijda NH 2004).

449 Our results are relevant for two longstanding debates in the literature, one on the nature  
450 and existence of abstract emotion representations and basic categories and the other on processes  
451 of multisensory integration. Concerning the first one, our results have implications for the debate  
452 on the existence of basic emotions (Ekman P 2016). Interestingly, modality specificity has rarely  
453 been considered as part of the issue as the basic emotion debate largely focusses on facial  
454 expressions. The present results might be viewed as evidence in favor of the view that basic  
455 emotions traditionally understood as specific representations of a small number of emotions with  
456 an identifiable brain correlate (Ekman P 2016) simply do not exist but that these are cognitive-  
457 linguistic constructions (Russell JA 2003). On the one hand, our results are consistent with  
458 critiques of basic emotions theories and meta-analysis (Lindquist KA *et al.* 2012) as we find no  
459 evidence for representations of emotions in general or specific emotions within or across  
460 modality and stimuli. Affective information processing thus appears not organized as  
461 categorically, neither by conceptual emotion category nor by modality, as was long assumed.  
462 Emotion representation, more so even than object representation, may possibly be sensory  
463 specific or idiosyncratic (Peelen MV and PE Downing 2017) and neural representations may  
464 reflect the circumstances under which specific types of signals are most useful or relevant rather  
465 than abstract category membership. This pragmatic perspective is consistent with the notion that

466 emotions are closely linked to action and stresses the need for more detailed ethological behavior  
467 investigations (de Gelder B 2016) .

468

469 While our study was not addressing issues of multisensory perception, our findings may  
470 have implications for theories of multisensory integration. As has often been noted, human  
471 emotion research by and far is still limited to the study of facial expressions. In line with the  
472 dominant view studies extending the scope of facial expression based theories have primarily  
473 been motivated to discover similarities across different modalities and stimulus types. Our group  
474 has initiated studies that go beyond the facial expression and found rapid and automatic influence  
475 of one type of expression on another (face and voice, (de Gelder B et al. 1999); face and body  
476 (Meeren HK et al. 2005); face and scene (Righart R and B de Gelder 2008; Van den Stock J et al.  
477 2013); body and scene, (Van den Stock J et al. 2014); auditory voice and tactile perception (de  
478 Borst AW and B de Gelder 2017). These original studies and subsequent ones (Müller VI et al.  
479 2012) investigated the impact of one modality on the other and targeted the area(s) where  
480 different signals converge. For example, Müller et al. (Müller VI *et al.* 2012) report posterior  
481 STS as the site of convergence of auditory and visual input systems and by implication, as the  
482 site of multisensory integration. Note that here and in many other studies the amygdala and its  
483 connectivity to face and voice areas emerges as a core structure involved in multisensory  
484 integration. Current studies did not yet raise the question of sensory specificity of the  
485 representations at stake in such cross-modal effects. Our findings stress the importance of  
486 sensory specific representation and indicate that aside from loci of integration based presumably  
487 on amodal representations, multisensory perception and integration seems to respect stimulus  
488 complementarity in the vertical plane rather than convergence between different emotion signals  
489 onto an abstract supramodal representation.

490

491 The motivation to include three stimulus categories led to some limitations of the current  
492 design because two separate scanning sessions were required to have the desired number of  
493 stimulus repetitions. To avoid that the comparison of representations of stimuli from two  
494 different sessions was biased by a scan session effect, we did not include any results that referred  
495 to differences or commonalities of stimuli from different sessions (e.g. bodies vs faces). A  
496 second limitation is that the voice stimuli were not controlled for low level acoustic features over

497 different emotions and therefore results from the decoding of emotion from the voice stimuli  
498 may partly have reflected these differences.

#### 499 *Conclusions*

500 Our results show that the brain correlates of observing emotional signals from the body,  
501 the face of the voice are specific for the modality as well as for the specific signal within the  
502 same modality. Our results underscore the importance of considering the specific contribution of  
503 each modality and each type of affective signal rather than only their higher order amodal  
504 similarity. We suggest that future research may look into the differences between the emotion  
505 signals and how they are complementary and not only at amodal similarity.

506

#### 507 *Acknowledgments*

508 We would like to thank Dr. Rebecca Watson for acquiring the MRI data. This work was  
509 supported by the European Research Council (ERC) FP7-IDEAS-ERC Grant agreement number  
510 295673 *Emobodies*, by the Future and Emerging Technologies (FET) Proactive Programme  
511 H2020-EU.1.2.2 Grant agreement 824160 *EnTimeMent* and by the Industrial Leadership  
512 Programme H2020-EU.1.2.2 Grant agreement 825079 *MindSpaces*.

