1   **Using intrinsic and contextual information associated with automated signal detections**

2   **to improve call recogniser performance: a case study using the cryptic and critically**

3   **endangered Night Parrot (*Pezoporus occidentalis*)**

4

5   Running Title: POST-PROCESSING TO IMPROVE RECOGNISER PERFORMANCE

6

7   *Nicholas P. Leseberg[A,B], William N. Venables[C], Stephen A. Murphy,[A,B], James E.M. Watson[A,B,D]*

8   [A]School of Earth and Environmental Sciences, The University of Queensland, Brisbane 4072,

9   Queensland, Australia

10  [B]Green Fire Science, The University of Queensland, Brisbane 4072, Queensland, Australia

11  [C]School of Mathematics and Physics, The University of Queensland, Brisbane 4072,

12  Queensland, Australia

13  [D]Centre for Biodiversity and Conservation Science, School of Biological Sciences, The

14  University of Queensland, Brisbane 4072, Queensland, Australia

15

16  Corresponding author contact details:

17  Nick Leseberg

18  School of Earth and Environmental Sciences

19  University of Queensland

20  St Lucia QLD 4072

21  Australia

22  Email: n.leseberg@uq.edu.au

23  Phone: 0488 636 010

24

25

26 **Abstract**

27 1. Rapid expansion in the collection of large acoustic datasets to answer ecological questions

28 has generated a parallel requirement for techniques that streamline analysis of these datasets.

29 In many cases, automated signal recognition algorithms, often termed 'call recognisers', are

30 the only feasible option for doing this. To date, most research has focused on what types of

31 recognisers perform best, and how to train these recognisers to optimise performance.

32 2. We demonstrate that once recogniser construction is complete and the data processed, further

33 improvements are possible using intrinsic and contextual information associated with each

34 detection. We initially construct a call recogniser for the Night Parrot (*Pezoporus occidentalis*)

35 using the R package monitoR, and scan a test dataset. We then examine a number of intrinsic

36 variables associated with each detection generated by the recogniser, and several contextual

37 variables, associated with the species' environment and ecology to determine if they might help

38 predict whether a given detection is a true positive (target signal) or false positive (non-target

39 signal). We test several logistic regression models incorporating different combinations of

40 intrinsic and contextual variables, selecting the best-performing model for application. We train

41 the model, using it to calculate the probability each detection is a true or false positive.

42 3. Substituting this model-derived probability for raw recogniser score improved the

43 recogniser's performance, reducing the number of detections requiring proofing by 60% to

44 achieve recall of 90%, and by 76% to achieve recall of 75%.

45 4. This technique is applicable to any recogniser output, regardless of the underlying algorithm.

46 Application requires an understanding of how the recogniser algorithm determines matches,

47 and knowledge of a species' ecology and environment. Because advanced programming skills

48 and expertise are not required to apply this technique, it will be particularly relevant to field

49 ecologists for whom building and operating call recognisers is an element of their research

50 toolbox, but not necessarily a focus.

53

54   **Introduction**

55   The increasing availability of technology to collect and analyse acoustic data, particularly

56   affordable automated recording units (ARUs), has seen a rapid expansion in this field of

57   research and its applications for ecology and conservation (Shonfield & Bayne, 2017; Teixeira,

58   Maron, & van Rensburg, 2019). The popularity of ARUs is largely due to their efficiency.

59   Particularly for long-term deployments, it is much cheaper to purchase, deploy, and maintain

60   an ARU than a human observer (Digby, Towsey, Bell, & Teal, 2013; Williams, O'Donnell, &

61   Armstrong, 2018). Unlike human observers, ARUs can be left in the field unattended for

62   extended periods, limited only by the availability of power and memory. As solar panels and

63   large capacity memory cards are now also relatively cheap, maintaining permanent acoustic

64   recording stations at remote sites has become feasible.

65

66   The easy collection of copious data has advantages and disadvantages. Large acoustic datasets

67   may contain powerful data (Magurran et al., 2010), but extracting that data can be challenging.

68   There are several techniques available to efficiently analyse large acoustic datasets, the most

69   suitable contingent on the nature of the signal of interest (Joshi, Mulder, & Rowe, 2017;

70   Towsey et al., 2018). Increasingly, research has focused on techniques that automate the signal

71   extraction process. This is typically performed using a signal detection algorithm, hereafter

72   termed 'call recogniser' (Potamitis, Ntalampiras, Jahn, & Riede, 2014; Priyadarshani,

73   Marsland, & Castro, 2018). For infrequent signals within large datasets, a call recogniser may

74   be the only feasible solution.

75

76 There are several options for researchers wanting to construct a call recogniser. They vary in

77 complexity, from commercial off-the-shelf programs such as Kaleidoscope (Wildlife Acoustics

78 Inc., Concord, Massachusetts, USA), to more recently, advanced machine learning algorithms

79 (Koops, van Balen, & Wiering, 2014; Salamon & Bello, 2017), acoustic indices (Towsey,

80 Wimmer, Williamson, & Roe, 2014), and wavelet based approaches (Priyadarshani, Marsland,

81 Juodakis, Castro, & Listanti, 2020). Although the computational processes behind each differ,

82 the basic premise remains the same; a computer is trained to detect and evaluate acoustic

83 signals by comparing them to a known target signal. Potential signals are classified depending

84 on their similarity to the target signal, with the user controlling the threshold at which a match

85 is declared.

86

87 Understanding the impact of this threshold is critical to understanding the performance of a

88 call recogniser. Setting a high threshold increases the precision of the recogniser, meaning a

89 higher proportion of matches will represent actual detections, or true positives. However, this

90 increases the likelihood of false negatives; target signals that do not meet the threshold, for

91 example soft or distant calls. This reduces the recogniser's recall, or ability to identify all target

92 signals within a dataset. Conversely, reducing the threshold ensures that more lower-scoring

93 target signals are returned as matches, but simultaneously returns more lower-scoring non-

94 target signals, or false positives. This increases the recogniser's recall, but also increases the

95 proportion of non-target signals in the resulting dataset, thereby decreasing precision. This false

96 positive / false negative trade-off is a well-known classification problem, with threshold choice

97 driven by the relative cost of false positive or false negative errors.

98

99 Besides an obvious focus on which computational techniques create the most successful

100 recognisers, research has also focused on the properties of training data that achieve the best

101 results (Knight & Bayne, 2018; Priyadarshani et al., 2018). Because a call recogniser's output

102 is dependent on how closely the signal of interest compares to the training data, efforts to

103 improve a specific type of recogniser's performance have largely focused on this aspect of their

104 development. However, little research has focused on how post-processing could be used to

105 derive improvements in overall performance. Typically, the output of a recogniser is a list of

106 potential 'detections', each with associated intrinsic information derived from the call

107 recognition process, for example a 'score' reflecting how similar the detection is to the training

108 data. There is also any number of contextual variables associated with each detection, such as

109 time-of-day and geographic location, that are known to affect detectability (Horton, Stepanian,

110 Wainwright, & Tegeler, 2015). Patterns in both intrinsic and contextual data could provide

111 clues to predict whether a detection is actually a signal of interest.

112

113 In this paper we outline a novel method to develop a model that uses both intrinsic and

114 contextual information associated with a call recogniser's raw output to generate an improved

115 output. We intentionally present a detailed description of the process, because one of our aims

116 is to demystify the process of automated call recognition for field ecologists, thereby

117 encouraging them to perform their own analyses. Broadly, our process was to first construct a

118 call recogniser for the Night Parrot (*Pezoporus occidentalis*), then investigate relationships

119 between the intrinsic and contextual variables associated with the recogniser's output to

120 establish if any could be incorporated into a model that predicts whether a detection is a true

121 positive or false positive. Following a model development and selection process, we selected

122 the best-performing model and tested whether this model improved recogniser performance.

