

1 Human temporal voice areas are sensitive to chimpanzee vocalizations

2

3 Leonardo Ceravolo^{1*}, Coralie Debracque^{1*}, Thibaud Gruber^{1‡} & Didier Grandjean^{1‡}

4

5 * shared first authors

6 ‡ shared last authors

7

8 ¹University of Geneva, Geneva, Switzerland

9

10

11

12 Corresponding author:

13 Leonardo Ceravolo, PhD

14 Boulevard Pont-d'Arve 40

15 CH-1205 Geneva

16 Switzerland

17 e-mail: leonardo.ceravolo@unige.ch

18

19 **Abstract**

20 In recent years, research on voice processing, particularly the study of temporal voice areas (TVA), was
21 dedicated almost exclusively to human voice. To characterize commonalities and differences regarding
22 primate vocalization representations in the human brain, the inclusion of closely related primates,
23 especially chimpanzees and bonobos, is needed. We hypothesized that commonalities would depend on
24 both phylogenetic and acoustic proximity, with chimpanzees ranking the closest to *Homo*. Presenting
25 human participants with four primate species vocalizations (rhesus macaques, chimpanzees, bonobos
26 and humans) and taking into account acoustic distance or removing voxels explained solely by
27 vocalization low-level acoustics, we observed within-TVA enhanced left and right anterior superior
28 temporal gyrus activity for chimpanzee compared to all other species, and chimpanzee compared to
29 human vocalizations. Our results provide evidence for a common neural basis in the TVA for the
30 processing of phylogenetically and acoustically close vocalizations, namely those of humans and
31 chimpanzees.

32

33 **Introduction**

34 The study of the cerebral mechanisms underlying speech and voice processing has gained steam since
35 the early 2000s with the emergence of functional magnetic resonance imaging (fMRI)¹. Voice-sensitive
36 areas, generally referred to as ‘temporal voice areas’ (TVA), have been highlighted along the upper,
37 superior part of the temporal cortex². Since then, great effort has been put into better characterizing these
38 TVA, with a specific focus on their spatial compartmentalization into functional subparts³⁻⁵. Repetitive
39 transcranial magnetic stimulations over the right mid TVA lead to persistent voice detection impairment
40 in a simple voice/non-voice discrimination task⁶ and a rather large body of literature is aligned with the
41 crucial role of the TVA in voice perception and processing^{3,7-9}. Subparts of the TVA have also been
42 directly linked to social perception¹⁰, vocal emotion processing^{11,12}, voice identity^{13,14} and gender¹⁵
43 perception. The developmental axis of voice processing has also been studied in infants, revealing the
44 existence of TVA as early as 7 but not 4 month-olds in the human brain¹⁶ while *in utero* fetuses have
45 been shown to be already able to recognize their parents’ voice¹⁷. With the constant development of
46 brain imaging and analysis techniques¹⁸, it is realistic to expect successful, though non-invasive, *in utero*
47 ‘task-related’ voice perception fMRI results in the near future. Along the evolutionary axis, evidence
48 for TVA or more generally voice-sensitive brain areas have emerged most notably for dogs¹⁹ and
49 monkeys^{20,21} (*Macaca mulatta*), raising the questions of whether TVA are species-specific²² and to
50 which extent human and non-human primates share neural mechanisms enabling them to process
51 conspecific vocalizations²³. Less attention has however been devoted to paradigms presenting animal
52 vocalizations to humans, and no study to date has ever reported human TVA activations for the
53 processing of such auditory material, namely other animals’ vocalizations. Human processing of animal
54 vocalizations has been studied using both monkey and cat material but no specific activations related to
55 any of the species was observed²⁴. Other studies have focused more specifically on phylogenetic
56 distance, including as stimuli human great ape (chimpanzee, *Pan troglodytes*) and old-world-monkey
57 (rhesus macaque, *Macaca mulatta*) vocalizations. Such studies could not identify species-specific brain
58 activations in spite of the correct discrimination of chimpanzee affective vocalizations²⁵, and observed
59 below²⁵ vs. above²⁶ chance discrimination of affective macaque vocalizations by human participants.

60 This scarce literature motivated the present study that aims at a reliable investigation of species-specific
61 TVA activations in humans asked to categorize phylogenetically close and distant species' vocalizations
62 while undergoing fast fMRI scanning. The importance of between-species acoustic differences and
63 distance, especially fundamental frequency was also of major interest^{27,28}. We therefore included
64 vocalizations of our closet sister taxon, *Pan* (chimpanzees; bonobos, *Pan paniscus*), whose estimated
65 split with *Homo* is only 6-8 million years ago as well as phylogenetically more distant species
66 (cercopithecidae: rhesus macaques, with an estimated split with *Homo* 25 million years ago). In fact,
67 any claim of human uniqueness for recruiting the TVA remains on hold and should be tested in light of
68 these closely related species. Bonobo vocalizations are of particular interest, as this species is thought
69 to have experienced evolutionary changes in their communication in part due to a neoteny process
70 involving acoustic modifications (i.e., fundamental frequency)²⁷ even though they are as
71 phylogenetically close to humans as chimpanzees²⁹. Whether such changes would affect the abilities of
72 human participants to recognize their calls should therefore be investigated in comparison to chimpanzee
73 and rhesus macaque vocalizations. We therefore predicted: i) acoustic proximity for human and
74 chimpanzee vocalizations, while more distance would separate those of bonobo and macaque
75 vocalizations; ii) an overlap between brain networks of *Homo* and the *Pan* branch (chimpanzee, bonobo)
76 but not the cercopithecidae (rhesus macaque) vocalizations; iii) shared and localized brain activations
77 for the categorization of human and chimpanzee vocalizations extending to the TVA, depending on both
78 phylogenetic proximity and acoustic distance. These hypotheses involve: a) a control of low-level
79 acoustic differences, namely vocalization mean fundamental frequency and energy—included as trial-
80 level covariates of 'no-interest' in the first neuroimaging statistical model; b) the inclusion of a measure
81 of acoustic distance—included as trial-level covariate of 'interest' in a second neuroimaging statistical
82 model.

83 **Material and Methods**

84 *Species categorization task*

85 *Participants*

86 Twenty-five right-handed, healthy, either native or highly proficient French-speaking participants took
87 part in the study. One participant was excluded because he had no correct response at all and may have
88 fallen asleep, while another participant was excluded due to incomplete scanning and technical issues,
89 leaving us with twenty-three participants (10 female, 13 male, mean age 24.65 years, SD 3.66). All
90 participants were naive to the experimental design and study, had normal or corrected-to-normal vision,
91 normal hearing and no history of psychiatric or neurologic incidents. Participants gave written informed
92 consent for their participation in accordance with ethical and data security guidelines of the University
93 of Geneva. The study was approved by the Ethics Cantonal Commission for Research of the Canton of
94 Geneva, Switzerland (CCER) and was conducted according to the Declaration of Helsinki.

