

1 **Counting with DNA in metabarcoding studies: how should we convert sequence reads to**
2 **dietary data?**

3

4 Bruce E. Deagle* (corresponding author: Bruce.Deagle@aad.gov.au)

5 Austen C. Thomast†

6 Julie C. McInnes*

7 Laurence J. Clarke‡*

8 Eero J. Vesterinen¹¶

9 Elizabeth L. Clare²

10 Tyler R. Kartzinel³

11 J. Paige Eveson§

12

13 *Australian Antarctic Division, Channel Highway, Kingston, Tasmania, Australia

14 †Science Department, Smith-Root Inc., Vancouver, Washington, USA

15 ‡ Antarctic Climate & Ecosystems Cooperative Research Centre, University of
16 Tasmania, Tasmania, Australia

17 ¹ Biodiversity Unit and Department of Biology, University of Turku, Turku, Finland

18 ¶ Department of Agricultural Sciences, University of Helsinki, Helsinki, Finland

19 ² School of Biological and Chemical Sciences, Queen Mary University of London, London, UK

20 ³ Department of Ecology and Evolutionary Biology, Brown University, Providence, Rhode
21 Island, USA

22 §CSIRO Oceans and Atmosphere, GPO Box 1538, Hobart, Tasmania, Australia

23

24 **Abstract**

25 Advances in DNA sequencing technology have revolutionised the field of molecular
26 analysis of trophic interactions and it is now possible to recover counts of food DNA barcode
27 sequences from a wide range of dietary samples. But what do these counts mean? To obtain
28 an accurate estimate of a consumer's diet should we work strictly with datasets
29 summarising the frequency of occurrence of different food taxa, or is it possible to use the
30 relative number of sequences? Both approaches are applied in the dietary metabarcoding
31 literature, but occurrence data is often promoted as a more conservative and reliable option
32 due to taxa-specific biases in recovery of sequences. Here, we point out that diet summaries
33 based on occurrence data overestimate the importance of food consumed in small
34 quantities (potentially including low-level contaminants) and are sensitive to the count
35 threshold used to define an occurrence. Our simulations indicate that even with recovery
36 biases incorporated, using relative read abundance (RRA) information can provide a more
37 accurate view of population-level diet in many scenarios. The ideas presented here highlight
38 the need to consider all sources of bias and to justify the methods used to interpret count
39 data in dietary metabarcoding studies. We encourage researchers to continue to addressing
40 methodological challenges, and acknowledge unanswered questions to help spur future
41 investigations in this rapidly developing area of research.

42 **1. Introduction**

43 Many recent studies documenting trophic interactions make use of metabarcoding,
44 an approach which combines high-throughput sequencing (HTS) with DNA barcoding to
45 characterise organisms in complex mixtures (Nielsen *et al.* 2017). When HTS first became
46 available the potential applications in diet studies were clear and the methods were quickly
47 embraced by the community (Deagle *et al.* 2009; Valentini *et al.* 2009). In a comprehensive
48 review of DNA-based diet analysis by King *et al.* (2008) the possibility of using HTS was only
49 briefly mentioned as a 'Future Direction', and just four years later another review paper
50 focussed entirely on this approach (Pompanon *et al.* 2012). While many underlying technical
51 and biological details vary between dietary metabarcoding studies, the general workflow is
52 now well defined. It involves extraction of DNA from faecal samples or stomach contents,
53 PCR amplification of DNA barcode markers from food taxa of interest, and then DNA
54 sequencing for taxonomic classification of the recovered sequences. The workflow has been
55 applied to determine diet in a range of animals, from invertebrates to large mammalian
56 herbivores and carnivores (representative studies summarised in Table 1).

57 The rapid adoption of HTS to characterise complex mixtures of DNA is not unique to
58 dietary studies; over the last ten years the technology has produced a wealth of new genetic
59 data providing insight into almost all areas of biology (Goodwin *et al.* 2016). One feature of
60 HTS is that it provides counts of DNA sequences in each sample and therefore it has the
61 potential not only to provide a qualitative list, but also to quantify what DNA is present. The
62 interpretation of sequence read counts as a numerical representation of sample
63 composition is common in many HTS applications. For example, studies sequencing
64 transcripts to determine differences in gene expression (Finotello & Di Camillo 2015),
65 profiling microbe communities (Vandeputte *et al.* 2017) or measuring epigenetic variation
66 (Schield *et al.* 2016) all rely on sequence read counts. However, decisions about how to
67 interpret read counts is certainly not routine and the validity of interpretations is sometimes
68 questioned even in fields where the practice is well established (e.g. Edgar 2017; Olova *et al.*
69 2017). These debates are constructive, and should motivate researchers to test the
70 underlying assumptions and justify their interpretations, but can inadvertently give rise to
71 the false impression that count data are always misleading.

72 The reality is that metabarcoding studies always use sequence counts to some
73 extent. In dietary investigations, count data are used either to record the occurrence of food
74 species within samples based on a threshold number of sequences (i.e. presence/absence of
75 taxa), or to calculate the percentage of DNA belonging to each food species as a proxy for
76 relative biomass consumed (i.e. relative abundance of taxa; Figure 1). The conversion of
77 sequence counts to occurrence data is often considered a more conservative approach than
78 using proportional data. In their introduction to the Molecular Ecology Special Issue on
79 ‘Molecular Detection of Trophic Interactions’, Symondson & Harwood (2014) pointed out
80 that authors of many metabarcoding papers “*now simply record numbers of predators*
81 *testing positive for a target prey or plant species, providing a pragmatic and useful surrogate*
82 *for truly quantitative information*”. This sentiment, that focusing only on occurrence data is
83 a reliable and safe option, is now common in the literature and this step in the analysis
84 pipeline is often uncritically applied as the default option. Using counts as an indication of
85 biomass in the sample is more controversial. Indeed, the difficulties of obtaining an accurate
86 biomass signature from sequence counts include both technical and biological biases that
87 affect barcode marker recovery rates from different taxa (Amend *et al.* 2010; Deagle *et al.*
88 2009; Pompanon *et al.* 2012). Therefore in the best-case scenario sequence read counts can
89 only provide a rough estimate of proportional abundance. Still, to accept the notion that
90 relative sequence counts provide no meaningful information would mean that, within one
91 sample, a few DNA sequences from one food taxon is equivalent to 10,000 sequences from
92 another. Most molecular ecologists would interpret these disparate counts to mean that
93 there are differences in template DNA abundance (as long as methods used to collect the
94 data are reasonable) and that there is some biological basis for that difference. Ignoring this
95 difference may inhibit ecological understanding.

96 Here, we review the approaches taken to interpret sequence count data in dietary
97 metabarcoding studies and consider their implications. We point out that converting
98 sequence read counts to occurrence information can introduce strong biases and thus we
99 suggest it is not always a “conservative” approach. We also assess the scale of biases in
100 recovery of sequences from different food taxa in study systems where it has been
101 examined. Using simulations we explore the impact of these biases on data summaries
102 (both based on occurrence and read counts). In this light, we evaluate factors that impact
103 data summaries in dietary metabarcoding and consider where using sequence count data as

104 an indication of relative biomass within samples might be justified to provide a more
105 nuanced picture of animal diet.

106 The issues we consider on how best to summarise dietary data have implications for
107 all metabarcoding studies (Taberlet *et al.* 2018) and similar issues have been considered
108 extensively in traditional diet studies (e.g. Barrett *et al.* 2007; Laake *et al.* 2002). In HTS-
109 based diet studies the ideas are most relevant when the underlying objective is to estimate
110 the true diet of a particular consumer (i.e. the relative biomass contributions of alternative
111 diet species). This may not be a clearly stated goal, but is often implicit in outcomes of
112 dietary metabarcoding studies. Approaches for summarising sequence counts may be of less
113 concern in studies aiming to provide a list of taxa consumed by a particular species (niche
114 breadth), a summary of trophic interactions in a food web, or an indicator of dietary
115 differences between sites. Throughout the paper we will refer to the two general
116 approaches of summarising sequence count data as ‘occurrence’ (i.e. presence/absence of
117 taxa) and ‘relative read abundance’ (RRA; i.e. proportional summaries of counts). We focus
118 mainly on dietary studies using DNA extracted from faecal material. The use of HTS to
119 identify food in stomach contents is common in invertebrates, and also fish, but the
120 material recovered is in various states of digestion and the sequence counts are less likely to
121 contain a meaningful quantitative signal compared to the more consistent signal seen in
122 faecal material (Deagle *et al.* 2013; Nakahara *et al.* 2015).

123

124 **2. Current Practice**

125 Non-dietary metabarcoding studies use a range of approaches to interpret sequence
126 count data, and these vary depending on the targeted organisms. Recent papers published
127 in *Molecular Ecology* on bacterial/archaeal communities all make use of RRA, although half
128 of these studies also presented summaries based on taxon occurrences (Table S1). There is
129 widespread acknowledgement of taxon-specific biases in recovery of the bacterial/archaeal
130 barcode markers, but RRA is accepted as a flawed, but useful, measure of these diverse
131 communities that cannot be easily characterized by other means (Forney *et al.* 2004;
132 Ibarbalz *et al.* 2014). There is no clear consensus in metabarcoding of eukaryotic
133 communities: RRA is sometimes used exclusively (often the case in studies of fungi),

134 whereas metazoan studies use either occurrence data only or both metrics in tandem
135 (recent examples listed in Table S1).

