1    **Freshwater diatom biomonitoring through benthic kick-net metabarcoding**

2    **Short title: Freshwater diatom biomonitoring through metabarcoding**

3    Victoria **Carley** Maitland[1], Chloe Victoria Robinson[1], Teresita M. Porter[1,2], Mehrdad

4    Hajibabaei[1]*.

5

6    [1]Centre for Biodiversity Genomics & Department of Integrative Biology, University of

7    Guelph, Guelph, Ontario, N1G 2W1

8    [2]Great Lakes Forestry Centre, Natural Resources Canada, 1219 Queen Street East,

9    Sault Ste. Marie, ON Canada

10

11    * Corresponding author

12    Email: mhajibab@uoguelph.ca

13

14

15

16

17

18

19

20

21

# Abstract

Biomonitoring is an essential tool for assessing ecological conditions and informing management strategies. The application of DNA metabarcoding and high throughput sequencing has improved data quantity and resolution for biomonitoring of taxa such as macroinvertebrates, yet, there remains the need to optimise these methods for other taxonomic groups. Diatoms have a longstanding history in freshwater biomonitoring as bioindicators of water quality status. However, periphyton scraping, a common diatom sampling practice, is time-consuming and thus costly in terms of labour. This study examined whether the benthic kick-net technique used for macroinvertebrate biomonitoring could be applied to bulk-sample diatoms for metabarcoding. To test this approach, we collected samples using both conventional microhabitat periphyton scraping and bulk-tissue kick-net methodologies in parallel from replicated sites with different habitat status (good/fair). We found there was no significant difference in community assemblages between conventional periphyton scraping and kick-net methodologies, but there was significant difference between diatom communities depending on site quality ($P = 0.029$). These results show the diatom taxonomic coverage achieved through DNA metabarcoding of kick-net is suitable for ecological biomonitoring applications. The shift to a more robust sampling approach and capturing diatoms and macroinvertebrates in a single sampling event has the potential to significantly improve efficiency of biomonitoring programmes.

## Introduction

45

46        As climate change and other anthropogenic impacts continue to alter the

47    environment, there is an increasing need for comprehensive ecological assessment.

48    Rapid and robust biomonitoring is essential for informing management plans and

49    mitigating further environmental degradation [1–3]. Freshwater biomonitoring typically

50    involves sampling a range of aquatic taxa, with particular focus on biological indicator

51    taxa, to assess environmental conditions based on diversity, richness, structure and

52    function of the existing communities [3–5].

53

54    Traditionally, biomonitoring data is generated through morphological taxonomic

55    classifications, however there has been a recent shift towards DNA-based identification

56    using metabarcoding [6] coupled with high throughput sequencing [7]. In aquatic

57    systems such as wadable streams, a combination of bulk-tissue benthic sampling using

58    kick-net methodology with DNA metabarcoding, facilitates rapid data collection whilst

59    maintaining data integrity [8–10]. The metabarcoding approach has been employed for

60    numerous biomonitoring studies involving macroinvertebrates [11,12] for assessing

61    freshwater health [5,10,13].

62

63    In addition to benthic macroinvertebrates, diatoms (members of Bacillariophyta) are also

64    ideal biomonitoring target taxa for assessing freshwater system conditions [14–16].

65    These single-celled algae have a short generation time which allows for rapid

66    responses to physical, chemical and biological changes in the environment [14,15,17].

67    Similar to macroinvertebrates, the high diversity and ubiquity of diatoms is used to

3

68    create biotic indices that can accurately report freshwater quality [16,18,19]. Studies

69    have shown that diatoms respond more readily to the presence of heavy metal

70    pollutants compared to macroinvertebrates, which are generally more sensitive to shifts

71    in hydrological conditions [17,20–22]. Monitoring only one of these taxonomic groups to

72    assess overall ecosystem health could potentially cause gaps in knowledge that could

73    subvert subsequent management strategies. Hence, diatoms are being used in a

74    number of national and regional biomonitoring programmes.

75

76    Current methods for diatom sampling are time-consuming and laborious, which could

77    hamper widespread use of diatoms for extensive freshwater biomonitoring [23,24]. The

78    conventional diatom collection method involves the scraping of periphyton (a

79    combination of algae, cyanobacteria, microbes, and detritus) from numerous substrates

80    within littoral habitats [23–26]. These samples are then fixed and visualised using light

81    microscopy [27–30]. From here, microscopy standards and keys are followed [29–31] to

82    enable identification of diatoms to different taxonomic ranks. Within recent years, there

83    has been the shift towards DNA metabarcoding-based identification of diatoms

84    [15,16,32,33]. This involves the manual homogenized of periphyton scrapings into

85    single samples, which are then processed via standard diatom metabarcoding

86    procedures [34,35]. Alternative sampling methods, such as collection through the

87    benthic kick-net technique, have not been tested for diatom biomonitoring applicability,

88    however it is expected that this technique would drastically reduce time spent collecting

89    samples. The ability to study diatom and macroinvertebrate assemblages from a single

90    sample would allow biomonitoring programs to achieve an intensive appraisal of

91    freshwater conditions. In a rapidly changing world, streamlining current methodology to

92    obtain as much data in as little time as possible is crucial.

93

94    Because DNA-based analysis of environmental samples such as contents of a kick-net

95    sample can provide a broad spectrum of organisms in the habitat sampled, we

96    hypothesized that kick-net metabarcoding will provide diatom biodiversity comparable to

97    commonly used scraping method. Specifically, we aimed to 1) investigate the feasibility

98    of kick-net sampling for capturing community assemblages of freshwater diatoms

99    versus conventional periphyton scraping using a high throughput sequencing coupled

100   metabarcoding approach and 2) compare diatom community assemblages across a

101   known habitat quality scale (Good and Fair) using both conventional and kick-net

102   sampling to investigate presence of diatom indicator groups.

103

# Methods

## Field Sampling

106        Samples were collected in November 2019 from Grand River tributaries across

107   four study sites in Waterloo, Ontario (Fig. 1). Status and location data were provided by

108   Dougan & Associates based on a 2018 benthos biomonitoring project for the City of

109   Waterloo (S1 Table). The four selected sites were a subset of the sites from this project

110   and were chosen based on accessibility and habitat quality. Hilsenhoff Biotic Index

111   ranges (weighted by species) informed the habitat quality scale [36] which categorized

112   sites into 'Good' (4.51-5.50) and 'Fair' (5.51-6.50).

113

5

114    Collection occurred in riffles, starting with a benthic kick-net sample, followed by

115    subsequent periphyton scrapings of microhabitats representative of the reach (S2

116    Table). Periphyton scraping refers to the sampling of sediment, rock, macrophytes and

117    leaf litter. Three replicates of each sampling type were collected at each site. Kick-net

118    collection followed the Canadian Aquatic Biomonitoring Network [CABIN] protocol [37].

119    Effort was standardized to three minutes. The sampler moved up stream in a zig-zag

120    pattern to encompass all microhabitats within the reach. Periphyton scraping samples

121    were comprised of five specimens per microhabitat type to account for variability within

122    the microhabitat [23]. Negative controls, consisting of molecular grade water, were

123    collected prior to the collection of each rock sample (n= 9) to ensure the toothbrushes

124    used for scraping biofilms from rocks had been adequately sterilised (S3 Table). All

125    other samples were collected using manufacture-sealed sterile equipment. All samples

126    were collected in 1L sample jars and placed in a cooler to transport back to the lab.

