



HAL
open science

Applying machine learning to primate bioacoustics: review and perspectives

Jules Cauzinille, Benoît Favre, Ricard Marxer, Arnaud Rey

► **To cite this version:**

Jules Cauzinille, Benoît Favre, Ricard Marxer, Arnaud Rey. Applying machine learning to primate bioacoustics: review and perspectives. *American Journal of Primatology*, inPress, 10.1002/ajp.23666 . hal-04658068

HAL Id: hal-04658068

<https://hal.science/hal-04658068v1>

Submitted on 22 Jul 2024

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.



Distributed under a Creative Commons Attribution - NonCommercial 4.0 International License

Applying machine learning to primate bioacoustics: review and perspectives

Jules Cauzinille^{1, 2, 3}, Benoit Favre^{1, 3}, Ricard Marxer^{4, 3}, and Arnaud Rey^{2, 3}

¹Aix-Marseille University, CNRS, LIS, Marseille, France

²Aix-Marseille University, CNRS, CRPN, Marseille, France

³Aix-Marseille University, ILCB, Marseille, France

⁴Université de Toulon, Aix-Marseille University, CNRS, LIS, Toulon France

Author Note

Jules Cauzinille  <https://orcid.org/0000-0002-8604-1801>

Correspondence concerning this article should be addressed to Jules Cauzinille, ILCB, Aix Marseille University, 5 Av. Pasteur, 13100 Aix-en-Provence, France.

E-mail: jules.cauzinille@lis-lab.fr, Phone : +33679585636

Abstract

This paper provides a comprehensive review of the use of computational bioacoustics as well as signal and speech processing techniques in the analysis of primate vocal communication. We explore the potential implications of machine learning and deep learning methods, from the use of simple supervised algorithms to more recent self-supervised models, for processing and analyzing large datasets obtained within the emergence of passive acoustic monitoring approaches. In addition, we discuss the importance of automated primate vocalization analysis in tackling essential questions on animal communication and highlighting the role of comparative linguistics in bioacoustic research. We also examine the challenges associated with data collection and annotation and provide insights into potential solutions. Overall, this review paper runs through a set of common or innovative perspectives and applications of machine learning for primate vocal communication analysis and outlines opportunities for future research in this rapidly developing field.

Keywords: computational bioacoustics, primate vocal communication, passive acoustic monitoring, deep learning.

Applying machine learning to primate bioacoustics: review and perspectives

1 Introduction

Acoustic communication can be observed in many animal species and offers a diverse set of cues in the study of their behavior as well as a prolific insight for the monitoring of their activity. Primates are certainly no exception in this context, and the study of their vocalizations has been of great interest for the scientific community in recent years. Directly following their success in processing speech and audio, deep learning (DL) models were introduced to the field of computational bioacoustics through the ever growing availability of datasets allowed by technical advances in data recording, sharing and storage. This led to the now widely spread Passive Acoustic Monitoring (PAM) approach and, in turn, to an increasing need for efficient automated workflows in addition to hand-made annotations and analysis. This can be seen as a slight change of paradigm in the way primate vocal communication is treated and understood by researchers, which is also fairly recent and prone to evolve. In fact, deep learning methods developed for speech, audio or image processing as we currently understand them were seldom mentioned in computational bioacoustics reviews (Ganchev, 2017) until Stowell (2019, 2022). Fifteen years after the exploratory perspective paper from Zimmermann et al. (1995), one of the first studies mentioning the direct use of artificial neural networks applied to primate vocalizations was carried out by Pozzi et al. (2010). Besides that, simpler machine learning approaches developed for the processing of large unsegmented PAM datasets were not explored, to our knowledge, before the work by Kalan et al. (2015). The computational analysis of primate vocal communication systems is thus a young and rapidly growing field of study.

We hereby present a concise survey of the latest trends and approaches in machine learning applied to primate vocal communication research. In this perspective, we carefully

25 selected experiments, mostly published in the last three years, with additional earlier
26 papers that we deemed interesting for a contextualized discussion. After introducing the
27 key concepts and the general approach found in the literature, we describe three different
28 approaches and the type of results they can provide to better understand primate vocal
29 communication or to be used as monitoring tools. We then put forward different
30 perspectives following the development of high-performance weakly **supervised** acoustic
31 models and their potential use in primate communication research. Finally, we discuss data
32 availability and ongoing efforts in collecting and sharing new exploitable datasets.

33 **2 General considerations and methods**

34 **2.1 Passive Acoustic Monitoring**

35 Directly observing primates in the scope of studying their communicative behavior
36 in the wild can prove to be challenging depending on their species and the natural
37 conditions encountered in their habitat. An essential limiting factor is their common
38 tendency to flee on contact and the fact that human presence may affect their natural
39 behavior upon direct observation (Crofoot et al., 2010). One solution to this problem is to
40 focus on habituated or captive animals, but recent advances in technology also allowed
41 researchers to resort to more passive behavioral data collection methods such as camera
42 traps, drone technology or Passive Acoustic Monitoring (PAM). The emergence of
43 high-storage and energy-efficient recording hardware and the rapid development of machine
44 learning software is now turning PAM into a promising scientific tool to indirectly monitor
45 either wild or captive animals. This approach can be summarized as follows: one or
46 several, usually synchronized, microphones are placed at specific locations to record
47 soundscapes over large spatiotemporal scales. This can be applied to any animal species
48 displaying acoustic signals (Sugai et al., 2018) from insects such as mosquitoes (Kiskin
49 et al., 2021) to birds (Pérez-Granados & Traba, 2021) and mammals including cetaceans

50 (Zimmer, 2011), deers (Enari et al., 2017) or primates (Do Nascimento et al., 2021).

51 A drawback of PAM experimental setups is that the collected data is inherently
52 unimodal. However, it can be coupled with visual data from cameras, or an array of
53 additional information (time of the day, meteorological conditions, location or expert
54 annotations). PAM may be restricted in terms of modality, but it also presents some
55 advantages. Regarding primate monitoring activities, for instance, it has been shown that
56 PAM data is more valuable for the detection of primates than visual recordings from
57 camera traps, as shown by Enari et al. (2019) for Japanese macaques (*Macaca fuscata*) and
58 by Crunchant et al. (2020) for chimpanzees (*Pan troglodytes*). The approach may also
59 allow researchers to rely on the extensive work, methods and software technologies
60 developed for signal and speech processing. This makes PAM data a reliable source of
61 answers for an array of ecological questions (Ross et al., 2023). In the perspective of
62 primatology, computational bioacoustics can thus lead to significant discoveries regarding
63 primate communication and vocal behavior. It allows researchers to process and filter large
64 collections of sounds by relying on so-called machine learning methods which help them
65 automatically analyze and interpret acoustic signals. Applied to primate vocalizations,
66 these methods show impressive results in an array of essential tasks such as the denoising
67 of recordings, the selection and segmentation of said recordings to extract meaningful or
68 primate-only vocalizations from lengthy recordings, the detection and classification of
69 species, individuals or specific types of calls, etc. (see Figure 1). We hereby present the
70 different trends, approaches and benefits related to the application of machine learning to
71 PAM datasets. To get a deeper understanding at the functioning of machine and deep
72 learning algorithms, one may refer to specialized reviews such as Pichler and Hartig (2023).

73 **2.2 Machine learning for bioacoustics**

74 Machine learning is the implementation of artificial intelligence through algorithms
75 and computer models trained to autonomously make predictions from data. This process

76 may be summarized as a search for parameters in mathematical functions that minimize
77 the difference between predicted outcomes and actual (human) observations, enabling
78 systems to generalize and make accurate decisions on new, unseen data. Deep learning is a
79 subset of machine learning involving the use of artificial neural networks characterized by
80 their large number of parameters interacting with data in hierarchical, complex and
81 unforeseeable ways. For a more comprehensive exploration of machine learning, readers are
82 encouraged to refer to James et al. (2023), which provides both introductory and in-depth
83 coverage of the topic. Deep learning (DL) is thus particularly well-suited to handle large
84 datasets containing complex patterns and unstructured information, such as those
85 encountered in PAM datasets. For bioacoustics, machine and deep learning can be seen as
86 a way to automatize, reproduce or enhance human annotations with computer models.
87 These are built through the selection and development of algorithms adapted to the task at
88 hand, and always rely on annotated data for their training and evaluation. Training a
89 machine learning model consists in optimizing a system to reproduce a labeling process
90 through trial and error by maximizing correct predictions and minimizing mistakes. Its
91 performance is thus heavily dependent on the quality of the training data and on its ability
92 to generalize the labeling procedure to previously unseen data. Experiments involving the
93 use of machine and DL models trained on acoustic recordings of primates often follow
94 standardized workflows. These, in turn, resemble methods developed for the study of other
95 species and can show strong parallels with computational linguistics and speech processing
96 research. As stated by Stowell (2022), “classification is indeed the main use of deep
97 learning seen in computational bioacoustics.” Although the task of classifying sound can be
98 divided in various sub-tasks ranging from segmentation to labeling (see Section 3), it
99 usually implies training a computational model on a given set of data for a specific task
100 and evaluating its ability to perform the task when presented a different subset of the data.
101 This common approach, which can be found in most of the experiments we hereby review,
102 unfolds as follows (see Figure 2 for a visual description):

- 103 1. Acoustic data is collected from either captive or wild settings. This implies technical
104 subtleties in terms of location of the recordings, number of acoustic sensors deployed,
105 specific recording configurations such as sampling rate or frequency range, and the
106 amount of data which can be obtained.
- 107 2. The data is then turned into processable inputs. It may be segmented into short
108 clips, specifically tailored for the use of some **DL** models. It may also be transformed
109 through feature extraction to facilitate its processing by the computer models (see
110 Section 2.2).
- 111 3. Depending on the task at hand, annotation is needed to provide target labels which
112 will be learned and predicted by the model. These can include several classes such as
113 species, individual identities or call types, but also binary labels such as presence and
114 absence of vocalizations, as well as dimensional labels.
- 115 4. The model architecture mostly depends on the task it needs to perform but also on
116 the type of data it will be trained on. For instance, acoustic data in the form of
117 spectrograms is efficiently processed by Convolutional Neural Networks (**CNN**). A
118 **CNN** is an artificial neural network constructed as a stack of layers (the so-called
119 *convolutional filters*) which efficiently extract meaningful information from audio or
120 image inputs by recognizing patterns in the data. **CNNs** are by far the most popular
121 choice in computational bioacoustics. With the advent of new **DL** models developed
122 in the scope of processing longer sequences of speech or audio, some experiments now
123 rely on state-of-the-art architectures such as Recurrent Neural Networks (RNN), or
124 the popular **transformer** networks (Lin et al., 2022) which leverage attention
125 mechanisms and further exploit longer and variable-rate relations in the data. Much
126 simpler machine learning models such as Multi-Layer Perceptrons (MLP), Hidden
127 Markov Models (HMM), Support Vector Machines (**SVM**) or **clustering** algorithms
128 (see Pichler and Hartig (2023) for a typology of existing algorithms) can also show

129 interesting results. These simpler models may also be used as baselines, a voluntarily
130 simple reference model used for comparison purposes.

131 5. The dataset is divided into *train*, *development* and *test* sets. The train set will be
132 used during training to present examples of target labels associated with input data
133 to the model, from which it will learn to extract cues and informative features. The
134 development set is also used during training to make design choices that cannot be
135 optimized using the machine learning method. It provides information in order to
136 select from different models or tune functional aspects of the learning algorithm itself
137 (often referred to as hyper-parameters). Finally, the test set, unseen by the model
138 during training, will be used at the inference step to evaluate said performances. It is
139 usually taken from a separate pool of data (different microphone, location or
140 vocalizing individual) to ensure true generalisability of the model and avoid biased
141 evaluations.

142 6. The evaluation of the model requires the selection of appropriate metrics depending
143 on the task at hand. These are chosen to be as informative as possible in the context
144 of the experiment and must reflect the prediction performance but also potential
145 flaws in terms of false positives or false negatives. It must be chosen according to the
146 way labels are balanced in the dataset as well as its size. To evaluate the automatic
147 segmentation of acoustic data over time, for instance, authors will often rely on
148 accuracy (the number of correctly predicted segments divided by the total number of
149 segments). They may also use the **F1-score** to account for the balance between false
150 negatives and false positives, which is ignored by the accuracy metric.

151 A great majority of computational bioacoustics experiments rely on spectral
152 representations of sound prior to their automatic processing by machine learning models or
153 their manual annotation. These spectral representations, relying on the short-time Fourier
154 transform, encode the temporal evolution of acoustic energy over a range of frequencies in

155 a sound signal. Although several types of spectral transformations may be found in the
156 literature, some can be turned into images to ease the manual analysis of acoustic material
157 (as in the spectrogram from Figure 4). This type of sound representation is often used to
158 speed up the manual annotation of bioacoustic data by human annotators. It can also be
159 directly employed as the input of vision-based DL models which can reach high
160 performances on a variety of tasks by processing sound as images. No particular
161 representation of sound has been proven to work best across all species and tasks, and the
162 use of a given method must be carefully justified because of its important implications on
163 the performances of a computational model. Within spectral representations, a first
164 distinction can be made between linear and logarithmic spectrograms, the latter being
165 designed to mimic the way human ears process sound by emphasizing discriminability of
166 lower frequencies and de-emphasizing it in higher ones. In this perspective, bioacoustic
167 researchers often rely on mel-spectrograms, which tend to show promising results when
168 used as features for animal vocalizations processing. Nonetheless, and despite their
169 popularity, spectral representations are not always a preferred solution and other acoustic
170 representation methods exist. As can be seen in Kiskin et al. (2020), wavelets can also
171 show great benefits compared to the short-time Fourier transform, due to their ability to
172 capture both fine details and broad trends in acoustic data. Although a less conventional
173 solution, the authors show the advantages of wavelets when facing weak and noisy signals
174 (such as mosquito sounds) and their ability to perform better across datasets compared to
175 spectral solutions (with bird species classification).

176 With the advent of more powerful computational technologies, the use of the raw
177 waveform also presents an array of advantages compared to spectral-based representations.
178 Although directly using a waveform as input usually requires larger datasets and more
179 computing power to efficiently train DL models, the approach allows researchers to bypass
180 yet another manual pre-processing step. This also means a model can learn to extract any
181 informative cues without the risk of losing information through spectral transformations of

182 its training dataset. Additionally, training classifiers directly on the waveform can greatly
183 simplify a classification pipeline, going from the acoustic data to the bioacoustic predictions
184 in a straightforward manner, usually referred to as an “end-to-end” approach. Although it
185 was not used on primate vocalizations to our knowledge, end-to-end approaches are gaining
186 popularity in bioacoustics, an sound processing in general, with successful applications such
187 as in Bravo Sanchez et al. (2021) on birds and Xie, Hu, et al. (2021) on frogs.

188 Finally, as we will further discuss in Section 4, acoustic representations may be
189 extracted by pre-trained upstream models (see Figure 3). This is referred to as *pre-trained*
190 *representation learning*, an increasingly popular solution in recent speech and audio
191 processing research. The approach relies on the pre-training of large generalistic foundation
192 models to enhance the performance of smaller task-specific ones. Once the pre-trained
193 representations are learned and extracted, they may be used as traditional features
194 containing useful information for an array of tasks.

195 A second optional step, directly preceding feature extraction, is *signal enhancement*,
196 or *denoising*, a process which consists in filtering out non-informational signal from raw
197 data. This signal processing method is quite popular in bioacoustics where clear recording
198 conditions are rarely encountered (Xie, Colonna, & Zhang, 2021). The amount of noise
199 which can affect the performances of computational models in processing primate
200 vocalizations greatly depends on the recording location or microphone sensitivity, and may
201 stem from an array of acoustic sources, from anthropogenic noise (vehicles, speech...) and
202 natural soundscapes (rain, wind...) to other species or conspecifics vocalizations. Although
203 denoising can be an essential tool, it is not always beneficial, as it may deprive the signal
204 from essential information which could potentially be extracted by a computational model.
205 The general approach consists in using simple fixed signal processing tools to perform the
206 so-called signal enhancement directly on a spectrogram. Such tools have been extensively
207 studied and engineered and are widely available through public softwares such as
208 **noisereduce**, a spectral gating algorithm developed by Sainburg et al. (2020). A more

209 refined option from which computational bioacoustics could greatly benefit is DL-based
210 denoising, a popular area of research in speech processing (Germain et al., 2019).

