

# Linking phylogenetic similarity and pollution sensitivity to develop ecological assessment methods: a test with river diatoms

François Keck      Agnès Bouchez      Alain Franc      Frédéric Rimet

February 5, 2016

## Abstract

1. Diatoms include a great diversity of taxa and are recognized as powerful bioindicators of freshwater quality. However using diatoms for bioassessment is costly and time consuming, because most of the indices necessitate species-level identification. Simplifying diatoms-based assessment protocols has focused the attention of water-managers and researchers in recent years.
2. The increasing availability of genomic data and phylogenies can benefit in the development of bioassessment methods making use of these tools, where a clade plays the role of a species if relevant. Indeed, the null hypothesis is that closely related species are more likely to exhibit similar environmental sensitivity because of phylogenetic constraints and inheritance. Such patterns have been reported recently for sensitivity to a variety of pollutants for two important groups of bioindicators used for freshwater monitoring: benthic macroinvertebrates and diatoms.
3. We introduce a method to extract clusters of species sharing similar traits and being phylogenetically related. We apply this method on the general pollution sensitivity (IPS specific sensitivity value; Coste, 1982) of 262 species of diatoms and, by tuning the method settings; we generate different clade-based derivatives of the traditional IPS index.
4. Finally, we estimate traditional and derived IPS scores for 2119 natural communities of diatoms in eastern France to compare and assess the performances of these new indices.
5. *Synthesis and applications.* We show that phylogenetic approaches offer a scope for simplification without an important loss of information and we discuss the potential of their use in biomonitoring.

This is a post-print version of an article originally published in *Journal of Applied Ecology*. ([link editor](#)). Please cite: Keck F., Bouchez A., Franc A. & Rimet F. (In press) Linking phylogenetic similarity and pollution sensitivity to develop ecological assessment methods: a test with river diatoms. *Journal of Applied Ecology*.

# Contents

<b>1 Introduction</b>	<b>2</b>
<b>2 Material and Methods</b>	<b>4</b>
2.1 Phylogenetic tree reconstruction . . . . .	4
2.2 Phylogenetically constrained clustering . . . . .	4
2.3 Defining new indices based on phylogenetic clusters . . . . .	4
2.4 Developing IPS <sub>P</sub> indices . . . . .	6
2.5 Comparing IPS <sub>P</sub> indices performances . . . . .	6
2.6 Statistical Packages . . . . .	7
<b>3 Results</b>	<b>7</b>
<b>4 Discussion</b>	<b>7</b>
4.1 Phylogenetic clustering – methodological discussion . . . . .	7
4.2 Phylogenetically based indices – potential for applications . . . . .	13
<b>5 Acknowledgments</b>	<b>14</b>
<b>References</b>	<b>14</b>

## 1 Introduction

Diatoms have been traditionally recognized as a good candidate group to monitor freshwater ecosystems, because they exhibit an important diversity and their community composition is strongly structured by numerous environmental factors including growth stimulating nutrients (Patrick, 1961; Lange-Bertalot, 1979). From this premise, the first environmental quality indices, based on diatoms assemblages, were developed about 50 years ago (*e.g.* Zelinka and Marvan, 1961) and nowadays diatoms are part of routine bioassessment standard methods in freshwater monitoring (Stevenson et al., 2010).

With an estimation of 100 000 extant species, diatoms constitute one of the most diverse algal classes (Mann and Vanormelingen, 2013). Taxonomic diversity is important for biomonitoring, because it promotes assemblage diversity and allows ecological assessment at a fine level (Birks, 2010). However, the extreme diversity of diatoms also constitutes a challenge for applied biomonitoring. Indices are traditionally developed by skilled diatomists and are usually derived at species-level to maximize their performance. Moreover, they may include several hundreds of species. Extending the use of such complex protocols – for example at a national network scale – is costly and requires training of many operators with continuous intercalibration (Prygiel et al., 2002; Kahlert et al., 2009). In addition, there is still the risk of imprecise or wrong identifications, which can lead to a biased estimation of environmental quality and ultimately lead managers to take unsatisfactory decisions (Besse-Lototskaya et al., 2006).

Simplifying and standardizing diatoms-based assessment protocols focused the attention of many researchers in recent years. Two main pathways have been explored: (i) reducing the number of species included in the indices by focusing on

most abundant and key species (Lenoir and Coste, 1996; Lavoie et al., 2009) and (ii) reducing the taxonomic resolution (Kelly et al., 1995; Chessman et al., 1999; Grouns, 1999; Hill et al., 2001; Wunsam et al., 2002; Raunio and Soinenen, 2007; Rimet and Bouchez, 2012).

With the increasing availability of genetic data and phylogenies (Benson et al., 2008; Wheeler et al., 2008), the idea arose that the development of bioassessment methods could also benefit from phylogenetic statistical approaches. Carew et al. (2011) first formulated the concept of phylogenetic redundancy in freshwater monitoring by analyzing links between mayflies and chironomids pollution sensitivity and phylogeny. The central idea is that closely related species are more likely to exhibit similar sensitivity because of phylogenetic constraints and inheritance. This hypothesis is commonly tested in the literature by measuring and testing the presence of the phylogenetic signal (“the tendency for related species to resemble each other more than they resemble species drawn at random from the tree”; Blomberg and Garland, 2002). The presence of such a signal may have direct consequences on biomonitoring, because it opens up interesting possibilities of simplification by using larger clades instead of species. Interestingly, the phylogenetic signal has been assessed for sensitivity to pollution on two important groups of bioindicators used for freshwater monitoring: benthic macroinvertebrates and diatoms (Ibáñez et al., 2010). Phylogenetic signal has been found significant for macroinvertebrates sensitivity to various metals (Buchwalter et al., 2008; Poteat et al., 2013; Poteat and Buchwalter, 2014) and to general pollutants (Carew et al., 2011). For diatoms, significant phylogenetic signal was found for sensitivities to different herbicides (Larras et al., 2014) and for general ecological preferences (Keck et al., 2016).

Demonstrating the presence of phylogenetic signal is essential, but is only the first step in making proposals for biomonitoring tools based on phylogenetic knowledge. A second step is to develop methods to extract informative groups of species based on phylogenetic signal to derive simpler indices and test their ability to predict environment quality. Thus, we introduce a simple distance-based method to extract clusters of species sharing similar traits, but also are phylogenetically related. It is classical in ecology to compare two distances within a set of individuals, typically through a Mantel test to compare phenetic or genetic distance with geographic distance (Sokal, 1979; Fortin and Gurevitch, 2001; Vignieri, 2005). Here, we go one step further, by building clusters of species based both on traits values and phylogenetic proximity meaning that, two distantly related species cannot be included in the same cluster even if they exhibit similar trait values.

