

1

2 Received Date: 25-Sep-2015

3 Revised Date: 17-Jan-2016

4 Accepted Date: 28-Jan-2016

5 Article Type: Article

6 A framework for inferring biological communities from environmental DNA

7

8 Andrew Olaf Shelton^{1*}

9 James Lawrence O'Donnell²

10 Jameal F. Samhouri¹

11 Natalie Lowell²

12 Gregory D. Williams³

13 Ryan P. Kelly²

14

15 ¹ Conservation Biology Division, Northwest Fisheries Science Center, National Marine
16 Fisheries Service, National Oceanic & Atmospheric Administration, Seattle, WA 98112

17

18 ² University of Washington, School of Marine and Environmental Affairs, 3707 Brooklyn
19 Ave NE, Seattle, WA 98105

20

21 ³ Pacific States Marine Fisheries Commission, Under contract to Northwest Fisheries
22 Science Center, National Marine Fisheries Service, National Oceanic & Atmospheric
23 Administration, Seattle, WA 98112

24

25 * Corresponding author: ole.shelton@noaa.gov

26

27 **Keywords:** environmental DNA, community surveys, Bayesian statistics, ecosystem
28 assessment, multinomial-Poisson transformation, quantitative PCR.

29 **Abstract**

This is the author manuscript accepted for publication and has undergone full peer review but has not been through the copyediting, typesetting, pagination and proofreading process, which may lead to differences between this version and the [Version of Record](#). Please cite this article as [doi: 10.1002/EAP.1336](https://doi.org/10.1002/EAP.1336)

This article is protected by copyright. All rights reserved

30 Environmental DNA (eDNA)—genetic material recovered from an environmental
31 medium such as soil, water, or feces—reflects the membership of the ecological
32 community present in the sampled environment. As such, eDNA is a potentially rich
33 source of data for basic ecology, conservation, and management, because it offers the
34 prospect of quantitatively reconstructing whole ecological communities from easily-
35 obtained samples. However, like all sampling methods, eDNA sequencing is subject to
36 methodological limitations that can generate biased descriptions of ecological
37 communities. Here, we demonstrate parallels between eDNA sampling and traditional
38 sampling techniques, and use these parallels to offer a statistical structure for framing the
39 challenges faced by eDNA and for illuminating the gaps in our current knowledge.
40 Although the current state of knowledge on some of these steps precludes a full estimate
41 of biomass for each taxon in a sampled eDNA community, we provide a map that
42 illustrates potential methods for bridging these gaps. Additionally, we use an original
43 dataset to estimate the relative abundances of taxon-specific template DNA prior to PCR,
44 given the abundance of DNA sequences recovered post-PCR-and-sequencing, a critical
45 step in the chain of eDNA inference. While we focus on the use of eDNA samples to
46 determine the relative abundance of taxa within a community, our approach also applies
47 to single-taxon applications (including applications using qPCR), studies of diversity, and
48 studies focused on occurrence. By grounding inferences about eDNA community
49 composition in a rigorous statistical framework, and by making these inferences explicit,
50 we hope to improve the inferential potential for the emerging field of community-level
51 eDNA analysis.

52 **Introduction:**

53 A central aim of ecology is to understand the distribution and abundance of
54 organisms, which requires estimates of the occurrence, density, or biomass of the
55 organisms in natural populations. Whether counting individuals in a habitat, in a
56 population, or across an assemblage, making inferences about an entire community from
57 an observed subset of individuals is fundamental to ecological science. Unfortunately all
58 sampling techniques are potentially subject to bias, undermining accuracy and confidence
59 in estimates of critical ecological parameters. Visual surveys may overlook or misidentify
60 cryptic species, surveys that capture individuals with nets or traps may under-represent

61 small or elusive prey, and quadrat-sampling methods for non-mobile flora and fauna can
62 underestimate the abundance of rare species or miss landscape-scale patterns.
63 Fortunately there is a large and sophisticated literature dedicated to examining and
64 improving efficacy and reducing bias for a range of sampling problems for terrestrial,
65 marine, and aquatic systems (Cochran 1977, Royle and Nichols 2003, Cotter and Pilling
66 2007, Elith and Leathwick 2009). In this paper, we contribute to this literature by
67 developing a general statistical framework as well as specific statistical sampling
68 methods for the emerging field of environmental DNA.

69 Recent advances in molecular biotechnology have resulted in the emergence of a
70 new survey tool, whereby the DNA present in an environmental medium (such as soil or
71 water; hereafter environmental DNA or eDNA), can be used to infer the presence of
72 organisms nearby (Jerde et al. 2011, Yoccoz 2012). There are currently two distinct
73 molecular approaches for eDNA. In the first, the amount of a known DNA sequence -
74 presumably from a single taxon - is determined from quantitative polymerase chain
75 reaction (qPCR; Thomsen et al. 2012, Nathan et al. 2014). The second approach is to
76 amplify some suitable region from all genomes present in a sample using PCR, and
77 sequence the resulting products (amplicons) using massively parallel sequencing
78 technologies, without *a priori* knowledge of the organisms present or their genetic
79 sequences (e.g. Ventner et al. 2004, de Barba et al 2015, Leray and Knowlton 2015).
80 While the qPCR approach is being used in several applications to monitor rare or
81 invasive species (Lodge et al. 2012; Turner et al. 2014), such methods can involve
82 extensive development for each taxon of interest, and cannot easily provide insight into
83 community-level patterns. Sequencing methods could feasibly provide relative
84 abundance data for a suite of species in the community, as the relative proportions of
85 taxon-specific DNA sequences observed may reflect the relative proportions of DNA in
86 the environment (Yoccoz 2012). While attempts have been made to link sequence counts
87 to biomass (e.g. Evans et al. 2015), no such study has yet evaluated the complex chain of
88 processes and associated uncertainty linking these two states (Iversen et al. *In Press*).
89 Thus, one barrier to the widespread adoption of the sequencing approach is the lack of
90 formal methods for linking this new data type (counts of DNA sequences) to the

91 underlying pattern of interest (the abundance or biomass of taxa comprising a
92 community; Yoccoz 2012).

93 Conceptually, using eDNA to infer the biomass or abundance in a community is
94 largely analogous to traditional non-molecular methods. Figure 1 illustrates how eDNA
95 and traditional sampling attempt to provide information about the same quantity: the
96 biomass of each species in the environment. Both eDNA and traditional sampling aim to
97 make inferences about distinct stages that are potentially measurable (latent states;
98 represented by boxes in Fig. 1), and processes which transform one stage to the next
99 (arrows in Fig. 1).

100 Before turning to sampling methods for eDNA, we first describe a general
101 theoretical framework in terms of traditional sampling of a marine fish community, with
102 the goal of quantifying the biomass of each taxon. Common sampling methods for fish
103 communities include using a variety of net technologies (trawl, gillnets, cast nets, seines,
104 etc.), systems using baited hooks, and visual surveys. Importantly, the process of
105 inference from data using any of these methods can be conceptualized using the diagram
106 in Figure 1. We use fish communities as an example with which we are familiar, and for
107 which there is a long history of explicitly modeling uncertainty, but the larger point
108 applies to all ecological sampling.

109 For example, for a sample collected using a trawl net, the total biomass of fish
110 taxon i at a particular location, l , and a particular time, t , $B_{i,l,t}$ is a function of the biomass
111 or counts observed in the net, $F_{i,l,t}$ (Fig. 1). Given that we only observe $F_{i,l,t}$ the process
112 of estimating $B_{i,l,t}$ from $F_{i,l,t}$ can be written as a conditional quantity, $[B_{i,l,t}|F_{i,l,t}]$. For
113 expositional purposes, we simplify notation by assuming a single sample time and
114 location, $[B_i|F_i]$. As Fig. 1 shows, the biomass in the environment (B_i) is not connected
115 to the biomass captured by the net (F_i) by a single process but rather a chain of distinct
116 processes. A full description of the sampling process would explicitly include each step.
117 For example, researchers commonly extract a subsample of individuals (E_i) from the full
118 catch of the net (D_i) to determine the taxon-specific count (F_i), which itself may be
119 influenced by taxonomic identification errors or other processes (Fig. 1).

120 From this conceptual framing it should be clear that our estimate of the taxon's
121 biomass B_i is influenced by at least three sets of processes: (1) the sampling approach to

122 obtain the collection D_i , (2) the methods used to reduce the full collection to the
123 subsample E_i , and (3) the identification and enumeration methods that result in the taxon
124 specific count F_i . Statistically, we can expand the inference of interest $[B_i|F_i]$ into three
125 conditionally independent processes (for general discussion of conditional modeling see
126 Clark 2007, Cressie and Wikle 2011)

$$127 \quad [B_i|F_i] = [B_i|D_i][D_i|E_i][E_i|F_i]. \quad (1)$$

128 Thus, any estimate from sampling data must implicitly or explicitly make assumptions or
129 estimate these three components. For example, the second term on the right side, $[D_i|E_i]$
130 describes the proportion of the total catch taken in a subsample. If the entire catch is
131 included, $[D_i|E_i] = 1$, and this term can be dropped from the model.

132 While accounting for $[D_i|E_i]$ is relatively straightforward, other terms in eq. 1 are
133 more difficult. Indeed, determining how biomass present in the environment corresponds
134 to the total catch in the net, $[B_i|D_i]$, is a classic and persistently difficult problem that has
135 been explored extensively in ecology (Royle and Nichols 2003, Elith and Leathwick
136 2009) and fisheries (see the fisheries concepts of “catchability”, and “selectivity”;
137 Beverton and Holt 1957 section 8, Arreguín-Sánchez 1996, Venebles and Dichmont
138 2004). For our hypothetical marine fish example, the mesh size, design, and deployment
139 of the net, among other characteristics, will interact with the true density of each species
140 to determine which are captured (the quantity $[B_i|D_i]$ in eq. 1; Beverton and Holt 1957
141 section 8, Arreguín-Sánchez 1996). Similar challenges face the determination of $[E_i|F_i]$;
142 individual skill and experience will affect the efficacy and accuracy of taxonomic
143 identification. Our purpose here is merely to note that such complexities plague virtually
144 all sampling problems—whether terrestrial or marine, from the poles to the equator.

145 The basic inferential framework introduced above (eq. 1, Fig. 1) readily applies to
146 the problem of reconstructing ecological communities from eDNA. Below, we outline the
147 processes connecting ecological communities to observations of eDNA, and briefly
148 summarize the state of knowledge about each process. We then construct a statistical
149 model for analyzing community eDNA data that accounts for some of the processes that
150 can potentially bias inference from eDNA data and provide a worked example for
151 applying these methods to a marine eDNA dataset. We end by briefly discussing further
152 methodological needs for eDNA data and making recommendations for best practices.