127    Upon arrival at the lab, samples (n=45) were preserved using 100% ethanol and stored

128    in a -20°C freezer until processing.

129

## Sample Validation and Extraction

131        To account for potential false negatives [38], diatom presence in the samples

132    was confirmed using microscopy. A small amount of ethanol used to preserve the

133    samples was placed on a slide and observed under a compound microscope at 100X

134    magnification. Visual inspection confirmed the presence of diatoms in each sample type

135    (S1 Fig.), however no taxonomic information was taken as morphological identification

136    was beyond the scope of this study.

6

137

138    Once diatom presence was validated, samples were homogenized using standard

139    blenders decontaminated by washing with ELIMINase® (VWR, Canada) then rinsing

140    with deionized water before treating with UV light for 30 minutes. Homogenate was

141    subsequently transferred to 50 mL Falcon tubes, where one tube was set aside and

142    centrifuged at 2400 rpm for two minutes. Supernatant was removed and residual pellets

143    were incubated at 70 ºC until fully dried. Next, approximately 300 mg dried tissue was

144    subsampled into PowerBead tubes and DNA extractions were completed using the

145    DNeasy Power Soil kit (Qiagen, CA) following the manufacturer's protocol. The only

146    exception being that 50 µL of buffer C6 (TE) was used for final elution. Negative

147    controls containing no tissue were also included with each batch of extractions. All

148    negative controls failed to amplify and therefore were not sequenced.

149

## DNA Amplification, Library Preparation and Sequencing

151        Amplification targeted the 312 base pair long region of the chloroplast gene

152    ribulose bisphosphate carboxylase large chain (rbcL) using five diatom specific primers.

153    Following the methods of Rivera et al. [39], forward primers Diat_rbcL_708F_1 (5'-

154    AGGTGAAG- TAAAAGGTTCWTACTTAAA-3'), Diat_rbcL_708F_2 (5'-AGGT-

155    GAAGTTAAAGGTTCWTAYTTAAA-3') and Diat_rbcL_708F_3 (5'-AGGTGAAAC-

156    TAAAGGTTCWTACTTAAA-3') were combined in an equimolar mix. Two reverse

157    primers, Diat_rbcL_R3_1 (5'-CCTTCTAATTTACC- WACWACTG-3') and

158    Diat_rbcL_R3_2 (5'-CCTTCTAATTTACCWA-CAACAG-3'), were also combined and

159    used for amplification. Each reaction used the following reagents: 17.5 µL HyPure$^{TM}$

7

160    molecular biology grade water, 2.5 µL 10X reaction buffer (200 mM Tris-HCl, 500 mM

161    KCl, pH 8.4), 1 µL MgCl$_2$ (50 mM), 05. µL dNTPs mix (10 mM), 0.5 µL of both forward

162    (10 mM) and reserve (10 mM) equimolar mixes, 0.5 µL Invitrogen's Platinum Taq

163    polymerase (5 U) and 2 µL of DNA. Final reaction volume totaled 25 µL.

164

165    PCR protocol largely followed Rivera et al. [39] with minor adjustments. Instead of thirty

166    cycles of denaturation at 95°C for 45 seconds, annealing at 55°C for 45 seconds and

167    extension at 72°C for 45 seconds [39], this study increased the number of cycles to

168    thirty-five. PCR amplification was also performed in two-steps, with the second PCR

169    using 2 µL of amplicons from the first PCR instead of DNA, and Illumina-tailed primers.

170    All PCRs were completed in Eppendorf Mastercycler ep gradient S thermal cycler.

171    Successful amplification was confirmed using 1.5% agarose gel electrophoresis before

172    purifying second PCR amplicons with the MinElute Purification kit (Qiagen). The next

173    step was quantifying purified samples with a QuantIT PicoGreen daDNA assay kit and

174    using these values to normalize all samples to 3 ng/µL. Samples were then indexed and

175    pooled before purifying with AMpure magnetic beads. QuantIT PicoGreen daDNA assay

176    kit was once again used to quantify the library and Bioanalyzer was used to determine

177    fragment length. The library was diluted to 4 nM and 10% PhiX was added before being

178    sequenced using Illumina MiSeq with a V3 MiSeq sequencing kit (300 X 2; MS-102-

179    2003).

180

181    **Bioinformatic Processing**

8

182    Illumina MiSeq paired-end reads were processed using the SCVURL rbcL

183    metabarcode pipeline-1.0.2 pipeline available from

184    https://github.com/terrimporter/SCVURL_rbcL_metabarcode_pipeline

185    . SCVURL is an automated snakemake [40] bioinformatic pipeline that runs in a conda

186    [41] environment. SeqPrep v1.3.2 [42] was used to pair raw reads requiring a minimum

187    Phred score of 20 to ensure 99% base-calling accuracy. CUTADAPT v2.6 was used to

188    trim primers from sequences, leaving a minimum fragment length of at least 150 base

189    pairs [43]. Global exact sequence variant (ESV) [44] analysis was performed on the

190    primer-trimmed reads. Reads were dereplicated using the 'derep_fulllength' command

191    with the 'sizein' and 'sizeout' options of VSEARCH v2.14.1 [45]. VSEARCH was also

192    used to denoise the data using the unoise3 algorithm [46]. These steps were taken to

193    remove sequences with errors, chimeric sequences, PhiX carry-over and rare reads

194    (singletons or doubletons) [47]. ESVs were classified using the rbcL diatom Classifier

195    available from https://github.com/terrimporter/rbcLdiatomClassifier. Reference rbcL

196    sequences were downloaded from the INRA diatom project [48]and reformatted to train

197    the naive Bayesian classifier to make rapid, accurate taxonomic assignments [49].  This

198    method makes assignments to the species rank and produces a statistical measure of

199    confidence for each taxon up to the domain rank to help reduce false positive taxonomic

200    assignments.  We used 0.60 cutoff at the family rank (99% accuracy) and 0.20 cutoff at

201    the genus rank (95% accuracy). The accuracy of the method assumes that target taxa

202    are present in the reference database.

203    

204    **Statistical Analysis**

205     RStudio was used to analyze the data [50]. To account for variable reads within

206     the library each sample was normalized to the 15th percentile using the 'rrarefy' function

207     in the vegan package [51,52].

208

209     ESV richness across the various sampling and status categories was calculated to

210     assess differences between the methods and sites. A non-metric multi-dimensional

211     (NMDS) analysis on Sorensen dissimilarities (binary Bray-Curtis) was conducted using

212     the vegan 'metaMDS' function to determine if sampling method or site status created

213     variation in community structure [5]. A scree plot was run using the 'dimcheckMDS'

214     command from the goeveg package to determine the number of dimensions (k=2) to

215     use with vegan metaMDS function[53]. Shephard's curve and goodness of fit

216     calculations were calculated using the vegan 'stressplot' and 'goodness' functions. The

217     vegan 'vegdist' command was used to build a Sorensen dissimilarity matrix. We

218     checked for heterogeneous distribution of dissimilarities using the 'betadisper' function.

219     We used the 'adonis' function to perform a permutational analysis of variance

220     (PERMANOVA). PERMANOVA was performed on conventional sampling methods

221     (periphyton scraping) and kick-net methods, as well as site status to test for significant

222     interactions between the categories [54].

