

1 **Early vocal production and functional flexibility in wild infant chimpanzees**

2 Guillaume Dezecache<sup>a,b,c,d,\*</sup>, Klaus Zuberbühler<sup>a,b,e</sup>, Marina Davila-Ross<sup>c</sup> & Christoph D.  
3 Dahl<sup>a,f,g,\*</sup>

4  
5 <sup>a</sup>Institute of Biology, University of Neuchâtel, Neuchâtel, Switzerland;

6 <sup>b</sup>Budongo Conservation Field Station, Masindi, Uganda;

7 <sup>c</sup>Department of Psychology, University of Portsmouth, Portsmouth, England, United  
8 Kingdom;

9 <sup>d</sup>Université Clermont Auvergne, CNRS, LAPSCO, Clermont-Ferrand, France;

10 <sup>e</sup>School of Psychology and Neuroscience, University of St Andrews, St Andrews, Scotland,  
11 United Kingdom;

12 <sup>f</sup>Graduate Institute of Mind, Brain and Consciousness, Taipei Medical University, Taipei,  
13 Taiwan;

14 <sup>g</sup>Brain and Consciousness Research Center, Taipei Medical University Shuang-Ho Hospital,  
15 New Taipei City, Taiwan

16  
17 \*Corresponding authors:

18 Guillaume Dezecache <[guillaume.dezecache@gmail.com](mailto:guillaume.dezecache@gmail.com)>

19 Christoph D. Dahl <[christoph.d.dahl@gmail.com](mailto:christoph.d.dahl@gmail.com)>

20  
21 **ABSTRACT**

22 How did human language evolve from earlier forms of communication? One way to address  
23 this question is to compare prelinguistic human vocal behavior with nonhuman primate calls.  
24 Here, an important finding has been that, prior to speech, human infant vocal behavior exhibits  
25 functional flexibility, the capacity of producing protophones that are not tied to one specific  
26 function. Nonhuman primate vocal behavior, by contrast, is comparably inflexible, with  
27 different call types tied to specific functions. Our research challenges the generality of this  
28 claim, with new findings of flexible vocal behavior in infant chimpanzees. We used artificial  
29 intelligence consisting of automated feature extraction and supervised learning algorithms to  
30 analyze grunt and whimper vocalizations from free-ranging infants during their first year of life.  
31 We found that grunt production was highly flexible occurring in positive, neutral and negative  
32 circumstances, as already shown in human infants. We also found acoustic variants of grunts  
33 produced in different affective contexts, suggesting gradation within this vocal category. By  
34 contrast, the second most common call type of infant chimpanzees, the whimpers, was produced  
35 in only one affective context in line with standard models of nonhuman primate vocal behavior.  
36 We concluded that the most common chimpanzee vocalization, the grunt, qualifies as  
37 functionally flexible, suggesting that evolution of vocal functional flexibility occurred before  
38 the split between the Homo and Pan lineages.

39  
40 Keywords: language evolution, vocal flexibility, wild chimpanzees, *Pan troglodytes*, grunts

## 41 INTRODUCTION

42 At some point in evolutionary history, there must have been a transition from primate-like  
43 inflexible to human-like flexible acoustic communication, which may have coincided with the  
44 origins of speech. The evolutionary history of this transition continues to be vividly debated  
45 (Fitch, 2018), with a large range of comparative evidence from animal communication systems,  
46 with the consensus view that direct evolutionary homologies are generally absent in the primate  
47 order (Rendall & Owren, 2002). More recently, however, some vocal and neural equipment has  
48 been identified in different primate species that allow for the production of speech-like sounds  
49 (Boë et al., 2017; Fitch, Boer, Mathur, & Ghazanfar, 2016; Lieberman, 2017) and for limited  
50 control over vocal fold oscillation (Lameira & Shumaker, 2019).

51  
52 A similar point of contention is whether there are fundamental differences in the ontogenetic  
53 trajectories between non-human primate and human vocal behavior prior to speech. By the age  
54 of one-month, human and great ape infants share parts of their vocal repertoire, such as crying  
55 or laughter, suggesting a common evolutionary root and biological function insofar as, in both  
56 species, the calls possess an illocutionary quality that expresses intuitively identifiable internal  
57 states to caregivers and other listeners (Jhang & Oller, 2017; Oller et al., 2013).

58  
59 However, human infants also produce a range of other sounds (called ‘protophones’), such as  
60 squeals, vocants and growls, that are not tied to particular affective states (Jhang & Oller, 2017;  
61 Oller et al., 2013). Some of these pre-linguistic sounds may convey relatively specific meanings,  
62 in the sense that infants produce them to express specific behavioral intentions rather than  
63 communicating specific affective states (Kersken, Zuberbühler, & Gomez, 2017). This apparent  
64 decoupling of signal structure and function in young infants, termed ‘vocal functional  
65 flexibility’, has been identified as a major evolutionarily precursor to language (Oller et al.,  
66 2013). Because of its early ontogenetic onset, vocal functional flexibility is said to be more  
67 foundational to human speech than other building blocks of the language faculty, such as proto-  
68 syntax or vocal elaboration (Oller et al., 2013). Vocal functional flexibility, in this view, is a  
69 prerequisite for speech development, and a major evolutionary departure from the functionally  
70 inflexible vocal behavior of non-human primates (Waal & Pollick, 2011).

71  
72 In one relevant study, Clay et al. (Clay, Archbold, & Zuberbühler, 2015) examined ‘peep’ calls  
73 in mature bonobos (*Pan paniscus*), their most common vocalizations, and found that they are  
74 produced in a variety of contexts, ranging from seemingly positive (food provisioning) to  
75 neutral (travel and resting) to negative (agonistic and alarm) situations. Based on these findings,  
76 the authors concluded that bonobos have the capability to produce sounds that are not  
77 affectively biased (Clay et al., 2015). Their peeps were, however, attributed to broad behavioral  
78 contexts (such as feeding or travelling) with no focus on more specific and transient behaviors  
79 that may infer affective contexts, such as when individuals suddenly experience aggression  
80 during travelling and feeding bouts. As such, the bonobo data are indicative of their peeps  
81 occurring across broad behavioral contexts but ultimately remain inconclusive in regards to  
82 whether vocal functional flexibility is indeed present in species other than humans. A second  
83 study, also on bonobos (Oller et al., 2019), suggests protophone-like vocal behavior with infants  
84 producing calls that occur in both low or moderate arousal situations, implying no affective  
85 binding. This conclusion has been preliminary, however, for the affective quality of the contexts  
86 surrounding vocalizations has proven difficult to discern.

87  
88 Here, we directly addressed the functional flexibility hypothesis by examining infant  
89 chimpanzee (*Pan troglodytes schweinfurthii*) vocal behavior at a very early age (< 12 months)  
90 and their affective context of occurrence. Examination of early vocal production is critical for

91 a more direct comparison with findings on human infants (Oller et al., 2013) and to test  
92 hypotheses about the evolutionary origins of functionally flexible vocal behavior. We took  
93 advantage of recent developments of applying machine learning techniques to the study of  
94 animal communication. We focused on two call types, the grunts and the whimpers, as they are  
95 acoustically distinct vocalizations that are common in young infants. Grunt calls are of  
96 particular importance as they develop into a central component of the vocal repertoire of  
97 chimpanzees and contribute to a variety of vocal sequences produced by juveniles, sub-adults  
98 and adults (Crockford & Boesch, 2005). For example, grunts complement panting elements  
99 during laughter and when encountering dominant individuals ('pant-grunts'). They are also  
100 produced upon encountering a food patch or when joining a foraging party ('rough grunts').  
101 Finally, they are routinely produced throughout resting or in relaxed social activities (Goodall,  
102 1986).

103  
104 Like in humans (McCune, Vihman, Roug-Hellichius, Delery, & Gogate, 1996), grunts are  
105 produced from the first days of life in chimpanzees. Their ontogenetic development has already  
106 been studied to some degree in chimpanzees, which has shown some flexibility in usage  
107 (Laporte & Zuberbühler, 2011). Two types of grunts can be distinguished, although no study  
108 has yet offered an acoustical validation of the existence of these diverse types. First, uh-grunts  
109 are short, tonal sounds, resembling human vowels {u}, {o} and {a}, sometimes produced in  
110 short series (staccato-grunts) (Kojima, 2003; Plooij, 1984). The second type are the so-called  
111 'effort' grunts, which represent the majority of grunting behavior in immature chimpanzees and  
112 are also present in humans and other mammals (McCune et al., 1996). Effort grunts are  
113 relatively soft and noisy and occur mostly during locomotor activities (Plooij, 1984). However,  
114 in adults, they are not mere byproducts of locomotion since they are sometimes emitted in the  
115 absence of movements, suggesting that the calls go through an ontogenetic transition from mere  
116 by-products of mechanical efforts to functionally active communicative signals in adults.

117  
118 Another common vocal utterance produced by chimpanzee infants is whimpers (Dezecache,  
119 Zuberbühler, Davila-Ross, & Dahl, 2019; Levréro & Mathevon, 2013; Plooij, 1984). They are  
120 short, tonal and often produced in series with an upward shift in fundamental frequency.  
121 Contrarily to grunts, whimpers preferentially occur in aversive contexts, likely homologous to  
122 human crying or distress calls in other mammals (Plooij, 1984). Previous research (e.g., (Plooij,  
123 1984)) has suggested the presence of whimper subtypes (single, serial and human-like  
124 whimpers), but we are not aware of any systematic acoustical analysis that would justify this  
125 nomenclature.

