

Derivation of Typical Assemblages

Jonathan Jupke

Oktober 07, 2020

Contents

1	Introduction	3
2	Harmonizing taxa names	3
3	What is the optimal taxonomic level?	3
4	Can we represent the stream types with our data?	5
5	What is a typical Assemblage?	5
6	Redundancy between typical assemblages	10
7	Characteristics of typical assemblages	11
8	Seasonal Typical Assemblages	12
9	Patterns and overlap in seasonal assemblages	14
10	Notes for Traits	17
	References	17

```
# TODO discuss PCoA plots shortly
# TODO add taxa lists
# TODO add to intro
# TODO add image of map
# TODO add link to maps
# TODO add figures
# TODO based on what criteria where the seasonal communities selected.
```

1 Introduction

Here we describe the methods we used to derive typical assemblages of macro-invertebrates and diatoms for selected river types across Europe. We also describe and briefly discuss the results.

2 Harmonizing taxa names

International diatom occurrence data sets require extensive harmonization because of the taxonomic resolution differing between data sets, different working groups using different nomenclatures, identification errors, and ongoing changes to the accepted nomenclature (Kahlert *et al.* 2020). Harmonization can reduce overall taxonomic resolution but also improve the detection of large-scale spatio-temporal patterns (Lee *et al.* 2019). We compared all our data sets against a series of databases that contain accepted names, synonyms with links to the respective accepted names and suggestions for grouping contentious taxa in larger complexes. If a taxon name was found in one of the databases the name was accepted, changed into the accepted name in case it was a synonym, or grouped into the respective complex. Once a taxon was found in a database, it would not be included in queries of subsequent databases. However, if the accepted name differed from the original one, the accepted name would be queried through all previous databases again. The results were also controlled visually for consistency. The following databases were used in the same order:

1. Table S2 from (Kahlert *et al.* 2020)
2. The taxon list associated with the OMNIDA software (Lecointe *et al.* 1993)
3. The German list of freshwater organisms (Mauch *et al.* 2017)
4. The diat.barcode database (Rimet *et al.* 2019)
5. The website algaebase.org (Guiry 2020)
6. The global biodiversity information platform (gbif) (GBIF.org 2020)

The harmonization of macro-invertebrate data required less effort, and was achieved with gbif (GBIF.org 2020) through the taxize R package (Scott Chamberlain & Eduard Szocs 2013).

3 What is the optimal taxonomic level?

One result of the last progress review for Get Real (held on the 29.04.2020) was that taxa in the TA should be included on their respective optimal taxonomic level instead of using one level (e.g. Genus) for all. *Oligochaetes*, for example, are usually only determined to subclass level, which should not prevent them to be part of a TA in the case that they are common in a given river type. Thus, the question arises: given a data set, what is the optimal taxonomic level to represent a specific taxon?

To establish the optimal level, we used a hierarchical approach. First, we removed all

observations from Phyla and Classes that were not present in all data sets. We assumed that these represented differences in sampling rather than in communities. That left us with the classes Clitellata (Annelida), Insecta, Malacostraca (Arthropoda), Bivalvia and Gastropoda (Mollusca).

In the following, a higher taxonomic level refers to levels with higher resolution, i.e. species is the highest taxonomic level and kingdom the lowest. For each taxon, we calculated the percentage of observations that are represented at each higher level. For example, 4.12% of observations from the order *Lepidoptera* are at the species level, 74.77% at the genus level, 7.75% at the family level, and 13.35% at the order level. Now given a threshold X, which is to be determined, we would call a taxon optimally represented at a certain taxonomic level if less than X% are represented by higher levels. For example, *Lepidoptera* would be represented on order level if $X > 4.12 + 74.77 + 7.75 = 86.64\%$. As there are no theoretical grounds on which to base such a threshold value we searched for noticeable patterns in the data (Figure 1). The most noticeable jump occurs between 85 and 86%. It occurs because for $X > 86$ *Chironomidae* are represented at the family level. Hence, we used 85% as threshold. Observations that were missed by this procedure, e.g. observations of *Chironomidae* at the family level, were included at their respective level.

