



## Flexibility in wild infant chimpanzee vocal behavior

Guillaume Dezecache, Klaus Zuberbühler, Marina Davila-Ross, Christoph Dahl

### ► To cite this version:

Guillaume Dezecache, Klaus Zuberbühler, Marina Davila-Ross, Christoph Dahl. Flexibility in wild infant chimpanzee vocal behavior. *Journal of Language Evolution*, 2021, 6 (1), pp.37-53. 10.1093/jole/lzaa009 . hal-03053559

**HAL Id: hal-03053559**

**<https://hal.science/hal-03053559>**

Submitted on 11 Dec 2020

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

# Flexibility in wild infant chimpanzee vocal behavior

Guillaume Dezecache<sup>1,2,3,4,\*</sup>, Klaus Zuberbühler<sup>1,2,5</sup>, Marina Davila-Ross<sup>3</sup>, and Christoph D. Dahl<sup>1,6,7,\*</sup>

<sup>1</sup>Institute of Biology, University of Neuchâtel, Neuchâtel, Switzerland, <sup>2</sup>Budongo Conservation Field Station, Masindi, Uganda, <sup>3</sup>Department of Psychology, University of Portsmouth, Portsmouth, UK, <sup>4</sup>Université Clermont Auvergne, CNRS, LAPSCO, Clermont-Ferrand, France, <sup>5</sup>School of Psychology and Neuroscience, University of St Andrews, St Andrews, UK, <sup>6</sup>Graduate Institute of Mind, Brain and Consciousness, Taipei Medical University, Taipei, Taiwan and <sup>7</sup>Brain and Consciousness Research Center, Taipei Medical University Shuang-Ho Hospital, New Taipei City, Taiwan

\*Corresponding authors: guillaume.dezecache@gmail.com; christoph.d.dahl@gmail.com

## 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, in 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 toward 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. In 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 reflects 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.

**Key words:** language evolution; vocal flexibility; affective decoupling; wild chimpanzees; Pan troglodytes; grunts

## 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 and Owren 2002). More recently, however, some vocal and neural equipments have 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 and Shumaker 2019). If the facial and gestural displays have undeniably played a crucial role in the evolution of language (Pollick and Waal 2007; Arbib et al. 2008; Waal and 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 nonhuman primate and human vocal behavior prior to speech. By the age of 1-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 (Oller et al. 2013; Jhang and Oller 2017). 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). In contrast to those vocal types, human infants use laughter and cries in single affective contexts (positive and negative, respectively).

This capacity to produce one vocal unit under a variety of affective state (such that vocal production is independent from the experience of a specific affective state—a capacity termed 'affective decoupling') later allows infants to use those sounds as they wish, and to express a variety of needs on different occasions (Oller and Griebel 2004; Oller et al., 2013). This decoupling between vocal production and affective experience is foundational for the ability to produce sounds that can later fulfil a variety of functions ('vocal functional flexibility'), that is, they can be used to invite a variety of responses in others (Oller and Griebel 2004). For

instance, a given utterance (such as 'the train is arriving') produced with neutral (a mere statement), negative (annoyance), and positive (exultation) tones have the same syntactic structure and semantic content, but are produced under antagonist affective states and cause vastly different responses in receivers. Without affective decoupling and vocal functional flexibility, basic speech acts cannot take place.

This decoupling of signal structure and affect in young infants' vocal repertoire has thus been identified as a major evolutionarily precursor to language (Oller et al. 2013). Because of their early ontogenetic onset, affective decoupling and vocal functional flexibility may be more foundational to human speech than other building blocks of the language faculty, such as proto-syntax or vocal elaboration (Oller et al. 2013). These capacities, in this view, are prerequisites for speech development, and major evolutionary departures from the affect-bound and functionally inflexible vocal behavior of nonhuman primates (Waal and Pollick 2011). In contrast to their primate cousins, humans would have evolved in a social ecology conducive to the development of such vocal flexibility. Notably, it is proposed that functionally flexible vocalizations of young human infants have evolved in humans in relation to allo-maternity (Hrdy 2007; Burkart et al. 2009; Kramer 2010; Schaik and Burkart 2010) or altriciality (Locke 2006) and associated pressures on young infants to signal their needs and attract caregivers (Locke, 2006; Zuberbühler 2012; Ghazanfar et al. 2019). Other species living in cooperative breeding systems (such as the marmosets (Burkart et al. 2007)) may display vocal functional flexibility in their vocal repertoire.

For long, primate (but also animal) communication systems have been characterized as affectively biased, such that affect has been proposed to be both necessary and sufficient for vocal production to occur. According to Hammerschmidt & Fischer (2008: 103), there could even exist '[a] correspondence between non-verbal vocalizations in humans and non-human primates [such] that they both function to communicate the affective state of the signaler'. In fact, a number of researchers have held the view that the equivalents of animal vocalizations are nonverbal affective expressions in humans, such as laughing, screaming, and crying (Gruber and Grandjean 2017; Marler 1980 for a discussion). Examining the neural pathways of vocal production in squirrel monkeys, Jürgens (Jürgens 1976, 1979) concluded that vocal production was mediated by affect. More recently in marmoset infants, Zhang & Ghazanfar (2016) found that fluctuations in cardiac rhythm shape respiratory patterns, which in turn contribute to vocal

production, thereby attributing a central role to affect in early vocal production of this species (Tchernichovski and Oller 2016). The vocal repertoire of species phylogenetically closer to humans (such as the chimpanzees) did not evade this conclusion. In her authoritative book on the behavior of Gombe chimpanzees, Goodall (1986) wrote that ‘chimpanzee vocalizations are closely bound to emotion’ and that ‘the production of a sound in the absence of the appropriate emotional state seems to be an almost impossible task for a chimpanzee’ (p. 125). As a consequence, she proposed a mapping between call and affect when describing the vocal repertoire of chimpanzees, with, for instance, a one-to-one correspondence between the experience of annoyance and the production of ‘soft barks’ (p. 127).

To which degree is vocal production affectively bound in other species? To which extent can the developmental trajectory seen in humans (with early independence between certain sounds and particular affective states (Oller et al. 2013)) also be observed in other primates? In fact, are affective decoupling and vocal functional flexibility unique to human vocal ontogeny? In one relevant study, Clay et al. (2015) examined ‘peep’ calls in adult bonobos (*Pan paniscus*), their most common vocalizations, and found that they are produced in a variety of contexts, ranging from seemingly positive (food provisioning) to neutral (travel and resting) and negative (agonistic and alarm) situations. Based on these findings, the authors concluded that bonobos have the capability to produce sounds that are not affectively biased (Clay et al. 2015), suggestive of affective decoupling in vocal production. Their peeps were, however, attributed to broad behavioral contexts (such as feeding or travelling) with no focus on more specific and transient behaviors that may help infer affective contexts, such as when individuals suddenly experience aggression during travelling and feeding bouts. In fact, peeps could well be bound to the expression of one particular affect, which could be common in both feeding and travelling contexts for example. As such, the bonobo data are indicative of their peeps occurring across broad behavioral contexts (‘contextual flexibility’) but may ultimately remain inconclusive in regards to whether affective decoupling and vocal functional flexibility are indeed present in species other than humans.

Similarly, the flexibility with which some call types are expressed in some primate species is only ‘suggestive’ of ‘vocal functional flexibility’ (the use of one vocal type to convey various illocutionary forces on different utterances), and may only correspond to contextual flexibility (the use of one call type in different contexts, with

core commonalities in the illocutionary force conveyed by all utterances). For example, Guinea baboons use a number of calls (e.g. grunts, roar grunts, barks, and wahoos) in a diversity of contexts (Maciej et al. 2013). Guinea baboons’ grunts are used not only in foraging and travelling contexts, but also in affiliative, infant handling and greeting contexts. Does that mean that Guinea baboons show functional flexibility when producing grunts? It is a possibility. However, and in the absence of a methodological focus on potential affective states experienced by the animal, a temporary conclusion is that Guinea baboon grunts likely are ‘contextually flexible’. The possibility that they also are ‘not affectively bound’ (i.e., not bound to the experience and expression of a particular affective state) or ‘functionally flexible’ (i.e., not assigned to the fulfilment of one particular function) awaits empirical confirmation. Indeed, the grunts of Guinea baboons could well be expressed under one particular affective state, and used to meet one single function in a variety of contexts (e.g. maintaining contact with other troop members). One should examine whether inferences about the affective state of animals (e.g. Guinea baboons) can be conducted using the behavioral contexts employed to describe the contextual occurrence of their vocalizations, and whether such analysis suggests that these vocalizations qualify as affectively decoupled.

