



This document is a postprint version of an article published in Science of The Total Environment© Elsevier after peer review. To access the final edited and published work see <https://doi.org/10.1016/j.scitotenv.2021.149029>

Document downloaded from:



1 DNA metabarcoding reveals differences in distribution patterns and
2 ecological preferences among genetic variants within some key freshwater
3 diatom species

4 Javier Pérez-Burillo^{1,2}, Rosa Trobajo^{1*}, Manel Leira^{3,4}, François Keck^{5,6}, Frédéric
5 Rimet^{7,8}, Javier Sigró² & David G. Mann^{1,9}

6 ¹IRTA-Institute for Food and Agricultural Research and Technology, Marine and Continental Waters
7 Programme. Ctra de Poble Nou Km 5.5, E43540, Sant Carles de la Ràpita, Tarragona, Spain

8 ²Center for Climate Change (C3), Departament de Geografia, Universitat Rovira i Virgili, C/ Joanot
9 Martorell 15, E43500, Vila-seca, Tarragona, Spain

10 ³BioCost Research Group, Facultade de Ciencias and Centro de Investigacións Científicas Avanzadas
11 (CICA), Universidade de A Coruña, 15071, A Coruña, Spain

12 ⁴Biodiversity and Applied Botany Research Group, Departamento de Botánica, Facultade de Biología, Universidade
13 de Santiago de Compostela, 15782 Santiago de Compostela, Spain;

14 ⁵Eawag: Swiss Federal Institute of Aquatic Science and Technology, Dübendorf, Switzerland

15 ⁶Department of Evolutionary Biology and Environmental Studies, University of Zürich, Zürich,
16 Switzerland

17 ⁷INRAE, UMR Carrtel, 75 av. de Corzent, FR-74203 Thonon les Bains cedex, France

18 ⁸University Savoie Mont-Blanc, UMR CARRTEL, FR-73370 Le Bourget du Lac, France

19 ⁹Royal Botanic Garden Edinburgh, Edinburgh, EH3 5LR, Scotland, UK

20 *corresponding author: rosa.trobajo@irta.cat

21

22

23

24 **Abstract**

25 Our study evaluates differences in the distribution and ecology of genetic variants within
26 several ecologically important diatom species that are also key for Water Framework Directive
27 monitoring of European rivers: *Fistulifera saprophila* (FSAP), *Achnantheidium minutissimum*
28 (ADMI), *Nitzschia inconspicua* (NINC) and *Nitzschia soratensis* (NSTS). We used DADA2 to infer
29 amplicon sequence variants (ASVs) of a short *rbcl* barcode in 531 environmental samples from
30 biomonitoring campaigns in Catalonia and France. ASVs within each species showed different
31 distribution patterns. Threshold Indicator Taxa ANalysis revealed three ecological groupings of
32 ASVs in both ADMI and FSAP. Two of these in each species were separated by opposite
33 responses to calcium and conductivity. Boosted regression trees additionally showed that both
34 variables greatly influenced the occurrence of these groupings. A third grouping in FSAP was
35 characterized by a negative response to total organic carbon and hence was better
36 represented in waters with higher ecological status than the other FSAP ASVs, contrasting with
37 what is generally assumed for the species. In the two *Nitzschia* species, our analyses confirmed
38 earlier studies: NINC preferred higher levels of calcium and conductivity. Our findings suggest
39 that the broad ecological tolerance of some diatom species results from overlapping
40 preferences among genetic variants, which individually show much more restricted
41 preferences and distributions. This work shows the importance of studying the ecological
42 preferences of genetic variants within species complexes, now possible with DNA
43 metabarcoding. The results will help reveal and understand biogeographical distributions and
44 facilitate the development of more accurate biological indexes for biomonitoring programmes.

45

46 **Keywords.** ASV, environmental DNA, Water Framework Directive, *rbcl*, ecological preferences,
47 species distribution

48

49 **1. Introduction**

50 Diatoms play a crucial role in aquatic systems due, amongst other things, to their importance
51 in food webs and biogeochemical cycling and their great contribution to carbon fixation
52 (Armbrust, 2009; Mann, 1999; Smetacek, 1999). They are also widely used as ecological
53 indicators in palaeoenvironmental studies and biomonitoring programmes. For example, in
54 European rivers, it is compulsory (within the European Union) to monitor benthic diatom
55 communities (Water Framework Directive [WFD], Directive 2000/60/EC, 2000) because of their
56 rapid and specific response to environmental changes, great diversity, and ubiquitous
57 distribution, and the availability of information on the ecological preferences of many species.
58 However, it has become evident in the last two decades that many of these species are
59 complexes of genetic variants (e.g. Pinseel et al. 2017; Souffreau et al. 2013). These often show
60 scarcely discernible or no morphological differences (they are “cryptic”) and therefore it is
61 difficult or impossible to determine their geographical distributions and ecological preferences
62 using traditional methods based on microscopical identifications. Therefore, the significance of
63 this intraspecific variation is still not clear: although It is suggested that closely related diatoms
64 often share a similar ecology (Keck et al., 2016a, b, 2018b), it is also evident that they can
65 differ (Pinseel et al., 2017; Poulíčková et al. 2008, 2017; Rynearson et al., 2006).

66 DNA metabarcoding has recently been developed for biomonitoring the ecological
67 status of rivers (e.g. Kelly et al. 2020; Mortágua et al., 2019; Pérez-Burillo et al. 2020; Rivera et
68 al., 2020; Vasselon et al., 2017b) and it has proved as well to be a reliable and efficient method
69 for surveying species diversity from environmental samples (Deiner et al., 2017; Malviya et al.
70 2016; Piredda et al., 2017). DNA metabarcoding also offers a way to study the significance of
71 genetic variants within species, especially following the development of bioinformatic
72 pipelines such as DADA2 (Callahan et al., 2017), which use a denoising algorithm to remove
73 sequencing artifacts and generate ‘amplicon sequence variants’ (ASVs); these are believed to
74 be real DNA sequences that were present in the original environmental samples. A recent

75 example of using an ASV approach in diatoms was by Tapolczai et al. (2021), where they
76 assessed the responses of river diatom communities to agricultural land use in Hungary; in
77 some cases, they reported different ecological preferences among ASVs from the same
78 species. However, despite the clear potential of ASV-based metabarcoding approaches, there
79 do not appear to have been any studies to date that have used a large dataset to examine the
80 ecology and distribution of genetic variants and hence to elucidate their significance.

81 The aim of this work was therefore to study the distribution and ecological preferences
82 of different ASVs within selected species complexes of diatoms. For this we chose two groups
83 that are ecologically important and have been shown to be key for the WFD (Pérez-Burillo et
84 al., 2020): *Achnantheidium minutissimum* sensu lato (said by Potapova and Hamilton, 2007, to
85 be “one of the most frequently occurring diatoms in freshwater benthic samples globally”) and
86 *Fistulifera saprophila*. Both are very small-celled species that are difficult to treat
87 morphologically. In addition, we selected *Nitzschia inconspicua*, because Sanger sequencing
88 has already demonstrated a complex pattern of genetic and physiological variation within it
89 (Rovira et al., 2015), and *N. soratensis*, which is so similar to *N. inconspicua* in the light
90 microscope that identifying the two species and determining their ecological separation is
91 highly challenging (Kelly et al., 2015). More specifically, we asked 1) do genetic variants (ASVs)
92 within a species complex have similar geographical distributions within the study area? 2) Do
93 ASVs within a species complex have the same ecological preferences or do they differ? 3) If
94 there are differences in the ecological preferences of genetic variants within a species, do
95 these correlate with their phylogeny?

96 To answer these questions, we used a large molecular dataset extracted from
97 environmental samples collected in several river biomonitoring campaigns in contiguous areas
98 of France and Catalonia (NE Spain). For evaluating ecological preferences and the spatial
99 distributions of ASVs, we performed Threshold Indicator Taxa Analyses (TITAN) and Boosted
100 Regression Trees (BRT) analyses, since both methods have been successfully applied in

101 morphological and metabarcoding studies addressing stressor-response and species
102 distribution models (Lanzén et al., 2020; Smucker et al., 2020; Soininen et al., 2018, Wagenhoff
103 et al., 2017).

104

105 **2. Material and Methods**

106 *2.1 Study site and diatom sampling*

107 The dataset used in this study consisted of 610 benthic diatom samples collected from both
108 Catalan and French biomonitoring networks. Samples were originally taken as a part of the
109 2017 Catalan biomonitoring programme and two French monitoring campaigns held in 2016
110 and 2017. The hydrographic area of Catalonia is divided into internal and interregional
111 hydrographic basins. The internal basins comprise a total of eleven main rivers, the basins of
112 the rivers Llobregat and Ter being the most extensive, and the interregional basins cover the
113 Catalan sections of the rivers Ebro, Garona (Garonne in French) and Xúquer. The French
114 monitoring network area corresponds to seven main basins (Adour–Garonne, Artois–Picardie,
115 Loire–Bretagne, Rhin–Meuse, Rhône–Méditerranée, Corse, and Seine-Normandie) of which
116 the largest belong to the rivers Loire, Rhône, Seine and Garonne (Supplementary Fig. 1).

117 All Catalan sites were sampled for periphyton between April and July of 2017 following
118 standard procedures (CEN, 2014). French sites were sampled between February and December
119 and between February and October, for the campaigns held in 2016 and 2017 respectively,
120 and followed French NFT 90 354 (AFNOR, 2007) and European (CEN, 2014) standards. At each
121 site, diatoms were collected from at least five stones by brushing their upper surfaces using a
122 toothbrush. The resulting samples were preserved by adding $\geq 90\%$ ethanol (to a final
123 concentration of 70%) and used for DNA metabarcoding analysis following the
124 recommendations of the technical report of the European Committee for Standardization
125 (CEN, 2018).

126

127 2.2 Physicochemical and biotic parameters

128 Physicochemical parameters that constituted the environmental dataset used for French and
129 Catalan river sites were obtained from the online “Naiades”
130 (<http://www.naiades.eaufrance.fr/>) and “SDIM” (<http://aca-web.gencat.cat/sdim21/>) water
131 quality datasets. Environmental parameters selected in this study were ammonium (NH_4^+ ;
132 mg/L), bicarbonates (HCO_3^- ; mg/L), calcium (mg/L), total organic carbon (TOC; mg/L),
133 conductivity ($\mu\text{S}/\text{cm}$), nitrates (NO_3^- ; mg/L), orthophosphates (PO_4^{3-} ; mg/L), pH, sulphates
134 (SO_4^{2-} ; mg/L), water temperature ($^\circ\text{C}$) and altitude (m) (Table 1). The measures selected for
135 these parameters corresponded to the mean of all the records available for a period of 80 days
136 preceding and 10 days following the biological sampling. The diatom indices IBD (“Indice
137 Biologique Diatomées”) and IPS (“Indice de Polluosensibilité spécifique”) were retrieved
138 respectively for French and Catalan rivers sites analysed.

139

140 2.3 DNA extraction, PCR amplification and high-throughput sequencing (HTS)

141 The procedures for DNA extraction, PCR amplification and HTS for French and Catalan rivers
142 are described in Rivera et al. (2020) and Pérez-Burillo et al. (2020), respectively. Briefly, DNA
143 extraction of French samples from the 2016 campaign was performed using GenElute TM-LPA
144 protocol, while the Macherey–Nagel NucleoSpin® soil kit (MN-Soil) protocol was followed for
145 DNA extraction of Catalan and French samples from the 2017 campaigns. A short *rbcl* region of
146 312 bp constituted the DNA marker and this was amplified by PCR using an equimolar mix of
147 the modified versions of the primers Diat_rbcl_708F (forward) and R3 (reverse) given by
148 Vasselon et al. (2017b). Four Illumina Miseq runs were performed for sequencing separately
149 the French (3 runs) and Catalan (1 run) samples. In order to prepare the HTS libraries using a 2-
150 step PCR strategy, half of P5 and P7 Illumina adapters were included to the 5' end of the

151 forward and reverse primers respectively. Adapter sequences used were
152 CTTTCCTACACGACGCTCTCCGATCT (P5) and GGAGTTCAGACGTGTGCTCTCCGATCT (P7) for
153 French samples and TCGTCGGCAGCGTCAGATGTGTATAAGAGACAG (P5) and
154 GTCTCGTGGGCTCGGAGATGTGTATAAGAGACA (P7) for Catalan samples.

155 PCR1 reactions for each DNA sample were performed in triplicate using 1 µL of the
156 extracted DNA in a final volume of 25 µL. Conditions and the reaction mix of the PCR1 followed
157 the procedure described in Vasselon et al. (2017b) For each sample, the three PCR1 replicates
158 were pooled and sent for sequencing to “Plateforme Génome Transcriptome” (PGTB,
159 Bordeaux, France) or “GenoToul Genomics and Transcriptomics” (GeT-PlaGe, Auzeville,
160 France), where the PCR1 products were purified and used as template for a second round of
161 PCR (PCR2), with Illumina tailed primers targeting the half of P5 and P7 adapters. Finally, all
162 generated amplicons were dual indexed and pooled for sequencing on an Illumina MiSeq
163 platform using the V3 and V2 paired-end sequencing kits (250 bp × 2) for the French and
164 Catalan samples respectively.