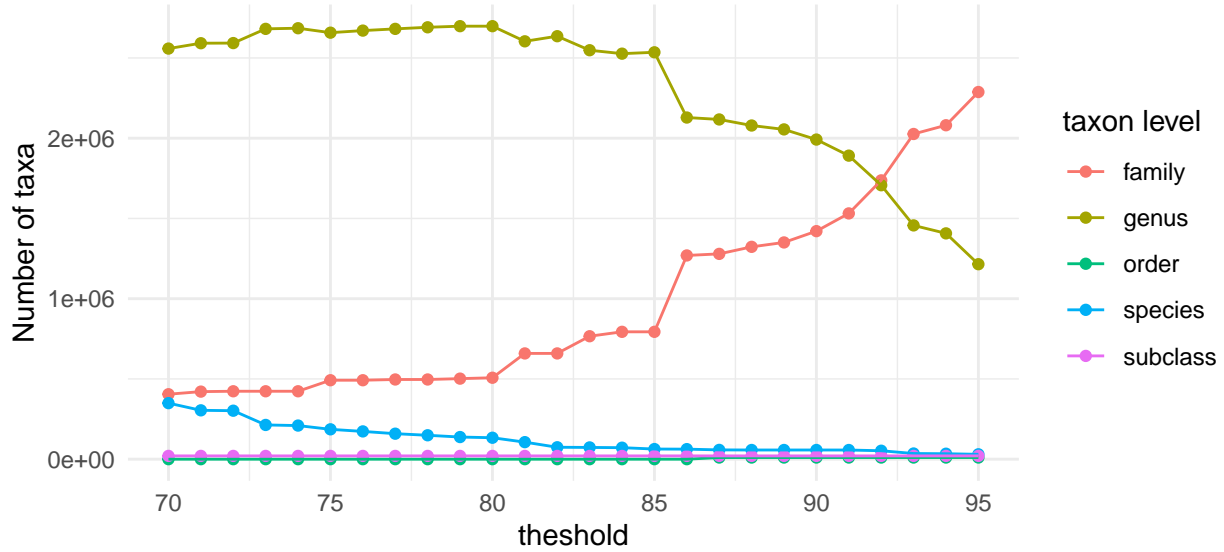


Figure 1: Number of taxa represented at a given taxonomic level as function of the threshold value

For the diatoms, we employed 75% as threshold, because Gomphonema, which is the fourth most common genus in our data set, had 81.43% observations at the species level. The taxonomic resolution was higher than in the macroinvertebrate data and the lowest resolution is the genus level.

Table 1: The ratings for all river types for macro-invertebrates and diatoms

Rating	Taxon	River.Types
high	macro-invertebrates	4, 5, 9, 10, 11, 12, 13, 16
high	diatoms	
medium	macro-invertebrates	1, 2, 3, 8, 14, 15, 18
medium	diatoms	1, 2, 3, 4, 5, 6, 8, 9, 12, 14, 16, 17, 18, 19
low	macro-invertebrates	6, 7, 17, 19, 20
low	diatoms	7, 10, 11, 13, 15, 20

4 Can we represent the stream types with our data?

We determined visually whether a our data set contained enough sampling sites in a given river type to derive meaningful TAs. The degree of representation for river type was graded in a three-tier system: High, medium, and low. A high degree of representation indicates, that we have many sampling locations, which are distributed across the instances of a river type which fall within the countries considered in GetReal. A low degree indicates the opposite, i.e. few and spatially clustered sites. A medium rating implies that we either have many sampling sites, but these only extend over parts of the countries or few sites that extend over most of the countries. The ratings for all river types for both macro-invertebrates and diatoms are shown in table 1.

For each river type x taxon combination, you can find maps with the associated sampling sites [for macro-invertebrates](#) and [for diatoms](#).

Further analyses were conducted for all stream types with a high or medium degree of representation. More information on the river types is available in Lyche Solheim *et al.* (2019). In general, we have fewer sampling sites for diatoms than for macro-invertebrates which entails that the representation of stream types is generally lower.