513

## 514 REFERENCES

- 515 Barrett LF. 2017. How emotions are made: The secret life of the brain: Houghton Mifflin Harcourt.
- 516 Barsalou LW, Simmons WK, Barbey AK, Wilson CD. 2003. Grounding conceptual knowledge in modality-  
517 specific systems. *Trends in cognitive sciences* 7:84-91.
- 518 Campanella S, Belin P. 2007. Integrating face and voice in person perception. *Trends in cognitive*  
519 *sciences* 11:535-543.
- 520 Caruana F, Gerbella M, Avanzini P, Gozzo F, Pelliccia V, Mai R, Abdollahi RO, Cardinale F, Sartori I, Lo  
521 Russo G, Rizzolatti G. 2018. Motor and emotional behaviours elicited by electrical stimulation of the  
522 human cingulate cortex. *Brain : a journal of neurology* 141:3035-3051.
- 523 Cisek P, Kalaska JF. 2010. Neural mechanisms for interacting with a world full of action choices. *Annual*  
524 *review of neuroscience* 33:269-298.
- 525 de Borst AW, de Gelder B. 2017. fMRI-based Multivariate Pattern Analyses Reveal Imagery Modality and  
526 Imagery Content Specific Representations in Primary Somatosensory, Motor and Auditory Cortices.  
527 *Cerebral cortex* 27:3994-4009.
- 528 de Gelder B. 2006. Towards the neurobiology of emotional body language. *Nature reviews Neuroscience*  
529 7:242-249.
- 530 de Gelder B. 2016. Emotional body perception in the wild. LF, Barrett, M, Lewis, JM Haviland-  
531 Jones, (Eds), *Handbook of emotions* (4th ed, pp 483-494): New York, NY: Guilford Press Publications.
- 532 de Gelder B, Bocker KB, Tuomainen J, Hensen M, Vroomen J. 1999. The combined perception of emotion  
533 from voice and face: early interaction revealed by human electric brain responses. *Neuroscience letters*  
534 260:133-136.
- 535 de Lima MAX, Baldo MVC, Canteras NS. 2019. Revealing a Cortical Circuit Responsive to Predatory  
536 Threats and Mediating Contextual Fear Memory. *Cerebral cortex* 29:3074-3090.
- 537 Dricu M, Fruhholz S. 2016. Perceiving emotional expressions in others: Activation likelihood estimation  
538 meta-analyses of explicit evaluation, passive perception and incidental perception of emotions.  
539 *Neuroscience and biobehavioral reviews* 71:810-828.
- 540 Ekman P. 1992. An argument for basic emotions. *Cognition & emotion* 6:169-200.
- 541 Ekman P. 2016. What scientists who study emotion agree about. *Perspectives on Psychological Science*  
542 11:31-34.
- 543 Ekman P, Cordaro D. 2011. What is meant by calling emotions basic. *Emotion review* 3:364-370.
- 544 Frijda NH`editor. Year Published|. Title|, Conference Name|; Year of Conference Date|; Conference  
545 Location| Place Published|:Publisher|. Pages p|.
- 546 Gerdes AB, Wieser MJ, Alpers GW. 2014. Emotional pictures and sounds: a review of multimodal  
547 interactions of emotion cues in multiple domains. *Front Psychol* 5:1351.
- 548 Grezes J, Pichon S, De Gelder B. 2007. Perceiving fear in dynamic body expressions. *NeuroImage* 35:959-  
549 967.
- 550 Hausfeld L, De Martino F, Bonte M, Formisano E. 2012. Pattern analysis of EEG responses to speech and  
551 voice: influence of feature grouping. *NeuroImage* 59:3641-3651.
- 552 Hausfeld L, Valente G, Formisano E. 2014. Multiclass fMRI data decoding and visualization using  
553 supervised self-organizing maps. *NeuroImage* 96:54-66.
- 554 Karle KN, Ethofer T, Jacob H, Bruck C, Erb M, Lotze M, Nizielski S, Schutz A, Wildgruber D, Kreifelts B.  
555 2018. Neurobiological correlates of emotional intelligence in voice and face perception networks. *Social*  
556 *cognitive and affective neuroscience* 13:233-244.
- 557 Klasen M, Kenworthy CA, Mathiak KA, Kircher TT, Mathiak K. 2011. Supramodal representation of  
558 emotions. *The Journal of neuroscience : the official journal of the Society for Neuroscience* 31:13635-  
559 13643.

