A sperm whale cautionary tale about estimating acoustic cue rates for deep divers

Tiago A. Marques,[1,2,a] Carolina S. Marques,[2] and Kalliopi C. Gkikopoulou[1,3]

*[1] Centre for Research into Ecological and Environmental Modelling, The Observatory, University of St Andrews, St Andrews, KY16 9LZ, Scotland, +44 1334 461842, +44 1334 461800*

*[2] Centro de Estatística e Aplicações, Departamento de Biologia Animal, Faculdade de Ciências da Universidade de Lisboa, Portugal*

*[3] Sea Mammal Research Unit, Scottish Ocean Institute, University of St Andrews, St Andrews, Fife, KY17 9LZ, UK*

1  Passive acoustic density estimation has been gaining traction in recent years. Cue

2  counting uses detected acoustic cues to estimate animal abundance. A cue rate,

3  the number of acoustic cues produced per animal per unit time, is required to

4  convert cue density into animal density. Cue rate information can be obtained

5  from animal borne acoustic tags. For deep divers, like beaked whales, data have

6  been analyzed considering deep dive cycles as a natural sampling unit, based on

7  either weighted averages or generalized estimating equations. Using a sperm

8  whale DTAG (sound-and-orientation recording tag) example we compare

9  different approaches of estimating cue rate from acoustic tags, illustrating that

10  both approaches used before, might introduce biases and suggest that the natural

11  unit of analysis should be the whole duration of the tag itself.
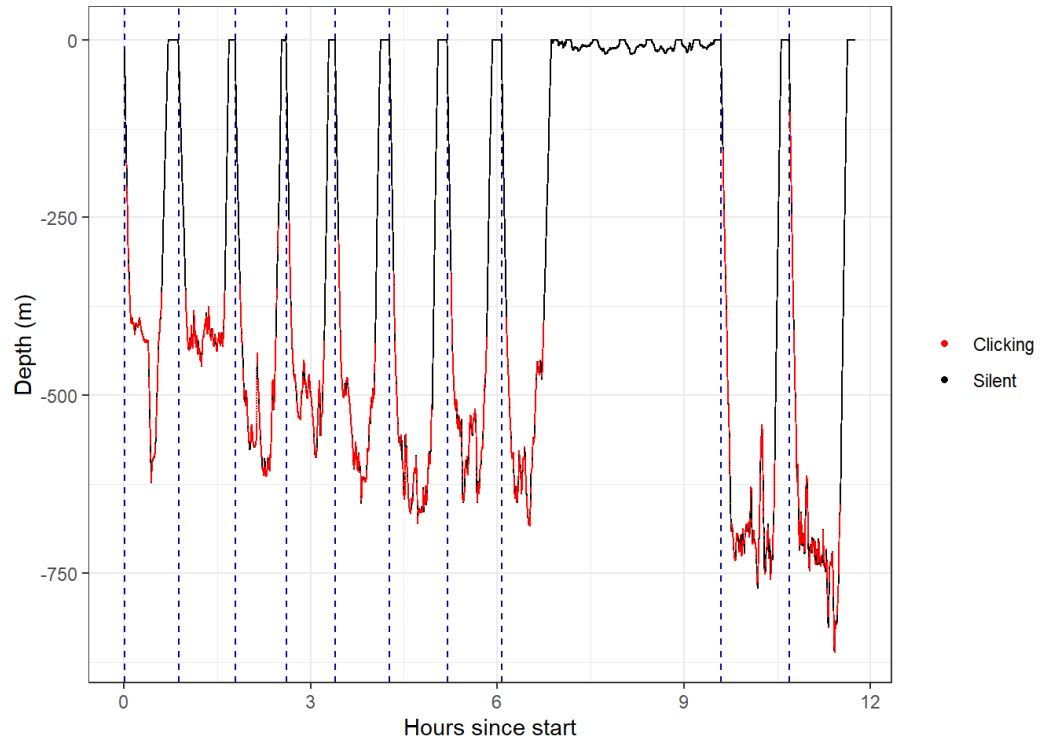
[a] tiago.marques@st-andrews.ac.uk

## 1. INTRODUCTION

14

15       Using the sounds produced by animals to estimate their abundance via

16 passive acoustic monitoring (PAM) is becoming increasingly popular for species that

17 are otherwise hard to detect visually. A prime example of such species are deep

18 diving cetaceans which spend prolonged periods at depth, making them hard to

19 survey visually. Therefore, it is not surprising that beaked whales Ziphiidae and

20 sperm whales *Physeter macrocephalus* densities have been estimated using PAM

21 methods (e.g. Barlow & Taylor, 2006, Lewis et al. 2007, Marques et al. 2009, Ward et

22 al. 2012).

23       Cue counting is an indirect PAM approach to estimate animal abundance,

24 where instead of counting the animals themselves, we count cues produced by the

25 animals. For PAM the cues are sounds of the species of interest, in the sperm whale

26 case usually those will be regular echolocation clicks. Cue counting was originally

27 developed in the 1980s within the realms of the IWC (International Whaling

28 Commission) for estimating baleen whale abundance from whale blows (e.g. Hiby &

29 Ward, 1996). If a cue is effectively instantaneous - as for whale blows, a short

30 duration sound or the onset of a long sound - then the only required multiplier to

31 convert an estimate of the density of cues into a density of animals is the cue

32 production rate. We define the cue production rate, or cue rate for short, as the

33 mean number of cues produced per animal per unit time. Naturally, one would like

34 to obtain a cue rate that is valid for the time when and place where the main survey

35 occurs (e.g. Marques et al. 2013). Otherwise, mismatches might potentially induce a

36  bias in cue rate, which will translate into a corresponding bias in the estimated animal

37  density.