In this paper, we apply this method to get different sets of clusters from 262 diatom species. Clustering is based on the phylogeny and on the general pollution sensitivity of the species (IPS specific sensitivity value; Coste, 1982). We use these sets of clusters to develop derivatives of the traditional IPS index. Finally we estimate traditional and derived IPS scores of 2119 samples to compare and assess the performances of these new indices.

## 2 Material and Methods

### 2.1 Phylogenetic tree reconstruction

We used the phylogenetic tree reconstructed in Keck et al. (2016). This phylogeny is based both on the nuclear gene coding for the small subunit 18S rRNA and the chloroplast *rbcL* gene coding for the RuBisCO enzyme. The tree includes 549 diatoms species for which genetic information is available for at least one of these markers. Reconstruction was done with RAxML 7.2.8 (Stamatakis, 2006) using a partitioned Maximum Likelihood analysis with a GTR+I+G evolutionary model (see Keck et al., 2016, for details). The tree was dated in relative time using a semi-parametric method based on penalized likelihood (Sanderson, 2002).

### 2.2 Phylogenetically constrained clustering

We present here a simple co-clustering method for a set of  $n$  species in a phylogeny with one or more associated trait values (as illustrated in Figure 1A). The method is based both on the pairwise trait distance matrix  $\mathbf{T}$  and the pairwise phylogenetic distance matrix  $\mathbf{P}$ . We consider the graph  $G = (V, E)$  where  $V$  denotes the vertices (the species) and  $E$  the set of edges connecting the vertices. Here, the graph  $G$  is defined by its adjacency binary matrix  $\mathbf{A}$ , an  $n \times n$  matrix where  $\mathbf{A}_{ij} = 1$  if there is an edge joining species  $i$  with species  $j$  and  $\mathbf{A}_{ij} = 0$  otherwise. A variety of rules can be used to decide whether there is an edge or not between two vertices. Here, we propose a linear rule given in Equation 1, for which a graphical illustration is given in Figure 1B.

$$A_{ij} = \begin{cases} 1 & \text{if } -\frac{t}{p}\mathbf{P}_{ij} \geq \mathbf{T}_{ij}; i \neq j \\ 0 & \text{otherwise} \end{cases} \quad (1)$$

where  $t$  and  $p$  are the upper bounds to be considered for respectively trait and phylogenetic distances (see Figure 1B) and must be manually set. Thus, the higher the values of  $t$  and  $p$ , the lower the trait and phylogenetic constraints are, respectively.

Once the adjacency matrix is given, we compute the connected components of the associated graph, which define the clusters (Figure 1C). Note, however, that different strategies are possible (*e.g.* selection of cliques, use of community detection algorithms) which will be discussed later.

### 2.3 Defining new indices based on phylogenetic clusters

We chose to work with the IPS index (indice de polluo-sensibilité; Coste, 1982). The IPS index is a weighted average autecological index based on a modified version of the Zelinka and Marvan (1961) equation (Equation 2) where  $a_i$  is the proportional abundance of the taxon  $i$ ,  $v_i$  is its indicator value and  $s_i$  its pollution sensitivity.

$$IPS = \frac{\sum_{i=1}^n a_i \times v_i \times s_i}{\sum_{i=1}^n a_i \times v_i} \quad (2)$$

We then have simplified the index by defining some clusters as explained above, and averaging values of  $v_i$  and  $s_i$  and summing the values of  $a_i$ , over clusters of

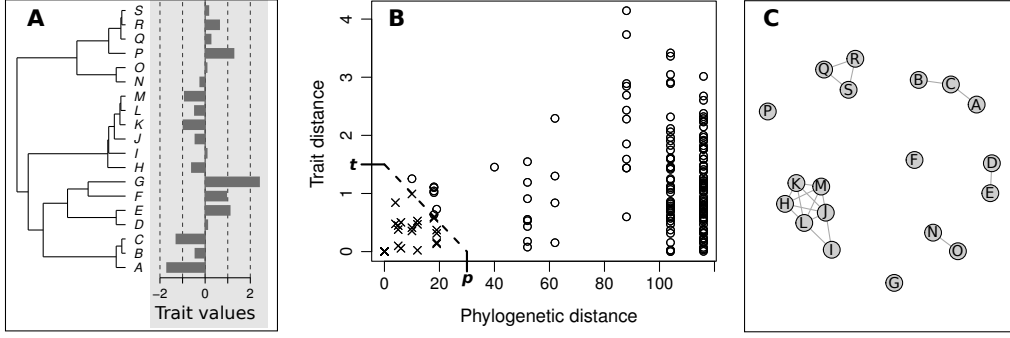


Figure 1: Phylogenetically constrained clustering process. **A.** The process is illustrated with a dataset of 19 species (identified by letters A–S). The trait data are simulated under a Brownian Motion model of trait evolution and are centred. **B.** Pairs of species in function of their phylogenetic (patristic) distance and trait (Euclidean) distance. The selected pairs (following Equation 1,  $p$  and  $t$  values) are represented with crosses while non-selected pairs are represented with circles. The dashed line illustrates the selection limit. **C.** A graph where species are connected according to previously selected pairs, unveiling 8 disconnected components (clusters).

species. Let us denote by  $\gamma$  a cluster. Then, aggregated IPS index is defined as in Equation 3.

$$IPS_P = \frac{\sum_{\gamma} a_{\gamma} \times v_{\gamma} \times s_{\gamma}}{\sum_{\gamma} a_{\gamma} \times v_{\gamma}} \quad (3)$$

where

$$a_{\gamma} = \sum_{i \in \gamma} a_i, \quad v_{\gamma} = \frac{1}{n_{\gamma}} \sum_{i \in \gamma} v_i, \quad s_{\gamma} = \frac{1}{n_{\gamma}} \sum_{i \in \gamma} s_i, \quad n_{\gamma} = \#\{i : i \in \gamma\}$$

As the grain for species sensitivity variation is coarser (in  $IPS_P$ , all species belonging to the same cluster are assumed to share a same value for  $s$  and  $v$ ), the estimate of IPS will be coarser. In order to evaluate the discrepancy between IPS and  $IPS_P$ , we have calculated the error made by using  $IPS_P$  instead of IPS, by a first order development of  $IPS_P$ . We show now that the discrepancy induced by using the coarse index is minimized when the clusters are such that the discrepancy between average value and individual values within each cluster are bounded from above. This is precisely the role of  $t$  in the above calculation. Indeed, the error made by using  $IPS_P$  instead of IPS depends on two types of terms given in Equation 4 (see Appendix S1 for details).

$$\Delta_{v,i} = \sum_{i \in \gamma} a_i (v_i - v_{\gamma}), \quad \Delta_{s,i} = \sum_{i \in \gamma} a_i (s_i - s_{\gamma}) \quad (4)$$

Each of these terms is a combination of two terms: the abundances in the environmental sample  $a_i$ , and the discrepancy between the species  $s_i$  and  $v_i$  and the

cluster  $s_\gamma$  and  $v_\gamma$  it belongs to. Each term remains small, *i.e.* the approximation is acceptable, if (i) the species is ill positioned in a group, *i.e.* the term  $|x_i - x_\gamma|$  is high, with  $x = s$  or  $x = v$  has a low abundance, or (ii) the discrepancy is acceptable. Let us note that for each cluster  $\gamma$ , we have  $\sum_{i \in \gamma} (s_i - s_\gamma) = \sum_{i \in \gamma} (v_i - v_\gamma) = 0$ , which means that the error terms are expected to be small. The detailed calculations are available in Appendix S1.

