1    **Flexibility in wild infant chimpanzee vocal behavior**

2    Guillaume Dezecache[a,b,c,d,*], Klaus Zuberbühler[a,b,e], Marina Davila-Ross[c] & Christoph D.

3    Dahl[a,f,g,*]

4

5    [a]Institute of Biology, University of Neuchâtel, Neuchâtel, Switzerland;

6    [b]Budongo Conservation Field Station, Masindi, Uganda;

7    [c]Department of Psychology, University of Portsmouth, Portsmouth, England, United

8    Kingdom;

9    [d]Université Clermont Auvergne, CNRS, LAPSCO, Clermont-Ferrand, France;

10   [e]School of Psychology and Neuroscience, University of St Andrews, St Andrews, Scotland,

11   United Kingdom;

12   [f]Graduate Institute of Mind, Brain and Consciousness, Taipei Medical University, Taipei,

13   Taiwan;

14   [g]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>, LAPSCO CNRS 17 Rue Paul

19   Collomp, 63000 Clermont-Ferrand, France;

20   Christoph D. Dahl <christoph.d.dahl@gmail.com>

21

**Flexibility in wild infant chimpanzee vocal behavior**

**ABSTRACT**

How did human language evolve from earlier forms of communication? One way to address this question is to compare prelinguistic human vocal behavior with nonhuman primate calls. An important finding has been that, prior to speech and from early on, human infant vocal behavior exhibits functional flexibility, or the capacity to produce sounds that are not tied to one specific function. This is reflected in human infants' use of single categories of protophones (precursors of speech sounds) in various affective circumstances, such that a given call type can occur in and express positive, neutral, or negative affective states, depending on the occasion. Nonhuman primate vocal behavior, by contrast, is seen as comparably inflexible, with different call types tied to specific functions and sometimes to specific affective states (e.g., screams mostly occur in negative circumstances). As a first step towards addressing this claim, we examined the vocal behavior of six wild infant chimpanzees during their first year of life. We found that the most common vocal signal, grunts, occurred in a range of contexts that were deemed positive, neutral and negative. Using automated feature extraction and supervised learning algorithms, we also found acoustic variants of grunts produced in the affective contexts, suggesting gradation within this vocal category. By contrast, the second most common call type of infant chimpanzees, the whimpers, was produced in only one affective context, in line with standard models of nonhuman primate vocal behavior. Insofar as our affective categorization reflect infants' true affective state, our results suggest that the most common chimpanzee vocalization, the grunt is not affectively bound. Affective decoupling is a prerequisite for chimpanzee grunts (and other vocal categories) to be deemed 'functionally flexible'. If later confirmed to be a functionally flexible vocal type, this would indicate that the evolution of this foundational vocal capability occurred before the split between the Homo and Pan lineages.

## 1. INTRODUCTION

At some point in evolutionary history, there must have been a transition from primate-like to human-like acoustic communication, which may have coincided with the origins of speech. The evolutionary history of this transition continues to be vividly debated (Fitch, 2018), with a large range of comparative evidence from animal communication systems, and the consensus view that direct evolutionary homologies are generally absent in the primate order (Rendall & Owren, 2002). More recently, however, some vocal and neural equipment has been identified in different primate species that allow for the production of speech-like sounds (Boë et al., 2017; Fitch et al., 2016; Lieberman, 2017) and for some degree of control over vocal fold oscillation (Lameira & Shumaker, 2019). If the facial and gestural displays have undeniably played a crucial role in the evolution of language (Arbib et al., 2008; Pollick & Waal, 2007; Waal & Pollick, 2011), vocal production appears so strongly different in humans and other primates that the quest for evolutionary precursors of human vocal production has been and will continue to be particularly fruitful.

One key divergence between humans and other species, according to Oller and colleagues (2013), lie in the ontogenetic trajectories between non-human primate and human vocal behavior prior to speech. By the age of one month old (and possibly earlier, see Oller et al., 2019), human infants produce sounds that are not tied to the expression of one particular need, such that they can hold different illocutionary force on different occasions, and cause appropriate reactions in caregivers (Jhang & Oller, 2017; Oller et al., 2013). This is reflected in infants' use of squeals, vocants and growls in circumstances associated with positive, negative, or neutral affective states, such that those sounds are not bound to the experience of one particular type of affect (Oller et al., 2013). By contrast to those vocal types, human infants use laughter and cries in single affective contexts (positive and negative, respectively).

73 This capacity to produce one vocal unit under a variety of affective state (such that vocal

74 production is independent from the experience of a specific affective state – a capacity termed

75 'affective decoupling') later allows infants to use those sounds as they wish, and to express a

76 variety of needs on different occasions (Oller et al., 2013; Oller & Griebel, 2004). This

77 decoupling between vocal production and affective experience is foundational for the ability to

78 produce sounds that can later fulfil a variety of functions ('vocal functional flexibility'), that is,

79 they can be used to invite a variety of responses in others (Oller & Griebel, 2004). For instance,

80 a given utterance (such as 'the train is arriving') produced with neutral (a mere statement),

81 negative (annoyance) and positive (exultation) tones have the same syntactic structure and

82 semantic content, but are produced under antagonist affective states and cause vastly different

83 responses in receivers. Without affective decoupling and vocal functional flexibility, basic

84 speech acts cannot take place.

85

86 This decoupling of signal structure and affect in young infants' vocal repertoire  has thus been

87 identified as a major evolutionarily precursor to language (Oller et al., 2013). Because of their

88 early ontogenetic onset, affective decoupling and vocal functional flexibility may be more

89 foundational to human speech than other building blocks of the language faculty, such as proto-

90 syntax or vocal elaboration (Oller et al., 2013). These capacities, in this view, are prerequisites

91 for speech development, and major evolutionary departures from the affect-bound and

92 functionally inflexible vocal behavior of non-human primates (Waal & Pollick, 2011). By

93 contrast to their primate cousins, humans would have evolved in a social ecology conducive to

94 the development of such vocal flexibility. Notably, it is proposed that functionally flexible

95 vocalizations of young human infants have evolved in humans in relation to allo-maternity

96 (Burkart et al., 2009, 2009, 2009; Hrdy, 2007; Kramer, 2010; Schaik & Burkart, 2010) or

97 altriciality (Locke, 2006) and associated pressures on young infants to signal their needs and

98    attract caregivers (Ghazanfar et al., 2019; Locke, 2006; Zuberbühler, 2012). Other species

99    living in cooperative breeding systems (such as the marmosets (Burkart et al., 2007)) may

100   display vocal functional flexibility in their vocal repertoire.

101

102   For long, primate (but also animal) communication systems have been characterized as

103   affectively-biased, such that affect has been proposed to be both necessary and sufficient for

104   vocal production to occur. According to Hammerschmidt and Fischer, there could even exist

105   "[a] correspondence between non-verbal vocalizations in humans and non-human primates

106   [such] that they both function to communicate the affective state of the signaler."

107   (Hammerschmidt & Fischer, 2008, p. 103). In fact, a number of researchers have held the view

108   that the equivalents of animal vocalizations are non-verbal affective expressions in humans,

109   such as laughing, screaming and crying (see Gruber & Grandjean, 2017 and Marler, 1980 for a

110   discussion).  Examining the neural pathways of vocal production in squirrel monkeys, Jürgens

111   (Jürgens, 1976; Jürgens, 1979) concluded that vocal production was mediated by affect. More

112   recently in marmoset infants, Zhang & Ghazanfar (Zhang & Ghazanfar, 2016) found that

113   fluctuations in cardiac rhythm shape respiratory patterns, which in turn contribute to vocal

114   production, thereby attributing a central role to affect in early vocal production of this species

115   (Tchernichovski & Oller, 2016). The vocal repertoire of species phylogenetically closer to

116   humans (such as the chimpanzees) did not evade this conclusion. In her authoritative book on

117   the behavior of Gombe chimpanzees, Goodall (Goodall, 1986) wrote that 'chimpanzee

118   vocalizations are closely bound to emotion' and that 'the production of a sound in the *absence*

119   of the appropriate emotional state seems to be an almost impossible task for a chimpanzee' (p.

120   125). As a consequence, she proposed a mapping between call and affect when describing the

121   vocal repertoire of chimpanzees, with, for instance, a one-to-one correspondence between the

122   experience of annoyance and the production of 'soft barks' (p .127).

123

124  To which degree is vocal production affectively bound in other species? To which extent can

125  the developmental trajectory seen in humans (with early independence between certain sounds

126  and particular affective states (Oller et al., 2013)) also be observed in other primates? In fact,

127  are affective decoupling and vocal functional flexibility unique to human vocal ontogeny? In

128  one relevant study, Clay et al. (Clay et al., 2015) examined 'peep' calls in mature bonobos (*Pan*

129  *paniscus*), their most common vocalizations, and found that they are produced in a variety of

130  contexts, ranging from seemingly positive (food provisioning) to neutral (travel and resting)

131  and negative (agonistic and alarm) situations. Based on these findings, the authors concluded

132  that bonobos have the capability to produce sounds that are not affectively biased (Clay et al.,

133  2015), suggestive of affective decoupling in vocal production. Their peeps were, however,

134  attributed to broad behavioral contexts (such as feeding or travelling) with no focus on more

135  specific and transient behaviors that may help infer affective contexts, such as when individuals

136  suddenly experience aggression during travelling and feeding bouts. In fact, peeps could well

137  be bound to the expression of one particular affect, which could be common in both feeding

138  and travelling contexts for example. As such, the bonobo data are indicative of their peeps

139  occurring across broad behavioral contexts ('contextual flexibility') but may ultimately remain

140  inconclusive in regards to whether affective decoupling and vocal functional flexibility are

141  indeed present in species other than humans.

142

143  Similarly, the flexibility with which some call types are expressed in some primate species is

144  only *suggestive* of 'vocal functional flexibility' (the use of one vocal type to convey various

145  illocutionary forces on different utterances), and may only correspond to contextual flexibility

146  (the use of one call type in different contexts, with core commonalities in the illocutionary force

147  conveyed by all utterances). For example, Guinea baboons use a number of calls (e.g., grunts,

148  roar grunts, barks and wahoos) in a diversity of contexts (Maciej et al., 2013). Guinea baboons'

149  grunts are used in foraging and travelling contexts, but also affiliative, infant handling and

150  greeting contexts. Does that mean that Guinea baboons show functional flexibility when

151  producing grunts? It is a possibility. However, and in the absence of a methodological focus on

152  potential affective states experienced by the animal, a temporary conclusion is that Guinea

153  baboon grunts likely are 'contextually flexible'. The possibility that they also are not affectively

154  bound (i.e., not bound to the experience and expression of a particular affective state) or

155  'functionally flexible' (i.e., not assigned to the fulfilment of one particular function) awaits

156  empirical confirmation, for grunts in Guinea baboons could well be expressed under one

157  particular affective state, and be used to meet one single function in a variety of contexts (e.g.,

158  maintaining contact with other troop members). In fact, a first step could be made by examining

159  whether inferences about the affective state of animals (e.g., Guinea baboons) can be conducted

160  using the behavioral contexts employed to describe the contextual occurrence of their

161  vocalizations, and whether such analysis suggests that these vocalizations qualify as affectively

162  decoupled.