38  Cue production rates reported in the literature for deep divers, namely

39  beaked whales, have considered deep dive cycles (DDCs) as a natural unit for

40  analysis. A DDC can be defined as the period corresponding from the time an

41  animal starts the descent for a deep foraging dive until the time it does the same for

42  the next DDC. DDCs are highlighted in an example sperm whale sound-and-

43  movement tag data from DTAG deployments (Johnson & Tyack, 2003) in Fig. 1.

44  DDCs might be more or less clearly defined units depending on a given species'

45  behaviour. Marques et al. (2009), considering Blainville's beaked whales (*Mesoplodon*

46  *densirostris*), estimated the cue rate from a weighted average of the number of cues per

47  unit time per DDC, where the weights were the durations of the DDCs. Warren et

48  al. (2017), working with DTAG data from both Blainville's beaked whales and

49  Cuvier's beaked whales (*Ziphius cavirostris*), considered a generalized estimation

50  equation (GEE) approach, using as response variable the number of clicks per DDC,

51  with DDC duration as an offset. Note this offset is equivalent to weighting by DDC

52  duration.

Fig. 1. An example sperm whale depth profile with the corresponding echolocation clicks overlaid. The 10 full DDCs available from this tag are highlighted. The eighth deep dive cycle is longer than the others, with the whale taking its time at the surface before submerging into the nineth deep dive.

We hypothesize that using the DDC as the unit for analysis when estimating cue rates might lead to biased inferences in the case where the DDC duration is correlated to the cue rate during a deep dive, which is likely the case by the very definition of a DDC. We compare different ways of calculating cue rates, and their associated precisions, from tag data, including averages and weighted averages, per DDC and per tag, and contrast these with regression modelling approaches to evaluate if previously used GEEs perform as expected. We illustrate the potential issues with a unique dataset of 104 sperm whale DTAGs and provide practical guidance for estimating acoustic cue rates from animal borne tags.

## 2. METHODS

We first describe the DTAG dataset used, then the methods used to extract the individual echolocation regular foraging clicks, considered the acoustic cue of interest, and finally the different analysis options to estimate cue rates from the acoustic data.

**A. Data collection**

We consider a sample of 104 DTAGs deployed on sperm whales, at 8 different sites and covering 13 different years. For additional details including counts of tags per year-site combination see the Supplementary Material. We focus on estimating a pooled cue rate for the species assuming the sample of tags would be representative for that purpose. We address potential issues in doing so in the discussion.

The DDCs were defined as periods starting at the moment the whale submerges into a deep foraging dive all the way till the next time it submerges for the subsequent deep foraging dive. For simplicity we considered deep foraging dives to be dives deeper than 100 meters. We note explicitly that a deep dive foraging cycle might include portions of time where the whale is at the surface not engaged in any deep foraging behaviour, say while resting or socializing at the surface (cf. eighth DDC in Fig. 1). This will be one of the reasons for why care must be taken when using such DDCs as sampling units.

## B. Data processing

For each tag, the sound files were processed to extract the times of emission for each regular echolocation click emitted by the tagged whale. Depending on the tags, custom built MATLAB functions to do so included either findclicks or findallclicks and findmissedclicks functions. These tools have been developed by Mark Johnson (freely available at: https://soundtags.wp.st-andrews.ac.uk/dtags/dtag-3/). For easier reference the functions are hosted also at https://github.com/TiagoAMarques/DeepDiverCueRates (folder: *click_extraction_matlab*).

To identify the timing of regular clicks from the tagged whale in the on-animal DTAG sound recording, a semi-supervised click detector was used, as described in Johnson et al. (2006). The sound files were processed sequentially in 15 second chunks through a supervised click detector to identify likely clicks from the tagged animal, using a 4-pole Butterworth band-pass filter (with cut-off frequencies at 3 and 20 kHz) and a level threshold based on the envelope of the click. An experienced analysist went through manual inspection of spectrograms (512 sample fast Fourier transform (FFT) with a Hamming window and 50% overlap) formed by 15s intervals of the sound recording (as described in Warren et al., 2017), accepting/rejecting the click identifications or, if needed, manually adding potential missed clicks. Clicks from the tagged animal were distinguished from those of other nearby whales in two ways (Johnson et al., 2006). Clicks from the tagged whale have both low-frequency energy that is absent in clicks recorded from non-tagged whales (Johnson et al., 2009) as well as a consistent angle of arrival on the tag, $\theta$ , computed from $\theta=\sin^{-1}(\tau c/d)$, where c

110   is the speed of sound in seawater, d is the hydrophone separation (0.025m) and $\tau$ is

111   the time delay between the two hydrophone signals, measured by cross-correlation.

112   The arrival angle of clicks from the tagged whale, when corrected for the tag

113   orientation on the whale, will be consistently close to zero, as the sound source from

114   the tagged animal is directly in front of the tag, while those from other whales will vary

115   widely as the focal and non-focal whales maneuver. The DTAG click extraction

116   process, which is a semi-supervised click identification ensures that all clicks produced

117   by the focal whale are identified.