223

224     To maintain a balanced design during statistical testing, we pooled all periphyton

225     sampling into one sample type (conventional) and maintained kick-net samples as a

226     separate sample type. The Jaccard index was calculated to assess the overall

227     similarities between the sites, collection methods and site status. Nestedness and

228    turnover of between kick-net and conventional samples were calculated using R

229    package betapart function 'beta-pair' [55] followed by vegan function 'betadisper'. The

230    number of diatom family ESVs detected from kick-net or pooled conventional samples

231    was also plotted. A dendrogram of diatom families detected was plotted using

232    RAWGraphs (app.rawgraphs.io) and color-coded to show the samples the families were

233    detected in [56]. Lastly, the frequency of ESVs detected from diatom families was

234    visualized using a heatmap generated using geom_tile (ggplot) in R, plotting individual

235    sample types for each site, split into two plots according to site status.

236

## Results

238        After bioinformatic processing, we generated 4,272 ESVs (2,166,157 reads).

239    After taxonomic filtering (removal of non-diatom phyla), a total of 3,940 diatom ESVs

240    (2,125,984 reads) were retained for data analysis. Read coverage per sample after

241    normalisation (15$^{th}$ percentile cut-off) was 37,735.

242

243    Since the rarefaction curves plateau, this indicated that the sequencing depth was

244    sufficient to capture the ESV diversity in our PCRs (S2 Fig.). In terms of the top 10

245    orders identified, the order Naviculales represented 30.6% of ESVs (30% of reads) and

246    Bacillariales represented 18.6% of ESVs (15.4% of reads; S3 Fig.).

247

## Taxonomic Coverage

249    In terms of taxonomic assignment, we identified a total of 1 phyla (Bacillariophyta), 4

250    classes, 23 orders, 44 families and 77 genera at the 95% correct assignment level. ESV

11

251    richness varied across different sampling methods (Fig. 2). Mean overall ESV richness

252    was used to calculate alpha diversity which displayed very similar values for all

253    sampling methods across the four sites (S4 Table). Averaged across sites, kick-net

254    samples produced the lowest mean ESV richness (225 ± 85), with sediment samples

255    producing the highest ESV richness (317 ± 92).

256

257    Through investigating diatom families, a majority of families detected were present in all

258    microhabitats and kick-net samples (Fig. 3). Two families (Coscinodiscaceae and

259    Orthoseriaceae) were solely present in leaf litter samples and two families

260    (Entomoneidaceae and Diadesmidaceae) were present only in sediment samples (Fig.

261    3).

262

263    In terms of diatom genera, some of the confidently identified genera represented by

264    more than 2 sequence variants, identified from kick-net and conventional samples,

265    included: *Nitzschia* (Bacillariales), *Polypedilum* (Chironomidae), *Navicula* (Naviculales),

266    *Amphora* (Thalassiophysales) and *Ulnaria* (Licmophorales; Fig. 4).

267

## Diatom Diversity by Method and Site Status

269         NMDS plots showed that replicates clustered close together for site and status,

270    with overlap observed between sampling methods and replicates (Fig. 5). When pooling

271    conventional periphyton samples (i.e. macrophyte, leaf litter, rock, and sediment) at

272    each site, there remained overlap between kick-net and conventional samples and

273    samples also remained clustered by site and status (S4 Fig). PERMANOVA of the

274    pooled samples, shows that analyzing data from kick-net or conventional samples

275    (method) explains 13% of the variation in Bray Curtis dissimilarities (p-value = 0.776),

276    sampling site (site) explains 58% of the variation (p-value = 0.009) and habitat quality

277    status (status) explains 22% of the variation observed (p-value = 0.029; S5 Table).  The

278    Jaccard index for kick-net compared with conventional samples is 0.53, indicating

279    samples are 53% similar, whereas the Jaccard index for fair compared to good site

280    quality status samples is 0.20, indicating samples are only 20% similar. In terms of beta

281    diversities of communities aggregated by the treatments of "kick-net" and

282    "conventional", there was no significant difference between turnover. For beta

283    diversities of communities aggregated by site status, there was a significant difference

284    between nestedness ($P < 0.05$) but not for turnover ($P = 0.06$).  Fair samples appear to

285    be significantly nested within good samples. These results further indicate that site

286    status has a significant effect on the sampled community composition whereas

287    conventional versus kick-net sampling methods do not.

288

289    For individual sample types (i.e. kick-net, macrophyte, leaf litter, rock, and sediment),

290    the heatmap shows that kick-net samples are largely representative of the diversity of

291    families detected within each conventional periphyton sampling method (Fig.6.). In

292    some cases, kick-net samples failed to detect diatom families which were present in

293    conventional periphyton samples (e.g. Sellaphoraceae and Diadesmidaceae in Clair15)

294    and conversely, kick-net samples also detected families which were not detected in

295    conventional periphyton samples (e.g. Eunotiaceae and Neidiaceae in Clair12; Fig. 6).

296    Similar assemblages of diatoms communities were detected across both fair and good

13

297    quality sites, with the main difference observed between fair and good sites being the

298    number of reads produced for families such as Thalassiosiraceae which was detected

299    with a high number of reads (1000+) in fair sites and a lower number of reads (10-100)

300    in good sites (Fig. 6).

301

## Discussion

302

303        The demand for high-quality, reproducible ecological data is increasing in

304    conjunction with the degradation of ecosystems globally [57]. There is a need to further

305    streamline existing biomonitoring methodologies without sacrificing the quality of data

306    produced [4,7,54,58]. With diatom assemblages providing a unique insight into the

307    water quality status of lentic and lotic systems, fast-tracking diatom data collection for

308    ecological assessments is a priority [39]. We have demonstrated that kick-net

309    methodology with DNA metabarcoding provides sufficient taxonomic coverage to

310    potentially be utilised as a for assessing diatom biodiversity in freshwater systems.

311

312    Kick-net sampling technique, whereby a zig-zag path is taken across the reach,

313    provided sufficient representation of existing diatom community assemblages within

314    site-specific microhabitats. Samples derived from the kick-net technique were highly

315    comparable with conventional samples in terms of diatom taxa detected, despite the

316    kick-net approach being more passive compared to direct periphyton scraping. Specific

317    diatom taxa are known to have ecological preferences for different freshwater

318    microhabitats [59,60]. For watershed-level health estimates, it is beneficial to be able to

319    efficiently detect the diversity of diatom taxa present without directly sampling each

14

320 microhabitat within a reach. We have demonstrated that kick-net methodology can

321 sufficiently capture the existing diatom biodiversity, ground truthed by comparing

322 assemblages detected with periphyton scrapings.

323

324 Ultimately, the detection of bioindicator species is a key variable to consider when

325 comparing biomonitoring methods, as these taxa are pivotal for detecting subtle

326 differences in freshwater health [3,5,14]. Naviculaceae contains diatom species

327 sensitive to herbicide exposure, which is a family we observed in all sites and with all

328 collection methods [61]. Additionally, the bioindicator family Stephanodiscaceae, (a

329 known tolerant taxon) [62], has a higher read abundance in 'Fair' sites compared to

330 'Good' in both conventional and kick-net sample types. Despite the direct sampling

331 approach of periphyton rock scraping, this methodology failed to detect this family at

332 one of the sites where kick-net samples were successful at detecting this benthic family.

333 Rock scrapings are commonly used as the sole collection method for diatoms

334 [14,39,63,64], which suggests that the kick-net approach facilitates the detection of taxa

335 which otherwise may be missed from conventional sampling.

