

Accepted Manuscript

Decoding the neural signatures of emotions expressed through sound

Matthew Sachs, Assal Habibi, Antonio Damasio, Jonas Kaplan

PII: S1053-8119(18)30165-4

DOI: [10.1016/j.neuroimage.2018.02.058](https://doi.org/10.1016/j.neuroimage.2018.02.058)

Reference: YNIMG 14759

To appear in: *NeuroImage*

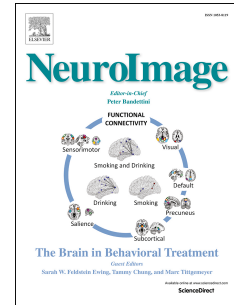
Received Date: 11 October 2017

Revised Date: 23 February 2018

Accepted Date: 27 February 2018

Please cite this article as: Sachs, M., Habibi, A., Damasio, A., Kaplan, J., Decoding the neural signatures of emotions expressed through sound, *NeuroImage* (2018), doi: 10.1016/j.neuroimage.2018.02.058.

This is a PDF file of an unedited manuscript that has been accepted for publication. As a service to our customers we are providing this early version of the manuscript. The manuscript will undergo copyediting, typesetting, and review of the resulting proof before it is published in its final form. Please note that during the production process errors may be discovered which could affect the content, and all legal disclaimers that apply to the journal pertain.



**Decoding the Neural Signatures of Emotions
Expressed Through Sound**

Matthew Sachs¹, Assal Habibi¹, Antonio Damasio¹, & Jonas Kaplan¹

¹Brain and Creativity Institute, University of Southern California
3620A McClintock Avenue
Los Angeles, CA 90089-2921
Telephone: (213) 821-2377

Correspondence should be addressed to Matthew Sachs:

Brain and Creativity Institute,

University of Southern California, 3620A McClintock Avenue

Los Angeles, California, 90089-2921, USA.

E-mail: msachs@usc.edu

1 Abstract

2 Effective social functioning relies in part on the ability to identify emotions from auditory stimuli
3 and respond appropriately. Previous studies have uncovered brain regions engaged by the
4 affective information conveyed by sound. But some of the acoustical properties of sounds that
5 express certain emotions vary remarkably with the instrument used to produce them, for example
6 the human voice or a violin. Do these brain regions respond in the same way to different
7 emotions regardless of the sound source? To address this question, we had participants (N = 38,
8 20 females) listen to brief audio excerpts produced by the violin, clarinet, and human voice, each
9 conveying one of three target emotions—happiness, sadness, and fear—while brain activity was
10 measured with fMRI. We used multivoxel pattern analysis to test whether emotion-specific
11 neural responses to the voice could predict emotion-specific neural responses to musical
12 instruments and vice-versa. A whole-brain searchlight analysis revealed that patterns of activity
13 within the primary and secondary auditory cortex, posterior insula, and parietal operculum were
14 predictive of the affective content of sound both within and across instruments. Furthermore,
15 classification accuracy within the anterior insula was correlated with behavioral measures of
16 empathy. The findings suggest that these brain regions carry emotion-specific patterns that
17 generalize across sounds with different acoustical properties. Also, individuals with greater
18 empathic ability have more distinct neural patterns related to perceiving emotions. These results
19 extend previous knowledge regarding how the human brain extracts emotional meaning from
20 auditory stimuli and enables us to understand and connect with others effectively.

21

22 *Keywords:* fMRI, music, voice, emotions, multivoxel pattern analysis

23

24

25 **1. Introduction**

26 The capacity to both convey and perceive emotions through sounds is crucial for successful
27 social interaction. For example, recognizing that a person is distressed based on vocal
28 expressions alone can confer certain advantages when it comes to communicating and
29 connecting with others. Intriguingly, emotions can be recognized in non-vocal sounds as well.
30 Music can convey emotions even when not mimicking the human voice, despite the fact that an
31 ability to express emotions through music does not serve as clear an evolutionary function as
32 vocal expressions of emotions (Frühholz et al., 2014). And yet, the capability to consistently and
33 reliably discern musical emotions appears to be universal, even in individuals with no musical
34 training (Fritz et al., 2009). Studying the neural overlap of expressions of emotions in both vocal
35 and musical stimuli therefore furthers our understanding of how auditory information becomes
36 emotionally relevant in the human brain.

37 Previous univariate neuroimaging studies that have examined this neural overlap have
38 reported activity in the superior temporal gyrus (Escoffier, Zhong, Schirmer, & Qiu, 2013),
39 amygdala and hippocampus (Frühholz et al., 2014) during both musical and non-musical, vocal
40 expressions of emotions. While these results support the notion that musical and vocal patterns
41 recruit similar brain regions when conveying emotions, they do not clarify whether these regions
42 are responsive to a specific emotional category or are involved in emotion processing more
43 generally. Neither study addressed the neural activity patterns that are specific to a particular
44 emotion, but conserved across the two different domains of music and vocals. One particular
45 univariate study did attempt to answer this question, but only with the emotion of fear: the
46 researchers found that the amygdala and posterior insula were commonly activated in response to
47 fear expressed through non-linguistic vocalizations and musical excerpts, as well as through

48 facial expressions, (Aubé et al., 2013).

49 In general, however, univariate methods are not well suited for evaluating commonalities
50 in the processing of emotions across the senses because, due to spatial smoothing and statistical
51 limitations, they cannot assess information that may be located in fine-grained patterns of
52 activity dispersed throughout the brain (Kaplan et al., 2015). Multivoxel pattern analysis
53 (MVPA), which entails classifying mental states using the spatially-distributed pattern of activity
54 in multiple voxels at once, can provide a more sensitive measure of the brain regions that are
55 responsible for distinguishing amongst different emotions (Norman et al., 2006). In combination
56 with a searchlight analysis, in which classification is performed on local activity patterns within
57 a sphere that traverses the entire brain volume, MVPA can reveal areas of the brain that contain
58 information regarding emotional categories (Kriegeskorte et al., 2006; Peelen et al., 2010). This
59 multivariate approach has been used in various capacities to predict emotional states from brain
60 data (Saarimaki et al., 2015). Spatial patterns within the auditory cortex, for example, were used
61 to classify emotions conveyed through both verbal (Ethofer et al., 2009) and nonverbal (Kotz et
62 al., 2013) speech. However, it remains unclear whether the neural activity in these regions
63 correspond to a particular category of emotion or are instead only sensitive to the lower-level
64 acoustical features of sounds.

65 Multivariate cross-classification, in which a classifier is trained on brain data
66 corresponding to an emotion presented in one domain and tested on separate brain data
67 corresponding to an emotion presented in another, is a useful approach to uncovering
68 representations that are modality independent (see Kaplan, Man, & Greening, 2015 for review).
69 Previously, this approach has been used to demonstrate that emotions induced by films, music,
70 imagery, facial expressions, and bodily actions can be successfully classified across different

71 sensory domains (Peelen et al., 2010; Skerry and Saxe, 2014; Kragel and LaBar, 2015; Saarimaki
72 et al., 2015; Kim et al., 2017). Cross-modal searchlight analyses revealed that successful
73 classification of emotions across the senses and across sources could be achieved based on signal
74 recorded from the cortex lying within the superior temporal sulcus (STS), the posterior insula,
75 the medial prefrontal cortex (MPFC), the precuneus, and the posterior cingulate cortex (Kim et
76 al., 2010; Peelen et al., 2010; Saarimaki et al., 2015). While informative for uncovering regions
77 of the brain responsible for representing emotions across the senses, these studies did not address
78 how the brain represents emotions within a single sensory domain when expressed in different
79 ways. To our knowledge, there has been no existing research on the affect-related neural patterns
80 that are conserved across vocal and musical instruments, two types of auditory stimuli with
81 differing acoustical properties.

82 Additionally, the degree to which emotion-specific predictive information in the brain
83 might be modulated by individual differences remains unexplored. Empathy, for example, which
84 entails understanding and experiencing the emotional states of others, is believed to rely on the
85 ability to internally simulate perceived emotions (Lamm et al., 2007). Activation of the anterior
86 insula appears to be related to linking observed expressions of emotions with internal empathic
87 responses (Carr et al., 2003) and the degree of activation during emotion processing tasks is
88 shown to be positively correlated with measures of empathy (Singer et al., 2004; Silani et al.,
89 2008). Emotion-distinguishing activity patterns in the insula may therefore relate to individual
90 differences in the tendency to share in the affective states of others.

91 Here, we used MVPA and cross-classification on two validated datasets of affective
92 auditory stimuli, one of non-verbal vocalizations (Belin et al., 2008) and one of musical
93 instruments (Paquette et al., 2013), to determine if patterns of brain activity can distinguish

94 discrete emotions when expressed through different sounds. Participants were scanned while
95 listening to brief (0-4s) audio excerpts produced by the violin, clarinet, and human voice and
96 designed to convey one of three target emotions—happiness, sadness, and fear. The authors who
97 published the original dataset chose the violin and clarinet because both musical instruments can
98 readily imitate the sounds of the human voice, but are from two different classes (strings and
99 woodwinds respectively; Paquette et al., 2013). These three target emotions were used because
100 (1) they constitute what are known as “basic” emotions, which are believed to be universal and
101 utilitarian (Ekman, 1992), (2) they can be reliably produced and conveyed on the violin and
102 clarinet (Hailstone et al., 2009) and (3) they are also present in both the vocal and musical
103 datasets.