136 In dietary metabarcoding studies, it is common to only interpret sequence data after
137 conversion to taxon occurrences (representative studies summarised in Table 1). This
138 conversion is done in various ways. During initial processing of sequence reads, most
139 researchers discard rare sequences to avoid incorporation of background sequencing errors
140 (e.g. Quéméré *et al.* 2013). After this a summary table of remaining sequence reads in each
141 sample is produced and sequences are assigned taxonomy (often with similar sequences
142 being clustered). Then, when converting these read counts to occurrence data, a threshold
143 number of reads is often required for each taxon to be tallied as an occurrence. Sequencing
144 depth can vary considerably between samples, so in practice a threshold percentage of
145 reads is often used (e.g. 1% of food sequences McInnes *et al.* 2017b), or sequencing depth
146 can be rarefied to a common level (O'Rorke *et al.* 2016). These approaches normalize
147 detection across samples, so that more sequences are required for an occurrence to be
148 recorded in samples with higher read depths.

149 Once occurrences are recorded in individual samples, several metrics can be used to
150 summarise the diet across samples. Those considered here are percent frequency of
151 occurrence (%FOO), percent of occurrence (POO) and weighted percent of occurrence
152 (wPOO) (Figure 1; see Box 1 for details).

153 Some dietary metabarcoding studies present RRA data along with occurrence
154 summaries, although relatively few have relied solely on information obtained from RRA
155 (Table 1). In almost all of these studies, the number of sequences obtained per sample are
156 converted to percentages (Figure 1a), because the absolute counts are dependent on
157 several factors unrelated to the overall importance of the sample (amount of starting
158 material used, DNA extraction efficiency, standardization of samples before HTS, etc.). To
159 provide an overall diet summary, sample-specific RRA values can be averaged across
160 samples; when doing so, each sample is given equal weight (Box 1; Figure 1b). The RRA of
161 taxa in each sample will vary depending on genetic marker, laboratory protocol, and
162 bioinformatic filtering strategy (Alberdi *et al.* 2017; Deagle *et al.* 2013). Ensuring laboratory
163 methods are robust (i.e. focussing on samples with sufficient target DNA and checking
164 replicates) and using a standardised bioinformatics pipeline without excessive filtering can

165 help ensure RRA data are reproducible and precise (Alberdi *et al.* 2017; Deagle *et al.* 2013;
166 McInnes *et al.* 2017a; Murray *et al.* 2015).

167

Box 1: Some metrics used to summarise sequence data in dietary metabarcoding

168

169

170

Occurrence Data

171

172

173

174

175

176

177

Frequency of occurrence (FOO) is the number of samples that contain a given food item, most often expressed as a percent (%FOO). Percent of occurrence (POO) is simply %FOO rescaled so that the sum across all food items is 100%. Weighted percent of occurrence (wPOO) is similar to POO, but rather than giving equal weight to all occurrences, this metric weights each occurrence according to the number of food items in the sample (e.g., if a sample contains 5 food items, each will be given weight 1/5). Mathematical expressions are as follows:

178

$$\%FOO_i = \frac{1}{S} \sum_{k=1}^S I_{i,k} \times 100\%$$

179

180

$$POO_i = \frac{\sum_{k=1}^S I_{i,k}}{\sum_{i=1}^T \sum_{k=1}^S I_{i,k}}$$

181

182

$$wPOO_i = \frac{1}{S} \sum_{k=1}^S \frac{I_{i,k}}{\sum_{i=1}^T I_{i,k}}$$

183

184

185

186

where T is the number of food items (taxa), S is the number of samples, and I is an indicator function such that $I_{i,k} = 1$ if food item i is present in sample k , and 0 if not.

187

188

189

190

191

192

193

Many metabarcoding diet studies make use of both %FOO and POO (e.g. Xiong *et al.* 2017). POO provides a convenient view since each food taxon contributes a percentage of total diet (unlike %FOO which does not sum to 100%). In POO summaries samples with a high number of food taxa have a stronger influence, whereas in wPOO each sample is weighted equally (i.e. lower weighting to food taxa in a mixed meal) and this may be more biologically realistic (wPOO is the same as split-sample frequency of occurrence; see Tollit *et al.* 2017 and references within).

194

Read Abundance Data

195

Using the sequence counts, relative read abundance (RRA_i) for food item i is calculated as:

196

$$RRA_i = \frac{1}{S} \sum_{k=1}^S \frac{n_{i,k}}{\sum_{i=1}^T n_{i,k}} \times 100\%$$

197

198

199

where $n_{i,k}$ is the number of sequences of food item i in sample k .

200

201

202

203 **3. Does converting read counts to occurrence data solve our problems?**

204 It is often assumed that because conversion to occurrence data moderates the
205 impact of taxa-specific bias in marker signal, it provides a trustworthy, or at least
206 conservative, view of diet. While it is true that occurrence-based summaries of diet are less
207 affected by recovery bias, it is not necessarily the case that they provide a more accurate
208 representation of overall diet. Our simulations suggest POO summaries are highly consistent
209 but generally less accurate representation of overall diet compared to RRA summaries even
210 when moderate taxa-specific recovery biases are present (see Box 2 for details).

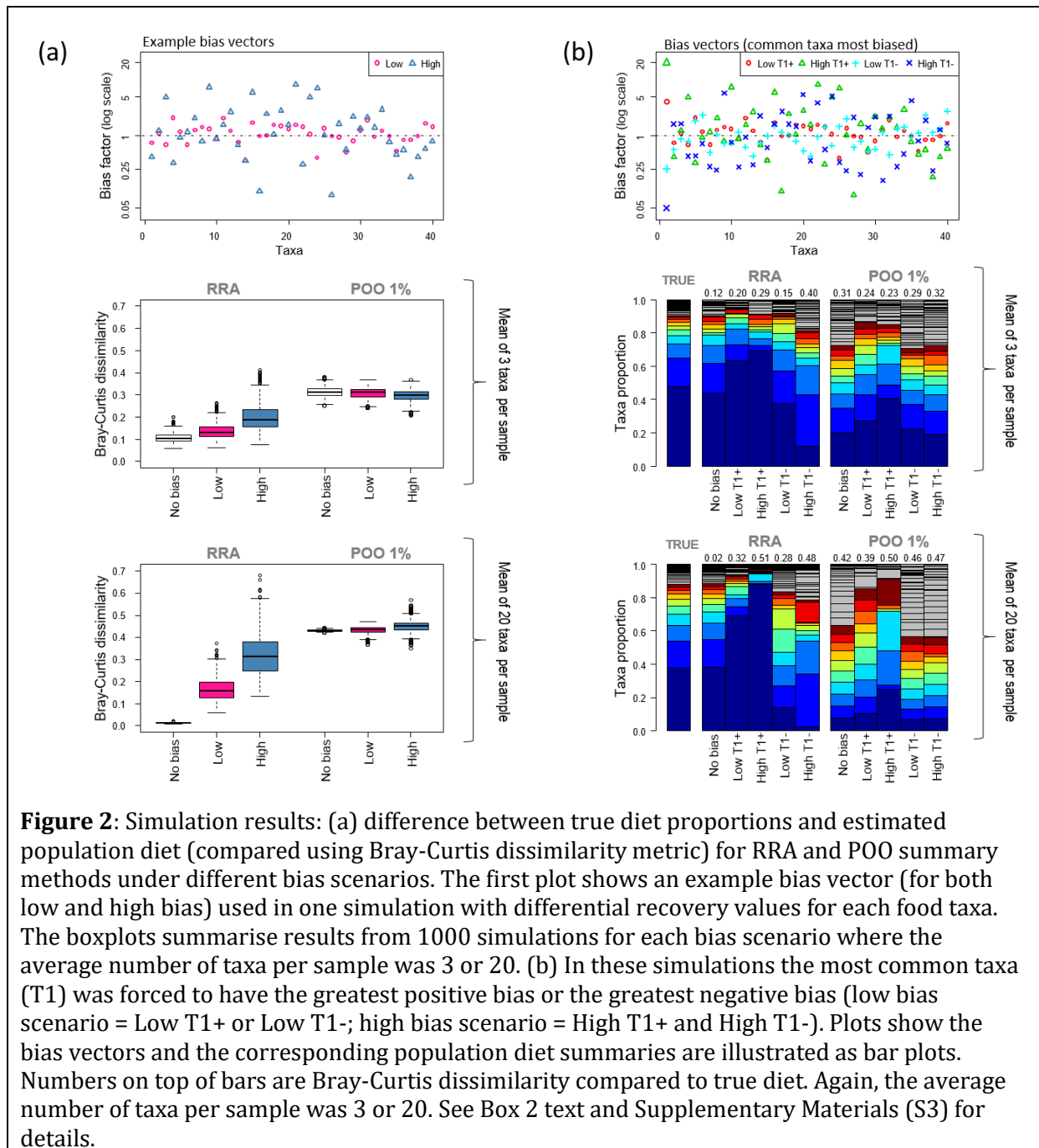
211

212 **Box 2: Simulations evaluating approaches for summarising population-level diet**
213 **composition**

214 To compare how effectively occurrence and RRA methods reconstruct population-level
215 diet we simulated HTS read counts for samples derived from a population with a fixed diet and
216 investigated the impact of taxa-specific sequence recovery biases (Figure 2). Our simulation
217 results are for a population with 40 food taxa in its diet, occurring in exponentially declining
218 abundance. Sequencing was simulated for 100 scat samples assuming a mean of either 3 or 20
219 food taxa per sample, and assuming different sequence recovery bias scenarios: no bias, low
220 bias or high bias. The biases introduce positive or negative biases of up to 4x and 20x (low and
221 high biases respectively) relative to a standard. In high bias scenario a 50:50 mixture could lead
222 to 400 fold recovery bias) Diet summaries were made using: (1) RRA; (2) POO with a 1%
223 minimum sequence threshold. For further details see Supplementary Material (S3).