## 2.4 Developing $\text{IPSP}$ indices

We carried out clustering analyses using the method described above for species for which both phylogenetic and IPS data ( $s$  and  $v$  values) were available. IPS data were retrieved from OMNIDIA (Lecointe et al., 1993). We used a phylogenetic distance matrix  $\mathbf{P}$ , based on the number of nodes separating two species  $i$  and  $j$  and a trait distance matrix  $\mathbf{T}$ , based on the pairwise Euclidean distance of IPS pollution sensitivity ( $\mathbf{T}_{ij} = \sqrt{(s_i - s_j)^2}$ ). To make things more interpretable, both  $\mathbf{P}$  and  $\mathbf{T}$  were divided by their respective maximum values so that all distances range between 0 and 1.

Since there is no rule to set  $t$  and  $p$  values, we tested different settings. A full grid of  $10^4$  combinations of  $t$  and  $p = \{0.01, 0.02, 0.03, \dots, 0.99, 1\}$  was processed. However, for clarity, we report results for a set of representative combinations of  $t = \{0.2, 0.4, 0.6, 0.8\}$  and  $p = \{0.05, 0.1, 0.15\}$  giving 12 different graphs and as many different sets of clusters. Then, we developed a series of phylogenetically IPS-derived indices using Equation 3 (referred as  $\text{IPSP}_{[t,p]}$  with  $t$  and  $p$  indicating the trait and phylogenetic constraints applied for the clustering).

## 2.5 Comparing $\text{IPSP}$ indices performances

To assess the performances of  $\text{IPSP}$  indices we used a database of 2119 diatom community samples collected in rivers and streams in eastern France between 2001 and 2008. For each of them, 400 diatom frustules were counted and identified at species level. Details about this database are given in Rimet and Bouchez (2012). These count data were used to compute the  $\text{IPSP}_{\text{standard}}$  value of each sample, which constitutes the reference index value. Thus we can compute different statistics to compare the ability of the different  $\text{IPSP}$  to recover the information contained in  $\text{IPSP}_{\text{standard}}$ . First, the Pearson correlation index is computed as a measure of the dependence between the samples scores as estimated by  $\text{IPSP}$  and by  $\text{IPSP}_{\text{standard}}$ . Second, the residual sum of square (RSS) is used as a measure of the discrepancy between the scores of  $\text{IPSP}$  and the scores of  $\text{IPSP}_{\text{standard}}$ . Finally, it is common to use IPS scores to classify samples in 5 levels of water quality:  $[0; 7[$  = Very Poor;  $[7; 11[$  = Poor;  $[11; 13.5[$  = Fair;  $[13.5; 16[$  = Good;  $[16; 20]$  = Very Good (Prygiel et al., 1996). In particular, these thresholds are currently used by managers to take decisions for environmental restoration. We reported the percentage of good classification and percentages of misclassification (over and under estimates) of samples by  $\text{IPSP}$  compared to  $\text{IPSP}_{\text{standard}}$ .

## 2.6 Statistical Packages

We performed all the statistical analyses with R 3.0.2 software (R Development Core Team, 2013). Phylogenies were handled with the `ape` package (Paradis et al., 2004) and the `phylobase` package (Hackathon et al. 2013). Phylogenetic distances were computed with the `adephylo` package (Jombart et al., 2010). Phylogenetic clustering was performed with the `phylosignal` package <sup>1</sup>.

## 3 Results

The dataset includes 262 taxa which were found both in the phylogenetic tree and the IPS database (Figure 2). The full grid approach generated  $10^4$  sets of cluster. We investigated the effects of phylogenetic and trait constraints on the number of clusters produced (Figure 3A) and the relationship between the number of cluster of an  $IPSP$  index and its ability for sample classification (Figure 3B).

The subset of 12 combinations of  $t$  and  $p$  values tested produced contrasting sets of clusters. The most restrictive graph ( $t = 0.2$ ,  $p = 0.05$ ; *i.e.* low trait and phylogenetic distances) is composed of 196 connected components (*i.e.* clusters) while the most relaxed graph ( $t = 0.8$  and  $p = 0.15$ ; *i.e.* high trait and phylogenetic distances) is composed of 9 components. Other graphs have a number of connected components ranging between these two extremes (Table 1).

Since the number of cluster on which they are based varies greatly, the capacity of the different  $IPSP$  indices to reflect the information of  $IPSP_{\text{standard}}$  varies also markedly. The correlation between  $IPSP$  indices and  $IPSP_{\text{standard}}$  is high ( $> 0.9$ ) as long as the number of clusters remains high ( $> 68$ ; Table 1 and Figure 4). The highest correlation (0.938) is achieved with  $IPSP_{[0.2,0.1]}$  and  $IPSP_{[0.2,0.15]}$  (*i.e.* low trait and moderate to high phylogenetic distances).

The residual sum square (RSS) ranged between 2990 and 4357 as long as the number of clusters remains high ( $> 68$ ). Under this threshold, the error increases drastically (Table 1 and Figure 4). The lowest RSS is achieved with  $IPSP_{[0.2,0.1]}$  (2990).

More than 73% of the samples are correctly classified as long as the number of clusters remains above 68. For indices based on very few clusters (16 and 9), the number of misclassified samples falls under the number of correctly classified samples. The best percentage of classification is achieved with  $IPSP_{[0.2,0.1]}$  (80.6% of good classification). Strong overestimations of water quality ( $\geq 2$  classes) are rare overall while strong underestimations ( $\geq 2$  classes) appear to be more frequent when the number of clusters is very low. Overall, in case of misclassification,  $IPSP$  indices are more likely to underestimate water quality than overestimate it.