A second study, also on bonobos (Oller et al. 2019), suggests protophone-like vocal behavior with bonobo infants producing calls that occur in both low or moderate arousal situations, implying no affective binding. This conclusion has been preliminary, however, for the affective quality of the contexts surrounding vocalizations (a reliable marker of illocutionary force and needs in human infants) has proven difficult to discern.

Here, we intended to provide a first evaluation of affective decoupling in infant chimpanzees’ (*Pan troglodytes schweinfurthii*) vocal behavior at a very early age (<12 months). Given the recent studies in both immature and mature bonobos, focusing on the other closest living relatives, the chimpanzees, is crucial to test hypotheses about the evolutionary origins of functionally flexible vocal behavior. What’s more, examination of ‘early’ vocal production is critical for a more direct comparison with findings on human infants (Oller et al. 2013). We focused on two call types, the grunts and the whimpers, as they are acoustically very distinct vocalization categories that are common in young infants (Plooij 1984). Finally, we tried to approach the affective dimension of the context of calling by focusing on transient behavioral cues (e.g. the infant escaping a situation) rather

than broader behavioral contexts (e.g. travelling context).

Grunt calls are of particular importance as they develop into a central component of the vocal repertoire of chimpanzees and contribute to a variety of vocal sequences produced by juveniles, sub-adults, and adults (Crockford and Boesch 2005). For example, grunts complement panting elements during laughter (Leavens, 2009) and when encountering dominant individuals ('pant-grunts') (Laporte and Zuberbühler 2010; Laporte and Zuberbühler 2011). They are also produced upon encountering a food patch or when joining a foraging party ('rough grunts') (Slocombe and Zuberbühler, 2005; Slocombe et al. 2010; Fedurek and Slocombe 2013; Schel et al. 2013; Watson et al., 2015). Finally, they are routinely produced throughout resting or in relaxed social activities (Goodall 1986). Grunts are produced from the first days of life in chimpanzees. Their ontogenetic development has already been studied to some degree in chimpanzees, which has shown some flexibility in usage (Laporte and Zuberbühler 2011). It has been suggested that at least two types of grunts could be distinguished. First, uh-grunts are short, tonal sounds, resembling human vowels {u}, {o}, and {a} (and possible homologous to quasi-vowels in humans), sometimes produced in short series (staccato-grunts) (Kojima 2003; Plooi 1984). The second type are the so-called 'effort' grunts, which are common in immature chimpanzees (Plooi 1984) and are also present in adult chimpanzees, mature and immature humans and other mammals (McCune et al. 1996). The so-called 'effort grunts' are very soft and require the close presence of observers to be reliably heard (Plooi 1984). They received their name from their presence during locomotor activities. Despite Plooi's (1984) suggestion that they could be mere by-products of locomotor activities, he also noted they can occur in the absence of movements. So far, no study has yet offered an acoustical validation of the existence of these diverse types, such that we (and others, see Laporte and Zuberbühler 2011) cannot rely on this distinction.

Another common vocal utterance produced by chimpanzee infants is whimpers (Plooi, 1984; Levréro and Mathevon 2013; Dezechache et al. 2020). They are short, tonal, and often produced in series with an upward shift in fundamental frequency. Contrarily to grunts, whimpers preferentially occur in aversive contexts, likely homologous to human crying or distress calls in other mammals (Plooi 1984). Previous research (Plooi 1984) has suggested the presence of whimper subtypes (single, serial, and human-like whimpers), but again, we are not aware of any systematic acoustical analysis that would

justify this nomenclature. Whimpers are also present in the repertoire of adult chimpanzees, notably in alarm (Tsukahara 1993), food begging (Crockford and Boesch 2005; Slocombe and Newton-Fisher 2005), and physical separation (Crockford and Boesch 2005) contexts.

To start addressing the hypothesis that affective decoupling and vocal functional flexibility evolved before the split between *Pan* and *Homo* lineages, we examined the vocal behavior of six wild chimpanzee infants aged between 0 and 12 months old from the Sonso community of Budongo Forest, Uganda. We analyzed the extent to which vocal production of grunt-like and whimper-like vocalizations were occurring with the so-called positive, negative, or neutral behaviors, as a first step toward evaluating the affective quality of the vocalization contexts. We also took advantage of recent developments of machine learning techniques to the study of animal communication (Mielke and Zuberbühler 2013; Fedurek et al. 2016) to evaluate acoustical differences between calls produced with positive, negative, and neutral markers.

## 2. Methods

### 2.1 Ethics

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

### 2.2 Subjects and data collection

Data were collected in the Sonso community of the Budongo Forest Reserve, Uganda (Reynolds 2005) between February–June 2014, December 2014, and March–June 2015. This community comprises around seventy individuals well habituated to human observers. The natural behavior of  $N=7$  infants was video recorded continuously during focal animal sampling (Altmann 1974), using Panasonic HC X909/V700 cameras, with a Sennheiser MKE-400 shotgun microphone. Six of those infants produced enough calls to be further considered for data analysis (see Table 1 for details).

### 2.3 Behavioral data analysis

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

As of today, there is no definitive repertoire of infant chimpanzee vocal behaviors, only suggestive classifications (Plooi 1984; Plooi et al. 2014). The categories

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

ID	Sex	Min. age (in months)	Max. age (in months)	N, whimper-like vocalizations	N, grunt-like vocalizations	N, 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

Also given are the numbers of grunt-like and whimper-like vocal behaviors collected, as well as grunt-like vocalizations acoustically analyzed.

**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', that is, 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 while moving (locomotion or clear leaning of the body) toward 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 toward 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 toward 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/kin	Seeking contact with the mother/kin 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

used in this research are based on first author's assessment. This assessment proved reliable when confronted to an independent assessment with Derry Taylor, using vocalizations from infant and juvenile semi-wild chimpanzees from the Chimfunshi Wildlife Orphanage, Zambia, collected by D.T. A total of 160 vocalizations were indeed classified as belonging to either the 'grunt', 'whimper', 'scream', or 'laughter' category. Agreement was excellent ( $k = 0.77$ ) and even better when considering only 'grunts' and 'whimpers' ( $k = 0.92$ ).

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

in Table 2). This list was established following data collection, with some inspiration from the behavioral categories established by Plooij (1984) during his study with the infant chimpanzees of the Gombe community between 1971 and 1973. As in the original human study (Oller et al. 2013), we reckoned the behavior of the infants could offer a reliable source of information unto their affective state, as a first step toward establishing affective descriptions of contexts. In fact, we originally aimed at mimicking their coding strategy, using categories appropriate to the study of wild infant chimpanzees. The affective quality of the infants' behavior was



classified as ‘positive’ if it showed one of the following four behaviors: (1) ‘play’; (2) giving or receiving ‘grooming’ (note that allo-grooming was never observed in our infants); (3) ‘feeding’; and (4) ‘social approach’. See [Table 2](#) for details.

The affective context was classified as ‘neutral’ if it showed one of the following behaviors: (5) ‘resting’; (6) ‘moving’; (7) ‘manipulating objects’ without playful postures; or (8) ‘greeting without approach’. See [Table 2](#) for details.

Infant behavior was classified as ‘negative’ if it showed one of the following behaviors: (9) ‘nuzzling’; (10) ‘begging’; (11) ‘hiding’; (12) ‘contact mother/kin’ was coded if infants were urgently seeking contact with the mother or a kin when contact was not already established between them; and (13) ‘escaping’. See [Table 2](#) for details.

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

In order to evaluate the evenness of the distributions of grunts and whimpers across affective contexts, we calculated, for each infant, and for grunts and whimpers separately, the dominance of one affect over the two others, using the Berger–Parker dominance index ([Morris et al. 2014](#)):

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

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

Dominance values (one per infant per call type) were compared between grunts and whimpers using a paired Wilcoxon signed-rank test. These analyses were carried out using R (version 3.6.1; [R Core Team, 2018](#)) and R Studio (version 1.2.1335; [RStudio Team, 2015](#)).