123

124 **Methods and Results**

125 *Study species and data collection*

126    The Night Parrot is a cryptic and extremely rare bird that formerly occurred throughout arid

127    central Australia (Higgins, 1999), but is now known from only a handful of sites. The species

128    is relatively sedentary, and predictably vocal (Leseberg et al., 2019; Murphy, Silcock, Murphy,

129    Reid, & Austin, 2017). They spend the day roosting in low, dense vegetation, as pairs or small

130    groups. The birds emerge at dusk to engage in a brief period of calling before leaving their

131    roost sites to feed. Birds occasionally return to their roost sites and call during the night, but

132    typically return for another brief period of calling just before dawn. Night Parrot vocalisations

133    are now relatively well known (Leseberg et al., 2019). Given this predictable calling behaviour,

134    acoustic monitoring has proven the most efficient technique for both monitoring the species at

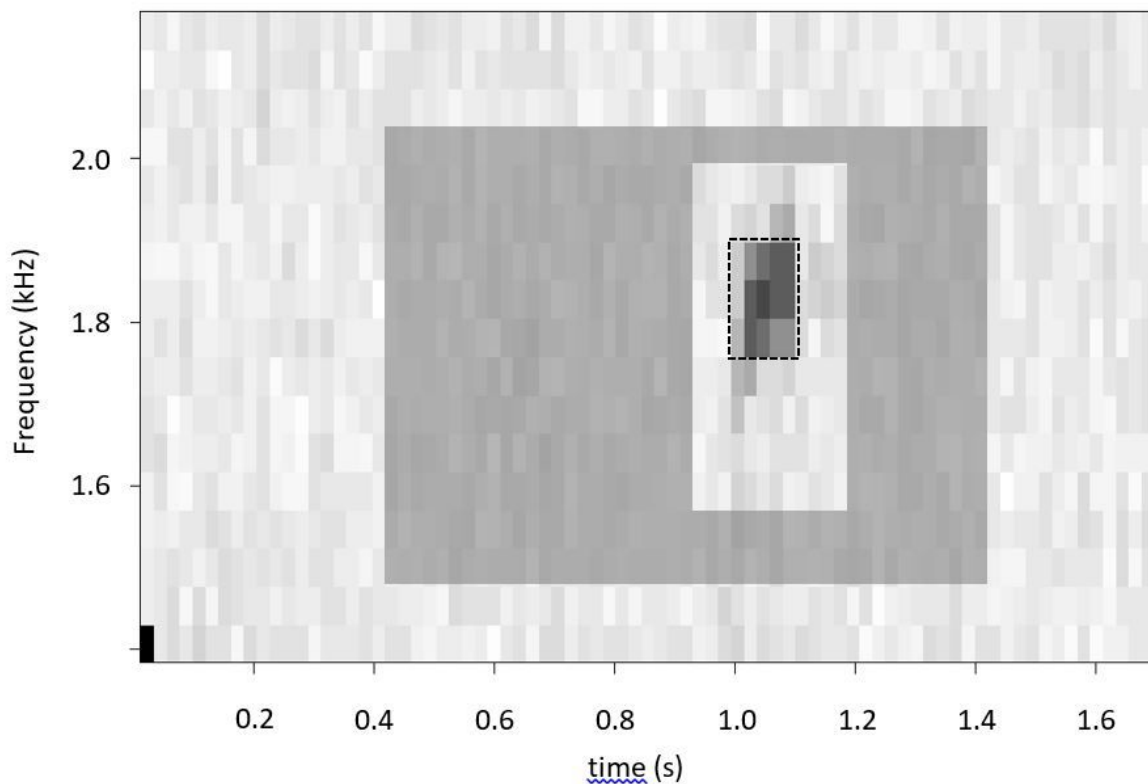135    known locations, and detecting it at new locations.

136

137    Since 2016, Night Parrot calling activity at three long-term stable roost sites in western

138    Queensland has been monitored using Song Meter 3 and Song Meter 4 ARUs (Wildlife

139    Acoustics Inc., Concord, Massachusetts, USA), fitted with standard external omnidirectional

140    microphones. ARUs recorded from sunset to sunrise, using the ARU's default gain settings.

141    Most ARUs recorded at sampling rates of 24000 Hz, or 48000 Hz, although some recorded at

142    16000 Hz. As required under the Nyquist-Shannon Sampling Theorem (Landau, 1967), these

143    sampling rates are greater than twice the peak frequency of all Night Parrot calls of interest to

144    this study.

145

146    *Call recogniser development and sound file analysis*

147    We used the R package monitoR (Katz, Hafner, & Donovan, 2016; R Core Team, 2018) to

148    build a call recogniser for the Night Parrot. R is a programming language accessible to users

149    without specialist programming skills, and in a comparison with recognisers using machine

150    learning methods and commercially available packages, monitoR performed well (Knight et

151 al., 2017). We used the technique outlined in Katz et al. (2016) to construct a series of binary

152 point templates. Templates are created by clipping an example call from a sound file and

153 creating a spectrogram (FFT transformation = Hann window, FFT size = 512, overlap = 0). A

154 selection of cells of the resulting spectrogram are then classified as 'on' or 'off'. 'On' cells are

155 selected to represent the expected region of strongest signal for the call, while 'off' cells are

156 placed strategically where no or little signal is expected (Fig. 1).



157

158 Figure 1. An example of a binary point matching template for the Night Parrot 'toot' call,

159 overlaid on the spectrogram of a 'toot' call. The central box with dotted outline represents the

160 'on' cells, and ideally contains most of the expected call energy. The shaded area represents

161 the 'off' cells.

162

163 Although Night Parrots have a variety of different calls, we focused on the bell-like and whistle

164 calls, as these are the calls most likely to be heard in and around roost sites (Leseberg et al.,

165 2019). These broad call types can be broken down further, and we constructed at least one

166 template for each of the ten specific call types known from the study area. We used example

167 calls extracted from the long-term monitoring dataset, adding further templates until testing

168 suggested the recogniser could detect most local variation within these call types. The final

169 recogniser used 31 different templates. Because monitoR requires template files and the sound

170 files that will be scanned to have the same sample rate, these were downsampled or upsampled

171 if required to a sampling rate of 24000 Hz. Qualitative testing confirmed that manipulating the

172 files in this way had no apparent effect on results.

173

174 Before analysis, each sound file is converted to a spectrogram using the same parameters as

175 were used to create the templates. Each template is then stepped along that spectrogram, and

176 for every step a similarity score is assigned based on the difference between the amplitude

177 detected in the 'on' cells, and the amplitude detected in the 'off' cells of the template. When

178 plotted against time this results in a series of peaks; the recogniser returns a list of these peaks

179 with their associated score. As some signals within the sound file are detected by more than

180 one template, a buffer of two seconds was applied so only the highest scoring peak within any

181 two-second period was returned. Because Night Parrot calls are generally short, temporally

182 discrete events, the risk of missing calls due to applying this buffer was low.