224 Overall results show that with these parameters RRA summaries were on average more
225 accurate but had higher variance than POO summaries. POO produced more consistent
226 estimates less impacted by recovery biases, but only outperformed RRA when the largest
227 recovery biases corresponded to the most common food items. Both methods were more
228 accurate when the number of food taxa per sample was small: with a small number of food taxa
229 per sample POO estimates provide more realistic enumeration of rare items and RRA estimates
230 are less impacted by sequence recovery biases (since biases are only expressed in the context of
231 other taxa in a sample).

232



233

234 **Figure 2:** Simulation results: (a) difference between true diet proportions and estimated
 235 population diet (compared using Bray-Curtis dissimilarity metric) for RRA and POO summary
 236 methods under different bias scenarios. The first plot shows an example bias vector (for both
 237 low and high bias) used in one simulation with differential recovery values for each food taxa.
 238 The boxplots summarise results from 1000 simulations for each bias scenario where the
 239 average number of taxa per sample was 3 or 20. (b) In these simulations RRA the most common taxa
 240 (T1) was forced to have the greatest positive bias or the greatest negative bias (low bias
 241 scenario = Low T1+ or Low T1-; high bias scenario = High T1+ and High T1-). Plots show the
 242 bias vectors and the corresponding population diet summaries are illustrated as bar plots.
 243 Numbers on top of bars are Bray-Curtis dissimilarity compared to true diet. Again, the average
 244 number of taxa per sample was 3 or 20. See Box 2 text and Supplementary Materials (S3) for
 245 details.

246

247 The primary drawback of occurrence datasets is that the importance of rare food
 248 taxa are often artificially inflated at the expense of food taxa eaten in large amounts,
 249 effectively flattening the rank-abundance species curves typically seen in dietary datasets
 250 (Figure 1; Box 2). This effect can be illustrated in metabarcoding data from a population-
 251 level diet study of killer whales (Figure 3). This study concluded that the whale population's
 252 diet consisted primarily of Chinook salmon (~80%) based on high RRA of this species in most
 253 samples (Ford *et al.* 2016). If we consider the killer whales' diet as occurrence (POO; each

254 food species occurrence given equal value), the view changes considerably because other
255 salmon species and halibut frequently detected at low levels become important prey. The
256 threshold level used to count an occurrence also impacts the relative importance of these
257 fish prey; a lower threshold increases the importance of rare diet items (Figure 3). These
258 different diet estimates have substantial implications when diet percentages are combined
259 with bioenergetics estimates and consumer population size to derive estimates of prey
260 consumption (Chasco *et al.* 2017). Another implication of rare-item inflation occurs in
261 studies of niche partitioning. Here, the conclusion that species feed on separate resources
262 may be inaccurate because separation may be driven primarily by partitioning of rare diet
263 items, which are given similar weight as shared important food. In contrast, the conclusion
264 that species overlap in their dietary niche is potentially less likely (i.e. requiring overlap in
265 both primary and rare food items), but may therefore be more biologically meaningful when
266 found (Clare 2014).

267 How much influence rare diet taxa have in overall diet estimates depends to some
268 extent on the foraging strategy of the focal species and food distribution. In cases where
269 small amounts of rare diet items are consumed in most feeding bouts, the importance of
270 these items could be strongly over-estimated in occurrence-based summaries (as seen in
271 the simulations with a high number of taxa per scat sample; Box 2). This may be the
272 situation for some large grazing herbivores that forage continuously across a grassland,
273 often eating relatively rare plant taxa in proportion to their availability (i.e., non-selective
274 feeding). In contrast, when rare diet items are eaten sporadically, their DNA would be
275 detected only occasionally and diet estimates would be more realistic. For instance, some
276 carnivores feed sporadically, individualistically, and in discrete foraging events such that
277 prey occurrences may provide a more meaningful indication of how often each taxon is
278 predated (Codron *et al.* 2016). The feeding ecology of a species is reflected to some extent
279 in the number of food taxa in individual faecal samples and this varies widely between
280 studies (Table 1). This value provides insight into the potential impact of rare-item inflation
281 bias. For example, in Figure 1, the zebra faecal samples have many food taxa per sample and
282 when summarised as occurrences, these have a predictably flat rank-abundance curve; this
283 curve would be generated regardless of the true amount of each plant consumed in each
284 meal (Box 2).

285 Summaries based on occurrences become less accurate when samples are pooled
286 (i.e. when sequence reads from individual scats are not identifiable; Clare *et al.* 2014; Deagle
287 *et al.* 2009; Ford *et al.* 2016) because rare diet taxa present in any one of the pooled
288 samples are weighted equally to taxa found in all of the pooled samples. The time period
289 over which food consumption is integrated in a faecal DNA sample (influenced by gut
290 passage time) can affect these data in a similar way, since longer integration will mean rare
291 taxa have a greater likelihood of being present in each sample.

292 The inflated importance of rare sequences in occurrence summaries could also
293 magnify some problems encountered in diet metabarcoding. There are occasions when
294 exogenous DNA can contaminate a sample of interest. This includes field-based
295 contamination from non-food eDNA (McInnes *et al.* 2017a), laboratory contamination (De
296 Barba *et al.* 2014), and misassignment of sequence-to-sample during HTS (i.e. tag-jumping;
297 Schnell *et al.* 2015). These problems will generally have less influence in RRA summaries
298 since the real food items should dominate unless samples are very poor quality. A similar
299 issue is the detection of secondary predation (i.e. DNA from gut contents of ingested prey).
300 Depending on the study system and research question, secondary predation may or may not
301 be a serious problem. However, occurrence-based datasets are expected to over-emphasise
302 these detections and ruling out secondary predation in occurrence summaries may require
303 information of RRA, examination of prey co-occurrence, or expert knowledge (Bowser *et al.*
304 2013; Hardy *et al.* 2017; McInnes *et al.* 2017b).

305

306 **4. Does RRA actually reflect food biomass?**

307 The relationship between proportions of biological material in a sample and
308 sequence reads recovered by HTS has been studied in many experiments by sequencing
309 artificial mixtures with known composition. These ‘mock communities’ are most relevant to
310 dietary metabarcoding studies when made from food tissues similar to what is being
311 consumed. Both mitochondrial and chloroplast DNA markers are present in multiple copies
312 in each cell and copy number varies between tissue types (e.g. leaves versus roots; Ma & Li
313 2015) and physiological state (e.g. juvenile vs. gravid adult; Veltri *et al.* 1990). Getting a
314 thoroughly homogeneous mix of tissues in a small volume suitable for DNA extractions is
315 challenging; therefore, mixtures made from DNA extracted separately for each taxa are
316 sometimes used (e.g. Ford *et al.* 2016; Krehenwinkel *et al.* 2017; Piñol *et al.* 2015). However,

317 mixing purified genomic DNA will miss differences in cell density, and differences in genome
318 size will confound results, making interpretation difficult (Piñol *et al.* 2015). Mixtures of PCR
319 products can identify technical biases (e.g. assessing PCR primers), but miss underlying
320 biological differences.

321 Conclusions from analyses of mock communities vary from no relationship to good
322 correlations between the composition of the mixture and sequence reads (Edgar 2017;
323 Kimmerling *et al.* 2018; Pornon *et al.* 2016). One reason for these different conclusions is
324 that the range of concentrations analysed varies considerably across studies, from equal
325 mixtures of a few taxa, to mixtures containing many taxa in very different abundances. A
326 positive relationship between RRA and sample composition across a broad range of
327 concentrations (often plotted on a log-log scale (e.g. Elbrecht & Leese 2015; Nichols *et al.*
328 2016)) might be missed over a smaller range. High variability between studies is also due to
329 biotic differences in target organisms and technical differences (e.g. different barcode
330 markers, PCR primers, sequencing platforms, etc.). This variation makes it difficult to
331 generalise, and considerable work is required to understand the reliability of RRA in any
332 system. Two taxonomic prey groups that have been the focus of several dietary
333 metabarcoding studies, and for which mock communities have been examined, are fish and
334 insects. These groups provide some insight into the expected scale of biases.

335 In metabarcoding of fish mixtures, conserved PCR primers are generally employed
336 and documented recovery biases are moderate. In their killer whale study, Ford *et al.* (2016)
337 analysed known percentages of DNA extracted from four fish species and the RRA of each
338 fish corresponded well to input (generally within 5% of expected values) providing
339 confidence in their conclusions. Using prey species of harbour seals Thomas *et al.* (2016)
340 carried out a detailed study on sequence recovery from blended tissue mixtures. Various
341 taxa (primarily fish; n=18) were sequenced in 50:50 tissue mixes with a control fish, and the
342 extent of deviations from the control fish measured. The recovered sequences varied from
343 20% to 60%, a 3-fold variation in marker recovery relative to the control. A recent study
344 looking at recovery of barcode markers from bulk samples of larval fish avoided marker
345 amplification by directly sequencing all DNA, then bioinformatically recovering relevant
346 marker sequences (Kimmerling *et al.* 2018). They found strong correspondence between
347 biomass in known mixtures and sequence counts, suggesting that without PCR amplification
348 biases, biological differences in mtDNA density between these fish are small. Even studies

349 looking at fish environmental DNA samples have found a relationship between fish density
350 and recovered sequence counts (Lacoursière-Roussel *et al.* 2016; Port *et al.* 2015; Thomsen
351 *et al.* 2016).