## 2.4 Acoustic analysis

Acoustic data analysis focused on grunts for they were the only vocal category for which at least two of the affective contexts were well represented. The acoustic structure of whimpers has been analyzed as part of

another study ([Dezecache et al. 2020](#)).  $N = 180$  grunts were extracted from independent vocal behaviors. For each affective context, sixty were randomly selected. Following extraction, we used MATLAB (MathWorks Inc., Natick, MA) for the acoustic data analysis, consisting of features extraction, feature selection, and call classification. We first pre-processed the audio files by applying a band pass filter from 50 to 4,000 Hz and normalized the signals using the following function:

$$\text{signal} = \left( \text{signal} - \text{mean}(\text{signal}) \right) / \max \left( \left| \text{signal} - \text{mean}(\text{signal}) \right| \right)$$

### 2.4.1 Feature extraction and selection

We first ran a feature extraction algorithm to reduce redundancy of information and computational efforts in classifying the grunts and to maximize the generalization ability of the classifier ([Tajiri et al. 2010](#)). A popular method is extraction of mel-frequency cepstral coefficients (MFCCs) ([Supplementary Fig. S2](#)). MFCCs represent the envelope of the short-time power spectrum, as determined by the shape of the vocal tract ([Logan 2000](#)). The idea behind the extraction of MFCCs is to obtain a comprehensive representation of the frequencies that compose an audio bout, while putting emphasis on certain frequency bands. Although a typical spectrogram linearly scales frequencies (i.e. each frequency bin is spaced an equal number of Hertz apart), the mel-frequency scale is a logarithmical spacing of frequencies. MFCCs are routinely used in speech recognition and are gaining prominence in the field of animal communication (see for instance, [Fedurek et al. 2016](#), in chimpanzees). The use of MFCCs to represent sounds can be considered to be a solution preferable to the selection of a limited set of parameters to describe acoustical phenomena (such as these related to the shape of the fundamental frequency) for it offers a more comprehensive representation of sounds. In the context of our work (the aim of which was to evaluate potential distinctiveness between grunts occurring in the so-called positive, neutral, and negative contexts), MFCCs appeared as the optimal solution to the problem of a false-negative conclusion.

We divided the calls into segments of 25 ms length and 10 ms steps between two successive segments. We warped twenty-six spectral bands and returned thirteen cepstra, which resulted in feature dimensions of thirteen values each. We then took the mean and co-variances of each cepstra over the collection of feature segments, resulting in a 13-value vector and a  $13 \times 13$ -value

matrix, respectively, and concatenated to 104-unit vectors (Mandel and Ellis 2005: 594–599) (Fig. 3). We applied feature scaling to [0–1] and mean normalization.

Second, we performed a feature selection procedure: too many feature dimensions are not useful for producing reliable classification systems, whereas low sample numbers can lead to over-fitting to noisy feature dimensions. We therefore selected a subset of the original feature dimensions and evaluated classification performance based on sequentially selected feature sets until there was no improvement in performance. At this end, we subdivided the entire dataset into a training (75%) and a test dataset (25%) and applied a *t*-test on each feature dimension, comparing values of given feature dimension sorted by predefined class labels (e.g. grunts occurring with negative (1) vs. positive (2) affects) and used *P*-values as a measure separability of the two classes. We plotted the *P*-values as an empirical cumulative distribution function (eCDF) to get an understanding of how well each feature separated the two classes and how many features contributed to a significant separation (5% level). We ran this procedure twenty times for each comparison and plotted the results individually (gray lines) and the mean of all repetitions (black line) (Fig. 2A). The classification routines were then independently run either on feature dimensions selected according to the discrimination power (decreasing order) (orange lines in Fig. 2B), as shown in the eCDF plots (Fig. 2A). Such procedure is referred to as a simple filter approach on feature selection, where general characteristics of the extracted features are taken into consideration when selecting feature dimensions, without subjecting them to a classifier. We also applied a more extensive procedure of feature selection by sequentially selecting feature dimensions by adding (forward search) feature dimensions, referred to as sequential feature selection (black lines in Fig. 2B). As part of this method, the algorithm searched the best feature dimensions (predictors) according to their individual classification performance in the given subset of data. For each candidate feature subset (predictor), the algorithm performed a ten-fold cross-validation procedure with different training and test subsets. After computing the mean performance values for each candidate feature subset, the algorithm chooses the candidate feature subset with minimal misclassification. For both methods, we systematically varied the number of features used for classification (*x*-axis in Fig. 2B). The selected features from a single run of the sequential search algorithm are illustrated in Fig. 2C. Scales reflect the feature-scaled and normalized values, as a result of

feature extraction, from which the grand means (i.e. for each feature dimensions across all data) were subtracted. This measure was used to visually highlight differences and was not used in further analyses.

## 2.4.2 Classification

We used support vector machine (SVM) with a radial basis function Kernel (Vert et al. 2004) for the classification of calls according to the class labels (so-called negative, neutral, and positive affective contexts). A classification procedure contains a training phase followed by a test phase. We separated training samples and labeled them according to an attribute of interest (e.g. negative (1) vs. positive (2) affective contexts). The algorithm then created a model that optimally separates the two classes. In the test phase, samples without attribute labels were fed into the model to measure its generalization performance. We used the SVM implementation from LIBSVM toolbox (Chang and Lin 2011). To evaluate how the classification results generalize to a novel and independent dataset, we ten-fold cross-validated the classification process and optimized the parameters *C* and gamma (Fedurek et al. 2016), with the *C* taking 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 single individuals contributed solely to the classification outcome, we ran a leave-one-out algorithm, where the procedure described above was re-run six times, excluding one of the individuals in each run. We applied one-sample *t*-tests to compare the classification scores with a 50% baseline condition. The 50% baseline results from the pairwise comparisons of affective contexts (positive, neutral, negative). To ensure samples were normally distributed (a key assumption behind the use of one-sample *t*-tests), we used Lilliefors test prior to each comparison at a significance level of 5%. In cases where data samples were not normally distributed, we used a one-sample Kolmogorov–Smirnov test. All reported *P*-values were adjusted for multiple comparisons using Bonferroni corrections.

## 2.4.3 Feature evaluation

To evaluate whether certain feature dimensions are particularly critical for the classification of grunts, we assessed whether feature dimensions have been repeatedly used by the classifier overall in the classification of grunts. We therefore considered the three types of comparisons, positive vs. neutral, positive vs. negative, and neutral vs. negative grunts, as well as the two feature evaluation algorithms (simple feature selection and sequential feature selection). Each comparison was ten-



fold cross-validated. We then calculated the empirical distribution of the ten features with best classification power, as determined by the feature selection algorithms (see above). Also, we determined a random distribution of ‘best features’ for each comparison by randomly selecting 10 out of 104 features. The frequency distribution across all comparisons was determined and 95% confidence intervals were calculated by running the procedure 1,000 times. We then traced back the significant feature dimensions to the underlying frequency bands in Hertz.

### 3. Results

#### 3.1 Types of vocal utterances

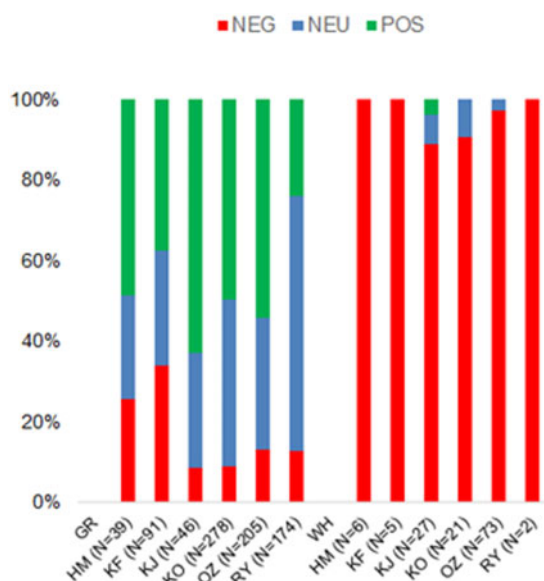
We inspected  $N=1,016$  vocal occurrences, of which  $N=967$  could be classified as either ‘grunts’ ( $N=833$ ) (corresponding to a rough, harsh and noisy sound) or ‘whimpers’ ( $N=134$ ) (usually a series of low-pitch tonal calls with increase in fundamental frequency throughout the series). Other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

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

**Grunts:** In total, 44.8% of grunt-like vocalizations co-occurred with contexts we classified as ‘positive’, 40.9% with ‘neutral’, and 14.3% with ‘negative’. When considering each individual separately, a similar picture emerged (Fig. 1), with most grunt-like vocalizations co-occurring with ‘positive’ and ‘neutral’ contexts. We found dominance to be relatively low in grunts, varying from 0.37 and 0.63 (mean = 0.53; SD = 0.10), suggesting a stable and relative evenness in the affective distribution of grunts, as defined by our coding system (Table 2).

