



Automatic detection for bioacoustic research: a practical guide to the state of the art and future directions

| | |
|-------------------------------|---|
| Journal: | <i>Biological Reviews</i> |
| Manuscript ID | Draft |
| Manuscript Type: | Original Article |
| Date Submitted by the Author: | n/a |
| Complete List of Authors: | <p>Kershenbaum, Arik; University of Cambridge, Girton College Akçay, Çağlar; Anglia Ruskin University, Behavioural Ecology Research Group Babu-Saheer, Lakshmi ; Anglia Ruskin University, Computing Informatics and Applications Research Group Barnhill, Alex; Friedrich-Alexander-Universität Erlangen-Nürnberg, Pattern Recognition Lab, Department of Computer Science Best, Paul; Université de Toulon, Aix Marseille Univ Cauzinille, Jules; Université de Toulon, Aix Marseille Univ Clink, Dena; Cornell Lab of Ornithology, K. Lisa Yang Center for Conservation Bioacoustics Dassow, Angela; Carthage College, Biology Dufourq, Emmanuel; African Institute for Mathematical Sciences; Stellenbosch University, Mathematical Sciences Growcott, Jonathan; University of Exeter - Cornwall Campus, Centre of Ecology and Conservation, College of Life and Environmental Sciences; Recanati-Kaplan Centre, Wildlife Conservation Research Unit Markham, Andrew; University of Oxford, Computer Science Marti-Domken, Barbara; University of Oviedo Marxer, Ricard; Université de Toulon, Aix Marseille Univ Muir, Jen; Anglia Ruskin University, Behavioural Ecology Research Group Reynolds, Sam; Anglia Ruskin University, Behavioural Ecology Research Group Root-Gutteridge, Holly; University of Lincoln, Dept of Life Sciences, College of Health and Science Sadhukhan, Sougata ; Bharati Vidyapeeth (Deemed to be University) Schindler, Loretta; Charles University, Zoology Smith, Bethany; Zoological Society of London Institute of Zoology Stowell, Dan; Tilburg University; Naturalis Biodiversity Center Wascher, Claudia; Anglia Ruskin University, Behavioural Ecology Research Group Dunn, Jacob; Anglia Ruskin University, Behavioural Ecology Research Group; University of Cambridge, Archaeology; University of Vienna, Behavioral and Cognitive Biology</p> |
| Keywords: | artificial intelligence, bioacoustics, machine learning, passive acoustic monitoring, acoustic monitoring, animal communication, automatic detection, classification, deep learning, neural networks |
| | |

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60



1
2
3 **1 Automatic detection for bioacoustic research: a practical guide to the state of the art and future**
4
5 **2 directions**
6

7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

3
4 Arik Kershenbaum¹, Çağlar Akçay², Lakshmi Babu-Saheer³, Alex Barnhill⁴, Paul Best⁵, Jules
5 Cauzinille⁵, Dena Clink⁶, Angela Dassow⁷, Emmanuel Dufourq^{8,9,10}, Jonathan Growcott^{11,12},
6 Andrew Markham¹³, Barbara Marti-Domken¹⁴, Ricard Marxer⁵, Jen Muir², Sam Reynolds²,
7 Holly Root-Gutteridge¹⁵, Sougata Sadhukhan¹⁶, Loretta Schindler¹⁷, Bethany R. Smith¹⁸, Dan
8 Stowell^{19,20}, Claudia A.F. Wascher², Jacob C. Dunn^{2,21,22*}

9
10 ¹Girton College and Department of Zoology, University of Cambridge, Cambridge, CB3 0JG,
11 United Kingdom

12 ²Behavioural Ecology Research Group, Anglia Ruskin University, East Road, Cambridge,
13 CB1 1PT, United Kingdom

14 ³Computing Informatics and Applications Research Group, Anglia Ruskin University, East
15 Road, Cambridge, CB1 1PT, United Kingdom

16 ⁴Pattern Recognition Lab, Department of Computer Science, Friedrich-Alexander-Universität
17 Erlangen-Nürnberg, 91058, Erlangen, Germany

18 ⁵Université de Toulon, Aix Marseille Univ, CNRS, LIS, ILCB, Toulon, France

19 ⁶K. Lisa Yang Center for Conservation Bioacoustics, Cornell Lab of Ornithology, Cornell
20 University, Ithaca, NY, USA

21 ⁷Biology Department, Carthage College, Kenosha, Wisconsin, United States

22 ⁸African Institute for Mathematical Sciences, South Africa

23 ⁹Department of Mathematical Sciences, Stellenbosch University, South Africa

24 ¹⁰National Institute for Theoretical and Computational Sciences, South Africa

- 1
2
3 25 ¹¹Centre of Ecology and Conservation, College of Life and Environmental Sciences,
4
5 26 University of Exeter, Cornwall Campus, TR10 9FE, United Kingdom
6
7
8 27 ¹²Wildlife Conservation Research Unit, Reanati-Kaplan Centre, Tubney House, Abingdon
9
10 28 Road Tubney, Abingdon OX13 5QL, United Kingdom
11
12 29 ¹³Department of Computer Science, University of Oxford, Oxford, OX1 3QD, United
13
14 30 Kingdom
15
16
17 31 ¹⁴University of Oviedo, Principality of Asturias, Spain
18
19 32 ¹⁵Dept of Life Sciences, College of Health and Science, University of Lincoln, Brayford
20
21 33 Pool, Lincoln, LN5 7TS
22
23
24 34 ¹⁶Institute of Environment Education and Research, Bharati Vidyapeeth Deemed University,
25
26 35 Pune, India
27
28 36 ¹⁷Department of Zoology, Faculty of Science, Charles University, Prague 128 44, Czech
29
30 37 Republic
31
32
33 38 ¹⁸Institute of Zoology, Zoological Society of London, London NW1 4RY, UK
34
35 39 ¹⁹Tilburg University, Tilburg, The Netherlands
36
37
38 40 ²⁰Naturalis Biodiversity Center, Leiden, The Netherlands
39
40 41 ²¹Department of Archaeology, University of Cambridge, Cambridge, CB2 3DZ, United
41
42 42 Kingdom
43
44
45 43 ²²Department of Behavioral and Cognitive Biology, University of Vienna, Vienna, Austria
46
47 44
48
49 45
50
51 46
52
53
54 47 *Correspondence: Dr Jacob C. Dunn, Behavioural Ecology Research Group, School of Life
55
56 48 Sciences, Anglia Ruskin University, Cambridge, CB1 1PT, United Kingdom; Telephone: +44
57
58 49 1223 698220; Email: jacob.dunn@aru.ac.uk
59
60

1
2
3 50 ABSTRACT
4

5 51 Recent years have seen a dramatic rise in the use of passive acoustic monitoring (PAM) for
6
7 52 biological and ecological applications, and a corresponding increase in the volume of data
8
9 53 generated. However, datasets are often becoming so sizable that analysing them manually is
10
11 54 burdensome and unrealistic. Fortunately, we have also seen a corresponding rise in
12
13 55 computing power and the capability of machine learning algorithms, which offer the
14
15 56 possibility of performing some of the analysis required for PAM automatically. Nonetheless,
16
17 57 the field of automatic detection of acoustic events is still in its infancy in biology and
18
19 58 ecology. In this review, we examine the trends in bioacoustic PAM applications, and their
20
21 59 implications for the burgeoning amount of data that needs to be analysed. We explore the
22
23 60 different methods of machine learning and other tools for scanning, analysing, and extracting
24
25 61 acoustic events automatically from large volumes of recordings. We then provide a step-by-
26
27 62 step practical guide for using automatic detection in bioacoustics.

28
29 63 One of the biggest challenges to greater use of automatic detection in bioacoustics is that
30
31 64 there is often a gulf in expertise between the biological sciences and the field of machine
32
33 65 learning and computer science. Therefore, this review first presents an overview of the
34
35 66 requirements for automatic detection in bioacoustics, intended to familiarise those from a
36
37 67 computer science background with the needs of the bioacoustics community, followed by an
38
39 68 introduction to the key elements of machine learning and artificial intelligence that a biologist
40
41 69 needs to understand to incorporate automatic detection into their research. We then provide a
42
43 70 practical guide to building an automatic detection pipeline for bioacoustic data, and conclude
44
45 71 with a discussion of possible future directions in the field.
46
47
48
49
50

51
52
53 72

54
55 73 KEYWORDS

56
57 74 Animal communication, Automatic detection, Artificial intelligence, Bioacoustics, Classification,
58
59 75 Deep learning, Machine learning, Neural networks, Passive acoustic monitoring
60

| | | | |
|----|----|---|----|
| 1 | | | |
| 2 | | | |
| 3 | 76 | | |
| 4 | | | |
| 5 | | | |
| 6 | 77 | CONTENTS | |
| 7 | | | |
| 8 | 78 | 1. INTRODUCTION | 7 |
| 9 | | | |
| 10 | | | |
| 11 | 79 | 1.1. Acoustic monitoring | 7 |
| 12 | | | |
| 13 | | | |
| 14 | 80 | 1.1.1. What is automatic detection? | 8 |
| 15 | | | |
| 16 | | | |
| 17 | 81 | 1.2. Scope of review | 9 |
| 18 | | | |
| 19 | 82 | 2. BACKGROUND OF AUTOMATIC DETECTION IN BIOACOUSTICS | 10 |
| 20 | | | |
| 21 | | | |
| 22 | 83 | 2.1. What is automatic detection and why do we need it?..... | 10 |
| 23 | | | |
| 24 | | | |
| 25 | 84 | 2.2. The current state of the art in automatic detection | 10 |
| 26 | | | |
| 27 | | | |
| 28 | 85 | 2.3. What do we aspire for from automatic detection?..... | 12 |
| 29 | | | |
| 30 | 86 | 3. PERSPECTIVES FROM BIOLOGICAL SCIENCES..... | 13 |
| 31 | | | |
| 32 | | | |
| 33 | 87 | 3.1. Overview of uses of automatic detection in the biological sciences | 13 |
| 34 | | | |
| 35 | | | |
| 36 | 88 | 3.1.1. Ecosystems and acoustic indices..... | 14 |
| 37 | | | |
| 38 | 89 | 3.1.2. Species occupancy and density | 15 |
| 39 | | | |
| 40 | | | |
| 41 | 90 | 3.1.3. Spatial analyses | 16 |
| 42 | | | |
| 43 | | | |
| 44 | 91 | 3.1.4. Species characteristics..... | 17 |
| 45 | | | |
| 46 | | | |
| 47 | 92 | 3.1.5. Populations and social groups..... | 19 |
| 48 | | | |
| 49 | 93 | 3.1.6. Individual characteristics..... | 19 |
| 50 | | | |
| 51 | | | |
| 52 | 94 | 3.2. Key challenges..... | 20 |
| 53 | | | |
| 54 | | | |
| 55 | 95 | 4. TECHNICAL PERSPECTIVES..... | 21 |
| 56 | | | |
| 57 | | | |
| 58 | 96 | 4.1. Perspectives from computer science..... | 21 |
| 59 | | | |
| 60 | | | |

| | | | |
|----|-----|--|----|
| 1 | | | |
| 2 | | | |
| 3 | 97 | 4.1.1. The role of computation in automatic detection..... | 21 |
| 4 | | | |
| 5 | | | |
| 6 | 98 | 4.1.2. State of the art in automatic detection methods | 23 |
| 7 | | | |
| 8 | | | |
| 9 | 99 | 4.1.3. Assessing pre-existing models | 27 |
| 10 | | | |
| 11 | 100 | 4.2. Conclusions of the technical constraints on the current uses, limitations and | |
| 12 | | | |
| 13 | | | |
| 14 | 101 | expectations of automatic detection | 28 |
| 15 | | | |
| 16 | 102 | 5. A PRACTICAL GUIDE TO AUTOMATIC DETECTION | 29 |
| 17 | | | |
| 18 | | | |
| 19 | 103 | 5.1. Define research questions..... | 30 |
| 20 | | | |
| 21 | | | |
| 22 | 104 | 5.2. Study design | 30 |
| 23 | | | |
| 24 | | | |
| 25 | 105 | 5.3. Start with a pilot study (if possible) | 31 |
| 26 | | | |
| 27 | | | |
| 28 | 106 | 5.4. Data collection and archiving..... | 31 |
| 29 | | | |
| 30 | | | |
| 31 | 107 | 5.5. Data annotation..... | 33 |
| 32 | | | |
| 33 | | | |
| 34 | 108 | 5.6. Choose your Detection Pipeline | 38 |
| 35 | | | |
| 36 | 109 | 5.6.1. Interfacing with your pipeline | 39 |
| 37 | | | |
| 38 | 110 | 5.6.2. Split your data | 40 |
| 39 | | | |
| 40 | | | |
| 41 | 111 | 5.6.3. Pick your feature representation..... | 41 |
| 42 | | | |
| 43 | | | |
| 44 | 112 | 5.6.4. Decide on feature transformation..... | 43 |
| 45 | | | |
| 46 | | | |
| 47 | 113 | 5.6.5. Decide on a method..... | 44 |
| 48 | | | |
| 49 | 114 | 5.7. Verifications - check your results..... | 48 |
| 50 | | | |
| 51 | | | |
| 52 | 115 | 5.7.1. When is a model good enough? Performance thresholds | 49 |
| 53 | | | |
| 54 | | | |
| 55 | 116 | 5.7.2. How harmful are mistakes (false positives vs false negatives)? | 49 |
| 56 | | | |
| 57 | | | |
| 58 | 117 | 5.7.3. Reproducibility and accessibility | 50 |
| 59 | | | |
| 60 | 118 | 5.7.4. Access to raw recordings | 51 |

| | | | |
|----|-----|--|----|
| 1 | | | |
| 2 | | | |
| 3 | 119 | 6. WAYS FORWARD..... | 51 |
| 4 | | | |
| 5 | | | |
| 6 | 120 | 6.1. Challenges | 51 |
| 7 | | | |
| 8 | | | |
| 9 | 121 | 6.1.1. Bioacoustic challenges..... | 51 |
| 10 | | | |
| 11 | 122 | 6.1.2. Computational challenges | 52 |
| 12 | | | |
| 13 | | | |
| 14 | 123 | 6.2. Future directions..... | 56 |
| 15 | | | |
| 16 | | | |
| 17 | 124 | 6.2.1. Accessibility..... | 56 |
| 18 | | | |
| 19 | 125 | 6.2.2. Foundation models..... | 57 |
| 20 | | | |
| 21 | | | |
| 22 | 126 | 6.2.3. Multi-modal detection..... | 58 |
| 23 | | | |
| 24 | | | |
| 25 | 127 | 6.2.4. Keeping a biologist in the loop | 59 |
| 26 | | | |
| 27 | | | |
| 28 | 128 | 7. CONCLUSIONS | 61 |
| 29 | | | |
| 30 | 129 | 7.1. Need for AD | 61 |
| 31 | | | |
| 32 | | | |
| 33 | 130 | 7.2. Cooperation between disciplines..... | 61 |
| 34 | | | |
| 35 | | | |
| 36 | 131 | 7.3. Deep neural networks..... | 62 |
| 37 | | | |
| 38 | | | |
| 39 | 132 | 7.4. Development pipelines | 62 |
| 40 | | | |
| 41 | 133 | 8. ACKNOWLEDGEMENTS..... | 62 |
| 42 | | | |
| 43 | | | |
| 44 | 134 | 9. BIBLIOGRAPHY..... | 63 |
| 45 | | | |
| 46 | | | |
| 47 | 135 | | |
| 48 | | | |
| 49 | 136 | | |
| 50 | | | |
| 51 | 137 | | |
| 52 | | | |
| 53 | | | |
| 54 | | | |
| 55 | | | |
| 56 | | | |
| 57 | | | |
| 58 | | | |
| 59 | | | |
| 60 | | | |

138 1. INTRODUCTION

139 1.1. Acoustic monitoring

140 The acoustic monitoring of captive and wild animals provides valuable data for many
141 purposes, including scientific research, conservation efforts, management decisions, and the
142 welfare of individual animals. Acoustic data can be collected using handheld microphones,
143 on-animal devices, or autonomous recording units (ARUs) placed in the field. Such data can
144 be collected over periods of time ranging from short, opportunistic, recordings, to long-term
145 deployments lasting months or years. The use of handheld microphones and ARUs are non-
146 invasive methods that do not require the capture of individual animals, and so reduce
147 disturbance and welfare impacts (Browning *et al.*, 2017; Soulsbury *et al.*, 2020; Ross *et al.*,
148 2023). Acoustic data can help with monitoring of elusive, cryptic, or nocturnal species that
149 are difficult to observe directly (Zwerts *et al.*, 2021), e.g., bats (Frick, 2013), wolves
150 (Harrington & Mech, 1982; Kershenbaum, Owens & Waller, 2019), or cetaceans (Zimmer,
151 2011). Additionally, where animals use long-distance vocalisations, ARUs are beneficial in
152 recording species over large spatial scales, e.g., crested argus pheasants (Vu *et al.*, 2023),
153 gibbons (Vu & Tran, 2019; Dufourq *et al.*, 2021), howler monkeys (Pérez-Granados &
154 Schuchmann, 2021), and wolves (Kershenbaum *et al.*, 2019). Such methods can offer
155 detection ranges in the order of several kilometres for some species, compared with tens of
156 metres for camera traps. However, as a passive technique, the obvious disadvantage of
157 acoustic monitoring is that the animal needs to be producing sound to be detected.

158 Whilst the collection of acoustic data offers many benefits and opportunities, it brings with it
159 certain challenges. First, the deployment and servicing of ARUs (e.g., replacing batteries and
160 memory storage cards) can be costly in terms of time and labour (Metcalf *et al.*, 2023b).
161 Second, although the tools for acoustic monitoring are now more widely available, cheaper in

1
2
3 162 cost, and include larger storage capacities and longer battery life (Hill *et al.*, 2019), this has
4
5 163 led to a very large increase in the quantity of data being stored, transferred, and analysed.
6
7 164 Third, a major challenge is distinguishing the sound(s) of interest from background sounds
8
9 165 which takes an enormous amount of researcher time, effort, and expertise, to recognise the
10
11 166 calls of species accurately and annotate the recordings reliably. All of this creates long delays
12
13 167 between data collection and the final results of a study, yet the need for real time results can
14
15 168 be pressing, especially in the field of conservation biology. Automatic detection can solve
16
17 169 many of these issues, as a tool to extract sounds of interest automatically, reducing or even
18
19 170 eliminating the need for manual analysis of the data.
20
21
22
23
24

25 171 1.1.1. What is automatic detection?

26
27 172 Automatic detection is the process of extracting acoustic signals from sound recordings
28
29 173 automatically, without human effort. Once detected, numerous properties of the acoustic
30
31 174 signal can be determined (with or without additional human effort). For example, the acoustic
32
33 175 signal could be classified as being produced by a particular species, its location determined,
34
35 176 and the identity of the animal inferred. The temporal and spectral properties (e.g.,
36
37 177 fundamental frequency, harmonics, modulation, etc) of the acoustic signal can be calculated
38
39 178 and used for additional processing or for inferring additional information about the sound
40
41 179 production. Some approaches implicitly combine the processes of automatic detection with
42
43 180 other tasks e.g., classification of the vocalising species, but fundamentally, the first step
44
45 181 within an automated bioacoustic processing pipeline is detection.
46
47
48
49
50

51 182 Throughout this paper, we will use the term “acoustic signal” to describe any sound or
52
53 183 acoustic event produced by an animal without regard to the purpose or intentionality of the
54
55 184 signal.
56
57
58
59
60

1
2
3 185 1.2. Scope of review
4
5

6 186 In this paper we set out to highlight and describe the emerging field of automatic detection of
7
8 187 acoustic signals as a highly interdisciplinary effort that requires expertise from both
9
10 188 biological and computer science to move forward. We present a review and tutorial that
11
12 189 addresses both the needs of the community of biologists using acoustic monitoring to answer
13
14 190 ecological, evolutionary, and conservation research questions, and the needs of computer
15
16 191 scientists developing new algorithms and implementations. As the overlap between these two
17
18 192 needs and the overlap between domain knowledge of these two groups is often small, this
19
20 193 review attempts to bridge that gap by addressing both groups simultaneously, enhancing the
21
22 194 missing knowledge of both. A reader from either field will find this review to be a useful
23
24 195 integration of both domains, providing new information to both without being inaccessible to
25
26 196 either. The review arose from an investigative workshop held in July 2023 at Girton College,
27
28 197 University of Cambridge, attended by 22 scientists from both the biological and computer
29
30 198 sciences.

31
32
33
34
35
36 199 By way of introduction, the review first presents the perspectives on automatic detection for
37
38 200 bioacoustics from the point of view of a biological researcher, aiming to instruct the
39
40 201 computer scientist in the needs of the end-user. Then, we present the perspective of the
41
42 202 computer scientist, aiming to instruct the biologist in the technologies available and their
43
44 203 limitations. There then follows a step-by-step guide to the practical implementation of
45
46 204 automatic detection, and finally a discussion of the potential future directions of the field.
47
48
49
50
51
52
53
54
55
56
57
58
59
60

205 2. BACKGROUND OF AUTOMATIC DETECTION IN BIOACOUSTICS

206 2.1. What is automatic detection and why do we need it?

207 To address the challenge of converting terabytes of acoustic recordings into useful
208 information, scientists have sought to develop techniques to automate the detection of
209 acoustic signals of interest. The traditional method of identifying the signals of interest from
210 longer acoustic recordings was to create a spectrogram and manually draw bounding boxes
211 around the signals of interest, incurring a significant cost in terms of time and expertise.
212 Fundamentally, the challenge is to replace the human annotator with computational methods
213 without a consequent loss in accuracy (Miller *et al.*, 2023). At its simplest, the aim of
214 automatic detection is to indicate segments or windows of audio which are likely to contain a
215 target sound of interest, substantially reducing the burden, even if the automated annotations
216 need then be checked by a human. The annotation label can simply be a binary label of
217 presence/absence of a sound, but this can also be further refined to classify by taxon, call-
218 type, number of individuals, etc., in increasing levels of precision and consequent difficulty
219 for both annotator and algorithm. For many species, it can be possible to identify an
220 individual through its unique vocal characteristics (Petso, Jamisola & Mpoeleng, 2021). In
221 addition to the class label, some systems also allow the position or bearing of the sound to be
222 estimated (Kershenbaum *et al.*, 2019; Smith *et al.*, 2021). Such information can then be used
223 in numerous downstream tasks such as occupancy monitoring, spatial habitat use, and
224 behavioural analysis, and automatic detection offers researchers the opportunity to scale to
225 larger spatiotemporal datasets.