104 After scanning, a classifier was trained to differentiate the spatial patterns of neural
105 activity corresponding to each emotion both within and across instruments. To understand the
106 contribution of certain acoustic features to our classification results, we compared cross-
107 instrument classification accuracy with fMRI data to cross-instrument classification accuracy
108 using acoustic features of the sounds alone. Then, a searchlight analysis was used to uncover
109 brain areas that represent the affective content that is shared across the two modalities, i.e. music
110 and the human voice. Finally, classification accuracies within a priori-defined regions of interest
111 in the auditory cortex, including the superior temporal gyrus and sulcus, as well as the insula
112 were correlated with behavioral measures of empathy. These regions were selected for further
113 investigation because of their well-validated roles in the processing of emotions from sounds
114 (Bamiou et al., 2003; Sander and Scheich, 2005) as well as across sensory modalities (Peelen et
115 al., 2010; Saarimaki et al., 2015). Based on previous results, we predict that BOLD signal in the
116 auditory and insular cortices will yield successful classification of emotions across all three

117 instruments. Moreover, given the known role of the insula in internal representations of observed
118 emotional states (Carr et al., 2003), we hypothesize that classification accuracies within the
119 insula will be positively correlated with empathy.

120

121 **2. Materials and Methods**

122 *2.1 Participants*

123 Thirty-eight healthy adult participants (20 females, mean age = 20.63, SD = 2.26, range = 18-31)
124 were recruited from the University of Southern California and surrounding Los Angeles
125 community. All participants were right-handed, had normal hearing and normal or corrected-to-
126 normal vision, and had no history of neurological or psychiatric disorders. All experimental
127 procedures were approved by the USC Institutional Review Board. All participants gave
128 informed consent and were monetarily compensated for participating in the study.

129

130 *2.2 Survey*

131 The Goldsmith Musical Sophistication Index (Gold-MSI; Mullensiefen, et al. 2014) was used to
132 evaluate past musical experience and degree of music training. The Gold-MSI contains 39 items
133 broken up into five subscales, each related to a separate component of musical expertise: *active*
134 *engagement, perceptual abilities, musical training, singing abilities, and emotions*. The scale
135 also contains a *general musical sophistication* score, which is the sum of responses to all items.
136 Each item is scored on a 7-point Likert scale from 1 = *completely disagree* to 7 = *completely*
137 *agree*.

138 Both cognitive and affective components of empathy were measured using the
139 Interpersonal Reactivity Index (Davis, 1983), which includes 28 items and four subscales:

140 fantasy and perspective taking (cognitive empathy) and empathic concern and personal distress
141 (affective empathy). Supplementary Table 1 summarizes the results obtained from the surveys.

142

143 *2.3 Stimuli*

144 Two validated, publically-available datasets of short, affective auditory stimuli were used: the
145 Music Emotional Bursts (MEB; Paquette, Peretz, & Belin, 2013) and the Montreal Affective
146 Voices (MAV; Belin, Fillion-Bilodeau, & Gosselin, 2008). Studying neural responses to
147 relatively short stimuli provided two main advantages: (1) As suggested Paquette et al., (2013),
148 these brief bursts of auditory emotions may mimic more primitive, and therefore more
149 biologically relevant, expressions of affect and (2) it allows us to maximize the number of trials
150 that can be presented to participants in the scanner, theoretically improving the training of the
151 classifier. The MEB contains 60 brief (1.64s on average) auditory clips designed to express 3
152 basic emotions (happiness, sadness, and fear) played on either the violin or clarinet. The dataset
153 contains 10 unique exemplars of each emotion on each instrument. The MAV is a set of brief
154 (1.35s on average) non-verbal vocalizations that reliably convey the same 3 emotions (happiness,
155 sadness, and fear) as well as several others (disgust, pain, surprise, and pleasure; these emotional
156 clips were not included in this study because they were not included in the MEB dataset). The
157 MAV dataset contains 10 unique exemplars of each emotion as well and includes both female
158 and male voices. Both the MEB and MAV also include a neutral condition, which were not
159 included in this study either. All clips from both datasets had been normalized so that the peak
160 signal value corresponded to 90% of the maximum amplitude (Belin et al., 2008; Paquette et al.,
161 2013). Combining these two stimulus datasets resulted in 90 unique stimuli: 30 featuring the
162 violin, 30 featuring the clarinet, and 30 featuring the human voice, with 30 clips for each of the

163 three emotions (happiness, sadness, and fear).

164 *2.4 Design and Procedure*

165 Stimuli were presented in an event-related design using MATLAB's PsychToolbox (Kleiner et
166 al., 2007) in 6 functional runs. During scanning, participants were instructed to be still, with their
167 eyes open and focused on a fixation point continually presented on a screen, and attend to the
168 audio clips when they heard them. Auditory stimuli were presented through MR-compatible
169 OPTOACTIVE headphones with noise-cancellation (Optoacoustics). An eye-tracking camera
170 was monitored to ensure that the participants were awake and alert during scanning.

171 During each functional run, participants listened to 45 audio clips, 5 clips for each trial
172 type (emotion x instrument). Each clip was followed by a rest-period that varied in length and
173 resulted in a total event length time (clip + rest) of 5s, regardless of the length of the clip. Five,
174 5-s rest events, in which no sound played, were also added as an additional condition, resulting
175 in a total functional run time of 250s (125 TRs, see Figure 1). Two unique orders of stimuli
176 presentation were created using a genetic algorithm (Kao et al., 2009), which takes into account
177 designed detection power and counterbalancing to generate an optimal design that is
178 pseudorandomized. One optimized order of stimuli presentation was used on odd-numbered runs
179 (1, 3 and 5) and the other order was used on even-numbered runs (2, 4, and 6). Over the course
180 of the 6 functional runs, each of the 90 audio stimuli were presented exactly 3 times.

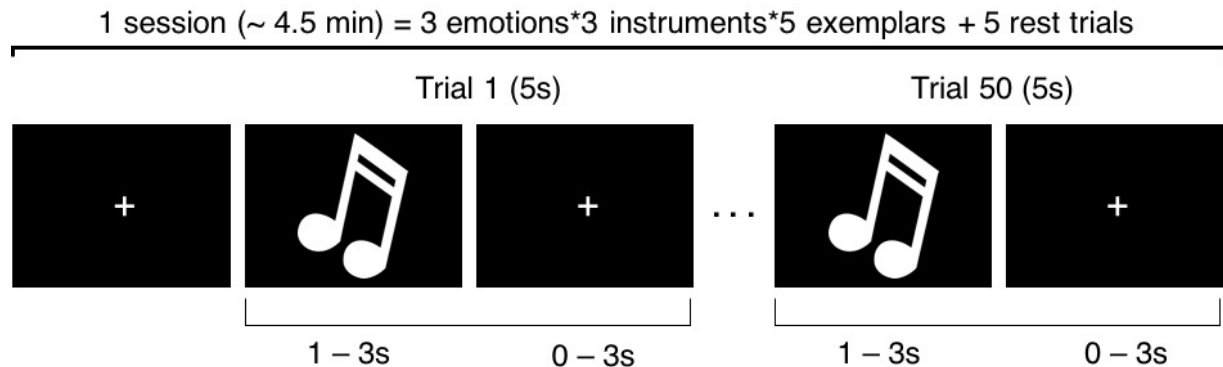
181 To validate the accuracy of the clips in terms of their ability to convey the intended
182 emotion, after scanning, participants listened to all 90 clips again in random order and selected
183 the single emotion, from the list of three, that they believed was being expressed in the clip. To
184 further describe their perceptions of each clip, participants also rated each clip for how intensely
185 it expressed each of the three target emotions using a scale ranging from 1 (not at all) to 5 (very

186 much).

187

188 **Figure 1:** Example of one functional session. In each session, participants listened to 45 clips
 189 and 5 rest trials, each of which lasted for a total of 5s. Each functional session lasted around
 190 4.5min and there were six functional scans in total.

191



192

193

194 2.5 Data Acquisition

195 Images were acquired with a 3-Tesla Siemens MAGNETOM Prisma System and using a 32-
 196 channel head coil. Echo-planar volumes were acquired continuously with the following
 197 parameters: repetition time (TR) = 2,000 ms, echo time (TE) = 25 ms, flip angle = 90°, 64 x 64
 198 matrix, in-plane resolution 3.0 x 3.0 mm, 41 transverse slices, each 3.0 mm thick, covering the
 199 whole brain. Structural T1-weighted magnetization-prepared rapid gradient echo (MPRAGE)
 200 images were acquired with the following parameters: TR = 2,530 ms, TE = 3.09 ms, flip angle =
 201 10°, 256 x 256 matrix, 208 coronal slices, 1 mm isotropic resolution.

202

203 2.6 Data processing

204 Data preprocessing and univariate analysis was done in FSL (FMRIB Software Library, Smith et
 205 al., 2004). Data were first preprocessed using brain extraction, slice-time correction, motion

206 correction, spatial smoothing with 5mm FWHM Gaussian kernel, and high-pass temporal
207 filtering. Each of the 9 trial types (emotion*instrument) was modeled with a separate regressor
208 derived from a convolution of the task design and a double gamma hemodynamic response
209 function. Six motion correction parameters were included in the design as nuisance regressors.
210 The functional data were registered to the high-resolution anatomical image of each subject and
211 to the standard Montreal Neurological Institute (MNI) brain using the FSL FLIRT tool
212 (Jenkinson and Smith 2001). Functional images were aligned to the high-resolution anatomical
213 image using a 7 degree-of-freedom linear transformation. Anatomical images were registered to
214 the MNI-152 brain using a 12 degree-of-freedom affine transformation. This entire procedure
215 resulted in one statistical image for each of the 9 trial types (3 emotions by 3 instrument) in each
216 run. Z-stat images were then aligned to the first functional run of that participant for within-
217 subject analysis.