5 What is a typical Assemblage?

We derived TAs based on a rule that considered:

1. The probability of site x belonging to stream type z given that species y is present (a measure of specificity, henceforth **A**)
2. The probability of species y being present given that site x belongs to stream type z (a measure of commonness, henceforth **B**)
3. The Species Indicator Value

The Species Indicator Value (Dufrêne & Legendre 1997; Cáceres & Legendre 2009) is the

weighted product of **A** and **B** (see Equation (1))

$$\sqrt{Indval} = \sqrt{A_g \times B} = \sqrt{\frac{\frac{n_p}{N_p}}{\sum_{k=1}^K \frac{n_k}{N_k}} \times \frac{n_p}{N_p}} \quad (1)$$

where N_p is the number of sites that belong to river type p and n_p the number of occurrences of the focal species in sites of type p . K is the number of river types. **A** is weighted by the total number of occurrences to account for unequal sample sizes. The statistical significance of the Indicator Value can be assessed with permutation-based pseudo- p -values, which we did with 999 permutations. Here, we are not interested in indicator species for each community, but TAs. Hence, simply continuing with those species that have a pseudo- p -value below some significance level would not serve our purpose. A species that occurs at each site, across all stream types, highlights the difference: while it would not be indicative of any stream type (low specificity) it should be part of each TA. Hence, we need additional criteria to derive the TAs which can be based on **A**, **B**, and the pseudo- p -value of the indicator value. We used the following rules:

For macro-invertebrates:

Species were considered typical if **B** > 0.25 or **B** > 0.20 and $p < 0.05$ or **A** > 0.80

Genera were considered typical if **B** > 0.50 or **B** > 0.33 and $p < 0.05$ or **A** > 0.95

Families were considered typical if **B** > 0.95 or **B** > 0.80 and $p < 0.01$ or **A** > 0.99

For diatoms:

Species where considered typical if **B** > 0.4 or **B** > 0.3 and $p < 0.05$ or **A** > 0.7.

Genera where considered typical if **B** > 0.8 or **B** > 0.6 and $p < 0.05$ or **A** > 0.95.

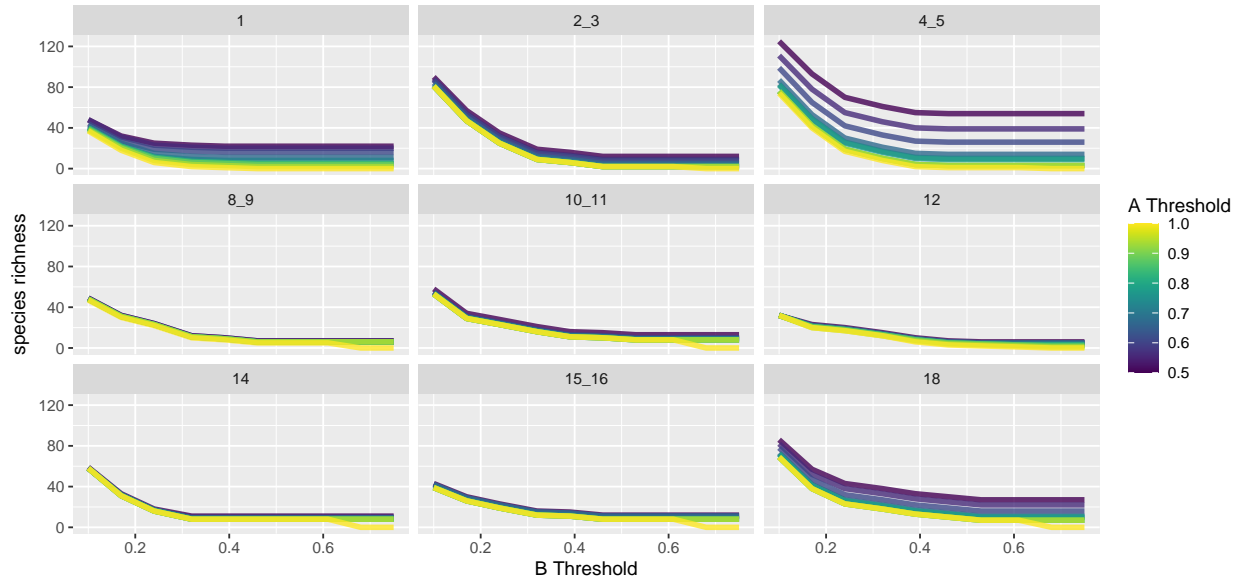
Note that there was no we did not systematically optimize these thresholds. Such procedures would require optimization criteria, but we are not aware of a criterion that would work in this context. We acknowledge that TA could be (i) very similar in composition or (ii) harbor strongly differing numbers of taxa. Thus, parametrizing the rules in a way that would (i) maximize dissimilarity between assemblages or (ii) maximize the mean assemblage richness would not lead to what we consider a typical assemblage. It would be possible to try a cross-validation-type approach where each taxon is scored based on the number of random-site-subsets it is included in, but such an approach would also entail making essentially arbitrary numerical assumptions. We think the use of subjectively defined thresholds is justified, as long as they are clearly and openly communicated, to be what we define as “typical assemblages”.