- 560 Kragel PA, LaBar KS. 2016. Decoding the Nature of Emotion in the Brain. *Trends in cognitive sciences*  
561 20:444-455.
- 562 Kret ME, Pichon S, Grezes J, de Gelder B. 2011. Similarities and differences in perceiving threat from  
563 dynamic faces and bodies. An fMRI study. *NeuroImage* 54:1755-1762.
- 564 Kret ME, Pichon S, Grèzes J, De Gelder B. 2011. Men fear other men most: gender specific brain  
565 activations in perceiving threat from dynamic faces and bodies—an fMRI study. *Frontiers in psychology*  
566 2:3.
- 567 Kriegeskorte N, Mur M, Ruff DA, Kiani R, Bodurka J, Esteky H, Tanaka K, Bandettini PA. 2008. Matching  
568 categorical object representations in inferior temporal cortex of man and monkey. *Neuron* 60:1126-  
569 1141.
- 570 Lee KH, Siegle GJ. 2012. Common and distinct brain networks underlying explicit emotional evaluation: a  
571 meta-analytic study. *Social cognitive and affective neuroscience* 7:521-534.
- 572 Lewis M, Haviland-Jones JM, Barrett LF. 2010. *Handbook of emotions*: Guilford Press.
- 573 Lindquist KA, Wager TD, Kober H, Bliss-Moreau E, Barrett LF. 2012. The brain basis of emotion: a meta-  
574 analytic review. *The Behavioral and brain sciences* 35:121.
- 575 Maddock RJ. 1999. The retrosplenial cortex and emotion: new insights from functional neuroimaging of  
576 the human brain. *Trends in neurosciences* 22:310-316.
- 577 Meeren HK, van Heijnsbergen CC, de Gelder B. 2005. Rapid perceptual integration of facial expression  
578 and emotional body language. *Proceedings of the National Academy of Sciences of the United States of*  
579 *America* 102:16518-16523.
- 580 Mobbs D, Trimmer PC, Blumstein DT, Dayan P. 2018. Foraging for foundations in decision neuroscience:  
581 insights from ethology. *Nature reviews Neuroscience* 19:419-427.
- 582 Müller VI, Cieslik EC, Turetsky BI, Eickhoff SB. 2012. Crossmodal interactions in audiovisual emotion  
583 processing. *NeuroImage* 60:553-561.
- 584 Ontivero-Ortega M, Lage-Castellanos A, Valente G, Goebel R, Valdes-Sosa M. 2017. Fast Gaussian Naive  
585 Bayes for searchlight classification analysis. *NeuroImage* 163:471-479.
- 586 Peelen MV, Atkinson AP, Vuilleumier P. 2010. Supramodal representations of perceived emotions in the  
587 human brain. *The Journal of neuroscience : the official journal of the Society for Neuroscience* 30:10127-  
588 10134.
- 589 Peelen MV, Downing PE. 2017. Category selectivity in human visual cortex: Beyond visual object  
590 recognition. *Neuropsychologia* 105:177-183.
- 591 Pichon S, de Gelder B, Grezes J. 2008. Emotional modulation of visual and motor areas by dynamic body  
592 expressions of anger. *Social neuroscience* 3:199-212.
- 593 Righart R, de Gelder B. 2008. Rapid influence of emotional scenes on encoding of facial expressions: an  
594 ERP study. *Social cognitive and affective neuroscience* 3:270-278.
- 595 Russell JA. 2003. Core affect and the psychological construction of emotion. *Psychological review*  
596 110:145.
- 597 Saarimäki H, Gotsopoulos A, Jaaskelainen IP, Lampinen J, Vuilleumier P, Hari R, Sams M, Nummenmaa L.  
598 2016. Discrete Neural Signatures of Basic Emotions. *Cerebral cortex* 26:2563-2573.
- 599 Schirmer A, Adolphs R. 2017. Emotion Perception from Face, Voice, and Touch: Comparisons and  
600 Convergence. *Trends in cognitive sciences* 21:216-228.
- 601 Van den Stock J, Vandenbulcke M, Sinke CB, de Gelder B. 2014. Affective scenes influence fear  
602 perception of individual body expressions. *Human brain mapping* 35:492-502.
- 603 Van den Stock J, Vandenbulcke M, Sinke CB, Goebel R, de Gelder B. 2013. How affective information  
604 from faces and scenes interacts in the brain. *Social cognitive and affective neuroscience* 9:1481-1488.
- 605 Whitehead JC, Armony JL. 2019. Multivariate fMRI pattern analysis of fear perception across modalities.  
606 *The European journal of neuroscience* 49:1552-1563.

