1  Modeling the ecological status response of rivers to multiple stressors using machine

2  learning: a comparison of environmental DNA metabarcoding and morphological data

3  Juntao Fan [a], Shuping Wang [a], Hong Li [b, c], Zhenguang Yan [a, *], Yizhang Zhang [a, d], Xin

4  Zheng [a], Pengyuan Wang [a]

5  [a] *State Key Laboratory of Environmental Criteria and Risk Assessment, Chinese*

6  *Research Academy of Environmental Sciences, Beijing, 100012, China*

7  [b] *Lancaster Environment Centre, Lancaster University, LA1 4YQ, UK*

8  [c] *UK Centre for Ecology & Hydrology, MacLean Building, Wallingford OX108 BB, UK*

9  [d] *Chinese Research Academy of Environmental Sciences Tianjin Branch, Tianjin,*

10  *300457, China*

11  * Corresponding author.

12   *E-mail address:* zgyan@craes.org.cn (Z. Yan).

13

14    **ABSTRACT**

15    Understanding the ecological status response of rivers to multiple stressors is a

16    precondition for river restoration and management. However, this requires the

17    collection of appropriate data, including environmental variables and the status of

18    aquatic organisms, and analysis via a suitable model that captures the nonlinear

19    relationships between ecological status and various stressors. The morphological

20    approach has been the standard data collection method employed for establishing the

21    status of aquatic organisms. However, this approach is very laborious and restricted to

22    a specific set of organisms. Recently, an environmental DNA (eDNA) metabarcoding

23    data approach has been developed that is far more efficient than the morphological

24    approach and potentially applicable to an unlimited set of organisms. However, it

25    remains unclear how well eDNA metabarcoding data reflects the impacts of

26    environmental stressors on aquatic ecosystems compared with morphological data,

27    which is essential for clarifying the potential applications of eDNA metabarcoding data

28    in the ecological monitoring and management of rivers. The present work addresses this

29    issue by modeling organism diversity based on three indices with respect to multiple

30    environmental variables in both the catchment and reach scales. This is done by

31    corresponding support vector machine (SVM) models constructed from eDNA

32    metabarcoding and morphological data on 24 sampling locations in the Taizi River

33    basin, China. According to the mean absolute percent error (MAPE) between the

34    measured diversity index values and the index values predicted by the SVM models,

35    the SVM models constructed from eDNA metabarcoding data (MAPE = 3.87) provide

36    more accurate predictions than the SVM models constructed from morphological data

37    (MAPE = 28.36), revealing that the eDNA metabarcoding data better reflects

38    environmental conditions. In addition, the sensitivity of SVM model predictions of the

39    ecological indices for both catchment-scale and reach-scale stressors is evaluated, and

40    the stressors having the greatest impact on the ecological status of rivers are identified.

41    The results demonstrate that the ecological status of rivers is more sensitive to

42    environmental stressors at the reach scale than to stressors at the catchment scale.

43    Therefore, our study is helpful in exploring the potential applications of eDNA

44    metabarcoding data and SVM modeling in the ecological monitoring and management

45    of rivers.

46    *Keywords*

49

## 1. Introduction

River ecosystems are impacted by multiple environmental variables at both the catchment scale and reach scale simultaneously, and any of these variables lying outside of their normal range can become a stressor. These natural and anthropogenic stressors always interact and are directly or indirectly impacting ecological status (Mori et al., 2019; Romero et al., 2018). For example, catchment scale stressors, such as increased impervious land use by humans, altere physical and chemical conditions of rivers such as increased nutrition through hydrological processes, affecting the structure and function of aquatic ecosystems (Bernhardt et al., 2012; Von Schiller et al., 2017). Here, aquatic communities play an important role in supporting ecosystem services, stability, and biodiversity, and their status can reflect the long-term cumulative effects of environmental stressors on aquatic ecosystems (Franzo and Del Negro, 2019). Therefore, biomonitoring is essential for assessing the impacts of human disturbance at the multiple scales of river basins. The standard approach that has been applied to river biomonitoring involves the sorting and morphological identification of aquatic communities, which is time-consuming and demands a high degree of taxonomic expertise (Pawlowski et al., 2018). However, the high-throughput amplicon sequencing of environmental DNA (eDNA) has recently provided a viable option for biomonitoring, which purified from substrates such as soil or water contains DNA fragments originating from organisms present in that environment (Cordier et al., 2017; Jarman et al., 2018; Mize et al., 2019; Visco et al., 2015). Moreover, a number of previous studies have shown that eDNA metabarcoding data can provide an accurate indication of environmental changes. For example, the relative abundance of operational taxonomic units (OTUs) indicative of plankton was demonstrated to have a significant negative correlation with river nutrient levels (Li et al., 2018a). The foraminifera diversity

75    inferred from eDNA metabarcoding data was found to have a significant positive

76    correlation with the biodiversity in the benthic zone impacted by fish farming activities

77    (He et al., 2019), and the distance from a wellhead in the ocean (Laroche et al., 2016).

78    Benthic macroinvertebrates diversity inferred from eDNA metabarcoding data were

79    also used to assess the freshwater quality (Fernandez et al., 2018; Hering et al., 2018).

80    In addition, previous studies have shown that, compared with morphological

81    classification, eDNA metabarcoding is a relatively simple and affordable method for

82    assessing biodiversity on a large temporal and spatial scale without the need for time-

83    consuming microscopy analysis by experts (He et al., 2019; Ji et al., 2013). Taxonomic

84    classification based on eDNA metabarcoding is usually more accurate than

85    morphological identification, particularly for species with similar morphology and

86    species with poor life cycle characteristics (He et al., 2019; Humbert et al., 2010).

87    Furthermore, eDNA metabarcoding data can be easily reanalyzed to make it suitable

88    for review by third parties (Ji et al., 2013). However, it remains unclear how well eDNA

89    metabarcoding data reflects the impacts of environmental stressors on aquatic

90    ecosystems in comparison with morphological identification data. Clarifying this issue

91    will illuminate potential applications of eDNA technology in the monitoring and

92    management of aquatic ecosystems.

93       Understanding the response of river ecosystems to multiple stressors and identifying

94    important stressors are prerequisites for conducting effective river restoration and

95    management (Meissner et al., 2019; Zhang, 2019). Developing this understanding

96    requires the analysis of biomonitoring data via a suitable model that captures the

97    relationships between the status of ecosystems and various stressors. However, the

98    interactions of multiple stressors produce a combined effect that can be equal to

99    (additive), greater than (synergistic), or less than (antagonistic) the sum of each single

100    effect (Piggott et al., 2015). Indeed, the response of aquatic ecosystems to multiple

101    stressors is typically nonlinear, which greatly complicates the development of accurate

102    models (Jones et al., 2017). The modeling of nonlinear responses can be conducted

103    using various methods, including mathematical/physical models, statistical models, and

104    data-driven models (Al-Mukhtar, 2019; Choubin et al., 2018; Park et al., 2015).

105    However, the complexity of relationships between ecological status and multiple

106    stressors limits the application of mathematical/physical models, and statistical models

107    also suffer from disadvantages, such as poor generalizability due to relatively small

108    sample sizes (Cui and Gong, 2018; Varoquaux, 2018). The development of machine

109    learning (ML) over the past few years has provided a new approach for quantifying

110    these nonlinear relationships (Torija and Ruiz, 2015). At present, ML models have been

111    widely used in the prediction of environmental or ecological indicators. For example, a

112    Bayesian belief network (BBN) was applied to model the combined effects of land use

113    change and climate change on the status of macroinvertebrates and fish in freshwater

114    bodies (Olson, 2018). In addition, artificial neural networks (ANNs), the support vector

115    machine (SVM) and generalized regression neural network, were used for predicting

116    chlorophyll-a concentrations in freshwater, and the results demonstrated that these data-

117    driven ML methods achieved better prediction performance than conventional

118    statistical methods (Marvuglia et al., 2015; Park et al., 2015). The SVM method is

119    particularly advantageous for modeling nonlinear response relationships because the

120    SVM is good for solving high-dimensional and nonlinear problems, while avoiding the

121    difficulties associated with determining the network structure and local minima of the

122    solutions, and provides good generalizability and relatively good prediction

123    performance under small sample size conditions (Vapnik, 1999). These advantages have

124    made SVM outperform other ML methods, e.g., standard ANNs, random forest (RF)

125    classifiers, and boosted trees (BT) classification, in the prediction of soil organic carbon,

126    clay content, and pH (Rossel and Behrens, 2010; Were et al., 2015) and chlorophyll-a

127    (Park et al., 2015) in some regions. Therefore, the SVM is well suited for modeling the