95 *Stimuli*

96 Seventy-two vocalizations of four primate species (human, chimpanzee, bonobo and rhesus macaque)
97 were used in this study (see Fig.1a). Therefore, eighteen human voices were selected and they were
98 expressed by two male and two female actors, obtained from a nonverbal validated stimuli set of Belin
99 and collaborators³⁰. The eighteen selected chimpanzee, bonobo and rhesus macaque vocalizations
100 contained single calls or call sequences produced by 6 to 8 different individuals in their natural
101 environment. All vocal stimuli were standardized to 750 milliseconds using PRAAT (www.praat.org)
102 but were not normalized in any way in order to preserve the naturality of the sounds³¹ and to allow for
103 low-level acoustic parameters of interest to be used in data modelling.

104 *Experimental procedure and paradigm*

105 Laying comfortably in a 3T scanner, participants listened to a total of seventy-two stimuli randomized
106 and played binaurally using MRI compatible earphones at 70 dB SPL. At the beginning of the
107 experiment, participants were instructed to identify the species that expressed the vocalizations using a
108 keyboard. For instance, the instructions could be “Human – press 1, Chimpanzee – press 2, Bonobo –

109 press 3 or Macaque – press 4”. The pressed keys were randomly assigned across participants. In a 3-5
110 second interval (jittering of 400 ms) after each stimulus, participants were asked to categorize the
111 species. If the participant did not respond during this interval, the next stimulus followed automatically.

112

113 *Temporal voice areas localizer task*

114 *Participants*

115 One-hundred and fifteen right-handed, healthy, either native or highly proficient French-speaking
116 participants (62 female, 54 male, mean age 25.34 years, SD 5.50) were included in this functional
117 magnetic resonance task. Among these participants, twenty-two out of the twenty-three who performed
118 the species categorization task were included (the temporal voice areas localizer task was not acquired
119 for one of them due to technical issues). All participants were naive to the experimental design and
120 study, had normal or corrected-to-normal vision, normal hearing and no history of psychiatric or
121 neurologic incidents. Participants gave written informed consent for their participation in accordance
122 with ethical and data security guidelines of the University of Geneva. The study was approved by the
123 Ethical Committee of the University of Geneva and was conducted according to the Declaration of
124 Helsinki.

125 *Stimuli and paradigm*

126 Auditory stimuli consisted of sounds from a variety of sources². Vocal stimuli were obtained from 47
127 speakers: 7 babies, 12 adults, 23 children and 5 older adults. Stimuli included 20 blocks of vocal sounds
128 and 20 blocks of non-vocal sounds. Vocal stimuli within a block could be either speech 33%: words,
129 non-words, foreign language or non-speech 67%: laughs, sighs, various onomatopoeia. Non-vocal
130 stimuli consisted of natural sounds 14%: wind, streams, animals 29%: cries, gallops, the human
131 environment 37%: cars, telephones, airplanes or musical instruments 20%: bells, harp, instrumental
132 orchestra. The paradigm, design and stimuli were obtained through the Voice Neurocognition
133 Laboratory website (<http://vnl.psy.gla.ac.uk/resources.php>). Stimuli were presented at an intensity that

134 was kept constant throughout the experiment 70 dB sound-pressure level. Participants were instructed
135 to actively listen to the sounds. The silent interblock interval was 8 s long.

136

137 ***Behavioral data analysis***

138 *Accuracy*

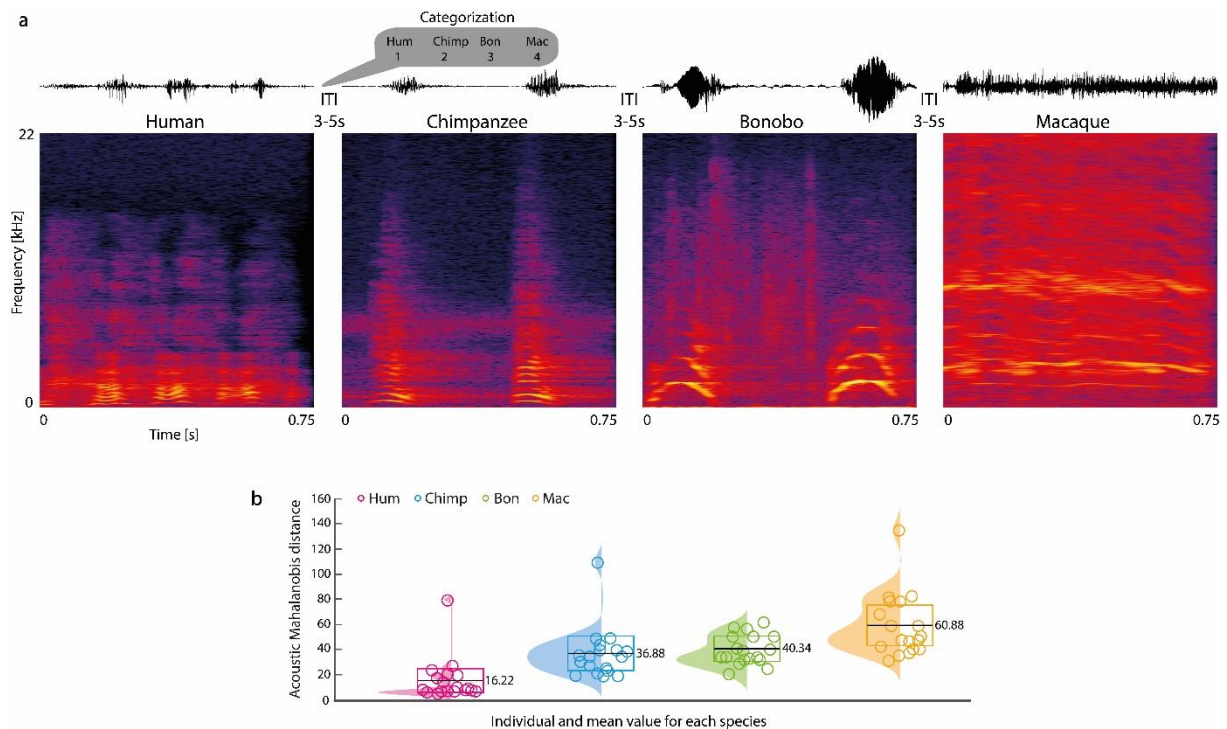
139 Behavioral data were exclusively used to exclude participants who had below chance level
140 categorization of human voices. Therefore, data from twenty-three participants mentioned in the *Species*
141 *Categorization Task - Participants* section above were analyzed using R studio software (R Studio
142 team³² Inc., Boston, MA, url: <http://www.rstudio.com/>). These data are reported in the supplementary
143 materials (Fig.S1) since they are not part of the questions of interest of this paper addressing neural
144 correlates of the species-specific processing of vocalizations within the temporal voice areas in human
145 participants.