226 2.2. The current state of the art in automatic detection

227 The use of automatic detection to accelerate acoustic monitoring has a long history (Acevedo
228 *et al.*, 2009; Aide *et al.*, 2013; Dufourq *et al.*, 2021; Oswald *et al.*, 2022). As an early

1
2
3 229 approach towards automation, simple techniques based on the energy within a particular
4
5 230 frequency range, characteristic to the target sound, have been used to detect signals of interest
6
7 231 (Morrissey *et al.*, 2006). However, this approach only works if the signal-to-noise ratio of the
8
9 232 target sound is sufficiently high, and if other sounds are not present in the same frequency
10
11 233 range which act to mask it. Subsequent techniques have used statistical modelling or classical
12
13 234 machine learning models such as hidden Markov models (HMMs) to detect calls that are
14
15 235 modulated in frequency and/or time (Duan *et al.*, 2013; Oswald *et al.*, 2022), by identifying
16
17 236 properties of the target sound beyond simply frequency range. Such models can provide more
18
19 237 robust and sensitive detections. More recently, there has been a strong push towards the use
20
21 238 of data-driven machine learning, exemplified by deep learning (DL), using techniques such as
22
23 239 convolutional neural networks (CNNs) (LeCun, Bengio & Hinton, 2015), recurrent neural
24
25 240 networks (RNNs) (Yu *et al.*, 2019) and more recently transformers (Lin *et al.*, 2022).
26
27 241 Transformers have been shown to obtain impressive detection accuracies, e.g. BirdNET
28
29 242 (Kahl *et al.*, 2021), and the BTO Acoustic Pipeline (Anon., 2023b).
30
31 243 There is, however, a highly fragmented landscape in the field of automatic detection – in
32
33 244 particular between the fields of computer science/machine learning, and bioacoustics/acoustic
34
35 245 ecology – and it can be very challenging for practitioners to know where to get started.
36
37 246 Should one build their own classifier from scratch, fine-tune an existing model, or simply use
38
39 247 an off-the-shelf pretrained model (Stowell, 2022a; Dufourq *et al.*, 2022b)? Good quality
40
41 248 detectors already exist in a relatively user-friendly format for birds, e.g. BirdNet (Kahl *et al.*,
42
43 249 2021); bats, e.g. BTO Acoustic Pipeline (Anon., 2023b), Kaleidoscope (Anon., 2023a);
44
45 250 cetaceans, e.g. PAMGuard (Gillespie *et al.*, 2009); rodents, e.g. DeepSqueak (Coffey, Marx
46
47 251 & Neumaier, 2019), MUPET (Van Segbroeck *et al.*, 2017). However, these detectors tend to
48
49 252 be known only by those using them in the field and are not straightforward to generalise to
50
51 253 other taxa without retraining or altering the model architecture or assumptions. There is also
52
53
54
55
56
57
58
59
60

1
2
3 254 an imbalance with some taxa being better represented than others in terms of the availability
4
5 255 of detectors. The process of building or fine-tuning a new deep learning model for a
6
7 256 practitioner's particular habitat and species of interest is non-trivial and involves several tasks
8
9
10 257 such as cloning repositories from Github, designing data-loaders, and training models on
11
12 258 specialised computing hardware such as Graphical Processing Unit (GPU) clusters. This
13
14 259 serves as a major barrier to widespread adoption of these new techniques unless a tame
15
16 260 computer scientist can be persuaded to assist in the process. In contrast, the more mature field
17
18 261 of automatic detection in camera trapping, e.g. WildlifeInsights, CameraTrapDetectorR
19
20 262 (Hendry & Mann, 2018); Camelot (Hendry & Mann, 2018); Agouti (Casaer *et al.*, 2019);
21
22 263 MegaDetector, can serve as an exemplar for deriving best practices, as existing tools are easy
23
24 264 to use for non-programmers, and easily generalised to different taxa.
25
26
27
28
29

30 265 2.3. What do we aspire for from automatic detection? 31

32
33 266 Despite the challenges associated with the automatic detection of acoustic signals, rapid
34
35 267 advances in machine learning are starting to bring this concept into reality. Although the
36
37 268 context under which acoustic data is recorded and its end use will differ, the common
38
39 269 requirement is for algorithms that take acoustic data as an input, and then detect and return
40
41 270 extracted sounds as the output. Some users may only require outputs of particular target
42
43 271 sounds, such as a specific species or anthropogenic sound, whereas others may require all
44
45 272 sounds to be classified. Ideally, the ultimate end goal of automatic detection for biologists
46
47 273 would be a universal, off-the-shelf algorithm capable of detecting and classifying all animal
48
49 274 vocalisations such that anybody, including those without any training in computational
50
51 275 methods, could process their acoustic data more efficiently and flexibly tailor it to their
52
53 276 particular use-case (Romero-Mujalli *et al.*, 2021). Where an off-the-shelf detector for a sound
54
55
56
57
58
59
60

1
2
3 277 of interest is not readily available, algorithms that are easy to train with a relatively small
4
5 278 amount of data and minimal annotation effort should be the aim.
6
7
8

9 279 3. PERSPECTIVES FROM BIOLOGICAL SCIENCES

10
11
12 280 In this section, we give, largely for the benefit of the reader from a computer science or other
13
14 281 non-biological background, an overview of the possible roles for bioacoustics in addressing
15
16 282 several important evolutionary, ecological, and conservation questions, highlighting the
17
18 283 potential benefit that automatic detection can provide.
19
20
21
22

23 284 3.1. Overview of uses of automatic detection in the biological sciences

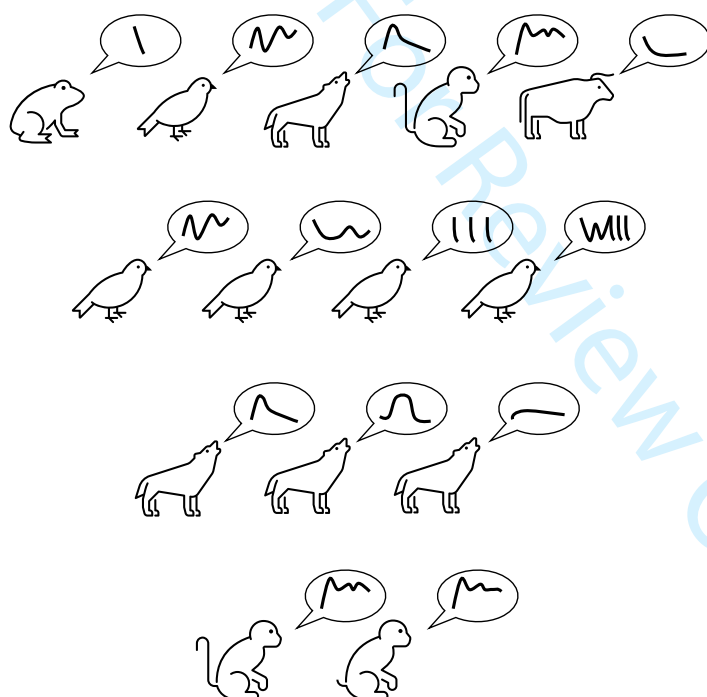
24
25 285 Detecting acoustically active animals through their acoustic signals can provide a wealth of
26
27 286 information that is important to conservation biology, ecology, evolutionary biology, animal
28
29 287 behaviour, and welfare (McLoughlin, Stewart & McElligott, 2019; Odom *et al.*, 2021; Erbe &
30
31 288 Thomas, 2022). Often, these areas of study can overlap: animals can produce sounds to
32
33 289 influence the behaviour of others in a wide range of contexts, e.g., to attract a mate or warn
34
35 290 off an intruder, or as a by-product of other behaviours, e.g., the sound of wings flapping or
36
37 291 footsteps.
38
39
40

41
42 292 Historically, conservation efforts and biodiversity surveys have been skewed towards species
43
44 293 that are easy to trap or track across the landscape, often depending on direct observation or
45
46 294 finding physical traces like scat or hair (Boakes *et al.*, 2010). However, the field of
47
48 295 bioacoustics allows us to survey remote or otherwise inaccessible areas, e.g., deep sea
49
50 296 environments, arctic and antarctic regions, and rainforests (Staaterman *et al.*, 2017), with
51
52 297 research often focusing on the loud and persistent calls of target species to detect their
53
54 298 presence. Like camera trapping, bioacoustics generates large datasets which challenge
55
56
57
58
59
60

299 analysis, but, unlike camera traps, the same event can be recorded in multiple places,
 300 multiplying the data to be assessed and analysed.
 301 Below, we provide a broad review of the use of acoustic data in the biological and ecological
 302 sciences, from measures of biodiversity at geographic scales to tracking the movements and
 303 behaviours of individual animals, and highlight how automatic detection can increase the
 304 efficiency and efficacy of monitoring.

305

306



Ecosystems and acoustic indices

Measuring acoustic variation and diversity across many different species in the environment.

Species repertoire

Measuring the range of different acoustic signals produced by a single species.

Populations and dialects

Measuring acoustic variation between different populations of the same species.

Individual identity

Identifying individuals by differences in their acoustic signals.

307

308 Figure 1: Hierarchy of acoustic signal specificity

309

310 3.1.1. Ecosystems and acoustic indices

311 Any multi-species soundscape will consist of a wide range of frequencies being used by
 312 different species in the same environment (Krause, 1993). To maximise the chance that their
 313 signal will be detected, animals usually avoid acoustic signal interference by vocalising in
 314 different frequency ranges or at different times, as described by the acoustic adaptation

1
2
3 315 hypothesis (Hansen, 1979; Rothstein & Fleischer, 1987). This ecological phenomenon makes
4
5 316 it possible to detect particular clades or species. It also means that estimates of biodiversity
6
7 317 can be made based on the number of different acoustic signals being produced at different
8
9 318 times/frequencies.

10
11
12 319 Acoustic indices provide a quantitative measure of acoustic complexity by analysing
13
14 320 variation in the frequency and timing of acoustic signals, rather than identifying individual
15
16 321 sounds. Such indices offer metrics for wildlife monitoring and assessment, characterising
17
18 322 biological communities through sound (Sueur *et al.*, 2014; Buxton *et al.*, 2018). While
19
20 323 acoustic indices are informative about the acoustic complexity or general biodiversity of a
21
22 324 landscape, they are less useful for deriving specific information about species or the
23
24 325 individuals vocalising.

25
26
27 326 Acoustic indices typically do not use automatic detection and classification of acoustic
28
29 327 signals, as, by their nature, they characterise the soundscape as a whole. However, automatic
30
31 328 detection of sound classes, for example distinguishing acoustic signals of anthropogenic
32
33 329 origin from those of biological origin, can improve the ability of acoustic indices to provide
34
35 330 useful indications of biological activity (Narasimhan, Fern & Raich, 2017; Fairbrass *et al.*,
36
37 331 2019; Clark *et al.*, 2023). Thus, effective automatic classification of acoustic signals may
38
39 332 become an important element of improving acoustic indices in future research.

40 41 42 333 3.1.2. Species occupancy and density

43
44
45
46 334 Occupancy modelling is the statistical analysis of the patterns and dynamics of a species in a
47
48 335 given space over time (MacKenzie *et al.*, 2003), which can be informed by acoustic signals
49
50 336 (Wood & Peery, 2022; Cole *et al.*, 2022). Bioacoustic occupancy monitoring can provide
51
52 337 critical information on the presence and absence of species and the dynamics of the
53
54 338 ecosystem, particularly for cryptic or elusive species.
55
56
57
58
59
60

1
2
3 339 Population density estimates model a species' abundance within a defined space. Density
4
5 340 estimates are an extremely important tool for assessing spatiotemporal population changes
6
7 341 that can be the result of declining prey numbers, land-use change, human-wildlife conflict
8
9 342 (Wolf & Ripple, 2016; Ogutu *et al.*, 2016; Rostro-García *et al.*, 2023), or other factors, and
10
11 343 bioacoustics data can provide an important tool for estimating the densities of animal
12
13 344 populations.
14
15
16
17

18 345 3.1.3. Spatial analyses

19
20
21 346 Population surveys and behavioural investigations often need to be able to determine the
22
23 347 location and/or movement patterns of animals. Bioacoustic surveys have been used in more
24
25 348 recent years to supplement or replace previous tracking methods (Frommolt & Tauchert,
26
27 349 2014). For example, the tracking of migratory species across their extensive ranges, where
28
29 350 radio/satellite telemetry is only useful if the individuals tagged with a transmitter survive
30
31 351 what may be a high mortality journey, can benefit from the application of bioacoustic
32
33 352 techniques. While telemetry is an effective method for learning about a species' movement, it
34
35 353 can also be highly invasive, can affect the behaviour of individuals being trapped, and is not
36
37 354 always suitable for all species/age groups, e.g. species that are too small to carry the weight
38
39 355 of a transmitter, or species in remote areas (Sharpe *et al.*, 2009).
40
41
42
43

44 356 Depending on the intended research goals, it may be sufficient simply to detect the
45
46 357 presence/absence of an animal within a recorder's range (macro-localisation), or one may
47
48 358 need to infer the exact position of an individual (micro-localisation). There are benefits and
49
50 359 limitations to each: macro-localisation can inform on occupancy, habitat suitability, territory
51
52 360 use, and migratory patterns. On the other hand, with a significant increase in the complexity
53
54 361 of use, advanced tools also allow a more targeted approach such as multilateration, where the
55
56 362 exact individual's location is calculated based on the time difference of arrival (TDOA) of an
57
58 363 acoustic signal to multiple recording devices (Mennill *et al.*, 2012). Such micro-localisation
59
60

1
2
3 364 allows avoiding double counting for density estimates, can inform on animal movement

4
5 365 speed and direction, as well as providing fine grained territory boundaries, but requires

6
7 366 additional downstream processing to carry out the localisation analysis.

8
9
10 367 Estimation of a focal animal's home range and territory provides wildlife managers with a

11
12 368 boundary for their activity (Powell, 2000), permitting the study of intraspecific dynamics and

13
14 369 spatial distribution of individuals across a landscape (Burgos & Zuberogitia, 2020), which

15
16 370 can be important for conservation action. Real-time automatic detection combined with

17
18 371 localisation reduces the research effort required for follow-up visual observation and can

19
20 372 obviate the need for visual observation entirely.

21
22
23
24
25 373 3.1.4. Species characteristics

26
27 374 Automatic detection of acoustic signals is complicated by the fact that there are relatively few

28
29 375 species that, like the American toad (*Anaryxus americanus*), produce a single call (Bee,

30
31 376 2012), while many species produce multiple call-types, e.g., the northern mockingbird

32
33 377 (*Mimus polyglottos*) produces hundreds of different song types (Derrickson, 1988). Thus,

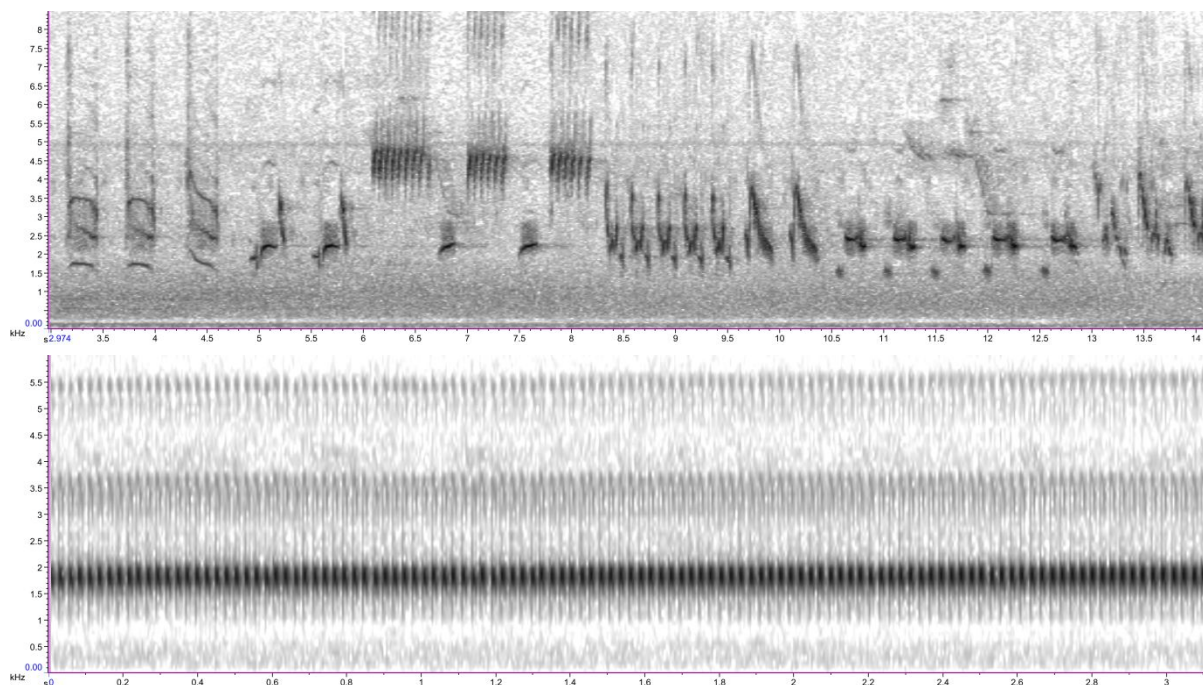
34
35 378 while it is relatively easy to link a croak to the presence of a toad, it can be more challenging

36
37 379 to capture all the potential acoustic signals of the mockingbird. This is further complicated if

38
39 380 the species' different call-types need to be classified beyond simple detection (e.g., as contact

40
41 381 calls vs alarm calls).

42
43
44
45
46 382



383

384 Figure 2. Top: the varied mimicry of the northern mockingbird (*Mimus polyglottus*),
385 composed of varied songs of other species, which would be difficult to detect in a general
386 way. Bottom: the call of the American toad (*Anaxyrus americanus*), which only produces a
387 single call, repeated for long periods. The spectrograms show time on the x-axis and
388 frequency on the y-axis.

389

390 Collectively, all the distinct call types a species produces can be defined as the vocal
391 repertoire. The size of the repertoire may be thought of as a simple proxy for vocal
392 complexity (Bouchet, Blois-Heulin & Lemasson, 2013; Manser *et al.*, 2014), and the
393 structure of the repertoire (e.g., how often call-types are used and interpretations of the
394 potential uses) are important for describing a species' behavioural ecology. Therefore, both
395 general acoustic signal detection ("the target species made a sound in some way") and
396 specific call-type detection ("the leopard-specific alarm call has been produced") are useful to
397 different studies and these analyses can be nested. Comparisons of vocal complexity between
398 species, species groups, and taxa (Kershenbaum *et al.*, 2021; Leighton & Birmingham, 2021)
399 may enable research into broad evolutionary or ecological questions, such as cognitive

1
2
3 400 abilities, adaptive advantages of cognitive skills, or the evolution of language (McComb &
4
5 401 Semple, 2005; Dunn & Smaers, 2018).
6
7 402 The more varied and less stereotyped calls are, the larger the challenge to automatic
8
9 403 detection. However, the implications of variability within a single call type on the
10
11 404 performance of automatic detection and classification have not been adequately investigated.
12
13
14
15

16 405 3.1.5. Populations and social groups

17
18 406 The same species can show variation in their vocalisations among social groups and or across
19
20 407 geographic regions. Research into these differences can offer unique insight into either
21
22 408 phylogenetic patterns, including speciation (Meyer *et al.*, 2012; Riesch *et al.*, 2012; Heaphy
23
24 409 & Cain, 2021), historic geographical patterns (Laiolo *et al.*, 2001; Kershenbaum *et al.*, 2012;
25
26 410 Hebets *et al.*, 2021), or differences between social groups (Ford, 1991; Velásquez *et al.*,
27
28 411 2013; Garland, Castellote & Berchok, 2015; Kershenbaum *et al.*, 2016b). Automatic
29
30 412 detection can scan through long-term recordings to unveil temporal and cultural variations of
31
32 413 vocal behaviours, for example in whales (McDonald, Hildebrand & Mesnick, 2009; Garland
33
34 414 *et al.*, 2011; Best *et al.*, 2022).
35
36
37
38
39
40

41 415 3.1.6. Individual characteristics

42
43 416 It may be important to identify individual animals and/or characterise the traits or states of
44
45 417 individuals of a target species, such as age, sex, body size, emotional valence/arousal, and
46
47 418 physiology. Acoustic signals can potentially encode all of this information. Examples of the
48
49 419 benefits of individual identification include gaining insights into the evolution and ecology of
50
51 420 a species, such as life history stages and social structure (Clutton-Brock & Sheldon, 2010);
52
53 421 facilitating conservation efforts, for example tracking movement of critically endangered
54
55 422 species in the landscape (McCloughlin *et al.*, 2019); and improving management in captivity,
56
57 423 for example measuring vocal activity as an indicator of welfare in zoo housed animals
58
59
60

1
2
3 424 (Castellote & Fossa, 2006). A diverse range of species' calls have been found to encode
4
5 425 individual identity from birds (Fox, Roberts & Bennamoun, 2008; Martin *et al.*, 2022) to
6
7 426 cattle (Green *et al.*, 2019), cetaceans (Kershenbaum, Sayigh & Janik, 2013; Böttcher *et al.*,
8
9 427 2018), and frogs (Qian *et al.*, 2023).

10
11
12 428 Individual identification provides an open scope for spatiotemporal monitoring of species
13
14 429 without tagging (Aide *et al.*, 2013), while also offering the opportunity for population
15
16 430 estimation using mark-capture recapture methods, which rely on individual identification
17
18 431 (Marques *et al.*, 2013b; Buxton *et al.*, 2018).

19
20
21 432 Acoustic signals can be used in a wide range of species to assess the intensity (high to low)
22
23 433 and valence (positive to negative) of emotional arousal of animals, which in turn can be used
24
25 434 as an estimate of welfare in animals in captivity (Volodina & Volodin, 1999; Clark & Dunn,
26
27 435 2022) and farms (Manteuffel, Puppe & Schön, 2004). Inferring emotional arousal from
28
29 436 acoustic signals also allows for the assessment of 'positive welfare' in animals (Laurijs *et al.*,
30
31 437 2021), and it is possible to monitor farm animals for the onset of disease, e.g., pigs (*Sus*
32
33 438 *domesticus*) (Exadaktylos *et al.*, 2008; McLoughlin *et al.*, 2019) and chickens (*Gallus*
34
35 439 *domesticus*) (Mao *et al.*, 2022).

40 41 440 3.2. Key challenges

42
43
44 441 As outlined above, many studies in ecology and evolution require relatively precise
45
46 442 identification of the type of acoustic signal, for example different call types, the source of the
47
48 443 sound, individual identification, or the localisation of the source of the sound in space.

49
50
51 444 Despite the huge potential of automatic detection to answer these challenges, the field is still
52
53 445 facing significant barriers during implementation in biological studies, ranging from
54
55 446 limitation in infrastructure, lack of training, inaccessibility of methods, and practical
56
57 447 limitations in the field. For example, field recordings are often not of optimal recording
58
59
60

1
2
3 448 quality and have a low signal to noise ratio. Even under ideal conditions, acoustic signals
4
5 449 themselves may be highly varied and irregular, with low stereotypy and a high degree of
6
7 450 variability between individuals and groups, or geographical dialects (Nelson, 2000), all of
8
9 451 which can present a challenge for automatic detection. The broad implementation of
10
11 452 automatic detection requires that the model is robust to the variation presented.
12
13

14 453 The training of models requires data to be robustly identified and correctly attributed to the
15
16 454 study species or individuals, often produced by visual observation of the callers. Collecting
17
18 455 these data can be challenging as, for instance, individuals may remain visually cryptic, or call
19
20 456 only at certain times. Thus, ground-truthing data requires high quality, reliably identified call
21
22 457 datasets which can be difficult to obtain but are essential. Furthermore, generalising data
23
24 458 from captive animals or in unique circumstances might give rise to misleading results. Thus,
25
26 459 robust identification of large datasets is rare but essential and should be a focus for future
27
28 460 research.
29
30
31
32
33
34

35 461 4. TECHNICAL PERSPECTIVES

36 37 38 462 4.1. Perspectives from computer science

39 40 41 42 463 4.1.1. The role of computation in automatic detection

43
44 464 Advanced computational methods can provide solutions to a wide range of bioacoustics
45
46 465 problems. For example, acoustic signals of interest can be merely detected (i.e., the start and
47
48 466 end times identified), or additional information can be extracted, such as classification of
49
50 467 signal type, or location of the sound source. If different types of acoustic signal are present,
51
52 468 they can be grouped into multiple classes, which might represent different species, or
53
54 469 different call types within a single species. Even when a single type of acoustic signal is
55
56 470 present, the task of counting the number of such events or sub-elements of the events is often
57
58
59
60

1
2
3 471 non-trivial (e.g., the different notes in a birdsong, or the individual barks of a dog). Therefore,
4
5 472 the role of automatic detection and automatic processing of bioacoustic data is a broad field,
6
7
8 473 with many possible applications.
9

10 474 Computational methods can help with any task which can be clearly defined. One way to
11
12 475 define the task is through explicit rules (an engineering approach), for example, to specify
13
14 476 that a target acoustic signal occurs solely and uniquely in a certain range of frequencies.
15
16
17 477 Alternatively, a set of examples can be provided to the algorithm (a machine-learning
18
19 478 approach), and the algorithm is trained to generalise those examples to detect successfully
20
21 479 when presented with novel examples. In the case of automatic detection, some tasks are
22
23
24 480 simple enough that a good method can be designed directly using the engineering approach:
25
26 481 this typically happens with situations of highly-stereotyped sounds, where template-matching
27
28 482 often works well (Barker, Herrera & West, 2014), or low-noise environments with few
29
30 483 interfering sounds, where energy-detection may work well (Hood, Flogeras & Theriault,
31
32
33 484 2016).

34
35 485 When the target sounds, or the background, are more complex—such as with recordings of
36
37 486 elaborate bird song or soundscapes with high levels of anthropogenic noise—then machine
38
39 487 learning (ML) is of benefit. As noisy problems can rarely be defined in a clear-cut
40
41
42 488 “engineering” way, ML attempts to reach a solution by generalising from a set of examples
43
44 489 instead. Although ML has been investigated for many years (Towsey *et al.*, 2012), it is the
45
46 490 era of deep learning that now makes many bioacoustic detection tasks achievable (Stowell,
47
48
49 491 2022b). It is still important to define the task to be solved clearly – by curating good datasets
50
51 492 for training and evaluating systems, and by specifying the input and output data formats.
52
53
54 493 Input data format, in bioacoustic applications, is generally some representation of the sounds
55
56 494 recorded, whereas the output format is defined by the nature of the “answer” that the system
57
58 495 is trained to supply, e.g., species presence, individual, call type, etc.
59
60

1
2
3 496 Data curation aside, the power of ML comes from having techniques that can “train”
4
5 497 (optimise) the system to achieve a particular goal, and so the output data format matters
6
7 498 because it is closely tied to this procedure of optimisation. If the output format is a yes/no
8
9
10 499 answer about species presence, this is the same format as a *binary classification* task in ML
11
12 500 and can be addressed directly by training a classifier (Stowell, 2022b), which takes sound as
13
14 501 input, and outputs a corresponding indicator: present/absent. Very often, however, the output
15
16 502 format wanted is more complex; for instance, given a long audio recording as input, we may
17
18 503 want to output a list of (predicted) events giving each event’s start and end time, and
19
20 504 optionally its frequency range as well. Note that this is quite similar to “object detection” in
21
22 505 image recognition, and indeed, most bioacoustic research uses spectrograms as a visual
23
24 506 representation of a sound, rather than working with the sound directly. In this case, we may
25
26 507 typically be looking for a list of “bounding boxes” along the time axis or in time-frequency,
27
28 508 leading some to directly adapt image object detection algorithms to spectrograms
29
30
31 509 (Kershenbaum & Roch, 2013; Venkatesh, Moffat & Miranda, 2022; Wu *et al.*, 2022).
32
33 510 When a ML model has been trained, better results may be obtained if the model is applied in
34
35 511 the same conditions as the training data, i.e., “in-domain” as opposed to “out-of-domain” data
36
37 512 (Best *et al.*, 2020). For example, conditions might be “in-domain” if they have the same
38
39 513 background conditions, microphone type, and sampling protocol as in the training data.
40
41
42
43
44
45

46 514 4.1.2. State of the art in automatic detection methods

47
48 515 No algorithm will generalise perfectly to all situations: the choice of training data represents
49
50 516 the choice of intended domain. Classic machine learning advice would be to avoid "out-of-
51
52 517 domain" situations. Yet many taxa do not benefit from such a large amount of prior work as
53
54 518 on birds. Could we nevertheless make use of off-the-shelf models from similar tasks, or must
55
56 519 we start building a large new dataset?
57
58
59
60

1
2
3 520 Happily, a recent widespread trend is “transfer learning”, in which one or more pretrained
4
5 521 models are used that have been trained on tasks that are different from (but usually related to)
6
7
8 522 the original domain: for example, we could consider models trained on human speech
9
10 523 recognition. The models are then re-used for the current application (i.e., animal
11
12 524 vocalisations), and it is often found that the original learning makes training the model on the
13
14
15 525 current data more effective (Zhuang *et al.*, 2021).