128    relationships between the ecological status of rivers and multiple stressors.
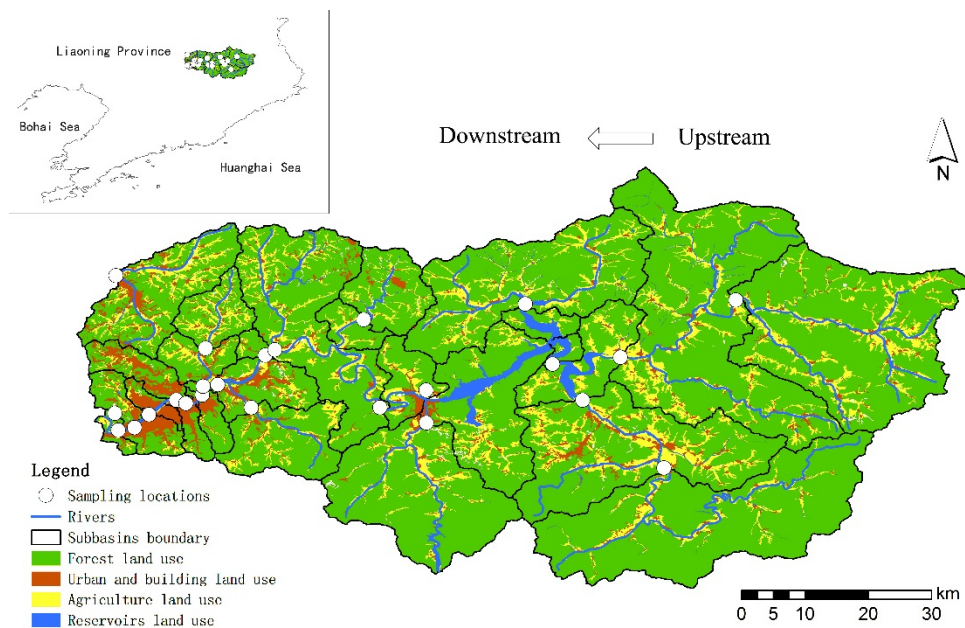
129    The present study compares the ability of eDNA metabarcoding data and

130    morphological identification data to reflect the nonlinear impact of multiple

131    environmental stressors on aquatic ecosystems by employing both sets of data in SVM

132    models corresponding to three ecological indices (i.e. observed species, Shannon

133    Wiener index, and Simpson index), which were commonly used in biodiversity

134    assessment inferred from eDNA metabarcoding or morphological data. As such, the

135    present work helps to explore the potential applications of eDNA technology in the

136    monitoring and management of aquatic ecosystems. In addition, the sensitivity of SVM

137    model predictions of the ecological indices to individual catchment-scale and reach-

138    scale stressors is evaluated, and the stressors having the greatest impact on the

139    ecological status of rivers are identified.

140

141    **2. Materials and methods**

142    *2.1. Study area*

143    The study area was the upstream area of the Taizi River basin (122°23'E–122°53'E,

144    40°28'N–41°39'N) in northeastern China. The location and characteristics of the study

145    area are illustrated in Fig. 1. A previous study demonstrated that the ecological status

146    of the Taizi River in this area was relatively good because the majority of the land in

147    the upstream area was covered by forests, and the intensity of human activities was

148    relatively low (Fan et al., 2015). The primary aquatic organisms of the Taizi River,

149    particularly those species most sensitive to environmental stressors, such as clean-type

fishes (*Lampetra morii* and *Odontobutis Obscurus*) and macrobenthos (*Epeorus melli* and *Cambaroides dauricus*), are mainly distributed in the upper reaches of the river. All of these organisms play an important role in maintaining the health of the aquatic ecosystem. However, the urbanization process in the region and the acceleration of human activities in recent years, such as agriculture and mining, have resulted in water shortages, the deterioration of water quality, habitat damage, loss of biodiversity, and the reduction of ecological functions.



**Fig. 1.** Map of the study area incorporating the upper area of the Taizi River basin at the time of sample collection in October, 2018. The 24 sampling sites and different types of land use in the sub-basins are indicated, and the location of the study area relative to the national boundary of China is shown in the inset.

*2.2. Ecological and environmental data collection*

The 24 sites sampled during October 2018 (Fig. 1) were located in the mainstem and tributaries of the upstream area of the Taizi River basin. Surface water was sampled using sterile bottles. One liter per site was used for eDNA metabarcoding analysis.

167   Three independent extractions of 300 mL were obtained from each one-liter water

168   sample within 6 h after sampling by filtering across a Millipore 0.22 μm hydrophilic

169   nylon membrane. The membrane discs containing captured eDNA were placed in 5.0

170   mL centrifugal tubes, and were instantly frozen and stored at −20°C until DNA

171   extraction. For morphological identification, phytoplankton samples were collected at

172   each sampling site by dragging a nylon mesh with a pore size of 64 μm under the water

173   surface for about 2 min. The water sample concentrated in the drip tube of the net was

174   collected in a 50 mL sample bottle and fixed using Lugol's solution.

175   Environmental variables considered include catchment-scale variables (i.e., land use

176   data) and reach-scale variables (i.e., physicochemical parameters). Land use data were

177   extracted from an analysis of Spot Image data obtained with a 2.5 m resolution. The

178   proportion of land use types (i.e., forest, agriculture, urban, and industrial) was

179   determined for the region of the catchment upstream of each sampling site contributing

180   to the sample characteristics and for a 250 m impact zone adjacent to the studied river

181   segment. Ten physicochemical indicators were selected, including electrical

182   conductivity (EC), dissolved oxygen (DO), pH, biological oxygen demand over 5 days

183   ($BOD_5$), permanganate index ($COD_{Mn}$), total phosphorus (TP), ammonia nitrogen

184   concentration ($NH_3$-N), total nitrogen (TN), suspended sediment (SS), and volatile

185   phenol (VP). The work of (Fan et al., 2015) and Chinese Quality Standards for Surface

186   Water Resources (Ministry of Water Resources, 1994) established thresholds not to be

187   exceed to assure high ecological status for these physicochemical parameters. These are

188   given as follows: EC = 400 μs/cm, DO = 7.5 mg/L, $BOD_5$ = 3 mg/L, $COD_{Mn}$ = 2 mg/L,

189   $NH_3$-N = 0.15 mg/L, TN = 0.2 mg/L (which was only considered in lake or reservoir

190   samples), TP = 0.02 mg/L, VP = 0.002 mg/L, SS = 20 mg/L, pH = 6.5~8.5.

191

192   *2.3. eDNA metabarcoding and morphological identification*

193   Phytoplankton is the target taxonomic group of eDNA metabarcoding and

194   morphological identification. Total eDNA was extracted using the cetyl

195   trimethylammonium bromide (CTAB) method combined with the Zymo DNA Clean &

196   Concentrator kit (Zymo Research Corp, Irvine, USA) (Yuan et al., 2015). The

197   concentration of eDNA was determined using a NanoDrop One microvolume

198   ultraviolet-visible (UV-vis) spectrophotometer (Thermo Fisher Scientific, Carlsbad,

199   USA). The eDNA was used as templates for the polymerase chain reaction (PCR)

200   method with 18S rRNA gene primers 18SV9F (5'-CCCTGCCNTTTGTACACAC-3')

201   and 18SV9FR (5'-CCTTCNGCAGGTTCACCTAC-3') (Amaral-Zettler et al., 2009; De

202   Vargas et al., 2015). The 18S rRNA gene primers were used because phytoplankton

203   diversity including the cryptic diversity in environmental samples can be indicated by

204   sequencing of 18S rRNA gene, and the SILVA datasets offered the 18S primer

205   opportunity to assess distribution patterns of phytoplankton species (Treusch et al.,

206   2012). The purified PCR products were added with 8-base sequence tags corresponding

207   to each sample. High throughput sequencing was conducted using a MiSeq sequencing

208   platform (Illumina, San Diego, USA). All low-quality sequencing data points with

209   adaptors, ambiguous bases, low complexity, and those having average quality scores

210   less than 20 were discarded using the UPARSE pipeline (Edgar, 2013). The OTUs were

211   determined at the ≥97% identity level (Edgar, 2013). Taxonomic annotation analysis

212   was performed using the Qiime2 pipeline (Caporaso et al., 2010) with respect to the

213   SILVA-119 reference database. The remaining high-quality data were transformed to

214   relative proportions before conducting subsequent statistical analysis.

215   For morphological identification, samples were concentrated and precipitated, and

216   the sample volume was adjusted to 20–50 mL. The concentrated sample was then

217     shaken uniformly, and 0.1 mL of the sample was immediately placed in a counting box

218     for morphological identification. The phytoplankton taxa in each sample were

219     identified under a 10 × 40 microscope. However, if a high concentration of diatoms

220     were observed, the sample was sealed and identified under a 10 × 100 microscope. The

221     specimens were identified to species level through microscopy and taxonomic experts

222     consultation. The reference used to identify phytoplankton is the Freshwater Alage of

223     China – Systematics, Taxonomy and Ecology (Hu and Wei, 2006).

224

225     *2.4. SVM model development*

226     The ecological status of the samples was evaluated according to the obtained eDNA

227     metabarcoding and morphology identification data based on three widely used

228     ecological indices, i.e., observed species (Kefford et al., 2011), Shannon Wiener index

229     (Strong, 2016), and Simpson index (Keylock, 2005). The abbreviations and ecological

230     significance of each of these indices are listed in Table 1. The values for these ecological

231     indices obtained from the eDNA metabarcoding and morphology identification data

232     were employed as the response/dependent variables in their respective SVM models.