607 Winston JS, O'Doherty J, Dolan RJ. 2003. Common and distinct neural responses during direct and  
608 incidental processing of multiple facial emotions. *NeuroImage* 20:84-97.

609 Zhan M, Goebel R, de Gelder B. 2018. Ventral and Dorsal Pathways Relate Differently to Visual  
610 Awareness of Body Postures under Continuous Flash Suppression. *eNeuro* 5.

611

612

613 **Supplementary results**

614 *Univariate*

615 The F-map for emotion showed regions with different activation levels across emotions in  
616 insular cortex, the superior temporal sulcus (STS), retrosplenial cortex, angular gyrus, and  
617 superior medial occipital cortex, see Fig. S2.

618 We ran separate ANOVAs with factor ‘emotion’ for each stimulus type condition (voices, faces,  
619 bodies; see Fig. S2), to identify regions that respond differentially across emotions within a  
620 specific stimulus type. For the face condition, this revealed clusters in STS and retrosplenial  
621 cortex. For the body condition, we observed regions in inferior temporal cortex and the  
622 intraparietal sulcus (IPS). Finally, the ANOVA for the voice condition (for session 1 and session  
623 2 separately) revealed auditory cortex, premotor cortex, IPS and angular gyrus.

624 We then tested for an interaction effect between the factors stimulus type and emotion. This  
625 revealed clusters in the retrosplenial cortex, the dorsolateral prefrontal cortex (dLPFC), auditory  
626 cortex and insula, see Fig. S3.