336

337

## Conclusion

339 Overall, this study found that benthic kick-net methodology enables a robust and

340 detailed assessment of freshwater diatom communities. This methodology is a scalable

341 option for generating a holistic insight into the health of freshwater systems. The high

342 similarity of diatom taxa detected between methods and significant differences between

15

343    diatom communities detected in sites of differing habitat quality, demonstrates that this

344    rapid method can provide accurate, fine-resolution taxonomic results. Future research

345    should examine the duo-analyses approach of macroinvertebrate and diatom

346    communities from a single kick-net sample, to determine reproducibility of multi-taxa

347    targeting with this method. Additionally, future studies should consider exploring the use

348    of multiple markers (i.e. rbcL cpDNA versus 18S rRNA gene), to address level of

349    taxonomic resolution that can be obtained with these markers commonly used for

350    diatom DNA barcoding.

351

352

# Supporting Information

354    **S1 Table. Information on study sites, including GPS coordinates and site status.**

355    **S2 Table. Outline of collections methods used in this study.** Samples for periphyton

356    scraping were taken from a depth no greater than 1m (King et al. 2006).

357    **S3 Table. Summary table of decontamination and sterilisation procedures**

358    **undertaken for the equipment in this study.**

359    **S4 Table. Mean ESV values (replicates pooled) for each sample type across the**

360    **four sites.** Based on normalised data.

361    **S5 Table. rbcL exact sequence variants (ESVs) are not significantly different**

362    **between sampling methods (kick-net versus conventional periphyton sampling).**

363    No significant beta dispersion was detected within groups (method, site, status). Only

364    significant difference detected was rbcL ESVs between sites and status. Summary of

365    PERMANOVA results based on a Sorensen dissimilarity matrix of rbcL ESVs.

366    Significant p-values are bolded.

367    **S1 Fig. Example of confirmation of diatom presence from preservative of kick-net**

368    **sample.** Image: CBG Photography Group.

369    **S2 Fig. All samples show that ESV sampling reached saturation.** Samples were

370    color-coded by site or method as shown in the legend.  The vertical dashed line

371    indicates the 15th percentile of sampling read depth, which is the number of reads that

372    would be used in any future analysis based on normalized data.

373    **S3 Fig. Naviculales is the most abundant diatom order detected.**  Results for the

374    top 10 orders are shown with respect to proportion of ESVs and reads recovered.

375    Based on raw unnormalized data.

376    **S4 Fig. Non-metric multi-dimensional scaling plots of microhabitat samples**

377    **pooled show clustering by due to site and status.** Specifically, a) depicts overlap

378    between the binary Bray Curtis (Sorensen) dissimilarities between different sampling

379    approaches, b) sample site clustering c) clustering based on habitat quality status

380    (stress = 0.012, $R^2$ = 0.98).  Based on rarefied data.

381

# Acknowledgements

383         We would like to extend thanks to Michael Wright for providing advice on

384    laboratory procedures, Josip Rudar for assistance with bioinformatics and Genevieve

385    Johnson for help collecting samples. We would also like to acknowledge Jessica

386    Robinson, Jaclyn McKeown, Allison Brown and Monica Young from the Centre for

387    Biodiversity Genomics' Collections department for assisting with imaging and use of

17

## Author Contributions and Competing Interests

394     V.C.M. collected samples, conducted molecular and genomic analyses and wrote

395    the manuscript with help from all authors, C.V.R and M.H. designed the study, C.V.R

396    contributed to sampling, bioinformatic processing and statistical analyses, T.M.P trained

397    the classifier, contributed to bioinformatic processing and advised on data analysis. All

398    authors helped to write/edit the manuscript.  The authors have declared that no

399    competing interests exist.

400

## References

402    1.    Bayramoglu S, Chakir R, Lungarska A. Impacts of Land Use and Climate Change

403          on Freshwater Ecosystems in France. Environ Monit Assess. 2020;25: 147–172.

404          doi:10.1007/bf02837416

405    2.    Karaouzas I, Smeti E, Vourka A, Vardakas L, Mentzafou A, Tornés E, et al.

406          Assessing the ecological effects of water stress and pollution in a temporary river

407          - Implications for water management. Sci Total Environ. 2018;618: 1591–1604.

408          doi:10.1016/j.scitotenv.2017.09.323

409    3.    Lefrançois E, Apothéloz-Perret-Gentil L, Blancher P, Botreau S, Chardon C,

18

410   Crepin L, et al. Development and implementation of eco-genomic tools for aquatic

411   ecosystem biomonitoring: the SYNAQUA French-Swiss program. Environ Sci

412   Pollut Res. 2018;25: 33858–33866. doi:doi.org/10.1007/s11356-018-2172-2

413   4.   Keck F, Vasselon V, Tapolczai K, Rimet F, Bouchez A. Freshwater biomonitoring

414   in the Information Age. Front Ecol Environ. 2017;15: 266–274.

415   doi:10.1002/fee.1490

416   5.   Hajibabaei M, Porter TM, Wright M, Rudar J. COI metabarcoding primer choice

417   affects richness and recovery of indicator taxa in freshwater systems. PLoS One.

418   2019;14: 1–18. doi:10.1371/journal.pone.0220953

419   6.   Taberlet P, Coissac E, Pompanon F, Brochmann C, Willerslev E. Towards next-

420   generation biodiversity assessment using DNA metabarcoding. Mol Ecol.

421   2012;21: 2045–2050. doi:10.1111/j.1365-294X.2012.05470.x

422   7.   Baird DJ, Hajibabaei M. Biomonitoring 2.0: A new paradigm in ecosystem

423   assessment made possible by next-generation DNA sequencing. Mol Ecol.

424   2012;21: 2039–2044. doi:10.1111/j.1365-294X.2012.05519.x

425   8.   Borisko JP, Kilgour BW, Stanfield LW, Jones FC. An Evaluation of Rapid

426   Bioassessment Protocols for Stream Benthic Invertebrates in Southern Ontario ,

427   Canada. Water Qual Res J Canada. 2007;42: 184–193.

428   9.   Brua RB, Culp JM, Benoy GA. Comparison of benthic macroinvertebrate

429   communities by two methods: Kick- and U-net sampling. Hydrobiologia. 2011;658:

430   293–302. doi:10.1007/s10750-010-0499-x

431   10.   Hajibabaei M, Porter TM, Robinson C V., Baird DJ, Shokralla S, Wright MTG.

432   Watered-down biodiversity? A comparison of metabarcoding results from DNA

433    extracted from matched water and bulk tissue biomonitoring samples. PLoS One.

434    2019;14: 1–16. doi:10.1371/journal.pone.0225409

435  11.  Emilson CE, Thompson DG, Venier LA, Porter TM, Swystun T, Chartrand D, et al.

436    DNA metabarcoding and morphological macroinvertebrate metrics reveal the

437    same changes in boreal watersheds across an environmental gradient. Sci Rep.

438    2017;7: 1–11. doi:10.1038/s41598-017-13157-x

439  12.  Leese F, Altermatt F, Bouchez A, Ekrem T, Hering D, Meissner K, et al. DNAqua-

440    Net: Developing new genetic tools for bioassessment and monitoring of aquatic

441    ecosystems in Europe. Res Ideas Outcomes. 2016;2: e11321.

442    doi:10.3897/rio.2.e11321

443  13.  Gibson JF, Shokralla S, Curry C, Baird DJ, Monk WA, King I, et al. Large-Scale

444    Biomonitoring of Remote and Threatened Ecosystems via High-Throughput

445    Sequencing. PLoS Genet. 2015;10: 1–15. doi:10.5061/dryad.vm72v

446  14.  Chonova T, Kurmayer R, Rimet F, Labanowski J, Vasselon V, Keck F, et al.