16
17 526 A common approach to transfer learning, known as fine-tuning, consists of modifying only a
18
19 527 small subset of parameters and adapting the inputs and/or outputs format. The modification
20
21 528 requires training the model on a new set of examples, made up of audio recordings and
22
23
24 529 corresponding annotations. This procedure is computationally much lighter than performing
25
26 530 the process from scratch. It also requires fewer labels since it exploits many of the regularities
27
28
29 531 in the initial data set. As a rule of thumb, one may try to choose a base model which has been
30
31 532 trained on similar target sounds or background noise, e.g. BirdNet (Kahl *et al.*, 2021). Yet we
32
33 533 have observed successful attempts in adapting models from significantly different acoustic
34
35 534 data, even from different frequency ranges (Çoban *et al.*, 2020; Sethi *et al.*, 2020; Leroux *et*
36
37
38 535 *al.*, 2021; Sarkar & -Doss, 2023).

39
40 536 When using transfer learning (also known as “pretrained” models), special care must be
41
42 537 taken. The model must be applied to acoustic data that closely resemble the data on which it
43
44
45 538 has been trained. The user must reflect on details such as matching sampling rates,
46
47 539 normalisations, SNR-levels, and duration of the input audio segments. Usually, the producers
48
49 540 of such models will have trained models on diverse data to ensure generalisation. However,
50
51 541 optimal performance is achieved when staying within the region of operation for which the
52
53
54 542 model was designed.

55
56 543 In this paradigm, the algorithm trained on a different system can be considered to perform a
57
58 544 role similar to the role of the spectrographic representation in aiding human interpretation of
59
60

1
2
3 545 sounds. In the same way that a spectrogram or filterbank takes a sound waveform and
4
5 546 presents it in a different format (and one where the important features are easy to detect by
6
7 547 eye), so a model trained on a different species, for example, cannot detect the target species
8
9 548 well, but may nonetheless produce an output (known as extracted acoustic features) that can
10
11 549 be used as the input to train another model, which will then be more successful in finding the
12
13 550 focal species. In the ML literature the resulting features are often referred to as embeddings
14
15 551 or latent representations. Unlike traditional acoustic features like a spectrogram these
16
17 552 embeddings are often difficult to interpret on their own. They are the result of a large
18
19 553 composition of complex functions whose parameters have been optimised to solve a
20
21 554 particular task such as classifying an acoustic scene or discriminating from a given set of
22
23 555 videos the one that matches a particular sound.

24
25
26 556 Despite this, fine-tuning alone may not be sufficient to obtain the desirable level of accuracy.
27
28 557 We may then further adapt the model to our specific needs by retraining all of its parameters
29
30 558 on the acoustic data of interest. One must take into consideration that these models have been
31
32 559 designed with a large number of parameters (317 million parameters for the large version of
33
34 560 HuBERT for instance; (Hsu *et al.*, 2021), to be optimised on thousands of hours of audio.
35
36 561 When trained on a small number of examples this may quickly lead to overfitting, where the
37
38 562 model will work as expected on the data presented during training but will fail to produce
39
40 563 satisfactory predictions for unseen audio examples.

41
42 564 Even when many hours of field recordings are available, it is not clear if the acoustic data
43
44 565 will be sufficiently diverse to produce acoustic features that will be performant enough for
45
46 566 downstream tasks such as the detection of vocalisations. In other words, if a bioacoustic
47
48 567 dataset does not contain any useful (or additional) information which could be reemployed in
49
50 568 the downstream detection tasks, then retraining the pre-trained model might not improve
51
52 569 performance. Furthermore, re-training these models on large amounts of data is usually a
53
54
55
56
57
58
59
60

1
2
3 570 tedious task which calls for the expertise of trained computer scientists and access to costly
4
5 571 computational resources such as GPU clusters.
6
7
8 572 The approach of adapting transfer learning models to automatic bioacoustic detection, can
9
10 573 still be carried out by pretraining models on bioacoustic data directly, instead of human
11
12 574 speech or generic audio. It has been shown to yield interesting results in the downstream
13
14 575 detection performances for a variety of species (Hagiwara, 2022), but much work still needs
15
16 576 to be done in this area. The success of this method relies on the availability of large datasets
17
18 577 which could allow for the pretraining of a single, large-scale, multispecies foundation model.
19
20 578 As is the case in the speech processing and image recognition domains, making such a model
21
22 579 available to the bioacoustic community could then allow for efficient user-friendly classifiers
23
24 580 to be trained in few-shot learning contexts within a unified pipeline.
25
26
27
28 581 An alternative to the transfer learning approaches is to use smaller models, with fewer
29
30 582 parameters that may be trained entirely on the target audio data. For example, an algorithm
31
32 583 called TweetyNet is designed for detecting/segmenting bird vocalisations in a laboratory
33
34 584 context, based on a CNN to be trained specifically for each target bird; the package includes
35
36 585 an interface to simplify that training process (Cohen *et al.*, 2022); DeepSqueak can do the
37
38 586 same for rodent vocalisations (Coffey *et al.*, 2019). Those algorithms directly train the CNN
39
40 587 as a classifier/detector. Another approach used by many in the bioacoustics community is to
41
42 588 train a so-called ‘auto-encoder’ on the dataset of interest to extract deep feature
43
44 589 representations from unlabelled data. This unsupervised approach consists in optimising a
45
46 590 neural network to compress an audio snippet into a numerical vector which is decompressed
47
48 591 to reconstruct the original sound. This technique has been applied to call categorisation in a
49
50 592 variety of species (Sainburg, Thielk & Gentner, 2020; Best *et al.*, 2023).
51
52
53
54 593 Even using such methods, it is common that bioacoustic datasets are not large enough to train
55
56 594 an ML detector well, or that some categories/contexts are underrepresented in the training
57
58
59
60

1
2
3 595 data. It is thus common (and recommended) to use ‘data augmentation’ to assist with this:
4
5 596 ‘new’ training examples can be created by small modifications of existing ones. This has
6
7 597 been widely investigated and found to improve performance, to a similar extent as the use of
8
9 598 pretrained networks (Lostanlen *et al.*, 2018).

10 599 The bioacoustics community often faces complex scenarios with sound events overlapping
11
12 600 both in time and frequency (e.g., dawn chorus) or with highly non-stationary background
13
14 601 noise (e.g., urban scenes). These require more advanced and specific solutions that tackle the
15
16 602 problem of working with mixtures of sounds. Data-augmentation techniques serve this
17
18 603 purpose by artificially constructing similar data for which annotations can be created by
19
20 604 design (Jansson *et al.*, 2017; Zhang *et al.*, 2018; Wisdom *et al.*, 2020). These approaches
21
22 605 have been applied to improve performance on up to ten simultaneously-calling bird species in
23
24 606 a simulation study (Parrilla & Stowell, 2022) and in real recordings with significantly fewer
25
26 607 simultaneous calls (Denton, Wisdom & Hershey, 2021; Bermant, 2021).

34 608 4.1.3. Assessing pre-existing models

35
36 609 The fast pace at which the ML community publishes new pretrained models renders them
37
38 610 outdated quickly. The availability of accessible learning resources for some models makes
39
40 611 them a go-to solution for many practitioners, despite having been superseded by other
41
42 612 options. Model publishers should document their work in a way approachable by non-experts
43
44 613 if they aspire to have an important impact on the bioacoustic community. On the other hand,
45
46 614 users of these models may consult the latest benchmarks and challenges that target diverse
47
48 615 applications of audio ML representations. For instance, HEAR (Turian *et al.*, 2022)
49
50 616 benchmarked multiple state-of-the-art methods on a varied set of tasks in speech, music and
51
52 617 environmental sounds. More recently BEANS (Hagiwara *et al.*, 2022) proposes a benchmark
53
54 618 specific to bioacoustics where representations are tested on detection and classification tasks
55
56 619 of several species.

1
2
3 620 4.2. Conclusions of the technical constraints on the current uses, limitations and
4
5 621 expectations of automatic detection
6
7
8 622 Automatic detection has been used for density estimation (McDonald & Fox, 1999; Marques
9
10 623 *et al.*, 2013a), occupancy (Dawson & Efford, 2009), species presence (Obrist *et al.*, 2010),
11
12 624 trends (Abrahams & Geary, 2020), and phenology, e.g., the start of breeding, or daily onset of
13
14 625 song (Willacy, Mahony & Newell, 2015; Oliver *et al.*, 2018). This technology can be used in
15
16 626 conjunction with other non-invasive monitoring methods such as camera traps, scat surveys,
17
18 627 hair collection, and human observation (Long, 2008), providing additional information and
19
20 628 allowing monitoring of otherwise cryptic species that might elude detection. There should be
21
22 629 ongoing conversations between biologists and computer scientists, bidirectional and iterative,
23
24 630 improving the survey quality, accuracy, and algorithm usability over time. Biologists can
25
26 631 provide the ground-truthing and validation of the use of automatic detection, while computer
27
28 632 scientists can develop the system and work with them to iteratively improve the automatic
29
30 633 detection system.
31
32
33
34
35 634 While we have argued for the widespread use of automatic detection systems, there are
36
37 635 limitations, and these should be considered at the start of a project. Some of these are self-
38
39 636 evident: signals that do not rise above background noise will be lost as undetectable. Also,
40
41 637 signals can be difficult to separate if they overlap with either intraspecific, interspecific, or
42
43 638 unrelated sounds, as in the dawn chorus when birds sing with many overlapping, very similar
44
45 639 elements, making extraction/detection of a single unit difficult. Dataset sizes (for both
46
47 640 training and deployment) may be a limiting factor. We have referred to data augmentation
48
49 641 and denoising to synthetically account for data limitations. These and other tools (e.g., data
50
51 642 imputation, generative deep learning) are often helpful, but the results are unlikely to be as
52
53 643 reliable or unbiased as they would be with a large representative dataset. They should not be
54
55 644 relied upon as a silver bullet when recordings are rarely observed, noisy or otherwise hard to
56
57
58
59
60

1
2
3 645 analyse. Just like with human annotation, automatic detection will always be subject to some
4
5 646 level of bias and inaccuracy; one advantage of automatic systems is that these factors can be
6
7 647 numerically evaluated. Automatic detection model predictions are only ever as good as the
8
9 648 input training data. Annotations which are not accurate or have not been conducted
10
11 649 appropriately for the intended application may worsen the efficacy of the model.
12
13 650 Furthermore, density estimation relies on the choice of robust thresholds for confidence in
14
15 651 attribution of sounds. There can be an accumulation of errors over time if the thresholds are
16
17 652 chosen either to be too low or too high, discarding weak identifications wrongly, or placing
18
19 653 too much confidence in others. Finally, all acoustic detection relies on the sound event
20
21 654 occurring, and often species may choose to not vocalise or create a sound and thus can be
22
23 655 missed. What is not heard cannot be counted. However, despite these caveats, we believe that
24
25 656 automatic detection and PAM offer the opportunity to collect and analyse data that cannot be
26
27 657 processed by other means, providing an exciting and valuable new tool for the biological
28
29 658 sciences.

30 31 32 33 34 35 36 37 659 5. A PRACTICAL GUIDE TO AUTOMATIC DETECTION

38
39 660 We now present a practical guide for using automatic detection. There are many decisions
40
41 661 that we must make when designing a study that uses automatic detection, and our goal is to
42
43 662 help practitioners optimise these decisions. We realise that some of these decisions may be
44
45 663 constrained by access to financial resources, lack of training in bioacoustics, limited technical
46
47 664 skills in coding and machine learning, and/or lack of access to high-speed internet for cloud
48
49 665 storage and computing. These limitations may be particularly pronounced for researchers in
50
51 666 the Global South. We acknowledge that there is still much to be done to make these tools and
52
53 667 approaches accessible for all.
54
55
56
57
58
59
60

1
2
3 668 This guide is developed to help users implement an ‘off-the-shelf’ automatic detection
4
5 669 approach, or for developing or adapting their own approach. We strongly advocate that
6
7 670 practitioners implement a pilot study to ensure the approach they plan to use is feasible
8
9 671 before embarking on a large-scale endeavour. Importantly, even with the most sophisticated
10
11 672 automated approach, a substantial amount of human investment is needed to create training
12
13 673 datasets, evaluate detector performance, and verify the detections.
14
15
16
17

18 674 5.1. Define research questions

20
21 675 The most important thing to consider when using automatic detection is the specific research
22
23 676 question. For example, if you are interested in detecting the presence or absence of a rare
24
25 677 signal (e.g., a gunshot or the presence of an endangered species) then you will want to use an
26
27 678 approach that will ensure high recall (i.e., high probability of detection) and you may tolerate
28
29 679 a relatively high number of false positives. Alternatively, if you are interested in subsequently
30
31 680 classifying individuals from the detections, you may prefer to focus on retaining high signal
32
33 681 to noise ratio (SNR) calls and will tolerate lower recall with higher precision. Your research
34
35 682 question will influence every decision you make in the automatic detection workflow,
36
37 683 including study design, data collection and the analytical approach. For guidance on study
38
39 684 design, we point readers to (Sugai *et al.*, 2020).
40
41
42
43
44
45

46 685 5.2. Study design

47
48 686 Depending on the nature of the research question, researchers will need to determine their
49
50 687 study design, including hardware needs, recording schedule and whether the processing of
51
52 688 data will be carried out in real time or at a later date (or, ‘offline’). For instance, for the
53
54 689 detection of a single species, researchers may deploy ARUs over the landscape for a period of
55
56 690 time and then download the data onto a hard drive to be processed offline. The recording
57
58 691 schedule also needs to be determined according to the research goal. We refer the readers to
59
60

1
2
3 692 more extensive discussions of this issue for further details (e.g., (Browning *et al.*, 2017;
4
5 693 Metcalf *et al.*, 2023a). Real-time processing is an emerging area, but due to the limitations of
6
7
8 694 placing power-efficient computation in the field, real-time automatic detection typically is
9
10 695 more bespoke and less accurate than offline processing.

14 696 5.3. Start with a pilot study (if possible)

16 697 Given the costs (both financial and in human labour) of implementing projects that use
17
18 698 automatic detection, we strongly advocate that researchers start with a small scale setup to
19
20 699 test out their planned approach. For a large-scale PAM study, deploying a few recorders over
21
22
23 700 a smaller spatial scale and a shorter time period may provide enough acoustic data to get
24
25 701 started with automatic detection. If the signals are relatively rare (e.g., gunshots) perhaps
26
27
28 702 finding online repositories or datasets of samples would be necessary. A well-designed pilot
29
30 703 study will help researchers make informed decisions about annotations, choosing an
31
32 704 automated detector, and reporting and interpreting their results.

36 705 5.4. Data collection and archiving

39 706 Data storage and archiving remains challenging, since the large data volume of the raw audio
40
41 707 in many projects often goes beyond the limits of free or easily-available services.
42
43 708 Furthermore, (Metcalf *et al.*, 2023a) recommend backing up audio data in multiple copies,
44
45 709 and also making use of cloud storage. Simply storing the audio is typically only part of the
46
47
48 710 issue: you and your collaborators will also need to access it, for example to visualise or to
49
50 711 apply an algorithm to the dataset, which means that speed of upload and download
51
52 712 (bandwidth) may be an equal or greater concern. Cost of storage and bandwidth are often
53
54 713 significant questions. Arbimon (Ganchev, 2020) is one project that aims to store and share
55
56 714 large volumes of wildlife audio on behalf of others.
57
58
59
60

1
2
3 715 Reducing data sizes can be achieved in many ways, including audio file compression and
4
5 716 data subsampling. Lossless compression (such as FLAC) can reduce file size without losing
6
7 717 information; lossy compression (such as MP3 or AAC) will discard at least some information
8
9 718 from the signal, but might still support reliable analysis (Heath *et al.*, 2021), depending on the
10
11 719 research question. An alternative strategy very relevant in automatic detection is to keep only
12
13 720 the audio corresponding to the positive detections: for rarely occurring sounds this will
14
15 721 greatly reduce the storage requirements, while keeping the detected audio clips available for
16
17 722 inspection or re-analysis. However, any missed (false-negative) sound events will be
18
19 723 irretrievably lost. This would prohibit future interrogation of the raw data for other potential
20
21 724 uses.

22
23
24 725 Good-quality metadata including time, date, location and more, is crucial for the success and
25
26 726 reproducibility of any project. This can be stored in the audio files (as “RIFF tags”) or
27
28 727 separately (Metcalf *et al.*, 2023a). Research and other publicly-shared data should be
29
30 728 “FAIR”- findable, accessible, interpretable, reusable (Wilkinson *et al.*, 2016) – and
31
32 729 publishing metadata in standardised formats is key to this. The Biodiversity Information
33
34 730 Standards (TDWG) group maintains the metadata standards Darwin Core (Darwin Core Task
35
36 731 Group, 2009) and Audiovisual Core (GBIF/TDWG Multimedia Resources Task Group,
37
38 732 2013) which help with this through a lightweight approach of specifying common field
39
40 733 names and their definitions (such as “Capture Device”, “Taxon Coverage”, “Locality”, “Start
41
42 734 Timestamp”). By using such standards, researchers can ensure that their metadata will be
43
44 735 understood by others and be findable. It also enables a next generation of methods that could
45
46 736 automatically generalise across multiple available datasets, since the metadata are
47
48 737 compatible.

738 5.5. Data annotation

739 A well-annotated dataset is critical to the performance of a ML-based automated detector.

740 When creating annotations, many decisions must be made, including which program will be

741 used, the specific approach, as well as (often subjective) decisions regarding specifics about

742 the granularity, or what ‘counts’ as an annotation, for example individual vocalisation bouts

743 or whole sequences. There have been calls to standardise annotation approaches in

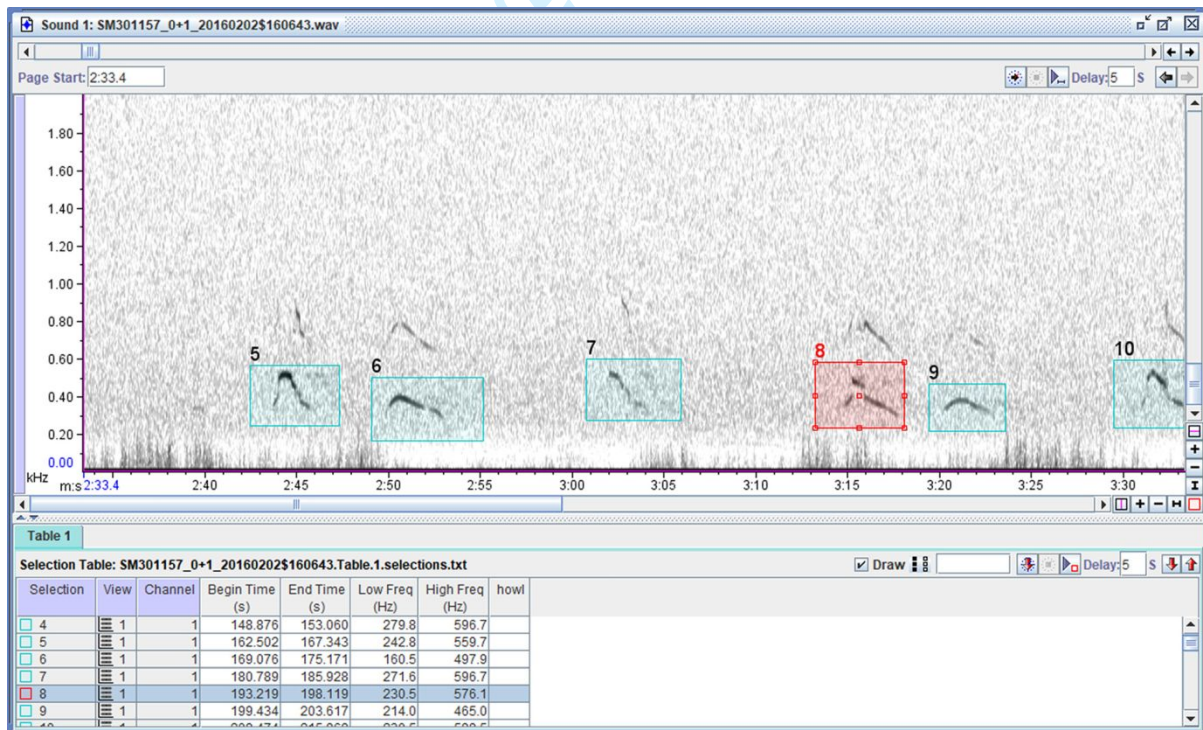
744 bioacoustics (Nicholson, 2023), similar to what has been done for human speech (Gibbon,

745 Moore & Winski, 1998) and music (Humphrey *et al.*, 2014). However, to our knowledge a

746 standardised protocol does not yet exist, perhaps due to the diversity of signal types and

747 research questions across bioacoustics and/or a lack of communication among fields. Here,

748 we aim to provide some guidance for annotating a dataset for automatic detection (Figure 4).



749 Figure 4. Example annotation of acoustic signals, in this case, wolf howls. Taken from

750 (Kershenbaum *et al.*, 2019), showing a spectrogram generated using Raven Pro.

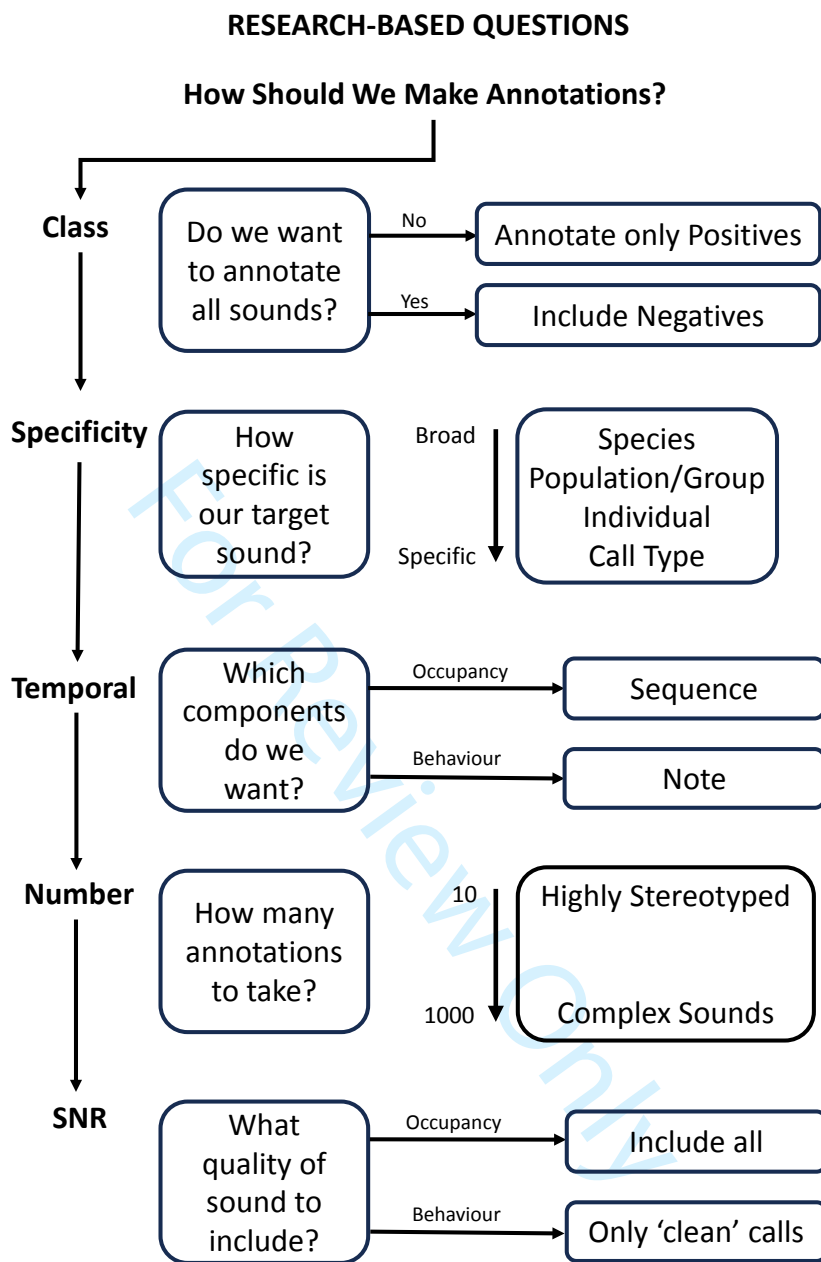
751

752

1
2
3 753 Due to the relatively large amount of human investment required to get high-quality
4
5 754 annotations, researchers often ask themselves how many annotations are needed. This
6
7
8 755 generally depends on the research question, and it is often recommended to annotate as many
9
10 756 signals as possible, however there are more specific questions that can help guide these
11
12 757 decisions. The first concerns the classes or discrete types of signals in your dataset. For
13
14 758 example, will you annotate every bird species in a long-term recording? Will you annotate a
15
16 759 single call type from a single species? Or will you annotate all the notes or elements in a
17
18 760 sequence from a single individual? In addition, one must decide whether to annotate the
19
20 761 “negative class” (oftentimes the “noise” or “absence” category). If doing exhaustive
21
22 762 annotation where all the signals of interest are annotated, then it can be assumed that anything
23
24 763 that is not annotated is the “negative class”. However, strategically annotating other
25
26 764 “distractor/noise” sound events may improve detector performance, especially sounds
27
28 765 occurring within the target frequency range which are loud or easily confused with the target
29
30 766 signal. These “noise” labels can help with error analysis and with the training of an
31
32 767 algorithm.
33
34
35
36
37 768 Decisions about the temporal scale of the annotations must also be made. A common
38
39 769 approach is to annotate the smallest acoustic unit, e.g., note or syllable (Kershenbaum *et al.*,
40
41 770 2016a); however this method can be very time-consuming for large datasets. For vocal
42
43 771 sequences that are comprised of multiple acoustic units (e.g., gibbon vocalisations) another
44
45 772 approach is to annotate particular call types or phrases within the longer sequence, e.g.,
46
47 773 annotate only the female gibbon contribution to the duet (Clink *et al.*, 2023).
48
49
50
51 774 The number of annotations needed will be influenced by the research question and the choice
52
53 775 of the automatic detection approach (see below) but may also be limited by external factors
54
55 776 such as funding support for analysts. It is important to consider the diversity of signal types
56
57 777 as well as background noise, and to work to include a distribution of annotations or samples
58
59
60