183

184 *Recogniser performance assessment*

185 To evaluate recogniser performance, 90 ten-minute field recordings known to contain Night

186 Parrot calls were extracted from the long-term monitoring dataset. We used field recordings to

187 ensure measured performance reflected what could be achieved on actual field recordings

188 rather than a manufactured test dataset (Potamitis et al., 2014). We used recordings from nights

189 that were either calm or with light winds, as wind noise significantly reduces both ARU and
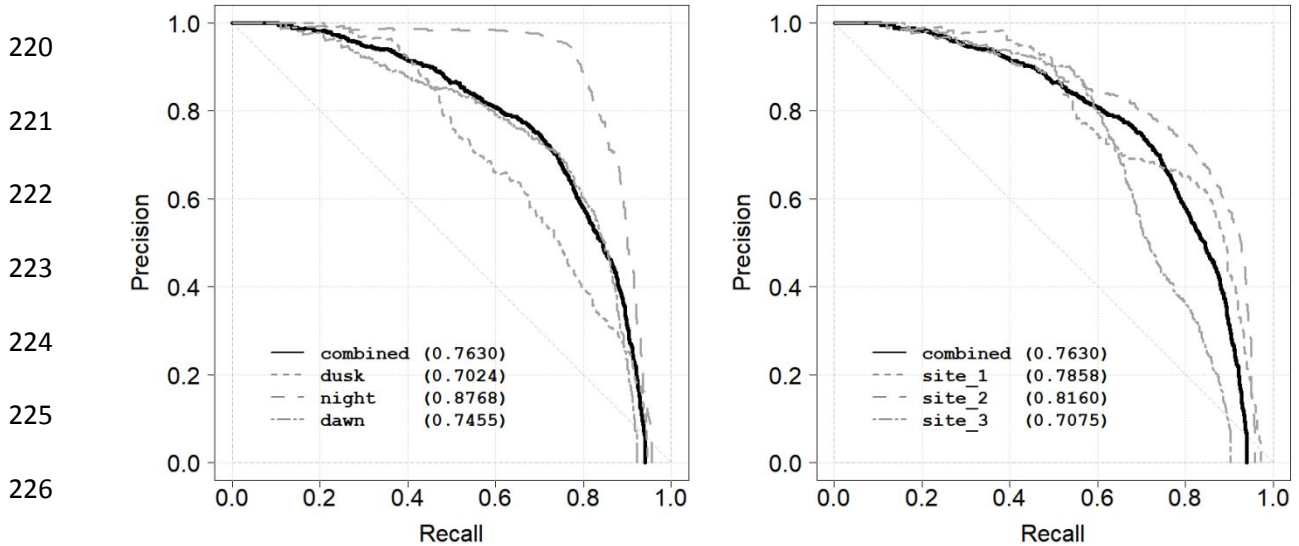
190    recogniser performance. While this imposes a limitation on the future data the results of this

191    research can be applied to, based on the species' ecology and our experience at the study site,

192    this limitation is not onerous, and is one we are willing to accept to improve efficiency. To

193    avoid overfitting, none of the field recordings contained calls that were used to train the

194    recogniser. The dataset was balanced across the three long-term stable roost sites, and three

195    discrete periods of the night: dusk, night, and dawn. Recordings for the dusk period occurred

196    within one hour of sunset, recordings for the dawn period occurred within one hour of sunrise,

197    and recordings for the night period included any time in between the defined dusk and dawn

198    periods. Using audio-editing public domain software Audacity (version 2.3.0,

199    http://audacity.sourceforge.net/), each clip was viewed in a spectrogram (spectrogram settings:

200    y-axis = 0-4000 Hz, x-axis = 30 secs, FFT transformation = Hann window, FFT size = 256),

201    and listened to at a consistent volume using a set of high-quality noise-cancelling headphones

202    (Sennheiser PXC480). 1850 definite Night Parrot calls were detected, ranging from loud calls

203    made in close proximity to the recorder, to faint, distant calls, that could not be seen on a

204    spectrogram and were only detectable by manual listening.

205

206    Each 10-minute recording was then analysed using the call recogniser, with the threshold score

207    set to zero, so all peaks in the similarity score were returned as 'detections'. It is important to

208    note that a 'detection' in this sense is a return from the recogniser representing a prospective

209    detection; it may or may not be an actual detection. The recogniser returned 31437 detections

210    from the 900-minute dataset. These detections were compared to the manually extracted data,

211    and each classified as either a true positive (an actual Night Parrot call) or false positive (not a

212    Night Parrot call). The recogniser did not detect 110 of the 1850 calls in the dataset. These

213    were added to the dataset and classified as false negatives. We assessed baseline performance

214    by producing a precision-recall curve, and calculating the area under the curve (AUC) (Fig. 2).

215     A precision-recall curve plots recall for each value of precision as the classification threshold

216     is reduced, allowing assessment of the trade-off between the two parameters. AUC of the

217     precision-recall curve is the recommended univariate statistic for comparing call recognisers

218     (Knight et al., 2017).

219



228     Figure 2.   Precision-recall curves calculated using raw recogniser scores, including separate

229     curves for each period (left) and site (right). The figures in brackets give the area under the

230     curve (AUC) for each curve. A higher AUC indicates better recogniser performance.

231

232     *Identification of potential intrinsic and contextual variables*

233     We next considered what intrinsic and contextual information could be used to assess the

234     likelihood that any given detection was a true positive detection. From the raw recogniser

235     output we extracted the following intrinsic variables for each detection: the score associated

236     with that detection (*score*); which template resulted in the detection (*template*); and, the parent

237     call type of that template (*call_class*). *Score* is the recogniser's most easily interpreted raw

238     output, with obvious predictive value.

239

240    A comparison of success rates for different values of *call_class* suggested these could have

241    predictive value. The Night Parrot calls incorporated into this recogniser are generally either

242    short or long. Short single notes are common components of other bird and insect calls

243    occurring in the study area, increasing the probability that templates for short calls will generate

244    false positives. Conversely, longer Night Parrot calls are relatively unique in the study area,

245    meaning their templates are less likely to generate false positives (Table 1).

246

247    Table 1.  Success rates for different categories of call templates, with recogniser threshold set

248    to zero. Three letter codes represent the different Night Parrot call types incorporated into the

249    recogniser. Short call templates, particularly the '1di' template, generate most false positives.

250    Most of the long call templates perform well.

| | Short Calls | | | | | Long Calls | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | ddt | too | 1di | 2di | 3nt | 1tr | 2tr | 2wh | 4wh | how |
| TRUE POS. | 50 | 287 | 647 | 25 | 5 | 33 | 13 | 567 | 6 | 107 |
| FALSE POS. | 388 | 4140 | 22053 | 2128 | 156 | 46 | 54 | 521 | 138 | 73 |

251

252

253    For each detection we clipped a 1.1 second segment of the original file that captured the precise

254    time of that detection, then used R package 'seewave' (Sueur, Aubin, & Simonis, 2008) to

255    calculate the difference between the maximum amplitude and mean amplitude within the

256    frequency range of the template that triggered the detection. Binary point matching compares

257    sound energy within a series of designated 'on' and 'off' cells for each template. Loud sounds

258    within the same frequency range as the binary point template can result in high sound energy

259    flooding both the 'on' and 'off' cells, and if slightly more energy is detected in the 'on' cells

260    this will trigger a detection. Typically though, it will receive a relatively low *score*. We

261    reasoned that if there was a large difference between the maximum and mean amplitude within

262    the template's frequency range, and the detection received only a moderate *score*, this was

11

263    likely to represent an example of excess sound energy flooding the template, and therefore a

264    false positive. If a large difference in the maximum and mean amplitude within the template's

265    frequency range resulted in a high *score*, the sound energy probably closely matched the 'on'

266    cells of the template, and was more likely to represent a true positive. A plot of amplitude

267    difference (*amp_diff*) against *score* confirmed this relationship (Fig. 3).