146

147 ***Acoustic Mahalanobis distances***

148 To quantify the impact of acoustic similarities in human recognition of affective vocalizations of other
149 primates, we extracted 88 acoustic parameters from all vocalizations using the extended Geneva
150 Acoustic parameters set defined as the optimal acoustic indicators related to voice analysis (GeMAPS)³³.
151 This set of acoustical parameters was selected based on: i) their potential to index affective physiological
152 changes in voice production, ii) their proven value in former studies as well as their automatic
153 extractability, and iii) their theoretical significance. Then, to assess the acoustic distance between
154 vocalizations of all species, we ran a General Discriminant Analysis model (GDA). More precisely, we
155 used the 88 acoustical parameters in a GDA in order to discriminate our stimuli based on the different
156 species (human, chimpanzee, bonobo, and rhesus macaque). Excluding the acoustical variables with the
157 highest correlations ($r > .90$) to avoid redundancy of acoustic parameters, we retained 16 acoustic
158 parameters.

159 We subsequently computed Mahalanobis distances to classify the 96 stimuli on these selected acoustical
160 features. A Mahalanobis distance is a generalized pattern analysis comparing the distance of each
161 vocalization from the centroids of the different species vocalizations. This analysis allowed us to obtain
162 an acoustical distance matrix used to test how the acoustical distances were differentially related to the
163 different species (see Fig.1bc).



164

165 **Fig.1: Timecourse of the species categorization task with stimuli example and acoustic distance**
166 **data.** a, Detail of the timecourse of four trials of the species categorization task in non-representative
167 order, including waveform and spectrogram graphs for one example stimulus of each species. b, Scatter
168 plot of the acoustic Mahalanobis distance data of each stimulus for each species including mean
169 (numbers represent exact mean value) and box plots of the standard error of the mean in addition to
170 distribution fit. ITI: inter trial interval; Hum: human; Chimp: chimpanzee; Bon: bonobo; Mac: macaque.

171

172 *Imaging data acquisition*

173 *Species categorization task*

174 Structural and functional brain imaging data were acquired by using a 3T scanner Siemens Trio,
175 Erlangen, Germany with a 32-channel coil. A 3D GR\IR magnetization-prepared rapid acquisition
176 gradient echo sequence was used to acquire high-resolution ($0.35 \times 0.35 \times 0.7 \text{ mm}^3$) T1-weighted
177 structural images (TR = 2400 ms, TE = 2.29 ms). Functional images were acquired by using fast fMRI,
178 with a multislice echo planar imaging sequence with 79 transversal slices in descending order, slice

179 thickness 3 mm, TR = 650 ms, TE = 30 ms, field of view = 205 x 205 mm², 64 x 64 matrix, flip angle
180 = 50 degrees, bandwidth 1562 Hz/Px. In total for this task, 636 functional volumes of 79 slices were
181 acquired for each participant for a total of 50244 slices per participant. For our whole sample of twenty-
182 three participants, 14628 volumes were acquired for a grand total of 1'155'612 slices.

183 *Temporal voice areas localizer task*

184 Structural and functional brain imaging data were acquired by using a 3T scanner Siemens Trio,
185 Erlangen, Germany with a 32-channel coil. A magnetization-prepared rapid acquisition gradient echo
186 sequence was used to acquire high-resolution (1 x 1 x 1 mm³) T1-weighted structural images TR = 1,900
187 ms, TE = 2.27 ms, TI = 900 ms. Functional images were acquired by using a multislice echo planar
188 imaging sequence with 36 transversal slices in descending order, slice thickness 3.2 mm, TR = 2,100
189 ms, TE = 30 ms, field of view = 205 x 205 mm², 64 x 64 matrix, flip angle = 90°, bandwidth 1562
190 Hz/Px. In total for this task, 230 functional volumes of 36 slices were acquired for each participant for
191 a total of 8280 slices per participant. For our whole sample of one hundred and fifteen participants,
192 26450 volumes were acquired for a grand total of 952'200 slices.

193

194 ***Wholebrain data analysis***

195 *Species categorization task region-of-interest analysis within the temporal voice areas*

196 Functional images were analyzed with Statistical Parametric Mapping software (SPM12, Wellcome
197 Trust Centre for Neuroimaging, London, UK). Preprocessing steps included realignment to the first
198 volume of the time series, slice timing, normalization into the Montreal Neurological Institute³³ (MNI)
199 space using the DARTEL toolbox³⁴ and spatial smoothing with an isotropic Gaussian filter of 8 mm full
200 width at half maximum. To remove low-frequency components, we used a high-pass filter with a cutoff
201 frequency of 128 s. Two general linear models were used to compute first-level statistics, in which each
202 event was modeled by using a boxcar function and was convolved with the hemodynamic response
203 function, time-locked to the onset of each stimulus. In model 1, separate regressors were created for all
204 trials of each species (Species factor: human, chimpanzee, bonobo, macaque vocalizations) and two

205 covariates of no-interest each (mean fundamental frequency and mean energy of each species) for a total
206 of 12 regressors. Finally, six motion parameters were included as regressors of no interest to account for
207 movement in the data and our design matrix therefore included a total of 18 columns plus the constant
208 term. The species regressors were used to compute simple contrasts for each participant, leading to
209 separate main effects of human, chimpanzee, bonobo and macaque vocalizations. Covariates were set
210 to zero in order to model them as no-interest regressors. In model 2, separate regressors were created
211 for all trials of each species (Species factor: human, chimpanzee, bonobo, macaque vocalizations) and
212 one covariate of interest for each species (acoustic distance for each species relative to human voice
213 stimuli) for a total of 8 regressors. Finally, six motion parameters were included as regressors of no
214 interest to account for movement in the data and our design matrix therefore included a total of 14
215 columns plus the constant term. The species regressors were used to compute simple contrasts for each
216 participant, leading to separate main effects of human, chimpanzee, bonobo and macaque vocalizations
217 including acoustic distance (the covariate was set to one in order to model it as ‘of interest’ regressor).
218 For each model, each of their respective four simple contrasts were then taken to two flexible factorial
219 second-level analyses. For both of these second-level analyses there were two factors: the Participants
220 factor (independence set to yes, variance set to unequal) and the Species factor (independence set to no,
221 variance set to unequal). For these analyses and to be consistent, we only included participants who were
222 above chance level (25%) in the species categorization task (N=18). Brain region labelling was defined
223 using xjView toolbox (<http://www.alivelearn.net/xjview>). All neuroimaging activations were
224 thresholded in SPM12 by using a voxelwise false discovery rate (FDR) correction at $p < .05$ and an
225 arbitrary cluster extent of $k > 10$ voxels to remove very small clusters of activity.