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

**Grunts vs. whimpers:** When comparing the distributional evenness of grunts vs. whimpers, we found dominance to be statistically higher in whimpers than in



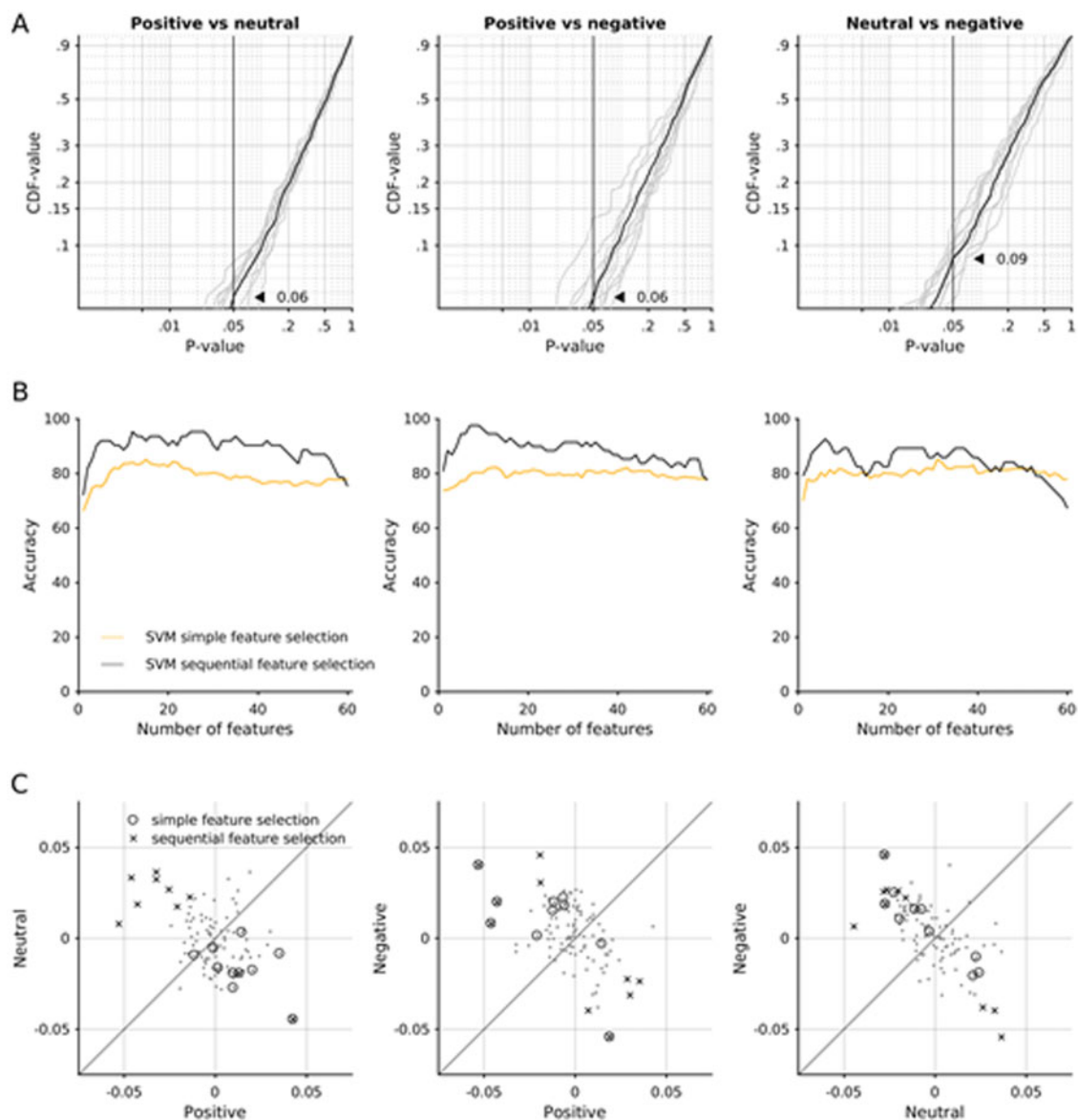
**Figure 1.** Proportion of grunt-like (GR) and whimper-like (WH) vocal behaviors recorded with negative (NEG), neutral (NEU), and positive (POS) affective categories of behaviors, for each individual separately. Numbers between brackets indicate the number of GR and WH calls contributed by each individual.

grunts (paired Wilcoxon signed-rank test:  $V=21$ ,  $P=0.031$ ).

#### 3.3 Acoustic variants of grunts

We classified the  $N=180$  grunts ( $N=60$  per affective contexts) according to their association with the so-called positive, neutral, negative contexts in order to test for the presence of acoustic variants. In the first step, we followed the feature extraction procedure by extracting the means and covariances of MFCCs for each call, and compared these values according to the calls’ associations (e.g. positive vs. negative) using  $t$ -tests. We displayed the resulting  $P$ -values in an eCDF (Fig. 2A). We found that 5–10% of all features showed significant differences between the class labels at a 5%-significance level. In other words, 5–10 of 104 feature dimensions had strong discrimination power to distinguish between grunts pertaining to the various affective contexts.

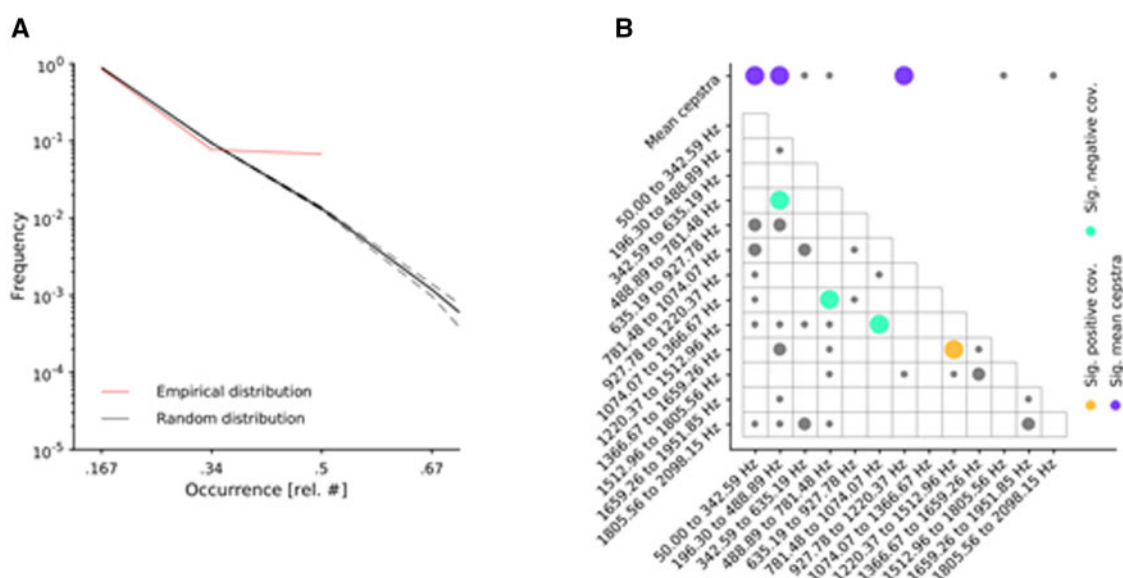
With the simple feature selection algorithm, the SVM correctly discriminated between classes at up to 80% (positive vs. neutral:  $M=78.99$ ,  $SD=3.53$ ,  $t(59)=63.69$ ,  $P<0.001$ ; positive vs. negative:  $M=79.58$ ,  $SD=1.83$ ,  $t(59)=125.37$ ,  $P<0.001$ ; neutral vs. negative:  $M=80.44$ ,  $SD=2.06$ ,  $t(59)=114.26$ ,  $P<0.001$ ; or orange lines in Fig. 2B). A substantial improvement was found when sequentially selecting feature dimensions:



**Figure 2.** Feature selection and classification performances. The columns represent the comparisons of affects during which the vocal utterance occurred. (a) For each feature dimensions the discrimination power of the two classes (e.g. positive vs. neutral) was evaluated using a *t*-test. *P*-values are shown as an eCDF. Gray lines show the results of individual runs of evaluation; black lines show the means of individual runs. Indicated with arrow heads are the proportions of feature dimensions that significantly discriminate between the two classes tested. (b) The classification performances are shown for the SVM classifier relying on feature dimensions extracted through a simple feature selection (orange lines) and a sequential feature selection procedure (black lines). (c) Feature selection outcomes are shown for simple (circles) and sequential feature selection (blue x-s) as overlays on all feature dimensions (gray dots).