163

164  A second study, also on bonobos (Oller et al., 2019), suggests protophone-like vocal behavior

165  with bonobo infants producing calls that occur in both low or moderate arousal situations,

166  implying no affective binding. This conclusion has been preliminary, however, for the affective

167  quality of the contexts surrounding vocalizations (a reliable marker of illocutionary force and

168  needs in human infants) has proven difficult to discern.

169

170  Here, we intended to provide a first evaluation of affective decoupling in infant chimpanzees'

171  (*Pan troglodytes schweinfurthii*) vocal behavior at a very early age (< 12 months). Given the

172  recent studies in both immature and mature bonobos, focusing on the other closest living

173    relatives, the chimpanzees, is crucial to test hypotheses about the evolutionary origins of

174    functionally flexible vocal behavior. What's more, examination of *early* vocal production is

175    critical for a more direct comparison with findings on human infants (Oller et al., 2013). We

176    focused on two call types, the grunts and the whimpers, as they are acoustically very distinct

177    vocalization categories that are common in young infants (Plooij, 1984). Finally, we tried to

178    approach the affective dimension of the context of calling by focusing on transient behavioral

179    cues (e.g., the infant escaping a situation) rather than broader behavioral contexts (e.g.,

180    travelling context).

181

182    Grunt calls are of particular importance as they develop into a central component of the vocal

183    repertoire of chimpanzees and contribute to a variety of vocal sequences produced by juveniles,

184    sub-adults and adults (Crockford & Boesch, 2005). For example, grunts complement panting

185    elements during laughter (Leavens, 2009) and when encountering dominant individuals ('pant-

186    grunts') (Laporte & Zuberbühler, 2011; Laporte & Zuberbühler, 2010). They are also produced

187    upon encountering a food patch or when joining a foraging party ('rough grunts') (Fedurek &

188    Slocombe, 2013; Schel et al., 2013; Slocombe et al., 2010; Slocombe & Zuberbühler, 2005;

189    Watson et al., 2015). Finally, they are routinely produced throughout resting or in relaxed social

190    activities (Goodall, 1986). Grunts are produced from the first days of life in chimpanzees. Their

191    ontogenetic development has already been studied to some degree in chimpanzees, which has

192    shown some flexibility in usage (Laporte & Zuberbühler, 2011). It has been suggested that at

193    least two types of grunts could be distinguished. First, uh-grunts are short, tonal sounds,

194    resembling human vowels {u}, {o} and {a} (and possible homologous to quasi-vowels in

195    humans), sometimes produced in short series (staccato-grunts) (Kojima, 2003; Plooij, 1984).

196    The second type are the so-called 'effort' grunts, which are common in immature chimpanzees

197    (Plooij, 1984) and are also present in adult chimpanzees, mature and immature humans and

9

198    other mammals (McCune et al., 1996). So-called 'effort grunts', are very soft and require the

199    close presence of observers to be reliably heard (Plooij, 1984). They received their name from

200    their presence during locomotor activities. Despite Plooij's suggestion that they could be mere

201    by-products of locomotor activities, he also noted they can occur in the absence of movements

202    (Plooij, 1984). So far, no study has yet offered an acoustical validation of the existence of these

203    diverse types, such that we (and others, see Laporte & Zuberbühler, 2011) cannot rely on this

204    distinction.

205

206    Another common vocal utterance produced by chimpanzee infants is whimpers (Dezecache et

207    al., 2020; Levréro & Mathevon, 2013; Plooij, 1984). They are short, tonal and often produced

208    in series with an upward shift in fundamental frequency. Contrarily to grunts, whimpers

209    preferentially occur in aversive contexts, likely homologous to human crying or distress calls

210    in other mammals (Plooij, 1984). Previous research (e.g., Plooij, 1984) has suggested the

211    presence of whimper subtypes (single, serial and human-like whimpers), but again, we are not

212    aware of any systematic acoustical analysis that would justify this nomenclature. Whimpers are

213    also present in the repertoire of adult chimpanzees, notably in alarm (Tsukahara, 1993), food

214    begging (Crockford & Boesch, 2005; Slocombe & Newton-Fisher, 2005), and physical

215    separation (Crockford & Boesch, 2005) contexts.

216

217    To start addressing the hypothesis that affective decoupling and vocal functional flexibility

218    evolved before the split between *Pan* and *Homo* lineages, we examined the vocal behavior of

219    six wild chimpanzee infants aged between 0-12 months old from the Sonso community of

220    Budongo Forest, Uganda. We analyzed the extent to which vocal production of grunt-like and

221    whimper-like vocalizations were occurring with so-called positive, negative or neutral

222    behaviors, as a first step towards evaluating the affective quality of the vocalization contexts.

10

223    We also took advantage of recent developments of machine learning techniques to the study of

224    animal communication (Fedurek et al., 2016; Mielke & Zuberbühler, 2013) to evaluate

225    acoustical differences between calls produced with positive, negative and neutral markers.

226

227    2.  **METHODS**

228    *2.1 Ethics*

229    Permission to conduct the study was obtained from the Ugandan Wildlife Authority (UWA)

230    and the Uganda National Council for Science and Technology (UNCST).

231

232    *2.2 Subjects and data collection*

233    Data were collected in the Sonso community of the Budongo Forest Reserve, Uganda

234    (Reynolds, 2005) between February-June 2014, December 2014 and March-June 2015. This

235    community comprises around 70 individuals well habituated to human observers. The natural

236    behavior of N = 7 infants was video recorded continuously during focal animal sampling

237    (Altmann, 1974), using Panasonic HC X909/V700 cameras, with a Sennheiser MKE-400

238    shotgun microphone. Six of those infants produced enough calls to be further considered for

239    data analysis (see Table 1 for details).

240

241    *2.3 Behavioral data analysis*

242    Videos were inspected for the presence of infant vocalizations. We defined vocal behavior as

243    the occurrence of single sound units or series of sounds produced by the infant's vocal apparatus,

244    separated by a least 5 seconds of silence.

245

246    As of today, there is no definitive repertoire of infant chimpanzee vocal behaviors, only

247    suggestive classifications (Plooij, 1984; Plooij et al., 2014). The categories used in this research

248 are based on First Author's assessment. This assessment proved reliable when confronted to an

249 independent assessment with Derry Taylor, using vocalizations from infant and juvenile semi-

250 wild chimpanzees from the Chimfunshi Wildlife Orphanage, Zambia, collected by DT. One

251 hundred-and-sixty vocalizations were indeed classified as belonging to either the 'grunt',

252 'whimper', 'scream' or 'laughter' category. Agreement was excellent (k = 0.77) and even better

253 when considering only 'grunts' and 'whimpers' (k = 0.92).

254

255 For each vocal occurrence, we coded infant behavior from a list of mutually exclusive behaviors

256 (summarized in Table 2). This list was established following data collection, with some

257 inspiration from the behavioral categories established by Plooij during his study with the infant

258 chimpanzees of the Gombe community between 1971 and 1973 (Plooij, 1984). As in the

259 original human study (Oller et al., 2013), we reckoned the behavior of the infants could offer a

260 reliable source of information unto their affective state, as a first step towards establishing

261 affective descriptions of contexts. In fact, we originally aimed at mimicking their coding

262 strategy, using categories appropriate to the study of wild infant chimpanzees. The affective

263 quality of the infants' behavior was classified as 'positive' if it showed one of the following

264 four behaviors: (1) 'play' (2) giving or receiving 'grooming' (note that allo-grooming was never

265 observed in our infants); (3) 'feeding', and (4) 'social approach'. See Table 2 for details.

266

267 The affective context was classified as 'neutral' if it showed one of the following behaviors:

268 (5) 'resting'; (6) 'moving'; (7) 'manipulating objects' without playful postures, or (8) 'greeting

269 without approach'. See Table 2 for details.

270

271 Infant behavior was classified as 'negative' if it showed one of the following behaviors: (9)

272 'nuzzling'; (10) 'begging'; (11) 'hiding'; (12) 'contact mother/kin' was coded if infants were

273    urgently seeking contact with the mother or a kin when contact was not already established

274    between them; (13) 'escaping'. See Table 2 for details.

275

276    We performed intra-coder reliability tests on the affective contexts coded as positive, neutral

277    and negative. For this, we randomly selected 200 video clips (around 19% of the coded dataset

278    composed of the 7 infants), which were coded independently during two coding sessions more

279    than a year apart (November 2015 and February 2017), so that the second coding was, notably,

280    naïve. We found strong agreement between the two coding sessions (k = 0.73).

281

282    In order to evaluate the evenness of the distributions of grunts and whimpers across affective

283    contexts, we calculated, for each infant, and for grunts and whimpers separately, the dominance

284    of one affect over the two others, using the Berger-Parker Dominance index (see Morris et al.,

285    2014):

286

287    $$dominance = N_{max} / N$$

288

289    where $N_{max}$ is the number of calls in the most abundant affective context; N the total number of

290    calls across all affective contexts. Dominance values range from 1 / number of affects (=

291    equiprobability of calls across affects; here 1 / 3 = 0.33) to 1 (= complete dominance of one

292    affective context over the others).

293

294    Dominance values (one per infant per call type) were compared between grunts and whimpers

295    using a paired Wilcoxon Sign-Ranked test. These analyses were carried out using R (version

296    3.6.1; R Core Team, 2018) and R Studio (version 1.2.1335; RStudio Team, 2015).