1
2
3 778 across sites, times of day, groups, individuals, etc. A higher number of annotations (and
4
5 779 therefore more available samples for training data) will likely improve detector performance
6
7 780 and may be necessary in cases where the signals of interest are highly variable. In some
8
9
10 781 cases, such as the use of transfer learning, a smaller number of training samples (~ 25) may
11
12 782 be sufficient (Dufourq *et al.*, 2022a), but even in these cases, only a test set in the order of
13
14 783 one hundred examples would enable a reliable evaluation of the model. Researchers also need
15
16 784 to make decisions about which target signals to include in their annotations, such as whether
17
18 785 to include low SNR acoustic signals, signals that substantially overlap with non-target
19
20 786 signals, or signals that are abnormal in structure.
21
22
23 787 A common way to do annotations is by visualising spectrograms in a graphical user interface
24
25 788 (GUI) such as Raven Pro (K. Lisa Yang Center for Conservation Bioacoustics, 2014),
26
27 789 SonicVisualizer (Cannam, Landone & Sandler, 2010) or Praat (Boersma & Weenink, 2007)
28
29 790 and creating bounding boxes around the signal(s) of interest. Other possibilities include the
30
31 791 use of an energy or coherence (Wijers *et al.*, 2021) detector to identify all signals above a
32
33 792 certain threshold in a given frequency range and then labelling these detections, applying an
34
35 793 unsupervised clustering algorithm and labelling the batches of samples that have been
36
37 794 grouped together, or the use of DL approaches to identify the start and stop times of signals
38
39 795 of interest automatically, e.g., TweetyNet (Cohen *et al.*, 2022). However, one must be
40
41 796 cautious about mass semi-automated annotations, since these may introduce non-obvious bias
42
43 797 that can affect the conclusions of the study. We recommend including random sampled
44
45 798 manual inspection steps in the procedure. It is important to document your annotation
46
47 799 protocol, including the decisions you made and why you made them, in a way that can be
48
49 800 reproduced by others. We suggest including these protocols as online supporting material in
50
51 801 publications. In addition, it is crucial to check both intra- and inter-observer reliability for
52
53 802 creating annotations (Nguyen Hong Duc *et al.*, 2021).
54
55
56
57
58
59
60

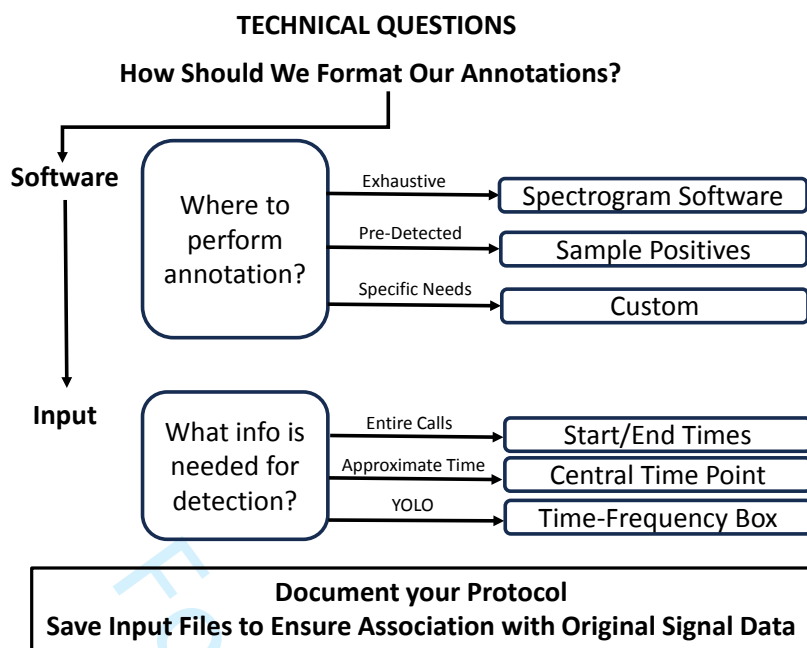
803



804

805 Figure 5a. A flowchart for designing research questions in relation to automatic detection.

806

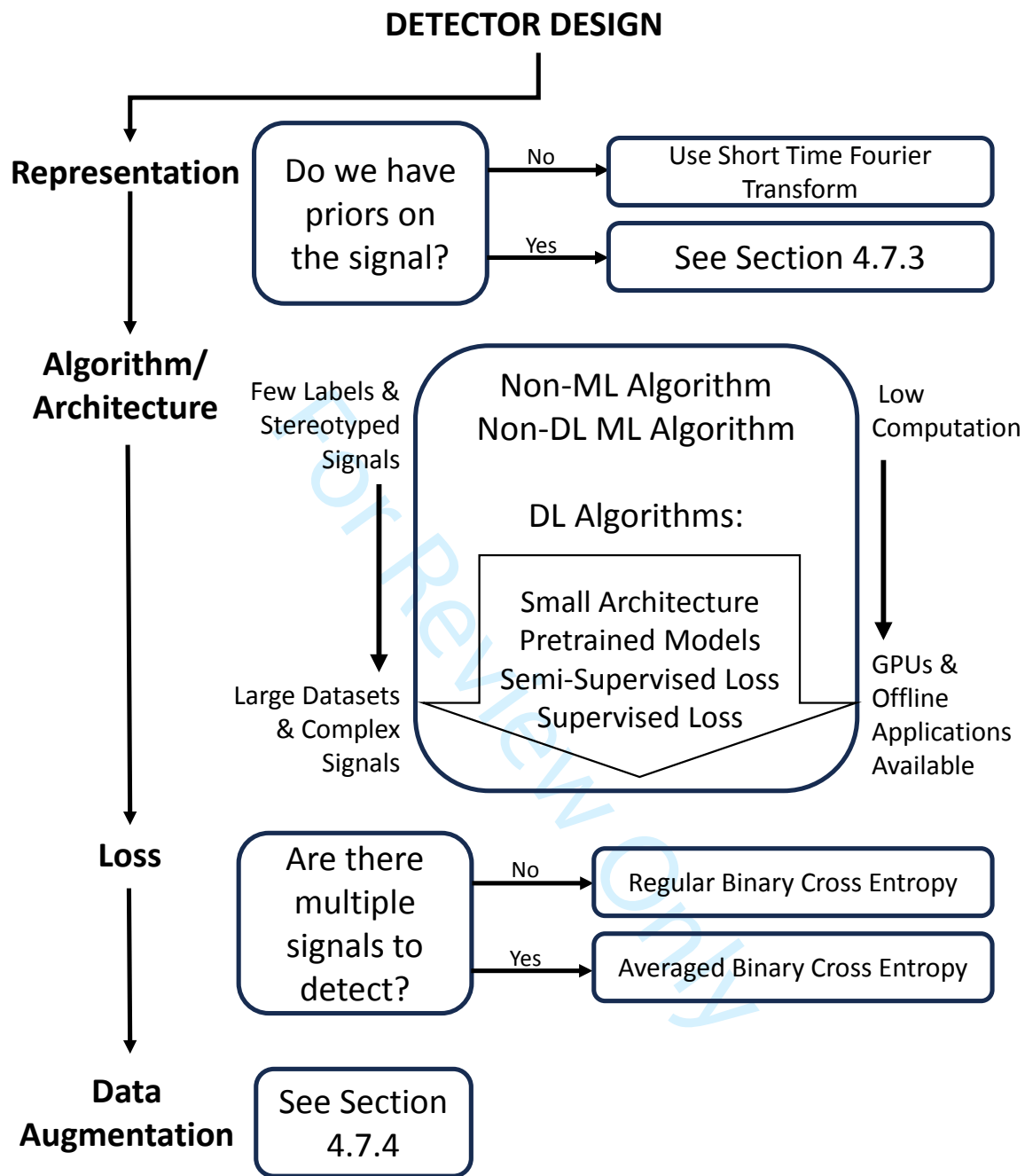


807

808 Figure 5b. A flowchart showing annotation decisions for automatic detection.

809

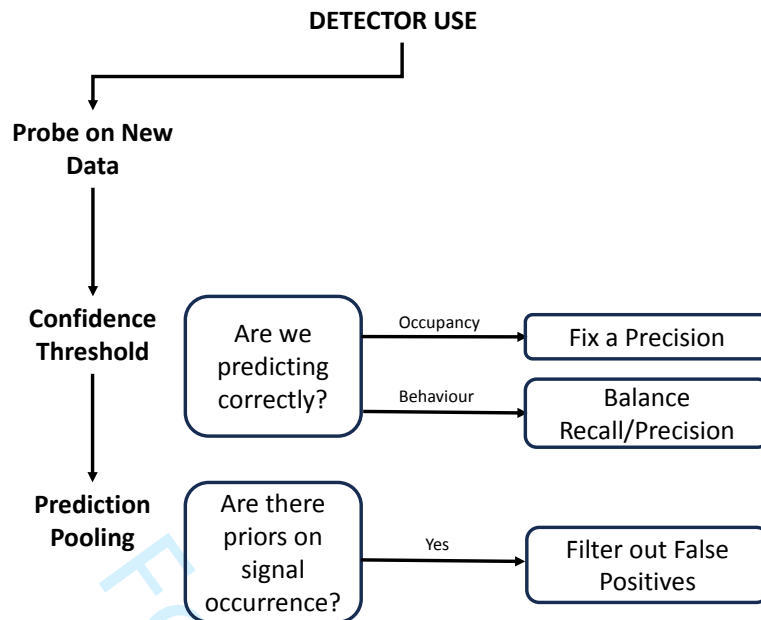
810 5.6. Choose your Detection Pipeline



811

812 Figure 5c. A flowchart showing the decisions necessary in automated detector design.

813



814

815 Figure 5d. A flowchart showing the constraints of end use on automated detector design.

816

817 5.6.1. Interfacing with your pipeline

818 Selecting an automatic detection approach depends on factors such as technical familiarity,
 819 desired granularity, and budgetary constraints. Products such as Kaleidoscope
 820 (<https://www.wildlifeacoustics.com/>), PAMGuard (<https://www.pamguard.org>), and Arbimon
 821 (<https://rfcx.org/ecoacoustics>) provide easy to use interfaces for systems that can perform
 822 automatic detection on audio samples originating from a wide variety of environmental
 823 samples. These tools come equipped with traditional approaches rooted in standard signal
 824 processing techniques but are limited in their ability to utilise modern advances in machine
 825 learning. Conversely, modern DL frameworks, such as TensorFlow, PyTorch, and Keras
 826 (Stowell, 2022b), as well as the models built with them, rarely come with an easy to use
 827 interface which makes them less accessible. Commercial approaches offering cloud-based
 828 machine learning as a service (MLaaS) solutions, such as those from Amazon, IBM, or
 829 Microsoft, allow easier access to these advanced methods, but can be prohibitively expensive.

1
2
3 830 Practitioners must decide whether easy-to-use tools are sufficient for the problem at hand, or
4
5 831 whether it would be advantageous to exploit the often superior performance of DL methods,
6
7 832 which require more investment of time, money or both. The complexity of the research
8
9 833 question has a significant influence on the selection but may be outweighed by the need to
10
11 834 invest further in expertise or funding.

12
13
14 835 In the case of any automatic detection approach, the pipeline must be evaluated in the context
15
16 836 of the research questions which necessitates dividing the data to properly evaluate
17
18 837 performance and generalisability, the choice of an appropriate detection mechanism, and the
19
20 838 selection of relevant, comparable, and appropriate metrics.

21 22 23 24 25 839 5.6.2. Split your data

26
27
28 840 As for most machine learning tasks, datasets should be split into train, test, and validation
29
30 841 subsets to ensure the true generalisability and comparability of a model's performance. This
31
32 842 means that an amount of data (usually around 10 to 20% of the total dataset) need to be kept
33
34 843 unseen during training and validation of the model, for which the remaining 80-90% of the
35
36 844 data are used. This helps to avoid model overfitting, which would cause the model to learn
37
38 845 only the characteristics of the training data, without the ability to generalise to new data, and
39
40 846 would bias performance scores (Gareth James *et al.*, 2013, p176).

41
42
43
44 847 The validation (or development) set is used for hyperparameter tuning. This is especially
45
46 848 useful in the case of DL models which involve empirical testing of optimal values and setups
47
48 849 for elements such as optimisers, learning-rates, or early stopping. The best performing model,
49
50 850 as determined using the validation set is then applied to the test set. Finally, the best
51
52 851 performing model on the validation set is applied to the test set. The test set should not be
53
54 852 used to fit the values of such hyperparameters or to compare model architectures since it
55
56 853 would no longer serve for generalisation assessment; it is kept for final performance
57
58 854 evaluation. Creating an effective test dataset may include the selection of a separate
59
60

1
2
3 855 microphone entry, specific time frames, separate recording locations, or subsets of
4
5 856 vocalisations from an individual which were not included in the training set, amongst others.
6
7 857 The general idea here is to separate the prediction capabilities of the computer model from
8
9 858 recording specificities and data related biases. We always want to ensure that an automatic
10
11 859 detection model is generalisable rather than specifically trained for a single recording setup,
12
13 860 location, or individual.
14

15 861 To provide an example, in the case of creating a presence/absence detection model, one
16
17 862 should not use annotations from the same file for training and testing. But instead, certain
18
19 863 audio files should be used to create the train dataset, and independent files should be used to
20
21 864 test the detection model. Furthermore, the model should be applied to entire testing audio
22
23 865 files and not only to parts of the test file that have been annotated, as this might result in an
24
25 866 overly optimistic evaluation of the model and potential false positives would be missed.
26
27
28
29
30

31 32 867 5.6.3. Pick your feature representation 33

34 868 Depending on the automatic detection approach, acoustic data may be transformed through
35
36 869 feature extraction to ease the automatic detection process. In the computational bioacoustics
37
38 870 literature, an array of such feature extraction methods can be found, each presenting their
39
40 871 own advantages and limitations.
41

42 872 In bioacoustics, the dominant approach is undoubtedly spectral representations such as
43
44 873 spectrograms or mel-spectrograms. This type of representation usually allows for
45
46 874 interpretable visualisation of acoustic data and provides an easy route to use popular vision-
47
48 875 based models such as CNNs for object detection and image classification. Despite this, some
49
50 876 information from the raw waveform may get lost when computing these representations. This
51
52 877 is especially the case for transient signals such as odontocetes' clicks which are poorly
53
54 878 represented by Fourier transforms (Jiang *et al.*, 2018). CNNs developed for spectrograms
55
56 879 cannot be used directly for waveforms, because the data is of different dimensionality;
57
58
59
60

1
2
3 880 however there have been a lot of recent developments in DL methods applied directly to
4
5 881 waveforms and so this is increasingly feasible (Baevski *et al.*, 2020).
6
7
8 882 DL methods now allow for high-dimensional inputs such as whole spectrograms, with the
9
10 883 succession of layers extracting higher level features and information. However, historically
11
12 884 users were the ones responsible for selecting relevant features to represent signals. In this
13
14 885 context, MFCCs were often used, and given to a classification algorithm such as a support
15
16 886 vector machine (Mitrovic, Zeppelzauer & Breiteneder, 2006). For relatively simple use cases
17
18 887 e.g., stereotyped signals and low background noise, this approach might suffice in bringing
19
20 888 satisfactory performances.
21
22
23 889 Recently, as stated in Section 0, extracting pretrained latent representations as features is also
24
25 890 being adopted as a promising solution. This approach may imply additional effort on the part
26
27 891 of the user and raises an array of questions on pretraining datasets, selected model
28
29 892 architectures or the need for higher computational power. It can also prove successful in
30
31 893 easing the downstream learning process or allowing for smaller annotated datasets in few-
32
33 894 shot learning perspectives.
34
35
36 895 Despite the advantage of using such abstract representations, using traditional engineered
37
38 896 features such as fundamental frequency, call duration or number of notes may still prove to
39
40 897 be effective depending on the task at hand. These can also be combined with features
41
42 898 extracted from the time domain such as energy and zero-crossing rates. These can then allow
43
44 899 for the use of simpler algorithms which may be easier to implement and require little
45
46 900 computational power and training time.
47
48
49 901 Overall, there is no such thing as the perfect feature extraction method for bioacoustics.
50
51
52 902 Comparing different feature representations should always be the preferred approach and can
53
54 903 be carried out on the previously mentioned validation set, ideally in a pilot study.
55
56
57
58
59
60

1
2
3 904 5.6.4. Decide on feature transformation
4

5
6 905 Prior to feature extraction, specifically in the case of noisy recordings characterised by low
7
8 906 SNR, some detectors may benefit from denoising, i.e., the automatic removal of background
9
10 907 noise from the acoustic signal of interest. An extensive overview of recent approaches can be
11
12 908 found in (Xie, Colonna & Zhang, 2021) with accessible open-source solutions. Some of these
13
14 909 methods are built on light-weight algorithms such as spectral-gating (Sainburg, 2019), others
15
16 910 involve the use of DL with CNNs, Noise-2-Noise-based approaches (Bergler *et al.*, 2020), or
17
18 911 denoising-autoencoder models (Vickers *et al.*, 2021; Yang *et al.*, 2021).
19

20 912 Although it is useful in some applications, this pre-processing step is not always
21
22 913 recommended and must be used with caution as it may result in a loss of information. In
23
24 914 some cases, noise can also be directly handled by the detector itself, especially when using
25
26 915 noise-resilient DL architectures or when stationary noise is not overlapping the target signals.
27
28 916 In cases where noise reduction is applied prior to training, the evaluation and test datasets
29
30 917 will need to be put through the same process, to ensure that training and testing data have
31
32 918 comparable characteristics and contain similar acoustic information. When building a noise
33
34 919 resilient model, one may also resort to multi-condition training approaches. This method can
35
36 920 imply adding noisy corrupted versions of the data to the training set or including both the
37
38 921 original and the denoised versions of the data during training to help with model robustness
39
40 922 to noisy acoustic contexts. This approach is fairly common in speech processing (Yin *et al.*,
41
42 923 2015) but needs further exploration in bioacoustics.
43
44

45 924 Depending on the amount of training data available, data augmentation techniques may be
46
47 925 used to artificially increase the variability of the data on which models are optimised. The
48
49 926 choice of which augmentation technique to use depends on the application. One should aim
50
51 927 to apply transformations that cover the range of variations found in real signals. However,
52
53 928 care must be taken to avoid transformations that could invalidate the annotations. For
54
55
56
57
58
59
60

1
2
3 929 instance, in a bird call detector reversing sounds could be a tempting simple transformation,
4
5 930 yet this could result in artificially making a bird call more similar to that of another species.
6
7 931 Simple transformations may also create artefacts that can complicate the modelling, for
8
9 932 example a pitch shift of a howl may also unrealistically shift the background noise.
10
11 933 Commonly used techniques include stretching or compressing the duration of acoustic
12
13 934 signals, shifting their pitch, making small volume modifications, or adding a variety of noise
14
15 935 or mixing with other audio events via some linear or non-linear combination (e.g., taking one
16
17 936 presence event and mixing it with one or more absence events). These transformations may
18
19 937 also be combined to produce more variation.
20
21 938 Recently generative deep-learning methods, such as Generative Adversarial Networks
22
23 939 (GANs) have been proposed in order to generate synthetic examples (Wang, She & Ward,
24
25 940 2022; Bergler *et al.*, 2022a).

31 32 941 5.6.5. Decide on a method

33 34 35 942 5.6.5.1. Deep learning or not

36
37 943 As mentioned above, the choice of a detection mechanism is dependent at least partially on
38
39 944 the complexity of the problem. If the signals are well defined, have high SNR, are highly
40
41 945 stereotyped, and the research question involves simple segmentation and can be done offline,
42
43 946 a package such as Kaleidoscope or Arbimon may be more than adequate.
44
45 947 Using machine (deep) learning may be advantageous in situations requiring a more complex
46
47 948 analysis, such as call type classification, or where robustness to environmental noise is
48
49 949 necessary (Aodha *et al.*, 2018; Stowell, 2022a). However, in situations where access to either
50
51 950 a large amount of computing resources or the training / expertise to use them effectively is
52
53 951 limited, the use of DL may not be possible. Additionally, it must be considered where the
54
55 952 detection mechanism will be deployed. If access to a large computing cluster is readily
56
57
58
59
60

1
2
3 953 available but the end result must function on a small device for field deployment, then a large
4
5 954 and complex model may not work. Conversely, if the final model will only be used offline
6
7
8 955 using minimal computing resources (budget GPU), then the model choice becomes somewhat
9
10 956 more flexible. Different machine learning approaches are given in Table 1, together with their
11
12 957 requirements and example studies.
13

14 958
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

For Review Only

959 Table 1. Different types of machine learning techniques

| Learning Type | Labelled Data Requirements | Metrics | Visualisations | Examples |
|---|---|---|--|---|
| Supervised (Segmentation, Classification) | Large amount of labelled data | Accuracy, Precision, Recall, F-Score, AUROC, mAP, UAR | Confusion Matrix ROC-Curve, PR-Curve | (Bergler <i>et al.</i> , 2022b) |
| Unsupervised or self-supervised (Clustering) | Labelled Data Not Necessary | Reconstruction Loss (MAE, MSE), homogeneity, completeness | Reconstructions, Dim-Reduction (t-SNE, UMAP) | (Cuevas <i>et al.</i> , 2017) |
| Semi-Supervised Learning | Some Labelled Data - Large amount of Unlabelled data (optional) | Both Supervised & Unsupervised | Both Supervised & Unsupervised | (Bermant <i>et al.</i> , 2019; Saeed, Grangier & Zeghidour, 2021; Leroux <i>et al.</i> , 2021; Hagiwara <i>et al.</i> , 2022) |

1
2
3
4 961 5.6.5.2. Choose your evaluation metrics
5

6
7 962 The evaluation of the automatic detection mechanism depends primarily on the type of task to
8
9 963 be performed. A fully supervised detection / classification task is typically evaluated using
10
11 964 metrics such as accuracy, precision, recall, F-score, or area under the receiver operating
12
13 965 characteristic curve (AUROC) (Lever, Krzywinski & Altman, 2016). These all provide
14
15 966 different insights and can help evaluate how the model is performing. For example, precision
16
17 967 indicates the fraction of relevant results (true positives) that are found among all detected
18
19 968 events, whereas recall indicates the fraction of signals in the dataset which were effectively
20
21 969 found. Typically, a balance must be decided which metrics are most important for a particular
22
23 970 task. For example, recall may be an important score to consider when detecting rare
24
25 971 phenomena where missing a single detection of an underrepresented class may prove costly.
26
27 972 Wrong choice of metrics may bias the results, for example, in the case of highly unbalanced
28
29 973 datasets, i.e., when the acoustic object to be detected is rather underrepresented in the dataset
30
31 974 compared to negative labels, accuracy may be very high despite low performances on the
32
33 975 small number of positive test samples.
34
35 976 Visualising results from supervised training methods can involve a confusion matrix, which
36
37 977 is a table that shows the ground truth values on one axis and predicted values on the other,
38
39 978 allowing visual analysis of model performance which is easy to digest. Another option is the
40
41 979 receiver operating characteristic curve (ROC curve), which plots the trade-off between true
42
43 980 positive rate (TPR) and false positive rate (FPR) at all confidence thresholds, enabling the
44
45 981 analyst to more easily choose a prediction threshold which suits their needs. The area under
46
47 982 the ROC curve (AUROC) gives a summary of the model's performance across threshold and
48
49 983 is agnostic of threshold choice.
50
51 984 A similar visualisation to the ROC curve is the Precision-Recall (PR) curve, which also
52
53 985 highlights the balance between missing out events (false negative) and making false alarms
54
55
56
57
58
59
60

1
2
3 986 (false positive). The area under the PR curve is commonly referred to as the mean average
4
5 987 precision (mAP). The important difference between PR and ROC curves is that the precision
6
7 988 gives the proportion of correct detection among all detections and the FPR indicates the
8
9 989 proportion of wrong detections among all negative examples. In the case of highly
10
11 990 unbalanced datasets (e.g., 1% of positive examples), the FPR can be rather optimistic as
12
13 991 compared to the precision, and thus the mAP might come out to be significantly lower than
14
15 992 the AUROC. Detailed discussions on possible performance metrics can be found in (Davis &
16
17 993 Goadrich, 2006; Hildebrand *et al.*, 2022).

18
19 994 Useful metrics for unsupervised learning are harder to identify, as it depends on the research
20
21 995 question. If labelled data are available, they can be used to assess the quality of a clustering
22
23 996 attempt by measuring completeness (across how many clusters are samples with the same
24
25 997 label) or homogeneity (the proportion of samples in a cluster with the same label).