447    Benthic Diatom Communities in an Alpine River Impacted by Waste Water

448    Treatment Effluents as Revealed Using DNA Metabarcoding. Front Microbiol.

449    2019;10: 1–17. doi:10.3389/fmicb.2019.00653

450  15.  Rivera SF, Vasselon V, Jacquet S, Bouchez A, Ariztegui D, Rimet F.

451    Metabarcoding of lake benthic diatoms : from structure assemblages to

452    ecological assessment. Hydrobiologia. 2018;807: 37–51. doi:10.1007/s10750-

453    017-3381-2

454  16.  Vasselon V, Rimet F, Tapolczai K, Bouchez A. Assessing ecological status with

455    diatoms DNA metabarcoding : Scaling-up on a WFD monitoring network (

456     Mayotte island , France ). Ecol Indic. 2017;82: 1–12.

457     doi:10.1016/j.ecolind.2017.06.024

458  17.  Pandey LK, Bergey EA, Lyu J, Park J, Choi S, Lee H, et al. The use of diatoms in

459     ecotoxicology and bioassessment: Insights , advances and challenges. Water

460     Res. 2017;118: 39–58. doi:10.1016/j.watres.2017.01.062

461  18.  Bailet B, Bouchez A, Franc A, Frigerio J-M, Keck F, Karjalainen S-M, et al.

462     Molecular versus morphological data for benthic diatoms biomonitoring in

463     Northern Europe freshwater and consequences for ecological status.

464     Metabarcoding and Metagenomics. 2019;3: 21–35. doi:10.3897/mbmg.3.34002

465  19.  Visco JA, Apothéloz-Perret-Gentil L, Cordonier A, Esling P, Pillet L, Pawlowski J.

466     Environmental Monitoring: Inferring the Diatom Index from Next-Generation

467     Sequencing Data. Environ Sci Technol. 2015;49: 7597–7605.

468     doi:10.1021/es506158m

469  20.  Blanco S, Bécares E. Chemosphere Are biotic indices sensitive to river

470     toxicants? A comparison of metrics based on diatoms and macro-invertebrates.

471     Chemosphere. 2010;79: 18–25. doi:10.1016/j.chemosphere.2010.01.059

472  21.  Muñoz I, Sabater S. Integrating chemical and biological status assessment:

473     Assembling lines of evidence for the evaliuation of river ecosystem risk. Acta Biol

474     Colomb. 2014;19: 25–33.

475  22.  Sharifinia M, Mahmoudifard A, Gholami K, Namin JI, Ramezanpour Z. Benthic

476     diatom and macroinvertebrate assemblages, a key for evaluation of river health

477     and pollution in the Shahrood River, Iran. Limnology. 2016;17: 95–109.

478     doi:10.1007/s10201-015-0464-5

479    23.    King L, Clarke G, Bennion H, Kelly M, Yallop M. Recommendations for sampling

480           littoral diatoms in lakes for ecological status assessments. J Appl Phycol.

481           2006;18: 15–25. doi:10.1007/s10811-005-9009-3

482    24.    Aloi JE. Review of Recent Freshwater Periphyton Field Methods. Can J Fish

483           Aquat Sci. 1990;47: 656–670.

484    25.    Smucker NJ, Vis ML. Contributions to diatom diversity and distributional patterns

485           in streams□: implications for conservation. Biodivers Conserv. 2011;20: 643–661.

486           doi:10.1007/s10531-010-9972-0

487    26.    Weilhoefer CL, Pan Y. A comparison of diatom assemblages generated by two

488           sampling protocols. J North Am Benthol Soc. 2007;26: 308–318.

489    27.    Rimet F, Bouchez A. Biomonitoring river diatoms: Implications of taxonomic

490           resolution. Ecol Indic. 2012;15: 92–99. doi:10.1016/j.ecolind.2011.09.014

491    28.    Winter JG, Duthie HC. Stream epilithic, epipelic and epiphytic diatoms: Habitat

492           fidelity and use in biomonitoring. Aquat Ecol. 2000;34: 345–353.

493           doi:10.1023/A:1011461727835

494    29.    Álvarez-Blanco I, Cejudo-Figueiras C, Bécares E, Blanco S. Spatiotemporal

495           changes in diatom ecological profiles: Implications for biomonitoring. Limnology.

496           2011;12: 157–168. doi:10.1007/s10201-010-0333-1

497    30.    Smucker NJ, Vis ML. Diatom biomonitoring of streams: Reliability of reference

498           sites and the response of metrics to environmental variations across temporal

499           scales. Ecol Indic. 2011;11: 1647–1657. doi:10.1016/j.ecolind.2011.04.011

500    31.    Lange-Bertalot H, Krammer K. Bacillariaceae, Epithemiaceae, Surirellaceae. J.

501           Kramer; 1987.

502   32.   Vasselon V, Domaizon I, Rimet F, Kahlert M, Bouchez A. Application of high-

503         throughput sequencing (HTS) metabarcoding to diatom biomonitoring: Do DNA

504         extraction methods matter? Freshw Sci. 2017;36: 162–177. doi:10.1086/690649

505   33.   Kermarrec L, Franc A, Rimet F, Chaumeil P, Frigerio JM, Humbert JF, et al. A

506         next-generation sequencing approach to river biomonitoring using benthic

507         diatoms. Freshw Sci. 2014;33: 349–363. doi:10.1086/675079

508   34.   Rimet F, Vasselon V, Keszte BA, Bouchez A. Do we similarly assess diversity

509         with microscopy and high-throughput sequencing? Case of microalgae in lakes.

510         Org Divers Evol. 2018;18: 51–62. doi:10.1007/s13127-018-0359-5

511   35.   Tapolczai K, Keck F, Bouchez A, Rimet F, Kahlert M, Vasselon V. Diatom DNA

512         Metabarcoding for Biomonitoring: Strategies to Avoid Major Taxonomical and

513         Bioinformatical Biases Limiting Molecular Indices Capacities. Front Ecol Evol.

514         2019;7: 1–15. doi:10.3389/fevo.2019.00409

515   36.   Gazedam E, Gharabaghi B, Jones FC, Whiteley. Evaluation of the Qualitative

516         Habitat Evaluation Index as a Planning and Design Tool for Restoration of Rural

517         Ontario Waterways. Can Water Resour J. 2011;36: 149–158.

518   37.   Environment Canada. Canadian aquatic Biomonitoring Network -Field Manual:

519         Wadeable Streams. 2012.

520   38.   Guillera-Arroita G, Lahoz-Monfort J., Rooyen AR Van, Weeks AR, Tingley R.