226 *Temporal voice areas localizer task*

227 Functional images were analyzed with Statistical Parametric Mapping software (SPM12, Wellcome
228 Trust Centre for Neuroimaging, London, UK). Preprocessing steps included realignment to the first
229 volume of the time series, slice timing, normalization into the Montreal Neurological Institute³³ (MNI)
230 space using the DARTEL toolbox³⁴ and spatial smoothing with an isotropic Gaussian filter of 8 mm full
231 width at half maximum. To remove low-frequency components, we used a high-pass filter with a cutoff

232 frequency of 128 s. A general linear model was used to compute first-level statistics, in which each
233 block was modeled by using a block function and was convolved with the hemodynamic response
234 function, time-locked to the onset of each block. Separate regressors were created for each condition
235 (vocal and non-vocal; condition factor). Finally, six motion parameters were included as regressors of
236 no interest to account for movement in the data. The condition regressors were used to compute simple
237 contrasts for each participant, leading to a main effect of vocal and non-vocal at the first-level of
238 analysis: [1 0] for vocal, [0 1] for non-vocal. These simple contrasts were then taken to a flexible
239 factorial second-level analysis in which there were two factors: Participants factor (independence set to
240 yes, variance set to unequal) and the Condition factor (independence set to no, variance set to unequal).
241 All neuroimaging activations were thresholded in SPM12 by using a voxelwise family-wise error (FWE)
242 correction at $p < .05$. Activation outline for vocal > nonvocal was precisely delineated and overlaid on
243 brain displays of the species categorization task.

244

245 **Results**

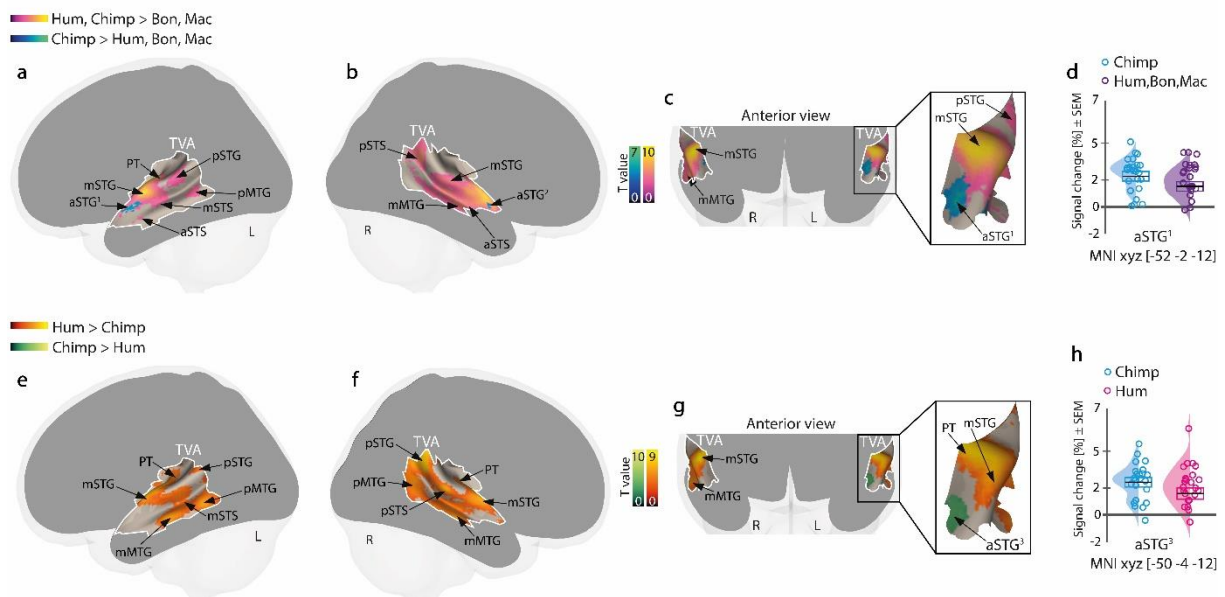
246 *Neuroimaging data within the temporal voice areas*

247 We adopted a region-of-interest approach to uncover functional changes relative to species
248 categorization and processing within the temporal voice areas, as delineated in our hypotheses. Low-
249 level acoustics were used in two distinct models, namely vocalization mean energy and mean
250 fundamental frequency (covariates of no-interest at the trial level, model 1) and a measure of acoustic
251 distance (covariate of interest, model 2). We were particularly interested in brain activity while
252 processing vocalizations of our closest relative (both acoustically and phylogenetically), the
253 chimpanzee. The present study did not aim at uncovering wholebrain results underlying the processing
254 of each species' vocalizations (see Fig.S2 and Fig.S3), although statistics presented in this section were
255 computed with a voxelwise approach on the wholebrain for higher data reproducibility and
256 generalizability.

257 *Model 1: Effects of species processing with vocalization mean energy and mean fundamental frequency*
258 *as covariates of no-interest at the trial level*

259 In this model, we wanted to remove from species' processing brain activations the part of variance
260 correlated with low-level acoustics of no-interest, namely mean voice energy and fundamental
261 frequency. Brain activations common to human and chimpanzee vocalizations using the [human,
262 chimpanzee > bonobo, macaque] contrast led to enhanced signal in the bilateral posterior, mid and
263 anterior superior temporal cortex (Fig.2abcd, Table 1). Brain activity specific to chimpanzee
264 vocalizations ([chimpanzee > human, bonobo, macaque]) led to enhanced activity in a cluster of the left
265 anterior STG located within the temporal voice areas (Fig.2c). A similar result was observed when
266 directly contrasting chimpanzee to human vocalizations ([chimpanzee > human]) in a slightly more
267 medial area of the anterior STG, also located again within the voice-sensitive areas (Fig.2g, Table 1).
268 Enhanced activity for human relative to chimpanzee vocalizations ([human > chimpanzee]) was
269 observed in large parts of the anterior, mid and posterior superior and middle temporal cortex (Fig.2efg,
270 Table 1). No voxels reached significance either at the wholebrain level or within the TVA for both the
271 [bonobo > human, chimpanzee, macaque] and the [macaque > human, chimpanzee, bonobo] contrasts.

Model 1: Species processing, mean of vocalization fundamental frequency and energy as covariates of no-interest (whole brain voxelwise $p < .05$ FDR, $k > 10$)



272 **Fig.2: Wholebrain results when selectively contrasting processing of chimpanzee to other species'**
273 **vocalizations with mean fundamental frequency and energy as trial-level covariates of no-interest.**
274 **abc**, Enhanced brain activity for human and chimpanzee compared to bonobo and macaque
275 vocalizations (purple to yellow) on a sagittal view, overlaid with activity specific to chimpanzee
276 vocalizations (blue to green) on a sagittal view, overlaid with activity specific to human vocalizations (orange to red) on a sagittal view, overlaid with activity specific to human vocalizations (green to blue) on a sagittal view.