118       For further analysis of our sperm whale dataset, we consider two datasets

119   derived from the above procedure: (1) The dataset with all the DDCs and (2) a

120   dataset corresponding to removing DDCs lasting over one hour. As an example, in

121   Figure1 this would correspond to exclude the 8$^{th}$ DDC. We refer to these as the

122   complete and the reduced datasets, respectively. These allow us to illustrate different

123   aspects of the analysis, as if these were two different datasets with slightly different

124   characteristics, the latter representing species with deep dive behavior similar to a

125   sperm whale, but without the long resting periods at the surface.

126   **C. Data analysis**

127       We focus on the estimation of a pooled (pooling across both years and

128   locations) cue production rate $r$ for sperm whales, defined as the number of cues

129   produced per time unit, per animal:

132
$$r = \frac{\#\ sounds}{time\ animal}. \qquad (1)$$

130   Note that in general we ignore, in wording and notation alike, the fact that this is a

131   measure per animal. Here we consider sounds to be regular echolocation clicks and

133 the time unit seconds, for convenience, but the above expression could be used for

134 any arbitrary time period, meaning in particular it could be calculated by DDC, by tag

135 or by any arbitrary time period (e.g. per 5 minutes). Then one can use those sampling

136 units to average across a sample to get a mean cue rate for the population of interest,

137 accounting for possible non-independence, as required. The population average cue

138 rate, using a standard mean based on DDCs, is then estimated by

139
$$\hat{r}^{ds} = \frac{\sum_{j=1}^{n_d} r_j^d}{n_d} = \frac{\sum_{i=1}^{n} \sum_{k=1}^{n_{di}} \frac{c_{ik}}{t_{ik}}}{n_d} , \qquad (2)$$

140 where the superscripts $d$ are used for deep $Dive$ and $s$ for $Standard$ mean, $r_j^d$

141 represents the cue rate for the $j^{th}$ dive cycle ($j=1,2,...,n_d$), $c_{ik}$ and $t_{ik}$ represent

142 respectively the number of clicks in, and the duration of, the $k^{th}$ DDC of the $i^{th}$ whale

143 ($k=1,2,...,n_{di}$), and $n_{di}$ is the number of deep dives recorded for whale i. On the other

144 hand, one could $Weight$ (note superscript $w$ below) for the DDC duration, as was

145 done in Marques et al. (2009), leading to

146
$$\hat{r}^{dw} = \frac{\sum_{i=1}^{n_d} \sum_{k=1}^{n_{di}} \frac{c_{ik}}{t_{ik}} t_{ik}}{\sum_{i=1}^{n_d} t_{ik}} = \frac{\sum_{i=1}^{n_d} c_{ik}}{\sum_{i=1}^{n_d} t_{ik}}. \qquad (3)$$

147 Note that this otherwise apparently more complex estimator (than the

148 standard average, given the weights) actually reduces to a simpler expression, the

149 total number of detected cues, across all tags, divided by the total recording time,

150 again across all tags. Both of the above consider the DDCs as the sampling unit, as

151 has been done before in the literature. We note that, strictly speaking, that is the

152 definition of pseudoreplication (Hurlbert, 1984), where the independence came in as

153 a stated assumption, for a proof of concept of PAM density estimation (DE) in

154 Marques et al. (2009). On the other hand, if we consider the $n$ tags as the sampling

9

155 units (i.e., the animals, superscript *a*), we have a standard average (superscript *s*)

156 estimator as

$$\hat{r}^{as} = \frac{\sum_{i=1}^{n} r_i^a}{n} = \frac{\sum_{i=1}^{n} \frac{c_i}{t_i}}{n} \qquad (4)$$

157

158 where $r_i^a$ represents the cue rate for the i[th] whale (i=1,2,...,n), $c_i$ and $t_i$

159 represent respectively the number of clicks in, and the duration of, the i[th] tag. The

160 corresponding weighted average (superscript *w*) version, now weighting by tag

161 duration, is

$$\hat{r}^{aw} = \frac{\sum_{i=1}^{n} \frac{c_i}{t_i} t_i}{\sum_{i=1}^{n} t_i} = \frac{\sum_{i=1}^{n} c_i}{\sum_{i=1}^{n} t_i}. \qquad (5)$$

162

163 For each of the approaches we also estimate the corresponding precision and

164 95% confidence intervals. The variance of a standard mean is straightforward and

165 present in any introductory statistics book, and the variance for the weighted mean

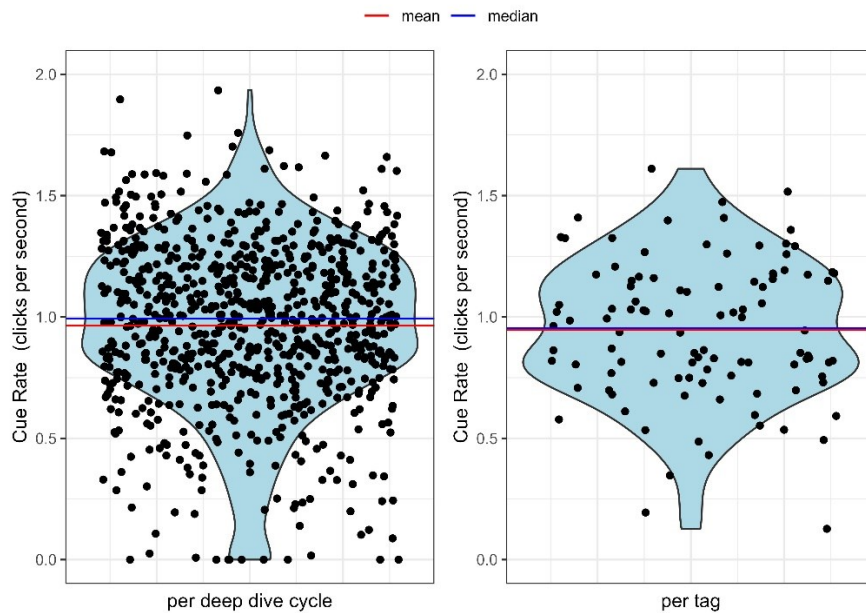166 was calculated considering the methods proposed by Gatz & Smith (1995).