## 4 Discussion

### 4.1 Phylogenetic clustering – methodological discussion

The idea of simplifying bioassessment methods using phylogenetics has been raised in the last few years (Carew et al., 2011; Larras et al., 2014), but no study pro-

---

<sup>1</sup><https://cran.r-project.org/web/packages/phylosignal/>



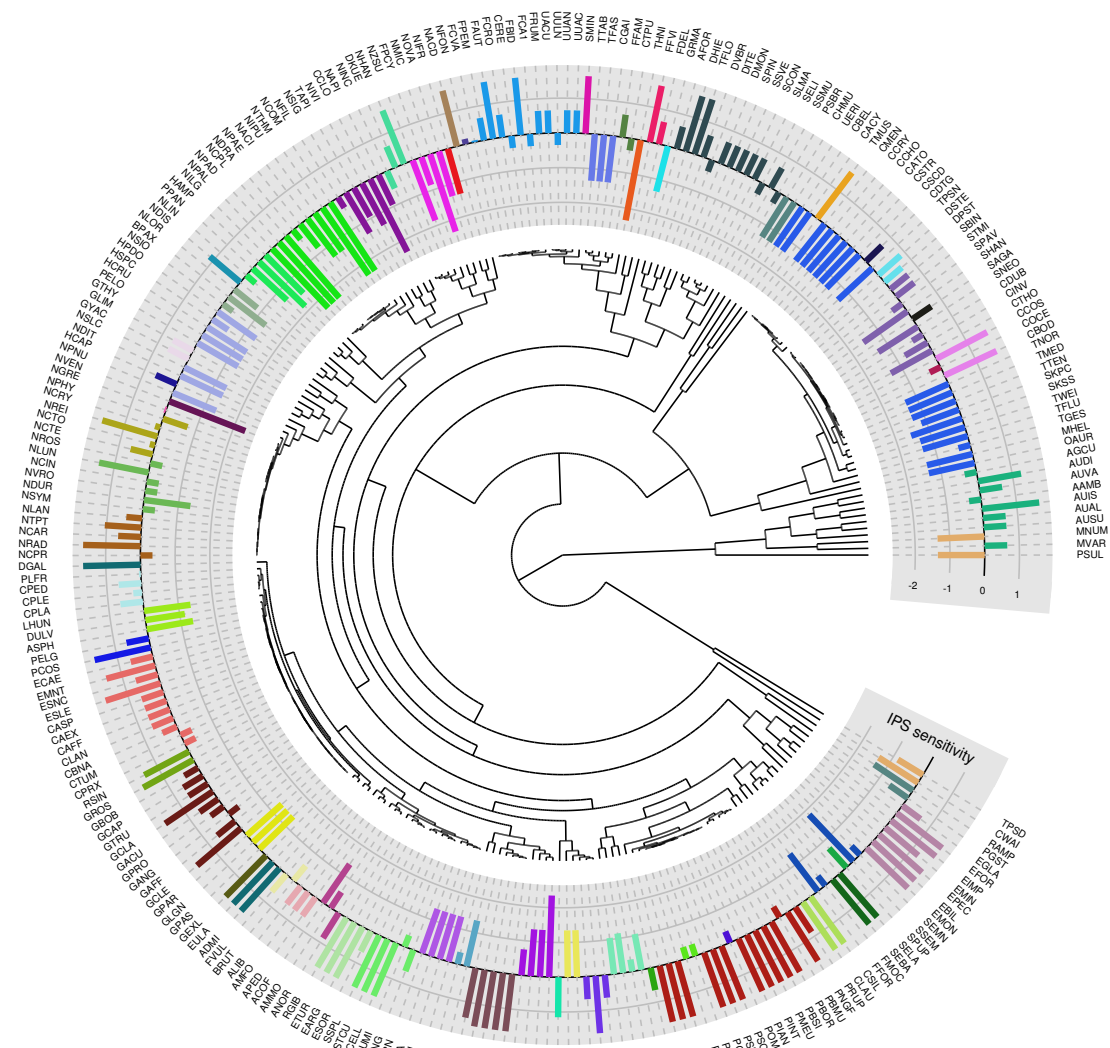


Figure 2: Phylogenetic tree of 262 diatoms species and their respective  $IPS_{\text{standard}}$  sensitivity value ( $s$ ). The colors delineate 68 clusters based on  $t = 0.6$  and  $p = 0.1$ . Diatoms names are reported using 4-letter codes (Lecointe et al., 1993, see Table S1 for corresponding Linnaean names).



Clustering constraints					Quality classification (% of samples) compared to $IPS_{\text{standard}}$					
					Overestimate		Exact	Underestimate		
$t$	$p$	Number of clusters	Correlation $IPS_{\text{standard}}$	RSS $IPS_{\text{standard}}$	2	1	0	1	2	3
0.2	0.05	196	0.935	3099	0.1	9.7	79.5	10	0.5	0.1
0.4	0.05	187	0.936	3081	0.1	8.2	79.8	11.2	0.5	0.1
0.6	0.05	157	0.929	3335	0.2	9.5	77.5	12.3	0.4	0.1
0.8	0.05	153	0.925	3989	0.2	6.5	73.9	18.1	1.2	0.1
0.2	0.1	126	0.938	2990	0.1	8	80.6	10.7	0.5	0.1
0.4	0.1	89	0.928	3594	0.2	7.4	76.1	15.4	0.8	0
0.2	0.15	86	0.938	3179	0.1	7.5	77.3	14.2	0.7	0.1
0.6	0.1	68	0.907	4357	0.4	9.9	73.4	15	1.2	0
0.8	0.1	51	0.863	7922	0.7	9	63.2	25	2	0
0.4	0.15	32	0.904	8684	0.1	2.2	55.3	39	3.1	0.2
0.6	0.15	16	0.774	17317	0.2	5.3	37.5	48.1	8.3	0.6
0.8	0.15	9	0.591	31120	0.9	5.9	23.5	37.4	30.2	2.1

Table 1: Comparison of the 12  $IPS_p$  indices. Each index is based on a set of clusters generated by a pair of  $t$  and  $p$  values. Performances of the indices are assessed by comparing with results of  $IPS_{\text{standard}}$  for 2119 diatom samples.

