

1 **Within-Community-Sampling Power Analysis to Detect Richness Change**

2
3 Eden Tekwa¹, Jake Lawlor¹, Matt Lemay², Ryan R.E. Stanley³, Nicholas W. Jeffery³,
4 Emily Rubidge⁴, Jennifer Sunday^{1*}

5 6 Affiliations:

- 7 1. Department of Biology, McGill University, 1205 Dr. Penfield, Montreal,
8 PQ, Canada.
- 9 2. Hakai Institute, street address, Hariott Bay, British Columbia
- 10 3. Fisheries and Oceans Canada, Bedford Institute of Oceanography, 1 Challenger
11 Dr., Dartmouth, NS, Canada
- 12 4. Fisheries and Oceans Canada, Institute of Ocean Sciences, 9860 West Saanich
13 Rd., Sidney, BC, Canada

14
15 Corresponding Author: Jennifer.Sunday@mcgill.ca

16 17 **Abstract**

18 Reliable biodiversity monitoring requires understanding how sampling effort influences
19 detectability of meaningful changes in species richness. Increased sampling within
20 independent units has been shown to reduce measurement error, while sampled
21 richness estimates are often subject to bias. However, robust methods for quantifying
22 the relationship between sampling effort and the power to detect biodiversity change
23 remain limited. Here we present a simulation-based power analysis framework that
24 explicitly links sampling effort within communities to the ability to detect richness change
25 between them. The approach uses empirical pilot data to simulate sampling and applies
26 a dimensionless coverage metric to translate simulated effort into real sample sizes
27 within communities. This produces direct relationships among power, effect size, and
28 coverage, enabling estimation of the sampling required to achieve a specified
29 probability of detecting richness change in the correct direction. We also provide a
30 variant designed for before–after monitoring scenarios in which pilot data are available
31 from a single community and the post-impact state is unknown. We demonstrate the
32 application of this framework using a field-based environmental DNA biodiversity
33 dataset. Additionally, we show that sample size recommendations rapidly converge
34 when we apply the initial recommendation and reanalyze the next survey. The method
35 is implemented in the R package *BioDivPoweR*, enabling broad application for
36 designing efficient biodiversity monitoring programs and evaluating emerging sampling
37 technologies.

38
39 **Keywords:** biodiversity, species richness, monitoring, sampling effort, quality control,
40 environmental DNA

41

42

43 **Introduction**

44

45 Biodiversity underpins the functioning of ecosystems [1] that ultimately support and
46 sustain human well-being [2]. As such, monitoring biodiversity change is essential for
47 maintaining ecosystem functions and the well-being they provide. However, reliably
48 detecting and attributing causes to these changes have been one of the most difficult
49 technical problems in biodiversity science [3–5]. Enhancing the effectiveness of
50 biodiversity monitoring programs will be vital for achieving international goals and
51 targets[6], such as the commitment among nations to effectively protect 30% of areas
52 under the Global Biodiversity Framework [7]. At more local scales, the ability to reliably
53 detect changes in species richness or abundance-weighted diversity indicators is
54 fundamental to evaluate ecological conditions, inform management decisions, and
55 consider attaching financial incentives to biodiversity outcomes [8].

56

57 Despite the importance of biodiversity change assessments, we lack tools for
58 understanding how much monitoring effort is required to robustly assess policy-relevant
59 biodiversity changes within a region of interest. In a recent meta-analysis of species
60 richness change associated with implementation of Marine Protected Areas (MPAs),
61 about one half of studies showed no increase in diversity compared to reference sites
62 [9] (from Figure S3 therein). Not observing an increase could have two very different
63 explanations: (1) there could be no biologically-meaningful biodiversity differences, or
64 (2) studies could lack the sampling effort needed to detect differences that were
65 present. An estimate of power to detect a known global effect with increasing
66 independent replicates advised caution in interpreting stable local richness as evidence
67 that biodiversity loss is not occurring [3]. When we lack independent replicates, such as
68 in a monitoring program at a local scale, we may still be interested in evaluating our
69 power to detect an effect with different within-community sampling designs. Yet no
70 analysis exists to our knowledge on how to advise local managers on how much
71 sampling effort is needed to robustly detect change of a given effect size. Here we
72 review why this is the case and propose a user-friendly method to achieve this.

73

74 ***What is a power analysis?***

75 Statistical power is the rate of positive identification of an effect (such as
76 differences in body size between two populations) when in fact there is an effect (i.e.,
77 the true positive rate, or 1-false negative rate). In prospective power analysis [10], we
78 ask how many independent replicates are needed to detect a difference between 2 or
79 more normally-distributed random variables, with an acceptable power level (1-B, often
80 set to 80%), a prespecified significance level α (usually 0.05), and the population effect

81 size, to be detected with probability $1-B$ [10]. For example, statistical power for detecting
82 a difference in body size between the two treatment groups can be estimated as a
83 function of the true difference in body sizes (effect), user-set α (probably of type I error),
84 and number of independent treatments (sample size)[11]. In ecology, examples of
85 detection objectives amenable to standard power analyses include detecting differences
86 or trends in ecological values when effect means are estimated using multiple
87 independent samples [12].

88

89 ***Why is power in biodiversity change difficult to assess?***

90 Biodiversity monitoring poses a special challenge to assessing statistical power
91 for two reasons: biased measurement error, and a tendency to not replicate
92 independent samples. Natural communities are typically composed of species that
93 cannot be found in every location because of spatial heterogeneity and rarity (reviewed
94 in [4]). As such, biodiversity is typically assessed by drawing ‘samples’ within a unit of
95 study, and as more samples are drawn, the total number of species encountered is
96 pooled (note, they are not averaged), bringing the measurement increasingly close to a
97 maximum and correct (or *census*) species richness. Sampled richness is therefore not
98 typically a mean across independent samples, but a cumulative metric of a single
99 independent unit (i.e. a place or single community), and is necessarily downward-biased
100 relative to true richness (Fig. 1b), at a scale that diminishes with increased observation
101 effort (such as area, volume, or numbers of samples). In addition, in richness estimates,
102 each individual’s species identity is only used toward a community’s richness if it
103 belongs to a previously unrecorded species. In other words, the weight of the
104 individual’s datapoint is dependent on other individuals. When comparing two or more
105 communities, a central challenge is determining how many samples within each
106 community are needed to robustly detect a difference between them in the correct
107 direction, but because samples contribute to a tally and not a mean, this question is not
108 addressed with a traditional power analysis.

109 Sometimes in biodiversity studies, independent treatments are replicated (e.g.,
110 multiple communities representing each treatment and controls) [13], and we may then
111 apply a traditional power analysis to assess differences between mean richness of two
112 groups. However, even in these cases, measurement error within each replicate can
113 affect outcomes. For example, a survey may repeatedly observe 5 species in
114 community type A and 3 species in community type B. However, if the survey always
115 misses 3 rare species in community type A and 10 rare species in community type B
116 (sampling bias, [4]), the survey misidentifies community type A as richer even though
117 the species estimates appear to be precise. Spatial heterogeneity and rarity of species
118 occurrences calls for attention to the role of sampling effort to assure that values from
119 each replicate are meaningful to test the effect at the scale of interest [14], yet the role
120 of measurement error has been little discussed in ecological power analysis [12].

121 Furthermore, in biodiversity monitoring applications, we are often only interested
122 in detecting temporal change between two communities without replication, such as
123 comparing before vs. after an environmental impact, or an affected vs. reference area
124 [15]. In these cases, it is even more important to consider the role of observation effort
125 on the rate of detecting true differences between biodiversity metrics. To focus efforts
126 on biologically meaningful change, we may consider a threshold of biodiversity change
127 that calls for management action, such as a 10% difference [12,16]. In either case, the
128 question of how many samples are needed in each community to detect a biodiversity
129 difference of a set size with a desired rate of correct detections (power) is a question of
130 quality assurance against measurement error, i.e. we are discussing the quality of
131 information contributing to essentially two data points (richness) that we then rank.