165 The influence of DNA extraction methods on the diatom inventory produced by DNA
166 metabarcoding has been evaluated by Vasselon et al. (2017a). In this study, the authors
167 evaluated 5 different extraction methods, including the two methods used in our study. They
168 found some slight differences in relative abundance between methods for some particular
169 species, but the differences in community composition caused were far less than the
170 differences attributable to habitat. Importantly for the current analyses, the slight differences
171 did not affect species richness (“Regardless of taxonomic level (OTU or species), the taxonomic
172 composition of the community represented in the extracts was not affected by DNA extraction
173 methods...”: Vasselon et al. 2017a). Furthermore, (1) the taxa contributing more than >1% of
174 the dissimilarities between diatom communities obtained with the GenElute and MN-Soil
175 protocols did not include either *Achnantheidium* or *Fistulifera* (Vasselon et al. 2017a, Table 3),
176 and (2) all our comparisons were made among ASVs belonging to the same species complex

177 and hence of diatoms with very similar physical characteristics (frustule shape, size and
178 robustness). Thus, we can expect that the different extraction methods used for the French
179 2016 and 2017 datasets will not have greatly affected either the presence/absence of ASVs
180 (especially since a high threshold of abundance was set for inclusion in the analyses) or relative
181 abundances across the combined dataset. Comparisons across wider ranges of species (e.g. all
182 *Achnantheidium*, all *Navicula*, etc) might have been more seriously affected.

183 2.4. Bioinformatic analysis

184 The sequencing facilities performed the demultiplexing of all the samples providing two fastq
185 files per sample, one corresponding to forward reads (R1) and one to reverse reads (R2). All
186 the demultiplexed Miseq reads were treated together using the R package *DADA2*, following
187 the method described by Callahan et al. (2016). Primers were removed from the R1 and R2
188 reads using *cutadapt* (Martin, 2011). The resulting R1 and R2 reads were truncated to 200 and
189 170 nucleotides respectively, based on to their quality profile (median quality score < 30) and
190 those reads with ambiguities or an expected error (maxEE) higher than 2 were discarded. The
191 *DADA2* denoising algorithm was applied to determine an error rates model in order to infer
192 amplicon sequence variants (ASVs); ASVs detected as chimeras were discarded using the
193 function “*removeBimeraDenovo*” implemented in *DADA2*. Finally, the taxonomic affiliation of
194 the ASVs was determined using the database “A ready-to-use database for *DADA2*:
195 *Diat.barcode_rbcL_312bp_DADA2*” (Chonova et al., 2020), which is derived from the curated
196 diatom reference library *Diat.barcode* v9 (Rimet et al., 2019, available at
197 https://www6.inra.fr/cartel-collection_eng/Barcoding-database and at
198 <https://data.inrae.fr/file.xhtml?persistentId=doi:10.15454/TOMBYZ/IEGUXB&version=10.0>),
199 and the naïve Bayesian classifier method (Wang et al., 2007); 50% was set as the minimum
200 confidence threshold (the default in *DADA2*). In this study we focused on ASVs that were
201 assigned by the pipeline to *Nitzschia inconspicua*, *N. soratensis*, *Achnantheidium minutissimum*
202 and *Fistulifera saprophila*. Of these, we retained for subsequent analyses only those with \geq

203 1000 reads and occurring in ≥ 2 samples with environmental data available, in order to remove
204 rare ASVs and residual sequencing artifacts. The ASVs were numbered according to the rank
205 order of their abundance; so, for example, *A. minutissimum* ASV6 was the sixth most abundant
206 sequence in the whole dataset.

207

208 2.5. Phylogenetic analyses

209 Phylogenetic analyses were performed in order to 1) elucidate the phylogeny of the different
210 ASVs obtained from *Nitzschia inconspicua*, *N. soratensis*, *Achnantheidium minutissimum* and
211 *Fistulifera saprophila*, and 2) assess the taxonomic assignation obtained after executing the
212 bioinformatics analyses by examining the phylogenetic relatedness between the ASVs and
213 curated reference sequences from Diat.barcode v9 (together with some other, more recent
214 sequences present in GenBank: <https://ncbi.nlm.nih.gov/>). For this purpose, maximum
215 likelihood trees were constructed using ASVs and the reference sequences. A first tree
216 included reference sequences and ASVs classified into *N. inconspicua* and *N. soratensis* species,
217 while a second and a third used those ASVs and reference sequences classified into *A.*
218 *minutissimum* and *F. saprophila* respectively. All three analyses were performed using
219 raxmlGUI with the GRT-Gamma model, with 1000 replicates for the bootstrap analyses.
220 Reference sequences and ASVs used for building each of the three trees were previously
221 aligned using the Muscle alignment algorithm (Edgar, 2004) in MegaX software (Kumar et al.,
222 2018). All the three trees calculated were drawn using iTOL (<https://itol.embl.de>) (Letunic et
223 al., 2019).

224

225 2.6. Statistical analyses

226 2.6.1 Spatial variables

227 In order to study spatial distribution patterns of ASVs, Moran's eigenvector maps (MEMs) were
228 used on sampling sites' latitude and longitude to generate explanatory variables that represent
229 spatial patterns at different scale and can be used in canonical analysis. (Dray et al.,2006).
230 MEMs are produced by the diagonalization of a spatial weighting matrix, which is obtained as
231 the Hadamard product of a connectivity matrix by a similarity matrix. The connectivity matrix
232 was based on Gabriel's graph geometrical connection scheme due to the non-regular
233 distribution of the sampling sites (Legendre and Legendre, 2012). The R package *adespatial*
234 (Dray et al., 2020) was used for calculating MEMs.

235

236 2.6.2 Redundancy analyses

237 ASV abundance data were Hellinger transformed and all environmental variables except pH
238 were standardized following $X_{st} = (X - \mu)/SD$. Variance inflation factors (VIFs) were calculated to
239 check the presence of collinearities among environmental variables and those variables with
240 VIF >10 were removed to avoid the impact of collinearity. Forward selection with two stopping
241 criteria (alpha significance level and adjusted coefficient of multiple determination, Blanchet et
242 al., 2008) was applied separately on environmental and MEMs sets of variables. Two
243 redundancy analyses (RDA) models were performed in order to analyse separately the
244 relationships between the selected environmental and spatial variables (MEMs) and the ASVs.
245 R packages *adespatial* (Dray et al., 2020) and *vegan* (Oksanen et al., 2020) were used for
246 performing forward selection and RDA models respectively.

247

248 2.6.3 TITAN analyses

249 Threshold indicator taxa analyses (TITAN) were conducted in order to characterize ASV-specific
250 responses for each environmental variable. TITAN handles multiple response variables (ASVs)
251 but only one explanatory variable (i.e. environmental variables) at each analysis and it detects

252 change points, which are the values of the environmental gradient where the greatest change
253 in taxon abundance and frequency occurs within the observed samples. TITAN standardizes
254 the magnitude of responses as z scores in order to facilitate cross-taxa comparison. Z scores
255 reflect the type of response, positive (+ z scores) or negative (- z scores), of a particular taxon
256 (ASV in this case) along the environmental gradient and the sum of the z scores (sum z) gives
257 information about the assemblage responses, either negative (sum -z) or positive (sum +z),
258 along the gradient, the maximum z score occurring at the point at which change in assemblage
259 composition is greatest (Baker & King, 2010). We conducted TITAN analyses for each of the
260 environmental parameters and using ASV relative abundance. Number of permutations was
261 set to 250, number of bootstrap replicates used was 500, the minimum number of
262 observations required on each side of a candidate change point was 5 and the TITAN filtering
263 metrics of uncertainty “purity” and “reliability”, used to separate reliable responders from
264 stochastic noise along the gradient, were set to 0.95. Z scores obtained for those ASVs whose
265 responses fulfilled purity and reliability criteria for at least 4 environmental variables were
266 hierarchically clustered and visualized through heatmaps in order distinguish groups of ASVs
267 with similar response patterns for environmental data. For that, Euclidean distance and ward.
268 D functions (Ward, 1963) were used to compute dissimilarity distance and hierarchical
269 clustering respectively. On the other hand, Kruskal–Wallis (Hollander and Wolfe, 1973) tests
270 with post hoc Dunn’s test (Dunn, 1964) were performed to determine environmental data
271 statistically significant ($p < 0.05$) among the sites where species and ecological groupings
272 occurred. We used the implementation available in the R packages *TITAN2* (Baker et al., 2019),
273 *gplots* (Warnes et al., 2020), *stats* (R core team, 2020) and *dunn.test* (Dinno, 2017) to conduct
274 the TITAN analyses, heatmaps, Kruskal–Wallis test and Dunn’s test respectively.

275

276 2.6.4 Boosted regression tress

277 Relationships of the groups of ASVs, defined after TITAN analysis, with environmental variables
278 were additionally evaluated using boosted regression trees (BRT). BRT is a machine learning
279 model that uses a boosting algorithm to combine large numbers of decision trees for
280 improving model accuracy (Elith et al., 2008). BRT handles multiple explanatory variables
281 (environmental variables) but only one response variable (groups of ASVs in our case). It
282 estimates the relative importance of environmental variables on the basis of the number of
283 times that a variable is selected and the extent to which it improves the model (Friedman,
284 2001). Partial dependence plots generated by BRT show the marginal effect of each predictor
285 on the response variable while accounting for the average effects of the other variables used
286 in the model. Thus, these plots are useful for comparing the relationship and influence of each
287 explanatory variable on the response variable. BRT analyses were conducted using the
288 Bernoulli family of presence/absence ASVs reads, a bag fraction of 0.5, a learning rate of 0.001
289 and a tree complexity of 3. BRT models were evaluated using a 10-fold cross validation
290 procedure (i.e. 90% of data is used for training and 10% for validation). The *dismo* (Hijmans et
291 al. 2020) R package was used to perform BRT analyses.

292

293 **3. Results**

294 **3.1 Metabarcoding data**

295 30,251,272 reads were obtained by Miseq Illumina sequencing of a total of 610 samples from
296 Catalan and French rivers. After quality filtering steps 25,452,802 reads and 6403 ASVs were
297 obtained, of which 148, 83, 29 and 14 were classified into *Achnanthydium minutissimum*,
298 *Fistulifera saprophila*, *Nitzschia inconspicua* and *N. soratensis*, respectively. After filtering to
299 remove ASVs having < 1000 reads and occurring in < 2 samples with environmental data, the
300 molecular inventory of the four species consisted of 531 samples and a total of 75 ASVs, of
301 which 45, 18, 9 and 3 belonged, respectively, to *Achnanthydium minutissimum*, *Fistulifera*

302 *saprophila*, *Nitzschia inconspicua* and *N. soratensis* (Supplementary table 1; Supplementary
303 data). We checked the 75 ASVs using MegaX and with reference to a matrix of available *rbcl*
304 sequences of diatoms. There was no evidence of sequencing artifacts, i.e. no indels, nor stop
305 codons, nor implausible amino-acids such as substitutions that have no parallel in other
306 diatoms or involve changes in the type of amino-acid (polar vs non-polar vs basic vs acidic).

307

308 3.2 Geographical distribution of ASVs within the study area

309 Most of the 20 most abundant ASVs from *A. minutissimum* (all with >10,000 reads in the
310 dataset) were widely distributed in both Catalan and French rivers, as shown in Fig. 1 and
311 Supplementary Fig. 2. However, despite their apparent ability to colonize sites across the
312 whole geographical region surveyed, individual ASVs seemed to show contagious distributions
313 (aka clumped distributions). For example, ASV70 dominated the *A. minutissimum* assemblage
314 in Mediterranean river sites from the south-central and north-east of Catalonia, but was much
315 less important than ASVs 6 and 7 over most parts of France, being an important ASV there only
316 in scattered sites, e.g. in Normandy and along the Loire (Fig. 1); ASV119 was an important
317 component at three sites located in the Pyrenees, but occurred also in the Jura and Alps
318 regions of eastern France (Fig. 1). Some of the 20 most abundant ASVs appeared to be
319 restricted to one or other country. Thus, ASVs 153 and 219 were only detected in Catalan
320 rivers (Fig. 1; Supplementary Fig. 2), and ASV269 only in French rivers (Supplementary Fig. 2).
321 Interestingly, one of the ASVs restricted to Catalonia, ASV219, was closely related – 1 bp
322 difference in the 263 bp alignment– to ASV7, which dominated the *A. minutissimum*
323 complement in the same central area of Catalonia where ASV219 occurred; likewise, ASV153
324 was generally found alongside or replacing ASV6 in central Catalan sites (Fig 1; Supplementary
325 Fig. 2), the two ASVs again differing by only 1 bp. Out of the 25 less abundant ASVs of *A.*

326 *minutissimum* considered here, 11 were common in both countries, 3 ASVs were only recorded
327 in Catalan rivers and 11 were only in French rivers (Supplementary Table 1).