126  
127 To address the hypothesis that vocal flexibility in grunts evolved before the split between *Pan*  
128 and *Homo* lineages, we examined the vocal behavior of six wild chimpanzee infants aged  
129 between 0-12 months old from the Sonso community of Budongo Forest, Uganda. We analyzed  
130 the extent to which vocal production of grunt-like and whimper-like vocalizations were  
131 affectively biased, i.e., occurring in positive, negative or neutral situations.

## 132 133 **RESULTS**

### 134 ***Types of vocal utterances***

135 We inspected N = 1,016 vocal occurrences, of which N = 967 could be classified as either  
136 'grunts' (N = 833) (corresponding to a rough, harsh and noisy sound) or 'whimpers' (N = 134)  
137 (usually a series of low-pitch tonal calls with increase in fundamental frequency throughout the  
138 series). Other types of calls were identified as 'hoos' (n = 23), 'pants' (n = 15), 'screams' (n =  
139 2), 'squeaks' (n = 2) or 'barks' (n = 4). 'Laughter' (defined as grunting and panting) was  
140 uncommon (n = 3).

141

### 142 ***Functional flexibility***

143 Grunts: 44.8% of grunt-like vocalizations co-occurred with contexts classified as ‘positive’,  
144 40.9% with ‘neutral’, and 14.3% with ‘negative’. When considering each individual separately,  
145 a similar picture emerged (Figure 1), with most grunt-like vocalizations co-occurring with  
146 ‘positive’ and ‘neutral’ affective contexts. We sought to evaluate the evenness of the  
147 distribution of grunts across affective contexts and did so by calculating, for each infant, the  
148 numerical dominance of one affective context over the others (from  $1/3$  [= equiprobability of  
149 all 3 affective contexts] to 1 [= complete dominance of one of the affective contexts over the  
150 two others]). We found dominance to be relatively low in grunts, varying from 0.37 and 0.63  
151 (mean = 0.53; SD = 0.10), suggesting a stable and relative evenness in the affective distribution  
152 of grunts.

153

154 Whimpers: 94.8% of whimpers co-occurred with negatively classified contexts, and rarely with  
155 neutral (4.5%) or positive (0.7%) affective contexts. Inspection of individual distributions  
156 revealed the same pattern with whimper-like vocalizations systematically co-occurring with  
157 negatively classified contexts (Figure 1). The dominance of one affective context over the  
158 others in whimpers was relatively high, ranging from 0.89 to 1 (mean = 0.96; SD = 0.05),  
159 indicating low evenness in the affective distribution of whimpers.

160

161 Grunts vs. Whimpers: When comparing the distributional evenness of grunts vs. whimpers, we  
162 found dominance to be statistically higher in whimpers than in grunts (paired Wilcoxon signed  
163 rank test:  $V = 21$ ,  $p = .031$ ).

164

### 165 ***Acoustic variants of grunts***

166 We then classified  $N=180$  grunts ( $N=60$  per affective context) according to their association  
167 with positive, neutral, negative contexts in order to test for the presence of acoustic variants. In  
168 the first step, we followed a feature extraction procedure by extracting the means and  
169 covariances of mel frequency cepstral coefficients (MFCCs) for each call, and compared these  
170 values according to the calls’ associations (e.g. positive vs negative) using  $t$ -tests. This approach  
171 provides a general idea of how well positive, neutral and negative calls can be separated. We  
172 displayed the resulting  $p$ -values in an empirical cumulative distribution function (eCDF)  
173 (Figure 2A). We found that 5-10% of all features showed significant differences between the  
174 class labels at a 5%-significance level. In other words, 5-10 of 104 feature dimensions had  
175 strong discrimination power to distinguish between grunts pertaining to the various affective  
176 contexts.

177

178 In the second step, the feature selection procedure, we systematically varied the number of  
179 feature dimensions to be considered into the classification process (x-axis in Figure 2B). The  
180 feature dimensions that went into the classification process were determined by means of two  
181 methods: (1) simple filter feature selection method and (2) sequential feature selection method.

182

183 With the simple feature selection algorithm, the SVM correctly discriminated between classes  
184 at up to 80% (positive vs neutral:  $M = 78.99$ ,  $SD = 3.53$ ,  $t(59) = 63.69$ ,  $p < .001$ ; positive vs  
185 negative:  $M = 79.58$ ,  $SD = 1.83$ ,  $t(59) = 125.37$ ,  $p < .001$ ; neutral vs negative:  $M = 80.44$ ,  $SD$   
186  $= 2.06$ ,  $t(59) = 114.26$ ,  $p < .001$ ; red lines in Figure 2B). A substantial improvement was found  
187 when sequentially selecting feature dimensions: SVM correctly classified samples at up to 95%  
188 (positive vs neutral:  $M = 89.56$ ,  $SD = 4.84$ ,  $t(59) = 143.42$ ,  $p < .001$ ; positive vs negative:  $M =$   
189  $88.72$ ,  $SD = 4.49$ ,  $t(59) = 153.11$ ,  $p < .001$ ; neutral vs negative:  $M = 84.27$ ,  $SD = 5.23$ ,  $t(59) =$   
190  $124.91$ ,  $p < .001$ ; blue lines in Figure 2B). For all comparisons chance levels were 50% due to

191 the two-class comparisons applied. We further illustrated the simple feature selection outcomes  
192 by highlighting the feature dimensions selected (red circles in Figure 2C) among the feature  
193 dimensions not selected (black dots). Further, the features selected via the sequential feature  
194 selection are marked with blue x's. It becomes evident that the sequential feature selection  
195 yields better performance through sequential combinations of feature dimensions that, on  
196 average, fall more distal to the diagonal midline than the feature dimensions selected by the  
197 simple feature selection process. Sequential feature selection, to a large extent, included feature  
198 dimensions not selected by the simple feature selection method.

199

200 We further ensured that each individual was not contributing solely to the classification results  
201 of various contrasts. To test this, we repeated the classification process for the number of  
202 individuals ( $N = 6$ ) and excluded one individual at each time. As can be seen in Supplementary  
203 Figure 1, the classification performance did not improve nor deteriorate systematically,  
204 suggesting no effect due to caller identity (the average  $t$ -value of one-sample  $t$ -tests is  $97.52 \pm$   
205  $30.25$  (SD); all  $p$ -values were smaller than .001).

206

207 The use of means and covariances of cepstra yielded relatively high-performance scores in the  
208 classification routines at low computational loads. To assess whether certain feature dimensions  
209 (means and covariances of cepstra) occurred above chance across all comparisons, we  
210 determined the empirical distribution of occurrences of feature dimensions and contrasted it  
211 with a random distribution. While the use of the same feature dimension in up to 33% of the  
212 comparisons was not significantly different in the empirical distribution from the random  
213 distribution, the use of the same feature dimension in 50% of comparisons was significantly  
214 increased in the empirical distribution (Figure 3A). To describe the frequency bands explaining  
215 significant variances between classes of calls, we traced back the frequency bands underlying  
216 the significant feature dimensions, i.e., covariances of cepstra, and determined the sign of the  
217 covariances. We found a negative covariances between the following frequency bands (Figure  
218 3B): (1) band 2 (196.30 to 488.89 Hz) and band 4 (488.89 to 927.78 Hz), (2) band 4 (488.89 to  
219 927.78 Hz) and band 8 (1074.07 to 1366.67 Hz), band 6 (781.48 to 1074.07 Hz) and band 9  
220 (1220.37 to 1512.96 Hz). We found a positive covariance between the frequency bands 9  
221 (1220.37 to 1512.96 Hz) and 10 (1366.67 to 1659.26 Hz). Mean cepstra were significantly  
222 contributing in the frequency bands from (1) 50 to 342.59 Hz, (2) 196.30 to 488.89 Hz, (3)  
223 927.78 to 1220.37 Hz.



## 224 **DISCUSSION**

225 Oller and colleagues (Jhang & Oller, 2017; Oller et al., 2013; Oller, Griebel, & Warlaumont,  
226 2016; Oller & Griebel, 2004) posit that speech emerged from pre-linguistic vocalizations that  
227 are free of predetermined biological function, a precursor called ‘vocal functional flexibility’.  
228 Modern human infants regularly vocalize in such a way, in supposed contrast to the relative  
229 inflexibility of vocalizations in non-human primates (e.g., (Pollick & Waal, 2007; Waal &  
230 Pollick, 2011)). Functionally flexible vocalizations of young human infants, according to this  
231 theory, have evolved in humans in relation to allo-maternity (Burkart, Hrdy, & Van Schaik,  
232 2009; Hrdy, 2007, 2009; Kramer, 2010; Schaik & Burkart, 2010) or altriciality (Locke, 2006)  
233 and associated pressures on young infants to signal their needs and attract caregivers (Ghazanfar,  
234 Liao, & Takahashi, 2019; Locke, 2006; Zuberbühler, 2012).

235  
236 In the current study, we focused on the grunt-like and whimper-like calls of young chimpanzee  
237 infants, using novel coding strategies and state-of-the-art acoustic analysis tools. By contrast to  
238 previous studies, we elaborated a workable coding system which provides insight into the  
239 affective state of the animal, without solely relying on the broad behavioral contexts. We found  
240 that grunt-like calls are functionally flexible vocal units, produced frequently by chimpanzee  
241 infants in both positive and neutral situations, and less commonly also in negative situations.  
242 Importantly, the presence of grunts in contexts of low-to-mild arousal is consistent with the  
243 hypothesis of vocal functional flexibility (Oller et al., 2019), and so is the finding that grunts  
244 occur in similar proportion in positive and neutral contexts (Oller et al., 2013). On the other  
245 hand, whimper-like vocalizations seem to be confined to situations involving negative affective  
246 states in the infants. Their near absence in positive and neutral situations suggests that they  
247 represent a functionally rigid vocalization that has evolved for a narrow range of biological  
248 purposes, similar to cries in humans (Oller et al., 2013), to which they may functionally  
249 correspond (Goodall, 1986). Grunts, more generally, are a promising class of calls, insofar as  
250 their functional flexibility is in line with the ubiquity of this vocal category in a diversity of  
251 contexts in other primate species (Cheney & Seyfarth, 1982, 2018; Range & Fischer, 2004;  
252 Rendall, Seyfarth, Cheney, & Owren, 1999; Salmi, Hammerschmidt, & Doran-Sheehy, 2013).  
253 This being said, data suggest that production of sounds that are typically uttered under high  
254 stress (e.g., alarm calls) can also occur in the absence of the triggering stimulus (Lameira &  
255 Call, 2018), a pattern also suggested by the use of alarm calls to deceive conspecifics during  
256 foraging activities (Møller, 1988; Wheeler, 2009).