521         Dealing with false-positive and false-negative errors about species occurrence at

522         multiple levels. Methods Ecol Evol. 2017;8: 1081–1091. doi:10.1111/2041-

523         210X.12743

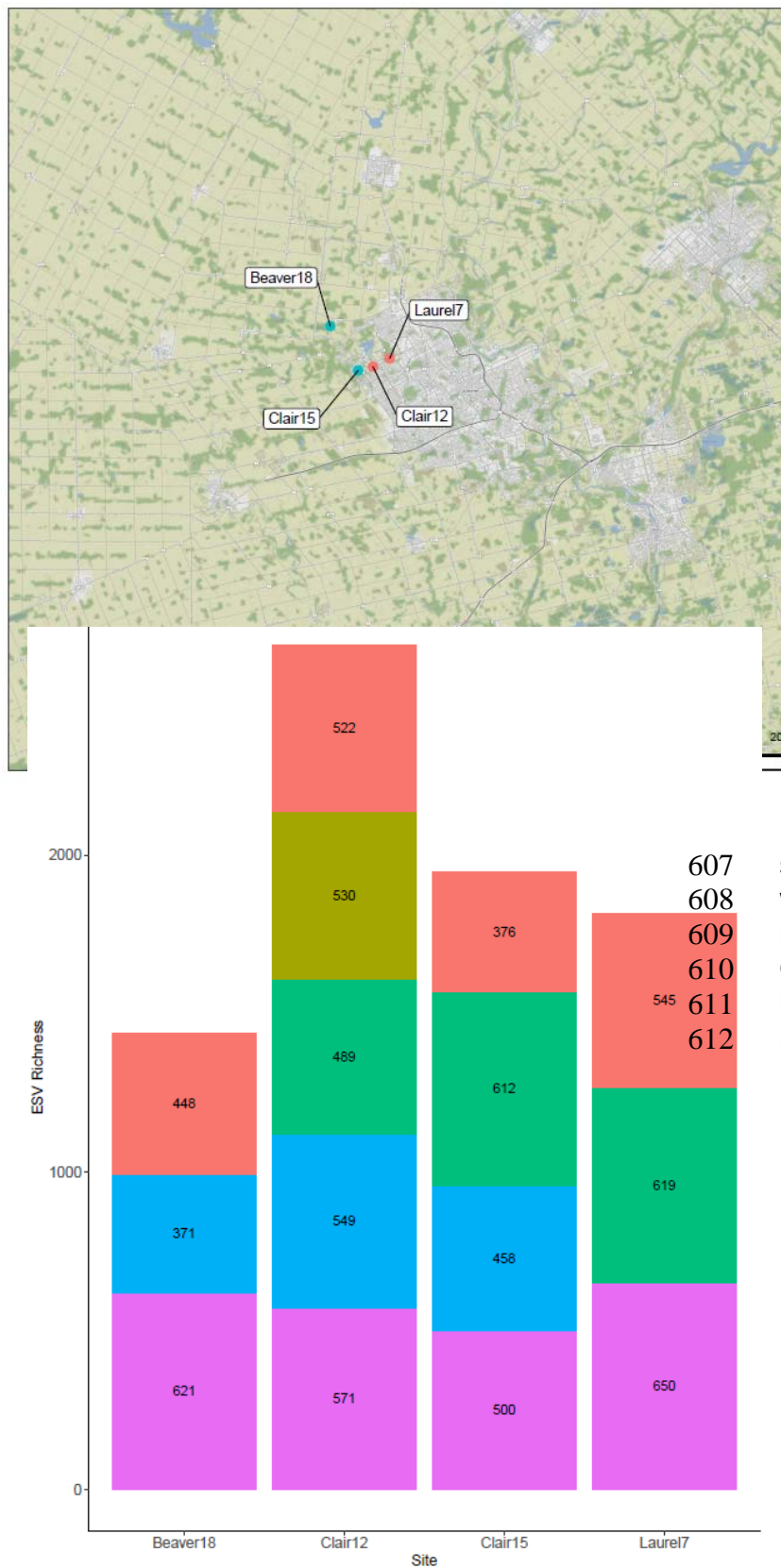524   39.   Rivera SF, Vasselon V, Jacquet S, Bouchez A, Ariztegui D, Rimet F.

23

525        Metabarcoding of lake benthic diatoms: from structure assemblages to ecological

526        assessment. Hydrobiologia. 2018;807: 37–51. doi:10.1007/s10750-017-3381-2

527   40.   Köster J, Rahmann S. Snakemake - A scalable bioinformatics workflow engine.

528        Bioinformatics. 2012;28: 2520–2522. doi:10.1093/bioinformatics/bty350

529   41.   Anaconda. Conda. 2016.

530   42.   Seqprep.

531   43.   Martin M. Cutadapt removes adapter sequences from high-throughput

532        sequencing reads. EMBnet.journal. 2011;17: 10–12.

533   44.   Callahan BJ, McMurdie PJ, Holmes SP. Exact sequence variants should replace

534        operational taxonomic units in marker-gene data analysis. ISME J. 2017;11:

535        2639–2643. doi:10.1038/ismej.2017.119

536   45.   Rognes T, Flouri T, Nichols B, Quince C, Mahé F. VSEARCH: A versatile open

537        source tool for metagenomics. PeerJ. 2016;2016: 1–22. doi:10.7717/peerj.2584

538   46.   Edgar RC. UNOISE2: improved error-correction for Illumina 16S and ITS

539        amplicon sequencing. bioRxiv. 2016. doi:10.1101/081257

540   47.   Reeder J, Knight R. The "rare biosphere": A reality check. Nat Methods. 2009;6:

541        636–637. doi:10.1038/nmeth0909-636

542   48.   Rimet F, Chaumeil P, Keck F, Kermarrec L, Vasselon V, Kahlert M, et al. R-

543        Syst□:: diatom□: an open-access and curated barcode database for diatoms and

544        freshwater monitoring. Database. 2016;2016: 1–21.

545        doi:10.1093/database/baw016

546   49.   Wang Q, Garrity GM, Tiedje JM, Cole JR. Naïve Bayesian Classifier for Rapid

547        Assignment of rRNA Sequences into the New Bacterial Taxonomy. Appl Environ

548   Microbiol. 2007;73: 5261 LP – 5267. doi:10.1128/AEM.00062-07

549 50. Team Rs. RStudio: Integrated Development for R. In: RStudio, Inc., Boston, MA.

550   2016.

551 51. Weiss S, Xu ZZ, Peddada S, Amir A, Bittinger K, Gonzalez A, et al. Normalization

552   and microbial differential abundance strategies depend upon data characteristics.

553   Microbiome. 2017;5: 1–18. doi:10.1186/s40168-017-0237-y

554 52. Oksanen J, Blanchet FG, Friendly M, Kindt R, Legendre P, McGlinn D, et al.

555   vegan: Community Ecology Package. R package version 2.5-2. 2018.

556 53. Goral F, Schellenberg J. goeveg: Functions for Community Data and Ordinations.

557   2018.

558 54. Robinson C V, Porter TM, Wright MTG, Hajibabaei M. Propylene glycol-based

559   antifreeze as an effective preservative for DNA metabarcoding of benthic

560   arthropods. bioRxiv. 2020; 2020.02.28.970475. doi:10.1101/2020.02.28.970475

561 55. Baselga A, Orme CDL. betapart: An R package for the study of beta diversity.

562   Methods Ecol Evol. 2012;3: 808–812. doi:10.1111/j.2041-210X.2012.00224.x

563 56. Mauri M, Elli T, Caviglia G, Uboldi G, Azzi M. RAWGraphs: A visualisation

564   platform to create open outputs. Proceedings of the 12th Biannual Conference on

565   Italian Sigchi Chapter. New York, NY, USA: Association for Computing

566   Machinery; 2017. pp. 1–5. doi:10.1145/3125571.3125585

567 57. Buss DF, Carlisle DM, Chon TS, Culp J, Harding JS, Keizer-Vlek HE, et al.

568   Stream biomonitoring using macroinvertebrates around the globe: a comparison

569   of large-scale programs. Environ Monit Assess. 2015;187. doi:10.1007/s10661-

570   014-4132-8

571   58.   Vasselon V, Domaizon I, Rimet F, Kahlert M, Bouchez A. Application of high-

572         throughput sequencing ( HTS ) metabarcoding to diatom biomonitoring□: Do DNA

573         extraction methods matter□? Freshw Sci. 2017;36: 162–177.