297

298    *2.4 Acoustic analysis*

299    Acoustic data analysis focused on grunts for they were the only vocal category for which at

300    least two of the affective contexts were well represented. The acoustic structure of whimpers

301    has been analyzed as part of another study (see Dezecache et al., 2020). N = 180 grunts were

302    extracted from independent vocal behaviors. For each affective context, 60 were randomly

303    selected. Following extraction, we used MATLAB (MathWorks Inc., Natick, MA, USA) for

304    the acoustic data analysis, consisting of features extraction, feature selection and call

305    classification. We first pre-processed the audio files by applying a band pass filter from 50 to

306    4000 Hz and normalized the signals using the following function:

307

308              *signal = (signal - mean(signal)) / max(abs(signal - mean(signal)))*

309

310    *2.4.1 Feature extraction and selection*

311    We first ran a feature extraction algorithm to reduce redundancy of information and

312    computational efforts in classifying the grunts and to maximize the generalization ability of the

313    classifier (Tajiri et al., 2010). A popular method is extraction of mel frequency cepstral

314    coefficients (MFCCs) (Supplementary Figure 2). MFCCs represent the envelope of the short-

315    time power spectrum, as determined by the shape of the vocal tract (Logan, 2000). The idea

316    behind the extraction of MFCCs is to obtain a comprehensive representation of the frequencies

317    that compose an audio bout, while putting emphasis on certain frequency bands. While a typical

318    spectrogram linearly scales frequencies (i.e., each frequency bin is spaced an equal number of

319    Hertz apart), the mel-frequency scale is a logarithmical spacing of frequencies. MFCCs is

320    routinely used in speech recognition and is gaining prominence in the field of animal

321    communication (see for instance Fedurek et al., 2016 in chimpanzees). The use of MFCCs to

322    represent sounds can be considered to be a solution preferable to the selection of a limited set

14

323    of parameters to describe acoustical phenomena (such as these related to the shape of the

324    fundamental frequency) for it offers a more comprehensive representation of sounds. In the

325    context of our work (the aim of which was to evaluate potential distinctiveness between grunts

326    occurring in so-called positive, neutral and negative contexts), MFCCs appeared as the optimal

327    solution to the problem of a false negative conclusion.

328

329    We divided the calls into segments of 25ms length and 10ms steps between two successive

330    segments. We warped 26 spectral bands and returned 13 cepstra, which resulted in feature

331    dimensions of 13 values each. We then took the mean and co-variances of each cepstra over

332    the collection of feature segments, resulting in a 13-value vector and a 13 x 13-value matrix,

333    respectively, and concatenated to 104-unit vectors (Mandel & Ellis, 2005, p. 594-599) (Figure

334    3). We applied feature scaling to [0 to 1] and mean normalization.

335

336    Second, we performed a feature selection procedure: too many feature dimensions are not

337    useful for producing reliable classification systems, whereas low sample numbers can lead to

338    over-fitting to noisy feature dimensions. We therefore selected a subset of the original feature

339    dimensions and evaluated classification performance based on sequentially selected feature sets

340    until there was no improvement in performance. At this end, we subdivided the entire data set

341    into a training (75%) and a test data set (25%) and applied a $t$-test on each feature dimension,

342    comparing values of given feature dimension sorted by predefined class labels (e.g., grunts

343    occurring with negative (1) vs. positive (2) affects) and used $p$-values as a measure separability

344    of the two classes. We plotted the $p$-values as an empirical cumulative distribution function

345    (eCDF) to get an understanding of how well each feature separated the two classes and how

346    many features contributed to a significant separation (5%-level). We ran this procedure 20 times

347    for each comparison and plotted the results individually (gray lines) and the mean of all

15

348    repetitions (black line) (Figure 2A). The classification routines were then independently run

349    either on feature dimensions selected according to the discrimination power (decreasing order)

350    (orange lines in Figure 2B), as shown in the eCDF plots (Figure 2A). Such procedure is referred

351    to as a simple filter approach on feature selection, where general characteristics of the extracted

352    features are taken into consideration when selecting feature dimensions, without subjecting

353    them to a classifier. We also applied a more extensive procedure of feature selection by

354    sequentially selecting feature dimensions by adding (forward search) feature dimensions,

355    referred to as sequential feature selection (black lines in Figure 2B). As part of this method, the

356    algorithm searched the best feature dimensions (predictors) according to their individual

357    classification performance in the given subset of data. For each candidate feature subset

358    (predictor), the algorithm performed a 10-fold cross-validation procedure with different

359    training and test subsets. After computing the mean performance values for each candidate

360    feature subset, the algorithm chooses the candidate feature subset with minimal

361    misclassification. For both methods, we systematically varied the number of features used for

362    classification (x-axis in Figure 2B). The selected features from a single run of the sequential

363    search algorithm are illustrated in Figure 2C. Scales reflect the feature-scaled and normalized

364    values, as a result of feature extraction, from which the grand means (i.e. for each feature

365    dimensions across all data) were subtracted. This measure was used to visually highlight

366    differences and was not used in further analyses.

367

368    *2.4.2 Classification*

369    We used support vector machine (SVM) with a radial basis function (RBF) Kernel (Vert et al.,

370    2004) for the classification of calls according to the class labels (so-called negative, neutral and

371    positive affective contexts). A classification procedure contains a training phase followed by a

372    test phase. We separated training samples and labelled them according to an attribute of interest

373    (e.g., negative (1) vs. positive (2) affective contexts). The algorithm then created a model that

374    optimally separates the two classes. In the test phase, samples without attribute labels were fed

375    into the model to measure its generalization performance. We used the SVM implementation

376    from LIBSVM toolbox (Chang & Lin, 2011). To evaluate how the classification results

377    generalize to a novel and independent data set, we 10-fold cross-validated the classification

378    process and optimized the parameters C and gamma (Fedurek et al., 2016), with the C taking

379    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

380    single individuals contributed solely to the classification outcome, we ran a leave-one-out

381    algorithm, where the procedure described above was re-run six times, excluding one of the

382    individuals in each run. We applied one-sample *t*-tests to compare the classification scores with

383    a 50% baseline condition. The 50% baseline results from the pairwise comparisons of affective

384    contexts (positive, neutral, negative). To ensure samples were normally distributed (a key

385    assumption behind the use of one-sample *t*-tests), we used Lilliefors test prior to each

386    comparison at a significance level of 5%. In cases where data samples were not normally

387    distributed, we used a one-sample Kolmogorov-Smirnov test. All reported *p*-values were

388    adjusted for multiple comparisons using Bonferroni corrections.

389

390    *2.4.3 Feature evaluation*

391    To evaluate whether certain feature dimensions are particularly critical for the classification of

392    grunts, we assessed whether feature dimensions have been repeatedly used by the classifier

393    overall in the classification of grunts. We therefore considered the three types of comparisons,

394    positive vs neutral, positive vs negative and neutral vs negative grunts, as well as the two feature

395    evaluation algorithms (simple feature selection and sequential feature selection). Each

396    comparison was ten-fold cross-validated. We then calculated the empirical distribution of the

397    ten features with best classification power, as determined by the feature selection algorithms

17

398 (see above). Also, we determined a random distribution of "best features" for each comparison

399 by randomly selecting 10 out of 104 features. The frequency distribution across all comparisons

400 were determined and 95% confidence intervals were calculated by running the procedure 1,000

401 times. We then traced back the significant feature dimensions to the underlying frequency bands

402 in Hertz.

403

404 ### 3. RESULTS

405 *3.1 Types of vocal utterances*

406 We inspected N = 1,016 vocal occurrences, of which N = 967 could be classified as either

407 'grunts' (N = 833) (corresponding to a rough, harsh and noisy sound) or 'whimpers' (N = 134)

408 (usually a series of low-pitch tonal calls with increase in fundamental frequency throughout the

409 series). Other types of calls were identified as 'hoos' (n = 23), 'pants' (n = 15), 'screams' (n =

410 2), 'squeaks' (n = 2), 'barks' (n = 4) and 'laughter' (defined as grunting and panting) (n = 3).

411

412 *3.2 Distribution of grunts and whimpers across so-called affective contexts*

413 Grunts: 44.8% of grunt-like vocalizations co-occurred with contexts we classified as 'positive',

414 40.9% with 'neutral', and 14.3% with 'negative'. When considering each individual separately,

415 a similar picture emerged (Figure 1), with most grunt-like vocalizations co-occurring with

416 'positive' and 'neutral' contexts. We found dominance to be relatively low in grunts, varying

417 from 0.37 and 0.63 (mean = 0.53; SD = 0.10), suggesting a stable and relative evenness in the

418 affective distribution of grunts, such as defined by our coding system (see Table 2).

419

420 Whimpers: 94.8% of whimpers co-occurred with negatively classified contexts, and rarely with

421 neutral (4.5%) or positive (0.7%) affects. Inspection of individual distributions revealed the

422 same pattern with whimper-like vocalizations systematically co-occurring with negatively

18

423     classified contexts (Figure 1). The dominance of one affective context over the others in

424     whimpers was relatively high, ranging from 0.89 to 1 (mean = 0.96; SD = 0.05), indicating low

425     evenness in the affective distribution of whimpers.

426

427     Grunts vs. Whimpers: When comparing the distributional evenness of grunts vs. whimpers, we

428     found dominance to be statistically higher in whimpers than in grunts (paired Wilcoxon signed

429     rank test: V = 21, *p* = .031).

430

431     *3.3 Acoustic variants of grunts*

432     We classified the N = 180 grunts (N = 60 per affective contexts) according to their association

433     with so-called positive, neutral, negative contexts in order to test for the presence of acoustic

434     variants. In the first step, we followed the feature extraction procedure by extracting the means

435     and covariances of MFCCs for each call, and compared these values according to the calls'

436     associations (e.g. positive vs. negative) using *t*-tests. We displayed the resulting *p*-values in an

437     empirical cumulative distribution function (eCDF) (Figure 2A). We found that 5-10% of all

438     features showed significant differences between the class labels at a 5%-significance level. In

439     other words, 5-10 of 104 feature dimensions had strong discrimination power to distinguish

440     between grunts pertaining to the various affective contexts.

441

442     With the simple feature selection algorithm, the SVM correctly discriminated between classes

443     at up to 80% (positive vs. neutral: M = 78.99, SD = 3.53, *t*(59) = 63.69, *p* < .001; positive vs.

444     negative: M = 79.58, SD = 1.83, *t*(59) = 125.37, *p* < .001; neutral vs. negative: M = 80.44, SD

445     = 2.06, *t*(59) = 114.26, *p* < .001; orange lines in Figure 2B). A substantial improvement was

446     found when sequentially selecting feature dimensions: SVM correctly classified samples at up

447     to 95% (positive vs. neutral: M = 89.56, SD = 4.84, *t*(59) = 143.42, *p* < .001; positive vs.

448    negative: M = 88.72, SD = 4.49, $t(59)$ = 153.11, $p < .001$; neutral vs. negative: M = 84.27, SD

449    = 5.23, $t(59)$ = 124.91, $p < .001$; black lines in Figure 2B). For all comparisons chance levels

450    were 50% due to the two-class comparisons applied. We, therefore, used one-sampled $t$-tests.

451    The classification scores in all (but one) comparisons fulfilled the requirement of normal

452    distribution. The first comparison (feature-selection algorithm, positive vs. neutral) was not

453    conform with a normal distribution and was, thus, re-evaluated using a one-sampled

454    Kolmogorov-Smirnov test, resulting in the following values (ks = 0.17; $p < .001$).

455

456    We further illustrated the simple feature selection outcomes by highlighting the feature

457    dimensions selected (circles in Figure 2C) among the feature dimensions not selected (gray

458    dots). Further, the features selected via the sequential feature selection are marked with x's. The

459    sequential feature selection yields better performance through sequential combinations of

460    feature dimensions that, on average, fall more distal to the diagonal mid-line than the feature

461    dimensions selected by the simple feature selection process. Sequential feature selection, to a

462    large extent, included feature dimensions not selected by the simple feature selection method.