218

219

220 *2.7 Multivoxel pattern analysis*

221 Multivoxel pattern analysis (MVPA) was conducted using the PyMVPA toolbox
222 (<http://www.pymvpa.org/>) in Python. A linear support vector machine (SVM) classifier was
223 trained to classify the emotion of each trial type. Leave-one run out cross-validation was used to
224 evaluate classification performance (i.e. 6-fold cross-validation with 45 data points in the
225 training dataset and 9 data points in the testing dataset for each fold). Classification was
226 conducted both within each instrument as well as across instruments (training the classifier on a
227 subset of data from two of the instruments and testing on a left-out subset from another
228 instrument) using a mask of the participant's entire brain. In addition to training on two
229 instruments and testing on the third, we also ran cross-instrument classification for every

230 pairwise combination of training on one instrument and testing on another (6 combinations in
231 total). Feature selection on the whole brain mask was employed on the training data alone using
232 a one-way ANOVA and keeping the top 5% most informative voxels (mean 3,320 voxels after
233 feature selection, SD = 251). Within participant classification accuracy was computed by
234 averaging the accuracy of predicting the emotion across each of the 6 folds. One-sample t-tests
235 on the population of participant accuracies were performed to determine if the achieved
236 accuracies were significantly above theoretical chance (33%).

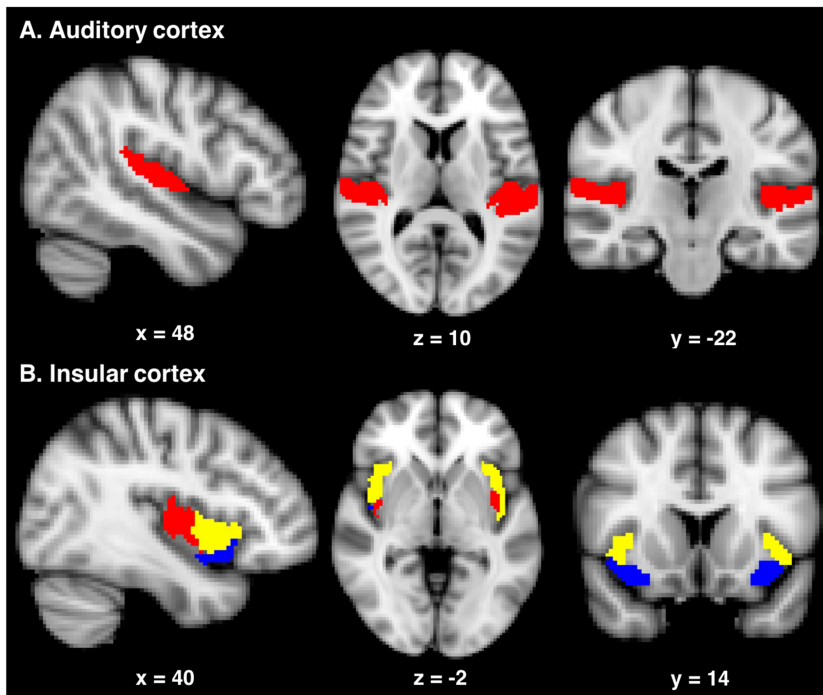
237

238 *2.8 Region of interest classification*

239 In addition to whole brain analysis, we performed a region of interest (ROI) analysis
240 focusing on a-priori ROIs in the auditory cortex and insular cortex. These two ROIs were chosen
241 because of their well-known roles in the processing of emotions from sounds (Bamiou et al.,
242 2003; Sander and Scheich, 2005). For the auditory cortex, we used the Harvard-Oxford Atlas
243 planum temporale mask, which is defined as the superior surface of the superior temporal gyrus,
244 as well as the Heschl's gyrus mask, merged and thresholded at 25 (Figure 2A). For the insula, we
245 used masks of the dorsal anterior, ventral anterior, and posterior insula described in Deen, Pitskel,
246 & Pelphrey (2011) that were defined by the results of a cluster analysis of functional
247 connectivity patterns (Figure 2B). Within and across instrument classification was conducted in
248 exactly the same way as described above. For the region of interest analysis, feature selection
249 was not used, that is, all voxels within the specified anatomical region were used.

250

251 **Figure 2:** Regions of interest. **A**, The auditory cortex was defined using the Harvard-Oxford
 252 Atlas by merging the planum temporale mask with Heschl's gyrus mask, both thresholded at 25.
 253 **B**, Three major subdivisions, the dorsal anterior, ventral anterior, and posterior, were identified
 254 based on the results from a previous study using cluster analysis of functional connectivity
 255 patterns.



256

257 2.9 Whole-brain searchlight analysis

258 A searchlight analysis for classifying emotions was conducted both within and across
 259 modalities (Kriegeskorte et al., 2006). For each subject, the classification accuracy was
 260 determined for spheres with radius 3 voxels throughout the entire brain. A sphere of that size was
 261 chosen to roughly match the size of the anatomical regions of interest, large enough to not be
 262 biased by individual variation in any one voxel and yet small enough to adhere to known
 263 anatomical boundaries. These accuracies were then mapped to the center voxel of the sphere and
 264 warped to standard space. The searchlight analysis was conducted both within instruments and
 265 across instruments. For the within instrument searchlights, the SVM classifier was trained on
 266 data from all but one of the six runs and tested on the left-out run (leave-one-run-out cross

267 validation). To evaluate the significance of clusters in the overlapped searchlight accuracy maps,
268 nonparametric permutation testing was performed using FSL's Randomise tool (Winkler et al.,
269 2014), which models a null distribution of expected accuracies at chance. The searchlight
270 accuracy maps were thresholded using threshold-free cluster enhancement (TFCE; Smith &
271 Nichols, 2009).

272 For the cross-instrument searchlight analysis, the classifier was trained on data from every
273 combination of two instruments and tested on data from the left-out instrument, resulting in a
274 total of three cross-instrument accuracy maps. The three cross-instrument searchlights were also
275 overlaid to determine the regions of overlap. To determine the significance of cross-
276 classification searchlight, we used a more complex nonparametric method than what was used to
277 determine significance of the within-instrument searchlight maps. As described in Stelzer, Chen,
278 & Turner (2013), this method involves random permutation tests on the subject level combined
279 with bootstrapping at the group level. While Randomise with TFCE, which was used to
280 determine significance of the within-instrument searchlight maps, does provide excellent control
281 of type 1 errors, the Stelzer et al., (2013) can provide a more accurate estimation of the group
282 level statistics because it models the null distribution of searchlight maps on both the individual
283 subject level and group level. However, because within-subject permutation testing and across-
284 subject bootstrapping is computationally intensive, we only used this method for determining
285 significance thresholds for the cross-modality searchlight maps, not for the within-instrument
286 searchlights maps. We believe this decision is justified because a) within modality classification
287 in auditory cortex is already well known and does not require a higher standard of proof, b)
288 successful cross-modal classification implies successful within modality classification, and c) the
289 cross-modal searchlights constitute the most direct test of our hypotheses.

290 To achieve this, we randomly permuted the class labels 50 times and performed whole-
291 brain cross searchlight analyses to create 50 single subject chance accuracy maps. One permuted
292 accuracy map per subject was selected at random (with replacement) to create a pooled group
293 accuracy map. This procedure was repeated 10,000 times to create a distribution of pooled group
294 accuracy maps. Next, a threshold accuracy was found for each voxel by determining the
295 accuracy that corresponded to a p-value of 0.001 in the voxel-wise pooled group accuracy map.
296 Clusters were then defined as a group of contiguous voxels that survived these voxel-wise
297 accuracy thresholds and cluster sizes were recorded for each of the 10,000 permuted group
298 accuracy maps to create a histogram of cluster sizes at chance. Finally, cluster-sizes from the
299 chance distribution were compared to cluster-sizes from the original, group accuracy maps to
300 determine significance. An FDR-method using Benjamini-Hochberg procedure was used to
301 correct for multiple comparisons at the cluster level (Heller et al., 2006).

302

303 *2.10 Multiple regression with personality measures*

304 Individual scores on the empathy subscales of the IRI were correlated with classification
305 accuracy within the four ROIs for both within and across classification to determine if the degree
306 of emotion-specific predictive information within these regions is associated with greater
307 emotional empathy. Age, gender, and music sophistication, as measured by the Gold-MSI, were
308 included in the model as regressors of no interest. Additionally, behavioral accuracy of correctly
309 identifying the intended emotions of the sound clips collected outside of the scanner were
310 correlated with performance of the classifier.

311

312 *2.11 Acoustic features of sound clips*

313 For extracting acoustic features from the sound clips believed to be relevant to emotional
314 expression, we used MIRToolbox, a publically available MATLAB toolbox primarily used for
315 music information retrieval (Lartillot et al., 2007), but well suited for extracting relevant
316 acoustical information from non-musical and vocal stimuli as well (Linke and Cusack, 2015;
317 Rigoulot et al., 2015). These included: spectral centroid, spectral brightness, spectral flux,
318 spectral rolloff, spectral entropy, spectral spread, and spectral flatness for evaluating timbral
319 characteristics, RMS energy for evaluating dynamics, mode and key clarity for evaluating tonal
320 characteristics, and fluctuation entropy and fluctuation centroid for evaluating rhythmic
321 characteristics of the clips (Alluri et al., 2012). We additionally added the acoustic features
322 published in Paquette et al., (2013), which included duration, mean fundamental frequency, max
323 fundamental frequency, and min fundamental frequency. We then evaluated how these features
324 varied by instrument and by emotion and conducted a classification analysis based on acoustic
325 features alone to predict the intended emotion of the sound clip.