233     The catchment-scale variables and reach-scale variables were input to the respective

234     SVM models as the independent variables.

235

236     **Table 1**

237     List of ecological indices with abbreviations and ecological significance.

| Ecological index/Response variables | Abbreviations in eDNA metabarcoding | Abbreviations in morphological identification | Ecological significance |
|---|---|---|---|
| Observed species | Species _E | Species _M | Number of species or OTU observed. |
| Shannon Wiener index | Shannon _E | Shannon _M | The species/OTUs richness and evenness of the community, but predominantly sensitive to richness. |

11

| | | | Richness increases with increasing index value. |
| Simpson index | Simpson _E | Simpson_M | The species/OTUs richness and evenness of the community, but predominantly sensitive to evenness. Evenness increases with increasing index value. |

238

239    The SVM was applied for nonlinear regression analysis to establish the response of

240    the ecological indices to the multiple environmental variables. Here, the input data were

241    mapped initially into a higher-dimensional feature space via a kernel function (i.e., a

242    linear kernel, polynomial kernel, radial basis kernel, and Gaussian kernel), and then

243    linear regression was performed in the high-dimensional feature space to obtain the

244    nonlinear regression effect in the original space (Balfer and Bajorath, 2015; Bouboulis

245    et al., 2015). The specific kernel function applied was selected by cross-validation

246    (Piette and Moore, 2018).

247    The regression performance of the SVM depends on the appropriate selection of

248    parameter values, including cost ($c$), epsilon ($\varepsilon$), and gamma ($\gamma$), where both $c$ and $\varepsilon$ are

249    employed to establish the penalty coefficient, which represents the error tolerance of

250    the regression analysis, and $\gamma$ determines the distribution of the data after it is mapped

251    to the new feature space. Here, the number of support vectors decreases with increasing

252    $\gamma$, which affects the speed of training and prediction. The values of these parameters are

253    optimized using a loop traversal algorithm (Cherkassky and Ma, 2004). Normalization

254    was applied to all independent variables to ensure that the indicator values were

255    comparable.

256    The generalization ability of the model was verified by 8-fold cross validation, where

257    the dataset was divided into 8 subsets, and each subset was employed as the testing set

258    once, while the remaining 7 subsets were used as the training set. Accordingly, this

259     process was repeated 8 times. The prediction error of each model was evaluated based

260     on the mean absolute percent error (MAPE), which is calculated for $n$ samples as

261     follows:

262 $$\text{MAPE} = \sum_{t=1}^{n} \left| \frac{Observed_t - Predicted_t}{Observed_t} \right| \times \frac{100}{n} \tag{1}$$

263     where $Observed_t$ is the observed value and $Predicted_t$ is the predicted value. Then, the

264     model with the smallest MAPE value was selected as the optimal model.

265       Sensitivity analysis was applied to determine the environmental variables that most

266     greatly influenced the model predictions of the ecological indices. This was conducted

267     using the one-factor-at-a-time (OAT) approach. Here, the MAPE values of the model

268     predictions were obtained with one environmental variable omitted at a time, while the

269     other environmental variables were held constant. Then, the impact of each

270     environmental variable on the model prediction was evaluated according to the absolute

271     value of the difference between the MAPE obtained with and without that variable,

272     which is denoted herein as ΔMAPE. Accordingly, the sensitivity of the ecological index

273     predictions to an environmental variable increases with increasing ΔMAPE.

274

275 **3. Results**

276 *3.1. Environmental conditions in catchment and reach scales*

277       All the environmental variables have become stressors, which are marked with "+"

278     in Table 2. Spatial analysis showed that almost all sites were under the selected

279     catchment-scale stressors, and the downstream sites (e.g., s19, s15 and s22) were under

280     more reach-scale stressors than the upstream sites (Table 2). We note that the proportion

281     of forest land use in the catchment scale (0.268–0.910) is greater than that in the 250 m

282     buffer zone (0.092–0.566). However, the proportion of agriculture land use in the 250

283     m buffer zone (maximum value of 0.596) is greater than that in the catchment scale

284  (maximum value of 0.265), which indicates that agricultural disturbance is greater in

285  the riparian zone than at the catchment scale, while urban and industrial disturbances

286  have opposite behaviors. Table 2 also indicates that, TN and VP were the reach-scale

287  variables with the highest number of sites exceeding the thresholds.

288

289  **Table 2**

290  List of environmental variables included in the modeling and spatial distribution of sites

291  with corresponding stressors. Stressors, i.e., catchment-scale variables impacted by any

292  artificial land use types (i.e., agriculture, urban and industrial land use), and reach-scale

293  variables with values less than or greater than the threshold values representing high

294  environmental status established by the work of (Fan et al., 2015) and Chinese quality

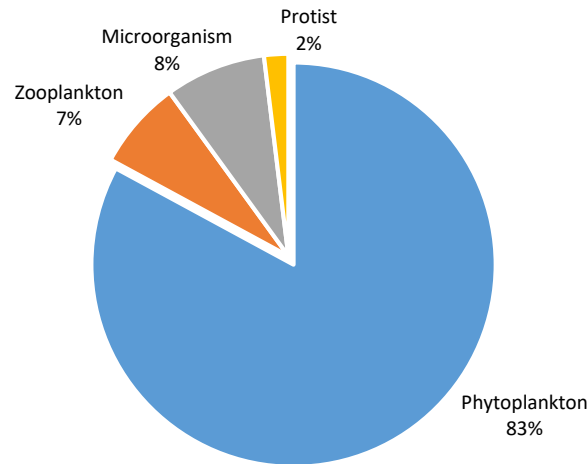295  standards for surface water resources (Ministry of Water Resources, 1994), are marked

296  with "+".

| Environmental variables | | Abbreviations (Units) | Ranges | Sites with corresponding stressors |
|---|---|---|---|---|
| Catchment-scale variables | | | | |
| Forest land use (catchment scale) | + | F_cat (proportion) | 0.268–0.910 | All sites |
| Forest land use (250 m buffer zone) | + | F_buf (proportion) | 0.092–0.566 | All sites |
| Agriculture land use (catchment scale) | + | A_cat (proportion) | 0.018–0.265 | All sites |
| Agriculture land use (250 m buffer zone) | + | A_buf (proportion) | 0.000–0.596 | All sites except s19, s21, s16, s20 |
| Urban and industrial land use (catchment scale) | + | U_cat (proportion) | 0.016–0.646 | All sites |
| Urban and industrial land use (250 m buffer zone) | + | U_buf (proportion) | 0.011–0.520 | All sites |
| Reach-scale variables | | | | |
| Electrical conductivity | + | EC (μs/cm) | 142.47–655.33 | s19, s13, s14, s21, s15, s17, |

14

| | | | | s20, s22 |
|---|---|---|---|---|
| Dissolved oxygen | + | DO (mg/L) | 7.02–14.26 | s22 |
| pH | + | pH | 7.84–8.98 | s19, s12, s10, s02, s20 |
| Permanganate index | + | $COD_{Mn}$ (mg/L) | 0.48–5.72 | All sites except s10, s03, s11 |
| Five-day biochemical oxygen demand | + | $BOD_5$ (mg/L) | 0.75–8.41 | s04, s14, s15, s16, s18, s24 |
| Ammonia nitrogen | + | $NH_3$-N (mg/L) | 0.12–3.87 | All sites except s05, s03, s07, s06, s01, s02 |
| Total nitrogen | + | TN (mg/L) | 1.55–6.75 | All sites |
| Total phosphorus | + | TP (mg/L) | 0.004–0.223 | s19, s04, s14, s06, s21, s23, s15, s09, s16, s24, s20, s22 |
| Suspended sediment | + | SS (mg/L) | 1.56–35.33 | s15, s01, s08 |
| Volatile phenol | + | VP (mg/L) | 0.004–0.112 | All sites |

297

*3.2. Ecological status derived from eDNA metabarcoding and morphological data*

A total of 67 18S rRNA gene libraries were analyzed according to the methodology presented in Subsection 2.3, which resulted in a total of 2,305,498 high-quality sequences, and a total of 6,635 OTUs. The number of OTUs in each sample was distributed between 477–2,661 (Table S1). The result of taxonomic group distribution of OTUs showed that approximately 83% eukaryotic sequences were annotated as phytoplankton (Fig. 2), which confirmed that the phytoplankton can be indicated by sequencing of 18S rRNA gene. Therefore, the eDNA metabarcoding data and morphological data are comparable in this study.
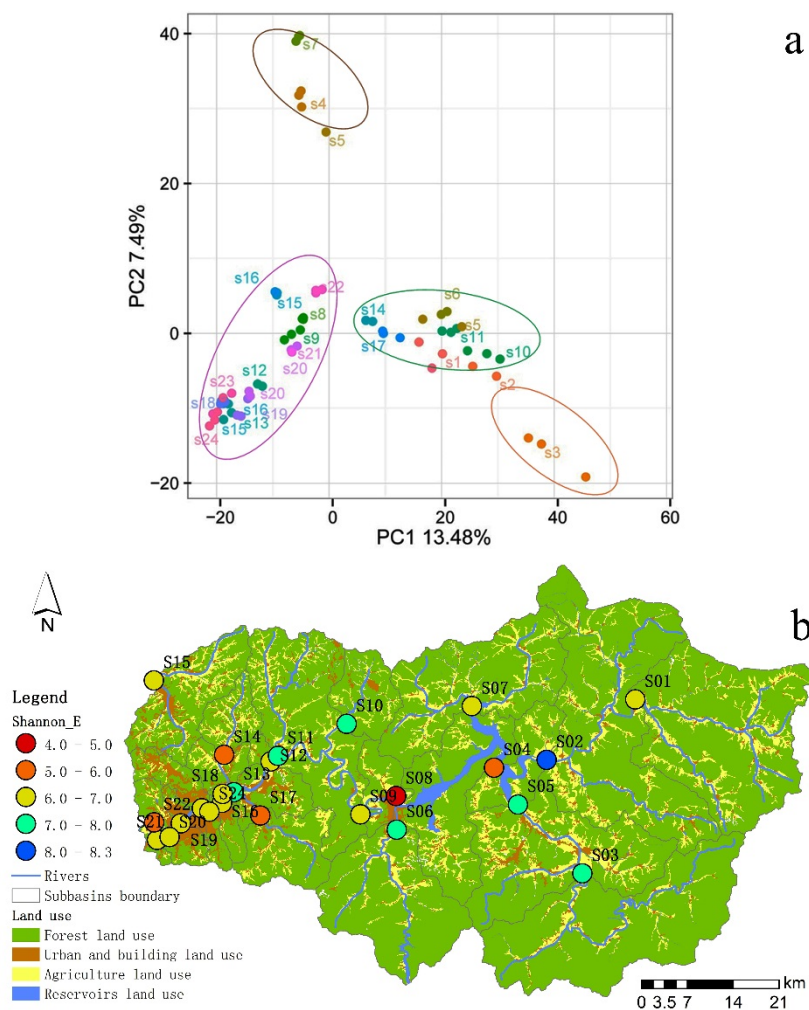
307

**Fig. 2.** Percentage of the sequences assigned to each of taxonomic groups.

309    An analysis of the relative abundances of the top 15 orders and families of organisms

310    for the three replications of the 67 samples were shown in Fig. S1 and S2, respectively.