SVM correctly classified samples at up to 95% (positive vs. neutral:  $M = 89.56$ ,  $SD = 4.84$ ,  $t(59) = 143.42$ ,  $P < 0.001$ ; positive vs. negative:  $M = 88.72$ ,  $SD = 4.49$ ,  $t(59) = 153.11$ ,  $P < 0.001$ ; neutral vs. negative:  $M = 84.27$ ,  $SD = 5.23$ ,  $t(59) = 124.91$ ,  $P < 0.001$ ; black

lines in Fig. 2B). For all comparisons, chance levels were 50% due to the two-class comparisons applied. We, therefore, used one-sampled *t*-tests. The classification scores in all (but one) comparisons fulfilled the requirement of normal distribution. The first comparison



**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 nonsignificant feature dimensions.

(feature-selection algorithm, positive vs. neutral) was not conform with a normal distribution and was, thus, re-evaluated using a one-sampled Kolmogorov–Smirnov test, resulting in the following values ( $k_s = 0.17$ ;  $P < 0.001$ ).

We further illustrated the simple feature selection outcomes by highlighting the feature dimensions selected (circles in Fig. 2C) among the feature dimensions not selected (gray dots). Further, the features selected via the sequential feature selection are marked with x's. The sequential feature selection yields better performance through sequential combinations of feature dimensions that, on average, fall more distal to the diagonal mid-line than the feature dimensions selected by the simple feature selection process. Sequential feature selection, to a large extent, included feature dimensions not selected by the simple feature selection method.

We further ensured that each individual was not contributing solely to the classification results of various contrasts. As can be seen in Supplementary Fig. S1, the classification performance did not improve nor deteriorate systematically when one individual was removed at a time, suggesting no effect due to caller identity (the average  $t$ -value of one-sample  $t$ -tests is  $97.52 \pm 30.25$  (SD); all  $P$ -values were smaller than 0.001).

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

routes at low computational loads. To assess whether certain feature dimensions (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 (Fig. 3A).

To describe the frequency bands explaining significant variances between classes of calls, we traced back the frequency bands underlying the significant feature dimensions, that is, covariances of cepstra, and determined the sign of the covariances. We found negative covariances between the following frequency bands (Fig. 3B): (1) band 2 (196.30–488.89 Hz) and band 4 (488.89–927.78 Hz); (2) band 4 (488.89–927.78 Hz) and band 8 (1,074.07–1,366.67 Hz), and band 6 (781.48–1,074.07 Hz) and band 9 (1,220.37–1,512.96 Hz). We found a positive covariance between the frequency bands 9 (1,220.37–1,512.96 Hz) and 10 (1,366.67–1,659.26 Hz). Mean cepstra were significantly contributing in the frequency bands from (1) 50 to 342.59 Hz, (2) 196.30 to 488.89 Hz, and (3) 927.78 to 1,220.37 Hz.

#### 4. Discussion

Oller and colleagues (Oller and Griebel 2004; Oller et al. 2013, 2016; Jhang and Oller 2017) posited 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 termed ‘affective decoupling’. Human infants regularly vocalize in such a way, in supposed contrast to the relative inflexibility of vocalizations in nonhuman primates (Pollick and 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 are not tied to the experience and expression of one particular affective state (Oller and Griebel 2004; Oller et al. 2013). In contrast primate (and more largely, ‘animal’) vocal behavior is said to be affectively bound, with particular calls being used to express particular affective state, ultimately constraining their signaling function. The view that primate vocalizations are mere read-outs of affective states is widespread in the animal communication literature (Marler 1980; Goodall 1986; Hammerschmidt and Fischer 2008; Gruber and Grandjean 2017).

In the current study, we looked at one of our closest living relative species, the chimpanzees. We focused on the grunt-like and whimper-like calls of young chimpanzee infants, using a novel coding strategy and state-of-the-art acoustic analysis. We elaborated a workable coding system, which was meant to provide insights into the affective state of infant chimpanzees, as seen in Oller et al. (2013), and to allow a first comparison between human and chimpanzee infants. We found that grunt-like calls were produced frequently by chimpanzee infants within both positive and neutral contexts, and less commonly with negative affective contexts. Importantly, the presence of grunts in contexts of low-to-mild arousal is consistent with the hypothesis of vocal functional flexibility (Oller et al. 2019), and so is the finding that grunts occur in similar proportion with contexts we deemed positive and neutral (Oller et al. 2013).

On the other hand, whimper-like vocalizations seemed to be confined to behaviors and contexts we associated with negative affective states. Their near absence in positive and neutral contexts suggests that they represent an affectively bound vocalization that has evolved to signal a narrow range of needs and one single (negative) affective valence, similar to cries in human infants (Oller et al. 2013), to which they may functionally correspond (Goodall 1986). Our results therefore

suggest that chimpanzee grunts are not bound to one affective context. They may qualify as a functionally flexible vocal unit, consistent with 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 humans, the finding has been that infants use protophones with a diversity of affects, with mothers reacting consequently, showing that infant calls are fully functionally flexible (Oller et al. 2013). In these studies, the mothers’ behavior could be examined, although protophones were not always socially directed (Oller and Griebel 2004). Protocols where mothers may be asked to interact with toddlers may yield to responsiveness from the mothers, regardless of the infant’s affective state (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 was 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 establish their functionality (Fischer et al. 2013; Zuberbühler 2014; Fischer 2016), this would require either playing the infants’ calls in its own presence (which is ethically inappropriate) or playing the calls of another infant to a mother (which may not trigger any reaction at all in the nongenetically related mother). Another possibility is that the sounds we examined are not meant to be fully functional, and could be considered to be vegetative sounds. The fact that they may not appear socially directed should, however, not speak against the hypothesis that they are affectively decoupled, for the fact that a given vocal unit

is independent from one particular affective valence is orthogonal with the fact that it is social directed or not. Our results are compatible with grunts being a functionally flexible call type in young chimpanzees, but do not yet demonstrate this, for the reactions of the mothers (and therefore, the function of the calls) could not be directly assessed.

Grunts (and other close calls (Oller and Griebel 2004)) are the promising class of vocalizations to investigate the evolutionary origins of vocal functional flexibility. In a number of species (such as the vervet monkeys (Cheney and Seyfarth 1982), western gorillas (Salmi et al. 2013), sooty mangabeys (Range and Fischer 2004), chacma (Meise et al. 2011), Guinea (Faraut et al. 2019; Maciej et al. 2013), and olive baboons (Ey and Fischer 2011; Silk et al. 2018)), grunts are used flexibly and can occur in a variety of contexts. So far, such evidence speaks in favor of grunts being a contextually flexible vocal unit (i.e. a vocal unit whose function can be fulfilled in a diversity of contexts). Future research should try delving into the affective state animals likely experience and express when producing grunts, to confirm whether these also displays affective decoupling (i.e. the independence between grunt production and the experience of one particular affective valence) and functional flexibility (i.e. the capacity of grunts to fulfil a variety of functions on different occasions). If the term ‘functional flexibility’ could appear misleading, its use in the field of child development should encourage the animal communication community to employ it, such that more fruitful cross-disciplinary work can best take place.

Our second main finding was systematic acoustic differences between grunts given with the so-called positive, neutral, and negative behaviors, which enabled us to segregate acoustic variants of grunts into these categories. Acoustical differences linked to the affect surrounding vocal production are common in humans as in other animals (Williams and Stevens 1972; Banse and Scherer 1996; Briefer 2012; Aucouturier et al. 2016; Arias et al. 2018; Ponsot et al. 2018, Goupil et al. 2019). Our data suggest that there is inter-gradation between grunt types, with differences in acoustics relating to differences in contexts. Grunts, in other words, represent a coherent and unified call type that can manifest itself in acoustic variants in relation to the affective contexts in which they are produced. It is possible that grunts acoustically vary with arousal of the animal (as seen in other primate species (Rendall 2003)), although positive and negative circumstances could, in principle, be equally arousing.