132

133 ***Biodiversity measurement quality assurance***

134 We propose that the question of how to draw samples from communities to
135 assess species richness is analogous to how to draw blood from patients in a
136 standardized way - a quality assurance step. The objective of quality assurance in this
137 case is to obtain sufficient within-subject sampling to increase the probability of
138 detecting a change in the correct direction. For the case of drawing blood in a single
139 human subject (analogous to a community), we can test how the needle, volume of
140 blood drawn, and blood processing affect test results in controlled experiments where
141 the conditions and expected outcomes can be precisely defined. This can be formally
142 justified by a power analysis on the number of blood particles one needs to draw in
143 order to compare two blood samples. Note, in this case the power analysis is framed
144 on two subjects, not on two treatment groups with multiple replicates. It still makes
145 sense to talk about power for two data points if the binary empirical decision on the
146 observed state (person one has higher iron concentration than person two) is simply
147 whether the point readings are higher (without a t-test). It is the same concept when we
148 envision a power analysis for how to sample two communities to detect a biodiversity
149 difference; we are still able to ask about power (e.g., probability of detecting a difference
150 in the correct direction) across different sampling qualities, even without replication
151 needed for a t-test across groups. The quality assurance power analysis is more difficult
152 for biodiversity science, because we generally do not know, or have control over, the
153 true states of ecological communities.

154

155 ***Both kinds of power analysis are needed in biodiversity science***

156 A quality assurance power analysis can be useful when we only have two
157 communities, or even a single community of interest in which we aim to detect change
158 over time. A quality assurance power analysis can also be useful to improve inference
159 when we have true replication of communities, by assessing how well each replicate
160 captures true richness. In the latter case, analyses can additionally leverage the broader

161 suite of traditional statistical tools, including hypothesis testing and power analyses. For
162 simplicity, we will use “power analysis” at times in this text to mean quality assurance
163 power analysis for detecting differences in two communities, not two groups of
164 communities. “Sample size” will be used to refer to the number of samples drawn from a
165 single community, rather than the number of communities sampled.

166

167 ***Proposed Solution***

168 We develop a method for conducting a quality-assurance power analysis by
169 using simulation to evaluate the effects of progressively reduced sampling efforts
170 starting from a pilot biodiversity dataset. In most biological communities, no practical
171 method approaches a true census. Therefore, we use simulated communities where
172 true biodiversity states are assumed and then subsampled, starting with an artificial
173 census.

174 The approach we propose is to computationally resample empirical community
175 composition datasets. While down-sampling these datasets to evaluate the relationship
176 between sampling effort and power is straightforward and has been done in previous
177 studies (e.g. [17]), the key challenge lies in translating the resulting power-effort
178 relationships to the real system, which exists at a broader spatial scale. Specifically, we
179 need to convert the abscissa variable (unit of effort) from the number of samples in each
180 community, which is a unit-specific measure, to a dimensionless value that represents a
181 percentage of the total sampling effort. This can be achieved using empirical
182 relationships between effort and percent coverage [18,19]. Using coverage as a
183 universal measure of sampling effort allows the simulated relationships between effect
184 size, power and sampling effort from pilot data can be rescaled to the real communities,
185 including effort levels beyond the pilot dataset, through the principle of scale-invariance
186 (analogous to modeling turbulence in fluid mechanics [20]). We explore this framework
187 to make full use of natural and detection-based variation in species occurrences across
188 pilot samples, in a manner analogous to how traditional power analyses use pilot data to
189 estimate sample variance around a mean.

190 We focus on two simple but common scenarios encountered by biodiversity
191 practitioners. The first is estimating differences between two “treatments” (e.g., affected
192 vs. control), where pilot data from two communities are used to estimate power-effort-
193 effect size relationships. The second is a biodiversity monitoring scenario where the
194 goal is to detect compositional changes across time between present and future states.
195 In the latter case, a range of possible future states and effect sizes are inferred based
196 on the present state and assuming future states to be within the range of random
197 reassembly of the observed community. To demonstrate the application of the power
198 analysis, we use pilot eDNA metabarcoding data on marine biodiversity collected from
199 nearshore coastal regions of the northeast Pacific Ocean. We also provide and describe

200 the implementation of the analysis in an openly accessible R package whereby users
201 can enter their own pilot data for power analysis in any system following our method.

202

203 **Methods**

204 ***Overview of Quality Assurance Power Analysis for Species Richness***

205 Our methods are organized into four subsections. First, we derive a relationship
206 between coverage and the number of samples so that coverage becomes the universal
207 measure of sampling effort. Second, we develop the computational method that uses a
208 pilot study of species occupancies from two pre-defined communities (a two-community
209 comparison scenario) to link effect size, coverage, and power. We repeat this for the
210 scenario when we have pilot samples from only one defined community (scenario for a
211 present-future comparison). Power in these cases was defined as the probability of
212 detecting a difference in the same direction as the true simulated difference, so the
213 statistic is the sign (+/-) of the difference in richness. Third, we develop analyses for the
214 same two scenarios, but this time assess the probability of detecting a difference that
215 exceeds a predefined minimum effect size (akin to *minimum effect power analysis*
216 [12,16]). Finally, we demonstrate the application with an eDNA metabarcoding dataset,
217 and run all four power analyses to assess how estimates of sample coverage needed to
218 monitor for changes between two habitats at a given level of power is predicted by pilot
219 study sample size.

220 All methods are implemented in the R-package *BioDivPoweR* [21], enabling
221 general application to user-supplied datasets, allowing users to evaluate their own pilot
222 data.

223

224 ***Developing a relationship between sample size and sample coverage***

225 To design a quality-assurance power analysis for detecting biodiversity change,
226 we first transform sampling effort to be comparable in both communities and
227 transferable from pilot studies to the hypothetical real study. If sampling coverage of
228 communities being compared is not the same, the comparison of richness (including
229 direction of richness differences between them) may be biased and not directly
230 comparable [22].

231 Here, we derive a relationship between coverage and sample number for each
232 community, where a sample is a population of individuals from any species collected as
233 a single unit of effort (e.g., one eDNA sample filter from the ocean or one quadrat from a
234 terrestrial survey). Coverage (Eq. 1) is defined as the portion of species captured (C ,
235 [18]), and a sample is an independent population draw (k) within each community. S is
236 the observed number of species and S' is an asymptotic richness estimator. If the same
237 percent of species is captured in two communities (i.e., coverage is equal), we can
238 determine which community has higher diversity.

239

240 (1) $C(k)=S(k)/S'$

241

242 We can use an asymptotic estimator of S' to estimate the relationship between
243 coverage C and species richness observed S . Here we use the Ω asymptotic estimator
244 that bases estimates on the entire observed pattern of species occupancies [19], that is
245 more accurate than the *Chao2* estimator of S' based on the number of species
246 occupying 1 or 2 samples [19].

247 Let's define $C_i(k)$ as the observation probability of species i after k samples. We
248 first make the assumption that each sample is a random independent draw of
249 community composition. $C(k)$ can be understood as the average observation probability
250 across species after k samples. For each species i , C_i is a function of P_i (observed
251 occupancy, or portion of samples where species i is observed). P_i is an observation,
252 whereas C_i is the estimated occupancy we would observe if k were to increase beyond
253 the original study's sample size. The across-species average probability C ranges from
254 0 to 1 and is an increasing function of sample size k . $\langle \rangle_i$ is the mathematical operator
255 for averaging over i . As sample size increases, the average probability approaches one.

256

257 (2) $C(k) = \langle C_i(k) \rangle_i = \langle 1 - (1 - P_i)^k \rangle_i$

258

259 In general, C in Eq. 2 is the average of a series of exponential functions
260 belonging to i , where the exponent is i and the base is P_i . C approaches one as k
261 approaches infinity. Unlike non-mechanistic statistical extrapolation on the rarefaction
262 curve that does not asymptote, Eq. 2 implies that the number of species captured (CS'
263 where S' is the true number of species - a constant) necessarily approaches a finite
264 number as k approaches infinity. The inverse of Eq. 2 allows us to rescale any function
265 of coverage (C) below 100% to the actual number of samples (k) by solving k given C .
266 Generally, every marginal increase in coverage (C) corresponds to an increasing
267 increment in sample size (k), meaning that a large increase in samples is required to
268 gain an increase in coverage when the coverage is already high.

269 Coverage (C) serves us to rarify to equal coverage when we directly compare
270 communities, and also to provide a dimensionless measure of sampling effort. We
271 therefore apply simulated relationships between power and coverage to the whole
272 community, by assuming that the pattern of occupancies of unobserved species is the
273 same as that of observed species. Since coverage is unitless, applying simulated *in*
274 *silico* relationships between power and coverage to real ecosystems is dimensionally
275 consistent.

