1 **COI metabarcoding primer choice affects richness and recovery of indicator taxa**

2 **in freshwater systems**

3

4 Mehrdad Hajibabaei[1*], Teresita M. Porter[1,2], Michael Wright[1], Josip Rudar[1]

5

6 [1]Centre for Biodiversity Genomics @ Biodiversity Institute of Ontario & Department of

7 Integrative Biology, University of Guelph, Guelph, Ontario, Canada, N1G 2W1.

8 [2]Natural Resources Canada, Canadian Forest Service, Great Lakes Forestry Centre,

9 1219 Queen St. East, Sault Ste. Marie, Ontario, Canada, P6A 2E5.

10

11 *Corresponding author (mhajibab@uoguelph.ca)

12

15

**Abstract**

DNA-based biodiversity analysis has gained major attention due to the use of high throughput sequencing technology in approaches such as mixed community or environmental DNA metabarcoding.  Many cytochrome c oxidase subunit I (COI) primer sets are now available for such work.  The purpose of this study is to look at how COI primer choice affects the recovery of arthropod richness, beta diversity, and recovery of site indicator taxa in benthos kick-net samples typically used in freshwater biomonitoring.  We examine 6 commonly used COI primer sets, on samples collected from 6 freshwater sites.  Richness is sensitive to primer choice and the combined use of additional multiple COI amplicons recovers higher richness.  Thus, to recover maximum richness, multiple primer sets should be used with COI metabarcoding.  Samples consistently cluster by site regardless of amplicon choice or PCR replicate.  Thus, for broadscale community analyses, overall beta diversity patterns are robust to COI marker choice.  Additionally, the recovery of traditional freshwater bioindicator assemblages such as Ephemeroptera, Trichoptera, Plectoptera, and Diptera may not fully capture the diversity of broadscale arthropod site indicators that can be recovered from COI metabarcoding.  Based on these results, studies that use different COI amplicons may not be directly comparable.  This work will help future biodiversity and biomonitoring studies develop not just standardized, but optimized workflows that maximize taxon-detection or order taxa along gradients.

**Introduction**

2

39

40      DNA-based biodiversity analysis has gained major attention due to the use of

41  high throughput sequencing technology in approaches such as mixed community or

42  environmental DNA metabarcoding (Hajibabaei, Shokralla, Zhou, Singer, & Baird, 2011;

43  Taberlet, Coissac, Pompanon, Brochmann, & Willerslev, 2012). Data generation

44  typically involves DNA extraction from an environmental sample such as water, soil or

45  collected biomass (e.g. benthic kicknet, malaise trap) followed by PCR amplification of

46  one or more taxonomic markers such as the COI DNA barcode region and subsequent

47  high throughput sequencing and bioinformatic analysis of marker gene sequences.

48  Resulting sequences are then assigned to sequence clusters (Operation Taxonomic

49  units, OTU; Exact Sequence Variants, ESV) and/or taxonomic names (Callahan,

50  McMurdie, & Holmes, 2017).  Sequence clusters and taxonomic lists obtained are used

51  in various statistical analyses for assessing different aspects of biodiversity such as

52  species richness or distribution, community composition, and functional diversity(Porter

53  & Hajibabaei, 2018c).  In practice, these questions are often geared towards identifying

54  assemblages or specific target taxa. Biodiversity information gained can contribute to

55  ecological investigations and applications such as biomonitoring as part of

56  environmental assessment programs (Baird & Hajibabaei, 2012; Leese et al., 2018).

57

58      A major step in obtaining sequence data from environmental samples involves

59  PCR amplification of target marker gene(s).  An important consideration in this multi-

60  template PCR step is the choice of primer sets.  It has been shown that primers can

61  differentially bind to template DNA and can result in both qualitative and quantitative

3

62    biases (Suzuki & Giovannoni, 1996; Polz & Cavanaugh, 1998; L. J. Clarke, Soubrier,

63    Weyrich, & Cooper, 2014).  Most previous metabarcoding studies use a single primer

64    set for generating sequence data from communities, but there is ample evidence that

65    multiple amplifications with different primer sets can provide better biodiversity coverage

66    from environmental samples (ex. J. Gibson et al., 2014).

67

68        Current methods for biomonitoring especially in freshwater typically rely on key

69    taxonomic groups that are considered ecological bioindicators. For example,

70    Ephemeroptera, Plecoptera, Trichoptera are known to be sensitive to water pollution

71    whereas Chironomidae (Diptera) have been shown to be tolerant to high levels of

72    pollution (Buss, Baptista, Silveira, Nessimian, & Dorvillé, 2002; Bonada, Prat, Resh, &

73    Statzner, 2006).  Here we collectively refer to this assemblage as the EPTC.  Because

74    of the difficulties associated with morphological identification of larval samples from

75    benthos, samples are generally identified to family or genus level. Sorting and

76    identifying individual samples from benthos poses a serious challenge in executing

77    large-scale biomonitoring programs.  With the advancement of genomics methods such

78    as DNA metabarcoding, sequence data is generated from whole communities without

79    the need to isolate individuals.  Therefore, the analysis can go beyond the target

80    assemblages such as EPTC.

81
82        The objective of this study was to test the performance of several newly

83    published COI metabarcode primers to detect freshwater benthic invertebrates.  We

84    wanted to determine the impact of primer choice on several components of diversity:

85    richness, beta diversity, and recovery of bioindicators.  We tested a total of 6 partial COI

4

86    metabarcode amplicons, including the two amplicons we have used routinely for

87    freshwater invertebrate monitoring.

88

89    **Methods**

90    *Field methods*

91    Six benthic invertebrate communities were sampled from shallow streams across

92    the City of Waterloo (Ontario, Canada) using a modified travelling kick-and-sweep

93    technique outlined in the Ontario Benthos Biomonitoring Network protocol (Jones,

94    Somers, Craig, & Reynoldson, 2007) (Table S1). Briefly, wetted width was measured

95    and used to calculate the number of return trips required to sample a 10m transect of

96    the stream specifically targeting a riffle habitat. Prior to sampling D-nets were

97    decontaminated by soaking them in a 10% bleach solution for 15 min, rinsing with

98    tapwater, and drying them overnight, A clean 500 μm mesh D-net was held downstream

99    to the person sampling, with the opening of the net facing the person sampling.

100    Substrate was disturbed by kicking the substrate at a constant effort for 3 minutes

101    across the 10 m transect dislodging invertebrates and allowing the flowing water to

102    guide the dislodged macroinvertebrates into the net. The samples were transferred from

103    the net to a clean 1 L polyethylene bottle, preserved with 80% ethanol and stored at -

104    20°C until further processing in the lab.

105

106    *Molecular Biology Methods*

107    **DNA Extraction.** Samples were homogenized separately in a clean blender

108    (decontaminated thoroughly with Eliminase solution (Decon Labs: King of Prussia, PA,

5

109 USA) (Black and Decker Model: BL2010BGC), distributing 50 mL of the homogenate

110 into 6 sterile conical tubes for each sample. Samples were centrifuged at 2400 x g for

111 2min to collect homogenate at the bottom of the tube, and excess preservative ethanol

112 was removed. Samples were covered and incubated at 65°C until residual ethanol was

113 evaporated (roughly 5-8 hours). DNA was extracted using Qiagen's DNeasy PowerSoil

114 kit (Toronto, Canada. Product Ref: 12888). Samples were lysed overnight (~15 hr).

