



LJMU Research Online

Lameira, AR, Hardus, ME, Mielke, A, Wich, SA and Shumaker, RW

Vocal fold control beyond the species-specific repertoire in an orang-utan

<http://researchonline.ljmu.ac.uk/id/eprint/3838/>

Article

Citation (please note it is advisable to refer to the publisher's version if you intend to cite from this work)

**Lameira, AR, Hardus, ME, Mielke, A, Wich, SA and Shumaker, RW (2016)
Vocal fold control beyond the species-specific repertoire in an orang-utan.
Scientific Reports, 6. ISSN 2045-2322**

LJMU has developed **LJMU Research Online** for users to access the research output of the University more effectively. Copyright © and Moral Rights for the papers on this site are retained by the individual authors and/or other copyright owners. Users may download and/or print one copy of any article(s) in LJMU Research Online to facilitate their private study or for non-commercial research. You may not engage in further distribution of the material or use it for any profit-making activities or any commercial gain.

The version presented here may differ from the published version or from the version of the record. Please see the repository URL above for details on accessing the published version and note that access may require a subscription.

For more information please contact researchonline@ljmu.ac.uk

<http://researchonline.ljmu.ac.uk/>

1 **Vocal fold control beyond the species-specific repertoire in an orang-utan**

2
3
4 **Authors:** Adriano R. Lameira^{1,2*}, Madeleine E. Hardus³, Alexander Mielke⁴, Serge A. Wich^{5,6},
5 Robert W. Shumaker^{7,8,9}

6
7
8 **Affiliations:**

9 ¹Evolutionary Anthropology Research Group, Department of Anthropology, Durham University,
10 Dawson Building, South Road, Durham, DH1 3LE, UK

11 ²Pongo Foundation, Papenhoeflaan 91, 3421XN Oudewater, the Netherlands

12 ³Independent researcher

13 ⁴Department of Primatology, Max Planck Institute for Evolutionary Anthropology, Deutscher
14 Platz 6, 04103 Leipzig, Germany

15 ⁵Research Centre in Evolutionary Anthropology, and Palaeoecology, School of Natural Sciences
16 and Psychology, Liverpool John Moores University, Byrom Street, Liverpool L3 3AF, UK

17 ⁶Institute for Biodiversity and Ecosystem Dynamics, University of Amsterdam, Sciencepark 904,
18 Amsterdam 1098, the Netherlands

19 ⁷Indianapolis Zoo, Indianapolis, IN 46222, USA

20 ⁸Krasnow Institute for Advanced Studies, George Mason University, Fairfax, VA 22030, USA

21 ⁹Anthropology Department, Indiana University, 107 S. Indiana Ave., Bloomington, IN 47405,
22 USA

23
24 *Correspondence to: adriano.lameira@durham.ac.uk

25
26
27 **KEYWORDS:** Voice, vocal behavior, call repertoire, great ape vocalization, learned call,
28 vocal learning, social learning, speech evolution

29 **ABSTRACT**

30 Vocal fold control was critical to the evolution of spoken language, much as it today allows us to
31 learn vowel systems. It has, however, never been demonstrated directly in a non-human primate,
32 leading to the suggestion that it evolved in the human lineage after divergence from great apes.
33 Here, we provide the first evidence for real-time, dynamic and interactive vocal fold control in a
34 great ape during an imitation “do-as-I-do” game with a human demonstrator. Notably, the orang-
35 utan subject skilfully produced “wookies” – an idiosyncratic vocalization exhibiting a unique
36 spectral profile among the orang-utan vocal repertoire. The subject instantaneously matched
37 human-produced wookies as they were randomly modulated in pitch, adjusting his voice
38 frequency up or down when the human demonstrator did so, readily generating distinct low vs.
39 high frequency sub-variants. These sub-variants were significantly different from spontaneous
40 ones (not produced in matching trials). Results indicate a latent capacity for vocal fold exercise
41 in a great ape (*i*) in real-time, (*ii*) up and down the frequency spectrum, (*iii*) across a register
42 range beyond the species-repertoire and, (*iv*) in a co-operative turn-taking social setup. Such
43 ancestral capacity likely provided the neuro-behavioural basis of the more fine-tuned vocal fold
44 control that is a human hallmark.

45
46
47

48 **INTRODUCTION**

49 Spoken languages are learned anew with every human generation. Great apes, however, our
50 closest relatives, are traditionally thought to be incapable of vocal learning^{1,2,cf.3,4} – the capacity
51 to expand their vocal repertoire with new calls learned from others⁵. This apparent paradox has
52 led to the suggestion that human vocal capacities have no imitative precursor in nonhuman
53 species⁶. The evolution of speech – the predominant means of expression of human language⁵ –
54 is hence currently hotly debated, as evidence seemingly challenges the importance of shared
55 ancestry for the emergence of speech within the primate lineage, even though shared ancestry
56 represents one of the founding pillars of Darwin’s theory of natural selection⁶.

57 Historical “great ape language projects” have trained captive individuals in the attempt to
58 teach them new word-like utterances^{7,8}. Results were, however, virtually null^{1,6}. One major
59 limitation of these landmark studies was the fact that detailed descriptions of the great ape vocal
60 repertoire were, for the most part, unavailable at that time. Importantly, scientists had no
61 verifiable catalogue or database to compare and gauge exhibited vocal flexibility. Ultimately,
62 great apes’ vocal skills were directly compared with humans’, rather than objectively against
63 their own natural vocal preferences, predispositions, and limitations.

64 This critical drawback has been addressed recently: new databases on the natural vocal
65 behaviour of great apes have allowed recognizing vocal learning of new voiceless consonant-like
66 calls^{3,9}, notably requiring supralaryngeal control of the vocal tract. A modern-day and informed
67 approach to great ape vocal repertoire could, therefore also clarify whether (besides
68 supralaryngeal control) vocal learning can also involve vocal fold control. This capacity would
69 permit volitional voice modulation⁵, enabling individuals to expand their repertoire with new
70 voiced vowel-like calls. Together with consonants, vowels represented the building blocks for
71 spoken language. Being able to socially learn new voiceless *and* voiced calls would have, thus,
72 effectively set the evolution of an ancestral hominid articulatory system on a course towards a
73 vocal system fundamentally similar to modern speech. The evolutionary implications of the

74 presence of vocal fold control (or volitional voice modulation⁵) in great apes warrants, therefore,
75 revisiting the “unsuccessful” protocols of previous historical studies under a new lens.

76 Thus far, great apes have been shown to exercise vocal fold control in some degree in
77 “species-specific” voiced calls (or “vocalizations”), i.e. that are typically produced by the
78 species^{10–14}. Other studies have shown that a number of individual-specific and population-
79 specific voiced calls in great apes do not conform to genetic and ecological divergence^{9,15,16},
80 suggesting that vocal fold control may play indeed an active role in shaping the composition of
81 the voiced repertoire of great apes. Together, these data confirm that it is imperative for our
82 understanding on the evolution of spoken language to assess the extent to which human vocal
83 fold skills elaborated upon those present in great apes^{17,18}.

84 Here, we report a novel orang-utan vocalization, coined “wookie,” idiosyncratic to the
85 vocal repertoire of an adolescent captive male – named Rocky. Our working hypothesis posed
86 that the study subject produced wookies through volitional control over the vocal folds. If this
87 hypothesis was in fact correct, the two major predictions followed. First, vocal fold activity and
88 acoustic profile of the wookie should be clearly different from those of other orang-utan calls.
89 Second, the study subject should be able to adapt vocal fold action in response to random stimuli
90 under rigorous controlled experimental settings (e.g. to rule out arousal-based mechanisms). The
91 calls produced in this fashion should be perceptually distinct according with their respective
92 stimuli.

93 To test the first prediction and verify the novelty of wookies, we evaluated wookies’
94 acoustic profile in light of the known orang-utan call repertoire. Specifically, we measured and
95 assessed parameters describing vocal fold activity and supralaryngeal manoeuvring between
96 wookies and its most similar call-type in the orang-utan repertoire. To test the second prediction,
97 we brought the subject’s putative vocal fold control under scrutiny by presenting him with a
98 imitative “do-as-I-do” game paradigm^{19,20}. Under this paradigm, a human demonstrator produced
99 wookie-approximations with varying acoustic features as an implicit request towards the subject
100 to produce vocalizations of matching features. Subject’s vocal responses were recorded and
101 compared with the human models and between themselves. Our results show that a nonhuman
102 great ape can achieve levels of volitional voice control qualitatively comparable to those
103 manifested in humans; notably, real-time, dynamic and interactive vocal fold control beyond the
104 species-specific vocal repertoire.

107 **Methods**

108 *Orang-utan wookies and the species-specific repertoire*

110 *Data Collection*

111 To test the first prediction of this study and verify the idiosyncrasy of wookies and their novelty
112 among the known orang-utan repertoire, we recorded spontaneous wookies from Rocky
113 (studbook ID: 3331) during interactions with the human experimenter (MEH) between April and
114 May 2012 at the Indianapolis Zoo, where he is currently housed. We used a ZOOM H4Next
115 Handy recorder via the inbuilt mic standing on a miniature tripod at approximately ~0.5m
116 distance from the subject. Recordings were collected at a sampling rate of 24bit/48,000kHz and
117 saved in wav format. These settings obtained high quality audio recording and are standard for
118 the collection of orang-utan call behaviour in captivity and the wild. The original version of
119 wookies has been produced by Rocky for at least the last 6.5 years. It was apparent when the

120 experimenters first met Rocky when he was 3.5 years old. It is unclear how he originally learned
121 the vocalization and no recordings are available from earlier years. Wookies are produced by the
122 subject to gather attention from caretakers^{16,21}. Recordings from the known orang-utan call
123 repertoire available from previous work²² were used in order to draw a comparison with
124 wookies.

125 126 *Data analyses*

127 In order to verify the novelty of wookies in relation to the remaining orang-utan call repertoire,
128 we assessed the largest database ever assembled of orang-utan calls²², currently spanning more
129 than 12,000 observation hours across 9 wild and 6 captive populations, and comprising more
130 than 120 individuals. We compared wookies produced spontaneously (i.e. not given in response
131 to human wookie-versions) with the spectrally most similar vocalization known to be produced
132 by orang-utans – the grumph²². Grumphs were the only vocalization presently described in the
133 orang-utan repertoire to exhibit a complete overlap in frequency range with wookies (grumphs:
134 86 – 1723Hz, wookies: 99.6 – 1418Hz). Both calls were the only orang-utan vocalizations to fall
135 below 100Hz and simultaneously reach above 350Hz²² (Fig. 1). Wookies were produced with
136 ingressive air-flow, whereas grumphs were presumably produced with egressive air-flow (as
137 various other orang-utan calls)²². Nevertheless, we decided to conduct a comprehensive acoustic
138 comparison in order to verify, with confidence, wookies' idiosyncrasy and prevent claims of
139 novelty strictly based on one immeasurable articulatory feature (i.e. air-flow direction). For this
140 comparative analysis, grumphs were sampled from wild adolescent males of similar age as
141 Rocky in order to control for the largest number of potentially confounding factors as possible;
142 primarily, sex and body size variation. In order to control for potential geographic variation in
143 grumph acoustics, all wild adolescent males were sampled from the same population (i.e.
144 Ketambe Forest, Aceh, Sumatra, Indonesia).

