1    Vocal complexity in the long calls of Bornean orangutans

2

3    *Erb, W.M.[1,2], Ross, W.[1], Kazanecki, H.[1], Mitra Setia, T.[3,4], Madhusudhana, S.[1,5], Clink, D.J.[1]

4    *Corresponding author

5    Email: erbivorous@gmail.com

6

7    [1] K. Lisa Yang Center for Conservation Bioacoustics, Cornell Lab of Ornithology, Cornell

8    University, Ithaca, New York, 14850, USA

9    [2] Department of Anthropology, Rutgers University, New Brunswick, New Jersey, 08901, USA

10    [3] Fakultas Biologi, Universitas Nasional Jakarta, Jakarta, Indonesia

11    [4] Primate Research Center, Universitas Nasional, Jakarta, Indonesia

12    [5] Centre for Marine Science and Technology, Curtin University, Perth, WA 6102, Australia

13

14    Short title: Orangutan Long Call Classification

15

16    **ABSTRACT**

17    Vocal complexity is central to many evolutionary hypotheses about animal communication. Yet,

18    quantifying and comparing complexity remains a challenge, particularly when vocal types are highly

19    graded. Male Bornean orangutans (*Pongo pygmaeus wurmbii*) produce complex and variable "long call"

20    vocalizations comprising multiple sound types that vary within and among individuals. Previous

21    studies described six distinct call (or pulse) types within these complex vocalizations, but none

22    quantified their discreteness or the ability of human observers to reliably classify them. We studied

23    the long calls of 13 individuals to: 1) evaluate and quantify the reliability of audio-visual classification

24    by three well-trained observers, 2) distinguish among call types using supervised classification and

25  unsupervised clustering, and 3) compare the performance of different feature sets. Using 46 acoustic

26  features, we used machine learning (i.e., support vector machines, affinity propagation, and fuzzy c-

27  means) to identify call types and assess their discreteness. We also used Uniform Manifold

28  Approximation and Projection (UMAP) to visualize the separation of pulses using both extracted

29  features and spectrograms. We found low inter-observer reliability and poor classification accuracy

30  using supervised approaches, indicating that pulse types were not discrete. We propose a new pulse

31  type classification scheme that is highly reproducible across observers and exhibits high classification

32  accuracy using support vector machines. Although the low number of call types suggests long calls

33  are fairly simple, the continuous gradation of sounds seems to greatly boost the complexity of this

34  system. This work responds to calls for more quantitative research to define call types and measure

35  the gradedness of animal vocal systems and highlights the need for a more comprehensive

36  framework for studying vocal complexity vis-à-vis graded repertoires.

37

38  **HIGHLIGHTS**

39  ● We used audio-visual (AV) analysis and machine-learning to discriminate pulse types.

40  ● AV and support vector machines (SVM) did not support the six published pulse types.

41  ● Hard and soft clustering algorithms showed a mixture of discrete and graded pulses.

42  ● We propose three pulse types that show high reproducibility and classification accuracy.

43  ● More work is needed to investigate the role of graded signals in vocal complexity.

44

45  **KEYWORDS**

46  acoustic communication; affinity propagation; fuzzy clustering; graded signals; machine learning;

47  supervised classification; support vector machines; Uniform Manifold Approximation and

48  Projection (UMAP); unsupervised clustering; vocal repertoire

## INTRODUCTION

Vocal complexity, or the diversity of sounds in a species' repertoire, is central to many evolutionary hypotheses about animal communication (Bradbury & Vehrencamp, 2011; Fischer et al., 2016; Freeberg et al., 2012; McComb & Semple, 2005). This complexity has been hypothesized to be shaped by a range of factors including predation pressure, sexual selection, habitat structure, and social complexity (Bradbury & Vehrencamp, 2011; Fischer et al., 2016). Two common measures of vocal complexity are: 1) the diversity (or number) of call types as well as 2) their discreteness. For instance, within black-capped chickadee (*Poecile atricapillus*) groups, individuals flexibly increase the diversity of note types when they are in larger groups, presumably increasing the number of potential messages that can be conveyed (Freeberg et al., 2012). When comparing across species, similar themes emerge in rodents and primates. Sciurid species with a greater diversity of social roles have more alarm call types (Blumstein & Armitage, 1997) and primate species in larger groups with more intense social bonding have larger vocal repertoires (McComb & Semple, 2005). Further, it has been proposed that while discrete repertoires facilitate signal recognition in dense habitats, graded calls allow more complexity in open habitats where intermediate sounds communicate arousal and can be linked with visual signals (Marler et al., 1975).

Yet, quantifying vocal complexity in a standardized manner remains a challenge for comparative analyses. A primary aspect of this challenge is related to the identification and quantification of discrete call types, which is particularly vexing in repertoires comprising intermediate calls and in species that exhibit significant inter-individual variation (Fischer et al., 2016). The most common approaches to identifying call types are: 1) manual (visual or audio-visual) classification of spectrograms by a human observer and 2) automated (quantitative or algorithmic) using features that are either manually or automatically measured from spectrograms (Kershenbaum et al., 2016). Audio-visual classification involves one or more observers inspecting spectrograms

73    visually while simultaneously listening to the sounds. This method has been applied to the

74    vocalizations of numerous taxa (e.g., manatees, *Trichechus manatus latirostris*: Brady et al., 2020; spear-

75    nosed bats, *Phyllostomus discolor*: Lattenkamp, 2019; humpback whales, *Megaptera novaeangliae*:

76    Madhusudhana et al., 2019; New Zealand kea parrots, *Nestor notabilis*: Schwing et al., 2012). Audio-

77    visual classification studies often rely on a single expert observer and only rarely quantify within- or

78    between-observer reliability (reviewed in Jones et al., 2001). On one hand, when classification is

79    done by a single observer, the study risks idiosyncratic or irreproducible results. On the other hand,

80    when multiple observers are involved, the study risks inconsistent assessments among scorers. To

81    assess the reproducibility of a human-based classification scheme, it is critical to evaluate the

82    consistency of scores within and/or among the human observers using inter-rater reliability (IRR)

83    statistics such as Cohen's kappa (Hallgren, 2012).

84         To compare and classify acoustic signals, researchers must make decisions about which

85    features to estimate, as analyses of the waveform are generally too computationally costly. A

86    commonly used approach for many classification problems is feature selection, in which a suite of

87    selected time- and frequency-based characteristics of sounds are measured and compiled from

88    manually annotated spectrograms (Odom et al., 2021). There is little standardization concerning the

89    selection of acoustic variables across studies, which often include a combination of qualitative and

90    quantitative measurements that are manually and/or automatically (i.e., using a sound analysis

91    program, such as Raven Pro 1.6) extracted. As an alternative to feature selection, some researchers

92    use automated approaches wherein the spectral content of sounds is measured using spectrograms,

93    cepstra, multi-taper spectra, wavelets, or formants (reviewed in Kershenbaum et al., 2016).

94         Once features have been manually or automatically extracted, multivariate analyses can be

95    used to classify or cluster sounds using supervised or unsupervised algorithms, respectively. In the

96    case of supervised classification, users manually label a subset of representative sounds which are

97      used to train the statistical model that will subsequently be used to automatically identify those

98      sound types in an unlabeled set of data (Cunningham et al., 2008). In contrast to supervised

99      classification, clustering is an unsupervised machine learning approach in which an algorithm divides

100     a dataset into several groups or clusters such that observations in the same group are similar to each

101     other and dissimilar to the observations in different groups (Greene et al., 2008). Thus, in the case of

102     unsupervised clustering, the computer – rather than the human observer – learns the groupings and

103     assigns labels to each value (Alloghani et al., 2020).

104             Enumeration of call types in a repertoire is especially challenging when there are

105     intermediate forms that fall between categories. These so-called graded call types have been well

106     documented across primate taxa (Fischer et al., 2016; Hammerschmidt & Fischer, 1998). An

107     alternative to "hard clustering" of calls into discrete categories (e.g., k-means, k-medoids, affinity

108     propagation), "soft clustering" (e.g., fuzzy c-means) allows for imperfect membership by assigning

109     probability scores for membership in each cluster, thereby making it possible to identify call types

110     with intermediate values (Cusano et al., 2021; Fischer et al., 2016). So-called fuzzy clustering can be

111     used in tandem with hard clustering by also quantifying the degree of ambiguity (or gradedness)

112     exhibited by particular sounds and continuities across call types. Thus, soft clustering provides a

113     means of quantifying gradedness in repertoires and can enable the identification of intermediate

114     members.

115             Across studies of animal vocal complexity, there is notable variation in the number and type

116     of feature sets used, ranging from fewer than 10 to more than 100 parameters that are manually

117     and/or automatically extracted. Table 1 provides a summary of 15 studies across mammalian and

118     avian taxa that used supervised classification and unsupervised clustering approaches to identify call

119     types across a range of mammalian and avian taxa. Though most studies paired audio-visual

120     classification with an unsupervised clustering method, a few also included discriminant function

121    analysis (DFA) to quantify the differences among the human-labeled call types and/or computer-

122    identified clusters. Authors relied on a broad range of unsupervised clustering algorithms, though

123    hierarchical agglomerative clustering was the most used method. Studies that aimed to provide an

124    accurate classification of different call types often relied on a combination of supervised

125    classification and unsupervised clustering methods to ensure results were robust and repeatable.

126    However, those that compared feature sets or clustering methods often reported a lack of agreement

127    on the number of clusters identified, highlighting the difficulty of the seemingly straightforward task

128    of identifying and quantifying call types.

129　**Table 1.** Review of studies using supervised classification and unsupervised clustering approaches to identify vocal types.

| Authors (Date) | Taxon | Goals | N Features | Classification (N observers) | Clustering Method |
|---|---|---|---|---|---|
| Wadewitz et al. (2015) | Chacma baboon *(Papio ursinus)* | Compare hard & soft clustering, evaluate influence of features | 9, 38, 118 (+ 19 PCA factors) | A/V * | K-means, Hierarchical agglomerative (Ward's), Fuzzy c-means |
| Fuller (2014) | Blue monkey *(Cercopithecus mitis stulmanni)* | Catalog vocal signals | 18 PCA factors | A/V (1), DFA ** | Hierarchical agglomerative |
| Fournet et al., (2015) | Humpback whale *(Megaptera novaeangliae)* | Catalog non-song vocalizations | 15 | A/V (1), DFA | Hierarchical agglomerative |
| Brady et al. (2020) | Florida manatee *(Trichechus manatus latirostris)* | Catalog vocal repertoire | 17 | A/V (1) | Maximum likelihood, CART |
| Hammerschmidt & Fischer (2019**)** | Chacma *(Papio ursinus)*, olive *(P. anubis)*, and Guinea baboon *(P. papio)* | Catalog & compare vocal repertoires, Compare A/V to clustering | 9 | A/V (multiple), DFA ** | Two-step cluster analysis |
| Sadhukhan et al. (2019) | Indian wolf *(Canis lupus pallipes)* | Catalog harmonic vocalizations | 8 | DFA | Hierarchical agglomerative |
| Hedwig et al. (2019) | African forest elephant *(Loxodonta cyclotis)* | Catalog vocal repertoire | 23 | DFA ** | PCA |
| Huijser et al. (2020) | Sperm whale *(Physeter macrocephalus)* | Catalog coda repertoires | 2 | A/V (1) | K-means, Hierarchical agglomerative |
| Vester et al. (2017) | Long-finned pilot whale *(Globicephala melas)* | Catalog vocal repertoire | 14 | A/V (2), DFA | Two-step cluster analysis |

| Soltis et al. (2012) | Key Largo woodrat *(Neotoma floridana smalli)* | Catalog vocal repertoire | 6 | A/V* | Multidimensional scaling analysis (MDS) |
|---|---|---|---|---|---|
| Elie & Theunissen (2016) | Zebra finch *(Taeniopygia guttata)* | Catalog vocal repertoire, determine distinguishing features | 22, 25 (MFCCs) | A/V (1), Fisher LDA, Random Forest | PCA, Gaussian mixture |
| Janik (1999) | Bottlenose dolphin *(Tursiops truncatus)* | Compare A/V to clustering | 20 | A/V (5) | K-means, Hierarchical agglomerative |
| Cusano et al. (2021) | Humpback whale *(Megaptera novaeangliae)* | Differentiate discrete vs. graded call types | 25 | A/V* | Fuzzy k-means |
| Garland et al. (2015) | Beluga whale *(Delphinapterus leucas)* | Catalog vocal repertoire | 12 | A/V* | CART, Random Forest |
| Thiebault et al. (2019) | Cape gannet *(Morus capensis)* | Catalog repertoire of foraging calls | 12 | A/V* | Random Forest |

130   * Study did not report # of observers

131   ** leave-one-out

132   In the present study, we examine vocal complexity in the long calls of Bornean orangutans

133 (*Pongo pygmaeus wurmbii*) by evaluating how the choice of classification or clustering methods and

134 feature inputs affect the number of call types we recognize. Orangutans are semi-solitary great apes

135 who exhibit a promiscuous mating system in which solitary adult males range widely in search of

136 fertile females (Spillmann et al., 2017). Flanged males (i.e., adult males who have fully developed

137 cheek pads, throat sacs, and body size approximately twice that of adult females) emit loud

138 vocalizations, or long calls, which travel up to a kilometer and serve to attract female mates and

139 repel rival males (Mitra Setia & van Schaik, 2007) In this social setting, long calls thus hold an

140 important function for coordination among widely dispersed individuals.