277 vocalizations (dark blue to green). **d**, Percentage of signal change for each individual and species in the
278 left anterior superior temporal gyrus (aSTG¹). Box plots represent mean value (black line) and the
279 standard error of the mean with distribution fit. **efg**, Direct comparison between human and chimpanzee
280 vocalizations (human > chimpanzee: dark red to yellow; chimpanzee > human: dark green to yellow)
281 on a sagittal render. **h**, Percentage of signal change in a more medial part of the anterior superior
282 temporal gyrus (aSTG³) when contrasting chimpanzee to human vocalizations for each individual and
283 species with box plots representing mean value (black line) and the standard error of the mean with
284 distribution fit. Brain activations are independent of low-level acoustic parameters for all species
285 (fundamental frequency 'F0' and mean energy of vocalizations). Data corrected for multiple comparison
286 using wholebrain voxelwise false discovery rate (FDR) at a threshold of $p < .05$. Percentage of signal
287 change extracted at cluster peak including 9 surrounding voxels, selecting among these the ones
288 explaining at least 85% of the variance using singular value decomposition. Hum: human; Chimp:
289 chimpanzee; Bon: bonobo; Mac: macaque. TVA: temporal voice areas. 'a' prefix: anterior; 'm' prefix:
290 mid; 'p' prefix: posterior; MTG: middle temporal gyrus; STG: superior temporal gyrus; STG: superior
291 temporal gyrus; STS: superior temporal sulcus; PT: planum temporale; L: left; R: right.

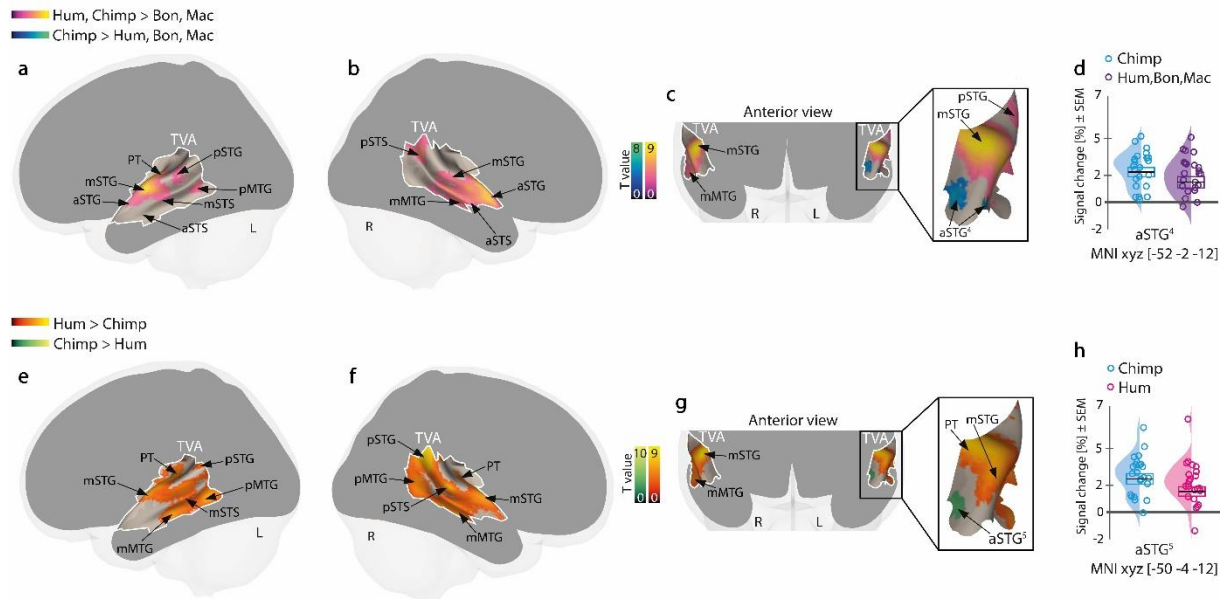
292

293 *Model 2: Effects of species processing with vocalization acoustic distance from human voice, per*
294 *species, as covariate of interest at the trial level*

295 In this second model, we wanted to add to species' processing brain activations the part of variance
296 correlated with acoustic distance between each species and the human voice. Brain activations specific
297 to human and chimpanzee vocalizations using the [human, chimpanzee > bonobo, macaque] contrast
298 led to enhanced signal in the bilateral posterior, mid and anterior superior temporal cortex (Fig.3**abcd**,
299 Table 2). Brain activity specific to chimpanzee vocalizations ([chimpanzee > human, bonobo, macaque])
300 led to enhanced activity in a cluster of the left anterior STG located within the temporal voice areas
301 (Fig.3**c**). A similar result was observed when directly contrasting chimpanzee to human vocalizations
302 ([chimpanzee > human]) in a slightly more medial area of the anterior STG, also located again within
303 the voice-sensitive areas (Fig.3**g**, Table 2). Enhanced activity for human relative to chimpanzee
304 vocalizations ([human > chimpanzee]) was observed in large parts of the anterior, mid and posterior
305 superior and middle temporal cortex (Fig.3**efg**, Table2). Again using this model, no voxels reached
306 significance either at the wholebrain level or within the TVA for both the [bonobo > human, chimpanzee,
307 macaque] and the [macaque > human, chimpanzee, bonobo] contrasts.

308

Model 2: Species processing, inter-species vocalization acoustic distance as covariate of interest (whole brain voxelwise $p < .05$ FDR, $k > 10$)



309
 310 **Fig.3: Wholebrain results when selectively contrasting processing of chimpanzee to other species'**
 311 **vocalizations with acoustic distance as trial-level covariate of interest. abc,** Enhanced brain activity
 312 for human and chimpanzee compared to bonobo and macaque vocalizations (purple to yellow) on a
 313 sagittal view, overlaid with activity specific to chimpanzee vocalizations (dark blue to green). **d,**
 314 Percentage of signal change for each individual and species in the left anterior superior temporal gyrus
 315 ($aSTG^4$). Box plots represent mean value (black line) and the standard error of the mean with distribution
 316 fit. **efg,** Direct comparison between human and chimpanzee vocalizations (human > chimpanzee: dark
 317 red to yellow; chimpanzee > human: dark green to yellow) on a sagittal render. **h,** Percentage of signal
 318 change in a more medial part of the anterior superior temporal gyrus ($aSTG^5$) when contrasting
 319 chimpanzee to human vocalizations for each individual and species with box plots representing mean
 320 value (black line) and the standard error of the mean with distribution fit. Brain activations are dependent
 321 of acoustic Mahalanobis distance between each species, see Methods for details. Data corrected for
 322 multiple comparison using wholebrain voxelwise false discovery rate (FDR) at a threshold of $p < .05$.
 323 Percentage of signal change extracted at cluster peak including 9 surrounding voxels, selecting among
 324 these the ones explaining at least 85% of the variance using singular value decomposition. Hum: human;
 325 Chimp: chimpanzee; Bon: bonobo; Mac: macaque. TVA: temporal voice areas. 'a' prefix: anterior; 'm'
 326 prefix: mid; 'p' prefix: posterior; MTG: middle temporal gyrus; STG: superior temporal gyrus; STG:
 327 superior temporal gyrus; STS: superior temporal sulcus; PT: planum temporale; L: left; R: right.