276

277 ***Inferring power from data in two communities (Scenario 1)***

278 Our first power analysis scenario assesses relationships between power, sample
279 size, and effect size from a pilot dataset in which multiple samples are taken from two

280 'treatment' communities we wish to compare. The two treatments could represent
281 different habitat types, sites experiencing different environmental stressors or regulatory
282 regimes, two time points from the same community, etc. Importantly, replicate samples
283 in this design are not independent replicates of the treatment but instead multiple draws
284 from two community types (i.e., two treatments, no replication) that we wish to compare.
285 Random samples from within each community provide estimates of occupancy across
286 samples that inform power. These are represented as multiple samples in red and blue
287 communities in Fig. 1b and Fig. 2.

288 From the dataset (first row of boxes in Fig. 2), we first use Eq. 2 to compute the
289 relationship between sample size and coverage for each community. We can also
290 determine rarified richness differences from the pilot communities for later comparison.

291

292 *Estimating power as a function of sample coverage using simulation*

293 From each community, we create multiple bootstraps (resample with
294 replacement) through a Monte Carlo procedure, and consider these as an array of
295 simulated communities (second row of boxes in Fig. 1). We randomly pair bootstraps of
296 each of two communities, and for each pair calculate the rarefied sample sizes using
297 Eq. 2. These numbers will be different for each bootstrap pair and different from the
298 rarefaction in the original communities.

299 For each bootstrapped pair of communities, individuals are subsampled (without
300 replacement) according to the number of samples that correspond to coverage values
301 ranging from zero to the maximum rarefied coverage (unique for each pair). We then
302 make the assumption that the maximum rarefied coverage is 100% for the subsequent
303 steps. For example, the rarefied coverage of 55% achieved by one bootstrapped pair is
304 relabelled as 100%, and all other rarefied coverage values below 55% are multiplied by
305 $100/55$. The richness difference is measured (in terms of \log_2 of the $\text{richness}_2/\text{richness}_1$
306 ratio, or fold difference for short) at each subsampled step. Using the assumption that
307 the observed difference at maximum coverage is the true difference (i.e., not from
308 samples but from a census, third row of boxes in Fig. 1), the coverage is rescaled to
309 pool results from different bootstrap pairs. We record these richness difference values
310 for the next step.

311

312 *Compiling results*

313 The bootstrap pairs of communities are next binned according to their richness
314 differences at maximum coverage as in a histogram (Fig. 2a), grouping bootstraps with
315 assumed-to-be similar 'true' effect sizes, and the continuous range of bins that contain a
316 minimum number of replicates desired for power analyses are extracted from each
317 (e.g., 40; coloured portion of bars in Fig. 2a). This range of effect size is proportional to
318 the heterogeneity in species compositions and represents a likely range of expected
319 richness for communities similar to the original data source. Smaller bin sizes increase

320 the effect size resolution of the final power analysis, while holding computation load
321 constant, thereby decreasing the range of effect sizes with sufficient replicates.

322 Then, for each effect size bin, we compare the measured richness difference at
323 different coverage values including full coverage with the again assumed to be 'true'
324 richness difference. The proportion of simulated differences in the correct direction is
325 taken as the statistical power or the probability of detecting the correct direction of
326 differences. The results taken together provide a 3-dimensional relationship between
327 coverage, effect size, and power to detect change in the correct direction (Fig. 1 bottom
328 plots, example in Fig. 2b).

329 We then rescale sample coverage to sample size from each community using
330 our dimensionless scalar (Eq. 2) that relates sample size to coverage for each
331 community. This allows for a specific relationship between power, effect size, and
332 sample size (Fig. 1 bottom plots, upper x-axis).

333 One particularly useful view of this relationship comes if we hold power constant
334 (e.g. typically 0.8 or 80% correct detections), and express the relationship between
335 coverage (or sample size in each community) and the minimum detectable effect size
336 (or amount of biodiversity change that is detectable in the correct direction at a given
337 sample size; bottom right of Fig. 1). This relationship provides potentially critical
338 information for planning monitoring efforts: *How much sampling is needed to detect a*
339 *particular effect size?* For every possible level of coverage, we can extract the minimum
340 effect size that achieves a given power (e.g., 0.8), and fit a linear regression to rank of
341 coverage to describe the general relationship between coverage and \log_2 of effect size.
342 We transformed coverage to exaggerate small increments at higher coverage values,
343 according to the non-linear sequence of coverage values: $([250:250:10000].^0.5)/100$.
344 This sequence tends to generate a linear relationship between detectable effect size
345 and coverage increments.

346

347 ***Inferring power from data in one community (Scenario 2)***

348 When empirical effect sizes between treatments are unavailable, the
349 relationships between power, coverage, and effect size may still be of interest. If no
350 treatment is defined in the initial dataset (i.e., sample data are from a single community
351 in one time point) we may still be interested in detecting changes over time. For this, we
352 can use a similar procedure described in the previous section, except that samples are
353 drawn from the original single dataset into two simulated communities leading to a
354 richness difference between them. The logic of this application is illustrated in Figure
355 S1. Here, the simulated community pairs will on average have an effect size of zero.
356 However, the variation generated from the between-sample compositional
357 heterogeneities can be understood as an estimate of how biodiversity could change.
358 This approach assumes that the community is a closed system within a particular time
359 frame of interest for tracking environmental changes, and that species retain their

360 occupancy distributions. Here, biodiversity changes must be selected from existing
361 variations in the original community, which is compatible with short term stability in
362 species occupancy patterns [23,24]. Therefore, our proposal to feed the same
363 community data to two simulated communities is not only constrained by data limitation
364 but has plausible biological justification.

365

366 ***Detecting effect sizes above a conservation-relevant minimum effect***

367 Because large changes in species richness could be important for triggering
368 conservation action, we also applied a “minimum effect” power analysis to estimate how
369 many samples are needed to detect a change greater than a predefined benchmark to
370 trigger conservation intervention [12,16]. We used a threshold of 0.1 log₂ fold change in
371 richness as our action-oriented benchmark (~7% increase in richness). Note, this is not
372 the same as asking how many samples are needed to detect a difference greater than 0
373 and in the correct direction if they were truly different by 0.1 log₂ fold change (that
374 would be attainable in our above analysis), instead, this shows how many samples
375 would be needed to detect a difference greater than 0.1 log₂ fold change in the correct
376 direction if they were truly different by at least that amount. We applied this to both
377 scenarios above, and expected a greater number of samples to be required for this
378 more stringent detection.

379

380 ***Applying our power analysis to biodiversity data derived from eDNA***

381 We apply our simulation under both scenarios to a dataset of marine community
382 assemblages, sampled using environmental DNA (eDNA) metabarcoding in the central
383 coast of British Columbia, Canada (methods and data are described in [25]). Briefly,
384 triplicate 1L-seawater samples were collected at 208 locations across a temperate
385 marine region roughly 100km² in scale, focussed on nearshore habitats of kelp,
386 seagrass, and rocky reef, and nearby deep-water sites. Samples were collected at
387 depth using a Niskin bottle and filtered using 0.22µM Sterivex filters. DNA
388 metabarcoding was implemented to target fish communities by amplifying a short region
389 of the 12S rRNA gene (MiFish-U; [26]). The resulting libraries were sequenced using
390 Illumina MiSeq V3 chemistry. Raw sequences were clustered into ASVs using the
391 DADA2 pipeline [27] and annotated using BLAST searches of the complete NCBI
392 nucleotide database (complete methods are available with the data package [25]).
393 Vegetation cover of each sample was estimated from video footage for every sample
394 and scored; we used a coarse category of high vs. low vegetation cover for our two-
395 treatment analysis, and used all of the data as pilot for our one-treatment analysis. We
396 refer to this dataset as the *Coastal Highly Density eDNA* dataset. Although this dataset
397 has three replicated samples within each of 208 sites, for the sake of our simulation we
398 pooled samples within sites, and treated every individual site as an independent
399 “sample” of the whole community.