153 Throughout, we focus on the use of eDNA for community sampling and highlight the
154 inferential and empirical connections between traditional and eDNA sampling methods.

155

156 **Conceptual models for eDNA**

157 Here we derive a model structure to estimate the relative amount of biomass
158 present in a community for some set of taxa of interest, by sampling eDNA. While we
159 develop the framework in the context of estimating abundance for multiple species from
160 sequenced DNA, both models of occurrence (e.g. Ji et al. 2014) and of single species
161 abundance (e.g. Jerde et al. 2011) are special cases in our framework, as will be discussed
162 later. Our general approach also applies to qPCR methodologies. We focus on the
163 detection and quantification of taxa that are not directly sampled. For example, if we
164 collected a liter of water from the environment, we focus primarily on inferring the
165 abundance of fishes, invertebrates, and mammals from individual cells (and
166 accompanying DNA) contained in that water sample. While similar methods could be
167 applied to bacteria and other microorganisms that can be directly measured and
168 sequenced from a small sample, we do not specifically address such cases here; direct
169 sequencing rather than PCR based approaches may be more appropriate for small,
170 abundant taxa (Yu et al. 2012).

171 To derive a general model for eDNA we need to explicitly consider the data in
172 hand and the process that led to the observation of the data. We assume that a researcher
173 has collected a sample of seawater—although soil, fecal, or other samples are essentially
174 equivalent—for the purposes of recovering eDNA from an ecological community. After
175 filtering the sample, extracting total DNA, and amplifying the DNA of interest using
176 oligonucleotide PCR primers, we observe counts of unique DNA sequences from a high-
177 throughput sequencer (e.g., Illumina, 454, Ion Torrent). Note that there are many possible
178 molecular methods by which the data can be derived. For all cases, though, the number of
179 observed DNA sequences for each type is a function of: 1) the true, but unknown, density
180 of DNA of each taxa present in the water, 2) the amount of DNA captured on the filter
181 and subsequent DNA extraction, 3) the primer set and its interaction with the DNA
182 sequence of each taxon present, 4) the number of PCR cycles performed, 5) the error rate
183 of the sequence analyzer, and myriad other factors. In short, the observed counts of DNA

184 sequences are a complicated stochastic realization of the true amount of DNA present in
185 the environment for each taxa. While eDNA protocols can be designed to minimize such
186 stochastic forces, they cannot be eliminated altogether. Defensible ecological inference
187 therefore depends upon identifying and estimating the parameters that may substantially
188 influence observed counts of DNA sequences.

189 By analogy with the net sampling example, the process by which biomass is
190 translated into DNA sequences matched to taxonomic groups is probabilistic (Fig. 1).
191 Specifically, the biomass of each taxon must be translated through several intermediate
192 states before it is observed as counts of DNA sequences. For taxon i , let W_i be the density
193 of DNA in the environment, X_i be the amount of DNA collected from the environment, Y_i
194 be the DNA present after DNA extraction, and Z_i be the DNA sequences recovered. We
195 acknowledge that there are other reasonable ways of parsing the process of generating
196 and making inference from eDNA (i.e. the framework we discuss here is extendable, and
197 additional states could be added to Fig. 1). However the latent states in Fig. 1 are intuitive
198 and, potentially, directly measurable with existing technologies.

199 As in eq. 1, the amount of biomass B_i estimated from eDNA sampling is the
200 product of four conditionally independent steps,¹

$$201 \quad [B_i|Z_i] = [B_i|W_i][W_i|X_i][X_i|Y_i][Y_i|Z_i] \quad (2)$$

202 Information about each link in this inferential chain is required to properly infer B_i from
203 observed counts of DNA sequences that emerge from a DNA sequencer Z_i . Such
204 information can be some combination of prior information about the processes
205 connecting these latent states, direct observations of the states, and biologically justified
206 assumptions about each component. There are two corollaries of this point: *i*) any
207 inferences made about B_i from eDNA make implicit and/or explicit assumptions about
208 the other components on the right side of eqn. 2; and, *ii*) if there is no information about
209 any of the components on the right side of eq. 2 (or researchers are unwilling or unable to
210 make assumptions about these components), it will be impossible to make inference
211 about B_i from eDNA observations alone. A parallel problem arises frequently in

¹ For the remainder of the manuscript, we let capital Roman letters denote random variables, lowercase roman letters denote realizations of random variables, and Greek letters denote parameters. Bold lowercase denote vectors and bold uppercase are matrices.

212 fisheries; biologists are unwilling to assert that the actual biomass is mirrored by
213 observed catches (i.e. the connection between B and D in Fig. 1 cannot be bridged).
214 Therefore survey catches are frequently used as indices of abundance not estimates of
215 absolute abundance (Kimura and Zenger 1997, Cotter and Pilling 2007). Despite not
216 reflecting actual abundance, such indices play a critical role in fisheries, wildlife
217 sciences, and management (Branch et al. 2010, Jannot and Holland 2013). The
218 formulation of eq. 2 also serves to point out where information is missing and to motivate
219 future research on poorly understood topics (Yoccoz 2012, Pedersen et al. 2015).

220 Other structures for Figure 1 are reasonable and we encourage investigators to
221 modify the chain of inference represented in Figure 1 to meet their specific sampling
222 needs. We view Figure 1 not as a rigid form for analyzing eDNA but as a framework
223 which can be modified to suit individual purposes and clarify thinking about the
224 inferences that can and cannot be drawn from available eDNA data. We expect improved
225 and more complex analytical structures to be developed for eDNA as the technology and
226 its use evolve.

227 An important point of Figure 1 is that the traditional sampling and eDNA arms of
228 the figure are only connected through the true biomass, B , represented at the top of the
229 figure. This structure serves to remind investigators that that directly comparing eDNA
230 and traditional sampling data is fraught with difficulty and can only be logically done
231 with a full sampling model for both how counts of OTUs observed from a sequencer (Z)
232 connect to biomass (B) and how traditional sampling observations connect to biomass.
233 Alternatively one could make strong assumptions about the connection between Z and B .
234 Indeed the most difficult step for both eDNA and traditional sampling in marine
235 environments is the first step in each pathway (between B and the density of DNA in the
236 environment W , and between B and collected individuals in a traditional sample, D ; Fig.
237 1). To date, we know of no eDNA study which has explicitly linked B and W under field
238 conditions and very few that have linked them under controlled laboratory conditions
239 (e.g. Takahara et al. 2012, Thomsen et al. 2012). To date, most researchers have either
240 asserted that the proportion of sequences observed from environmental samples mirror
241 the abundance (either count or biomass) of physically collected individuals or,
242 alternatively, concluded proportions of sequences are proportional to abundance based on

243 visual inspection (for example, see de Vargas et al. 2015, their Figs. W2B and W2C and
244 accompanying text). While these correlations may accurately reflect a functional link
245 between individuals in the environment and eDNA, we would point out that a complex
246 and diverse set of processes that separate D and Y mean that there are large number of
247 ways to arrive at spurious correlations between these two states. It is therefore desirable
248 to explicitly assess each link in the inferential chain linking observed DNA sequences to
249 biomass or some other biological/ecological parameter of interest.

250 While eq. 2 is instructive to broadly frame eDNA problems, the processes that
251 connect the latent states must be detailed to make this model useful in practice.
252 Specifically, the rates of transition between the states presented in Fig. 1 are controlled by
253 parameters that do not appear in eq. 2; we introduce those parameters here. Let θ_i be a set
254 of species-specific parameters associated with transition from B_i to W_i (e.g. DNA
255 shedding (Klymus et al. 2015, Iversen et al. *In Press*) and degradation (Thomsen et al.
256 2012, Strickler et al. 2015; Fig. 1), ϕ_i be taxon-specific parameters associated with
257 transition from W_i to X_i (e.g. the small scale patchiness of DNA in the water), ψ_i be
258 taxon-specific parameters associated with DNA filtering and extraction (the transition
259 from X_i to Y_i), and ξ_i define taxon-specific parameters associated with PCR amplification
260 and sequencing driving the transition from Y_i to Z_i (e.g. the match of a primer sequence to
261 the DNA input to template DNA sequence). Eq. 2 can be rewritten to include these
262 parameters for all taxa simultaneously,

$$263 \quad [B|Z, \theta, \phi, \psi, \xi] = [B|W, \theta][W|X, \phi][X|Y, \psi][Y|Z, \xi] \quad (3)$$

264 To connect these equations to empirical observations, they must be matched to
265 appropriate likelihood functions; we demonstrate in detail how to do so in the section
266 “*Statistical methods for community eDNA*” below.

267 It bears noting that the current state of knowledge with respect to eDNA limits our
268 ability to estimate all terms on the right-hand side of eq. 3, although at least some data are
269 available from which to begin such estimation. Here we briefly summarize the state of
270 knowledge with respect to each term in eq. 3 (Fig. 1).

- 271 1) *Processes in the transition from biomass, B , to DNA present in the environment, W (θ)*
- 272 • DNA shedding rates are positively correlated with biomass and influenced by diet
273 (Takahara et al. 2012, Kelly et al. 2014, Klymus et al. 2015, Evans et al. 2015) and

274 ambient eDNA density varies by species (Thomsen et al. 2012). Small DNA
275 fragments (ca. 100 base pairs) degrade within a few days in the marine environment
276 (Thomsen et al. 2012) but in some cases DNA signals are detectable for weeks to
277 months (Barnes et al. 2014, Strickler et al. 2015). DNA shedding and degradation
278 rates likely differ among taxa and among life-stages (Maruyama et al. 2014, Iversen
279 et al. *In Press*), though these differences are not well studied.

- 280 • In aquatic environments transported DNA does not appear to accumulate downstream
281 from the organism shedding it (Laramie et al. 2015) but rather remains at similar
282 concentrations downstream over short distances (Pilliod et al. 2014). DNA may be
283 moved over longer distances by bulk flow (Deiner and Altermatt 2014) or by mobile
284 predators that transport prey DNA in their gut and deposit it in their feces (Merkes et
285 al. 2014).

286 2) Processes in the transition from DNA in the environment, W , to DNA collected on a
287 filter, X (ϕ), and from DNA collected on a filter, X , to DNA present after extraction, Y
288 (ψ)

- 289 • Although methods for capturing eDNA influence the amount of useful sequence data,
290 they likely do not cause taxon-specific biases (Feinstein et al. 2009, Turner et al.
291 2014, Deiner et al. 2015). However, pre- and post-processing sample storage and
292 DNA extraction methods can produce taxon-specific biases (Carrigg et al. 2007,
293 Deiner et al. 2015).