141   Long calls are complex and variable vocalizations comprising multiple call (or pulse) types

142 that vary within and among individuals (Askew & Morrogh-Bernard, 2016; Spillmann et al., 2010).

143 Long calls typically begin with a bubbly introduction of soft, short sounds that build into a climax of

144 high-amplitude frequency-modulated pulses followed by a series of lower-amplitude and -frequency

145 pulses that gradually transition to soft and short sounds, similar to the introduction (cf. MacKinnon,

146 1977, Table S1). Although Davilla Ross and Geissmann (2007) first attempted to classify and name

147 the different elements of these calls, they noted a "wide variety of call elements do not belong to any

148 of these note types" (Davila Ross & Geissmann, 2007 p. 309).

149   Spillmann and colleagues (2010) presented the most detailed description of orangutan long

150 calls in which they identified six different pulse types (Table 2), but thus far there has been no

151 attempt to systematically classify pulses or quantify how discrete these sounds are. Further, no

152 studies have described the process for or the number of observers classifying sound types nor the

153 reliability of classifications within or among observers. Thus, it is presently unclear how well pulse

154 types can be discriminated by human observers or quantitative classification tools, thereby limiting

155 our ability to repeat, reproduce, and replicate these studies.

156 **Table 2.** Names and descriptions of sound labels used in previous studies, using Spillmann et al.

157 (2010) labels as reference.

| Sound Type | MacKinnon 1974 | Davila Ross & Geissmann 2007 | Spillmann et al. 2010 |
|---|---|---|---|
| Grumbles | bubbly introduction | bubbling | "preceding bubbling-like elements that are low in loudness" |
| Bubbles | n/a | bubbling | "low amplitude, looks like a cracked sigh" |
| Roar | "climax of full roars" | roar | "more rounded and lower in frequency" |
| Low Roar | n/a | n/a | "half the fundamental frequency at the highest point than roar" |
| Volcano Roar | n/a | n/a | "sharp tip and higher frequency than roar" |
| Huitus | n/a | huitus | "high amplitude with steeply ascending and descending part that are not connected" |
| Intermediary | n/a | intermediary | "low amplitude, frequency modulation starts with a rising part followed by a falling part that changes again into a rising and ends with a falling part" |
| Sigh | "tails off gradually into a series of sighs" | sigh | "low amplitude, starts with a short rising part and changes in a long falling part" |

158

159 The present study aims to evaluate vocal complexity in orangutan long calls to compare

160 different approaches to identifying the number of discrete calls and estimating the degree of

161 gradedness in a model vocal system. Specifically, the objectives of our study are to: 1) evaluate and

162 quantify the reliability of manual audio-visual (AV) classification by three well-trained observers, 2)

163 classify and cluster call types using supervised classification (support vector machines) and

164 unsupervised hard (affinity propagation) and soft (fuzzy c-means) clustering methods, and 3)

165    compare the results using different feature sets (i.e., feature engineering, complete spectrographic

166    representations). Based on these findings, we will make explore and assess alternative classification

167    systems for identifying discrete and graded call types in this system.

168

169    **METHODS**

170    **Ethical Note**

171        This research was approved by the Institutional Animal Care and Use Committee of Rutgers,

172    the State University of New Jersey (protocol number 11-030 granted to Erin Vogel). Permission to

173    conduct the research was granted to WME by the Ministry of Research and Technology of the

174    Republic of Indonesia (RISTEK Permit #137/SIP/FRP/SM/V/2013-2015). The data included in

175    the present study comprise recordings collected during passive observations of wild habituated

176    orangutans at distances typically exceeding 10 m. The population has been studied since 2003 and

177    individual orangutans were not disturbed by observers in the execution of this study.

178

179    **Study Site and Subjects**

180        We conducted our research at the Tuanan Orangutan Research Station in Central

181    Kalimantan, Indonesia ($2^0$ 09' 06.1" S; $114^0$ 26' 26.3" E). Tuanan comprises approximately 900

182    hectares of secondary peat swamp forest that was selectively logged prior to the establishment of the

183    study site in 2003 (see Erb et al., 2018 for details). For the present study, data were collected

184    between June 2013 and May 2016 by WME and research assistants (see Acknowledgments) during

185    focal observations of adult flanged males. Whenever flanged males were encountered, our field team

186    followed them until they constructed a night nest and we returned to the nest before dawn the next

187    morning to continue following the same individual. All subjects were individually recognized on the

188    basis of unique facial features, scars, and broken or missing digits. Individuals were followed

189    continuously for five days, unless they were lost or left the study area. During 316 partial- and full-

190    day focal observations, we recorded 1,013 long calls from 23 known individuals.

191

192    **Long Call Recording**

193    During observations, we used all-occurrences sampling (Altmann, 1974) of long calls noting:

194    time, GPS location, stimulus (preceded within 15 minutes by another long call, tree fall, approaching

195    animal, or other loud sounds), and any accompanying movements or displays. Recordings of long

196    calls were made opportunistically, using a Marantz PMD-660 solid-state recorder (44,100 Hz

197    sampling frequency, 16 bits: Kanagawa, Japan) and a Sennheiser directional microphone (K6 power

198    module and ME66 recording head: Wedemark, Germany). Observers made voice notes at the end of

199    each recording noting the date and time, orangutan's name, height(s), distance(s), and movement(s),

200    as well as the gain and microphone directionality (i.e., directly or obliquely oriented).

201

202    **Long Call Analysis**

203    For the present study, we selected a subset of recordings from 13 males from whom we had

204    collected at least 10 high-quality long call recordings. When more than 10 long call recordings were

205    available for a given individual, we randomly selected 10 of his recordings, stratified by study year, to

206    balance our dataset across individuals and years. The final dataset comprised 130 long calls, 10 from

207    each of 13 males.

208    Prior to annotating calls, we used Adobe Audition 14.4 to downsample recordings to 5,100

209    Hz (cf. Hammerschmidt & Fischer, 2019). We then generated spectrograms in Raven Pro 1.6 (K.

210    Lisa Yang Center for Conservation Bioacoustics, 2019) with a 512-point (92.9 ms) Hann window (3

211    dB bandwidth = 15.5 Hz), with 90% overlap and a 512-point DFT, yielding time and frequency

212    measurement precision of 9.25 ms and 10.8 Hz. Three observers (WME, WR, HK) annotated calls
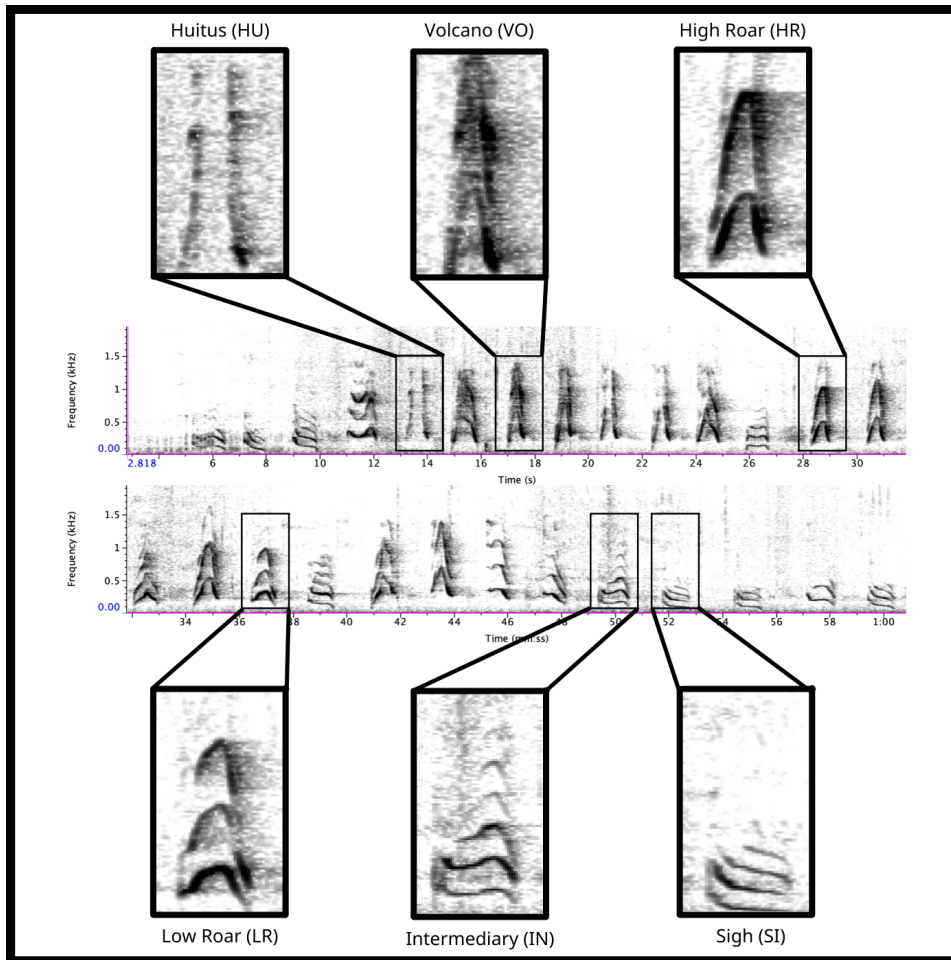
213    by drawing selections that tightly bounded the start and end of each pulse (Fig. S1) and assigned call

214    type labels using the classification scheme outlined in Table 2. Except for huitus pulses (for which

215    the rising and falling sounds are broken by silence), we operationally defined a pulse as the longest

216    continuous sound produced on a single exhalation. Because most long calls are preceded and/or

217    followed by a series of short bubbling sounds, we used a threshold duration of $\geq 0.2$ seconds to

218    differentiate pulses from these other sounds. Most selections were drawn with a fixed frequency

219    range from 50 Hz to 1 kHz; however, in cases where the maximum fundamental frequency exceeded

220    1 kHz (e.g., huitus and volcano roars), selections were drawn from 50 Hz to 1.5 kHz. Occasionally,

221    we manually reduced the frequency range of selections if there were disturbing background sounds,

222    but only if this did not affect measures of the fundamental frequency contour or high-energy

223    harmonics. We noted whether selections were tonal (i.e., the fundamental frequency contour was

224    fully or partially visible) and whether they contained disturbing background noises such as birds,

225    insects, or breaking branches.