How exactly functionally flexible vocalizations produced by human infants transition into speech sounds has been described in previous studies (Oller et al. 1976; Elbers and Ton 1985; de Boysson-Bardies and Vihman 1991; de Boysson-Bardies 1993; Oller 2000; Boysson-Bardies 2001; Nathani et al. 2006). Chimpanzee infants may produce grunts in ways consistent with the functional flexibility hypothesis but they of course never produce speech sounds and, historically, have failed to acquire human speech utterance even after extensive training (Hayes and Hayes 1951). Instead, infant chimpanzee grunts may gradually develop into call variants with seemingly relatively narrow biological functions (Slocombe and Zuberbühler 2005; Slocombe and Zuberbühler 2010; Laporte and Zuberbühler 2011; Watson et al. 2015), with clear acoustical boundaries notably between grunts used to greet conspecifics (‘pant-grunts’ (Laporte and Zuberbühler 2011)) and those produced upon encountering food (‘rough’ or ‘food grunts’ (Slocombe and Zuberbühler 2005)). It is possible that the acoustic boundaries we identified between the grunts produced across affective states (under our nomenclatures and coding system) are the foundation of acoustic diversification in adults, although the categories used here (for instance, feeding and social approach are together considered ‘positive’) are not consistent with the vocal differentiation seen in adults (the grunts produced in feeding vs. social approach situations are acoustically distinct in adults (Goodall 1986; Crockford, 2020). Alternatively, those calls may simply disappear and be absent from the adult repertoire, one causal factor being the relative absence of social reinforcement (including contingent vocal responses (Ghazanfar et al. 2019)) associated with grunt production, as compared with the frequent maternal reactions to distress calls (Dezecache et al. 2020).

Our tentative to explore the affective state of the infant may be seen as preliminary, insofar as the categories we have used do not represent read-outs of physiological states. This being said, the acoustical differentiation we found spoke in favor of the appropriate character of our affective distinctions. Ideally, other cues should be considered, such as the infants’ facial expressions or the mothers’ behavior. This approach would, however, face considerable challenges. We found that infant facial movements are extremely fast and fluid, which prevented us from reliable coding particularly in the wild. For this reason, the behavioral context of the infant alone (although imperfect and probably still questionable) was the most relevant available cue to approach the affective dimension of the situation. Although we must again acknowledge the limitations pertaining to



the fact that judgments of infants' affect were made based on the infants' behavioral contexts and done so by a human observer, the results of the acoustic analysis are providing support for the approach used to categorize affect in the present work. Future studies should investigate the affective impact of other communicative signals used by infants, such as gesture and facial behavior, and their combinations (Fröhlich et al. 2018; Fröhlich and Hobaiter 2018).

Besides the limitations pertaining to our coding system (and its shortcomings with respect to the production of inferences regarding infants' affective states), one other limitation of this study is the small sample size, as we could only collect enough data from six infants. One particular difficulty with collecting data from such young chimpanzee infants is that some of their calls (notably a large part of their grunts) are very soft (a point also acknowledged by Plooij (1984)) and can only be heard from close, limiting the number of individuals whose mothers are unwary enough of continuous and long-lasting human observational efforts. We could not use already published data, because, to the best of our knowledge, no previous studies on the vocal behavior of wild infant chimpanzees (such as Plooij 1984; Laporte and Zuberbühler 2011; Plooij et al. 2014) used a coding system amenable to inferences about the affective state of the infant.

In latest research, the comparative volubility (quantity of sounds produced in a given period of time) of human infants and other animals (Ghazanfar and Takahashi 2014; Takahashi et al. 2015; Oller et al. 2019), and the privileged function of protophone-like vocalizations to increasingly elicit social interactions and vocal turn-taking with caregivers (Yoo et al. 2018; Oller et al. 2019). In humans, affectively flexible vocalizations appear to occur more often than affectively bound vocalizations (such as crying) (Oller et al. 2019). They occur in solitary contexts where infants invest in vocal practice and exploration. They also occur in interactive contexts, so as to elicit and regulate social interactions with caregivers. Caregivers appear to detect the functional difference between protophones (as potentially interactive calls) and other calls (such as cries), where caregiver intervention is solicited (Yoo et al. 2018). Comparison with bonobo infants suggested much higher rate of production of affectively flexible vocalizations and much higher vocal investment in social interactions in human infants (Oller et al. 2019). Whether human infants also are comparably more 'talkative' than their chimpanzee counterparts is a question we need to be exploring. This should be preferably investigated in captive or semi-captive settings, where true calling rate can

be assessed, for video monitoring is less likely to be interrupted and for levels of ambient noise could be comparatively less problematic. Such problems have already been acknowledged by Oller et al. (2019) regarding previous report on grunting behavior in wild chimpanzees (Laporte and Zuberbühler 2011). Data from the vocal development of one captive chimpanzee indicated lower volubility than in humans (Kojima 2003). Future studies should evaluate this fact with a larger sample.

Our study suggests that, insofar as one can delve into the affective state of infants using our coding system, chimpanzees may possess a feature that is fundamental to the development of speech in humans, the ability to produce vocalizations that are not strongly bound to the experience and expression of one particular affective valence. However, we should expect that future research will reveal further examples. For instance, coo calls in several macaque species (Hsu et al. 2005; Owren and Casale 1994), wahoos of baboons (Maciej et al. 2013), or grunts of a number of primate species seem to be given in a variety of contexts, a precondition for affective decoupling in vocal production, itself a prerequisite for vocal functional flexibility. More largely, close calls appear to be excellent candidates (Oller and Griebel 2004). Importantly, methodologically efforts to infer the affective states of the animals should be made in order for affective decoupling to be hypothesized.

Future research will also have to address the question of how selection favored acoustic diversification of functionally flexible vocal behavior into speech in humans. The main driver for this transition, it has been argued, may have been the highly cooperative breeding system of humans, with infants regularly looked after by individuals other than the mother, which requires infants to become more active agents in forming social bonds from a much younger age than in great ape infants (Zuberbühler 2012; Ghazanfar et al. 2019).

Cooperative breeding, in this view, may thus have transformed a functionally flexible vocal system into the uniquely human way of using vocal signals to interact socially. Another complementary reasoning is that humans' high altriciality selected for the most vocal individuals, capable of attracting caregivers (Locke 2006). The relative contribution of both factors through mapping the phylogenetic distribution of affective decoupling and vocal functional flexibility remains to be investigated.

## Data availability

Data are accessible from both corresponding Authors.



## Supplementary material

Supplementary material is available at the OUP webpage.

## Conflict of interest statement

No conflict of interest.

## Acknowledgments

We thank UWA and UNCST for permission to conduct the study, Geoffrey Muhanguzi, Caroline Asiimwe and Sam Adué for their support in the field, Derry Taylor, the editor and two anonymous reviewers for critical comments on earlier versions of the article. We are grateful to the Royal Zoological Society of Scotland for providing core funding to the Budongo Conservation Field Station.

## Funding

The research was supported by a Fyssen Fellowship, British Academy Newton International Fellowship (NF171514), and CAP2025 funding awarded to G.D., funding from the European Union's Seventh Framework Programme for research, technological development, and demonstration (Grant agreement no. 283871) awarded to KZ the Swiss National Science Foundation (PZ00P3\_154741) awarded to C.D.D., the Startup-funding of Taipei Medical University (108-6402-004-112) awarded to C.D.D., and (310030\_185324) awarded to KZ.

## References

Altmann, J. (1974) 'Observational Study of Behavior: Sampling Methods', *Behaviour*, 49/3–4: 227–66

Arbib, M. A. et al. (2008) 'Primate Vocalization, Gesture, and the Evolution of Human Language', *Current Anthropology*, 49/6: 1053–76

Arias, P., Belin, P., and Aucouturier, J.-J. (2018) 'Auditory Smiles Trigger Unconscious Facial Imitation', *Current Biology*, 28/14: R782–3.

Aucouturier, J.-J. et al. (2016) 'Covert Digital Manipulation of Vocal Emotion Alter Speakers' Emotional States in a Congruent Direction', *Proceedings of the National Academy of Sciences*, 113/4: 948–53

Banase, R., and Scherer, K. R. (1996) 'Acoustic Profiles in Vocal Emotion Expression', *Journal of Personality and Social Psychology*, 70/3: 614–36

Boë, L.-J. et al. (2017) 'Evidence of a Vocalic Proto-System in the Baboon (*Papio Papio*) Suggests Pre-Hominin Speech Precursors', *Plos One*, 12/1: e0169321 <https://doi.org/10.1371/journal.pone.0169321>

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

Briefer, E. F. (2012) 'Vocal Expression of Emotions in Mammals: Mechanisms of Production and Evidence', *Journal of Zoology*, 288/1: 1–20.