463

464    We further ensured that each individual was not contributing solely to the classification results

465    of various contrasts. As can be seen in Supplementary Figure 1, the classification performance

466    did not improve nor deteriorate systematically when one individual was removed at a time,

467    suggesting no effect due to caller identity (the average $t$-value of one-sample $t$-tests is 97.52 +/-

468    30.25 (SD); all $p$-values were smaller than .001).

469

470    The use of means and covariances of cepstra yielded relatively high-performance scores in the

471    classification routines at low computational loads. To assess whether certain feature dimensions

472    (means and covariances of cepstra) occurred above chance across all comparisons, we

determined the empirical distribution of occurrences of feature dimensions and contrasted it with a random distribution. While the use of the same feature dimension in up to 33% of the comparisons was not significantly different in the empirical distribution from the random distribution, the use of the same feature dimension in 50% of comparisons was significantly increased in the empirical distribution (Figure 3A).

To describe the frequency bands explaining significant variances between classes of calls, we traced back the frequency bands underlying the significant feature dimensions, i.e., covariances of cepstra, and determined the sign of the covariances. We found negative covariances between the following frequency bands (Figure 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 927.78 Hz) and band 8 (1074.07 to 1366.67 Hz), band 6 (781.48 to 1074.07 Hz) and band 9 (1220.37 to 1512.96 Hz). We found a positive covariance between the frequency bands 9 (1220.37 to 1512.96 Hz) and 10 (1366.67 to 1659.26 Hz). Mean cepstra were significantly contributing in the frequency bands from (1) 50 to 342.59 Hz, (2) 196.30 to 488.89 Hz, (3) 927.78 to 1220.37 Hz.

## 4. DISCUSSION

Oller and colleagues (Jhang & Oller, 2017, 2017; Oller et al., 2013, 2016; Oller & Griebel, 2004) posit that speech emerged from pre-linguistic vocalizations that are free of predetermined biological function, a precursor called 'vocal functional flexibility'. One capacity foundational to vocal functional flexibility is the ability to use sounds that are not affectively-bound, a capacity we call 'affective decoupling'. Modern human infants regularly vocalize in such a way, in supposed contrast to the relative inflexibility of vocalizations in non-human primates (e.g., Pollick & Waal, 2007). Indeed, human infants can use sounds ('protophones') that can be uttered into a diversity of affective circumstances on diverse occasions, such that these sounds

21

498      are not tied to the experience and expression of one particular affective state (Oller et al., 2013;

499      Oller & Griebel, 2004). By contrast primate (and more largely, 'animal') vocal behavior would

500      be affectively bound, with particular calls being used to express particular affective state,

501      ultimately constraining their signaling function. The view that primate vocalizations are read-

502      outs of the affective states of the animal has otherwise long been held in the literature (Goodall,

503      1986; Gruber & Grandjean, 2017; Hammerschmidt & Fischer, 2008; Marler, 1980).

504

505      In the current study, we specifically looked at one of our closest living relative species, the

506      chimpanzees. We focused on the grunt-like and whimper-like calls of young chimpanzee

507      infants, using novel coding strategies and state-of-the-art acoustic analysis tools. We elaborated

508      a workable coding system, which was meant to provide first insights into the affective state of

509      infant chimpanzees, as seen in Oller et al. (Oller et al., 2013), and so as to allow for a first

510      comparison between human and chimpanzee infants. We found that grunt-like calls are

511      produced frequently by chimpanzee infants with both contexts we deemed positive and neutral,

512      and less commonly also with the so-called negative affective context. Importantly, the presence

513      of grunts in contexts of low-to-mild arousal is consistent with the hypothesis of vocal functional

514      flexibility (Oller et al., 2019), and so is the finding that grunts occur in similar proportion with

515      contexts we deemed positive and neutral (Oller et al., 2013).

516

517      On the other hand, whimper-like vocalizations seem to be confined to behaviors and contexts

518      we associated with negative affective states in the infants. Their near absence with positive and

519      neutral contexts suggests that they represent an affectively bound vocalization that has evolved

520      to signal a narrow range of needs and one single (negative) affective valence, similar to cries in

521      humans (Oller et al., 2013), to which they may functionally correspond (Goodall, 1986). Our

522      results therefore suggest that grunts are not bound to one particular affective context in

chimpanzees. They may also further qualify as a functionally flexible vocal unit, consistent with the observations of the circumstances of production of squeals, vocants and growls in young human infants (Oller et al., 2013). This, however, requires further examination, notably by improving our capacity to produce inferences about animals' transient affective states, and measuring whether recipients respond to these calls in a way consistent with the affect they are meant to convey.

Indeed, vocal functional flexibility requires not only affective decoupling (or the independence between particular vocalization and one affective dimension) but also evidence for consistent functionality. In human infants, the findings have been that infants use protophones with a diversity of affects, with mothers reacting consequently, showing that infants' calls are indeed fully functionally flexible (Oller et al., 2013). In these studies, the mothers' behavior could be examined, although protophones are not always socially directed (Oller & Griebel, 2004). Protocols where mothers may be asked to interact with toddlers may yield to responsiveness from the mothers whichever the affective state of the infant is (Yoo et al., 2018), which is critical in determining the function of the calls. In the course of spontaneous behavior, though, we expected little intervention from the chimpanzee mothers, except in situations where the infant is in danger. In our sample, responsiveness of the mother (tentatively defined in pilot coding as being either proactive, protective or neutral by the observer) was relatively low, a pattern which might be due to differences in mothering style between chimpanzees and humans, or a difference between our own study (where no particular demand is put on the mother) and others (where mothers may be interacting with their infant, e.g., Oller et al., 2013). This leaves us with the impossibility to conclude on whether mothers would react in ways consistent with the affective dimension of the vocal production, as seen in the human studies. Although playback of infant grunts to the mother may appear like a methodological possibility to further

548  establish their functionality (Fischer et al., 2013; Fischer, 2016; Zuberbühler, 2014), this would

549  require either playing the infants' calls in its own presence (which is ethically inappropriate) or

550  playing the calls of another infant to a mother (which may not trigger any reaction at all in the

551  non-genetically related mother). Another possibility is that the sounds we examined are not

552  meant to be fully functional, and could be considered to be vegetative sounds. The fact that they

553  may not appear socially directed should, however, not speak against the hypothesis that they

554  are affectively decoupled, for the fact that a given vocal unit is independent from one particular

555  affective valence is orthogonal with the fact that it is social directed or not. Our results are

556  compatible with grunts being a functionally flexible call type in young chimpanzees, but do not

557  yet demonstrate this, for the reactions of the mothers (and therefore, the function of the calls)

558  could not be directly assessed.

559

560  Grunts (and other close calls (Oller & Griebel, 2004)) are a promising class of vocalizations to

561  investigate the evolutionary origins of vocal functional flexibility. In a number of species (such

562  as the vervet monkeys (Cheney & Seyfarth, 1982), western gorillas (Salmi et al., 2013), sooty

563  mangabeys (Range & Fischer, 2004), chacma (Meise et al., 2011), Guinea (Faraut et al., 2019;

564  Maciej et al., 2013) and olive baboons (Ey & Fischer, 2011; Silk et al., 2018)), grunts are used

565  flexibly and can occur in a variety of contexts. So far, such evidence speaks in favor of grunts

566  being a contextually flexible vocal unit (that is, a vocal unit whose function can be fulfilled in

567  a diversity of contexts). Future research should try delving into the affective state animals likely

568  experience and express when producing grunts, to confirm whether these also displays affective

569  decoupling (i.e., the independence between grunt production and the experience of one

570  particular affective valence) and functional flexibility (i.e., the capacity of grunts to fulfil a

571  variety of functions on different occasions). If the term 'functional flexibility' could appear

572  misleading, its use in the field of child development should encourage the animal

24

573 communication community to employ it, such that more fruitful cross-disciplinary work can

574 best take place.

575

576 Our second main finding was systematic acoustic differences between grunts given with so-

577 called positive, neutral and negative behaviors, which enabled us to segregate acoustic variants

578 of grunts into these categories. Acoustical differences linked to the affect surrounding vocal

579 production are common in humans as in other animals (Arias et al., 2018; Aucouturier et al.,

580 2016; Banse & Scherer, 1996; Briefer, 2012; Goupil et al., 2019; Ponsot et al., 2018; Williams

581 & Stevens, 1972). Our data suggest that there is inter-gradation between grunt-types, with

582 differences in acoustics relating to differences in contexts. Grunts, in other words, represent a

583 coherent and unified call type that can manifest itself in acoustic variants in relation to the

584 affective contexts in which they are produced. It is possible that grunts acoustically vary with

585 arousal of the animal (as seen in other primate species (Rendall, 2003)), although positive and

586 negative circumstances could, in principle, be equally arousing.

587

588 How exactly functionally flexible vocalizations produced by human infants transition into

589 speech sounds has been described in previous studies (Boysson-Bardies, 2001; de Boysson-

590 Bardies, 1993; de Boysson-Bardies & Vihman, 1991; Elbers & Ton, 1985; Nathani et al., 2006;

591 Oller, 2000; Oller et al., 1976). Chimpanzee infants may produce grunts in ways consistent with

592 the functional flexibility hypothesis but they of course never produce speech sounds and,

593 historically, have failed to acquire human speech utterance even after extensive training (Hayes

594 & Hayes, 1951). Instead, infant chimpanzee grunts may gradually develop into call variants

595 with seemingly relatively narrow biological functions (Laporte & Zuberbühler, 2011;.

596 Slocombe & Zuberbühler, 2010; Slocombe & Zuberbühler, 2005; Watson et al., 2015), with

597 clear acoustical boundaries notably between grunts used to greet conspecifics ('pant-grunts'

598 (Laporte & Zuberbühler, 2011)) and those produced upon encountering food ('rough' or 'food

599 grunts' (Slocombe & Zuberbühler, 2005)). It is possible that the acoustic boundaries we

600 identified between the grunts produced across affective states (under our nomenclatures and

601 coding system) are the foundation of acoustic diversification in adults, although the categories

602 used here (for instance, feeding and social approach are together considered 'positive') are not

603 consistent with the vocal differentiation seen in adults (the grunts produced in feeding vs. social

604 approach situations are acoustically distinct in adults (Crockford, in press; Goodall, 1986)).

605 Alternatively, those calls may simply disappear and be absent from the adult repertoire, one

606 causal factor being the relative absence of social reinforcement (including contingent vocal

607 responses (Ghazanfar et al., 2019)) associated with grunt production, as compared to the

608 frequent maternal reactions to distress calls (Dezecache et al., 2020).