294 3) Processes included in the transition from DNA present after extraction, Y , to DNA
295 present after sequencing, Z (ξ)

- 296 • PCR amplification of multi-taxon DNA samples introduces sequence-specific biases
297 due to differential primer binding strength (Lee et al. 2012); to a lesser degree the
298 number of PCR cycles may exacerbate these biases (Polz and Cavanaugh 1998, Sipos
299 et al. 2007).
- 300 • To improve cost efficiency by increasing sample throughput, a unique nucleotide
301 sequence (a “tag”) can be adjoined to the 5’ end of PCR primers. While these tags
302 allow multiple samples to be pooled for simultaneous (multiplex) sequencing, they
303 can introduce sequence-specific bias by changing primer binding strength (Berry et
304 al. 2011). In effect, these additions simply lengthen the primer sequence.

- 305 • High throughput sequencing platforms are thought to be relatively free from
306 sequence-specific biases, though low nucleotide diversity can degrade sequence
307 quality (Fadrosh et al. 2014). Further, the bioinformatic protocols used to process raw
308 sequence data can influence the inferred number of reads for a given taxon (Schloss et
309 al. 2011).
- 310 • Lastly, the taxonomic information DNA provides varies among loci, taxa, and
311 environments (Soergel et al. 2012), and nucleotide sequence repositories (e.g.
312 Genbank) are incomplete and both geographically and taxonomically biased
313 (Puillandre et al. 2009; Hijmans et al. 2000), limiting our ability to confidently
314 connect identified DNA sequences with specific taxa.

315 The above list is not a complete set of hurdles faced by eDNA methods and we expect
316 additional challenges will arise in the future. However, the model structure and logical
317 process of dividing the production of eDNA into conditionally independent processes is
318 general and broadly applicable to eDNA problems.

319

320 **Statistical methods for community eDNA**

321 As discussed above, methods for eDNA are not sufficiently well developed at
322 present to make full inference about density or biomass in an ecological community from
323 eDNA. Similar challenges confront estimation of density and biomass based on
324 traditional sampling methods (Burnham et al. 1980, Hankin and Reeves 1988, Kéry and
325 Royle 2010), but do not prevent researchers from making the best approximations
326 possible given existing knowledge and data. In this section we provide a statistical
327 framework for estimating the final term in eq. 3, $[Y|Z, \xi]$, in a community context. Once
328 we have an estimate of Y , if we can assume that the transitions from Y all the way to B do
329 not have taxon-specific biases, our approach allows statistically-justified inferences about
330 the relative abundance of taxa within a sampled community. As the processes related to
331 sampling eDNA become increasingly well understood, the other three terms in eq. 3 can
332 be modeled using a logic similar to the one detailed below.

333 For a sample of seawater that has been filtered, has had its total DNA extracted,
334 amplified by PCR, and has been processed by a high-throughput sequencer, our empirical
335 observations will be counts of unique DNA sequences. DNA sequences may be classified

336 into types on the basis of their similarity with respect to a user-specified threshold. These
 337 are most often referred to as operational taxonomic units (OTUs), and hereafter we refer
 338 to them as OTUs. For simplicity, we initially treat each unique DNA sequence observed
 339 as an OTU, and later discuss how to combine distinct OTUs into groups. The results of a
 340 single sequencing run can be written as $\mathbf{Z}=\mathbf{z}$, where \mathbf{z} is a realization of the random
 341 variable \mathbf{Z} and is vector of length I . Each entry in the vector, z_i , then contains the counts
 342 of the i^{th} OTU.

343 Using Bayes' theorem, we write the posterior probability of \mathbf{Y} , given our
 344 observations and parameters as proportional to the likelihood of the observations,
 345 $[\mathbf{Z} = \mathbf{z}|\mathbf{Y}, \xi]$, and the prior probability of the parameters $[\xi]$,

$$346 \quad [\mathbf{Y}|\mathbf{Z} = \mathbf{z}, \xi] \propto [\mathbf{Z} = \mathbf{z}|\mathbf{Y}, \xi][\xi] \quad (4)$$

347 A logical sampling model for counts with many possible categories is a multinomial
 348 model. We replace the general parameter notation ξ with $\boldsymbol{\pi} = \{\pi_1, \pi_2, \dots, \pi_I\}$ which
 349 represents the proportion of each OTU sequence present in the collected sample. Then we
 350 can write the likelihood as

$$351 \quad [\mathbf{Z} = \mathbf{z}|\mathbf{Y}, \boldsymbol{\pi}] \sim \text{Multinomial}(\boldsymbol{\pi}, n) \quad (5)$$

352 where n is the total number of DNA sequences observed. With a single sequencing run
 353 we have single set of observed counts, \mathbf{z} . However if we have M total observations of
 354 DNA sequences from a single DNA extraction – potentially from multiple independent
 355 PCR reactions or sequencing runs – we have,

$$356 \quad [\mathbf{Z} = \mathbf{z}_1, \dots, \mathbf{z}_M|\mathbf{Y}, \boldsymbol{\pi}] \sim \text{Multinomial}(\boldsymbol{\pi}, n_1, n_2, \dots, n_M) \quad (6)$$

357 This equation states that each \mathbf{z} is a sample from a shared process (i.e. there is a single
 358 true proportion of DNA from each taxon in the eDNA sample and we have M
 359 observations of this process). Variation among the observations of \mathbf{z} can be attributed to
 360 stochastic processes occurring during PCR and sequencing, and the model described here
 361 can be generalized to include these effects as individually modeled parameters if desired.

362 Because a multinomial distribution can be written as a combination of
 363 independent Poisson distributions (the multinomial-Poisson transformation; Baker 1994),
 364 it is convenient to write the number of sequenced DNA fragments observed for each
 365 OTU as an independent Poisson random variable,

366
$$\begin{aligned} z_{im} &\sim \text{Poisson}(e^{\lambda_i}) \\ \lambda_i &= \beta_i \end{aligned} \quad (7)$$

367 Here, β_i , the OTU-specific fixed effects, and λ_i are identical, but we use this notation to
 368 allow later elaboration in circumstances where additional processes are thought to
 369 influence λ_i . The proportion of DNA associated with each OTU can then be found by
 370 calculating

371
$$\pi_i = \frac{e^{\beta_i}}{\sum_i e^{\beta_i}} \quad (8)$$

372 Note that Eq. 6 provides identical inference to eqs. 7 and 8 (Baker 1994).

373 The model formulation in eq. 7 assumes that each observed DNA fragment is
 374 sampled independently from a multinomial distribution. Due to the compounding process
 375 of sequential amplification in PCR, counts of DNA sequences from a sequencer are not
 376 truly independent observations of the extracted DNA. One method to deal with such non-
 377 independence is to allow for overdispersion in the Poisson parameter λ . With $m =$
 378 $1, 2, \dots, M$ replicate observations, we can write the observed species counts as an over-
 379 dispersed Poisson and estimate the amount of over-dispersion, σ^2 ,

380
$$\begin{aligned} z_{im} &\sim \text{Poisson}(e^{\lambda_{im}}) \\ \lambda_{im} &= \beta_i + \eta_{im} \\ \eta_{im} &\sim N(0, \sigma^2) \end{aligned} \quad (9)$$

381 This is a simple random effects model, but one that allows great flexibility in modeling
 382 count data. Note that in the case that only a single OTU is present, eq. 9 simplifies to a
 383 log-linear model of the DNA counts and thus the single OTU version of this model is
 384 appropriate for qPCR data. When we observe more than one OTU, we can still produce
 385 estimates of the proportion of DNA from each taxa across all of our observations (eq. 8).
 386 After specifying prior parameters, we can use standard Bayesian Markov chain Monte
 387 Carlo (MCMC) methods to estimate the model and provide uncertainty bounds (Gelman
 388 et al. 2003). Likelihood optimization methods are also available. A further benefit of the
 389 structure is the possibility of multiple random effects that can represent multiple sources
 390 of variation in the observed counts. We present a more complicated example in the online
 391 supplement. We note that the above model is similar to other models for sequencing data
 392 proposed in a different context for other applications (Love et al. 2014).

393

394 *Addressing primer bias: a framework and a simulated example*

395 Equation 9 implicitly makes assumptions that eDNA data almost certainly violate.
396 Importantly, eq. 9 assumes all OTUs present in the DNA extraction will be amplified
397 equally well by PCR, and will subsequently appear in the count data emerging from the
398 sequencer, yet PCR primers are intentionally designed to amplify specific taxonomic
399 groups (e.g. vertebrates) to the exclusion of others (e.g., Riaz et al. 2011). Even within a
400 target group of taxa, intra-group genetic variability in the primer binding site can cause
401 variation in template-primer mismatch, resulting in unequal amplification among
402 templates and thus bias in the observed sequences (e.g. Hong et al. 2009). Estimating the
403 extent of amplification bias due to this interaction requires detailed information about
404 both the primer set and the template (target) sequence for all taxa of interest. Generally, a
405 way to incorporate a series of covariates—such as would describe these OTU-specific
406 effects—is to construct a matrix of covariates, \mathbf{H} , and estimated coefficients, $\boldsymbol{\gamma}$, given
407 available information about primer mismatches with existing sequence data from target
408 taxa. Accordingly, the second line of eq. 9 can be modified to accommodate variation in
409 primer specificity to become:

$$410 \quad \lambda_{im} = \beta_i + \boldsymbol{\gamma}\mathbf{H}_{im} + \eta_{im} \quad (10)$$

411 where $\boldsymbol{\gamma}$ defines how covariates shared across taxa (e.g. the quality of match between the
412 primer and taxa DNA) will affect the observed number of DNA sequences for each
413 taxon. Also, note that the researcher-specified design matrix \mathbf{H} includes the subscript m .
414 This indicates multiple PCR or sequencing runs conducted using distinct methods on a
415 single sample can be used jointly to improve the reconstruction of the ecological
416 community of interest. For example, if two or more independent analyses were carried
417 out on the same DNA extraction—such as in the case of multi-locus eDNA studies—the
418 results could be formally combined into a single analysis. Furthermore, such
419 methodological variation will help inform how changing primer specificity, the PCR
420 reaction parameters, or other methods affect the inference about the proportion of DNA
421 associated with each OTU. We illustrate an application of these methods below in
422 “*Understanding marine invertebrate communities using eDNA*”.