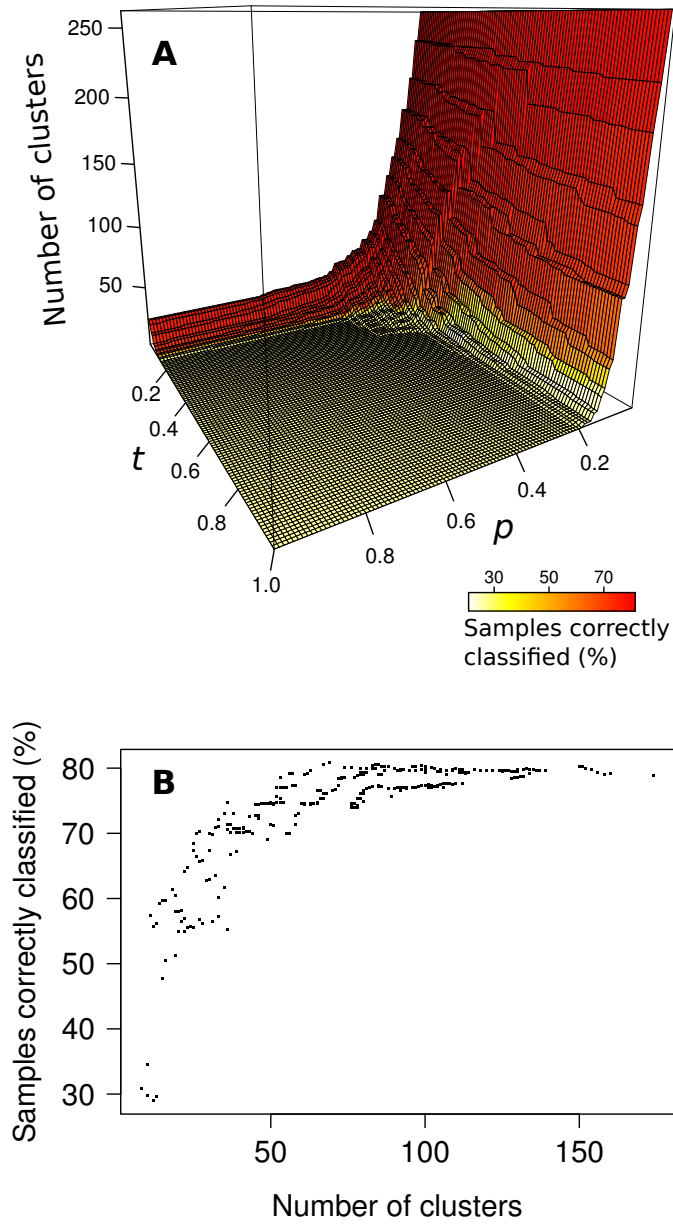


Figure 3: **A.** Relation between the pair of  $t$  and  $p$  values and the number of clusters produced by the method. The color gradient indicates the percentage of samples correctly classified by the  $\text{IPSP}_p$  developed from the corresponding set of cluster. **B.** Relation between the number of clusters produced by the method and the percentage of samples correctly classified by the  $\text{IPSP}_p$  developed from the corresponding set of cluster. Data presented only for  $p < 0.2$ .

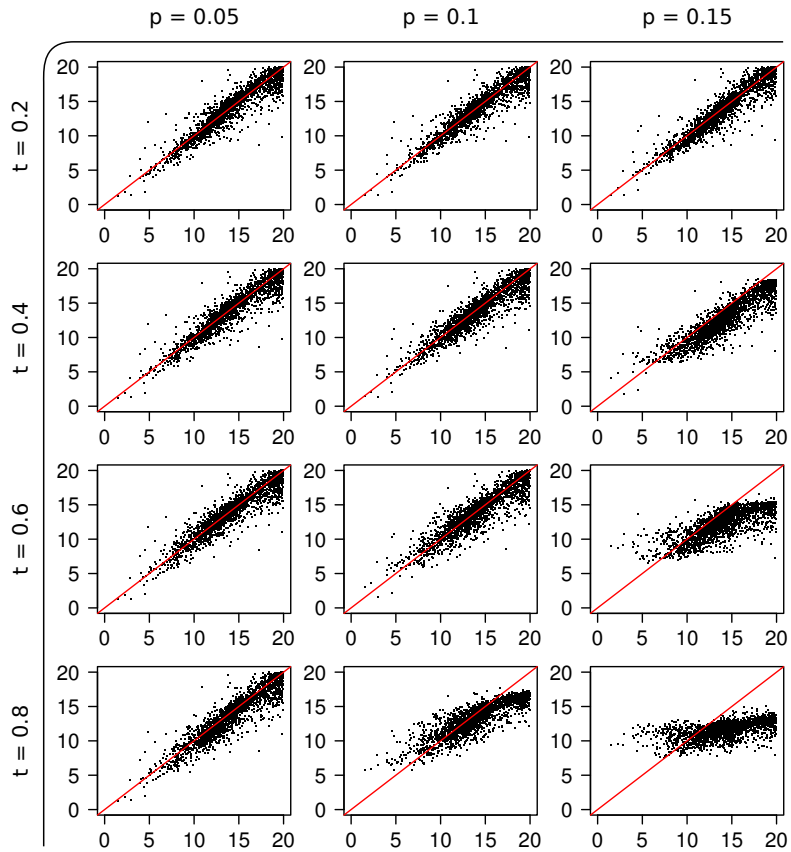


Figure 4: Relation between samples scores estimated with  $IPS_{\text{standard}}$  (horizontal axis) and  $IPS_{P[t,p]}$  (vertical axis) for the 20 tested pairs of  $t$  (trait constraint) and  $p$  (phylogenetic constraint) values. The solid red line represents the full equivalence between  $IPS_{\text{standard}}$  and  $IPS_P$ .

posed a phylogenetically based biomonitoring tool. We introduced a simple and general approach to develop such tools and tested it in order to simplify a popular biomonitoring diatomic index: the IPS (Coste, 1982).

We proposed a simple method, which allows clustering species taking into account both their phylogenetic proximities and trait similarities. Clusters generated by this method are not necessarily monophyletic clades. The method has many declinations possible, since each step is independent and adaptable. First, a different phylogenetic distance matrix ( $\mathbf{P}$ ) can be used. Here we used the number of internal nodes separating two species, but patristic distance (length of branches separating two species) or more complex distances (Pavoine et al., 2008; Pavoine and Ricotta, 2013) can be considered, as well as transformation of these distances (*e.g.* square root of patristic distance; Hardy and Pavoine, 2012). Second, we applied the method on a single trait (species sensitivity), but since clustering is based on the Euclidean distance of trait ( $\mathbf{T}$ ), it can be easily extended to a multivariate framework. Third, different rules can be used to select which pairs of species are connected by an edge in the graph. We used a simple rule based on a linear equation (Equation 1, Figure 1B), but other options can be developed (*e.g.* rectangular and elliptical selections are included in the R package `phylosignal`). Finally, different cluster extraction approaches can be tested. In particular, for complex data, clusters can be detected using community detection algorithms (Newman and Girvan, 2004) and clusters validity can be assessed with statistics derived from graph theory like measures of density and connectivity (Van Steen, 2010). Since the method is extremely general and flexible, this gives an opportunity to fit to a large variety of data. Clustering can be applied to any kind of trait. For example in freshwater biomonitoring, other indices could be clustered like the trophic diatom index (Kelly and Whitton, 1995), the global periphyton indices (Rott et al., 1997; Rott et al., 1999) or the Brettum index for lakes monitoring (Brettum, 1989), but the method could also be applied directly on species preferences through a multivariate approach (see Keck et al., 2016).