211 **3 Tasks and applications**

212 We hereby describe three main categories of tasks which can be tackled through the
213 use of machine learning for primate vocalization analysis. For an overview of these
214 categories and their applications in primate bioacoustics, see Figure 1.

215 **3.1 Detection and segmentation**

216 The most practical application of machine learning in bioacoustics, when facing
217 large unlabeled recordings of natural soundscapes, is the detection of animal vocalizations
218 among ambient noise. As we previously mentioned when introducing the PAM approach,
219 the ever-increasing storage and battery life capabilities of microphones may result in
220 recordings lasting several hours or days. Primates, however, are not constantly vocalizing
221 and their calls usually span specific segments of time which need to be identified and
222 extracted for their subsequent analysis. The manual segmentation of recordings (i.e.
223 annotating start and end times of primate calls among a continuous audio clip) is an
224 essential step in processing PAM data. An efficient way to carry out this segmentation is to
225 directly inspect spectrograms of the recording in a specialized software such as PRAAT¹ or
226 Raven Pro². Although it results in precise annotations, manual segmentation may prove to
227 be quite time-consuming depending on the length of the audio files and the nature of the
228 recorded vocalizations in terms of frequency ranges, unit-rates, distances to the microphone
229 and amounts of background noise. Automatic detection and segmentation were proposed as

¹ Boersma, Paul & Weenink, David (2023). Praat: doing phonetics by computer [Computer program]. Version 6.3.17, retrieved 10 September 2023 from <http://www.praat.org/>

² K. Lisa Yang Center for Conservation Bioacoustics at the Cornell Lab of Ornithology. (2023). Raven Pro: Interactive Sound Analysis Software (Version 1.6.4) [Computer software]. Ithaca, NY: The Cornell Lab of Ornithology. Available from <https://ravensoundsoftware.com/>

230 an answer to this issue. It may take at least three different forms:

- 231 • binary detection: the machine learning algorithm is given a segment of audio as input
232 and outputs the probability of this segment containing a call. This may be referred
233 to as “occupancy” or “presence” prediction.
- 234 • time-wise segmentation: the task can still be developed as a binary one but it results
235 in a more fine-grained annotation of the input file with start and end time-codes of
236 each call. This is often solved by making an occupancy prediction in short windows
237 (10 or 50 ms) and merging consecutive positive decisions into a single segment.
- 238 • time and frequency-wise detection: directly inspired by image object detection, this
239 task usually implies the use of spectrograms. The model is constructed as an object
240 detector and outputs time and frequency boundaries of the target call as in Figure 4.
241 To our knowledge, this approach was never explored for primate vocalizations and is
242 scarcely applied to other species as well, although it could be used to identify various
243 simultaneously vocalizing species in single segments.

244 Automatic segmentation is undoubtedly one of the most studied aspects of
245 automatic detection in bioacoustics. The following examples exclusively focus on detecting
246 gibbons, but similar directions are being explored on other primate species (Anders et al.,
247 2021; Bonafos et al., 2023). Recently, approaches involving the use of deep learning are
248 being predominantly adopted in automatic audio recognition and tend to replace the use of
249 hand-crafted features and simpler machine learning algorithms. This scientific trend is
250 widely adopted across bioacoustics, specifically through the use of **CNN**-based solutions
251 from spectral inputs, as can be seen in the evolution of the DCASE challenge over the
252 years (Mesaros et al., 2017, 2019). Although rarely explicitly compared with more simple
253 statistical approaches, **CNNs** present an array of advantages in these tasks. Their primary
254 benefit lies in their ability to generalize predictions across varying recording conditions.
255 They also allow efficiently tackling tasks with noisy and unbalanced annotated datasets of

256 limited size (Anders et al., 2021), as is often the case with PAM recordings of primates in
257 the wild. As we will see, they may also be coupled with so-called RNNs to account for the
258 sequentiality of primates' vocalisations. This makes them especially effective at detecting
259 primates with temporally dynamic calls, as is the case with gibbons. Finally, CNNs being a
260 very popular option in the deep learning community, extensive research, publicly available
261 resources and off-the-shelf solutions can be accessed with little expertise to develop fast
262 and efficient models. All these advantages made CNNs a go-to solution in bioacoustic
263 detection (Stowell, 2022), progressively replacing the use of simpler statistical algorithms.
264 As we will see throughout this review, machine learning for primate bioacoustics follows a
265 similar trend. Nonetheless, the efficiency and advantages of deep learning solutions come
266 with some drawbacks, specifically in the interpretability of a model's predictions as well as
267 in their need for higher computational resources and larger datasets (see Figure 5). In
268 recent bioacoustics papers, the specific reasons for choosing large deep learning models over
269 lightweight statistical solutions are rarely explicated. However, they are often implicitly
270 shown through comparison with simpler baseline performances.

271 An illustrative example in this perspective is a model developed by Dufourq et al.
272 (2021) applied to the highly endangered Hainan Gibbons (*Nomascus hainanus*) from the
273 *Bawangling National Nature Reserve*. The proposed model is a common one in bioacoustic
274 event detection as it relies on the popular CNN architecture. In this case, the model is
275 designed to differentiate between two classes of sounds : *non-primate background noise* and
276 *primate vocalizations*. It is trained on mel-spectrogram representations of short sound
277 segments which were previously labeled as such. The authors additionally resort to data
278 augmentation, which consists in increasing the size and diversity of a dataset by slightly
279 modifying it to create additional synthetic data. This method permits rendering the
280 models robust to certain transformations that we know should not affect the system's
281 prediction. Here, each segment is shifted in time to double the size of the initial dataset.
282 The authors evaluate two types of architectures, namely a one-dimensional CNN leveraging

283 temporal patterns, and a two-dimensional one which captures frequency as well as time
284 from the input spectrograms. One last step consists in post-processing the model
285 predictions by removing unrealistic detections (such as isolated or very short calls). The
286 2-D CNN with data-augmentation paired with the post-processing step achieves 99.37%
287 accuracy (compared to 97.60% without post-processing and 92.32% without
288 data-augmentation as well)/. This shows that the CNN approach can be highly efficient in
289 facilitating the segmentation of large PAM recordings of Gibbons, with an eight hours long
290 test recording taking six minutes on average to be processed by the model. In addition to
291 its high performance, this automated procedure also provides exhaustive quantitative
292 information about Hainan gibbon’s vocal behavior, including their preferred vocalization
293 times, the amount of calls they produce in a day and their geographical distribution over
294 the study site. By coupling the detections with additional metadata such as meteorological
295 information and environmental parameters, this approach could lead to many more
296 interesting observations.

297 As is common in machine learning experiments, the work by Dufourq et al. (2021)
298 constitutes a baseline performance which was promptly improved by Ruan et al. (2022)
299 with a slightly different architecture. This “baseline” can be considered as a performance
300 milestone aimed at being improved upon and was included in one of the only bioacoustics
301 benchmarks available to date: *the Benchmark of Animal Sounds* by Hagiwara et al. (2023).
302 Ruan et al. (2022) approach relies on deep learning solutions, namely Residual Networks (a
303 former state-of-the-art model, known for its high performances in image classification),
304 SpecAugment (the random masking and warping of portions of the input spectrogram to
305 improve the generalization ability of the model) and label smoothing (a method with
306 similar results based on the addition of noise to the label distribution during training).
307 These methods were originally developed for speech and image classification and allow the
308 authors to propose *BPDnet*, a model with high performance on Hainan Gibbon’s presence
309 detection carried out on the same dataset as used by Dufourq et al. (2021) without the

310 need for manual intervention and post-processing. When compared with the baseline
311 experimental setup, this new model improves the F1-score by 0.16 without post-processing
312 and by 0.09 with post-processing.

313 The approach selected by Dufourq et al. (2021) corresponds to the most
314 recommended and wide-spread approach for the segmentation of bioacoustic data. It relies
315 on years of research into the **CNN** architecture and tackles major limitations through
316 post-processing and data-augmentation in addition to the use of traditional spectral
317 representations. Comparatively, Y. Wang et al. (2022) work on the same dataset with a
318 rather innovative perspective relying on more complex state-of-the-art models. They
319 implemented two solutions: a **CNN** stacked with a Hidden Markov Model (HMM) and a
320 Convolutional Recurrent Neural Network (CRNN). The HMM architecture works as a
321 post-processing step and allows for a correction of the **CNN** decision from contextual
322 information of neighboring segments. The CRNN, in contrast, outputs decisions from
323 sequential information rather than from a fixed segment. It still relies on a **CNN** for feature
324 extraction which is subsequently passed to Gated Recurrent Units accounting for the
325 temporal information. These particular types of models have been shown to be more
326 effective at modeling long-term dependencies in sequential data (like sound), while also
327 being computationally more efficient than traditional RNNs. Here, the choice of this
328 architecture is motivated by the sequential nature of gibbon vocalizations which are known
329 to produce varying sequences of notes combined into phrases. The authors show that the
330 CRNN architecture improves performance compared to the CNN-HMM and that it is
331 resistant to low Sound to Noise Ratio (SNR), a metric used to account for the amount of
332 background noise in a given audio segment. Although the results are not comparable with
333 previously mentioned experiments because of distinct evaluation metrics and dataset
334 processing, CRNNs are a viable option for automatically processing large PAM recordings
335 of gibbons and their resistance to noise can be seen as a major asset in the bioacoustic
336 context.

337 Following both previously described frameworks, Tzirakis et al. (2020) conducted
338 somewhat similar experiments on large recordings of Müller gibbons (*Hylobates muelleri*)
339 from Malaysian Borneo. The authors built a publicly available model consisting of a 2-D
340 **CNN** followed by a RNN which also capture longer-term temporal dynamics of the input
341 signal. To validate their approach, the proposed architecture is compared to two publicly
342 available toolkits for audio representation and analysis: End2You (a simpler yet similar
343 model comprised of **CNNs** followed by Gated Recurrent Units) and openXBOW (based on
344 the bag of words approach from computational linguistics and several machine learning
345 models, including the Random Forest algorithm). The dataset was processed as a
346 collection of positive (gibbon’s presence) and negative (background noise) audio clips. The
347 results indicate that the author’s model reaches an Unweighted Average Recall of 93.3% on
348 the test set compared to 84.8% for the best openXBOW model.

349 Although we only discussed three experiments on a single taxon, we can see that
350 state-of-the-art segmentation of primate vocalization datasets mostly resort to closely
351 related approaches and tend to show high performance if sufficient annotated data is
352 provided. **CNN** solutions can be found successful for an array of other primates and
353 animals: Stowell (2022) surveyed 83 such experiments for a variety of tasks including
354 segmentation; and RNNs seem to be a viable option in further improving their
355 performance. This approach also shows the benefit of potentially relying on off-the-shelf
356 models made available by researchers in other fields. The limitation here is the availability
357 of annotated data itself and the potential complexity of usage of such publicly available
358 models for non-specialist practitioners.

359 **3.2 Identification and density estimation**

360 We have seen that segmentation and detection tasks relying on state-of-the-art deep
361 learning architectures can be very effective for primate species shown to communicate
362 vocally like gibbons, even in their noisy natural environment. Although the obtained

363 results and performances can be used in quantitative analysis of their vocal behavior and
364 are essential for qualitative studies of the segmented calls, they lack in conservation value.
365 In this perspective, identifying vocalizing individuals to estimate primate population
366 density from their vocalizations and studying individual vocal signatures are both complex
367 but potentially significant tasks. Few studies can be found exploring automatic density
368 estimation relying on **DL** models, although human based detection seems to be a good
369 option for this task. Using humans as acoustic detectors has indeed been proven successful
370 in estimating the density of yellow-cheeked gibbons (*Nomascus gabriellae*) in Cambodia by
371 Kidney et al. (2016). This means that acoustic data contains enough information to
372 develop similar experiments with computational models.

373 In a similar perspective, several studies were conducted in developing automatic
374 classifiers for caller identification. These often rely on recordings in captivity, during
375 mark-recapture events or by focal recordings of individuals to allow for an easier extraction
376 of the caller identity and the constitution of an annotated dataset. In fact, obtaining
377 identity annotations from PAM data is a difficult process due to low control possibilities
378 over the recordings. In captivity however, some solutions were explored by Bayestehtashk
379 et al. (2014) who recorded groups of captive rhesus macaques (*Macaca mulatta*) using
380 individual collars. The authors provide an interesting semi-automatic pipeline for the
381 constitution of an ID-labeled corpus. To process the obtained data from multiple collars,
382 they construct a segmentation model based on manually designed acoustic features often
383 used in music and speech from the OpenSmile toolkit (Eyben et al., 2010). Their best
384 performing model is a Support Vector Machine (**SVM**). **SVMs** are machine learning
385 algorithms designed to separate samples of different class labels in the feature space. They
386 are light machine learning models, easy to train, and show good performances on simple
387 tasks. In this case, the **SVM** architecture reaches 88.9% accuracy on a manually segmented
388 subset of the data. Once segmented, the obtained set of vocalizations is to be matched to
389 the recordings from each collar in order to detect the one which most likely emitted a

390 specific vocalization. The authors rely on Dynamic Time Warping (Müller, 2007) to
391 compute the acoustic similarity of each segment and show good results compared to
392 manually aligned data. The obtained dataset can then be used to train classifiers for the
393 automatic identification of individual monkeys.

394 In an opposite approach, highly territorial primates like Northern grey gibbons
395 (*Hylobitae funereus*) can provide interesting datasets to perform this type of task on wild
396 PAM recordings. By placing microphones inside individual group territories, Clink et al.
397 (2017) successfully identified the acoustic parameters contributing to individuality in
398 female's great calls. The authors manually extracted acoustic features from their PAM
399 dataset and computed a Mahalanobis acoustic distance measure between pairs of
400 vocalizations. They were able to discriminate between pairs of 33 females with a 95.7%
401 accuracy using linear Discriminant Function Analysis, a method consisting in searching for
402 linear combinations of the extracted features to separate the different individuals.

403 Machine learning algorithms such as Discriminant Function Analysis do not rely on
404 deep learning, contrary to the **CNNs** and RNNs discussed in the previous section. For
405 vocal signature classification, they seem to be a preferred approach with the advantage of
406 demanding less computational power and data all the while yielding competitive results.
407 We want to stress that resorting to deep learning solutions in bioacoustics is not always a
408 preferred approach, especially when facing scarce annotated data, as is the case for
409 identification of primate voice prints.

410 With a similar dataset of Northern grey gibbons and an approach involving Support
411 Vector Machines and Mel-Frequency Cepstral Coefficients (**MFCC**), Lakdari et al. (2024)
412 reached high performance on classifying female great calls by recognizing their emitter
413 from a pool of 12 individuals. They further examined the performance of their approach by
414 recording the calls in playback at varying distances to account for its resilience on low
415 sound to noise ratios. They find that **MFCCs** are outperforming other feature extraction
416 methods, namely acoustic indices or pre-trained **DL** models, when calls are recorded at

417 larger distances.

418 Another example of such solutions can be found in a study by Fedurek et al. (2016).
419 The authors examine chimpanzees (*Pan troglodytes schweinfurthii*) pant-hoots and attempt
420 at identifying the type of information acoustically embedded in them. Their experiment
421 also relies on Support Vector Machines trained on MFCCs. They find that all four phases
422 of the pant-hoot (introduction, build-up, climax and let-down) are associated with a
423 variety of information, including individual identity, which is more specifically encoded in
424 the introduction and climax. Despite these promising results, there seems to be a lack in
425 the implementation of state-of-the-art DL models for vocal signature classification in
426 primates. We have seen that this shortfall can be explained by difficulties in annotating
427 data accordingly. It may also be explained by the more consistent amount of work put into
428 identifying primates from visual data with face recognition (Guo et al., 2020; Schofield
429 et al., 2019). Yet, few experiments relying on complex and innovative DL architectures
430 from sound show promising results. We have mentioned Lakdari et al. (2024) who
431 compared MFCCs feature extraction with deep learning models pre-trained on birds,
432 speech or general sound. Leroux et al. (2021) also introduce transfer learning from DL
433 models pre-trained on speech for chimpanzees voice print recognition. Both approaches will
434 be further discussed in Section 4.