423 To illustrate the potential consequences of the effect of primer-template mismatch
424 on estimates of OTU composition, we simulated small changes to the quality of primer

425 match and used estimates of γ to show how they affected estimates in a simple three-
426 species community (Fig. 2, supplementary materials). Simulations show that a change in
427 primer-template match of as little as 5% (e.g. a 3 base-pair difference between a 60 bp
428 long template and the combined forward and reverse primer) can change estimates of
429 relative abundance (Fig. 2). The most important point of Fig. 2 is that because the
430 estimates are relative proportions that must sum to one, if one taxon has a biased
431 estimate, all of the other taxa's estimates are biased as well. A consequence of this
432 observation is that analyzing data derived from multi-species primers on a species-by-
433 species basis (i.e. treating the number of reads for each taxa independently in later
434 analyses) is likely to decrease statistical precision and introduce bias in the relationship
435 between the number of reads and virtually any other variable.

436

437 *Estimating the absolute concentration of DNA in an extraction*

438 Thus far, we have not provided direct estimates of the concentration of template
439 DNA in the sample, Y , but only estimates of the proportional abundance of each OTU, $\boldsymbol{\pi}$.
440 To generate estimates of DNA concentration, we need to incorporate additional
441 information about the absolute abundance of DNA from at least some of the OTUs to
442 scale the proportional abundance to true abundance. We can use the posterior estimates of
443 proportional abundance $\boldsymbol{\pi}$ in combination with posterior estimates of the density of DNA
444 from a single OTU, ω_1 , to scale the proportions to DNA densities for all OTUs. Current
445 methods using qPCR are adept at producing estimates of ω_1 (Jerde et al. 2011, Lodge et
446 al. 2012, Takahara et al. 2012). If we assume that ω_1 and $\boldsymbol{\pi}$ are derived from independent
447 methods, we can use draws from the posterior distributions of each to derive the posterior
448 distribution of Y . For the j^{th} OTU and g^{th} draw from the posterior distribution, we have

$$449 \quad Y_j^{(g)} = \omega_1^{(g)} \left(\frac{\pi_j^{(g)}}{\pi_1^{(g)}} \right) \quad (11)$$

450 After calculating Y for a large number of posterior draws, we can summarize Y using
451 standard descriptors (mean, standard deviation, etc.). This method is appealing because it
452 reflects the uncertainty in both $\boldsymbol{\pi}$ and the concentration of DNA derived from qPCR. It
453 also shows how qPCR and sequencing approaches are complementary data types that can
454 be combined and re-emphasizes how the structure presented in Figure 1 is applicable to a

455 wide variety of eDNA methods. We highlight the utility of this two-pronged validation
456 method for future applications.

457

458 *Detection probabilities and power analysis*

459 A trade-off between detection probability for any given taxon and breadth of the
460 community observed is common to surveys using both eDNA and non-molecular (i.e.,
461 traditional) methods. In many eDNA applications, the risk of false-negative detections (in
462 which a taxon is present, but not detected) is one of the most pressing issues (Yoccoz
463 2012, Yu et al. 2012, Ji et al. 2014.). Conveniently, the model outlined in eqs. 9 and 10
464 provides a method for determining the thresholds for detection. However, because the
465 PCR primers for community eDNA analyses will almost never be strictly taxon-specific,
466 the power analysis cannot be determined on a single-taxon basis but must always be
467 phrased in terms of a larger DNA community that is “observed” by a given PCR protocol.

468 The relative abundance of an arbitrary OTU, taxon “A”, can be fully defined by
469 four quantities: the true relative proportion of DNA from OTU A in the sample π_A ; the
470 estimated effect of covariates for that OTU, $\boldsymbol{\gamma}H_A$; the total number of DNA sequences
471 observed, n ; and the stochasticity in the PCR and sequencing process, σ^2 . Because for
472 the observed data, $n = \sum_i e^{\lambda_i}$ (eq. 8), we can combine eq. 8 and 10 and use the properties
473 of the log-normal distribution to show that for any true value of π_A , the median value of
474 λ_A , λ_A^* , will be

$$475 \quad \lambda_A^* = \log(\pi_A) + \log(n) + \boldsymbol{\gamma}H_A \quad (12)$$

476 Using the probability mass function of the Poisson distribution, the probability that the
477 observed number of DNA sequences for OTU A will exceed 0 at λ_A^* is,

$$478 \quad p(z_A > 0) = 1 - e^{-\lambda_A^*} \quad (13)$$

479 In this way, the detection probability can be approximated for a given primer, the number
480 of DNA sequences observed, and DNA community. This type of power analysis based on
481 the median estimate is likely sufficient for most applications, but it is important to
482 acknowledge that this approach ignores variability in the PCR process (σ^2) and
483 uncertainty in the estimate of $\boldsymbol{\gamma}$. However, simulation approaches could incorporate this
484 variability if desired. Importantly, eqs. 12 and 13 make explicit that analytical approaches
485 based on the occurrence data (Yu et al. 2012, Ji et al. 2014) are special cases of multi-

486 taxa count data. In its simplest form, occurrence data is simply the count data for each
487 OTU converted into two classes: $z_i = 0$ and $z_i > 0$. Other investigators have suggested
488 that OTUs below a certain threshold abundance should be excluded (e.g. OTUs below
489 0.005% of the total number of DNA reads is recommended by Bokulich et al. 2013).
490 Regardless of the exact cutoff used, this section demonstrates that the same biases that
491 plague estimating abundance from eDNA will also plague estimations of occurrence –
492 though signatures of bias will be more difficult to detect and estimate using occurrence
493 data.

494 We illustrate power curves in Fig. 3 to provide a graphical method for
495 understanding the detection probability of a taxon for a given primer, extracted DNA, and
496 number of DNA reads. Specifically, we compare three values of a single covariate
497 representing the match between the primer and taxon A 's DNA. $H_A = 0$ represents the
498 average match between the primer and the taxa observed in the sample, while $H_A = -0.15$
499 corresponds to A having a 15% better match to primer than average and $H_A = 0.15$
500 corresponds to A having a 15% worse match to primer than average (e.g. for a 20
501 basepair primer, 15% corresponds to a change of 3 basepair matches between primer and
502 template). For all three simulations, we used a slope parameter that reflect real-world
503 estimates of primer bias discussed below in “Understanding marine invertebrate
504 communities using eDNA” ($\gamma = -14$). An important result of Fig. 3 is that even when a
505 taxon is present in a sample, it may not be observed in the DNA counts emerging from
506 the sequencer. The probability of observing at least one instance of taxon A is
507 affected both by its true abundance (relative to other species amplified by the PCR
508 product) and the match between the DNA sequence and the PCR primer used.

509 Eq. 12 and Fig. 3 suggest that there are several intuitive and non-mutually
510 exclusive methods for increasing detection probability of a particular taxon: 1) increase
511 the number of sequences observed for each PCR (increase n); 2) decrease the number of
512 taxa amplified by the primer (decrease I and thereby increase the relative abundance of
513 the OTU of interest, π_A); 3) improve the efficiency of the primer for taxon A relative to
514 other taxon in the DNA community (i.e. modify H_A). In practice, a PCR primer that more
515 closely matches a particular taxon will likely contribute to both point 2 and 3. However,
516 increased primer specificity will always reduce the diversity of taxa detected in a single

517 sequencing run. Both highly specific and more general primers have important real world
518 applications (Simmons et al. *In Press*).

519

520 *Combining unique DNA sequences into biologically meaningful groups*

521 Genetic variation among individuals both within and across taxa can result in two
522 problematic scenarios: 1) high diversity within a taxon will result in it being represented
523 by more than one OTU in the sequence data or 2) low diversity across taxa will result in
524 many taxa being represented by a single OTU. An ideal PCR primer would target a locus
525 with high inter-taxon diversity and low intra-taxon diversity. Unfortunately, we know of
526 no such locus that can be used for a broad swath of taxa. For the case where a single
527 taxon is represented by multiple OTUs, we describe two approaches for obtaining
528 abundance estimates.

529 The first is to estimate the model treating each OTU separately (eq. 12), and
530 combine the output of the estimation procedure. Because each iteration of a Markov
531 chain provides a draw from the posterior distribution of the parameters, the draws can
532 simply be added together for the OTUs of interest, and the proportion of the resulting
533 taxon recalculated (Shelton et al. 2012). To provide a concrete, but fictitious, example,
534 suppose that OTU *A* and OTU *B* were both observed in a sequencing run. Both OTUs are
535 subsequently determined to represent unique sequences from woolley mammoth
536 (*Mammuthus primigenius*) and need to be combined to provide an estimate of the total
537 mammoth present in the extracted DNA sample. After estimating a Poisson model (e.g.
538 eq. 10) we can simply add the two estimated parameters for OTU *A* and OTU *B*
539 (β_A and β_B , respectively) such that $\beta_{mammoth} = \beta_A + \beta_B$ for each MCMC iteration. The
540 proportion of DNA attributable to mammoth would then be $\pi_{mammoth} = \frac{e^{\beta_{mammoth}}}{\sum_i e^{\beta_i}}$.

541 Using draws from the posterior distribution maintains the correlation structure and
542 uncertainty bounds of the proportional occurrence. However, this approach has the
543 downside of requiring parameter estimates and the collection of covariates to populate \mathbf{H}
544 for each OTU, slowing computation speed if there are large numbers of OTUs.

545 The second option is to group the OTUs into broader taxonomic groups before
546 they are included as input data for the model estimation. While the choice of method for
547 clustering sequence data into OTU counts is of general concern (Edgar et al. 2011, Yu et

548 al. 2012), this approach also requires that all OTUs within a group be assumed to have
549 shared covariates related to PCR. Continuing our previous mammoth example, the
550 primer-template mismatch might differ between OTU A and OTU B, and yet if their
551 counts were to be combined, information about their distinct matching characteristics
552 could not be directly incorporated in the model. A summary statistic such as the median
553 dissimilarity would have to be used instead. Depending on the details of the primers and
554 match quality, such averaging across covariates may or may not substantially influence
555 the result. Given these considerations, we advocate the first approach of combining taxa
556 after model estimation, unless speed is favored over accuracy or researchers are
557 sufficiently confident that grouped taxa do not differ in PCR or sequencing efficiency.

558

559 **Understanding marine invertebrate communities using eDNA**

560 To illustrate the utility of our statistical framework, we apply the above methods
561 to eDNA isolated, amplified, and sequenced from eleven, 1-L seawater samples collected
562 from a single location in Puget Sound, WA on June 26, 27, and 29, 2014 (Carkeek Park,
563 Seattle, WA, USA; 47°42'40.44"N, 122°22'20.10"W). Because we use this empirical
564 dataset here only to illustrate the application of statistical methods to counts of DNA
565 sequences emerging from a high throughput sequencer, we only outline the methods that
566 affect the statistical estimation. We provide detailed molecular protocols in the online
567 supplement for interested readers.