115 Following lysis, samples were extracted according to the manufacturer's protocol,

116 eluting with 30 µL molecular biology grade water. All extractions included a negative

117 control where no sample was included.

118 **Polymerase Chain Reaction**. The six amplicons from CO1 barcode region used in this

119 study are shown in Figure 1. The primers were aligned against the *Drosophila yakuba*

120 COI barcode region obtained from GenBank accession X03240 using Mesquite v3.10

121 (Maddison & Maddison, 2015). COI secondary structure, 6 alpha helices, from *Bos*

122 *taurus* were obtained from UniProt accession P00396. All samples were amplified for

123 six primer sets according to their published amplification regime (Table 2) with the

124 exception that a two-step PCR was used for all reactions (first PCR using untailed

125 primers, second PCR using Illumina adapter-tailed primers), even if a one-step PCR

126 was used in the original protocol. PCRs were run in duplicate with a negative control.

127 Amplification success was confirmed through gel electrophoresis (not pictured).

128 Amplicons were purified using a MinElute PCR Purification kit, quantified on a

129 TBS-380 Mini-Fluorometer (Turner Biosystems Sunnyvale California, United States)

130 using a Quant-iT PicoGreen dsDNA assay (Invitrogen Waltham Massachusetts, United

131 States Product Ref: P11496). The concentration of each purified sample was

6

132    normalized across samples and primer sets were pooled for each sample. Although all

133    primers were tested using the PowerSoil kit, samples extracted with the NucleoSpin

134    Tissue Kit and amplified for BR5 and F230R primer sets were also sequenced as a

135    comparison for their past use (J. Gibson et al., 2015). PCR replicates were sequenced

136    separately from each other. Pooled samples were indexed using Illumina's Nextera

137    index kit (San Diego, California, United States Product Ref: FC-121-1011). All indexed

138    samples were pooled, purified through magnetic bead purification, quantified using the

139    PicoGreen dsDNA assay, and average fragment length for the library was determined

140    on an Agilent Bioanalyzer 2100 (Santa Clara, California, United States. Product

141    ref:G2939BA) using the Agilent DNA 7500 assay chip (Product Ref: 5067-4627). The

142    library was diluted then sequenced using Illumina's MiSeq v3 sequencing chemistry kit

143    (2x300 cycle. Product Ref: MS-102-3003) on an Illumina MiSeq, comprising

144    approximately half of a sequencing run.

145

146    *Bioinformatic processing*

147         Raw sequences were processed with the SCVUC COI metabarcode pipeline

148    v2.1 available from Github at https://github.com/EcoBiomics-

149    Zoobiome/SCVUC_COI_metabarcode_pipeline . The acronym SCVUC stands for the

150    major programs or algorithms used for bioinformatic processing: "S" – SEQPREP, "C" –

151    CUTADAPT, "V" – VSEARCH, "U" – UNOISE, "C" – COI Classifier.  Briefly, this semi-

152    automated pipeline is described below.  Jobs were spread across multiple cores using

153    GNU Parallel (Tange, 2011).  Raw compressed fastq Illumina read files were paired

154    using SeqPrep specifying a minimum Phred score of 20 at the ends of the reads and an

155    overlap of at least 25 bp  (St. John, 2016).  The following steps were conducted

156    separately for each of the 6 amplicons tested in this study.  Primers were trimmed using

157    CUTADAPT v1.14 and reads were retained if they were at least 150 bp long after

158    trimming, had a minimum Phred score of 20 at the ends of the reads, and contained no

159    more than 3 N's.  CUTADAPT was also used to convert fastq files to FASTA files

160    (Martin, 2011).  The individual sample files were combined into a single file for global

161    ESV generation.  VSEARCH v2.4.2 was used to dereplicate the data (get the unique

162    reads) using the –derep_fulllength option (Rognes, Flouri, Nichols, Quince, & Mahé,

163    2016).  The USEARCH v10.0.240 unoise3 algorithm was used to denoise the reads

164    (Edgar, 2016).  This involved the removal of any contaminating PhiX reads (carry over

165    from Illumina sequencing), prediction and removal of sequences with errors, removal of

166    putative chimeric sequences, and removal of rare sequences.  We defined rare

167    sequences to be those clusters comprised of less than 3 reads (singletons and

168    doubletons) (Brown et al., 2015; Tedersoo et al., 2010).  We used this set of exact

169    sequence variants (ESVs) as a reference, and all primer trimmed reads were mapped to

170    this reference set with an identity of 1.0 (100% sequence similarity) to generate a

171    sample x ESV table.  The COI Classifier v3.2, that uses the Ribosomal Database

172    Project naïve Bayesian classifier v2.12 with a custom COI reference set, was used to

173    taxonomically assign the ESVs (Porter & Hajibabaei, 2018a; Wang, Garrity, Tiedje, &

174    Cole, 2007).  Taxonomic assignments were mapped to ESVs detected in each sample

175    with a custom Perl script.  The final taxonomy table for each primer was concatenated.

176

177    *Data analysis*

178    The final taxonomy table above was formatted in R v3.4.3 in RStudio v1.1.419

179    (RStudio Team, 2016; R Core Team, 2017).  Custom scripts are available from GitHub

180    at URL.  Data was summarized multiple taxonomic ranks.  High confidence taxonomic

181    assignments were retained by filtering for bootstrap support cutoffs >= 0.30 at the genus

182    rank and >= 0.20 at the family rank.  Using these cutoffs ensures that 95-99% of the

183    taxonomic assignments are correct, assuming our query taxa are in the reference

184    database (Porter & Hajibabaei, 2018a).  We retained taxa at the species rank with a

185    bootstrap support cutoff >= 0.70.  Assuming our query species are present in the

186    reference database, this should ensure that at least 95% of species level assignments

187    are correct.  To check whether we had sufficient sequencing depth, we used the

188    package VEGAN v2.5-2 to plot rarefaction curves using the 'rarecurve' function

189    (Oksanen et al., 2018).  Curves that reach a plateau show saturated sequencing.  To

190    account for variable library sizes, reads/library were rarefied down to the 15$^{th}$ percentile

191    library size using the 'rrarefy' function (S. Weiss et al., 2017).

192    We compared average richness recovered from each amplicon using the VEGAN

193    'specnum' function and total richness was plotted with ggplot2 (Wickham, 2009).

194    Richness data was checked for normality using visual distribution plots (ggdensity and

195    ggqqplot, not shown) as well as using the Shapiro-Wilk test of normality (W=0.97,

196    p=0.36) and this data was treated as normally distributed in comparisons (Shapiro &

197    Wilk, 1965).  We compared average richness using paired t-tests with the Holm

198    adjustment for multiple comparisons.

199    There is uncertainty in how to interpret read abundance from arthropod

200    metabarcoding studies due to unexpected template to product ratios after PCR due to

201    stochasticity and GC content (Polz & Cavanaugh, 1998) as well as the effect of primer

202    bias and body size variation across life stages and species that can vary by orders of

203    magnitude and affect recovery (Elbrecht & Leese, 2015).  As a result, we chose to

204    transform read abundance into presence-absence data for all subsequent analyses.

205    We checked for correlations in the presence-absence of ESVs recovered from DNA

206    extractions processed with soil or tissue kits as well as between two PCR replicates

207    using the 'cor' and 'corrplot' functions in R (Wei & Simko, 2017).