145 To acoustically compare wookies with orang-utan grumphs, acoustic measures were
146 conducted with Praat, using “voice report” standard settings, except for voicing threshold in the
147 pitch settings, which was set to 0.15. Seven acoustic parameters describing vocal fold oscillation
148 were measured: duration, median pitch, mean pitch, pitch standard deviation, minimum pitch,
149 maximum pitch and pitch amplitude. Complementary, three acoustic parameters describing
150 supralaryngeal action were measured: first, second and third formant. Because these parameters
151 directly express the position of the tongue and jaw during vocal production, they were used to
152 assess whether wookies also involved different oral manoeuvres, besides different oscillation
153 patterns at the vocal folds.

154 Statistical analyses were conducted using nonparametric tests with IBM SPSS Statistics
155 21 (SPSS, Inc.). To compare the differences between wookies and grumphs, one would typically
156 use a Mann-Whitey U test for each parameter. However, because different individuals
157 contributed with several calls to our dataset, this condition violated the assumption of data
158 independence for conducting Mann-Whitney U tests. As such, we opted to conduct Kruskal
159 Wallis tests between individuals for each parameter, while correcting for multiple testing using
160 Bonferroni correction. We expected that Kruskal Wallis test results would show the following.
161 For each parameter, our study subject should be different from all other individuals, while all
162 other individuals should not be different between themselves, since wookies only derived from
163 our study subject whereas grumphs derived from all the remaining individuals. For these
164 analyses, we included our subject and the other adolescent males for whom a sample size larger

165 than one was available (i.e. 2 individuals with 24 and 12 calls). This operation resulted in the
166 exclusion of three adolescent males for which one grumph recording was available.

167
168

169 *Orang-utan vocal fold action in match trials*

170
171

171 *Data Collection*

172 To test the second prediction of this study, experimental testing was conducted with Rocky
173 during April and May 2012 at the Indianapolis Zoo. The zoo's committee provided ethical
174 approval and permission to conduct research, and the methods were carried out in accordance
175 with the approved guidelines. "Do-as-I-do" paradigm was selected for match trials because this
176 paradigm has been successfully used previously to invoke voluntary call responses in captive
177 orang-utans^{19,20}. Human demonstrator used protective gloves and a facial mask at all times and
178 interacted with Rocky always through enclosed mesh. Rocky was rewarded during trial sessions
179 with customary food snacks (i.e. raisins and dried plums) or drinks, prepared and provided by
180 full-time orang-utan caretakers at the zoo. Caretakers assured the items used differed in no
181 noticeable way in terms of the subject's food preferences and food rewards did not vary within
182 trial sessions.

183 Under the "do-as-I-do" test paradigm, the human demonstrator presented Rocky with
184 random sequences (Runs test, $Z = -4.751$, $p < 0.001$) of human wookie-versions varying in
185 frequency (Hz) – low vs. high wookies. 513 trials were presented (272 low, 241 high), divided
186 through 13 sessions (~49 trials/session, ~472 seconds/session) over the course of 5 days. The
187 subject typically responded to the model signal within approximately 500ms.

188 Trial sessions were recorded at ~0.5m distance from the subject with a ZOOM H4Next
189 Handy recorder via the inbuilt mic standing on a miniature tripod. Recordings were collected at a
190 sampling rate of 24bit/48,000kHz and saved in wav format. These settings obtained high quality
191 audio recordings. Rocky only joined trial sessions voluntarily and never refused to participate.
192 Rocky was never food deprived during trials sessions and trial sessions never interfered with
193 normal feeding times or working schedule at the orang-utan enclosure so as to prevent imposing
194 any stress. Rocky was tested when he and his cohort (four other orang-utans) were housed in
195 their individual quarters.

196 During trial sessions, only the first reply immediately after the human model was
197 considered for analyses, unless the human demonstrator verbally instructed (repeating the call
198 model or saying the name of the variant to be matched, "low" or "high") the focal to repeat, in
199 which case we considered the call produced after the last instruction provided by the human
200 demonstrator, or the last call produced by the focal before the human demonstrator verbally
201 closed the bout (e.g. by saying "yes" or "very good"). We did not consider calls when overlap
202 between human model and orang-utan match reply did not allow suitable extraction of acoustic
203 parameters from both calls (i.e. focal was too quick to reply).

204 We intentionally selected a human demonstrator with no previous voice training or music
205 experience. Because our main aim was fundamentally evolutionary, we deliberately avoided
206 using a demonstrator with vocal skills well beyond those potentially present in a human ancestor.
207 We mandated model calls to be as "raw" and naturally sounding as much as possible. No *a priori*
208 guidelines were given to the human demonstrator before match trials and no acoustic treatment
209 was given to her utterances. Moreover, we purposefully did not obstruct the human demonstrator
210 from deploying her natural behaviour during the interaction (e.g. occasional approximation to the

211 subject, occasional arm movement). Crucially, this decision allowed the demonstrator to keep the
212 subject engaged and cooperative during the tests. Nevertheless, we were adamant about
213 providing no training sessions, opportunities or time to the subject before the match trials, and
214 the subject was presented a human demonstrator with whom he was not familiar. These factors
215 confidently assured that our subject did not develop conditioned responses.

216

217 *Data analyses*

218 In order to compare the acoustic profile and general vocal fold oscillation between human- and
219 orang-utan-produced wookies, we selected and analyzed call maximum frequency (Hz). This
220 parameter was also used to compare the subject's wookie sub-variants between each other
221 (spontaneous, high and low). Maximum frequency is the frequency at which maximum energy
222 (dB) occurs within a call. For this reason, maximum frequency contributes disproportionately to
223 pitch and, in the case of wookies, it represented one of the best proxies available for pitch
224 (Spearman test between maximum frequency and mean pitch, $r = 0.341$, $N_{\text{spontaneous wookies}} = 124$, p
225 > 0.001). Moreover, maximum frequency was equal to the fundamental frequency (F_0) 93.4% of
226 500 measured cases. Therefore, maximum frequency provided one of the most reliable measures
227 of the oscillation rate of the vocal folds and its perception. In order to assess the subject's level of
228 accuracy during the task, we also conducted the same test but analysing low and high wookies
229 separately.

230 Besides maximum frequency, we measured duration and maximum power (dB) within
231 each call. Because all recordings were conducted at a constant distance from the study subject,
232 maximum power could be used as a proxy of glottal air pressure during call production. This
233 measure allowed us, thus, to monitor the contribution of abdominal action (generating air current
234 within the vocal tract) during the production of wookies exhibiting different maximum
235 frequencies.

236 Maximum frequency, duration and maximum power were extracted from recordings
237 using Raven Pro software package (version 1.5, Ithaca, NY: The Cornell Lab of Ornithology)
238 and Hann type spectrogram grip spacing at 2.93Hz. The use of other important parameters
239 characterizing vocal fold oscillation (e.g. harmonics-to-noise ratio) was hampered because these
240 parameters are particularly susceptible to recording settings²⁰.

241 Nonparametric statistical analyses were conducted using IBM SPSS Statistics 21 (SPSS,
242 Inc.). Spearman binomial correlation test was used to assess a potential effect of human model
243 calls on the responses produced by the study subject. Wilcoxon signed ranks test was used to
244 identify potential differences between wookie subvariants produced by the study subject.
245 Discriminant function analyses were used to assess whether wookie subvariants produced by the
246 study subject could be distinguished perceptually. Discriminant function analyses were
247 conducted both by setting prior probabilities (i.e. chance probability of correct assignment) equal
248 between all groups and by computing prior probabilities based on group size. Because our data
249 set for these analyses derived from the same individual, this did not require conducting a
250 permuted discriminant function analysis. A permuted analysis would have otherwise allowed
251 controlling for a possible confounding variable. For instance, if several individuals had
252 contributed wookie subvariants, the permuted analysis would have allowed controlling for
253 individual variation while assessing the capacity to correctly distinguish wookie subvariants.

254 Because receivers sense acoustic signals holistically instead of attending to one or few
255 acoustic parameters separately²³, we tested whether low and high wookies produced by Rocky
256 were overall perceptually distinct from each other by using automated classification algorithms,

257 combined with artificial neural networks (ANN) and mel frequency cepstral coefficients
258 (MFCC)²⁴, a classification method that scans and analyses signals based on their general acoustic
259 profile. These analyses allowed assessing the differences between wookie sub-variants while
260 taking in consideration their complete acoustic profile simultaneously, other than one acoustic
261 parameter at a time. For both feature extraction and network analyses, Matlab R2012b (The
262 MathWorks, Inc., Natick, MS, U.S.A.) was used. The MFCCs in this study were computed using
263 the ‘melcepst’-routine available in the toolbox Voicebox. We optimized both MFCC and ANN
264 according to published guidelines²⁴. To acquire a MFCC, each call was sliced into seven frames
265 using a Hamming window, two-thirds frame overlap and 16 mel-spaced filters²⁴. We used 10
266 hidden layer neurons and 100 iterations to obtain an optimal ANN²⁴. To increase the reliability of
267 the results, every call was tested against seven neural networks, and the condition proposed by
268 the majority of the networks was considered final²⁴. Calls were tested using a leave-one-out
269 procedure²⁴.

270 Lastly, we conducted Spearman binomial correlation tests between maximum frequency,
271 duration and maximum power of the subject’s wookies in order to investigate general production
272 dynamics. With these analyses, we were particularly interested in examining to what extent low
273 and high wookies could have been produced strictly by means of changes in glottal air pressure
274 generated by abdominal control (other than by vocal fold control).
275

276

277

277 **Results**

278 *Orang-utan wookies and the species-specific repertoire*

279

280 A number of acoustic parameters was measured characterizing the oscillation pattern of the vocal
281 folds with high accuracy. Significant differences were detected within our sample comprised by
282 our study subject ($n_{\text{wookies}} = 124$) and other adolescent males ($n_{\text{grumphs}} = 36$, $n_{\text{individuals}} = 2$,
283 $n_{\text{grumphs/ind}} = 24$, 12) with regards to duration (Kruskal Wallis test, $df = 2$, $X^2 = 62.080$, $p <$
284 0.001), median pitch ($X^2 = 29.404$, $p < 0.001$), mean pitch ($X^2 = 56.899$, $p < 0.001$), pitch
285 standard deviation ($X^2 = 20.592$, $p < 0.001$), minimum pitch ($X^2 = 26.508$, $p < 0.001$), maximum
286 pitch ($X^2 = 62.201$, $p < 0.001$), and pitch amplitude ($X^2 = 20.540$, $p < 0.001$). Post hoc pairwise
287 comparisons between individuals revealed that, for all parameters, our study subject was (with
288 the exception of two out of 14 pairwise comparisons) always significantly different from the
289 remaining individuals (duration: $p < 0.001$ and $p < 0.001$; median pitch: $p < 0.001$ and $p = 0.002$;
290 mean pitch: $p < 0.001$ and $p < 0.001$; pitch standard deviation: $p < 0.001$ and $p = 0.054$; minimum
291 pitch: $p < 0.001$ and $p = 0.004$; maximum pitch: $p < 0.001$ and $p < 0.001$; pitch amplitude: $p <$
292 0.001 and $p = 0.133$). At the same time, the remaining individuals showed always no significant
293 differences between each other (duration: $p = 0.539$; median pitch = 1.000; mean pitch: 1.000;
294 pitch standard deviation: 0.124; minimum pitch: $p = 1.000$; maximum pitch: $p = 0.884$; pitch
295 amplitude: $p = 0.051$). Overall, wookies were significantly longer and exhibited lower pitch
296 values than grumphs (Fig. 2 and Table S1 in Supplementary material).