226         Our selected feature set comprised 25 extracted measurements made in Raven (Table S1) as

227    well as an additional 19 measurements estimated using the R package *warbleR* (Araya-Salas & Smith-

228    Vidaurre, 2017). Prior to analyzing sounds in warbleR, we filtered out all pulse selections that were

229    atonal or contained disturbing background noise, resulting in 2,270 clips. Two additional

230    measurements (minimum and maximum) of the fundamental frequency (F0) were made using the

231    "freq_ts" function in *warbleR* with the following settings: wavelength = 512, Hanning window, 70%

232    overlap, 50 - 1,500 Hz, threshold = 85%. We then saved printed spectrograms depicting the F0

233    contours for each. One observer (WME) visually screened the minimum and maximum values of

234    the F0 contours and scored them as accurate or inaccurate. After removing those pulses for which

235    one or both F0 measures were inaccurate, the final full dataset comprised 1,033 pulses from 117

236    long calls for which all 46 parameters were measured.

237

238    **Audio-Visual Analysis**

239        To assess the inter-rater reliability (IRR) of the audio-visual analysis, we randomly selected

240    300 pulses (saved as individual .wav files). We included this step to remove any bias that may be

241    introduced by information about the position or sequence of a pulse-type within a long call. Using

242    the spectrograms and descriptions of pulse types published by Spillmann and colleagues (Spillmann

243    et al., 2010), three observers (WME, WR, HK) labeled each sound as one of six pulse types (Fig. 1).

244    Prior to completing this exercise, all observers had at least six months' experience classifying pulse

245    types, which involved routine feedback and three-way discussion. We used the R package *irr* (Gamer

246    et al., 2012) to calculate Cohen's kappa (a common statistic for assessing IRR for categorical

247    variables) for each pair of observers, and averaged these values to provide an overall estimate of IRR

248    (Light's kappa) across all pulse types (cf. Hallgren, 2012; Light, 1971).

**Figure 1. Spectrogram depicting long call pulse types.** Pulses include HU = huitus, VO = volcano, HR = (high) roar, LR = low roar, IN = intermediary, SI = sigh. Spectrograms produced in Raven Pro 1.6.

**Supervised Classification**

For the supervised classification analysis, one observer (WME) manually classified all pulses (N=1,033). We then used support vector machines (SVM) in the R package *e1071* (Meyer et al., 2021) to evaluate how well pulse types could be discriminated using a supervised machine learning approach. SVMs are commonly used for supervised classification and have been successfully applied to the classification of primate calls (Clink & Klinck, 2020; Fedurek et al., 2016; Turesson et al.,

260  2016). We used the sigmoidal kernel as previous research using SVM has found the most robust

261  results using this kernel type (Clink & Klinck, 2020) and we estimated the best values for the gamma

262  and cost parameters using the "tune" function. Following this, we calculated our classification

263  accuracy using 10 iterations of leave-one-out cross-validation. Lastly, we used SVM recursive feature

264  elimination to rank variables in order of their importance for classifying call types (cf. Clink et al.,

265  2018). For each of the top five most influential variables identified by recursive feature elimination,

266  we used Kruskal-Wallis nonparametric tests due to the non-normal distribution of the residuals

267  when applying linear models. We followed these with Dunn's test of multiple comparisons to

268  examine differences among pulse types and unsupervised clusters (described below) – applying the

269  Benjamini-Hochberg adjustment to control the false discovery rate – using the R package *FSA*

270  (Ogle et al., 2022).

271

272  **Unsupervised Clustering**

273  For the unsupervised analysis, we used both hard- and soft-clustering approaches. For hard

274  clustering, we used affinity propagation, which has the advantage that it does not require the user to

275  identify the number of clusters a priori; further, because all data points are considered

276  simultaneously, the results are not influenced by the selection of an initial set of points (Frey &

277  Dueck, 2007). Using the R package *apcluster* (Bodenhofer et al., 2011), we systematically varied the

278  value of 'q' in 0.25 increments from 0 to 1. By comparing the mean silhouette coefficient for each of

279  the cluster solutions (Wang et al., 2008), we found that q = 0 produced the optimal number of

280  clusters and thus we report the results from this model. We used silhouette coefficients to quantify

281  the stability of the resulting clusters (cf. Clink & Klinck, 2020).

282  For the soft clustering analysis, we used C-means fuzzy clustering. In this analysis, each pulse

283  is assigned a membership value (m ranges from 0 = none to 1 = full accordance) for each of the

284    clusters. We first determined the optimal number of clusters (c) by evaluating measures of internal

285    validation and stability generated in the R package *clValid* (Brock et al., 2008) when c varied from 2

286    (the minimum) to 7 (one more than the previously described number of pulse types). We then

287    systematically varied the fuzziness parameter μ from 1.1 to 5 (i.e., 1.1, 1.5, 2, 2.5, etc.: cf. Zhou et al.,

288    2014)) using the R package 'cluster' (Maechler et al., 2021). When μ = 1, clusters are tight and

289    membership values are binary; however, as μ increases, cases can show partial membership to

290    multiple clusters, and the clusters themselves thereby become fuzzier and can eventually merge,

291    leading to fewer clusters (Fischer et al., 2016). We used measures of internal validity (connectivity,

292    silhouette width, and Dunn index) and stability (average proportion of non-overlap = APN, average

293    distance = AD, average distance between means = ADM, and figure of merit = FOM) to evaluate

294    the cluster solutions in the R package *clValid* (Brock et al., 2008). Once we had identified the best

295    solution, we calculated typicality coefficients to assess the discreteness of each pulse, wherein higher

296    values indicate pulses that are well separated from other clusters and lower values indicate pulses

297    that are intermediate between classes (cf. Cusano et al., 2021; Wadewitz et al., 2015).

298            Non-linear dimensionality reduction techniques have recently emerged as fruitful

299    alternatives to traditional linear techniques (e.g., principal component analysis) for classifying animal

300    sounds (Sainburg et al., 2020). Uniform Manifold Approximation and Projection (UMAP) is a state-

301    of-the-art unsupervised machine learning algorithm (McInnes et al., 2018) that has been applied to

302    visualizing and quantifying structures in animal vocal repertoires (Sainburg et al., 2020). Like

303    ISOMAP and t-SNE, UMAP constructs a topology of the data and projects that graph into a lower-

304    dimensional embedding (McInnes et al., 2018; Sainburg et al., 2020) UMAP has been shown to

305    preserve more global structure while achieving faster computation times (McInnes et al., 2018) and

306    has been effectively applied to meaningful representations of acoustic diversity (reviewed in

307 Sainburg et al., 2020). This approach removes any a priori assumptions about which acoustical

308 features are most salient or easily measured by humans.

309   We applied UMAP separately to the 46-feature set and to time-frequency representations of

310 extracted pulses. In the latter case, we used as inputs power density spectrograms of 0.9-s duration

311 audio clips centered at the temporal midpoint of annotated pulses. The chosen duration was fixed

312 irrespective of the selection duration. This means that, for short selections, the spectrograms also

313 included sounds outside of the original selection. Short-time Fourier transforms of the clips were

314 computed, using SciPy's (https://scipy.org/) *spectrogram* function, with a Hann window of 50 ms and

315 50% frame overlap (20 Hz frequency resolution, 25 ms time resolution). Spectral levels were

316 converted to the decibel scale by applying $10 \times \log_{10}$. The bandwidth of the resulting spectrograms

317 was limited to 50-1000 Hz prior to UMAP computation to suppress the influence of low-frequency

318 noise on clustering. We used the *UMAP* function from the Python package *umap-learn* (McInnes et

319 al., 2018) to compute the low-dimensional embeddings. Finally, we calculated Hopkin's statistic of

320 clusterability on the resultant UMAP using the R package *factoextra* (Kassambara & Mundt, 2020).
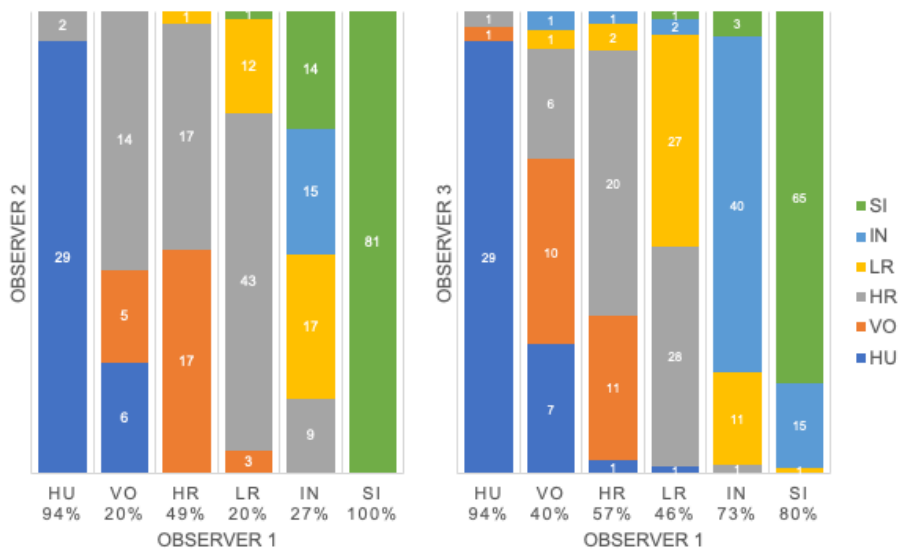
321   Finally, we reviewed the outputs of our unsupervised clustering approaches to assess the

322 putative number of pulses and graded variants. To identify a simple, data-driven, repeatable method

323 for manually classifying pulse types, we began by pooling the typical pulses that belonged to each of

324 the clusters identified by fuzzy clustering. Because F0 is a highly salient feature in long call

325 spectrograms, our approach focused on the shape and height (or maximum frequency) of this

326 feature. Using our revised definitions, we repeated the 1) audio-visual analysis and calculated IRR

327 using manual labels from the same 300 pulses reviewed by the same three observers as before, and

328 2) SVM classification of 500 randomly selected pulses scored by a single observer (WME) following
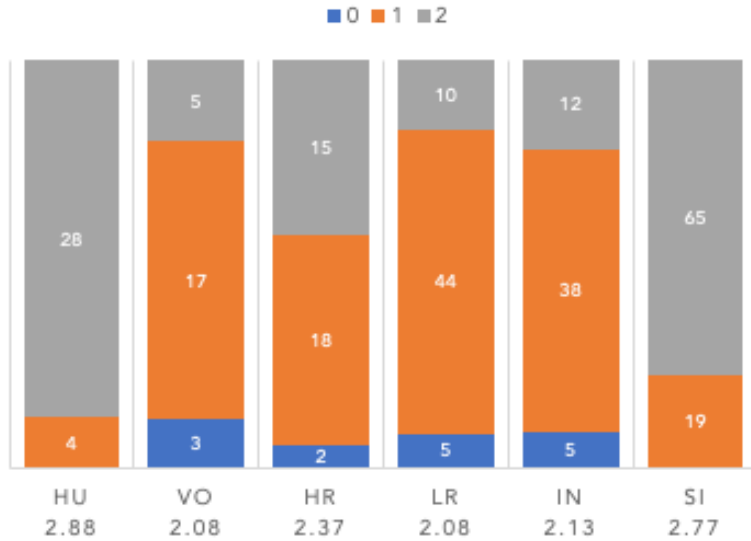
329 the methods described above.

330

331     **RESULTS**

332     **Audio-visual analysis**

333            Based on manual labels from three observers using audio-visual classification methods, we

334     calculated Light's kappa $\varkappa = 0.599$ (i.e., the arithmetic mean of Cohen's Kappa for observers 1-2 =

335     0.48, 1-3 = 0.60, and 2-3 = 0.60), which indicated only moderate agreement among observers

336     (Landis & Koch, 1977). Classification agreement varied widely by pulse type (Fig. 2, Table 3).

337     Whereas huitus and sigh pulse types showed high agreement among observers (mean 2.88 and 2.77,

338     respectively, where 3 indicates full agreement), low roar and volcano pulse types showed very low

339     agreement (mean 2.08).



340

341

342 **Figure 2. Audio-visual classification agreement across observers.** Stacked barplots indicating

343 (top) classification agreement by pulse type between observer 1-2 and observer 1-3 and (bottom) the

344 number of observers who agreed on the pulse types assigned by observer 1; the average agreement

345 index is indicated below each pulse type and demonstrates high agreement for HU and SI ($\geq$2.77),

346 but low agreement for VO and LR (2.08).