268

269
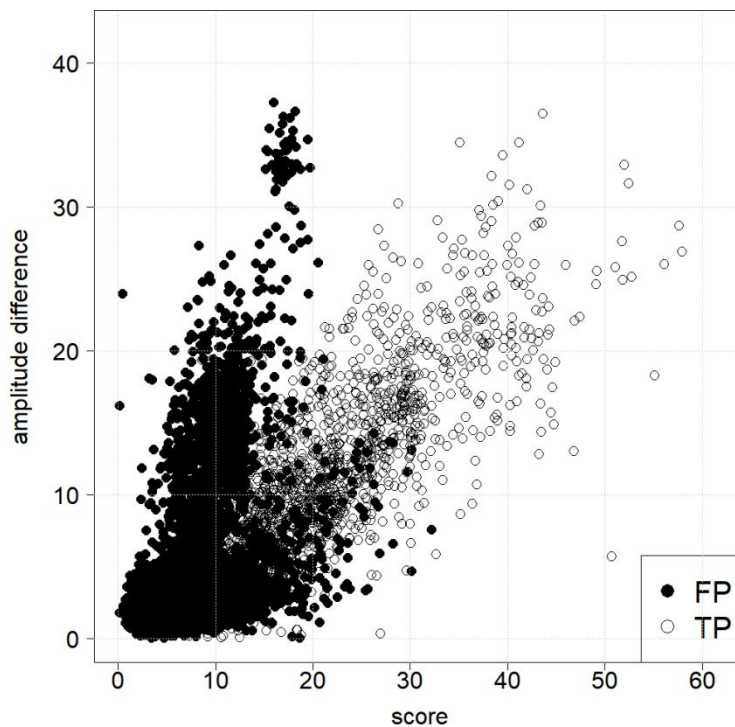
270

271

272

273

274

275

276

277

278

279



280    Figure 3.  Plot of the relationship between amplitude difference and score for each detection,

281    categorised by detection classification (true positive or false positive). As predicted, detections

282    with a higher amplitude difference but moderate to low score are mostly false positives.

283

284    We next considered potential contextual variables. All detections were classified according to

285    which *period* ('dusk', 'night' and 'dawn'), and which *site* they were recorded from ('site_1',

286    'site_2', 'site_3'). Precision-recall curves were plotted and AUC calculated for each *period*

287    and *site*, then compared to the recogniser's baseline precision-recall curve, to explore their

288  influence on recogniser performance (Fig. 2). Recogniser performance varied between periods,

289  performing best during the night, and most poorly at dusk. This is expected, given the

290  likelihood of false positives is reduced during the night when diurnal birds are not calling.

291  There was no apparent effect of site on recogniser performance. For each detection we also

292  noted which model of ARU (*ARU_type*) and which specific ARU (*machine*) recorded the

293  detection, and in which of the 90 test files (*file*) the detection occurred.

294

295  *Model development procedure*

296  Our aim was to determine whether a model-derived probability calculated using intrinsic and

297  contextual variables could be substituted for the recogniser's initial *score* value, and achieve

298  better results. We chose a generalised linear mixed-effects model structure, to enable inclusion

299  of both fixed and random effects. As our response variable was binary (true positive or false

300  positive), models were fitted assuming a binomial response distribution, and a logit link

301  function (logistic regression) using the lme4 package (Bates & Sarkar, 2007).

302

303   As the practical purpose of this model is to facilitate the process of sifting through recogniser

304  outputs, the process of model building can be more informal than for research purposes that

305  involve *a priori* questions. The approach to selecting the final model was to initially generate

306  a comprehensive set of possible fixed and random effects and compare candidate models

307  containing main effects and interactions for the fixed effect terms, together with the random

308  effects. We then assessed the performance of the candidate models via summary statistics and

309  selected the most promising ones for further development. We determined which variables and

310  variable combinations were critical to those models' performance. Finally, we re-evaluated the

311  refined models before selecting the best performing model as the final model. Model selection

312  was completed using the entire performance dataset.

313

*Fixed and random effects selection*

As the aim was to apply the model developed using the performance dataset to any data collected at the study site, we limited fixed effects to those whose complete range of variation was represented in the performance dataset, and which could be determined *a priori* from the resulting raw recogniser output. Factors whose variation was not entirely represented in the performance dataset were included as random effects, and not used in predictions. For example, as *ARU_type* for any data collected at the study site will be either SM3 or SM4, and both were adequately represented in the performance dataset, this could be included as a fixed effect. However, more than 80 individual ARUs have been used at the study site, and only a portion of these were represented in the performance dataset. As this portion represents a random sample from the set of possible ARUs, *machine* (representing the specific ARU used) is included as a random effect. This still allowed the variance associated with this factor to be captured and an allowance made for it in the training phase, but only that level of variance determined during the training phase can be used when the model is applied to future data collected from any machine.

Data exploration revealed interactions were needed between *score* and both *period* and *amp_diff*, so these were initially included as a three-way interaction fixed effect. Because the relationship between a detection's *score* and the probability that the detection is a true positive is curved in the logistic scale, *score* was fitted as a quadratic term. Also included as fixed effects were *call_class* and *ARU_type*. As factors whose level will very likely be new for future datasets, *site*, *file*, and *machine* were all included as random effects. The factor *template* can be established *a priori* from the raw results, but as it contains 31 levels and is nested within

337  *call_class* its predictive power is likely to be limited. However, understanding its impact on

338  model performance may still be important, so it was included as a random effect.

339

340  We initially tested a series of 16 models. Each model included all fixed effects, but varied in

341  the combination of random effects. All possible combinations of the four random effects were

342  tested, including a model with no random effects. Models were compared using both Akaike's

343  Information Criterion (AIC) and Bayesian Information Criterion (BIC). AIC and BIC are

344  statistics for comparing relative model performance, with the primary difference being that

345  BIC penalises more heavily for model complexity (Burnham & Anderson, 2004). Four models

346  stood out as having much lower AIC than the other 12 (Table 2). These four models also had

347  a much lower BIC than the other 12 models. Examining the variance components for each

348  random effect revealed that *file* and *template* were the source of most variation in each of the

349  four best-ranked models, with the contribution of both *machine* and *site* limited (Table 3).

350  Therefore, we retained *file* and *template* as random effects.

351

352  We next ran the model including all fixed effects and our chosen random effects, before

353  examining the significance of resulting individual fixed effect coefficients (Table 4). These

354  suggest that the three-way interaction between *period*, *score* and *amp_diff* is not substantially

355  influencing model performance, but that each of the two way interactions between these

356  variables should be retained. *Call_class* has an effect on model performance, but not

357  consistently across classes. Calls that are short have less influence on the model than calls

358  which are long. To investigate this, we created two new variables based on call length. The

359  variable *call_length_1* categorised detections based on the template that detects the call as

360  either short or long, while *call_length_2* categorised all detections based on the template that

361   detects the call as either short, medium, or long. The influence of *ARU_type* is significant, but

362   marginally so.