311    However, approximately 70% of sequences cannot be assigned to genus level because

312    the limitation of reference information in the SILVA database. Analysis of top 15

313    families of organisms indicated that the Mediophyceae, Ochromonadales and

314    Chlorodendrales accounted for approximately 17.5%, 9.9% and 5.4% of all taxa,

315    respectively. Analysis of variance (ANOVA) results indicated that no significant

316    differences were observed for the relative family abundances among the sample

317    replications ($p > 0.05$).

318    The OTU compositions of the different samples were analyzed according to beta

319    diversity to reflect differences between samples using principal component analysis

320    (PCA). Here, PCA uses variance decomposition to reflect the differences between

321    multiple sets of data on a two-dimensional coordinate graph, where the coordinate axes

322    are two eigenvalues that reflect the variance to the greatest extent. As such, samples

323    with similar compositions were clustered in the PCA graph, as shown in Fig. 3A based

324    on the sampling locations illustrated in Fig. 3B, which also showed the Shannon_E

325    values for the individual sampling locations and the land use types of the study area.
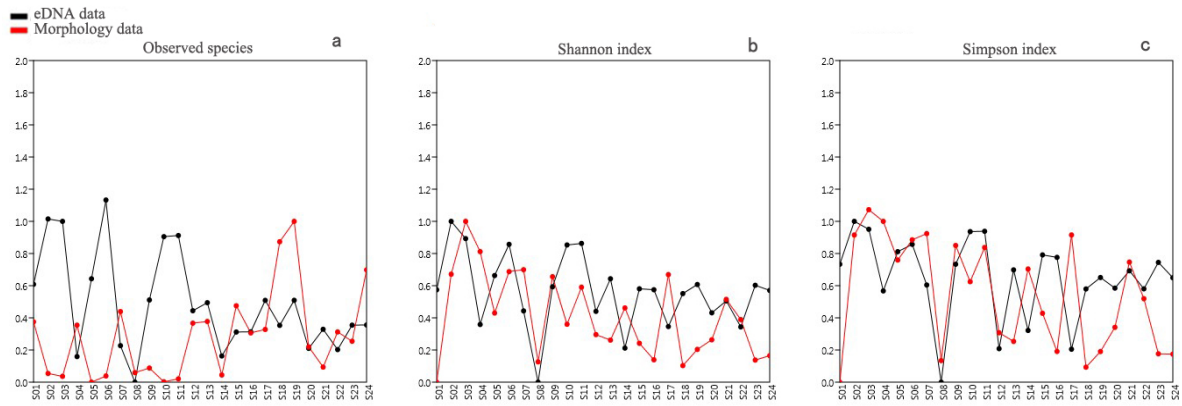
16

326  The results indicated that significant differences exist between the sampling sites of

327  upstream tributaries (e.g., s3, s2 and s1) and the sampling sites of the middle and lower

328  mainstem, while differences were also observed between the urban (e.g., s22, s21, and

329  s19) and mountainous sections (e.g., s6, s2, and s5) of the mainstem. However, the some

330  sites were impacted by the reservoir located in the mainstem of upstream (e.g. s04 and

331  s08). The spatial distribution of Shannon_E values presented the same pattern, where

332  the Shannon_E value tended to gradually decrease with increasing disturbance from

333  human activity from the upstream to the downstream regions, as reflected by increasing

334  urban and industrial land use.

335



336

337  **Fig. 3.** (a) Principal component analysis graph for all samples based on the beta

338 diversity derived from eDNA metabarcoding data and (b) the spatial distribution of

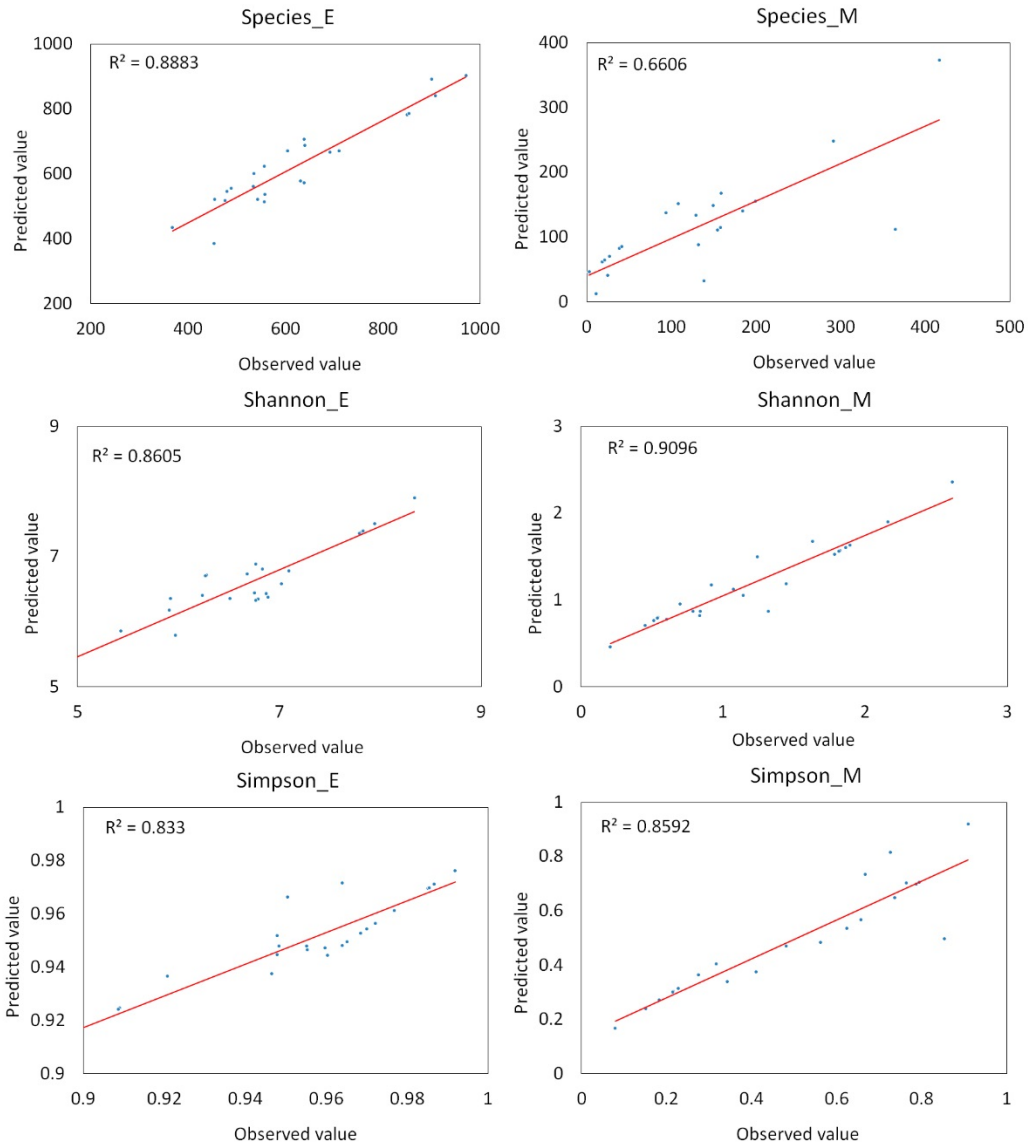339 ecological status based on the Shannon Wiener index.

340

341     The diversity values measured according to the observed species, Shannon Wiener

342 index, and Simpson index derived from eDNA metabarcoding and morphological

343 identification data were normalized and compared, and the results were given in Fig.

344 4A, B, and C, respectively, for sample locations s01–s24. The results in Fig. 4A

345 indicated that in most sites, the observed species values obtained based on eDNA

346 metabarcoding were higher than the values based onmorphological data, because OTUs

347 contained a greater number of taxa information. This difference decreased in the

348 Shannon Wiener and Simpson index values, which demonstrated that the data obtained

349 by the two methods reflect similar richness and evenness characteristics of community

350 composition in most sampling sites (Fig. 4B and C). Fig. 4A also showed that 8 sites

351 out of 24 were higher for morphological data than eDNA metabarcoding data, and most

352 of these sites are located in the downstream of study area (e.g. s15, s18, s19 and s24),

353 where a large number of *Cyclotella meneghiniana* were detected in morphological data.