568 *Summary of molecular methods*

569 To test the effect of primer mismatch on template-specific PCR efficiency, we
570 amplified each environmental sample using two different sets of primers, which in each
571 direction shared a common core 22bp region targeting the 16S region of the
572 mitochondrion, but differed by an index sequence on the 5' end (see Table S1 for the
573 primer sequences used). These index sequences have been demonstrated to cause
574 differential amplification efficiency among template DNA in a mixed-template PCR
575 (Berry et al. 2011), and thus provide an opportunity to test the efficacy of our framework
576 for estimating biomass and uncertainty in the face of bias. PCR, library preparation,
577 sequencing, and bioinformatics protocols are described in the supplementary material.

578 The experimental design yielded sequence data from six PCR products per
579 environmental sample: three sequencing replicates arising from each of two distinct
580 primer sets. In total, we observed over 10.5 million individual DNA reads representing
581 27,973 unique OTUs. For the purpose of this example, we model only 9 of the most
582 common OTUs and focus on estimating the proportional DNA contribution for these 9
583 OTUs and a tenth “Other” category which encompasses all remaining OTUs. We
584 investigate only 10 OTUs for illustration purposes, though this approach is directly
585 applicable to a much larger set of OTUs. We present the raw data and models for
586 estimating these models for these nine OTUs in the supplementary materials.

587 *Statistical modeling of OTU counts*

588 To estimate the proportion of each of these 9 OTUs on each sampling occasion,
589 we use a version of eq. 10 that adds a subscript t corresponding to each sample time and
590 includes m observed DNA replicates for each time. Then the full model is

$$\begin{aligned} z_{itm} &\sim \text{Poisson}(e^{\lambda_{itm}}) \\ \lambda_{itm} &= \beta_{it} + \boldsymbol{\gamma} \mathbf{H}_{itm} + \eta_{itm} \\ \eta_{itm} &\sim N(0, \sigma^2) \end{aligned} \quad (14)$$

592 Again, β_{it} indicates the count of OTU i at time t , $\boldsymbol{\gamma} \mathbf{H}_{itm}$ controls the fixed effect of PCR
593 and sequencing bias on the observed number of OTU counts for each replicate, with $\boldsymbol{\gamma}$
594 estimated regression coefficients and the covariate matrix \mathbf{H}_{itm} supplied by the
595 investigator on the basis of available information about target-taxon sequences in the
596 primer region. Finally, η_{itm} provides for additional error not accounted for by either the
597 fixed taxon effect β_{it} or the other fixed effects. While it is possible to include a large
598 variety of potential covariates in $\boldsymbol{\gamma} \mathbf{H}_{itm}$ for illustration purposes we include only a single
599 covariate, the total genetic distance between the OTUs’ primer binding sites and the
600 primers, γ , at both forward and reverse priming sites. Thus \mathbf{H} is a design matrix with a
601 single column corresponding to the proportion of nucleotide mismatches between the
602 primers and each template (OTU primer binding site). A value of 0 would indicate no
603 difference between the primer and the template, while 0.10 would indicate 10% of base
604 pairs do not match between the primer and the OTU. Distance calculations were
605 performed using the function `dist.dna` in the R package `ape` (Paradis et al. 2004). To
606 derive estimates of the design matrix \mathbf{H} we assessed the quality of match between the

607 primer and each taxon's DNA. For the nine focal OTUs in the dataset, we performed a
608 BLAST search of NCBI's nucleotide database (GenBank) to identify the likely sequence
609 of the primer binding sites given existing sequence information for taxa in GenBank
610 matching the OTU sequences (see below). We centered the covariate values in \mathbf{H} before
611 analysis by subtracting each value by the average across all OTU-primer pairs. The
612 process of centering makes β_{it} the intercept for each OTU in this generalized linear
613 model. We assumed the "Other" category had a covariate value of 0, (i.e. $H_{Other,t,m} = 0$)
614 corresponding to the average amplification value of the "Other" category. Centering the
615 covariates also means that when we calculate the proportional contribution of each OTU,
616 we can calculate the proportion of each OTU in the sample as $\pi_{it} = \frac{e^{\beta_{it}}}{\sum_i e^{\beta_{it}}}$. This produces
617 estimates of proportional composition of each OTU at a standardized match between the
618 primer and substrate for all OTUs.

619 We estimated eq. 14, using Just Another Gibbs Sampler (JAGS; Plummer 2003)
620 implemented in R (R Core Team 2014) using the R2jags package (Su and Yajima 2014).
621 We used non-informative prior distributions for each parameter. Specifically we let
622 $\gamma \sim Normal(0,1000)$, $\beta_- \sim Normal(0,1000)$, and $\sigma^2 \sim Uniform(0,1000)$. We ran three
623 replicate MCMC chains using a 100,000 iteration burn-in and 10,000 monitoring
624 iterations. We confirmed appropriate model mixing and convergence using visual
625 inspection of trace plots and Gelman-Rubin diagnostics as implemented by the R package
626 "coda" (Plummer et al. 2006).

627

628 *Results*

629 Using eq. 14, we estimated the proportional composition for nine focal OTUs and
630 the "Other" category for all eleven time periods (Fig. 4, Fig. 5). Our model estimated a
631 large amount of overdispersion in the observed count data ($\sigma^2 = 8.34[0.68]$; posterior
632 mean[sd]) indicating that there remains a substantial effect of unknown and unmodeled
633 factors on variation among samples. The large estimated overdispersion translates into
634 large uncertainty in the estimated proportional composition (Fig. 4). Our estimates are
635 statistically well-justified and reflect the uncertainty present in our observations, but
636 suggest that methodological improvements will be required to provide more precise

637 estimates of the marine community. Across all times, OTUs 3, 5, and 7 were particularly
638 frequent. Both OTU 3 and 7 correspond to the mussel, *Mytilus trossulus*, while OTU 5
639 corresponds to acorn barnacles (suborder Balanomorpha; likely *Balanus glandula*), both
640 of which are among the most commonly observed species at our study site. We found no
641 dramatic patterns of OTU relative abundance over time or with respect to an important
642 covariate, tidal height (Fig. 5). However, the large degree of uncertainty limits our power
643 to detect strong effects of time or environmental factors.

644 Among our nine focal OTUs—which, again, represent sequences amplified and
645 recovered from environmental samples—the variance in primer-template mismatch was
646 substantial. Across all primer-template pairs the mean proportional mismatch was 0.193
647 (range: 0.11-0.28), indicating that, on average 10.81 out of a total 56 base pairs were
648 mismatched. We estimated, as expected, that the effect of decreasing match between the
649 primer and substrate was strongly negative, $\gamma = -14.3[6.11]$ (posterior mean[sd])
650 indicating OTUs with a poor match between the primer binding site and primer were
651 underrepresented in the observed DNA counts. Our estimated effect of primer quality is
652 similar to experimental results exploring the effect of primer mismatch on preferential
653 PCR amplification (Polz 1998, Sipos 2007, Wright et al. 2014). We emphasize that there
654 are a great many possible other covariates that could be used in this type of analysis.

655

656 **Discussion and conclusions**

657 eDNA is an exciting emerging method for describing ecological communities.
658 Given the enormous potential for eDNA applications in the environmental sciences,
659 recent reviews of eDNA methods have stressed the need for improved molecular and
660 statistical techniques for eDNA (Yu et al. 2012, Yoccoz et al. 2012, Schmidt et al. 2013,
661 Ji et al. 2014). Conceptually, the challenges posed by eDNA are largely analogous to
662 those faced by traditional sampling techniques (Fig. 1). Both conventional and eDNA
663 sampling ultimately attempt to make inferences about the same quantity: the biomass or
664 density of each species in the environment. It should also be clear that traditional
665 sampling methods suffer from a parallel set of sampling problems to eDNA and, as noted
666 earlier, our current inability to estimate abundance or biomass from eDNA samples alone
667 is not a fatal flaw for eDNA data. A specific topic that deserves special consideration in

668 future work is understanding the spatial and temporal spread of eDNA under natural
669 conditions and how the scale of inference from eDNA sampling matches (or, potentially,
670 does not match) the spatial and temporal inference available from traditional sampling
671 methods.

672 While we have framed our analysis in terms of biomass, we note that an
673 equivalent structure is necessary for estimation of count data and for deriving most
674 community metrics of interest as well. Estimated species richness is the number of
675 species with biomass greater than 0 while Shannon diversity is species richness weighted
676 by the relative biomass (or count) of each species. Both richness and Shannon diversity –
677 and indeed virtually all community and diversity metrics – are directly derived from
678 estimates of occurrence and abundance of individual species. Thus this framework
679 provides a pathway for investigating communities as well as individual taxa.

680 In closing we offer a few recommendations to ensure that eDNA study designs—
681 and the resulting datasets—are adequate to develop a meaningful estimate of the target
682 biological community structure.

683 Foremost, it should be clear from the framework we discuss here that sample
684 replication (in space, time, laboratory treatment, etc.) is critical to partitioning variance
685 among steps in the eDNA analytical chain. Because real-world constraints on time and
686 funding generally prevent replication at every step, we emphasize that replication is most
687 important at the step or steps that are likely to introduce the greatest amount of variance
688 or where the variance attributable to that step is of special interest. For example, if one
689 has data demonstrating that eDNA capture, extraction, and sequencing are likely to
690 introduce little systematic bias, but that PCR primer choice has an unknown and
691 potentially large effect, PCR is the most important target for replication and independent
692 analysis. Samples treated separately can then subsequently be combined using
693 hierarchical models, where this would provide analytical benefit (see online supplement).
694 Note additionally that we advocate avoiding pooling samples and then running analyses
695 on the pooled output whenever possible; there is information in the variability among
696 replicated outputs of molecular methods.

697 Second, because taxa are not equally abundant in a sampled environment, and
698 because taxa are not equally likely to amplify with a given set of PCR primers, eDNA

699 community surveys are necessarily an uneven reflection of taxa present, even for a
700 specifically targeted groups. The same issues arise with traditional sampling methods, as
701 alternative survey methods have different but non-negligible selectivity issues (Beverton
702 and Holt 1957 section 8, Arreguín-Sánchez 1996, Venebles and Dichmont 2004).