347

348 **Table 3.** Mean values for A/V agreement index, SVM pulse classification accuracy, typicality

349 coefficient, and frequency measures by pulse type.

| Pulse | A/V index | SVM | Typicality | Center | Peak | Mean peak | 3rd quart | 1st quart |
|---|---|---|---|---|---|---|---|---|
| HU | 2.88 | 77% | 0.90 | 443.3 | 421.0 | 436.4 | 585.3 | 370.3 |
| VO | 2.08 | 41% | 0.98 | 483.1 | 442.3 | 505.6 | 592.6 | 376.7 |
| HR | 2.37 | 61% | 0.94 | 440.0 | 409.9 | 450.1 | 533.8 | 358.7 |
| LR | 2.08 | 54% | 0.81 | 266.3 | 252.2 | 271.8 | 312.0 | 231.4 |
| IN | 2.13 | 52% | 0.84 | 249.7 | 242.7 | 244.6 | 288.8 | 225.5 |
| SI | 2.77 | 81% | 0.97 | 203.0 | 201.1 | 194.6 | 239.1 | 172.5 |

350

351

352

**Supervised classification using extracted feature set: support vector machines**

353

354     We tested the performance of SVM for the classification of orangutan long call pulse types

355  using our full acoustic feature dataset. Using leave-one-out cross-validation, we found the average

356  classification accuracy of pulse types was 64.8% (range: 64.28 – 65.44 $\pm$ 0.10 SE). SVM classification

357  accuracy was higher than IRR agreement scores for most pulse types, though human observers were

358  better at discriminating huitus and sigh pulses (Fig. 3). Classification accuracy was highly variable

359  across pulse types. Whereas sighs and huituses were classified with the highest accuracy (81 and

360  77%, respectively), volcanoes were classified with the lowest accuracy (41%: Fig. 3, Table 3).

361     Recursive feature elimination revealed that center frequency, peak frequency, mean peak

362  frequency, and third and first frequency quartiles were the most influential variables (Table 3). In all

363  five influential features, high roars, huituses, and volcanoes overlapped, and in four of five features,

364  intermediaries overlapped low roars (Fig. S2, Table S2). All other pairwise comparisons of pulse

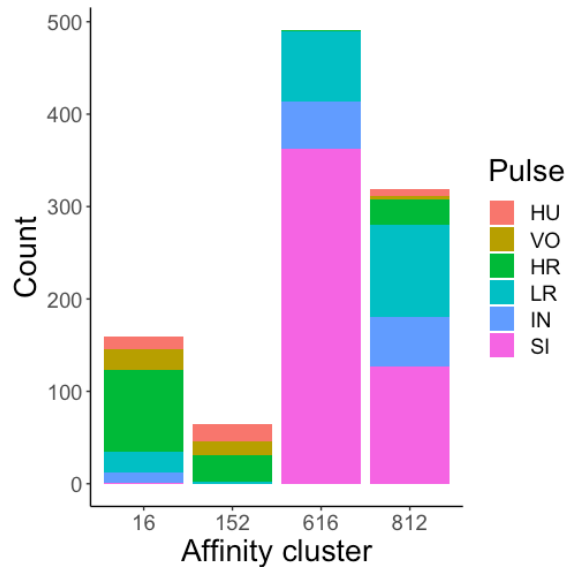365  types showed significant differences in all features.

366

**Figure 3. Barplot of classification accuracy for original pulse scheme.** Comparison of

classification accuracy of audio-visual classification (AV), calculated as the average agreement

between three observer pairs compared to supervised machine learning classification (SVM).

370

**Unsupervised clustering using extracted feature set: hard and soft clustering**

Affinity propagation resulted in four clusters with an average silhouette coefficient of 0.32

(range: -0.22 – 0.61). Of these four clusters, two (clusters 616 and 152: Fig. 4) had relatively high

silhouette coefficients (0.45 and 0.29, respectively) and separated the higher-frequency pulses (i.e.,

HU, VO, and HR pulses) from lower-frequency ones (i.e., LR, IN, and SI). The remaining two

clusters had low silhouette coefficients (cluster 16 = 0.19, cluster 812 = 0.21) and both contained

377     calls from all six pulse types (Fig. 4). We analyzed the separation of unsupervised clusters using the

378     influential features identified from recursive feature elimination (Fig. S2). Two of the four clusters

379     (16 and 152) overlapped in four of five features. These clusters primarily comprised high roars,

380     volcanoes, and huituses.

381



382

383     **Figure 4. Stacked barplots of affinity propagation clusters** showing the number of calls in each

384     cluster classified by pulse type.

385

386        In a final approach to clustering our extracted feature set, we used c-means fuzzy clustering

387     to provide another estimate of the number of clusters in our dataset and quantify the degree of

388     gradation across pulse types. All three internal validity measures (connectivity, Dunn, and silhouette)

389     and three of four stability measures (APN, AD, and ADM) indicated that the two-cluster solution

390     was optimal. Only FOM indicated a 3-cluster solution was marginally more stable (0.855 for 2 vs.

391     0.860 for 3 clusters). We found that mu = 1.1 yielded the highest average silhouette width (0.312);

392     silhouette widths decreased as mu increased.

393          Typicality coefficients were high overall (mean: 0.92 + 0.006 SE, Fig. 5) but varied widely by

394      pulse type. Whereas volcanoes and sighs had the highest typicality coefficients (0.98 and 0.97,

395      respectively) and intermediaries and low roars had the lowest coefficients (0.84 and 0.81,

396      respectively, Table 3). Pairwise comparisons of typicality coefficients showed that typicality

397      coefficients for low roars and intermediaries were significantly lower than those of all other pulse

398      types but did not significantly differ between these two pulses (Fig. S2, Table S2).

399          We determined the thresholds for typical (>0.976) and atypical calls (<0.855) (cf. Wadewitz

400      et al., 2015). Overall, 69% of calls were 'typical' for their cluster and 17% were 'atypical'; however,

401      pulse types varied greatly (Fig. 6). Whereas sighs and volcanoes had a high proportion of typical calls

402      (85% and 80% respectively), low roars and intermediaries had a high proportion of atypical calls

403      (44% and 40% respectively).

404          Typical calls were found in both clusters (Fig. 6). Typical calls in cluster one included high

405      roars, huituses, low roars, and volcanoes and those in cluster two included sighs, low roars, and

406      intermediaries. Whereas typical sighs, huituses, and volcanoes were found in only one cluster (and

407      only 1-2 intermediaries and high roars were typical for a secondary cluster), 24% of low roars

408      belonged to a secondary cluster. Overall, cluster one comprised 189 typical and 99 atypical calls

409      (53% and 28% of 353 calls, respectively) and cluster two comprised 526 typical and 75 atypical calls

410      (77% and 11% of 680 calls, respectively), indicating that calls in cluster two were better separated

411      from other call types than those in cluster one. We compared typical calls in each cluster and found

412      that calls in different clusters significantly differed from each other in all five influential features (Fig.
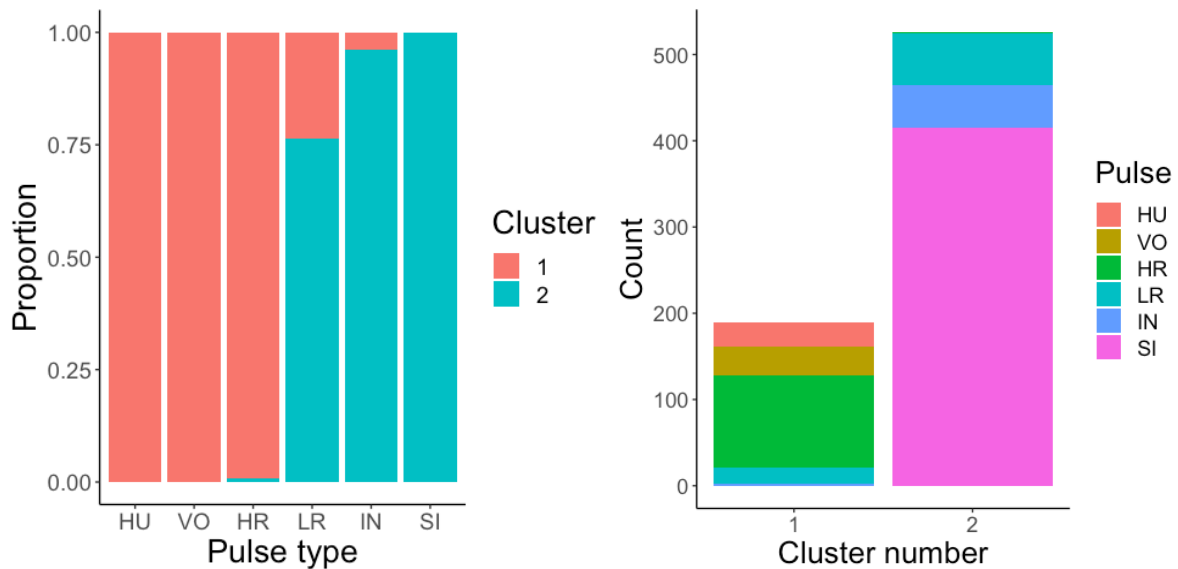
413      S2, Table S2a).

414

**Figure 5. Typicality coefficients for each pulse type** a) Histogram showing the distribution of

coefficients and b) boxplot showing typicality values for each pulse type. Typicality thresholds were

calculated following (Wadewitz et al., 2015). Typical calls were those whose typicality coefficients

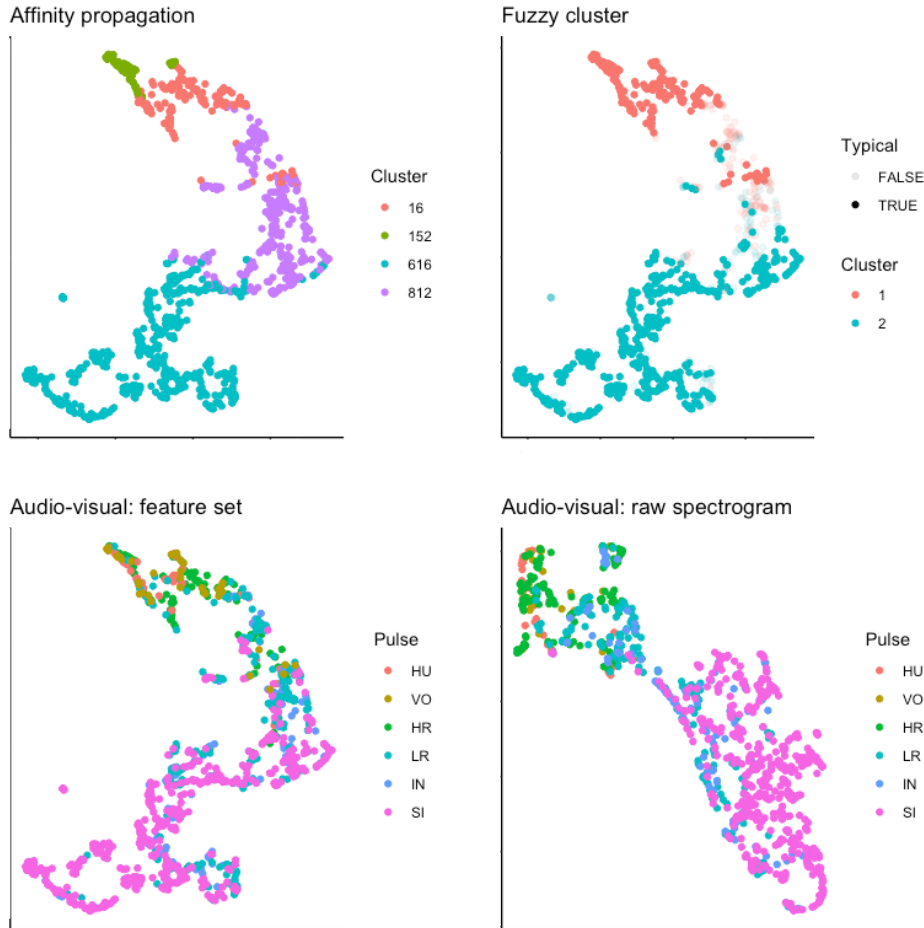exceeded 0.976 and atypical calls were those below 0.855.