Burkart, J. M. et al. (2007) 'Other-regarding Preferences in a Non-Human Primate: Common Marmosets Provision Food Altruistically', *Proceedings of the National Academy of Sciences*, 104/50: 19762–6.

——, Hrdy, S. B., and Van Schaik, C. P. (2009) 'Cooperative Breeding and Human Cognitive Evolution', *Evolutionary Anthropology*, 18/5: 175–86.

Chang, C.-C., and Lin, C.-J. (2011) 'LIBSVM: A Library for Support Vector Machines', *ACM Transactions on Intelligent Systems and Technology (TIST)*, 2/3: 27.

Cheney, D. L., and Seyfarth, R. M. (1982) 'How Vervet Monkeys Perceive Their Grunts: Field Playback Experiments', *Animal Behaviour*, 30/3: 739–51.

Clay, Z., Archbold, J., and Zuberbühler, K. (2015) 'Functional Flexibility in Wild Bonobo Vocal Behaviour', *PeerJ*, 3: e1124.

Crockford, C. (2020). 'Why Does the Chimpanzee Vocal Repertoire Remain Poorly Understood? And What can be Done About It', in Boesch C. and Wittig R. (eds) *The Tai Chimpanzees: 40 Years of Research*. Cambridge: Cambridge University Press.

——, and Boesch, C. (2005) 'Call Combinations in Wild Chimpanzees', *Behaviour*, 142/4: 397–421.

de Boysson-Bardies, B. (1993). 'Ontogeny of Language-Specific Syllabic Productions' in de Boysson-Bardies B., de Schonen S., Jusczyk P., McNeillage P., and Morton J. (eds) *Developmental Neurocognition: Speech and Face Processing in the First Year of Life*, pp. 353–63. Netherlands: Springer.

——, and Vihman, M. M. (1991) 'Adaptation to Language: Evidence from Babbling and First Words in Four Languages', *Language*, 67/2: 297–319.

Dezecache, G. et al. (2020) 'A Machine Learning Approach to Infant Distress Calls and Maternal Behaviour of Wild Chimpanzees', *Animal Cognition*, 10.1007/s10071-020-01437-5.

Elbers, L., and Ton, J. (1985) 'Play Pen Monologues: The Interplay of Words and Babbles in the First Words Period', *Journal of Child Language*, 12/3: 551–65.

Ey, E., and Fischer, J. (2011). 'Keeping in Contact: Flexibility in Calls of Olive Baboons' in Sommer V. and Ross C. (eds) *Primates of Gashaka*, pp. 413–436. New York: Springer.

Faraut, L. et al. (2019) 'How Life in a Tolerant Society Affects the Usage of Grunts: Evidence from Female and Male Guinea Baboons', *Animal Behaviour*, 153: 83–93.

Fedurek, P., and Slocombe, K. E. (2013) 'The Social Function of Food-Associated Calls in Male Chimpanzees', *American Journal of Primatology*, 75/7: 726–39.

——, Zuberbühler, K., and Dahl, C. D. (2016) 'Sequential Information in a Great Ape Utterance', *Scientific Reports*, 6: 38226.

- Fischer, J. (2016) 'Playback Experiments' in Fuentes A. (ed) *The International Encyclopedia of Primatology*, pp. 1–2. Hoboken, NY: Wiley Blackwell.
- , Noser, R., and Hammerschmidt, K. (2013) 'Bioacoustic Field Research: A Primer to Acoustic Analyses and Playback Experiments with Primates', *American Journal of Primatology*, 75/7: 643–63
- Fitch, W. T. (2018) 'The Biology and Evolution of Speech: A Comparative Analysis', *Annual Review of Linguistics*, 4/1: 255–79
- et al. (2016) 'Monkey Vocal Tracts Are Speech-Ready', *Science Advances*, 2/12: e1600723.
- Fröhlich, M., and Hobaiter, C. (2018) 'The Development of Gestural Communication in Great Apes', *Behavioral Ecology and Sociobiology*, 72/12: 194.
- , Wittig, R. M., and Pika, S. (2018). *The Ontogeny of Intentional Communication in Chimpanzees in the Wild Developmental Science*, 22/1: e12716.
- Ghazanfar, A. A., Liao, D. A., and Takahashi, D. Y. (2019) 'Volition and Learning in Primate Vocal Behaviour', *Animal Behaviour*, 151: 239–47
- , and Takahashi, D. Y. (2014) 'The Evolution of Speech: Vision, Rhythm, Cooperation', *Trends in Cognitive Sciences*, 18/10: 543–53.
- Goodall, J. (1986). *The Chimpanzees of Gombe: Patterns of Behavior*. Cambridge, MA: Harvard University Press.
- Goupil, L. et al. (2019) *Influence of Vocal Feedback on Emotions Provides Causal Evidence for the Self-Perception Theory*, <https://doi.org/10.1101/510867>.
- Gruber, T., and Grandjean, D. (2017) 'A Comparative Neurological Approach to Emotional Expressions in Primate Vocalizations', *Neuroscience & Biobehavioral Reviews*, 73: 182–90
- Hammerschmidt, K., and Fischer, J. (2008). 'Constraints in Primate Vocal Production'. In Oller K. and Griebel U. (eds) *Evolution of Communicative Flexibility: Complexity, Creativity and Adaptability in Human and Animal Communication*. Massachusetts: MIT Press.
- Hayes, K. J., and Hayes, C. (1951) 'The Intellectual Development of a Home-Raised Chimpanzee', *Proceedings of the American Philosophical Society*, 95/2: 105–9
- Hrdy, S. B. (2007). 'Evolutionary Context of Human Development: The Cooperative Breeding Model' in Salmon C. A. and Shackelford T. K. (eds) *Family Relationships: An Evolutionary Perspective*. New York: Oxford University Press.
- Hsu, M. J., Chen, L.-M., and Agoramoorthy, G. (2005) 'The Vocal Repertoire of Formosan Macaques, *Macaca cyclops*: Acoustic Structure and Behavioral Context', *Zoological Studies*, 44/2: 275.
- Jhang, Y., and Oller, D. K. (2017) 'Emergence of Functional Flexibility in Infant Vocalizations of the First 3 Months', *Frontiers in Psychology*, 8, 615.
- Jürgens, U. (1976) 'Reinforcing Concomitants of Electrically Elicited Vocalizations', *Experimental Brain Research*, 26/2: 203–14
- (1979) 'Vocalization as an Emotional Indicator', *Behaviour*, 69/1–2: 88–117
- Kojima, S. (2003). *A Search for the Origins of Human Speech: Auditory and Vocal Functions of the Chimpanzee*. Kyoto: Kyoto University Academic Press.
- Kramer, K. L. (2010) 'Cooperative Breeding and its Significance to the Demographic Success of Humans', *Annual Review of Anthropology*, 39/1: 417–36.
- Lameira, A. R., and Shumaker, R. W. (2019) 'Orangutans Show Active Voicing Through a Membranophone', *Scientific Reports*, 9: 12289.
- Laporte, M. N. C., and Zuberbühler, K. (2011) 'The Development of a Greeting Signal in Wild Chimpanzees', *Developmental Science*, 14/5: 1220–34.
- Laporte, M. N., and — (2010) 'Vocal Greeting Behaviour in Wild Chimpanzee Females', *Animal Behaviour*, 80/3: 467–73.
- Leavens, D. A. (2009) 'Animal Communication: Laughter is the Shortest Distance between Two Apes', *Current Biology*, 19/13: R511–13.
- Levréro, F., and Mathevon, N. (2013) 'Vocal Signature in Wild Infant Chimpanzees', *American Journal of Primatology*, 75/4: 324–32.
- Lieberman, P. (2017) 'Comment on "Monkey Vocal Tracts Are Speech-Ready"', *Science Advances*, 3/7: e1700442.
- Locke, J. L. (2006) 'Parental Selection of Vocal Behavior', *Human Nature*, 17/2: 155–68.
- Logan, B. (2000) 'Mel Frequency Cepstral Coefficients for Music Modeling', *Ismir*, 270: 1–11.
- Maciej, P. et al. (2013) 'Vocal Communication in a Complex Multi-Level Society: Constrained Acoustic Structure and Flexible Call Usage in Guinea Baboons', *Frontiers in Zoology*, 10/1: 58.
- Mandel, M. I., and Ellis, D. P. (2005). Song-level features and support vector machines for music classification. *Proceedings of the 6th International Conference on Music Information Retrieval (ISMIR)*, pp. 594–599.
- Marler, P. (1980). 'Primate Vocalization: Affective or Symbolic?' in Bourne G. H. (ed) *Speaking of Apes*, pp. 221–229. Boston: Springer.
- McCune, L. et al. (1996) 'Grunt Communication in Human Infants (*Homo sapiens*)', *Journal of Comparative Psychology*, 110/1: 27
- Meise, K. et al. (2011) 'Sources of Acoustic Variation: Implications for Production Specificity and Call Categorization in Chacma Baboon (*Papio Ursinus*) Grunts', *The Journal of the Acoustical Society of America*, 129/3: 1631–41
- Mielke, A., and Zuberbühler, K. (2013) 'A Method for Automated Individual, Species and Call Type Recognition in Free-Ranging Animals', *Animal Behaviour*, 86/2: 475–82
- Morris, E. K. et al. (2014) 'Choosing and Using Diversity Indices: Insights for Ecological Applications from the German Biodiversity Exploratories', *Ecology and Evolution*, 4/18: 3514–24