326 Because we found a main effect of emotion label on duration of the clips (i.e. fear clips
327 were significantly shorter than sad clips) and we do not believe that this difference reflects a
328 meaningful difference amongst emotions, we added an additional regressor of no interest where
329 the height of the regressor reflected the duration of each clip in a separate GLM analysis. MVPA
330 and searchlight analysis were then repeated with this model for comparison.

331

332 **3. Results**

333 *3.1 Behavioral results*

334 Behavioral ratings of the sound clips outside of the scanner were collected for 37 out of the
335 38 participants. Overall, participants correctly labeled 85% of the clips (SD = 17%). Averaged

336 correct responses for each emotion and instrument are presented in Table 1. The between-within
337 ANOVA on accuracy scores for each of the clips showed a significant interaction between
338 emotion and instrument ($F(4,144) = 64.30, p < 0.0001$). Post-hoc follow-up t-tests showed that
339 the fear condition on the clarinet was more consistently labeled incorrectly (mean accuracy =
340 55%, $SD = 16\%$) than the happy condition on the clarinet (mean accuracy = 97%, $SD = 5\%$; $t(36)$
341 = 15.99, $p < 0.0001$, paired t-test) as well as the fear condition in the violin (mean accuracy =
342 90%, $SD = 13\%$; $t(36) = 13.83, p < 0.0001$, paired t-test) or voice (mean accuracy = 94%, $SD =$
343 7%; $t(36) = 13.27, p < 0.0001$, paired t-test).

344 For intensity ratings of the clips, we calculated the average intensity of each emotion for
345 each participant. Again, an interaction between emotion and instrument was found for the
346 intensity ratings ($F(4,144) = 38.73, p < 0.0001$, ANOVA). Fear clips on the clarinet were rated as
347 significantly less intense than fear clips on the violin ($t(36) = 14.49, p < 0.0001$, paired t-test) and
348 voice ($t(36) = 13.16, p < 0.0001$, paired t-test), whereas sad clips on the voice were rated as
349 significantly more intense than sad clips on the violin ($t(36) = 9.84, p < 0.0001$, paired t-test) or
350 clarinet ($t(36) = 9.90, p < 0.0001$, paired t-test; see Table 1 for average ratings of intensity and
351 average accuracy). Overall, the intensity ratings provide further information for the accuracy
352 scores: fear on the clarinet was the most difficult to identify and was rated as significantly less
353 intense. Participant intensity ratings were not found to be related to the performance of the brain-
354 based classifier.

355
356
357
358
359
360
361
362

363 **Table 1.** Behavioral ratings of intensity and emotion label for stimuli by instrument and emotion.
 364 Accuracy is calculated as the number of clips correctly identified with the intended emotion label.

	Happy		Sad		Fear		Total	
	Acc	Intensity	Acc	Intensity	Acc	Intensity	Acc	Intensity
All	0.92	3.86 (0.58)	0.85	3.81 (0.68)	0.79	3.52 (0.87)	0.85	3.86 (0.58)
Voice	0.82	4.17 (0.45)	0.81	4.35 (0.48)	0.90	4.04 (0.56)	0.84	4.19 (0.51)
Violin	0.96	3.54 (0.61)	0.89	3.54 (0.56)	0.94	3.87 (0.69)	0.93	3.64 (0.64)
Clarinet	0.97	3.85 (0.50)	0.85	3.55 (0.64)	0.55	2.67 (0.61)	0.79	3.36 (0.77)

365

366 3.2 Multivariate results

367 MVPA applied to the whole brain to predict the emotion of each clip showed above chance
 368 (0.33) accuracies using data from all instruments ($M = 0.43$, $SD = 0.08$, $t(37) = 7.58$, $p < 0.0001$).
 369 Above chance accuracy was also obtained using data collected from each instrument individually
 370 (*clarinet*: $M = 0.39$, $SD = 0.12$, $t(37) = 3.06$, $p = 0.004$; *violin*: $M = 0.37$, $SD = 0.11$, $t(37) = 2.18$,
 371 $p = 0.04$; *voice*: $M = 0.43$, $SD = 0.15$, $t(37) = 4.14$, $p = 0.0002$). Within instrument classification
 372 accuracy was also significantly above chance in both the auditory cortex ($M = 0.49$, $SD = 0.06$,
 373 $t(37) = 15.06$, $p < 0.0001$) and all three regions of the insula (*dorsal anterior*: $M = 0.38$, $SD =$
 374 0.07 , $t(37) = 4.06$, $p = 0.0002$; *ventral anterior*: $M = 0.38$, $SD = 0.07$, $t(37) = 4.41$, $p < 0.0001$;
 375 *posterior*: $M = 0.38$, $SD = 0.08$, $t(37) = 3.68$, $p = 0.0007$, see Figure 3). Confusion matrices for
 376 within instrument classification in the whole brain and auditory cortex are provided in the
 377 supplementary materials (Supplementary Figure 1) as well as additional measures of
 378 classification performance, including sensitivity, specificity, positive predictive value, and
 379 negative predictive value (Supplementary Table 2).

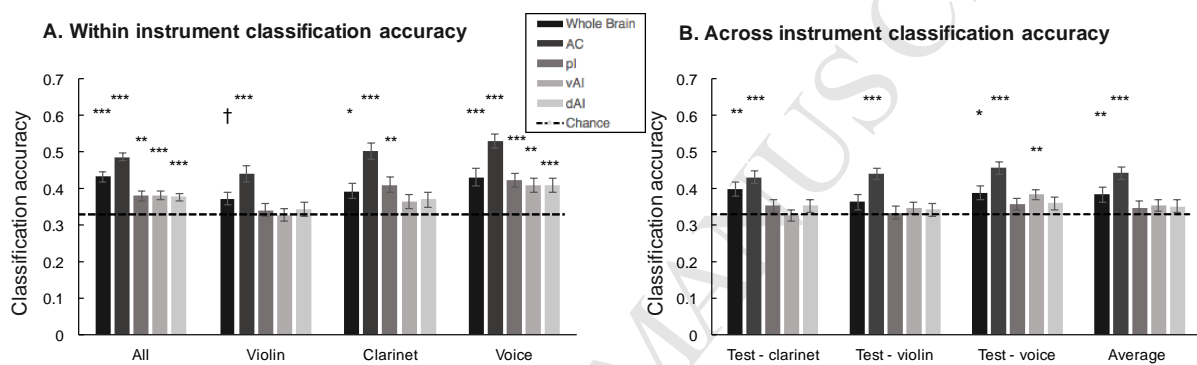
380 Cross-classification accuracies, in which the classifier was trained on data from two

381 instruments and tested on data from the left-out third instrument, also showed successful
382 classification for each combination of training and testing (3 in total). Classification accuracy
383 averaged across the 3 combinations of training and testing was significantly greater than chance
384 in the whole brain ($M = 0.38$, $SD = 0.09$, $t(37) = 3.56$, $p = 0.001$) as well as the region of interest
385 in the auditory cortex ($M = 0.44$, $SD = 0.08$, $t(37) = 8.83$, $p < 0.0001$), but not in the three insula
386 ROIs. Graphs of the accuracies for each combination of training and testing in both the whole
387 brain analysis and ROI analysis are presented in Figure 3. Confusion matrices for cross
388 instrument classification in the whole brain and auditory cortex are provided in Supplementary
389 Figure 2.

390 We additionally conducted cross-classification for all 6, pairwise combinations of training
391 on one instrument and testing on one other instrument. Overall, the pairwise cross-instrument
392 classification accuracies were significantly above chance in the auditory cortex. The results are
393 presented in Supplementary Figure 3.

394
395
396
397
398
399
400
401
402
403
404
405
406
407
408
409
410
411
412

413 **Figure 3.** Classification accuracies for MVPA decoding of emotions in auditory stimuli using
 414 fMRI data from the whole brain and four region of interest. **A**, Classification accuracies in the
 415 whole brain, auditory cortex (AC), posterior insula (pI), dorsal anterior insula (dAI), and ventral
 416 anterior insula (vAI) with all three instruments (violin, clarinet, voice) as well as within each
 417 instrument individually. **B**, Cross-instrument classification accuracies in the whole brain,
 418 auditory cortex (AC), posterior insula (pI), dorsal anterior insula (dAI), and ventral anterior
 419 insula (vAI), leaving out data from one instrument and training on the other two. Error bars
 420 represent indicate error. p values are calculated based on a one-sample t-test comparing
 421 classification with chance (0.33,dotted line). †p < 0.05, uncorrected; *p < 0.05;**p < 0.01,***p
 422 < 0.001, corrected for multiple comparisons across the four ROIs.
 423
 424
 425



426
 427
 428

429 3.3 Searchlight results

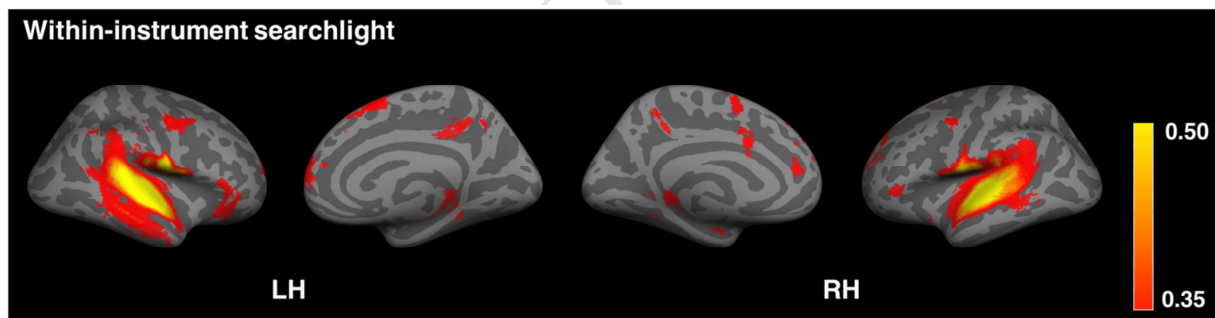
430 The whole-brain, within instrument searchlight analysis revealed that successful
 431 classification of the emotions of the musical clips could be found bilaterally in the primary and
 432 secondary auditory cortices, including the cortices lying within the superior temporal gyrus and
 433 sulcus, as well as the bilateral posterior insular cortices, parietal operculum, precentral gyrus,
 434 inferior frontal gyrus, right middle temporal gyrus, the right medial prefrontal cortex, right
 435 superior frontal gyrus, right precuneus, and right supramarginal gyrus (Figure 4). Center
 436 coordinates and accuracies for significant regions in the within instrument searchlight analysis
 437 are presented in Supplementary Table 2.