257  
258 Our second finding was systematic acoustic differences between grunts given in positive,  
259 neutral and negative situations, which enabled us to segregate acoustic variants of grunts into  
260 these categories. Acoustical differences linked to the affective context surrounding vocal  
261 production are common in humans as in other animals (Arias, Belin, & Aucouturier, 2018;  
262 Aucouturier et al., 2016; Banse & Scherer, 1996; Briefer, 2012; Goupil, Johansson, Hall, &  
263 Aucouturier, 2019; Ponsot, Burred, Belin, & Aucouturier, 2018; Williams & Stevens, 1972).  
264 Our data suggest that there is inter-gradation between grunt-types, with differences in acoustics  
265 relating to differences in contexts. Grunts, in other words, represent a coherent and unified call  
266 type that can manifest itself in acoustic variants in relation to the affective contexts in which  
267 they are produced.

268  
269 How exactly functionally flexible vocalization produced by human infants transition into  
270 speech sounds has been described in previous studies (Boysson-Bardies, 2001; de Boysson-  
271 Bardies, 1993; de Boysson-Bardies & Vihman, 1991; Elbers & Ton, 1985; Nathani, Ertmer, &  
272 Stark, 2006; Oller, 2000; Oller, Wieman, Doyle, & Ross, 1976). Although chimpanzee infants  
273 produce grunts in ways consistent with the functional flexibility hypothesis, they of course

274 never produce speech sounds and, historically, have failed to acquire human speech utterance  
275 even after extensive training (Hayes & Hayes, 1951). Instead, infant chimpanzee grunts may  
276 gradually develop into context-specific call variants with seemingly relatively narrow  
277 biological functions (Laporte & Zuberbühler, 2011; Slocombe & Zuberbühler, 2010; Slocombe  
278 & Zuberbühler, 2005; Watson et al., 2015), with clear acoustical boundaries notably between  
279 grunts used in greeting situations ('pant-grunts') and those produced upon encountering food.  
280 It is possible that the acoustic boundaries we identified between the grunts produced across  
281 affective states are the foundation of acoustic diversification in adults, although the categories  
282 used to define the affective states here (for instance, feeding and social approach are together  
283 considered 'positive') are not consistent with the vocal differentiation seen in adults (the grunts  
284 produced in feeding vs. social approach situations are acoustically distinct in adults (Crockford,  
285 in press; Goodall, 1986)). Alternatively, those calls may simply disappear and be absent from  
286 the adult repertoire, one causal factor being the relative absence of social reinforcement  
287 (including contingent vocal responses (Ghazanfar et al., 2019)) associated with grunt  
288 production, as compared to the frequent maternal reactions to distress calls (Dezecache et al.,  
289 2019).

290  
291 Our tentative to further explore the affective state of the infant by considering other cues, such  
292 as the infants' facial expressions or the mothers' behavior, faced considerable challenges. We  
293 found that infant facial movements are extremely fast and fluid, which prevented us from  
294 reliable coding particularly in the wild. For this reason, the behavioral context of the infant  
295 alone was the most relevant available cue to approach the affective dimension of the situation.  
296 While we must yet acknowledge the limitations pertaining to the fact that judgments of infants'  
297 affect were made based on the infants' behavioral contexts and done so by a human observer,  
298 the results of the acoustic analysis are providing important support for the approach used to  
299 categorize affect in the present work. Future studies should investigate the affective impact of  
300 other communicative signals used by infants (Fröhlich & Hobaiter, 2018; Fröhlich, Wittig, &  
301 Pika, 2018).

302  
303 Another hint to the affective dimension of the situation is the mothers' behavior. Protocols  
304 where mothers may be asked to interact with toddlers may yield to responsiveness from the  
305 mothers whichever the affective state of the infant is (Yoo, Bowman, & Oller, 2018). In the  
306 course of spontaneous behavior, though, we may expect little intervention from the chimpanzee  
307 mothers, except in situations where the infant is in danger. In our sample, responsiveness of the  
308 mother (tentatively defined in pilot coding as being either proactive, protective or neutral by  
309 the observer) was relatively low (proactive or protective less than half of the time), a pattern  
310 which might be due to differences in mothering style between chimpanzees and humans, or a  
311 difference between our own study (where no particular demand is put on the mother) and others  
312 (where mothers may be interacting with their infant, e.g., (Oller et al., 2013)).

313  
314 Although playback of infant grunts to the mother may appear like a methodological possibility  
315 to further establish maternal assessment of the affective state of the infant or to be able to see  
316 whether mothers respond differently to positive, neutral and negative grunts (Fischer, Noser, &  
317 Hammerschmidt, 2013; Fischer, 2016; Zuberbühler, 2014), this would require either playing  
318 the infants' calls in its own presence (which is ethically inappropriate) or playing the calls of  
319 another infant to a mother (which may not trigger any reaction at all in the non-genetically  
320 related mother).

321  
322 In latest research, the comparative volubility (quantity of sounds produced in a given period of  
323 time) of human infants and other animals (Ghazanfar et al., 2019; Ghazanfar & Takahashi,

324 2014; Oller et al., 2019; Takahashi et al., 2015), and the privileged function of protophone-like  
325 vocalizations to elicit social interactions and vocal turn-taking with caregivers (Oller et al.,  
326 2019; Yoo et al., 2018). In humans, non-affectively bound vocalizations appear to occur more  
327 often than affectively bound vocalizations (such as crying) (Oller et al., 2019). They occur in  
328 solitary contexts where infants invest in practice and vocal exploration. They also occur in  
329 interactive contexts, so as to elicit and regulate social interactions with caregivers. Caregivers  
330 appear to detect the functional difference between protophones (as potentially interactive calls)  
331 and other calls (such as cries), where caregiver intervention is solicited (Yoo et al., 2018).  
332 Comparison with bonobo infants suggested much higher rate of production of non-affectively  
333 bound vocalizations and much higher vocal investment in social interactions in human infants  
334 (Oller et al., 2019). Whether human infants also are comparably more talkative than their  
335 chimpanzee counterparts is a question we need to be exploring. This should be preferably  
336 investigated in captive or semi-captive settings, where true calling rate can be assessed, for  
337 video monitoring is less likely to be interrupted and for levels of ambient noise could be  
338 comparatively less problematic. Such problems have already been acknowledged by (Oller et  
339 al., 2019) regarding previous report on the flexible development of grunting behavior in wild  
340 chimpanzees as well as their rate of occurrence (Laporte & Zuberbühler, 2011). Data from the  
341 vocal development of one captive chimpanzee indicate lower volubility than in humans  
342 (Kojima, 2003). Future studies should evaluate this fact with a larger sample.

343  
344 In conclusion, our study suggests that chimpanzees may possess a feature that is fundamental  
345 to the development of speech in humans, the ability to produce vocalizations that are not  
346 strongly bound to one particular affective context, but are produced in a functionally flexible  
347 manner. However, we should expect that future research will reveal further examples. For  
348 instance, coo calls in several macaque species (Hsu, Chen, & Agoramorthy, 2005; Owren &  
349 Casale, 1994) or grunts in vervet monkeys (Seyfarth & Cheney, 1984) also seem to be given in  
350 a variety of contexts. Future research will equally have to address the question of how selection  
351 favored acoustic diversification of functionally flexible vocal behavior into speech in humans.  
352 The main driver for this transition, it has been argued, may have been the highly cooperative  
353 breeding system of humans, with infants regularly looked after by individuals other than the  
354 mother, which requires infants to become more active agents in forming social bonds from a  
355 much younger age than in great ape infants (Ghazanfar et al., 2019; Zuberbühler, 2012).

356  
357 Cooperative breeding, in this view, may thus have transformed a functionally flexible vocal  
358 system, evolved prior to the split between humans and apes, into the uniquely human way of  
359 using vocal signals to interact socially. Another complementary reasoning is that humans' high  
360 altriciality selected for the most vocal individuals, capable of attracting caregivers (Locke,  
361 2006). Future studies should clarify the relative contribution of both factors through mapping  
362 the phylogenetic distribution of vocal functional flexibility.

## 363 364 **METHODS**

### 365 *Ethics*

366 Permission to conduct the study was obtained from the Ugandan Wildlife Authority (UWA)  
367 and the Uganda National Council for Science and Technology (UNCST).

### 368 369 *Subjects and data collection*

370 Data were collected in the Sonso community of the Budongo Forest Reserve, Uganda  
371 (Reynolds, 2005) between February-June 2014, December 2014 and March-June 2015. This  
372 community comprises around 70 individuals well habituated to human observers. The natural  
373 behavior of N=7 infants was video recorded continuously during focal animal sampling, using



374 Panasonic HC X909/V700 cameras, with a Sennheiser MKE-400 shotgun microphone. Six of  
375 those infants produced enough calls to be further considered in the statistical analysis (see Table  
376 1 for details).