297

298 In addition, we compared in the same way the first, second, and third formant (F1, F2,
299 F3) between our subject and other adolescent males to assess differences in supralaryngeal
300 maneuvering during vocal production. Significant differences within our sample of individuals
301 were found for F1 (Kruskal Wallis test, $df = 2$, $X^2 = 11.964$, $p < 0.001$), but neither for F2 nor F3
302 ($X^2 = 0.470$, $p = 0.791$; $X^2 = 2.307$, $p = 0.316$, respectively). Post hoc pairwise comparisons
303 between individuals revealed that our study subject was significantly different from the

303 remaining individuals for F1 ($p = 0.037$ and $p = 0.019$), but the remaining individuals were not
304 different between each other ($p = 1.000$). Overall, tongue body (F2) and tip (F3) positioning was
305 relatively similar between the two calls types but wookies (presenting a higher F1) involved a
306 wider opening of the mouth during call production than that required for grumph production²⁵.

307 These analyses encompassed multiple testing. Correction of significance level was
308 therefore required. Even though Bonferroni correction represents an over-conservative method
309 ($0.05/10 = 0.005$)²⁶, this adjustment did not affect our results on vocal focal action, since all our
310 tests yielding significant differences provided p values smaller than 0.001. The only significant
311 difference dissolved by Bonferroni correction concerned F1 between our subject and the
312 remaining adolescent males. Essentially, this result indicates that differences in vocal fold action
313 provided the most reliable and consistent way of distinguishing wookies versus grumphs,
314 whereas differences in supralaryngeal action were less secure.

315
316

317 *Orang-utan vocal fold action in match trials*

318

319 Maximum call frequency (Hz) of human-wookies and orang-utan-wookies showed a significant
320 positive correlation (Spearman, $r = 0.574$, $N = 513$, $p < 0.001$) (Fig. 3). When testing for low and
321 high wookies separately, a significant correlation between human-wookies and orang-utan-
322 wookies was also reached for high wookies (Spearman, $r = 0.141$, $p = 0.029$).

323 Maximum frequency differences between low and high wookies produced by Rocky
324 significantly differed from each other (Wilcoxon Signed Ranks test, $Z = -10.409$, $p < 0.001$),
325 with low and high wookies exhibiting a median frequency of 126Hz and 161.1Hz, respectively, a
326 difference nearly equivalent to a four-note interval on a standard musical octave (B–E) (Fig. 4,
327 Table S2 in Supplementary material). Low and high frequency wookies produced by the subject
328 also significantly differed in maximum frequency from spontaneous wookies ($n = 124$) (low vs.
329 spontaneous wookies: Wilcoxon Signed Ranks test, $Z = -4.405$, $p < 0.001$; high vs. spontaneous
330 wookies: $Z = -3.101$, $p = 0.002$), with spontaneous wookies exhibiting an intermediary median
331 frequency of 134.8Hz (Fig. 4, Table S2 in Supplementary material). Bonferroni correction of our
332 significance value ($0.05:3=0.0167$) did not affect our results.

333 Discriminant function analysis, based on maximal frequency alone, attained 50.1% of
334 corrected assignments between low, high, and spontaneous wookies (49.6% using leave-one-out
335 procedure), performing significantly above chance (Wilks' Lambda Chi-square, $X^2 = 47.128$, df
336 $= 2$, $p < 0.001$; Binomial test, chance probability = 0.333, $p < 0.001$). Correct assignments
337 decreased slightly to 48.0% (48.0% using leave-one-out procedure), but remained well above
338 chance, when computing chance levels according to group size (low wookies: 42.6%; high:
339 38.0%; spontaneous: 19.4%). Percentage of correct assignments to the three sub-variants
340 increased to 69.5% (69.3% using leave-one-out procedure) when supplementing duration and
341 maximum power to the analyses (Fig. 5). In these conditions, maximum frequency (together with
342 maximum power) held the largest absolute correlation with the first discriminant function, which
343 explained 79.4% of the total observed variation. Percentage of correct assignments increased to
344 72.5% (72.1% using leave-one-out procedure) when computing chance levels according to group
345 size.

346 These results were corroborated when ascribing the classification of low and high
347 wookies to an automated process scanning the vocalizations' general acoustic profile. The mean
348 (25%; 75% percentiles) percentage of correct assignments per session was 87.82% (84.82%;

349 95.12%). Altogether, these results confirmed that low and high wookies were perceptually
350 distinct, and thus, that they could potentially encode biologically pertinent differences.

351 Maximum frequency, duration, and maximum power of Rocky's wookies showed
352 significant positive correlations (Spearman, $n = 639$, maximum frequency x duration: $r = 0.116$,
353 $p = 0.003$; maximum frequency x maximum power: $r = 0.134$, $p = 0.001$). Bonferroni correction
354 of our significance value ($0.05:3=0.025$) did not affect these results. Graphical examination of
355 Rocky's phonetogram (Fig. 6) showed that at any given sound pressure level Rocky was capable
356 of generating a frequency range wider than 100Hz. This effect was particularly visible in high
357 frequency wookies, with Rocky producing most of the calls around 160 dB but spanning well
358 above 200Hz. At the same time, Rocky was able to produce any specific frequency tone across a
359 range of more than 20dB.

360

361

362 **DISCUSSION**

363 *Orang-utan wookies and the species-specific repertoire*

364

365 Our results validated our first prediction, showing that wookies represent an acoustically distinct
366 voiced call within the orang-utan call repertoire. Wookies exhibit features of air-flow, vocal fold
367 action and jaw position unique to Rocky and described here for the first time in the *Pongo* genus.
368 These results confirm the capacity of orang-utans to learn and acquire new calls into their
369 individual repertoires, both in the form of voiceless consonant-like calls^{3,4,9,15,20} and voiced
370 vowel-like calls^{9,15}.

371 Because our analyses focused on an idiosyncratic vocalization, there were inevitable
372 limitations in our statistical analyses. However, after conducting procedures that contemplated
373 the potential of confounding effects, results were always highly significant. Together with the
374 observation that wookies and their closest counterpart in the known orang-utan repertoire exhibit
375 opposite air-flow directions, our analyses allow determining with confidence that wookies are a
376 novel vocalization based on parameters describing vocal fold oscillation.

377 Despite an N of 1, our study allows reevaluating current assumptions on great ape vocal
378 capacities as well as reformulating some of the basic premises of a general theory of spoken
379 language evolution. By demonstrating vocal learning beyond the species-specific repertoire in a
380 great ape, our results unveil a fundamental parallel with human spoken languages. Namely, the
381 two vocal systems, separated by approximately 10mya²⁷, can be assumed homologous regarding
382 open-endedness and the voiced/voiceless nature of their two building blocks.

383

384 *Orang-utan vocal fold action in match trials*

385

386 Our results validated our second prediction, indicating that the subject modulated vocal fold
387 oscillation according to the model-calls provided by the human demonstrator under controlled
388 settings. The subject adjusted his voice frequency up or down when the human model did so. For
389 this, the subject produced significantly different vocal sub-variants that stood in average outside
390 his normal spectrum of wookie vocalizations. Human demonstrations, thus, effectively guided
391 the subject's vocal output. Moreover, results suggest that the subject attended, was sensitive to
392 and coordinated his vocal responses according to the spectral dispersion of sub-variants beyond
393 the low/high dichotomy and down to a scale of tens of Hz. Manual and automated procedures

394 demonstrated that his low vs. high wookies exhibited clear perceptible differences, allowing
395 discerning the two with high accuracy.

396 Correlation between wookies' acoustic parameters produced by the subject indicated that
397 high frequency wookies were simultaneously louder and longer. That is, high wookies were
398 partly underlined by higher airflow pressure exciting the vocal folds. Accordingly, the
399 production of wookie sub-variants by our subject resulted from the synchronized exercise of the
400 vocal folds and the abdominal musculature generating glottal airflow (e.g. diaphragm). The
401 action of abdominal muscles may have partially alleviated the degree of vocal fold control
402 required to obtain the observed dynamic production across frequencies during match-trials. This
403 positive acoustic interdependence between frequency and glottal air pressure also characterizes,
404 however, overall human vocal production, including people with musical training²⁸, and is a
405 phenomenon predicted to be common among animal vocal communication systems.
406 Nevertheless, different wookies produced by Rocky with equal frequencies exhibited wide
407 differences in acoustic power, and vice versa. These observations would have been theoretically
408 impossible if Rocky had not exercised some degree of direct control over vocal fold oscillation,
409 and instead had only resorted to abdominal action to produce modulations at the level of vocal
410 fold oscillation. The subject's phonetogram attests that vocal fold control was effective and
411 moderately autonomous from abdominal control.

412 Our match trials were conducted in constant settings in one-to-one interactions between
413 the subject and the human demonstrator. Food rewards were part of the subject's daily diet and
414 were always kept constant within sessions. Accordingly, we can ascertain that the subject's
415 performance and vocal output was not affected by the influence of other orang-utans, physical
416 surroundings or food-driven arousal. Thus, the different wookie sub-variants produced by the
417 subject were unrelated to specific changes in context and can be considered to have conveyed no
418 change in function or informational content.

419 Any possible biasing effects deriving from the natural behaviour of the human
420 demonstrator can also be excluded in light of our results. For example, the demonstrator
421 occasionally approached the subject and moved her arm during low and high vocal models,
422 respectively. The subject could have hypothetically used these supplementary visual cues to
423 know which response was "correct" (instead of directly mimicking the demonstrator's voice
424 modulation), or these cues could have somehow affected the subject's arousal in a coherent way
425 with correct responses ("clever Hans effect"). Such interpretations can, however, be dismissed at
426 least for three reasons. First, the subject neither necessarily gazed directly at the human
427 demonstrator to produce a correct response, nor did human supplementary cues ensured subject's
428 correct responses (see supplementary video). Second, the subject never raised his arm in
429 response to the similar movement by the demonstrator. Thus, he attended to human *acoustic*
430 signals, not other cues. Third, in case the subject's arousal had been affected, one would expect
431 an increase in subject's arousal when interacting with a human. However, subject's low calls
432 were lower than his spontaneous calls. Overall, visual cues or arousal offer no parsimonious
433 explanation for our results.

434

435 *Implications for spoken language evolution*

436

437 Our findings imply the functional presence of direct pathways between the primary motor cortex
438 and the nucleus ambiguus (site of the laryngeal motor-neurons in medulla oblongata) in the ape
439 brain, as observed in an chimpanzee by Kuypers²⁹, allowing some degree of vocal fold control

440 autonomous from context and individual's affective state. Specifically, our analyses indicate that
441 vocal fold control pathways and respective firing in the ape brain integrate with pathways
442 innervating other musculatures engaged in vocal production (namely, abdominal muscles).
443 Several motor maneuvers are brought together synergistically to generate a particular acoustic
444 output.

445 Contrarily to the notion that spoken language emerged abruptly sometime along the
446 genus *Homo*³⁰, our findings amplify the spoken language evolution timeline at least five-fold
447 (assuming speech evolution onset in *Homo* paleodemes, from 2 mya onwards) and up to 50-fold
448 (assuming speech emergence in *H. sapiens*, 200kya)³¹. Full articulatory range and excellent
449 vocal control as observed today in humans may be relatively recent within the human lineage.
450 However, the presence of learned consonant- and vowel-like calls, potentially as far as 10 mya
451 within our lineage, allows considering gradual forces and progression in stages towards full-
452 blown language. This intriguing possibility raises caution in the inference of the vocal capacities
453 of extinct hominoidae from the fossil record without complementary assessment of the vocal
454 capacities of extant great apes.