352 Many studies have sequenced DNA from insect mock communities; however, rather
353 than considering if read counts are proxies for input biomass, the focus of these studies has
354 generally been to test if taxa can be detected (Alberdi *et al.* 2017; Clarke *et al.* 2014;
355 Elbrecht & Leese 2015; Yu *et al.* 2012). The reason for this focus is that insect communities
356 tend to be complex, with many rare taxa, and the recovery biases large. In studies by Yu *et al.*
357 *et al.* (2012) and Clarke *et al.* (2014), a paltry 43-76% of species known to be present in mock
358 communities were recovered. A study that included a mixture containing equal amounts of
359 purified DNA from 12 arthropod species (10 insects, 2 spiders), reported RRA values for half
360 of the species were that were more than 100 times lower than expected (i.e. expected 8%
361 and recovered at <0.08% (Piñol *et al.* 2015)). Another arthropod study found consistent
362 relationships between percentages of DNA and RRA; however, the slope of the correlation
363 deviated from the expected value of 1 in different insect orders and with different DNA
364 markers, which was attributed to copy number variation (Krehenwinkel *et al.* 2017). Even a
365 change in PCR primers used to amplify a marker from the same gene can produce very
366 different results (Alberdi *et al.* 2017). Because of the generally poor correlation between
367 biomass and read counts most diet studies looking at insectivorous predators focus on
368 occurrence data (Table 1), but methodological improvements may change this (Jusino *et al.*
369 2017).

370 Diet studies incorporate more complexity than analysis of mock communities due to
371 potential differential digestion of food taxa. Relatively few captive feeding experiments have
372 examined how well dietary DNA counts reflect known diet, but studies have been carried
373 out on herbivores (sheep, deer) and marine predators (penguins, seals). These have
374 focussed on simple diets (~2-6 diet items) and results generally show that comparisons
375 between major and minor diet components are reflected in RRA. For example, the diet of
376 sheep fed two plants in ratios of 0:100, 25:75, 50:50, 75:25, 100:0 had a good correlation
377 with the percentages of DNA marker sequences amplified from rumen content (Willerslev *et al.*
378 *et al.* 2014). In a study on captive deer, >90% of the diet was made up of three plant species
379 with two other species fed in low amounts. In this case >90% of sequences came from the
380 three dominant taxa, but considering just these taxa, the correlation between what went in

381 and what came out was poor (Nakahara *et al.* 2015). Similarly, in faecal samples from
382 captive penguins fed pilchards as the majority of their diet, sequence reads from pilchards
383 were most common in the data; however, the three other fish species fed in mass ratios
384 45:35:20 produced sequences counts of 60:6:34 (Deagle *et al.* 2010).

385 Detailed captive feeding studies examining quantitative prey DNA recovery have
386 been carried out on captive seals and sea lions (Bowles *et al.* 2011; Deagle & Tollit 2007;
387 Thomas *et al.* 2014). Early studies used quantitative PCR rather than HTS and found the
388 amount of marker DNA recovered provided a reasonable indication of biomass ingested
389 (Bowles *et al.* 2011; Deagle & Tollit 2007). A trial with harbour seals by Thomas *et al.* (2014)
390 compared HTS data from food tissue (affected by biological and technical biases) with faecal
391 DNA (affected by digestion as well). The scale of bias introduced by digestion was generally
392 smaller than biases observed in undigested fish tissue mix. Since digestion bias may be in
393 the same or opposite direction to tissue biases, the overall effect is expected to increase
394 variance in prey-specific recovery biases compared to tissue mixes. These seal studies all
395 excluded prey hard parts from DNA extractions, but in other systems where this may not be
396 feasible, digestion biases could be larger. For example, faeces from insectivorous animals
397 often contain relatively undigested hard body parts (i.e. exoskeleton). The impact on DNA
398 recovery is difficult to assess: hard fragments will contain undigested DNA, but the DNA may
399 not be extracted as efficiently as DNA present from soft bodied prey (Clare 2014).

400 Another approach to understanding how much of a signal is present in counts from
401 DNA sequences is to compare results with other methods of diet analysis. In a study of large
402 mammalian herbivores, Kartzinel *et al.* (2015) found a nearly one-to-one correlation
403 between estimates of C₄ grass (family Poaceae) consumption based on stable isotopes
404 analyses and RRA based on metabarcoding of the chloroplast marker (trnL-P6). The use of
405 alternative proxies for diet composition can also reveal complexities. Craine *et al.* (2015)
406 used similar protocols to Kartzinel *et al.* (2015) but found C₄ grass RRA to be under-
407 represented compared to measures based on stable isotopes. They suggested that
408 chloroplast density scales with foliar nitrogen concentrations so that RRA values could
409 reflect dietary sources of protein, and thus may deviate from dietary sources of biomass as
410 represented by carbon stable isotopes. If RRA values based on this marker occasionally
411 reflect an animal's source of protein more closely than its source of carbon (i.e., biomass),
412 this knowledge can enable count data to still be interpreted appropriately.

413 Several studies have used traditional morphological analysis of food remains to help
414 cross-validate RRA data (Soininen *et al.* 2009; Thomas *et al.* 2017). Thomas *et al.* (2017)
415 analysed DNA and prey hard parts in >1000 seal faecal samples, and while there were minor
416 differences between methods in prey recovery and taxonomic resolution, both methods
417 provided a highly similar picture of population-level diet (Thomas *et al.* 2017; Table S2).
418 Cross-validation has the problem that all methods of diet determination are biased, so if
419 there is disagreement the correct answer may be unclear (Soininen *et al.* 2009). However,
420 congruence between datasets is reassuring and known biases can be taken into account
421 when making conclusions (e.g. jellyfish are digested quickly, so occurrence in faecal DNA but
422 not stomach contents is credible; Jarman *et al.* 2013; McInnes *et al.* 2017b). Large
423 differences in results between methods warrant further investigation; multiple lines of
424 independent evidence provide the strongest support for any conclusion.

425 Overall, assessing recovery bias between food taxa is complex, specific to a study
426 system, and can require significant effort. In some cases, broad correlations are likely, but
427 this cannot be taken for granted and very large biases may occur (e.g. Pawluczyk *et al.*
428 2015).

429

430 **5. A view of the way forward in interpreting sequence counts**

431 What should be considered best practice given the potential biases we have outlined
432 in diet metabarcoding studies? First of all, we should take a step back and remember that
433 getting estimates of the true diet of any species using any method is a challenging
434 proposition – all methods of diet analysis have biases. A well-designed metabarcoding diet
435 study may provide as accurate an estimate as any other approach, while also providing high
436 taxonomic resolution, the opportunity to detect rare foods and the potential to solve
437 otherwise intractable problems (e.g. liquid feeding). We should also remember that other
438 classic experimental design issues, such as collecting appropriate sample sizes and getting
439 representative samples, will potentially have a bigger impact on study outcomes than the
440 diet estimation method. Furthermore, dietary metabarcoding has a huge variety of
441 applications, many of which do not require highly accurate dietary proportions.

442 Still, we will inevitably come to a point in dietary metabarcoding studies where we
443 need to decide how to interpret sequence counts. It is often the case that the overarching
444 views of population-level diet are consistent regardless of how sequence counts are

445 summarised (i.e. when commonly occurring food items are also represented by the highest
446 number of sequences). This is most likely to be the case when faecal samples contain a
447 limited number of food taxa (in the extreme case where there is only one taxon per sample,
448 occurrence and RRA estimates are identical and recovery biases have no impact). However,
449 some outcomes will depend on how we consider counts. Occurrence summaries are less
450 affected by differential recovery of markers from food taxa, but tend to put much more
451 weight on food consumed in small quantities and potential contaminants. RRA can
452 potentially provide a weighting of food present in a sample based on biomass, but
453 differential recovery of markers (especially from dominant food taxa) can impact data
454 summaries. Our strongest recommendation is that if one approach is relied on heavily,
455 some justification should be given for its use, and potential biases inherent in the method
456 should be acknowledged and taken into account when drawing conclusions.

457

458 **5.1 Using occurrence data**

459 Many future diet studies will have almost no information on the scale of biases in
460 the recovery of sequences from specific food taxa. The use of occurrence data may be a
461 sensible approach, but careful consideration of the impact of this choice is still required and
462 the bioinformatics steps taken to produce this dataset should be documented. We
463 recommend converting counts to percentages (excluding non-food sequences from total
464 count) and then defining a minimum sequence percentage threshold to determine
465 occurrences. This will limit the impact of variation in read depth. The threshold is a trade-off
466 between maximizing inclusion of real diet sequences and excluding low-level background
467 noise (secondary predation, contamination, sequencing errors). A 1% threshold may be
468 suitable for many situations, but when diets are extremely diverse with potentially large
469 recovery biases (e.g. some bats species), then a much lower threshold may be justified (e.g.
470 0.01% in Alberdi *et al.* 2017). In these cases, ensuring contaminant sequences do not
471 influence results requires additional vigilance (De Barba *et al.* 2014; Nguyen *et al.* 2015).
472 Given that many of the issues we have raised regarding the use of occurrence data stem
473 from the disproportionate influence of rarer sequences, it may seem advantageous to use a
474 higher minimum sequence threshold (e.g. >5% constitutes occurrence). While this type of
475 summary can provide insight, rare taxa that make up a small percentage of sequences in
476 many samples would be missed completely (Alberdi *et al.* 2017) and taxa-specific biases in

477 recovery also have a larger impact on these high threshold occurrence summaries (see
478 simulations in Supplementary Material S3 comparing different threshold levels). Since the
479 purported benefit of occurrence-based approaches is to record food taxa even when there
480 is strong bias against them, thresholds higher than 1% cannot be generally recommended.