208         Indicator species can be used as a proxy to indicate differences among sites or

209    conditions (De Cáceres & Legendre, 2009).  For example, in freshwater systems, the

210    diversity of EPTC taxa have been used as water quality indicators (Emilson et al.,

211    2017).  To test whether the recovery of indicator taxa depends on the COI amplicon

212    used for the analysis, we performed indicator species analyses using the

213    INDICSPECIES package in R and the 'multipatt' function with default settings (De

214    Cáceres & Legendre, 2009).  The six sites were used as groups for the analysis at

215    multiple ranks and significant site indicators were selected if the resulting p-value was

216    <= 0.05. We tested how often traditional EPTC are recovered with COI metabarcode

217    data by repeating the analysis using all arthropod ESVs and just the ESVs assigned to

218    EPTC.

219         To test whether sample clusters are affected by COI amplicon choice or PCR

220    replicate, we used non-metric multidimensional scaling.  Plots were created using the

221    vegan 'metaMDS' function using the default settings with two dimensions (scree plot not

222    shown) and dissimilarities were calculated using the method 'bray' for binary data

223    (Sorensen dissimilarity) and plotted with ggplot.  Goodness of fit was calculated using

224    the VEGAN 'goodness' function.  To check whether we had homogenous dispersion of

225    dissimilarities, an assumption of permutational multivariate analysis of variance

226    (PERMANOVA), we created a dissimilarity matrix with the VEGAN 'vegdist' function,

227    then calculated beta dispersion using the 'betadisper' function in R.  We tested for

228    significant heterogeneity using analysis of variance (ANOVA) in R.  We checked for

229    significant interactions among sites, amplicons, and replicates as well as the

230    significance of group clusters with PERMANOVA using the VEGAN 'adonis' function

231    with 999 permutations.

232

233    **Results**

234        A total of 9,980,584 x 2 paired-end reads were generated for this study and they

235    have been deposited to the NCBI SRA: accession numbers # (Table S2).  After pairing

236    and primer trimming we retained a total of 7,619,108 reads.  A summary of ESV counts

237    for all taxa are shown in Table S3.  About 23% of raw reads were retained in the

238    denoised set of ESVs whereas the difference was removed during denoising as putative

239    sequence errors, chimeras, PhiX contamination, or rare singletons and doubletons.

240    About 24% of all ESVs were taxonomically assigned to Arthropoda taxa and the final

241    Arthropoda ESV counts are shown in Table 2.  When 6 COI primer pairs are compared,

242    F230R ESVs contained the highest proportion of Arthropoda ESVs (43.9%) and

243    contained the highest proportion of raw reads (4.7%). About 13% of raw reads were

244    retained in this final set of Arthropoda ESVs.  Out of all the Arthropoda taxonomic

245    assignments, 11% of unique species, 15% of genera, and 26% of families were

246    considered high confidence assignments (Table S4).  The proportion of raw reads

11

247     represented in these high confidence Arthropoda assignments was 7% for species, 8%

248     for genera, and 10% for families.

249         Rarefaction curves show that at each rank, all samples reached a plateau,

250     indicating that we had sufficient sequencing coverage for these samples (Figure S1).

251     The proportion of taxa that are arthropods and the proportion of arthropods that are

252     EPTC are shown in Figure S2.  The average Arthropoda richness was not significantly

253     different across the pairwise amplicon comparisons (pairwise t test, p > 0.05) and there

254     was substantial variation in richness across sites (Figure S3).  The total number of

255     unique Arthropoda taxa were compared across COI amplicons (Figure S4) and the

256     amplicon that detects the most unique taxa varied depending on the taxonomic

257     resolution of the results.  At the ESV rank, the ml-jg amplicon recovered the highest

258     richness.  We also note that the presence-absence of ESVs are positively correlated

259     whether tissue or soil DNA kits are used for extraction (Figure S5) and across 2 PCR

260     replicates (Figure S6).

261         To test the effect of using multiple COI amplicons on richness, we pooled

262     increasing numbers of combined amplicons.  We show that using a multi-amplicon

263     approach can detect greater richness than using any single amplicon alone (Figure 2).

264     In this study, ESV richness increases linearly as amplicons are added whereas species

265     richness reaches a plateau when at least 4 amplicons are combined.  In some cases,

266     multiple combinations of amplicons recover equivalent richness.  Due to limitations in

267     the underlying reference sequence database, it is likely that species richness will also

268     increase as additional reference taxa are added so that more ESVs can be assigned

269     with high-confidence (Porter & Hajibabaei, 2018b).

12

270    We also looked at how the recovery of broad-based indicator taxa from across

271    the Arthropoda and more traditional freshwater indicator taxa from the EPTC varied with

272    amplicon choice (Figure 3).  Overall, more site indicator taxa were recovered from

273    across the Arthropoda compared with limiting analyses to the EPTC.  Generally, the

274    amplicon that recovers the greatest number of site indicators varies according to the

275    taxonomic resolution of the analysis.  At the ESV rank, BF1R2 recovers the greatest

276    number of broadscale site indicator taxa and F230R specifically recovers the greatest

277    number of EPTC indicator taxa.  Since the subset of indicator taxa presented for the

278    species to family ranks only represents the portion of the ESVs assigned with high

279    confidence, rank specific results may change over time as reference databases better

280    represent local taxa (Porter & Hajibabaei, 2018b).  We further explored the identities of

281    non-traditional freshwater indicator taxa by looking at the taxonomic distribution of the

282    broadscale indicator species and how this varied for each amplicon (Figure 4).  Indicator

283    taxa from the Elmidae (riffle beetles), Limoniidae (crane flies), and Simuliidae (black

284    flies) were detected in addition to the expected indicator species from the

285    Ephemeroptera and Trichoptera.

286    To investigate the effect of amplicon choice on beta diversity we looked at how

287    sites cluster with respect to COI amplicon choice and PCR replicates.  We compared all

288    the data at the ESV rank (Figure 5).  Samples cluster by site (stress = 0.154, linear $R^2$ =

289    0.912).  We found significant heterogeneity of beta diversity among sites (p-value <

290    0.05), but since we had a balanced design, proceeded to use PERMANOVA to test the

291    significance of groupings (Anderson & Walsh, 2013).  There were no significant

292    interactions among sites, makers, or PCR replicates.  Amplicon choice and PCR

13

293     replicate did not explain any significant variation in beta diversity among samples, but

294     sites explained 76% of the variation among samples ($R^2 = 0.76$, p-value = 0.001).

295

296     **Discussion**

297

298          As showcased in recent literature, DNA metabarcoding has gained significant

299     popularity in various ecological studies where biodiversity in a habitat or a sample is

300     investigated (Bik et al., 2012; Yu et al., 2012; J. Gibson et al., 2014, 2015; Creer et al.,

301     2016; Leese et al., 2016; Bush et al., 2017; Porter & Hajibabaei, 2018c; Bush et al.,

302     Submitted).  In this study we show that the optimal choice of amplicon(s) ought to be

303     based on the objective of the study: optimizing richness, optimizing the differentiation of

304     samples based on sites/conditions, or optimizing the detection of target taxa.  Here we

305     show the impact of using varied primer sets all of which have been used in recent

306     metabarcoding studies of freshwater benthic macroinvertebrates.