609

610 Our tentative to explore the affective state of the infant may be seen as preliminary, insofar as

611 the categories we have used do not represent read-outs of physiological states. This being said,

612 the acoustical differentiation we found speak in favor of the appropriate character of our

613 affective distinctions. Ideally, other cues should be considered, such as the infants' facial

614 expressions or the mothers' behavior. This approach would however face considerable

615 challenges. We found that infant facial movements are extremely fast and fluid, which

616 prevented us from reliable coding particularly in the wild. For this reason, the behavioral

617 context of the infant alone (although imperfect and probably still questionable) was the most

618 relevant available cue to approach the affective dimension of the situation. While we must again

619 acknowledge the limitations pertaining to the fact that judgments of infants' affect were made

620 based on the infants' behavioral contexts and done so by a human observer, the results of the

621 acoustic analysis are providing support for the approach used to categorize affect in the present

622 work. Future studies should investigate the affective impact of other communicative signals

623     used by infants, such as gesture and facial behavior, and their combinations (Fröhlich et al.,

624     2018; Fröhlich & Hobaiter, 2018).

625

626     Besides the limitations pertaining to our coding system (and its shortcomings with respect to

627     the production of inferences regarding infants' affective states), one other limitation of this

628     study is the small sample size, as we could only collect enough data from 6 infants. One

629     particular difficulty with collecting data from such young chimpanzee infants is that some of

630     their calls (notably a large part of their grunts) are very soft (a point also acknowledged by

631     Plooij (Plooij, 1984)) and can only be heard from close, limiting the number of individuals

632     whose mothers are unwary enough of continuous and long-lasting human observational efforts.

633     We could not use already published data, because, to the best of our knowledge, no previous

634     studies on the vocal behavior of wild infant chimpanzees (such as Laporte & Zuberbühler, 2011,

635     Plooij et al., 2014 or Plooij, 1984) used a coding system amenable to inferences about the

636     affective state of the infant.

637

638     In latest research, the comparative volubility (quantity of sounds produced in a given period of

639     time) of human infants and other animals (Ghazanfar & Takahashi, 2014; Oller et al., 2019;

640     Takahashi et al., 2015), and the privileged function of protophone-like vocalizations to

641     increasingly elicit social interactions and vocal turn-taking with caregivers (Oller et al., 2019;

642     Yoo et al., 2018). In humans, non-affectively bound vocalizations appear to occur more often

643     than affectively bound vocalizations (such as crying) (Oller et al., 2019). They occur in solitary

644     contexts where infants invest in practice and vocal exploration. They also occur in interactive

645     contexts, so as to elicit and regulate social interactions with caregivers. Caregivers appear to

646     detect the functional difference between protophones (as potentially interactive calls) and other

647     calls (such as cries), where caregiver intervention is solicited (Yoo et al., 2018). Comparison

648 with bonobo infants suggested much higher rate of production of non-affectively bound

649 vocalizations and much higher vocal investment in social interactions in human infants (Oller

650 et al., 2019). Whether human infants also are comparably more 'talkative' than their

651 chimpanzee counterparts is a question we need to be exploring. This should be preferably

652 investigated in captive or semi-captive settings, where true calling rate can be assessed, for

653 video monitoring is less likely to be interrupted and for levels of ambient noise could be

654 comparatively less problematic. Such problems have already been acknowledged by Oller and

655 colleagues (2019) regarding previous report on the flexible development of grunting behavior

656 in wild chimpanzees as well as their rate of occurrence (Laporte & Zuberbühler, 2011). Data

657 from the vocal development of one captive chimpanzee indicated lower volubility than in

658 humans (Kojima, 2003). Future studies should evaluate this fact with a larger sample.

659

660 Our study suggests that, insofar as one can delve into the affective state of infants using our

661 coding system, chimpanzees may possess a feature that is fundamental to the development of

662 speech in humans, the ability to produce vocalizations that are not strongly bound to the

663 experience and expression of one particular affective valence. However, we should expect that

664 future research will reveal further examples. For instance, coo calls in several macaque species

665 (Hsu et al., 2005; Owren & Casale, 1994), wahoos of baboons (Maciej et al., 2013) or grunts

666 of a number of primate species seem to be given in a variety of contexts, a precondition for

667 affective decoupling in vocal production, itself a prerequisite for vocal functional flexibility.

668 More largely, close calls appear to be excellent candidates (Oller & Griebel, 2004). Importantly,

669 methodologically efforts to infer the affective states of the animals should be made in order for

670 affective decoupling to be hypothesized.

671

672 Future research will have to address the question of how selection favored acoustic

673 diversification of functionally flexible vocal behavior into speech in humans. The main driver

674 for this transition, it has been argued, may have been the highly cooperative breeding system

675 of humans, with infants regularly looked after by individuals other than the mother, which

676 requires infants to become more active agents in forming social bonds from a much younger

677 age than in great ape infants (Ghazanfar et al., 2019; Zuberbühler, 2012).

678

679 Cooperative breeding, in this view, may thus have transformed a functionally flexible vocal

680 system into the uniquely human way of using vocal signals to interact socially. Another

681 complementary reasoning is that humans' high altriciality selected for the most vocal

682 individuals, capable of attracting caregivers (Locke, 2006). The relative contribution of both

683 factors through mapping the phylogenetic distribution of affective decoupling and vocal

684 functional flexibility remains to be investigated.

685

686 **REFERENCES**

687 Altmann, J. (1974). Observational study of behavior: Sampling methods. *Behaviour*,

688      *49*(3–4), 227–266.

689 Arbib, M. A., Liebal, K., Pika, S., Corballis, M. C., Knight, C., Leavens, D. A., Maestripieri, D.,

690      Tanner, J. E., Arbib, M. A., & Liebal, K. (2008). Primate vocalization, gesture, and

691      the evolution of human language. *Current Anthropology*, *49*(6), 1053–1076.

692 Arias, P., Belin, P., & Aucouturier, J.-J. (2018). Auditory smiles trigger unconscious facial

693      imitation. *Current Biology*, *28*(14), R782–R783.

694      https://doi.org/10.1016/j.cub.2018.05.084

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

696      Covert digital manipulation of vocal emotion alter speakers' emotional states in a

697   congruent direction. *Proceedings of the National Academy of Sciences*, *113*(4),

698   948–953. https://doi.org/10.1073/pnas.1506552113

699 Banse, R., & Scherer, K. R. (1996). Acoustic profiles in vocal emotion expression. *Journal*

700   *of Personality and Social Psychology*, *70*(3), 614–636.

701   https://doi.org/10.1037/0022-3514.70.3.614

702 Boë, L.-J., Berthommier, F., Legou, T., Captier, G., Kemp, C., Sawallis, T. R., Becker, Y., Rey,

703   A., & Fagot, J. (2017). Evidence of a Vocalic Proto-System in the Baboon (Papio

704   papio) Suggests Pre-Hominin Speech Precursors. *PLOS ONE*, *12*(1), e0169321.

705   https://doi.org/10.1371/journal.pone.0169321

706 Boysson-Bardies, B. de. (2001). *How Language Comes to Children: From Birth to Two*

707   *Years*. MIT Press.

708 Briefer, E. F. (2012). Vocal expression of emotions in mammals: Mechanisms of

709   production and evidence. *Journal of Zoology*, *288*(1), 1–20.

710   https://doi.org/10.1111/j.1469-7998.2012.00920.x

711 Burkart, J. M., Fehr, E., Efferson, C., & Schaik, C. P. van. (2007). Other-regarding

712   preferences in a non-human primate: Common marmosets provision food

713   altruistically. *Proceedings of the National Academy of Sciences*, *104*(50), 19762–

714   19766. https://doi.org/10.1073/pnas.0710310104

715 Burkart, J. M., Hrdy, S. B., & Van Schaik, C. P. (2009). Cooperative breeding and human

716   cognitive evolution. *Evolutionary Anthropology*, *18*(5), 175–186.

717 Chang, C.-C., & Lin, C.-J. (2011). LIBSVM: A library for support vector machines. *ACM*

718   *Transactions on Intelligent Systems and Technology (TIST)*, *2*(3), 27.

719   https://doi.org/10.1145/1961189.1961199

720 Cheney, D. L., & Seyfarth, R. M. (1982). How vervet monkeys perceive their grunts: Field

721   playback experiments. *Animal Behaviour*, *30*(3), 739–751.

722    Clay, Z., Archbold, J., & Zuberbühler, K. (2015). Functional flexibility in wild bonobo vocal

723          behaviour. *PeerJ*, *3*, e1124. https://doi.org/10.7717/peerj.1124

724    Crockford, C. (in press). Why Does the Chimpanzee Vocal Repertoire Remain Poorly

725          Understood? And What Can Be Done About It. In *The Tai Chimpanzees: 40 years of*

726          *Research. Eds: Boesch C. and Wittig R.* Cambridge University Press.

727    Crockford, C., & Boesch, C. (2005). Call combinations in wild chimpanzees. *Behaviour*,

728          *142*(4), 397–421. https://doi.org/10.1163/1568539054012047

729    de Boysson-Bardies, B. (1993). Ontogeny of Language-Specific Syllabic Productions. In B.

730          de Boysson-Bardies, S. de Schonen, P. Jusczyk, P. McNeilage, & J. Morton (Eds.),

731          *Developmental Neurocognition: Speech and Face Processing in the First Year of Life*

732          (pp. 353–363). Springer Netherlands. https://doi.org/10.1007/978-94-015-

733          8234-6_29

734    de Boysson-Bardies, B., & Vihman, M. M. (1991). Adaptation to Language: Evidence from

735          Babbling and First Words in Four Languages. *Language*, *67*(2), 297–319. JSTOR.

736          https://doi.org/10.2307/415108

737    Dezecache, G., Zuberbühler, K., Davila-Ross, M., & Dahl, C. D. (2020). A machine learning

738          approach to infant distress calls and maternal behaviour of wild chimpanzees.

739          *Animal Cognition*. https://doi.org/10.1007/s10071-020-01437-5

740    Elbers, L., & Ton, J. (1985). Play pen monologues: The interplay of words and babbles in

741          the first words period. *Journal of Child Language*, *12*(3), 551–565.

742          https://doi.org/10.1017/S0305000900006644

743    Ey, E., & Fischer, J. (2011). Keeping in contact: Flexibility in calls of olive baboons. In

744          *Primates of Gashaka* (pp. 413–436). Springer.

745    Faraut, L., Siviter, H., Dal Pesco, F., & Fischer, J. (2019). How life in a tolerant society

746        affects the usage of grunts: Evidence from female and male Guinea baboons.

747        *Animal Behaviour*, *153*, 83–93.

748    Fedurek, P., & Slocombe, K. E. (2013). The social function of food-associated calls in male

749        chimpanzees. *American Journal of Primatology*, *75*(7), 726–739.