481 The sequencing depth required per sample is directly related to the minimum
482 threshold; in diverse and/or potentially highly biased situations warranting a very low
483 threshold (e.g. 0.01%), high numbers of reads per sample would be needed (e.g. >10000).
484 Lower read depth is sufficient with a 1% threshold and increasing replication (biological or
485 technical) would be preferable to having redundant sequences within samples. Even when
486 sequence counts are not used directly, it is important these data are available as
487 supplementary material (and ideally the sequence reads archived) with appropriate
488 explanatory files outlining potential biases. This allows others to revisit the data and will
489 allow insight in future comparative meta-analyses.

490 Summaries of data based only on occurrence information will remain appropriate in
491 many situations. This includes dietary metabarcoding studies that use DNA from food
492 remains in gut contents since differences in time since ingestion will have a major impact on
493 relative number of reads recovered per taxon (Egeter *et al.* 2015; Greenstone *et al.* 2014). In
494 studies using faecal samples, occurrence summaries will often be appropriate when food is
495 clearly differentially digested, the sequence recovery bias is known to be high (e.g. many
496 animals with an insectivorous diet), or this bias is unknown and results cannot be cross-
497 validated. Note, that this appropriateness may differ between dietary analyses of relatively
498 similar consumers. For example, most bat diet studies only analyse occurrence data, but the
499 bat *Myotis daubentonii* (Figure 1) has relatively low diet richness compared to other bats
500 and RRA may be suitable (Vesterinen *et al.* 2016).

501

502 **5.2 Using RRA**

503 Incorporation of RRA into analyses will have the most benefit when individual faecal
504 samples contain several food taxa and the same food taxa occur across many samples. In
505 these cases, occurrence summaries may provide very inaccurate summaries (Box 2).
506 Unfortunately RRA-based summaries from these types of samples can be most affected by
507 recovery biases (Box 2) and careful decisions about how to interpret data are required.
508 When there is uncertainty surrounding which method will be more accurate, presentation

509 of results summarised with both methods is recommended. Conclusions relying heavily on
510 RRA should include justification as to why the counts are expected to contain a roughly
511 accurate signature. One way to justify interpretations based on RRA is through cross-
512 validation of the diet data with alternative methods, and this is recommended whenever
513 possible. Alternatively, mock community experiments and/or feeding trials can be carried
514 out, but this is feasible in a limited number of situations. In study systems where diet is
515 relatively well known, examining biases in a small number of dominant food taxa can ensure
516 they are not drastically over or underestimated and will lend support to using RRA
517 information. The dominant diet items have by far the strongest impact on RRA diet
518 summaries as significant shifts in percentages of these species will adjust percentages of all
519 food taxa (i.e. unit sum constrained data must sum to 100%). One question that inevitably
520 arises is at what point does “semi-quantitative” RRA information stop being useful? Our
521 simulations indicate that even in scenarios with 20x overestimation of some food and 20x
522 underestimation of others (i.e. in 50:50 mixtures this could lead to 400 fold recovery bias) the
523 population-level RRA summaries often still provides a more accurate view of diet compared
524 to POO (Figure 2). But the limits of usefulness will depend on the application. It is probable
525 that comparisons between closely related food taxa will provide more reliable RRA data,
526 because biological differences should be smaller and technical biases less pronounced (e.g.
527 animal COI primer binding sites will be more conserved, or length differences in the plant
528 trnL-P6 marker will be low). However, it is risky to make generalizations and to transfer
529 specific methodological findings between study systems.

530 Further refinements to increase confidence in RRA dietary metabarcoding data are
531 possible. Because conversion to occurrence datasets has been seen as a necessary remedy
532 for biases in sequence recovery, there has been less incentive for researchers to test new
533 protocols and evaluate markers on their ability to obtain accurate RRA data. While it is
534 sensible to use standard DNA barcode markers, by ignoring information in RRA during
535 marker development we might have inadvertently imposed limitations on the field.
536 Fortunately, we are starting to move towards a point where markers used in different
537 applications are better understood and alternative less-biased approaches are being
538 explored (e.g. the use of multiple target markers (Stat et al. 2017) or PCR-free approaches
539 (Srivathsan et al. 2016)). Inclusion of control materials in sequencing runs can also ensure
540 consistency between experiments (Hardwick *et al.* 2017). For the most accurate diet

541 estimates, correction factors can be developed to take into account known biological
542 differences between taxa in mixtures (e.g. gene copy number differences; Angly et al. 2014;
543 Vasselon et al. 2018). Such species-specific correction factors have been developed for fish,
544 with the intent of applying them in field-collected seal diet samples (Thomas *et al.* 2016).

545 While the effort needed to justify the RRA approach may be challenging, the
546 possibility of obtaining more accurate diet estimates will make it worthwhile in many
547 situations. We have seen such effort undertaken in papers addressing broad ecological
548 questions (Kartzinel *et al.* 2015; Willerslev *et al.* 2014), and in diet studies of marine
549 predators, where population consumption have significant fisheries management
550 implications (Ford *et al.* 2016; Thomas *et al.* 2017). This approach should also be possible in
551 monitoring programs, such as those carried out on seabird diet (Jarman *et al.* 2013;
552 Sydeman *et al.* 2017), where the long-term investment warrants the development of robust
553 DNA-based methods that provide the best possible data.

554

555 **5.3 Outstanding issues**

556 There are a number of issues in the diet metabarcoding literature that have an
557 impact on both occurrence and RRA summaries that have yet to be clearly addressed. One
558 of these is the impact of collecting data with markers that have low taxonomic resolution
559 (McInnes *et al.* 2017b) or collating data at higher taxonomic levels to increase certainty in
560 taxonomic assignment (Biffi *et al.* 2017). Depending on how broad the grouping are,
561 occurrence summaries may not be very informative as many occurrences are potentially
562 pooled. For RRA it is unclear whether pooling counts from multiple taxa will nullify fine-scale
563 stochasticity in recovery biases, or magnify lineage-specific biases. A related issue is how to
564 summarise data from diet metabarcoding studies using multiple markers. When markers are
565 targeting the same food taxa, either additive (i.e. include detections by any marker) or
566 restrictive strategies (only include food detected by all markers) could be logically applied in
567 occurrence and RRA summaries (Alberdi *et al.* 2017). The situation is even more complex
568 when a “universal” primer set is used to define the broad diet and group-specific primers
569 subsequently improve taxonomic resolution for particular groups (e.g. a marker targeting all
570 plants together with several that offer greater resolution for specific plant families). Errors
571 based on the universal marker will be propagated when attempting to incorporate data
572 from the other primer sets (i.e. if the grass family is estimated to be 20% of a diet instead of

573 the true 40%, then the perceived importance of each grass species is reduced). This problem
574 can be avoided to some extent by reporting each component separately, but this provides
575 an unsatisfactory synthesis for omnivorous and other species with a very diverse diet that
576 can only be characterised with several markers (De Barba *et al.* 2014). Studies that use a
577 marker capturing only one component of the diet need to be very clear that the results
578 comprise an unknown amount of the total diet.

579 Simulations such as the ones outlined in this paper can help establish which
580 scenarios are most sensitive to biases (from either occurrence or RRA). When informed by
581 experimental work to assign an error range to each parameter, and combined with
582 sensitivity analysis, this can identify which sources of bias have the largest impact on
583 conclusions. Some of the details we have focussed on may be inconsequential for many
584 studies and we have not considered the effect of alternate summaries on downstream
585 applications. For example, it would be very interesting to see how switching between
586 occurrence and RRA datasets affects outputs in the context of food web studies (Roslin &
587 Majaneva 2016).

588 The ultimate test for how to deal with sequence counts in HTS diet analyses will
589 remain in empirical studies. We hope this opinion piece will be a starting point to highlight
590 the need to consider all sources of bias and to justify the methods used when confronting
591 count data in metabarcoding studies. We also hope that this critique is not discouraging to
592 researchers approaching this new and rapidly developing area of research, as no single
593 study should be rightly expected to address all issues arising from DNA-based diet analyses.
594 Instead, our aim is to encourage researchers to continue to addressing methodological
595 challenges, and acknowledge unanswered questions to help spur future investigations. As
596 the field matures, we envisage publication standards will emerge to provide the most robust
597 diet data and provide an accurate indication of the uncertainty associated with dietary
598 assessments.

599

600 **Data Accessibility**

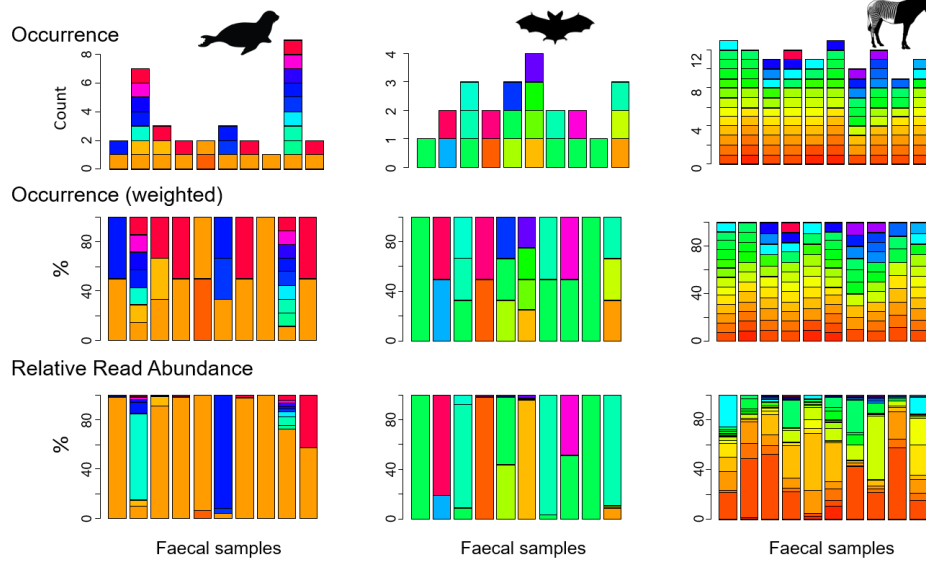
601 All data in figures is either publically accessible or will be deposited in Dryad along with R
602 scripts to produce the figures (including simulations).