419

420



421

422     **Figure 6. Stacked barplots of typical calls** a) the proportion of each pulse type that was typical

423     for each cluster and b) the number of typical calls in each cluster classified by pulse type.

424

425     **UMAP visualization of extracted features and spectrograms**

426           We used UMAP to visualize the separation of individual pulses using our extracted feature

427     set, comparing the cluster results from affinity propagation and fuzzy clustering with manual

428     classification (Fig. 7). We also used UMAP to visualize the separation of pulses based on the power

429     density spectrograms (Fig. 7). For both datasets, it appears that there are two loose and incompletely

430     separated clusters as well as a smaller number of pulses that grade continuously between the two

431     clusters. The Hopkins statistic of clusterability for the extracted feature set was 0.940 and 0.957 for

432     the power spectrograms, both of which indicate strong clusterability of calls.

433

**Figure 7. UMAP projection of 46-feature dataset.** Colors indicate four clusters identified using unsupervised affinity propagation (upper left), two clusters and typical calls identified by fuzzy clustering (upper right), six pulse types labeled by human observer using the extracted feature set (lower left), and raw power density spectrograms (lower right).

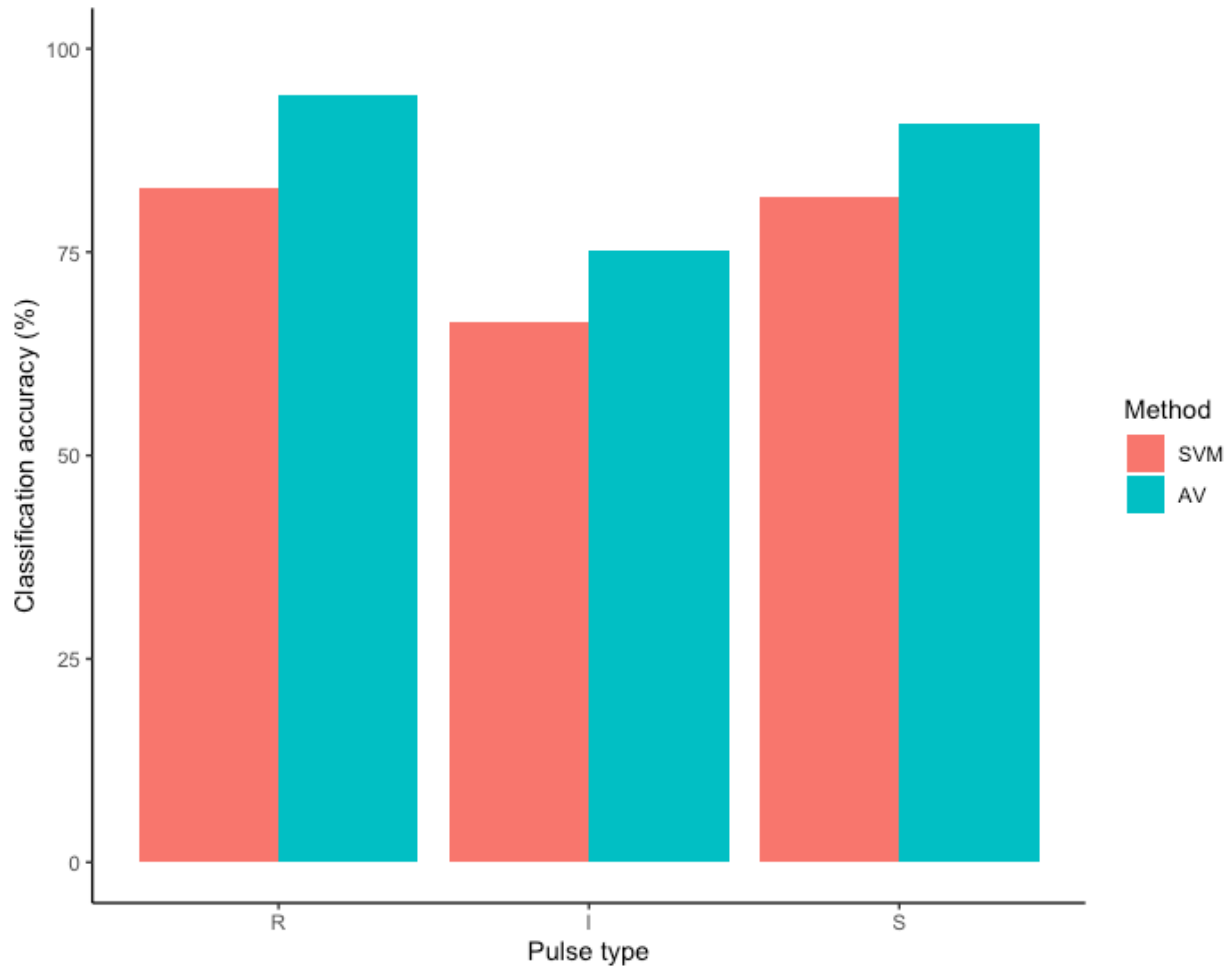**Identification and evaluation of a new classification scheme**

Collectively, our unsupervised clustering approaches showed broad agreement for a two-cluster solution with graded pulses occurring along a spectrum between the two classes. In fuzzy typical cluster 1, the mean value for F0 max was 764.3 Hz $\pm$ 351.5 SD Hz (range = 320-1,500); whereas for those pulses belonging to fuzzy typical cluster 2, the mean value of F0 max was 225.3 $\pm$

444   SD 67.9 SD Hz (range = 80-440). Pulses that were not typical for either cluster had a mean F0 max

445   of 345.8 $\pm$ 159.9 SD Hz. Based on these patterns, as well as the shape of the F0 (a feature that was

446   commonly used to distinguish among pulse types in previous studies), we distinguished among

447   pulses as follows: **Roar (R)** = F0 ascends and reaches its maximum (>350 Hz) at or near the

448   midpoint of the pulse before descending, **Sigh (S)** = F0 descends and reaches its  maximum

449   (typically, but not always < 350 Hz) at start of the pulse (i.e., no ascending portion of F0), and

450   **Intermediate (I)** = either a) maximum F0 value occurs at the start of the pulse but with an

451   ascending portion later in pulse, or b) F0 ascends and reaches its maximum (<350 Hz) at or near the

452   midpoint of the pulse.

453         These revised definitions yielded Light's kappa $\varkappa$ = 0.838 (i.e., the arithmetic mean of

454   Cohen's Kappa for observers 1-2 = 0.84, 1-3 = 0.86, and 2-3 = 0.78), indicating near-perfect

455   agreement among observers (Landis & Koch, 1977). Classification agreement varied only slightly by

456   pulse type, with roars showing the highest agreement among observers (mean 2.92, where 3

457   indicates full agreement), and intermediaries and sighs showing slightly lower agreement (mean 2.79

458   and 2.72, respectively). Using leave-one-out cross-validation, we found the average classification

459   accuracy of pulse types using SVM was 82.1% (range: 80.8 – 85.0 $\pm$ 0.47 SE). SVM classification

460   accuracy was lower than IRR agreement scores for most pulse types (Fig. 8) but both roars and sighs

461   were classified with high agreement using both methods.

462

**Figure 8. Barplot of classification accuracy for revised pulse scheme.** Comparison of classification accuracy of audio-visual classification (AV), calculated as the average agreement between three observer pairs compared to supervised machine learning classification (SVM).

**DISCUSSION**

Here we present an extensive qualitative and quantitative assessment of the vocal complexity of the long-call vocalizations of Bornean orangutans. Relying on a large dataset comprising 46 acoustic measurements from 1,033 pulses from 117 long calls recorded from 13 males, we compared the ability of human observers and supervised and unsupervised machine-learning techniques to discriminate unique call (or pulse) types. Three human observers performed relatively well at

474    discriminating two pulse types – huitus and sigh – but our inter-rater reliability score (i.e., Light's

475    kappa) showed only moderate agreement across the six pulse types. Although support vector

476    machines (SVM) performed better than human observers in classifying most pulse types (except for

477    huitus and sigh pulses), the overall accuracy was less than 65%. Like humans, SVM's were best at

478    discriminating huitus and sigh pulse types but performed relatively poorly for the others. Poor

479    classification accuracy across audio-visual and supervised machine learning approaches indicates that

480    these six pulse types are not discrete. This finding suggests that attempting comparisons of different

481    pulse types (cf. Davila Ross & Geissmann, 2007; Spillmann et al., 2010) across observers or studies

482    is not advisable, since these classes are not reliably reproduced by different observers or well

483    separated by a robust set of acoustic features.

484        Having demonstrated that these six pulse types were not well discriminated, we turned to

485    unsupervised clustering to characterize and classify the diversity of pulses comprising orangutan long

486    calls. Whereas hard clustering, such as affinity propagation, seeks to identify a set of high-quality

487    exemplars and corresponding clusters (Frey & Dueck, 2007), soft, or fuzzy, clustering is an

488    alternative or complementary approach to evaluate and quantify the discreteness of call types within

489    a graded repertoire (Cusano et al., 2021; Fischer et al., 2016; Wadewitz et al., 2015). Although the

490    hard and soft unsupervised techniques yielded different clustering solutions – four clusters for

491    affinity propagation and two for fuzzy c-means – both methods showed relatively poor separation

492    across pulse types. Importantly, both hard and soft clustering solutions separated high-frequency

493    pulses (i.e., HU, VO, and HR) from low-frequency ones (i.e., LR, IN, SI), but low roars and

494    intermediaries showed low typicality coefficients and occurred in both fuzzy clusters. Together, the

495    results of unsupervised clustering support our interpretation of the manual and supervised

496    classification analysis in demonstrating that orangutan long calls contain a mixture of discrete and

497    graded pulse types.

498  We used a final approach, UMAP, to visualize the separation and quantify the clusterability

499 of call types. Because the number and type of features selected can have a strong influence on the

500 cluster solutions and their interpretations (Fischer et al., 2016; Wadewitz et al., 2015), we compared

501 the results of our extracted 46-feature dataset with raw power density spectrograms as inputs. Both

502 datasets yielded similar and high Hopkin's statistic values, indicating strong clusterability of calls. At

503 the same time, both datasets generated a V-shaped cloud of points showing two large loose clusters

504 with a spectrum of points lying along a continuum between them.

505  Based on our comprehensive evaluation of orangutan pulse types, we have proposed a

506 revised approach to the classification of orangutan pulses that we hope provides improved

507 reproducibility for future researchers. We recommend using the following terms to categorize the

508 range of pulse types comprising orangutan long calls: 1) 'Roar' for high-frequency pulses, 2) 'Sigh'

509 for low-frequency pulses, and 3) 'Intermediate' for graded pulses that fall between these two

510 extremes. We have provided detailed descriptions of each of these pulse types and demonstrated

511 that they can be easily and reliably identified among different observers and exhibit high

512 classification accuracy using SVM.

513  Thus, we find that orangutan calls can be clustered into three pulse types (two discrete and

514 one graded). The low diversity of call types suggests that these vocalizations are not particularly

515 complex. Like the long-calls of other apes (chimpanzees, *Pan troglodytes schweinfurthii*: Arcadi, 1996;

516 Marler & Hobbett, 1975; gibbons, *Hylobates* spp: Marshall & Marshall, 1976), orangutan long calls

517 typically comprise an intro and/or build-up phase (quiet, staccato grumbles, not analyzed in the

518 present study), climax (high-energy, high-frequency roars), and a let-down phase (low-energy sighs).

519 The low number of discrete pulse types could be interpreted as support for the hypothesis that long-

520 distance signals have been selected to facilitate signal recognition in dense and noisy habitats (Marler

521     et al., 1975). Yet, there is a spectrum of intermediate call types that yield a continuous gradation of

522     sounds across phases and pulses, that seems to greatly boost the complexity of this signal.

523          Unfortunately, only a handful of studies have quantified the gradedness of animal vocal

524     systems (but see Cusano et al., 2021; Fischer et al., 2017; Taylor et al., 2021; Wadewitz et al., 2015).