307          As predicted, different primer sets produced varied richness results.  For

308     example, while the ml-jg amplicon produced the highest overall ESV richness,

309     combinations of amplicons together detected even greater richness.  Moreover, even

310     though ml-jg maximizes ESV richness, at the species rank the best choice is F230R, yet

311     at the genus rank BR5 optimizes richness.  The decision to present the results of a

312     study at various taxonomic ranks is often based on the desire to include all the data

313     (ESV rank), or to present results at a fine level of taxonomic resolution (species rank),

314     or to present results based on previous knowledge.  For example, 94% of North

315     American freshwater specimens identified by morphology are represented by a DNA

14

316   sequence so it may be desirable to present results at the genus rank (Curry, Gibson,

317   Shokralla, Hajibabaei, & Baird, 2018).  These observations have important implications

318   for choosing primers, especially when considering the level of standardization required

319   in biomonitoring programs.  While the use of a single primer set is desirable, richness

320   based on metabarcoding is sensitive to the number of combined amplicons, primer

321   choice, and the taxonomic resolution of the results.  Based on results from this study

322   and elsewhere, primer binding biases during amplification steps can have tangible

323   impacts on results and using multiple primer sets will aid in increasing taxonomic

324   coverage (J. Clarke et al., 2009; Bellemain et al., 2010; Hajibabaei, Spall, Shokralla, &

325   van Konynenburg, 2012; J. Gibson et al., 2014).  For the sake of flexibility and forward

326   compatibility, aside from the storage of raw data, we encourage data reporting at the

327   ESV rank.  When reports are summarized to other taxonomic ranks, we encourage

328   disclaimer statements that results are limited by the taxonomic coverage of current

329   reference databases that may improve in the future (Porter & Hajibabaei, 2018b).

330         Our study provides important insights with regards to use of varied PCR primer

331   sets and replicates.  Contrary to measures of alpha diversity (above), beta diversity

332   measures do not seem to be affected by primer sets or PCR replicates when ESVs are

333   used for the spatial analysis.  In other words, spatial separation of sites based on these

334   varied parameters are robust as used in biomonitoring applications.  While the use of

335   standard primer sets may be desirable, it may not be required as beta diversity is less

336   sensitive to primer choice and technical replicates.

337         Another important and widespread use of metabarcoding data is in determining

338   ecosystem status or "biomonitoring" where the state of the ecosystem is derived from

339    bioindicator assemblages such as EPTC (Buss et al., 2002).  The recovery of

340    bioindicators from metabarcoding data in this study varied with amplicon choice.  For

341    example, the BF1R2 amplicon detects the greatest number of broadscale Arthropoda

342    site indicators but the F230R amplicon detects the greatest number of traditional EPTC

343    site indicators.  This has implications for benchmarking studies that compare

344    metabarcoding against the results based on the use of traditional bioindicators.  With

345    the application of DNA-based methods, our ability to detect a broad range of taxa has

346    improved such that it may not be necessary to limit bioindicator reporting to just the

347    traditional bioindicators (Bush et al., Submitted).

348        Note that even though equimolar amounts of each amplicon were combined for

349    sequencing, variable numbers of reads were obtained across amplicons.  This may be

350    caused by variable amplification efficiency during library preparation or slight differences

351    in the number of transferred amplicons when they are pooled prior to library preparation.

352    Since the recovery of variable library sizes is such a common occurrence, it is important

353    to normalize library size across samples prior to conducting data analysis.  It has been

354    shown that there is a trade-off between the use of rarefaction (removal of sequences

355    such that each sample can be compared at a common library size) to reduce the false

356    positive rate, and a loss of sensitivity because of the removal of sequences (McMurdie

357    & Holmes, 2014; S. J. Weiss et al., n.d.).  This has implications for beta diversity

358    analyses, where false positives can occur when samples cluster by sequencing depth

359    obscuring real differences, especially for samples with very small library sizes.

360    Common normalization methods include rarefaction down to the smallest library size

361    and working with proportions (ESV reads per sample / total reads per sample).  A

362    simulation study showed that rarefaction combined with the analysis of presence-

363    absence data worked best to cluster samples when groups are substantially different (S.

364    Weiss et al., 2017).  For differential abundance testing, however, methods that take into

365    consideration the compositional nature of metabarcode datasets (log ratio

366    transformation) may be more appropriate (Gloor, Macklaim, Pawlowsky-Glahn, &

367    Egozcue, 2017; S. Weiss et al., 2017; S. J. Weiss et al., n.d.).

368

369    **Conclusions**

370        This study analyzed how arthropod richness, beta diversity, and recovery of site

371    indicator taxa vary with COI amplicon choice.  We show how richness is sensitive to

372    primer choice and the combined use of multiple COI amplicons; beta diversity is robust

373    to primer choice and PCR replicates when ESVs are used for this analysis; and that

374    limiting analyses to the traditional bioindicator assemblages may not capture the

375    diversity of broadscale arthropod site indicators that can be recovered from COI

376    metabarcoding.  We note that the proportion of raw reads retained in the final set of

377    Arthropod ESVs was relatively low (13%) and this proportion varied noticeably among

378    different primer sets.  There are several reasons why raw reads are filtered out of the

379    final dataset: the removal of chimeras during PCR, removal of sequences with errors

380    generated during PCR or sequencing, and/or to the removal of the substantial 'tail' of

381    rare taxa as commonly seen in such studies reflecting either genuine rare novelty or

382    artefacts (Kunin, Engelbrektson, Ochman, & Hugenholtz, 2010; Brown et al., 2015).

383    Regardless, we suggest that method efficiency may also be gauged using a similar

384    measure of raw data that can be retained in final analyses.

385

## Data Accessibility

Raw reads were submitted to the NCBI SRA: #. FASTA files of the final ESVs are available as supplementary material. The SCVUC semi-automated bioinformatic pipeline is available from GitHub (#). The custom scripts used to generate figures are also available from GitHub (#).

391

## Acknowledgements

397

## References

Anderson, M. J., & Walsh, D. C. I. (2013). PERMANOVA, ANOSIM, and the Mantel test in the face of heterogeneous dispersions: What null hypothesis are you testing? *Ecological Monographs*, *83*(4), 557–574. doi:10.1890/12-2010.1

Baird, D. J., & Hajibabaei, M. (2012). Biomonitoring 2.0: a new paradigm in ecosystem assessment made possible by next-generation DNA sequencing. *Molecular Ecology*, *21*(8), 2039–2044.

Bellemain, E., Carlsen, T., Brochmann, C., Coissac, E., Taberlet, P., & Kauserud, H. avard. (2010). ITS as an environmental DNA barcode for fungi: an in silico approach reveals potential PCR biases. *BMC Microbiology*, *10*(1), 189.

Bik, H. M., Porazinska, D. L., Creer, S., Caporaso, J. G., Knight, R., & Thomas, W. K. (2012). Sequencing our way towards understanding global eukaryotic biodiversity. *Trends in Ecology & Evolution*, *27*(4), 233–243. doi:10.1016/j.tree.2011.11.010

Bonada, N., Prat, N., Resh, V. H., & Statzner, B. (2006). DEVELOPMENTS IN AQUATIC INSECT BIOMONITORING: A Comparative Analysis of Recent Approaches. *Annual Review of Entomology*, *51*(1), 495–523. doi:10.1146/annurev.ento.51.110104.151124

415    Brown, S. P., Veach, A. M., Rigdon-Huss, A. R., Grond, K., Lickteig, S. K., Lothamer, K., …
416        Jumpponen, A. (2015). Scraping the bottom of the barrel: are rare high throughput
417        sequences artifacts? *Fungal Ecology*, *13*, 221–225. doi:10.1016/j.funeco.2014.08.006

418    Bush, A., Compson, Z., Monk, W., Porter, T. M., Steeves, R., Emilson, E., … Baird, D. J.
419        (Submitted). Studying ecosystems with DNA metabarcoding: lessons from aquatic
420        biomonitoring. *Methods in Ecology and Evolution*.