354 *Cyclotella meneghiniana* is a typical indicator of water pollution (Duong et al., 2008).

355 This was proved by Fig. 4B and C, which showed that the Shannon and Simpson indices

356 derived from morphological data were relatively low at these downstream sites.

357 However, the ecological indices derived from eDNA data showed better consistency at

358 these sites, which indicated that the difference between eDNA metabarcoding and

359 morphological data may become larger in polluted river sections.

360

361

**Fig. 4.** Comparison of the three ecological index values derived from eDNA metabarcoding and morphology identification data.

*3.2. Predictive performances and sensitivity analysis of SVM models*

After optimizing the model parameters ($c = 10000$, $\varepsilon = 0.2$, and $\gamma = 0.025$) according to the methodology presented in Subsection 2.4, the nonlinear regression analysis results obtained by the SVM models for the three indices (Species_E, Shannon_E, Simpson_E) derived from eDNA metabarcoding data and the three indices (Species_M, Shannon_M, Simpson_M) derived from morphological identification data are presented in Fig. 5. The results indicated that, with the exception of Species_M (squared correlation coefficient $R^2 = 0.66$), the SVM models achieved good prediction performance, with $R^2$ values that were all greater than 0.80.

**Fig. 5.** Nonlinear regression fitting plots of the support vector machine (SVM) models for the measured values and predicted values of the three ecological indices.

The minimum values of MAPE for all samples (MAPE_ALL) and the minimum values of MAPE for the test samples (MAPE_TEST) obtained by 8-fold cross-validation indicated the accuracy of different models (Table 3). The results indicated that the MAPE_ALL values of the three most accurate SVM models obtained from eDNA metabarcoding data were in the order of Species_E > Shannon_E > Simpson_E, and the MAPE_ALL values of the three most accurate SVM models obtained from

20

morphology identification data exhibited an equivalent pattern. Nevertheless, the SVM

models constructed from the eDNA metabarcoding data had MAPE values that were

much smaller than those of the models constructed from the morphological

identification data whether based on MAPE_ALL or MAPE_TEST values. This

indicated that the models constructed from eDNA metabarcoding data were more

accurate than those constructed from the morphological identification data.

**Table 3**

Results of model selection using 8-fold cross-validation for each ecological index given

in terms of the minimum values of MAPE for all samples (MAPE_ALL), and the

minimum values of MAPE for the test samples (MAPE_TEST).

| Ecological index | MAPE_ALL | MAPE_TEST |
|---|---|---|
| Index derived from eDNA metabarcoding data | | |
| Species _E | 9.06 | 6.72 |
| Shannon_E | 5.14 | 4.14 |
| Simpson_E | 1.33 | 0.75 |
| Index derived from morphology identification data | | |
| Species _M | 183.96 | 49.57 |
| Shannon_M | 25.61 | 15.50 |
| Simpson_M | 25.37 | 20.00 |

The sensitivity of each ecological index to multiple stressors were varying (Table 4).

For Species_E, the largest value of $\Delta$MAPE = 1.12 was obtained for SS, indicating that

the Species_E prediction was most sensitive to this variable. For Shannon_E, the largest

value of $\Delta$MAPE = 0.47 was obtained for SS, indicating that the Shannon_E prediction

was most sensitive to this variable. For Simpson_E, the largest value of $\Delta$MAPE = 0.05

was obtained for DO, indicating that the Simpson_E prediction was most sensitive to

402 this variable. Likewise, we can determine that the Species_M prediction was most

403 sensitive to DO (ΔMAPE = 21.79), the Shannon_M prediction was most sensitive to

404 VP (ΔMAPE = 2.17), and the Simpson_M prediction was most sensitive to VP

405 (ΔMAPE = 2.13). We also note from Table 4 that the magnitudes of the ΔMAPE values

406 for the ecological indices obtained from DNA metabarcoding data are much smaller

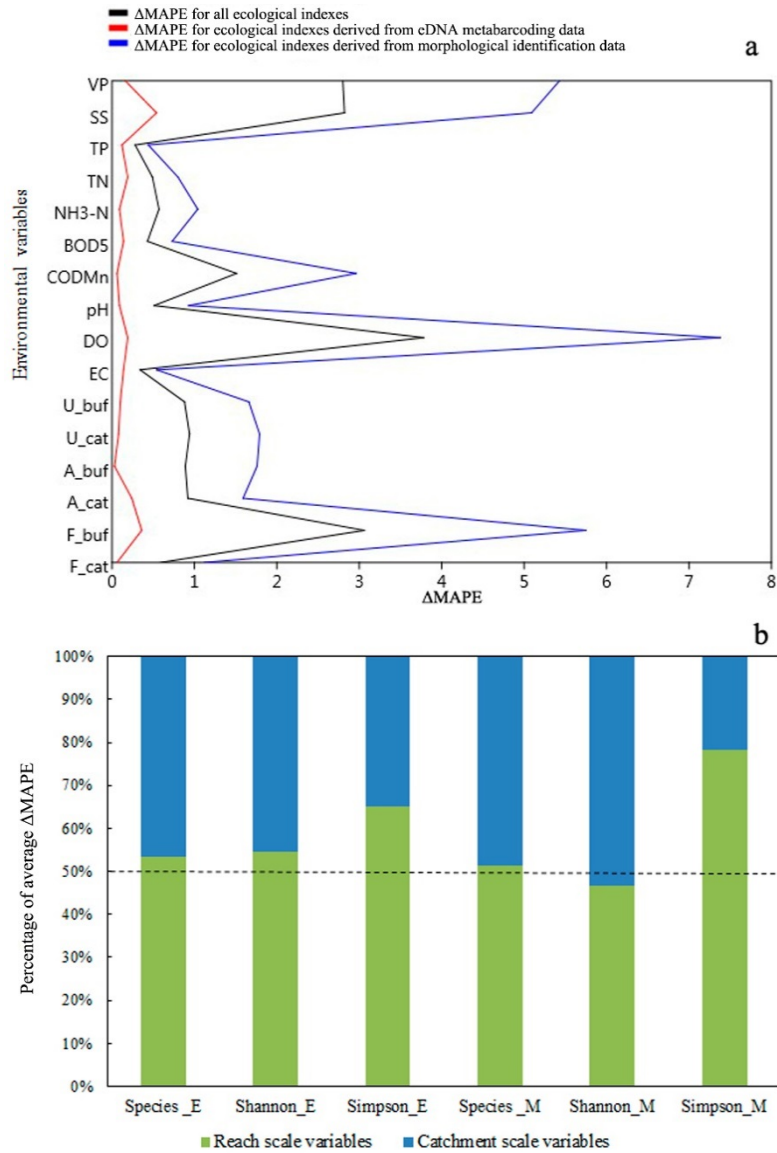407 than those obtained from morphological identification data.

408
409 **Table 4**

410 Results of sensitivity analysis based on the change in MAPE values (ΔMAPE) for all

411 samples with respect to the individual environmental variables.

| Environmental variables | Species_E | Shannon_E | Simpson_E | Species_M | Shannon_M | Simpson_M |
|---|---|---|---|---|---|---|
| | ΔMAPE | | | | | |
| Catchment-scale variables | | | | | | |
| F_cat | 0.13 | 0.04 | 0 | 3.19 | 0.16 | 0.02 |
| F_buf | 0.62 | 0.44 | 0.03 | 15.47 | 1.59 | 0.18 |
| A_cat | 0.62 | 0.09 | 0 | 3.69 | 0.78 | 0.31 |
| A_buf | 0 | 0.07 | 0.01 | 4.67 | 0.55 | 0.07 |
| U_cat | 0.1 | 0.11 | 0.04 | 4.87 | 0.07 | 0.42 |
| U_buf | 0.29 | 0.01 | 0 | 4.2 | 0.62 | 0.16 |
| Reach-scale variables | | | | | | |
| EC | 0.36 | 0.02 | 0.03 | 1.34 | 0.15 | 0.14 |
| DO | 0.16 | 0.35 | 0.05 | 21.79 | 0.06 | 0.29 |
| pH | 0.15 | 0.08 | 0.03 | 1.84 | 0.68 | 0.28 |
| $COD_{Mn}$ | 0.14 | 0.04 | 0 | 8.27 | 0.13 | 0.47 |
| $BOD_5$ | 0.24 | 0.16 | 0.02 | 0.68 | 0.11 | 1.39 |
| $NH_3$-N | 0.22 | 0.02 | 0.03 | 2.01 | 0.32 | 0.8 |
| TN | 0.34 | 0.19 | 0.04 | 0.63 | 1.12 | 0.64 |

| | | | | | | |
|---|---|---|---|---|---|---|
| TP | 0.28 | 0.07 | 0.01 | 0.8 | 0.05 | 0.47 |
| SS | 1.12 | 0.47 | 0.02 | 14.22 | 0.72 | 0.34 |
| VP | 0.35 | 0.12 | 0.02 | 11.99 | 2.17 | 2.13 |

412

413    The variations in the ΔMAPE values for the ecological indices obtained from DNA

414    metabarcoding and morphological identification data are more clearly shown in Fig.