525     Consequently, we are still lacking a comprehensive framework through which to quantify and

526     interpret vocal complexity vis-à-vis graded repertoires (Fischer et al., 2017). Future research will

527     explore the production of graded call types across individuals and call types to examine the sources

528     of variation and the potential role of graded call types in orangutan communication.

529          In summary, we evaluated a range of supervised and unsupervised approaches to classifying

530     and clustering sounds in animal vocal repertoires. We used a combination of traditional audio-visual

531     methods and modern machine learning techniques that relied on human eyes and ears, a set of 46

532     features measured from spectrograms, and raw power density spectrograms to triangulate diverse

533     datasets and methods to answer a simple question: how many pulse types exist within orangutan

534     long calls, how can they be distinguished, and how graded are they? While each approach has its

535     strengths and limitations, taken together, they can lead to a more holistic understanding of call types

536     within graded repertoires and contribute to a growing body of literature documenting the graded

537     nature of animal communication systems.

538    **SUPPLEMENTARY MATERIALS**

539    **Table S1.** Table describing features measured in Raven Pro and warbleR (Specan and freq_ts**)**

| No | Program | Feature | Description |
|---|---|---|---|
| 1 | Raven | Delta.Time.s | difference between Begin Time and End Time for the selection (s) |
| 2 | Raven | Freq.5%.Hz | frequency at which summed energy exceeds 5% of total energy |
| 3 | Raven | Freq.95%.Hz | frequency at which summed energy exceeds 95% of total energy |
| 4 | Raven | Agg.Entropy.bits | aggregate entropy measures the disorder in a sound by analyzing the energy distribution (pure tone ~ 0) |
| 5 | Raven | Avg.Entropy.bits | average entropy measurement describes the amount of disorder for a typical spectrum within the selection |
| 6 | Raven | BW.50% | difference between the 25% and 75% frequencies (Hz) |
| 7 | Raven | BW.90% | difference between the 5% and 95% frequencies (Hz) |
| 8 | Raven | Center.Freq | frequency that divides the selection into two frequency intervals of equal energy (Hz) |
| 9 | Raven | Center.Time.Rel. | proportion of selection at which 50% of the sound energy has an earlier time |
| 10 | Raven | Dur.50% | difference between the 25% and 75% times (s) |
| 11 | Raven | Dur.90% | difference between the 5% and 95% times (s) |
| 12 | Raven | Freq.25% | frequency at which summed energy exceeds 25% of total energy (Hz) |
| 13 | Raven | Freq.75% | frequency at which summed energy exceeds 75% of total energy (Hz) |
| 14 | Raven | Peak.Freq | frequency at which Peak Power occurs within the selection (Hz) |
| 15 | Raven | PFC.Avg.Slope | Mean of the Peak Frequency Contour Slope Series of numbers (Hz/ms) |
| 16 | Raven | PFC.Max.Freq | Maximum of the Peak Frequency Contour Series of numbers (Hz) |
| 17 | Raven | PFC.Max.Slope | Maximum of the Peak Frequency Contour Slope Series of numbers (Hz/ms) |
| 18 | Raven | PFC.Min.Freq | Minimum of the Peak Frequency Contour Series of numbers (Hz) |
| 19 | Raven | PFC.Min.Slope | Minimum of the Peak Frequency Contour Slope Series of numbers (Hz/ms) |
| 20 | Raven | PFC.Num.Inf.Pts | Number of times the slope changes sign in Peak Frequency Contour Slope Series of numbers |
| 21 | Raven | Peak.Time.Rel. | proportion of selection at first time in a selection at which amplitude equal to Peak Amplitude occurs |
| 22 | Raven | Time.25%.Rel. | proportion of selection at which 25% of the sound energy has an earlier time |
| 23 | Raven | Time.5%.Rel. | proportion of selection at which 5% of the sound energy has an earlier time |
| 24 | Raven | Time.75%.Rel. | proportion of selection at which 75% of the sound energy has an earlier time |

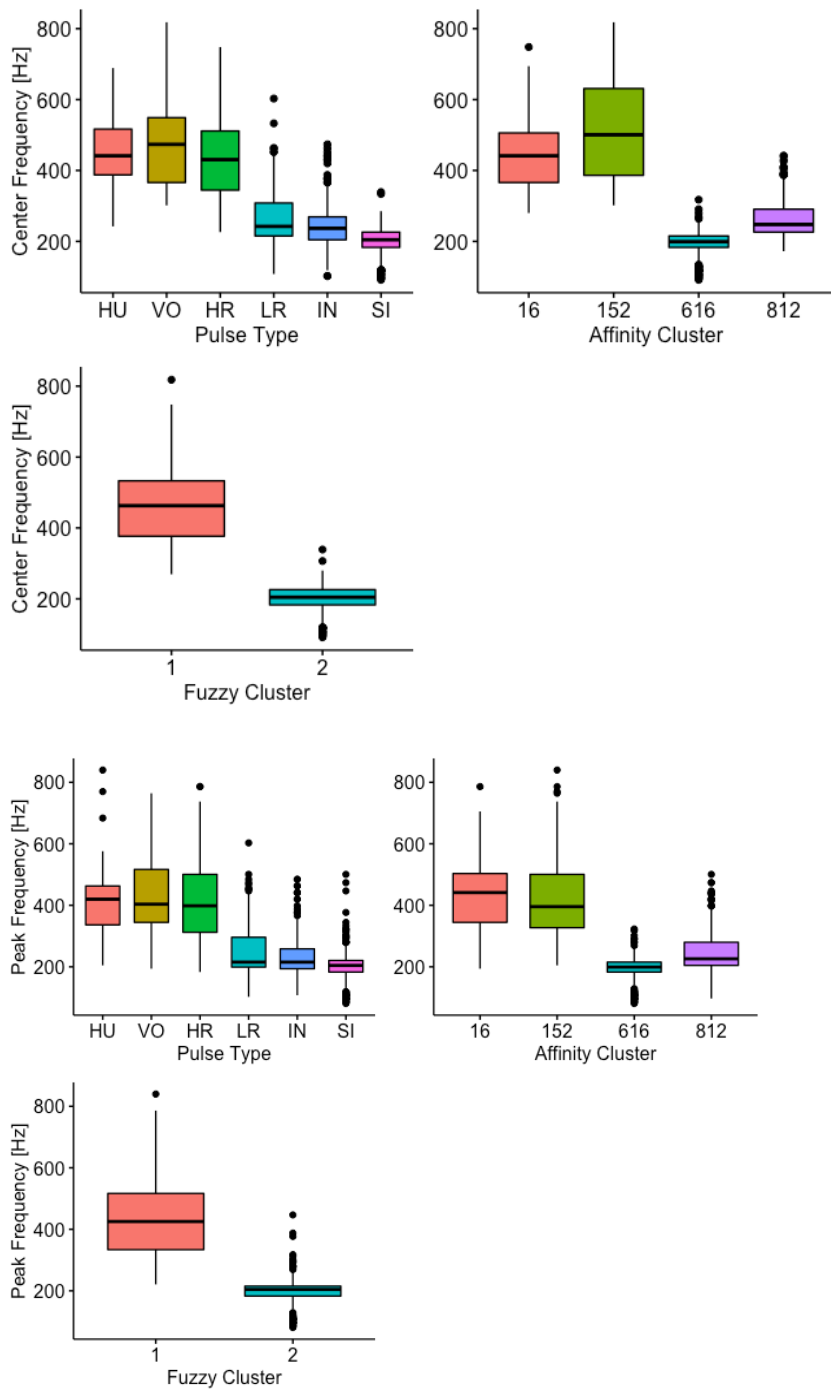| 25 | Raven | Time.95%.Rel. | proportion of selection at which 95% of the sound energy has an earlier time |
|----|-------|---------------|------------------------------------------------------------------------------|
| 26 | specan | meanfreq | mean of frequency spectrum (kHz) |
| 27 | specan | sd | standard deviation of frequency (kHz) |
| 28 | specan | skew | skewness: asymmetry of the spectrum |
| 29 | specan | kurt | kurtosis: peakedness of the spectrum |
| 30 | specan | sp.ent | energy distribution of the frequency spectrum (pure tone ~ 0) |
| 31 | specan | time.ent | energy distribution on the time envelope (pure tone ~ 0) |
| 32 | specan | entropy | spectrographic entropy: product of time x spectral entropy |
| 33 | specan | sfm | spectral flatness (pure tone ~ 0) |
| 34 | specan | meandom | average of dominant frequency measured across the acoustic signal |
| 35 | specan | mindom | minimum of dominant frequency measured across the acoustic signal |
| 36 | specan | maxdom | maximum of dominant frequency measured across the acoustic signal |
| 37 | specan | dfrange | range of dominant frequency measured across the acoustic signal |
| 38 | specan | modindx | modulation index: cumulative difference between adjacent dominant frequencies / dominant frequency range |
| 39 | specan | startdom | dominant frequency measurement at the start of the signal |
| 40 | specan | enddom | dominant frequency measurement at the end of the signal |
| 41 | specan | dfslope | slope of the change in dominant frequency through time |
| 42 | specan | meanpeakf | frequency with highest energy from the mean frequency spectrum |
| 43 | specan | Freq_IQR | interquartile frequency range. Frequency range between 'freq.Q25' and 'freq.Q75' (kHz) |
| 44 | specan | Time_IQR | interquartile time range. Time range between 'time.Q25' and 'time.Q75' (s) |
| 45 | freq_ts | F0_min | frequency at which F0 contour is at its minimum value (kHz) |
| 46 | freq_ts | F0_max | frequency at which F0 contour reaches its maximum value (kHz) |

540

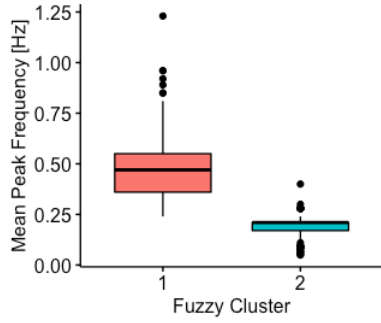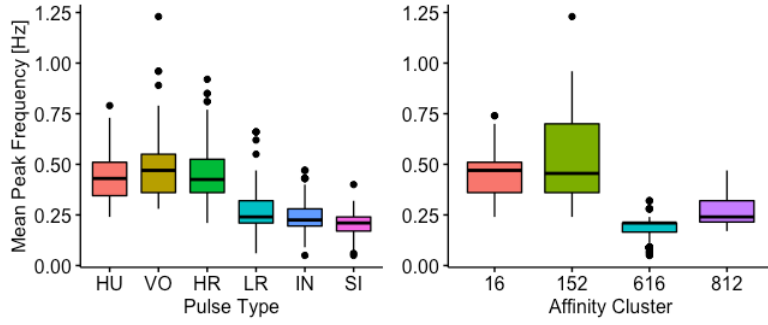541

542     **Figure S1.** Example of annotated spectrogram

544 **Figure S2.** Boxplots of features that differed across human-labeled pulses (upper left), affinity

545 propagation clusters (upper right), and typical calls in fuzzy clusters (lower left) for each of the

546 following influential features: a) center frequency, b) peak frequency, c) mean peak frequency, d)