377  
378 ***Behavioral data analysis***

379 Videos were inspected for the presence of infant vocalizations. We defined vocal behavior as  
380 the occurrence of single sound units or series of sounds produced by the infant's vocal apparatus,  
381 separated by a least 5 seconds of silence.

382  
383 As of today, there is no definitive repertoire of infant chimpanzee vocal behaviors. The  
384 categories used in this research are based on GD's assessment. This assessment proved reliable  
385 when confronted to an independent assessment with Derry Taylor, using vocalizations from  
386 infant and juvenile semi-wild chimpanzees from the Chimfunshi Wildlife Orphanage, Zambia,  
387 collected by DT. One hundred-and-sixty vocalizations were indeed classified as belonging to  
388 either the 'grunt', 'whimper', 'scream' or 'laughter' category. Agreement was excellent ( $k =$   
389  $0.77$ ), even better when considering only 'grunts' and 'whimpers' ( $k = 0.92$ ).

390  
391 For each vocal occurrence, we coded infant behavior from the following list of mutually  
392 exclusive behavioral contexts (summarized in Table 2). The internal state of the infant was  
393 classified as 'positive' if it showed one of the following four behaviors: (1) 'play' (showing  
394 relaxed, joyous movements without obvious purpose, either as 'social play' (accompanied by  
395 tell-tale behavior such as embracing and gentle biting) or 'solitary play'; (2) giving or receiving  
396 'grooming' (note that allo-grooming was never observed in our infants); (3) 'feeding'  
397 (breastfeeding or swallowing an edible element), and (4) 'social approach' (greeting a  
398 conspecific while moving towards it).

399  
400 The infant's internal state was classified as 'neutral' if it showed one of the following behaviors:  
401 (5) 'resting' (remaining with a limited area, sometimes moving within); (6) 'moving'; (7)  
402 'manipulating objects' (such as leaves, branches, rocks) without playful postures, or (8)  
403 'greeting without approach' (calling upon the approach of a conspecific without showing nor  
404 approach or avoidance behavior towards it).

405  
406 Infant behavior was classified as 'negative' if it showed one of the following behaviors: (9)  
407 'nuzzling' (unsuccessfully trying to access the mother's nipple); (10) 'begging' (attempting to  
408 access food other than breast milk); (11) 'hiding' (increased gripping or seek for contact with  
409 the mother when contact was already established between them); (12) 'contact mother/kin' was  
410 coded if infants were urgently seeking contact with the mother or a kin when contact was not  
411 already established between them; (13) 'escaping' (when the infant shows escape movements  
412 away from an activity it is involved in). Escaping could also include moment of discomfort  
413 when the infant is pressed against the belly of the mother or stuck in a bad position.

414  
415 We performed intra-coder reliability tests on the affective contexts coded as positive, neutral  
416 and negative. For this, we randomly selected 200 video clips (around 19% of the coded dataset  
417 composed of the 7 infants), which were coded independently during two coding sessions more  
418 than a year apart (November 2015 and February 2017), so that the second coding was, notably,  
419 naïve. We found strong agreement between the two coding sessions ( $k = 0.73$ ).

420  
421 ***Statistical analyses***

422 In order to evaluate the evenness of the distributions of grunts and whimpers across affective  
423 contexts, we calculated, for each infant, and for grunts and whimpers separately, the dominance  
424 of one affective context over the two others, using the Berger-Parker Dominance index:

$$425 \text{dominance} = N_{\max} / N$$

426 where  $N_{\max}$  is the number of calls in the most abundant affective context;  $N$  the total number of  
427 calls across all affective contexts. Dominance values range from 1 / number of affective  
428 contexts (= equiprobability of calls across affective contexts; here  $1 / 3 = 0.33$ ) to 1 (= complete  
429 dominance of one of the affective contexts over the others).

430 Dominance values were compared between grunts and whimpers using a paired Wilcoxon Sign-  
431 Ranked test. Analyses were carried out using R (version 3.6.1 (R Core Team, 2018)) and R  
432 Studio (version 1.2.1335 (RStudio Team, 2015)).

### 433 *Acoustic analyses*

434 Acoustic data analysis focused on grunts for they were the only vocal category for which at  
435 least two of the affective contexts were well represented. The acoustic structure of whimpers  
436 has been analyzed as part of another study.  $N=180$  grunts were extracted. For each affective  
437 context, 60 were randomly selected. Following extraction, we used MATLAB (MathWorks  
438 Inc., Natick, MA, USA) for the acoustic data analysis, consisting of features extraction, feature  
439 selection and call classification. We first pre-processed the audio files by applying a band pass  
440 filter from 50 to 4000 Hz and normalized the signals using the following function:

$$441 \text{signal} = (\text{signal} - \text{mean}(\text{signal})) / \max(\text{abs}(\text{signal} - \text{mean}(\text{signal})))$$

### 442 *Feature extraction and selection*

443 We first ran a feature extraction algorithm to reduce redundancy of information and  
444 computational efforts in classifying the calls and to maximize the generalization ability of the  
445 classifier (Tajiri, Yabuwaki, Kitamura, & Abe, 2010). A popular method is extraction of mel  
446 frequency cepstral coefficients (MFCCs) (Supplementary Figure 2), a procedure that adapts  
447 function parameters to the primate auditory system (Fedurek, Zuberbühler, & Dahl, 2016;  
448 Mielke & Zuberbühler, 2013). While a typical spectrogram linearly scales frequencies (i.e.,  
449 each frequency bin is spaced an equal number of Hertz apart), the mel-frequency scale is a  
450 logarithmical spacing of frequencies. We divided the calls into segments of 25ms length and  
451 10ms steps between two successive segments. We warped 26 spectral bands and returned 13  
452 cepstra, which resulted in feature dimensions of 13 values each. We then took the mean and co-  
453 variances of each cepstra over the collection of feature segments, resulting in a 13-value vector  
454 and a 13 x 13-value matrix, respectively, and concatenated to 104-unit vectors ((Mandel & Ellis,  
455 2005), p. 594-599) (Figure 3). We applied feature scaling to [0 to 1] and mean normalization.

456 Second, we performed a feature selection procedure, a crucial part in statistical learning: too  
457 many feature dimensions are not useful for producing reliable classification systems, whereas  
458 low sample numbers can lead to over-fitting to noisy feature dimensions. We therefore selected  
459 a subset of the original feature dimensions and evaluated classification performance based on  
460 sequentially selected feature sets until there was no improvement in performance. At this end,  
461 we subdivided the entire data set into a training (75%) and a test data set (25%) and applied a  
462  $t$ -test on each feature dimension, comparing values of given feature dimension sorted by  
463 predefined class labels (e.g. grunts occurring in negative (1) vs. positive (2) contexts) and used  
464  $p$ -values as a measure separability of the two classes. We plotted the  $p$ -values as an empirical  
465 cumulative distribution function (eCDF) to get an understanding of how well each feature  
466

472 separated the two classes and how many features contributed to a significant separation (5%-  
473 level). We ran this procedure 20 times for each comparison and plotted the results individually  
474 (grey lines) and the mean of all repetitions (black line) (Figure 2A). The classification routines  
475 were then independently run either on feature dimensions selected according to the  
476 discrimination power (decreasing order) (red lines in Figure 2B), as shown in the eCDF plots  
477 (Figure 2A). Such procedure is referred to as a simple filter approach on feature selection, where  
478 general characteristics of the extracted features are taken into consideration when selecting  
479 feature dimensions, without subjecting them to a classifier. We also applied a more extensive  
480 procedure of feature selection by sequentially selecting feature dimensions by adding (forward  
481 search) feature dimensions, referred to as sequential feature selection (blue lines in Figure 2B).  
482 As part of this method, the algorithm searched the best feature dimensions (predictors)  
483 according to their individual classification performance in the given subset of data. For each  
484 candidate feature subset (predictor), the algorithm performed a 10-fold cross-validation  
485 procedure with different training and test subsets. After computing the mean performance  
486 values for each candidate feature subset, the algorithm chooses the candidate feature subset  
487 with minimal misclassification. For both methods, we systematically varied the number of  
488 features used for classification (x-axis in Figure 2B). The selected features from a single run of  
489 the sequential search algorithm are illustrated in Figure 2C. Scales reflect the feature-scaled  
490 and normalized values, as a result of feature extraction, from which the grand means (i.e. for  
491 each feature dimensions across all data) were subtracted. This measure was used to visually  
492 highlight differences and was not used in further analyses.

493

#### 494 ***Classification***

495 We used support vector machine (SVM) with a radial basis function (RBF) Kernel (Vert, Tsuda,  
496 & Schölkopf, 2004) for the classification of calls according to the class labels (negative, neutral  
497 and positive contexts). A classification procedure contains a training phase followed by a test  
498 phase. We therefore separated training samples and labelled them according to an attribute of  
499 interest (e.g. negative (1) vs. positive (2) contexts). The algorithm then created a model that  
500 optimally separates the two classes. In the test phase, samples without attribute labels were fed  
501 into the model to measure its generalization performance. We used the SVM implementation  
502 from LIBSVM toolbox (Chang & Lin, 2011). To evaluate how the classification results  
503 generalize to a novel and independent data set, we 10-fold cross-validated the classification  
504 process and optimized the parameters C and gamma (Fedurek et al., 2016), with the C taking  
505 values in a range of ( $2^{-1}$ ,  $2^3$ ) and gamma in a range of [ $2^{-4}$ ,  $2^1$ ]. In addition, to ensure that no  
506 single individuals contributed solely to the classification outcome, we ran a leave-one-out  
507 algorithm, where the procedure described above was re-run six times, excluding one of the  
508 individuals in each run.