- Nathani, S., Ertmer, D. J., and Stark, R. E. (2006) 'Assessing Vocal Development in Infants and Toddlers', *Clinical Linguistics & Phonetics*, 20/5: 351–69
- Oller, D. K. (2000). *The Emergence of the Speech Capacity*. Mahwah, NJ: Psychology Press.
- et al. (2013) 'Functional Flexibility of Infant Vocalization and the Emergence of Language', *Proceedings of the National Academy of Sciences of the United States of America*, 110/16: 6318–23
- , and Griebel, U. (2004). 'Contextual Freedom in Human Infant Vocalization and the Evolution of Language' in *Evolution of Communicative Flexibility: Complexity, Creativity and Adaptability in Human and Animal Communication*, p. 135. Cambridge, MA: MIT Press.
- et al. (2019) 'Language Origins Viewed in Spontaneous and Interactive Vocal Rates of Human and Bonobo Infants', *Frontiers in Psychology*, 10: 729.
- , Griebel, U., and Warlaumont, A. S. (2016) 'Vocal Development as a Guide to Modeling the Evolution of Language', *Topics in Cognitive Science*, 8/2: 382–92
- et al. (1976) 'Infant Babbling and Speech', *Journal of Child Language*, 3/1: 1–11
- Owren, M. J., and Casale, T. M. (1994) 'Variations in Fundamental Frequency Peak Position in Japanese Macaque (*Macaca fuscata*) Coo Calls', *Journal of Comparative Psychology*, 108/3: 291
- Plooij, F. X. (1984). *The Behavioral Development of Free-Living Chimpanzee Babies and Infants*. Norwood, NJ: ALEX Publishing Corporation.
- et al. (2014) 'Longitudinal Recordings of the Vocalizations of Immature Gombe Chimpanzees for Developmental Studies', *Scientific Data*, 1/1: 1–10
- Pollick, A. S., and Waal, F. B. M. d. (2007) 'Ape Gestures and Language Evolution', *Proceedings of the National Academy of Sciences*, 104/19: 8184–9
- Ponsot, E. et al. (2018) 'Cracking the Social Code of Speech Prosody Using Reverse Correlation', *Proceedings of the National Academy of Sciences*, 115/15: 3972–7
- R Core Team. (2018). *R: A Language and Environment for Statistical Computing*. Vienna: R Foundation for Statistical Computing.
- Range, F., and Fischer, J. (2004) 'Vocal Repertoire of Sooty Mangabeys (*Cercocebus Torquatus Atys*) in the Tai National Park', *Ethology*, 110/4: 301–21
- Rendall, D. (2003) 'Acoustic Correlates of Caller Identity and Affect Intensity in the Vowel-like Grunt Vocalizations of Baboons', *The Journal of the Acoustical Society of America*, 113/6: 3390–402
- , and Owren, M. J. (2002). *Animal Vocal Communication: Say What? In the Cognitive Animal: Empirical and Theoretical Perspectives on Animal Cognition*, pp. 307–313. Cambridge, MA: MIT Press.
- Reynolds, V. (2005). *The Chimpanzees of the Budongo Forest: Ecology, Behaviour, and Conservation*. Oxford: Oxford University Press.
- Team R. (2015). *RStudio: Integrated Development for R*. Boston, MA: RStudio, Inc., p. 42.
- Salmi, R., Hammerschmidt, K., and Doran-Sheehy, D. M. (2013) 'Western Gorilla Vocal Repertoire and Contextual Use of Vocalizations', *Ethology*, 119/10: 831–47
- Schaik, C. P. V., and Burkart, J. M. (2010). 'Mind the Gap: Cooperative Breeding and the Evolution of Our Unique Features' in Kappeler P.M. and Silk J. B. (eds) *Mind the Gap*, pp. 477–496. Springer: Berlin, Heidelberg.
- Schel, A. M. et al. (2013) 'Chimpanzee Food Calls Are Directed at Specific Individuals', *Animal Behaviour*, 86/5: 955–65.
- Silk, J. B. et al. (2018) 'To Grunt or Not to Grunt: Factors Governing Call Production in Female Olive Baboons, *Papio Anubis*', *PLoS One*, 13/11: e0204601.
- Slocombe, K. E., and Zuberbühler, K. (2010). Vocal Communication in Chimpanzees. *The Mind of the Chimpanzee: Ecological and Experimental Perspectives*, pp. 192–207. Chicago: University of Chicago Press.
- et al. (2010) 'Production of Food-Associated Calls in Wild Male Chimpanzees is Dependent on the Composition of the Audience', *Behavioral Ecology and Sociobiology*, 64/12: 1959–66.
- , and Newton-Fisher, N. E. (2005) 'Fruit Sharing between Wild Adult Chimpanzees (*Pan troglodytes Schweinfurthii*): A Socially Significant Event? ', *American Journal of Primatology*, 65/4: 385–91
- , and Zuberbühler, K. (2005) 'Functionally Referential Communication in a Chimpanzee', *Current Biology*, 15/19: 1779–84
- Tajiri, Y. et al. (2010). 'Feature Extraction Using Support Vector Machines' in Wong K. W., Mendis B. S. U., and Bouzerdoum A. (eds) *Neural Information Processing. Models and Applications*, pp. 108–115. Springer.
- Takahashi, D. Y. et al. (2015) 'The Developmental Dynamics of Marmoset Monkey Vocal Production', *Science*, 349/6249: 734–8
- Tchernichovski, O., and Oller, D. K. (2016) 'Vocal Development: How Marmoset Infants Express Their Feelings', *Current Biology*, 26/10: R422–4.
- Tsukahara, T. (1993) 'Lions Eat Chimpanzees: The First Evidence of Predation by Lions on Wild Chimpanzees', *American Journal of Primatology*, 29/1: 1–11
- Vert, J.-P., Tsuda, K., and Schölkopf, B. (2004). 'A Primer on Kernel Methods', *Kernel Methods in Computational Biology*, 47: 35–70.
- Waal, F. B. M. d., and Pollick, A. S. (2011) 'Gesture as the Most Flexible Modality of Primate Communication' in Gibson K. R. and Tallerman M. (eds) *The Oxford Handbook of Language Evolution*, <https://doi.org/10.1093/oxfordhb/9780199541119.013.0006>.
- Watson, S. K. et al. (2015) 'Vocal Learning in the Functionally Referential Food Grunts of Chimpanzees', *Current Biology*, 25/4: 495–9
- Williams, C. E., and Stevens, K. N. (1972) 'Emotions and Speech: Some Acoustical Correlates', *The Journal of the Acoustical Society of America*, 52/4B: 1238–50

- Yoo, H., Bowman, D. A., and Oller, D. K. (2018) 'The Origin of Protoconversation: An Examination of Caregiver Responses to Cry and Speech-Like Vocalizations', *Frontiers in Psychology*, 9, Article 1510.
- Zhang, Y. S., and Ghazanfar, A. A. (2016) 'Perinatally Influenced Autonomic System Fluctuations Drive Infant Vocal Sequences', *Current Biology*, 26/10: 1249–60
- Zuberbühler, K. (2012). 'Cooperative Breeding and the Evolution of Vocal Flexibility' in Gibson K. R. and Tallerman M. (eds) *The Oxford Handbook of Language Evolution*, pp. 71–81. Oxford: Oxford University Press.
- (2014) 'Experimental Field Studies with Non-Human Primates', *Current Opinion in Neurobiology*, 28: 150–6