627 **Supplementary Figure 1**

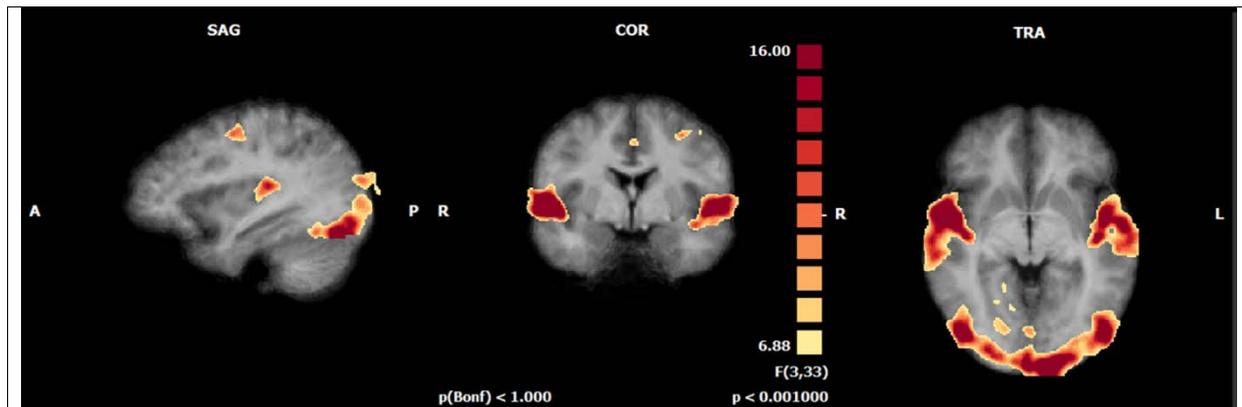
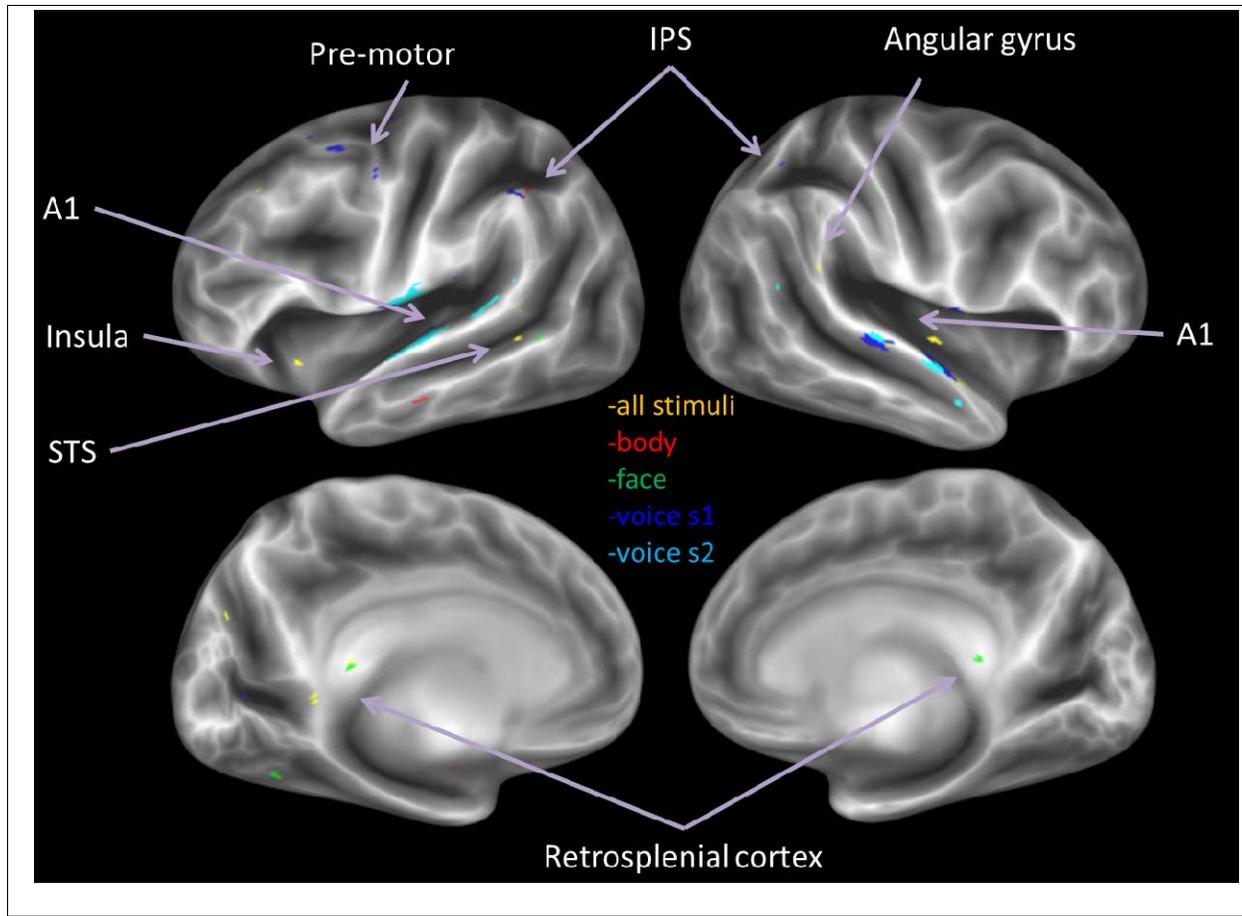