509

#### 510 ***Feature evaluation***

511 To evaluate whether certain feature dimensions are particularly critical for the classification of  
512 grunts, we assessed whether feature dimensions have been repeatedly used by the classifier  
513 overall in the classification of grunts. We therefore considered the three types of comparisons,  
514 positive vs neutral, positive vs negative and neutral vs negative grunts, as well as the two feature  
515 evaluation algorithms (simple feature selection and sequential feature selection). Each  
516 comparison, as described above, was ten-fold cross-validated. We then calculated the empirical  
517 distribution of the ten features with best classification power, as determined by the feature  
518 selection algorithms (see above). Also, we determined a random distribution of “best features”  
519 for each comparison by randomly selecting 10 out of 104 features. The frequency distribution  
520 across all comparisons were determined and 95% confidence intervals were calculated by

521 running the procedure 1,000 times. We then traced back the significant feature dimensions to  
522 the underlying frequency bands in Hertz.

523

## 524 REFERENCES

525 Arias, P., Belin, P., & Aucouturier, J.-J. (2018). Auditory smiles trigger unconscious facial  
526 imitation. *Current Biology*, *28*(14), R782–R783.

527 <https://doi.org/10.1016/j.cub.2018.05.084>

528 Aucouturier, J.-J., Johansson, P., Hall, L., Segnini, R., Mercadié, L., & Watanabe, K. (2016).

529 Covert digital manipulation of vocal emotion alter speakers' emotional states in a  
530 congruent direction. *Proceedings of the National Academy of Sciences*, *113*(4),

531 948–953. <https://doi.org/10.1073/pnas.1506552113>

532 Banse, R., & Scherer, K. R. (1996). Acoustic profiles in vocal emotion expression. *Journal*  
533 *of Personality and Social Psychology*, *70*(3), 614–636.

534 <https://doi.org/10.1037/0022-3514.70.3.614>

535 Boë, L.-J., Berthommier, F., Legou, T., Captier, G., Kemp, C., Sawallis, T. R., ... Fagot, J.

536 (2017). Evidence of a Vocalic Proto-System in the Baboon (*Papio papio*) Suggests  
537 Pre-Hominin Speech Precursors. *PLOS ONE*, *12*(1), e0169321.

538 <https://doi.org/10.1371/journal.pone.0169321>

539 Boysson-Bardies, B. de. (2001). *How Language Comes to Children: From Birth to Two*  
540 *Years*. Cambridge, MA: MIT Press.

541 Briefer, E. F. (2012). Vocal expression of emotions in mammals: mechanisms of  
542 production and evidence. *Journal of Zoology*, *288*(1), 1–20.

543 <https://doi.org/10.1111/j.1469-7998.2012.00920.x>

544 Burkart, J. M., Hrdy, S. B., & Van Schaik, C. P. (2009). Cooperative breeding and human  
545 cognitive evolution. *Evolutionary Anthropology*, *18*(5), 175–186.

- 546 Chang, C.-C., & Lin, C.-J. (2011). LIBSVM: a library for support vector machines. *ACM*  
547 *Transactions on Intelligent Systems and Technology (TIST)*, 2(3), 27.  
548 <https://doi.org/10.1145/1961189.1961199>
- 549 Cheney, D. L., & Seyfarth, R. M. (1982). How vervet monkeys perceive their grunts: Field  
550 playback experiments. *Animal Behaviour*, 30(3), 739–751.  
551 [https://doi.org/10.1016/S0003-3472\(82\)80146-2](https://doi.org/10.1016/S0003-3472(82)80146-2)
- 552 Cheney, D. L., & Seyfarth, R. M. (2018). Flexible usage and social function in primate  
553 vocalizations. *Proceedings of the National Academy of Sciences*, 115(9), 1974–  
554 1979. <https://doi.org/10.1073/pnas.1717572115>
- 555 Clay, Z., Archbold, J., & Zuberbühler, K. (2015). Functional flexibility in wild bonobo vocal  
556 behaviour. *PeerJ*, 3, e1124. <https://doi.org/10.7717/peerj.1124>
- 557 Crockford, C. (in press). Why Does the Chimpanzee Vocal Repertoire Remain Poorly  
558 Understood? And What Can Be Done About It. In *The Tai Chimpanzees: 40 years of*  
559 *Research*. Eds: Boesch C. and Wittig R. Cambridge University Press.
- 560 Crockford, C., & Boesch, C. (2005). Call combinations in wild chimpanzees. *Behaviour*,  
561 142(4), 397–421. <https://doi.org/10.1163/1568539054012047>
- 562 de Boysson-Bardies, B. (1993). Ontogeny of Language-Specific Syllabic Productions. In B.  
563 de Boysson-Bardies, S. de Schonen, P. Jusczyk, P. McNeilage, & J. Morton (Eds.),  
564 *Developmental Neurocognition: Speech and Face Processing in the First Year of Life*  
565 (pp. 353–363). [https://doi.org/10.1007/978-94-015-8234-6\\_29](https://doi.org/10.1007/978-94-015-8234-6_29)
- 566 de Boysson-Bardies, B., & Vihman, M. M. (1991). Adaptation to Language: Evidence from  
567 Babbling and First Words in Four Languages. *Language*, 67(2), 297–319.  
568 <https://doi.org/10.2307/415108>



- 569 Dezechache, G., Zuberbühler, K., Davila-Ross, M., & Dahl, C. D. (2019). Machine learning  
570 reveals adaptive maternal responses to infant distress calls in wild chimpanzees.  
571 *BioRxiv*, 835827. <https://doi.org/10.1101/835827>
- 572 Elbers, L., & Ton, J. (1985). Play pen monologues: the interplay of words and babbles in  
573 the first words period. *Journal of Child Language*, 12(3), 551–565.  
574 <https://doi.org/10.1017/S0305000900006644>
- 575 Fedurek, P., Zuberbühler, K., & Dahl, C. D. (2016). Sequential information in a great ape  
576 utterance. *Scientific Reports*, 6, 38226. <https://doi.org/10.1038/srep38226>
- 577 Fischer, J. (2016). Playback Experiments. *The International Encyclopedia of Primatology*,  
578 1–2. <https://doi.org/10.1002/9781119179313.wbprim0140>
- 579 Fischer, J., Noser, R., & Hammerschmidt, K. (2013). Bioacoustic field research: a primer  
580 to acoustic analyses and playback experiments with primates. *American Journal*  
581 *of Primatology*, 75(7), 643–663. <https://doi.org/10.1002/ajp.22153>
- 582 Fitch, W. T. (2018). The Biology and Evolution of Speech: A Comparative Analysis.  
583 *Annual Review of Linguistics*, 4(1), 255–279. [https://doi.org/10.1146/annurev-](https://doi.org/10.1146/annurev-linguistics-011817-045748)  
584 [linguistics-011817-045748](https://doi.org/10.1146/annurev-linguistics-011817-045748)
- 585 Fitch, W. T., Boer, B. de, Mathur, N., & Ghazanfar, A. A. (2016). Monkey vocal tracts are  
586 speech-ready. *Science Advances*, 2(12), e1600723.  
587 <https://doi.org/10.1126/sciadv.1600723>
- 588 Fröhlich, M., & Hobaiter, C. (2018). The development of gestural communication in great  
589 apes. *Behavioral Ecology and Sociobiology*, 72(12), 194.  
590 <https://doi.org/10.1007/s00265-018-2619-y>
- 591 Fröhlich, M., Wittig, R. M., & Pika, S. (2019). The ontogeny of intentional communication  
592 in chimpanzees in the wild. *Developmental Science*, e12716. doi:  
593 [10.1111/desc.12716](https://doi.org/10.1111/desc.12716)

- 594 Ghazanfar, A. A., Liao, D. A., & Takahashi, D. Y. (2019). Volition and learning in primate  
595 vocal behaviour. *Animal Behaviour*, *151*, 239–247.  
596 <https://doi.org/10.1016/j.anbehav.2019.01.021>
- 597 Ghazanfar, A. A., & Takahashi, D. Y. (2014). The evolution of speech: vision, rhythm,  
598 cooperation. *Trends in Cognitive Sciences*, *18*(10), 543–553.  
599 <https://doi.org/10.1016/j.tics.2014.06.004>
- 600 Goodall, J. (1986). *The chimpanzees of Gombe: Patterns of behavior*. Cambridge, MA:  
601 Harvard University Press.
- 602 Goupil, L., Johansson, P., Hall, L., & Aucouturier, J.-J. (2019). *Influence of Vocal Feedback*  
603 *on Emotions Provides Causal Evidence for the Self-Perception Theory*.  
604 <https://doi.org/10.1101/510867>
- 605 Hayes, K. J., & Hayes, C. (1951). The Intellectual Development of a Home-Raised  
606 Chimpanzee. *Proceedings of the American Philosophical Society*, *95*(2), 105–109.
- 607 Hrdy, S. (2007). Evolutionary context of human development: the cooperative breeding  
608 model. In *Family Relationships: An Evolutionary Perspective*. Oxford: Oxford  
609 University Press.
- 610 Hrdy, S. (2009). *Mothers and others*. Cambridge, MA: Harvard University Press.
- 611 Hsu, M. J., Chen, L.-M., & Agoramorthy, G. (2005). The vocal repertoire of Formosan  
612 macaques, *Macaca cyclopis*: acoustic structure and behavioral context. *Zoological*  
613 *Studies*, *44*(2), 275.
- 614 Jhang, Y., & Oller, D. K. (2017). Emergence of Functional Flexibility in Infant  
615 Vocalizations of the First 3 Months. *Frontiers in Psychology*, *8*.  
616 <https://doi.org/10.3389/fpsyg.2017.00300>