438 Three whole-brain across instrument searchlight analyses were conducted where the

439 classifier was trained on data from two of the instruments and tested on the held-out third
440 instrument. All three searchlights showed significant classification bilaterally in primary auditory
441 cortex, including Heschl's gyrus, and the superior temporal gyrus and sulcus, as well as the
442 posterior insula and parietal operculum (Figure 5). Several other brain regions showed
443 significant classification in one or more of the searchlight analyses, but not all three. These
444 included the right middle and inferior frontal gyri and precentral gyrus (leaving out the violin
445 and the voice) and the MPFC (leaving out clarinet). Center coordinates and accuracies for each
446 significant region of the cross-instrument searchlight analysis are presented in Supplementary
447 Table 3.

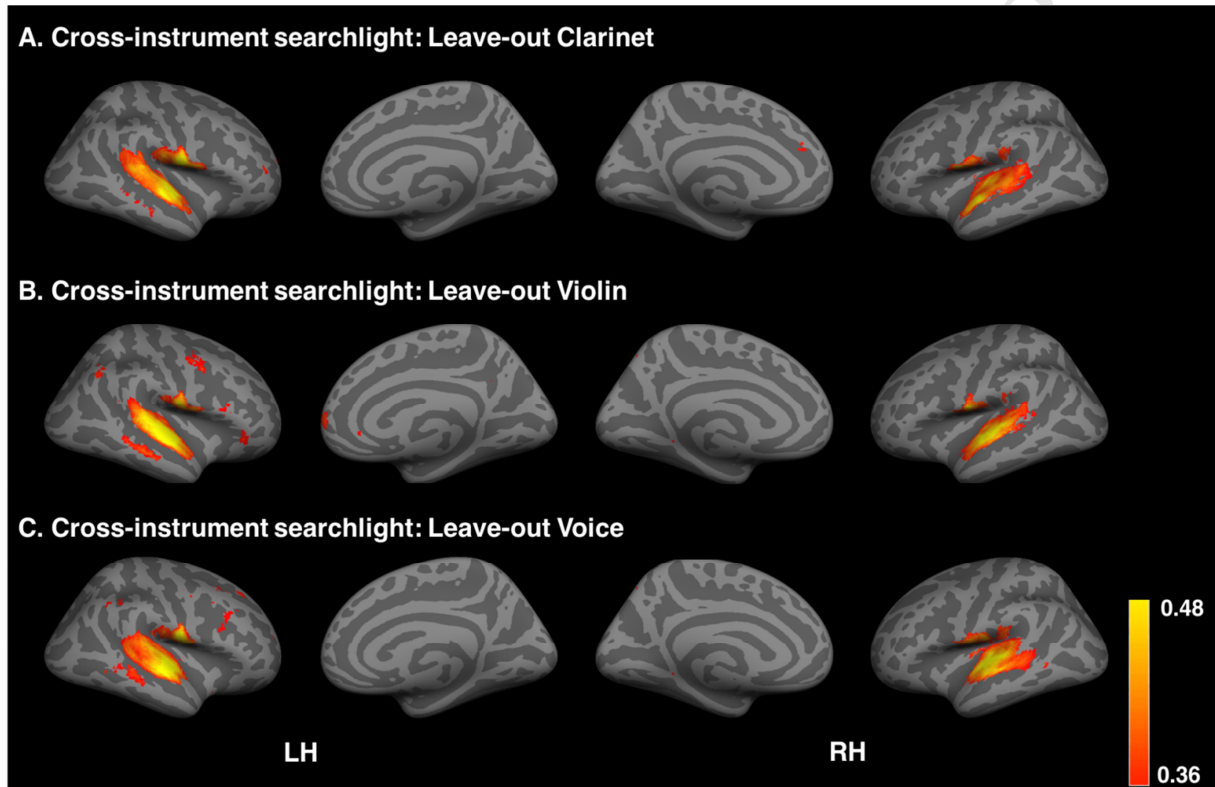
448
449
450

451 **Figure 4:** Within instrument whole-brain searchlight results using data from all instruments and
452 leave-one-run out cross validation. Red-yellow colors represent classification accuracy.
453 Significant clusters determined by permutation testing. All images are thresholded to show
454 clusters that reached a FDR-corrected significance level at $\alpha = 0.05$



455
456
457
458
459
460
461
462
463
464
465
466

467 **Figure 5:** **A**, Cross-instrument whole-brain searchlight results training on data collected during
 468 violin and voice clips, testing on clarinet clips. **B**, Cross-instrument whole-brain searchlight
 469 results training on data collected during clarinet and voice clips, testing on violin clips. **C**,
 470 Cross-instrument whole-brain searchlight results training on data collected during violin and
 471 clarinet clips, testing on voice clips. Red-yellow colors represent classification accuracy.
 472 Significant clusters were determined by permutation testing. All images are thresholded to show
 473 clusters that reached a FDR-corrected significance level at $\alpha = 0.05$
 474



475

476

477 3.4 Multiple regression results

478 Measures of the four subscales of the IRI were modeled in a multiple regression to predict
 479 the classification accuracies in each of the four regions of interest (auditory cortex and three
 480 subcomponents of the insula) with age and gender added as covariates of no interest. In this
 481 model, empathic concern was positively correlated with both the within and cross classification
 482 accuracies in the dorsal anterior insula (*Within*: $\beta = 0.08$, $p = 0.0101$, *Cross*: $\beta = 0.08$, $p =$
 483 0.0158). The significance of the regression coefficient between empathic concern and within-

484 instrument accuracy in the dorsal anterior insula survived correction for multiple comparisons,
485 though the regression with cross-instrument accuracies did not (Bonferroni correction with four
486 regions of interest, $\alpha = 0.0125$). No other predictors were significantly correlated with
487 accuracy in the regions of interest.

488 Scores corresponding to the five subscales of the MSI were additionally modelled in a
489 separate multiple regression as covariates of interest to predict classification accuracies in the
490 four ROIs. No significant correlations were found between musical experience and classification
491 accuracy in either the auditory cortex or insula. Additionally, no significant correlations were
492 found between behavioral accuracies of correctly identifying the intended emotion of the clip
493 (collected outside of the scanner) and classification accuracy.

494

495 *3.6 Acoustic features classification and duration*

496 Duration was significantly different for the three emotions according to a one-way
497 ANOVA ($F(2,87) = 110.30$, $p < 0.0001$). Sad clips ($M = 2.39s$, $SD = 0.54$) were significantly
498 longer than the happy ($M = 1.48s$, $SD = 0.39$; $t(58) = 7.45$, $p_{\text{adjust}} = 1.2 \times 10^{-12}$) and fear clips (M
499 $= 0.81s$, $SD = 0.25$; $t(58) = 14.38$, $p_{\text{adjust}} < 0.0001$). Because duration is not an acoustic feature
500 directly related to the expression of an emotion, and because it differed significantly by emotion,
501 we wanted to ensure that the classifier was not only classifying based on stimuli length rather
502 than its emotional content. We therefore added the length of each clip as an additional parametric
503 regressor in the lower-level GLM models and redid both within and cross instrument
504 classification with the z-stat images obtained from this analysis. The average within instrument
505 accuracy using duration as a regressor was 44% ($SD = 7\%$) within the whole brain and the
506 average across instrument accuracy using duration as a regressor was 39% ($SD = 9\%$), which

507 were both statistically significant according to a one-way t-test against theoretical chance (*within*:
508 $t(37) = 9.12, p = 5.33 \times 10^{-11}$; *across*: $t(37) = 4.22, p = 0.0002$). No significant differences were
509 found between the classification accuracies when duration was added as a regressor (*within*: $t(37)$
510 $= 0.61, p = 0.55$, *across*: $t(37) = 0.53, p = 0.60$, paired t-test). Because of these, we did not re-
511 compute the searchlight analysis using the results from the GLM analysis with duration modelled.