363

364   Table 2.  Summary statistics for all random effects models, ranked by AIC. There is strong

365   support for the top four models, warranting further inspection of each component's variation

366   within these models.

| Random effects | AIC | BIC | Deviance | log lik. | Resid. df |
|---|---|---|---|---|---|
| file + template + site | 2520.79 | 2779.82 | 2174.11 | -1229.40 | 31406 |
| machine + file + template | 2520.96 | 2779.98 | 2174.03 | -1229.48 | 31406 |
| machine + file + template + site | 2522.75 | 2790.14 | 2174.13 | -1229.38 | 31405 |
| file + template | 2528.47 | 2779.14 | 2172.61 | -1234.23 | 31407 |
| file + site | 2716.34 | 2967.01 | 2436.64 | -1328.17 | 31407 |
| machine + file | 2716.37 | 2967.04 | 2436.45 | -1328.18 | 31407 |
| machine + file + site | 2718.31 | 2977.34 | 2436.58 | -1328.16 | 31406 |
| file | 2722.61 | 2964.93 | 2434.75 | -1332.31 | 31408 |
| machine + template | 2730.59 | 2981.26 | 2561.90 | -1335.30 | 31407 |
| machine + template + site | 2732.03 | 2991.06 | 2561.98 | -1335.01 | 31406 |
| template + site | 2740.30 | 2990.97 | 2581.41 | -1340.15 | 31407 |
| template | 2840.21 | 3082.53 | 2699.26 | -1391.10 | 31408 |
| machine | 2955.80 | 3198.11 | 2873.55 | -1448.90 | 31408 |
| machine + site | 2957.47 | 3208.15 | 2873.65 | -1448.74 | 31407 |
| site | 2965.63 | 3207.95 | 2892.61 | -1453.82 | 31408 |
| fixed effects only | 3066.66 | 3300.62 | 3010.66 | -1505.33 | 31409 |

367

368   We tested a series of nine models, including all possible combinations of the following fixed

369   effects: *score*, *period* and *amp_diff* as either a three-way, or three separate two-way

370   interactions; template category as either *call_class*, *call_length_1* or *call_length_2*; and, with

371   or without *ARU_type*. The random effects for *file* and *template* were retained for all models.

372   The three best models had an AIC value no larger than one unit above the model with the

373   minimum AIC (Table 5). However, the third ranked of these models had a much lower BIC

374   than the other two, with $\Delta$BIC > 30 between this model and the next ranked model by BIC.

375 Given there was not clear support for one of these three models using AIC, we contend that the

376 best-ranked model using BIC could be considered preferable. We selected this model for use

377 in practice.

378

379 Table 3. Variance of each random effects component within each of the top four models used

380 for random effects testing. The contribution of both *machine* and *site* are limited in each case,

381 supporting the decision to retain only *file* and *template* for model simplicity.

| *file + template + site* | |
| --- | --- |
| Component | Std dev. |
| file | 1.2177 |
| template | 1.2545 |
| site | 0.6554 |

| *machine + file + template* | |
| --- | --- |
| Component | Std dev. |
| file | 1.2113 |
| template | 1.2492 |
| machine | 0.5789 |

| *machine + file + template + site* | |
| --- | --- |
| Component | Std dev. |
| file | 1.2127 |
| template | 1.2538 |
| machine | 0.2847 |
| site | 0.5536 |

| *file + template* | |
| --- | --- |
| Component | Std dev. |
| file | 1.3584 |
| template | 1.2222 |

382

383 *Model testing*

384 To test the model, we partitioned the performance dataset, using one third of the files, balanced

385 by site and period, to train the model. The remaining files were set aside to test the model. After

386 training, the model was used to predict whether each detection in the test dataset was a true

387 positive. Because we would not know file in advance for a future dataset, this random effect

388 was predicted using the estimate from model training. The predicted probability for each

389 detection was then then substituted for raw recogniser score, and the precision-recall curves

390 replotted (Fig. 4).

391

392    Table 4.  Significance of the fixed effect coefficients for the model incorporating all fixed

393    effects. Of particular note are the consistent differences between short calls ('ddt', '1di', '2di',

394    '3nt', 'too') and long calls ('1tr', '2tr', '2wh', '4wh', 'how').

| Fixed effect | Estimate | Std. Error | z value | Pr(>\|z\|) |
|---|---|---|---|---|
| (Intercept) | -7.271 | 0.702 | -10.356 | 0.000 |
| *period* 'dusk' | 1.623 | 0.543 | 2.986 | 0.003 |
| *period* 'night' | 0.970 | 0.588 | 1.650 | 0.099 |
| *score*$^2$(1) | 133.608 | 56.883 | 2.349 | 0.019 |
| *score*$^2$(2) | -494.164 | 59.590 | -8.293 | 0.000 |
| *amp_diff* | 0.801 | 0.086 | 9.258 | 0.000 |
| *ARU_type* 'SM4' | -1.367 | 0.416 | -3.288 | 0.001 |
| *call_class* '1tr' | 5.001 | 1.469 | 3.404 | 0.001 |
| *call_class* '2di' | 0.054 | 0.787 | 0.068 | 0.945 |
| *call_class* '2tr' | 5.767 | 1.868 | 3.088 | 0.002 |
| *call_class* '2wh' | 3.352 | 0.763 | 4.391 | 0.000 |
| *call_class* '3nt' | 1.972 | 1.198 | 1.645 | 0.100 |
| *call_class* '4wh' | 3.824 | 1.449 | 2.638 | 0.008 |
| *call_class* 'ddt' | 2.254 | 1.348 | 1.673 | 0.094 |
| *call_class* 'how' | 5.211 | 1.323 | 3.938 | 0.000 |
| *call_class* 'too' | 1.597 | 0.984 | 1.623 | 0.105 |
| *period* 'dusk': *score*$^2$(1) | 173.151 | 71.366 | 2.426 | 0.015 |
| *period* 'night': *score*$^2$(1) | -63.386 | 107.565 | -0.589 | 0.556 |
| *period* 'dusk': *score*$^2$(2) | 198.939 | 73.734 | 2.698 | 0.007 |
| *period* 'night': *score*$^2$(2) | -214.750 | 102.485 | -2.095 | 0.036 |
| *period* 'dusk':*amp_diff* | -0.588 | 0.094 | -6.288 | 0.000 |
| *period* 'night':*amp_diff* | -0.428 | 0.106 | -4.024 | 0.000 |
| *score*$^2$(1):*amp_diff* | -7.443 | 5.904 | -1.261 | 0.207 |
| *score*$^2$(2):amp_diff | 34.896 | 5.888 | 5.926 | 0.000 |
| *period* 'dusk': *score*$^2$(1):*amp_diff* | 9.761 | 7.708 | 1.266 | 0.205 |
| *period* 'night': *score*$^2$(1):*amp_diff* | 27.259 | 12.111 | 2.251 | 0.024 |
| *period* 'dusk': *score*$^2$(2):*amp_diff* | -8.322 | 8.106 | -1.027 | 0.305 |
| *period* 'night': *score*$^2$(2):*amp_diff* | 5.248 | 11.050 | 0.475 | 0.635 |

395

396

397

398

399 Table 5. Summary statistics for the final set of nine models. Only fixed effects for each model

400 are shown; the random effects for each model were file and template. There is strong support

401 for each of the top three models by AIC, but the third of these (in bold) has much stronger

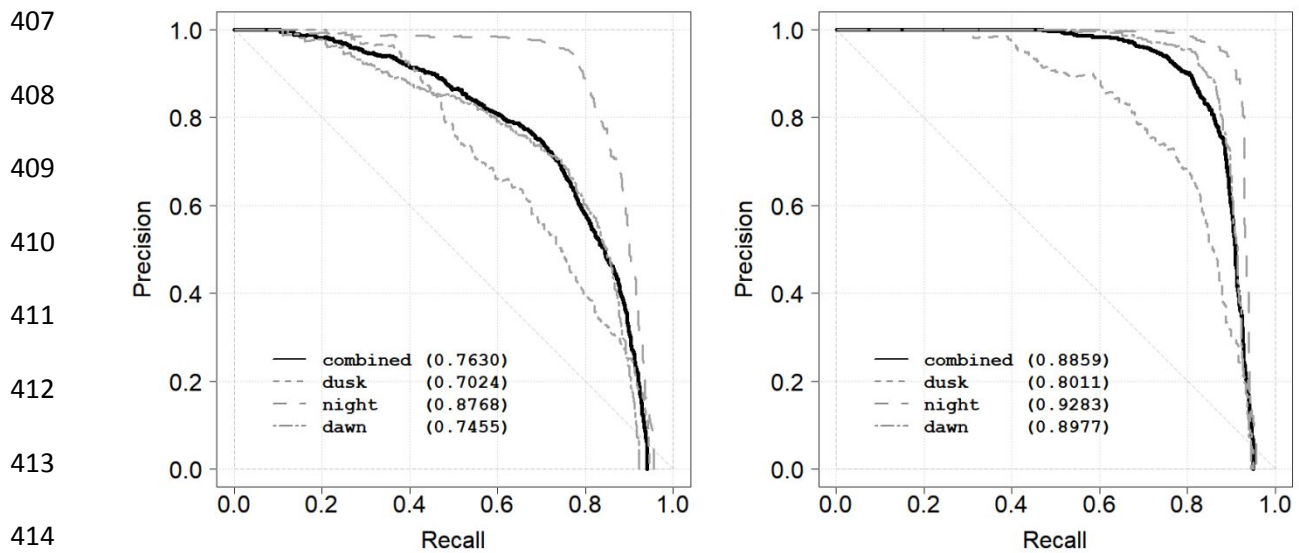402 support by BIC and was selected as the final model.