26
27 998 Visualisation for unsupervised clustering results are often done by reducing the
28
29 999 dimensionality of the reductions to either two or three dimensions using t-Stochastic
30
31 1000 Neighbour Embedding (t-SNE) (Maaten & Hinton, 2008), Uniform Manifold Approximation
32
33 1001 and Projection (UMAP) (McInnes, Healy & Melville, 2020), or a similar method.

34 35 1002 5.7. Verifications - check your results

36
37 1003 The verification of model performance on the test data should involve quantitative and
38
39 1004 qualitative evaluations. Quantitative metrics give the performance in terms of comparable
40
41 1005 values like the F1-score, accuracy, precision, or recall. Whilst the qualitative metrics would
42
43 1006 help to understand the practical implications of the model. Qualitative analysis involves
44
45 1007 manually checking or visualising the predictions. This may involve plotting automatic
46
47 1008 segmentation results on spectrograms to visually account for the precision of detected time
48
49 1009 frames. It may also be carried out through a simple manual inspection of a subset of results.
50
51
52
53
54
55
56
57
58
59
60

1
2
3 1010 Careful manual analysis of the signals with missed detections or false alarms could help to
4
5 1011 identify the characteristics that trigger the models and help to improve the models further by
6
7
8 1012 adding the specific variations needed in the training data or clean up train data (especially
9
10 1013 wrong annotations or mislabelled data).
11
12

14 1014 5.7.1. When is a model good enough? Performance thresholds

16 1015 Understanding the performance thresholds and being realistic about the task is a pragmatic
17
18 1016 way of approaching the problem. It is important to understand that machine learning models
19
20
21 1017 are statistical in nature and may never provide 100% performance even with perfect data or
22
23 1018 models. Understanding the limitations of the model and the desirable performance in the real
24
25 1019 world scenario can help set the thresholds for performance, for example, trade-off between
26
27 1020 false positives and missed detections (Karnan, Akila & Krishnaraj, 2011). In some scenarios
28
29 1021 it may not be even practically feasible to achieve a desirable performance due to factors like
30
31 1022 overlapping sounds, environment noise or very low SNR. But understanding and defining the
32
33 1023 problem based on a trade-off between what is feasible with the acoustic data and what is
34
35 1024 desirable could help define performance thresholds and build practical models. For example,
36
37 1025 defining the range of distance within which the target species needs to be detected.
38
39
40
41
42

43 1026 5.7.2. How harmful are mistakes (false positives vs false negatives)?

45 1027 The use case for automatic detection will influence how much (section 4.8.1) and what kind
46
47 1028 of errors are acceptable. For instance, if doing an analysis on vocal behaviour, missing a call
48
49 1029 in a sequence might strongly distort results. Conversely, if occupancy trends are aimed for,
50
51 1030 missing one call in a sequence is insignificant, and imperfect detection can be incorporated
52
53 1031 into occupancy models (Bailey, MacKenzie & Nichols, 2014). Recall is thus more or less
54
55 1032 important depending on the type of study being conducted.
56
57
58
59
60

1
2
3 1033 In general, false positives are undesirable, but a certain amount might be acceptable (Shiu *et*
4
5 1034 *al.*, 2020). In any case, converting the precision into the number of false positives per hour
6
7 1035 allows an unambiguous interpretation by the user and the planning of how to deal with false
8
9 1036 alarms.

10
11
12 1037 Additionally, prior knowledge on vocal behaviour such as the sequence regularities might
13
14 1038 allow filtering out of false positives. Such priors can be used to reduce confidence thresholds
15
16 1039 and increase the recall, but with the risk of imposing too strong priors and missing out on
17
18 1040 uncommon sequences.

21 22 23 1041 5.7.3. Reproducibility and accessibility

24
25 1042 We also expect automated vocalisation detection systems to be made available to other users,
26
27 1043 thus broadening the contribution to the field of bioacoustics (especially to users without a
28
29 1044 strong computer science background). For this purpose, code for detection systems should be
30
31 1045 shared in comprehensive and accessible ways, such as version control repositories, and
32
33 1046 should be well documented with detailed user manuals (Braga *et al.*, 2023). An easy way to
34
35 1047 make a detection model available to the community is also to follow common APIs that will
36
37 1048 allow their integration into pre-existing interfaces, such as ARISE (Hogeweg & Stowell,
38
39 1049 2023) or Raven Pro (K. Lisa Yang Center for Conservation Bioacoustics, 2014).

40
41
42 1050 Besides publishing code for experiments to be reproducible, datasets used for training and
43
44 1051 testing should be made available to the community for building new systems and comparing
45
46 1052 them using standard annotation protocols (see Section 5.5). Indeed, public benchmarking
47
48 1053 datasets exist (Joly *et al.*, 2015; Politis *et al.*, 2020) but cover only a relatively small set of
49
50 1054 species targeted by bioacoustic studies.

1055 5.7.4. Access to raw recordings

1056 Ideally, additional to labelled training dataset, raw recordings (as opposed to cut-out
1057 snapshots) are of potential value to the machine learning community (to train self-supervised
1058 models for instance) and to the research community in general to reuse the data for other
1059 tasks or to create new annotated datasets from previously recorded data. But it might not be
1060 always feasible to make this readily accessible in public repositories due to storage and other
1061 constraints. We encourage researchers to store the raw recordings locally and share them on
1062 demand with the community or with interested parties.

1063 6. WAYS FORWARD

1064 We now consider some important ways forward for automatic detection for bioacoustics,
1065 including best practices which should be implemented now, the challenges still to be
1066 overcome, and the future direction of the field.

1067 6.1. Challenges

1068 6.1.1. Bioacoustic challenges

1069 Although automatic detection has already brought large improvements to the field of
1070 bioacoustics, challenges remain which are closely related to the nature of animal sound
1071 and/or the desired uses of such data. For instance, population density estimates rely on
1072 detections being reliable, without double-counting individuals' vocalisations when they are
1073 picked up by multiple devices (Kimura *et al.*, 2010; Marin-Cudraz *et al.*, 2019), and are
1074 further improved if calls can be localised and attributed to an identified individual (Nijman,
1075 2007; Knight & Bayne, 2019; Hedley *et al.*, 2021; Law *et al.*, 2021). Moreover, in most cases
1076 population density cannot be estimated without knowing the detection range of the system
1077 (Metcalf *et al.*, 2023a). The detection range of the acoustic signal will depend on multiple

1
2
3 1078 factors including source level and frequency range of the signal, characteristics of the habitat
4
5 1079 including ambient noise levels, vegetation and topography, along with specifications of the
6
7
8 1080 ARU (Hauptert, Sèbe & Sueur, 2022). However, detection range is often difficult to estimate,
9
10 1081 especially in forest environments or areas with extreme topography, and in many cases is
11
12 1082 ignored or assumed to be consistent across studies when this may not be the case. When
13
14 1083 species of interest are near the limit of the detection range of the device, recordings of vocal
15
16 1084 signals may become attenuated or missed. This might cause problems in some tasks which try
17
18
19 1085 to capture specific aspects of the vocalisation, for example to infer behaviour, caller identity
20
21 1086 or communication patterns, rather than generic tasks which look at occupancy (Spillmann *et*
22
23 1087 *al.*, 2017).
24
25
26 1088 Even when accurately focusing on our target species' vocal signals, animals might engage in
27
28 1089 simultaneous vocalisations or choruses (Torti *et al.*, 2018), which makes a simple
29
30 1090 timestamped detection system insufficient for acoustic behaviour analysis. Also, it can be
31
32 1091 difficult to distinguish vocalisations of similar species if they share characteristics, e.g., dog
33
34 1092 barks and coyote barks share a number of similarities which make it difficult to determine
35
36 1093 which species produced the rapid-fire sequence of noisy barks, though there are some
37
38 1094 quantitative differences (Feddersen-Petersen, 2000).
39
40
41
42

43 1095 6.1.2. Computational challenges

44
45
46 1096 Computational challenges in this field include questions of algorithms, datasets,
47
48 1097 computational efficiency, computing platforms and more.

49
50
51 1098 One overarching challenge within machine learning in the broad, and with particular
52
53 1099 relevance to automatic detection, is the ability to generalise. For example, a model well-
54
55 1100 trained for a particular species can perform poorly with even slight variations in recording
56
57 1101 devices, ambient noise, or operating environments. This could lead to low accuracy without
58
59 1102 further testing and adjustment. Creating scalable models that have the flexibility to add new
60

1
2
3 1103 species to the training dataset, to increase the number of vocally active species which can be
4
5 1104 detected, is still a challenging task. Transferring knowledge from models built with data from
6
7 1105 one species to a new species without further training data is even more desirable. We also
8
9 1106 note that many models are highly task specific - the data specification, annotations, model
10
11 1107 architectures, and systems are highly optimised for best performance. For example, a system
12
13 1108 used to determine the occupancy of a species may not be suitable for individual
14
15 1109 identification, understanding communication, or behaviour patterns which superficially
16
17 1110 appear to be related but are subtly different tasks. It is not immediately clear to a user how far
18
19 1111 to trust in the generalisation of a detector.
20
21
22

23
24 1112 Acquiring generic datasets that can address multiple tasks, such as population density
25
26 1113 estimation and behavioural characteristics, poses a significant challenge due to the limitations
27
28 1114 in data collection strategies. Typically, data collection is initially planned to address specific
29
30 1115 tasks, which makes it difficult to acquire datasets that can be scaled to any given task. This is
31
32 1116 a challenge as it is essential to streamline and optimise the recordings to collect only data of
33
34 1117 interest to a particular task to increase storage and computational efficiency. But, at the same
35
36 1118 time, the data collected might not include the context or information that was needed to use it
37
38 1119 for a new task. A lack of generic, benchmark datasets has significant implications for the
39
40 1120 standardisation of methods in the field and the appropriate evaluation of research.
41
42
43

44 1121 In bioacoustics as in other fields, DL comes with very limited interpretability, an issue known
45
46 1122 as the 'black box problem'. This amplifies the problem that conclusions drawn about DL
47
48 1123 models will be specific to the dataset they were tested on, which significantly hinders the
49
50 1124 process of finding a consensus for the best architecture or training procedure to be used. In
51
52 1125 certain cases, it is also unclear as to how different research studies split their datasets and
53
54 1126 conduct model evaluation. As it stands, little to no standards on the best approaches exist and
55
56 1127 without these best practices put in place, authors will implement their own approaches within
57
58
59
60

1
2
3 1128 their research. The best opportunity to overcome issues such as these is firstly to encourage
4
5 1129 further development of public access or benchmark datasets, and secondly to probe models
6
7
8 1130 on their detailed behaviour regarding these datasets (Alain & Bengio, 2018). Within the
9
10 1131 current literature, the approach that authors have taken to implement their machine learning
11
12 1132 testing methodologies and model evaluation differs drastically. In most cases, comparisons
13
14 1133 are not made to existing results on datasets that are publicly available, instead, most studies
15
16 1134 present their findings related to their proposed method on the dataset that was collected for
17
18
19 1135 the study. These observations are quite different to what has been observed within the
20
21 1136 computer vision and natural language processing literature whereby most studies will
22
23
24 1137 compare their proposed method to various baselines and existing state-of-the-art methods on
25
26 1138 the same datasets. Consequently, a comparison between research studies within bioacoustics
27
28 1139 is not feasible and determining the state-of-the-art is non-trivial. Various initiatives exist that
29
30 1140 provide bioacoustic benchmark datasets and standardised public evaluations, including
31
32
33 1141 automatic detection in particular, though these are neither as large nor as widely-used as in
34
35 1142 mainstream ML application domains (Stowell *et al.*, 2019; Ferrari *et al.*, 2020; Hagiwara *et*
36
37 1143 *al.*, 2022).

38
39
40 1144 Training machine learning models, particularly deep neural networks, is computationally
41
42 1145 intensive. Specifically, computers, workstations, or servers with a large amount of CPU
43
44 1146 (central processing unit) and GPU may be needed, to speed up the training or just to make it
45
46
47 1147 achievable in reasonable time. Furthermore, certain deep neural networks require a large
48
49 1148 amount of GPU RAM to load the model into memory given the large number of trainable
50
51 1149 neural network parameters that need optimisation. The issue of access to computational
52
53 1150 power can exacerbate inequalities between people, institutions, and countries. However, the
54
55
56 1151 good news is that the widespread use of pretrained models can massively decrease the
57
58 1152 amount of computation needed: most researchers should not need to train a model from
59
60

1
2
3 1153 scratch. This helps to reduce inequalities as well as the carbon footprint incurred through a
4
5 1154 move to ML methods.
6
7 1155 In conjunction with computation, data storage requirements have skyrocketed with the
8
9 1156 amount of data being collected from PAM and necessities to store, share and create backups
10
11 1157 of these very large datasets. In certain cases, practitioners have had to ship hard drives
12
13 1158 physically across the world to share acoustic datasets, and in other cases practitioners share
14
15 1159 large datasets via cloud-based solutions. It is unlikely that storing all audio for all projects is
16
17 1160 feasible, and yet discarding audio takes away the possibility of reanalysis or new uses.
18
19 1161 Bioacoustics will benefit from the development of mixed schemes with well-designed
20
21 1162 heuristics to store some audio in full resolution (e.g., detected audio clips) and the remainder
22
23 1163 in highly compressed formats which are still reusable (e.g., embeddings or low-bitrate lossy
24
25 1164 compression).
26
27 1165 There are other considerations that arise from the large data volumes that are required both
28
29 1166 for training automatic detection systems, and for investigating biological questions using
30
31 1167 bioacoustics. Logistical challenges in maintaining the data collection devices include
32
33 1168 changing batteries, calibration of microphones, and general wear and tear. Sometimes the
34
35 1169 devices need to be deployed in remote, difficult-to-access, or even dangerous locations,
36
37 1170 which makes the maintenance even more challenging. Therefore, the effort required to gather
38
39 1171 the volume of data needed for training automatic detection models needs to be considered
40
41 1172 carefully. However, artificial intelligence being a rapidly evolving field means that new
42
43 1173 techniques and models may ease (or indeed exacerbate) the problems of providing enough
44
45 1174 data.
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

1
2
3 1175 6.2. Future directions
4
5

6
7 1176 6.2.1. Accessibility
8

9 1177 The extent to which automatic detection for bioacoustics is accessible to a wide range of
10 1178 researchers across different fields and geographical regions is patchy and insufficient. Future
11
12 1179 developments in the field must include increasing the ease with which researchers can
13
14 1180 implement and customise the technology. Usable, stable, and open-source tool kits with an
15
16 1181 associated GUI, and potentially a cloud-based solution, can aid the entry of practitioners from
17
18 1182 a non-machine-learning background and reduce the learning curve. Standards-based
19
20 1183 interoperability and component-based approaches will help ensure that solutions remain well-
21
22 1184 maintained and usable.
23
24
25

26
27 1185 To move to the next generation of automatic detection, we look forward to further work
28
29 1186 developing the scale, reliability, and generality of machine learning methods in bioacoustics.
30
31 1187 But even considering the current state of the art, the barrier to entry for practitioners, students
32
33 1188 and researchers who are new to the field of machine learning is high (Broll & Whitaker,
34
35 1189 2017; Schultze, Gruenefeld & Boll, 2020). This barrier is potentially even higher for
36
37 1190 newcomers in machine learning for bioacoustics than those entering the field of machine
38
39 1191 learning for computer vision or natural language processing. For the latter two, there are large
40
41 1192 quantities of educational material, including blog posts, online tutorials, books, videos, and
42
43 1193 software repositories. The number of research laboratories, and researchers from tertiary
44
45 1194 educational institutions working on automatic detection for PAM or bioacoustics in general is
46
47 1195 not evenly distributed between the Global North and South, and thus, the ability to train
48
49 1196 students may differ between regions. There is a pressing need for more educational material
50
51 1197 to become available so that those entering the field can rapidly learn the necessary skills to
52
53 1198 facilitate progress, and as such, we encourage researchers and practitioners to create and
54
55 1199 share open-access educational material.
56
57
58
59
60

1
2
3 1200 Complementary to educational materials is of course that systems themselves should be more
4
5 1201 accessible and user-friendly. The required use of Python or R (let alone libraries such as
6
7 1202 Tensorflow, and repositories such as Github, etc) acts as a barrier to many potential users,
8
9 1203 and so projects that develop good interfaces are to be celebrated. However, the pace of
10
11 1204 change in ML methods is fast, as well as the diversity of platforms (e.g., mobile devices), so
12
13 1205 it is risky to advocate a single graphical interface. The solution is to rely on component-based
14
15 1206 approaches and well-documented standards; as long as user interfaces can use standards-
16
17 1207 based methods to “talk to” algorithms and datasets, and each of these components can be
18
19 1208 replaced, substituted and improved, we provide a good substrate that makes it easy for
20
21 1209 interface developers to add value to the work (Darwin Core Task Group, 2009; GBIF/TDWG
22
23 1210 Multimedia Resources Task Group, 2013). For all these components, the community needs to
24
25 1211 consider their maintenance models (open source or commercial, free or subscription-based)
26
27 1212 and the ongoing maintenance of core components should not be left to chance.
28
29
30
31
32
33

34 1213 6.2.2. Foundation models

35
36
37 1214 As with the maturation of machine learning in fields such as image or speech recognition, we
38
39 1215 expect animal vocalisation detection models to progressively standardise, not only in terms of
40
41 1216 model architectures but also in data representation. Indeed, pretrained models created from
42
43 1217 large datasets with a variety of species or taxa can yield rather generic embeddings, allowing
44
45 1218 good performances when fine-tuning for a specific task, even when relatively few labels are
46
47 1219 available (see Section 0). Fields such as text processing and image recognition are beginning
48
49 1220 to move to a scale where “foundation models” emerge, meaning DL models which are trained
50
51 1221 across massive and highly varied datasets, whose scales lead to emergent generalisation
52
53 1222 behaviour and which can be reused for a wide range of downstream tasks (Bommasani *et al.*,
54
55 1223 2022). The same could happen for bioacoustics and automatic detection: although the size of
56
57 1224 the benefit is hard to foresee, large-scale highly generalised models could indeed overcome
58
59
60

1
2
3 1225 the significant limitation in bioacoustics that many custom tasks do not come with strong
4
5 1226 training datasets. An alternative approach is few-shot learning, recently explored to
6
7 1227 generalise robustly from as few as five examples (Nolasco *et al.*, 2023). Such methods
8
9 1228 indicate that “one big dataset” is not necessarily the main objective for the field. These trends
10
11 1229 may converge, with the many public bioacoustic datasets forming a richly structured
12
13 1230 pretraining curriculum for systems to generalise well from simple examples.
14
15
16
17

18 1231 6.2.3. Multi-modal detection

19
20
21 1232 Some challenges posed by automatic bioacoustic detection, including difficulties in
22
23 1233 separating individual emitters, precisely assessing population density, double counting, or
24
25 1234 missing detections, could potentially be eased by multi-modal approaches. In fact,
26
27 1235 incorporating additional modalities such as images, video or GPS data, may result in
28
29 1236 complementary information missing from the acoustic data and enhance the detector’s
30
31 1237 performance, which can then enable uses such as abundance estimation (Akamatsu *et al.*,
32
33 1238 2013) and activity tracking (Li *et al.*, 2020; Morrison & Novikova, 2023). Automatic
34
35 1239 multimodal approaches can also allow tackling complex and innovative behavioural
36
37 1240 questions for species known to communicate in multimodal ways, such as primates
38
39 1241 (Slocombe, Waller & Liebal, 2011; Liebal & Oña, 2018) and spiders (Uetz & Roberts, 2002;
40
41 1242 Hebets, 2005). Multimodal data thus presents many advantages for automatic bioacoustic
42
43 1243 detection, all the while raising an array of limitations and adding a certain degree of
44
45 1244 complexity to machine learning solutions. Recording multimodal data is a first important
46
47 1245 challenge which can be partly addressed through the increasing availability of new efficient
48
49 1246 hardware solutions, such as lightweight, inexpensive camera traps and drones. The automatic
50
51 1247 processing of non-acoustic data is also being investigated and numerous machine learning
52
53 1248 models exist as promising solutions (Akamatsu *et al.*, 2013). Yet, the simplicity, diversity and
54
55 1249 quantity of information contained in bioacoustic data seem to make it a superior solution in
56
57
58
59
60

1
2
3 1250 most detection tasks (Enari *et al.*, 2019), at least as long as vision-based machine learning and
4
5 1251 visual recording hardware / large data storage and processing don't show significant
6
7
8 1252 improvements.

9
10
11 1253 6.2.4. Keeping a biologist in the loop

12
13
14 1254 Some of the ML models and systems are designed without the full domain knowledge or
15
16 1255 context of the problem being addressed. There needs to be close collaboration between the
17
18 1256 ML engineer designing the systems and training models, and biological scientists, as domain
19
20 1257 experts, who can validate the solutions and performance of the models. The process pipeline
21
22 1258 needs to be designed such that domain experts should closely monitor every stage from the
23
24 1259 methodology for data collection, design of data collection devices, data annotation techniques
25
26 1260 or methodology, data splits, model architecture (including inputs and outputs), and
27
28 1261 performance metrics and performance threshold values. It is also worth noting that the very
29
30 1262 same biologists may also be the ideal audience for the commercialisation of foundational
31
32 1263 models once they become available and the technologies and methods are easily accessible.
33
34 1264 The system should be iteratively improved with the active feedback from experts in the field
35
36 1265 or through the knowledge of the domain expert. This in turn maps to the process flow
37
38 1266 standardisation discussed in earlier sections.

39
40
41
42
43 1267 Since bioacoustic tasks deal with big datasets, demanding high computational power, there
44
45 1268 needs to be considerations on the environmental impacts of data storage, data transfer,
46
47 1269 computation power in terms of model training, validation or deployment in the real world.
48
49 1270 Training machine learning models is computationally very expensive and the use of GPUs
50
51 1271 results in large amounts of energy consumption. This raises the question of sustainability with
52
53 1272 respect to the research being conducted. Various independent researchers training similar
54
55 1273 models on the same datasets would result in a suboptimal use of resources. Energy
56
57 1274 consumption may be reduced by training smaller models (from model pruning, or

1
2
3 1275 “distillation”) or by sharing models. There are options of cloud storage or cloud computations
4
5 1276 (*Aide et al.*, 2013) which could benefit from the usage of green data centres in remote
6
7 1277 locations (Ministry of Local Government and Modernisation, 2021) that have green
8
9
10 1278 infrastructure for energy production (through renewable energy sources) and are perhaps less
11
12 1279 harmful to the environment rather than local GPUs or server solutions.

13
14 1280 It is also important to think of low footprint, low power usage models and systems in real
15
16 1281 world deployment for data collection or final deployment. Currently, many research studies
17
18 1282 are applying automatic detection algorithms on data that were collected in the past. We,
19
20 1283 however, anticipate that the field will move towards real-time algorithms which require
21
22 1284 systems that consume less energy in comparison to modern GPUs. To achieve this, more
23
24 1285 efforts are required within model compression, for these models to be embedded into small
25
26 1286 devices during data collection or deployment in the field.

27
28
29
30 1287 Automatic detection holds large opportunities for advances in the field of conservation and
31
32 1288 welfare, and drawing on the domain knowledge of biologists not currently involved in
33
34 1289 bioacoustics can open up new research directions. The advantages of processing large
35
36 1290 amounts of acoustic data seem clear to those currently involved in the field, but the wider
37
38 1291 biological community should be involved to find new fundamental research questions in the
39
40 1292 field of ecology and evolution (Clutton-Brock & Sheldon, 2010; De Frenne *et al.*, 2018), for
41
42 1293 example around species occurrence (Sebastián-González *et al.*, 2015; Rice *et al.*, 2021;
43
44 1294 Sattar, 2023), population density (Marques *et al.*, 2013b) and diversity (Kotera & Phillott,
45
46 1295 2022), habitat use (Brookes, Bailey & Thompson, 2013; Kotila *et al.*, 2023), phenology
47
48 1296 (*Dede et al.*, 2014; Monczak *et al.*, 2017), and the early detection of invasive species (Juanes,
49
50 1297 2018). Such questions offer opportunities for research into major conservation challenges like
51
52 1298 biodiversity loss and the effects of climate change (Sugai & Llusia, 2019; Ross *et al.*, 2023).
53
54 1299 Presently, studies driven by existing bioacoustics practitioners mostly focus on occurrence, or
55
56
57
58
59
60

1
2
3 1300 spatial or temporal distribution of a single species, whereas the advancement of automatic

4
5 1301 detection potentially allows for a focus on multiple species and to map biodiversity and

6
7 1302 potentially the functioning of whole ecosystems (Ross *et al.*, 2018).

8
9 1303 Another example of how biologists and ecologists can steer the direction in which automatic

10
11 1304 detection may be developed in the future is to identify research questions without current

12
13 1305 technological solutions. For example, although detecting signs of poor animal welfare in

14
15 1306 captivity has been the subject of many studies (Zhang *et al.*, 2022; Mao *et al.*, 2022), there

16
17 1307 are comparably very few studies investigating the of wild animals (Mcloughlin *et al.*, 2019).

18
19 1308 This is surprising given the great potential acoustic monitoring of threatened species could

20
21 1309 provide, for example on species' reproduction, or social behaviour (Teixeira, Maron &

22
23 1310 Rensburg, 2019; Greggor *et al.*, 2021).