435 Nonetheless, a parallel task involving multi-label classification and voice prints with
436 interesting machine learning solutions is primate species identification. As we have seen,
437 most bioacoustic studies carried out on primates focus on single datasets from one species
438 of interest. However, wild environments may host various cohabiting species which often
439 end up overlapping in single PAM recordings. A first interesting study in this regard was
440 conducted by Mielke and Zuberbühler (2013) with a combination of classification tasks for
441 species, call type and caller identification. It relies on a MLP trained on a dataset of
442 Stuhlman’s blue monkey (*Cercopithecus mitis stuhlmanni*) vocalizations. This particular
443 species allows for identity labeling because each group hosts a single male which also

444 happens to be the only producer of “pyow” calls. Additionally, other species’ calls found in
445 the same environment were added for the species discrimination task (olive baboons, *Papio*
446 *anubis*; redtail monkeys, *Cercopithecus ascanius schmidti*, and guereza colobus monkeys,
447 *Colobus guereza occidentalis*). After extracting MFCCs and training various MLPs with
448 distinct hyperparameters, male identity classification resulted in 73% accuracy in average
449 and species recognition resulted in 96% accuracy for the four classes. Despite the promising
450 results, we must point out that substantial manual work had to be allocated for the
451 pre-processing, segmentation and identification of the calls prior to the automated
452 classification. Other early experiments by Kalan et al. (2015) and Heinicke et al. (2015)
453 also showed interesting approaches with simpler algorithms including SVMs and Gaussian
454 Mixture Models. Both papers focus on the identification of chimpanzees (*Pan troglodytes*
455 *verus*), diana monkeys (*Cercopithecus diana*), red colobus (*Procolobus badius*) and king
456 colobus (*Colobus polykomos*). Both SVMs and Gaussian Mixture Models were trained on
457 MFCCs and other spectral information extracted PAM recordings. The algorithms show
458 relatively low results with less than 5% of detected segments being true-positives for the
459 best model.

460 A more recent approach involving Kernel Extreme Learning Machine was adopted
461 by Zwerts et al. (2021). This particular type of model is a supervised learning algorithm
462 using a kernel function to map input data into a high-dimensional space and allows for the
463 learning of complex and non-linear relationships between input features and output targets.
464 It was trained on MFCC representations of vocalizations from captive chimpanzees (*Pan*
465 *troglodytes*), mandrills (*Mandrillus sphinx*), red-capped mangabeys (*Cercocebus torquatus*)
466 and a mixed group of guenons (*Cercopithecus sp.*), with an additional class of background
467 noise. The performances of the model are above chance (25% for the four species) with
468 76.7% accuracy in a four class setup and 69.7% accuracy with the addition of the noise
469 class, but stay relatively low compared to more recent approaches.

470 As we have mentioned with other experiments, the publication of this new dataset

471 and baseline model may be seen as a benchmark which was promptly integrated as part of
472 the INTERSPEECH 2021 Computational Paralinguistics Challenge (Schuller et al., 2021).
473 In an attempt to tackle the species identification problem with state-of-the-art
474 architectures (similar to what we have seen in Section 3.1), Pellegrini (2021) compared
475 several DL models including CNNs, MobileNet and ResNets. They also revolved to data
476 augmentation methods like SpecAugment and MixUp (another technique for data
477 augmentation relying on the blending of pairs of training examples). The main difference
478 between each model lies in the definition of their convolutional blocks. The first two
479 models are standard CNN architectures with 6 and 10 layers respectively. The
480 MobileNetV1 model relies on *depthwise separable convolutions* to reduce computational
481 costs and gain efficiency, usually in the scope of being used in mobile and embedded
482 devices. Finally, two CNN ResNet models make use of *residual connections*, allowing for a
483 deeper network architecture to be trained without suffering from the vanishing gradient
484 problem which affects models with many stacked layers, such as standard CNNs. The
485 results show good improvement compared to Zwerts et al. (2021) baseline with an
486 unweighted average recall of 92.5% achieved by the 10 layer CNN, closely followed by the
487 large ResNet model. In this case, the 10 layer CNN is preferable to ResNet as it achieves
488 better performance with a much smaller model size. The authors also note that the most
489 common confusion made by their models regards the background noise class versus the
490 primate vocalizations one. This confirms the importance of ongoing efforts in resolving
491 “low level” tasks, such as the identification of primate vocalizations among natural noise.

492 **3.3 Vocal repertoires and clustering**

493 Despite this, “high level” tasks can be found in computational bioacoustics
494 literature, with many relying on machine and deep learning-based solutions. One such task
495 with great scientific value for primatology is the discovery or the classification of call types
496 (i.e., the categories of calls produced by a species). This task may be carried out through

497 different approaches, including **supervised** classification (each class corresponding to a
498 predefined call type) and **unsupervised clustering** (the grouping of similar acoustic objects
499 into undefined call type categories). Each of these has been explored for various primate
500 species using an array of machine learning algorithms. The relative success of one
501 approach, especially in **unsupervised** contexts, is often seen as a form of validation of
502 predefined expert descriptions of a species vocal repertoire. Call type classification thus
503 serves the purpose of automatically processing large amounts of data while potentially
504 questioning human bias in the definition of vocal repertoires and fostering replicable results
505 across studies.

506 We hereby refer to “**unsupervised**” approaches to account for all experiments
507 involving the training of a model with little to no reliance on expert labels and
508 annotations. In the context of call type discovery, for instance, this means that a **clustering**
509 model is trained on unlabeled acoustic samples and should discover its own typography of
510 calls in order to categorize them. The related “semi-supervised” approach is one where a
511 limited amount of information is given to the model prior to **clustering**, such as the number
512 of categories to be discovered. Evaluating the results of such **clustering** approaches is a
513 highly debated topic in computer science, as no single solution can objectively quantify the
514 validity of a set of clusters compared to another. Von Luxburg et al. (2012) review the
515 different issues related to **clustering** and its evaluation. Although not centered around
516 bioacoustics, the paper draws inherent limits specific to the idea of automatic clustering:

- 517 • Evaluating **clustering** results is not problem-independent and must be related to the
518 end-user intent and their scientific scope.
- 519 • As is always the case with high-dimensional data, selecting the features on which
520 **clustering** will be carried out can greatly modify the output typology.
- 521 • A given **clustering** output can be found to be qualitatively reasonable for a specific
522 research question but may be meaningless to others.

- 523 • Computing internal **clustering** quality scores (centrality of the clusters, likelihood
524 scores, silhouette values, etc.) can be informative at the algorithm level but does not
525 provide an objective and domain-specific evaluation of the results.
- 526 • Comparing the results of **unsupervised clustering** with predefined categories of calls
527 should not be seen as undisputable proof of the validity of said “expert” categories,
528 as both may be biased in different ways.

529 **Clustering** primate call-types should thus be seen as an exploratory approach, and
530 experiments using it as a confirmatory solution to predefined human vocal repertoires
531 should be taken with care.

532 To our knowledge, the first paper mentioning the use of Artificial Neural Networks
533 for primate vocalization analysis, over and above the preliminary work of Zimmermann
534 et al. (1995), is, in fact, aimed at call type classification on black lemurs (*Eulemur macaco*).
535 Pozzi et al. (2010) compare the performances of **supervised** neural networks, statistical
536 models and **clustering** algorithms in recognizing a set of predefined call types. They show
537 that basic artificial neural networks trained to classify seven call types from which spectral
538 (F0 and formants) and temporal (duration) acoustic features were extracted, can show high
539 performances with a general accuracy of approximately 94%. Statistical analysis with
540 Discriminant Function Analysis and K-means **clustering** showed slightly lower
541 performances, with large disparities in classification accuracies for some call-types,
542 potentially due to the unbalanced classes context. The authors thus give a first example of
543 some advantages presented by deep learning methods compared to statistical approaches.
544 They mention their ability to handle noisy recordings, to generalize human annotations to
545 unseen data, and the reusability of a model’s weights once it has been successfully trained.
546 The authors also mention a set of limitations that can still be found in such experiments.
547 These include the over-fitting problem where neural networks learn dataset-specific
548 information related to individuals or to their sex rather than universal cues generalizable to
549 the entire species. They also mention the problem of biases in the manual annotation of

550 datasets, which may greatly affect the evaluation of **clustering**.

551 As previously mentioned, statistical and machine learning algorithms that do not
552 involve neural networks can show promising results in the analysis of call types. Turesson
553 et al. (2016) investigated the use of seven different such models in addition to **DL** ones for
554 the categorization of common marmoset (*Callithrix jacchus*) calls. The automatic
555 identification of call types appears as an essential tool for marmosets, as they produce large
556 amounts of characteristically complex and overlapped vocalizations on which annotation is
557 rather tedious and time-consuming. The authors collected a dataset from captive monkeys'
558 recordings with approximately 30 examples for each of the 11 call types investigated. They
559 chose linear predictive coding as a feature extractor. This technique, traditionally used to
560 encode the timbre of human voice signals in speech compression for telephony, consists of
561 modeling the spectral envelope of sound samples as the weighted sum of previous samples
562 from a given acoustic sample. The extracted features are used as an input to train several
563 classifiers, namely an Optimum-Path Forest, an MLP, an **SVM**, a k-Nearest Neighbors
564 **clustering** algorithm, etc. In addition, various proportions of the training set were tested to
565 understand performance trade-offs relative to training size. The **SVM**, k-Nearest Neighbors
566 and Optimum-Path Forest were found to be the best performing algorithms in both the
567 smallest and largest training set sizes, with Optimum-Path Forest being parameter-free and
568 requiring less computational resources. This suggests that simple statistical algorithms can
569 show high performances when facing limited amounts of clean data, although we could
570 argue that larger training datasets would increase performances in general and might be in
571 favor of other more complex deep learning-based models.

572 The task of automatically discriminating between different types of calls may prove
573 useful in quickly processing large amounts of data but it can also be used to infer new
574 properties of primate communicative systems, especially when tackled with **clustering**. Erb
575 et al. (2023) adapted different models to the classification of Bornean orangutans (*Pongo*
576 *pygmaeus wurmbii*) pulse-types to investigate problematic elements in the specie's

577 predefined vocal repertoire. They collected a dataset of focal recordings from 23 individual
578 males. Comparing human annotations and the **unsupervised** predictions of an **SVM** as well
579 as soft and hard **clustering** algorithms, they showed that a set of six pulse-types gives
580 rather poor results in terms of inter-annotator agreement as well as automatic predictions
581 in this specific experimental setup. This negative result allows them to propose a new
582 repertoire comprised of only three pulse-types, which shows higher classification accuracy
583 and reproducibility. Finally, they highlight the importance of graded categories of signals,
584 in opposition to strictly separated call types, in the typology of orangutan call types. This
585 type of experiment shows how automatic **clustering** and classification, although not
586 sufficient to objectively refute a predefined vocal repertoire, can still be used as an
587 exploratory tool to identify its potential biases. Similarly, Wadewitz et al. (2015) question
588 the discreteness of chacma baboons' (*Papio ursinus*) call type categories by investigating
589 the results of a "fuzz" **clustering** algorithm. They argue that labeling primate vocal
590 repertoires as being either fully discrete or fully graded may be considered an
591 oversimplification. Hard **clustering** (found in K-means algorithms, for example) assigns
592 each call to a single cluster or call-type. Fuzzy **clustering** (such as the C-means algorithm),
593 allows separating different classes of calls in a gradual manner rather than a sharp one.
594 Each call is given a membership value, ranging from 0 to 1, assigning it to each cluster.
595 Intermediate membership values characterize calls ambiguously pertaining to multiple
596 clusters. The authors find that, although hard K-mean **clustering** shows good alignment
597 with predefined human labeled call types on chacma baboons, fuzzy **clustering** gives
598 additional information regarding the atypicality of some of the species' calls. Again, when
599 used as an exploratory tool, fuzzy **clustering** may question unforeseen biases in the
600 constitution of a species vocal repertoire.

601 We have seen that call-type classification and **clustering** can be carried out for
602 different reasons with different algorithms, but the outcome of call-type **clustering** also
603 greatly depends on the choice of acoustic features it is built on. A preferred approach is the

604 dimensionality reduction of spectral features, as in an experiment proposed by Sainburg
605 et al. (2020) relying on the Uniform Manifold Approximation and Projection algorithm
606 (UMAP), among others. This popular method has the benefits of being relatively efficient
607 in discovering significant sound features and can result in informative visual
608 representations of clusters. UMAP is one of the many dimensionality reduction algorithms
609 used to classify vocalizations by taking high-dimensional features (like deep, spectral, or
610 acoustic features) and mapping them to a lower-dimensional space while maintaining the
611 underlying distances between different sounds. By reducing the dimensionality of
612 spectrograms with said algorithm, the authors show the implications of automatic
613 **unsupervised clustering** in a variety of topics related to primate vocalization analysis,
614 including the discreteness of macaques vocal signatures (see Figure 6) or the apparent
615 continuity of gibbon (*Hylobates sp.*) syllables (see Figure 7).

616 This approach should both be extended to other primates and explored through the
617 use of different algorithms and input features. In fact, UMAP may struggle with capturing
618 the global structure of acoustic data, particularly when dealing with complex and highly
619 varied vocalizations, as is the case for primates. In addition, UMAP's performance may be
620 strongly affected by the presence of outliers and noise in the data (as is often the case with
621 PAM recordings), potentially leading to distorted representations and the absence of
622 interpretable results.

623 As an alternative, Best et al. (2023) were inspired by deep representation learning
624 and extended this methodological framework to an array of animal species, showing once
625 again the great flexibility of the approach in validating vocal repertoires and alleviating
626 their manual annotation. Contrary to the latter experiment, the features extracted prior to
627 **clustering** are derived from a **self-supervised CNN**-based auto-encoder trained to *encode*
628 informative components of a spectrogram through a bottleneck approach in order to
629 subsequently *decode* input signals with minimal loss of information. This method, inspired
630 by speech and image processing techniques, yields significant results in the **unsupervised**

631 **clustering** of call-types for a variety of taxa ranging from birds to marine mammals. The
632 authors showed the benefits of working with UMAP and **clustering** algorithms based on
633 deep representations of sound rather than spectral or handcrafted features. As was
634 previously mentioned, the evaluation of **clustering** results through their comparison with
635 pre-defined expert categories is rather exploratory and does not prove the objective validity
636 of a given algorithm or feature extraction method (none of said expert categories).
637 Nevertheless, **clustering** solutions can be compared in terms of their alignment with human
638 typologies to provide interesting insights on their ability to extract information deemed
639 important by expert labelers. The authors thus demonstrate the ability of models
640 pre-trained on non-bioacoustic datasets to extract features that encode sufficient
641 information for an efficient **unsupervised clustering** of call-types. Although their
642 Autoencoder architecture yields better results in most datasets, models like wav2vec
643 (Schneider et al., 2019) and OpenL3 (Cramer et al., 2019), a model trained on audio/video
644 correspondence from YouTube data, also show comparable agreement scores between found
645 clusters and expert labels. This innovative paradigm, i.e., using large **DL** models
646 pre-trained on non-bioacoustic data for bioacoustic tasks, seems to be an increasingly
647 popular one in a variety of experiments, although it stays quite seldom explored for
648 primate vocalizations. In the next section, we will discuss recent papers making use of this
649 approach and see the potential implications and perspectives it may offer for the study of
650 primate vocal communication.

651 **4 Transfer learning and promising approaches**

652 Fairly recently, the advent of so-called “Pre-Trained Models” (PTM) has
653 undoubtedly revolutionized the use of deep learning for text, image and speech processing.
654 For Natural Language Processing, PTMs such as BERT (Devlin et al., 2019) or GPT
655 (Radford et al., 2018) have become a milestone in the field of artificial intelligence with
656 large language models like ChatGPT showing impressive applications way beyond

657 computer science research by leveraging an ever-increasing access to high computational
658 power and large amounts of data.