328 In *Fistulifera* too, some ASVs showed wide patterns of distribution. Eleven out of the
329 total of 18 *F. saprophila* ASVs were detected in both Catalan and French rivers and most of
330 them were recorded at numerous sites across the study area (Supplementary Figs. 3 and 4;
331 Supplementary Table 1), Nevertheless, contagion was obvious. For example, ASVs 16 and 74
332 dominated in Catalonia, though they also occurred in scattered sites in France, mostly in
333 eastern parts, while ASV43 showed no obvious pattern in Catalonia, but was restricted to the
334 upper regions of the Rhone catchment in France (Supplementary Fig. 3). Three ASVs were
335 recorded only in French rivers and four in Catalan rivers (Supplementary Table 1). ASVs
336 exclusively recorded in French rivers (ASVs 187, 198 and 233) were much more abundant than
337 those only recorded in Catalan rivers (ASVs 643, 823 and 983) and were widely distributed
338 within France (Supplementary Figs. 3 and 4; Supplementary Table 1). For example, while there
339 was a noticeable concentration of ASV233 in the Garonne catchment of SE France, this ASV
340 also occurred in scattered locations in almost all the other major French basins surveyed
341 (Supplementary Fig. 3). One of the ASVs occurring only in Catalonia, ASV643, was restricted to
342 just two sites in the NE, where it formed c. 50 to more than 75% of the *Fistulifera* reads
343 (Supplementary Fig. 4). Conversely, *Fistulifera* ASVs 823 and 983 occurred in 13 and 17 sites,
344 respectively, but were always rare (Supplementary Fig. 4; Supplementary Table 1). *Fistulifera*
345 ASV234 tended to dominate any assemblage where it was present and rarely occurred with
346 any of the other common ASVs (Supplementary Fig. 3).

347 In the case of *N. inconspicua*, 7 out of the 9 ASVs analysed were detected in both
348 Catalan and French rivers, most of them (ASVs 53, 56, 113, 273, 463) being broadly distributed
349 (Supplementary Fig. 5; Supplementary Table 1). The remaining 2 ASVs were detected in only
350 one country, but were represented there by 10 (ASV572 in France) and 8 sites (ASV615 in
351 Catalonia) respectively, and were not restricted to a single catchment (Supplementary Fig. 5;

352 Supplementary Table 1); indeed, ASV572 spanned the whole of France, from the extreme
353 north to almost the most southerly site sampled (Supplementary Fig. 5). In France, ASVs 53
354 and 113 tended not to co-occur; for example, ASV53 picked out the course of the Loire but
355 ASV113 the Garonne. In Catalonia, the pattern seemed to differ, the two cooccurring quite
356 frequently (Supplementary Fig. 5). One of the rarer ASVs, ASV463, though only found at 14
357 sites (Supplementary Table 1), nevertheless occurred over a wide range, from the southern
358 part of Catalonia to the Jura mountains in eastern France (Supplementary Fig. 5). The 3 ASVs
359 identified as *N. soratensis* had a higher occurrence in French rivers, though all of them were
360 also identified in Catalan rivers (Supplementary Fig. 6; Supplementary Table 1). However, in
361 Catalonia ASV288 was found only at one site in the extreme north. There was no obvious
362 pattern in the geographical distribution of ASV94 and ASV117 in France (Supplementary Fig. 6).

363

364 3.3 ASVs identified in our study area in relation to previously sequenced haplotypes

365 The ASVs detected in our dataset extended the known genetic diversity of each species,
366 although the most abundant ASVs generally found a match among the Sanger sequences
367 already available. For example, of the five most abundant *A. minutissimum* ASVs, only one
368 (ASV70) represented a haplotype not in GenBank or Diat.barcode version 9 (Fig. 2) and this
369 ASV differed by only 1 bp from the most similar haplotypes. Likewise, the three most abundant
370 *F. saprophila* ASVs (ASVs 16, 43 and 48) all matched sequences already in Diat.barcode version
371 9 (Supplementary Fig. 7). Perhaps the most surprising 'newcomers' were ASV164 (24961
372 reads) and the clade of ASVs 156, 272 and 956 (together >42000 reads; Fig. 1), all within *A.*
373 *minutissimum*: none of these seem to have been found before, despite their high abundance
374 and wide distribution.

375 Among the matches between ASVs and Sanger sequences were several that were not
376 surprising, given the source of the clones previously sequenced. For example, in *N.*

377 *inconspicua*, ASVs were found that were 100% similar to Sanger-sequenced clones isolated
378 from southern Catalonia; these were ASVs 615, 113 and 273, which are identical to *N.*
379 *inconspicua* genotypes G2, G3 and G4 of Rovira et al. (2015; see also Mann et al. 2021 for a
380 four-gene analysis) (Supplementary Fig. 8). In other cases, however, the ASVs extended the
381 range of a haplotype. For example, in the *A. minutissimum* complex, ASV253 was 100%
382 identical to Sanger sequences from eastern N America and Portugal, and ASV77 and ASV256
383 were each identical to clones isolated from Siberia (Fig. 2). The Sanger sequences themselves
384 provide a further example of a widely distributed haplotype, with 100% *rbcl* identity between
385 clones isolated from Montana and Hawaii (GenBank accessions KJ658384 and KJ658385). Of
386 the *F. saprophila* ASVs, ASV48 was 100% identical to sequences from Luxembourg and South
387 Korea (Supplementary Fig. 7), while *N. inconspicua* ASV463 was identical to a *N. inconspicua*
388 isolated from the tropical Indian Ocean island of Mayotte (Supplementary Fig. 8).

389 In contrast, some *rbcl* Sanger sequences in the four species were not represented in
390 our HTS dataset. These included a clade of three *N. inconspicua* haplotypes from the islands of
391 La Réunion and Mayotte, c. 1400 km apart in the Indian Ocean (from clones TCC474, 510 and
392 571; all were at least 10 bp different from any *inconspicua* ASV in our dataset); and two *N.*
393 *inconspicua* haplotypes isolated from saline habitats in S Catalonia (G5 and G6, from 31–40
394 PSU and 5 PSU, respectively: Rovira et al. 2015) (Supplementary Fig. 8). In *Fistulifera* a ‘tropical
395 clade’ of haplotypes from Mayotte and S Japan had no parallel in our dataset, nor *F. alcalina*,
396 recently described from Florida, USA (Supplementary Fig. 7).

397

398 3.4 Redundancy analysis

399 Given the non-uniform distributions of the ASVs of all four species in the study area, we
400 examined their occurrence and abundance in relation to environmental variables. Those
401 selected by forward selection ($p < 0.05$) were altitude, calcium, conductivity, HCO_3^- , pH, PO_4^{3-} ,

402 SO_4^{2-} , TOC and water temperature. An RDA model that included these variables explained 15%
403 of the constrained variance, the first two axes accounting respectively for 7.3% and 4.6%
404 (Supplementary Fig 9). 69 MEMs were selected by forward selection and an RDA model that
405 included these selected MEMs explained a total of 39% of the constrained variance, of which
406 the first and second axes accounted for 11.3% and 10% respectively. This indicates an
407 important degree of spatial structuring of the ASVs assemblages.

408

409 3.5 Responses to environmental data

410 3.5.1 *Achnanthydium minutissimum* (ADMI)

411 Z scores obtained by TITAN analyses performed on ASVs of *Achnanthydium minutissimum* were
412 hierarchically clustered and visualized through a heatmap plot. Three main groups of ASVs
413 ADMI EG1, ADMI EG2 and ADMI EG3 (= *A. minutissimum* Ecological Groupings 1, 2 and 3) could
414 be distinguished on the basis of the magnitude (given by z score) and type (either positive or
415 negative) of their responses (Fig. 3; Supplementary Table 2). ADMI EG1 constituted a group
416 formed by 7 ASVs which, shared a positive response to altitude, calcium, conductivity, NH_4^+ ,
417 pH, SO_4^{2-} and a negative response to water temperature (Fig. 3; Supplementary Table 2). In
418 contrast to the positive response showed by ADMI EG1, the 7 ASVs that constituted ADMI EG3
419 group were characterized by an often negative response to altitude, calcium, conductivity,
420 NH_4^+ , pH, SO_4^{2-} and the response was especially strong for calcium and conductivity (Fig. 3;
421 Supplementary Table 2). Assemblage changes points to calcium and conductivity differed
422 between positive and negative responders (Supplementary Fig.10). Kruskal–Wallis and post-
423 hoc Dunn’s test indicated that the tree groupings were distributed at waters with significantly
424 different levels of calcium, conductivity, pH, NH_4^+ and SO_4^{2-} (Table 3).

425 BRT analyses indicated that calcium importantly influenced the occurrence of ADMI
426 EG3 and ADMI EG1 since for these groups it was the variable with the highest and second

427 highest relative importance respectively (Table 2). As in the TITAN analysis, partial dependence
428 plots generated by BRT models indicated a positive relationship of ADMI EG1 with both
429 calcium and conductivity but a negative relationship of ADMI EG3 with both variables (Fig. 4).
430 These plots showed that the response to calcium and conductivity largely increased in ADMI
431 EG1 group from 0 to 120 mg/L and from 0 to 700 $\mu\text{S}/\text{cm}$ respectively but decreased in ADMI
432 EG3 group from 35 to 55 mg/L and from 200 to 400 $\mu\text{S}/\text{cm}$ respectively (Fig. 4). BRT models
433 explained 47% and 44% of the total deviance and 30% and 27% of cross-validated deviance for
434 ADMI EG1 and ADMI EG3 groups respectively.

435 In contrast to the ADMI EG1 and ADMI EG3 groups, the ADMI EG2 group, formed by 18
436 ASVs was characterized by a positive response to altitude and a negative response to NO_3^- ,
437 NH_4^+ , PO_4^{3-} , TOC and water temperature (Fig. 3). The magnitude of response to altitude and
438 TOC was especially strong in some ASVs (Fig. 3; Supplementary Table 2).

439 BRT models indicated that altitude, conductivity, TOC and water temperature were the
440 four variables that most influenced the occurrence of ASVs in the ADMI EG2 group (Table 2).
441 Partial dependence plots showed a positive relationship of the grouping with altitude and a
442 negative with TOC and conductivity (Fig. 4). These plots depicted a large increase in the
443 response to altitude from 0 to 200 m and a decrease in the response to TOC from 1.5 to 5 mg/L
444 (Fig. 4). The BRT model based on the ADMI EG2 group explained 47% of deviance and 26% of
445 cross-validated deviance.

446

447 3.5.2 *Fistulifera saprophila* (FSAP)

448 According to the heatmap based on TITAN z scores obtained for ASVs of *F. saprophila*, three
449 ecological groupings were distinguished: FSAP EG1, FSAP EG2 and FSAP EG3 (Fig. 3). FSAP EG1
450 group was formed by 7 ASVs, most of them showing a positive response to calcium,
451 conductivity, NO_3^- , NH_4^+ , pH, SO_4^{2-} , PO_4^{3-} and TOC. Out of these variables, the strongest

452 responses (high z scores) in the group were to conductivity, NH_4^+ and PO_4^{3-} (Fig. 3;
453 Supplementary Table 2)

454 FSAP EG2 was comprised by 4 ASVs and they were characterized by a negative
455 response to altitude, calcium and conductivity and by a positive response to TOC and water
456 temperature (Fig. 3; Supplementary Table 2). In contrast to the responses shown by ASVs from
457 FSAP EG1 and FSAP EG2 groups, the ASVs from FSAP EG3 group were characterized by being
458 the only ASVs of *Fistulifera saprophila* that responded negatively to SO_4^{2-} , PO_4^{3-} and TOC (Fig.
459 3; Supplementary Table 2). With respect to SO_4^{2-} and TOC, assemblage change points differed
460 between positive and negative responders (Supplementary Fig. 11). Kruskal–Wallis and post-
461 hoc Dunn’s test indicated that FSAP EG3 ASVs were distributed in river sites with statistically
462 different values of TOC, PO_4^{3-} , NO_3^- and diatom indexes (i.e. IPS and IBD) (Table 3).

463 BRT analyses indicated that altitude and TOC importantly influenced the occurrence of
464 FSAP EG3, since these variables were respectively the first and second variables with the
465 highest relative importance. PO_4^{3-} was the most important variable in FSAP EG1 models but
466 altitude in FSAP EG2 models (Table 2). Partial dependence plots showed that the response of
467 FSAP EG1 and G3 groups to TOC decreased from 1 mg/L to 4-4.5 mg/L TOC. After this gradient,
468 the response of FSAP EG1 to TOC largely increased from 4.5 mg/L to 5 mg/L whereas there was
469 not any response for FSAP EG3 after 4 mg/L TOC (Fig. 4).