24 25 26 27 28 29 30 1311 7. CONCLUSIONS

31 32 33 34 1312 7.1. Need for AD

35
36 1313 Automatic detection is no longer an optional capability in bioacoustics. Increasing data

37
38 1314 volumes, the need for near real-time analysis, and the expanding range of questions that

39
40 1315 biologists want to answer using passive acoustics mean that opening up the capabilities of

41
42 1316 this promising technology require parallel new developments in the field of machine learning.

43 44 45 46 47 1317 7.2. Cooperation between disciplines

48
49 1318 Mature fields in machine learning, such as image or voice recognition, are not immediately

50
51 1319 transferrable to automatic detection in bioacoustics. Close cooperation between biologist

52
53 1320 practitioners and machine learning developers will help advance solution creation by (a)

54
55 1321 enabling developers with an understanding of the problems facing bioacoustics practitioners,

1
2
3 1322 and (b) inform biologists what can and cannot be provided by the state of the art in machine
4
5 1323 learning.

6
7
8
9 1324 7.3. Deep neural networks

10
11
12 1325 Despite this, impressive advances in machine learning, particularly deep neural networks,
13
14 1326 hold out the potential for very significant developments that would cut processing time and
15
16 1327 enable a new wave of bioacoustics applications.

17
18
19
20 1328 7.4. Development pipelines

21
22
23 1329 Application development pipelines are of necessity problem-specific, however, certain
24
25 1330 guidelines and workflows (Section 4) should smooth the integration of solutions constrained
26
27 1331 both by the biological features of the problem, and by the available machine learning
28
29 1332 capabilities.

30
31
32
33
34 1333 8. ACKNOWLEDGEMENTS

35
36
37 1334 This paper arose from an investigative workshop, “Automatic detection for bioacoustics”,
38
39 1335 organised by JD and AK, supported with funding from the Cambridge Centre for Data-
40
41 1336 Driven Discovery and Accelerate Programme for Scientific Discovery, made possible by a
42
43 1337 donation from Schmidt Futures.

44
45
46 1338 ED is supported by a research chairship from the African Institute for Mathematical Sciences
47
48 1339 South Africa. ED's work was carried out with the aid of a grant from the International
49
50 1340 Development Research Centre, Ottawa, Canada, www.idrc.ca, and with financial support
51
52 1341 from the Government of Canada, provided through Global Affairs Canada (GAC),
53
54 1342 www.international.gc.ca.

55
56
57 1343 AM is supported by NERC NE/W005468/1. We thank the Agence Nationale de la Recherche
58
59 1344 (ANR) for funding RM grant ANR-20-CE23-0012-01 MIM, and PB's postdoc grant ANR-

1
2
3 1345 20-CHIA-0014. PB and JC's participation in the workshop were funded by the Institute of
4
5 1346 Convergence ILCB (ANR-16-CONV-0002), which has benefited from the support of the
6
7 1347 French government (France 2030) and is managed by the Excellence Initiative of Aix-
8
9 1348 Marseille University (A*MIDEX).

10
11
12 1349

13
14
15 1350

16
17
18
19 1351 9. BIBLIOGRAPHY

20
21 1352 ABRAHAMS, C. & GEARY, M. (2020) Combining bioacoustics and occupancy modelling for
22
23 1353 improved monitoring of rare breeding bird populations. *Ecological Indicators* **112**,
24
25 1354 106131.

26
27
28
29 1355 ACEVEDO, M.A., CORRADA-BRAVO, C.J., CORRADA-BRAVO, H., VILLANUEVA-RIVERA, L.J. &
30
31 1356 AIDE, T.M. (2009) Automated classification of bird and amphibian calls using
32
33 1357 machine learning: A comparison of methods. *Ecological Informatics* **4**, 206–214.
34
35 1358 Elsevier B.V.

36
37
38
39 1359 AIDE, T.M., CORRADA-BRAVO, C., CAMPOS-CERQUEIRA, M., MILAN, C., VEGA, G. &
40
41 1360 ALVAREZ, R. (2013) Real-time bioacoustics monitoring and automated species
42
43 1361 identification. *PeerJ* **1**, e103.

44
45
46
47 1362 AKAMATSU, T., URA, T., SUGIMATSU, H., BAHL, R., BEHERA, S., PANDA, S., KHAN, M., KAR,
48
49 1363 S.K., KAR, C.S., KIMURA, S. & SASAKI-YAMAMOTO, Y. (2013) A multimodal
50
51 1364 detection model of dolphins to estimate abundance validated by field experiments.
52
53 1365 *The Journal of the Acoustical Society of America* **134**, 2418–2426.

54
55
56
57 1366 ALAIN, G. & BENGIO, Y. (2018) Understanding intermediate layers using linear classifier
58
59 1367 probes. arXiv. <http://arxiv.org/abs/1610.01644> [accessed 5 July 2023].

- 1
2
3 1368 ANON. (2023a) Kaleidoscope Pro Analysis Software. *Wildlife Acoustics*.
4
5 1369 <https://www.wildlifeacoustics.com/products/kaleidoscope-pro> [accessed 12 August
6
7 1370 2023].
8
9
10
11 1371 ANON. (2023b) BTO Acoustic Pipeline. *BTO - British Trust for Ornithology*.
12
13 1372 <https://www.bto.org/our-science/products-and-technologies/bto-acoustic-pipeline>
14
15 1373 [accessed 12 August 2023].
16
17
18
19 1374 AODHA, O.M., GIBB, R., BARLOW, K.E., BROWNING, E., FIRMAN, M., FREEMAN, R., HARDER,
20
21 1375 B., KINSEY, L., MEAD, G.R., NEWSON, S.E., PANDOURSКИ, I., PARSONS, S., RUSS, J.,
22
23 1376 SZODORAY-PARADI, A., SZODORAY-PARADI, F., ET AL. (2018) Bat detective—Deep
24
25 1377 learning tools for bat acoustic signal detection. *PLOS Computational Biology* **14**,
26
27 1378 e1005995. Public Library of Science.
28
29
30
31 1379 BAEVSKI, A., ZHOU, Y., MOHAMED, A. & AULI, M. (2020) wav2vec 2.0: A framework for
32
33 1380 self-supervised learning of speech representations. *Advances in Neural Information*
34
35 1381 *Processing Systems (NeurIPS)* **33**, 12449–12460.
36
37
38
39 1382 BAILEY, L.L., MACKENZIE, D.I. & NICHOLS, J.D. (2014) Advances and applications of
40
41 1383 occupancy models. *Methods in Ecology and Evolution* **5**, 1269–1279.
42
43
44
45 1384 BARKER, D.J., HERRERA, C. & WEST, M.O. (2014) Automated detection of 50-kHz ultrasonic
46
47 1385 vocalizations using template matching in XBAT. *Journal of Neuroscience Methods*
48
49 1386 **236**, 68–75.
50
51
52
53 1387 BEE, M.A. (2012) Sound source perception in anuran amphibians. *Current Opinion in*
54
55 1388 *Neurobiology* **22**, 301–310.
56
57
58
59
60

- 1
2
3 1389 BERGLER, C., BARNHILL, A., PERRIN, D., SCHMITT, M., MAIER, A. & NÖTH, E. (2022a)
4
5 1390 ORCA-WHISPER: An Automatic Killer Whale Sound Type Generation Toolkit
6
7 1391 Using Deep Learning. In *Interspeech 2022* pp. 2413–2417. ISCA.
8
9
10
11 1392 BERGLER, C., SCHMITT, M., MAIER, A., SMEELE, S., BARTH, V. & NÖTH, E. (2020) ORCA-
12
13 1393 CLEAN: A Deep Denoising Toolkit for Killer Whale Communication. In *Interspeech*
14
15 1394 *2020* pp. 1136–1140. ISCA.
16
17
18
19 1395 BERGLER, C., SMEELE, S.Q., TYNDEL, S.A., BARNHILL, A., ORTIZ, S.T., KALAN, A.K.,
20
21 1396 CHENG, R.X., BRINKLØV, S., OSIECKA, A.N., TOUGAARD, J., JAKOBSEN, F.,
22
23 1397 WAHLBERG, M., NÖTH, E., MAIER, A. & KLUMP, B.C. (2022b) ANIMAL-SPOT
24
25 1398 enables animal-independent signal detection and classification using deep learning.
26
27 1399 *Scientific Reports* **12**, 21966. Nature Publishing Group.
28
29
30
31 1400 BERMANT, P.C. (2021) BioCPPNet: automatic bioacoustic source separation with deep neural
32
33 1401 networks. *Scientific Reports* **11**, 23502. Nature Publishing Group.
34
35
36
37 1402 BERMANT, P.C., BRONSTEIN, M.M., WOOD, R.J., GERO, S. & GRUBER, D.F. (2019) Deep
38
39 1403 Machine Learning Techniques for the Detection and Classification of Sperm Whale
40
41 1404 Bioacoustics. *Scientific Reports* **9**, 12588. Nature Publishing Group.
42
43
44
45 1405 BEST, P., FERRARI, M., POUPARD, M., PARIS, S., MARXER, R., SYMONDS, H., SPONG, P. &
46
47 1406 GLOTIN, H. (2020) Deep Learning and Domain Transfer for Orca Vocalization
48
49 1407 Detection. In *2020 International Joint Conference on Neural Networks (IJCNN)* pp.
50
51 1408 1–7.
52
53
54
55 1409 BEST, P., MARXER, R., PARIS, S. & GLOTIN, H. (2022) Temporal evolution of the
56
57 1410 Mediterranean fin whale song. *Scientific Reports* **12**, 13565. Nature Publishing Group.
58
59
60

- 1
2
3 1411 BEST, P., PARIS, S., GLOTIN, H. & MARXER, R. (2023) Deep audio embeddings for
4
5 1412 vocalisation clustering. *PLOS ONE* **18**, e0283396. Public Library of Science.
6
7
8
9 1413 BOAKES, E.H., MCGOWAN, P.J.K., FULLER, R.A., CHANG-QING, D., CLARK, N.E., O'CONNOR,
10
11 1414 K. & MACE, G.M. (2010) Distorted Views of Biodiversity: Spatial and Temporal Bias
12
13 1415 in Species Occurrence Data. *PLOS Biology* **8**, e1000385. Public Library of Science.
14
15
16 1416 BOERSMA, P. & WEENINK, D. (2007) PRAAT: Doing phonetics by computer (Version
17
18 1417 5.3.51).
19
20
21
22 1418 BOMMASANI, R., HUDSON, D.A., ADELI, E., ALTMAN, R., ARORA, S., VON ARX, S., BERNSTEIN,
23
24 1419 M.S., BOHG, J., BOSSELUT, A., BRUNSKILL, E., BRYNJOLFSSON, E., BUCH, S., CARD,
25
26 1420 D., CASTELLON, R., CHATTERJI, N., ET AL. (2022) On the Opportunities and Risks of
27
28 1421 Foundation Models. arXiv. <http://arxiv.org/abs/2108.07258> [accessed 7 August 2023].
29
30
31
32 1422 BØTTCHER, A., GERO, S., BEEDHOLM, K., WHITEHEAD, H. & MADSEN, P.T. (2018) Variability
33
34 1423 of the inter-pulse interval in sperm whale clicks with implications for size estimation
35
36 1424 and individual identification. *The Journal of the Acoustical Society of America* **144**,
37
38 1425 365–374.
39
40
41
42
43 1426 BOUCHET, H., BLOIS-HEULIN, C. & LEMASSON, A. (2013) Social complexity parallels vocal
44
45 1427 complexity: a comparison of three non-human primate species. *Frontiers in*
46
47 1428 *Psychology* **4**.
48
49
50 1429 BRAGA, P.H.P., HÉBERT, K., HUDGINS, E.J., SCOTT, E.R., EDWARDS, B.P.M., SÁNCHEZ REYES,
51
52 1430 L.L., GRAINGER, M.J., FOROUGHIRAD, V., HILLEMANN, F., BINLEY, A.D., BROOKSON,
53
54 1431 C.B., GAYNOR, K.M., SHAFIEI SABET, S., GÜNCAN, A., WEIERBACH, H., ET AL. (2023)
55
56 1432 Not just for programmers: How GitHub can accelerate collaborative and reproducible
57
58 1433 research in ecology and evolution. *Methods in Ecology and Evolution* **14**, 1364–1380.
59
60

- 1
2
3 1434 BROLL, B. & WHITAKER, J. (2017) DeepForge: An open source, collaborative environment
4
5 1435 for reproducible deep learning.
6
7
8
9 1436 BROOKES, K.L., BAILEY, H. & THOMPSON, P.M. (2013) Predictions from harbor porpoise
10
11 1437 habitat association models are confirmed by long-term passive acoustic monitoring.
12
13 1438 *The Journal of the Acoustical Society of America* **134**, 2523–2533.
14
15
16
17 1439 BROWNING, E., GIBB, R., GLOVER-KAPFER, P. & JONES, K.E. (2017) Passive acoustic
18
19 1440 monitoring in ecology and conservation. Report, WWF-UK.
20
21
22 1441 BURGOS, G. & ZUBEROGOITIA, I. (2020) A telemetry study to discriminate between home
23
24 1442 range and territory size in Tawny Owls. *Bioacoustics* **29**, 109–121. Taylor & Francis.
25
26
27
28 1443 BUXTON, R.T., MCKENNA, M.F., CLAPP, M., MEYER, E., STABENAU, E., ANGELONI, L.M.,
29
30 1444 CROOKS, K. & WITTEMYER, G. (2018) Efficacy of extracting indices from large-scale
31
32 1445 acoustic recordings to monitor biodiversity. *Conservation Biology* **32**, 1174–1184.
33
34
35
36 1446 CANNAM, C., LANDONE, C. & SANDLER, M. (2010) Sonic visualiser: an open source
37
38 1447 application for viewing, analysing, and annotating music audio files. In *Proceedings*
39
40 1448 *of the 18th ACM international conference on Multimedia* pp. 1467–1468. Association
41
42 1449 for Computing Machinery, New York, NY, USA.
43
44
45
46 1450 CASAER, J., MILOTIC, T., LIEFTING, Y., DESMET, P. & JANSEN, P. (2019) Agouti: A platform
47
48 1451 for processing and archiving of camera trap images. *Biodiversity Information Science*
49
50 1452 *and Standards* **3**, e46690.
51
52
53
54 1453 CASTELLOTE, M. & FOSSA, F. (2006) Measuring acoustic activity as a method to evaluate
55
56 1454 welfare in captive beluga whales (*Delphinapterus leucas*). *Aquatic Mammals* **32**, 325–
57
58 1455 333.
59
60

- 1
2
3 1456 CLARK, F.E. & DUNN, J.C. (2022) From Soundwave to Soundscape: A Guide to Acoustic
4
5 1457 Research in Captive Animal Environments. *Frontiers in Veterinary Science* **9**.
6
7
8
9 1458 CLARK, M.L., SALAS, L., BALIGAR, S., QUINN, C.A., SNYDER, R.L., LELAND, D.,
10
11 1459 SCHACKWITZ, W., GOETZ, S.J. & NEWSAM, S. (2023) The effect of soundscape
12
13 1460 composition on bird vocalization classification in a citizen science biodiversity
14
15 1461 monitoring project. *Ecological Informatics* **75**, 102065.
16
17
18
19 1462 CLINK, D.J., KIER, I., AHMAD, A.H. & KLINCK, H. (2023) A workflow for the automated
20
21 1463 detection and classification of female gibbon calls from long-term acoustic
22
23 1464 recordings. *Frontiers in Ecology and Evolution* **11**.
24
25
26
27 1465 CLUTTON-BROCK, T. & SHELDON, B.C. (2010) Individuals and populations: the role of long-
28
29 1466 term, individual-based studies of animals in ecology and evolutionary biology. *Trends*
30
31 1467 *in Ecology & Evolution* **25**, 562–573.
32
33
34
35 1468 ÇOBAN, E.B., PIR, D., SO, R. & MANDEL, M.I. (2020) Transfer Learning from Youtube
36
37 1469 Soundtracks to Tag Arctic Ecoacoustic Recordings. In *ICASSP 2020 - 2020 IEEE*
38
39 1470 *International Conference on Acoustics, Speech and Signal Processing (ICASSP)* pp.
40
41 1471 726–730.
42
43
44
45 1472 COFFEY, K.R., MARX, R.E. & NEUMAIER, J.F. (2019) DeepSqueak: a deep learning-based
46
47 1473 system for detection and analysis of ultrasonic vocalizations.
48
49 1474 *Neuropsychopharmacology* **44**, 859–868.
50
51
52
53 1475 COHEN, Y., NICHOLSON, D.A., SANCHIONI, A., MALLABER, E.K., SKIDANOVA, V. &
54
55 1476 GARDNER, T.J. (2022) Automated annotation of birdsong with a neural network that
56
57 1477 segments spectrograms. *eLife* **11**, e63853. eLife Sciences Publications, Ltd.
58
59
60

- 1
2
3 1478 COLE, J.S., MICHEL, N.L., EMERSON, S.A. & SIEGEL, R.B. (2022) Automated bird sound
4
5 1479 classifications of long-duration recordings produce occupancy model outputs similar
6
7 1480 to manually annotated data. *Ornithological Applications* **124**, duac003.
8
9
10
11 1481 CUEVAS, A., VERAGUA, A., ESPAÑOL-JIMÉNEZ, S., CHIANG, G. & TOBAR, F. (2017)
12
13 1482 Unsupervised blue whale call detection using multiple time-frequency features. In
14
15 1483 *2017 CHILEAN Conference on Electrical, Electronics Engineering, Information and*
16
17 1484 *Communication Technologies (CHILECON)* pp. 1–6.
18
19
20
21 1485 DARWIN CORE TASK GROUP (2009) Darwin Core. <https://www.tdwg.org/standards/dwc/>
22
23 1486 [accessed 5 July 2023].
24
25
26
27 1487 DAVIS, J. & GOADRICH, M. (2006) The relationship between Precision-Recall and ROC
28
29 1488 curves. In *Proceedings of the 23rd international conference on Machine learning* pp.
30
31 1489 233–240. Association for Computing Machinery, New York, NY, USA.
32
33
34
35 1490 DAWSON, D.K. & EFFORD, M.G. (2009) Bird population density estimated from acoustic
36
37 1491 signals. *Journal of Applied Ecology* **46**, 1201–1209. WILEY-BLACKWELL.
38
39
40
41 1492 DE FRENNE, P., VAN LANGENHOVE, L., VAN DRIESSCHE, A., BERTRAND, C., VERHEYEN, K. &
42
43 1493 VANGANSBEKE, P. (2018) Using archived television video footage to quantify
44
45 1494 phenology responses to climate change. *Methods in Ecology and Evolution* **9**, 1874–
46
47 1495 1882.
48
49
50
51 1496 DEDE, A., ÖZTÜRK, A.A., AKAMATSU, T., TONAY, A.M. & ÖZTÜRK, B. (2014) Long-term
52
53 1497 passive acoustic monitoring revealed seasonal and diel patterns of cetacean presence
54
55 1498 in the Istanbul Strait. *Journal of the Marine Biological Association of the United*
56
57 1499 *Kingdom* **94**, 1195–1202.
58
59
60

- 1
2
3 1500 DENTON, T., WISDOM, S. & HERSHEY, J.R. (2021) Improving Bird Classification with
4
5 1501 Unsupervised Sound Separation. arXiv. <http://arxiv.org/abs/2110.03209> [accessed 6
6
7 1502 July 2023].
8
9
10
11 1503 DERRICKSON, K.C. (1988) Variation in Repertoire Presentation in Northern Mockingbirds.
12
13 1504 *The Condor* **90**, 592–606.
14
15
16 1505 DUAN, S., ZHANG, J., ROE, P., WIMMER, J., DONG, X., TRUSKINGER, A. & TOWSEY, M. (2013)
17
18 1506 Timed Probabilistic Automaton: A Bridge between Raven and Song Scope for
19
20 1507 Automatic Species Recognition. *Proceedings of the AAAI Conference on Artificial*
21
22 1508 *Intelligence* **27**, 1519–1524.
23
24
25
26 1509 DUFOURQ, E., BATIST, C., FOQUET, R. & DURBACH, I. (2022a) Passive acoustic monitoring of
27
28 1510 animal populations with transfer learning. *Ecological Informatics* **70**, 101688.
29
30
31
32 1511 DUFOURQ, E., BATIST, C., FOQUET, R. & DURBACH, I. (2022b) Passive acoustic monitoring of
33
34 1512 animal populations with transfer learning. *Ecological Informatics* **70**, 101688.
35
36
37
38 1513 DUFOURQ, E., DURBACH, I., HANSFORD, J.P., HOEPFNER, A., MA, H., BRYANT, J.V., STENDER,
39
40 1514 C.S., LI, W., LIU, Z., CHEN, Q., ZHOU, Z. & TURVEY, S.T. (2021) Automated detection
41
42 1515 of Hainan gibbon calls for passive acoustic monitoring. *Remote Sensing in Ecology*
43
44 1516 *and Conservation* **7**, 475–487.
45
46
47
48 1517 DUNN, J.C. & SMAERS, J.B. (2018) Neural Correlates of Vocal Repertoire in Primates.
49
50 1518 *Frontiers in Neuroscience* **12**.
51
52
53
54 1519 ENARI, H., ENARI, H.S., OKUDA, K., MARUYAMA, T. & OKUDA, K.N. (2019) An evaluation of
55
56 1520 the efficiency of passive acoustic monitoring in detecting deer and primates in
57
58 1521 comparison with camera traps. *Ecological Indicators* **98**, 753–762.
59
60

- 1
2
3 1522 ERBE, C. & THOMAS, J.A. (eds) (2022) *Exploring Animal Behavior Through Sound: Volume*
4
5 1523 *1: Methods*. Springer Nature.
6
7
8
9 1524 EXADAKTYLOS, V., SILVA, M., AERTS, J.-M., TAYLOR, C.J. & BERCKMANS, D. (2008) Real-
10
11 1525 time recognition of sick pig cough sounds. *Computers and Electronics in Agriculture*
12
13 1526 **63**, 207–214.
14
15
16
17 1527 FAIRBRASS, A.J., FIRMAN, M., WILLIAMS, C., BROSTOW, G.J., TITHERIDGE, H. & JONES, K.E.
18
19 1528 (2019) CityNet—Deep learning tools for urban ecoacoustic assessment. *Methods in*
20
21 1529 *Ecology and Evolution* **10**, 186–197.
22
23
24
25 1530 FEDDERSEN-PETERSEN, D.U. (2000) Vocalization of European wolves (<i>Canis lupus
26
27 1531 *lupus*; L.) and various dog breeds (<i>Canis lupus; f. fam.).
28
29 1532 *Archives Animal Breeding* **43**, 387–398.
30
31
32
33 1533 FERRARI, M., GLOTIN, H., MARXER, R. & ASCH, M. (2020) DOCC10: Open access dataset of
34
35 1534 marine mammal transient studies and end-to-end CNN classification. In *2020*
36
37 1535 *International Joint Conference on Neural Networks (IJCNN)* pp. 1–8.
38
39
40
41 1536 FORD, J.K.B. (1991) Vocal traditions among resident killer whales (*Orcinus orca*) in coastal
42
43 1537 waters of British Columbia. *Canadian Journal of Zoology* **69**, 1454–1483. NRC
44
45 1538 Research Press.
46
47
48
49 1539 FOX, E.J.S., ROBERTS, J.D. & BENNAMOUN, M. (2008) Call-Independent Individual
50
51 1540 Identification in Birds. *Bioacoustics* **18**, 51–67. Taylor & Francis.
52
53
54
55 1541 FRICK, W.F. (2013) Acoustic monitoring of bats, considerations of options for long-term
56
57 1542 monitoring. *Therya* **4**, 69–70. Centro de Investigaciones Biológicas del Noroeste.
58
59
60

- 1
2
3 1543 FROMMOLT, K.-H. & TAUCHERT, K.-H. (2014) Applying bioacoustic methods for long-term
4
5 1544 monitoring of a nocturnal wetland bird. *Ecological Informatics* **21**, 4–12.
6
7
8
9 1545 GANCHEV, T.D. (2020) Chapter 8 - Ubiquitous computing and biodiversity monitoring. In
10
11 1546 *Advances in Ubiquitous Computing* (ed A. NEUSTEIN), pp. 239–259. Academic Press.
12
13
14 1547 GARETH JAMES, DANIELA WITTEN, TREVOR HASTIE, & ROBERT TIBSHIRANI (2013) *An*
15
16 1548 *Introduction to Statistical Learning: With Applications in R*. Springer, New York,
17
18 1549 NY.
19
20
21
22 1550 GARLAND, E.C., CASTELLOTE, M. & BERCHOK, C.L. (2015) Beluga whale (*Delphinapterus*
23
24 1551 *leucas*) vocalizations and call classification from the eastern Beaufort Sea population.
25
26 1552 *The Journal of the Acoustical Society of America* **137**, 3054–3067.
27
28
29
30 1553 GARLAND, E.C., GOLDIZEN, A.W., REKDAHL, M.L., CONSTANTINE, R., GARRIGUE, C.,
31
32 1554 HAUSER, N.D., POOLE, M.M., ROBBINS, J. & NOAD, M.J. (2011) Dynamic Horizontal
33
34 1555 Cultural Transmission of Humpback Whale Song at the Ocean Basin Scale. *Current*
35
36 1556 *Biology* **21**, 687–691. Elsevier.
37
38
39
40 1557 GBIF/TDWG MULTIMEDIA RESOURCES TASK GROUP (2013) Audiovisual Core Multimedia
41
42 1558 Resources Metadata Schema. <https://www.tdwg.org/standards/ac/> [accessed 5 July
43
44 1559 2023].
45
46
47
48 1560 GIBBON, D., MOORE, R. & WINSKI, R. (eds) (1998) *Vol 1 Spoken Language System and*
49
50 1561 *Corpus Design*. De Gruyter Mouton, Berlin, Boston.
51
52
53
54 1562 GILLESPIE, D., GORDON, J., MCHUGH, R., MCLAREN, D., MELLINGER, D.K., REDMOND, P.,
55
56 1563 THODE, A., TRINDER, P. & DENG, X.Y. (2009) PAMGUARD: Semiautomated, open
57
58 1564 source software for real-time acoustic detection and localisation of cetaceans. In *The*
59
60