415    6A. We note that, among all six ecological indices, DO, SS, and VP are the three

416    environmental variables in the reach scale that most greatly affect the index value

417    predictions. These are followed by F_buf, the variable in the catchment scale. A

418    comparison of the average ΔMAPE values obtained for the environmental variables

419    shown in Fig. 6B indicate that, with the exception of Shannon_M, the environmental

420    variables at the reach scale have a greater impact on the ecological indices than those

421    at the catchment scale.

**Fig. 6.** (a) ΔMAPE values for each environmental variable in the sensitivity analysis and (b) a comparison of sensitivities between catchment and reach scale environmental variables.

**4. Discussion**

*4.1. SVM model development and validation*

The SVM models increase our understanding of the non-linear relationships between ecological status and multiple stressors on the one hand and the sensitivity of the ecological status to each stressor on the other. More importantly, the MAPE and high

432 $R^2$ values obtained by the SVM models demonstrate quantitatively that eDNA

433 metabarcoding data provide modeling results that were more indicative of

434 environmental degradation compared with morphological identification data. However,

435 we must note that OTUs contained more taxa information than species, which may

436 increase the uncertainty of the model comparison. However, the greater the number of

437 OTUs does not necessarily mean the better model performance, because a larger data

438 may also bring noise for modeling (Lu et al., 2018). In many biodiversity surveys and

439 assessments, the concept of OUTs diversity has been roughly equated with the concept

440 of species diversity (Caron and Hu, 2019), because OUTs use 3% sequence difference

441 to distinguish species, which is an accepted standard in molecular biology techniques

442 (Schloss and Handelsman, 2005). Previous study also showed that eDNA

443 metabarcoding and morphological macroinvertebrate metrics are positively correlated

444 and indicate the same key gradients in stream condition (Emilson et al., 2017).

445    Furthermore, this kind of uncertainty can be reduced by using some ecological

446 indices (e.g. the Simpson and Shannon-Weiner indices), which represent the relative

447 diversity of taxa. These indices have all been normalized before modeling, and the

448 results showed that the normalized values of these ecological indices are relatively

449 consistent in most sampling sites (Fig. 4). Therefore, the uncertainty of the model due

450 to different classification levels can be reduced. In addition, our results were obtained

451 with a relatively small training dataset, and increasing the number of samples or

452 applying a larger sampling area can lead to the process of refining our predictive models.

453 This is supported by a previous study, which has shown that the accuracy and stability

454 of predictions increased exponentially with increasing sample size regardless of the

455 type of ML algorithm adopted (Cui and Gong, 2018).

456

*4.2. Ecological response derived from eDNA metabarcoding and morphological identification data to multiple stressors*

Although the ecological indices obtained from morphological identification data exhibited good response relationships with multi-scale stressors, the models constructed from eDNA metabarcoding data provided better accuracy, as shown in Table 3. This is because the effective eDNA sequencing information includes a large number of intact and fragmentary organisms, and even includes the DNA information of many historically existing organisms. This is supported by a previous study, which found that the DNA information of some species may exist in water for up to one month after the removal of DNA release sources (Li et al., 2018b). In addition, it has been shown that eDNA metabarcoding data provide more integral information regarding biology, including the taxa and even the potential bioindicators of pollution, for example, the OTUs that dominate eDNA datasets in high mercury concentration do not need to be assigned taxonomically, which are typically overlooked in morphological identification (Frontalini et al., 2018). In conclusion, the biodiversity information contained in eDNA data is massive, and the large volume of data available may alleviate model prediction uncertainties caused by sample size limitations to some extent.

It is worth noting that eDNA metabarcoding data are not able to provide some information available from morphological identification data. For example, eDNA metabarcoding data provides no information regarding the morphological deformations of target organisms, which are often found in highly polluted environments, and are commonly used as evidence for heavy metal pollution (Yanko et al., 1998). Therefore, eDNA metabarcoding data cannot replace morphological analysis when studying the response of a particular species, but methods to detect change population of one or multiple organisms to environmental stressors by eDNA metabarcoding are developing,

482   such as screening for functional genes, which may enable the eDNA metabarcoding to

483   assess toxicological information (Zhang et al., 2019). This will widen the application

484   of eDNA metabarcoding in environmental sciences.

485

486   *4.3. Sensitivity differences with respect to catchment and reach scale stressors*

487     The sensitivity analysis results indicated that DO, SS, VP, and F_buf have the greatest

488   impact on the model predictions of diversity indices, and that environmental variables

489   at the reach scale are more influential than that at the catchment scale, as shown in Table

490   4 and Fig. 6. This greater sensitivity of ecological status to reach-scale stressors can be

491   explained by noting that disturbances in land use at the catchment scale affect the

492   aquatic ecological status by generating non-point source pollutants, such as fertilizers,

493   pesticides, and sewage irrigation, that enter water bodies, resulting in increased

494   nutrition, bacteria, toxicity, and harmful substances, which means that the changes at

495   the reach scale affect the ecological status of rivers directly (Meador and Goldstein,

496   2003; Piggott et al., 2012).

497     The DO is directly decreased under these degradating conditions (Mineau et al., 2015)

498   Moreover, DO has been shown to be a key variable impacting the status of many aquatic

499   species because it can affect the tolerance limit of organisms (Marshall and Elliott,

500   1998). In addition, sites with the highest DO level have also been shown to have the

501   highest aquatic species diversity (Wilhm and McClintock, 1978). A previous study has

502   also demonstrated that SS is critical to phytoplankton communities because

503   phytoplankton growth requires photosynthesis, and light intensity in the photic zone

504   has a significant negative correlation with SS (Van Duin et al., 2001). Finally, we note

505   that urban and industrial land use in the urban section of the study area increased

506   significantly since the work of (Fan et al., 2015), and this can be expected to have

507    released toxic chemicals from industrial pollution, such as VP, into water bodies. In this

508    regard, the photosynthetic activity parameters of algae have been shown to have a

509    negative dose-response relationship to phenol toxicity (Kottuparambil et al., 2014).

510    Therefore, VP represents another critical environmental variable impacting the status

511    of many aquatic species.

512

**5. Conclusion**

513

514      The present study compared the ability of eDNA metabarcoding data and

515    morphological identification data to reflect the nonlinear impact of multiple

516    environmental stressors on aquatic ecosystems by employing both sets of data in SVM

517    models corresponding to three ecological indices (i.e. observed species, Shannon

518    Wiener index, and Simpson index). Analysis of the environmental variables at the

519    catchment and reach scales of the study area indicated that most of the variables

520    exceeded their natural thresholds at some of the sampling sites, and became a complex

521    of simultaneously interacting stressors affecting the ecological status of the river. The

522    SVM models constructed from eDNA metabarcoding data (MAPE = 3.87) provided

523    more accurate predictions than the SVM models constructed from morphological

524    identification data (MAPE = 28.36), revealing that the eDNA metabarcoding data better

525    reflected ecological conditions. As such, the present work helps to explore the potential

526    applications of eDNA technology in the monitoring and management of aquatic

527    ecosystems. In addition, the sensitivity of SVM model predictions of aquatic ecosystem

528    diversity to catchment-scale and reach-scale stressors was evaluated, and the stressors

529    having the greatest impact on the ecological status of rivers were identified. These

530    results indicated that the model predictions were more sensitive to the environmental

531    variables at the reach scale than those at the catchment scale. In addition, DO, SS, VP,

and F_buf were found to be the most influential variables impacting the ecological status of the river.

**References**

Al-Mukhtar, M., 2019. Random forest, support vector machine, and neural networks to modelling suspended sediment in Tigris River-Baghdad. Environmental Monitoring and Assessment 191 (11), 673. https://doi.org/10.1007/s10661-019-7821-5.

Amaral-Zettler, L.A., McCliment, E.A., Ducklow, H.W. and Huse, S.M., 2009. A method for studying protistan diversity using massively parallel sequencing of V9 hypervariable regions of small-subunit ribosomal RNA genes. PLoS One 4 (7), e6372. https://doi.org/10.1371/journal.pone.0006372.

Balfer, J. and Bajorath, J., 2015. Systematic artifacts in support vector regression-based compound potency prediction revealed by statistical and activity landscape analysis. PLoS One 10 (3), e0119301. https://doi.org/10.1371/journal.pone.0119301.

557    Bernhardt, E.S., Lutz, B.D., King, R.S., Fay, J.P., Carter, C.E., Helton, A.M., Campagna,

558        D. and Amos, J., 2012. How many mountains can we mine? Assessing the regional

559        degradation of central Appalachian rivers by surface coal mining. Environmental

560        Science & Technology 46 (15), 8115–8122. https://doi.org/10.1021/es301144q.

561    Bouboulis, P., Theodoridis, S., Mavroforakis, C. and Evaggelatou-Dalla, L., 2015.

562        Complex support vector machines for regression and quaternary classification. IEEE

563        Transactions on Neural Networks and Learning Systems 26 (6), 1260–1274.

564        https://doi.org/10.1109/TNNLS.2014.2336679.

565    Caporaso, J.G., Kuczynski, J., Stombaugh, J., Bittinger, K., Bushman, F.D., Costello,

566        E.K., Fierer, N., Pena, A.G., Goodrich, J.K., Gordon, J.I., Huttley, G.A., Kelley, S.T.,

567        Knights, D., Koenig, J.E., Ley, R.E., Lozupone, C.A., McDonald, D., Muegge, B.D.,

568        Pirrung, M., Reeder, J., Sevinsky, J.R., Turnbaugh, P.J., Walters, W.A., Widmann, J.,

569        Yatsunenko, T., Zaneveld, J. and Knight, R., 2010. QIIME allows analysis of high-

570        throughput community sequencing data. Nature Methods 7 (5), 335–336.