574         doi:10.1086/690649.

575   59.   Castro E, Siqueira T, Melo AS, Bini LM, Landeiro VL, Schneck F. Compositional

576         uniqueness of diatoms and insects in subtropical streams is weakly correlated

577         with riffle position and environmental uniqueness. Hydrobiologia. 2019;842: 219–

578         232. doi:10.1007/s10750-019-04037-8

579   60.   Yang X, Lv H, Li W, Guo M, Zhang X. Effect of water motion and microhabitat

580         preferences on spatio-temporal variation of epiphytic communities: a case study

581         in an artificial rocky reef system, Laoshan Bay, China. Environ Sci Pollut Res.

582         2018;25: 12896–12908. doi:10.1007/s11356-018-1349-z

583   61.   Debenest T, Pinelli E, Coste M, Silvestre J, Mazzella N, Madigou C, et al.

584         Sensitivity of freshwater periphytic diatoms to agricultural herbicides. Aquat

585         Toxicol. 2009;93: 11–17. doi:10.1016/j.aquatox.2009.02.014

586   62.   Reavie ED, Cai M. Consideration of species-specific diatom indicators of

587         anthropogenic stress in the Great Lakes. PLoS One. 2019;14: 1–15.

588         doi:10.1371/journal.pone.0210927

589   63.   Jakovljević OS, Popović SS, Vidaković DP, Stojanović KZ, Krizmanić J. The

590         application of benthic diatoms in water quality assessment (Mlava River, Serbia).

591         Acta Bot Croat. 2016;75: 199–205. doi:10.1515/botcro-2016-0032

592   64.   Srivastava P, Grover S, Verma J, Khan AS. Applicability and efficacy of diatom

593          indices in water quality evaluation of the Chambal River in Central India. Environ

594          Sci Pollut Res. 2017;24: 25955–25976. doi:10.1007/s11356-017-0166-0

**Fig. 1. Map of sample sites located within the Waterloo region (Ontario, Canada)** Scale bar shown in km, site habitat status indicated in legend.

27

617

618

619

620

621

622

623

624

**Fig. 2. ESV richness varies across different sample types.** Methods refer to the

different sampling approaches analyzed (i.e. Kick-net, Macrophyte, Leaf Litter, Rock

and Sediment). Replicates are pooled. Based on rarefied data.

628

28

629

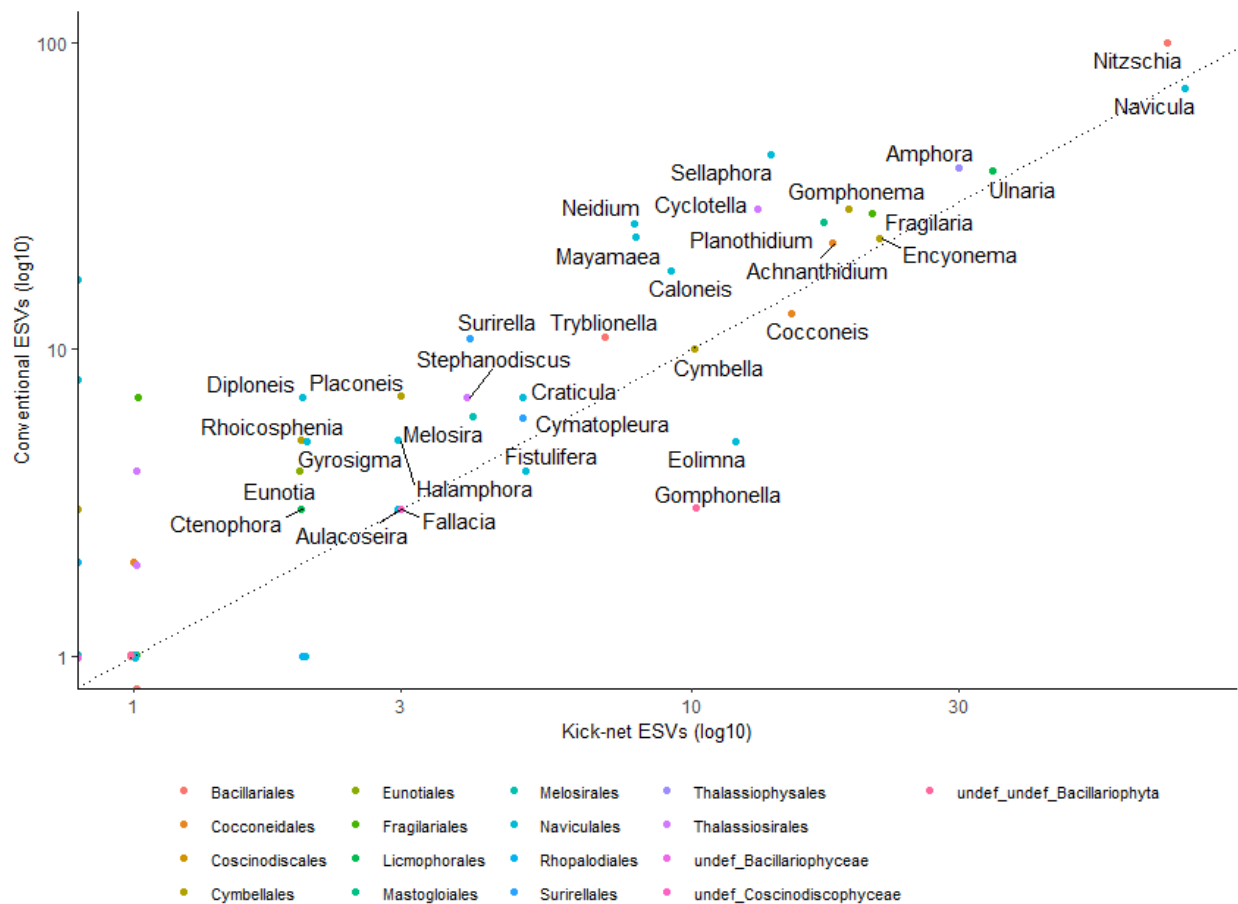**Fig. 3. A majority of diatom families were detected in both microhabitat and kick-net samples.**

632

633

634

635

636

637

**Fig. 4. Number of ESVs detected from genera detected from kick-net versus conventionally sampled diatoms are similar.** The points are color-coded for the orders detected in this study. A 1:1 correspondence line (dotted) is also shown. A log10 scale is shown on each axis to improve the spread of points with small values. Based on rarefied data.