167 For a direct comparison with the methods from Warren et al. (2017), we also

168 estimate the cue rate based on regression models. For these, the unit of analysis

169 considered was always the DDC. We consider the same GEE inspired approach:

170 modelling the number of clicks per DDC, using the DDC duration as an offset,

171 considering tag as the grouping variable (akin to a random effect), with a log link

172 function, an "independence" correlation matrix and robust standard errors. The

173 GEE treats the deep dives as the sampling unit but inflates standard errors on the

174 estimated cue rate via the correlation structure induced by the grouping variable.

175 Therefore, in terms of how the correlation structure is dealt with, it represents a half-

176 way house between treating deep dives or tags as the sampling unit, while

177 considering a regression model instead of an analytical formula for the average. The

178 intercept of this model corresponds to an estimate of the cue rate. To obtain the

179 standard error on the response scale we used a delta method approximation (Powell,

180 2007). To obtain 95% confidence intervals we assumed normality on the link scale

181 and back-transformed to the response scale. For comparison we also implement the

182 same GEE analysis without the offset and without both offset and grouping variable,

183 and use corresponding generalized linear mixed models (GLMM) with whale as a

184 random effect, instead of the GEEs. When not considering the offset we modelled

185 the rate directly and considered a Gamma distribution for the response. For 9 DDCs

186 (corresponding to 2% of the DDCs) there were 0 clicks. To avoid issues with the

187 Gamma not coping with the response variable being exactly zero we replaced these

188 observations by 0.5 (or $1/20^{th}$ of the observed minimum positive count of 10 clicks

189 per DDC). The practical impact of this tweak is negligible, but fitting with the

190 Gamma family becomes possible. Good reviews on GLMMs and GEEs in Ecology

191 are Bolker et al. (2009) and Pekár & Brabec (2017), respectively.

192     Analysis was implemented using R (R Core Team, 2022). The GEE model

193 was implemented using the geeglm function in the geepack package (Højsgaard et al.

194 2006) and the GLMM model via glmer in the R package lme4 (Bates et al. 2015). All

195 of the code to reproduce the statistics and figures in the paper is provided as

196 supplementary pdf file. This pdf is generated via an RMarkdown dynamic report. All

197 the data and original .Rmd file that allows one to reproduce or update the analysis
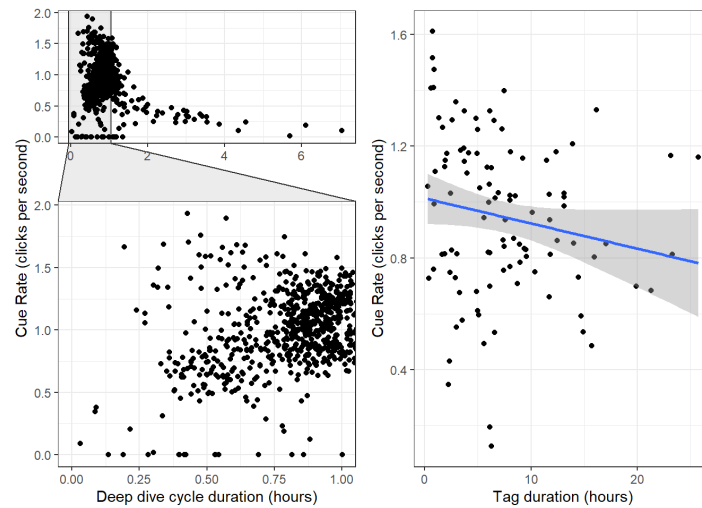
198 are shared as a github repository:

199 https://github.com/TiagoAMarques/DeepDiverCueRates.

## 3. RESULTS

We had 826 DDCs across a total of 104 tag deployments from sperm whales, with a median number of 7 DDC per tag, ranging from 1 to 31 DDC per tag. Tag durations ranged from 0.3 to 25.8 hours. The observed cue rates per tag varied between 0.13 and 1.61 clicks per second, with a median value of 0.86, while in the case of DDC these ranged between 0 and 1.93, with a median value of 0.93. (Figure 2). In the case of our dataset, the cue rate per DDC tended to increase with DDC duration for the reduced dataset, but in the full dataset, the longer DDC presented long periods without vocalizations, and hence cue rates tended to decrease with DDC duration (Figure 3).



Fig. 2. Violin plots of the observed cue rates per deep dive cycle (left) and per tag (right). The x-coordinate is a non-interpretable jitter for improved data visualization. There are 826 points in the left plot corresponding to the cue rates for

214  all the deep dive cycles available for the 104 tag deployments on whales in the right

215  plot.