470 These plots reflected a positive relationship of FSAP EG1 and FSAP EG2 with SO_4^{2-} and
471 a negative one of FSAP EG3 with SO_4^{2-} . This was observed in the increasing response of both
472 FSAP EG1 and EG2 (though intermittently in the former case) along the gradient between 0 to
473 500 mg/L and in the large decreasing response of FSAP EG3 from 0 to 20 mg/L (Fig. 4). Partial
474 dependence plots also indicated a negative relationship of ASVs from FSAP EG2 with altitude,
475 since the plot depicted a large decrease in the response from 200 to 600m (Fig. 4). In contrast,
476 the response increased from 0 to 400 m for the FSAP EG1 group, while in the case of FSAP EG3,

477 the response increased from 0 to 600 m and partially and gradually decreased from 700 to
478 1030 m (Fig. 4). BRT models explained 53.4%, 40.8%, 41.7% of the deviance for FSAP EG1, FSAP
479 EG2 and FSAP EG3 respectively, and 37.2%, 24.3% and 21.6% of cross-validated deviance for
480 FSAP EG1, FSAP EG2 and FSAP EG3 respectively.

481

482 3.5.3 *Nitzschia* species

483 Based on TITAN analysis of *Nitzschia inconspicua* (NINC) and *N. soratensis* (NSTS), two
484 ecological groupings were defined (Fig. 5). All the 5 ASVs in the first group corresponded to *N.*
485 *inconspicua* species and they were characterized by a marked positive response to NO_3^- , NH_4^+ ,
486 SO_4^{2-} , PO_4^{3-} and TOC and by a very strong positive response to conductivity. The second group
487 comprised all three ASVs from *N. soratensis*, which, unlike the *N. inconspicua* ASVs, showed a
488 negative response to calcium, conductivity, pH, SO_4^{2-} , the responses to the first two being
489 especially strong (Fig. 5; Supplementary Table 2). Sum z scores for calcium differed between
490 ASVs from NINC and NSTS (Supplementary Fig. 12). A Kruskal–Wallis test showed that both
491 species were distributed in waters with significant differences levels of calcium, conductivity,
492 NH_4^+ , NO_3^- , pH, PO_4^{3-} , SO_4^{2-} and TOC (Table 3).

493 BRT models were performed separately for the group of ASVs from *N. inconspicua* and
494 the group of ASVs from *N. soratensis*. These models highlighted the importance of calcium for
495 explaining the distribution of ASVs from *N. soratensis*, since it was the variable with the highest
496 relative importance in the model (Table 2). In the case of the ASVs from *N. inconspicua*, the
497 two variables with the highest relative importance were conductivity and PO_4^{3-} respectively
498 (Table 2). Partial dependence plots (Fig. 6) indicated that the relationship of calcium with *N.*
499 *inconspicua* was positive but it was negative with *N. soratensis*. The models depicted an
500 increase in the response, though not continuously, from 10 mg/L to 150 mg/L for *N.*
501 *inconspicua* ASVs and a decrease from 50 to 70 mg/L for *N. soratensis* ASVs (Fig. 6). BRT

502 models explained 52.8% and 48.4% of deviance and 38.3% and 27.9% of cross-validated
503 deviance in *N. inconspicua* and *N. soratensis* ASVs respectively.

504

505 3.6. Relationship between phylogeny and geographical–ecological groupings

506 Phylogenetic trees of the *A. minutissimum* complex showed very little correlation between the
507 phylogeny and the ecological groupings (Fig. 2), although bootstrap for the tree nodes was
508 low. Five out of the seven ASVs that comprised the ecological grouping ADMI EG3 were placed
509 in the major Clade B and all the ASVs from the ADMI EG1 and ADMI EG2 groupings that passed
510 TITAN uncertainty criteria, except for ASV156 and ASV164, were classified into a second major
511 clade (Clade C). More specifically, all the ASVs from subclade d and all the ASVs except ASV
512 219 from the subclade h of the major clade C, belonged to the same ecological grouping, ADMI
513 EG2. However, some important exceptions showed that preferences are not always clade-
514 specific and must be determined at the ASV level: thus, ASVs 156 and 272 belong to the same
515 clade and differ by just two base-pairs, but belong to different ecological groupings (2 and 3,
516 respectively).

517 In the case of *F. saprophila* complex, the ASVs from the different ecological groupings
518 were scattered across the phylogenetic tree, without following any clear pattern
519 (supplementary Fig. 7).

520

521 **4. Discussion**

522 4.1 High diversity within species is captured by a short *rbcl* barcode

523 *RbcL* metabarcoding has been successfully applied for studying diatom species diversity (e.g.,
524 Rimet et al., 2018, Stoof-Leichsenring et al., 2020) and is especially useful for species that are
525 difficult to identify based on their morphological characteristics, such as those studied here –

526 *A. minutissimum*, *F. saprophila*, *Nitzschia inconspicua* and *N. soratensis*. An extra dimension is
527 given by the use of bioinformatics pipelines that generate amplicon sequence variants (ASVs)
528 as opposed to OTUs, since it is possible not only to identify species but also to detect and
529 quantify genetic diversity within them. Despite its short length, the 312-bp *rbcL* barcode we
530 used revealed substantial genetic diversity within the species studied, even when analysis was
531 restricted to the commoner ASVs, with ≥ 1000 reads and occurring in at least 2 samples with
532 environmental data. These comprised 45 ASVs identified as belonging to the *A. minutissimum*
533 complex and 18 of *F. saprophila*. However, it must be underlined that the total numbers of
534 ASVs obtained for these two species were much higher: 148 for *A. minutissimum* and 76 for *F.*
535 *saprophila* when ASVs having < 1000 reads and occurring in < 2 samples are also considered.

536 Interpretation of the low abundance ASVs is not straightforward, because both PCR
537 and Illumina sequencing generate errors. Despite the variety of quality and filtering steps
538 implemented in the various commonly used pipelines for HTS data analyses (Bailet et al.,
539 2020), it is impossible to be sure in all cases which ASVs are real though rare genetic variants
540 and which are artefactual. Clearly this can introduce a major bias in biodiversity studies (Turon
541 et al. 2019; Tsuji et al., 2019). A partial solution in the case of *rbcL*, if no matching Sanger
542 sequence is available, is to see whether the same ASV is present in different datasets
543 generated in different Illumina runs. Another is to assess each ASV by reference to the amino-
544 acids encoded: changes that are unlikely, based on amino-acid substitution matrices (e.g.
545 BLOSUM-62: Styczynski et al. 2008) can be tentatively discarded as artefactual. In this study,
546 the most common sequence of *A. minutissimum* that must be artefactual is ASV2237, the 72nd
547 most abundant sequence assigned to the species and represented by 114 reads; this contains a
548 stop codon and so cannot be functional. However, the least abundant *A. minutissimum* ASV
549 analysed (ASV6401), represented by just one read in the whole dataset, had an amino-acid
550 sequence identical to that of 8 of the 10 most abundant ASVs and cannot be discounted as an
551 error. These results illustrate, therefore, the importance of assessing the validity of sequences

552 even after denoising; it is dangerous to rely only on the abundances, since moderately
553 abundant sequences may nevertheless be artefacts. Conversely, rare sequences or even
554 singletons (i.e. sequences detected with only 1 read) are not necessarily artefacts but can be
555 reliable, as noted in other studies (e.g. Alberdi et al. 2017).

556 The reliability of DNA metabarcoding studies also depends on successful taxonomic
557 assignment of the sequences generated and for this it is important to choose an appropriate
558 confidence threshold. This issue has already been addressed in some studies (Rivera et al.
559 2020; Zizka et al., 2020) and in particular, for the short region of 312pb of the *rbcL* marker,
560 non-strict confidence thresholds have been demonstrated for benthic diatom biomonitoring
561 purposes (Rivera et al., 2020). We chose to set a similarity threshold of 50% (the default in
562 DADA2) in order to catch the maximum number of ASVs assigned to the studied species
563 because there is a risk of losing important ecological information when real ASVs are discarded
564 from a dataset, as has been shown in the taxonomy-free approach developed by Tapolczai et
565 al. (2021). In our dataset, although some of the ASVs' taxonomic assignments had low
566 bootstrap support (i.e. the percentage of times that the sequence was classified into the same
567 taxonomy was low), phylogenetic analyses that included curated reference sequences
568 indicated that all the abundant ASVs used in this study were properly classified into the
569 relevant species complex. Our results indicate that it is advisable to use a non-strict similarity
570 threshold to capture high diversity, provided that other analyses can guarantee the reliability
571 of the taxonomic assignment.

572

573 4.2. Wide geographical distributions of ASVs suggest dispersal is not a major constraint

574 The spatial structuring of ASVs suggested by MEMs analyses is congruent with the fact that
575 different ASVs have different geographical distributions, which ultimately could imply dispersal
576 constraints or different environmental preferences, or both. Although individual ASVs tended

577 to be abundant only in particular regions, in most cases the most abundant ASVs were
578 nevertheless found across more or less the whole region surveyed: only a few abundant ASVs
579 were restricted to one or other of France and Catalonia. Furthermore, in several cases the
580 ASVs matched Sanger-sequenced clones isolated from locations far from the study area, even
581 on different continents. It seems therefore that the ASVs of the species studied here are
582 dispersed quite effectively. Hence, when a ASV of the four species is *not* found in the France–
583 Catalonia dataset, there is a *prima facie* case that the appropriate environmental conditions do
584 not occur there, or at least, not in rivers. Examples are the Indian Ocean clade of *N.*
585 *inconspicua* (TCC clones 474, 510 and 571) and the tropical clade of *Fistulifera*. The species
586 considered here could therefore be argued to conform to the ubiquitous dispersal hypothesis
587 (e.g. Finlay, 2002), like some previous examples that have been sampled extensively, including
588 *Sellaphora capitata* (Evans and Mann 2009) and *S. bisexualis* (Mann et al., 2009), in which
589 identical or extremely similar haplotypes enjoy very wide ranges, despite evidence from
590 microsatellite data (in *S. capitata*) of genetic differentiation between populations separated by
591 only some 10s of km (Vanormelingen et al., 2015). In *N. palea* too, particular haplotypes have
592 extremely wide distributions (Trobajo et al., 2010), even though overall there is evidence of a
593 positive relationship between genetic and geographical distances (Rimet et al. 2014),
594 suggesting that dispersal is not fully effective in preventing genetic divergence.

595 Extra factors that need to be taken into account in interpreting the spatial structuring
596 observed in some ASVs are i) spatial structuring of key environmental variables and ii) the
597 possibility that important variables were not measured. Spatial structuring of the environment
598 was particularly obvious in the case of calcium, conductivity and sulphates, whose levels were
599 generally higher in Catalan rivers than French ones (Table 1). This could partly explain why
600 ASVs characterized by a strong positive response to calcium and conductivity often
601 predominated in Catalan rivers or were restricted there (e.g. ASV153; ASV219; ASVs from *N.*
602 *inconspicua*), whereas ASVs that showed a strong negative response were often better

603 represented in France (e.g. ASV269; ASVs from *N. soratensis*). Unmeasured environmental
604 parameters - such as substrate composition, dissolved oxygen, turbidity, water flow, channel
605 width or metals concentration – may also be influential (cf. Castro et al., 2019; Dalu et al.,
606 2017; Keck et al., 2018a) accounting for the low amount of variance explained by the RDA
607 model built from environmental data.

608 Overall, our results support the idea that individuals can disperse over long distances
609 while stochastic events of colonization and extinction possibly combined with fine scale
610 environmental variation are likely to generate local patchiness, outlining the importance of
611 considering spatial scale when studying diatom biogeographical patterns (Keck et al. 2018a).

612

613 4.3 Ecological preferences differ among ASVs in *A. minutissimum* and *F. saprophila*.

614 Our findings evidence the existence of different ecological preferences among different
615 populations and lineages of both *A. minutissimum* and *F. saprophila*, and importantly, that
616 these preferences are correlated with variations in the short *rbcl* barcode. Clearly, base
617 substitutions in *rbcl* within species (most of which do not in fact affect the amino-acid
618 composition and structure of RuBisCO) are unrelated to the causes of ecotypic differentiation
619 in the four diatom species studied; they are instead useful markers that can be used in
620 metabarcoding datasets to explore the existence and distributions of ecotypes.