512 We also conducted a classification analysis using the acoustic features of the sound clips
513 only. These included 12 features related to timbre, rhythm, and tonality as described in (Alluri et
514 al., 2012) as well as fundamental frequency and duration (Paquette et al. 2013). The linear SVM
515 classifier could successfully classify the emotion label of the sound clip 82% of the time using all
516 data and 60% on average when training and testing on data from separate instruments (cross
517 classification). The duration of the sound clips was determined to be the most important feature
518 used by the SVM. When duration was removed, classification accuracy was 72% when using
519 data from all instruments and 57% on average when training and testing across all three
520 instruments (see Table 2). After removing duration, the most important features for classification
521 were fluctuation centroid and spectral flux. Fluctuation centroid is a measure of rhythmic
522 changes in sounds and is calculated by taking the mean (center of gravity) of the fluctuation
523 spectrum, which conveys the periodicities contained in a sound wave's envelope (Alluri et al.,
524 2012). A one-way ANOVA revealed a significant main effect of emotion $F(2,87) = 29.93, p <$
525 0.0001 on fluctuation centroid. Fear clips ($M = 2,977; SD = 1,220$) were significantly higher
526 than both sad ($M = 1,325; SD = 554; t(58) = 6.76, p_{\text{adjust}} < 0.0001$) and happy clips ($M = 2,432;$
527 $SD = 632.21; t(58) = , p_{\text{adjust}} < 0.0001$). Spectral flux is a measure of how the variance in the
528 audio spectrum changes over time and therefore conveys both spatial and temporal components
529 of sound (Alluri et al., 2012). It is highly correlated with fluctuation centroid with our sound

530 clips. A one-way ANOVA revealed that fearful clips ($M = 175.03$, $SD = 86.05$) had a
 531 significantly higher spectral flux than both sad ($M = 59.70$, $SD = 33.32$; $t(58) = 6.85$, $p < 0.001$)
 532 and happy clips ($M = 93.20$, $SD = 28.25$; $t(58) = 4.95$, $p < 0.001$).

533 The results from the acoustic classification provide information regarding how the fMRI-
 534 based classifier is able to decode the auditory emotions and suggests that differences in neural
 535 responses to changes in rhythm and timbre between the emotions might contribute to the
 536 classifier's performance.

537

538 **Table 2:** Classification of emotion of stimuli using acoustic features (duration and fundamental
 539 frequency)

		Happy	Sad	Fear	Total
Within classification	w/ duration	0.73	0.83	0.90	0.82
	w/out duration	0.53	0.87	0.77	0.72
Cross-classification: Test on voice	w/ duration	0.50	0.20	0.90	0.57
	w/out duration	0.50	0.40	0.70	0.53
Cross-classification: Test on violin	w/ duration	0.50	0.70	0.30	0.50
	w/out duration	0.50	0.30	0.30	0.37
Cross-classification: Test on clarinet	w/ duration	0.60	1.00	0.60	0.73
	w/out duration	0.80	1.00	0.60	0.80

540

541 4. Discussion

542 By using multivariate cross-classification and searchlight analyses with different types of
 543 auditory stimuli that convey the same three emotions, we identified higher-level neural regions
 544 that process the affective information of sounds produced from various sources. Using fMRI data
 545 collected from the entire brain, above-chance classification of emotions expressed through
 546 auditory stimuli was found both within and across instruments. Searchlight analyses revealed

547 that the primary and secondary auditory cortices, including the superior temporal gyrus (STG)
548 and sulcus (STS), extending into the parietal operculum and posterior insula, exhibit emotion-
549 specific and modality-general patterns of neural activity. This is supported by the fact that BOLD
550 signal in these regions could differentiate the affective content when the classifier was trained on
551 data from one instrument and tested on data from another instrument. Furthermore, within and
552 cross-modal classification performance within a region spatially confined to the dorsal anterior
553 portion of the insula was positively correlated with a behavior measure of empathy. To our
554 knowledge, this is the first study to report the emotion-related spatial patterns that are shared
555 across both musical instruments and vocal sounds as well as to link the degree of predictive
556 information within these spatial patterns with individual differences.

557 The findings confirm the role of the cortices in the STG and the STS regions in perceiving
558 emotions conveyed by auditory stimuli. Significant classification of vocal expressions of
559 emotions was previously reported in the STG (Kotz et al., 2013) and the region is active when
560 processing acoustical (Salimpoor et al., 2015) and affective components of music (Koelsch,
561 2014). The left STG was also found to code for both lower-level acoustic aspects as well as
562 higher-level evaluative judgments of nonlinguistic vocalizations of emotions (Bestelmeyer et al.,
563 2014). It has been suggested that the STG and STS bilaterally may be involved in tracking the
564 changing acoustic features of sounds as they evolve over time (Schonwiesner et al., 2005). The
565 STS in particular appears to integrate audio and visual information during the processing of non-
566 verbal affective stimuli (Kreifelts et al., 2009). Both facial and vocal expressions of emotions
567 activate the STS (Escoffier et al., 2013; Wegrzyn et al., 2015). Multivariate neuroimaging studies
568 have proposed that supramodal mental representations of emotions lie in the STS (Peelen et al.,
569 2010). Furthermore, aberrations in both white-matter volume (von dem Hagen et al., 2011) and

570 task-based functional activity (Alaerts et al., 2014) in these regions were associated with emotion
571 recognition deficits in individuals with autism spectrum disorder (ASD). Our findings suggest
572 that discrete emotions expressed through music are represented by similar patterns of activity in
573 the auditory cortex as when expressed through the human voice. This confirms the role of the
574 STG and STS in processing the perceived affective content of a range of sounds, both musical
575 and non-musical, that is not purely dependent on lower-level acoustic features.

576 While the peak of the searchlight accuracy maps was located in the auditory cortex, the
577 significant results extend into the parietal operculum. Because these two regions are adjacent and
578 because the neuroimaging data are spatially smoothed both in the preprocessing steps and in the
579 searchlight analysis, we cannot be certain that the significant classification accuracy found in the
580 parietal operculum indicates that this region is additionally involved in representing emotions
581 from sounds. Nonetheless, the idea that cross-modal representation of emotions could be located
582 in this region is consistent with previous research. The inferior portion of the somatosensory
583 cortex, which is located in the parietal operculum (Eickhoff et al., 2006), has been shown to be
584 engaged during vicarious experiences of perceived emotions (Straube and Miltner, 2011).
585 Furthermore, patients with lesions in the right primary and secondary somatosensory cortices
586 performed poorly in emotion recognition tasks (Adolphs et al., 2000) and reported reduced
587 intensity of subjective feelings in response to music (Johnsen et al., 2009). Transcranial magnetic
588 stimulation applied over the right parietal operculum region was also shown to impede the ability
589 to detect the emotions of spoken language (Rijn et al., 2005). Using multivariate methods, Man,
590 Damasio, Meyer, & Kaplan, (2015) found activity in the parietal operculum could be used to
591 reliably classify objects when presented aurally, visually, and tactilely, suggesting that this
592 region contains modality invariant representations of objects and may therefore serve as a

593 convergence zone for information coming from multiple senses. Taken together, the fact that we
594 find significant predictive affective information in the parietal operculum may suggest that the
595 ability to recognize the emotional content of sounds relies on an internal simulation of the
596 actions and sensations that go into producing such sounds.

597 The searchlight accuracy maps additionally extended into the posterior portion of the insula.
598 The insula is believed to be involved in mapping bodily state changes associated with particular
599 feeling states (Damasio et al., 2013; Immordino-Yang et al., 2014). A range of subjectively-
600 labeled feeling states could be decoded from brain activity in the insula, suggesting that the
601 physiological experience that distinguishes one emotion from another is linked to distinct spatial
602 patterns of activity in the insula (Saarimaki et al., 2015). Studies have shown that the region is
603 largely modality invariant, activated in response to facial expressions of emotions (Wegrzyn et
604 al., 2015), perceptual differences between emotions conveyed through non-speech vocalizations
605 (Bestelmeyer et al., 2014), multimodal presentations of emotions (Schirmer and Adolphs, 2017)
606 and by a wide range of emotions conveyed through music (Baumgartner et al., 2006; Park et al.,
607 2013). The insula's role in auditory processing may be to allocate attentional resources to salient
608 sounds (Bamiou et al., 2003) as evidenced by cases in which patients with lesions that include
609 the insula but not Heschl's gyrus develop auditory agnosia (Fifer et al 1995). The function of the
610 insula in processing emotions expressed across the senses is further substantiated by the
611 observation that a patient with a lesion in the insula showed an impaired ability to recognize the
612 emotion disgust when expressed in multiple modalities (Calder et al., 2000). Developmental
613 disorders characterized by deficits in emotional awareness and experience may be linked to
614 aberrant functioning of the insula, as decrease insular activity was observed in ASD children
615 observing emotional faces (Dapretto et al., 2006) and altered resting-state functional connectivity

616 between the posterior insula and somatosensory cortices was observed in adults with ASD
617 (Ebisch et al., 2011). The fact that emotions conveyed through auditory stimuli could be
618 classified based on activity in the posterior insula in our study provides further evidence for the
619 hypothesis that perceiving and recognizing an emotion entails recruiting the neural mechanisms
620 that represent the subjective experience of that same emotion.

621 Despite this finding, classification accuracy within a region of interest in the *dorsal*
622 *anterior* portion of the insula was significantly positively correlated with empathy. The anterior
623 insula was not one of the significant regions found in the searchlight analysis. These two results
624 might be explained in the context of previous functional and structural imaging studies that
625 suggest that subdivisions of the insular cortex are associated with specific functions (Deen et al.,
626 2011). According to such accounts, the posterior insula, which is structurally connected to the
627 somatosensory cortices, is more directly involved in affective processing of visceral sensations
628 (Kurth et al., 2010) and interoceptive awareness (Craig, 2009), whereas the dorsal anterior insula,
629 which is connected to the cognitive control network and the ACC, is more directly involved in
630 socio-emotional abilities such as empathy and subjective awareness (Craig, 2009). This is
631 evidenced by the fact that the anterior insula is activated when both observing and imitating the
632 emotions of others (Carr et al., 2003). Measures of empathic concern, a subtype of affective
633 empathy referring to the tendency to feel sympathy or concern for others, have been shown to be
634 positively correlated with anterior insula activity when viewing emotional pictures (Silani et al.,
635 2008) as well as when observing loved-ones in pain (Singer et al., 2004). Given that
636 classification accuracy obtained from data within the dorsal anterior insula specifically, was
637 correlated with empathic concern, our results provide further evidence for the unique role of this
638 subdivision in enabling the emotional resonance that is essential to understanding the feelings of

639 others. We speculate that individuals who readily behave empathically might have more finely
640 tuned representations of emotions in the dorsal anterior insula when processing affective
641 information.