The method does not provide an optimal pair of  $p$  and  $t$  constraining values. This can limit the ease of use, but is also a source of flexibility. Since the clustering algorithm we propose is not computationally intensive, it can be easy to test thousands of settings. Thus, a practitioner developing a new index can pick up the pair of phylogenetic and trait constraints which fit the best with his/her own needs in terms of trade-off between simplification and precision. Representing the relationship between the number of clusters and the efficiency of indices (Figure 4) may be a good way to support the decision process.

Overall, the results must always be interpreted carefully and we stress the importance to make a detailed analysis of how  $t$  and  $p$  influence the clustering outcomes. An identified issue is the linkage effect: if there is an edge between species A and species B and an edge between species B and C, then A, B and C will be included in the same connected component (*i.e.* cluster), even if A and C are not connected. A way to overcome this problem might be to use more sophisticated method to extract clusters from the graph, as discussed above. Another point which needs attention is that if  $t$  or  $p$  values are too high, the method will converge to phylogenetic-only clustering (*i.e.* clusters strictly based on phylogenetic distances) or trait-only clustering (*i.e.* clusters strictly based on traits distances), respectively. For example

in our dataset, when the number of clusters is very low and the performance of the index is very high, this is due to the phylogenetic constraint, which is non-existent ( $p > 0.25$ ; high phylogenetic distance). Therefore, the results converge to a trait-only clustering, which is definitely not the aim here.

## 4.2 Phylogenetically based indices – potential for applications

The tests we conducted showed that the number of clusters can be reduced without an important loss of information. These results tend to confirm that biomonitoring with diatoms can be simplified using taxonomic levels higher than the species level as previously suggested by other authors (Kelly et al., 1995; Chessman et al., 1999; Grouns, 1999; Hill et al., 2001; Wunsam et al., 2002; Raunio and Soininen, 2007; Rimet and Bouchez, 2012). This is achieved for the first time using a phylogenetic approach in order to take account phylogenetic redundancy (Carew et al., 2011).

It is important to note that the phylogeny of diatoms is far from complete. Only 262 species have been included in the clustering method, whereas IPS computation is based on more than 5000 species with 909 of them present in the samples of our test dataset. Including more species in the phylogenetic tree could produce more clusters, but also it will probably increase significantly the performance of IPS<sub>P</sub> indices. In particular, some missing taxa are important for biomonitoring like *Achnantheidium subatomus*, *A. subatomoides*, *A. daonense*, which are indicators of pristine rivers of relatively low conductivity. These missing species can probably explain the tendency of IPS<sub>P</sub> to underestimate the water quality. On the other hand, pollution tolerant taxa are better represented in the current phylogeny. Significant progress has been made in our understanding of diatom phylogeny in recent years (Theriot et al., 2010; Theriot et al., 2011; Medlin, 2011). Large scale phylogenetic trees, including many more species and based on many more markers, will be made available, making phylogenetic approaches more robust and relevant. Biomonitoring methods based on phylogenies can be easily updated as new data are made available.

The use of phylogenetic approaches aims principally to simplify biomonitoring by avoiding phylogenetic redundancy (Carew et al., 2011). In this paper, we try to address this issue by extracting clusters of phylogenetically related species sharing similar pollution sensitivities. Ideally, these clusters would provide the best compromise between simplicity and efficiency. However, it seems difficult to import a tool developed with phylogeny in the traditional biomonitoring workflows based on classical microscopic counts, because there are several incongruencies between morphology and DNA-phylogeny (Kermarrec et al., 2011; Zimmermann et al., 2014). Moreover taxonomical classification is rarely matching the IPS<sub>P</sub> clusters proposed, and it would ask the technician counting diatoms to learn these new clusters even if many of them are intuitive (*e.g.* *Ulnaria* group, *Nitzschia lanceolatae* group). This is probably hardly applicable for already trained diatomists. However, in some cases, clusters match the taxonomy, especially at genus level. For example, in Figure 2, the genera *Eunotia* and *Stauroneis* are identified as two clusters with high sensitivity while all *Entomoneis* species are detected in the same non sensitive (low sensitivity values) cluster. Such results can be interesting to develop biomonitoring tools based on a mixture of taxonomical levels as suggested by Jones (2008) for macroinvertebrates. For diatoms, tools based both on species and genus levels exist (Kelly and Whitton, 1995), but could undergo new developments with phylogenetic approaches.

Finally, these approaches seem to be much better adapted for next generation biomonitoring, so-called biomonitoring 2.0, based on metabarcoding and high throughput sequencing methods (Baird and Hajibabaei, 2012), which aims to use DNA-barcodes to assess environmental quality. Since this approach is based on molecular characters, it is much more straightforward to integrate phylogenetic considerations. In metabarcoding, one of the difficulties is the taxonomic assignment of metabarcode sequences (Coissac et al., 2012). Assigning DNA sequences to clusters of species, rather than species would be more flexible and probably would be achieved more easily. Another common issue is the lack of data in taxon-stressor response libraries. The use of phylogenetic methods to infer taxa traits from their phylogenetic position could offer a solution to this problem (Keck et al., 2016) and a complete modeling framework has been proposed by Guénard et al. (2013). However, as a first step in a biomonitoring context, it would be simple to infer trait values of unknown sampled species if they fall within a given cluster. Thus, increasing information on traits and taxa – thanks to metabarcoding associated together with phylogenetically based methods – should significantly enhance the efficiency of environmental monitoring.

## 5 Acknowledgments

This work was funded by ONEMA (French National Office for Water and Aquatic Ecosystems) in the context of the 2013-2015 “Phylogeny and Bioassessment” program.