421    Bush, A., Sollmann, R., Wilting, A., Bohmann, K., Cole, B., Balzter, H., … Yu, D. W. (2017).
422        Connecting Earth observation to high-throughput biodiversity data. *Nature Ecology &*
423        *Evolution*, *1*(7), 0176. doi:10.1038/s41559-017-0176

424    Buss, D. F., Baptista, D. F., Silveira, M. P., Nessimian, J. L., & Dorvillé, L. F. (2002). Influence of
425        water chemistry and environmental degradation on macroinvertebrate assemblages in a
426        river basin in south-east Brazil. *Hydrobiologia*, *481*(1), 125–136.

427    Callahan, B. J., McMurdie, P. J., & Holmes, S. P. (2017). Exact sequence variants should replace
428        operational taxonomic units in marker-gene data analysis. *The ISME Journal*, *11*, 2639–
429        2643.

430    Clarke, J., Wu, H.-C., Jayasinghe, L., Patel, A., Reid, S., & Bayley, H. (2009). Continuous base
431        identification for single-molecule nanopore DNA sequencing. *Nature Nanotechnology*,
432        *4*(4), 265–270. doi:10.1038/nnano.2009.12

433    Clarke, L. J., Soubrier, J., Weyrich, L. S., & Cooper, A. (2014). Environmental metabarcodes for
434        insects: *in silico* PCR reveals potential for taxonomic bias. *Molecular Ecology Resources*,
435        *14*(6), 1160–1170. doi:10.1111/1755-0998.12265

436    Creer, S., Deiner, K., Frey, S., Porazinska, D., Taberlet, P., Thomas, W. K., … Bik, H. M. (2016).
437        The ecologist's field guide to sequence-based identification of biodiversity. *Methods in*
438        *Ecology and Evolution*, *7*(9), 1008–1018. doi:10.1111/2041-210X.12574

439    Curry, C. J., Gibson, J. F., Shokralla, S., Hajibabaei, M., & Baird, D. J. (2018). Identifying North
440        American freshwater invertebrates using DNA barcodes: are existing COI sequence
441        libraries fit for purpose? *Freshwater Science*, *37*(1), 178–189. doi:10.1086/696613

442    De Cáceres, M., & Legendre, P. (2009). Associations between species and groups of sites:
443        indices and statistical inference. *Ecology*, *90*(12), 3566–3574. doi:10.1890/08-1823.1

444    Edgar, R. C. (2016). UNOISE2: improved error-correction for Illumina 16S and ITS amplicon
445        sequencing. *BioRxiv*. doi:10.1101/081257

446    Elbrecht, V., & Leese, F. (2015). Can DNA-Based Ecosystem Assessments Quantify Species
447        Abundance? Testing Primer Bias and Biomass—Sequence Relationships with an
448        Innovative Metabarcoding Protocol. *PLOS ONE*, *10*(7), e0130324.
449        doi:10.1371/journal.pone.0130324

450    Elbrecht, V., & Leese, F. (2017). Validation and Development of COI Metabarcoding Primers for
451        Freshwater Macroinvertebrate Bioassessment. *Frontiers in Environmental Science*, *5*, 11.
452        doi:10.3389/fenvs.2017.00011

453    Emilson, C. E., Thompson, D. G., Venier, L. A., Porter, T. M., Swystun, T., Chartrand, D., …
454        Hajibabaei, M. (2017). DNA metabarcoding and morphological macroinvertebrate
455        metrics reveal the same changes in boreal watersheds across an environmental
456        gradient. *Scientific Reports*, *7*(1). doi:10.1038/s41598-017-13157-x

457    Folmer, O., Black, M., Hoeh, W., Lutz, R., & Vrijenhoek, R. (1994). DNA primers for amplification
458        of mitochondrial cytochrome c oxidase subunit I from diverse metazoan invertebrates.
459        *Molecular Marine Biology and Biotechnology*, *3*(5), 294–299.

460    Geller, J., Meyer, C., Parker, M., & Hawk, H. (2013). Redesign of PCR primers for mitochondrial
461        cytochrome c oxidase subunit I for marine invertebrates and application in all-taxa biotic
462        surveys. *Molecular Ecology Resources*, *13*(5), 851–861. doi:10.1111/1755-0998.12138

463    Gibson, J., Shokralla, S., Curry, C., Baird, D. J., Monk, W. A., King, I., & Hajibabaei, M. (2015).
464        Large-Scale Biomonitoring of Remote and Threatened Ecosystems via High-Throughput
465        Sequencing. *PLOS ONE*, *10*(10), e0138432. doi:10.1371/journal.pone.0138432

466    Gibson, J., Shokralla, S., Porter, T. M., King, I., Konynenburg, S. van, Janzen, D. H., … Hajibabaei,
467        M. (2014). Simultaneous assessment of the macrobiome and microbiome in a bulk
468        sample of tropical arthropods through DNA metasystematics. *Proceedings of the*
469        *National Academy of Sciences*, *111*(22), 8007–8012. doi:10.1073/pnas.1406468111

470    Gloor, G. B., Macklaim, J. M., Pawlowsky-Glahn, V., & Egozcue, J. J. (2017). Microbiome
471        Datasets Are Compositional: And This Is Not Optional. *Frontiers in Microbiology*, *8*.
472        doi:10.3389/fmicb.2017.02224

473    Hajibabaei, M., Shokralla, S., Zhou, X., Singer, G. A. C., & Baird, D. J. (2011). Environmental
474        Barcoding: A Next-Generation Sequencing Approach for Biomonitoring Applications
475        Using River Benthos. *PLOS ONE*, *6*(4), e17497. doi:10.1371/journal.pone.0017497

476    Hajibabaei, M., Spall, J. L., Shokralla, S., & van Konynenburg, S. (2012). Assessing biodiversity of
477        a freshwater benthic macroinvertebrate community through non-destructive
478        environmental barcoding of DNA from preservative ethanol. *BMC Ecology*, *12*, 28.
479        doi:10.1186/1472-6785-12-28

480    Jones, C., Somers, K. M., Craig, B., & Reynoldson, T. B. (2007). *Ontario Benthos Biomonitoring*
481        *Network: Protocol Manual*. Toronto: Queens Printer for Ontario.