750    Fedurek, P., Zuberbühler, K., & Dahl, C. D. (2016). Sequential information in a great ape

751        utterance. *Scientific Reports*, *6*, 38226. https://doi.org/10.1038/srep38226

752    Fischer, J. (2016). Playback Experiments. *The International Encyclopedia of Primatology*,

753        1–2. https://doi.org/10.1002/9781119179313.wbprim0140

754    Fischer, J., Noser, R., & Hammerschmidt, K. (2013). Bioacoustic field research: A primer

755        to acoustic analyses and playback experiments with primates. *American Journal*

756        *of Primatology*, *75*(7), 643–663. https://doi.org/10.1002/ajp.22153

757    Fitch, W. T. (2018). The Biology and Evolution of Speech: A Comparative Analysis.

758        *Annual Review of Linguistics*, *4*(1), 255–279. https://doi.org/10.1146/annurev-

759        linguistics-011817-045748

760    Fitch, W. T., Boer, B. de, Mathur, N., & Ghazanfar, A. A. (2016). Monkey vocal tracts are

761        speech-ready. *Science Advances*, *2*(12), e1600723.

762        https://doi.org/10.1126/sciadv.1600723

763    Fröhlich, M., & Hobaiter, C. (2018). The development of gestural communication in great

764        apes. *Behavioral Ecology and Sociobiology*, *72*(12), 194.

765    Fröhlich, M., Wittig, R. M., & Pika, S. (2018). The ontogeny of intentional communication

766        in chimpanzees in the wild. *Developmental Science*, e12716.

767    Ghazanfar, A. A., Liao, D. A., & Takahashi, D. Y. (2019). Volition and learning in primate

768        vocal behaviour. *Animal Behaviour*, *151*, 239–247.

769        https://doi.org/10.1016/j.anbehav.2019.01.021

770 Ghazanfar, A. A., & Takahashi, D. Y. (2014). The evolution of speech: Vision, rhythm,

771   cooperation. *Trends in Cognitive Sciences*, *18*(10), 543–553.

772   https://doi.org/10.1016/j.tics.2014.06.004

773 Goodall, J. (1986). *The chimpanzees of Gombe: Patterns of behavior*. Harvard University

774   Press.

775 Goupil, L., Johansson, P., Hall, L., & Aucouturier, J.-J. (2019). *Influence of Vocal Feedback*

776   *on Emotions Provides Causal Evidence for the Self-Perception Theory*.

777   https://doi.org/10.1101/510867

778 Gruber, T., & Grandjean, D. (2017). A comparative neurological approach to emotional

779   expressions in primate vocalizations. *Neuroscience & Biobehavioral Reviews*, *73*,

780   182–190.

781 Hammerschmidt, K., & Fischer, J. (2008). Constraints in primate vocal production. In

782   *Evolution of communicative flexibility: Complexity, creativity and adaptability in*

783   *human and animal communication*. MIT press.

784 Hayes, K. J., & Hayes, C. (1951). The Intellectual Development of a Home-Raised

785   Chimpanzee. *Proceedings of the American Philosophical Society*, *95*(2), 105–109.

786 Hrdy, S. (2007). Evolutionary context of human development: The cooperative breeding

787   model. In *Family Relationships: An Evolutionary Perspective*. Oxford University

788   Press.

789 Hsu, M. J., Chen, L.-M., & Agoramoorthy, G. (2005). The vocal repertoire of Formosan

790   macaques, Macaca cyclopis: Acoustic structure and behavioral context. *Zoological*

791   *Studies*, *44*(2), 275.

792 Jhang, Y., & Oller, D. K. (2017). Emergence of Functional Flexibility in Infant

793   Vocalizations of the First 3 Months. *Frontiers in Psychology*, *8*.

794   https://doi.org/10.3389/fpsyg.2017.00300

795 Jürgens, U. (1976). Reinforcing concomitants of electrically elicited vocalizations.

796     *Experimental Brain Research*, *26*(2), 203–214.

797 Jürgens, Uwe. (1979). Vocalization as an emotional indicator. *Behaviour*, *69*(1–2), 88–

798     117.

799 Kojima, S. (2003). *A search for the origins of human speech: Auditory and vocal functions*

800     *of the chimpanzee*. Kyoto University Academic Press.

801 Kramer, K. L. (2010). Cooperative Breeding and its Significance to the Demographic

802     Success of Humans. *Annual Review of Anthropology*, *39*(1), 417–436.

803     https://doi.org/10.1146/annurev.anthro.012809.105054

804 Lameira, A. R., & Shumaker, R. W. (2019). Orangutans show active voicing through a

805     membranophone. *Scientific Reports.*, *9*, 12289. https://doi.org/10.1038/s41598-

806     019-48760-7

807 Laporte, M. N. C., & Zuberbühler, K. (2011). The development of a greeting signal in wild

808     chimpanzees. *Developmental Science*, *14*(5), 1220–1234.

809     https://doi.org/10.1111/j.1467-7687.2011.01069.x

810 Laporte, M. N., & Zuberbühler, K. (2010). Vocal greeting behaviour in wild chimpanzee

811     females. *Animal Behaviour*, *80*(3), 467–473.

812 Leavens, D. A. (2009). Animal communication: Laughter is the shortest distance between

813     two apes. *Current Biology*, *19*(13), R511–R513.

814 Levréro, F., & Mathevon, N. (2013). Vocal Signature in Wild Infant Chimpanzees.

815     *American Journal of Primatology*, *75*(4), 324–332.

816     https://doi.org/10.1002/ajp.22108

817 Lieberman, P. (2017). Comment on "Monkey vocal tracts are speech-ready." *Science*

818     *Advances*, *3*(7), e1700442. https://doi.org/10.1126/sciadv.1700442

819     Locke, J. L. (2006). Parental selection of vocal behavior. *Human Nature*, *17*(2), 155–168.

820         https://doi.org/10.1007/s12110-006-1015-x

821     Logan, B. (2000). Mel frequency cepstral coefficients for music modeling. *Ismir*, *270*, 1–

822         11.

823     Maciej, P., Ndao, I., Hammerschmidt, K., & Fischer, J. (2013). Vocal communication in a

824         complex multi-level society: Constrained acoustic structure and flexible call

825         usage in Guinea baboons. *Frontiers in Zoology*, *10*(1), 58.

826     Mandel, M. I., & Ellis, D. P. (2005). Song-level features and support vector machines for

827         music classification. *Proceedings of the 6th International Conference on Music*

828         *Information Retrieval (ISMIR)*, 594–599.

829     Marler, P. (1980). Primate vocalization: Affective or symbolic? In *Speaking of apes* (pp.

830         221–229). Springer.

831     McCune, L., Vihman, M. M., Roug-Hellichius, L., Delery, D. B., & Gogate, L. (1996). Grunt

832         communication in human infants (Homo sapiens). *Journal of Comparative*

833         *Psychology*, *110*(1), 27.

834     Meise, K., Keller, C., Cowlishaw, G., & Fischer, J. (2011). Sources of acoustic variation:

835         Implications for production specificity and call categorization in chacma baboon

836         (Papio ursinus) grunts. *The Journal of the Acoustical Society of America*, *129*(3),

837         1631–1641.

838     Mielke, A., & Zuberbühler, K. (2013). A method for automated individual, species and call

839         type recognition in free-ranging animals. *Animal Behaviour*, *86*(2), 475–482.

840     Morris, E. K., Caruso, T., Buscot, F., Fischer, M., Hancock, C., Maier, T. S., Meiners, T.,

841         Müller, C., Obermaier, E., & Prati, D. (2014). Choosing and using diversity indices:

842         Insights for ecological applications from the German Biodiversity Exploratories.

843         *Ecology and Evolution*, *4*(18), 3514–3524.

844 Nathani, S., Ertmer, D. J., & Stark, R. E. (2006). Assessing vocal development in infants

845 and toddlers. *Clinical Linguistics & Phonetics*, *20*(5), 351–369.

846 https://doi.org/10.1080/02699200500211451

847 Oller, D. K. (2000). *The Emergence of the Speech Capacity*. Psychology Press.

848 https://doi.org/10.4324/9781410602565

849 Oller, D. K., Buder, E. H., Ramsdell, H. L., Warlaumont, A. S., Chorna, L., & Bakeman, R.

850 (2013). Functional flexibility of infant vocalization and the emergence of

851 language. *Proceedings of the National Academy of Sciences of the United States of*

852 *America*, *110*(16), 6318–6323. https://doi.org/10.1073/pnas.1300337110

853 Oller, D. K., & Griebel, U. (2004). Contextual freedom in human infant vocalization and

854 the evolution of language. In *Evolution of communicative flexibility: Complexity,*

855 *creativity and adaptability in human and animal communication* (p. 135). MIT

856 Press.

857 Oller, D. K., Griebel, U., Iyer, S. N., Jhang, Y., Warlaumont, A. S., Dale, R., & Call, J. (2019).

858 Language Origins Viewed in Spontaneous and Interactive Vocal Rates of Human

859 and Bonobo Infants. *Frontiers in Psychology*, *10*.

860 https://doi.org/10.3389/fpsyg.2019.00729

861 Oller, D. K., Griebel, U., & Warlaumont, A. S. (2016). Vocal Development as a Guide to

862 Modeling the Evolution of Language. *Topics in Cognitive Science*, *8*(2), 382–392.

863 https://doi.org/10.1111/tops.12198

864 Oller, D. K., Wieman, L. A., Doyle, W. J., & Ross, C. (1976). Infant babbling and speech.

865 *Journal of Child Language*, *3*(1), 1–11.

866 https://doi.org/10.1017/S0305000900001276

867    Owren, M. J., & Casale, T. M. (1994). Variations in fundamental frequency peak position

868        in Japanese macaque (Macaca fuscata) coo calls. *Journal of Comparative*

869        *Psychology*, *108*(3), 291. https://doi.org/10.1037/0735-7036.108.3.291

870    Plooij, F. X. (1984). *The behavioral development of free-living chimpanzee babies and*

871        *infants.* Ablex.

872    Plooij, F. X., Van De Rijt-plooij, H., Fischer, M., & Pusey, A. (2014). Longitudinal

873        recordings of the vocalizations of immature Gombe chimpanzees for

874        developmental studies. *Scientific Data*, *1*(1), 1–10.

875    Pollick, A. S., & Waal, F. B. M. de. (2007). Ape gestures and language evolution.

876        *Proceedings of the National Academy of Sciences*, *104*(19), 8184–8189.

877        https://doi.org/10.1073/pnas.0702624104

878    Ponsot, E., Burred, J. J., Belin, P., & Aucouturier, J.-J. (2018). Cracking the social code of

879        speech prosody using reverse correlation. *Proceedings of the National Academy of*

880        *Sciences*, *115*(15), 3972–3977. https://doi.org/10.1073/pnas.1716090115

881    R Core Team. (2018). *R: A language and environment for statistical computing.* R

882        Foundation for Statistical Computing. https://www.R-project.org/

883    Range, F., & Fischer, J. (2004). Vocal repertoire of sooty mangabeys (Cercocebus

884        torquatus atys) in the Taï National Park. *Ethology*, *110*(4), 301–321.