659 These models often rely on a **self-supervised** learning pre-training step, consisting in
660 storing and extracting information from massive datasets which can then be reemployed in
661 a variety of downstream tasks with great performance benefits compared to more
662 traditional **supervised** approaches (X. Liu et al., 2023; Mohamed et al., 2022). In the
663 acoustic domain, **self-supervised** models have also shown impressive capabilities in
664 generalizing knowledge with performance gains across a wide range of domains. The
665 typical approach in this regard is to pre-train a model on large unannotated datasets
666 (which should be relatively close in nature to the target domain data) and to use the
667 learned representation for downstream tasks on smaller manually labeled datasets. This
668 process involves transfer learning, i.e., relying on the knowledge learned during the
669 pre-training task for a new, potentially different, downstream task (see Figure 3). This
670 approach was successfully carried out in a variety of domains including music or biomedical
671 signal processing (Banville et al., 2021; Wu et al., 2021).

672 Bioacoustic tasks and use-cases are no exception here. Researchers in animal vocal
673 communication progressively turned to this new paradigm in recent years by adapting
674 methods initially developed for speech and sound processing to the analysis of acoustic
675 data produced by animals. When it comes to primates, however, the success of
676 **self-supervised** and transfer-learning approaches is yet to be confirmed and widely adopted.
677 However, several such experiments can be found relying on a variety of parallel approaches,
678 each showing its own benefits. We hereby discuss three main solutions that arise from
679 using transfer learning for automatic primate vocalizations processing.

680 **4.1 Retraining**

681 In the field of machine learning, transfer learning refers to the *pre-training* of a
682 model on a given dataset or task and the development of a downstream model aimed at

683 performing a specific downstream task on a different (annotated) dataset. In bioacoustics,
684 however, a slightly different popular approach, partly relying on knowledge transfer, is the
685 *retraining* of a model initially developed for similar tasks but on a different species. This
686 retraining approach is not to be confused with the pre-training of a single foundation
687 model in which previously learned weights may be reused for several applications. For
688 primate vocalizations, a good example of knowledge transfer through retraining is the work
689 by Romero-Mujalli et al. (2021). Here, the authors show the benefits of retraining the
690 ultrasonic vocalization detector model *DeepSqueak* (originally developed for rodents) on a
691 gray mouse lemur (*Microcebus murinus*) vocalizations dataset. Both taxa, rodents and
692 gray mouse lemurs, show relative similarity in the frequency range and general spectral
693 dynamics of their vocal communication. This similarity, the simplicity of the retraining
694 approach, the efficiency of *DeepSqueak*'s Faster-RCNN and the user-friendly environment
695 of the software allows yielding competitive results by training *DeepSqueak* on lemur's
696 vocalization for their segmentation, classification and the **unsupervised clustering** of call
697 types. The approach reaches high accuracy in the detection of calls (with 91% of correctly
698 identified calls from a training set containing $\approx 2,000$) with interesting insights on the
699 effects of recording quality and inter-individual variation.

700 In a second part of the experiment, the authors also turn to transfer learning
701 through pre-training. After having trained *DeepSqueak* on a gray mouse lemur dataset,
702 they test its robustness in the detection of calls from Goodmans mouse lemurs (*M.*
703 *lehilahytsara*), a closely related species which was never seen by the model during its
704 training, achieving very similar results. This is thus a first example of how a model trained
705 on one species or one dataset can be leveraged in processing a second species or dataset
706 with the assumption that information extracted from the first task can be efficiently
707 reemployed in the latter.

708 4.2 Pre-training

709 Surprisingly, relying on pre-trained models knowledge from taxonomically related
710 primates is not a preferred approach in computational bioacoustics. Rather, most transfer
711 learning experiments are built upon speech-based models with the underlying assumption
712 that human and non-human primates share, at least, some vocal characteristics and that
713 models are sufficiently resistant to such a domain shift. In addition to this, pre-trained
714 **self-supervised** models for speech have been extensively explored in recent years and
715 state-of-the-art solutions are now publicly available and easy to access through dedicated
716 APIs like [HuggingFace](#), [S3PRL](#) or [SpeechBrain](#). We hereby give some of the few examples
717 of how large speech-based PTM such as HuBERT (Hsu et al., 2021), wav2vec2 (Schneider
718 et al., 2019) or DeepTone can be used to efficiently process primate vocalizations, either as
719 frozen feature extractors replacing **mel-spectrograms** and engineered features or as
720 foundation models aimed at offering a unified solution to multiple tasks and species.

721 An essential part in the development of transfer learning models is the comparison
722 of their performance with more traditional approaches, as was done by Jiang et al. (2023).
723 Here, the authors train a Long Short-Term Memory (LSTM) model and a **transformer**
724 model on sound event detection: the segmentation and the automatic identification of call
725 sequences from continuous vocalizations of bonobos (*Pan paniscus*), chimpanzees (*Pan*
726 *troglodites*) and orangutans (*Pongo pygmaeus*). They focus their experiment on comparing
727 performances across three feature extraction processes as input to the models: the raw
728 waveform, spectrograms and wav2vec **embeddings**. Wav2vec (Schneider et al., 2019) is a
729 speech based PTM, relying on **self-supervised** representation learning from raw audio and
730 initially developed for speech recognition. The model consists in a multi-layer **CNN** trained
731 on a noise contrastive binary classification task: it learns to extract informative
732 representations of short sound frames of 30 ms from their context by differentiating them
733 from other, randomly sampled, sound frames. This learned “latent representation”, also
734 known as pre-trained **embedding**, can be seen as a fixed-length vector of 768 elements (for

735 wav2vec LARGE), supposedly encoding essential information from the input audio data.
736 The encoded information was proven to be useful for its initial intended purpose of
737 automatic speech recognition, but it may also incorporate other acoustic properties from
738 speech such as the identity of a speaker or voice print, language information or even
739 emotional expressivity from prosodical content (Y. Wang et al., 2021).

740 In the paper by Jiang et al. (2023), the assumption is that a wav2vec representation,
741 although initially trained on speech, can also encode enough acoustic information to
742 distinguish between great apes call types or to differentiate them from background forest
743 noise. The authors find that training an LSTM on these **embeddings** yields better results
744 compared to spectrograms or the raw waveform. They show the benefits of balancing
745 classes when facing small annotated datasets and give an example of how pre-trained
746 representation can also be used in *zero-shot* classification contexts by training a model on
747 the orangutan dataset and using it to classify bonobo's call types without further training.
748 These results thus show how the use of speech-based models is a promising solution for
749 zero or few-shot learning from small primate datasets, even for cross-species classification.

750 Leroux et al. (2021) give an example of transfer learning from speech for primate
751 vocal signature classification. They formulate a hypothesis for the existence of an
752 acoustically encoded unique individual signature across call types in chimpanzees (*Pan*
753 *troglydytes*) and test it through automatic classification of individuals. In doing so, the
754 authors train several shallow classifiers on top of DeepTone Identity **embeddings**, a model
755 pre-trained on 10,000 unique utterances from human IDs, and compare performances with
756 **MFCC** inputs (a spectral representation often used for speaker identification). Despite the
757 classifier relying on a simple **SVM** architecture and being trained on a rather unbalanced
758 dataset of calls from three individual chimpanzees including three different call types, they
759 reach 80% accuracy with consistently higher performances from DeepTone **embeddings**
760 compared to **MFCCs**. Additionally, the transfer learning approach tends to show higher
761 results compared to the spectral approach in low training data contexts, reaching a

762 maximally higher accuracy when using only 40 training examples. Again, this shows the
763 value of pre-trained **embeddings** for few-shot learning when annotated data is scarce. The
764 approach also gives a hint into the potential acoustic similarity between great apes
765 vocalizations and human speech as well as the great generalizability and domain transfer
766 abilities of speech based PTMs. In contrast, Lakdari et al. (2024) show that **MFCCs** can
767 outperform pre-trained **embeddings** from wav2vec when used as input to an **SVM** for
768 gibbon vocal identity classification. Yet, their experiment relies on modified versions of said
769 **embeddings** (with **embeddings** averaged on several dimensions) which may result in an
770 important loss of information prior to classification. We think that these modifications and
771 the use of a single model may impair fair comparisons and do not properly reflect the
772 abilities of pre-trained speech models (Jiang et al., 2023; Leroux et al., 2021; Sarkar &
773 Doss, 2023).

774 **4.3 Pretext tasks and pre-training data**

775 We have seen that relying on learned representations from **self-supervised** PTMs can
776 boost primate vocalization classification performances compared to using the raw waveform
777 or spectral representations, even when said PTMs were initially trained on speech. With
778 this in mind, an open question remains on the influence of different PTM architectures and
779 pre-training datasets on the performances of downstream classifiers. Currently, speech
780 processing state-of-the-art models are mostly **self-supervised** PTMs, and recent years have
781 seen the emergence of innovative architectures frequently improving benchmark
782 performance with new pre-training datasets, larger numbers of parameters or different
783 “pretext tasks”.

784 These pretext tasks are proxy tasks used to pre-train models on raw acoustic
785 datasets without the need for human supervision (hence the **self-supervised** nature of these
786 models). They consist in generating supervision from the data itself, requiring informative
787 data representation learning from the model which will automatically learn to capture

788 structural information and acoustic patterns to reach low losses and higher predictive
789 performances during training. A first example of such a pretext task is masked modeling,
790 an approach introduced for the textual language model BERT (Devlin et al., 2019), which
791 consists in masking portions of the data (either text, images or sound) and reconstructing
792 said portions from their surrounding context. This approach was successfully implemented
793 for speech in HuBERT by Hsu et al. (2021). As we have already seen for wav2vec,
794 contrastive predictive coding is another pretext task consisting in predicting future sound
795 frames from previous ones, and yields similar results compared to HuBERT.

796 Many such examples exist in the literature and could result in different performances
797 gains when adapted to bioacoustic classification. This was tested by Sarkar and Doss
798 (2023) who compared 11 speech-based PTMs, all trained on similar speech datasets (i.e.,
799 Librispeech for 10 of them and Libri-Light for Modified-CPC) on a common marmoset
800 (*Callithrix jacchus*) caller detection task. These models include wav2vec and HuBERT as
801 well as other state-of-the-art models including APC (Chung & Glass, 2020), Mockingjay
802 (A. T. Liu et al., 2020) or WavLM (Chen et al., 2022), each presenting some specificity in
803 their architecture, sizes and pretext tasks. As in both previously mentioned experiments,
804 all PTM weights are kept frozen (the models are not further trained on unlabeled data)
805 and used as feature extractors for downstream classifiers: **SVMs** and an LSTM for binary
806 caller classification. The downstream models thus predict if two calls are uttered from the
807 same individual or not. As can be seen in Figure 8, the authors test the performance of the
808 downstream model in terms of Area Under the Curve and compare it with PTM size and
809 pretext task (also referred to as the pre-training objective). The autoregressive
810 reconstruction implemented in APC (Chung & Glass, 2020) and its vector quantized
811 variant VQ-APC (Chung et al., 2020) seem to perform slightly better despite smaller model
812 sizes. Surprisingly, Data2vec (Baevski et al., 2022) which was the most successful masked
813 model for several speech tasks at the time of the experiment, performs lower than the rest,
814 thus showing weaker representation learning capabilities in a domain adaptation context.

815 Following the question of model architectures and pretext tasks, the nature of the
816 pre-training dataset could be considered as an essential part of the process, for the reason
817 that PTMs are inherently conditioned to capture knowledge dependent on their training
818 data. We should point out that the idea of using pre-training models from speech to
819 perform bioacoustic tasks is not solely related to a theoretical similarity between speech
820 and animal vocalizations. The approach can also be explained by the extensive availability
821 of speech data in recent years, when the size of a pre-training dataset is an essential
822 prerequisite to the success of *self-supervised* models. Yet, the effect of the nature of a
823 pre-training dataset on bioacoustic tasks performances remains an open question.
824 Although an intuitive answer to this second question would be that the closest in domain a
825 pre-training dataset is to the downstream one, the preliminary results recently showcased
826 by Hagiwara (2023) seem to indicate a more complicated situation. With their
827 *self-supervised* model AVES, the authors go a step further from using speech-based PTMs
828 as feature extractors and test performance gains in terms of pre-training data for an array
829 of downstream tasks (classification and detection on marine and terrestrial mammals,
830 amphibians, birds and primates). Heavily inspired by the HuBERT architecture, they
831 entirely retrain the masked modeling *transformer* on several curated datasets including
832 animal vocalizations, speech and general sound. They propose different pre-training sets by
833 filtering audioset and VGGsound: two collections of several millions of 10 second audio
834 clips drawn from YouTube videos with corresponding categories. By filtering said
835 categories, they build 4 distinct data subsets :

- 836 • *core*: a configuration containing 153 hours of general sounds
- 837 • *bio*: the core configuration with added sounds corresponding to the animal label in
838 VGG sound (360 hours)
- 839 • *non-bio*: a similarly sized control dataset containing random sounds from all
840 categories except the animal one (360 hours)

- 841 • *all*: a dataset containing all types of sounds on top of the core configuration, making
842 up to 5,054 hours of audio.

843 To further test the performance gains of their models, they compare them with
844 VGGish and ResNet (both PTMs developed for general purpose audio-tagging) which were
845 further trained in a **supervised** manner on the tasks at hand. The authors find that the
846 AVES version trained on bioacoustic data (*bio*) outperforms other PTMs, including the
847 **supervised** topline from VGGish and ResNet on most tasks. Although these results seem
848 promising and show the validity of the approach, they must be taken carefully as the
849 bioacoustic pre-training set only increases performance by a small margin. Furthermore,
850 the primate detection task, carried out on Müller gibbons (*Hylobitae muelleri*), shows
851 slightly lower mean average precision compared to the **supervised** and **unsupervised**
852 versions of ResNet. This result may be explained by the scarcity of vocalizations contained
853 in the gibbon dataset and might not entirely reflect the advantages of AVES which can be
854 seen in most of the other tasks.

855 In a broader perspective, the authors compute t-scores to compare the average
856 results obtained from the four pre-training datasets. Surprisingly, *bio* and *non-bio* reach
857 very close performance and improve upon *all* despite their much smaller sizes. This
858 indicates that selecting reduced curated datasets may give better results in a pre-training
859 configuration rather than opting for very large and miscellaneous collections of sounds. The
860 authors thus show their model’s ability to generalize well across domains. This might be
861 seen as a counterargument towards the need for a specific bioacoustic pre-training dataset
862 for transfer learning from **self-supervised** models. It could also mean speech-based models
863 are not successful in bioacoustics because of some acoustic resemblance between speech and
864 animal vocalizations but rather because of their ability to transfer knowledge across
865 acoustic domains. In any case, such assumptions will need to be further tested to account
866 for the many technical limitations which might also explain these counter-intuitive results.

867 Lastly, the results obtained by Ghani et al. (2023), although not specifically tailored

868 for primate vocalization analysis, give a good example of transfer learning across species.
869 In their experiments, large models pre-trained for bird sound classification (namely
870 BirdNet 2.3 and Perch) are compared to general audio tagging models pre-trained on
871 AudioSet (YAMNet, VGGish and AudioMAE). This comparison is carried out through
872 probing: a method consisting in training simple linear layers on the pre-trained **embeddings**
873 to understand how much of the information needed for the downstream task they are able
874 to linearly encode. This gives a better account for the ability of a PTM to capture
875 information for a given task, as the downstream model (a simple linear probe) adds very
876 little knowledge to what was effectively captured by the PTM. In this case, results show
877 that both bird-based models outperform the general event-detection ones by a good margin
878 in detecting and classifying bird sounds as well as other animals such as frogs, cetaceans
879 and bats. This is also the case in few-shot learning, as both models are still on the topline
880 when downstream datasets are reduced in size, thus showing that pre-training models on
881 bird sounds may be a viable option for few-shot learning on other scarcely annotated
882 species. The authors state that this performance gain may be explained by the rich and
883 diverse sounds produced by birds which occupy a broad range both temporally and in the
884 spectral domain with great frequency, harmonic and rhythmical complexity. Added to this,
885 the large amount of publicly available bird song datasets, in par with what can be found
886 for speech compared to the scarcity of primate recordings, makes it another viable option
887 for bioacoustic transfer learning across species. We think that, in addition to model testing
888 and the development of primate-only PTMs, a good amount of work still needs to be put
889 into understanding the influence of pre-training datasets for automatic primate
890 vocalization analysis.