## References

- Baird, D. J. and M. Hajibabaei (2012). “Biomonitoring 2.0: a new paradigm in ecosystem assessment made possible by next-generation DNA sequencing”. In: *Molecular Ecology* 21.8, pp. 2039–2044.
- Benson, D. A., I. Karsch-Mizrachi, D. J. Lipman, J. Ostell, and D. L. Wheeler (2008). “GenBank”. In: *Nucleic Acids Research* 36 (Database Issue), pp. D25–D30.
- Besse-Lototskaya, A., P. F. M. Verdonschot, and J. A. Sinkeldam (2006). “Uncertainty in diatom assessment: sampling, identification and counting variation”. In: *Hydrobiologia* 566.1, pp. 247–260.
- Birks, H. J. B. (2010). “Numerical methods for the analysis of diatom assemblage data”. In: *The Diatoms: Applications for the Environmental and Earth Sciences*. Ed. by J. P. Smol and E. F. Stoermer. 2nd. Cambridge, UK: Cambridge University Press, pp. 23–54.
- Blomberg, S. P. and T. Garland (2002). “Tempo and mode in evolution: phylogenetic inertia, adaptation and comparative methods”. In: *Journal of Evolutionary Biology* 15.6, pp. 899–910.
- Brettum, P. (1989). *Algen als Indikatoren für die Gewässerqualität in norwegischen Binnenseen*. Norway: Norsk Institutt for vannforskning (NIVA), p. 102.
- Buchwalter, D. B., D. J. Cain, C. A. Martin, L. Xie, S. N. Luoma, and T. Garland Jr (2008). “Aquatic insect ecophysiological traits reveal phylogenetically based differences in dissolved cadmium susceptibility”. In: *Proceedings of the National Academy of Sciences* 105.24, pp. 8321–8326.

- Carew, M. E., A. D. Miller, and A. A. Hoffmann (2011). “Phylogenetic signals and ecotoxicological responses: potential implications for aquatic biomonitoring”. In: *Ecotoxicology* 20.3, pp. 595–606.
- Chessman, B., I. Growns, J. Currey, and N. Plunkett-Cole (1999). “Predicting diatom communities at the genus level for the rapid biological assessment of rivers”. In: *Freshwater Biology* 41.2, pp. 317–331.
- Coissac, E., T. Riaz, and N. Puillandre (2012). “Bioinformatic challenges for DNA metabarcoding of plants and animals”. In: *Molecular Ecology* 21.8, pp. 1834–1847.
- Coste, M. (1982). *Étude des méthodes biologiques d’appréciation quantitative de la qualité des eaux*. Cemagref, p. 218.
- Fortin, M.-J. and J. Gurevitch (2001). “Mantel tests: spatial structure in field experiments”. In: *Design and Analysis of Ecological Experiments*. Ed. by S. M. Scheiner and J. Gurevitch. 2nd. New York: Oxford University Press, pp. 308–326.
- Growns, I. (1999). “Is genus or species identification of periphytic diatoms required to determine the impacts of river regulation?” In: *Journal of Applied Phycology* 11.3, pp. 273–283.
- Guénard, G., P. Legendre, and P. Peres-Neto (2013). “Phylogenetic eigenvector maps: a framework to model and predict species traits”. In: *Methods in Ecology and Evolution* 4.12, pp. 1120–1131.
- Hackathon et al. (2013). *phylobase: Base package for phylogenetic structures and comparative data*. Version 0.6.5.2.
- Hardy, O. J. and S. Pavoine (2012). “Assessing phylogenetic signal with measurement error: A comparison of Mantel tests, Blomberg et al.’s K, and phylogenetic distograms”. In: *Evolution* 66.8, pp. 2614–2621.
- Hill, B. H., R. J. Stevenson, Y. Pan, A. T. Herlihy, P. R. Kaufmann, and C. B. Johnson (2001). “Comparison of correlations between environmental characteristics and stream diatom assemblages characterized at genus and species levels”. In: *Journal of the North American Benthological Society* 20.2, pp. 299–310.
- Ibáñez, C., N. Caiola, P. Sharpe, and R. Trobajo (2010). “Ecological indicators to assess the health of river ecosystems”. In: *Handbook of Ecological Indicators for Assessment of Ecosystem Health*. Ed. by S. E. Jørgensen, L. Xu, and R. Costanza. 2nd. Boca Raton, Florida: CRC Press, pp. 447–464.
- Jombart, T., F. Balloux, and S. Dray (2010). “adephylo: new tools for investigating the phylogenetic signal in biological traits”. In: *Bioinformatics* 26.15, pp. 1907–1909.
- Jones, F. C. (2008). “Taxonomic sufficiency: the influence of taxonomic resolution on freshwater bioassessments using benthic macroinvertebrates”. In: *Environmental Reviews* 16, pp. 45–69.
- Kahlert, M. et al. (2009). “Harmonization is more important than experience—results of the first Nordic–Baltic diatom intercalibration exercise 2007 (stream monitoring)”. In: *Journal of Applied Phycology* 21.4, pp. 471–482.
- Keck, F., F. Rimet, A. Franc, and A. Bouchez (2016). “Phylogenetic signal in diatom ecology: perspectives for aquatic ecosystems biomonitoring”. In: *Ecological Applications*.



- Kelly, M. G., C. J. Penny, and B. A. Whitton (1995). "Comparative performance of benthic diatom indices used to assess river water quality". In: *Hydrobiologia* 302.3, pp. 179–188.
- Kelly, M. G. and B. A. Whitton (1995). "The trophic diatom index: a new index for monitoring eutrophication in rivers". In: *Journal of Applied Phycology* 7.4, pp. 433–444.
- Kermarrec, L., L. Ector, A. Bouchez, F. Rimet, and L. Hoffmann (2011). "A preliminary phylogenetic analysis of the Cymbellales based on 18S rDNA gene sequencing". In: *Diatom Research* 26.3, pp. 305–315.
- Lange-Bertalot, H. (1979). "Pollution tolerance of diatoms as a criterion for water quality estimation". In: *Nova Hedwigia* 64, pp. 285–304.
- Larras, F., F. Keck, B. Montuelle, F. Rimet, and A. Bouchez (2014). "Linking Diatom Sensitivity to Herbicides to Phylogeny: A Step Forward for Biomonitoring?" In: *Environmental Science & Technology* 48.3, pp. 1921–1930.
- Lavoie, I., P. J. Dillon, and S. Campeau (2009). "The effect of excluding diatom taxa and reducing taxonomic resolution on multivariate analyses and stream bioassessment". In: *Ecological Indicators* 9.2, pp. 213–225.
- Lecointe, C., M. Coste, and J. Prygiel (1993). "'Omnidia': software for taxonomy, calculation of diatom indices and inventories management". In: *Hydrobiologia* 269-270.1, pp. 509–513.
- Lenoir, A. and M. Coste (1996). "Development of a practical diatom index of overall water quality applicable to the French National Water Board Network". In: *Use of Algae for Monitoring Rivers II*. International symposium, Volksbildungsheim Grilhof Vill, AUT, 17-19 September 1995. Universität Innsbruck: Whitton, B.A., Rott, E., Eds., pp. 29–43.
- Mann, D. G. and P. Vanormelingen (2013). "An inordinate fondness? The number, distributions, and origins of diatom species". In: *Journal of Eukaryotic Microbiology* 60.4, pp. 414–420.
- Medlin, L. K. (2011). "A review of the evolution of the diatoms from the origin of the lineage to their populations". In: *The Diatom World*. New York, USA: Seckbach, J., Kociolek, P., Eds., pp. 93–118.
- Newman, M. E. J. and M. Girvan (2004). "Finding and evaluating community structure in networks". In: *Physical review E* 69.2, p. 026113.
- Paradis, E., J. Claude, and K. Strimmer (2004). "APE: analyses of phylogenetics and evolution in R language". In: *Bioinformatics* 20.2, pp. 289–290.
- Patrick, R. (1961). "A study of the numbers and kinds of species found in rivers in eastern United States". In: *Proceedings of the Academy of Natural Sciences of Philadelphia* 113.10, pp. 215–258.
- Pavoine, S. and C. Ricotta (2013). "Testing for Phylogenetic Signal in Biological Traits: The Ubiquity of Cross-Product Statistics". In: *Evolution* 67.3, pp. 828–840.
- Pavoine, S., S. Ollier, D. Pontier, and D. Chessel (2008). "Testing for phylogenetic signal in phenotypic traits: new matrices of phylogenetic proximities". In: *Theoretical Population Biology* 73.1, pp. 79–91.
- Poteat, M. D. and D. B. Buchwalter (2014). "Phylogeny and Size Differentially Influence Dissolved Cd and Zn Bioaccumulation Parameters among Closely Related Aquatic Insects". In: *Environmental Science & Technology* 48.9, pp. 5274–5281.