328

329 Discussion

330 The present study provides evidence of the sensitivity of the TVA to chimpanzee vocalizations,
 331 materialized by chimpanzee-specific enhanced activity in the left and right anterior STG. Second, our
 332 results highlight shared brain networks for the processing of both human voices and chimpanzee calls
 333 involving posterior, mid and anterior parts of bilateral superior temporal gyrus. Therefore, our results
 334 suggest that vocalizations expressed by another great ape species can also recruit subparts of the human
 335 temporal cortex normally dedicated to the processing of human voices, namely the anterior TVA.

336 Because we controlled our analyses for low-level acoustics and acoustic distance, we importantly
337 demonstrate that similar TVA activity for the processing of human voices and chimpanzee vocalizations
338 directly relate to both phylogenetic and acoustic proximity.

339 Often linked to the processing of conspecific vocalizations only (e.g., in humans²; macaques^{20,21}; and
340 dogs¹⁹), the present study questions the current view of TVA ‘selectivity’ showing that both human
341 voices and chimpanzee calls enhance activity in the anterior TVA. Indeed, our neuroimaging analyses
342 revealed the specific involvement of the left anterior STG when processing chimpanzee vocalizations.
343 No specific results were observed for bonobo and macaque vocalizations, respectively. Furthermore,
344 such anterior STG activity was also observed when a direct comparison between chimpanzee calls and
345 human voice was made. This result adds to the specificity of subparts of the anterior TVA for the
346 processing of species with human-like phylogeny and acoustics. Differences at the level of processing
347 complexity between the two vocalizations could explain such observations. In fact, previous studies
348 have shown the role of the left anterior STG and anterior STS in the conceptual representation of social
349 context through the human voice³⁵⁻³⁷. Hence, our data could suggest that the anterior part of the left
350 superior temporal cortex is recruited to process the social context of vocal stimuli expressed by human
351 and chimpanzee species. Yet, this processing would be more automated for the perception of human
352 voice due to our high exposition and expertise as humans, as opposed to chimpanzee calls that we do
353 not encounter on a daily basis. For this reason, processing chimpanzee vocalizations and their context
354 could trigger enhanced activity in the anterior superior temporal cortex, especially when compared to
355 human voice³⁷.

356 Importantly, our data stress the importance of acoustic proximity between human and chimpanzee
357 vocalizations: activity in the anterior STG and more generally in the anterior TVA would in fact depend
358 on phylogenetic *and* acoustic proximity. If phylogenetic proximity was the only actor at play, bonobo
359 calls should also trigger activity in the TVA, since they are similarly close to humans as chimpanzees
360 as far as phylogeny is concerned. Concerning macaque calls, since they are both phylogenetically and
361 acoustically distant from humans, the absence of TVA activity specific to this species was expected.
362 This interpretation is strongly supported by the inclusion of acoustic Mahalanobis distance for each

363 species compared to humans as covariate of interest. Additionally, previous research showed a higher
364 pitch in young bonobo screams in comparison to chimpanzee and human baby cries³⁸, giving steam to
365 the crucial role of acoustic Mahalanobis distance in our results. Therefore, it seems reasonable to
366 hypothesize that TVA activity is not human-specific^{2,6} *per se* but that it would instead be sensitive to
367 the vocalizations of other primate species, provided that such vocalizations share sufficient acoustic (and
368 phylogenetic) proximity with the human vocal signal.

369 We already mentioned that the interaction between phylogeny and acoustic distance or proximity would
370 be at the origin of TVA enhancement for the processing of chimpanzee but not bonobo vocalizations.
371 However, the absence of similar results for bonobo calls also support the evolutionary divergence of this
372 peculiar species. In fact, according to the self-domestication hypothesis, bonobos would have evolved
373 differently compared to chimpanzees due to selection against aggression³⁹. Interestingly, differentiation
374 in the evolutionary pathway of bonobos has affected both their behavior²⁹ and morphology. For instance,
375 research has shown a shorter larynx in bonobos in comparison to chimpanzees resulting in a higher
376 fundamental frequency in their calls²⁷, contributing to their greater acoustic distance from human or
377 chimpanzee vocalizations. Putting into perspective the self-domestication hypothesis and our
378 neuroimaging data, we can suppose that the calls of our common ancestor together with the other great
379 apes 8 million year ago⁴⁰ would be close to the ones currently expressed by chimpanzees.

380 Taken together, our data allow us to draw the conclusion that both phylogenetic and acoustic proximity
381 of primate vocalizations seem necessary to trigger activity in the human temporal voice areas. For this
382 reason, anterior TVA activity was observed solely for the processing of chimpanzee but not bonobo or
383 macaque vocalizations. Contrary to what was reported in recent years, we claim that the human TVA
384 are also involved in the processing of heterospecific vocalizations, provided they share sufficient
385 phylogenetic and acoustic proximity. Finally, our findings support a critical evolutionary continuity
386 between the structure of human and chimpanzee vocalizations.

387

388 **Acknowledgements**

389 We thank the Swiss National Science foundation (SNSF) for supporting this interdisciplinary project
390 (grants CR13II_162720 / 1 to DG-TG), the National Centre of Competence in Research (NCCR)
391 (51NF40-104897 - DG) hosted by the Swiss Center for Affective Sciences, as well as the Fondation
392 Ernst et Lucie Schmidheiny supporting CD.

393

394 **References**

- 395 1 Ogawa, S., Lee, T.-M., Kay, A. R. & Tank, D. W. Brain magnetic resonance imaging with
396 contrast dependent on blood oxygenation. *proceedings of the National Academy of Sciences* **87**,
397 9868-9872 (1990).
- 398 2 Belin, P., Zatorre, R. J., Lafaille, P., Ahad, P. & Pike, B. Voice-selective areas in human auditory
399 cortex. *Nature* **403**, 309-312 (2000).
- 400 3 Kriegstein, K. V. & Giraud, A.-L. Distinct functional substrates along the right superior
401 temporal sulcus for the processing of voices. *Neuroimage* **22**, 948-955 (2004).
- 402 4 Pernet, C. R. *et al.* The human voice areas: Spatial organization and inter-individual variability
403 in temporal and extra-temporal cortices. *Neuroimage* **119**, 164-174 (2015).
- 404 5 Aglieri, V., Chaminade, T., Takerkart, S. & Belin, P. Functional connectivity within the voice
405 perception network and its behavioural relevance. *NeuroImage* **183**, 356-365 (2018).
- 406 6 Bestelmeyer, P. E., Belin, P. & Grosbras, M.-H. Right temporal TMS impairs voice detection.
407 *Current Biology* **21**, R838-R839 (2011).
- 408 7 Frühholz, S., Trost, W., Grandjean, D. & Belin, P. Neural oscillations in human auditory cortex
409 revealed by fast fMRI during auditory perception. *NeuroImage* **207**, 116401 (2020).
- 410 8 Latinus, M. & Belin, P. Human voice perception. *Current Biology* **21**, R143-R145 (2011).
- 411 9 Zäske, R., Hasan, B. A. S. & Belin, P. It doesn't matter what you say: FMRI correlates of voice
412 learning and recognition independent of speech content. *cortex* **94**, 100-112 (2017).