891 5 Discussion

892 5.1 Available datasets

893 As an encouragement for researchers to partake in further testing of the many
894 options we have surveyed so far, we draw a non-exhaustive list of some publicly available
895 primate vocalization datasets. These can be used either as pre-training data for transfer
896 learning models or as manually annotated datasets for **supervised** downstream
897 classification. Some also include open-source models and code to be used as inspiration or
898 as baselines for performance evaluation. See Table 1 for a list of the previously-mentioned
899 papers with code and dataset availability.

900 5.2 General lack of publicly available data

901 As can be seen in Table 1, few annotated datasets of primate vocalization recordings
902 are made publicly available (compared to other taxa, or to the amount of public speech
903 datasets). Despite this, our list is not exhaustive and more unpublished data could be
904 found in addition to ongoing efforts still being carried out to this day in the recording of
905 new datasets and their annotation. This lack of published data could be said to hinder
906 research efforts into primate vocalization analysis and similar issues can be found in the
907 whole field of bioacoustics (Baker & Vincent, 2019). Furthermore, within primate related
908 research, some species are clearly underrepresented for various reasons, including the lack
909 of interest put into the study of apparently poorly complex vocal systems, the remoteness
910 of their habitat or the scarcity of endangered species. This is why the development of
911 efficient machine learning solutions for the processing of primate vocalizations should
912 always be made in parallel with annotation and recording projects as well as a substantial
913 amount of work put into their deposition as supplementary material into public websites.
914 The annotation itself should be thoroughly documented with an emphasis on reduced bias
915 and reproducible methods. Finally, codes and models need to be published in open-source

916 (as is often the case). This encourages their reuse or adaptation, especially when
917 developing large pre-trained foundation models which weights can be reemployed without
918 the need for time and energy-consuming re-training.

919 **5.3 Ethical and environmental concerns**

920 In the field of computational bioacoustics applied to primates, we can observe a
921 general tendency towards the use of PAM paired with deep learning approaches. This may
922 be seen as a promising direction in terms of ethical and environmental concerns. In fact,
923 PAM is considered as a non-invasive solution to the study of animal communication, and
924 its automatic processing with machine learning methods leads to great opportunities for
925 conservation and monitoring projects. Yet, several drawbacks should also be mentioned.

926 First of all, the lack of control over the elements recorded during PAM may lead to
927 privacy concerns when human speech is picked up by the acoustic sensors. This problem, in
928 turn, can be easily circumvented with similar machine learning methods as the ones used
929 for the animal vocalization analysis. As we have mentioned before, speech processing
930 methods for the automatic detection of speech have shown impressive results in the recent
931 years and their implementation is strongly facilitated by the availability of user friendly
932 open-source models. Employing such models to filter out speech, especially when facing
933 recording of animals in captivity should be included as a preprocessing step in such
934 experiments (Janetzky et al., 2021).

935 A second limit is the well-known environmental impact of AI, although
936 computational bioacoustics stays a relatively niche domain of study compared to the
937 research effort put into computer science for image or natural language processing. This
938 important issue has been thoroughly addressed by specific reviews and studies
939 (Van Wynsberghe, 2021). We should mention that promising solutions include the reuse of
940 weights from pre-trained models which limits training time, and the development of
941 foundation models for transfer learning as we discussed in Section 4.

942 Finally, the malicious use of automatic primate monitoring tools should always
943 remain an important concern. As stated in Piel et al. (2022): “The ability of remote
944 sensing tools to incidentally (or deliberately, in the case of poachers) reveal the location,
945 movement and behavior of individuals raises concerns about informed consent, privacy,
946 civil liberties, and fear of arrest”. This problem is unfortunately embedded in the advocacy
947 for open-source models, and no single solution exists. More work needs to be carried out in
948 evaluating the impact of open-source animal detection models and datasets, as was done by
949 Lennox et al. (2020) for biotelemetric data sharing. In contrast however, the detection of
950 poaching activity may be tackled through machine learning solutions and seems to be an
951 actively addressed problem in recent computational bioacoustics studies with tasks like
952 gunshot, chainsaw or illegal cattle farming detection from acoustic data (Pérez-Granados &
953 Schuchmann, 2023; Sethi et al., 2020).

954 **6 Conclusion**

955 Primate vocal communication research has seen a significant shift with the advent of
956 machine learning and artificial neural networks, inspired by bioacoustics studies on other
957 taxa, or sound and speech processing. The use of passive acoustic monitoring and the
958 availability of large annotated datasets have paved the way for innovative automated
959 workflows, reducing our reliance on manual annotations and analysis and questioning some
960 aspects of human bias. This paradigm shift has seen the emergence of new automated
961 tasks with important scientific implications and is starting to turn into valuable monitoring
962 and conservation tools. We have provided a concise survey of the recent directions taken in
963 computational bioacoustics, highlighting emerging approaches and the valuable insights
964 they offer for the study of primate communication. We have discussed recent interests for
965 state-of-the-art deep learning models and their prolific application to primate bioacoustic
966 research, all the while reaffirming the validity and greater transparency of earlier statistical
967 approaches. Looking ahead, the development of high-performance weakly supervised

968 transfer learning models holds promise for further advancements in the field, but challenges
969 remain in understanding the behavior of these *black box* models for their subsequent use as
970 scientific tools. Challenges also remain in terms of data availability, and in turning existing
971 solutions into user-friendly light-weight models accessible to primatologists. As the field of
972 computational bioacoustics evolves, we can expect deep learning research to bring further
973 exciting developments that will continue to enrich our understanding of primate vocal
974 communication and its role in their social interactions and ecological contexts.

975 **Author contributions**

976 Conceptualization: Jules Cauzinille, Ricard Marxer, Arnaud Rey, Benoit Favre;
977 writing and preparation of original version: Jules Cauzinille; revision and editing: Ricard
978 Marxer, Arnaud Rey, Benoit Favre. All authors have read and agreed to the published
979 version of the manuscript

980 **Acknowledgments**

981 This work, carried out within the Institute of Convergence ILCB
982 (ANR-16-CONV-0002), has benefited from support from the French government (France
983 2030), managed by the French National Agency for Research (ANR) and the Excellence
984 Initiative of Aix-Marseille University (A*MIDEX). This work was also supported by the
985 HEBBIAN (ANR-23-CE28-0008) and COMPO (ANR-23- CE23-0031) ANR projects.

986 **Financial disclosure**

987 This research received no external funding. The authors state that the views
988 expressed in the submitted article are their own and do not constitute an official position of
989 the institution or funder.

990 **Conflict of interest**

991 The authors declare no conflict of interests.

992 **Glossary**

993 **clustering** is a machine learning approach which consists in grouping similar objects into
994 different subsets (or clusters). Although many algorithms exist, they generally work
995 by partitioning a dataset into said clusters according to some similarity and/or
996 dissimilarity metric. In the case of primate vocalizations, this approach is usually
997 employed to examine repertoires of call-types. 8, 22–28, 30

998 **CNN** (Convolutional Neural Networks) are a popular **DL** architecture in bioacoustics,
999 specifically designed to extract meaningful information from image inputs (such as
1000 spectrograms) by recognizing patterns in the data. 7, 13–17, 19, 21, 22, 27, 31

1001 **DL** (Deep Learning) models refer to a category of machine learning algorithms based on
1002 neural networks and capable of learning complex patterns from data such as primate
1003 vocalizations. 3, 6, 7, 9–11, 17, 19–21, 24, 28

1004 **embedding** (or latent representation) refers to a numerical vector outputted by a deep
1005 learning model and supposed to encode relevant informative features from input
1006 soundframes. Embeddings are usually extracted from a pre-trained model and used
1007 as input for downstream models in the transfer learning approach (see Figure 3).
1008 31–33, 37

1009 **F1-score** is a metric used to evaluate the performances of a machine learning model in
1010 binary classification tasks. Ranging from 0 (worst) to 1 (best), it is computed as a
1011 balance between precision and recall, making it particularly useful in scenarios with
1012 imbalanced classes or to account for false positives and false negatives. 8

1013 **mel-spectrogram** provides a detailed visual representation of the frequency content of an
1014 audio signal. They can be considered as the conversion of a traditional spectrogram
1015 into the Mel scale, which is better aligned with human auditory perception. 9, 14, 31

1016 **MFCC** (Mel-Frequency Cepstral Coefficients) are a spectral feature extraction method
1017 commonly used in sound processing. They allow capturing essential characteristics of
1018 an audio signal by converting the frequency domain into the so-called “Mel” scale.
1019 This makes MFCCs quite difficult to interpret visually but a very successful feature
1020 extraction method for machine learning models. 19–21, 32, 33

1021 **self-supervised** learning is a pre-training approach relying on pseudo-labels found within
1022 the data itself, without human interventions (as in designing models that will predict
1023 future sound frames given a context). 27–29, 31, 33, 35, 36

1024 **supervised** learning is a machine learning approach involving the use of annotated or
1025 labeled data as training material for a given algorithm. Weak supervision is more
1026 closely related to the **unsupervised** approach and involves using little expert
1027 knowledge during this process (a pre-defined number of call-types during
1028 unsupervised call-type clustering for example). 4, 21, 22, 24, 29, 36, 38, 40

1029 **SVM** (Support Vector Machines) are a class of “simple” machine learning algorithms
1030 designed to linearly separate datapoints into different categories. They are especially
1031 popular in bioacoustics where the use of more complex deep learning models might
1032 not be necessary to reach acceptable performances on a given task. 8, 18, 21, 25,
1033 32–34

1034 **transformer** networks are a type of neural network specifically designed to model
1035 long-range dependencies in sequential data such as sound or text. Their popularity
1036 grew in recent years within natural language processing research but their application
1037 to bioacoustic data remains underexplored. 7, 31, 35

1038 **unsupervised** learning is a machine learning approach which does not rely on expert
1039 annotations. The data used for unsupervised learning is unlabeled and the machine
1040 learning algorithms must categorize it according to the structure of the data itself or
1041 by recognizing specific patterns in said data. [22](#), [23](#), [25](#), [27](#), [28](#), [30](#), [36](#)

1042 **References**

1043 Anders, F., Kalan, A. K., Köhl, H. S., & Fuchs, M. (2021). Compensating class imbalance
1044 for acoustic chimpanzee detection with convolutional recurrent neural networks.
1045 Ecological Informatics, *65*, 101423.

1046 <https://doi.org/https://doi.org/10.1016/j.ecoinf.2021.101423>

1047 Baevski, A., Hsu, W.-N., Xu, Q., Babu, A., Gu, J., & Auli, M. (2022). Data2vec: A general
1048 framework for self-supervised learning in speech, vision and language.
1049 International Conference on Machine Learning, 1298–1312.

1050 <https://doi.org/10.48550/arXiv.2202.03555>

1051 Baker, E., & Vincent, S. (2019). A deafening silence: A lack of data and reproducibility in
1052 published bioacoustics research? Biodiversity Data Journal.

1053 <https://doi.org/10.3897/BDJ.7.e36783>

1054 Banville, H., Chehab, O., Hyvärinen, A., Engemann, D.-A., & Gramfort, A. (2021).
1055 Uncovering the structure of clinical eeg signals with self-supervised learning.

1056 Journal of Neural Engineering, *18*(4). <https://doi.org/10.1088/1741-2552/abca18>

1057 Bayestehtashk, A., Shafran, I., Coleman, K., & Robertson, N. (2014). Detecting
1058 vocalizations of individual monkeys in social groups.

1059 2014 36th Annual International Conference of the IEEE Engineering in Medicine and Biology Soci
1060 4775–4779. <https://doi.org/10.1109/EMBC.2014.6944692>

1061 Best, P., Paris, S., Glotin, H., & Marxer, R. (2023). Deep audio embeddings for vocalisation
1062 clustering. Plos one, *18*(7), e0283396. <https://doi.org/10.1371/journal.pone.0283396>

1063 Bonafos, G., Pudlo, P., Freyermuth, J.-M., Legou, T., Fagot, J., Tronçon, S., & Rey, A.
1064 (2023, October).

1065 Detection and classification of vocal productions in large scale audio recordings
1066 [working paper or preprint]. <https://doi.org/10.48550/arXiv.2302.07640>

1067 Bravo Sanchez, F. J., Hossain, M. R., English, N. B., & Moore, S. T. (2021). Bioacoustic
1068 classification of avian calls from raw sound waveforms with an open-source deep

1069
1070
1071
1072
1073
1074
1075
1076
1077
1078
1079
1080
1081
1082
1083
1084
1085
1086
1087
1088
1089
1090
1091
1092
1093
1094

learning architecture [Publisher: Nature Publishing Group]. Scientific Reports,
11(1), 15733. <https://doi.org/10.1038/s41598-021-95076-6>

Chen, S., Wang, C., Chen, Z., Wu, Y., Liu, S., Chen, Z., Li, J., Kanda, N., Yoshioka, T.,
Xiao, X., et al. (2022). Wavlm: Large-scale self-supervised pre-training for full stack
speech processing. IEEE Journal of Selected Topics in Signal Processing, 16(6),
1505–1518. <https://doi.org/10.1109/JSTSP.2022.3188113>

Chung, Y.-A., & Glass, J. (2020). Generative pre-training for speech with autoregressive
predictive coding.
ICASSP 2020-2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP 2020)
3497–3501. <https://doi.org/10.48550/arXiv.1910.12607>

Chung, Y.-A., Tang, H., & Glass, J. (2020). Vector-Quantized Autoregressive Predictive
Coding. Proc. Interspeech 2020, 3760–3764.
<https://doi.org/10.21437/Interspeech.2020-1228>

Clink, D. J., Bernard, H., Crofoot, M. C., & Marshall, A. J. (2017). Investigating
Individual Vocal Signatures and Small-Scale Patterns of Geographic Variation in
Female Bornean Gibbon (*Hylobates muelleri*) Great Calls.
International Journal of Primatology, 38(4), 656–671.
<https://doi.org/10.1007/s10764-017-9972-y>

Cramer, A. L., Wu, H.-H., Salamon, J., & Bello, J. P. (2019). Look, listen, and learn more:
Design choices for deep audio embeddings.
ICASSP 2019 - 2019 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP 2019)
3852–3856. <https://doi.org/10.1109/ICASSP.2019.8682475>

Crofoot, M., Lambert, T., Kays, R., & Wikelski, M. (2010). Does watching a monkey
change its behaviour? quantifying observer effects in habituated wild primates using
automated telemetry. Animal Behaviour, 80, 475–480.
<https://doi.org/10.1016/j.anbehav.2010.06.006>

1095 Crunchant, A.-S., Borchers, D., Kühl, H., & Piel, A. (2020). Listening and watching: Do
1096 camera traps or acoustic sensors more efficiently detect wild chimpanzees in an open
1097 habitat? Methods in Ecology and Evolution, 11(4), 542–552.

1098 <https://doi.org/10.1111/2041-210X.13362>

1099 Devlin, J., Chang, M.-W., Lee, K., & Toutanova, K. (2019). BERT: Pre-training of deep
1100 bidirectional transformers for language understanding.

1101 Proceedings of the 2019 Conference of the North American Chapter of the Association for Computa-
1102 4171–4186. <https://doi.org/10.18653/v1/N19-1423>

1103 Do Nascimento, L. A., Pérez-Granados, C., & Beard, K. H. (2021). Passive acoustic
1104 monitoring and automatic detection of diel patterns and acoustic structure of
1105 howler monkey roars. Diversity, 13(11). <https://doi.org/10.3390/d13110566>

1106 Dufourq, E., Durbach, I., Hansford, J. P., Hoepfner, A., Ma, H., Bryant, J. V.,
1107 Stender, C. S., Li, W., Liu, Z., Chen, Q., et al. (2021). Automated detection of
1108 hainan gibbon calls for passive acoustic monitoring.