- 617 Kersken, V., Zuberbühler, K., & Gomez, J.-C. (2017). Listeners can extract meaning from  
618 non-linguistic infant vocalisations cross-culturally. *Scientific Reports*, 7,  
619 srep41016. <https://doi.org/10.1038/srep41016>
- 620 Kojima, S. (2003). *A search for the origins of human speech: Auditory and vocal functions*  
621 *of the chimpanzee*. Kyoto: Kyoto University Academic Press.
- 622 Kramer, K. L. (2010). Cooperative Breeding and its Significance to the Demographic  
623 Success of Humans. *Annual Review of Anthropology*, 39(1), 417–436.  
624 <https://doi.org/10.1146/annurev.anthro.012809.105054>
- 625 Lameira, A. R., & Call, J. (2018). Time-space–displaced responses in the orangutan vocal  
626 system. *Science Advances*, 4(11), eaau3401.  
627 <https://doi.org/10.1126/sciadv.aau3401>
- 628 Lameira, A. R., & Shumaker, R. W. (2019). Orangutans show active voicing through a  
629 membranophone. *Scientific Reports*, 9, 12289. [https://doi.org/10.1038/s41598-](https://doi.org/10.1038/s41598-019-48760-7)  
630 [019-48760-7](https://doi.org/10.1038/s41598-019-48760-7)
- 631 Laporte, M. N. C., & Zuberbühler, K. (2011). The development of a greeting signal in wild  
632 chimpanzees. *Developmental Science*, 14(5), 1220–1234.  
633 <https://doi.org/10.1111/j.1467-7687.2011.01069.x>
- 634 Levréro, F., & Mathevon, N. (2013). Vocal Signature in Wild Infant Chimpanzees.  
635 *American Journal of Primatology*, 75(4), 324–332.  
636 <https://doi.org/10.1002/ajp.22108>
- 637 Lieberman, P. (2017). Comment on “Monkey vocal tracts are speech-ready.” *Science*  
638 *Advances*, 3(7), e1700442. <https://doi.org/10.1126/sciadv.1700442>
- 639 Locke, J. L. (2006). Parental selection of vocal behavior. *Human Nature*, 17(2), 155–168.  
640 <https://doi.org/10.1007/s12110-006-1015-x>

- 641 Mandel, M. I., & Ellis, D. P. (2005). Song-level features and support vector machines for  
642 music classification. *Proceedings of the 6th International Conference on Music*  
643 *Information Retrieval (ISMIR)*, 594–599.
- 644 McCune, L., Vihman, M. M., Roug-Hellichius, L., Delery, D. B., & Gogate, L. (1996). Grunt  
645 communication in human infants (*Homo sapiens*). *Journal of Comparative*  
646 *Psychology*, 110(1), 27.
- 647 Mielke, A., & Zuberbühler, K. (2013). A method for automated individual, species and call  
648 type recognition in free-ranging animals. *Animal Behaviour*, 86(2), 475–482.  
649 <https://doi.org/10.1016/j.anbehav.2013.04.017>
- 650 Møller, A. P. (1988). False Alarm Calls as a Means of Resource Usurpation in the Great Tit  
651 *Parus major*. *Ethology*, 79(1), 25–30. [https://doi.org/10.1111/j.1439-](https://doi.org/10.1111/j.1439-0310.1988.tb00697.x)  
652 [0310.1988.tb00697.x](https://doi.org/10.1111/j.1439-0310.1988.tb00697.x)
- 653 Nathani, S., Ertmer, D. J., & Stark, R. E. (2006). Assessing vocal development in infants  
654 and toddlers. *Clinical Linguistics & Phonetics*, 20(5), 351–369.  
655 <https://doi.org/10.1080/02699200500211451>
- 656 Oller, D. K. (2000). *The Emergence of the Speech Capacity*.  
657 <https://doi.org/10.4324/9781410602565>
- 658 Oller, D. K., Buder, E. H., Ramsdell, H. L., Warlaumont, A. S., Chorna, L., & Bakeman, R.  
659 (2013). Functional flexibility of infant vocalization and the emergence of  
660 language. *Proceedings of the National Academy of Sciences of the United States of*  
661 *America*, 110(16), 6318–6323. <https://doi.org/10.1073/pnas.1300337110>
- 662 Oller, D. K., & Griebel, U. (2004). Contextual freedom in human infant vocalization and  
663 the evolution of language. In *Evolution of communicative flexibility: Complexity,*  
664 *creativity and adaptability in human and animal communication* (p. 135).  
665 Cambridge, MA: MIT Press.

- 666 Oller, D. K., Griebel, U., Iyer, S. N., Jhang, Y., Warlaumont, A. S., Dale, R., & Call, J. (2019).  
667 Language Origins Viewed in Spontaneous and Interactive Vocal Rates of Human  
668 and Bonobo Infants. *Frontiers in Psychology, 10*.  
669 <https://doi.org/10.3389/fpsyg.2019.00729>
- 670 Oller, D. K., Griebel, U., & Warlaumont, A. S. (2016). Vocal Development as a Guide to  
671 Modeling the Evolution of Language. *Topics in Cognitive Science, 8*(2), 382–392.  
672 <https://doi.org/10.1111/tops.12198>
- 673 Oller, D. K., Wieman, L. A., Doyle, W. J., & Ross, C. (1976). Infant babbling and speech.  
674 *Journal of Child Language, 3*(1), 1–11.  
675 <https://doi.org/10.1017/S0305000900001276>
- 676 Owren, M. J., & Casale, T. M. (1994). Variations in fundamental frequency peak position  
677 in Japanese macaque (*Macaca fuscata*) coo calls. *Journal of Comparative*  
678 *Psychology, 108*(3), 291. <https://doi.org/10.1037/0735-7036.108.3.291>
- 679 Plooij, F. X. (1984). *The behavioral development of free-living chimpanzee babies and*  
680 *infants*. Norwood, NJ: Ablex.
- 681 Pollick, A. S., & Waal, F. B. M. de. (2007). Ape gestures and language evolution.  
682 *Proceedings of the National Academy of Sciences, 104*(19), 8184–8189.  
683 <https://doi.org/10.1073/pnas.0702624104>
- 684 Ponsot, E., Burred, J. J., Belin, P., & Aucouturier, J.-J. (2018). Cracking the social code of  
685 speech prosody using reverse correlation. *Proceedings of the National Academy of*  
686 *Sciences, 115*(15), 3972–3977. <https://doi.org/10.1073/pnas.1716090115>
- 687 R Core Team. (2018). *R: A language and environment for statistical computing*. Retrieved  
688 from <https://www.R-project.org/>



- 689 Range, F., & Fischer, J. (2004). Vocal Repertoire of Sooty Mangabeys (*Cercocebus*  
690 *torquatus atys*) in the Taï National Park. *Ethology*, *110*(4), 301–321.  
691 <https://doi.org/10.1111/j.1439-0310.2004.00973.x>
- 692 Rendall, D., Seyfarth, R. M., Cheney, D. L., & Owren, M. J. (1999). The meaning and  
693 function of grunt variants in baboons. *Animal Behaviour*, *57*(3), 583–592.  
694 <https://doi.org/10.1006/anbe.1998.1031>
- 695 Rendall, D., & Owren, M. J. (2002). Animal vocal communication: Say what? In *The*  
696 *cognitive animal: Empirical and theoretical perspectives on animal cognition* (pp.  
697 307–313). Cambridge, MA, US: MIT Press.
- 698 Reynolds, V. (2005). *The chimpanzees of the Budongo Forest: ecology, behaviour, and*  
699 *conservation*. Retrieved from  
700 [http://books.google.fr/books?hl=fr&lr=&id=C6hzM5lQJ6YC&oi=fnd&pg=PR11&](http://books.google.fr/books?hl=fr&lr=&id=C6hzM5lQJ6YC&oi=fnd&pg=PR11&dq=budongo+reynolds&ots=0OfjtMfycP&sig=X1c6kEzGs8ZlzhXdNH3aK3foPrA)  
701 [dq=budongo+reynolds&ots=0OfjtMfycP&sig=X1c6kEzGs8ZlzhXdNH3aK3foPrA](http://books.google.fr/books?hl=fr&lr=&id=C6hzM5lQJ6YC&oi=fnd&pg=PR11&dq=budongo+reynolds&ots=0OfjtMfycP&sig=X1c6kEzGs8ZlzhXdNH3aK3foPrA)
- 702 RStudio Team. (2015). RStudio: integrated development for R. *RStudio, Inc., Boston, MA*  
703 *URL [Http://Www. Rstudio. Com](http://www.rstudio.com), 42.*
- 704 Salmi, R., Hammerschmidt, K., & Doran-Sheehy, D. M. (2013). Western Gorilla Vocal  
705 Repertoire and Contextual Use of Vocalizations. *Ethology*, *119*(10), 831–847.  
706 <https://doi.org/10.1111/eth.12122>
- 707 Schaik, C. P. van, & Burkart, J. M. (2010). Mind the Gap: Cooperative Breeding and the  
708 Evolution of Our Unique Features. In *Mind the Gap* (pp. 477–496).  
709 [https://doi.org/10.1007/978-3-642-02725-3\\_22](https://doi.org/10.1007/978-3-642-02725-3_22)
- 710 Seyfarth, R. M., & Cheney, D. L. (1984). The acoustic features of vervet monkey grunts.  
711 *The Journal of the Acoustical Society of America*, *75*(5), 1623–1628.  
712 <https://doi.org/10.1121/1.390872>