However, we conducted a sensitivity analysis to see how much varying the parameters of the rules would alter the results. We altered the threshold values of **A** and **B**. The rules above contain two distinct **B** Threshold: B_1 which does not consider the pseudo- p -value (**B** > 0.25 for macro-invertebrate species and **B** > 0.40 for diatom species) and B_2 which does take the pseudo- p -value into account ($p < 0.05$ and **B** > 0.2 macro-invertebrate species and **B** > 0.30

for diatom species). In the following simulations, the B_2 was always taken to be 25% below B_1 . Henceforth, when referring to the threshold for B , we refer to B_1 . For species, we varied the threshold for B in ten steps between 0.10 and 0.75 and that for \mathbf{A} in ten steps between 0.5 and 1.0. For lower taxonomic levels these thresholds were raised. For genera, the threshold values of \mathbf{A} and \mathbf{B} were raised by a factor of 1.25 and 2 respectively. All levels family and lower taxonomic levels were grouped in “families or lower” (fol). For fol, the thresholds were raised by factors of 1.5 and 3 respectively. The taxon richness and uniqueness scores of each TA were computed for all 100 combinations of these parameters and each taxonomic level. Please note that results are only shown and discussed for the non-redundant TAs (see section 6). Taxa richness decreased with increasing \mathbf{A} and \mathbf{B} threshold in macro-invertebrates and diatoms (Figure 2A and Figure 3A), while the uniqueness scores increased with \mathbf{B} thresholds but decreased with \mathbf{A} thresholds ((Figure 2B and Figure 3B)). Uniqueness scores decreased noticeably with very high \mathbf{A} thresholds (> 0.9), indicating that taxa that are specific to certain river types are an important driver of TA differentiation. Note that graphs are only shown for all taxa levels combined. Plots for each taxon level separately are available for [macro-invertebrates](#). However, the general patterns visible in Figure 2 and Figure 3, hold for them as well.

Species richness decreased with increasing \mathbf{A} and \mathbf{B} threshold (Figure 1 and Figure 2), while the uniqueness scores increased with \mathbf{B} thresholds but decreased with \mathbf{A} thresholds (Figure 3 and Figure 4). The rate of change in species richness along gradients in \mathbf{A} and \mathbf{B} threshold differed markedly between TAs but seemed to be correlated with overall species richness, i.e. more species-rich TA lost species more quickly than less species-rich ones. The TAs of RT03 and RT16 serve as examples at both extremes of our data set. Uniqueness scores decreased noticeably with very high \mathbf{A} thresholds (> 0.9), indicating that taxa that are specific to certain river types are an important driver of TA differentiation.