- Poteat, M. D., T. Garland, N. S. Fisher, W.-X. Wang, and D. B. Buchwalter (2013). “Evolutionary Patterns in Trace Metal (Cd and Zn) Efflux Capacity in Aquatic Organisms”. In: *Environmental Science & Technology* 47.14, pp. 7989–7995.
- Prygiel, J., P. Carpentier, S. Almeida, M. Coste, J.-C. Druart, L. Ector, D. Guillard, M.-A. Honoré, R. Iserentant, and P. Ledeganck (2002). “Determination of the biological diatom index (IBD NF T 90–354): results of an intercomparison exercise”. In: *Journal of Applied Phycology* 14.1, pp. 27–39.
- Prygiel, J., L. Lévêque, and R. Iserentant (1996). “Un nouvel Indice Diatomique Pratique pour l’évaluation de la qualité des eaux en réseau de surveillance”. In: *Journal of Water Science* 9.1, pp. 97–113.
- Raunio, J. and J. Soininen (2007). “A practical and sensitive approach to large river periphyton monitoring: comparative performance of methods and taxonomic levels.” In: *Boreal environment research* 12.1.
- R Development Core Team (2013). *R: A Language and Environment for Statistical Computing*. Vienna, Austria: R Foundation for Statistical Computing.
- Rimet, F. and A. Bouchez (2012). “Biomonitoring river diatoms: Implications of taxonomic resolution”. In: *Ecological Indicators* 15.1, pp. 92–99.
- Rott, E., P. G. Hofmann, K. Pall, P. Pfister, and E. Pipp (1997). “Indikationslisten für Aufwuchsalgen in österreichischen Fließgewässern. Teil 1: Saprobienliste”. In: *Bundesministerium für Land-und Forstwirtschaft, Wasserwirtschaftskataster, Wien*.
- Rott, E., E. Pipp, P. Pfister, H. Van Dam, K. Ortler, N. Binder, and K. Pall (1999). “Indikationslisten für Aufwuchsalgen in österreichischen Fließgewässern. Teil 2: Trophie-indikation sowie geochemische Präferenz; taxonomische und toxikologische Anmerkungen”. In: *Bundesministerium für Land-und Forstwirtschaft, Wasserwirtschaftskataster, Wien*.
- Sanderson, M. J. (2002). “Estimating Absolute Rates of Molecular Evolution and Divergence Times: A Penalized Likelihood Approach”. In: *Molecular Biology and Evolution* 19.1, pp. 101–109.
- Sokal, R. R. (1979). “Testing Statistical Significance of Geographic Variation Patterns”. In: *Systematic Zoology* 28.2, pp. 227–232.
- Stamatakis, A. (2006). “RAxML-VI-HPC: maximum likelihood-based phylogenetic analyses with thousands of taxa and mixed models”. In: *Bioinformatics* 22.21, pp. 2688–2690.
- Stevenson, R. J., P. Yangdong, and H. Van Dam (2010). “Assessing environmental conditions in rivers and streams with diatoms”. In: *The Diatoms: Applications for the Environmental and Earth Sciences*. Ed. by J. P. Smol and E. F. Stoermer. 2nd ed. Cambridge University Press, pp. 55–85.
- Theriot, E. C., M. Ashworth, E. Ruck, T. Nakov, and R. K. Jansen (2010). “A preliminary multigene phylogeny of the diatoms (Bacillariophyta): challenges for future research”. In: *Plant Ecology and Evolution* 143.3, pp. 278–296.
- Theriot, E. C., E. Ruck, M. Ashworth, T. Nakov, and R. K. Jansen (2011). “Status of the pursuit of the diatom phylogeny: Are traditional views and new molecular paradigms really that different?” In: *The Diatom World*. Ed. by J. Seckbach and J. Kociolek. New York, USA: Springer, pp. 119–142.
- Van Steen, M. (2010). *Graph Theory and Complex Networks: An Introduction*. Maarten van Steen. 300 pp.

- Vignieri, S. N. (2005). "Streams over mountains: influence of riparian connectivity on gene flow in the Pacific jumping mouse (*Zapus trinotatus*)". In: *Molecular Ecology* 14.7, pp. 1925–1937.
- Wheeler, D. L. et al. (2008). "Database resources of the National Center for Biotechnology Information". In: *Nucleic Acids Research* 36 (Database Issue), pp. D13–D21.
- Wunsam, S., A. Cattaneo, and N. Bourassa (2002). "Comparing diatom species, genera and size in biomonitoring: a case study from streams in the Laurentians (Quebec, Canada)". In: *Freshwater Biology* 47.2, pp. 325–340.
- Zelinka, M. and P. Marvan (1961). "Zur präzisierung der biologischen klassifikation der reinheit fließender gewässer". In: *Archiv für Hydrobiologie* 57.3, pp. 389–407.
- Zimmermann, J., N. Abarca, N. Enk, O. Skibbe, W.-H. Kusber, and R. Jahn (2014). "Taxonomic Reference Libraries for Environmental Barcoding: A Best Practice Example from Diatom Research". In: *PLoS ONE* 9.9, e108793.