- 713 Slocombe, K. E., & Zuberbühler, K. (2010). Vocal communication in chimpanzees. *The*  
714 *Mind of the Chimpanzee: Ecological and Experimental Perspectives*. University of  
715 *Chicago Press, Chicago*, 192–207.
- 716 Slocombe, K. E., & Zuberbühler, K. (2005). Functionally Referential Communication in a  
717 Chimpanzee. *Current Biology*, *15*(19), 1779–1784.  
718 <https://doi.org/10.1016/j.cub.2005.08.068>
- 719 Tajiri, Y., Yabuwaki, R., Kitamura, T., & Abe, S. (2010). Feature Extraction Using Support  
720 Vector Machines. In K. W. Wong, B. S. U. Mendis, & A. Bouzerdoum (Eds.), *Neural*  
721 *Information Processing. Models and Applications* (pp. 108–115).  
722 [https://doi.org/10.1007/978-3-642-17534-3\\_14](https://doi.org/10.1007/978-3-642-17534-3_14)
- 723 Takahashi, D. Y., Fenley, A. R., Teramoto, Y., Narayanan, D. Z., Borjon, J. I., Holmes, P., &  
724 Ghazanfar, A. A. (2015). The developmental dynamics of marmoset monkey vocal  
725 production. *Science*, *349*(6249), 734–738.  
726 <https://doi.org/10.1126/science.aab1058>
- 727 Vert, J.-P., Tsuda, K., & Schölkopf, B. (2004). A primer on kernel methods. In *Kernel*  
728 *methods in computational biology* (Vol. 47, pp. 35–70). Cambridge, MA: MIT Press.
- 729 Waal, F. B. M. de, & Pollick, A. S. (2011). Gesture as the most flexible modality of primate  
730 communication. *The Oxford Handbook of Language Evolution*.  
731 <https://doi.org/10.1093/oxfordhb/9780199541119.013.0006>
- 732 Watson, S. K., Townsend, S. W., Schel, A. M., Wilke, C., Wallace, E. K., Cheng, L., ...  
733 Slocombe, K. E. (2015). Vocal Learning in the Functionally Referential Food  
734 Grunts of Chimpanzees. *Current Biology*, *25*(4), 495–499.  
735 <https://doi.org/10.1016/j.cub.2014.12.032>
- 736 Wheeler, B. C. (2009). Monkeys crying wolf? Tufted capuchin monkeys use anti-predator  
737 calls to usurp resources from conspecifics. *Proceedings of the Royal Society B:*

- 738 *Biological Sciences*, 276(1669), 3013–3018.
- 739 <https://doi.org/10.1098/rspb.2009.0544>
- 740 Williams, C. E., & Stevens, K. N. (1972). Emotions and Speech: Some Acoustical
- 741 Correlates. *The Journal of the Acoustical Society of America*, 52(4B), 1238–1250.
- 742 <https://doi.org/10.1121/1.1913238>
- 743 Yoo, H., Bowman, D. A., & Oller, D. K. (2018). The Origin of Protoconversation: An
- 744 Examination of Caregiver Responses to Cry and Speech-Like Vocalizations.
- 745 *Frontiers in Psychology*, 9. <https://doi.org/10.3389/fpsyg.2018.01510>
- 746 Zuberbühler, K. (2012). Cooperative breeding and the evolution of vocal flexibility. In
- 747 *The Oxford Handbook of Language Evolution*. Oxford: Oxford University Press.
- 748 Zuberbühler, K. (2014). Experimental field studies with non-human primates. *Current*
- 749 *Opinion in Neurobiology*, 28, 150–156.
- 750 <https://doi.org/10.1016/j.conb.2014.07.012>

751

752 **ACKNOWLEDGMENTS**

753 We thank UWA and UNCST for permission to conduct the study, Geoffrey Muhanguzi,

754 Caroline Asimwe and Sam Adué for their support in the field, Derry Taylor for critical

755 comments on previous versions of the manuscript.

756 We are grateful to the Royal Zoological Society of Scotland for providing core funding to the

757 Budongo Conservation Field Station.

758 The research was supported by a Fyssen Fellowship and British Academy Newton International

759 Fellowship awarded to GD, funding from the European Union's Seventh Framework

760 Programme for research, technological development and demonstration (grant agreement no

761 283871), and the Swiss National Science Foundation (PZ00P3\_154741) awarded to CDD.

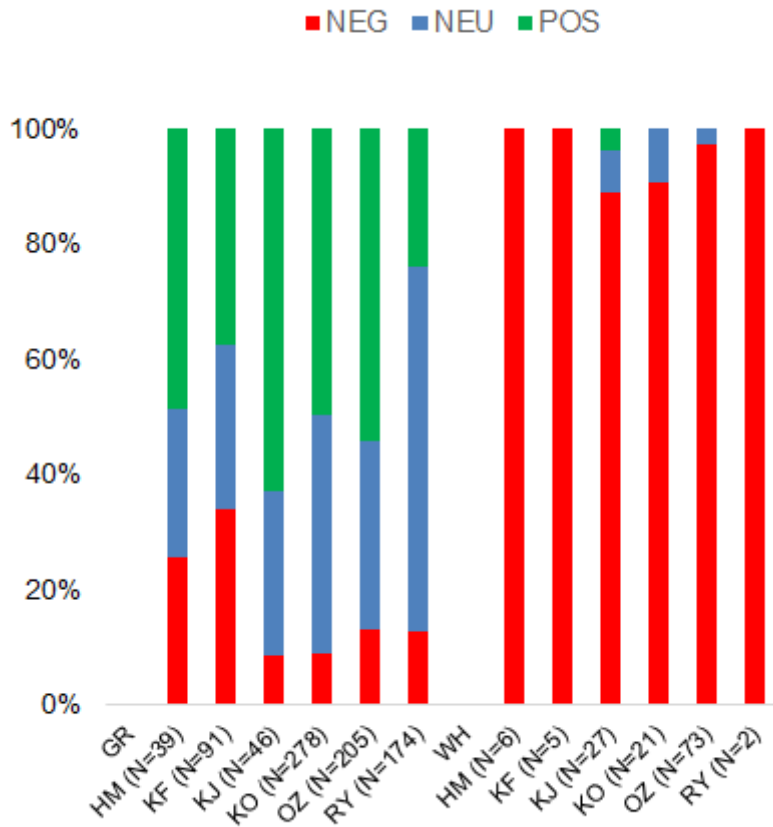
762

763 **CONFLICTS OF INTEREST**

764 No conflicts of interest.

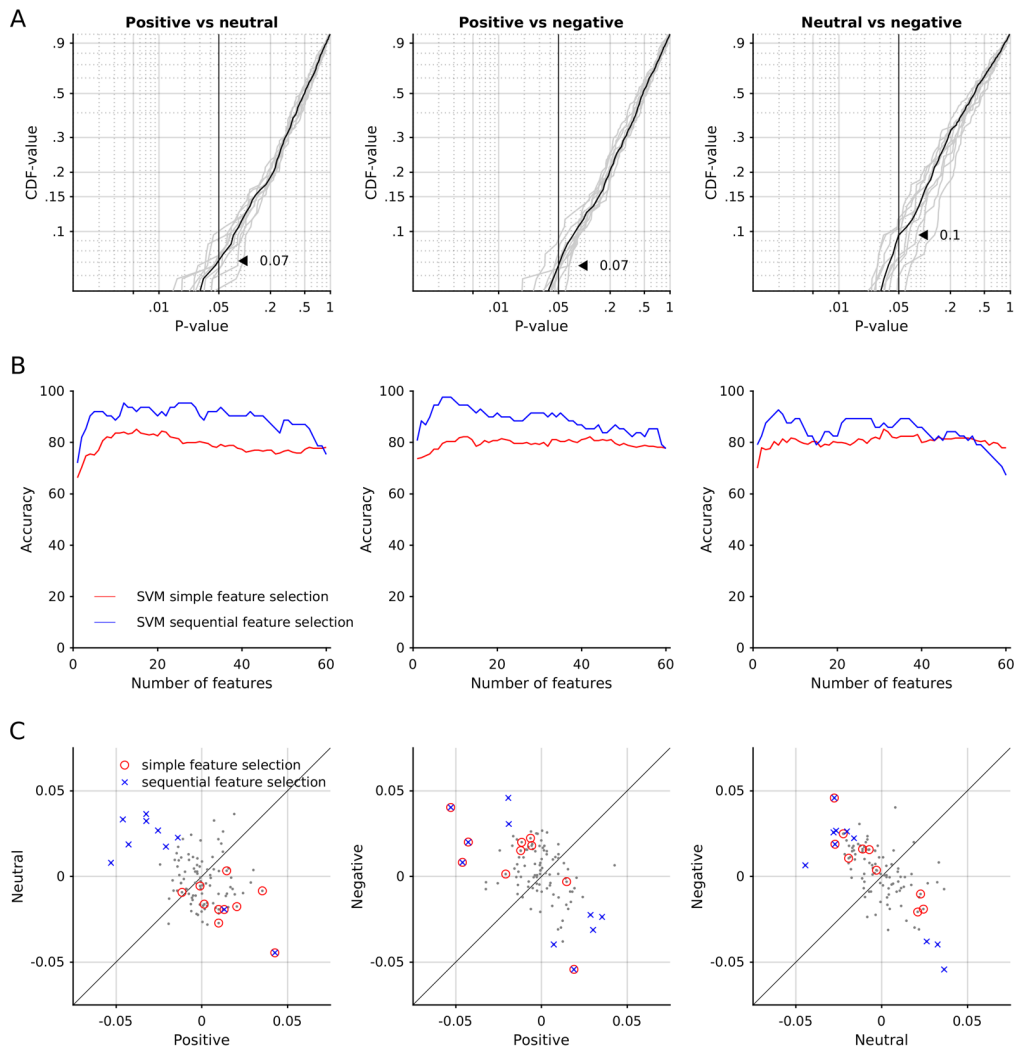
765 **FIGURES**

766 **Figure 1** Proportion of grunt-like (GR) and whimper-like (WH) vocal behaviors recorded in  
767 negative (NEG), neutral (NEU) and positive (POS) affective contexts, for each individual  
768 separately. Number between brackets indicate the number of GR and WH calls contributed by  
769 each individual.  
770



771  
772  
773

774 **Figure 2** Feature selection and classification performances. The columns represent the  
775 comparisons of affective contexts during which the vocal utterance occurred.  
776 A. For each feature dimensions the discrimination power of the two classes (e.g. positive vs.  
777 neutral) was evaluated using a t-test. P-values are shown as an empirical cumulative distribution  
778 function (eCDF). Gray lines show the results of individual runs of evaluation; black lines show  
779 the means of individual runs. Indicated with arrow heads are the proportions of feature  
780 dimensions that significantly discriminate between the two classes tested.  
781 B. The classification performances are shown for the SVM classifier relying on feature  
782 dimensions extracted through a simple feature selection (red lines) and a sequential feature  
783 selection procedure (blue lines).  
784 C. Feature selection outcomes are shown for simple (red circles) and sequential feature  
785 selection (blue x-s) as overlays on all feature dimensions (black dots).