A Macro-invertebrates – Richness – All Levels



B Macro-invertebrates – Uniqueness – All levels

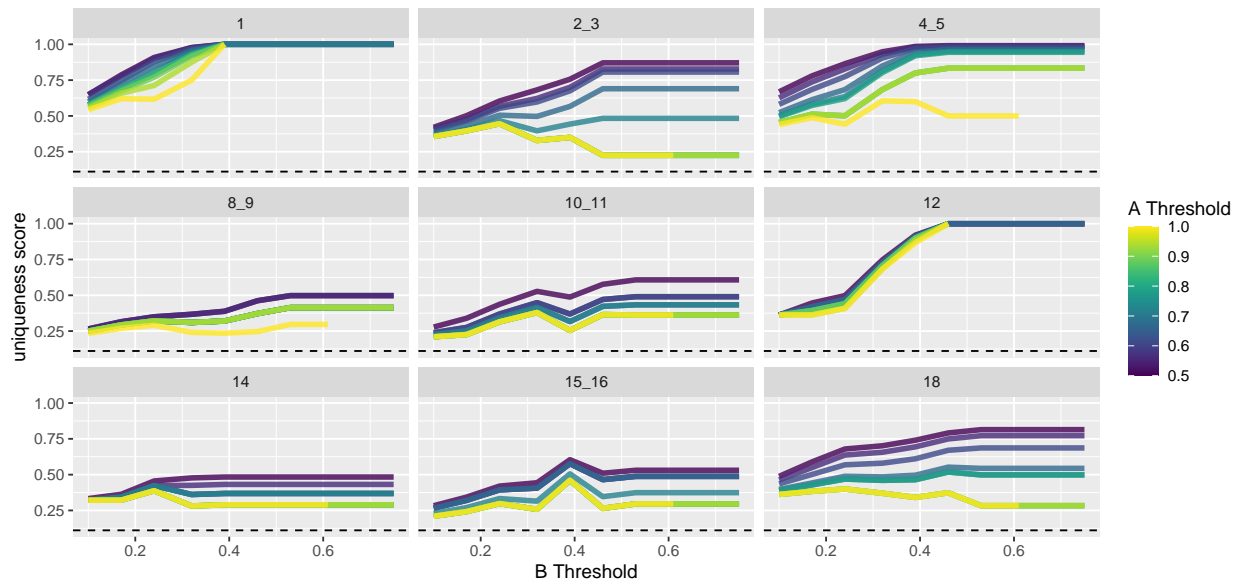
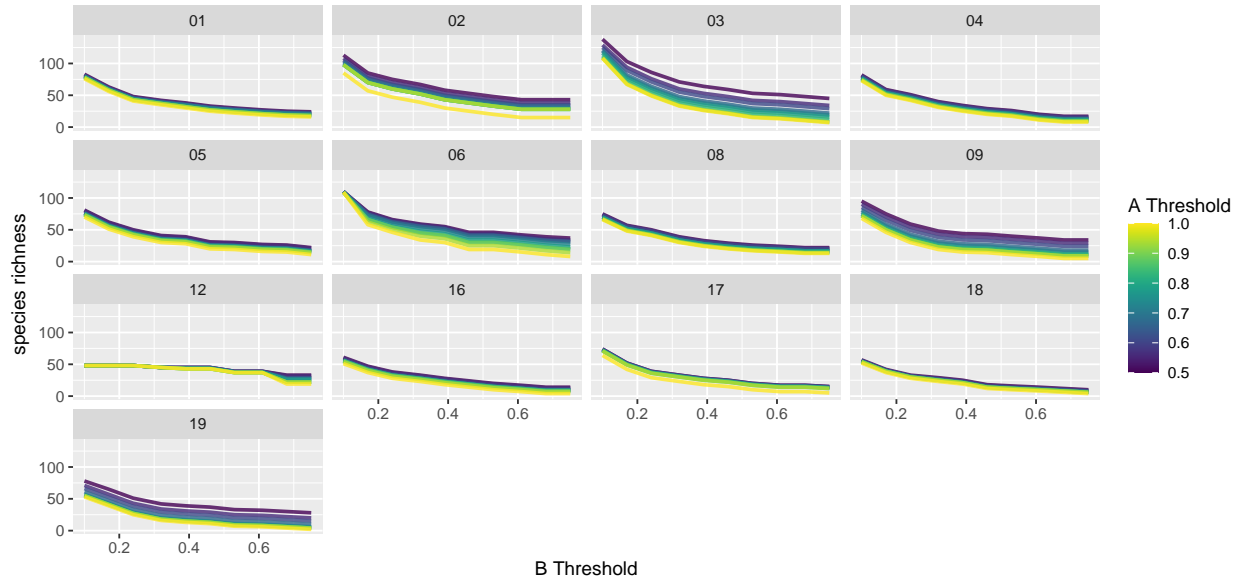


Figure 2: Changes in taxon richness along a changing B threshold. Line color indicates the A threshold