482    Kunin, V., Engelbrektson, A., Ochman, H., & Hugenholtz, P. (2010). Wrinkles in the rare
483        biosphere: pyrosequencing errors can lead to artificial inflation of diversity estimates.
484        *Environmental Microbiology*, *12*(1), 118–123. doi:10.1111/j.1462-2920.2009.02051.x

485    Leese, F., Altermatt, F., Bouchez, A., Ekrem, T., Hering, D., Meissner, K., … Zimmermann, J.
486        (2016). DNAqua-Net: Developing new genetic tools for bioassessment and monitoring of

487     aquatic ecosystems in Europe. *Research Ideas and Outcomes*, *2*, e11321.
488     doi:10.3897/rio.2.e11321

489 Leese, F., Bouchez, A., Abarenkov, K., Altermatt, F., Borja, Á., Bruce, K., … Weigand, A. M.
490     (2018). Why We Need Sustainable Networks Bridging Countries, Disciplines, Cultures
491     and Generations for Aquatic Biomonitoring 2.0: A Perspective Derived From the
492     DNAqua-Net COST Action. In *Advances in Ecological Research* (Vol. 58, pp. 63–99).
493     Elsevier. doi:10.1016/bs.aecr.2018.01.001

494 Leray, M., Yang, J. Y., Meyer, C. P., Mills, S. C., Agudelo, N., Ranwez, V., … Machida, R. J. (2013).
495     A new versatile primer set targeting a short fragment of the mitochondrial COI region
496     for metabarcoding metazoan diversity: application for characterizing coral reef fish gut
497     contents. *Frontiers in Zoology*, *10*(1), 34. doi:10.1186/1742-9994-10-34

498 Maddison, W. P., & Maddison, D. R. (2015). *Mesquite* (Vol. Version 3.10). Retrieved from
499     http://mesquiteproject.org

500 Martin, M. (2011). Cutadapt removes adapter sequences from high-throughput sequencing
501     reads. *EMBnet. Journal*, *17*(1), pp–10.

502 McMurdie, P. J., & Holmes, S. (2014). Waste Not, Want Not: Why Rarefying Microbiome Data Is
503     Inadmissible. *PLOS Comput Biol*, *10*(4), e1003531. doi:10.1371/journal.pcbi.1003531

504 Oksanen, J., Blanchet, G. F., Friendly, M., Kindt, R., Legendre, P., McGlinn, D., … Wagner, H.
505     (2018). vegan: Community Ecology Package. R package version 2.5-2. Retrieved from
506     https://CRAN.R-project.org/package=vegan

507 Polz, M. F., & Cavanaugh, C. M. (1998). Bias in template-to-product ratios in multitemplate PCR.
508     *Applied and Environmental Microbiology*, *64*(10), 3724–3730.

509 Porter, T. M., & Hajibabaei, M. (2018a). Automated high throughput animal CO1 metabarcode
510     classification. *Scientific Reports*, *8*, 4226.

511 Porter, T. M., & Hajibabaei, M. (2018b). Over 2.5 million COI sequences in GenBank and
512     growing. *PLoS ONE*, *13*(9), e0200177. doi:10.1101/353904

513 Porter, T. M., & Hajibabaei, M. (2018c). Scaling up: A guide to high-throughput genomic
514     approaches for biodiversity analysis. *Molecular Ecology*, *27*(2), 313–338.
515     doi:10.1111/mec.14478

516 R Core Team. (2017). R: A Language and Environment for Statistical Computing. Retrieved from
517     https://www.R-project.org/

518 Rognes, T., Flouri, T., Nichols, B., Quince, C., & Mahé, F. (2016). VSEARCH: a versatile open
519     source tool for metagenomics. *PeerJ*, *4*, e2584. doi:10.7717/peerj.2584

520 RStudio Team. (2016). RStudio: Integrated Development for R. Retrieved from
521     http://www.rstudio.com/

522     Shapiro, S. S., & Wilk, M. B. (1965). An Analysis of Variance Test for Normality (Complete
523             Samples). *Biometrika*, *52*, 591–611.

524     St. John, J. (2016, Downloaded). SeqPrep. Retrieved from
525             https://github.com/jstjohn/SeqPrep/releases

526     Suzuki, M. T., & Giovannoni, S. J. (1996). Bias caused by template annealing in the amplification
527             of mixtures of 16S rRNA genes by PCR. *Applied and Environmental Microbiology*, *62*(2),
528             625–630.

529     Taberlet, P., Coissac, E., Pompanon, F., Brochmann, C., & Willerslev, E. (2012). Towards next-
530             generation biodiversity assessment using DNA metabarcoding. *Molecular Ecology*, *21*(8),
531             2045–2050.

532     Tange, O. (2011). GNU Parallel - The Command-Line Power Tool. ;;*Login: The USENIX Magazine*,
533             *February*, 42–47.

534     Tedersoo, L., Nilsson, R. H., Abarenkov, K., Jairus, T., Sadam, A., Saar, I., … Kõljalg, U. (2010). 454
535             Pyrosequencing and Sanger sequencing of tropical mycorrhizal fungi provide similar
536             results but reveal substantial methodological biases. *New Phytologist*, *188*(1), 291–301.
537             doi:10.1111/j.1469-8137.2010.03373.x

538     Vamos, E., Elbrecht, V., & Leese, F. (2017). Short COI markers for freshwater macroinvertebrate
539             metabarcoding. *Metabarcoding and Metagenomics*, *1*, e14625.
540             doi:10.3897/mbmg.1.14625

541     Wang, Q., Garrity, G. M., Tiedje, J. M., & Cole, J. R. (2007). Naive Bayesian Classifier for Rapid
542             Assignment of rRNA Sequences into the New Bacterial Taxonomy. *Applied and
543             Environmental Microbiology*, *73*(16), 5261–5267. doi:10.1128/AEM.00062-07

544     Wei, T., & Simko, V. (2017). R package "corrplot": Visualization of a Correlation Matrix (Version
545             0.84). Retrieved from https://github.com/taiyun/corrplot

546     Weiss, S. J., Xu, Z., Amir, A., Peddada, S., Bittinger, K., Gonzalez, A., … Knight, R. (n.d.). Effects of
547             library size variance, sparsity, and compositionality on the analysis of microbiome data.
548             doi:10.7287/peerj.preprints.1157v1

549     Weiss, S., Xu, Z. Z., Peddada, S., Amir, A., Bittinger, K., Gonzalez, A., … Knight, R. (2017).
550             Normalization and microbial differential abundance strategies depend upon data
551             characteristics. *Microbiome*, *5*, 27. doi:10.1186/s40168-017-0237-y

552     Wickham, H. (2009). *ggplot2: Elegant Graphics for Data Analysis*. New York: Springer-Verlag.
553             Retrieved from http://ggplot2.org

554     Yu, D. W., Ji, Y., Emerson, B. C., Wang, X., Ye, C., Yang, C., & Ding, Z. (2012). Biodiversity soup:
555             metabarcoding of arthropods for rapid biodiversity assessment and biomonitoring:
556             Biodiversity soup. *Methods in Ecology and Evolution*, *3*(4), 613–623. doi:10.1111/j.2041-
557             210X.2012.00198.x