- 1
2
3 1565 *Journal of the Acoustical Society of America* pp. 2547–2547. Acoustical Society of
4
5 1566 AmericaASA.
6
7
8
9 1567 GREEN, A., CLARK, C., FAVARO, L., LOMAX, S. & REBY, D. (2019) Vocal individuality of
10
11 1568 Holstein-Friesian cattle is maintained across putatively positive and negative farming
12
13 1569 contexts. *Scientific Reports* **9**. Nature Research.
14
15
16
17 1570 GREGGOR, A.L., MASUDA, B., GAUDIOSO-LEVITA, J.M., NELSON, J.T., WHITE, T.H., SHIER,
18
19 1571 D.M., FARABAUGH, S.M. & SWAISGOOD, R.R. (2021) Pre-release training, predator
20
21 1572 interactions and evidence for persistence of anti-predator behavior in reintroduced
22
23 1573 `alalā, Hawaiian crow. *Global Ecology and Conservation* **28**, e01658.
24
25
26
27 1574 HAGIWARA, M. (2022) AVES: Animal Vocalization Encoder based on Self-Supervision.
28
29 1575 arXiv. <http://arxiv.org/abs/2210.14493> [accessed 6 July 2023].
30
31
32 1576 HAGIWARA, M., HOFFMAN, B., LIU, J.-Y., CUSIMANO, M., EFFENBERGER, F. & ZACARIAN, K.
33
34 1577 (2022) BEANS: The Benchmark of Animal Sounds. arXiv.
35
36 1578 <http://arxiv.org/abs/2210.12300> [accessed 6 July 2023].
37
38
39
40 1579 HANSEN, P. (1979) Vocal learning: Its role in adapting sound structures to long-distance
41
42 1580 propagation, and a hypothesis on its evolution. *Animal Behaviour* **27**, 1270–1271.
43
44 1581 Elsevier Science, Netherlands.
45
46
47
48 1582 HARRINGTON, F.H. & MECH, L.D. (1982) An Analysis of Howling Response Parameters
49
50 1583 Useful for Wolf Pack Censusing. *The Journal of Wildlife Management* **46**, 686–693.
51
52 1584 Allen Press.
53
54
55
56 1585 HAUPERT, S., SÈBE, F. & SUEUR, J. (2022) Physics-based model to predict the acoustic
57
58 1586 detection distance of terrestrial autonomous recording units over the diel cycle and
59
60

- 1
2
3 1587 across seasons: Insights from an Alpine and a Neotropical forest. *Methods in Ecology*
4
5 1588 *and Evolution*. Wiley.
6
7
8
9 1589 HEAPHY, K. & CAIN, K. (2021) Song variation between sexes and among subspecies of New
10
11 1590 Zealand Fantail (*Rhipidura fuliginosa*). *Emu - Austral Ornithology* **121**, 198–210.
12
13 1591 Taylor & Francis.
14
15
16 1592 HEATH, B.E., SETHI, S.S., ORME, C.D.L., EWERS, R.M. & PICINALI, L. (2021) How index
17
18 1593 selection, compression, and recording schedule impact the description of ecological
19
20 1594 soundscapes. *Ecology and Evolution* **11**, 13206–13217.
21
22
23
24 1595 HEBETS, E.A. (2005) Attention-altering signal interactions in the multimodal courtship
25
26 1596 display of the wolf spider *Schizocosa uetzi*. *Behavioral Ecology* **16**, 75–82.
27
28
29
30 1597 HEBETS, E.A., BERN, M., MCGINLEY, R.H., ROBERTS, A., KERSHENBAUM, A., STARRETT, J. &
31
32 1598 BOND, J.E. (2021) Sister species diverge in modality-specific courtship signal form
33
34 1599 and function. *Ecology and Evolution* **11**, 852–871.
35
36
37
38 1600 HEDLEY, R.W., WILSON, S.J., YIP, D.A., LI, K. & BAYNE, E.M. (2021) Distance truncation
39
40 1601 via sound level for bioacoustic surveys in patchy habitat. *Bioacoustics* **30**, 303–323.
41
42
43
44 1602 HENDRY, H. & MANN, C. (2018) Camelot—intuitive software for camera-trap data
45
46 1603 management. *Oryx* **52**, 15–15.
47
48
49 1604 HILDEBRAND, J.A., FRASIER, K.E., HELBLE, T.A. & ROCH, M.A. (2022) Performance metrics
50
51 1605 for marine mammal signal detection and classification. *The Journal of the Acoustical*
52
53 1606 *Society of America* **151**, 414–427.
54
55
56
57
58
59
60

- 1
2
3 1607 HILL, A.P., PRINCE, P., SNADDON, J.L., DONCASTER, C.P. & ROGERS, A. (2019) AudioMoth:
4
5 1608 A low-cost acoustic device for monitoring biodiversity and the environment.
6
7 1609 *HardwareX* **6**, e00073.
8
9
10
11 1610 HOGEWEG, L. & STOWELL, D. (2023) An API for AI species recognition. *Arise*.
12
13 1611 <https://www.arise-biodiversity.nl/post/an-api-for-ai-species-recognition> [accessed 4
14
15 1612 July 2023].
16
17
18
19 1613 HOOD, J.D., FLOGERAS, D.G. & THERIAULT, J.A. (2016) Improved passive acoustic band-
20
21 1614 limited energy detection for cetaceans. *Applied Acoustics* **106**, 36–41.
22
23
24
25 1615 HSU, W.-N., BOLTE, B., TSAI, Y.-H.H., LAKHOTIA, K., SALAKHUTDINOV, R. & MOHAMED, A.
26
27 1616 (2021) HuBERT: Self-Supervised Speech Representation Learning by Masked
28
29 1617 Prediction of Hidden Units. arXiv. <http://arxiv.org/abs/2106.07447> [accessed 6 July
30
31 1618 2023].
32
33
34
35 1619 HUMPHREY, E.J., SALAMON, J., NIETO, O., FORSYTH, J., BITTNER, R.M. & BELLO, J.P. (2014)
36
37 1620 JAMS: A JSON ANNOTATED MUSIC SPECIFICATION FOR REPRODUCIBLE
38
39 1621 MIR RESEARCH.
40
41
42
43 1622 JANSSON, A., HUMPHREY, E., MONTECCHIO, N., BITTNER, R., KUMAR, A. & WEYDE, T.
44
45 1623 (2017) Singing voice separation with deep U-Net convolutional networks. conference,
46
47 1624 Suzhou, China. <https://ismir2017.smcnus.org/> [accessed 6 July 2023].
48
49
50
51 1625 JIANG, J., BU, L., WANG, X., LI, C., SUN, Z., YAN, H., HUA, B., DUAN, F. & YANG, J. (2018)
52
53 1626 Clicks classification of sperm whale and long-finned pilot whale based on continuous
54
55 1627 wavelet transform and artificial neural network. *Applied Acoustics* **141**, 26–34.
56
57
58
59
60

- 1
2
3 1628 JOLY, A., GOËAU, H., GLOTIN, H., SPAMPINATO, C., BONNET, P., VELLINGA, W.-P., PLANQUÉ,
4
5 1629 R., RAUBER, A., PALAZZO, S., FISHER, B. & MÜLLER, H. (2015) LifeCLEF 2015:
6
7 1630 Multimedia Life Species Identification Challenges. In *Experimental IR Meets*
8
9 1631 *Multilinguality, Multimodality, and Interaction* (eds J. MOTHE, J. SAVOY, J. KAMPS,
10
11 1632 K. PINEL-SAUVAGNAT, G. JONES, E. SAN JUAN, L. CAPELLATO & N. FERRO), pp. 462–
12
13 1633 483. Springer International Publishing, Cham.
- 14
15
16
17
18 1634 JUANES, F. (2018) Visual and acoustic sensors for early detection of biological invasions:
19
20 1635 Current uses and future potential. *Journal for Nature Conservation* **42**, 7–11.
- 21
22
23 1636 K. LISA YANG CENTER FOR CONSERVATION BIOACOUSTICS (2014) Bioacoustics Research
24
25 1637 Program. <https://ravensoundsoftware.com/> [accessed 5 July 2023].
- 26
27
28
29 1638 KAHL, S., WOOD, C.M., EIBL, M. & KLINCK, H. (2021) BirdNET: A deep learning solution
30
31 1639 for avian diversity monitoring. *Ecological Informatics* **61**, 101236.
- 32
33
34
35 1640 KARNAN, M., AKILA, M. & KRISHNARAJ, N. (2011) Biometric personal authentication using
36
37 1641 keystroke dynamics: A review. *Applied Soft Computing* **11**, 1565–1573.
- 38
39
40 1642 KERSHENBAUM, A., BLUMSTEIN, D.T., ROCH, M.A., AKÇAY, Ç., BACKUS, G., BEE, M.A.,
41
42 1643 BOHN, K., CAO, Y., CARTER, G., CÄSAR, C., COEN, M., DERUITER, S.L., DOYLE, L.,
43
44 1644 EDELMAN, S., FERRER-I-CANCHO, R., ET AL. (2016a) Acoustic sequences in non-
45
46 1645 human animals: a tutorial review and prospectus. *Biological Reviews* **91**, 13–52.
- 47
48
49
50 1646 KERSHENBAUM, A., DEMARTSEV, V., GAMMON, D.E., GEFFEN, E., GUSTISON, M.L., ILANY, A.
51
52 1647 & LAMEIRA, A.R. (2021) Shannon entropy as a robust estimator of Zipf's Law in
53
54 1648 animal vocal communication repertoires. *Methods in Ecology and Evolution* **12**, 553–
55
56 1649 564.
- 57
58
59
60

- 1
2
3 1650 KERSHENBAUM, A., ILANY, A., BLAUSTEIN, L. & GEFFEN, E. (2012) Syntactic structure and
4
5 1651 geographical dialects in the songs of male rock hyraxes. *Proceedings of the Royal*
6
7 1652 *Society B: Biological Sciences* **279**, 2974–2981. Royal Society.
8
9
10
11 1653 KERSHENBAUM, A., OWENS, J.L. & WALLER, S. (2019) Tracking cryptic animals using
12
13 1654 acoustic multilateration: A system for long-range wolf detection. *The Journal of the*
14
15 1655 *Acoustical Society of America* **145**, 1619–1628.
16
17
18
19 1656 KERSHENBAUM, A. & ROCH, M.A. (2013) An image processing based paradigm for the
20
21 1657 extraction of tonal sounds in cetacean communications. *The Journal of the Acoustical*
22
23 1658 *Society of America* **134**, 4435–4445.
24
25
26
27 1659 KERSHENBAUM, A., ROOT-GUTTERIDGE, H., HABIB, B., KOLER-MATZNICK, J., MITCHELL, B.,
28
29 1660 PALACIOS, V. & WALLER, S. (2016b) Disentangling canid howls across multiple
30
31 1661 species and subspecies: Structure in a complex communication channel. *Behavioural*
32
33 1662 *Processes* **124**, 149–157.
34
35
36
37 1663 KERSHENBAUM, A., SAYIGH, L.S. & JANIK, V.M. (2013) The Encoding of Individual Identity
38
39 1664 in Dolphin Signature Whistles: How Much Information Is Needed? *PLOS ONE* **8**,
40
41 1665 e77671. Public Library of Science.
42
43
44
45 1666 KIMURA, S., AKAMATSU, T., LI, S., DONG, S., DONG, L., WANG, K., WANG, D. & ARAI, N.
46
47 1667 (2010) Density estimation of Yangtze finless porpoises using passive acoustic sensors
48
49 1668 and automated click train detection. *The Journal of the Acoustical Society of America*
50
51 1669 **128**, 1435.
52
53
54
55 1670 KNIGHT, E.C. & BAYNE, E.M. (2019) Classification threshold and training data affect the
56
57 1671 quality and utility of focal species data processed with automated audio-recognition
58
59 1672 software. *Bioacoustics* **28**, 539–554.
60

- 1
2
3 1673 KOTERA, M.M. & PHILLOTT, A.D. (2022) Calls for conservation: A review of bioacoustics
4
5 1674 monitoring with case studies from India. *Asian Journal of Environment & Ecology*,
6
7 1675 142–150.
8
9
10
11 1676 KOTILA, M., SUOMINEN, K.M., VASKO, V.V., BLOMBERG, A.S., LEHIKONEN, A., ANDERSSON,
12
13 1677 T., ASPI, J., CEDERBERG, T., HÄNNINEN, J., INKINEN, J., KOSKINEN, J., LUNDBERG, G.,
14
15 1678 MÄKINEN, K., RONTTI, M., SNICKARS, M., ET AL. (2023) Large-scale long-term
16
17 passive-acoustic monitoring reveals spatio-temporal activity patterns of boreal bats.
18 1679
19
20 1680 *Ecography* **2023**, e06617.
21
22
23 1681 KRAUSE, B.L. (1993) The niche hypothesis: a virtual symphony of animal sounds, the origins
24
25 of musical expression and the health of habitats. *The Soundscape Newsletter* **6**, 6–10.
26 1682
27
28
29 1683 LAIOLO, P., ROLANDO, A., DELESTRADE, A. & SANCTIS, A. DE (2001) GEOGRAPHICAL
30
31 1684 VARIATION IN THE CALLS OF THE CHOUGHS. *The Condor* **103**, 287–297.
32
33 1685 American Ornithological Society.
34
35
36
37 1686 LAURIJS, K.A., BRIEFER, E.F., REIMERT, I. & WEBB, L.E. (2021) Vocalisations in farm
38
39 1687 animals: A step towards positive welfare assessment. *Applied Animal Behaviour*
40
41 1688 *Science* **236**, 105264.
42
43
44
45 1689 LAW, B., GONSALVES, L., BURGAR, J., BRASSIL, T., KERR, I., WILMOTT, L., MADDEN, K.,
46
47 1690 SMITH, M., MELLA, V., CROWTHER, M., KROCKENBERGER, M., RUS, A., PIETSCH, R.,
48
49 1691 TRUSKINGER, A., EICHINSKI, P., ET AL. (2021) Estimating and validating koala.
50
51 1692 *Wildlife Research* **49**, 438–448.
52
53
54
55 1693 LECUN, Y., BENGIO, Y. & HINTON, G. (2015) Deep learning. *Nature* **521**, 436–444.
56
57
58
59
60

- 1
2
3 1694 LEIGHTON, G.M. & BIRMINGHAM, T. (2021) Multiple factors affect the evolution of repertoire
4
5 1695 size across birds. *Behavioral Ecology* **32**, 380–385.
6
7
8
9 1696 LEROUX, M., AL-KHUDHAIRY, O.G., PERONY, N. & TOWNSEND, S.W. (2021) Chimpanzee
10
11 1697 voice prints? Insights from transfer learning experiments from human voices. arXiv.
12
13 1698 <http://arxiv.org/abs/2112.08165> [accessed 6 July 2023].
14
15
16
17 1699 LEVER, J., KRZYWINSKI, M. & ALTMAN, N. (2016) Classification evaluation. *Nature Methods*
18
19 1700 **13**, 603–604. Nature Publishing Group.
20
21
22 1701 LI, N., REN, Z., LI, D. & ZENG, L. (2020) Automated techniques for monitoring the behaviour
23
24 1702 and welfare of broilers and laying hens: towards the goal of precision livestock
25
26 1703 farming. *Animal* **14**, 617–625.
27
28
29
30 1704 LIEBAL, K. & OÑA, L. (2018) Different Approaches to Meaning in Primate Gestural and
31
32 1705 Vocal Communication. *Frontiers in Psychology* **9**.
33
34
35
36 1706 LIN, T., WANG, Y., LIU, X. & QIU, X. (2022) A survey of transformers. *AI Open* **3**, 111–132.
37
38
39 1707 LONG, R.A. (2008) *Noninvasive survey methods for carnivores*. In p. 385. Island Press.
40
41
42 1708 LOSTANLEN, V., SALAMON, J., FARNSWORTH, A., KELLING, S. & BELLO, J.P. (2018) Birdvox-
43
44 1709 Full-Night: A Dataset and Benchmark for Avian Flight Call Detection. In *2018 IEEE*
45
46 1710 *International Conference on Acoustics, Speech and Signal Processing (ICASSP)* pp.
47
48 1711 266–270.
49
50
51
52 1712 MAATEN, L. VAN DER & HINTON, G. (2008) Visualizing Data using t-SNE. *Journal of*
53
54 1713 *Machine Learning Research* **9**, 2579–2605.
55
56
57
58
59
60

- 1
2
3 1714 MACKENZIE, D.I., NICHOLS, J.D., HINES, J.E., KNUTSON, M.G. & FRANKLIN, A.B. (2003)
4
5 1715 Estimating Site Occupancy, Colonization, and Local Extinction When a Species Is
6
7 1716 Detected Imperfectly. *Ecology* **84**, 2200–2207.
- 9
10
11 1717 MANSER, M.B., JANSEN, D.A.W.A.M., GRAW, B., HOLLÉN, L.I., BOUSQUET, C.A.H., FURRER,
12
13 1718 R.D. & LE ROUX, A. (2014) Chapter Six - Vocal Complexity in Meerkats and Other
14
15 1719 Mongoose Species. In *Advances in the Study of Behavior* (eds M. NAGUIB, L.
16
17 1720 BARRETT, H.J. BROCKMANN, S. HEALY, J.C. MITANI, T.J. ROPER & L.W. SIMMONS),
18
19 1721 pp. 281–310. Academic Press.
- 22
23 1722 MANTEUFFEL, G., PUPPE, B. & SCHÖN, P.C. (2004) Vocalization of farm animals as a
24
25 1723 measure of welfare. *Applied Animal Behaviour Science* **88**, 163–182.
- 28
29 1724 MAO, A., GIRAUDET, C.S.E., LIU, K., DE ALMEIDA NOLASCO, I., XIE, Z., XIE, Z., GAO, Y.,
30
31 1725 THEOBALD, J., BHATTA, D., STEWART, R. & MCELLIGOTT, A.G. (2022) Automated
32
33 1726 identification of chicken distress vocalizations using deep learning models. *Journal of*
34
35 1727 *The Royal Society Interface* **19**, 20210921.
- 38
39 1728 MARIN-CUDRAZ, T., MUFFAT-JOLY, B., NOVOA, C., AUBRY, P., DESMET, J.-F., MAHAMOUD-
40
41 1729 ISSA, M., NICOLÈ, F., VAN NIEKERK, M.H., MATHEVON, N. & SÈBE, F. (2019)
42
43 1730 Acoustic monitoring of rock ptarmigan: A multi-year comparison with point-count
44
45 1731 protocol. *Ecological Indicators* **101**, 710–719.
- 48
49 1732 MARQUES, T.A., THOMAS, L., MARTIN, S.W., MELLINGER, D.K., WARD, J.A., MORETTI, D.J.,
50
51 1733 HARRIS, D. & TYACK, P.L. (2013a) Estimating animal population density using
52
53 1734 passive acoustics. *Biological reviews of the Cambridge Philosophical Society* **88**,
54
55 1735 287–309.

- 1
2
3 1736 MARQUES, T.A., THOMAS, L., MARTIN, S.W., MELLINGER, D.K., WARD, J.A., MORETTI, D.J.,
4
5 1737 HARRIS, D. & TYACK, P.L. (2013b) Estimating animal population density using
6
7 1738 passive acoustics. *Biological Reviews* **88**, 287–309.
8
9
10
11 1739 MARTIN, K., ADAM, O., OBIN, N. & DUFOUR, V. (2022) Rookognise: Acoustic detection and
12
13 1740 identification of individual rooks in field recordings using multi-task neural networks.
14
15 1741 *Ecological Informatics* **72**, 101818.
16
17
18
19 1742 MCCOMB, K. & SEMPLE, S. (2005) Coevolution of vocal communication and sociality in
20
21 1743 primates. *Biology Letters* **1**, 381–385. Royal Society.
22
23
24 1744 MCDONALD, M., HILDEBRAND, J. & MESNICK, S. (2009) Worldwide decline in tonal
25
26 1745 frequencies of blue whale songs. *Endangered Species Research* **9**, 13–21.
27
28
29
30 1746 MCDONALD, M.A. & FOX, C.G. (1999) Passive acoustic methods applied to fin whale
31
32 1747 population density estimation. *Journal of the Acoustical Society of America* **105**,
33
34 1748 2643–2651. ACOUSTICAL SOC AMER AMER INST PHYSICS.
35
36
37
38 1749 MCINNIS, L., HEALY, J. & MELVILLE, J. (2020) UMAP: Uniform Manifold Approximation
39
40 1750 and Projection for Dimension Reduction. arXiv. <http://arxiv.org/abs/1802.03426>
41
42 1751 [accessed 7 August 2023].
43
44
45
46 1752 MCLOUGHLIN, M.P., STEWART, R. & MCELLIGOTT, A.G. (2019) Automated bioacoustics:
47
48 1753 methods in ecology and conservation and their potential for animal welfare
49
50 1754 monitoring. *Journal of The Royal Society Interface* **16**, 20190225.
51
52
53
54 1755 MENNILL, D.J., BATTISTON, M., WILSON, D.R., FOOTE, J.R. & DOUCET, S.M. (2012) Field test
55
56 1756 of an affordable, portable, wireless microphone array for spatial monitoring of animal
57
58 1757 ecology and behaviour. *Methods in Ecology and Evolution* **3**, 704–712.
59
60

- 1
2
3 1758 METCALF, O., ABRAHAMS, C., ASHINGTON, B., BAKER, E., BRADFER-LAWRENCE, T.,
4
5 1759 BROWNING, E., CARRUTHERS-JONES, J., DARBY, J., DICK, J., ELDRIDGE, A., ELLIOTT,
6
7 1760 D., HEATH, B., HOWDEN-LEACH, P., JOHNSTON, A., LEES, A., ET AL. (2023a) Good
8
9 practice guidelines for long-term ecoacoustic monitoring in the UK. In pp. 1–82.
10 1761
11 Report, The UK Acoustics Network.
12 1762
13
14
15 1763 METCALF, O., ABRAHAMS, C., ASHINGTON, B., BAKER, E., BRADFER-LAWRENCE, T.,
16
17 1764 BROWNING, E., CARRUTHERS-JONES, J., DARBY, J., DICK, J., ELDRIDGE, A., & OTHERS
18
19 (2023b) Good practice guidelines for long-term ecoacoustic monitoring in the UK.
20 1765
21 The UK Acoustics Network.
22 1766
23
24
25 1767 MEYER, D., HODGES, J.K., RINALDI, D., WIJAYA, A., ROOS, C. & HAMMERSCHMIDT, K.
26
27 (2012) Acoustic structure of male loud-calls support molecular phylogeny of
28 1768
29 Sumatran and Javanese leaf monkeys (genus *Presbytis*). *Bmc Evolutionary Biology*
30 1769
31 **12**, 16–16. BIOMED CENTRAL LTD.
32 1770
33
34
35 1771 MILLER, B.S., MADHUSUDHANA, S., AULICH, M.G. & KELLY, N. (2023) Deep learning
36
37 algorithm outperforms experienced human observer at detection of blue whale D-
38 1772
39 calls: a double-observer analysis. *Remote Sensing in Ecology and Conservation* **9**,
40 1773
41 104–116.
42 1774
43
44
45 1775 MINISTRY OF LOCAL GOVERNMENT AND MODERNISATION (2021) Norwegian data centres -
46
47 sustainable, digital powerhouses. Plan, regjeringen.no. *Government.no*.
48 1776
49 [https://www.regjeringen.no/en/dokumenter/norwegian-data-centres-sustainable-](https://www.regjeringen.no/en/dokumenter/norwegian-data-centres-sustainable-digital-powerhouses/id2867155/)
50 1777
51 [digital-powerhouses/id2867155/](https://www.regjeringen.no/en/dokumenter/norwegian-data-centres-sustainable-digital-powerhouses/id2867155/) [accessed 7 August 2023].
52 1778
53
54
55
56
57
58
59
60

- 1
2
3 1779 MITROVIC, D., ZEPPELZAUER, M. & BREITENEDER, C. (2006) Discrimination and retrieval of
4
5 1780 animal sounds. In *2006 12th International Multi-Media Modelling Conference* p. 5
6
7 1781 pp.-.
- 8
9
10
11 1782 MONCZAK, A., BERRY, A., KEHRER, C. & MONTIE, E. (2017) Long-term acoustic monitoring
12
13 1783 of fish calling provides baseline estimates of reproductive timelines in the May River
14
15 1784 estuary, southeastern USA. *Marine Ecology Progress Series* **581**, 1–19.
- 16
17
18
19 1785 MORRISON, A. & NOVIKOVA, A. (2023) Monitoring technologies for animal welfare: A
20
21 1786 review of aspirations and deployments in zoos. In *Proceedings of the Future*
22
23 1787 *Technologies Conference (FTC) 2022, Volume 3* (ed K. ARAI), pp. 155–178. Springer
24
25 1788 International Publishing, Cham.
- 26
27
28
29 1789 MORRISSEY, R.P., WARD, J., DIMARZIO, N., JARVIS, S. & MORETTI, D.J. (2006) Passive
30
31 1790 acoustic detection and localization of sperm whales (*Physeter macrocephalus*) in the
32
33 1791 tongue of the ocean. *Applied Acoustics* **67**, 1091–1105.
- 34
35
36
37 1792 NARASIMHAN, R., FERN, X.Z. & RAICH, R. (2017) Simultaneous segmentation and
38
39 1793 classification of bird song using CNN. In *2017 IEEE International Conference on*
40
41 1794 *Acoustics, Speech and Signal Processing (ICASSP)* pp. 146–150.
- 42
43
44
45 1795 NELSON, D.A. (2000) Song overproduction, selective attrition and song dialects in the white-
46
47 1796 crowned sparrow. *Animal Behaviour* **60**, 887–898.
- 48
49
50
51 1797 NGUYEN HONG DUC, P., TORTEROTOT, M., SAMARAN, F., WHITE, P.R., GÉRARD, O., ADAM,
52
53 1798 O. & CAZAU, D. (2021) Assessing inter-annotator agreement from collaborative
54
55 1799 annotation campaign in marine bioacoustics. *Ecological Informatics* **61**, 101185.
- 56
57
58
59
60