455 Vocal control over laryngeal and supralaryngeal structures at the root of a 10 mya
456 timeline for spoken language evolution suggests that vocal evolution could have co-evolved with
457 cognition within the human lineage. Whereas monkey cognitive skills have been hitherto
458 assumed to surpass their vocal counterparts^{32,33}, the possibility that the two skillsets originally
459 exhibited even levels of sophistication in an ancestral hominid opens new considerations on
460 speech/language evolution. In this scenario, vocal control would have allowed the immediate
461 manifestation, or “verbalization,” of advanced cognition. Forces propelling cognitive processes
462 would have then compelled vocal progress by association, and vice versa. For instance, with the
463 emergence of theory of mind, individuals would have been able to exploit deceptive calls³⁴⁻³⁶,
464 effectively launching new communicative and social dynamics within a population where
465 acoustic deception was previously absent. If vocal and cognitive sophistication developed hand-
466 in-hand over the course of human evolution during the last 10 mya, then, the processes of speech
467 evolution and language evolution could be considered to have been one and the same. This
468 “speech-language co-evolution” hypothesis will require future examination but it may perhaps
469 expedite, for example, our understanding on the evolution of syntax and semantics. Because
470 vocal control allows a functional divide between a signal (signifier³⁵) and its functional use or
471 meaning (signified³⁵) – as suggested in our results – there would be few articulatory limitations
472 for the assemblage of vocal sequences and the attribution of their respective informational
473 content, so long as we had the required cognitive machinery to do so. In other words, in a
474 condition where vocal evolution kept close pace with cognition, a human ancestor (regardless
475 his/her position along human evolution timeline) would rarely have cognitive computations for
476 which there were no matching vocal counterparts.

477

478

479 CONCLUSION

480

481 We demonstrate real-time, dynamic and interactive vocal fold control beyond the vocal range of
482 the orang-utan genus. This study offers a new category of vocal learning in great apes, in
483 addition to previous cases describing gradual (over the course of months) and directional shift
484 (exclusively downwards in frequency) of a species-specific vocalization¹⁰. Orang-utans (and
485 possibly other great apes) possess a latent capacity for controlled deployment of vocal fold

486 oscillation, allowing the volitional production of novel vowel-like calls. Theoretically, together
487 with the capacity of great apes to socially learn voiceless consonant-like calls^{3,20}, this proto-
488 linguistic capacity constituted a crucial prerequisite for the onset of spoken language evolution.
489

490

491 References

492

- 493 1. Janik, V. M. & Slater, P. J. B. in *Advances in the Study of Behavior* (eds. Snowdon, C. T.
494 & Manfred, M.) **26**, 59–99; DOI:10.1016/s0065-3454(08)60377-0 (Academic Press,
495 1997).
- 496 2. Fitch, W., Huber, L. & Bugnyar, T. Social cognition and the evolution of language:
497 constructing cognitive phylogenies. *Neuron* **65**, 795–814;
498 DOI:10.1016/j.neuron.2010.03.011 (2010).
- 499 3. Lameira, A. R., Maddieson, I. & Zuberbühler, K. Primate feedstock for the evolution of
500 consonants. *Trends Cogn. Sci.* **18**, 60–62; DOI:10.1016/j.tics.2013.10.013 (2014).
- 501 4. Lameira, A. R. The forgotten role of consonant-like calls in theories of speech evolution.
502 *Behav. Brain Sci.* **37**, 559–560; DOI:10.1017/S0140525X1300407X (2014).
- 503 5. Pisanski, K., Cartei, V., McGettigan, C., Raine, J. & Reby, D. Voice Modulation: A
504 Window into the Origins of Human Vocal Control? *Trends Cogn. Sci.* **20**, 304–318;
505 DOI:10.1016/j.tics.2016.01.002 (2016).
- 506 6. Bolhuis, J. J. & Wynne, C. D. L. Can evolution explain how minds work? *Nature* **458**,
507 832–833; DOI:10.1038/458832a (2009).
- 508 7. Furness, W. H. Observations on the mentality of chimpanzees and orang-utans. *Proc. Am.*
509 *Philos. Soc.* **55**, 281–290; (1916).
- 510 8. Hayes, K. J. & Hayes, C. The intellectual development of a home-raised chimpanzee.
511 *Proc. Am. Philos. Soc.* **95**, 105–109; (1951).
- 512 9. Lameira, A. R. *et al.* Speech-Like Rhythm in a Voiced and Voiceless Orangutan Call.
513 *PLoS One* **10**, e116136; DOI:10.1371/journal.pone.0116136 (2015).
- 514 10. Watson, S. K. *et al.* Vocal Learning in the Functionally Referential Food Grunts of
515 Chimpanzees. *Curr. Biol.* **25**, 495–499; DOI:10.1016/j.cub.2014.12.032 (2015).
- 516 11. Crockford, C., Herbinger, I., Vigilant, L. & Boesch, C. Wild Chimpanzees Produce
517 Group-Specific Calls: a Case for Vocal Learning? *Ethology* **110**, 221–243;
518 DOI:10.1111/j.1439-0310.2004.00968.x (2004).
- 519 12. Marshall, A. J., Wrangham, R. W. & Arcadi, A. C. Does learning affect the structure of
520 vocalizations in chimpanzees? *Anim. Behav.* **58**, 825–830; DOI:10.1006/anbe.1999.1219
521 (1999).
- 522 13. Mitani, J. C. & Gros-Louis, J. Chorusing and Call Convergence in Chimpanzees: Tests of
523 Three Hypotheses. *Behaviour* **135**, 1041–1064; DOI:10.1163/156853998792913483
524 (1998).
- 525 14. Slocombe, K. E. & Zuberbühler, K. Chimpanzees modify recruitment screams as a
526 function of audience composition. *Proc. Natl. Acad. Sci. U. S. A.* **104**, 17228–17233;
527 DOI:10.1073/pnas.0706741104 (2007).
- 528 15. Wich, S. A. *et al.* Call cultures in orang-utans? *PLoS One* **7**, e36180;
529 DOI:10.1371/journal.pone.0036180 (2012).
- 530 16. Hopkins, W. D., Taglialatela, J. P. & Leavens, D. A. Chimpanzees differentially produce
531 novel vocalizations to capture the attention of a human. *Anim. Behav.* **73**, 281–286;

- 532 DOI:10.1016/j.anbehav.2006.08.004 (2007).
- 533 17. Taglialatela, J. P., Russell, J. L., Schaeffer, J. A. & Hopkins, W. D. Communicative
534 signaling activates 'Broca's' homolog in chimpanzees. *Curr. Biol.* **18**, 343–8;
535 DOI:10.1016/j.cub.2008.01.049 (2008).
- 536 18. Taglialatela, J. P., Russell, J. L., Schaeffer, J. A. & Hopkins, W. D. Chimpanzee Vocal
537 Signaling Points to a Multimodal Origin of Human Language. *PLoS One* **6**, e18852;
538 DOI:10.1371/journal.pone.0018852 (2011).
- 539 19. Wich, S. *et al.* A case of spontaneous acquisition of a human sound by an orangutan.
540 *Primates* **50**, 56–64; DOI:10.1007/s10329-008-0117-y (2009).
- 541 20. Lameira, A. R. *et al.* Orangutan (*Pongo* spp.) whistling and implications for the
542 emergence of an open-ended call repertoire: A replication and extension. *J. Acoust. Soc.*
543 *Am.* **134**, 1–11; DOI:10.1121/1.4817929 (2013).
- 544 21. Poss, S. R., Kuhar, C., Stoinski, T. S. & Hopkins, W. D. Differential use of attentional and
545 visual communicative signaling by orangutans (*Pongo pygmaeus*) and gorillas (*Gorilla*
546 *gorilla*) in response to the attentional status of a human. *Am. J. Primatol.* **68**, 978–992;
547 DOI:10.1002/ajp.20304 (2006).
- 548 22. Hardus, M. E. *et al.* in *Orangutans* (eds. S. Wich, T. Mitra Setia, S.S. Utami & Schaik, C.
549 P.) 49–60; (Oxford University Press, 2009).
- 550 23. Belin, P. Voice processing in human and non-human primates. *Philos. Trans. R. Soc.*
551 *Lond. B. Biol. Sci.* **361**, 2091–2107; DOI:10.1098/rstb.2006.1933 (2006).
- 552 24. Mielke, A. & Zuberbühler, K. A method for automated individual, species and call type
553 recognition in free-ranging animals. *Anim. Behav.*
554 DOI:http://dx.doi.org/10.1016/j.anbehav.2013.04.017 (2013).
- 555 25. Ladefoged, P., Harshman, R., Goldstein, L. & Rice, L. Generating vocal tract shapes from
556 formant frequencies. *J. Acoust. Soc. Am.* **64**, 1027; DOI:10.1121/1.382086 (1978).
- 557 26. Glickman, M. E., Rao, S. R. & Schultz, M. R. False discovery rate control is a
558 recommended alternative to Bonferroni-type adjustments in health studies. *J. Clin.*
559 *Epidemiol.* **67**, 850–857; DOI:10.1016/j.jclinepi.2014.03.012 (2014).
- 560 27. Hobolth, A., Dutheil, J. Y., Hawks, J., Schierup, M. H. & Mailund, T. Incomplete lineage
561 sorting patterns among human, chimpanzee, and orangutan suggest recent orangutan
562 speciation and widespread selection. *Genome Res.* **21**, 349–356;
563 DOI:10.1101/gr.114751.110 (2011).
- 564 28. Pabon, P., Ternström, S. & Lamarche, A. Fourier descriptor analysis and unification of
565 voice range profile contours: method and applications. *J. Speech. Lang. Hear. Res.* **54**,
566 755–76; DOI:10.1044/1092-4388(2010/08-0222) (2011).
- 567 29. Kuypers, M. G. J. M. Some projections from the peri-central cortex to the pons and lower
568 brain stem in monkeys and chimpanzee. *J. Comp. Neurol.* 211–255;
569 DOI:10.1002/cne.901100205 (1958).
- 570 30. Maclarnon, A. & Hewitt, G. Increased breathing control: Another factor in the evolution
571 of human language. *Evol. Anthropol. Issues, News, Rev.* **13**, 181–197;
572 DOI:10.1002/evan.20032 (2004).
- 573 31. Shea, J. J. *Homo sapiens* Is as *Homo sapiens* Was. *Curr. Anthropol.* **52**, 1–35;
574 DOI:10.1086/658067 (2011).
- 575 32. Seyfarth, R. M. & Cheney, D. L. Production, usage, and comprehension in animal
576 vocalizations. *Brain Lang.* **115**, 92–100; DOI:10.1016/j.bandl.2009.10.003 (2010).
- 577 33. Seyfarth, R. M., Cheney, D. L. & Bergman, T. J. Primate social cognition and the origins

- 578 of language. *Trends Cogn. Sci.* **9**, 264–266; DOI:10.1016/j.tics.2005.04.001 (2005).
579 34. Hardus, M. E., Lameira, A. R., van Schaik, C. P. & Wich, S. A. Tool use in wild orang-
580 utans modifies sound production: a functionally deceptive innovation? *Proc. R. Soc. B*
581 *Biol. Sci.* **276**, 3689–3694; DOI:10.1098/rspb.2009.1027 (2009).
582 35. Lameira, A. R. *et al.* Population-specific use of the same tool-assisted alarm call between
583 two wild orangutan populations (*Pongo pygmaeus wurmbii*) indicates functional
584 arbitrariness. *PLoS One* **8**, e69749; DOI:10.1371/journal.pone.0069749 (2013).
585 36. de Boer, B., Wich, S. A., Hardus, M. E. & Lameira, A. R. Acoustic models of orangutan
586 hand-assisted alarm calls. *J. Exp. Biol.* **218**, 907–914; DOI:10.1242/jeb.110577 (2015).
587
588
589
590

591 **Acknowledgments**

592 ARL was supported by a European Union COFUND/Durham Junior Research Fellowship. We
593 thank the Indianapolis Zoo for permission to conduct this study and logistical support. We thank
594 Tecumseh Fitch, Robert Barton and Roger Mundry for helpful comments on earlier versions of
595 the manuscript. The authors confirm that all data underlying the findings are fully available
596 without restriction. The majority of relevant data are within the paper. Remaining data is
597 available upon request to the corresponding author.
598

599 **Competing interests**

600 None of the authors have any competing interests.
601

602 **Author contribution statement**

603 ARL, MEH, SW and RS conceived the study and methodological protocol. ARL and MEH
604 conducted the experiments. ARL and AM conducted data analyses. ARL, MEH, SW and RS
605 wrote the manuscript.
606
607
608

609 **Figure legends**

610

611 Fig. 1. Spectrographic representation of two orang-utan grumphs followed by two wookies

612

613 Fig. 2. Boxplot per acoustic parameter of Rocky (producing wookies) and other adolescent males
614 (producing grumphs) (middle line represents the median, the box represents the interquartile
615 range (IQ) containing the middle 50% of the data, and the whiskers represent 1.5 times the IQ).