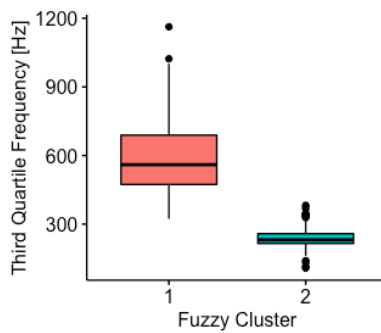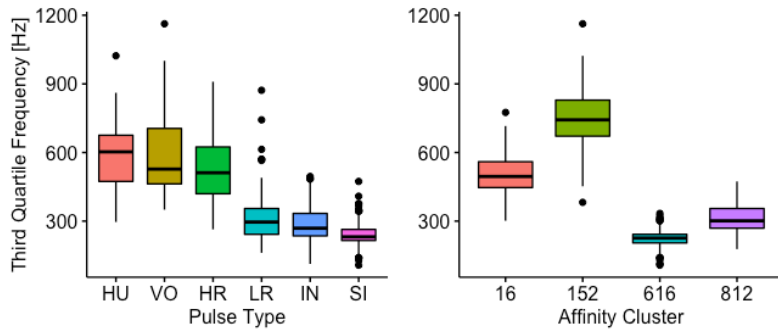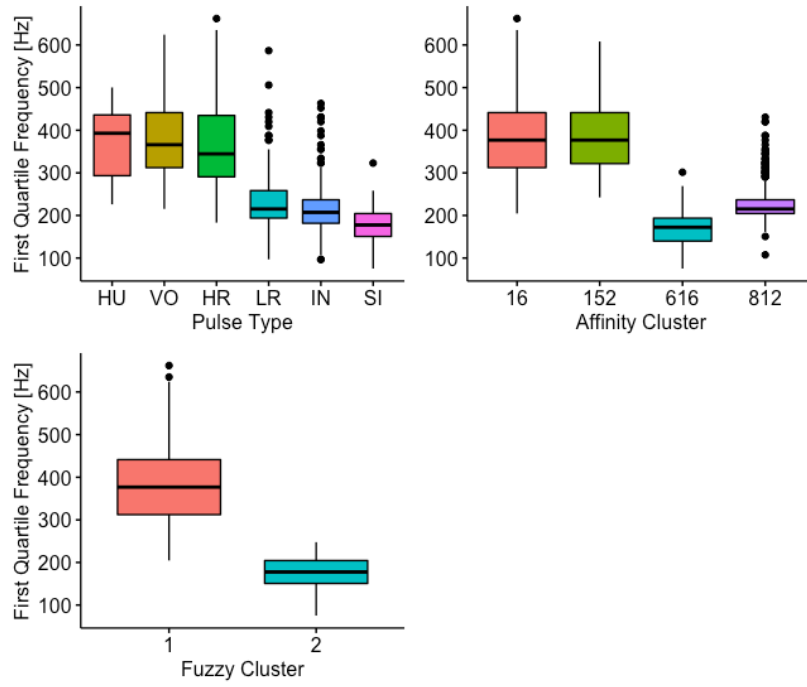547 third quartile frequency, e) first quartile frequency.

548



549

550



551

552

553    **Table S2a.** Table summarizing results of Kruskal-Wallis tests for differences among pulses or clusters identified by human observers,

554    affinity propagation, and fuzzy clustering for each of the top five influential variables.

| Variable | A/V | | | AFFINITY | | | FUZZY | | |
|---|---|---|---|---|---|---|---|---|---|
| | $\chi^2$ | df | p | $\chi^2$ | df | p | $\chi^2$ | df | p |
| Center | 557.81 | 5.00 | 0.00 | 738.53 | 3.00 | 0.00 | 417.16 | 1.00 | 0.00 |
| Peak | 425.31 | 5.00 | 0.00 | 588.78 | 3.00 | 0.00 | 406.49 | 1.00 | 0.00 |
| Mean peak | 528.18 | 5.00 | 0.00 | 677.78 | 3.00 | 0.00 | 421.95 | 1.00 | 0.00 |
| Third quart | 570.30 | 5.00 | 0.00 | 777.79 | 3.00 | 0.00 | 416.69 | 1.00 | 0.00 |
| First quart | 536.72 | 5.00 | 0.00 | 684.50 | 3.00 | 0.00 | 414.47 | 1.00 | 0.00 |

555

556    **Table S2b.** Table summarizing results of Dunn tests for pair-wise differences among pulses identified by human observers for each of the

557    top five influential variables.

| A/V | Center | | | Peak | | | Mean peak | | | Third quart | | | First quart | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Pair | Z | P.unadj | P.adj | Z | P.unadj | P.adj | Z | P.unadj | P.adj | Z | P.unadj | P.adj | Z | P.unadj | P.adj |
| HR-HU | -0.21 | 0.84 | 0.84 | -0.54 | 0.59 | 0.63 | 0.12 | 0.90 | 0.90 | -0.60 | 0.55 | 0.64 | -0.54 | 0.59 | 0.68 |
| HR-IN | 10.29 | 0.00 | 0.00 | 9.12 | 0.00 | 0.00 | 10.59 | 0.00 | 0.00 | 10.52 | 0.00 | 0.00 | 9.04 | 0.00 | 0.00 |
| HR-LR | 9.80 | 0.00 | 0.00 | 9.81 | 0.00 | 0.00 | 9.68 | 0.00 | 0.00 | 9.97 | 0.00 | 0.00 | 9.35 | 0.00 | 0.00 |
| HR-SI | 19.72 | 0.00 | 0.00 | 17.10 | 0.00 | 0.00 | 19.35 | 0.00 | 0.00 | 19.86 | 0.00 | 0.00 | 19.26 | 0.00 | 0.00 |
| HR-VO | -0.67 | 0.50 | 0.58 | -0.84 | 0.40 | 0.50 | -0.44 | 0.66 | 0.71 | -0.55 | 0.58 | 0.62 | -0.50 | 0.62 | 0.66 |
| HU-IN | 7.12 | 0.00 | 0.00 | 6.65 | 0.00 | 0.00 | 7.00 | 0.00 | 0.00 | 7.65 | 0.00 | 0.00 | 6.60 | 0.00 | 0.00 |
| HU-LR | 6.31 | 0.00 | 0.00 | 6.66 | 0.00 | 0.00 | 5.90 | 0.00 | 0.00 | 6.81 | 0.00 | 0.00 | 6.37 | 0.00 | 0.00 |
| HU-SI | 11.40 | 0.00 | 0.00 | 10.27 | 0.00 | 0.00 | 10.83 | 0.00 | 0.00 | 11.90 | 0.00 | 0.00 | 11.49 | 0.00 | 0.00 |
| HU-VO | -0.36 | 0.72 | 0.77 | -0.23 | 0.82 | 0.82 | -0.44 | 0.66 | 0.76 | 0.04 | 0.97 | 0.97 | 0.04 | 0.97 | 0.97 |
| IN-LR | -1.83 | 0.07 | 0.08 | -0.57 | 0.57 | 0.66 | -2.26 | 0.02 | 0.03 | -1.91 | 0.06 | 0.07 | -0.91 | 0.36 | 0.45 |
| IN-SI | 5.60 | 0.00 | 0.00 | 4.63 | 0.00 | 0.00 | 4.91 | 0.00 | 0.00 | 5.46 | 0.00 | 0.00 | 6.70 | 0.00 | 0.00 |
| IN-VO | -7.70 | 0.00 | 0.00 | -7.06 | 0.00 | 0.00 | -7.67 | 0.00 | 0.00 | -7.74 | 0.00 | 0.00 | -6.67 | 0.00 | 0.00 |
| LR-SI | 9.44 | 0.00 | 0.00 | 6.48 | 0.00 | 0.00 | 9.18 | 0.00 | 0.00 | 9.38 | 0.00 | 0.00 | 9.51 | 0.00 | 0.00 |
| LR-VO | -6.92 | 0.00 | 0.00 | -7.09 | 0.00 | 0.00 | -6.60 | 0.00 | 0.00 | -6.90 | 0.00 | 0.00 | -6.45 | 0.00 | 0.00 |
| SI-VO | -12.16 | 0.00 | 0.00 | -10.83 | 0.00 | 0.00 | -11.70 | 0.00 | 0.00 | -12.12 | 0.00 | 0.00 | -11.71 | 0.00 | 0.00 |

558

559    **Table S2c.** Table summarizing results of Dunn tests for pair-wise differences among clusters identified by affinity propagation for each of

560    the top five influential variables.

| AFFINITY | Center | | | Peak | | | Mean peak | | | Third quart | | | First quart | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Pair | Z | P.unadj | P.adj | Z | P.unadj | P.adj | Z | P.unadj | P.adj | Z | P.unadj | P.adj | Z | P.unadj | P.adj |
| 152-16 | 0.84 | 0.40 | 0.40 | -0.43 | 0.66 | 0.66 | 0.24 | 0.81 | 0.81 | 2.33 | 0.02 | 0.02 | 0.20 | 0.84 | 0.84 |
| 152-616 | 16.64 | 0.00 | 0.00 | 14.23 | 0.00 | 0.00 | 15.60 | 0.00 | 0.00 | 18.09 | 0.00 | 0.00 | 15.68 | 0.00 | 0.00 |
| 16-616 | 22.88 | 0.00 | 0.00 | 21.43 | 0.00 | 0.00 | 22.33 | 0.00 | 0.00 | 22.56 | 0.00 | 0.00 | 22.51 | 0.00 | 0.00 |
| 152-812 | 7.95 | 0.00 | 0.00 | 7.59 | 0.00 | 0.00 | 7.57 | 0.00 | 0.00 | 8.88 | 0.00 | 0.00 | 7.73 | 0.00 | 0.00 |
| 16-812 | 9.94 | 0.00 | 0.00 | 11.37 | 0.00 | 0.00 | 10.31 | 0.00 | 0.00 | 8.97 | 0.00 | 0.00 | 10.60 | 0.00 | 0.00 |
| 616-812 | -15.61 | 0.00 | 0.00 | -11.83 | 0.00 | 0.00 | -14.41 | 0.00 | 0.00 | -16.52 | 0.00 | 0.00 | -14.26 | 0.00 | 0.00 |

561

589

590

591

592     **REFERENCES**

593     Alloghani, M., Al-Jumeily, D., Mustafina, J., Hussain, A., & Aljaaf, A. J. (2020). A Systematic Review

594             on Supervised and Unsupervised Machine Learning Algorithms for Data Science. In

595             *Unsupervised and Semi-Supervised Learning: Supervised and Unsupervised Learning for Data Science* (pp.

596             3–21). Springer International Publishing.

597     Altmann, J. (1974). Observational study of behavior: Sampling methods. *Behaviour*, *49*, 227–267.

598     Araya-Salas, M., & Smith-Vidaurre, G. (2017). WarbleR: an r package to streamline analysis of

599             animal acoustic signals. *Methods in Ecology and Evolution*, *8*(2), 184–191.

600     Arcadi, A. C. (1996). Phrase structure of wild chimpanzee pant hoots: Patterns of production and

601             interpopulation variability. *American Journal of Primatology*, *39*(3), 159–178.

602     Askew, J. A., & Morrogh-Bernard, H. C. (2016). Acoustic characteristics of long calls produced by

603             male orangutans (Pongo pygmaeus wurmbii): Advertising individual identity, context, and

604             travel direction. *Folia Primatologica*, *87*(5), 305–319.

605     Blumstein, D. T., & Armitage, K. B. (1997). Does sociality drive the evolution of communicative

606             complexity? A comparative test with ground-dwelling sciurid alarm calls. *The American

607             Naturalist*, *150*(2), 179–200.

608     Bodenhofer, U., Kothmeier, A., & Hochreiter, S. (2011). APCluster: An R package for affinity

609             propagation clustering. *Bioinformatics*, *27*, 2463–2464.

610          https://doi.org/10.1093/bioinformatics/btr406

611   Bradbury, J. W., & Vehrencamp, S. L. (2011). *Principles of animal communication (2nd ed.).* Sinauer

612          Associates.

613   Brady, B., Hedwig, D., Trygonis, V., & Gerstein, E. (2020). Classification of Florida manatee

614          (*Trichechus manatus latirostris*) vocalizations. *The Journal of the Acoustical Society of America*, *147*(3),

615          1597–1606. https://doi.org/10.1121/10.0000849

616   Brock, G., Pihur, V., Datta, S., & Datta, S. (2008). clValid: An R package for Cluster Validation.

617          *Journal of Statistical Software*, *25*(4), 1–22.

618   Clink, D. J., Crofoot, M. C., & Marshall, A. J. (2018). Application of a semi-automated vocal

619          fingerprinting approach to monitor Bornean gibbon females in an experimentally

620          fragmented landscape in Sabah, Malaysia. *Bioacoustics*, 1–17.

621   Clink, D. J., & Klinck, H. (2020). Unsupervised acoustic classification of individual gibbon females

622          and the implications for passive acoustic monitoring. *Methods in Ecology and Evolution*.