885    Rendall, D. (2003). Acoustic correlates of caller identity and affect intensity in the vowel-

886        like grunt vocalizations of baboons. *The Journal of the Acoustical Society of*

887        *America*, *113*(6), 3390–3402.

888    Rendall, D., & Owren, M. J. (2002). Animal vocal communication: Say what? In *The*

889        *cognitive animal: Empirical and theoretical perspectives on animal cognition* (pp.

890        307–313). MIT Press.

891    Reynolds, V. (2005). *The chimpanzees of the Budongo Forest: Ecology, behaviour, and*

892          *conservation*. Oxford University Press.

893          http://books.google.fr/books?hl=fr&lr=&id=C6hzM5lQJ6YC&oi=fnd&pg=PR11&

894          dq=budongo+reynolds&ots=OOfJtMfycP&sig=X1c6kEzGs8ZlzhXdNH3aK3foPrA

895    RStudio Team. (2015). RStudio: Integrated development for R. *RStudio, Inc., Boston, MA*

896          *URL Http://Www. Rstudio. Com*, *42*.

897    Salmi, R., Hammerschmidt, K., & Doran-Sheehy, D. M. (2013). Western gorilla vocal

898          repertoire and contextual use of vocalizations. *Ethology*, *119*(10), 831–847.

899    Schaik, C. P. van, & Burkart, J. M. (2010). Mind the Gap: Cooperative Breeding and the

900          Evolution of Our Unique Features. In *Mind the Gap* (pp. 477–496). Springer,

901          Berlin, Heidelberg. https://doi.org/10.1007/978-3-642-02725-3_22

902    Schel, A. M., Machanda, Z., Townsend, S. W., Zuberbühler, K., & Slocombe, K. E. (2013).

903          Chimpanzee food calls are directed at specific individuals. *Animal Behaviour*,

904          *86*(5), 955–965.

905    Silk, J. B., Roberts, E. R., Städele, V., & Strum, S. C. (2018). To grunt or not to grunt:

906          Factors governing call production in female olive baboons, Papio anubis. *PloS*

907          *One*, *13*(11), e0204601.

908    Slocombe, K. E., & Zuberbühler, K. (2010). Vocal communication in chimpanzees. *The*

909          *Mind of the Chimpanzee: Ecological and Experimental Perspectives. University of*

910          *Chicago Press, Chicago*, 192–207.

911    Slocombe, Katie E., Kaller, T., Turman, L., Townsend, S. W., Papworth, S., Squibbs, P., &

912          Zuberbühler, K. (2010). Production of food-associated calls in wild male

913          chimpanzees is dependent on the composition of the audience. *Behavioral*

914          *Ecology and Sociobiology*, *64*(12), 1959–1966.

915 Slocombe, Katie E., & Newton-Fisher, N. E. (2005). Fruit sharing between wild adult

916       chimpanzees (Pan troglodytes schweinfurthii): A socially significant event?

917       *American Journal of Primatology: Official Journal of the American Society of*

918       *Primatologists*, *65*(4), 385–391.

919 Slocombe, Katie E., & Zuberbühler, K. (2005). Functionally Referential Communication in

920       a Chimpanzee. *Current Biology*, *15*(19), 1779–1784.

921       https://doi.org/10.1016/j.cub.2005.08.068

922 Tajiri, Y., Yabuwaki, R., Kitamura, T., & Abe, S. (2010). Feature Extraction Using Support

923       Vector Machines. In K. W. Wong, B. S. U. Mendis, & A. Bouzerdoum (Eds.), *Neural*

924       *Information Processing. Models and Applications* (pp. 108–115). Springer.

925       https://doi.org/10.1007/978-3-642-17534-3_14

926 Takahashi, D. Y., Fenley, A. R., Teramoto, Y., Narayanan, D. Z., Borjon, J. I., Holmes, P., &

927       Ghazanfar, A. A. (2015). The developmental dynamics of marmoset monkey vocal

928       production. *Science*, *349*(6249), 734–738.

929       https://doi.org/10.1126/science.aab1058

930 Tchernichovski, O., & Oller, D. K. (2016). Vocal Development: How Marmoset Infants

931       Express Their Feelings. *Current Biology*, *26*(10), R422–R424.

932       https://doi.org/10.1016/j.cub.2016.03.063

933 Tsukahara, T. (1993). Lions eat chimpanzees: The first evidence of predation by lions on

934       wild chimpanzees. *American Journal of Primatology*, *29*(1), 1–11.

935 Vert, J.-P., Tsuda, K., & Schölkopf, B. (2004). A primer on kernel methods. In *Kernel*

936       *methods in computational biology* (Vol. 47, pp. 35–70). MIT Press.

937 Waal, F. B. M. de, & Pollick, A. S. (2011). Gesture as the most flexible modality of primate

938       communication. *The Oxford Handbook of Language Evolution*.

939       https://doi.org/10.1093/oxfordhb/9780199541119.013.0006

940 Watson, S. K., Townsend, S. W., Schel, A. M., Wilke, C., Wallace, E. K., Cheng, L., West, V., &

941      Slocombe, K. E. (2015). Vocal Learning in the Functionally Referential Food

942      Grunts of Chimpanzees. *Current Biology*, *25*(4), 495–499.

943      https://doi.org/10.1016/j.cub.2014.12.032

944 Williams, C. E., & Stevens, K. N. (1972). Emotions and Speech: Some Acoustical

945      Correlates. *The Journal of the Acoustical Society of America*, *52*(4B), 1238–1250.

946      https://doi.org/10.1121/1.1913238

947 Yoo, H., Bowman, D. A., & Oller, D. K. (2018). The Origin of Protoconversation: An

948      Examination of Caregiver Responses to Cry and Speech-Like Vocalizations.

949      *Frontiers in Psychology*, *9*. https://doi.org/10.3389/fpsyg.2018.01510

950 Zhang, Y. S., & Ghazanfar, A. A. (2016). Perinatally influenced autonomic system

951      fluctuations drive infant vocal sequences. *Current Biology*, *26*(10), 1249–1260.

952 Zuberbühler, K. (2012). Cooperative breeding and the evolution of vocal flexibility. In

953      *The Oxford Handbook of Language Evolution*. Oxford University Press.

954 Zuberbühler, K. (2014). Experimental field studies with non-human primates. *Current

955      Opinion in Neurobiology*, *28*, 150–156.

956      https://doi.org/10.1016/j.conb.2014.07.012

957

958

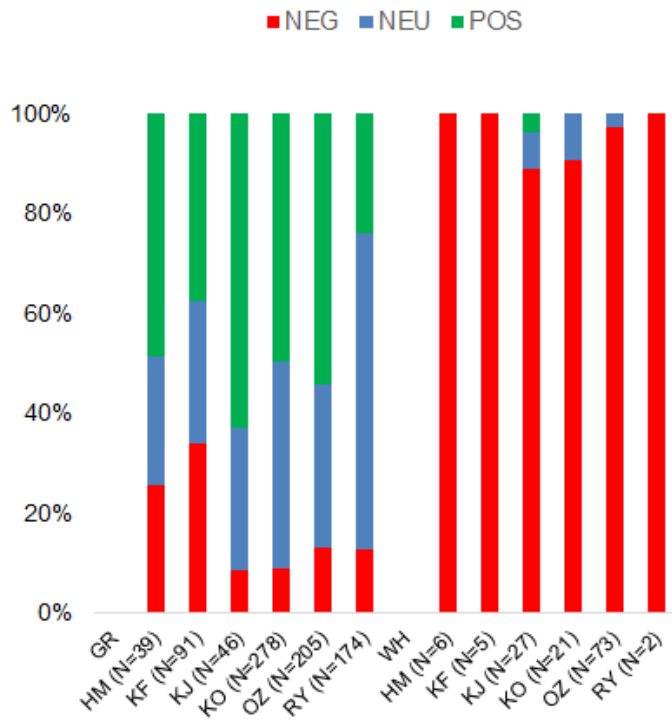959 **CONFLICTS OF INTEREST**

960 No conflicts of interest.

961

962 **FIGURES**

963 **Figure 1** Proportion of grunt-like (GR) and whimper-like (WH) vocal behaviors recorded with

964 negative (NEG), neutral (NEU) and positive (POS) affective categories of behaviors, for each

965    individual separately. Number between brackets indicate the number of GR and WH calls

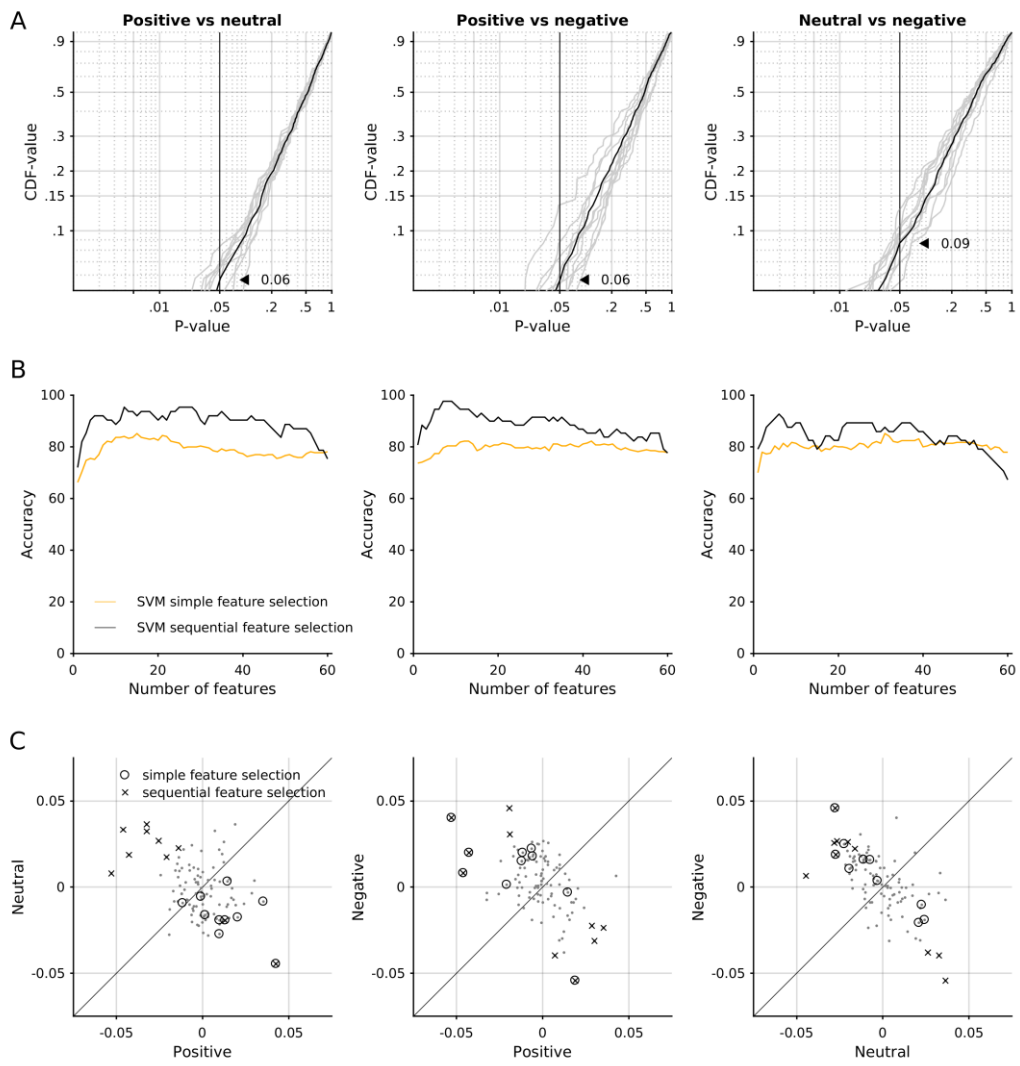966    contributed by each individual.