642 While we were mainly interested in identifying brain regions that conserve affect-related
643 information across sounds with differing acoustical properties, we recognize that certain acoustic
644 properties are also integral to specific emotional categories regardless of the source of the sound.
645 Previous results have shown that happiness, for example, is characterized by higher fundamental
646 frequencies and faster tempos when conveyed through both vocal expressions and musical pieces
647 (Juslin and Laukka, 2003). An earlier attempt to disentangle the neural processing of acoustic
648 changes from the neural processing of perceptual changes associated with two different emotions
649 conveyed through auditory stimuli acknowledged that the two are interrelated and that a
650 complete and straightforward separation of the two would be overly simplistic; indeed, the
651 researchers found evidence for both distinct and overlapping neural networks associated with
652 these two processes (Bestelmeyer et al., 2014). Because of this, we did not intend to control for
653 all potential acoustic differences between our stimuli, believing that these features may be
654 essential to that emotional category. Because the duration of the clips varied significantly by
655 emotion and is a feature not directly tied to affective expression, we did regress out the variance
656 explained by duration from our GLM model and showed that cross-instrument classification
657 performance did not change. Besides for duration, no other acoustic properties of the sounds
658 were regressed out of the signal. We therefore might expect that our classifier may be sensitive
659 to signal that is responsive to certain acoustic variations.

660 To make predictions about the types of acoustic variation that the classifier may be
661 sensitive to, we conducted classification using several audio features extracted for each clip. We

662 found that the acoustic-based classifier performance was also largely dependent on differences in
663 duration between the emotions, as evidenced by the attenuating in performance when durational
664 information was removed. While it is difficult to know what types of information the fMRI-
665 based classifier is using to make distinctions between the emotional states, the classifier trained
666 on acoustic features alone without duration can provide some hypotheses. Once duration was
667 removed, the most informative features for classification of emotions include a rhythmic feature
668 called fluctuation centroid. These results suggest that the fMRI-based classifier is not only
669 sensitive to BOLD signal corresponding to the duration and frequency of the sounds, and may be
670 capturing finely-tuned responses in the auditory cortex and insula that are sensitive to changes in
671 rhythm and timbre that are integral to conveying emotions through sound.

672 Using BOLD data from a mask of the entire brain, classification accuracies were between
673 38-43% for the whole-brain within instrument classification and 36-40% for the whole-brain
674 cross-instrument classification. Theoretically, if the classifier was guessing the emotional
675 category at random, just by chance, it would correctly identify the emotion 33% of the time.
676 While we recognize that the accuracies obtained here are not impressively high compared to
677 theoretical chance, they are statistically significant according to a one-way t-test corrected for
678 multiple comparisons and comparable to those reported in other multivariate cross-classification
679 fMRI studies (Skerry and Saxe, 2014; Kim et al., 2017). Furthermore, the cross instrument
680 classification accuracies should be interpreted in relation to the within instrument classification
681 accuracies, which may set the upper bound of possible performance of a cross-modal classifier
682 (Kaplan et al., 2015). We therefore would not expect cross instrument classification to perform
683 better than within instrument classification and the fact that the cross instrument accuracies are
684 still significantly above chance provides us with compelling evidence that unique spatial patterns

685 of BOLD signal throughout does contain some predictive information regarding the emotional
686 category.

687 Of additional note, upon inspection of the confusion matrices, the within-instrument
688 classification performance correctly identified fearful clips to a greater degree than the other two
689 emotions, despite the fact positive predictive value across all three was not drastically different.
690 This contrasts with the behavioral findings, in which fear was the emotion most difficult to
691 identify and label. This could indicate that while the fearful clips were easily distinguishable
692 from the other two emotions, these perceived differences may not necessarily adhere to the
693 concepts and features humans have learned to associate with the categorical label of fear. Further
694 exploration into the acoustic components and behavioral responses to fearful musical and vocal
695 clips will help to interpret these opposing findings.

696 In sum, our study reveals the emotion-specific neural information that is shared across
697 sounds from musical instruments and the human voice. The results support the idea that the
698 emotional meaning of sounds can be represented by unique spatial patterns of neural activity in
699 sensory and affect processing areas of the brain, representations that do not depend solely on the
700 specific acoustic properties associated with the source instrument. These findings therefore have
701 implications for scientific investigations of neurodevelopmental disorders characterized by an
702 impaired ability to recognize vocal expressions of emotions (Allen et al., 2013) and provide a
703 clearer picture of the remarkable ability of the human brain to instantaneously and reliably infer
704 emotions when conveyed nonverbally.

705

706

707

708 **Acknowledgements**

709 The authors would like to thank Hanna Damasio for her assistance and input regarding
710 neuroanatomical distinctions and all private donors to the Brain and Creativity Institute.

711

712 **Funding**

713 Funding for this work was provided by the Brain and Creativity Institute.

714

715

716

717

718

719

720

721

722

723

724

725

726

727

728

729

730

731 **References**

- 732 Adolphs R, Damasio H, Tranel D, Cooper G, Damasio AR (2000) A role for somatosensory
733 cortices in the visual recognition of emotion as revealed by three-dimensional lesion
734 mapping. *J Neurosci* 20:2683–2690.
- 735 Allen R, Davis R, Hill E (2013) The Effects of Autism and Alexithymia on Physiological and
736 Verbal Responsiveness to Music. *J Autism Dev Disord* 43:432–444.
- 737 Alluri V, Toiviainen P, Jääskeläinen IP, Glerean E, Sams M, Brattico E (2012) Large-scale brain
738 networks emerge from dynamic processing of musical timbre, key and rhythm. *Neuroimage*
739 59:3677–3689.
- 740 Aubé W, Angulo-Perkins A, Peretz I, Concha L, Armony JL (2013) Fear across the senses: Brain
741 responses to music, vocalizations and facial expressions. *Soc Cogn Affect Neurosci*
742 10:399–407.
- 743 Bamiou D, Musiek FE, Luxon LM (2003) The insula (Island of Reil) and its role in auditory
744 processing: Literature review. *Brain Res Rev* 42:143–154.
- 745 Baumgartner T, Lutz K, Schmidt CF, Jäncke L (2006) The emotional power of music: How
746 music enhances the feeling of affective pictures. *Brain Res* 1075:151–164.
- 747 Belin P, Fillion-Bilodeau S, Gosselin F (2008) The Montreal Affective Voices: a validated set of
748 nonverbal affect bursts for research on auditory affective processing. *Behav Res Methods*
749 40:531–539.
- 750 Bestelmeyer PEG, Maurage P, Rouger J, Latinus M, Belin P (2014) Adaptation to Vocal
751 Expressions Reveals Multistep Perception of Auditory Emotion. *J Neurosci* 34:8098–8105.
- 752 Calder AJ, Keane J, Manes F, Antoun N, Young AW (2000) Impaired recognition and
753 experience of disgust following brain injury. *Nat Neurosci* 3:1077–1078.

- 754 Carr L, Iacoboni M, Dubeau M, Mazziotta JC, Lenzi GL (2003) Neural mechanisms of empathy
755 in humans : A relay from neural systems for imitation to limbic areas. *Proc Natl Acad Sci*
756 100:5497–5502.
- 757 Craig AD (2009) How do you feel--now? The anterior insula and human awareness. *Nat Rev*
758 *Neurosci* 10:59–70.
- 759 Damasio A, Damasio H, Tranel D (2013) Persistence of feelings and sentience after bilateral
760 damage of the insula. *Cereb Cortex* 23:833–846.
- 761 Dapretto M, Davies MS, Pfeifer JH, Scott AA, Sigman M, Bookheimer SY, Iacoboni M (2006)
762 Understanding emotions in others: mirror neuron dysfunction in children with autism
763 spectrum disorders. *Nat Neurosci* 9:28–30.
- 764 Davis MH (1983) Measuring individual differences in empathy: Evidence for a multidimensional
765 approach. *J Pers Soc Psychol* 44:113–126.
- 766 Deen B, Pitskel NB, Pelphrey KA (2011) Three Systems of Insular Functional Connectivity
767 Identified with Cluster Analysis. *Cereb Cortex* 21:1498–1506.
- 768 Ebisch SJH, Gallese V, Willems RM, Mantini D, Groen WB, Romani GL, Buitelaar JK,
769 Bekkering H (2011) Altered Intrinsic Functional Connectivity of Anterior and Posterior
770 Insula Regions in High-Functioning Participants With Autism Spectrum Disorder. *Hum*
771 *Brain Mapp* 32:1013–1028.
- 772 Eickhoff SB, Schleicher A, Zilles K (2006) The Human Parietal Operculum . I .
773 Cytoarchitectonic Mapping of Subdivisions. *Cereb Cortex* 15:254–267.
- 774 Ekman P (1992) An argument for basic emotions. *Cogn Emot* 6:169–200.
- 775 Escoffier N, Zhong J, Schirmer A, Qiu A (2013) Emotional expressions in voice and music:
776 Same code, same effect? *Hum Brain Mapp* 34:1796–1810.