616

617 Fig. 3. Maximum frequency of human wookie demonstrations against maximum frequency of
618 Rocky's match wookies (linear fit line with intercept suppressed).

619

620 Fig. 4. Boxplot of the maximum frequency of low, spontaneous, and high wookie by Rocky
621 (middle line represents the median, the box represents the interquartile range (IQ) containing the
622 middle 50% of the data, and the whiskers represent 1.5 times the IQ).

623

624 Fig. 5. Graphic representation of first and second canonical discriminant functions, displaying
625 distribution and group centroids of Rocky's low frequency (1), high frequency (2), and
626 spontaneous wookies (3).

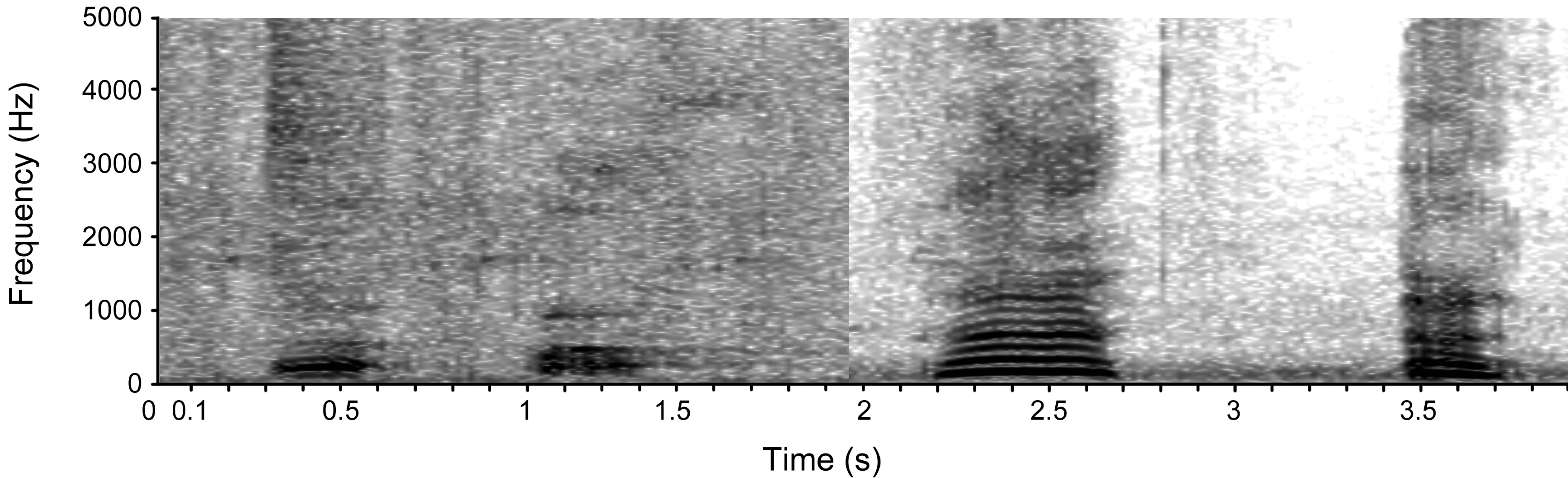
627

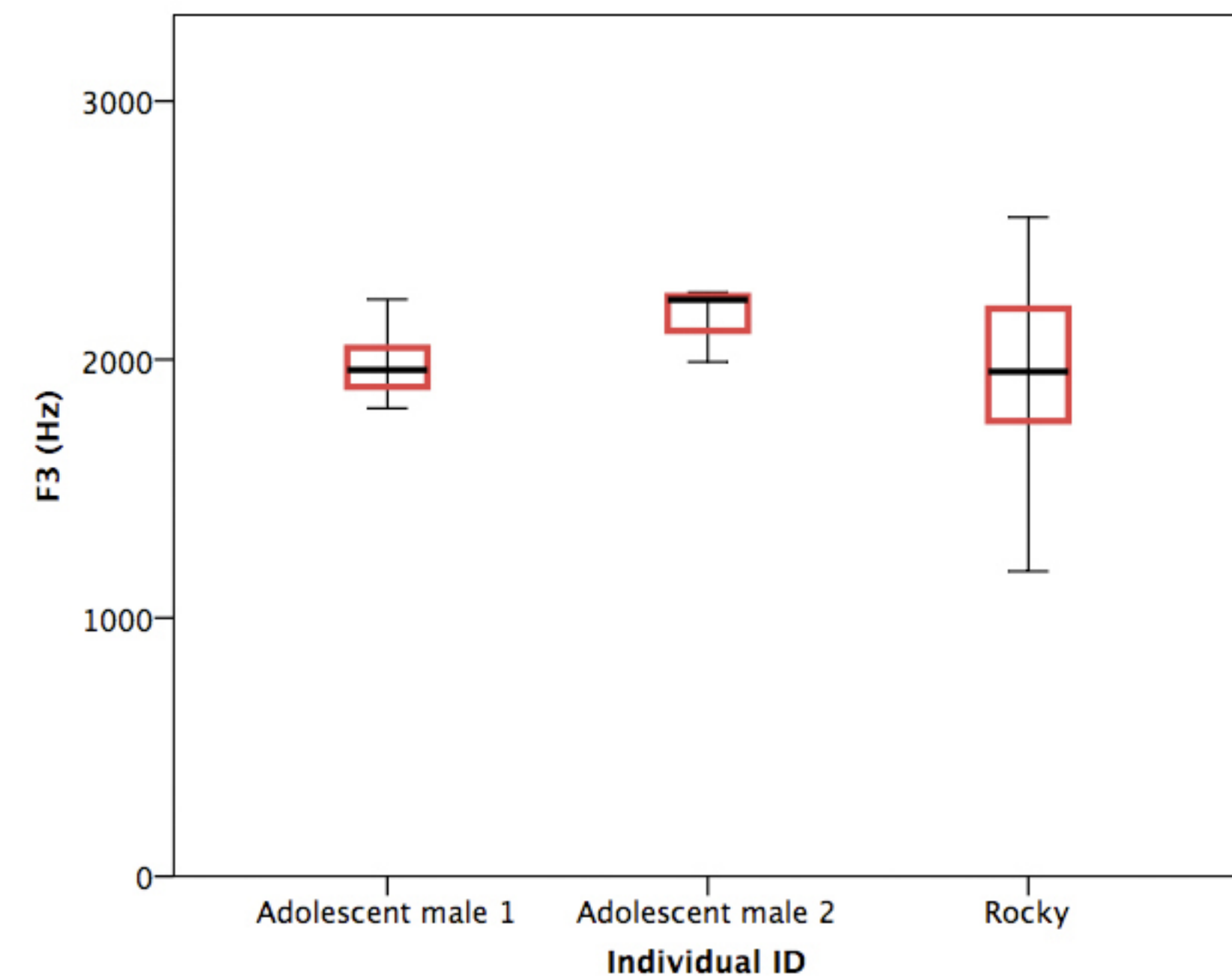
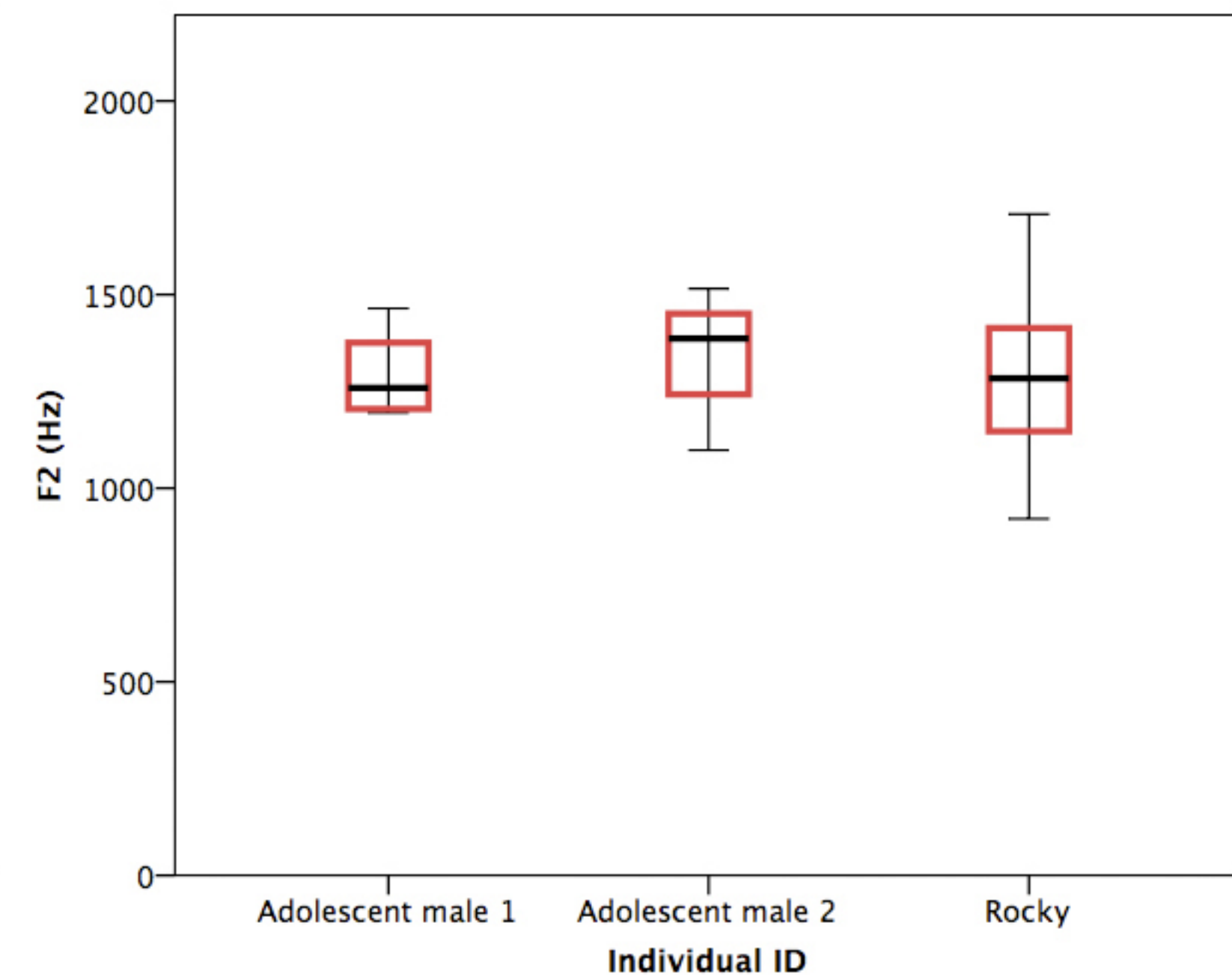
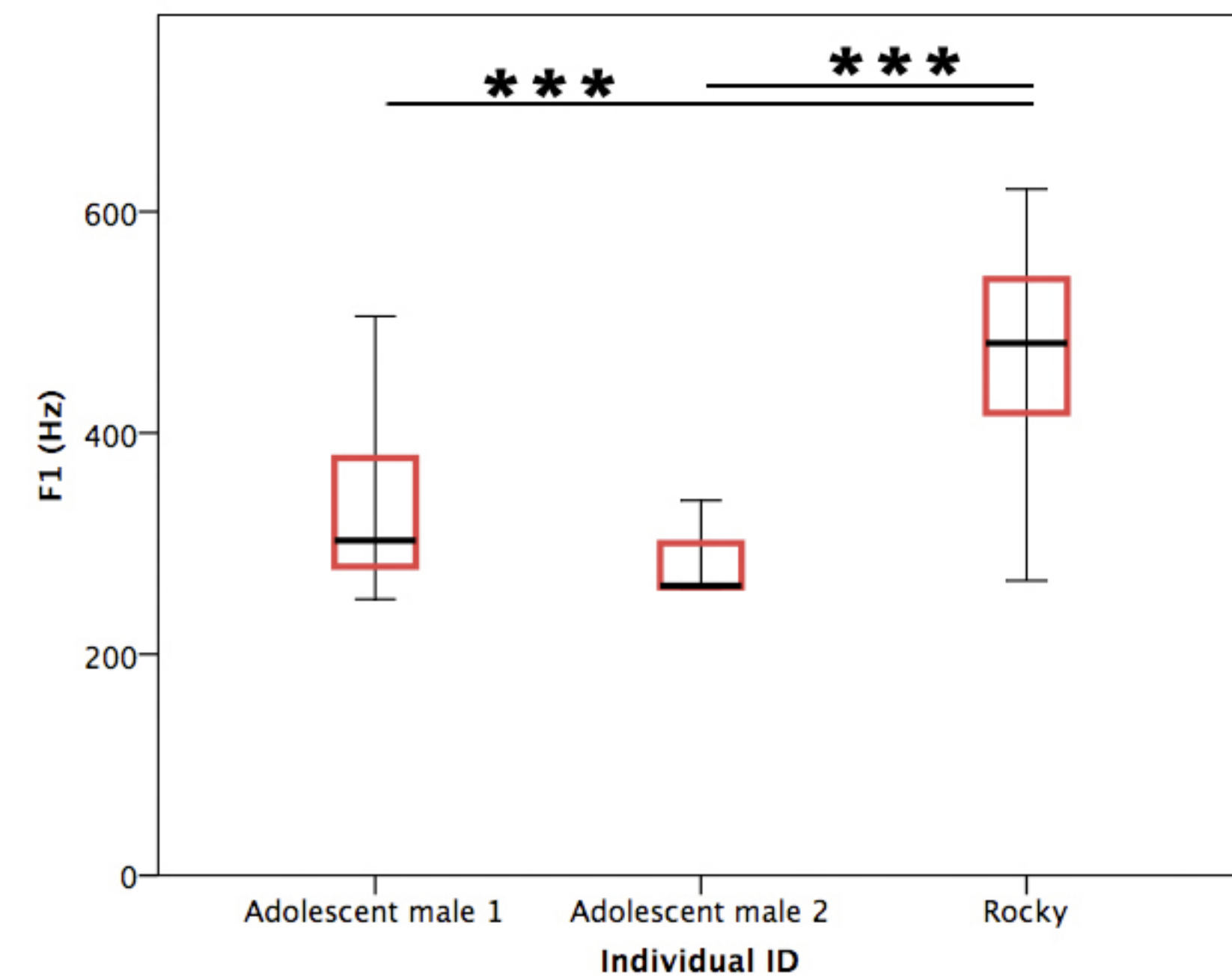