400

401 **Sensitivity Analyses**

402 We assessed, through simulation, the sensitivity of our power analysis to the size
403 of the initial pilot dataset and explored the possibility of convergence if the
404 recommendations of the power analysis were iteratively used. We assessed sensitivity
405 to pilot size by randomly drawing (without replacement) smaller subsets of the full
406 dataset (208 samples, from 5% to 100% of the total) and using it to estimate number of
407 samples needed to detect a 0.1 fold change at 0.8 power in the 1-treatment power
408 analysis. We assessed 'convergence' by drawing, for each initial pilot data simulation, a
409 new dataset that matched the newly recommended sample size, and using that for a
410 subsequent power analysis, to estimate the number of samples needed to detect a 0.1
411 fold change at 0.8 power. This simulates 'adaptive monitoring', or iteratively using new
412 samples to run a new power analysis, assuming stationarity in mean species occupancy
413 across samples and richness. We repeated the adaptive process 10 times to visualize
414 trajectories and possible convergence on a stable recommendation.

415

416 **Results**

417 **Application**

418 Our *Coastal Highly Density eDNA* fish dataset shows different compositional
419 characteristics. In 208 samples, the total richness of fish was 86 species. When the
420 dataset was divided into two treatments of high vs. low-vegetation, rarefied fish richness
421 difference between treatments was 11.4 species (81 in the high-vegetation sites, and an
422 average of 69.6 species in the rarefied low-vegetation sites - a fold difference of 0.22,
423 Fig. 3A). The ranges of effect sizes that contained enough replications to conduct our
424 quality-assurance power analysis by simulation were between -0.38 to 0.07 fold
425 difference in fish (coloured bins in Fig. 3A).

426 Using the datasets with 2 treatments (high vs. low vegetation), the minimum
427 detectable fish richness difference at a power of 0.8 decreased with coverage values
428 (Fig. 3C, high values on the y-axis correspond to smaller effect sizes). The power
429 analysis suggests that, assuming the richness difference between communities is 0.1 or
430 greater, a study of these communities would require sampling sufficient to reach 89%
431 coverage (or 98 samples in the high-vegetation community and 116 samples in the low-
432 vegetation community) to detect a true difference 80% of the time. For scale, this
433 assumed difference (0.1 log fold) is equivalent to a difference of 6 or 7 species between
434 the communities. Note, we consider only the application of these simulated relationships
435 to future samples or with assumed true effect sizes, to avoid the pitfalls of retroactive
436 power analysis [28]. Our analysis also suggests that with smaller sample sizes (e.g., 6-
437 10 samples), we would only be able to detect differences with 80% power when true
438 differences are large (e.g., log fold difference of 0.3, equivalent to a difference of ~19
439 species).

440 The 1-treatment analyses estimated relationships between coverage and
441 detectable effect size, providing estimates of sample sizes needed to detect differences
442 assuming different true (future) effect sizes. For fish, the full pilot dataset achieved a
443 coverage of 93% (Fig. 4C) with a detectable effect size of 0.09 fold difference (± 6
444 species from the total of 86, at a power of 0.8). In other words, the pilot study size can
445 detect changes of ± 6 starting from 86 species.

446 The minimum effect quality assurance power analysis provides
447 recommendations for sampling effort needed to detect changes above a user-defined
448 biologically-relevant threshold. When we set this to 0.1 log₂ fold-difference, or a 7.18%
449 increase, we found that, as expected, the effect size always needs to be greater than
450 this minimum level, in order to detect change of that scale or greater in 80% of the runs
451 when it was indeed that large (Fig. 5). Also as expected, the extent to which the true
452 effect needed to be greater than 0.1 decreased with sample size. As an example, for
453 the sample size of the pilot to have a power of 0.8 to detect a change greater than 0.1
454 log₂ change, the true effect size would need to be 0.19 log₂ fold change, or a 14%
455 difference (Fig. 5).

456

457 **Sensitivity Analysis**

458 We tested to see if our method converges on a stable coverage recommendation
459 from repeated application to subsequent samples as 'pilot datasets' for the next set of
460 samples. We used the fish dataset without treatments and arbitrarily chose 0.1 (~7%
461 increase) as the target effect size we want to detect, which is close to one standard
462 deviation from bootstrapped communities. Repeating the power analysis starting from
463 different pilot study sizes and successively adopting new recommendations showed
464 convergent coverage recommendations. Most pilot studies recommended sample sizes
465 between 95 and 145 in the test case (Fig. 6) after the second iteration. At very small
466 sample sizes (e.g., $k=6$), our analysis could not proceed due to cases in which the same
467 sample is drawn by chance 100% of the time within a bootstrap, which fixes
468 occupancies of all observed species to 100% (coverage $C(k)=1$ for all sample sizes k in
469 Eq.2), resulting in no possible control over coverage values. This illustrates that a
470 minimal number of samples is required for our power analysis, but it does not imply
471 these small studies have no statistical power. For example, Fig. 2 shows that with 6
472 samples, a fold difference of 0.35 can be detected, but such an analysis was only
473 possible from a larger pilot study. Our R package *BioDivPower* provides a warning
474 when the simulations could not generate coverage variations.

475

476 **Discussion**

477 We present a novel, non-parametric quality-assurance power analysis for
478 assessing sampling intensity needed within communities to robustly detect biodiversity
479 change between them. This addresses a unique problem in biodiversity detection, and

480 contrasts from traditional power analyses that rely on independent replicates of
481 treatment groups (e.g. Valdez et al). Our work supports previous findings of power
482 relationships increasing with sampling effort within communities [17], but provides a
483 method to apply the relationships from pilot data to the unsampled community to guide
484 future sampling. Our main analysis estimates the rate of detecting biodiversity
485 differences in the correct direction, but our modified *minimum-effect* power analysis
486 estimates sampling effort needed to detect a change greater than an ecologically-
487 relevant benchmark with sufficient power. This provides a practical approach to
488 designing field surveys to detect change relevant for conservation action [12]. Through
489 application to a large eDNA metabarcoding dataset, we found that the method yields
490 convergent results on the recommended sample coverage required to detect a fixed
491 change.

492 A key innovation that makes the power analysis applicable to real communities is
493 the use of coverage, which was used to both rarefy within measurements and scale
494 simulated findings for application. Our derivation shares the same mechanics as the
495 best available asymptotic richness estimator [19] which uses information from observed
496 species occupancies from the pilot study, and provides an analytical rarefaction curve
497 that is used to obtain measurable effect sizes at different levels of effort. We consider
498 coverage to be a dimensionless universal measure of sampling effort under the
499 assumption that the pattern of occupancies of unobserved species is the same as that
500 of observed species. We recognize that sampling biases stemming from rarity and
501 spatial heterogeneity will likely violate the above assumption, but we do not yet know
502 how to correct these second-order biases [19]. Therefore, we suggest an adaptive
503 monitoring approach where, after adopting the sample coverage recommended for the
504 next field survey, one repeats the power analysis and adopts the new recommendation
505 until a convergent pattern becomes clear (Fig. 5). In addition, our single-sample method
506 assumes stationarity in occupancy distributions and other features by which a treatment
507 can affect a species accumulation curve [29], and we therefore suggest conservatism
508 by erring on the side of more samples for planning before-after surveys where a pilot
509 dataset is not available for both treatments.

510 We caution that like all power analyses, our approach should not be used to
511 retroactively assess the power of a study to detect its own effect size [28,30]. This has
512 been shown to be invalid because the probability to detect an observed effect size is not
513 relevant after the effect has been observed (reviewed in [28]). Instead, our quality-
514 assurance power analysis makes it possible to estimate the sampling effort required to
515 detect a difference with a given probability under an assumed true effect in future
516 monitoring (i.e. a prospective power analysis [10,12]), or to estimate how large of an
517 effect would be necessary to detect a difference in the correct direction with an
518 acceptable error rate if the sample size is fixed (minimum detectable effect size [31]).
519 We also made it possible to assess sampling effort needed to detect a difference at

520 least as large as a given biologically-meaningful effect size (minimum effects power
521 analysis [12,16]). We expect our approach to be useful to guide monitoring planning,
522 especially with new technologies, when an effect size is in mind or a conservation action
523 is tied to an ecologically-relevant amount of change. With caution, our analysis could be
524 also used to ask if previous studies had sufficient power to detect a change of a given
525 *externally-expected* size, but refer to the previous literature on cautions against
526 incorrect retrospective power analysis [28,30].