| Fixed effects | AIC | BIC | Deviance | log lik. | Resid. df |
|---|---|---|---|---|---|
| period * score$^2$ * amp_diff + call_length_1 + ARU_type | 2522.12 | 2705.94 | 2171.67 | -1239.06 | 31415 |
| period * score$^2$ * amp_diff + call_length_2 + ARU_type | 2522.60 | 2714.78 | 2169.74 | -1238.30 | 31414 |
| **period * score$^2$ + score$^2$ * amp_diff + period * amp_diff + call_length_1 + ARU_type** | **2522.84** | **2673.25** | **2180.75** | **-1243.42** | **31419** |
| period * score$^2$ + score$^2$ * amp_diff + period * amp_diff + call_length_2 + ARU_type | 2524.54 | 2683.30 | 2179.03 | -1243.27 | 31418 |
| period * score$^2$ * amp_diff + call_class + ARU_type | 2528.47 | 2779.14 | 2172.61 | -1234.23 | 31407 |
| period * score$^2$ + score$^2$ * amp_diff + period * amp_diff + call_class + ARU_type | 2529.69 | 2746.94 | 2182.30 | -1238.84 | 31411 |
| period * score$^2$ + score$^2$ * amp_diff + period * amp_diff + call_length_1 | 2532.23 | 2674.28 | 2179.28 | -1249.12 | 31420 |
| period * score$^2$ + score$^2$ * amp_diff + period * amp_diff + call_length_2 | 2533.67 | 2684.07 | 2177.83 | -1248.83 | 31419 |
| period * score$^2$ + score$^2$ * amp_diff + period * amp_diff + call_class | 2538.85 | 2747.74 | 2181.04 | -1244.42 | 31412 |

403

404

405

406

Figure 4. Precision recall curves calculated for each period using raw recogniser scores (left), and model-derived probabilities (right). When using model-derived probabilities, the increase in AUC is evident overall, and across all periods, meaning this approach improves recogniser performance.

The precision-recall curves for the combined data, and for each *period*, demonstrate that substituting model-derived probability for raw score results in an increased AUC overall (AUC = 0.89 for model-derived probability, and AUC = 0.76 for raw score), meaning overall recogniser performance is improved. As expected, this improvement is modest for the night *period*, but marked for both the dusk and dawn *period*, with AUC improving by 0.10 and 0.15 respectively.

To quantify the practical improvements resulting from this modelling procedure, we investigated the number of detections requiring proofing to achieve a specific level of recall. Recall is of particular importance because the recall of a recogniser equals the probability that a species will be detected if it is available for detection, an important component of the overall probability of detection (Pollock et al., 2004). Furthermore, it is important for rare species

20

432　research because prioritising recall maximises the likelihood of detecting the species if it is

433　available in the acoustic dataset. This emphasis on recall manifests itself in the increased

434　number of detections that require proofing to achieve the increased level of recall.

435

436　We calculated the mean number of false positive detections requiring proofing per 10-minute

437　file in the test dataset to achieve a specific recall; a proxy for the amount of time an analyst

438　needs to spend proofing recogniser output. We initially calculated the score cut-off that

439　achieved a specified recall for both raw score, and for the model-derived probability. Because

440　model-derived probability incorporates *period* as a fixed effect in the calculation, cut-off scores

441　for a specific value of recall under the model-derived probability may vary between periods.

442　Accordingly, the model-derived probability cut-off for each recall threshold was calculated

443　separately for each *period* using only the test dataset to avoid overfitting. Using these data, we

444　also simulated for both raw score and model-derived probability, how many false positive

445　detections would need to be checked during a complete 12 hour night of acoustic data, with

446　one hour of 'dusk', ten hours of 'night' and one hour of 'dawn' recordings to be assessed.

447

448　The model-derived probability markedly reduced the number of false positives that needed

449　checking to achieve each level of recall tested (Table 6). This improvement is most pronounced

450　during the night period, and at lower levels of recall. However, even at 90% recall, if using the

451　model-derived probability as a substitute for score, the number of false positives that would

452　need checking during an entire night of acoustic data is 40% of what would need to be checked

453　if using the raw score.

454

455

Table 6. The mean number of false positives requiring proofing in a 10-minute recording for a set level of recall, using either raw recogniser score (Score), or the model-derived probability (MDP). The final three columns present the number of false positives that would need proofing if analysing a 12-hour night of recordings, with the '%' column representing the percentage of proofing, and therefore time required when using model-derived probability compared to raw score.

| Recall | Dusk | | Night | | Dawn | | 12-hour Night | | |
|--------|-------|------|-------|------|-------|------|-------|-------|-----|
| | Score | MDP | Score | MDP | Score | MDP | Score | MDP | % |
| 0.50 | 2.80 | 0.85 | 0.20 | 0.00 | 0.70 | 0.00 | 30.6 | 5.1 | 17 |
| 0.55 | 4.30 | 1.10 | 0.25 | 0.05 | 1.00 | 0.00 | 43.8 | 9.0 | 21 |
| 0.60 | 6.20 | 1.30 | 0.25 | 0.05 | 1.55 | 0.00 | 58.5 | 10.2 | 17 |
| 0.65 | 7.55 | 2.15 | 0.25 | 0.05 | 2.00 | 0.10 | 69.3 | 15.9 | 23 |
| 0.70 | 9.80 | 3.30 | 0.40 | 0.05 | 2.60 | 0.35 | 93.6 | 24.3 | 26 |
| 0.75 | 15.30 | 4.45 | 0.50 | 0.05 | 3.65 | 0.55 | 137.7 | 32.4 | 24 |
| 0.80 | 22.25 | 6.55 | 0.70 | 0.25 | 5.80 | 0.85 | 201.9 | 56.4 | 28 |
| 0.85 | 29.30 | 13.70 | 1.85 | 0.55 | 9.70 | 2.35 | 322.8 | 122.7 | 38 |
| 0.90 | 45.35 | 34.95 | 9.05 | 1.35 | 22.25 | 10.10 | 840.0 | 335.1 | 40 |

**Discussion**

The method we have outlined demonstrates that intrinsic and contextual information associated with a call recogniser's output can be used to improve the performance of that recogniser. This approach is compatible with any signal detection algorithm, not just binary point matching as is the case here. While the improvements are revealed through the AUC of the precision-recall curve, this representation is somewhat abstract. The practical benefits of this approach are more clearly demonstrated in the reduced effort required to achieve a specific recall. For practitioners using call recognisers to analyse large quantities of field recordings, the limiting factor is typically time, which manifests itself as the number of detections that can be manually proofed. However, while this technique does result in efficiencies, there are limitations.