1109 Remote Sensing in Ecology and Conservation, 7(3), 475–487.

1110 <https://doi.org/10.1101/2020.09.07.285502>

1111 Enari, H., Enari, H., Okuda, K., Yoshita, M., Kuno, T., & Okuda, K. (2017). Feasibility
1112 assessment of active and passive acoustic monitoring of sika deer populations.

1113 Ecological Indicators, 79, 155–162.

1114 <https://doi.org/https://doi.org/10.1016/j.ecolind.2017.04.004>

1115 Enari, H., Enari, H. S., Okuda, K., Maruyama, T., & Okuda, K. N. (2019). An evaluation
1116 of the efficiency of passive acoustic monitoring in detecting deer and primates in
1117 comparison with camera traps. Ecological Indicators, 98, 753–762.

1118 <https://doi.org/10.1016/j.ecolind.2018.11.062>

1119 Erb, W., Ross, W., Kazanecki, H., Mitra Setia, T., Madhusudhana, S., & Clink, D. (2023,
1120 April). Vocal complexity in the long calls of Bornean orangutans (preprint). Animal
1121 Behavior and Cognition. bioarXiv. <https://doi.org/10.1101/2023.04.05.535487>

1122
1123
1124
1125
1126
1127
1128
1129
1130
1131
1132
1133
1134
1135
1136
1137
1138
1139
1140
1141
1142
1143
1144
1145
1146
1147

Eyben, F., Wöllmer, M., & Schuller, B. (2010). Opensmile: The munich versatile and fast open-source audio feature extractor.

Proceedings of the 18th ACM International Conference on Multimedia, 1459–1462.

<https://doi.org/10.1145/1873951.1874246>

Fedurek, P., Zuberbühler, K., & Dahl, C. D. (2016). Sequential information in a great ape utterance. Scientific Reports, 6(1), 38226. <https://doi.org/10.1038/srep38226>

Ganchev, T. (2017). Computational bioacoustics: Biodiversity monitoring and assessment (Vol. 4). Walter de Gruyter GmbH & Co KG.

<https://doi.org/10.1515/9781614516316>

Germain, F. G., Chen, Q., & Koltun, V. (2019). Speech denoising with deep feature losses.

Proc. Interspeech 2019, 2723–2727. <https://doi.org/10.1109/Interspeech.2019.8782422>

Ghani, B., Denton, T., Kahl, S., & Klinck, H. (2023, July). Feature Embeddings from Large-Scale Acoustic Bird Classifiers Enable Few-Shot Transfer Learning

[arXiv:2307.06292 [cs, eess]]. Retrieved July 20, 2023, from

<http://arxiv.org/abs/2307.06292>

Guo, S., Xu, P., Miao, Q., Shao, G., Chapman, C. A., Chen, X., He, G., Fang, D.,

Zhang, H., Sun, Y., Shi, Z., & Li, B. (2020). Automatic identification of individual primates with deep learning techniques. iScience, 23(8), 101412.

<https://doi.org/10.1016/j.isci.2020.101412>

Hagiwara, M. (2023). Aves: Animal vocalization encoder based on self-supervision.

ICASSP 2023 - 2023 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP 2023)

1–5. <https://doi.org/10.1109/ICASSP49357.2023.10095642>

Hagiwara, M., Hoffman, B., Liu, J.-Y., Cusimano, M., Effenberger, F., & Zacarian, K.

(2023). Beans: The benchmark of animal sounds.

ICASSP 2023 - 2023 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP 2023)

1–5. <https://doi.org/10.1109/ICASSP49357.2023.10096686>

1148 Heinicke, S., Kalan, A. K., Wagner, O. J., Mundry, R., Lukashevich, H., & Köhl, H. S.
1149 (2015). Assessing the performance of a semiautomated acoustic monitoring system
1150 for primates (K. Jones, Ed.). Methods in Ecology and Evolution, 6(7), 753–763.
1151 <https://doi.org/10.1111/2041-210X.12384>

1152 Hsu, W.-N., Bolte, B., Tsai, Y.-H., Lakhotia, K., Salakhutdinov, R., & Mohamed, A.
1153 (2021). Hubert: Self-supervised speech representation learning by masked prediction
1154 of hidden units.
1155 IEEE/ACM Transactions on Audio, Speech, and Language Processing, PP, 1–1.
1156 <https://doi.org/10.1109/TASLP.2021.3122291>

1157 James, G., Witten, D., Hastie, T., Tibshirani, R., & Taylor, J. (2023).
1158 An introduction to statistical learning: With applications in python. Springer
1159 Nature. <https://doi.org/10.1007/978-3-031-38747-0>

1160 Janetzky, P., Davidson, P., Steininger, M., Krause, A., & Hotho, A. (2021). Detecting
1161 presence of speech in acoustic data obtained from beehives.
1162 Proceedings of the 6th Detection and Classification of Acoustic Scenes and Events 2021 Workshop
1163 26–30.

1164 Jiang, Z., Soldati, A., Schamberg, I., Lameira, A. R., & Moran, S. (2023). Automatic
1165 Sound Event Detection and Classification of Great Ape Calls Using Neural
1166 Networks. <https://doi.org/https://doi.org/10.48550/arXiv.2301.02214>

1167 Kalan, A. K., Mundry, R., Wagner, O. J., Heinicke, S., Boesch, C., & Köhl, H. S. (2015).
1168 Towards the automated detection and occupancy estimation of primates using
1169 passive acoustic monitoring. Ecological Indicators, 54, 217–226.
1170 <https://doi.org/10.1016/j.ecolind.2015.02.023>

1171 Kidney, D., Rawson, B. M., Borchers, D. L., Stevenson, B. C., Marques, T. A., &
1172 Thomas, L. (2016). An efficient acoustic density estimation method with human
1173 detectors applied to gibbons in cambodia. PLOS ONE, 11(5), 1–16.
1174 <https://doi.org/10.1371/journal.pone.0155066>

- 1175 Kiskin, I., Sinka, M., Cobb, A. D., Rafique, W., Wang, L., Zilli, D., Gutteridge, B.,
1176 Dam, R., Marinos, T., Li, Y., et al. (2021). Humbugdb: A large-scale acoustic
1177 mosquito dataset. arXiv e-prints. <https://doi.org/10.5281/zenodo.4904800>
- 1178 Kiskin, I., Zilli, D., Li, Y., Sinka, M., Willis, K., & Roberts, S. (2020). Bioacoustic
1179 detection with wavelet-conditioned convolutional neural networks.
1180 Neural Computing and Applications, *32*(4), 915–927.
1181 <https://doi.org/10.1007/s00521-018-3626-7>
- 1182 Lakdari, M. W., Ahmad, A. H., Sethi, S., Bohn, G. A., & Clink, D. J. (2024).
1183 Mel-frequency cepstral coefficients outperform embeddings from pre-trained
1184 convolutional neural networks under noisy conditions for discrimination tasks of
1185 individual gibbons. Ecological Informatics, *80*, 102457.
1186 <https://doi.org/10.1016/j.ecoinf.2023.102457>
- 1187 Lennox, R. J., Harcourt, R., Bennett, J. R., Davies, A., Ford, A. T., Frey, R. M.,
1188 Hayward, M. W., Hussey, N. E., Iverson, S. J., Kays, R., Kessel, S. T.,
1189 McMahon, C., Muelbert, M., Murray, T. S., Nguyen, V. M., Pye, J. D.,
1190 Roche, D. G., Whoriskey, F. G., Young, N., & Cooke, S. J. (2020). A Novel
1191 Framework to Protect Animal Data in a World of Ecosurveillance. BioScience,
1192 *70*(6), 468–476. <https://doi.org/10.1093/biosci/biaa035>
- 1193 Leroux, M., Al-Khudhairi, O. G., Perony, N., & Townsend, S. W. (2021, December).
1194 Chimpanzee voice prints? Insights from transfer learning experiments from human
1195 voices. <https://doi.org/https://doi.org/10.48550/arXiv.2112.08165>
- 1196 Lin, T., Wang, Y., Liu, X., & Qiu, X. (2022). A survey of transformers. AI Open, *3*,
1197 111–132. <https://doi.org/https://doi.org/10.1016/j.aiopen.2022.10.001>
- 1198 Liu, A. T., Yang, S.-w., Chi, P.-H., Hsu, P.-c., & Lee, H.-y. (2020). Mockingjay:
1199 Unsupervised speech representation learning with deep bidirectional transformer
1200 encoders.

1201
1202
1203
1204
1205
1206
1207
1208
1209
1210
1211
1212
1213
1214
1215
1216
1217
1218
1219
1220
1221
1222
1223
1224
1225
1226
1227

ICASSP 2020-2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP 2020), 6419–6423. <https://doi.org/10.1109/ICASSP40776.2020.9054458>

Liu, X., Zhang, F., Hou, Z., Mian, L., Wang, Z., Zhang, J., & Tang, J. (2023). Self-supervised learning: Generative or contrastive. IEEE Transactions on Knowledge; Data Engineering, *35*(01), 857–876. <https://doi.org/10.1109/TKDE.2021.3090866>

Mesaros, A., Heittola, T., Diment, A., Elizalde, B., Shah, A., Vincent, E., Raj, B., & Virtanen, T. (2017). Dcase 2017 challenge setup: Tasks, datasets and baseline system. DCASE 2017-Workshop on Detection and Classification of Acoustic Scenes and Events. <https://doi.org/doi:10.1109/TASLP.2019.2907016>

Mesaros, A., Heittola, T., & Virtanen, T. (2019). Acoustic scene classification in dcase 2019 challenge: Closed and open set classification and data mismatch setups. Workshop on Detection and Classification of Acoustic Scenes and Events. <https://doi.org/10.33682/m5kp-fa97>

Mielke, A., & Zuberbühler, K. (2013). A method for automated individual, species and call type recognition in free-ranging animals. Animal Behaviour, *86*(2), 475–482. <https://doi.org/10.1016/j.anbehav.2013.04.017>

Mohamed, A., Lee, H.-y., Borgholt, L., Havtorn, J., Edin, J., Igel, C., Kirchhoff, K., Li, S.-W., Livescu, K., Maaløe, L., Sainath, T., & Watanabe, S. (2022). Self-supervised speech representation learning: A review. IEEE Journal of Selected Topics in Signal Processing, *16*(6), 1179–1210. <https://doi.org/10.1109/JSTSP.2022.3207050>

Müller, M. (2007). Dynamic time warping. Information retrieval for music and motion, 69–84. https://doi.org/10.1007/978-3-540-74048-3_4

Pellegrini, T. (2021). Deep-Learning-Based Central African Primate Species Classification with MixUp and SpecAugment. Interspeech 2021, 456–460. <https://doi.org/10.21437/Interspeech.2021-1911>

- 1228 Pérez-Granados, C., & Schuchmann, K.-L. (2023). The sound of the illegal: Applying
1229 bioacoustics for long-term monitoring of illegal cattle in protected areas.
1230 Ecological Informatics, 74, 101981.
1231 <https://doi.org/https://doi.org/10.1016/j.ecoinf.2023.101981>
- 1232 Pérez-Granados, C., & Traba, J. (2021). Estimating bird density using passive acoustic
1233 monitoring: A review of methods and suggestions for further research. Ibis, 163(3),
1234 765–783. <https://doi.org/10.1111/ibi.12944>
- 1235 Pichler, M., & Hartig, F. (2023). Machine learning and deep learning a review for ecologists.
1236 Methods in Ecology and Evolution, 14(4), 994–1016.
1237 <https://doi.org/10.1111/2041-210X.14061>
- 1238 Piel, A. K., Cruncheon, A., Knot, I. E., Chalmers, C., Fergus, P., Mulero-Pázmány, M., &
1239 Wich, S. A. (2022). Noninvasive Technologies for Primate Conservation in the 21st
1240 Century. International Journal of Primatology, 43(1), 133–167.
1241 <https://doi.org/10.1007/s10764-021-00245-z>
- 1242 Pozzi, L., Gamba, M., & Giacoma, C. (2010). The use of artificial neural networks to
1243 classify primate vocalizations: A pilot study on black lemurs.
1244 American Journal of Primatology: Official Journal of the American Society of Primatologists,
1245 72(4), 337–348. <https://doi.org/10.1002/ajp.20786>
- 1246 Radford, A., Narasimhan, K., Salimans, T., Sutskever, I., et al. (2018). Improving language
1247 understanding by generative pre-training.
- 1248 Robakis, E., Watsa, M., & Erkenwick, G. (2018). Classification of producer characteristics
1249 in primate long calls using neural networks.
1250 The Journal of the Acoustical Society of America, 144(1), 344–353.
1251 <https://doi.org/10.1121/1.5046526>
- 1252 Romero-Mujalli, D., Bergmann, T., Zimmermann, A., & Scheumann, M. (2021). Utilizing
1253 DeepSqueak for automatic detection and classification of mammalian vocalizations:

- 1254 A case study on primate vocalizations. Scientific Reports, 11(1), 24463.
1255 <https://doi.org/10.1038/s41598-021-03941-1>
- 1256 Ross, S. R.-J., O'Connell, D. P., Deichmann, J. L., Desjonquères, C., Gasc, A.,
1257 Phillips, J. N., Sethi, S. S., Wood, C. M., & Burivalova, Z. (2023). Passive acoustic
1258 monitoring provides a fresh perspective on fundamental ecological questions.
1259 Functional Ecology, 37(4), 959–975. <https://doi.org/10.1111/1365-2435.14275>
- 1260 Ruan, W., Wu, K., Chen, Q., & Zhang, C. (2022). ResNet-based bio-acoustics presence
1261 detection technology of Hainan gibbon calls. Applied Acoustics, 198, 108939.
1262 <https://doi.org/10.1016/j.apacoust.2022.108939>
- 1263 Sainburg, T., Thielk, M., & Gentner, T. Q. (2020). Finding, visualizing, and quantifying
1264 latent structure across diverse animal vocal repertoires.
1265 PLoS computational biology, 16(10), e1008228.
1266 <https://doi.org/10.1371/journal.pcbi.1008228>
- 1267 Sarkar, E., & Doss, M. M. (2023, May). Can Self-Supervised Neural Networks Pre-Trained
1268 on Human Speech distinguish Animal Callers? [arXiv:2305.14035 [cs, eess]].
1269 <https://doi.org/https://doi.org/10.48550/arXiv.2305.14035>
- 1270 Schneider, S., Baevski, A., Collobert, R., & Auli, M. (2019). wav2vec: Unsupervised
1271 Pre-Training for Speech Recognition. Proc. Interspeech 2019, 3465–3469.
1272 <https://doi.org/10.21437/Interspeech.2019-1873>
- 1273 Schofield, D., Nagrani, A., Zisserman, A., Hayashi, M., Matsuzawa, T., Biro, D., &
1274 Carvalho, S. (2019). Chimpanzee face recognition from videos in the wild using deep
1275 learning. Science Advances, 5(9), eaaw0736. <https://doi.org/10.1126/sciadv.aaw0736>
- 1276 Schuller, B. W., Batliner, A., Bergler, C., Mascolo, C., Han, J., Lefter, I., Kaya, H.,
1277 Amiriparian, S., Baird, A., Stappen, L., Otth, S., Gerczuk, M., Tzirakis, P.,
1278 Brown, C., Chauhan, J., Grammenos, A., Hasthanasombat, A., Spathis, D., Xia, T.,
1279 ... Kaandorp, C. (2021). The INTERSPEECH 2021 Computational Paralinguistics

- 1280 Challenge: COVID-19 Cough, COVID-19 Speech, Escalation & Primates.
1281 <https://doi.org/10.48550/arXiv.2102.13468>
- 1282 Sethi, S., Jones, N., Fulcher, B., Picinali, L., Clink, D., Klinck, H., Orme, D., Wrege, P., &
1283 Ewers, R. (2020). Characterizing soundscapes across diverse ecosystems using a
1284 universal acoustic feature set. Proceedings of the National Academy of Sciences,
1285 117, 202004702. <https://doi.org/10.1073/pnas.2004702117>
- 1286 Stowell, D. (2019). State of the art in computational bioacoustics and machine learning:
1287 How far have we come? Biodiversity Information Science and Standards, 3, e37227.
1288 <https://doi.org/10.3897/biss.3.37227>
- 1289 Stowell, D. (2022). Computational bioacoustics with deep learning: A review and roadmap.
1290 PeerJ, 10. <https://doi.org/10.7717/peerj.13152>
- 1291 Sugai, L. S. M., Silva, T. S. F., Ribeiro, J., José Wagner, & Llusia, D. (2018). Terrestrial
1292 Passive Acoustic Monitoring: Review and Perspectives. BioScience, 69(1), 15–25.
1293 <https://doi.org/10.1093/biosci/biy147>
- 1294 Turesson, H. K., Ribeiro, S., Pereira, D. R., Papa, J. P., & de Albuquerque, V. H. C.
1295 (2016). Machine learning algorithms for automatic classification of marmoset
1296 vocalizations (M. Smotherman, Ed.). PLOS ONE, 11(9), e0163041.
1297 <https://doi.org/10.1371/journal.pone.0163041>
- 1298 Tzirakis, P., Shiarella, A., Ewers, R., & Schuller, B. W. (2020). Computer Audition for
1299 Continuous Rainforest Occupancy Monitoring: The Case of Bornean Gibbons Call
1300 Detection. Interspeech 2020, 1211–1215.
1301 <https://doi.org/10.21437/Interspeech.2020-2655>
- 1302 Van Wynsberghe, A. (2021). Sustainable ai: Ai for sustainability and the sustainability of
1303 ai. AI and Ethics, 1(3), 213–218. <https://doi.org/10.1007/s43681-021-00043-6>
- 1304 von Luxburg, U., Williamson, R. C., & Guyon, I. (2012, February). Clustering: Science or
1305 art? In I. Guyon, G. Dror, V. Lemaire, G. Taylor, & D. Silver (Eds.),

1306
1307
1308
1309
1310
1311
1312
1313
1314
1315
1316
1317
1318
1319
1320
1321
1322
1323
1324
1325
1326
1327
1328
1329
1330
1331

Proceedings of icml workshop on unsupervised and transfer learning (pp. 65–79, Vol. 27). PMLR.