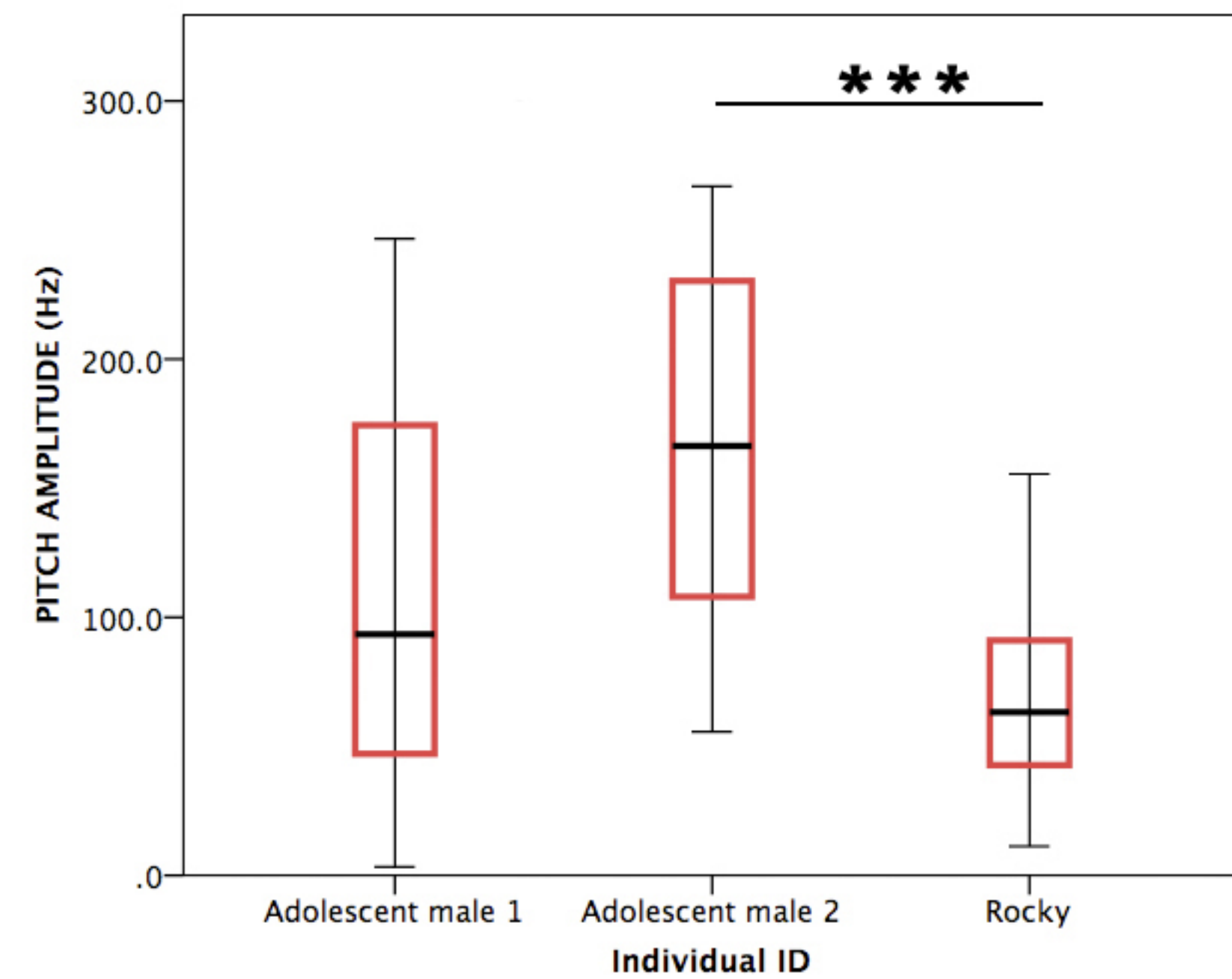
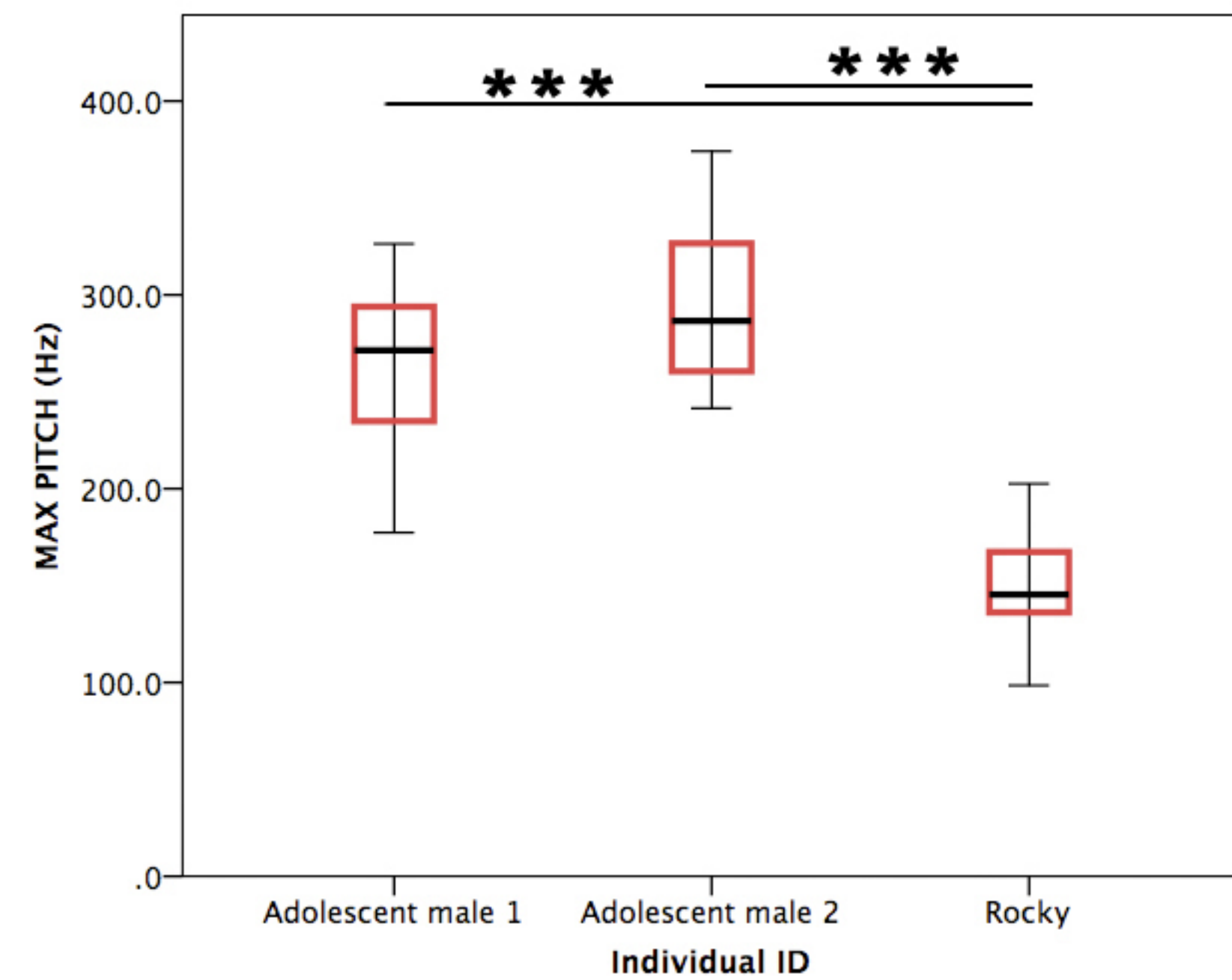
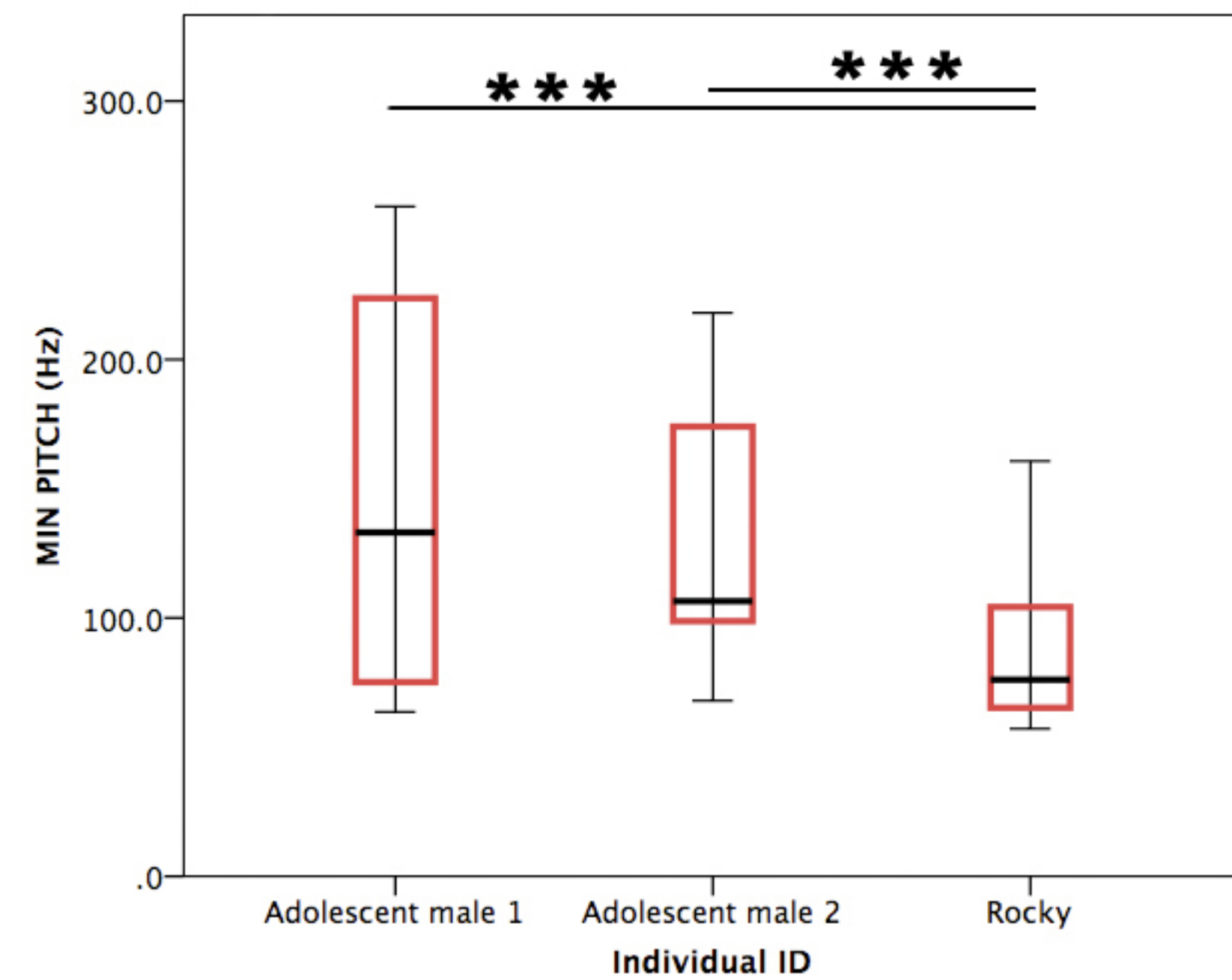
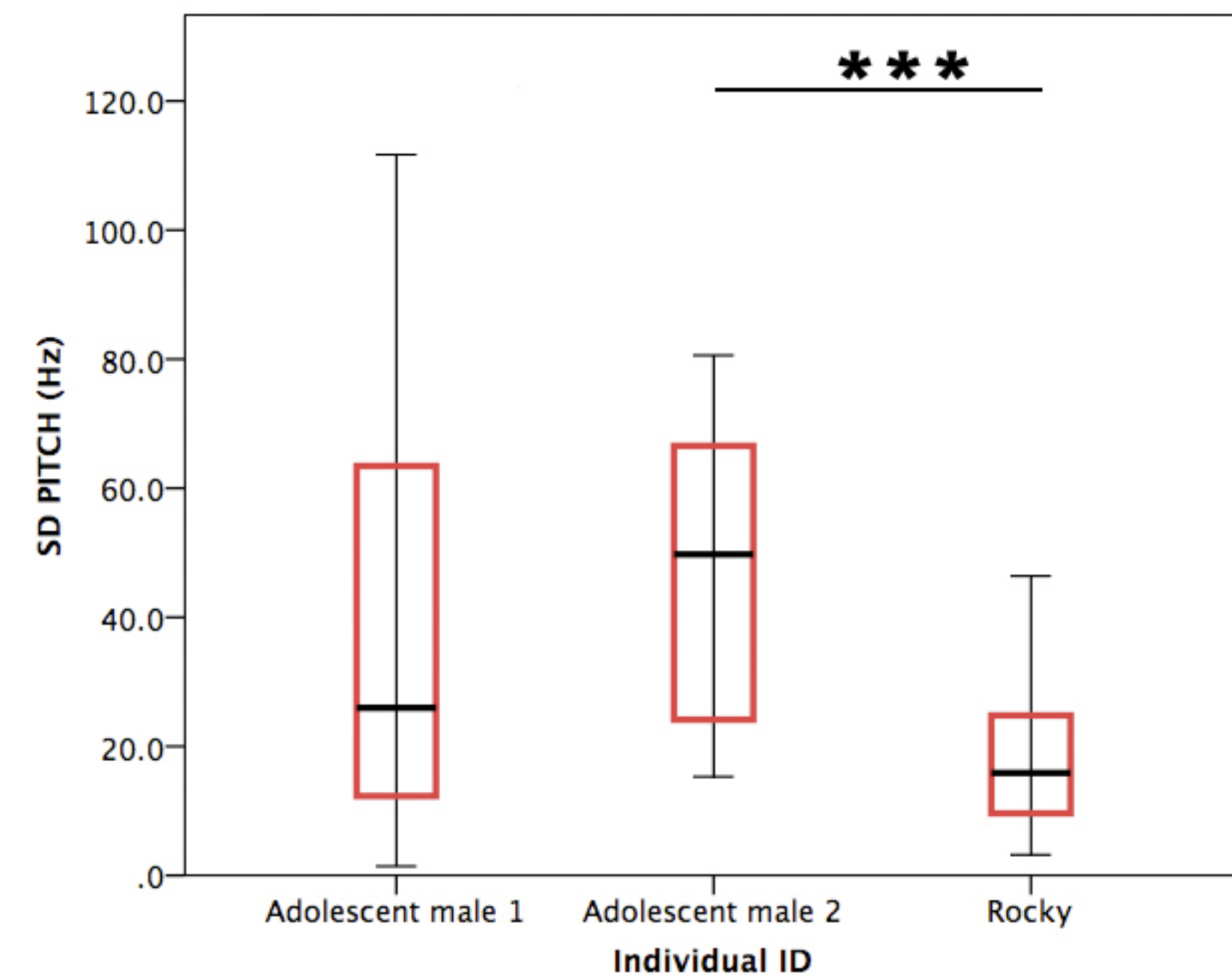
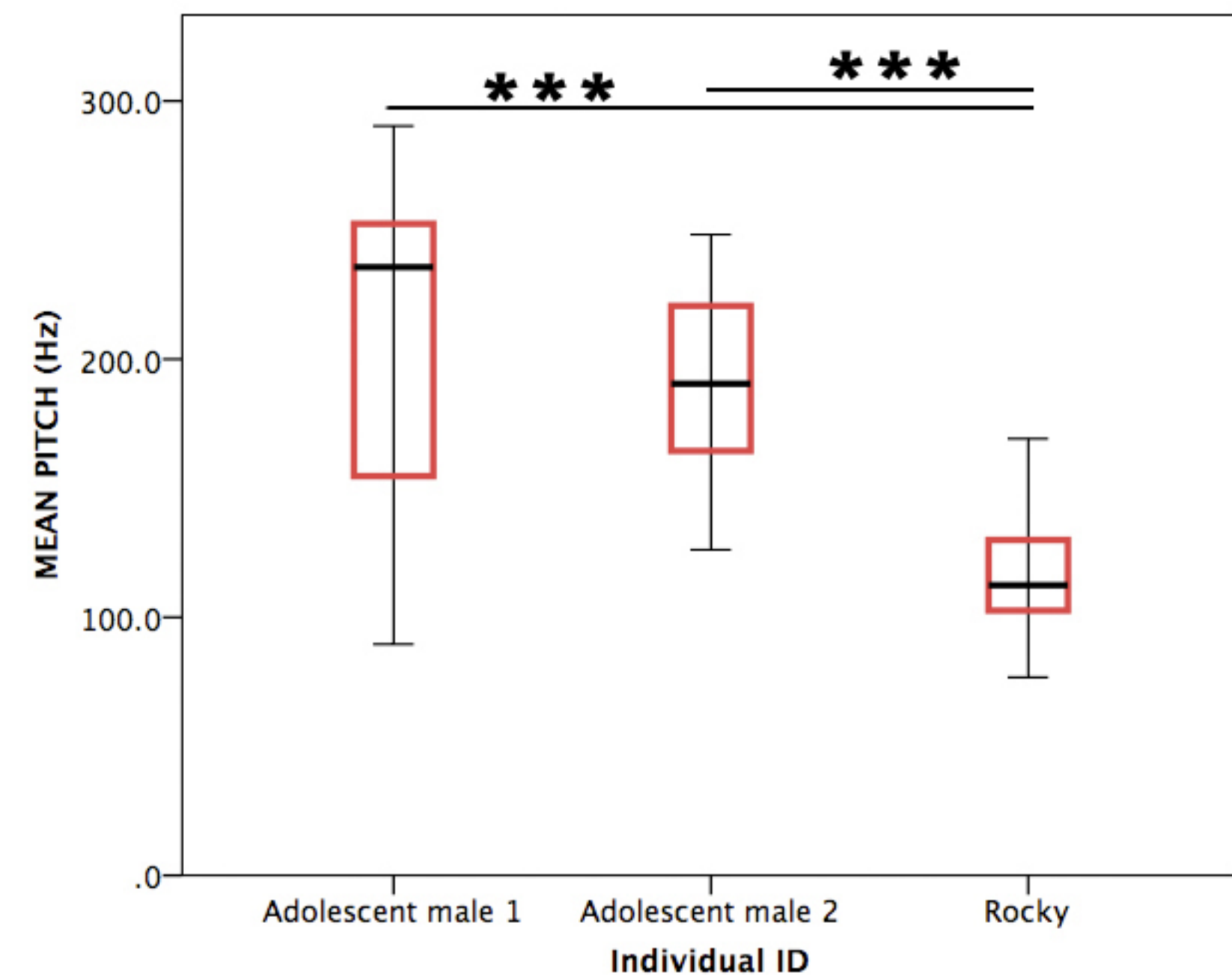
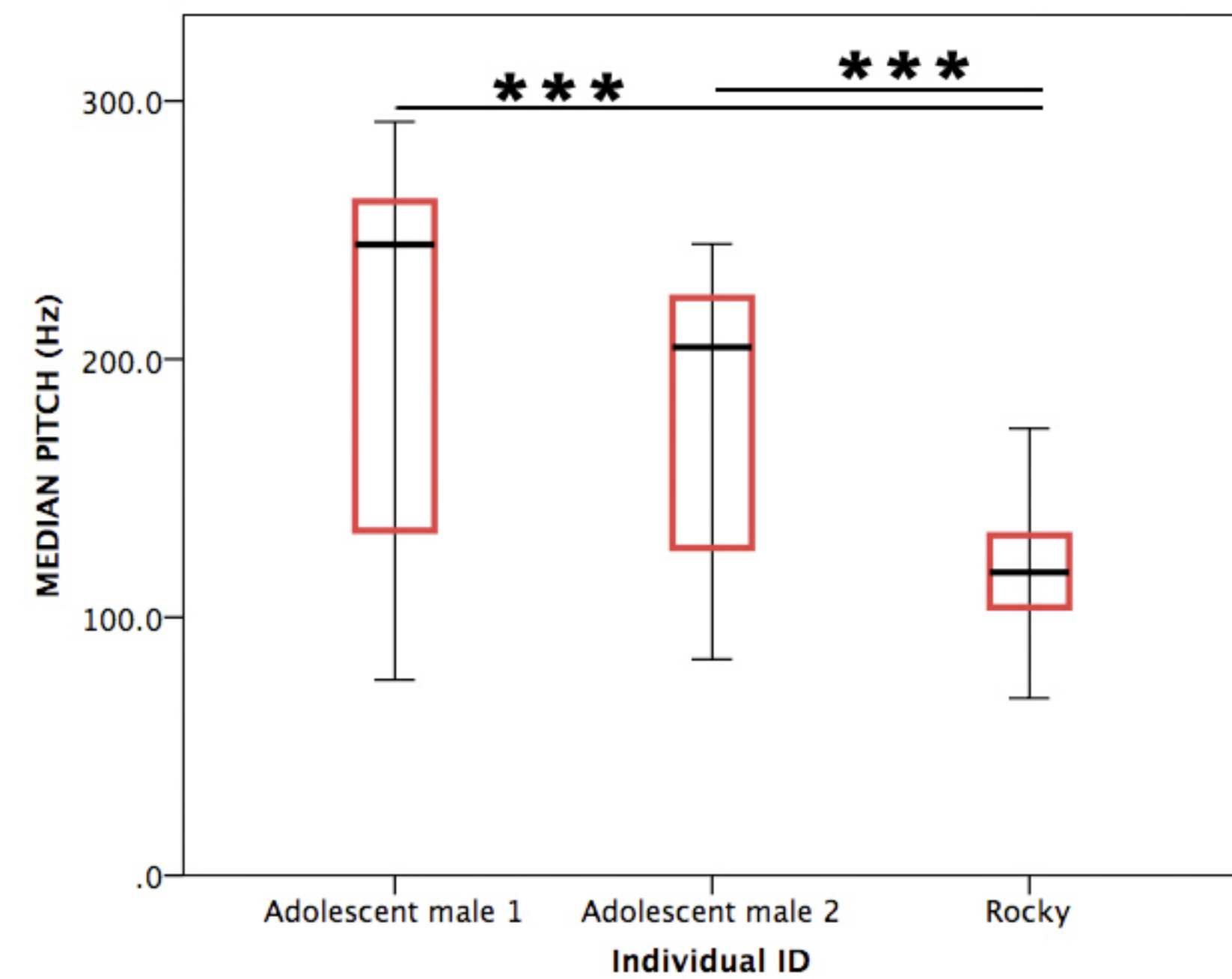
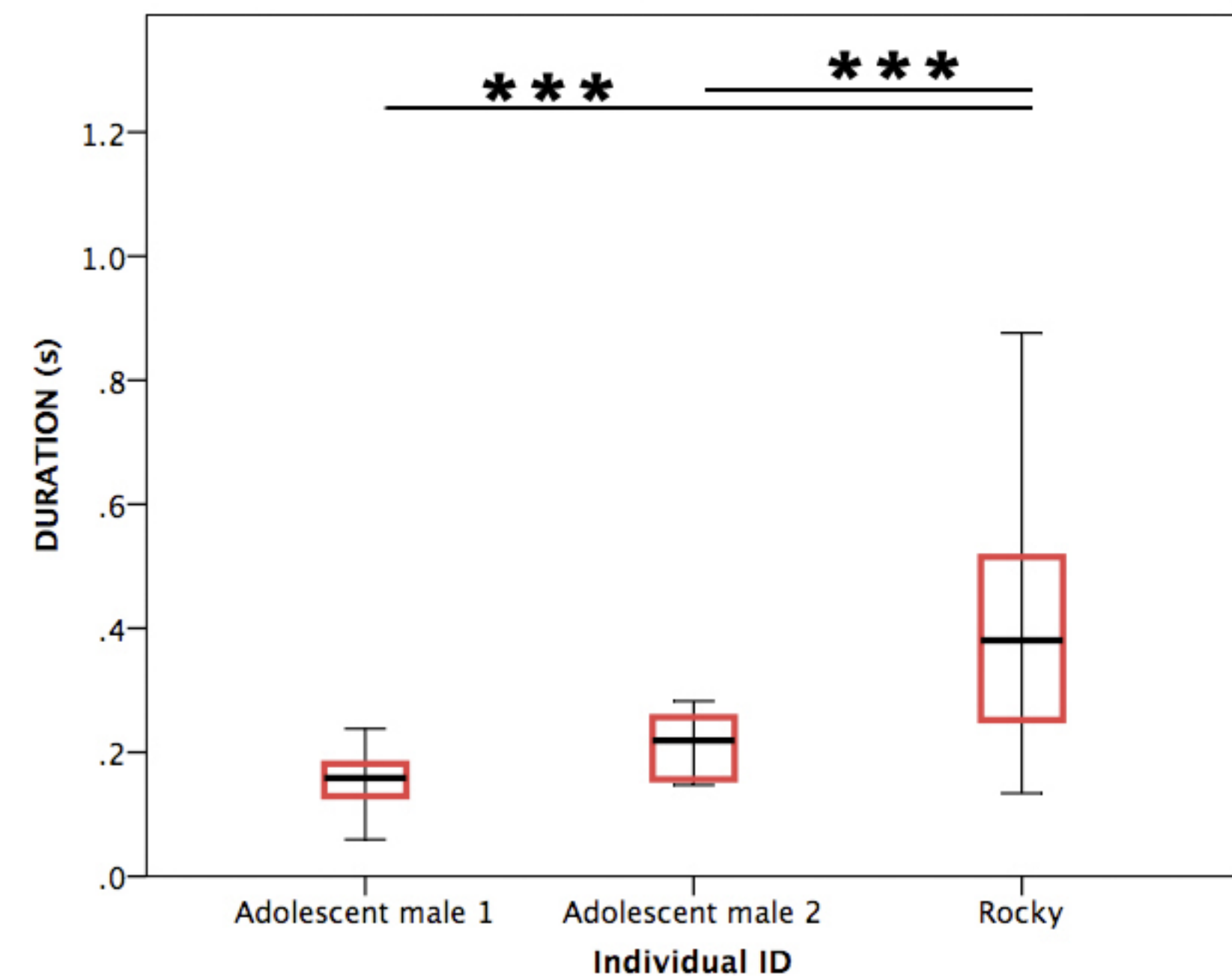
628 Fig. 6. Phonetogram displaying Rocky's wookies according to maximum frequency and
629 maximum power.