Figure S1: F-map ( $p < 0.001$  uncorrected) for stimulus type effect

628

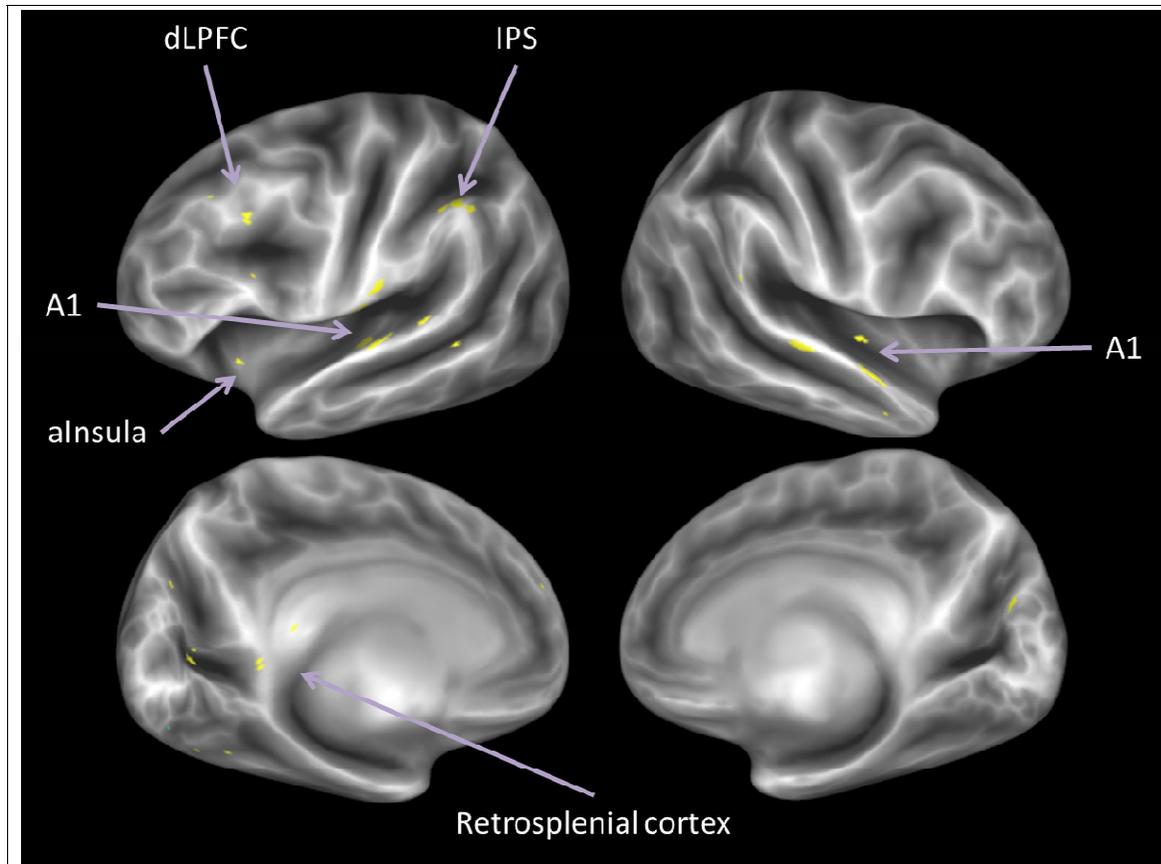
629 **Supplementary Figure 2**



**Figure S2.** F-map for the ANOVA with factor 'emotion', calculated either for all stimuli (orange) or for each stimulus type separately (bodies = red; faces = green; voices session 1 = dark blue; voices session 2 = light blue). The displayed results are label maps derived from the volume map thresholded at  $p < 0.001$  uncorrected.

630

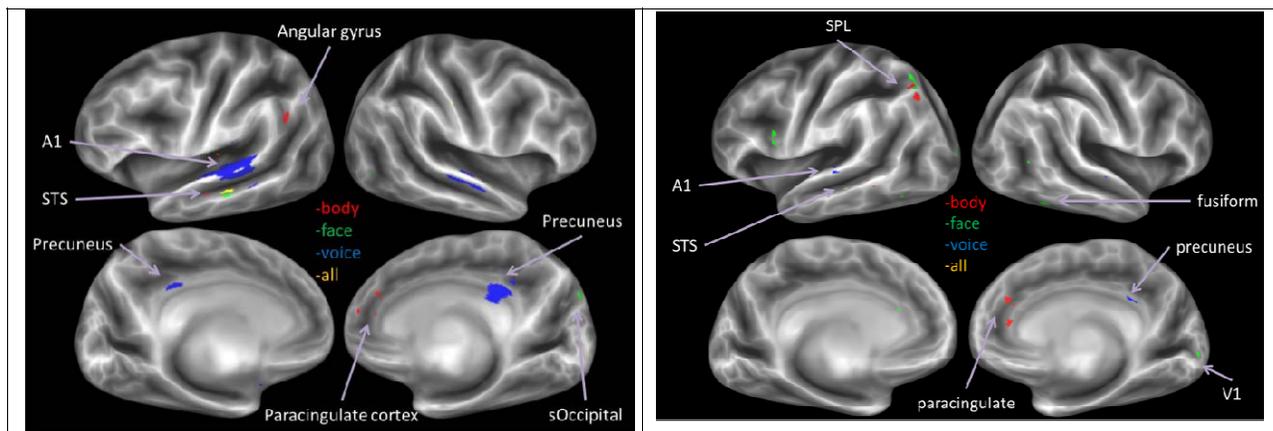
631 **Supplementary Figure3**



**Figure S3.** Interaction emotion and stimulus type. The displayed results are label maps derived from the volume map thresholded at  $p < 0.001$  uncorrected.