527 We built our quality-assurance power analysis with relatively simple scenarios in
528 mind, but expect that our general approach could be extended further. For example, we
529 focused on detected changes in species richness, i.e. Hill number $q=0$ [32]. We believe
530 it could be extended to abundance-weighted biodiversity estimators ($q>0$) in the future,
531 possibly by implementing coverage estimators designed for $q>0$ [18]. Similarly, our
532 method is focused on differences in richness across two treatments, but we believe the
533 concept can be extended for detecting trends across multiple levels of treatments, or
534 detecting trends across multiple timepoints. For example, a previous study found
535 relatively lower sampling effort is needed to detect trends across multiple time points
536 compared to detecting direction from two timepoints [17]. Our method is not currently
537 designed to detect differences in community compositions (beta diversity); to do this, we
538 would need a measure of effort analogous to coverage but specific to beta-diversity
539 estimators, where future work is needed.

540 We made our method accessible through the R-package *BioDivPowerR* [21],
541 which generates quality-assurance power analyses for any monitoring program that
542 have collected pilot data with species and presence/absence occupancies across
543 samples. The package provides relationships between sampling effort, true effect sizes,
544 and rate of detecting a change in the correct direction or above a use-set difference.
545 Additionally, the package can be used to compare the costs of sampling associated with
546 statistical power that can be used to compare methods or different monitoring
547 technologies; we believe this will help decisions on the relative value of changing
548 sampling technologies.

549 Our method addresses an urgent need for designing and assessing biodiversity
550 monitoring programs. Efficient and scalable monitoring systems are paramount,
551 particularly as investments are made in new sampling technologies and survey
552 programs to improve both the detection and attribution of biodiversity change. So far,
553 such investments have often proceeded under the assumption that increased sampling
554 effort and improved technology will enhance detectability; our framework provides a
555 new quantitative basis for evaluating how much improvement is actually achieved.

556 **Acknowledgements**

558 We thank Evan Morien, Carolyn Prentice, Rute Carvalho, Kristin Meagher Robinson,
559 Rosie Savage, Kyle Hall, Ben Millard-Martin, and Kate Sheridan for their contributions to
560 field collections and processing of the pilot eDNA dataset. This work was funded by a

561 Genomics Applications Partnership Program grant from Genome Canada, and a
562 Convergent Research Themes grant from the Computational and Data Systems
563 Institute at McGill University

564

565 **Author Contributions**

566 ET designed the quality-control power analysis and wrote the original paper draft; JL
567 contributed substantially to methods development and created the R-package; JS
568 contributed to methods development, writing, and organized the team; ER, ML, RS, and
569 NJ contributed substantially to project conceptualization, multiple rounds of feedback on
570 methods, and contributed to writing.

571

572 **Funding Statement**

573 Eden Tekwa - Genomics Applications Partnership Program grant from Genome Canada
574 Jennifer Sunday - supported in part by the NSERC Canada Research Chairs Program
575 Jake Lawlor - supported by a Convergent Research Themes grant from the
576 Computational and Data Systems Institute at McGill University, and by a R  al Decoste-
577 Ouranos Partnership Training Award from the Province of Quebec.

578

579 **Data availability**

580 Pilot species occupancy dataset from eDNA survey are available here:

581 <https://doi.org/10.21966/vdyq-r660>.

582 R-package and associated GitHub repository will be made available upon project
583 publication

584 Code to replicate analyses for this manuscript will be made available upon project
585 publication

586

587 **Competing Interests**

588 We have no competing interests.

589

590

591

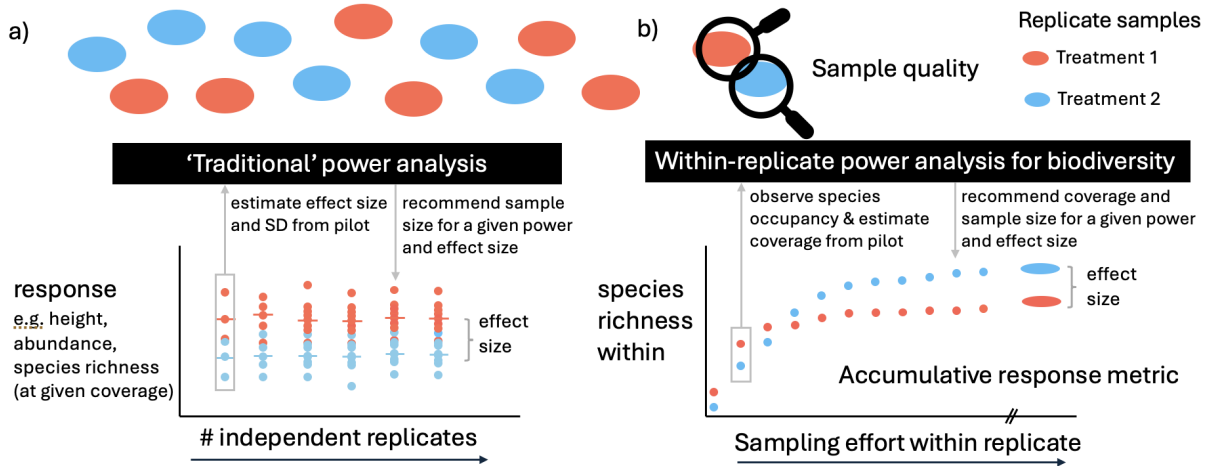
592 **References**

- 593 1. Gonzalez A *et al.* 2020 Scaling-up biodiversity-ecosystem functioning research. *Ecol. Lett.*
594 **23**, 757–776. (doi:10.1111/ele.13456)
- 595 2. Dasgupta P, Levin S. 2023 Economic factors underlying biodiversity loss. *Philos. Trans. R.*
596 *Soc. B Biol. Sci.* **378**, 20220197. (doi:10.1098/rstb.2022.0197)
- 597 3. Valdez JW, Callaghan CT, Junker J, Purvis A, Hill SLL, Pereira HM. 2023 The
598 undetectability of global biodiversity trends using local species richness. *Ecography* **2023**,
599 e06604. (doi:10.1111/ecog.06604)
- 600 4. Magurran, McGill. 2011 *Biological Diversity: Frontiers in Measurement and Assessment*.
601 Oxford University Press.
- 602 5. Cardinale BJ, Gonzalez A, Allington GRH, Loreau M. 2018 Is local biodiversity declining or
603 not? A summary of the debate over analysis of species richness time trends. *Biol. Conserv.*
604 **219**, 175–183. (doi:10.1016/j.biocon.2017.12.021)
- 605 6. Gonzalez A *et al.* 2023 A global biodiversity observing system to unite monitoring and guide
606 action. *Nat. Ecol. Evol.* **7**, 1947–1952. (doi:10.1038/s41559-023-02171-0)
- 607 7. Convention on Biological Diversity (CBD). 2022 Convention on Biological Diversity.
608 Kunming-Montreal Global Biodiversity Framework (Decision CBD/COP/DEC/15/4).
- 609 8. Wunder S *et al.* 2025 Biodiversity Credits: An Overview of the Current State, Future
610 Opportunities, and Potential Pitfalls. *Bus. Strategy Environ.* **34**, 8470–8499.
611 (doi:10.1002/bse.70018)
- 612 9. Jacquemont J, Blasiak R, Le Cam C, Le Gouellec M, Claudet J. 2022 Ocean conservation
613 boosts climate change mitigation and adaptation. *One Earth* **5**, 1126–1138.
614 (doi:10.1016/j.oneear.2022.09.002)
- 615 10. Cohen J. 1988 *Statistical Power Analysis for the Behavioral Sciences*. 0 edn. Routledge.
616 (doi:10.4324/9780203771587)
- 617 11. Steidl RJ, Hayes JP, Schauber E. 1997 Statistical Power Analysis in Wildlife Research. *J.*
618 *Wildl. Manag.* **61**, 270. (doi:10.2307/3802582)
- 619 12. Morrison LW. 2007 Assessing the Reliability of Ecological Monitoring Data: Power Analysis
620 and Alternative Approaches. *Nat. Areas J.* **27**, 83–91. (doi:10.3375/0885-
621 8608(2007)27[83:ATROEM]2.0.CO;2)
- 622 13. Nichols JD, Kendall WL, Boomer GS. 2019 Accumulating evidence in ecology: Once is not
623 enough. *Ecol. Evol.* **9**, 13991–14004. (doi:10.1002/ece3.5836)
- 624 14. Chase JM, McGill BJ, McGlinn DJ, May F, Blowes SA, Xiao X, Knight TM, Purschke O,
625 Gotelli NJ. 2018 Embracing scale-dependence to achieve a deeper understanding of
626 biodiversity and its change across communities. *Ecol. Lett.* **21**, 1737–1751.
627 (doi:10.1111/ele.13151)