558

559    **Table 1.  COI amplicons used in this study.**

560

561

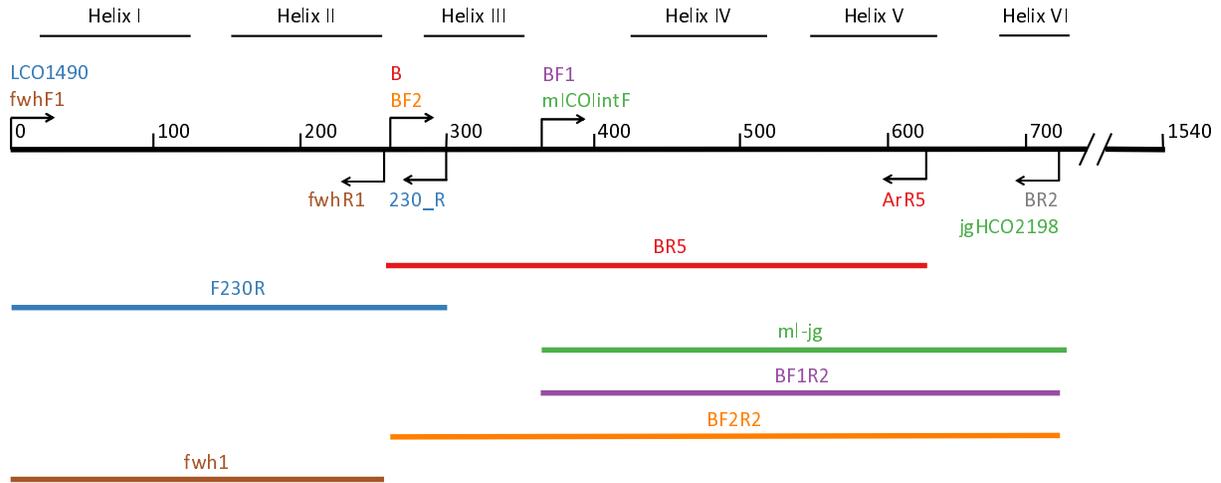| COI Amplicon | Primer | 5'-3' Primer sequence | Mode amplicon length (bp) | Primer reference | PCR conditions |
|---|---|---|---|---|---|
| BR5 | B | CCIGAYATRGCITTYCCICG | 310 | (Hajibabaei et al., 2012) | 95°C for 5min, 35 cycles of 94°C for 40s, 46°C for 1min, and 72°C for 30s, and a final extension at 72°C for 5min |
|  | ArR5* | GTRATIGCICCIGCIARIACIGG |  | (J. Gibson et al., 2014) |  |
| F230R | LCO1490 | GGTCAACAAATCATAAAGATATTGG | 229 | (Folmer, Black, Hoeh, Lutz, & Vrijenhoek, 1994) | 95°C for 5min, 35 cycles of 94°C for 40s, 46°C for 1min, and 72°C for 30s, and a final extension at 72°C for 5min |
|  | 230_R | CTTATRTTRTTTATICGIGGRAAIGC |  | (J. Gibson et al., 2015) |  |
| ml-jg | mlCOIintF | GGWACWGGWTGAACWGTWTAYCCYCC | 313 | (Leray et al., 2013) | 95°C for 1 min, 35 cycles of 94°C for 15 s, 46°C for 15 s, 72°C for 10s, and final extension at 72°C for 3 min |
|  | jgHCO2198 | TAIACYTCIGGRTGICCRAARAAYCA |  | (Geller, Meyer, Parker, & Hawk, 2013) |  |
| BF1R2 | BF1 | ACWGGWTGRACWGTNTAYCC | 316 | (Elbrecht & Leese, 2017) | 94 °C for 3 min; 40 cycles of 94 °C for 30 s, 50 °C for 30 s, and 65 °C for 2 min; and final extension at 65 °C for 5 min |
|  | BR2 | TCDGGRTGNCCRAARAAYCA |  |  |  |
| BF2R2 | BF2 | GCHCCHGAYATRGCHTTYCC | 421 | (Elbrecht & Leese, 2017) | 94 °C for 3 min; 40 cycles of 94 °C for 30 s, 50 °C for 30 s, and 65 °C for 2 min; and final extension at 65 °C for 5 min |
|  | BR2 | TCDGGRTGNCCRAARAAYCA |  |  |  |
| fwh1 | fwhF1 | YTCHACWAAYCAYAARGAYATYGG | 178 | (Vamos, Elbrecht, & Leese, 2017) | 95°C for 5 min, 34 cycles of 95°C for 30 s, 52°C for 30 s, 72°C for 2 min, and final extension at 72°C for 10 min |
|  | fwhR1 | ARTCARTTWCCRAAHCCHCC |  |  |  |

562

563

564 **Table 2: Arthropoda ESV and read counts vary by COI amplicon.**

565

|  | BR5 | F230R | ml-jg | BF1R2 | BF2R2 | fwh1 | Total |
|---|---|---|---|---|---|---|---|
| Arthropoda ESVs | 873 | 1,143 | 1,342 | 803 | 477 | 302 | 4,940 |
| Proportion of all ESVs assigned to Arthropoda (%) | 25 | 43.9 | 40.6 | 13.1 | 13.7 | 15.5 | 23.5 |
| Reads in Arthropoda ESVs | 187,353 | 467,910 | 285,933 | 147,697 | 24,375 | 167,129 | 1,280,397 |
| Proportion of raw reads in Arthropoda ESVs (%) | 1.9 | 4.7 | 2.9 | 1.5 | 0.2 | 1.7 | 12.8 |

566

567  **Figure 1.  Map of primers and amplicons tested in this study.**  The reference
568  sequence shown in black is *Drosophila yakuba*, cytochrome c oxidase region 1470-
569  3009 bp (1540 nt).  Secondary structure is shown for reference, comprised of 6 alpha
570  helices in the standard DNA barcode region shown here.
571
572



Reference sequence is *Drosophila yakuba* X03240, cytochrome c oxidase region 1470 – 3009 bp, 1540 nt
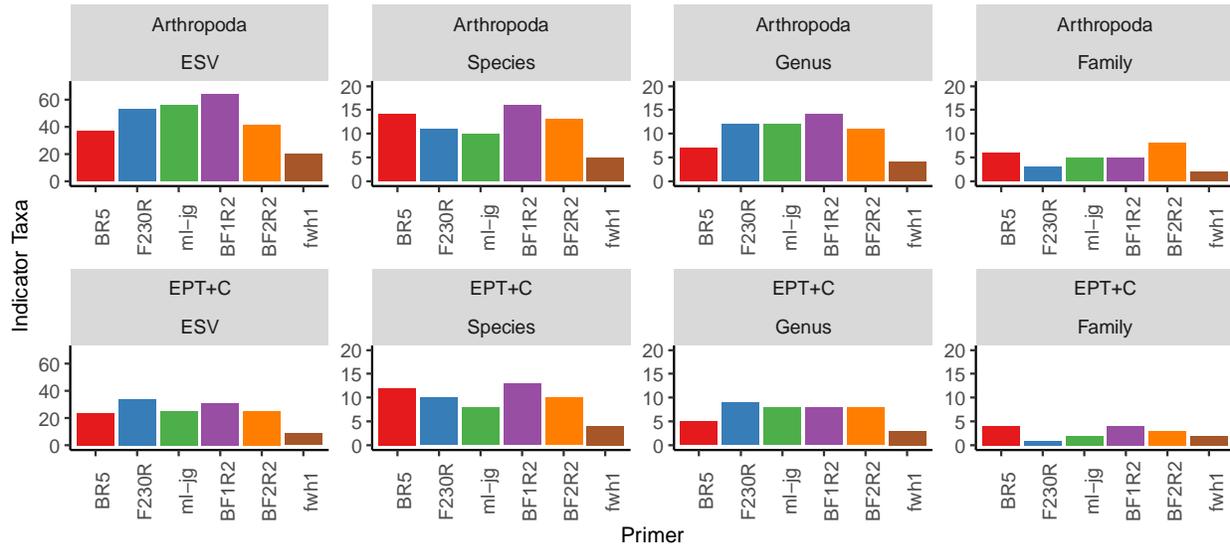
573
574