216

217

218      Fig. 3. Observed cue rates as a function of duration, both for deep dive

219  cycles (DDCs, left) and tags (right). The left bottom panel zooms in on DDCs

220  shorter than 1 hour. On the right plot a regression line is represented, with the grey

221  lines representing the corresponding 95% confidence interval on the regression line.
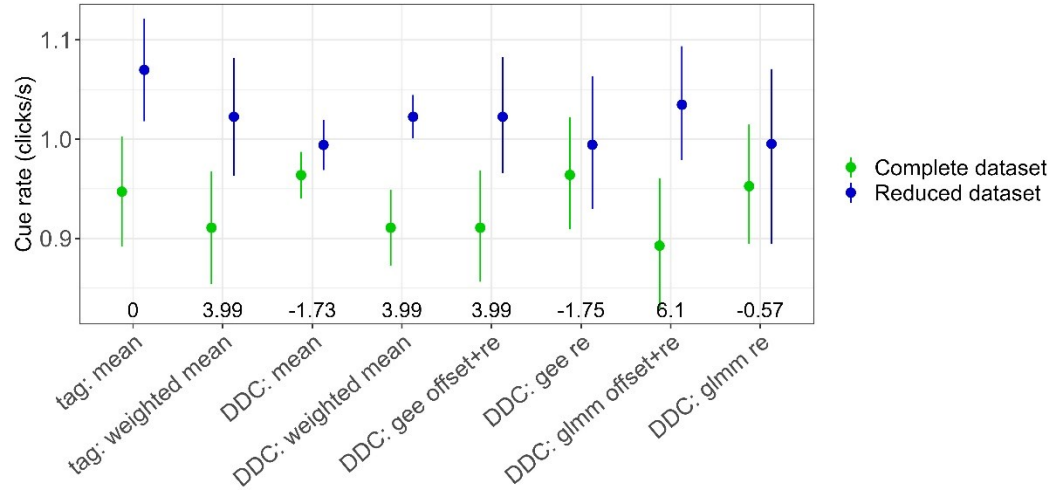
222

223      The results obtained for each of the analytic estimators and the regression

224  models are compared in Fig. 4. As expected, for any given method, the estimates

225  based on the reduced dataset were always higher than those obtained for the

226  complete dataset. This is a rather obvious consequence of the fact that, by

227  construction, the longer DDCs tend to include periods when the animal spent

228  considerable time at the surface, when they are mostly silent. The 8th DDC in Fig. 1

229  is a good example.

230        The sperm whale cue rate, considering an unweighted average based on the

231        full dataset at the tag level, was 0.947 clicks per second, with a 95% confidence

232        interval (CI) of 0.892,1.003. The same tag-level analysis, but weighted by tag

233        duration, estimated a lower cue rate: 0.911 clicks per second (95% CI 0.854, 0.968).

234        Corresponding unweighted estimates at the DDC level were lower in the reduced

235        dataset (the difference in green values from the 2nd analysis to the 3rd analysis in Fig.

236        4), but higher for the complete dataset (the difference in blue values from the 2nd

237        analysis to the 3rd analysis in Fig. 4). When weighed by DDC, the effect is opposite in

238        the reduced and the complete datasets, with the cue rate increasing in the former, but

239        decreasing in the latter. This is the same pattern observed when adding an offset to

240        the GEE with whale as grouping variable.

241        The analyses weighted by duration, either at the tag or DDC level, lead

242        necessarily to the same cue rate point estimate, with a marked difference in precision,

243        naturally higher for the DDC analysis.

244        Similar estimates are obtained for the GEE and the corresponding GLMM

245        counterparts.

246

Fig. 4 Estimated cue rates, and corresponding 95% confidence intervals, considering the different approaches. Left to right, the first two analyses are at the tag level, the remaining ones at the deep dive cycle (DDC) level. For the regression analysis, offset means the duration of the DDC was included as an offset and "re" denotes if the model included whale as a grouping variable (GEE) or as a random effect (GLMM). The number shown for each approach corresponds to the bias in animal density that a biased cue rate would induce, assuming as the truth for cue rate the unweighted average at the tag level, unbiased by design.

**4. DISCUSSION**

We presented pooled cue rate estimates for sperm whales across all the tag recordings available to us, for the purpose of comparing several methods to estimate cue rates, including averages and weighted averages, and a variety of regression models. Variability over time and space means that such overall mean cue rates might not be representative of any one place at any one time, and therefore we recommend these values are not used to inform any sperm whale PAM density

263    estimation exercises. This is the topic of a separate research thread we are currently

264    pursuing.

265        To inform a PAM density survey, we are interested in estimating a cue rate

266    for a population of whales. Therefore, the true variability that we are interested in is

267    the variability across whales. For that reason alone, one should expect that

268    approaches that consider the whale (i.e. tag record) as the sampling unit would be

269    preferable. Perhaps surprisingly, that was not considered by earlier attempts (e.g.

270    Marques et al. 2009, Warren et al. 2017). Here we consider the standard average at

271    the tag level as a gold standard, since that should be unbiased by design.

272    **A. To weight or not to weight: that is (and might remain) the question**

273        The weighted average at the tag level was lower than the unweighted average,

274    a reflection of the fact that longer tag records had slightly lower cue rates (cf. Fig. 3,

275    right panel). There is considerable overlap in the 95% CI between the weighted and

276    unweighted average at the tag level, and there is no obvious reason for why tags with

277    longer duration (hence, everything else being the same, also on average with a larger

278    number of DDCs) might have lower cue rates. We put forward a couple of possible

279    explanations. The first is that deep diving itself could represent a behavior that

280    promotes tag displacement and/or release. In such a case, animals spending longer

281    times at the surface and hence for which the tag would be more likely to stay on

282    longer would be oversampled, and hence the true cue rate would be biased low.