630

631

632





Orangutan maximum frequency (Hz)

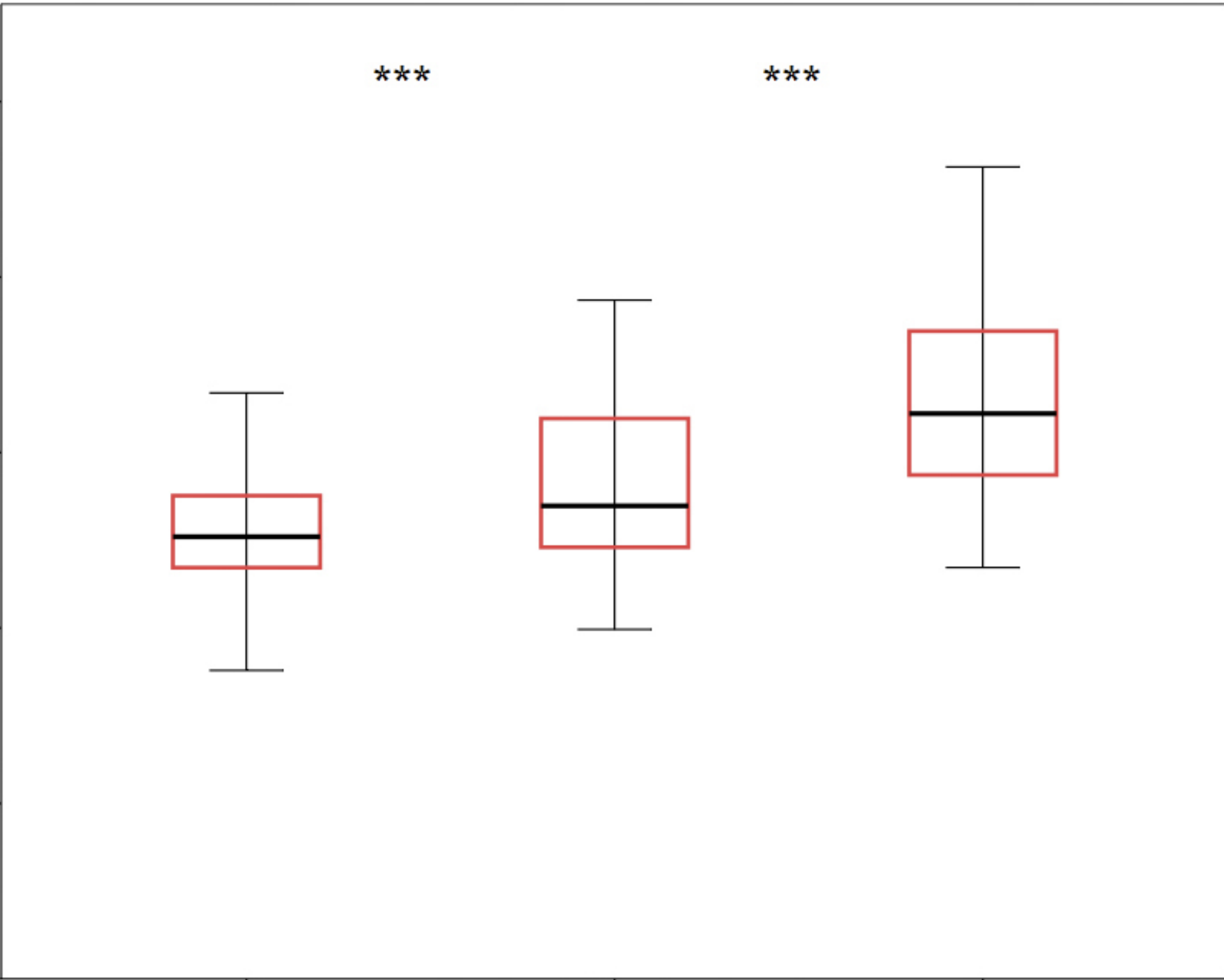
250
200
150
100
50
0

low

spontaneous

high

Wookie sub-variants



Function 2