571        https://doi.org/10.1038/nmeth.f.303.

572    Caron, D.A. and Hu, S.K., 2019. Are we overestimating protistan diversity in nature?

573        Trends in Microbiology 27 (3), 197–205. https://doi.org/10.1016/j.tim.2018.10.009.

574    Cherkassky, V. and Ma, Y.Q., 2004. Practical selection of SVM parameters and noise

575        estimation for SVM regression. Neural Networks 17 (1), 113–126.

576        https://doi.org/10.1016/S0893-6080(03)00169-2.

577    Choubin, B., Darabi, H., Rahmati, O., Sajedi-Hosseini, F. and Klove, B., 2018. River

578        suspended sediment modelling using the CART model: a comparative study of

579        machine learning techniques. Science of the Total Environment 615, 272–281.

580        https://doi.org/10.1016/j.scitotenv.2017.09.293.

581    Cordier, T., Esling, P., Lejzerowicz, F., Visco, J., Ouadahi, A., Martins, C., Cedhagen,

582    T. and Pawlowski, J., 2017. Predicting the ecological quality status of marine

583    environments from eDNA metabarcoding data using supervised machine learning.

584    Environmental Science & Technology 51 (16), 9118–9126.

585    https://doi.org/10.1021/acs.est.7b01518.

586  Cui, Z.X. and Gong, G.L., 2018. The effect of machine learning regression algorithms

587    and sample size on individualized behavioral prediction with functional connectivity

588    features. Neuroimage 178, 622–637.

589    https://doi.org/10.1016/j.neuroimage.2018.06.001.

590  De Vargas, C., Audic, S., Henry, N., Decelle, J., Mahé, F., Logares, R., Lara, E., Berney,

591    C., Le Bescot, N. and Probert, I., 2015. Eukaryotic plankton diversity in the sunlit

592    ocean. Science 348 (6237), 1261605. https://doi.org/10.1126/science.1261605.

593  Duong, T. T., Morin, S., Herlory, O., Feurtet-Mazel, A., Coste, M., 2008. Seasonal

594    effects of cadmium accumulation in periphytic diatom communities of freshwater

595    biofilms. Aquatic Toxicology 90, 19–28. https://doi.org/

596    10.1016/j.aquatox.2008.07.012.

597  Edgar, R.C., 2013. UPARSE: highly accurate OTU sequences from microbial amplicon

598    reads. Nature Methods 10 (10), 996–998. https://doi.org/10.1038/nmeth.2604.

599  Emilson, C.E., Thompson, D.G., Venier, L.A., Porter, T.M., Swystun, T., Chartrand, D.,

600    Capell, S. and Hajibabaei, M., 2017. DNA metabarcoding and morphological

601    macroinvertebrate metrics reveal the same changes in boreal watersheds across an

602    environmental gradient. Scientific Reports 7, 1–11. https://doi.org/10.1038/s41598-

603    017-13157-x.

604  Fan, J.T., Semenzin, E., Meng, W., Giubilato, E., Zhang, Y., Critto, A., Zabeo, A., Zhou,

605    Y., Ding, S. and Wan, J., 2015. Ecological status classification of the Taizi River

606    Basin, China: a comparison of integrated risk assessment approaches. Environmental

607    Science and Pollution Research 22 (19), 14738–14754.

608    https://doi.org/10.1007/s11356-015-4629-x.

609    Fernandez, S., Rodriguez, S., Martinez, J.L., Borrell, Y.J., Ardura, A. and Garcia-

610    Vazquez, E., 2018. Evaluating freshwater macroinvertebrates from eDNA

611    metabarcoding: a river Nalón case study. PLoS One 13 (8), e0201741.

612    https://doi.org/10.1371/journal.pone.0201741.

613    Franzo, A. and Del Negro, P., 2019. Functional diversity of free-living nematodes in

614    river lagoons: can biological traits analysis (BTA) integrate traditional taxonomic-

615    based approaches as a monitoring tool? Marine Environmental Research 145, 164–

616    176. https://doi.org/10.1016/j.marenvres.2019.02.015.

617    Frontalini, F., Greco, M., Di Bella, L., Lejzerowicz, F., Reo, E., Caruso, A., Cosentino,

618    C., Maccotta, A., Scopelliti, G. and Nardelli, M.P., 2018. Assessing the effect of

619    mercury pollution on cultured benthic foraminifera community using morphological

620    and eDNA metabarcoding approaches. Marine Pollution Bulletin 129 (2), 512–524.

621    https://doi.org/10.1016/j.marpolbul.2017.10.022.

622    He, X., Sutherland, T.F., Pawlowski, J. and Abbott, C.L., 2019. Responses of

623    foraminifera communities to aquaculture-derived organic enrichment as revealed by

624    environmental DNA metabarcoding. Molecular Ecology 28 (5), 1138–1153.

625    https://doi.org/10.1111/mec.15007.

626    Hering, D., Borja, A., Jones, J.I., Pont, D., Boets, P., Bouchez, A., Bruce, K., Drakare,

627    S., Hänfling, B. and Kahlert, M., 2018. Implementation options for DNA-based

628    identification into ecological status assessment under the European Water

629    Framework Directive. Water Research 138, 192–205.

630    https://doi.org/10.1016/j.watres.2018.03.003.

631    Hu, H.J. and Wei, Y.X. (2006) The Freshwater Algae of China, Systematics, Taxonomy

632    and Ecology, Science Press, Beijing.

633    Humbert, J.-F., Quiblier, C. and Gugger, M., 2010. Molecular approaches for

634        monitoring potentially toxic marine and freshwater phytoplankton species.

635        Analytical    and    Bioanalytical    Chemistry    397    (5),    1723–1732.

636        https://doi.org/10.1007/s00216-010-3642-7.

637    Jarman, S.N., Berry, O. and Bunce, M., 2018. The value of environmental DNA

638        biobanking for long-term biomonitoring. Nature Ecology & Evolution 2 (8), 1192–

639        1193. https://doi.org/10.1038/s41559-018-0614-3.

640    Ji, Y.Q., Ashton, L., Pedley, S.M., Edwards, D.P., Tang, Y., Nakamura, A., Kitching, R.,

641        Dolman, P.M., Woodcock, P. and Edwards, F.A., 2013. Reliable, verifiable and

642        efficient monitoring of biodiversity via metabarcoding. Ecology Letters 16 (10),

643        1245–1257. https://doi.org/10.1111/ele.12162.

644    Jones, F.C., Plewes, R., Murison, L., MacDougall, M.J., Sinclair, S., Davies, C., Bailey,

645        J.L., Richardson, M. and Gunn, J., 2017. Random forests as cumulative effects

646        models: a case study of lakes and rivers in Muskoka, Canada. Journal of

647        Environmental        Management        201,        407–424.

648        https://doi.org/10.1016/j.jenvman.2017.06.011.

649    Kefford, B.J., Marchant, R., Schäfer, R.B., Metzeling, L., Dunlop, J.E., Choy, S.C. and

650        Goonan, P., 2011. The definition of species richness used by species sensitivity

651        distributions approximates observed effects of salinity on stream macroinvertebrates.

652        Environmental        Pollution        159        (1),        302–310.

653        https://doi.org/10.1016/j.envpol.2010.08.025.

654    Keylock, C.J., 2005. Simpson diversity and the Shannon–Wiener index as special cases

655        of a generalized entropy. Oikos 109 (1), 203–207. https://doi.org/10.1111/j.0030-

656        1299.2005.13735.x.

657    Kottuparambil, S., Kim, Y.J., Choi, H., Kim, M.S., Park, A., Park, J., Shin, W. and Han,

658        T., 2014. A rapid phenol toxicity test based on photosynthesis and movement of the

659        freshwater flagellate, Euglena agilis Carter. Aquatic Toxicology 155, 9–14.

660        https://doi.org/10.1016/j.aquatox.2014.05.014.

661    Laroche, O., Wood, S.A., Tremblay, L.A., Ellis, J.I., Lejzerowicz, F., Pawlowski, J.,

662        Lear, G., Atalah, J. and Pochon, X., 2016. First evaluation of foraminiferal

663        metabarcoding for monitoring environmental impact from an offshore oil drilling

664        site.    Marine    Environmental    Research    120,    225–235.

665        https://doi.org/10.1016/j.marenvres.2016.08.009.

666    Li, F., Peng, Y., Fang, W., Altermatt, F., Xie, Y., Yang, J. and Zhang, X., 2018a.

667        Application of environmental DNA metabarcoding for predicting anthropogenic

668        pollution in rivers. Environmental Science & Technology 52 (20), 11708–11719.

669        https://doi.org/10.1021/acs.est.8b03869.

670    Li, M., Shan, X.J., Wang, W.J., Dai, F.Q., Lu, D. and Wu, H.H., 2018b. Study on the

671        retention time of environmental DNA of fenneropenaeus chinensis in water. Progress

672        in Fishery Sciences, https://doi.org/10.19663/j.issn2095-9869.20180906005.

673    Lu, X.J., Ming, L., Liu, W.B. and Li, H.X., 2018. Probabilistic regularized extreme

674        learning machine for robust modeling of noise data. IEEE Transactions on

675        Cybernetics 48 (8), 2368–2377. https://doi.org/10.1109/TCYB.2017.2738060.

676    Marshall, S. and Elliott, M., 1998. Environmental influences on the fish assemblage of

677        the Humber estuary, UK. Estuarine, Coastal and Shelf Science 46 (2), 175–184.