283    However, this is unlikely, and in fact one might even argue the bias would be the

284    other way around, with high pressure experienced by the tag during deep dives

285    meaning the suction cups would be less prone to displace. The second is that animals

286    available for tagging might be preferentially in a vocal mode (say because they are

287    found acoustically), and hence they might be, more often than not, engaged in

288    behaviors that have higher cue rates than the average animal in the population. On

289    the other hand, Warren et al. (2020) noted, using a subset of the tags used here, an

290    opposite pattern: if anything a time decaying effect detectable up to the 5th deep dive

291    with less buzz sounds (but no results are reported for echolocation clicks, the focus

292    here). There is no strong reason to prefer one estimate over the other; weighted

293    averages can have higher precision if the weights are sensible. But if they are not,

294    bias could creep in. Unless there is an alternative explanation that would not induce

295    bias for why longer tag deployments might have lower cue rates, we suggest that an

296    unweighted average at the tag level might be preferred. The weighted by DDC

297    approach might introduce bias, or in this case perhaps differences in the estimates

298    which might be driven by randomness rather than signal.

299        The weighted averages lead to identical point estimates, irrespective of

300    whether we consider an analysis at the tag level or at the DDC level (cf. equations 3

301    and 5), but the precision is higher at the DDC level. This will always be the case, the

302    analysis at the DDC level is strictly pseudoreplication, but naturally only the analysis

303    at the tag level is sensible. When Marques et al. (2009) estimated a cue rate to inform

304    a passive acoustic density estimation exercise for the first time, only 5 tags were

305    available, and hence considering the DDC as the independent sampling unit was an

306    attractive choice. That choice, and implicit assumption about independence across

307    DDCs, meant that instead of a 5-tag sample size, the authors considered a sample

308    size of 21 DDCs. But as with most assumptions - there are no free lunches in

309    statistics - that assumption comes at a cost. Here illustrated for the case of the sperm

310    whales, that cost is underestimating the true precision. In other words, we

311    underestimate the variance of the cue rate estimate, and therefore we would

312    underestimate the true variance on the corresponding estimated density.

313         For our sperm whale example, considering the unweighted average at the

314    DDC overestimates the cue rate. This is a consequence of most DDCs being below

315    1 hour, and therefore, having higher cue rates than all the DDCs together. This is

316    clearly evidenced when one looks at the contrast in behavior between the reduced

317    and complete datasets at the DDC level. Cue rate increases in the complete dataset,

318    but it decreases in the reduced dataset. When you compare just between DDC

319    averages, naturally you obtain a higher value for the weighted average in the reduced

320    dataset. Again, this happens because longer DDCs, with lower cue rates, were

321    removed. On the other hand, you obtain a lower value for the DDC weighted

322    average in the full dataset, since in that dataset there are a considerable large number

323    of longer DDCs, typically with lower cue rates, being then given larger weights.

324    There is a correlation between the cue rate by DDC and the DDC duration, with the

325    largest DDCs being associated with the lowest cue rates. This happens because these

326    correspond to DDCs where whales might spend a considerable amount of time at

327    the surface. In fact, for these instances, the definition we considered for a DDC can

328    be misleading: for long DDCs the whale might spend much more time doing

329    something else, like resting or socializing, than actually performing the deep foraging

330    dive than coins the DDC term used to define the period. Hence, when we weight by

331    DDC duration, compared to an unweighted tag analysis, we underestimate the cue

332    rate by about $100*(0.947-0.911)/0.947 = 3.84\%$, which would correspond to, all else

333    being equal, an upward bias in density of $100*(1/0.911-1/0.947)/(1/0.947)=3.99\%$.

334     While 4% might be a relatively small bias given the usual coefficient of variation of

335     abundance/density estimates, any bias that can be removed necessarily leads to

336     improved inferences. Additionally, this correlation between cue rate and the weights

337     (DDC duration) means that properly calculating the variance of the corresponding

338     weighted mean would require incorporating the covariance between the observations

339     and the weights, which is far from straightforward.

340        The decision of weighing or not by recording duration is unfortunately more

341     nuanced than one might hope. One can easily construct a scenario where such

342     weighting would be desirable. An example is when shorter duration tags do not

343     provide a reliable cue rate for the tagged animal. On the other hand, when all tags are

344     long enough to obtain a reliable individual cue rate per tag, weighting could induce

345     bias towards some animals with unusually long records. The decision will depend on

346     how variable animals are over time compared to the variability between animals.

347     Weighting becomes more relevant as variability within animals increases and across

348     animals decreases, but where to change from a standard average to a weighted

349     average given said ratio of variabilities and average tag duration might be a hard

350     question to answer.

351     **B. About regression models**

352        The analysis considering (1) the GEE regression model and (2) the weighted

353     average, considering DDC as the sampling unit, estimate the same quantity, and

354     hence we get the same point estimates for either dataset. We do note that the GEE

355     confidence intervals are wider than those for the weighted mean, reflecting lower

356 precision in the estimated means. This is a more sensible precision, since tag

357 deployments are the independent sampling units, not DDCs.

358     The pattern found in Fig. 3 illustrates that, while required given the definition

359 of what a cue rate is (cues per unit time), from a regression modelling perspective,

360 DDC duration might not be a sensible offset, since the relationship between it and

361 the cue rate is far from being proportional. In other words, an estimated coefficient

362 for the offset included in the model as a variable would not be 1, which is strictly

363 what an offset corresponds to. This is a reminder that use of offsets in regression

364 corresponds to an implicit, often unstated, assumption. Since this is an assumption

365 that is easy to test, by plotting the data as we did, we suggest in general should be

366 tested when using an offset.