621 In both species we found that two of the ecological groupings of ASVs were clearly
622 separated by their opposite responses to calcium and conductivity, while in the case of *F.*
623 *saprophila* a third ecological grouping (FSAP EG3) showed a preference for waters with low
624 organic pollution. It might be argued that the type of response shown by this grouping
625 corresponds better, within the genus *Fistulifera*, to *F. pelliculosa*, since this species is
626 considered to occur from oligo to mesotrophic habitats (Lange-Bertalot et al., 2017). The
627 morphology of FSAP EG3 cells is of course unknown. However, the two ASVs from this

628 grouping (i.e. ASV234 and ASV655) have probably been reliably assigned to *F. saprophila* since
629 phylogenetic analyses positioned these ASVs (which are not close relatives of each other)
630 within clades defined by curated reference sequences of *F. saprophila* (Supplementary Fig. 7).
631 We therefore treat the EG3 ASVs as belonging to *F. saprophila*. However, their ecological
632 preferences contrast starkly with the ecology often assumed for the species. Thus, Lange-
633 Bertalot et al. (2017) wrote that *F. saprophila* exhibits “large populations in heavily degraded,
634 highly eutrophic habitats with strong organic pollution up to polysaprobic conditions ... It is ...
635 one of the most pollution-tolerant diatoms.” A similar assessment was made by Gevrey et al.
636 (2004) and the IPS sensitivity value assigned by OMNIDIA (v5.5; Lecoq et al., 1993) is low
637 (IPSS=2). On the other hand, Lange-Bertalot et al. also noted that *F. saprophila* “can also be
638 found in moderately polluted water although in smaller numbers” and Zgrundo et al. (2013)
639 commented that the species is “a widely distributed taxon with broad ecological tolerances”.
640 Our data suggest that, if there is a ‘broad tolerance’, it may be because the species comprises
641 variants with contrasting requirements and tolerances, not because all *F. saprophila* can grow
642 across a wide range of water types. There are implications for biomonitoring, since the same
643 indicator values cannot be assigned to all the genetic varieties and metabarcoding assessments
644 should take this into account. The well-known tolerance of *F. saprophila* to a wide salinity
645 range, eutrophic conditions, and heavily degraded and organically polluted waters (Zgrundo et
646 al., 2013, Lange-Bertalot et al., 2017, Pniewski et al., 2010) must surely reflect the preferences
647 of the EG1 and EG2 groupings, not the EG3 ASVs. Moreover, the contrasting responses of the
648 EG1 and EG2 *rbcl* ASVs to conductivity suggest that the wide range of salinities recorded for
649 the species (e.g. Zgrundo et al., 2013) is also somewhat misleading, primarily reflecting
650 genotypic diversity rather than phenotypic plasticity.

651 In a lesser way, deviation from the ‘expected’ ecology was also observed in *A.*
652 *minutissimum*. Whereas one grouping of ASVs (ADMI EG3) was particularly restricted to low
653 nutrient concentrations (i.e. PO_4^{3-} , SO_4^{2-} , NH_4^+ and NO_3^-), as might be expected from the

654 characterization of *A. minutissimum* as an indicator of nutrient-poor, good quality waters (e.g.
655 Potapova and Charles, 2007, especially Appendix A), the other two groupings of ASVs tolerated
656 higher nutrient levels and would explain extension of the species complex into more nutrient-
657 rich waters, creating the impression of a broad ecological tolerance – hence the
658 characterization by Lange-Bertalot et al. (2017) “ecological amplitude apparently very wide”
659 (see also Potapova and Hamilton, 2007; Snoeijs and Balashova, 1998; Round, 2004). The idea
660 that *A. minutissimum* is a heterogeneous collection of lineages with different ecological
661 preferences is not new. For example, Potapova and Hamilton (2007) were able to distinguish
662 morphotypes within *A. minutissimum* and to associate them to some extent with different
663 preferences for conductivity, pH and nutrients. However, the morphological differences
664 between these variants (and between some of those documented by Pinseel et al., 2017), are
665 very subtle and distinguishing them in LM-based assessments is arguably impractical. The
666 metabarcoding approach not only aids identification but also allows vastly greater sampling of
667 *A. minutissimum* across natural communities.

668 Thus, our results for *A. minutissimum* and *F. saprophila* tell the same story, that while
669 overall the two species (i.e. all ASVs assigned to each of *A. minutissimum* and *F. saprophila*
670 taken together) have a very broad ecological tolerance, individual genetic variants (ASVs) do
671 not, and the perceived ecological preferences – and indicator value – of the species will differ
672 according to the types and relative abundances of the different ASVs present.

673

674 4.4 Ecological groupings of ASVs do not correspond well to phylogenetic groupings

675 The preferences we obtained for the ASVs are based on correlations between their relative
676 abundances in different samples and the environmental characteristics at the sites where the
677 samples were obtained, exactly as has been done previously with microscopical cell counts to
678 determine the preferences of morphologically defined species. These correlations likely reflect

679 adaptations of the ASVs to different ecological conditions and ASVs that are closely related
680 phylogenetically might be expected to share similar adaptations and belong to the same
681 ecological grouping (Keck et al., 2016a, b, 2018b). Overall, we did not find very strong evidence
682 of a correlation between ecological and phylogenetic group, though there were some trend
683 that could be observed in some cases. For instance, in *A. minutissimum*, the more distantly
684 related ASVs generally belonged to the groupings that differed most (i.e. ADMI EG1 and ADMI
685 EG3). And in *Fistulifera*, ASV74, which tolerated a high conductivity level (c. 9.000 $\mu\text{S}/\text{cm}$), was
686 closely related to a sequence (HQ337547) from clone CCMP543, isolated from a brackish pond
687 (in Massachusetts USA; this clone is often kept in fully marine medium), and clone TCC809,
688 isolated from the River Arão estuary in Portugal (Rimet et al., 2019). However, the *F.*
689 *saprophila* ASV recorded in the highest conductivity site in our dataset (c. 13.000 $\mu\text{S}/\text{cm}$) was
690 ASV445, which is not closely related to ASV74 and belongs to a clade whose other members
691 were recorded from freshwaters.

692

693 4.5. *Nitzschia inconspicua* and *N. soratensis* differ in their ecology but ASVs in each species
694 showed very similar preferences

695 Phylogenetic analyses show that *Nitzschia inconspicua* and *N. soratensis* are not close relatives
696 (Mann et al. 2021) but in the light microscope they are barely separable (Trobajo et al. 2013).
697 However, the value of differentiating between them in ecological and biomonitoring studies
698 has already been shown (Trobajo et al. 2013 and Kelly et al. 2015) and is further confirmed
699 here. Calcium and conductivity were the environmental parameters that most influenced the
700 occurrence of these species according to our data and the preference of *N. soratensis* for low
701 calcium and conductivity (see also Kelly et al. 2015) might explain why this species was
702 widespread in French rivers but scarcely detected in the Catalan ones.

703 In relation to ecological preferences, we found no differentiation between the ASVs in
704 *N. inconspicua* or *N. soratensis*, in contrast to *A. minutissimum* and *F. saprophila*. For
705 *inconspicua* this was surprising because Rovira et al. (2015) showed that this 'species' is
706 paraphyletic and comprises several very distantly related lineages. Furthermore, their
707 experimental work showed different salinity responses among *inconspicua* genotypes (Rovira
708 et al., 2015). However, the absence of the 'Indian Ocean' haplotypes from French and Catalan
709 rivers (section 3.3) may suggest ecological differentiation from the European ASVs and hence
710 that the structure of the *N. inconspicua* complex is not unlike that in *A. minutissimum* and *F.*
711 *saprophila*, containing populations adapted to different ecological conditions. This can only be
712 studied using molecular markers via a metabarcoding approach.

713

714 Conclusions

715 Our results show how intraspecific and cryptic diversity can be assessed and understood
716 through the application of DNA metabarcoding, leading to improvements in the knowledge of
717 dispersion patterns, phylogeny and ecological preferences of species and intraspecific variants
718 (see also De Luca et al., 2021; Rivera et al., 2018; Wattier et al., 2020; Zizka et al., 2020). This
719 approach is particularly appropriate for species or species complexes that are difficult to
720 distinguish on the basis of morphological characteristics and whose preferences are therefore
721 still not well-defined. There are many further examples in diatoms that would benefit greatly
722 from this approach, such as the *Cocconeis placentula* complex (Lange-Bertalot et al., 2017) and
723 *Planothidium* species (Jahn et al., 2017).

724 In relation to the questions we posed for this study, it is clear that genetic variants
725 within *Achnantheidium minutissimum* and *Fistulifera saprophila* are not distributed evenly
726 across the study area and it seems that this is at least partly due to differences in their
727 ecological preferences. Our data indicate that the broad ecological tolerances and wide

728 distributions claimed for some diatom species may well be the result of a continuum of
729 overlapping preferences among individual genetic variants, which can only be discriminated
730 using molecular markers. Importantly, however, there was little or no agreement between
731 ecological and phylogenetic groupings in *A. minutissimum* and *F. saprophila*, which shows that,
732 at least here, it is necessary to work at the lowest “taxonomic” level possible – ASVs – because
733 it cannot be assumed that clades of species and infraspecific variants share the same
734 ecological preferences and distributions.

735

736 **Acknowledgements**

737 We are very grateful to the Catalan Water Agency (ACA) for managing and organizing the river
738 survey and the following consultancies for taking DNA samples for us: Sorelló, Estudis del Medi
739 Aquàtic; CERM, Centre d'Estudis dels Rius Mediterranis -Universitat de Vic; GESNA Estudis
740 Ambientals; and Hidrologia i Qualitat de l'Aigua. We also thank the OFB (Office Français de la
741 Biodiversité), the French Water Agencies and the DREAL (Direction Régionale de
742 l'Environnement, de l'Aménagement et du Logement) who made possible the study in France;
743 and two anonymous reviewers for very constructive comments on the manuscript.

744 The authors also acknowledge support from the CERCA Programme/ Generalitat de
745 Catalunya. J. Pérez-Burillo acknowledges IRTA and Universitat Rovira i Virgili for his PhD grant
746 (2018PMF-PIPF-22). The Royal Botanic Garden Edinburgh is supported by the Scottish
747 Government's Rural and Environment Science and Analytical Services Division. This article was
748 also facilitated by COST Action DNAqua-Net (CA15219), supported by the COST (European
749 Cooperation in Science and Technology) program.

750

752 References

- 753 AFNOR (2007). NF T90-354. Qualité de l'eau - Détermination de l'Indice Biologique Diatomées
754 (IBD). AFNOR, 1–79.
- 755 Alberdi, A., Aizpurua, O., Gilbert, M., & Bohmann, K. (2017). Scrutinizing key steps for reliable
756 metabarcoding of environmental samples. *Methods in Ecology and Evolution*, 9(1),
757 134–147. <https://doi.org/10.1111/2041-210X.12849>
- 758 Armbrust, E.V. (2009). The life of diatoms in the world's oceans. *Nature* 459, 185–192.
759 <https://doi.org/10.1038/nature08057>.
- 760 Bailet, B., Apothéloz-Perret-Gentil, L., Baričević, A., Chonova, T., Franc, A., Frigerio, J.-M., Kelly,
761 M., Mora, D., Pfannkuchen, M., Proft, S., Ramon, M., Vasselon, V., Zimmermann, J. &
762 Kahlert, M. (2020). Diatom DNA metabarcoding for ecological assessment: Comparison
763 among bioinformatics pipelines used in six European countries reveals the need for
764 standardization. *Science of the Total Environment*, 745, 140948.
765 <https://doi.org/10.1016/j.scitotenv.2020.140948>.
- 766 Baker, M. E., & King, R. S. (2010). A new method for detecting and interpreting biodiversity and
767 ecological community thresholds. *Methods in Ecology and Evolution*, 1(1), 25–37.
768 <https://doi.org/10.1111/j.2041-210X.2009.00007.x>
- 769 Baker, M. E., King, R. S., & Kahle., D. (2019). TITAN2: Threshold Indicator Taxa Analysis. R
770 package, version 2.4. <https://CRAN.R-project.org/package=TITAN2>
- 771 Blanchet, F. G., Legendre, P., & Borcard, D. (2008). Forward selection of explanatory variables.
772 *Ecology*, 89(9), 2623–2632. <https://doi.org/10.1890/07-0986.1>
- 773 Callahan, B. J., McMurdie, P. J., Rosen, M. J., Han, A. W., Johnson, A. J. A., & Holmes, S. P.
774 (2016). DADA2: High resolution sample inference from Illumina amplicon data. *Nature*
775 *Methods*, 13, 581–583. <https://doi.org/10.1038/nmeth.3869>
- 776 Castro, E., Siqueira, T., Melo, A. S., Bini, L. M., Landeiro, V. L., & Schneck, F. (2019).
777 Compositional uniqueness of diatoms and insects in subtropical streams is weakly
778 correlated with riffle position and environmental uniqueness. *Hydrobiologia*, 842, 219–
779 232. <https://doi.org/10.1007/s10750-019-04037-8>
- 780 CEN (2014). CEN_EN 13946: Water Quality - Guidance for the Routine Sampling and
781 Preparation of Benthic Diatoms From Rivers and Lakes. pp. 1–22.
- 782 CEN (2018). CEN/TR 17245: Water Quality –Technical Report for the Routine Sampling of
783 Benthic Diatoms From Rivers and Lakes Adapted for Metabarcoding Analysis. CEN/ TC
784 230/WG23 – Aquatic Macrophyte and Algae. pp. 1–8.
- 785 Chonova, T., Keck, F., Bouchez, A., & Rimet, F. (2020). "A ready-to-use database for DADA2:
786 Diat.barcode_rbcL_263bp_DADA2 based on Diat.barcode v9". Portail Data INRAE, V2.
787 <https://doi.org/10.15454/QBLXP>
- 788 Dalu, T., Wasserman, R. J., Magoro, M. L., Mwedzi, T., Froneman, P. W., & Weyl, O. L. F. (2017).
789 Variation partitioning of benthic diatom community matrices: effects of multiple
790 variables on benthic diatom communities in an Austral temperate river system. *Science*
791 *of the Total Environment*, 601, 73–82. <https://doi.org/10.1016/j.scitotenv.2017.05.162>
- 792 Deiner, K., Bik, H. M., Mächler, E., Seymour, M., Lacoursière-Roussel, A., Altermatt, F., Creer,
793 S., Bista, I., Lodge, D. M., Vere, N., Pfrender, M. E., & Bernatchez, L. (2017).
794 Environmental DNA metabarcoding: Transforming how we survey animal and plant