- 777 Ethofer T, Van De Ville D, Scherer K, Vuilleumier P (2009) Decoding of Emotional Information
778 in Voice-Sensitive Cortices. *Curr Biol* 19:1028–1033.
- 779 Fritz T, Jentschke S, Gosselin N, Sammler D, Peretz I, Turner R, Friederici AD, Koelsch S (2009)
780 Universal Recognition of Three Basic Emotions in Music. *Curr Biol* 19:573–576.
- 781 Frühholz S, Trost W, Grandjean D (2014) The role of the medial temporal limbic system in
782 processing emotions in voice and music. *Prog Neurobiol* 123:1–17.
- 783 Hailstone JC, Omar R, Henley SMD, Frost C, Michael G, Warren JD, Hailstone JC, Omar R,
784 Henley SMD, Frost C, Hailstone JC, Omar R, Henley SMD, Frost C, Warren JD (2009) It's
785 not what you play , it's how you play it: Timbre affects perception of emotion in music. *Q J*
786 *Exp Psychol* 62:2141–2155.
- 787 Heller R, Stanley D, Yekutieli D, Nava R, Benjamini Y (2006) Cluster-based analysis of fMRI
788 data. *Neuroimage* 33:599–608.
- 789 Immordino-Yang MH, Yang X-F, Damasio H (2014) Correlations between social-emotional
790 feelings and anterior insula activity are independent from visceral states but influenced by
791 culture. *Front Hum Neurosci* 8:1–15.
- 792 Johnsen EL, Tranel D, Lutgendorf S, Adolphs R (2009) A neuroanatomical dissociation for
793 emotion induced by music. *Int J Psychophysiol* 72:24–33.
- 794 Juslin PN, Laukka P (2003) Communication of emotions in vocal expression and music
795 performance: different channels, same code? *Psychol Bull* 129:770–814.
- 796 Kao M, Mandal A, Lazar N, Stufken J (2009) NeuroImage Multi-objective optimal experimental
797 designs for event-related fMRI studies. *Neuroimage* 44:849–856.
- 798 Kaplan JT, Man K, Greening SG (2015) Multivariate cross-classification : applying machine
799 learning techniques to characterize abstraction in neural representations. *Front Hum*

- 800 Neurosci 9:1–12.
- 801 Kim J, Shinkareva S V, Wedell DH (2017) Representations of modality-general valence for
802 videos and music derived from fMRI data. *Neuroimage* 148:42–54.
- 803 Kim YE, Schmidt EM, Migneco R, Morton BG, Richardson P, Scott J, Speck J a, Turnbull D
804 (2010) Music Emotion Recognition : a State of the Art Review. *Inf Retr Boston*:255–266.
- 805 Kleiner M, Brainard DH, Pelli D (2007) What’s new in Psychtoolbox-3? *Percept 36 ECVF Abstr*
806 *Suppl.*
- 807 Koelsch S (2014) Brain correlates of music-evoked emotions. *Nat Rev Neurosci* 15:170–180.
- 808 Kotz SA, Kalberlah C, Bahlmann J, Friederici AD, Haynes JD (2013) Predicting vocal emotion
809 expressions from the human brain. *Hum Brain Mapp* 34:1971–1981.
- 810 Kragel PA, LaBar KS (2015) Multivariate neural biomarkers of emotional states are
811 categorically distinct. *Soc Cogn Affect Neurosci* 10:1437–1448.
- 812 Kriegeskorte N, Goebel R, Bandettini P (2006) Information-based functional brain mapping.
- 813 Kurth F, Zilles K, Fox PT, Laird AR, Eickhoff SB (2010) A link between the systems: functional
814 differentiation and integration within the human insula revealed by meta-analysis. *Brain*
815 *Struct Funct* 214:519–534.
- 816 Lamm C, Batson CD, Decety J (2007) The neural substrate of human empathy: effects of
817 perspective-taking and cognitive appraisal. *J Cogn Neurosci* 19:42–58.
- 818 Lartillot O, Lartillot O, Toivianen P, Toivianen P (2007) A matlab toolbox for musical feature
819 extraction from audio. *Int Conf Digit Audio ...*:1–8.
- 820 Linke AC, Cusack R (2015) Flexible Information Coding in Human Auditory Cortex during
821 Perception, Imagery, and STM of Complex Sounds. *J Cogn Neurosci* 27.
- 822 Man K, Damasio A, Meyer K, Kaplan JT (2015) Convergent and Invariant Object

- 823 Representations for Sight, Sound, and Touch. *Hum Brain Mapp* 36:3629–3640.
- 824 Mullensiefen D, Gingas B, Musil J, Steward L (2014) The Musicality of Non-Musicians : An
825 Index for Assessing Musical Sophistication in the General Population. *PLoS One* 9.
- 826 Norman K a, Polyn SM, Detre GJ, Haxby J V (2006) Beyond mind-reading: multi-voxel pattern
827 analysis of fMRI data. *Trends Cogn Sci* 10:424–430.
- 828 Paquette S, Peretz I, Belin P (2013) The “ Musical Emotional Bursts ”: a validated set of musical
829 affect bursts to investigate auditory affective processing. *Front Psychol* 4:1–7.
- 830 Park M, Hennig-Fast K, Bao Y, Carl P, Pöppel E, Welker L, Reiser M, Meindl T, Gutyrchik E
831 (2013) Personality traits modulate neural responses to emotions expressed in music. *Brain*
832 *Res* 1523:68–76.
- 833 Peelen M V, Atkinson AP, Vuilleumier P (2010) Supramodal Representations of Perceived
834 Emotions in the Human Brain. *J Neurosci* 30:10127–10134.
- 835 Rigoulot S, Pell MD, Armony JL (2015) Time course of the influence of musical expertise on the
836 processing of vocal and musical sounds. *Neuroscience* 290:175–184.
- 837 Rijn S Van, Aleman A, Diessen E Van, Berckmoes C, Vingerhoets G, Kahn RS (2005) What is
838 said or how it is said makes a difference : role of the right fronto-parietal operculum in
839 emotional prosody as revealed by repetitive TMS. *Eur J Neurosci* 21:3195–3200.
- 840 Saarimaki H, Gotsopoulos A, Jaaskelainen IP, Lampinen J, Vuilleumier P, Hari R, Sams M,
841 Nummenmaa L (2015) Discrete Neural Signatures of Basic Emotions. *Cereb Cortex*:1–11.
- 842 Salimpoor VN, Zald DH, Zatorre RJ, Dagher A, McIntosh AR (2015) Predictions and the brain :
843 how musical sounds become rewarding. *Trends Cogn Sci* 19:86–91.
- 844 Sander K, Scheich H (2005) Left Auditory Cortex and Amygdala , but Right Insula Dominance
845 for Human Laughing and Crying. *J Cogn Neurosci* 17:1519–1531.

- 846 Schirmer A, Adolphs R (2017) Emotion Perception from Face, Voice, and Touch: Comparisons
847 and Convergence. *Trends Cogn Sci* 21:216–228.
- 848 Schonwiesner M, RübSamen R, von Cramon DY (2005) Hemispheric asymmetry for spectral and
849 temporal processing in the human antero-lateral auditory belt cortex. *Eur J Neurosci*
850 22:1521–1528.
- 851 Silani G, Bird G, Brindley R, Singer T, Frith C, Frith U (2008) Levels of emotional awareness
852 and autism : An fMRI study. *Soc Neurosci* 3:97–112.
- 853 Singer T, Seymour B, Doherty JO, Kaube H, Dolan RJ, Frith CD (2004) Empathy for Pain
854 Involves the Affective but not Sensory Components of Pain. *Science* (80-) 303:1157–1162.
- 855 Skerry AE, Saxe R (2014) A Common Neural Code for Perceived and Inferred Emotion. *J*
856 *Neurosci* 34:15997–16008.
- 857 Smith SM, Jenkinson M, Woolrich MW, Beckmann CF, Behrens TEJ, Johansen-Berg H,
858 Bannister PR, De Luca M, Drobnjak I, Flitney DE, Niazy RK, Saunders J, Vickers J, Zhang
859 Y, De Stefano N, Brady JM, Matthews PM (2004) Advances in functional and structural
860 MR image analysis and implementation as FSL. *Neuroimage* 23:S208-19.
- 861 Smith SM, Nichols TE (2009) Threshold-free cluster enhancement : Addressing problems of
862 smoothing , threshold dependence and localisation in cluster inference. *Neuroimage* 44:83–
863 98.
- 864 Stelzer J, Chen Y, Turner R (2013) Statistical inference and multiple testing correction in
865 classification-based multi-voxel pattern analysis (MVPA): Random permutations and
866 cluster size control. *Neuroimage* 65:69–82.
- 867 Straube T, Miltner WHR (2011) Attention to aversive emotion and specific activation of the
868 right insula and right somatosensory cortex. *Neuroimage* 54:2534–2538.

869 von dem Hagen EAH, Nummenmaa L, Yu R, Engell AD, Ewbank MP, Calder AJ (2011) Autism
870 Spectrum Traits in the Typical Population Predict Structure and Function in the Posterior
871 Superior Temporal Sulcus. *Cereb Cortex* 21:492–500.

872 Wegrzyn M, Riehle M, Labudda K, Woermann F, Baumgartner F, Pollmann S, Bien CG, Kissler
873 J (2015) Investigating the brain basis of facial expression perception using multi-voxel
874 pattern analysis. *Cortex* 69:131–140.

875 Winkler AM, Ridgway GR, Webster MA, Smith SM, Nichols TE (2014) Permutation inference
876 for the general linear model. *Neuroimage* 92:381–397.

877

878

879

880

881

882

883

884

885

886

887

888

889

890

891

892

893

894

895

896

897
898
899
900
901
902
903
904
905
906
907
908
909

ACCEPTED MANUSCRIPT