A Diatoms – Richness – All Levels



B Diatoms – Uniqueness – All levels

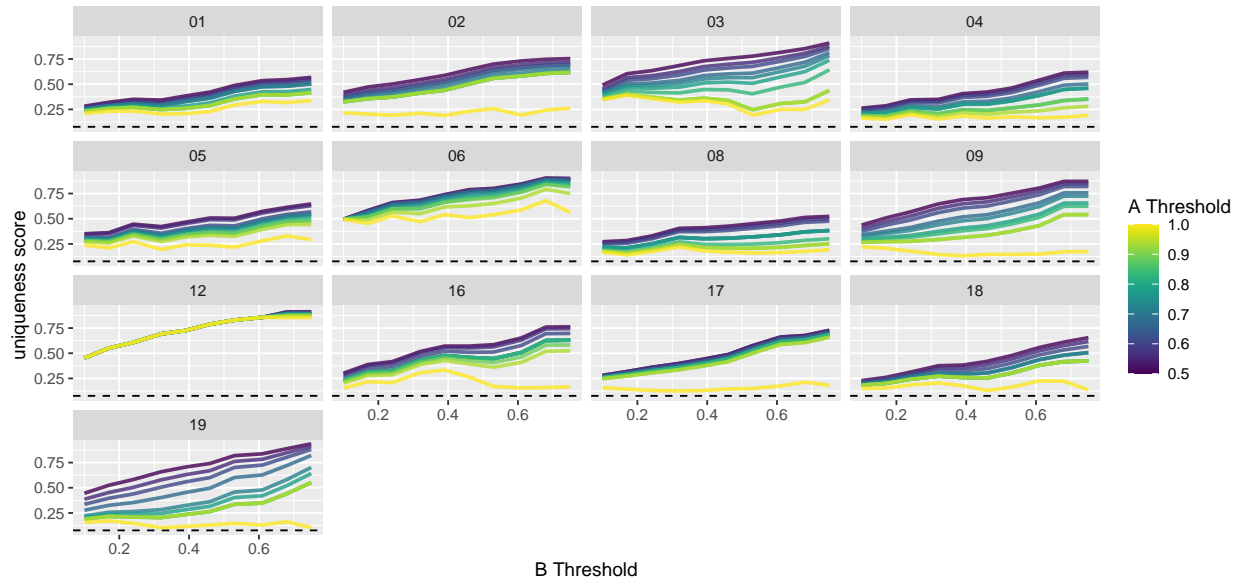


Figure 3: Changes in uniqueness score along a changing B threshold. Line color indicates the A threshold

Table 2: Overlap between different typical assemblages.

River.Type	Macroinvertebrate	Diatoms
RT1	RT2+4 (25%)	RT2 (62.9%)
RT2	RT3 (88.2%)	RT1 (46.8%)
RT3	RT2 (68.2%)	RT2 (48.7%)
RT4	RT3 (45.8%)	RT2 (74.1%)
RT5	RT4 (83.3%)	RT2 (62.5%)
RT6		RT12 (38.3%)
RT8	RT10 (77.8%)	RT2 (71.4%)
RT9	RT8 (76.5%)	RT5+8 (56%)
RT10	RT11+18 (65.2%)	
RT11	RT10 (88.2%)	
RT12	RT9 (50.0%)	RT6 (40.9%)
RT13	RT2 (87.5%)	
RT14	RT16+18 (69.2%)	
RT15	RT16 (85.7%)	
RT16	RT9+10+11+15 (57.1%)	RT18 (56.5%)
RT17		RT1 (63%)
RT18	RT10 (55.6%)	RT1 (71.4%)
RT19		RT17 (68.4%)

6 Redundancy between typical assemblages

We assessed to which degree the different TAs overlap (Table 2). The degree of overlap is the percentage of taxa in a TA that is also present in the most similar (largest overlap) TA. Again, choosing a threshold above which we consider two assemblages to be redundant is somewhat arbitrary. We proceeded with 75% but are open to other suggestions. This threshold leads to five redundant assemblages in macroinvertebrates and none in diatoms. Of the redundant TAs most belong to two river types that only differ in river size: RT02 and 03, 04 and 05, 08 and 09, as well as 10 and 11. The only exception is the combination of RT15 and 16. Both are high altitude river types that occur mainly in southern Europe, which differentiates them from the northern high altitude rivers in RT14. RT13 is also redundant with RT02 and 03 however joining it with these two river types led to a drastically reduced number of taxa in the TA, when compared to that of the combined river type RT02_03. Since RT13 represents an exceedingly rare river type we decided to omit it from the analysis and proceed with RT02_03 instead of RT02_03_13. The new TAs resulted in overall lower degrees of overlap, none of which exceeds the 75% threshold. The largest overlaps were between RT8_9 and RT10_11, with 70%.