795 communities. *Molecular Ecology*, 26(21), 5872–5895.
796 <https://doi.org/10.1111/mec.14350>

797 De Luca, D., Piredda, R., Sarno, D., & Kooistra, W. H. C. F. (2021). Resolving cryptic species
798 complexes in marine protists: phylogenetic haplotype networks meet global DNA
799 metabarcoding datasets. *ISME journal*. <https://doi.org/10.1038/s41396-021-00895-0>

800 Dinno, A. (2017). Dunn's Test of Multiple Comparisons Using Rank Sums. R package, version
801 1.3.5. <https://cran.r-project.org/web/packages/dunn.test/dunn.test.pdf>

802 Dray, S., Bauman, D., Blanchet, G., Borcard, D., Clappe, S., Guenard, G., Jombart, T., Larocque,
803 G., Legendre, P., Madi, N. & Wagner, H.H. (2020). Adespatial: Multivariate Multiscale
804 Spatial Analysis. R package, version 0.3-8. [https://CRAN.R-](https://CRAN.R-project.org/package=adespatial)
805 [project.org/package=adespatial](https://CRAN.R-project.org/package=adespatial)

806 Dray, S., Legendre, P., & Peres-Neto, P.R. (2006). Spatial modelling: a comprehensive frame-
807 work for principal coordinate analysis of neighbours matrices (PCNM). *Ecological*
808 *Modelling*, 196(3–4), 483–493. <https://doi.org/10.1016/j.ecolmodel.2006.02.015>

809 Dunn, O. J. (1964). Multiple comparisons using rank sums. *Technometrics*, 6(3), 241–252.

810 Edgar, R. C. (2004). MUSCLE: multiple sequence alignment with high accuracy and high
811 throughput. *Nucleic acids research*, 32(5), 1792–1797.
812 <https://doi.org/10.1093/nar/gkh340>

813 Elith, J., Leathwick, J. R., & Hastie, T. (2008). A working guide to boosted regression trees.
814 *Journal of Animal Ecology*, 77(4), 802–813. [https://doi.org/10.1111/j.1365-](https://doi.org/10.1111/j.1365-2656.2008.01390.x)
815 [2656.2008.01390.x](https://doi.org/10.1111/j.1365-2656.2008.01390.x)

816 European Commission (2016). Introduction to the new EU Water Framework Directive.
817 Available at [https://ec.europa.eu/environment/water/water-](https://ec.europa.eu/environment/water/water-framework/info/intro_en.htm)
818 [framework/info/intro_en.htm](https://ec.europa.eu/environment/water/water-framework/info/intro_en.htm)

819 Evans, K.M. & Mann, D.G. (2009). A proposed protocol for nomenclaturally effective DNA
820 barcoding of microalgae. *Phycologia* 48, 70–74. <https://doi.org/10.2216/08-70.1>

821 Finlay, B. J., Monaghan, E. B., & Maberly, S. C. (2002). Hypothesis: The Rate and Scale of
822 Dispersal of Freshwater Diatom Species is a Function of their Global Abundance.
823 *Protist*, 153(3), 261–273. <https://doi.org/10.1078/1434-4610-00103>

824 Friedman, J.H. (2001). Greedy function approximation: a gradient boosting machine. *Annals of*
825 *Statistics*, 29(5), 1189–1232. <https://doi.org/10.1214/aos/1013203451>

826 Gevrey, M., Rimet, F., Park, Y. S., Giraudel, J.-L., Ector, L., & Lek, S. (2004). Water quality
827 assessment using diatom assemblages and advanced modelling techniques.
828 *Freshwater Biology*, 49(2), 208–220. [https://doi.org/10.1046/j.1365-](https://doi.org/10.1046/j.1365-2426.2003.01174.x)
829 [2426.2003.01174.x](https://doi.org/10.1046/j.1365-2426.2003.01174.x)

830 Hijmans, R.J., Phillips, S., Leathwick, J., & Elith, J. (2020). dismo: Species Distribution Modeling.
831 R package version 1.3-3. <https://CRAN.R-project.org/package=dismo>

832 Hollander, M., & Wolfe, D.A. (1973). *Nonparametric statistical methods*, 2nd ed. New York, NY,
833 USA, Wiley.

834 Jahn, R., Abarca, N., Gemeinholzer, B., Mora, D., Skibbe, O., Kulikovskiy, M., Gusev, E., Kusber,
835 W.H., Zimmermann, J. (2017). *Planothidium lanceolatum* and *Planothidium*
836 *frequentissimum* reinvestigated with molecular methods and morphology: four new
837 species and the taxonomic importance of the sinus and cavum. *Diatom Research*,
838 32(1), 75–107. <https://doi.org/10.1080/0269249X.2017.1312548>.

- 839 Keck, F., Bouchez, A., Franc, A., & Rimet, F. (2016a). Linking phylogenetic similarity and
840 pollution sensitivity to develop ecological assessment methods: a test with river
841 diatoms. *Journal of Applied Ecology*, 53(3), 856–864. [https://doi.org/10.1111/1365-
842 2664.12624](https://doi.org/10.1111/1365-2664.12624)
- 843 Keck, F., Rimet, F., Franc, A., & Bouchez, A. (2016b). Phylogenetic signal in diatom ecology:
844 perspectives for aquatic ecosystems biomonitoring. *Ecological Applications*, 26(3),
845 861–872. <https://doi.org/10.1890/14-1966>
- 846 Keck, F., Franc, A., & Kahlert, M. (2018a). Disentangling the processes driving the
847 biogeography of freshwater diatoms: a multiscale approach. *Journal of Biogeography*,
848 45, 1582–1592. <https://doi.org/10.1111/jbi.13239>.
- 849 Keck, F., Vasselon, V., Rimet, F., Bouchez, A., & Kahlert, M. (2018b). Boosting DNA
850 metabarcoding for biomonitoring with phylogenetic estimation of operational
851 taxonomic units' ecological profiles. *Molecular Ecology Resources*, 18(6), 1299–1309.
852 <https://doi.org/10.1111/1755-0998.12919>
- 853 Kelly, M. G., Juggins, S., Mann, D. G., Sato, S., Glover, R., Boonham, N., Sapp, M., Lewis, E.,
854 Hany, U., Kille, P., Jones, T. & Walsh, K. (2020). Development of a novel metric for
855 evaluating diatom assemblages in rivers using DNA metabarcoding. *Ecological
856 Indicators*, 118, 106725. <https://doi.org/10.1016/j.ecolind.2020.106725>
- 857 Kelly, M. G., Trobajo, R., Rovira, L., & Mann, D. G. (2015). Characterizing the niches of two very
858 similar *Nitzschia* species and implications for ecological assessment, *Diatom Research*,
859 30(1), 27–33. <https://doi.org/10.1080/0269249X.2014.951398>
- 860 Kumar, S., Stecher, G., Li, M., Nnyaz, C., & Tamura, K. (2018). MEGA X: Molecular Evolutionary
861 Genetics Analysis across computing platforms. *Molecular Biology and Evolution*, 35(6),
862 1547–1549. <https://doi.org/10.1093/molbev/msy096>
- 863 Lange-Bertalot, H., Hofmann, G., Werum, M., & Cantonati, M. (2017). *Freshwater benthic
864 diatoms of Central Europe: over 800 common species used in ecological assessment.
865 English Edition with updated taxonomy and added species.* English Edition with
866 updated taxonomy and added species. Schmitten-Oberreifenberg, Koeltz Botanical
867 Books.
- 868 Lanzén, A., Mendibil, I., Borja, Á., & Alonso-Sáez, L. (2020). A microbial *mandala* for
869 environmental monitoring: Predicting multiple impacts on estuarine prokaryote
870 communities of the Bay of Biscay. *Molecular Ecology*, 1– 19.
871 <https://doi.org/10.1111/mec.15489>
- 872 Lecointe, C., Coste, M., & Prygiel, J. (1993). OMNIDIA—software for taxonomy, calculation of
873 diatom indexes and inventories management. *Hydrobiologia*, 269, 509–513.
874 <https://doi.org/10.1007/BF00028048>
- 875 Legendre, P., & Legendre, L. (2012). Chapter thirteen – Spatial Analysis. *Numerical ecology*, 3rd
876 ed. (pp. 785–858). Amsterdam, Elsevier Science BV.
- 877 Letunic I., & Bork P. (2019). Interactive Tree Of Life (iTOL) v4: recent updates and new
878 developments. *Nucleic Acids Research*, 47(W1), W256–W259.
879 <https://doi.org/10.1093/nar/gkz239>.
- 880 Malviya, S., Scalco, E., Audic, S., Vincent, F., Veluchamy, A., Poulain, J., Wincker, P., Iudicone,
881 D., de Vargas, C., Bittner, L., Zingone, A., & Bowler, C. (2016). Insights into global
882 diatom distribution and diversity in the world's ocean. *Proceedings of the National*

883 *Academy of Sciences of the USA*, 113(111), E1516–E1525.
884 <https://doi.org/10.1073/pnas.1509523113>

885 Mann, D.G. (1999). The species concept in diatoms. *Phycologia*, 38(6), 437–495.
886 <https://doi.org/10.2216/i0031-8884-38-6-437.1>

887 Mann, D. G., Evans, K. M., Chepurinov, V. A., & Nagai, S. (2009). Morphology and formal
888 description of *Sellaphora bisexualis*, sp. nov. (Bacillariophyta). *Fottea*, 9(2), 199–209.
889 <https://doi.org/10.5507/fot.2009.021>

890 Mann, D. G., Trobajo, R., Sato, S., Li, C., Witkowski, A., Rimet, F., Ashworth, M. P., Hollands, R.
891 M., & Theriot, E. C. (2021). Ripe for reassessment: A synthesis of available molecular
892 data for the speciose diatom family Bacillariaceae. *Molecular Phylogenetics and*
893 *Evolution*, 158, 106985. <https://doi.org/10.1016/j.ympev.2020.106985>

894 Martin, M. (2011) Cutadapt removes adapter sequences from high-throughput sequencing
895 reads. *EMBnet. Journal*, 17, 10–12.

896 Mortágua, A., Vasselon, V., Oliveira, R., Elias, C., Chardon, C., Bouchez, A., Rimet, F., João Feio,
897 M., & Almeida, S. F. (2019). Applicability of DNA metabarcoding approach in the
898 bioassessment of Portuguese rivers using diatoms. *Ecological Indicators*, 106, 105470.
899 <https://doi.org/10.1016/j.ecolind.2019.105470>.