632

633 **Supplementary Figure 4**



**Figure S4:** (Left) Decoder trained to classify emotion, for all stimuli and for each stimulus type separately each modality separately. (Right) Decoder trained to classify emotion excluding the

neutral stimuli, for all stimuli and for each stimulus type separately each modality separately. The displayed results are label maps derived from the volume map thresholded at  $p < 0.001$  uncorrected.

634

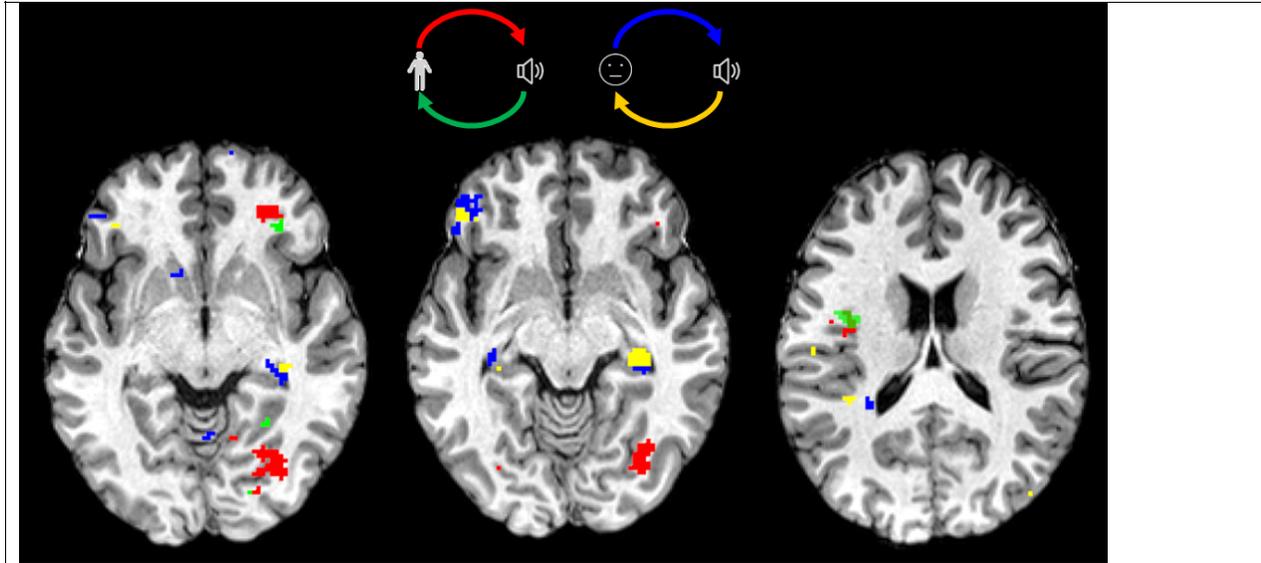
635 **Supplementary table 1:**

	cluster size	cluster (unc)	p	peak T	Peak p unc	x	y	z
<b>All emotion</b>								
planum temporale	22	0.505		4.237	0.001	-60	-28	14
<b>Emotion, no neutral</b>								
inferior frontal gyrus	26	0.111		5.868	0.000	-44	4	16
<b>Emotion from body only</b>								
paracingulate gyrus	76	0.213		7.364	0.000	2	50	8
<b>Emotion from body only, no neutral</b>								
frontal pole	22	0.150		6.957	0.000	26	58	18
intraparietal sulcus	34	0.079		6.048	0.000	-34	-62	48
superior temporal sulcus	25	0.127		5.521	0.000	-60	-52	4
paracingulate gyrus	39	0.062		4.915	0.000	2	50	8
<b>Emotion from face only</b>								
superior lateral occipital cortex	31	0.442		5.441	0.000	14	-80	40
precuneus / posterior cingulate gyrus	23	0.513		4.569	0.000	10	-36	28
<b>Emotion from face only, no neutral</b>								
ventrolateral prefrontal cortex	26	0.121		8.237	0.000	-48	6	18
temporooccipital cortex	68	0.018		6.286	0.000	46	-50	-6
occipital cortex / V3	22	0.151		5.055	0.000	-26	-98	10
superior parietal lobe	20	0.170		4.778	0.000	-24	-62	52
<b>Emotion from voice</b>								
planum temporale	749	0.001		6.462	0.000	-58	-34	6
posterior cingulate gyrus / precuneus	187	0.058		5.654	0.000	8	-34	30

planum temporale	210	0.046	5.206	0.000	62	-18	2
precentral gyrus	50	0.307	4.629	0.000	-30	2	38
<b>Emotion from voice, no neutral</b>							
planum temporale	21	0.148	4.773	0.000	-66	-24	4