603 **Author Contributions**

604 All Authors contributed ideas and to the writing of the paper

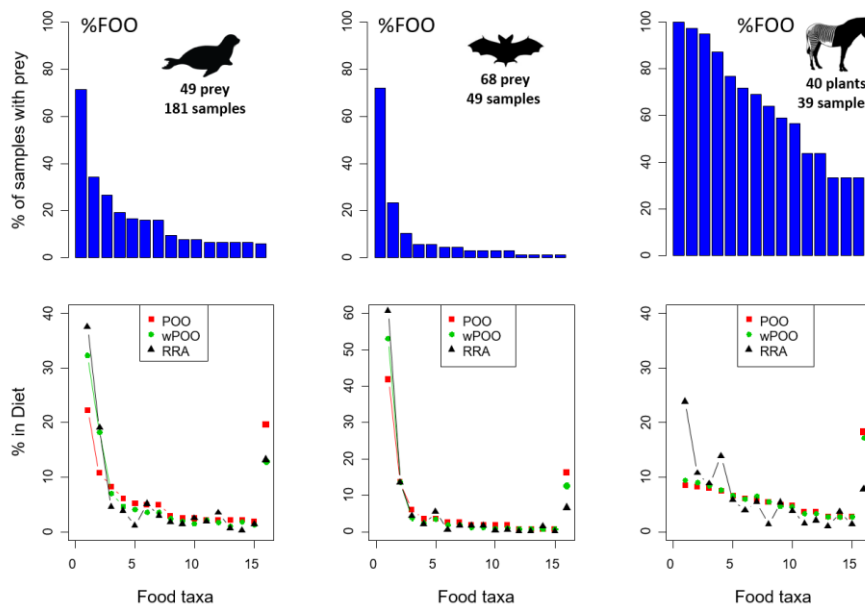
605

(a) Views of HTS data from 10 individual faecal samples



606

(b) Different population-level summaries of HTS datasets



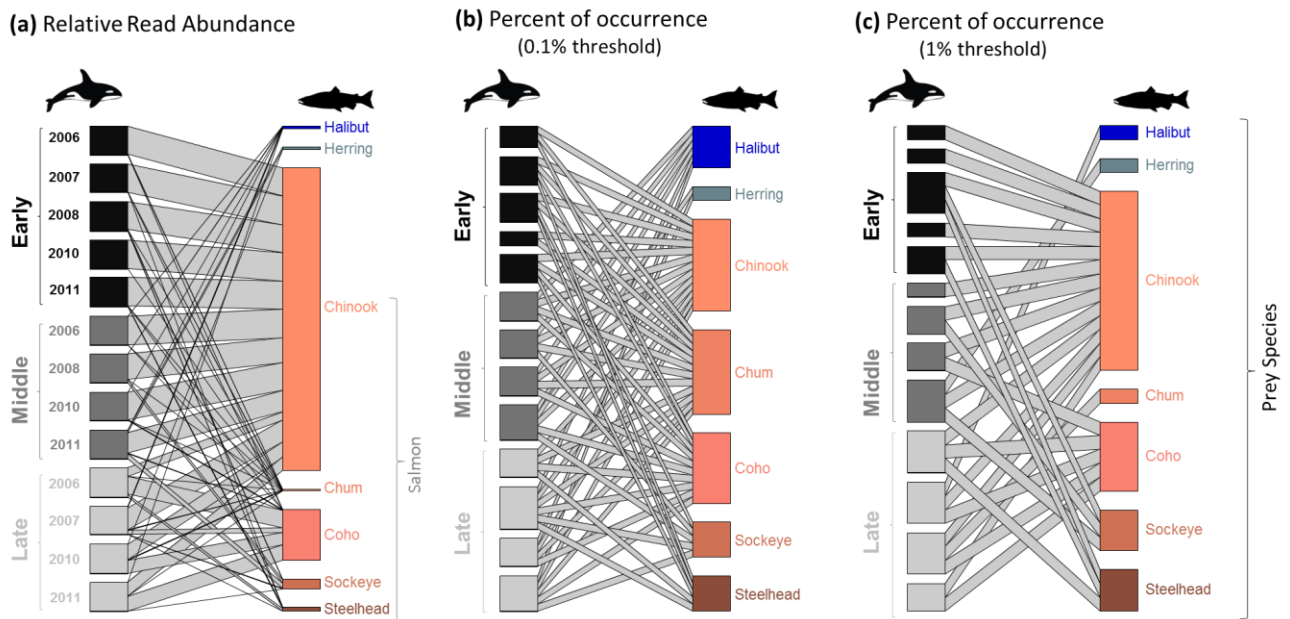
607

608

609 **Figure 1:** Information in faecal samples from dietary metabarcoding datasets of harbour seal
 610 (Thomas *et al.* 2017), an insectivorous bat (Vesterinen *et al.* 2016) and Grevy's zebra
 611 (Kartzinel *et al.* 2015). (a) Individual-level data in 10 faecal samples viewed using different
 612 metrics. Colours represent different food taxa. (b) Population-level summaries of these
 613 datasets showing the top 15 food taxa (%FOO ranking); 1% threshold used for occurrence in
 614 POO and wPOO calculations. In the lower plots the sum contribution of remaining food taxa
 615 are plotted at end. In each example population data include only collections from one site
 616 and samples with >50 food taxa reads.

617

618



619

620 **Figure 3:** Killer whale diet in the Salish Sea illustrated with bipartite graphs constructed from
621 data in Ford *et al.* (2016) using either (a) RRA (b) POO with a 0.1% threshold or (c) POO with
622 a 1% threshold. Samples (DNA from faecal material) are shown on left of each plot and were
623 pooled according to collection dates (Early, Middle, Late) in different years. The overall diet
624 calculated by the different methods is shown on the right of each plot (includes the seven
625 prey taxa with >1% of sequences in at least one sample). Line thickness shows contribution
626 of taxa in each sample to the overall diet.

627

628 **Table 1** Use of sequence counts in 20 metabarcoding diet studies carried out using faecal DNA collected from a range of different species. Representative
 629 studies across a range of focal taxa carried out by different research groups are shown rather than trying to summarise all dietary metabarcoding studies.

Focal Taxa	Reference	FOO†	RRA‡	Sample number	Number food taxa§	Taxa per sample¶	Marker	Target group	Sequences per sample*	Sequencer	Count data Available
Snail	O'Rorke <i>et al.</i> (2016)	N	Y	35	>50	NR	ITS	fungus	3500 (rarefied)	MiSeq	Yes
Snail	Waterhouse <i>et al.</i> (2014)	Y	N	60	26	4.7	16S	earthworms	1047	454	No
Pigeon	Ando <i>et al.</i> (2013)	Y	Y	48	44	6.7	trnL	plants	743	454	No
Albatross	McInnes <i>et al.</i> (2017b)	Y	Y	447	~20	NR	18S	metazoan	>100 prey	MiSeq	Yes
Puffin	Bowser <i>et al.</i> (2013)	Y	Y	129	~40	NA	CO1, 16S	metazoan	>50 prey	454	No
Sandpiper	Gerwing <i>et al.</i> (2016)	Y	N	164	132	NA	CO1, 16S	metazoan, fish/cephalopod/crustacea	721^	454	No
Desman (Rodent)	Biffi <i>et al.</i> (2017)	Y	N	383	156	5.8	CO1	arthropods	6910^	Ion Torrent	No
Bat	Clare <i>et al.</i> (2014)	Y	N	25 (pooled)	>158	NA	CO1	arthropods	>10000*	Ion Torrent	No
Bats	Burgar <i>et al.</i> (2014)	Y	N	64	>120	15	CO1	arthropods	230	454	No
Bat	Vesterinen <i>et al.</i> (2016)	Y	Y	82	59	NR	CO1	arthropods	995	Ion Torrent	Yes
Bat	Aizpurua <i>et al.</i> (2018)	Y	Y	79	>276	8.4	CO1, 16S	arthropods	>10000*	MiSeq	No
Seal	Thomas <i>et al.</i> (2017)	N	Y	1166	71	3.2	CO1, 16S	salmon, fish and cephalopods	1227	MiSeq	No (Available on request)
Seal	Hardy <i>et al.</i> (2017)	Y	N	112	115	3 to 6	16S, 12S	vertebrates, invertebrates	>10000*	MiSeq	Yes
Killer Whale	Ford <i>et al.</i> (2016)	N	Y	13 (pooled)	16	NA	16S	fish	>10000*	MiSeq	Yes (raw sequences)
Bear	De Barba <i>et al.</i> (2014)	Y	N	91	>84	NA	trnL, 12S, 16S, ITS	plants, vertebrates, invertebrates	>500	HiSeq	Yes

Cats	Xiong <i>et al.</i> (2017)	Y	N	103	40	3.6-4.1	16S	vertebrates	>10000*	HiSeq	No
Monkey	Quéméré <i>et al.</i> (2013)	Y	N	96	>130	13.9	trnL	plants	23793	Illumina	No
Deer	Erickson <i>et al.</i> (2017)	N	Y	12	>91	71	rbcl	plants	>10000*	MiSeq	No
Large herbivores	Kartzinel <i>et al.</i> (2015)	Y	Y	292	>110	NA	trnL, ITS	plants	>10000*	HiSeq	Yes
Ibex and Goat	Gebremedhin <i>et al.</i> (2016)	Y	Y	39	>50	NR	trnL	plants	>8000	454	Yes

630

631 † For this table Frequency Of Occurrence (FOO) refers to any use of presence/absence data

632 ‡ For this table Relative Read Abundance (RRA) refers to the use of sequence counts to weight taxa present in samples. This includes distance methods such as Bray-Curtis dissimilarity.