703 The methods we present for community eDNA data offer the ability to correct for
704 attributable biases and to be statistically honest about biases and variability that we do not
705 understand. However, real differences in DNA abundance and susceptibility to
706 amplification mean that for any given set of PCR primers there is a limited set of taxa
707 that can successfully be detected. This observation gives rise to three recommendations:

- 708 1. Using multiple markers offers the chance to broaden the scope of an eDNA survey
709 and to generate mutually reinforcing datasets that might be combined in the
710 framework we present here (Evans et al. 2015).
- 711 2. Community surveys that focus on the most common sequences generated—rather
712 than on the rare sequence “tails”—are more likely to be repeatable and robust to
713 statistical inference. At the same time, we acknowledge that some analyses –
714 particularly those focused on measures of biodiversity (e.g. Ji et al. 2014) - are
715 intrinsically interested in rare taxa. We think an increased focus on understanding
716 how the probability of detection may affect diversity estimates is an important area
717 for further research (Fig. 2; Schmidt et al. 2013).
- 718 3. Finally, a focus on the most common species (or most common sequences) found in
719 an environment has implications for primer design. Rather than accepting a very
720 broad set of sequence constraints on primer design (e.g., all metazoans), ensuring that
721 primers are likely to be good matches for the few dozen most common target species
722 in the sampled area is likely to yield a better range of acceptable primer sequences.
723 Increased specificity is more likely to lead to the intended results of a community
724 eDNA survey. Again, this approach is appropriate only when the interest is focused
725 on relatively common species, not on rare or unknown taxa in the community.

726 As we have suggested throughout this paper, we believe there is ample room for
727 cross-pollination between eDNA, both qPCR and sequencing based, and traditional
728 sampling approaches. Notably, the conceptual framework we outline suggests that it is
729 possible to construct models that jointly model data from traditional and eDNA sampling

730 to draw inference about natural populations. We also expect that methodological biases
731 inherent to eDNA and traditional sampling may often produce complementary, rather
732 than overlapping, estimates of community composition. Regardless, here we have shown
733 how to start toward this ultimate goal by providing a framework and detailed statistical
734 models for a particularly challenging aspect of eDNA work—calculating the relative
735 abundance of DNA from multi-species primers while accounting for variation in PCR.
736 However, multiple elements of the eDNA processing chain remain poorly described from
737 a quantitative perspective, and as future work clarifies biases introduced at each
738 experimental step, our framework provides a means of using such emerging information
739 to improve quantitative estimates of community biomass from eDNA.

740

741 *Acknowledgements*

742 This work was supported in part by a grant to RPK from the David and Lucile Packard
743 Foundation. Thanks also to the Helen R. Whiteley Center at Friday Harbor Laboratories
744 for supporting the writing workshop that substantially advanced this product. We thank
745 E. Buckner, E. Garrison, M. Klein, and A. Wong for help with field collections and
746 Seattle Parks and Recreation for access to Carkeek Park. We thank P. Levin, A. Stier, B.
747 Feist, K. Marshall, and S. Hennessey for comments on earlier versions of this manuscript.
748 K. Deiner and three anonymous reviewers improved this manuscript.

749 *Literature Cited*

- 750 Arreguín-Sánchez, F. 1996. Catchability: a key parameter for fish stock assessment.
751 *Reviews in Fish Biology and Fisheries* 6:221–242.
- 752 Baker, S.G. 1994. The multinomial-Poisson transformation. *The Statistician* 43:495.
- 753 Barnes, M.A., C.R. Turner, C.L. Jerde, M.A. Renshaw, W.L. Chadderton, and D.M.
754 Lodge. 2014. Environmental conditions influence eDNA persistence in aquatic
755 systems. *Environmental Science and Technology* 48:1819–1827.
- 756 Berry, D., K.B. Mahfoudh, M. Wagner, and A. Loy. 2011. Barcoded primers used in
757 multiplex amplicon pyrosequencing bias amplification. *Applied and Environmental*
758 *Microbiology* 77:7846–7849.
- 759 Beverton, R. J. H., and S. J. Holt. 1957. *On the dynamics of exploited fish populations.*
760 Chapman and Hall.

761 Bokulich, N.A., S. Subramanian, J.J. Faith, D. Gevers, J.I. Gordon, R. Knight, D.A.
762 Mills, and J.G. Caporaso. 2013. Quality-filtering vastly improves diversity estimates
763 from Illumina amplicon sequencing. *Nature Methods* 10:57–59.

764 Branch, T.A., R. Watson, E.A. Fulton, S. Jennings, C.R. McGilliard, G.T. Pablico, D.
765 Ricard, S.R. Tracey. 2010. The trophic fingerprint of marine
766 fisheries. *Nature* 468:431-435.

767 Burnham, K.P., D.R. Anderson, and J.L. Laake. 1980. Estimation of density from line
768 transect sampling of biological populations. *Wildlife Monographs* 3–202.

769 Camacho C., G. Coulouris, V. Avagyan, N. Ma, J. Papadopoulos, K. Bealer, and T.L.
770 Madden. 2009. BLAST+: architecture and applications. *BMC bioinformatics* 10:421.

771 Carrigg, C., O. Rice, S. Kavanagh, G. Collins, and V. O’Flaherty. 2007. DNA extraction
772 method affects microbial community profiles from soils and sediment. *Applied*
773 *Microbiology and Biotechnology* 77:955–964.

774 Clark, J.S. 2007. *Models for ecological data*. Princeton, Princeton, NJ.

775 Cochran, W.G. 1977. *Sampling techniques*. John Wiley and Sons Inc.

776 Cotter, A.J.R., and G.M. Pilling. 2007. Landings, logbooks and observer surveys:
777 improving the protocols for sampling commercial fisheries. *Fish and Fisheries*
778 8:123–152.

779 Cressie, N., and C.K. Wikle. 2011. *Statistics for spatio-temporal data*. John Wiley and
780 Sons Inc.

781 de Vargas, C., S. Audic, N. Henry, J. Decelle, F. Mahé, R. Logares, et al. 2015.
782 Eukaryotic plankton diversity in the sunlit ocean. *Science*, 348: 1261605.

783 Deiner, K. and F. Altermatt F. 2014. Transport distance of invertebrate environmental
784 DNA in a natural river. *PLoS ONE* 9: e88786. doi:10.1371/journal.pone.0088786

785 Deiner, K., J.-C. Walser, E. Mächler, and F. Altermatt. 2015. Choice of capture and
786 extraction methods affect detection of freshwater biodiversity from environmental
787 DNA. *Biological Conservation* 183:53–63.

788 Edgar R.C. 2010. Search and clustering orders of magnitude faster than BLAST.
789 *Bioinformatics* 26:2460-2461.

790 Edgar, R.C., B.J. Haas, J.C. Clemente, C. Quince, and R. Knight. 2011. UCHIME
791 improves sensitivity and speed of chimera detection. *Bioinformatics* 27:2194–2200.

792 Elith, J. and J.R. Leathwick. 2009. Species distribution models: ecological explanation
793 and prediction across space and time. *Annual Review of Ecology, Evolution, and*
794 *Systematics* 40:677-697.

795 Evans, N.T., B.P. Olds, M.A. Renshaw, C.R. Turner, Y. Li, C.J. Jerde, A.R. Mahon, M.E.
796 Pfrender, G.A. Lamberti, and D.M. Lodge. 2015. Quantification of mesocosm fish
797 and amphibian species diversity via environmental DNA metabarcoding. *Molecular*
798 *Ecology Resources*. doi: 10.1111/1755-0998.12433

799 Fadrosch, D.W., B. Ma, P. Gajer, N. Sengamalay, S. Ott, R.M. Brotman, and J. Ravel.
800 2014. An improved dual-indexing approach for multiplexed 16S rRNA gene
801 sequencing on the Illumina MiSeq platform. *Microbiome* 2:1–7.

802 Feinstein, L.M., W.J. Sul, and C.B. Blackwood. 2009 Assessment of bias associated with
803 incomplete extraction of microbial DNA from soil. *Applied Environmental*
804 *Microbiology*. 75:5428–5433.

805 Hankin, D.G. and G.H. Reeves. 1988. Estimating total fish abundance and total habitat
806 area in small streams based on visual estimation methods. *Canadian Journal of*
807 *Fisheries and Aquatic Sciences* 45:834–844.

808 Hijmans, R.J., K.A. Garrett, Z. Huamán, D.P. Zhang, M. Schreuder, and M. Bonierbale.
809 2000. Assessing the geographic representativeness of genebank collections: the case
810 of Bolivian wild potatoes. *Conservation Biology* 14:1755–1765.

811 Iversen, L.L., J. Kielgast, and K. Sand-Jensen. *In press*. Monitoring of animal
812 abundance by environmental DNA- An increasingly obscure perspective: A reply
813 to Klymus et al., 2015. *Biological Conservation*.
814 <http://dx.doi.org/10.1016/j.biocon.2015.09.024>

815 Ficetola, G.F., J. Pansu, A. Bonin, E. Coissac, C. Giguët-Covex, M. De Barba, L. Gielly,
816 C.M. Lopes, F. Boyer, F. Pompanon, G. Rayé and P. Taberlet. 2015. Replication
817 levels, false presences and the estimation of the presence/absence from eDNA
818 metabarcoding data. *Molecular Ecology Resources*. 15: 543–556

819 Gelman, A., J.B. Carlin, H.S. Stern, and D.B. Rubin. 2003. *Bayesian data Analysis*.
820 Second edition. Chapman and Hall/CRC.

821 Hong S., J. Bunge, C. Leslin, S. Jeon, and S.S. Epstein. 2009. Polymerase chain reaction
822 primers miss half of rRNA microbial diversity. *ISME* 3: 1365–73.

823 Jannot, J.E. and D.S. Holland 2013. Identifying ecological and fishing drivers of bycatch
824 in a U.S. groundfish fishery. *Ecological Applications* 23:1645–1658.

825 Jerde, C.L. A.R. Mahon, W.L. Chadderton, and D.M. Lodge 2011. ‘Sight-unseen’
826 detection of rare aquatic species using environmental DNA. *Conservation Letters*
827 4:150-157.

828 Ji, Y., L. Ashton, S.M. Pedley, D.P. Edwards, Y. Tang, A. Nakamura, R. Kitching, P.M.
829 Dolman, P. Woodcock, F.A. Edwards, T.H. Larsen, W.W. Hsu, S. Benedick, K.C.
830 Hamer, D.S. Wilcove, C. Bruce, X. Wang, T. Levi, M. Lott, B.C. Emerson, and D.W.
831 Yu. 2013. Reliable, verifiable and efficient monitoring of biodiversity via
832 metabarcoding. *Ecology Letters* 16:1245-1257.

833 Kéry, M., and J.A. Royle. 2010. Hierarchical modelling and estimation of abundance and
834 population trends in metapopulation designs. *Journal of Animal Ecology* 79:453–
835 461.

836 Kelly, R.P., J.A. Port, K.M. Yamahara, and L.B. Crowder. 2014 Using Environmental
837 DNA to Census Marine Fishes in a Large Mesocosm. *PLoS ONE* 9: e86175

838 Kimura, D.K. and H.H. Zenger, Jr. 1997. Standardizing sablefish (*Anoplopoma fimbria*)
839 long-line survey abundance indices by modeling the log-ratio of paired comparative
840 fishing cpues. *ICES Journal of Marine Science*, 54:48–59.