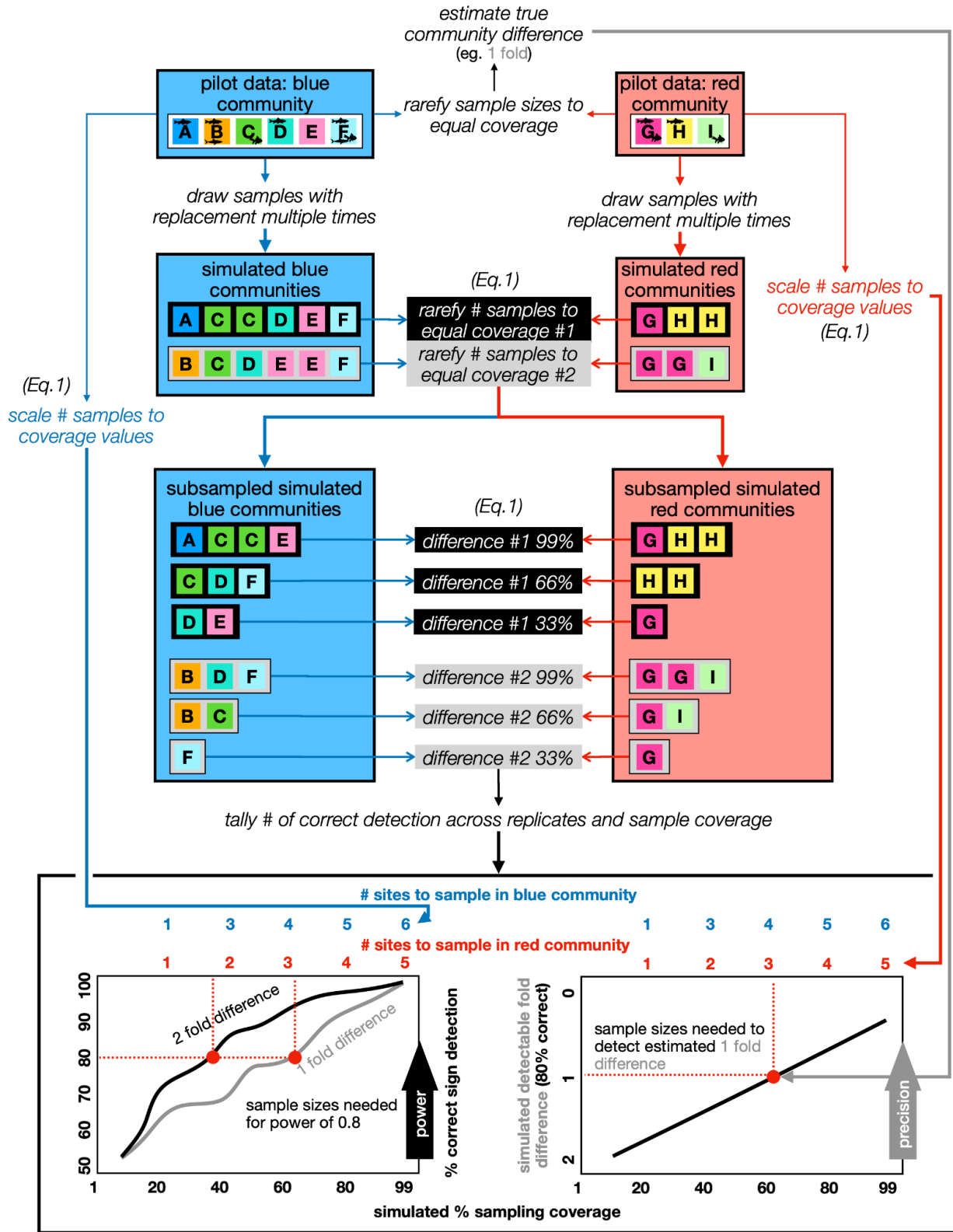
- 628 15. Davies GM, Gray A. 2015 Don't let spurious accusations of pseudoreplication limit our ability
629 to learn from natural experiments (and other messy kinds of ecological monitoring). *Ecol.*
630 *Evol.* **5**, 5295–5304. (doi:10.1002/ece3.1782)
- 631 16. Myers B, Murphy KR. 2023 *Statistical Power Analysis: A Simple and General Model for*
632 *Traditional and Modern Hypothesis Tests, Fifth Edition*. 5th edn. New York: Routledge.
633 (doi:10.4324/9781003296225)
- 634 17. Gwinn DC, Allen MS, Bonvechio KI, V. Hoyer M, Beesley LS. 2016 Evaluating estimators of
635 species richness: the importance of considering statistical error rates. *Methods Ecol. Evol.*
636 **7**, 294–302. (doi:10.1111/2041-210X.12462)
- 637 18. Chao A *et al.* 2020 Quantifying sample completeness and comparing diversities among
638 assemblages. *Ecol. Res.* **35**, 292–314. (doi:10.1111/1440-1703.12102)
- 639 19. Tekwa EW, Whalen MA, Martone PT, O'Connor MI. 2023 Theory and application of an
640 improved species richness estimator. *Philos. Trans. R. Soc. B Biol. Sci.* **378**, 20220187.
641 (doi:10.1098/rstb.2022.0187)
- 642 20. Moin P, Chan WHR. 2024 *Fundamentals of Turbulent Flows*. 1st edn. Cambridge University
643 Press. (doi:10.1017/9781009431385)
- 644 21. BioDivPowerR: Analyze Power to Detect Change in Biodiversity Sampling Schemes. R
645 package version 0.1.0, *Package DOI withheld for anonymity. Instead, a local,*
646 *anonymized version of the package is available as a compressed file within the project*
647 *repository uploaded as a supplementary file for reviewers (BioDivPower-*
648 *87C7_June02).*
- 649 22. Gotelli NJ, Colwell RK. 2001 Quantifying biodiversity: procedures and pitfalls in the
650 measurement and comparison of species richness. *Ecol. Lett.* **4**, 379–391.
651 (doi:10.1046/j.1461-0248.2001.00230.x)
- 652 23. Webb TJ, Noble D, Freckleton RP. 2007 Abundance–occupancy dynamics in a human
653 dominated environment: linking interspecific and intraspecific trends in British farmland and
654 woodland birds. *J. Anim. Ecol.* **76**, 123–134. (doi:10.1111/j.1365-2656.2006.01192.x)
- 655 24. Zuckerberg B, Porter WF, Corwin K. 2009 The consistency and stability of abundance–
656 occupancy relationships in large-scale population dynamics. *J. Anim. Ecol.* **78**, 172–181.
657 (doi:10.1111/j.1365-2656.2008.01463.x)
- 658 25. *DOI of data redacted for anonymity. Data are instead available within the anonymized*
659 *project R repository, uploaded as a supplementary file for review. Details of data*
660 *collection methods are also available with the online dataset; an anonymized version of*
661 *these methods are provided with the data within the anonymized project repository for*
662 *review (Anonymized eDNA data methods.doc, within BioDivPower_manuscript_figs-*
663 *master).*
- 664 26. Miya M *et al.* 2015 MiFish, a set of universal PCR primers for metabarcoding environmental
665 DNA from fishes: Detection of more than 230 subtropical marine species. *R. Soc. Open Sci.*
666 (doi:10.1098/rsos.150088)