- 413 10 Lahnakoski, J. M. *et al.* Naturalistic fMRI mapping reveals superior temporal sulcus as the hub
414 for the distributed brain network for social perception. *Frontiers in human neuroscience* **6**, 233
415 (2012).
- 416 11 Ethofer, T. *et al.* Emotional voice areas: anatomic location, functional properties, and structural
417 connections revealed by combined fMRI/DTI. *Cerebral cortex* **22**, 191-200 (2012).
- 418 12 Witteman, J., Van Heuven, V. J. & Schiller, N. O. Hearing feelings: a quantitative meta-analysis
419 on the neuroimaging literature of emotional prosody perception. *Neuropsychologia* **50**, 2752-
420 2763 (2012).
- 421 13 Latinus, M., Crabbe, F. & Belin, P. Learning-induced changes in the cerebral processing of
422 voice identity. *Cerebral Cortex* **21**, 2820-2828 (2011).
- 423 14 Latinus, M., McAleer, P., Bestelmeyer, P. E. & Belin, P. Norm-based coding of voice identity
424 in human auditory cortex. *Current Biology* **23**, 1075-1080 (2013).
- 425 15 Charest, I., Pernet, C., Latinus, M., Crabbe, F. & Belin, P. Cerebral processing of voice gender
426 studied using a continuous carryover fMRI design. *Cerebral Cortex* **23**, 958-966 (2013).
- 427 16 Grossmann, T., Oberecker, R., Koch, S. P. & Friederici, A. D. The developmental origins of
428 voice processing in the human brain. *Neuron* **65**, 852-858 (2010).
- 429 17 Kisilevsky, B. S. *et al.* Effects of experience on fetal voice recognition. *Psychological science*
430 **14**, 220-224 (2003).
- 431 18 Hüppi, P. S. Cortical development in the fetus and the newborn: advanced MR techniques.
432 *Topics in Magnetic Resonance Imaging* **22**, 33-38 (2011).
- 433 19 Andics, A., Gácsi, M., Faragó, T., Kis, A. & Miklósi, Á. Voice-sensitive regions in the dog and
434 human brain are revealed by comparative fMRI. *Current Biology* **24**, 574-578 (2014).
- 435 20 Perrodin, C., Kayser, C., Logothetis, N. K. & Petkov, C. I. Voice cells in the primate temporal
436 lobe. *Current Biology* **21**, 1408-1415 (2011).
- 437 21 Petkov, C. I. *et al.* A voice region in the monkey brain. *Nature neuroscience* **11**, 367-374 (2008).
- 438 22 Fecteau, S., Armony, J. L., Joanette, Y. & Belin, P. Is voice processing species-specific in
439 human auditory cortex? An fMRI study. *Neuroimage* **23**, 840-848 (2004).

- 440 23 Belin, P. Voice processing in human and non-human primates. *Philosophical Transactions of*
441 *the Royal Society B: Biological Sciences* **361**, 2091-2107 (2006).
- 442 24 Belin, P. *et al.* Human cerebral response to animal affective vocalizations. *Proceedings of the*
443 *Royal Society B: Biological Sciences* **275**, 473-481 (2008).
- 444 25 Fritz, T. *et al.* Human behavioural discrimination of human, chimpanzee and macaque affective
445 vocalisations is reflected by the neural response in the superior temporal sulcus.
446 *Neuropsychologia* **111**, 145-150 (2018).
- 447 26 Linnankoski, I., Laakso, M., Aulanko, R. & Leinonen, L. Recognition of emotions in macaque
448 vocalizations by children and adults. *Language & Communication* **14**, 183-192 (1994).
- 449 27 Grawunder, S. *et al.* Higher fundamental frequency in bonobos is explained by larynx
450 morphology. *Current Biology* **28**, R1188-R1189 (2018).
- 451 28 Slocombe, K. E. & Zuberbühler, K. Chimpanzees modify recruitment screams as a function of
452 audience composition. *Proceedings of the National Academy of Sciences* **104**, 17228-17233
453 (2007).
- 454 29 Gruber, T. & Clay, Z. A comparison between bonobos and chimpanzees: A review and update.
455 *Evolutionary Anthropology: Issues, News, and Reviews* **25**, 239-252 (2016).
- 456 30 Belin, P., Fillion-Bilodeau, S. & Gosselin, F. The Montreal Affective Voices: a validated set of
457 nonverbal affect bursts for research on auditory affective processing. *Behavior research*
458 *methods* **40**, 531-539 (2008).
- 459 31 Ferdenzi, C. *et al.* Voice attractiveness: Influence of stimulus duration and type. *Behavior*
460 *research methods* **45**, 405-413 (2013).
- 461 32 Team, R. RStudio: integrated development for R. *RStudio, Inc., Boston, MA URL* [http://www.](http://www.rstudio.com)
462 *rstudio.com* **42**, 14 (2015).
- 463 33 Collins, D. L., Neelin, P., Peters, T. M. & Evans, A. C. Automatic 3D intersubject registration
464 of MR volumetric data in standardized Talairach space. *Journal of computer assisted*
465 *tomography* **18**, 192-205 (1994).
- 466 34 Ashburner, J. A fast diffeomorphic image registration algorithm. *Neuroimage* **38**, 95-113
467 (2007).

- 468 35 Mellem, M. S., Jasmin, K. M., Peng, C. & Martin, A. Sentence processing in anterior superior
469 temporal cortex shows a social-emotional bias. *Neuropsychologia* **89**, 217-224 (2016).
- 470 36 Simmons, W. K., Reddish, M., Bellgowan, P. S. & Martin, A. The selectivity and functional
471 connectivity of the anterior temporal lobes. *Cerebral Cortex* **20**, 813-825 (2010).
- 472 37 Zahn, R. *et al.* Social concepts are represented in the superior anterior temporal cortex.
473 *Proceedings of the National Academy of Sciences* **104**, 6430-6435 (2007).
- 474 38 Kelly, T. *et al.* Adult human perception of distress in the cries of bonobo, chimpanzee, and
475 human infants. *Biological Journal of the Linnean Society* **120**, 919-930 (2017).
- 476 39 Hare, B., Wobber, V. & Wrangham, R. The self-domestication hypothesis: evolution of bonobo
477 psychology is due to selection against aggression. *Animal Behaviour* **83**, 573-585 (2012).
- 478 40 Perelman, P. *et al.* A molecular phylogeny of living primates. *PLoS Genet* **7**, e1001342 (2011).
- 479
- 480

481 **Tables**

482 **Table 1:** Activations, cluster size and coordinates for each contrast of interest of model 1 (mean
 483 of vocalization fundamental frequency and energy as trial-level covariates of no-interest) in the
 484 temporal voice areas, wholebrain voxelwise $p < .05$ FDR corrected, $k > 10$.