Wadewitz, P., Hammerschmidt, K., Battaglia, D., Witt, A., Wolf, F., & Fischer, J. (2015). Characterizing vocal repertoires hard vs. soft classification approaches. PLOS ONE, 10(4), 1–16. <https://doi.org/10.1371/journal.pone.0125785>

Wang, Y., Boumadane, A., & Heba, A. (2021). A fine-tuned wav2vec 2.0/hubert benchmark for speech emotion recognition, speaker verification and spoken language understanding. CoRR, abs/2111.02735. <https://doi.org/10.48550/arXiv.2111.02735>

Wang, Y., Ye, J., & Borchers, D. L. (2022). Automated call detection for acoustic surveys with structured calls of varying length. Methods in Ecology and Evolution, 13(7), 1552–1567. <https://doi.org/10.1111/2041-210X.13873>

Wu, H.-H., Kao, C.-C., Tang, Q., Sun, M., McFee, B., Bello, J. P., & Wang, C. (2021). Multi-task self-supervised pre-training for music classification. ICASSP 2021 - 2021 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP 2021), 556–560. <https://doi.org/10.1109/ICASSP39728.2021.9414405>

Xie, J., Colonna, J. G., & Zhang, J. (2021). Bioacoustic signal denoising: A review. Artif. Intell. Rev., 54(5), 3575–3597. <https://doi.org/10.1007/s10462-020-09932-4>

Xie, J., Hu, K., Guo, Y., Zhu, Q., & Yu, J. (2021). On loss functions and cnns for improved bioacoustic signal classification. Ecological Informatics, 64, 101331. <https://doi.org/10.1016/j.ecoinf.2021.101331>

Zimmer, W. M. X. (2011). Passive acoustic monitoring of cetaceans. Cambridge University Press. <https://doi.org/10.1017/CBO9780511977107>

Zimmermann, A., Zimmermann, E., Newman, J. D., & Jürgens, U. (1995). Artificial neural networks for analysis and recognition of primate vocal communication. In Current topics in primate vocal communication (pp. 29–46). Springer US. https://doi.org/10.1007/978-1-4757-9930-9_2

- 1332 Zwerts, J. A., Treep, J., Kaandorp, C., Meewis, F., Koot, A. C., Kaya, H., et al. (2021).
1333 Introducing a central african primate vocalisation dataset for automated species
1334 classification. INTERSPEECH 2021, 466–470.
1335 <https://doi.org/10.21437/Interspeech.2021-154>

Table 1
Summary of cited experiments with data and code availability.

Citation	Species	Task	Architecture	Features	Availability
Dufourq et al. (2021)	Hainan Gibbons	Segmentation	supervised CNN		Dataset + code
Ruan et al. (2022)	Hainan Gibbons	Segmentation	ResNet		Subset + code
Y. Wang et al. (2022)	Hainan Gibbons	Segmentation	CRNN		Dataset + code
Tzirakis et al. (2020)	Miller gibbons	Segmentation	CRNN		Code
Bayestehshak et al. (2014)	rhesus macaques	Segmentation	SVM	OpenSmile	
Robakis et al. (2018)	emperor tamarins + saddleback tamarins	Individual classification	Linear layers	Handcrafted	
Clink et al. (2017)	Miller gibbons	Individual discrimination	DFA	Handcrafted	Dataset
Fedurek et al. (2016)	chimpanzees	Individual classification	SVM	MFCC	
Mielke and Zuberbühler (2013)	Stuhman's blue monkeys + others	Individual classification + species recognition	MLP		
Zwerts et al. (2021)	chimpanzees + mandrills + red-capped mangabeys + guenons	species recognition	KELM	MFCC	Dataset + code
Pellegrini (2021)	chimpanzees + mandrills + red-capped mangabeys + guenons	species recognition	CNN + MobileNet + Resnet		Code
Pozzi et al. (2010)	black lemurs	call type classification	ANN + DFA + clustering		
Tureson et al. (2016)	common marmoset	call type classification	Seven models (SVM, kNN, OPE...)		Dataset + Code
Erb et al. (2023)	Bornean orangutans	call type clustering	SVM + clustering		
Sainburg et al. (2020)	gibbons + rhesus macaques + other taxa	call type clustering	UMAPS	spectrograms	Dataset
Best et al. (2023)	no primates	call type clustering	SSL CNN auto-encoder + UMAPS	spectrograms + PTMs	Code
Romero-Mujalli et al. (2021)	gray mouse lemurs + goodman's mouse lemurs	segmentation + classification + call type clustering	DeepSqueek (FastCNN)		Dataset + code
Jiang et al. (2023)	bonobos + chimpanzees + orangutans	segmentation + classification	CNN	waveform + spectrograms + wav2vec embeddings	Code
Leroux et al. (2021)	chimpanzees	Individual classification	ANNS	DeepTone	
Sarkar and Doss (2023)	common marmosets	individual discrimination	SVM, LSTM	11 speech PTMs	
Hagiwara (2023)	gibbons + other taxa	segmentation + other	AVES (bioacoustic HUBERT)		Dataset + code

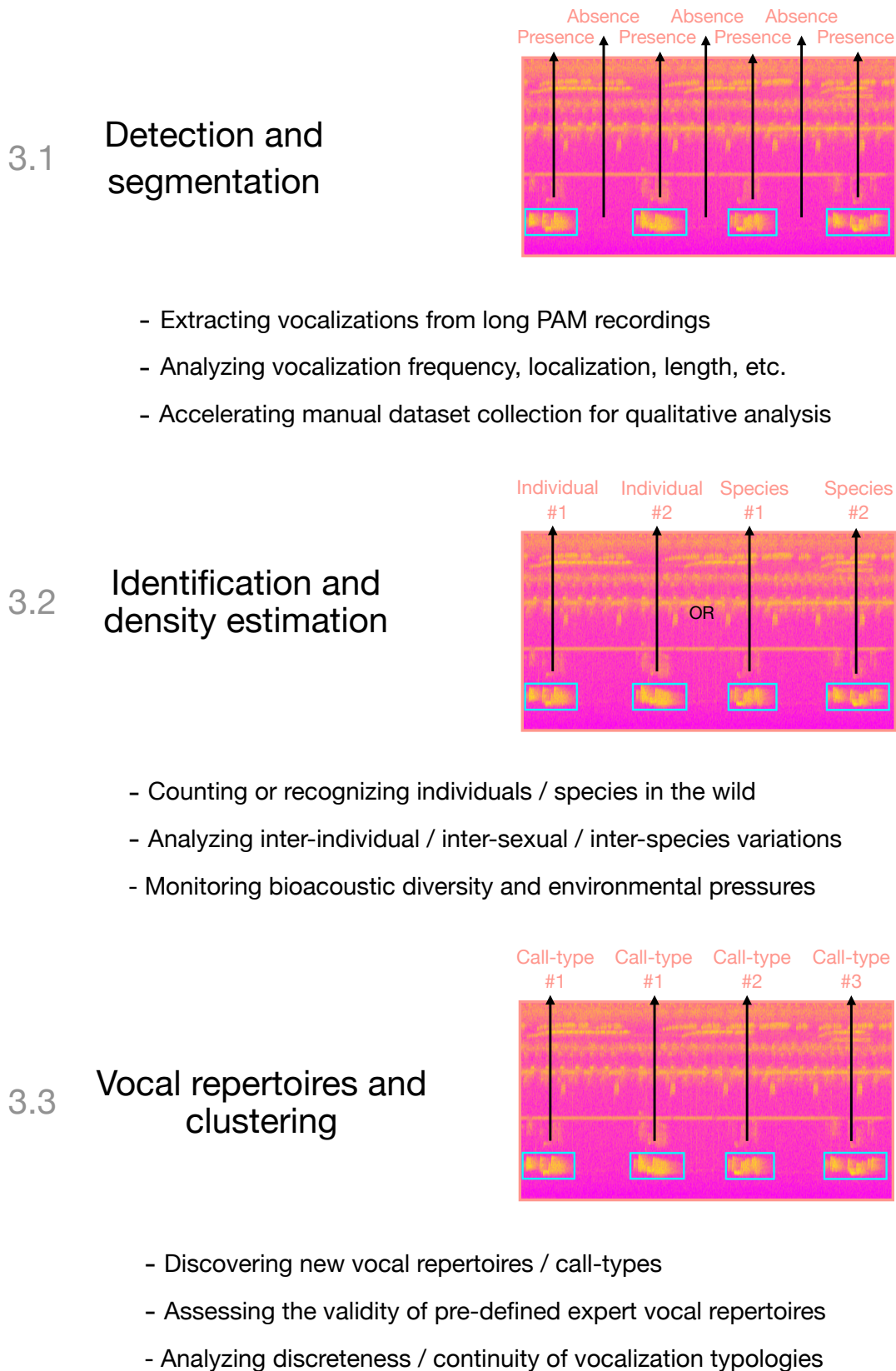


Figure 1

The three main categories of tasks tackled with machine learning for primate bioacoustics.

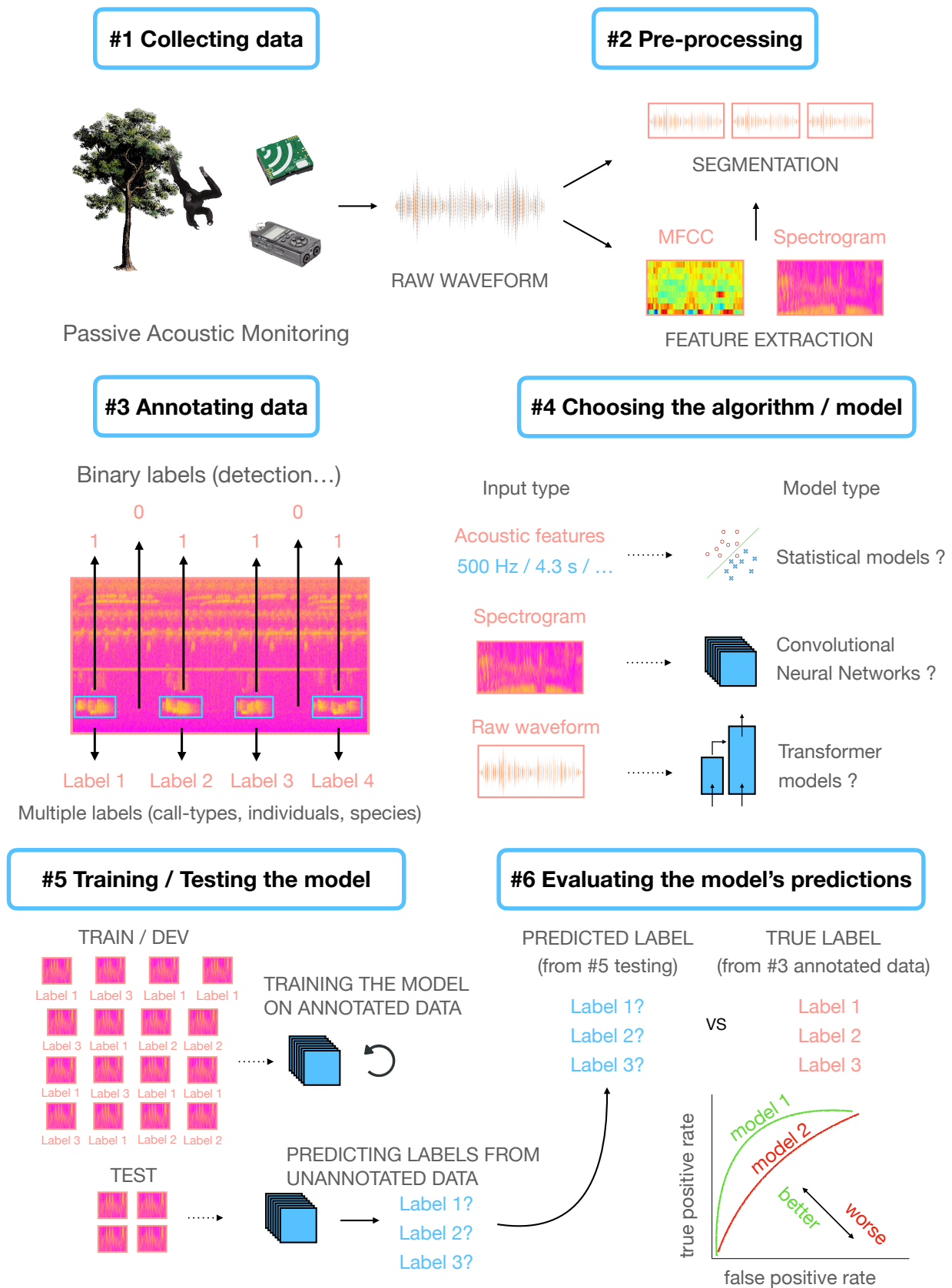
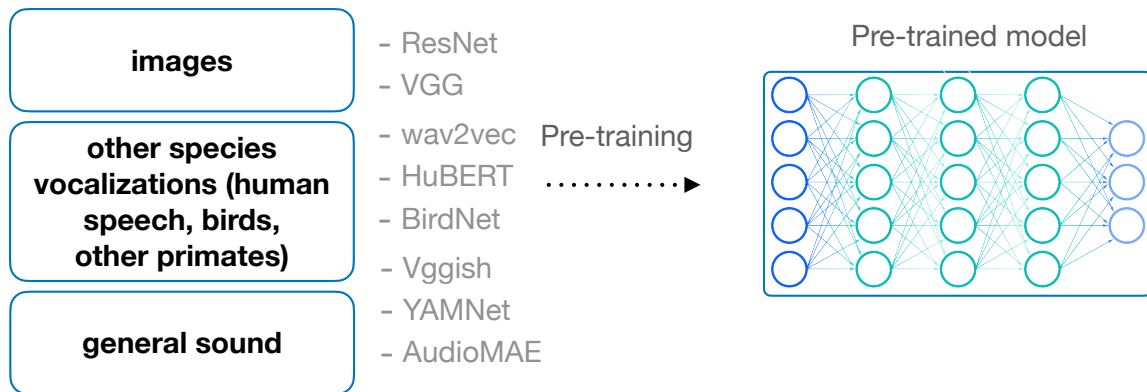


Figure 2
Machine learning workflow for primate bioacoustics.