- 667 27. Callahan BJ, McMurdie PJ, Rosen MJ, Han AW, Johnson AJA, Holmes SP. 2016 DADA2:
668 High-resolution sample inference from Illumina amplicon data. *Nat. Methods* **13**, 581–583.
669 (doi:10.1038/nmeth.3869)
- 670 28. Dormann CF, Ellison AM. 2025 *Statistics by Simulation: A Synthetic Data Approach*.
671 Princeton and Oxford: Princeton University Press.
- 672 29. Chase JM, Knight TM. 2013 Scale-dependent effect sizes of ecological drivers on
673 biodiversity: why standardised sampling is not enough. *Ecol. Lett.* **16**, 17–26.
674 (doi:10.1111/ele.12112)
- 675 30. O’Keefe DJ. 2007 Brief Report: Post Hoc Power, Observed Power, A Priori Power,
676 Retrospective Power, Prospective Power, Achieved Power: Sorting Out Appropriate Uses of
677 Statistical Power Analyses. *Commun. Methods Meas.* **1**, 291–299.
678 (doi:10.1080/19312450701641375)
- 679 31. Rotenberry JT, Wiens JA. 1985 Statistical Power Analysis and Community-Wide Patterns.
680 *Am. Nat.* **125**, 164–168. (doi:10.1086/284335)
- 681 32. Jost L. 2006 *Entropy and diversity*. *Oikos* **113**, 363–375. (doi:10.1111/j.2006.0030-
682 1299.14714.x)
- 683

684 **Figures**



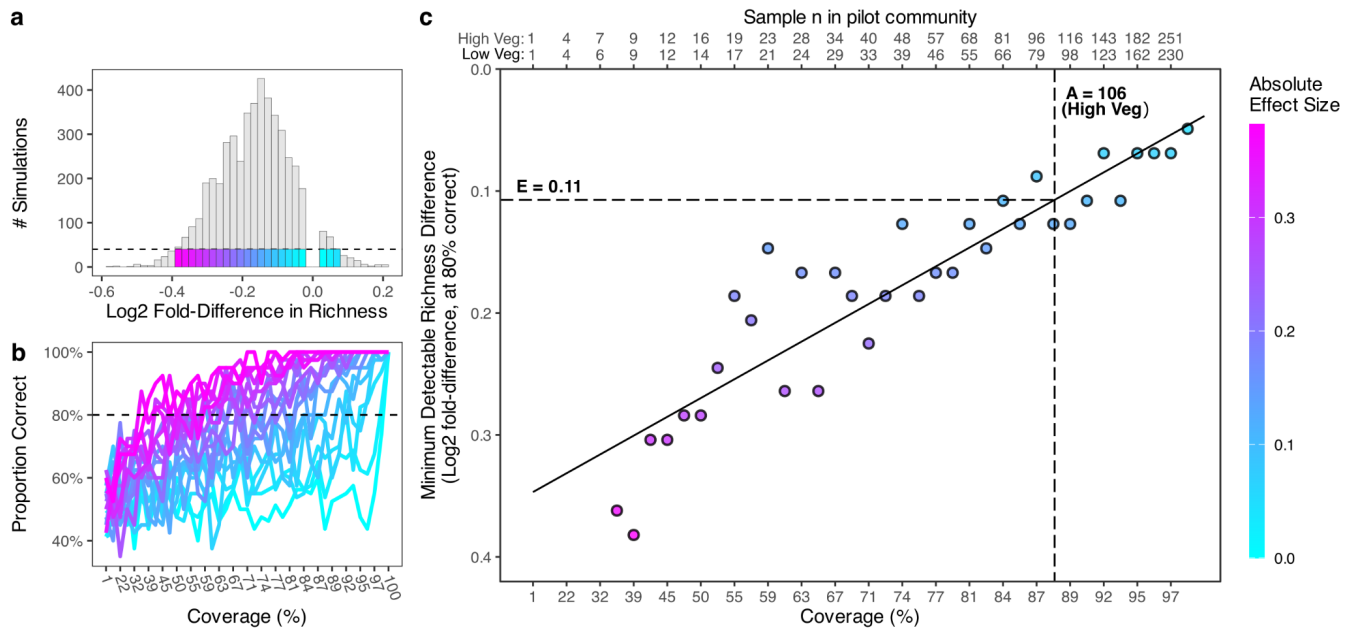
685
 686 **Fig 1. Traditional power analysis vs. within-replicate power analysis for**
 687 **biodiversity assessment.** (a) In traditional power analysis, a small number of
 688 independent replicates, across two treatments (red vs. blue) are used to estimate the
 689 mean effect and second moment of effect (SD, standard error), to assess
 690 recommended sample size at a given effect size and power. Points represent estimates
 691 within each replicate, and horizontal lines represent mean across replicates (sample
 692 mean); mean and standard error are invariant with increasing numbers of replicates.
 693 Response values are any metric assessed within a replicate, and include species
 694 richness at a given level of sampling, or sample coverage. (b) In the quality-assurance
 695 power analysis for biodiversity, we focus on sampling effort *within* replicates, i.e. sample
 696 quality. Here, we are interested in how sampling effort (such as volume, area, or
 697 number of sites visited) improves the estimated richness difference between
 698 communities. An accumulated response metric like species richness estimates the true
 699 (census) richness with bias that is reduced with increased coverage. In this power
 700 analysis, we assess occupancy and estimate coverage from a pilot study and use
 701 simulations from the data to recommend a sample size at a given power and effect size.



702
 703 **Figure 2. Power analysis from two-treatment data.** Coloured boxes labelled A to F
 704 represent independent replicated samples from pilot data, and blue vs. red communities
 705 represent two communities with different treatments that are being compared.

706 Workflows and calculations are described in methods. Bottom plots summarize
707 hypothetical simulation results for estimates of power/accuracy (LHS) and detectable
708 effect size/precision (RHS), each as a function of sampling intensity in the pilot area
709 given the same sampling method. LHS shows the percent of times that an observed
710 effect size (differences in bottom red and blue boxes) over a range of coverage (from
711 subsampling) is in the correct direction as the simulated "truth" (eg. 1 or 2 fold
712 difference, measured between the rarefied simulated communities). RHS shows the
713 minimum effect size that can be correctly detected at a given level of power (eg. 0.8).
714 Note here that the y-axis is reversed so that increases in precision (smaller detectable
715 effect sizes) are in the positive direction. All sample-to-coverage calculations, including
716 rarefactions, subsampling, and scaling to final applications, use Eq. 2 (samples vs.
717 coverage relationship).
718