623   Cunningham, P., Cord, M., & Delany, S. (2008). Supervised learning. In M. Cord & P. Cunningham

624          (Eds.), *Machine Learning Techniques for Multimedia.* Springer. https://doi.org/10.1007/978-3-

625          540-75171-7_2

626   Cusano, D. A., Noad, M. J., & Dunlop, R. A. (2021). Fuzzy clustering as a tool to differentiate

627          between discrete and graded call types. *JASA Express Letters*, *1*(6), 061201.

628   Davila Ross, M., & Geissmann, T. (2007). Call diversity of wild male orangutans: A phylogenetic

629          approach. *American Journal of Primatology*, *69*(3), 305–324.

630   Elie, J., & Theunissen, F. (2016). The vocal repertoire of the domesticated zebra finch: A data-driven

631          approach to decipher the information-bearing acoustic features of communication signals.

632          *Animal Cognition*, *19*(2), 285–315.

633   Erb, W. M., Barrow, E. J., Hofner, A. N., Utami-Atmoko, S. S., & Vogel, E. R. (2018). Wildfire

634      smoke impacts activity and energetics of wild Bornean orangutans. *Scientific Reports*, *8*(1),

635      7606. https://doi.org/10.1038/s41598-018-25847-1

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

637      *Scientific Reports*, *6*, 38226.

638   Fischer, J., Wadewitz, P., & Hammerschmidt, K. (2017). Structural variability and communicative

639      complexity in acoustic communication. *Animal Behaviour*.

640   Fournet, M. E., Szabo, A., & Mellinger, D. K. (2015). Repertoire and classification of non-song calls

641      in Southeast Alaskan humpback whales (Megaptera novaeangliae). *The Journal of the Acoustical*

642      *Society of America*, *137*(1), 1–10. https://doi.org/10.1121/1.4904504

643   Freeberg, T. M., Dunbar, R. I. M., & Ord, T. J. (2012). Social complexity as a proximate and ultimate

644      factor in communicative complexity. *Philosophical Transactions of the Royal Society London*

645      *Biological Sciences*, *367*(1597), 1785–1801.

646   Frey, B., & Dueck, D. (2007). Clustering by passing messages between data points. *Science*, *315*(5814),

647      972–976.

648   Fuller, J. L. (2014). The vocal repertoire of adult male blue monkeys (Cercopithecus mitis

649      stulmanni): A quantitative analysis of acoustic structure. *American Journal of Primatology*, *76*(3),

650      203–216. https://doi.org/10.1002/ajp.22223

651   Gamer, M., Lemon, J., Fellows, I., & Singh, P. (2012). *irr: Various coefficients of interrater reliability and*

652      *agreement*.

653   Garland, E., Castellote, M., & Berchok, C. (2015). Beluga whale (Delphinapterus leucas)

654      vocalizations and call classification from the eastern Beaufort Sea population. *The Journal of*

655      *the Acoustical Society of America*, *137*(6), 3054–3067.

656   Greene, D., Cunningham, P., & Mayer, R. (2008). Unsupervised learning and clustering. In M. Cord

657      & P. Cunningham (Eds.), *Machine Learning Techniques for Multimedia*. Springer.

658        https://doi.org/10.1007/978-3-540-75171-7_3

659    Hallgren, K. A. (2012). Computing inter-rater reliability for observational data: An overview and

660        tutorial. *Tutorials in Quantitative Methods for Psychology*, *8*(1), 23.

661    Hammerschmidt, K., & Fischer, J. (1998). The vocal repertoire of Barbary macaques: A quantitative

662        analysis of a graded signal system. *Ethology*, *104*(3), 203–216.

663    Hammerschmidt, K., & Fischer, J. (2019). Baboon vocal repertoires and the evolution of primate

664        vocal diversity. *Journal of Human Evolution*, *126*, 1–13.

665        https://doi.org/10.1016/j.jhevol.2018.10.010

666    Hedwig, D., Verahrami, A. K., & Wrege, P. H. (2019). Acoustic structure of forest elephant rumbles:

667        A test of the ambiguity reduction hypothesis. *Animal Cognition*, *22*(6), 1115–1128.

668        https://doi.org/10.1007/s10071-019-01304-y

669    Huijser, L. A. E., Estrade, V., Webster, I., Mouysset, L., Cadinouche, A., & Dulau-Drouot, V.

670        (2020). Vocal repertoires and insights into social structure of sperm whales (*Physeter*

671        *macrocephalus*) in Mauritius, southwestern Indian Ocean. *Marine Mammal Science*, *36*(2), 638–

672        657. https://doi.org/10.1111/mms.12673

673    Janik, V. M. (1999). Pitfalls in the categorization of behaviour: A comparison of dolphin whistle

674        classification methods. *Animal Behaviour*, *57*(1), 133–143.

675    Jones, A. E., ten Cate, C., & Bijleveld, C. C. J. H. (2001). The interobserver reliability of scoring

676        sonagrams by eye: A study on methods, illustrated on zebra finch songs. *Animal Behaviour*,

677        *62*(4), 791–801.

678    K. Lisa Yang Center for Conservation Bioacoustics. (2019). *Raven Pro: Interactive Sound Analysis*

679        *Software (Version 1.6.1) [Computer software]. Ithaca, NY: The Cornell Lab of Ornithology.* Available

680        from http://ravensoundsoftware.com

681    Kassambara, A., & Mundt, F. (2020). *factoextra: Extract and Visualize the Results of Multivariate Data*

682          *Analyses*. http://www.sthda.com/english/rpkgs/factoextra

683    Kershenbaum, A., Blumstein, D., Roch, M., Akçay, Ç., Backus, G., Bee, M., Bohn, K., Cao, Y.,

684          Carter, G., Cäsar, C., Coen, M., DeRuiter, S., Doyle, L., Edelman, S., Ferrer-i-Cancho, R.,

685          Freeberg, T., Garland, E., Gustison, M., Harley, H., … Zamora-Gutierrez, V. (2016).

686          Acoustic sequences in non-human animals: A tutorial review and prospectus. *Biol Rev Camb*

687          *Philos Soc*, *91*(1), 13–52.

688    Landis, J. R., & Koch, G. G. (1977). The measurement of observer agreement for categorical data.

689          *Biometrics*, *33*(1), 159–174.

690    Lattenkamp, E. Z. (2019). The Vocal Repertoire of Pale Spear-Nosed Bats in a Social Roosting

691          Context. *Frontiers in Ecology and Evolution*, *7*, 14.

692    Light, R. J. (1971). Measures of response agreement for qualitative data: Some generalizations and

693          alternatives. *Psychological Bulletin*, *76*(5), 365.

694    MacKinnon, J. (1977). A comparative ecology of Asian apes. *Primates*.

695    Madhusudhana, S. K., Chakraborty, B., & Latha, G. (2019). Humpback whale singing activity off the

696          Goan coast in the Eastern Arabian Sea. *Bioacoustics*, *28*(4), 329–344.

697          https://doi.org/10.1080/09524622.2018.1458248

698    Maechler, M., Rousseeuw, P., Struyf, A., Hubert, M., & Hornik, K. (2021). *cluster: Cluster Analysis*

699          *Basics and Extensions*. https://CRAN.R-project.org/package=cluster

700    Marler, P., & Hobbett, L. (1975). Individuality in a long-range vocalization of wild chimpanzees.

701          *Zeitschrift Für Tierpsychologie*, *38*(1), 97–109.

702    Marler, P., Kavanaugh, J. F., & Cutting, J. E. (1975). On the origin of speech from animal sounds. In

703          *On the origin of speech from animal sounds*. MIT Press.

704    Marshall, J., & Marshall, E. (1976). Gibbons and their territorial songs. *Science*, *193*(4249), 235–237.

705    McComb, K., & Semple, S. (2005). Coevolution of vocal communication and sociality in primates.

706      *Biol Lett*, *1*(4), 381–385.

707   McInnes, L., Healy, J., & Melville, J. (2018). UMAP: Uniform Manifold Approximation and

708      Projection for Dimension Reduction. *ArXiv*, 1802.03426v3.

709   Meyer, D., Dimitriadou, E., Hornik, K., Weingessel, A., & Leisch, F. (2021). *e1071: Misc functions of*

710      *the Department of Statistics, probability*. https://CRAN.R-project.org/package=e1071

711   Mitra Setia, T., & van Schaik, C. P. (2007). The response of adult orang-utans to flanged male long

712      calls: Inferences about their function. *Folia Primatologica*, *78*(4), 215–226.

713   Odom, K., Araya-Salas, M., Morano, J., Ligon, R., Leighton, G., Taff, C., Dalziell, A., Billings, A.,

714      Germain, R., Pardo, M., de Andrade, L., Hedwig, D., Keen, S., Shiu, Y., Charif, R., Webster,

715      M., & Rice, A. (2021). Comparative bioacoustics: A roadmap for quantifying and comparing

716      animal sounds across diverse taxa. *Biol Rev Camb Philos Soc.*

717   Ogle, D. H., Doll, J. C., Wheeler, P., & Dinno, A. (2022). *FSA: Fisheries Stock Analysis.*

718      https://github.com/fishR-Core-Team/FSA

719   Sadhukhan, S., Hennelly, L., & Habib, B. (2019). Characterising the harmonic vocal repertoire of the

720      Indian wolf (Canis lupus pallipes). *PLOS ONE*, *14*(10), e0216186.

721      https://doi.org/10.1371/journal.pone.0216186

722   Sainburg, T., Thielk, M., & Gentner, T. (2020). Finding, visualizing, and quantifying latent structure

723      across diverse animal vocal repertoires. *PLoS Comput Biol*, *16*(10), e1008228.

724   Schwing, R., Parsons, S., & Nelson, X. J. (2012). Vocal repertoire of the New Zealand kea parrot

725      Nestor notabilis. *Current Zoology.*

726   Soltis, J., Alligood, C., Blowers, T., & Savage, A. (2012). The vocal repertoire of the Key Largo

727      woodrat (Neotoma floridana smalli). *J Acoust Soc Am*, *132*(5), 3550–3558.

728   Spillmann, B., Dunkel, L. P., van Noordwijk, M. A., Amda, R. N. A., Lameira, A. R., Wich, S. A., &

729      van Schaik, C. P. (2010). Acoustic properties of long calls given by flanged male orang-utans

730    (*Pongo pygmaeus wurmbii*) reflect both individual identity and context. *Ethology*, *116*(5), 385–395.

731    https://doi.org/10.1111/j.1439-0310.2010.01744.x

732  Spillmann, B., Willems, E. P., van Noordwijk, M. A., Setia, T. M., & van Schaik, C. P. (2017).

733    Confrontational assessment in the roving male promiscuity mating system of the Bornean

734    orangutan. *Behavioral Ecology and Sociobiology*, *71*(1), 20.

735  Taylor, D., Dezecache, G., & Davila-Ross, M. (2021). Filling in the gaps: Acoustic gradation

736    increases in the vocal ontogeny of chimpanzees (Pan troglodytes). *Am J Primatol*, *83*(5),

737    e23249.

738  Thiebault, A., Charrier, I., Pistorius, P., & Aubin, T. (2019). At sea vocal repertoire of a foraging

739    seabird. *Journal of Avian Biology*, *50*(5).

740  Turesson, H., Ribeiro, S., Pereira, D., Papa, J., & de Albuquerque, V. (2016). Machine learning

741    algorithms for automatic classification of marmoset vocalizations. *PLoS One*, *11*(9),

742    e0163041.

743  Vester, H., Hallerberg, S., Timme, M., & Hammerschmidt, K. (2017). Vocal repertoire of long-

744    finned pilot whales (Globicephala melas) in northern Norway. *J Acoust Soc Am*, *141*(6), 4289.

745  Wadewitz, P., Hammerschmidt, K., Battaglia, D., Witt, A., Wolf, F., & Fischer, J. (2015).

746    Characterizing Vocal Repertoires—Hard vs. Soft Classification Approaches. *PLOS ONE*,

747    *10*(4), e0125785. https://doi.org/10.1371/journal.pone.0125785

748  Zhou, K., Fu, C., & Yang, S. (2014). Fuzziness parameter selection in fuzzy c-means: The

749    perspective of cluster validation. *Science China Information Sciences*, *57*(11), 1–8.