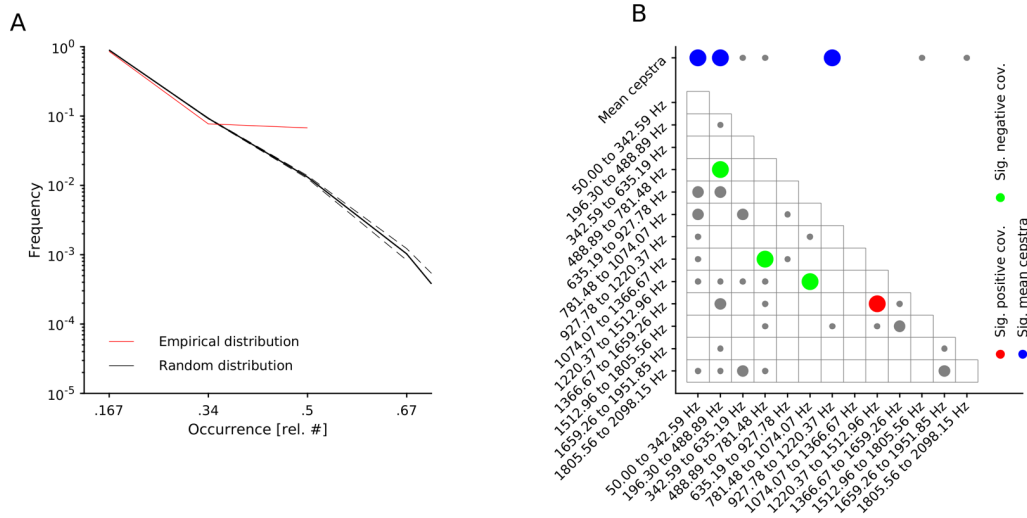
786  
787



788 **Figure 3** Overall feature importance.

789 A. The empirical distribution of feature dimensions across all comparisons.

790 B. Significant feature dimensions are shown in colors, according to their sign: in red positive  
 791 covariances, in green negative covariance. The means of cepstra are shown in blue. The marker  
 792 size indicates the occurrence: small = 1, medium-large = 2, large = 3 (significant). Gray-colored  
 793 markers are non-significant feature dimensions.



794  
 795

796 **TABLES**

797 **Table 1** List of focal animals, with their name (ID), sex and minimum and maximum age in  
 798 months. Also given are the number of grunt-like and whimper-like vocal behaviors collected,  
 799 as well as grunt-like vocalizations acoustically analyzed.

800

<i>ID</i>	<i>Sex</i>	<i>Min. Age</i> <i>(in</i> <i>months)</i>	<i>Max. Age</i> <i>(in</i> <i>months)</i>	<i>N</i>	<i>whimper-like</i> <i>vocalizations</i>	<i>N</i>	<i>grunt-like</i> <i>vocalizations</i>	<i>N</i>	<i>of grunt-like</i> <i>vocalizations used in</i> <i>acoustical analysis</i>
HM	F	3.41	6.85	6		39		10	
KF	M	<1	11.87	5		91		20	
KJ	M	6.98	10.52	27		46		7	
KO	M	3.08	8.46	21		278		67	
OZ	M	1.38	8.16	73		205		32	
RY	M	4.75	8.16	2		174		44	

801

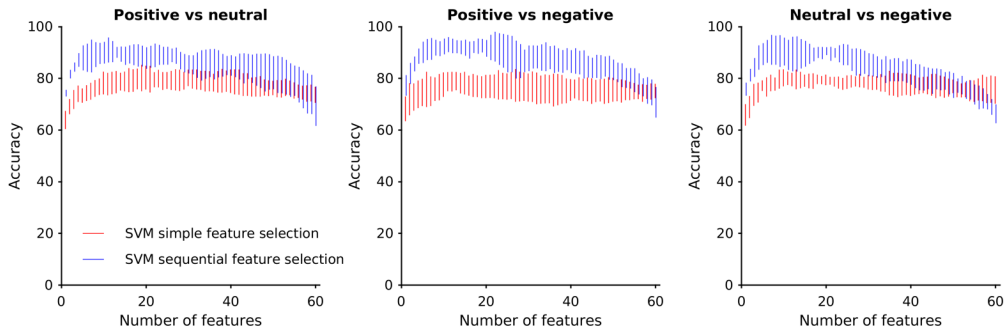
802 **Table 2** Affective coding of infant behavior

Behavior	Description	Affect
Play	Relaxed moving without obvious purpose. Can be solitary or social	POSITIVE
Grooming	Giving or receiving 'grooming'	POSITIVE
Feeding	Breastfeeding or swallowing an edible element	POSITIVE
Social approach	Greeting a conspecific whilst moving towards this individual	POSITIVE
Resting	Remaining within a limited area, may involve some degree of moving around, marked by relative idleness	NEUTRAL
Moving	Simple movements	NEUTRAL
Manipulating objects	Manipulating objects (leaves, branches, rocks)	NEUTRAL
Greeting without approach	Calling upon the approach of a conspecific without showing approach or avoidance behavior towards it	NEUTRAL
Nuzzling	Unsuccessfully trying to access the mother's nipple	NEGATIVE
Begging	Unsuccessfully attempting to access food other than breast milk	NEGATIVE
Hiding	Increased gripping or seeking contact with the mother when contact already established between them	NEGATIVE
Contact mother	Seeking contact with the mother when contact not established between them	NEGATIVE
Escaping	Showing escape movements, away from an activity involved in; could be associated with discomfort	NEGATIVE

803

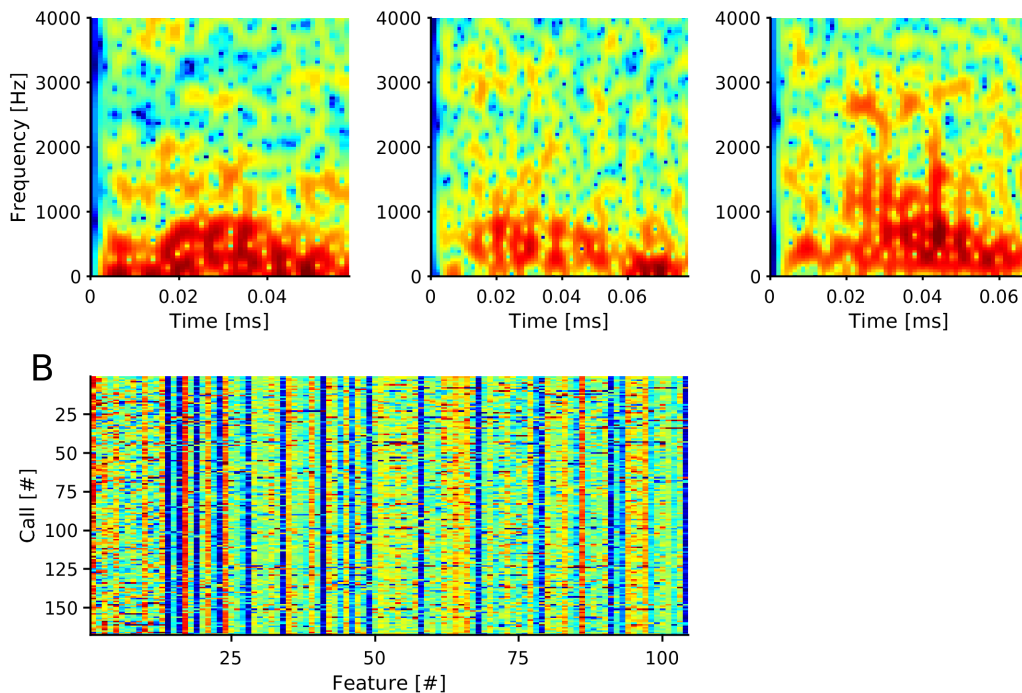
804 **SUPPLEMENTARY INFORMATION**

805 **Supplementary Figure 1** Leave-one-out method to account for subject effects. The accuracies  
806 of the three comparisons of grunt types are shown as function of number of features. These  
807 graphs illustrate the variability of accuracy caused by leaving out one of the 6 individuals per  
808 each separate classification procedure. The vertical bars indicate the minimum and maximum  
809 scores.



810  
811

812 **Supplementary Figure 2** MFCCs extracted from example calls and extracted feature matrix.  
813 A. Time-frequency spectra of three arbitrarily chosen calls.  
814 B. From each call 26 spectral bands and 13 cepstra were extracted. Feature vectors containing  
815 the means and covariances of cepstra are shown for each call. Means are shown as features 1  
816 to 13 on the x-axis, followed by covariances (91 values).



817