367     GEEs and GLMMs model different conceptual quantities. GEEs are often

368 referred to as marginal models, and provide population level averages, while

369 GLMMs, also known as conditional models, will provide individual level averages

370 (Fieberg et al. 2009). This is often ignored, especially if the interest of inferences is

371 on how covariates influence a response. However, in the case of cue rates from tag

372 data, the distinction might be crucial. We are actually interested on  the mean

373 estimated by the GLMM, a mean across whales, not by that of the GEE as

374 implemented, a mean across DDCs. At least for our example, differences in point

375 estimates across the two approaches were minimal, with slightly higher variances

376 estimated via the GLMM, but that might not be the case in general for other species

377 that might have a different diving and sound production behaviour.

378     One might wonder why bother with regression models if analytical

379 expressions of averages provide such similar results for mean cue rates. The power

380  of regression approaches truly emerges when additional covariates that cue rates

381  might depend on are available. That opens the door to model-based estimates of cue

382  rate, that can be predicted for the actual survey conditions. As an hypothetical

383  example, one could imagine cue rate being dependent on survey level covariates, like

384  study area bottom depth, because animals spend less or more time silent travelling to

385  the bottom to feed depending on bottom depth. Then a regression model allows one

386  to estimate the cue rate for the depths at which the survey sensors were placed.  One

387  might additionally be able to model cue rates as a function of animal level covariates

388  (e.g. sex), or even covariates that change over time within animal (e.g. animal depth).

389  While this might provide interesting biological information, for a PAM density

390  estimate we will require an average cue rate, which will have to be averaged over the

391  survey conditions. Hence, knowing that cue rate differs by sex or by the depth at

392  which an animal is not enough and might be of little use in practice: To use that

393  information we would need to know, for all the animals within the survey area, the

394  sex of the animals or the depths at which they were diving, to average across those

395  distributions to obtain the correct multiplier. A pragmatic approach in such cases,

396  which comes at the cost of a strong untestable assumption, is to assume that the

397  sample of animals we have tagged provides an unbiased sample over which, once

398  averaged across, we can estimate the corresponding average cue rate.

399  **C. Cue rates to inform PAM surveys**

400       As a conclusion, we highlight what we knew from the start. We want a cue

401  rate estimate to convert a cue density per time into an animal density. Therefore, we

402  need the average cue rate that applies for the time and place the survey took place.

403 Fundamentally, the variability in this cue rate estimate that we are interested is that

404 across animals, hopefully obtained from a random sample of animals. Therefore, an

405 analysis that focuses on the DDCs, where the duration of the DDCs is correlated

406 with the variable of interest, here cue rate, might be biased. We recommend that

407 researchers calculating cue rates from similar tag data take due care to consider

408 analysis that reflect the variability at the level of the individual whales. We encourage

409 researchers to be careful when considering inferences where DDCs might be

410 tempting to use as natural sampling units.

411     Cue rates are a required multiplier for cue counting approaches to estimate

412 animal density and abundance. Reliable methods to estimate the cue rate and its

413 corresponding variability are needed. For the sake of this paper, with the objective of

414 evaluating potential bias induced by the methods used, we pooled all the data and

415 focused on a pooled cue rate across space and time, as if the samples were a suitable

416 random sample for that purpose. Nonetheless, the sample was not balanced in either

417 time or space, so it might be biased for any given time or place. Even for a given

418 time and place we have tags for, one should consider carefully whether a sample of

419 tagged animals is, in general, a representative sample of the animals available. With a

420 small sample might be likely to get a few STRANGE animals (a la Webster and Rutz,

421 2020) that could compromise inferences. As an example, if animals to be tagged are

422 found acoustically, implying they would be in a "vocally active" mode, cue rates of

423 these animals could be potentially higher than the cue rates of other animals, biasing

424 cue rate up and correspondingly density estimates low. Obtaining a cue rate for a

425 new location or time period should require a good understanding about a species'

426 cue production and implications of potential spatio-temporal differences in cue rates.

427   While we considered here a pooled mean across all tags, collected at different times

428   and different places, we know ultimately a cue rate might be affected by a multitude

429   of factors, be it season, region, demography, etc. That is the objective of current

430   research and is unlikely to lead to an answer that fits all questions.

431       We did not consider any additional covariates in the regression models. A

432   model-based approach where cue rate can be predicted for the time and place one's

433   survey was conducted, conditional on observed covariates, might be desirable. For

434   cues like echolocation clicks the cue rate variability and factors affecting it might be

435   tamed, and a reasonable value obtained in such a way. At the other extreme, say for

436   social sounds, said variability might preclude obtaining estimates with acceptable

437   precision. In such cases, or when cue rate is density dependent, the only option to

438   use a cue counting approach for density estimation might be to use a sample of

439   tagged animals to estimate the cue rate for the place and time the survey took place.

440   This remains a considerable drawback in estimating animal abundance from passive

441   acoustic data via cue counting, and further work is required to understand drivers of

442   cue rate variability and to identify for which species and cues these might be stable

443   enough to lead to reliable estimates of abundance.