967



968

969

970

971 **Figure 2** Feature selection and classification performances. The columns represent the

972 comparisons of affects during which the vocal utterance occurred.

973 A. For each feature dimensions the discrimination power of the two classes (e.g. positive vs.

974 neutral) was evaluated using a t-test. P-values are shown as an empirical cumulative distribution

975 function (eCDF). Gray lines show the results of individual runs of evaluation; black lines show

976 the means of individual runs. Indicated with arrow heads are the proportions of feature

977 dimensions that significantly discriminate between the two classes tested.

978 B. The classification performances are shown for the SVM classifier relying on feature

979 dimensions extracted through a simple feature selection (orange lines) and a sequential feature

980 selection procedure (black lines).

981 C. Feature selection outcomes are shown for simple (circles) and sequential feature selection

982 (blue x-s) as overlays on all feature dimensions (gray dots).
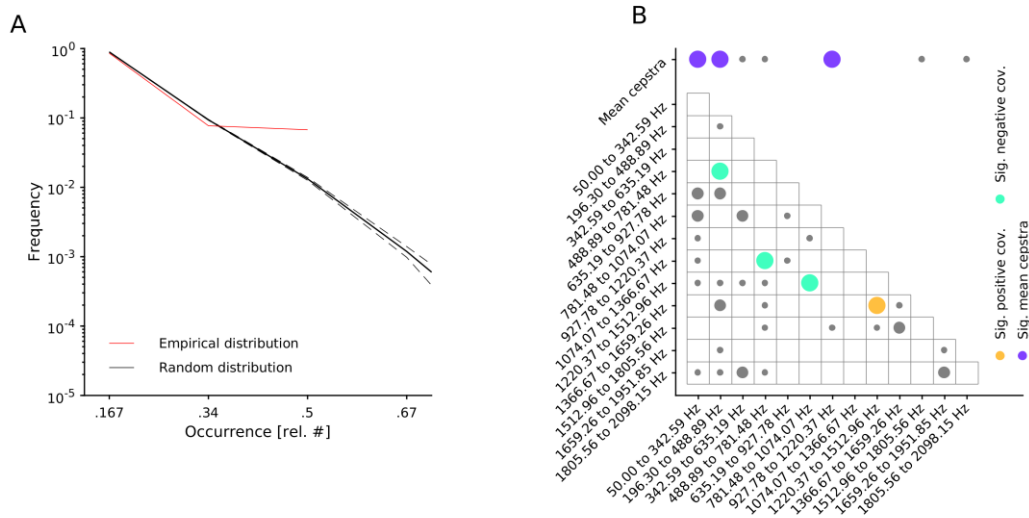
983

984

**Figure 3** Overall feature importance.

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

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

**TABLES**

995 **Table 1** List of focal animals, with their name (ID), sex and minimum and maximum age in

996 months. Also given are the number of grunt-like and whimper-like vocal behaviors collected,

997 as well as grunt-like vocalizations acoustically analyzed.

998

| ID | Sex | Min. Age (in months) | Max. Age (in months) | N whimper-like vocalizations | N grunt-like vocalizations | N of grunt-like vocalizations used in acoustical analysis |
|---|---|---|---|---|---|---|
| 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 |

999

1000   **Table 2** Affective coding of infant behavior

| Affect | Behavior | Description |
|---|---|---|
| POSITIVE | Play | Relaxed movements without obvious purpose. Can be solitary (shaking, biting and gnawing vegetation, swinging) or social (wrestling, gentle biting, gentle hitting, chasing or being chased). |
| POSITIVE | Grooming | Giving or receiving 'grooming', i.e., defined following Plooij (1984) as 'picking through the fur of another individual', using one's hands or lips. |
| POSITIVE | Feeding | Breastfeeding or swallowing an edible element |
| POSITIVE | Social approach | Greeting a conspecific whilst moving (locomotion or clear leaning of the body) towards this individual |
| NEUTRAL | Resting | Remaining within a limited area, may involve some degree of moving around, marked by relative idleness |
| NEUTRAL | Moving | Locomotion not directed towards a specific individual, and not involving play |
| NEUTRAL | Manipulating objects | Manipulating objects (leaves, branches, rocks) |
| NEUTRAL | Greeting without approach | Calling upon the approach of a conspecific without showing approach (as in Social approach) or avoidance behavior towards it |
| NEGATIVE | 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 movements meant to avoid or withdraw from a certain situation (play, grooming) or a physical position (such as moments |

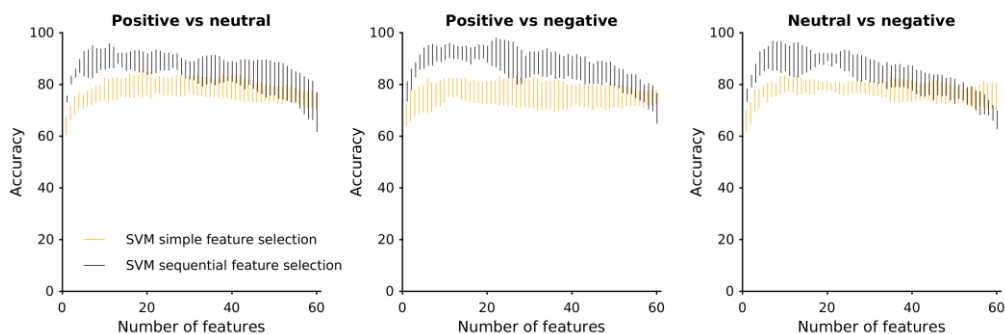| | | of discomfort when the infant is suddenly pressed against the belly of the mother) the infant is in |
|---|---|---|

1001

## SUPPLEMENTARY INFORMATION

**Supplementary Table 1** Number of calls per infant per affective category.

|  | Negative | Neutral | Positive | Grand Total |
|---|---|---|---|---|
| **Grunt-like** | **119** | **341** | **373** | **833** |
| HM | 10 | 10 | 19 | 39 |
| KF | 31 | 26 | 34 | 91 |
| KJ | 4 | 13 | 29 | 46 |
| KO | 25 | 115 | 138 | 278 |
| OZ | 27 | 67 | 111 | 205 |
| RY | 22 | 110 | 42 | 174 |
| **Whimper-like** | **127** | **6** | **1** | **134** |
| HM | 6 |  |  | 6 |
| KF | 5 |  |  | 5 |
| KJ | 24 | 2 | 1 | 27 |
| KO | 19 | 2 |  | 21 |
| OZ | 71 | 2 |  | 73 |
| RY | 2 |  |  | 2 |
| **Grand Total** | **246** | **347** | **374** | **967** |

**Supplementary Figure 1** Leave-one-out method to account for subject effects. The accuracies of the three comparisons of grunt types are shown as function of number of features. These graphs illustrate the variability of accuracy caused by leaving out one of the 6 individuals per each separate classification procedure. The vertical bars indicate the minimum and maximum scores.
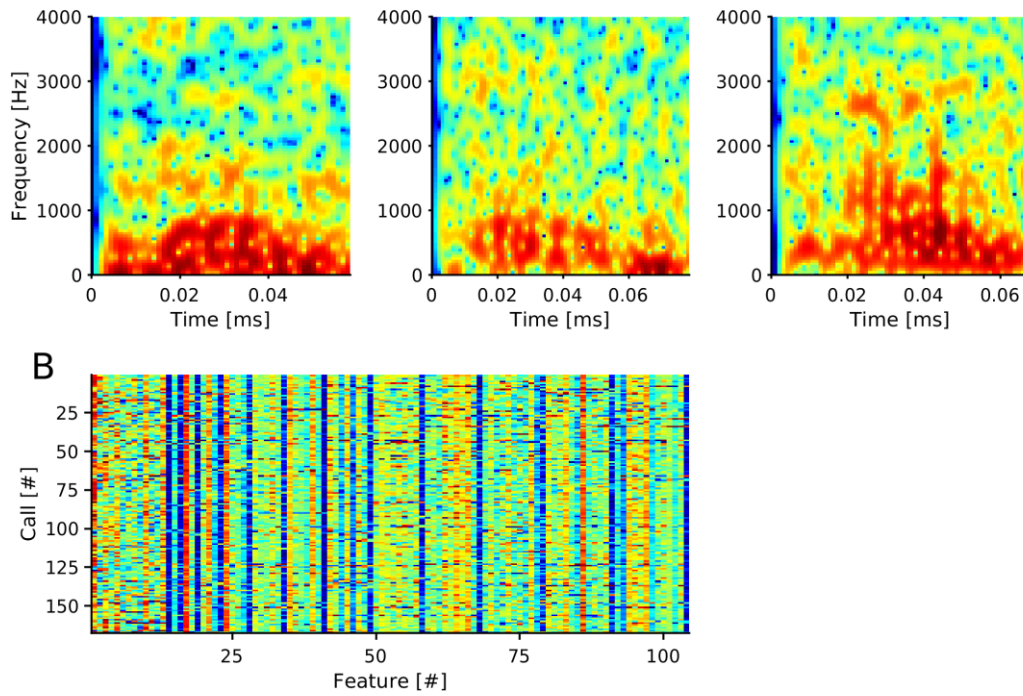
1012 **Supplementary Figure 2** MFCCs extracted from example calls and extracted feature matrix.

1013 A. Time-frequency spectra of three arbitrarily chosen calls.

1014 B. From each call 26 spectral bands and 13 cepstra were extracted. Feature vectors containing

1015 the means and covariances of cepstra are shown for each call. Means are shown as features 1

1016 to 13 on the x-axis, followed by covariances (91 values).



1017