- 1
2
3 1800 NICHOLSON, D. (2023) Crowsetta: A Python tool to work with any format for annotating
4
5 1801 animal vocalizations and bioacoustics data. *Journal of Open Source Software* **8**, 5338.
6
7
8
9 1802 NIJMAN, V. (2007) Effects of vocal behaviour on abundance estimates of rainforest
10
11 1803 Galliforms. *Acta Ornithologica* **42**, 186–190.
12
13
14 1804 NOLASCO, I., SINGH, S., MORFI, V., LOSTANLEN, V., STRANDBURG-PESHKIN, A., VIDAÑA-
15
16 1805 VILA, E., GILL, L., PAMUŁA, H., WHITEHEAD, H., KISKIN, I., JENSEN, F.H., MORFORD,
17
18 1806 J., EMMERSON, M.G., VERSACE, E., GROUT, E., ET AL. (2023) Learning to detect an
19
20 1807 animal sound from five examples. arXiv. <http://arxiv.org/abs/2305.13210> [accessed 5
21
22 1808 July 2023].
23
24
25
26
27 1809 OBRIST, M.K., PAVAN, G., SUEUR, J., RIEDE, K., LLUSIA, D. & MÁRQUEZ, R. (2010)
28
29 1810 Bioacoustics approaches in biodiversity inventories. *Manual on field recording*
30
31 1811 *techniques and protocols for all taxa biodiversity inventories* **8**, 68–99.
32
33
34
35 1812 ODOM, K.J., ARAYA-SALAS, M., MORANO, J.L., LIGON, R.A., LEIGHTON, G.M., TAFF, C.C.,
36
37 1813 DALZIELL, A.H., BILLINGS, A.C., GERMAIN, R.R., PARDO, M., DE ANDRADE, L.G.,
38
39 1814 HEDWIG, D., KEEN, S.C., SHIU, Y., CHARIF, R.A., ET AL. (2021) Comparative
40
41 1815 bioacoustics: a roadmap for quantifying and comparing animal sounds across diverse
42
43 1816 taxa. *Biological Reviews* **96**, 1135–1159.
44
45
46
47 1817 OGUTU, J.O., PIEPHO, H.-P., SAID, M.Y., OJWANG, G.O., NJINO, L.W., KIFUGO, S.C. &
48
49 1818 WARGUTE, P.W. (2016) Extreme Wildlife Declines and Concurrent Increase in
50
51 1819 Livestock Numbers in Kenya: What Are the Causes? *PLOS ONE* **11**, e0163249.
52
53 1820 Public Library of Science.
54
55
56
57 1821 OLIVER, R.Y., ELLIS, D.P.W., CHMURA, H.E., KRAUSE, J.S., PÉREZ, J.H., SWEET, S.K.,
58
59 1822 GOUGH, L., WINGFIELD, J.C. & BOELMAN, N.T. (2018) Eavesdropping on the Arctic:

- 1
2
3 1823 Automated bioacoustics reveal dynamics in songbird breeding phenology. *Science*
4
5 1824 *Advances* **4**, eaaq1084.
6
7
8
9 1825 OSWALD, J.N., ERBE, C., GANNON, W.L., MADHUSUDHANA, S. & THOMAS, J.A. (2022)
10
11 1826 Detection and Classification Methods for Animal Sounds. In *Exploring Animal*
12
13 1827 *Behavior Through Sound: Volume 1: Methods* (eds C. ERBE & J.A. THOMAS), pp.
14
15 1828 269–317. Springer International Publishing, Cham.
- 16
17
18
19 1829 PARRILLA, A.G.A. & STOWELL, D. (2022) Polyphonic sound event detection for highly dense
20
21 1830 birdsong scenes. arXiv. <http://arxiv.org/abs/2207.06349> [accessed 6 July 2023].
22
23
24 1831 PÉREZ-GRANADOS, C. & SCHUCHMANN, K.-L. (2021) Passive acoustic monitoring of the diel
25
26 1832 and annual vocal behavior of the Black and Gold Howler Monkey. *American Journal*
27
28 1833 *of Primatology* **83**, e23241.
- 29
30
31
32 1834 PETSO, T., JAMISOLA, R.S. & MPOELENG, D. (2021) Review on methods used for wildlife
33
34 1835 species and individual identification. *European Journal of Wildlife Research* **68**, 3.
- 35
36
37
38 1836 POLITIS, A., MESAROS, A., ADAVANNE, S., HEITTOLA, T. & VIRTANEN, T. (2020) Overview
39
40 1837 and Evaluation of Sound Event Localization and Detection in DCASE 2019.
41
42 1838 *arXiv:2009.02792 [cs, eess]*.
- 43
44
45
46 1839 POWELL, R. (2000) Animal home ranges and territories and home range estimators. *Research*
47
48 1840 *Techniques in Animal Ecology: Controversies and Consequences*, 65–110.
- 49
50
51 1841 QIAN, T., DENG, G., LI, Y. & YANG, D. (2023) Description of the advertisement call of
52
53 1842 *Boulenophrys nanlingensis* (Anura, Megophryidae), with a case of individual
54
55 1843 identification using its dorsum pattern. *Herpetozoa* **36**, 123–128. Pensoft Publishers.
- 56
57
58
59
60

- 1
2
3 1844 RICE, A., ŠIROVIĆ, A., TRICKEY, J.S., DEBICH, A.J., GOTTLIEB, R.S., WIGGINS, S.M.,
4
5 1845 HILDEBRAND, J.A. & BAUMANN-PICKERING, S. (2021) Cetacean occurrence in the
6
7 1846 Gulf of Alaska from long-term passive acoustic monitoring. *Marine Biology* **168**, 72.
8
9
10
11 1847 RIESCH, R., BARRETT-LENNARD, L.G., ELLIS, G.M., FORD, J.K.B. & DEECKE, V.B. (2012)
12
13 1848 Cultural traditions and the evolution of reproductive isolation: ecological speciation in
14
15 1849 killer whales? *Biological Journal of the Linnean Society* **106**, 1–17. WILEY-
16
17 1850 BLACKWELL.
18
19
20
21 1851 ROMERO-MUJALLI, D., BERGMANN, T., ZIMMERMANN, A. & SCHEUMANN, M. (2021)
22
23 1852 Utilizing DeepSqueak for automatic detection and classification of mammalian
24
25 1853 vocalizations: a case study on primate vocalizations. *Scientific Reports* **11**, 24463.
26
27 1854 Nature Publishing Group.
28
29
30
31 1855 ROSS, S.R.P. -J., O'CONNELL, D.P., DEICHMANN, J.L., DESJONQUÈRES, C., GASC, A.,
32
33 1856 PHILLIPS, J.N., SETHI, S.S., WOOD, C.M. & BURIVALOVA, Z. (2023) Passive acoustic
34
35 1857 monitoring provides a fresh perspective on fundamental ecological questions.
36
37 1858 *Functional Ecology* **37**, 959–975.
38
39
40
41 1859 ROSS, S.R.P.-J., FRIEDMAN, N.R., DUDLEY, K.L., YOSHIMURA, M., YOSHIDA, T. & ECONOMO,
42
43 1860 E.P. (2018) Listening to ecosystems: data-rich acoustic monitoring through
44
45 1861 landscape-scale sensor networks. *Ecological Research* **33**, 135–147.
46
47
48
49 1862 ROSTRO-GARCÍA, S., KAMLER, J.F., SOLLMANN, R., BALME, G., AUGUSTINE, B.C., KÉRY, M.,
50
51 1863 CROUTHERS, R., GRAY, T.N.E., GROENENBERG, M., PRUM, S. & MACDONALD, D.W.
52
53 1864 (2023) Population dynamics of the last leopard population of eastern Indochina in the
54
55 1865 context of improved law enforcement. *Biological Conservation* **283**, 110080.
56
57
58
59
60

- 1
2
3 1866 ROTHSTEIN, S.I. & FLEISCHER, R.C. (1987) Vocal Dialects and Their Possible Relation to
4
5 1867 Honest Status Signalling in the Brown-Headed Cowbird. *The Condor* **89**, 1–23.
6
7
8
9 1868 SAEED, A., GRANGIER, D. & ZEGHIDOUR, N. (2021) Contrastive Learning of General-Purpose
10
11 1869 Audio Representations. In *ICASSP 2021 - 2021 IEEE International Conference on*
12
13 1870 *Acoustics, Speech and Signal Processing (ICASSP)* pp. 3875–3879.
14
15
16
17 1871 SAINBURG, T. (2019) timsainb/noisereducer: v1.0. Zenodo. <https://zenodo.org/record/3243139>
18
19 1872 [accessed 5 July 2023].
20
21
22 1873 SAINBURG, T., THIELK, M. & GENTNER, T.Q. (2020) Latent space visualization,
23
24 1874 characterization, and generation of diverse vocal communication signals. bioRxiv.
25
26 1875 <https://www.biorxiv.org/content/10.1101/870311v2> [accessed 6 July 2023].
27
28
29
30 1876 SARKAR, E. & -DOSS, M.M. (2023) Can Self-Supervised Neural Representations Pre-Trained
31
32 1877 on Human Speech distinguish Animal Callers? arXiv. <http://arxiv.org/abs/2305.14035>
33
34 1878 [accessed 6 July 2023].
35
36
37
38 1879 SATTAR, F. (2023) A new acoustical autonomous method for identifying endangered whale
39
40 1880 calls: A case study of blue whale and fin whale. *Sensors* **23**, 3048.
41
42
43
44 1881 SCHULTZE, S., GRUENEFELD, U. & BOLL, S. (2020) Demystifying deep learning: Developing
45
46 1882 a learning app for beginners to gain practical experience. *Proceedings of the Mensch*
47
48 1883 *und Computer 2020 Workshop*.
49
50
51
52 1884 SEBASTIÁN-GONZÁLEZ, E., PANG-CHING, J., BARBOSA, J.M. & HART, P. (2015) Bioacoustics
53
54 1885 for species management: two case studies with a Hawaiian forest bird. *Ecology and*
55
56 1886 *Evolution* **5**, 4696–4705.
57
58
59
60

- 1
2
3 1887 SETHI, S.S., JONES, N.S., FULCHER, B.D., PICINALI, L., CLINK, D.J., KLINCK, H., ORME,
4
5 1888 C.D.L., WREGE, P.H. & EWERS, R.M. (2020) Characterizing soundscapes across
6
7 1889 diverse ecosystems using a universal acoustic feature set. *Proceedings of the National*
8
9 1890 *Academy of Sciences* **117**, 17049–17055. Proceedings of the National Academy of
10
11 1891 Sciences.
- 12
13
14
15 1892 SHARPE, F., BOLTON, M., SHELDON, R. & RATCLIFFE, N. (2009) Effects of color banding,
16
17 1893 radio tagging, and repeated handling on the condition and survival of Lapwing chicks
18
19 1894 and consequences for estimates of breeding productivity. *Journal of Field*
20
21 1895 *Ornithology* **80**, 101–110.
- 22
23
24
25 1896 SHIU, Y., PALMER, K.J., ROCH, M.A., FLEISHMAN, E., LIU, X., NOSAL, E.-M., HELBLE, T.,
26
27 1897 CHOLEWIAK, D., GILLESPIE, D. & KLINCK, H. (2020) Deep neural networks for
28
29 1898 automated detection of marine mammal species. *Scientific Reports* **10**, 607. Nature
30
31 1899 Publishing Group.
- 32
33
34
35 1900 SLOCOMBE, K.E., WALLER, B.M. & LIEBAL, K. (2011) The language void: the need for
36
37 1901 multimodality in primate communication research. *Animal Behaviour* **81**, 919–924.
- 38
39
40
41 1902 SMITH, B.R., ROOT-GUTTERIDGE, H., BUTKIEWICZ, H., DASSOW, A., FONTAINE, A.C.,
42
43 1903 MARKHAM, A., OWENS, J., SCHINDLER, L., WIJERS, M., KERSHENBAUM, A., SMITH,
44
45 1904 B.R., ROOT-GUTTERIDGE, H., BUTKIEWICZ, H., DASSOW, A., FONTAINE, A.C., ET AL.
46
47 1905 (2021) Acoustic localisation of wildlife with low-cost equipment: lower sensitivity,
48
49 1906 but no loss of precision. *Wildlife Research* **49**, 372–381. CSIRO PUBLISHING.
- 50
51
52
53 1907 SOULSBURY, C.D., GRAY, H.E., SMITH, L.M., BRAITHWAITE, V., COTTER, S.C., ELWOOD,
54
55 1908 R.W., WILKINSON, A. & COLLINS, L.M. (2020) The welfare and ethics of research
56
57 1909 involving wild animals: A primer. *Methods in Ecology and Evolution* **11**, 1164–1181.
58
59
60

- 1
2
3 1910 SPILLMANN, B., VAN SCHAİK, C.P., SETIA, T.M. & SADJADI, S.O. (2017) Who shall I say is
4
5 1911 calling? Validation of a caller recognition procedure in Bornean flanged male
6
7 1912 orangutan (*Pongo pygmaeus wurmbii*) long calls. *Bioacoustics* **26**, 109–120.
8
9
10
11 1913 STAATERMAN, E., OGBURN, M.B., ALTIERI, A.H., BRANDL, S.J., WHIPPO, R., SEEMANN, J.,
12
13 1914 GOODISON, M. & DUFFY, J.E. (2017) Bioacoustic measurements complement visual
14
15 1915 biodiversity surveys: preliminary evidence from four shallow marine habitats. *Marine*
16
17 1916 *Ecology Progress Series* **575**, 207–215.
18
19
20
21 1917 STOWELL, D. (2022a) Computational bioacoustics with deep learning: a review and roadmap.
22
23 1918 *PeerJ* **10**, e13152. PeerJ Inc.
24
25
26
27 1919 STOWELL, D. (2022b) Computational bioacoustics with deep learning: a review and roadmap.
28
29 1920 *PeerJ* **10**, e13152. PeerJ Inc.
30
31
32 1921 STOWELL, D., WOOD, M.D., PAMUŁA, H., STYLIANOU, Y. & GLOTIN, H. (2019) Automatic
33
34 1922 acoustic detection of birds through deep learning: The first Bird Audio Detection
35
36 1923 challenge. *Methods in Ecology and Evolution* **10**, 368–380.
37
38
39
40 1924 SUEUR, J., FARINA, A., GASC, A., PIERETTI, N. & PAVOINE, S. (2014) Acoustic Indices for
41
42 1925 Biodiversity Assessment and Landscape Investigation. *Acta Acustica united with*
43
44 1926 *Acustica* **100**, 772–781.
45
46
47
48 1927 SUGAI, L.S.M., DESJONQUÈRES, C., SILVA, T.S.F. & LLUSIA, D. (2020) A roadmap for survey
49
50 1928 designs in terrestrial acoustic monitoring. *Remote Sensing in Ecology and*
51
52 1929 *Conservation* **6**, 220–235.
53
54
55
56 1930 SUGAI, L.S.M. & LLUSIA, D. (2019) Bioacoustic time capsules: Using acoustic monitoring to
57
58 1931 document biodiversity. *Ecological Indicators* **99**, 149–152.
59
60

- 1
2
3 1932 TEIXEIRA, D., MARON, M. & RENSBURG, B.J. (2019) Bioacoustic monitoring of animal vocal
4
5 1933 behavior for conservation. *Conservation Science and Practice* **1**.
6
7
8
9 1934 TORTI, V., VALENTE, D., DE GREGORIO, C., COMAZZI, C., MIARETSSOA, L., RATSIMBAZAFY, J.,
10
11 1935 GIACOMA, C. & GAMBA, M. (2018) Call and be counted! Can we reliably estimate the
12
13 1936 number of callers in the indri's (Indri indri) song? *PLOS ONE* **13**, e0201664.
14
15
16 1937 TOWSEY, M., PLANITZ, B., NANTES, A., WIMMER, J. & ROE, P. (2012) A toolbox for animal
17
18 1938 call recognition. *Bioacoustics* **21**, 107–125. Taylor & Francis.
19
20
21
22 1939 TURIAN, J., SHIER, J., KHAN, H.R., RAJ, B., SCHULLER, B.W., STEINMETZ, C.J., MALLOY, C.,
23
24 1940 TZANETAKIS, G., VELARDE, G., McNALLY, K., HENRY, M., PINTO, N., NOUFI, C.,
25
26 1941 CLOUGH, C., HERREMANS, D., ET AL. (2022) HEAR 2021: Holistic Evaluation of
27
28 1942 Audio Representations. In *Proceedings of the NeurIPS 2021 Competitions and*
29
30 1943 *Demonstrations Track* (eds D. KIELA, M. CICCONE & B. CAPUTO), pp. 125–145.
31
32 1944 PMLR.
33
34
35
36
37 1945 UETZ, G.W. & ROBERTS, J.A. (2002) Multisensory Cues and Multimodal Communication in
38
39 1946 Spiders: Insights from Video/Audio Playback Studies. *Brain Behavior and Evolution*
40
41 1947 **59**, 222–230.
42
43
44
45 1948 VAN SEGBROECK, M., KNOLL, A.T., LEVITT, P. & NARAYANAN, S. (2017) MUPET—Mouse
46
47 1949 Ultrasonic Profile ExTraction: A Signal Processing Tool for Rapid and Unsupervised
48
49 1950 Analysis of Ultrasonic Vocalizations. *Neuron* **94**, 465–485.e5.
50
51
52
53 1951 VELÁSQUEZ, N.A., MARAMBIO, J., BRUNETTI, E., MÉNDEZ, M.A., VÁSQUEZ, R.A. & PENNA,
54
55 1952 M. (2013) Bioacoustic and genetic divergence in a frog with a wide geographical
56
57 1953 distribution. *Biological Journal of the Linnean Society* **110**, 142–155.
58
59
60

- 1
2
3 1954 VENKATESH, S., MOFFAT, D. & MIRANDA, E.R. (2022) You Only Hear Once: A YOLO-like
4
5 1955 Algorithm for Audio Segmentation and Sound Event Detection. *Applied Sciences* **12**,
6
7 1956 3293. Multidisciplinary Digital Publishing Institute.
- 8
9
10
11 1957 VICKERS, W., MILNER, B., RISCH, D. & LEE, R. (2021) Robust North Atlantic right whale
12
13 1958 detection using deep learning models for denoising). *The Journal of the Acoustical*
14
15 1959 *Society of America* **149**, 3797–3812.
- 16
17
18
19 1960 VOLODINA, E.V. & VOLODIN, I.A. (1999) Bioacoustics in zoos: A review of applications and
20
21 1961 perspectives. *International Zoo News* **46**, 208–213.
- 22
23
24 1962 VU, T.T., CHI, T.N., DOHERTY JR, P.F., NGUYEN, H.T., CLINK, D.J., DAC, M.N., THANH, H.D.
25
26 & TRONG, T.G. (2023) Using mobile smartphones and bioacoustics to monitor
27 1963 endangered bird species. *Ibis*. Wiley Online Library.
- 28
29 1964
30
31
32 1965 VU, T.T. & TRAN, L.M. (2019) An Application of Autonomous Recorders for Gibbon
33
34 1966 Monitoring. *International Journal of Primatology* **40**, 169–186. Springer US.
- 35
36
37
38 1967 WANG, Z., SHE, Q. & WARD, T.E. (2022) Generative Adversarial Networks in Computer
39
40 1968 Vision: A Survey and Taxonomy. *ACM Computing Surveys* **54**, 1–38.
- 41
42
43 1969 WIJERS, M., LOVERIDGE, A., MACDONALD, D.W. & MARKHAM, A. (2021) CARACAL: a
44
45 1970 versatile passive acoustic monitoring tool for wildlife research and conservation.
46
47 1971 *Bioacoustics* **30**, 41–57. Taylor & Francis.
- 48
49
50
51 1972 WILKINSON, M.D., DUMONTIER, M., AALBERSBERG, I.J.J., APPLETON, G., AXTON, M., BAAK,
52
53 1973 A., BLOMBERG, N., BOITEN, J.-W., DA SILVA SANTOS, L.B., BOURNE, P.E., BOUWMAN,
54
55 1974 J., BROOKES, A.J., CLARK, T., CROSAS, M., DILLO, I., ET AL. (2016) The FAIR Guiding
56
57
58
59
60

- 1
2
3 1975 Principles for scientific data management and stewardship. *Scientific Data* **3**, 160018.
4
5 1976 Nature Publishing Group.
6
7
8
9 1977 WILLACY, R.J., MAHONY, M. & NEWELL, D.A. (2015) If a frog calls in the forest: Bioacoustic
10
11 1978 monitoring reveals the breeding phenology of the endangered Richmond Range
12
13 1979 mountain frog (*Philoria richmondensis*). *Austral Ecology* **40**, 625–633.
14
15
16 1980 WISDOM, S., TZINIS, E., ERDOGAN, H., WEISS, R.J., WILSON, K. & HERSHEY, J.R. (2020)
17
18 1981 Unsupervised Sound Separation Using Mixture Invariant Training. arXiv.
19
20 1982 <http://arxiv.org/abs/2006.12701> [accessed 6 July 2023].
21
22
23
24 1983 WOLF, C. & RIPPLE, W.J. (2016) Prey depletion as a threat to the world's large carnivores.
25
26 1984 *Royal Society Open Science* **3**, 160252. Royal Society.
27
28
29
30 1985 WOOD, C.M. & PEERY, M.Z. (2022) What does 'occupancy' mean in passive acoustic
31
32 1986 surveys? *Ibis* **164**, 1295–1300.
33
34
35 1987 WU, S.-H., CHANG, H.-W., LIN, R.-S. & TUANMU, M.-N. (2022) SILIC: A cross database
36
37 1988 framework for automatically extracting robust biodiversity information from
38
39 1989 soundscape recordings based on object detection and a tiny training dataset.
40
41 1990 *Ecological Informatics* **68**, 101534.
42
43
44
45 1991 XIE, J., COLONNA, J.G. & ZHANG, J. (2021) Bioacoustic signal denoising: a review. *Artificial*
46
47 1992 *Intelligence Review* **54**, 3575–3597. Springer.
48
49
50
51 1993 YANG, W., CHANG, W., SONG, Z., ZHANG, Y. & WANG, X. (2021) Transfer learning for
52
53 1994 denoising the echolocation clicks of finless porpoise (*Neophocaena phocaenoides*
54
55 1995 sunameri) using deep convolutional autoencoders. *The Journal of the Acoustical*
56
57 1996 *Society of America* **150**, 1243–1250.
58
59
60

- 1
2
3 1997 YIN, S., LIU, C., ZHANG, Z., LIN, Y., WANG, D., TEJEDOR, J., ZHENG, T.F. & LI, Y. (2015)
4
5 1998 Noisy training for deep neural networks in speech recognition. *EURASIP Journal on*
6
7 *Audio, Speech, and Music Processing* **2015**, 2.
8 1999
9
10
11 2000 YU, Y., SI, X., HU, C. & ZHANG, J. (2019) A Review of Recurrent Neural Networks: LSTM
12
13 2001 Cells and Network Architectures. *Neural Computation* **31**, 1235–1270.
14
15
16 2002 ZHANG, H., CISSE, M., DAUPHIN, Y.N. & LOPEZ-PAZ, D. (2018) mixup: Beyond Empirical
17
18 Risk Minimization. arXiv. <http://arxiv.org/abs/1710.09412> [accessed 6 July 2023].
19 2003
20
21
22 2004 ZHANG, Z., ZHANG, H., HE, Y. & LIU, T. (2022) A review in the automatic detection of pigs
23
24 2005 behavior with sensors. *Journal of Sensors* **2022**, 1–17.
25
26
27 2006 ZHUANG, F., QI, Z., DUAN, K., XI, D., ZHU, Y., ZHU, H., XIONG, H. & HE, Q. (2021) A
28
29 2007 Comprehensive Survey on Transfer Learning. *Proceedings of the IEEE* **109**, 43–76.
30
31
32
33 2008 ZIMMER, W.M.X. (2011) *Passive Acoustic Monitoring of Cetaceans*. Cambridge University
34
35 2009 Press.
36
37
38
39 2010 ZWERTS, J.A., STEPHENSON, P.J., MAISELS, F., ROWCLIFFE, M., ASTARAS, C., JANSEN, P.A.,
40
41 2011 VAN DER WAARDE, J., STERCK, L.E.H.M., VERWEIJ, P.A., BRUCE, T., BRITAIN, S. &
42
43 2012 VAN KUIJK, M. (2021) Methods for wildlife monitoring in tropical forests: Comparing
44
45 2013 human observations, camera traps, and passive acoustic sensors. *Conservation*
46
47 *Science and Practice* **3**, e568.
48 2014
49
50
51 2015
52
53
54 2016
55
56
57
58
59
60