634 § Taxonomic level of assignments varies between studies, therefore the number of taxa is not directly comparable.

635 ¶ In some cases multiple markers were used, or multiple samples were pooled, making this value Not Applicable (NA). NR indicate the number of food taxa per sample was Not Reported.

637 † Most studies report mean number of food taxa sequences recovered per sample, but variance is not usually provided. The minimum number was reported in some cases.

638 ^ Unclear if these sequence counts include non-target DNA such as consumer DNA.

639 * The maximum value reported here was 10000 reads per sample.

640 References

641

- 642 Aizpurua O, Budinski I, Georgiakakis P, *et al.* (2018) Agriculture shapes the trophic niche of a bat
643 preying on multiple pest arthropods across Europe: evidence from DNA metabarcoding.
644 *Molecular Ecology* **27**, 815–825.
- 645 Alberdi A, Aizpurua O, Gilbert MTP, Bohmann K (2017) Scrutinizing key steps for reliable
646 metabarcoding of environmental samples. *Methods in Ecology and Evolution*.
- 647 Amend AS, Seifert KA, Bruns TD (2010) Quantifying microbial communities with 454 pyrosequencing:
648 does read abundance count? *Molecular Ecology* **19**, 5555-5565.
- 649 Ando H, Setsuko S, Horikoshi K, *et al.* (2013) Diet analysis by next-generation sequencing indicates
650 the frequent consumption of introduced plants by the critically endangered red-headed
651 wood pigeon (*Columba janthina nitens*) in oceanic island habitats. *Ecology and Evolution* **3**,
652 4057-4069.
- 653 Barrett RT, Camphuysen K, Anker-Nilssen T, *et al.* (2007) Diet studies of seabirds: a review and
654 recommendations. *ICES Journal of Marine Science* **64**, 1675-1691.
- 655 Biffi M, Gillet F, Laffaille P, *et al.* (2017) Novel insights into the diet of the Pyrenean desman
656 (*Galemys pyrenaicus*) using next-generation sequencing molecular analyses. *Journal of*
657 *Mammalogy*.
- 658 Bowles E, Schulte PM, Tollit DJ, Deagle BE, Trites AW (2011) Proportion of prey consumed can be
659 determined from faecal DNA using real-time PCR. *Molecular Ecology Resources* **11**, 530-540.
- 660 Bowser AK, Diamond AW, Addison JA (2013) From puffins to plankton: a DNA-based analysis of a
661 seabird food chain in the northern Gulf of Maine. *PLoS One* **8**, e83152.
- 662 Burgar JM, Murray DC, Craig MD, *et al.* (2014) Who's for dinner? High-throughput sequencing
663 reveals bat dietary differentiation in a biodiversity hotspot where prey taxonomy is largely
664 undescribed. *Molecular Ecology* **23**, 3605-3617.
- 665 Chasco BE, Kaplan IC, Thomas AC, *et al.* (2017) Competing tradeoffs between increasing marine
666 mammal predation and fisheries harvest of Chinook salmon. *Scientific Reports* **7**, 15439.
- 667 Clare EL (2014) Molecular detection of trophic interactions: emerging trends, distinct advantages,
668 significant considerations and conservation applications. *Evolutionary Applications*, 1144–
669 1157.
- 670 Clare EL, Symondson WOC, Fenton MB (2014) An inordinate fondness for beetles? Variation in
671 seasonal dietary preferences of night-roosting big brown bats (*Eptesicus fuscus*). *Molecular*
672 *Ecology* **23**, 3633-3647.
- 673 Clarke LJ, Soubrier J, Weyrich LS, Cooper A (2014) Environmental metabarcodes for insects: in silico
674 PCR reveals potential for taxonomic bias. *Molecular Ecology Resources* **14**, 1160-1170.
- 675 Codron D, Codron J, Sponheimer M, Clauss M (2016) Within-population isotopic niche variability in
676 savanna mammals: disparity between carnivores and herbivores. *Frontiers in Ecology and*
677 *Evolution* **4**, 15.
- 678 Craine JM, Towne EG, Miller M, Fierer N (2015) Climatic warming and the future of bison as grazers.
679 *Scientific Reports* **5**, 16738.
- 680 De Barba M, Miquel C, Boyer F, *et al.* (2014) DNA metabarcoding multiplexing and validation of data
681 accuracy for diet assessment: application to omnivorous diet. *Molecular Ecology Resources*
682 **14**, 306-323.
- 683 Deagle BE, Chiaradia A, McInnes J, Jarman SN (2010) Pyrosequencing faecal DNA to determine diet
684 of little penguins: is what goes in what comes out? *Conservation Genetics* **11**, 2039-2048.
- 685 Deagle BE, Kirkwood R, Jarman SN (2009) Analysis of Australian fur seal diet by pyrosequencing prey
686 DNA in faeces. *Molecular Ecology* **18**, 2022-2038.
- 687 Deagle BE, Thomas AC, Shaffer AK, Trites AW, Jarman SN (2013) Quantifying sequence proportions in
688 a DNA-based diet study using Ion Torrent amplicon sequencing: which counts count?
689 *Molecular Ecology Resources* **13**, 620-633.

- 690 Deagle BE, Tollit DJ (2007) Quantitative analysis of prey DNA in pinniped faeces: potential to
691 estimate diet composition? *Conservation Genetics* **8**, 743-747.
- 692 Edgar RC (2017) UNBIAS: An attempt to correct abundance bias in 16S sequencing, with limited
693 success. *bioRxiv*, 124149.
- 694 Egeter B, Bishop PJ, Robertson BC (2015) Detecting frogs as prey in the diets of introduced
695 mammals: a comparison between morphological and DNA-based diet analyses. *Molecular
696 Ecology Resources* **15**, 306-316.
- 697 Elbrecht V, Leese F (2015) Can DNA-based ecosystem assessments quantify species abundance?
698 Testing primer bias and biomass—sequence relationships with an innovative metabarcoding
699 protocol. *PLoS One* **10**, e0130324.
- 700 Erickson DL, Reed E, Ramachandran P, *et al.* (2017) Reconstructing a herbivore's diet using a novel
701 rbcl DNA mini-barcode for plants. *AoB PLANTS* **9**.
- 702 Finotello F, Di Camillo B (2015) Measuring differential gene expression with RNA-seq: challenges and
703 strategies for data analysis. *Briefings in functional genomics* **14**, 130-142.
- 704 Ford MJ, Hempelmann J, Hanson MB, *et al.* (2016) Estimation of a killer whale (*Orcinus orca*)
705 population's diet using sequencing analysis of DNA from feces. *PLoS One* **11**, e0144956.
- 706 Forney LJ, Zhou X, Brown CJ (2004) Molecular microbial ecology: land of the one-eyed king. *Current
707 opinion in microbiology* **7**, 210-220.
- 708 Gebremedhin B, Flagstad Ø, Bekele A, *et al.* (2016) DNA Metabarcoding Reveals Diet Overlap
709 between the Endangered Walia Ibex and Domestic Goats - Implications for Conservation.
710 *PLoS One* **11**, e0159133.
- 711 Gerwing TG, Kim J-H, Hamilton DJ, Barbeau MA, Addison JA (2016) Diet reconstruction using next-
712 generation sequencing increases the known ecosystem usage by a shorebird. *The Auk* **133**,
713 168-177.
- 714 Goodwin S, McPherson JD, McCombie WR (2016) Coming of age: ten years of next-generation
715 sequencing technologies. *Nature Reviews Genetics* **17**, 333.
- 716 Greenstone MH, Payton ME, Weber DC, Simmons AM (2014) The detectability half-life in arthropod
717 predator-prey research: what it is, why we need it, how to measure it, and how to use it.
718 *Molecular Ecology* **23**, 3799-3813.
- 719 Hardwick S, Deveson I, Mercer T (2017) Reference standards for next-generation sequencing. *Nature
720 Reviews Genetics* **18**, 473-484.
- 721 Hardy N, Berry T, Kelaher BP, *et al.* (2017) Assessing the trophic ecology of top predators across a
722 recolonisation frontier using DNA metabarcoding of diets. *Marine Ecology Progress Series*
723 **573**, 237-254.
- 724 Ibarbalz FM, Pérez MV, Figuerola EL, Erijman L (2014) The bias associated with amplicon sequencing
725 does not affect the quantitative assessment of bacterial community dynamics. *PLoS One* **9**,
726 e99722.
- 727 Jarman SN, McInnes JC, Faux C, *et al.* (2013) Adélie penguin population diet monitoring by analysis of
728 food DNA in scats. *PLoS One* **8**, e82227.
- 729 Jusino MA, Banik MT, Palmer JM, *et al.* (2017) An improved method for utilizing high-throughput
730 amplicon sequencing to determine the diets of insectivorous animals. *PeerJ PrePrints*.
- 731 Kartzinel TR, Chen PA, Coverdale TC, *et al.* (2015) DNA metabarcoding illuminates dietary niche
732 partitioning by African large herbivores. *Proceedings of the National Academy of Sciences*
733 **112**, 8019-8024.
- 734 Kimmerling N, Zuqert O, Amitai G, *et al.* (2018) Quantitative species-level ecology of reef fish larvae
735 via metabarcoding. *Nature ecology & evolution* **2**, 306.
- 736 King RA, Read DS, Traugott M, Symondson WOC (2008) Molecular analysis of predation: a review of
737 best practice for DNA-based approaches. *Molecular Ecology* **17**, 947-963.
- 738 Krehenwinkel H, Wolf M, Lim JY, *et al.* (2017) Estimating and mitigating amplification bias in
739 qualitative and quantitative arthropod metabarcoding. *Scientific Reports* **7**, 17668.