#1 Pre-training

Pre-training data:

- annotated (supervised)
- not annotated (unsupervised / self-supervised)



#2 Knowledge transfer

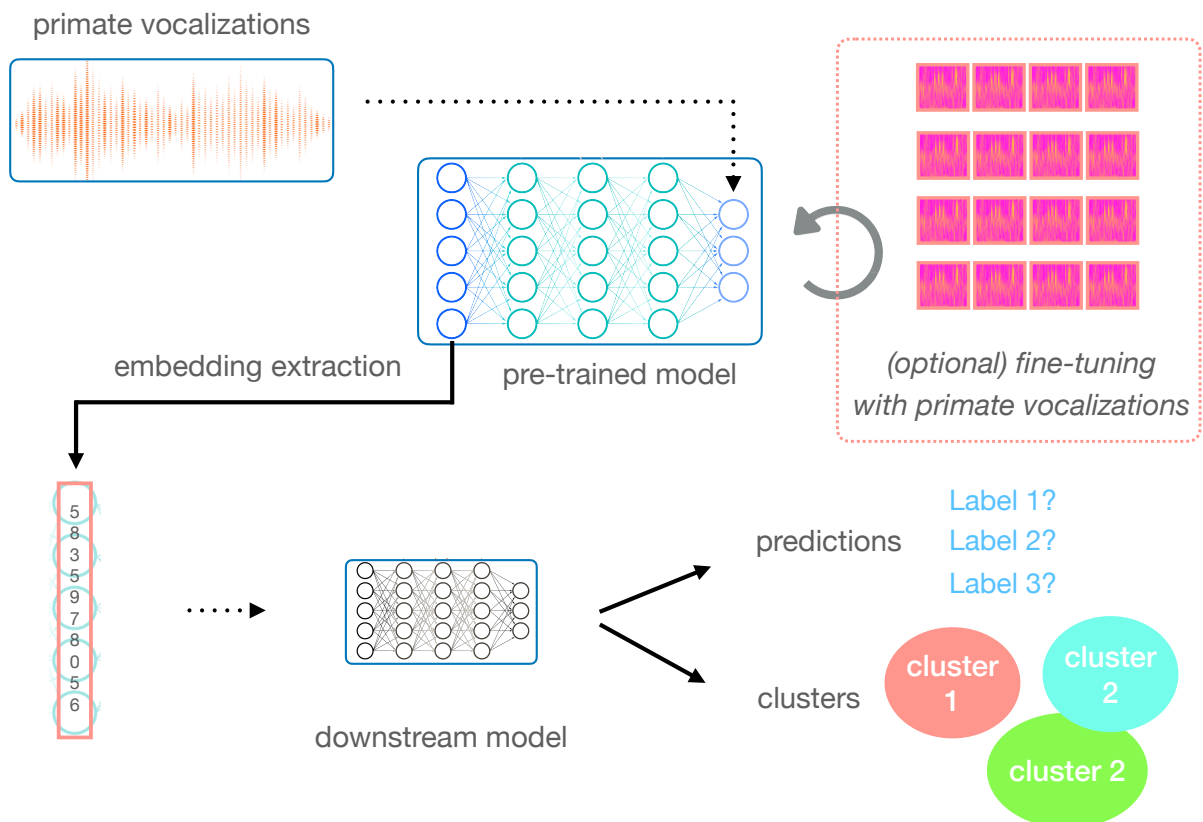


Figure 3
 The transfer learning approach: Using a pre-trained model to classify primate vocalizations (ResNet, VGG, wav2vec, etc. are examples of publicly available pre-trained models).

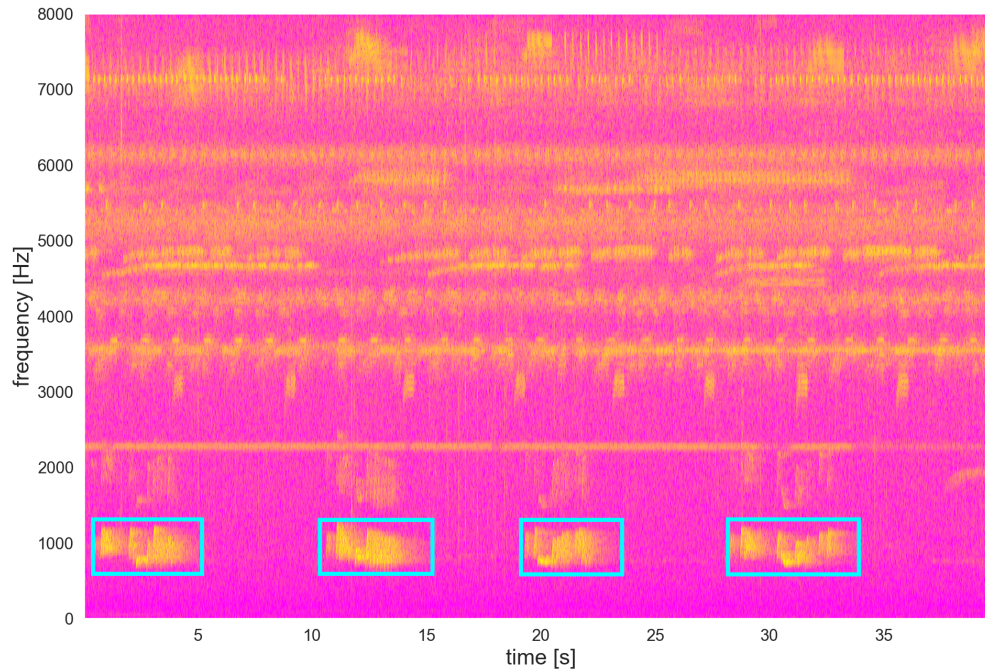


Figure 4

Spectrogram of a Müller Gibbon call. The blue boxes correspond to time and frequency boundaries of the calls. Data from Clink et al. (2017).

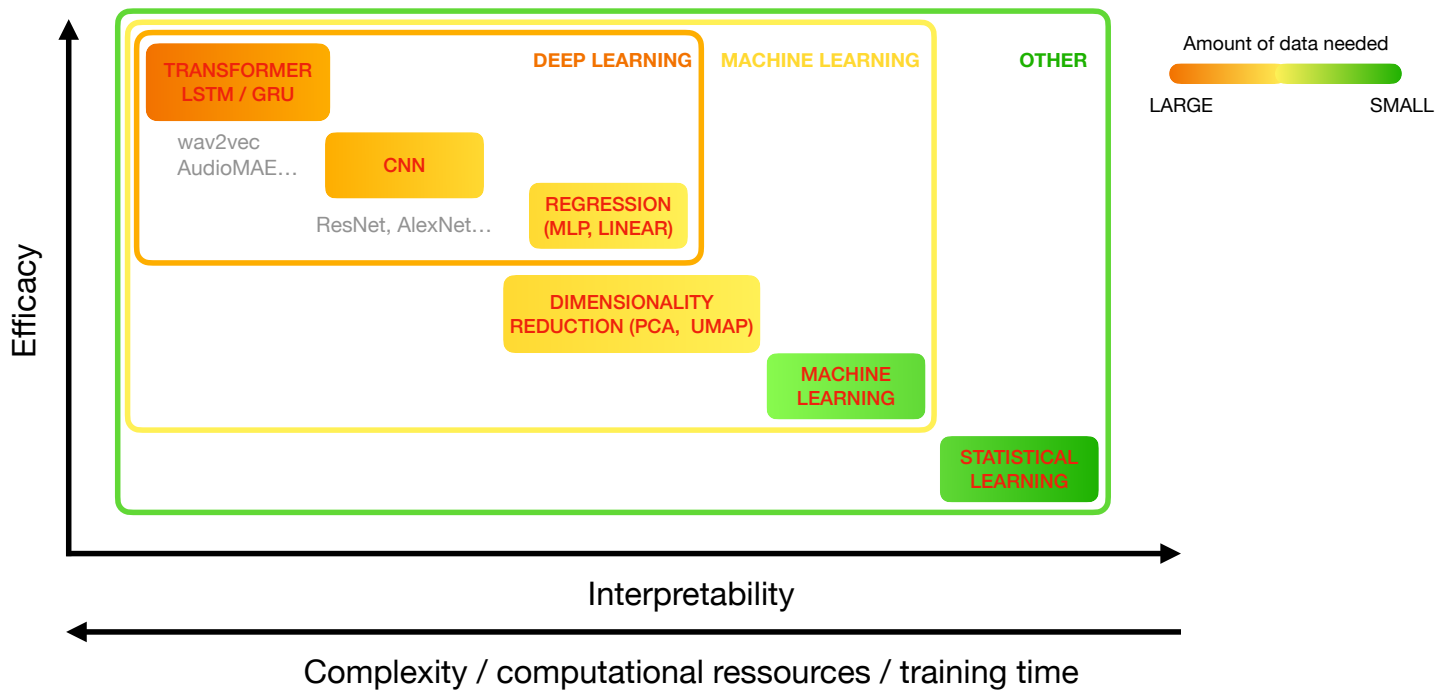


Figure 5

The efficacy / explainability tradeoff between different machine learning architectures (wav2vec, AudioMAE, ResNet and AlexNet are examples of popular deep learning architectures used in bioacoustics).

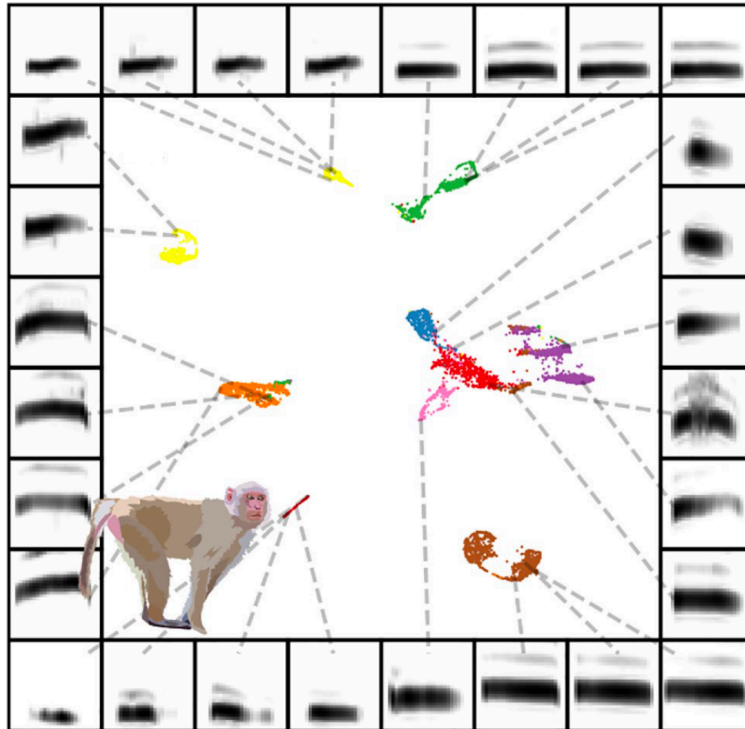


Figure 6

*Spectrograms of Rhesus macaques (*Macaca mulatta*) vocal elements discretized and embedded into a 2D UMAP space. Scatter plot points are colored by individual identity. Image from Sainburg et al. (2020).*

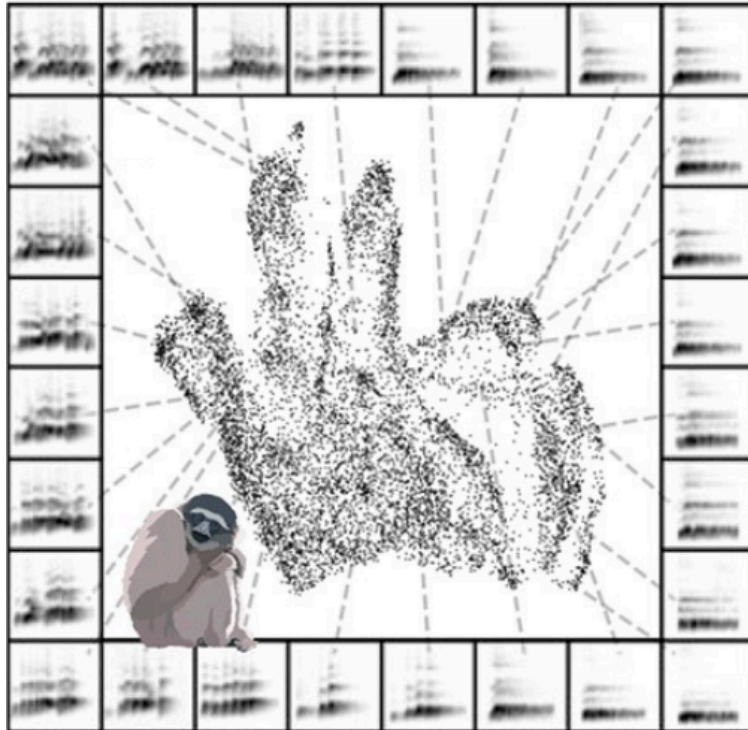


Figure 7

Gibbon syllable spectrograms embedded into a 2D UMAP space. Image from Sainburg et al. (2020).

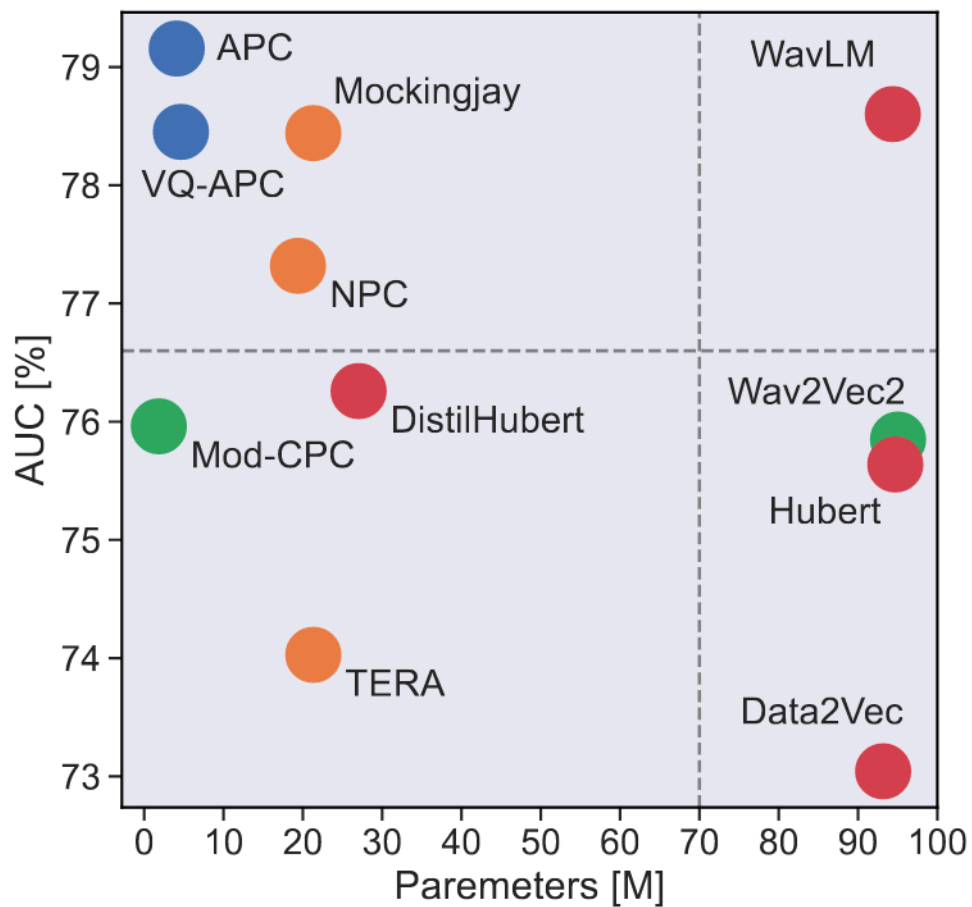


Figure 8

Model size against performance on a primate bioacoustics task. Model pre-training objective denoted as: Masked prediction (red). Autoregressive reconstruction (blue). Contrastive (green). Masked reconstruction (orange). AUC is the Area Under Curve evaluation metric. Figure from Sarkar and Doss (2023).