474

*Raw recogniser performance and improvement*

These improvements will only apply to detections within the recogniser's output; it does not change the recogniser's ability to detect false negatives. False negatives occur for two reasons. The recogniser may detect some other signal that occurs concurrently with the call of interest and achieves a higher score, meaning the call of interest is missed. Such events are difficult to overcome. Alternatively, a call of interest may not match the training data. Post-processing techniques, as outlined here, will not improve recogniser performance in that respect. This can only be overcome by updating the recogniser's training dataset to improve the probability the recogniser will detect that missed call. If new templates are added to the recogniser, the model selection process will need to be rerun, with sufficient training and test files added to model the impact of the new templates.

*Model application for different species and new sites*

Even though calls used to create this recogniser's templates were excluded from the training and test datasets, because the Night Parrot population at the study site is very small, it is likely calls from the same individuals were incorporated into the training and test datasets. There is a resultant risk of model overfitting. Additionally, the repertoire of this population is well-known (Leseberg et al., 2019), and the recogniser templates featured most of the variation that occurs at the study site. It is possible this combination of factors has exaggerated the success of our model. In scenarios where the subject species does not have such a consistent repertoire, because it has a larger number of individuals, a more dynamic population, or greater variation in its calls, this technique will still be applicable provided this variation is incorporated into the training and test datasets.

499    The properties of the general soundscape, including likely non-target calls that occur in the

500    dataset will also influence model applicability. For example, the model developed here could

501    be reasonably applied to other datasets from western Queensland, where Night Parrots are

502    known to have similar calls to those in this dataset (NL pers. obs.), and where the suite of likely

503    non-target species will also be similar. However, the model may not be as effective if applied

504    to a dataset from Western Australia, where the suite of Night Parrot and likely non-target

505    species are slightly different to western Queensland. Testing on an annotated dataset would

506    determine if the model does improve recogniser performance and by how much. Otherwise,

507    the model selection and training process would need to be rerun using a performance dataset

508    compiled from the new region of interest.

509

510    *Impact of model treatment of different call types*

511    The fixed effect *call_length_1* boosts the model-derived probability for longer calls, when

512    compared to shorter calls. In a scenario where shorter calls predominate at a site, this may affect

513    the recogniser's ability to detect birds at that site. It is likely that faint short calls are most

514    affected. Because an ARU established at a prospective long-term stable roost site will record a

515    variety of short calls over time, the probability of at least some calls being detected by the

516    recogniser is high. Additionally, over long periods at long-term stable roost sites, there is

517    typically a mix of long and short calls (SM, NL unpub. data), ensuring the recogniser will detect

518    birds if they are present. This may still be an issue if a short deployment limits the variety of

519    calls that occur within the dataset.

520

521    An additional consequence of the differing treatment of call types by the model will be the

522    distortion of potential distance effects. Researchers can extract distance information from

523    acoustic data, using signal strength, or variables closely related to signal strength such as the

524     call recogniser's raw score, as a proxy for distance from the recorder (Knight & Bayne, 2018;

525     Lambert & McDonald, 2014). This information is then used in distance-sampling procedures,

526     or for establishing survey effort parameters (Yip, Leston, Bayne, Solymos, & Grover, 2017).

527     The mechanics of this modelling technique will confound any attempts to use the model-

528     derived probabilities as a proxy for distance, because they are influenced by factors other than

529     signal strength, whereas raw score is typically heavily dependent on signal strength (Knight &

530     Bayne, 2018). For example, if ranked by model-derived probability, a faint long call is likely

531     to rank higher than if it were ranked by raw score alone. If model-derived probability is being

532     used as a proxy for distance from the recorder, this would be equivalent to the call being made

533     closer to the recorder, an incorrect assumption that could distort conclusions around that call's

534     likely distance from the recorder.

535

536     Depending on the aim of the distance-sampling approach, this issue could be overcome in

537     several ways, although each has limitations. Research could assess the relationship between

538     model outputs and distance, although this is likely to vary across call types, and for a species

539     like the Night Parrot would require a test dataset that would be almost impossible to collect.

540     Alternatively, signal strength or raw score for a given detection could be extracted after model

541     application to determine distance data, although this will mean the calls extracted will be

542     influenced by the model. Again, long calls are more likely to be extracted than short calls,

543     possibly interfering with subsequent conclusions. A final option could be to first sort data by

544     raw score, before applying the model to the subset of data whose raw score satisfies the distance

545     sampling criteria.

546

547     *Other parameters with potential predictive power*

548      The modelling approach applied here was successful using a relatively limited number of

549      parameters, some that were particular to the subject species' biology, such as *call_length_1*

550      and *period*, while others were generic, such as *amp_diff*, *ARU_type* and the random effects

551      *template* and *file*. It is likely that a number of other parameters could be incorporated to further

552      improve results. As Night Parrots call more frequently in response to local rain events (Murphy,

553      Austin, et al., 2017), a variable quantifying antecedent rainfall could be an obvious inclusion.

554      An emerging question in Night Parrot research is the merit of acoustic surveys at water points

555      and likely feeding sites, compared to current protocols that focus solely on roosting habitat. If

556      autecological research determines a consistent pattern of nocturnal activity, *site resource* (i.e.

557      water point, feeding site, roosting site) could be included as a fixed effect in the model.

558

559      The predictable calling behaviour and site fidelity of the Night Parrot make it particularly suited

560      to the approach we have outlined here, but with careful consideration, it will be applicable in

561      other scenarios. Intrinsic variables related to raw recogniser output can be developed that are

562      species specific, as call type was here, or recogniser specific, as *amp_diff* was in this case,

563      being relevant specifically to the binary point matching technique used in this recogniser. There

564      are likely to be similar variables that could be developed for the numerous other recogniser

565      algorithms. Improvements to the raw output for more advanced algorithms may not be as

566      significant as for the relatively basic binary point matching, but for field ecologists, any

567      reduction in the time required to proof recogniser returns will be beneficial. The contextual

568      variables that could be trialled will relate to a species' biology and might include long-term

569      seasonal and short-term weather effects, habitat or other environmental parameters at both the

570      local and landscape scale, and calling biology. The number of contextual parameters that could

571      be tested is limited only by a researcher's ability to compile a performance testing dataset that

572      satisfactorily represents the variation in each parameter.

573

574  This technique's biggest advantages are its simplicity, and compatibility with any recognition

575  algorithm. For the ecologist or practitioner, call recogniser development is daunting, with high

576  performing recognisers generally built using state-of-the-art techniques that in many cases

577  require advanced programming skills and research time. The foundation of the post-processing

578  technique we outline in this paper is a relatively straightforward procedure that can be

579  completed using graduate level statistics. For that reason, it will be of particular use to

580  practicing field ecologists looking to improve a simple recogniser, which may only be one part

581  of a broader research project. It may also be applied to any state-of-the-art recognition

582  algorithm to further improve results.

583

584  **Acknowledgements**

591

592  **Authors' contributions**

593  NL, SM, WV, and JW conceived the ideas and designed the methodology; NL, SM and JW

594  collected the data; NL and WV analysed the data; NL led the writing of the manuscript. All

595  authors contributed critically to the drafts and gave final approval for publication.

596

597  **Data availability**

598 We intend to make the outputs of the recogniser, and the code used to create and apply our

599 model available via Github.

600

601 **References**

602 Bates, D. M., & Sarkar, D. (2007). lme4: Linear mixed-effects models using S4 classes.

603 Burnham, K. P., & Anderson, D. R. (2004). Multimodel Inference:Understanding AIC and BIC

604     in Model Selection. *Sociological Methods & Research, 33*(2), 261-304.