- 740 Laake J, Browne P, DeLong R, Huber H (2002) Pinniped diet composition: a comparison of estimation
741 models. *Fishery Bulletin* **100**, 434-447.
- 742 Lacoursière-Roussel A, Côté G, Leclerc V, Bernatchez L (2016) Quantifying relative fish abundance
743 with eDNA: a promising tool for fisheries management. *Journal of Applied Ecology* **53**, 1148-
744 1157.
- 745 Ma J, Li X-Q (2015) Organellar genome copy number variation and integrity during moderate
746 maturation of roots and leaves of maize seedlings. *Current genetics* **61**, 591-600.
- 747 McInnes JC, Alderman R, Deagle BE, *et al.* (2017a) Optimised scat collection protocols for dietary
748 DNA metabarcoding in vertebrates. *Methods in Ecology and Evolution* **8**, 192-202.
- 749 McInnes JC, Alderman R, Lea M-A, *et al.* (2017b) High occurrence of jellyfish predation by black-
750 browed and Campbell albatross identified by DNA metabarcoding. *Molecular Ecology* **26**,
751 4831-4845.
- 752 Murray DC, Coghlan ML, Bunce M (2015) From benchtop to desktop: important considerations when
753 designing amplicon sequencing workflows. *PLoS One* **10**, e0124671.
- 754 Nakahara F, Ando H, Ito H, *et al.* (2015) The applicability of DNA barcoding for dietary analysis of sika
755 deer. *DNA Barcodes* **3**, 200-206.
- 756 Nguyen NH, Smith D, Peay K, Kennedy P (2015) Parsing ecological signal from noise in next
757 generation amplicon sequencing. *New Phytologist* **205**, 1389-1393.
- 758 Nichols RV, Åkesson M, Kjellander P (2016) Diet assessment based on rumen contents: A comparison
759 between DNA metabarcoding and macroscopy. *PLoS One* **11**, e0157977.
- 760 Nielsen JM, Clare EL, Hayden B, Brett MT, Kratina P (2017) Diet tracing in ecology: method
761 comparison and selection. *Methods in Ecology and Evolution*.
- 762 O'Rorke R, Holland BS, Cobian GM, Gaughen K, Amend AS (2016) Dietary preferences of Hawaiian
763 tree snails to inform culture for conservation. *Biological Conservation* **198**, 177-182.
- 764 Olova N, Krueger F, Andrews S, *et al.* (2017) Comparison of whole-genome bisulfite sequencing
765 library preparation strategies identifies sources of biases affecting DNA methylation data.
766 *bioRxiv*, 165449.
- 767 Pawluczyk M, Weiss J, Links MG, *et al.* (2015) Quantitative evaluation of bias in PCR amplification
768 and next-generation sequencing derived from metabarcoding samples. *Analytical and*
769 *bioanalytical chemistry* **407**, 1841-1848.
- 770 Piñol J, Mir G, Gomez-Polo P, Agustí N (2015) Universal and blocking primer mismatches limit the use
771 of high-throughput DNA sequencing for the quantitative metabarcoding of arthropods.
772 *Molecular Ecology Resources* **15**, 819-830.
- 773 Pompanon F, Deagle BE, Symondson WO, *et al.* (2012) Who is eating what: diet assessment using
774 next generation sequencing. *Molecular Ecology* **21**, 1931-1950.
- 775 Pornon A, Escaravage N, Burrus M, *et al.* (2016) Using metabarcoding to reveal and quantify plant-
776 pollinator interactions. *Scientific Reports* **6**, 27282.
- 777 Port J, O'Donnell J, Lowell N, Romero-Maraccini O, Kelly R (2015) Assessing the vertebrate
778 community of a kelp forest ecosystem using environmental DNA. *Molecular Ecology*.
- 779 Quéméré E, Hibert F, Miquel C, *et al.* (2013) A DNA metabarcoding study of a primate dietary
780 diversity and plasticity across its entire fragmented range. *PLoS One* **8**, e58971.
- 781 Schield DR, Walsh MR, Card DC, *et al.* (2016) EpiRADseq: scalable analysis of genomewide patterns
782 of methylation using next-generation sequencing. *Methods in Ecology and Evolution* **7**, 60-
783 69.
- 784 Schnell IB, Bohmann K, Gilbert MTP (2015) Tag jumps illuminated—reducing sequence-to-sample
785 misidentifications in metabarcoding studies. *Molecular Ecology Resources* **15**, 1289-1303.
- 786 Soininen EM, Valentini A, Coissac E, *et al.* (2009) Analysing diet of small herbivores: the efficiency of
787 DNA barcoding coupled with high-throughput pyrosequencing for deciphering the
788 composition of complex plant mixtures. *Frontiers in Zoology* **6**, 16.
- 789 Sydeman WJ, Piatt JF, Thompson SA, *et al.* (2017) Puffins reveal contrasting relationships between
790 forage fish and ocean climate in the North Pacific. *Fisheries Oceanography* **26**, 379-395.

- 791 Symondson WO, Harwood JD (2014) Special issue on molecular detection of trophic interactions:
792 Unpicking the tangled bank. *Molecular Ecology* **23**, 3601-3604.
- 793 Taberlet P, Bonin A, Zinger L, Coissac E (2018) *Environmental DNA: For Biodiversity Research and*
794 *Monitoring* Oxford University Press.
- 795 Thomas AC, Deagle BE, Eveson JP, Harsch CH, Trites AW (2016) Quantitative DNA metabarcoding:
796 improved estimates of species proportional biomass using correction factors derived from
797 control material. *Molecular Ecology Resources* **16**, 714-726.
- 798 Thomas AC, Jarman SN, Haman KH, Trites AW, Deagle BE (2014) Improving accuracy of DNA diet
799 estimates using food tissue control materials and an evaluation of proxies for digestion bias.
800 *Molecular Ecology* **23**, 3706-3718.
- 801 Thomas AC, Nelson BW, Lance MM, Deagle BE, Trites AW (2017) Harbour seals target juvenile
802 salmon of conservation concern. *Canadian Journal of Fisheries and Aquatic Sciences* **74**, 907-
803 921.
- 804 Thomsen PF, Møller PR, Sigsgaard EE, *et al.* (2016) Environmental DNA from seawater samples
805 correlate with trawl catches of subarctic, deepwater fishes. *PLoS One* **11**, e0165252.
- 806 Tollit D, Fritz L, Joy R, *et al.* (2017) Diet of endangered Steller sea lions in the Aleutian Islands: New
807 insights from DNA detections and bio-energetic reconstructions. *Canadian Journal of*
808 *Zoology*.
- 809 Valentini A, Miquel C, Nawaz MA, *et al.* (2009) New perspectives in diet analysis based on DNA
810 barcoding and parallel pyrosequencing: the trnL approach. *Molecular Ecology Resources* **9**,
811 51-60.
- 812 Vandeputte D, Kathagen G, D'hoë K, *et al.* (2017) Quantitative microbiome profiling links gut
813 community variation to microbial load. *Nature* **551**.
- 814 Veltri KL, Espiritu M, Singh G (1990) Distinct genomic copy number in mitochondria of different
815 mammalian organs. *Journal of cellular physiology* **143**, 160-164.
- 816 Vesterinen EJ, Ruokolainen L, Wahlberg N, *et al.* (2016) What you need is what you eat? Prey
817 selection by the bat *Myotis daubentonii*. *Molecular Ecology* **25**, 1581-1594.
- 818 Waterhouse BR, Boyer S, Wratten SD (2014) Pyrosequencing of prey DNA in faeces of carnivorous
819 land snails to facilitate ecological restoration and relocation programmes. *Oecologia* **175**,
820 737-746.
- 821 Willerslev E, Davison J, Moora M, *et al.* (2014) Fifty thousand years of Arctic vegetation and
822 megafaunal diet. *Nature* **506**, 47.
- 823 Xiong M, Wang D, Bu H, *et al.* (2017) Molecular dietary analysis of two sympatric felids in the
824 Mountains of Southwest China biodiversity hotspot and conservation implications. *Scientific*
825 *Reports* **7**.
- 826 Yu DW, Ji Y, Emerson BC, *et al.* (2012) Biodiversity soup: metabarcoding of arthropods for rapid
827 biodiversity assessment and biomonitoring. *Methods in Ecology and Evolution* **3**, 613-623.

828