900 Oksanen, J., Guillaume Blanchet, F., Friendly, M., Kindt, R., Legendre, P., McGlinn, D., Minchin,
901 P. R., O'Hara, B., Simpson, G. L., Solymos, P., Stevens, M. H. H., Szoecs, E., & Wagner, H.
902 (2020). Vegan: Community Ecology Package. R package, version 2.5-7. [https://CRAN.R-](https://CRAN.R-project.org/package=vegan)
903 [project.org/package=vegan](https://CRAN.R-project.org/package=vegan)

904 Pérez-Burillo, J., Trobajo, R., Vasselon, V., Rimet, F., Bouchez, A., & Mann, D. G. (2020).
905 Evaluation and sensitivity analysis of diatom DNA metabarcoding for WFD
906 bioassessment of Mediterranean rivers. *Science of the Total Environment*, 727, 138445.
907 <https://doi.org/10.1016/j.scitotenv.2020.138445>

908 Pinseel, E., Vanormelingen, P., Hamilton, P. B., Vyverman, W., Van de Vijver, B., & Kopalova, K.
909 (2017). Molecular and morphological characterization of the *Achnantheidium*
910 *minutissimum* complex (Bacillariophyta) in Petuniabukta (Spitsbergen, High Arctic)
911 including the description of *A. digitatum* sp. nov. *European Journal of Phycology*, 52(3),
912 264–280. <https://doi.org/10.1080/09670262.2017.1283540>

913 Piredda, R., Claverie, J.-M., Decelle, J., De Vargas, C., Dunthorn, M., Edvardsen, B., Eikrem, W.,
914 Forster, D., Kooistra, W. H. C. F., Logares, R., Massana, R., Montresor, M., Not, F.,
915 Ogata, H., Pawlowski, J., Romac, S., Sarno, D., Stoeck, T., & Zingone, A. (2018). Diatom
916 diversity through HTS-metabarcoding in coastal European seas. *Scientific Reports*, 8,
917 18059. <https://doi.org/10.1038/s41598-018-36345-9>

918 Pniewski, F., Friedl, T., & Latała, A. (2010). Identification of diatom isolates from the Gulf of
919 Gdańsk: testing of species identifications using morphology, 18S rDNA sequencing and
920 DNA barcodes of strains from the Culture Collection of Baltic Algae (CCBA).
921 *Oceanological and Hydrobiological Studies*, 39(3), 3-20.
922 <https://doi.org/10.2478/v10009-010-0031-7>

923 Potapova, M., & Charles, D. F. (2007). Diatom metrics for monitoring eutrophication in rivers of
924 the United States. *Ecological Indicators*, 7(1), 48–70.
925 <https://doi.org/10.1016/j.ecolind.2005.10.001>

- 926 Potapova, M., & Hamilton, P. B. (2007). Morphological and ecological variation within the
927 *Achnantheidium minutissimum* (Bacillariophyceae) species complex. *Journal of*
928 *Phycology*, 43(3), 561–575. <https://doi.org/10.1111/j.1529-8817.2007.00332.x>
- 929 Poulíčková, A., Letáková, M., Hašler, P., Cox, E., & Duchoslav, M. (2017). Species complexes
930 within epiphytic diatoms and their relevance for the bioindication of trophic status.
931 *Science of the Total Environment*, 599-600, 820–833.
932 <https://doi.org/10.1016/j.scitotenv.2017.05.034>
- 933 Poulíčková, A., Špačková, J., Kelly, M. G., Duchoslav, M., & Mann, D. G. (2008). Ecological
934 variation within *Sellaphora* species complexes (Bacillariophyceae): specialists or
935 generalists? *Hydrobiologia*, 614, 373–386. <https://doi.org/10.1007/s10750-008-9521-y>
- 936 R Core Team, R. (2020). *A language and environment for statistical computing*. Vienna, Austria:
937 *R Foundation for Statistical Computing*. Retrieved from <https://www.R-project.org/>
- 938 Rimet, F., Gusev, E., Kahlert, M., Kelly, M. G., Kulikovskiy, M., Maltsev, Y., Mann, D. G.,
939 Pfannkuchen, M., Trobajo, R., Vasselon, V., Zimmermann, J., & Bouchez, A. (2019).
940 Diat.barcode, an open-access curated barcode library for diatoms. *Scientific Reports*, 9,
941 15116, 1–12. <https://doi.org/10.1038/s41598-019-51500-6>.
- 942 Rimet, F., Trobajo, R., Mann, D. G., Kermarrec, L., Franc, A., Domaizon, I., & Bouchez, A. (2014).
943 When is sampling complete? The effects of geographical range and marker choice on
944 perceived diversity in *Nitzschia palea* (Bacillariophyta). *Protist*, 165(3), 245–259.
945 <https://doi.org/10.1016/j.protis.2014.03.005>
- 946 Rimet, F., Vasselon, V., A.-Keszte, B., & Bouchez, A. (2018). Do we similarly assess diversity with
947 microscopy and high-throughput sequencing? Case of microalgae in lakes. *Organism*
948 *Diversity & Evolution*, 18, 51–62. <https://doi.org/10.1007/s13127-018-0359-5>.
- 949 Rivera, S. F., Vasselon, V., Ballorain, K., Carpentier, A., Wetzel, C. E., Ector, L., Bouchez, A., &
950 Rimet, F. (2018). DNA metabarcoding and microscopic analyses of sea turtles biofilms:
951 Complementary to understand turtle behavior. *PLoS ONE*, 13(4): e0195770.
952 <https://doi.org/10.1371/journal.pone.0195770>
- 953 Rivera, S. F., Vasselon, V., Bouchez, A., & Rimet, F. (2020). Diatom metabarcoding applied to
954 large scale monitoring networks: Optimization of bioinformatics strategies using
955 Mothur software. *Ecological Indicators*, 109, 105775.
956 <https://doi.org/10.1016/j.ecolind.2019.105775>
- 957 Round, F. E. (2004). pH scaling and diatom distribution. *Diatom*, 20, 9–12.
958 https://doi.org/10.11464/diatom1985.20.0_9
- 959 Rovira, L., Trobajo, R., Sato, S., Ibáñez, C., & Mann, D. G. (2015). Genetic and physiological
960 diversity in the diatom *Nitzschia inconspicua*. *Journal of Eukaryotic Microbiology*, 62(6),
961 815–832. <https://doi.org/10.1111/jeu.12240>
- 962 Ryneerson, T. A., Newton, J. A., & Armbrust, E. V. (2006). Spring bloom development, genetic
963 variation, and population succession in the planktonic diatom *Ditylum brightwelli*.
964 *Limnology and Oceanography*. 51(3), 1249–1261.
965 <https://doi.org/10.4319/lo.2006.51.3.1249>
- 966 Snoeijs, P., & Balashova, N. (1998). *Intercalibration and distribution of diatom species in the*
967 *Baltic Sea*. Opulus Press, Uppsala.
- 968 Smucker, N. J., Pilgrim, E. M., Nietch, C. T., Darling, J. A., & Johnson B. R. (2020). DNA
969 metabarcoding effectively quantifies diatom responses to nutrients in streams.
970 *Ecological Applications*, 30(8), e02205. <https://doi.org/10.1002/eap.2205>

- 971 Smetacek, V. (1999). Diatoms and the ocean carbon cycle. *Protist*, 150(1), 25–32.
 972 [https://doi.org/10.1016/S1434-4610\(99\)70006-4](https://doi.org/10.1016/S1434-4610(99)70006-4)
- 973 Soininen, J., Jamoneau, A., Rosebery, J., Leboucher, T., Wang, J., Kokociński, M., & Passy, S. I.
 974 (2018). Stream diatoms exhibit weak niche conservation along global environmental
 975 and climatic gradients. *Ecography*, 42(2), 346–353.
 976 <https://doi.org/10.1111/ecog.03828>
- 977 Souffreau, C., Vanormelingen, P., Van de Vijver, B., Isheva, T., Verleyen, E., Sabbe, K. &
 978 Vyverman, W. (2013). Molecular evidence for distinct antarctic lineages in the
 979 cosmopolitan terrestrial diatoms *Pinnularia borealis* and *Hantzschia amphioxys*.
 980 *Protist*, 164(1), 101–115. <https://doi.org/10.1016/j.protis.2012.04.001>
- 981 Stoof-Leichsenring, K. R., Pestryakova, L. A., Epp, L. S., Herzsich, U. (2020). Phylogenetic
 982 diversity and environment form assembly rules for Arctic diatom genera—A study on
 983 recent and ancient sedimentary DNA. *Journal of Biogeography*. 47(5), 1166–1179.
 984 <https://doi.org/10.1111/jbi.13786>
- 985 Styczynski, M. P., Jensen, K. L., Rigoutsos, I., & Stephanopoulos, G. (2008). BLOSUM62
 986 miscalculations improve search performance. *Nature Biotechnology*, 26(3), 274–275.
 987 <https://doi.org/10.1038/nbt0308-274>
- 988 Tapolczai, K., Selmečzy, G. G., Szabó, B., B-Béres, V., Keck, F., Bouchez, A., Rimet, F., & Padisák,
 989 J. (2021). The potential of exact sequence variants (ESVs) to interpret and assess the
 990 impact of agricultural pressure on stream diatom assemblages revealed by DNA
 991 metabarcoding. *Ecological Indicators*, 122, 107322.
 992 <https://doi.org/10.1016/j.ecolind.2020.107322>.
- 993 Trobajo, R., Mann, D. G., Clavero, E., Evans, K. M., Vanormelingen, P., & McGregor, R. C.
 994 (2010). The use of partial *cox1*, *rbcL* and LSU rDNA sequences for phylogenetics and
 995 species identification within the *Nitzschia palea* complex (Bacillariophyceae). *European*
 996 *Journal of Phycology*, 45(4), 413–425. <https://doi.org/10.1080/09670262.2010.498586>
- 997 Trobajo, R., Rovira, L., Ector, L., Wetzel, C. E., Kelly, M., & Mann, D. G. (2013). Morphology and
 998 identity of some ecologically important small *Nitzschia* species. *Diatom Research*,
 999 28(1), 37–59. <https://doi.org/10.1080/0269249X.2012.734531>
- 1000 Tsuji, S., Miya, M., Ushio, M., Sato, H., Minamoto, T., & Yamanaka, H. (2019). Evaluating
 1001 intraspecific genetic diversity using environmental DNA and denoising approach: A
 1002 case study using tank water. *Environmental DNA*, 2(1), 42–52.
 1003 <https://doi.org/10.1002/edn3.44>
- 1004 Turon, X., Antich, A., Palacín, C., Præbel, K., & Wangensteen, O. S. (2019). From metabarcoding
 1005 to metaphylogeography: separating the wheat from the chaff. *Ecological Applications*,
 1006 30(2), e02036. <https://doi.org/10.1002/eap.2036>
- 1007 Vanormelingen, P., Evans, K. M., Mann, D. G., Lance, S., Debeer, A.-E., D’Hondt, S., Verstraete,
 1008 T., De Meester, L., & Vyverman, W. (2015). Genotypic diversity and differentiation
 1009 among populations of two benthic freshwater diatoms as revealed by microsatellites.
 1010 *Molecular Ecology*, 24(17), 4433–4448. <https://doi.org/10.1111/mec.13336>
- 1011 Vasselon, V., Domaizon, I., Rimet, F., Kahlert, M., Bouchez, A. (2017a). Application of
 1012 highthroughput sequencing (HTS) metabarcoding to diatom biomonitoring: do DNA
 1013 extraction methods matter? *Freshwater Science*, 36, 162–177.
 1014 <https://doi.org/10.1086/690649>

- 1015 Vasselon, V., Rimet, F., Tapolczai, K., & Bouchez, A. (2017b). Assessing ecological status with
 1016 diatoms DNA metabarcoding: scaling-up on a WFD monitoring network (Mayotte
 1017 island, France). *Ecological Indicators*, 82, 1–12.
 1018 <https://doi.org/10.1016/j.ecolind.2017.06.024>.
- 1019
- 1020 Wagenhoff, A., Liess, A., Pastor, A., Clapcott, J. E., Goodwin, E. O., & Young, R. G. (2017).
 1021 Thresholds in ecosystem structural and functional responses to agricultural stressors
 1022 can inform limit setting in streams. *Freshwater Science*, 36(1), 178–194.
 1023 <https://doi.org/10.1086/690233>
- 1024 Wang, Q., Garrity, G. M., Tiedje, J. M., & Cole, J. R. (2007). Naïve Bayesian classifier for rapid
 1025 assignment of rRNA sequences into the new bacterial taxonomy. *Applied and
 1026 Environmental Microbiology*, 73(16), 5261–5267. [https://doi.org/10.1128/AEM.00062-](https://doi.org/10.1128/AEM.00062-07)
 1027 [07](https://doi.org/10.1128/AEM.00062-07).
- 1028 Ward, J.H. (1963). Hierarchical Grouping to Optimize an Objective Function. *Journal of the
 1029 American Statistical Association*, 58, 236–244. [10.1080/01621459.1963.10500845](https://doi.org/10.1080/01621459.1963.10500845).
- 1030 Warnes, G. R., Bolker, B., Bonebakker, L., Gentleman, R., Huber, W., Liaw, A., Lumley, T.,
 1031 Maechler, M., Magnusson, A., Moeller, S., Schwartz, M., & Venables, B. (2020). Gplots:
 1032 Various R Programming Tools for Plotting Data. R package, version 3.1.1.
 1033 <https://CRAN.R-project.org/package=gplots>
- 1034 Wattier, R., Mamos, T., Copilaş-Ciocianu, D., Jelić, M., Ollivier, A., Chaumot, A., Danger, M.,
 1035 Felten, V., Piscart, C., Žganec, K., Rewicz, T., Wysocka, A., Rigaud, T., & Grabowski, M.
 1036 (2020). Continental-scale patterns of hyper-cryptic diversity within the freshwater
 1037 model taxon *Gammarus fossarum* (Crustacea, Amphipoda). *Scientific Reports*, 10,
 1038 16536. <https://doi.org/10.1038/s41598-020-73739-0>
- 1039 Zgrundo, A., Lemke, P., Pniewski, F., Cox, E.J., & Latala, A. (2013). Morphological and molecular
 1040 phylogenetic studies on *Fistulifera saprophila*. *Diatom Research*, 28(4), 431–443.
 1041 <https://doi.org/10.1080/0269249X.2013.833136>
- 1042 Zizka, V. M. A., Weiss, M., & Leese, F. (2020). Can metabarcoding resolve intraspecific genetic
 1043 diversity changes to environmental stressors? A test case using river
 1044 macrozoobenthos. *Metabarcoding and Metagenomics*, 4, e51925.
 1045 <https://doi.org/10.3897/mbmg.4.51925>
- 1046

1047 **Figures caption**

1048

1049 Fig 1. Spatial distribution of the 10 most abundant ASVs from *Achnanthydium minutissimum* in
1050 French and Catalan rivers. Segments in each circle represent the proportion of *A.*
1051 *minutissimum* reads recorded in each sample site.