605     doi:10.1177/0049124104268644

606 Digby, A., Towsey, M., Bell, B. D., & Teal, P. D. (2013). A practical comparison of manual

607     and autonomous methods for acoustic monitoring. *Methods in Ecology and Evolution,*

608     *4*, 675-683. doi:10.1111/2041-210X.12060

609 Higgins, P. J. (1999). *Night Parrot (Pezoporus occidentalis)* (Vol. 4, Parrots to Dollarbird).

610     South Melbourne: Oxford University Press.

611 Horton, K. G., Stepanian, P. M., Wainwright, C. E., & Tegeler, A. K. (2015). Influence of

612     atmospheric properties on detection of wood-warbler nocturnal flight calls.

613     *International Journal of Biometeorology, 59*, 1385-1394. doi:10.1007/s00484-014-

614     0948-8

615 Joshi, K. A., Mulder, R. A., & Rowe, K. M. (2017). Comparing manual and automated species

616     recognition in the detection of four common south-east Australian forest birds from

617     digital field recordings. *Emu, 117*(3), 233-246. doi:10.1371/journal.pone.0199396

618 Katz, J., Hafner, S. D., & Donovan, T. (2016). Tools for automated acoustic monitoring within

619     the R package monitoR. *Bioacoustics, 25*(2), 197-210.

620     doi:10.1080/09524622.2016.1138415

621 Knight, E. C., & Bayne, E. M. (2018). Classification threshold and training data affect the

622    quality and utility of focal species data processed with automated audio-recognition

623    software. *Bioacoustics*. doi:10.1080/09524622.2018.1503971

624 Knight, E. C., Hannah, K. C., Foley, G. J., Scott, C. D., Brigham, R. M., & Bayne, E. (2017).

625    Recommendations for acoustic recognizer performance assessment with application to

626    five common automated signal recognition programs. *Avian Conservation and*

627    *Ecology, 12*(2). doi:10.5751/ACE-01114-120214

628 Koops, H. V., van Balen, J., & Wiering, F. (2014). A deep Neural Network Approach to the

629    LifeCLEF 2014 Bird task. *CEUR Workshop Proceedings, 1180*, 634-642.

630 Lambert, K. T. A., & McDonald, P. G. (2014). A low-cost, yet simple and highly repeatable

631    system for acoustically surveying cryptic species. *Austral Ecology, 39*, 779-785.

632    doi:10.1111/aec.12143

633 Landau, H. J. (1967). Sampling, data transmission, and the Nyquist rate. *Proceedings of the*

634    *IEEE, 55*(10), 1701-1706. doi:10.1109/PROC.1967.5962

635 Leseberg, N. P., Murphy, S. A., Jackett, N. A., Greatwich, B. R., Brown, J., Hamilton, N., . . .

636    Watson, J. E. M. (2019). Descriptions of known vocalisations of the Night Parrot

637    *Pezoporus occidentalis*. *Australian Field Ornithology, 36*, 79-88.

638    doi:10.20938/afo36079088

639 Magurran, A. E., Baillie, S. R., Buckland, S. T., Dick, J. P., Elston, D. A., Scott, E. M., . . .

640    Watt, A. D. (2010). Long-term datasets in biodiversity research and monitoring:

641    assessing change in ecological communities through time. *Trends in Ecology and*

642    *Evolution, 25*(10), 574-582. doi:10.1016/j.tree.2010.06.016

643 Murphy, S. A., Austin, J. J., Murphy, R. K., Silcock, J., Joseph, L., Garnett, S. T., . . . Burbidge,

644    A. H. (2017). Observations on breeding Night Parrots (*Pezoporus occidentalis*) in

645    western Queensland. *Emu, 117*(2), 107-113. doi:10.1080/01584197.2017.1292404

646 Murphy, S. A., Silcock, J., Murphy, R. K., Reid, J. R. W., & Austin, J. J. (2017). Movements

647 and habitat use of the night parrot *Pezoporus occidentalis* in south‐western

648 Queensland. *Austral Ecology, 42*(7), 858-868. doi:10.1111/aec.12508

649 Pollock, K. H., Marsh, H., Bailey, L. L., Farnsworth, G. L., Simons, T. R., & Alldredge, M.

650 W. (2004). Separating Components of Detection Probability in Abundance Estimation:

651 An Overview with Diverse Examples. In W. L. Thompson (Ed.), *Sampling rare or*

652 *elusive species: concepts, designs, and techniques for estimating population*

653 *parameters*. Washington, USA: Island Press.

654 Potamitis, I., Ntalampiras, S., Jahn, O., & Riede, K. (2014). Automatic bird sound detection in

655 long real-field recordings: Applications and tools. *Applied Acoustics, 80*, 1-9.

656 doi:10.1016/j.apacoust.2014.01.001

657 Priyadarshani, N., Marsland, S., & Castro, I. (2018). Automated birdsong recognition in

658 complex acoustic environments: a review. *Journal of Avian Biology, 49*(5).

659 doi:10.1111/jav.01447

660 Priyadarshani, N., Marsland, S., Juodakis, J., Castro, I., & Listanti, V. (2020). Wavelet filters

661 for automated recognition of birdsong in long-time field recordings. *Methods in*

662 *Ecology and Evolution, 11*, 403-417. doi:10.1111/2041-210X.13357

663 R Core Team. (2018). R: A Language and Environment for Statistical Computing: R

664 Foundation for Statistical Computing, Vienna.

665 Salamon, J., & Bello, J. P. (2017). Deep Convolutional Neural Networks and Data

666 Augmentation for Environmental Sound Classification. *IEEE Signal Processing*

667 *Letters, 24*(3), 279-283. doi:10.1109/LSP.2017.2657381

668 Shonfield, J., & Bayne, E. M. (2017). Autonomous recording units in avian ecological research:

669 current use and future applications. *Avian Conservation and Ecology, 12*(1).

670 doi:10.5751/ACE-00974-120114

671    Sueur, J., Aubin, T., & Simonis, C. (2008). Seewave: A free modular tool for sound analysis

672        and synthesis. *Bioacoustics, 18*, 213-226. doi:10.1080/09524622.2008.9753600

673    Teixeira, D., Maron, M., & van Rensburg, B. J. (2019). Bioacoustic monitoring of animal vocal

674        behaviour for conservation. *Conservation Science and Practice*. doi:10.1111/csp2.72

675    Towsey, M., Wimmer, J., Williamson, I., & Roe, P. (2014). The use of acoustic indices to

676        determine avian species richness in audio-recordings of the environment. *Ecological*

677        *Informatics, 21*, 110-119. doi:https://doi.org/10.1016/j.ecoinf.2013.11.007

678    Towsey, M., Znidersic, E., Broken-Brow, J., Indraswari, K., Watson, D. M., Phillips, Y., . . .

679        Roe, P. (2018). Long-duration, false-colour spectrograms for detecting species in large

680        audio data-sets. *Journal of Ecoacoustics, 2*. doi:10.22261/JEA.IUSWUI

681    Williams, E. M., O'Donnell, C. F. J., & Armstrong, D. P. (2018). Cost-benefit analysis of

682        acoustic recorders as a solution to sampling challenges experienced monitoring cryptic

683        species. *Ecology and Evolution, 8*, 6839-6848. doi:10.1002/ece3.4199

684    Yip, D. A., Leston, L., Bayne, E. M., Solymos, P., & Grover, A. (2017). Experimentally

685        derived detection distances from audio recordings and human observers enable

686        integrated analysis of point count data. *Avian Conservation and Ecology, 12*(1).

687        doi:10.5751/ACE-00997-120111

688