841 Klymus, K., C.A. Richter, D. Chapman, and C. Paukert. 2015. Quantification of eDNA
842 shedding rates from invasive bighead carp *Hypophthalmichthys nobilis* and silver
843 carp *Hypophthalmichthys molitrix*. *Biological Conservation* 183:77–84

844 Laramie, M.B., D.S. Pilliod, and C.S. Goldberg. 2015. Characterizing the distribution of
845 an endangered salmonid using environmental DNA analysis. *Biological*
846 *Conservation*. 183:29-37.

847 Lee, C.K., C.W. Herbold, S.W. Polson, K.E. Wommack, S.J. Williamson, I.R.
848 McDonald, and S.C. Cary. 2012. Groundtruthing next-gen sequencing for microbial
849 ecology—biases and errors in community structure estimates from PCR amplicon
850 pyrosequencing. *PloS One*, 7:e44224.

851 Leray, M., and N. Knowlton. 2015. DNA barcoding and metabarcoding of standardized
852 samples reveal patterns of marine benthic diversity. *Proceedings of the National*
853 *Academy of Sciences*, 112:2076–2081.

854 Lodge, D.M., C.R. Turner, C.L. Jerde, M.A. Barnes, L. Chadderton, S.P. Egan, J.L.
855 Feder, A.R. Mahon, and M.E. Pfrender. 2012. Conservation in a cup of water:
856 estimating biodiversity and population abundance from environmental DNA.
857 *Molecular Ecology* **21**: 2555–2558.

858 Love, M.I., W. Huber, and S. Anders. 2014. Moderated estimation of fold change and
859 dispersion for RNA-seq data with DESeq2. *Genome Biology* 15:550.

860 Martin M. 2011. Cutadapt removes adapter sequences from high-throughput sequencing
861 reads. *EMBnet.journal* 17:10-12.

862 Maruyama A., K. Nakamura, H. Yamanaka, M. Kondoh, and T. Minamoto. 2014. The
863 release rate of environmental DNA from juvenile and adult fish. *PLoS ONE*
864 9:e114639.

865 Merkes C.M., S.G. McCalla, N.R. Jensen, M.P. Gaikowski, and J.J. Amberg. 2014
866 Persistence of DNA in carcasses, slime and avian feces may affect interpretation of
867 environmental DNA data. *PLoS ONE* 9(11): e113346.
868 doi:10.1371/journal.pone.0113346.

869 Nathan, L.M., M. Simmons, B.J. Wegleitner, C.L. Jerde, and A.R. Mahon. 2014.
870 Quantifying environmental DNA signals for aquatic invasive species across
871 multiple detection platforms. *Environmental Science and Technology* 48:12800-
872 12806.

873 O'Donnell, J.L., R.P. Kelly, N. Lowell, and J.A. Port. 2015. Indexed PCR primers induce
874 template-specific bias in large-scale DNA sequencing studies. *In revision*. *PLoS One*,
875 July 2015.

876 Paradis E., J. Claude, and K. Strimmer. 2004. APE: analyses of phylogenetics and
877 evolution in R language. *Bioinformatics* 20:289-290.

878 Pedersen, M.W., S. Overballe-Petersen, L. Ermini, C. Der Sarkissian, J. Haile, M.
879 Hellstrom, J. Spens, P.F. Thomsen, K. Bohmann, E. Cappellini, I.B. Schnell, N.A.
880 Wales, C. Carøe, P.F. Campos, A.M.Z. Schmidt, M.T.P. Gilbert, A.J. Hansen, L.
881 Orlando, and E. Willerslev. 2015. Ancient and modern environmental DNA.
882 *Philosophical Transactions of the Royal Society B* 370:20130383.

883 Pilliod, D.S., Goldberg, C. S., Arkle, R. S., and Waits, L. P. (2014). Factors influencing
884 detection of eDNA from a stream-dwelling amphibian. *Molecular Ecology Resources*

885 **14**:109–116.

886 Plummer, M. 2003. JAGS: A program for analysis of Bayesian graphical models using
887 Gibbs sampling.

888 Plummer, M., N. Best, K. Cowles, and K. Vines. 2006. CODA: Convergence Diagnosis
889 and Output Analysis for MCMC 6:7–11.

890 Polz, M.F., and C.M. Cavanaugh. (1998). Bias in template-to-product ratios in
891 multitemplate PCR. *Applied and Environmental Microbiology* 64: 3724–3730.

892 Puillandre, N., E.E. Strong, P. Bouchet, M.-C. Boisselier, A. Couloux, and S. Samadi.
893 2009. Identifying gastropod spawn from DNA barcodes: possible but not yet
894 practicable. *Molecular Ecology Resources* 9:1311–1321.

895 R Core Team. 2014. R: A language and environment for statistical computing. R
896 Foundation for Statistical Computing, Vienna, Austria.

897 Renshaw, M.A, B.P. Olds, C.L. Jerde, M.M McVeigh, and D.M. Lodge. 2015. The room
898 temperature preservation of filtered environmental DNA samples and assimilation
899 into a phenol–chloroform–isoamyl alcohol DNA extraction. *Molecular Ecology*
900 *Resources* 15:68-176.

901 Riaz T., W. Shehzad, A. Viari, F. Pompanon, P. Taberlet, E. Coissac. 2011. ecoPrimers:
902 inference of new DNA barcode markers from whole genome sequence analysis.
903 *Nucleic Acids Research* 39:e145. doi:10.1093/nar/gkr732.

904 Royle, J.A., and J.D. Nichols. 2003. Estimating abundance from repeated presence–
905 absence data or point counts. *Ecology* 84:777–790.

906 Schloss, P.D., D. Gevers, and S.L. Westcott. 2011. Reducing the effects of PCR
907 amplification and sequencing artifacts on 16S rRNA-based studies. *PLoS ONE* 6:
908 e27310. doi:10.1371/journal.pone.0027310

909 Schmidt, B.R., M. Kéry, S. Ursenbacher, O.J. Hyman, and J.P. Collins. 2013. Site
910 occupancy models in the analysis of environmental DNA presence/absence surveys:
911 a case study of an emerging amphibian pathogen. *Methods in Ecology and Evolution*
912 4:646–653

913 Shelton, A.O., E.J. Dick, D.E. Pearson, S. Ralston, and M. Mangel. 2012. Estimating
914 species composition and quantifying uncertainty in multispecies fisheries:
915 hierarchical Bayesian models for stratified sampling protocols with missing data.

916 Canadian Journal of Fisheries and Aquatic Sciences 69:231–246.

917 Simmons, M., A. Tucker, W.L. Chadderton, C.L. Jerde, A.R. Mahon. *In Press*. Active and
918 passive environmental DNA surveillance of aquatic invasive species. Canadian
919 Journal of Fisheries and Aquatic Sciences.

920 Sipos, R., Székely, A.J., Palatinszky, M., Révész, S., Márialigeti, K., and Nikolausz, M.
921 2007. Effect of primer mismatch, annealing temperature and PCR cycle number on
922 16S rRNA gene-targeting bacterial community analysis. *FEMS Microbiology*
923 *Ecology* 60:341–350.

924 Soergel, D.A.W., N. Dey, R. Knight, and S.E. Brenner. 2012. Selection of primers for
925 optimal taxonomic classification of environmental 16S rRNA gene sequences. *The*
926 *ISME Journal* 6:440–1444.

927 Strickler, K.M., A.K. Fremier, and C.S. Goldberg. Quantifying effects of UV-B,
928 temperature, and pH on eDNA degradation in aquatic microcosms. 2015. *Biological*
929 *Conservation* 183: 85-92.

930 Su, Y.-S., and M. Yajima. 2014. R2jags: A package for running jags from R. R package
931 version 0.04-03.

932 Takahara, T., T. Minamoto, H. Yamanaka, H. Doi, and Z. Kawabata. 2012. Estimation of
933 fish biomass using environmental DNA. *PLoS ONE*, 7:e35868

934 Thomsen, P.F., J. Kielgast, L.L. Iversen, P.R. Møller, M. Rasmussen, and E. Willerslev.
935 2012. Detection of a diverse marine fish fauna using environmental DNA from
936 seawater samples. *PloS One*, 7:e41732.

937 Thomsen, P.F., J. Kielgast, L.L. Iversen, C. Wiuf, M. Rasmussen, M.T.P. Gilbert, L.
938 Orlando, and E. Willerslev. 2012. Monitoring endangered freshwater biodiversity
939 using environmental DNA. *Molecular Ecology* 21: 2565-2573.

940 Turner, C.R., M.A. Barnes, C.C.Y. Xu, S.E. Jones, C.L. Jerde, and D.M. Lodge. 2014.
941 Particle size distribution and optimal capture of aqueous microbial eDNA. *Methods*
942 *in Ecology and Evolution* 5:676–684.

943 Venables, W.N. and C.M. Dichmont. 2004. GLMs, GAMs and GLMMs: an overview of
944 theory for applications in fisheries research. *Fish. Res.* 70:319–37.

945 Venter, J.C., K. Remington, J.F. Heidelberg, A.L. Halpern, D. Rusch, J.A. Eisen, D. Wu
946 et al. 2004. Environmental genome shotgun sequencing of the Sargasso Sea. *Science*

947 304:66-74.
948 Wright, E.S., L.S. Yilmaz, S. Ram, J.M. Gasser, G.W. Harrington, and D.R. Noguera.
949 2014. Exploiting extension bias in polymerase chain reaction to improve primer
950 specificity in ensembles of nearly identical DNA templates. *Environmental*
951 *Microbiology* 16:1354-1365.
952 Yoccoz, N.G. 2012. The future of environmental DNA in ecology. 2012. *Molecular*
953 *Ecology* 21:2031-2038.
954 Yu, D.W., J. Yinqiu, B.C. Emerson, X. Wang, C. Ye, C. Yang, and Z. Ding. 2012.
955 Biodiversity soup: metabarcoding of arthropods for rapid biodiversity assessment and
956 biomonitoring. *Methods in Ecology and Evolution* 3:613–623
957 Zhang J.J., K. Kobert, T. Flouri, and A. Stamatakis. 2014. PEAR: a fast and accurate
958 Illumina paired-end read merger. *Bioinformatics* 30:614-620.

eDNA:

$B_{i,l,t}$ Biomass of taxon i at location l , time t

Traditional Sampling:

Shedding rate
Degradation rate
Immigration/Emigration

θ

Detection probability
Sampling method biases

$W_{i,l,t}$ Density of DNA in the environment

$D_{i,l,t}$ Collected Individuals

Water collection
Soil sampling

ϕ

Sample preservation
Pre-sorting
Subsampling

$X_{i,l,t}$ Collected DNA (e.g. DNA on filter)