636

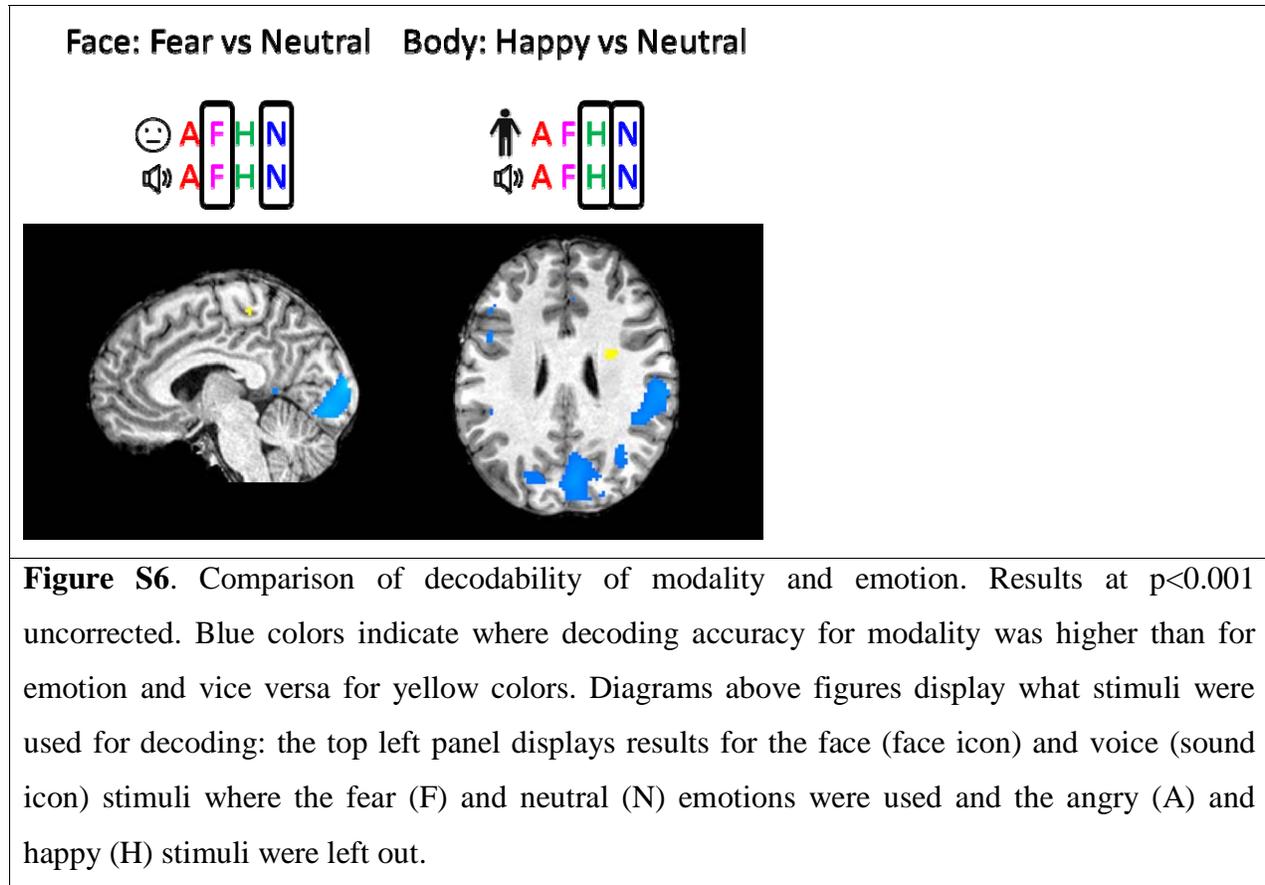
637 **Supplementary Figure 5**



**Figure S5.** Crossmodal decoding of emotion. Red colors are for regions where emotion could be decoded from training on the body stimuli and tested on the voice stimuli from the body session. Green colors are for training on the voice stimuli from the body session and testing on the body stimuli. The same logic applies to stimuli from the face session (see top diagram). Results are at  $t > 3.00$  ( $p < 0.01$ ) uncorrected.

638

639 **Supplementary Figure 6**



640

641