1052

1053 Fig 2. Maximum likelihood phylogenetic tree of *Achnanthydium minutissimum* ASVs obtained in
1054 this study and related reference sequences extracted from Diat.barcode v9 and GenBank. The
1055 tree was obtained using raxmlGUI and a GRT-Gamma model with 1000 replicates for the
1056 bootstrap analyses. The tree was drawn using iTOL. ASVs belonging to the different ecological
1057 groupings defined after TITAN analyses are represented: EG1 in red, EG2 in green and EG3 in
1058 blue. Black circles represent bootstrap support values of 50-100. Three major clades can be
1059 distinguished (A, B and C) and a number of subclades (Ca to Cl).

1060

1061 Fig 3. Heatmap dendrogram based on z score obtained by the different TITAN analyses
1062 performed on ASVs from a) *Achnanthydium minutissimum* and b) *Fistulifera saprophila*.
1063 Euclidean distance and ward.D functions were used to compute dissimilarity distance and
1064 hierarchical clustering respectively on ASVs z scores obtained for the different environmental
1065 variables. Only those ASVs with more than 3 responses that fulfilled purity and reliability
1066 criteria are represented. Red colour indicates positive responses while blue negative
1067 responses. Magnitude of response (z score) are given by the contrast of the colour; dark
1068 colours depict strong responses while light colours indicate weak responses. Chart in the
1069 upper-left corner indicates the correspondence between colour gradient and z-score.

1070

1071 Fig 4. Partial dependence plots generated by boosted regression trees analyses depicting the
1072 response of the ecological groupings of ASVs from *Achnanthydium minutissimum* to altitude
1073 (m), Calcium (mg/L), conductivity ($\mu\text{S}/\text{cm}$) and Total Organic Carbon (TOC, mg/L) and ecological
1074 groupings of ASVs from *Fistulifera saprophila* to Total Organic Carbon (TOC, mg/L), sulphates
1075 (mg/L) and altitude (m). The different groups of ASVs were defined after TITAN analyses. Y axis
1076 shows fitted function.

1077

1078 Fig 5. Heatmap based on z score obtained by the different TITAN analyses performed on ASVs
1079 from *Nitzschia inconspicua* (NINC) and *N. soratensis* (NSTS). Only those ASVs with more than 3
1080 responses that fulfilled purity and reliability criteria are represented. Red colour indicates
1081 positive responses while blue negative responses. Magnitude of response (z score) are given
1082 by the contrast of the colour; dark colours depict strong responses while light colours indicate
1083 weak responses. Chart in the upper-left corner indicates the correspondence between colour
1084 gradient and z-score.

1085

1086 Fig 6. Partial dependence plots generated by boosted regression trees analyses depicting the
1087 response of ASVs from *Nitzschia inconspicua* and *Nitzschia soratensis* to Calcium (mg/L). Y axis
1088 represents fitted function.

1089

1090

1091

1092

Tables

Table 1. Physicochemical parameters information from the 531 river sites studied.

Variable	Number of sampling sites with available data	Number of sampling sites with available data (Catalan rivers)	Number of sampling sites with available data (French rivers)	Average \pm standard deviation of number of records per sampling site within the 90-day period (Catalan rivers)	Average \pm standard deviation of number of records per sampling site within the 90-day period (French rivers)	Range (average \pm standard deviation) in Catalan rivers	Range (average \pm standard deviation) in French rivers
Altitude (m)	531	148	383	NA	NA	3.89 - 1243.97 (303.75 \pm 255.29)	0 - 1933 (255.7 \pm 314.39)
Ammonium (mg/L)	513	136	377	1 \pm 0.08	2.61 \pm 2.03	0.1 - 15.33 (0.69 \pm 2.21)	0.004 - 1.4 (0.07 \pm 0.13)
Bicarbonates (mg/L)	200	35	165	1 \pm 0.08	2.09 \pm 1.25	25 - 182 (54.18 \pm 39.11)	6.4 - 600 (184.37 \pm 95.8)
Calcium (mg/L)	335	148	187	1 \pm 0.08	2.3 \pm 1.81	2.5 - 673.33 (116.05 \pm 93.83)	0.7 - 333 (63.66 \pm 43.35)
Conductivity (μ S/cm)	336	136	200	1 \pm 0.08	3.85 \pm 3.25	99.5 - 13341.33 (1054.61 \pm 1382.76)	25.67 - 2377.67 (341.23 \pm 283.14)
Total organic carbon (mg/L)	514	136	378	1 \pm 0.08	2.69 \pm 2.33	0.5 - 10.65 (3.47 \pm 1.85)	0.2 - 15 (2.46 \pm 1.67)
Nitrates (mg/L)	502	123	379	1 \pm 0.08	2.63 \pm 2.12	2.5 - 76.45 (13.94 \pm 13.55)	0.48 - 47.27 (7.49 \pm 7.22)
Orthophosphates (mg/L)	515	136	379	1 \pm 0.08	2.57 \pm 1.90	0.1 - 9.73 (0.58 \pm 1.09)	0.01 - 2.53 (0.16 \pm 0.25)
pH	336	136	200	1 \pm 0.08	3.84 \pm 3.25	7.65 - 8.8 (8.19 \pm 0.23)	6.3 - 8.6 (7.83 \pm 0.42)
Sulphates (mg/L)	301	136	165	1 \pm 0.08	2.09 \pm 1.25	4 - 1500 (178.69 \pm 217.83)	1 - 416 (39.92 \pm 57.1)
Water temperature ($^{\circ}$ C)	330	130	200	1 \pm 0.00	4.26 \pm 6.44	5 - 26 (12.45 \pm 3.17)	7.13 - 24.5 (18.12 \pm 3.99)

Table 2. Relative importance (%) of each environmental variable resulting from the boosted regression tree models (with 10-fold cross validation of data) performed for the different groups of ASVs of *Achnanthydium minutissimum* (ADMI), *Fistulifera saprophila* (FSAP), *Nitzschia inconspicua* (NINC) and *Nitzschia soratensis* (NSTS). Groups of ASVs were defined on the basis of TITAN analyses.

Variable	ADMI G1	ADMI G2	ADMI G3	FSAP G1	FSAP G2	FSAP G3	NINC	NSTS
Orthophosphates	17.90	8.35	9.96	23.44	8.81	2.93	19.69	14.14
Calcium	13.81	5.01	21.75	8.57	5.20	3.89	4.77	20.49
Conductivity	11.89	11.79	7.43	16.28	4.38	3.83	27.49	14.14
Nitrates	8.50	6.41	13.35	8.57	7.98	7.93	13.68	12.35
Altitude	10.09	13.94	9.43	5.30	32.33	29.66	7.00	7.38
pH	8.38	10.03	1.52	2.98	4.66	1.72	2.87	6.85
Water temperature	10.79	22.27	6.88	6.25	11.22	9.14	5.04	3.77
TOC	5.90	12.80	12.12	3.52	7.34	20.79	4.28	8.98
Bicarbonates	7.38	2.26	3.63	6.05	8.05	6.14	1.53	5.24
Ammonium	2.93	4.72	10.89	11.43	7.24	5.89	3.32	3.36
Sulphates	2.45	2.41	3.03	5.46	2.75	8.24	10.28	3.24

Table 3. Range, average and standard deviation environmental parameters analysed in the sites were different defined ecological groupings occurred. ^a and ^b indicate species and ecological groupings with statistically significant differences (Kruskal–Wallis for *Nitzschia inconspicua* and *N. soratensis* and post-hoc Dunn’s test for groupings from *Achnanthydium minutissimum* and *Fistulifera saprophila*, $p < 0.05$).

Variable	ADMI G1	ADMI G2	ADMI G3	FSAP G1	FSAP G2	FSAP G3	NINC	NSTS
Orthophosphates	^a 0.01-3.35 (0.23±0.39)	^b 0.01- 4.1(0.22±0.40)	^{ab} 0.01- 2(0.13±0.22)	^a 0.01-9.73(0.50±0.95)	^a 0.01-4.3(0.28±0.54)	^a 0.01-4.3(0.22±0.53)	^a 0.01-9.73(0.48±0.24)	^a 0.01-4.10(0.24±0.42)
Calcium	^a 8.25-605 (108.09±76.54)	^a 1.90-477 (86.02±60.58)	^a 1.55-379 (44.12 ± 55.97)	^a 5.11- 673.33(108.52±88.99)	^a 3.2-477(81.41±69.77)	4.1-266(85.29±47.23)	^a 8.25- 673.33(111.02±87.92)	^a 0.7- 333(51.71±52.66)
Conductivity	^a 71.2-9371 (842.55±874.67)	^a 30.67- 2885(599.95±479.62)	^a 30-2885 (271.09±335.46)	^{ab} 83.01- 13341.33(1077.06±1452.74)	^a 48- 2738(550.23±452.36)	^b 30.67- 2377.67(556.84±480.72)	^a 77.5- 13341.33(914.78±1254.12)	^a 25.66- 2377.67(341.69±346.34)
Nitrates	^a 0.47-61.20 (10.74±11.07)	^b 0.47-76.45 (8.86±9.68)	^{ab} 0.47-53.50 (6.31±7.67)	^a 0.95-76.45(11.88±11.23)	^a 0.5-61.2(9.36±8.57)	^a 0.5-36.90(6.98±7.21)	^a 0.95- 76.45(11.54±10.81)	^a 0.5-53.50(6.94±7.27)
Altitude	^a 0-1476 (323±282.69)	^b 0-1933 (311.16 ±311.80)	^{ab} 0-1243.97 (193.38 ±220.99)	^a 0-1200 (277.53±228.95)	^a 0-1933(150.48±196.97)	^a 0- 1589(409.20±338.60)	0- 1042.78(183.22±184.16)	0- 1243.97(189.91±235.5)
pH	^a 7.07-8.8 (8.16±0.27)	^a 6.90-8.8 (8.05±0.34)	^a 6.3-8.6 (7.77±0.50)	^{ab} 7.21-8.8(8.15±0.28)	^a 7-8.8(7.98±0.35)	^b 6.80-8.6(8.03±0.30)	^a 7.33-8.8(8.11±0.28)	^a 6.43-8.6(7.87±0.42)
Water temperature	^a (5- 24.02)13.47±3.7	^a 5- 24.35(15.03±4.52)	^a 6-23.9 (16.93±4.70)	^a 6-26(13.63±3.72)	^{ab} 6-24.5(17.29±4.42)	^b 5-23.7(14.18±4.11)	^a 5-26(15.76±4.64)	^a 6-24.35(17.28±4.56)
TOC	0.2- 8.5(2.60±1.58)	^a 0.2-15(2.49±1.66)	^a 0.6-13.71 (2.86±1.75)	^a 0.2-15(3.25±2.00)	^b 0.2-15(3.03±1.56)	^{ab} 0.2- 10.65(1.95±1.56)	^a 0.5-10.65(3.39±1.62)	^a 0.6-15(3.00±1.71)
Bicarbonates	^a 25- 345(169.65±102.89)	^b 8.67- 350(161.64±94.81)	^{ab} 6.4- 600(117.61±127.64)	^a 21.25- 384(136.91±102.59)	^b 12.75- 350(139.77±93.16)	^{ab} 13.33- 384(199.47±81.13)	25-600(148.66±109.90)	12.2- 384(125.93±98.00)
Ammonium	^a 0.01-4.8(0.17 ±0.39)	^a 0.004- 12.1(0.15±0.63)	^a 0.05-1.2(0.06 ±0.11)	^a 0.01-15.33(0.55±1.93)	^a 0.01-15.27(0.20±0.98)	^a 0.01- 12.10(0.23±1.23)	^a 0.01-15.33(0.48±1.81)	^a 0.01-2.6(0.11±0.26)
Sulphates	^a 3.73-1500 (139.28±197.63)	^a 2.4-970 (96.11±137.11)	^a 1.2-538.5 (36.76±78.65)	^{ab} 4-1500(154.22±212.57)	^a 3.05- 970(103.71±151.58)	^b 2.80- 458(79.88±104.84)	^a 4-1500(160.54±209.45)	^a 1-135(32.54±25.85)
IPS	6.19- 19.95(13.88±3.53)	6.19- 19.95(13.75±3.60)	9.05- 19.65(15.23±3.42)	^a 6.19-18.89(12.42±3.16)	^b 6.52-18.57(12±2.99)	^{ab} 8.01- 19.95(14.88±3.57)	6.19-18.71(12.28±3.16)	6.52- 18.71(13.01±3.93)
IBD	^{ab} 10.9- 20(17.05±2.47)	^a 5.4- 20(15.56±3.33)	^b 8.2- 20(15.02±2.82)	^a 5.7-20(14.04±3.33)	^b 5.4-20(13.62±2.81)	^{ab} 8.7-20(16.86±3.11)	5.4-19.1(12.78±2.86)	5.4-20(13.33±3.21)