575 **Figure 2. ESV richness continues to increase as COI amplicons are added but**
576 **species - order richness reaches a plateau.** For the primer comparison experiment
577 that used the soil DNA extraction kit, we pooled the results from the 6 sites and show
578 the top COI amplicon combinations that detected the greatest richness. We report the
579 recovered richness when up to 6 amplicons are combined at the 1) ESV, 2) species, 3)
580 genus, 4) family, and 5) order ranks. ESV = exact sequence variant; A = BR5; B =
581 F230R; C = ml-jg; D = BF1R2; E = BF2R2; F = fwh1.
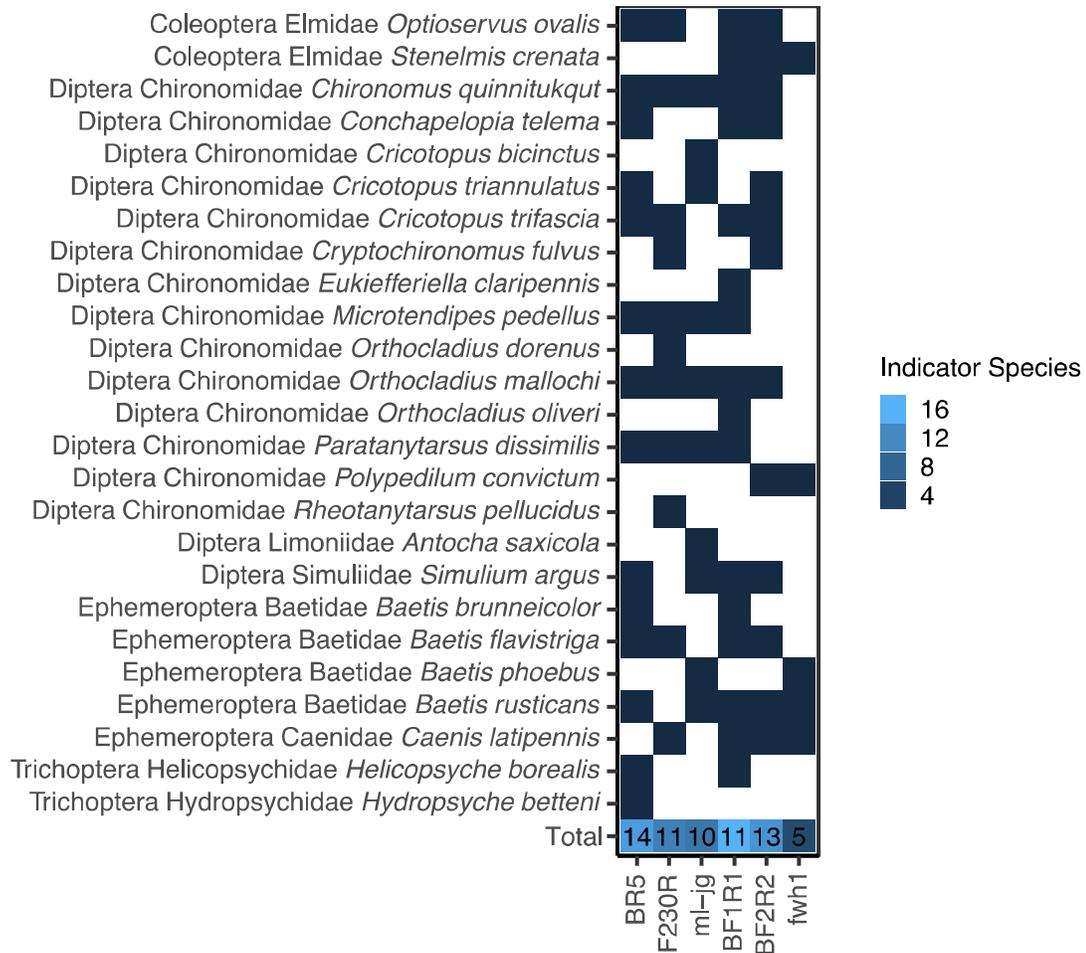582
583

**Figure 3. Ephemeroptera, Plecoptera, Trichoptera, and Chironomidae comprise a subset of total number of site indicator taxa drawn from across the Arthropoda.**

27

588 An indicator taxon analysis was used to determine the number of taxa that could
589 distinguish among 6 sampled sites. In the top panel, the number of broadscale indicator
590 taxa from across the Arthropoda are shown. In the bottom panel, the number of typical
591 freshwater indicator taxa from the EPT+C are shown. This analysis was based on
592 normalized data. ESV = exact sequence variant; EPTC = Ephemeroptera, Plecoptera,
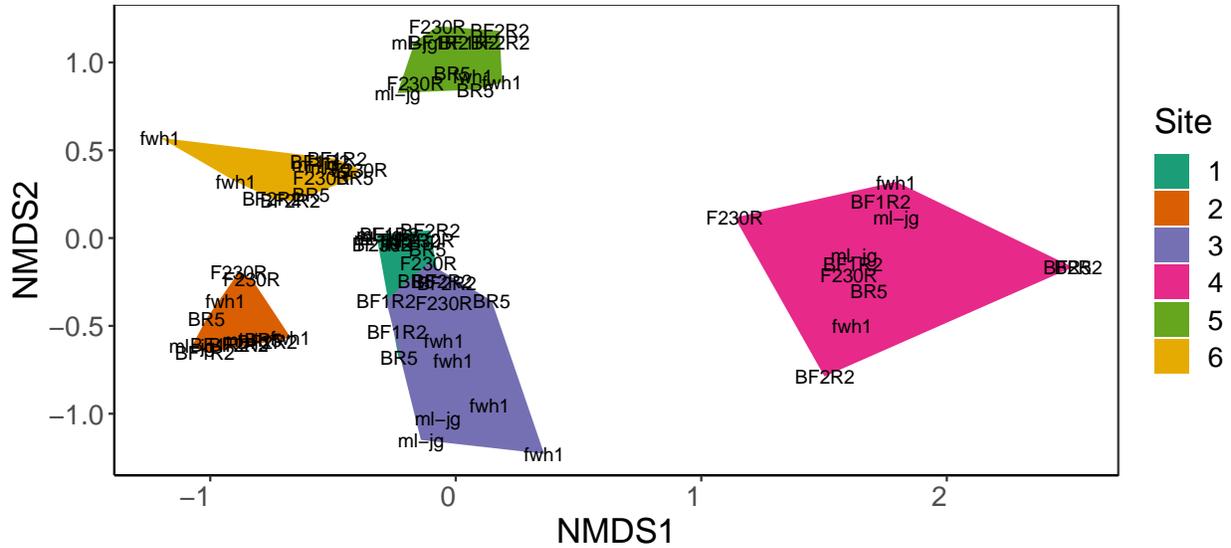593 Trichoptera, Chironomidae.

594



595
596
597

**Figure 4. Site indicator taxa chosen based on metabarcode sequencing are comprised of Coleoptera, Diptera, Ephemeroptera, and Trichoptera.** Presence is indicated by a dark square, absence by a white square. The total number of broadscale indicator taxa detected by each amplicon is shown in the bottom row according to the legend.

606 **Figure 5. Samples cluster mainly by site despite differences in amplicons and**
607 **replicates.** Results are based normalized data. COI amplicons are labelled directly in
608 the plot. Amplicons shown twice represent the two PCR replicates. Sites are grouped
609 by color according to the legend.
610



611
612
613
614