$E_{i,l,t}$ Subsampled Individuals

DNA extraction

ψ

Taxonomic identification
Measurement (count or mass)

$Y_{i,l,t}$ DNA present after DNA extraction

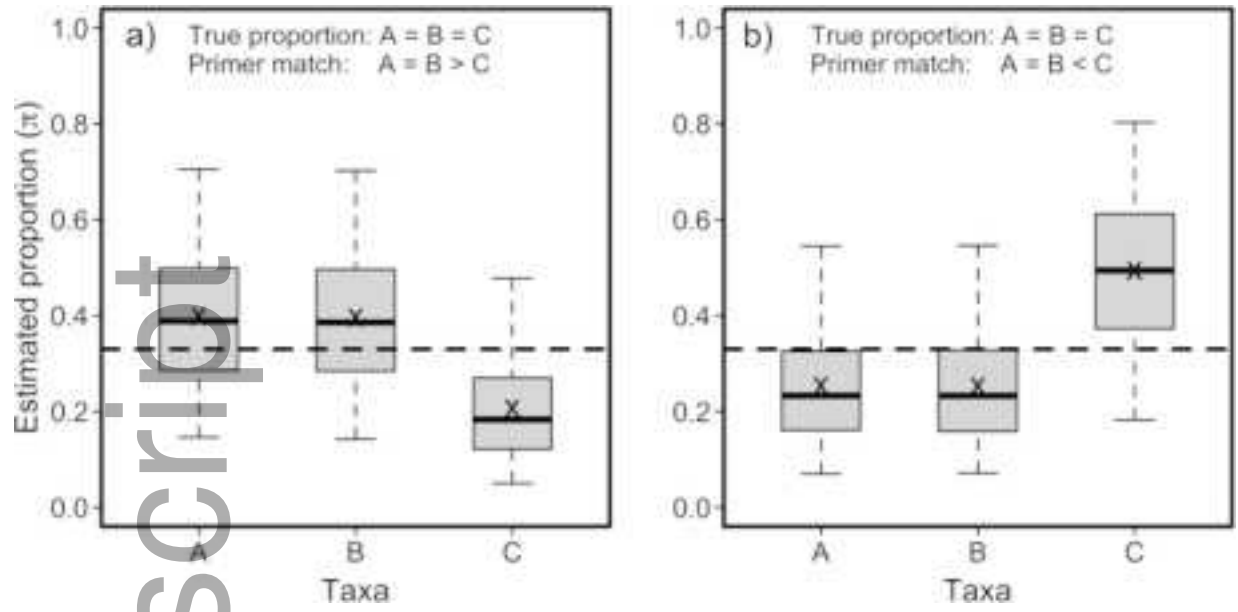
$F_{i,l,t}$ Taxon specific Count or Biomass

PCR amplification
DNA sequencing
Assigning OTUs to taxa

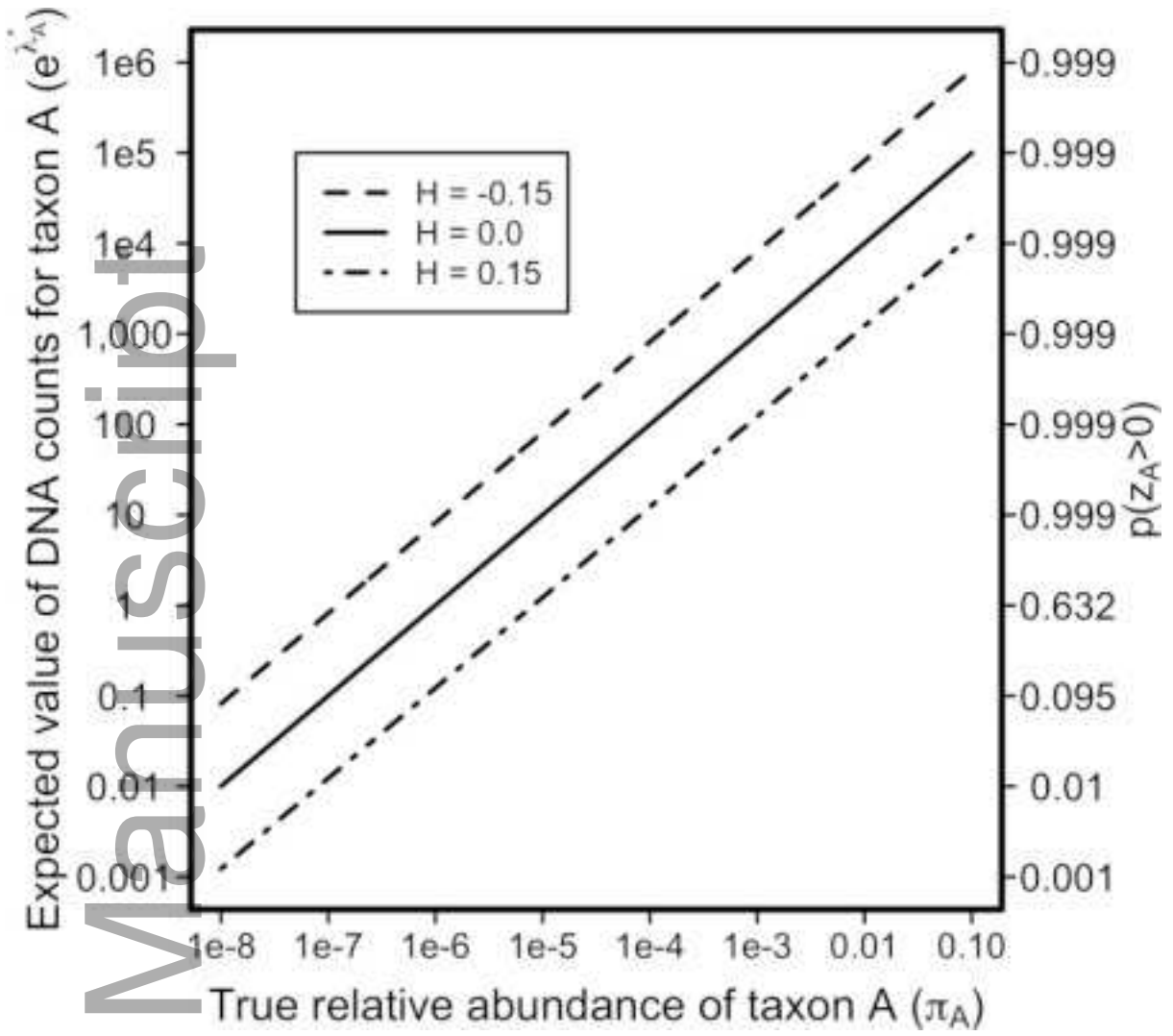
ξ

$Z_{i,l,t}$ DNA present after sequencing

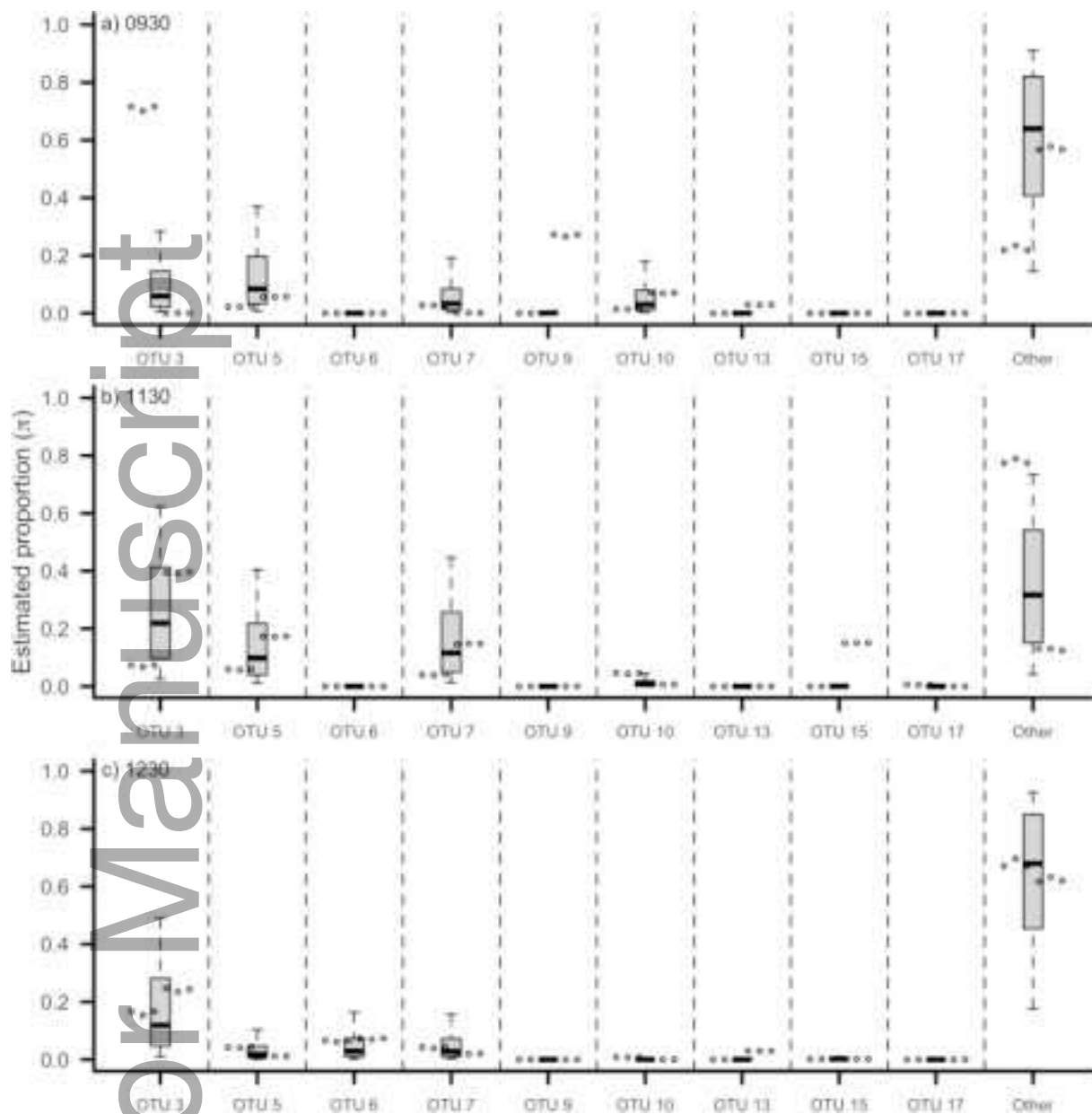
Author Manuscript



eap_1336_f2.jpg



eap_1336_f3.jpg



eap_1336_f4.jpg