444

445   **SUPPLEMENTARY MATERIAL**

446       See supplementary material at [URL will be inserted by AIP] for a dynamic

447   report that reproduces all the analysis and outputs statistics and figures on the paper

448   via RMarkdown. The data and the dynamic report are hosted at a github repository:

449     https://github.com/TiagoAMarques/DeepDiverCueRates/. We also host there the

450     custom-built MATLAB functions used to identify the clicks of the tagged whales.

451

## ACKNOWLEDGMENTS

## AUTHOR DECLARATIONS

467     The authors have no conflicts of interest to declare. The data collection was

468     carried out by the researchers contributing with their tags to this joint effort, all of

469     them operating under strict ethical permits. This manuscript was produced within

470 ACCURATE, and the project got a secondary data Ethics Approval from the School

471 of Biology Ethics Committee at the University of St Andrews.

472 **REFERENCES**

473 Barlow, J. & Taylor, B. L. 2006 Estimates of sperm whale abundance in the

474 northeastern temperate Pacific from a combined acoustic and visual

475 survey Marine Mammal Science 21, 429-445

476 Bates, D., Maechler, M., Bolker, B. & Walker, S. (2015). Fitting linear mixed-

477 effects models using lme4. Journal of Statistical Software, 67: 1-48

478 Bolker, B. M.; Brooks, M. E.; Clark, C. J.; Geange, S. W.; Poulsen, J. R.; Stevens,

479 M. H. H. & White, J.-S. S. 2009 Generalized linear mixed models: a

480 practical guide for ecology and evolution Trends in Ecology & Evolution

481 24, 127 – 135

482 Fieberg, J.; Rieger, R. H.; Zicus, M.C. & Schildcrout, J. S. 2009 Regression

483 modelling of correlated data in ecology: subject-specific and population

484 averaged response patterns Journal of Applied Ecology 46, 1018-1025

485 Gatz, D. F. & Smith, L. 1995 The standard error of a weighted mean

486 concentration--I. Bootstrapping vs other methods Atmospheric

487 Environment, 29, 1185-1193

488 Hiby, A. R. & Ward, A. J. 1986 Analysis of cue-counting and blow rate

489 estimation experiments carried out during the 1984/85 IDCR minke whale

490 assessment cruise Report of the International Whaling Commission 36,

491 473-476

492    Højsgaard, S., Halekoh, U. & Yan J. (2006) The R Package geepack for

493         Generalized Estimating Equations Journal of Statistical Software, 15: 1--

494         11

495    Hurlbert, Stuart H. 1984 Pseudoreplication and the design of ecological field

496         experiments Ecological Monographs 54, 187-211

497    Johnson, M. P. and Tyack, P. L. (2003) A digital acoustic recording tag for

498         measuring the response of wild marine mammals to sound IEEE Journal

499         of Oceanic Engineering 28,  3-12

500    Johnson, M., Madsen, P. T., Zimmer, W. M. X., De Soto, N. A., & Tyack, P. L.

501         (2006). Foraging Blainville's beaked whales (*Mesoplodon densirostris*)

502         produce distinct click types matched to different phases of

503         echolocation. Journal of Experimental Biology, 209(24), 5038-5050

504    Lewis, T., Gillespie, D.,  Lacey, C.,  Matthews, J.,  Danbolt, M.,  Leaper, R.,

505         McLanaghan, R. &  Moscrop, A. 2007 Sperm whale abundance estimates

506         from acoustic surveys of the Ionian Sea and Straits of Sicily in 2003

507         Journal of the Marine Biological Association of the United Kingdom 87,

508         353-357

509    Marques, T. A.; Thomas, L.; Ward, J.; DiMarzio, N. & Tyack, P. L. (2009)

510         Estimating cetacean population density using fixed passive acoustic

511         sensors: an example with Blainville's beaked whales The Journal of the

512         Acoustical Society of America, 125, 1982-1994

513    Marques, T. A.; Thomas, L.; Martin, S. W.; Mellinger, D. K.; Ward, J. A.;

514         Moretti, D. J.; Harris, D. & Tyack, P. L. (2013) Estimating animal

515       population density using passive acoustics Biological Reviews, 88, 287-

516       309

517    Pekár, S. & Brabec, M. 2017 Generalized estimating equations: A pragmatic and

518       flexible approach to the marginal GLM modelling of correlated data in the

519       behavioural sciences Ethology 124, 86-93

520    Powell, L. A. 2007 Approximating variance of demographic parameters using the

521       delta method: a reference for avian biologists. The Condor, 109, 949-954

522    Ward, J. A.; Thomas, L.; Jarvis, S.; DiMarzio, N.; Moretti, D.; Marques, T. A.;

523       Dunn, C.; Claridge, D.; Hartvig, E. & Tyack, P. (2012) Passive acoustic

524       density estimation of sperm whales in the Tongue of the Ocean, Bahamas

525       Marine Mammal Science, 28, E444-E455

526    Warren, V. E.; Marques, T. A.; Harris, D.; Tyack, P. L.; Thomas, L.; de Soto, N.

527       A.; Hickmott, L. & Johnson, M. P. (2017) Spatio-temporal variation in

528       click production rates of beaked whales: implications for passive acoustic

529       density estimation The Journal of the Acoustical Society of America, 141,

530       1962-1974

531    Warren, V.E., Miller, P.J. and Tyack, P.L. 2020. Short-term responses of sperm

532       whales *Physeter macrocephalus* to the attachment of suction cup

533       tags. Marine Ecology Progress Series 645, 219-234.

534    Webster, Michael M., and Christian Rutz. How STRANGE are your study

535       animals? Nature 582, 337-340

536