485 MNI coordinates

486	Region label	Hemisphere	X	Y	Z	T value	Cluster size (voxels)
487	Human, Chimpanzee > Bonobo, Macaque						
488	Superior temporal gyrus mid	L	-58	-12	0	9.21	1664
489	<i>Superior temporal sulcus post</i>	L	-58	-42	-2	4.30	
490	<i>Superior temporal sulcus mid</i>	L	-54	-24	-6	3.92	
491							
492	Superior temporal gyrus mid	R	56	-8	2	8.86	3051
493	<i>Superior temporal gyrus ant</i>	R	60	-2	-8	7.22	
494	<i>Superior temporal sulcus mid</i>	R	56	-10	-14	6.37	
495	<i>Superior temporal sulcus post</i>	R	58	-40	-4	3.75	
496							
497							
498	Chimpanzee > Human, Bonobo, Macaque						
499	Superior temporal gyrus ant ¹	L	-52	-2	-12	4.84	91
500	<i>Superior temporal gyrus mid</i>	L	-50	-8	-12	4.12	
501							
502	Superior temporal gyrus ant ²	R	54	0	-12	3.63	18
503							
504	Human > Chimpanzee						
505	Supramarginal gyrus	R	56	-42	28	8.41	5941
506	<i>Superior temporal gyrus mid</i>	R	56	-10	2	8.09	
507	<i>Superior temporal gyrus post</i>	R	58	-46	14	6.76	
508	<i>Superior temporal sulcus post</i>	R	66	-32	0	6.50	
509	<i>Middle temporal gyrus mid</i>	R	68	-20	-10	5.52	
510	<i>Superior temporal sulcus ant</i>	R	66	-14	-6	5.18	
511	<i>Middle temporal gyrus ant</i>	R	60	4	-20	4.45	
512							
513	Supramarginal gyrus	L	-56	-40	34	7.51	6109
514	<i>Superior temporal gyrus mid</i>	L	-50	-18	4	7.37	
515	<i>Superior temporal gyrus post</i>	L	-56	-52	18	6.53	
516	<i>Middle temporal gyrus mid</i>	L	-64	-22	-10	6.10	
517	<i>Middle temporal gyrus post</i>	L	-62	-44	-6	4.87	
518	<i>Superior temporal sulcus mid</i>	L	-54	-32	-4	4.76	
519							
520							
521	Chimpanzee > Human						
522	Superior temporal gyrus ant ³	L	-50	-4	-12	3.36	74
523							

524 ant: anterior; mid: central part; post: posterior.

525 ¹Figure 2 cluster label: aSTG¹

526 ²Figure 2 cluster label: aSTG²

527 ³Figure 2 cluster label: aSTG³

528 **Table 2:** Activations, cluster size and coordinates for each contrast of interest of model 2 (inter-
 529 species vocalization acoustic distance as trial-level covariate of interest) in the temporal voice
 530 areas, wholebrain voxelwise $p < .05$ FDR corrected, $k > 10$.

531 MNI coordinates

532	Region label	Hemisphere	X	Y	Z	T value	Cluster size (voxels)
533	Human, Chimpanzee > Bonobo, Macaque						
534	Superior temporal gyrus mid	L	-58	-12	0	9.35	1112
535	<i>Superior temporal gyrus mid</i>	<i>L</i>	<i>-50</i>	<i>-18</i>	<i>4</i>	<i>7.27</i>	
536	<i>Superior temporal sulcus ant</i>	<i>L</i>	<i>-54</i>	<i>-2</i>	<i>-12</i>	<i>4.19</i>	
537							
538	Superior temporal gyrus mid	R	56	-8	2	8.52	1619
539	<i>Superior temporal gyrus ant</i>	<i>R</i>	<i>58</i>	<i>-2</i>	<i>-10</i>	<i>6.66</i>	
540	<i>Superior temporal sulcus mid</i>	<i>R</i>	<i>56</i>	<i>-10</i>	<i>-14</i>	<i>6.09</i>	
541	<i>Superior temporal sulcus post</i>	<i>R</i>	<i>58</i>	<i>-40</i>	<i>-4</i>	<i>3.75</i>	
542							
543	Superior temporal gyrus post	R	58	-46	14	5.32	641
544							
545							
546	Chimpanzee > Human, Bonobo, Macaque						
547	Superior temporal gyrus ant ⁴	L	-52	-2	-12	4.74	71
548							
549							
550	Human > Chimpanzee						
551	Supramarginal gyrus	R	56	-42	28	8.91	6411
552	<i>Superior temporal gyrus mid</i>	<i>R</i>	<i>56</i>	<i>-8</i>	<i>2</i>	<i>8.62</i>	
553	<i>Superior temporal gyrus post</i>	<i>R</i>	<i>58</i>	<i>-44</i>	<i>20</i>	<i>7.80</i>	
554	<i>Middle temporal gyrus mid</i>	<i>R</i>	<i>68</i>	<i>-20</i>	<i>-10</i>	<i>5.76</i>	
555	<i>Superior temporal sulcus ant</i>	<i>R</i>	<i>58</i>	<i>4</i>	<i>-20</i>	<i>5.00</i>	
556							
557	Supramarginal gyrus	L	-62	-48	24	7.66	6527
558	<i>Superior temporal gyrus mid</i>	<i>L</i>	<i>-50</i>	<i>-18</i>	<i>6</i>	<i>7.50</i>	
559	<i>Middle temporal gyrus mid</i>	<i>L</i>	<i>-58</i>	<i>-32</i>	<i>-6</i>	<i>5.24</i>	
560	<i>Superior temporal gyrus ant</i>	<i>L</i>	<i>-54</i>	<i>0</i>	<i>2</i>	<i>5.18</i>	
561							
562							
563	Chimpanzee > Human						
564	Inferior temporal gyrus ant	L	-36	-2	-36	5.54	1310
565	<i>Hippocampus</i>	<i>L</i>	<i>-40</i>	<i>-24</i>	<i>-18</i>	<i>4.50</i>	
566	<i>Superior temporal gyrus ant</i> ⁵	<i>L</i>	<i>-48</i>	<i>-6</i>	<i>-14</i>	<i>3.29</i>	
567							

568 ant: anterior; mid: central part; post: posterior.

569 ⁴Figure 3 cluster label: aSTG⁴

570 ⁵Figure 3 cluster label: aSTG⁵

571