643

644
645
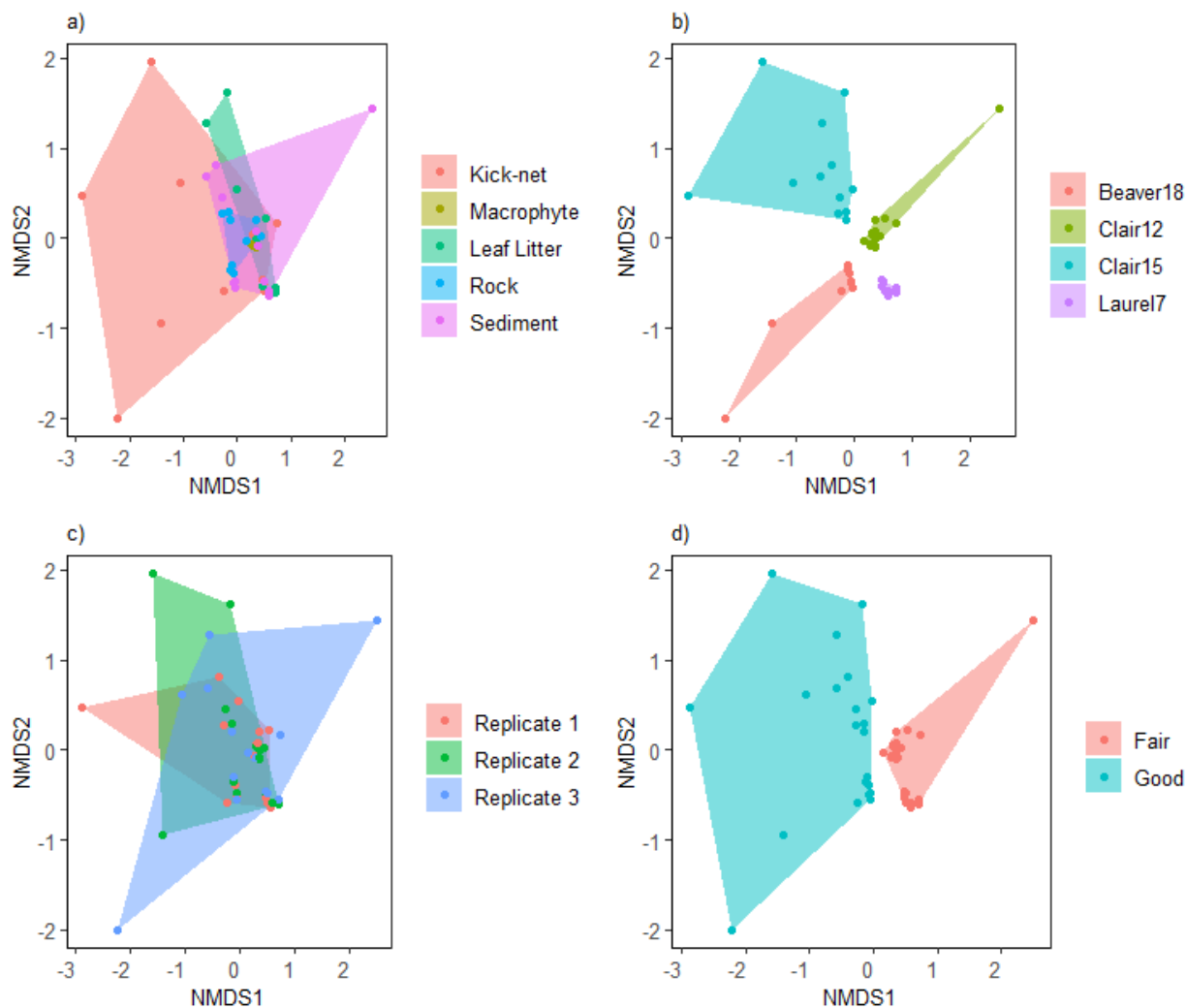646
647
648
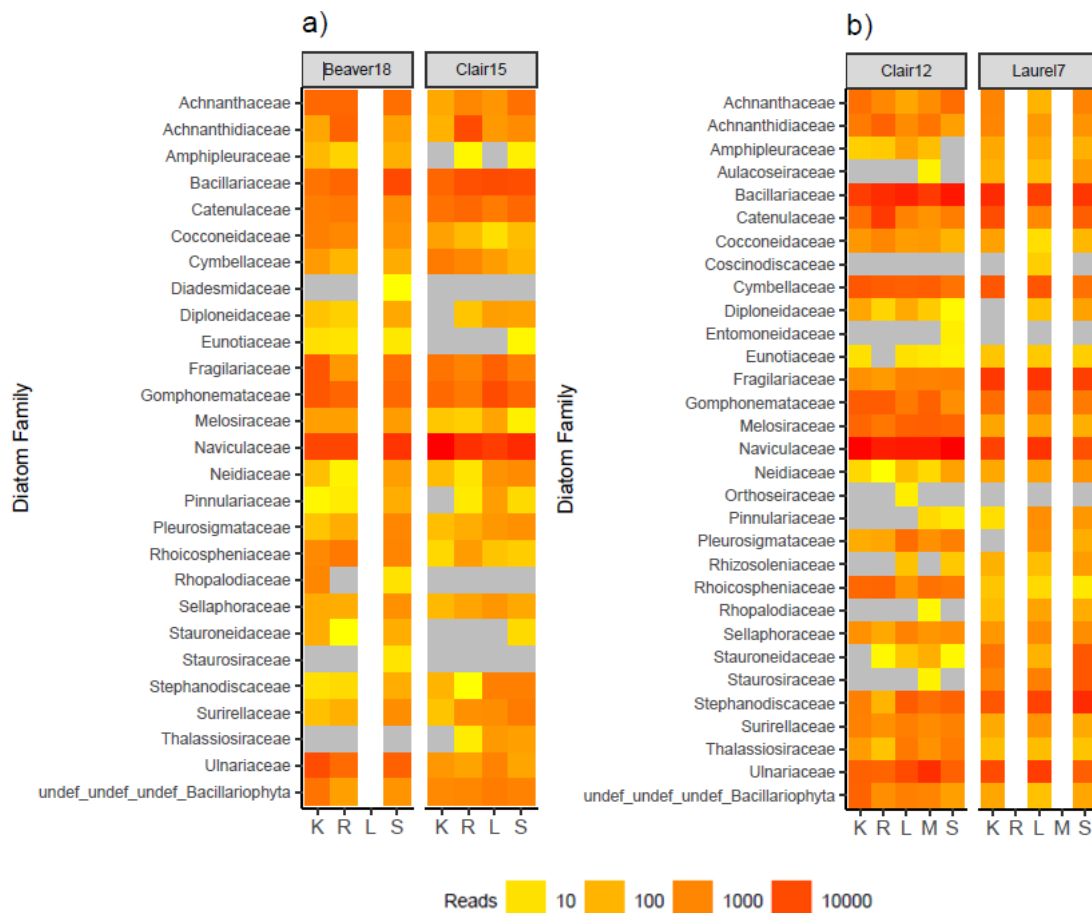
649

650

651

652

653 **Fig. 5. Non-metric multi-dimensional scaling plots show clustering mainly due to**

654 **site and status.** Specifically, a) binary Bray Curtis (Sorensen) dissimilarities

655 overlapping across different sampling approaches, b) clustering by site, c) overlap

656 between replicates, and d) clustering based on habitat quality status (stress = 0.111, $R^2$

657 = 0.98). Based on rarefied data.

658

667

**Fig. 6. Samples detect similar diatom families across sampling methods and site**

**status.** Only ESVs taxonomically assigned to families with high confidence (bootstrap

support >= 0.60 for 95% accuracy) are included. Part a) shows sites with a 'good'

quality status b) sites with a 'fair' quality status. Sampling methods: K = kick-net; R =

rock scraping; L = leaf litter; M = macrophyte; S = sediment. Empty lanes indicate the

corresponding microhabitat was not present at the site. For each site, three replicates

for each sampling method are pooled. Based on normalized data.

675

676

## **Data Availability**

677

678 Raw sequences will be available from NCBI SRA on acceptance. The SCVURL rbcL

679 metabarcode pipeline-1.0.2 is available from

680 https://github.com/terrimporter/SCVURL_rbcL_metabarcode_pipeline and the

681 rbcLdiatomClassifier v1 we used is available on GitHub at

682 https://github.com/terrimporter/rbcLdiatomClassifier.

683