678        https://doi.org/10.1006/ecss.1997.0268.

679    Marvuglia, A., Kanevski, M. and Benetto, E., 2015. Machine learning for toxicity

680        characterization of organic chemical emissions using USEtox database: learning the

681        structure of the input space. Environment International 83, 72–85.

682 https://doi.org/10.1016/j.envint.2015.05.011.

683 Meador, M.R. and Goldstein, R.M., 2003. Assessing water quality at large geographic

684 scales: relations among land use, water physicochemistry, riparian condition, and

685 fish community structure. Environmental Management 31 (4), 0504–0517.

686 https://doi.org/10.1007/s00267-002-2805-5.

687 Meissner, T., Sures, B. and Feld, C.K., 2019. Multiple stressors and the role of

688 hydrology on benthic invertebrates in mountainous streams. Science of the Total

689 Environment 663, 841–851. https://doi.org/10.1016/j.scitotenv.2019.01.288.

690 Mineau, M.M., Wollheim, W.M. and Stewart, R.J., 2015. An index to characterize the

691 spatial distribution of land use within watersheds and implications for river network

692 nutrient removal and export. Geophysical Research Letters 42 (16), 6688–6695.

693 Ministry of Water Resources, Quality Standards for Surface Water Resources.

694 http://www.mwr.gov.cn/zwgk/zfxxgkml/201301/t20130124_964289.html. (1994,

695 accessed 24 January 2013)

696 Mize, E.L., Erickson, R.A., Merkes, C.M., Berndt, N., Bockrath, K., Credico, J.,

697 Grueneis, N., Merry, J., Mosel, K. and Tuttle-Lau, M., 2019. Refinement of eDNA

698 as an early monitoring tool at the landscape-level: Study design considerations.

699 Ecological Applications 29 (6), e01951. https://doi.org/10.1002/eap.1951.

700 Mori, N., Debeljak, B., Skerjanec, M., Simcic, T., Kanduc, T. and Brancelj, A., 2019.

701 Modelling the effects of multiple stressors on respiration and microbial biomass in

702 the hyporheic zone using decision trees. Water Research 149, 9–20.

703 https://doi.org/10.1016/j.watres.2018.10.093.

704 Olson, J.R., 2018. Predicting combined effects of land use and climate change on river

705 and stream salinity. Philosophical Transactions of the Royal Society B: Biological

706 Sciences 374 (1764), 20180005. https://doi.org/doi:10.1098/rstb.2018.0005.

707     Park, Y., Cho, K.H., Park, J., Cha, S.M. and Kim, J.H., 2015. Development of early-

708       warning protocol for predicting chlorophyll-a concentration using machine learning

709       models in freshwater and estuarine reservoirs, Korea. Science of the Total

710       Environment 502, 31–41. https://doi.org/10.1016/j.scitotenv.2014.09.005.

711     Pawlowski, J., Kelly-Quinn, M., Altermatt, F., Apothéloz-Perret-Gentil, L., Beja, P.,

712       Boggero, A., Borja, A., Bouchez, A., Cordier, T. and Domaizon, I., 2018. The future

713       of biotic indices in the ecogenomic era: Integrating (e) DNA metabarcoding in

714       biological assessment of aquatic ecosystems. Science of the Total Environment 637,

715       1295–1310. https://doi.org/10.1016/j.scitotenv.2018.05.002.

716     Piette, E.R. and Moore, J.H., 2018. Improving machine learning reproducibility in

717       genetic association studies with proportional instance cross validation (PICV).

718       BioData Mining 11, 6. https://doi.org/10.1186/s13040-018-0167-7.

719     Piggott, J.J., Lange, K., Townsend, C.R. and Matthaei, C.D., 2012. Multiple stressors

720       in agricultural streams: a mesocosm study of interactions among raised water

721       temperature, sediment addition and nutrient enrichment. PLoS One 7 (11), e49873.

722       https://doi.org/10.1371/journal.pone.0049873.

723     Piggott, J.J., Townsend, C.R. and Matthaei, C.D., 2015. Reconceptualizing synergism

724       and antagonism among multiple stressors. Ecology and Evolution 5 (7), 1538–1547.

725       https://doi.org/10.1002/ece3.1465.

726     Romero, F., Sabater, S., Timoner, X. and Acuña, V., 2018. Multistressor effects on river

727       biofilms under global change conditions. Science of the Total Environment 627, 1–

728       10. https://doi.org/10.1016/j.scitotenv.2018.01.161.

729     Rossel, R.A.V. and Behrens, T., 2010. Using data mining to model and interpret soil

730       diffuse reflectance spectra. Geoderma 158 (1–2), 46–54.

731       https://doi.org/0.1016/j.geoderma.2009.12.025.

732     Schloss, P.D. and Handelsman, J., 2005. Introducing DOTUR, a computer program for

733       defining operational taxonomic units and estimating species richness. Applied and

734       Environmental       Microbiology       71       (3),       1501–1506.

735       https://doi.org/10.1128/AEM.71.3.1501-1506.2005.

736     Strong, W.L., 2016. Biased richness and evenness relationships within Shannon–

737       Wiener       index       values.       Ecological       Indicators       67,       703–713.

738       https://doi.org/10.1016/j.ecolind.2016.03.043.

739     Torija, A.J. and Ruiz, D.P., 2015. A general procedure to generate models for urban

740       environmental-noise pollution using feature selection and machine learning methods.

741       Science of the Total Environment 505, 680–693.

742     Treusch, A. H., Demirhilton, E., Vergin, K. L., Worden, A. Z., Carlson, C. A., Donatz,

743       M. G., ... & Giovannoni, S. J., 2012. Phytoplankton distribution patterns in the

744       northwestern Sargasso Sea revealed by small subunit rRNA genes from plastids. The

745       ISME Journal, 6(3), 481-492. https://doi.org/10.1038/ismej.2011.117

746     Van Duin, E.H.S., Blom, G., Los, F.J., Maffione, R., Zimmerman, R., Cerco, C.F.,

747       Dortch, M. and Best, E.P.H., 2001. Modeling underwater light climate in relation to

748       sedimentation, resuspension, water quality and autotrophic growth. Hydrobiologia

749       444 (1–3), 25–42. https://doi.org/10.1023/A:1017512614680.

750     Vapnik, V.N., 1999. An overview of statistical learning theory. IEEE Transactions on

751       Neural Networks 10 (5), 988–999. https://doi.org/10.1109/72.788640.

752     Varoquaux, G., 2018. Cross-validation failure: small sample sizes lead to large error

753       bars.       Neuroimage       180       (Part       A),       68–77.

754       https://doi.org/10.1016/j.neuroimage.2017.06.061.

755     Visco, J.A., Apotheloz-Perret-Gentil, L., Cordonier, A., Esling, P., Pillet, L. and

756       Pawlowski, J., 2015. Environmental monitoring: inferring the diatom index from

757    next-generation sequencing data. Environmental Science & Technology 49 (13),

758    7597–7605. https://doi.org/10.1021/es506158m.

759    Von Schiller, D., Acuna, V., Aristi, I., Arroita, M., Basaguren, A., Bellin, A., Boyero, L.,

760    Butturini, A., Ginebreda, A. and Kalogianni, E., 2017. River ecosystem processes: a

761    synthesis of approaches, criteria of use and sensitivity to environmental stressors.

762    Science of the Total Environment 596, 465–480.

763    https://doi.org/10.1016/j.scitotenv.2017.04.081.

764    Were, K., Bui, D.T., Dick, O.B. and Singh, B.R., 2015. A comparative assessment of

765    support vector regression, artificial neural networks, and random forests for

766    predicting and mapping soil organic carbon stocks across an Afromontane landscape.

767    Ecological Indicators 52, 394–403. https://doi.org/10.1016/j.ecolind.2014.12.028.

768    Wilhm, J. and McClintock, N., 1978. Dissolved oxygen concentration and diversity of

769    benthic macroinvertebrates in an artificially destratified lake. Hydrobiologia 57 (2),

770    163–166. https://doi.org/10.1007/BF00016460.

771    Yanko, V., Ahmad, M. and Kaminski, M., 1998. Morphological deformities of benthic

772    foraminiferal tests in response to pollution by heavy metals: implications for

773    pollution monitoring. The Journal of Foraminiferal Research 28 (3), 177–200.

774    Yuan, J., Li, M. and Lin, S., 2015. An improved DNA extraction method for efficient

775    and quantitative recovery of phytoplankton diversity in natural assemblages. PLoS

776    One 10 (7), e0133060. https://doi.org/10.1371/journal.pone.0133060.

777    Zhang, L., Calvo-Bado, L., Murray, A.K., Amos, G.C.A., Hawkey, P.M., Wellington,

778    E.M. and Gaze, W.H., 2019. Novel clinically relevant antibiotic resistance genes

779    associated with sewage sludge and industrial waste streams revealed by functional

780    metagenomic screening. Environment International 132, 105120.

781    https://doi.org/10.1016/j.envint.2019.105120.

782    Zhang, X., 2019. Environmental DNA shaping a new era of ecotoxicological research.

783    Environmental     Science     &     Technology     53     (10),     5605–5612.

784    https://doi.org/10.1021/acs.est.8b06631.

785

786

787