719
720



721

722

723

724

725

726

727

728

729

730

731

732

733

734

735

736

737

738

739

740

741

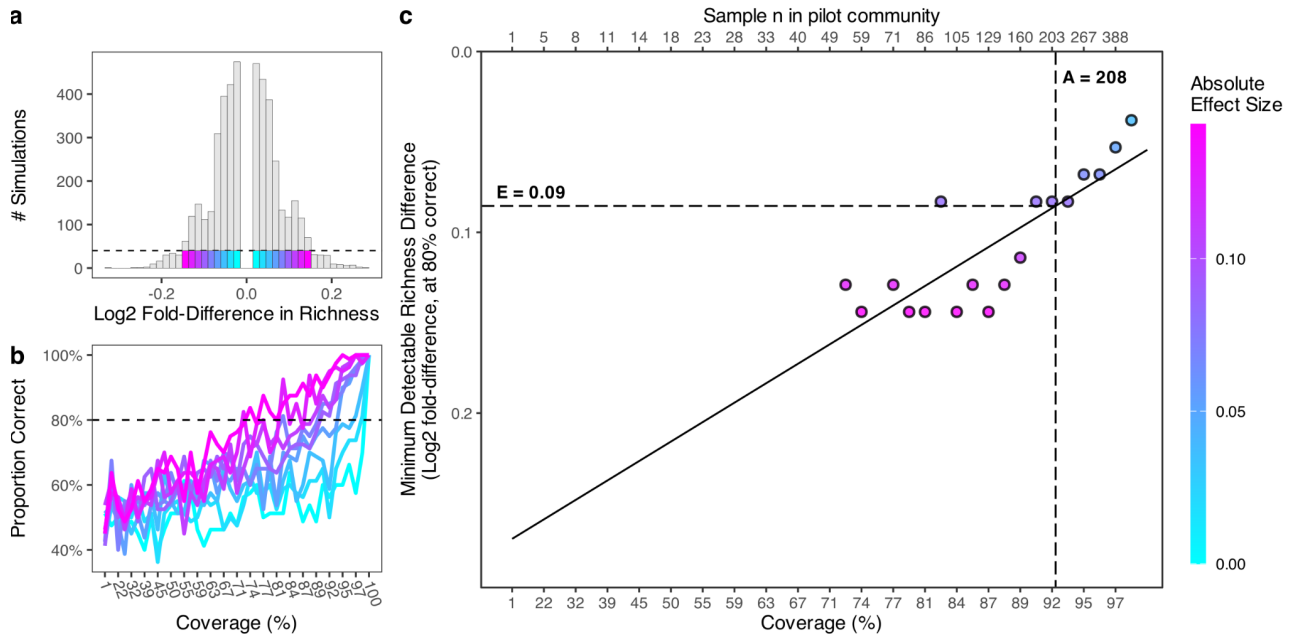
742

743

744

745

Figure 3. Power analysis results for two-treatment fish data. Simulations are based on data from two marine communities that we wish to compare: fish species at high-vegetation sites vs. fish species at low-vegetation sites ($n = 106$ & 102 , rarefied richness = 81 & 69.6). **A.** Simulated richness log-ratio from random draws of each community ($\log_2(\text{high-veg richness} / \text{low-veg richness})$), with zero-change simulations removed, provides a range of effect sizes with an equal number of simulations in each, colour coded from blue to magenta in the order of increasing absolute effect size; dashed line shows number of simulations retained (user-set, default is 40 replicates). Retained simulations are downsampled once to provide an equal number of replicates for each simulated effect size for subsequent power analysis. **B.** Relationship between power to detect change and richness coverage from simulations. Percent of simulations with correctly detected difference sign is shown as a function of down-sampled % coverage, assuming observed community is 100% coverage. Lines show results for each effect size bin; dashed line indicates a target power (user-set, default 80%). **C.** Minimum detectable effect size as a function of sample size and % coverage in each community, at 80% power. The upper x-axis shows sample size recommended from each community, while lower x-axis shows the corresponding equal sample coverage. The y-axis is reversed in scale to show smaller differences towards the top, indicating greater precision. The solid line is an ordinary least-squares regression fit to the data; vertical dashed line is achieved equal coverage from the original dataset; the horizontal dashed line is the minimum detectable effect size at achieved equal coverage. The regression estimates that the current sample size (A) can detect an effect size (E) higher than 0.11, 80% of the time in the correct direction.



746

747

748

749

750

751

752

753

754

755

756

757

758

759

760

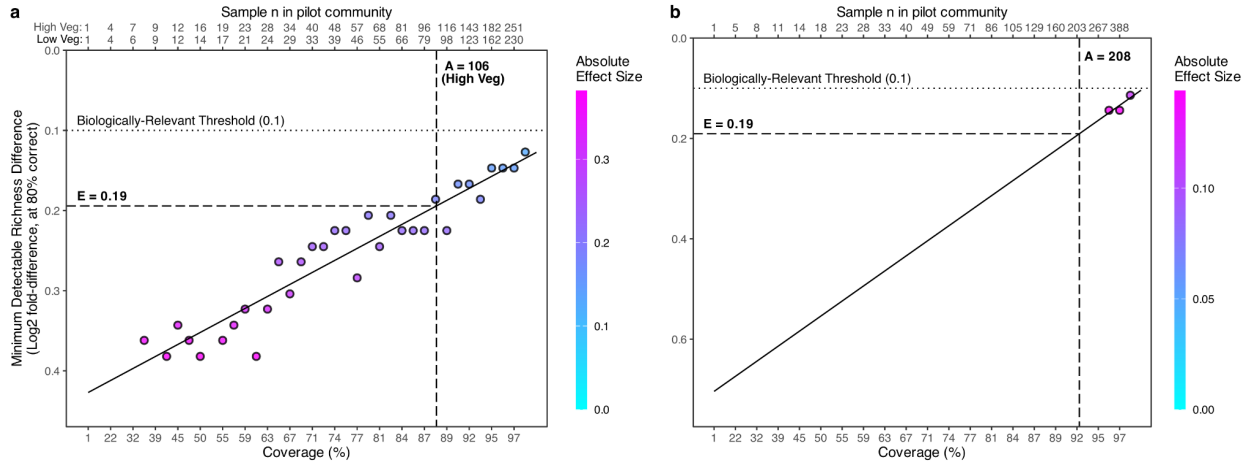
761

762

763

764

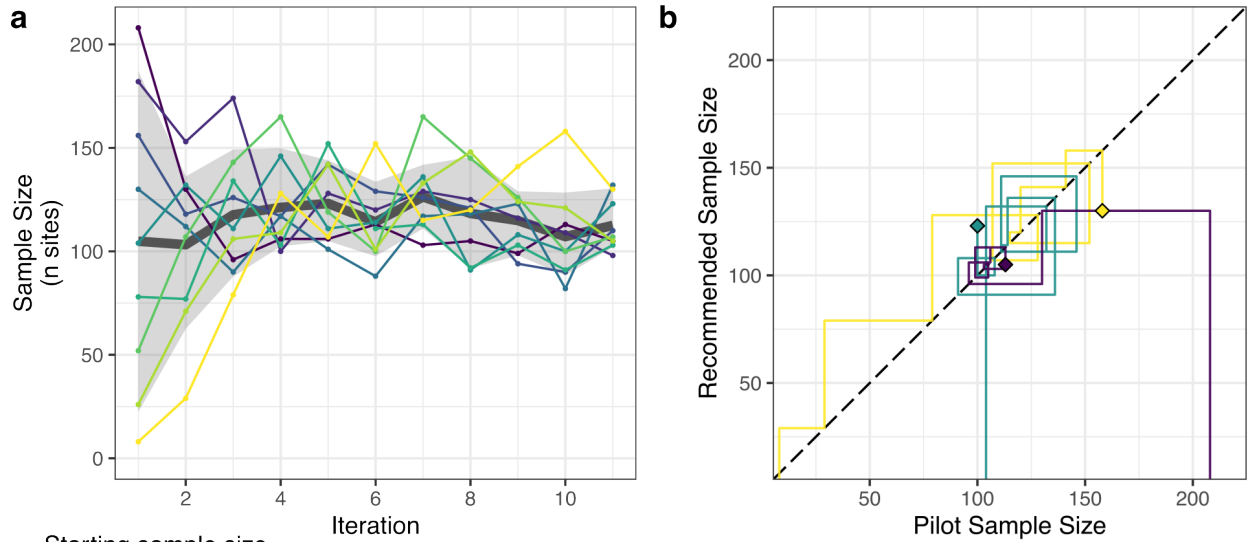
Figure 4. Power analysis results for single-treatment fish data. Simulations are based on data from a single community in time. Species richness is 86. **A.** Simulated rarefied richness fold difference ($\log_2(\text{richness at time 2} / \text{richness at time 1})$) colour coded from blue to magenta in the order of increasing effect size, with data from time points being two bootstraps from the original dataset. The range of differences (binned effect sizes) with 40 or more replicates (user-set, dashed line) is used for subsequent power analysis. **B.** Percent of simulations with correct difference sign across true effect size magnitude bins and coverage. **C.** Detectable effect size as a function of sample coverage at a power of 0.8. The fitted line is from an ordinary least-squares regression. The regression estimates that the current sample size (A) can detect an effect size (E) higher than 0.09, 80% of the time in the correct direction.



765
766
767
768
769
770
771
772
773
774
775
776

Figure 5. Minimum Effects Power Analysis for detecting relevant thresholds.

Recommendations for sampling effort to detect changes above a user-defined biologically-relevant threshold (here, 0.1 log₂ fold-difference, a 7.18% increase) for the two-treatment method (a) and the single-treatment method (b) using the example fish eDNA dataset. Plots mirror Figures 3c and 4c, respectively, but here power is set as a 0.8 probability of detecting change at least as great and in the correct direction as the user-defined minimum effect size for biological or conservation relevance. Because of measurement error, the effect size always needs to be greater than this minimum (threshold) level, and needs to be increasingly greater as sample size decreases (higher values are lower on y-axis as samples size is decreased on x-axis).



777
 778 **Figure 6. Convergence of sample size recommendations.** Recommendations from
 779 repeating simulated power analysis and adapting through iterations starting from 9
 780 different pilot study sizes. Simulations are from the 1-treatment dataset with an objective
 781 of detecting a log-fold difference of 0.1 with a power of 0.8. **A.** Recommended sample
 782 size through 11 adaptive iterations. Black trend shows the mean recommendation over
 783 iterations, and shade shows the 10th and 90th quantiles. **B.** Phase diagram relating pilot
 784 sample size (iteration t) to recommended sample size (iteration $t+1$) for three
 785 trajectories (origins o ; pilot sample sizes 8, 104, 208 sites) over 11 iterations.