7 Characteristics of typical assemblages

In all macro-invertebrate TAs, genus is the prevalent taxonomic level (figure 4 A). The numbers of species and fol are similar with both exceeding the other in four assemblages. The mean number of species was 3.2, mean number of genera 14.3, and the mean number of fol 2.4. For diatoms, species is the prevalent taxonomic level in all TAs (figure 4 B). Some assemblages consist entirely of species (i.e. RT04, 06, 09, 16, 17, and 19). The mean number of species per TA is 30 and the mean number of genera 0.5. RT03 has the most taxa rich TA with 49 taxa and RT19 the most taxa poor with 18.

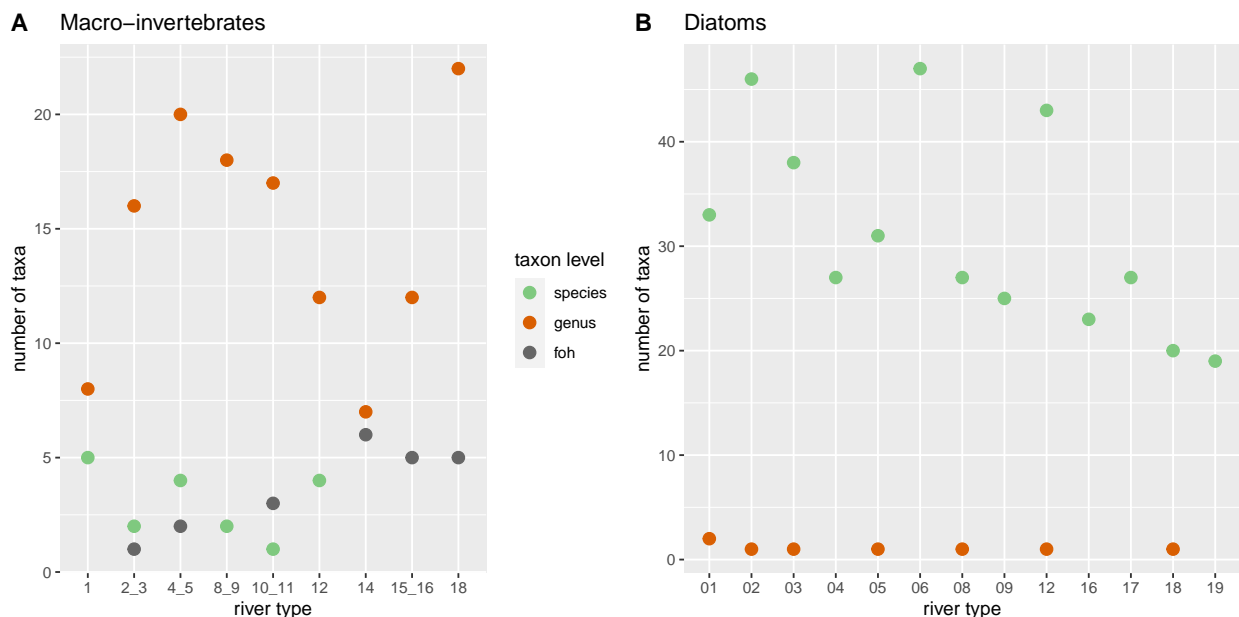


Figure 4: Numbers of taxa on each taxonomical level for all typical assemblages

We can express the uniqueness of a TA with the following score: Each taxon receives a taxon uniqueness score that is one divided by the number of TAs it occurs in. For each river type, we sum the taxon scores of all taxa up and divide it by the number of taxa in the river type's TA. If all taxa in the TA are unique to that TA the score is one. If all species occur in one other TA the score is 0.5. The minimal score depends on the number of TAs, as it is 1 divided by that number and it signals that all species in that TA occur in all other TAs. These scores are shown in figure ???. The dashed horizontal lines indicate the minimum score for each representation level.

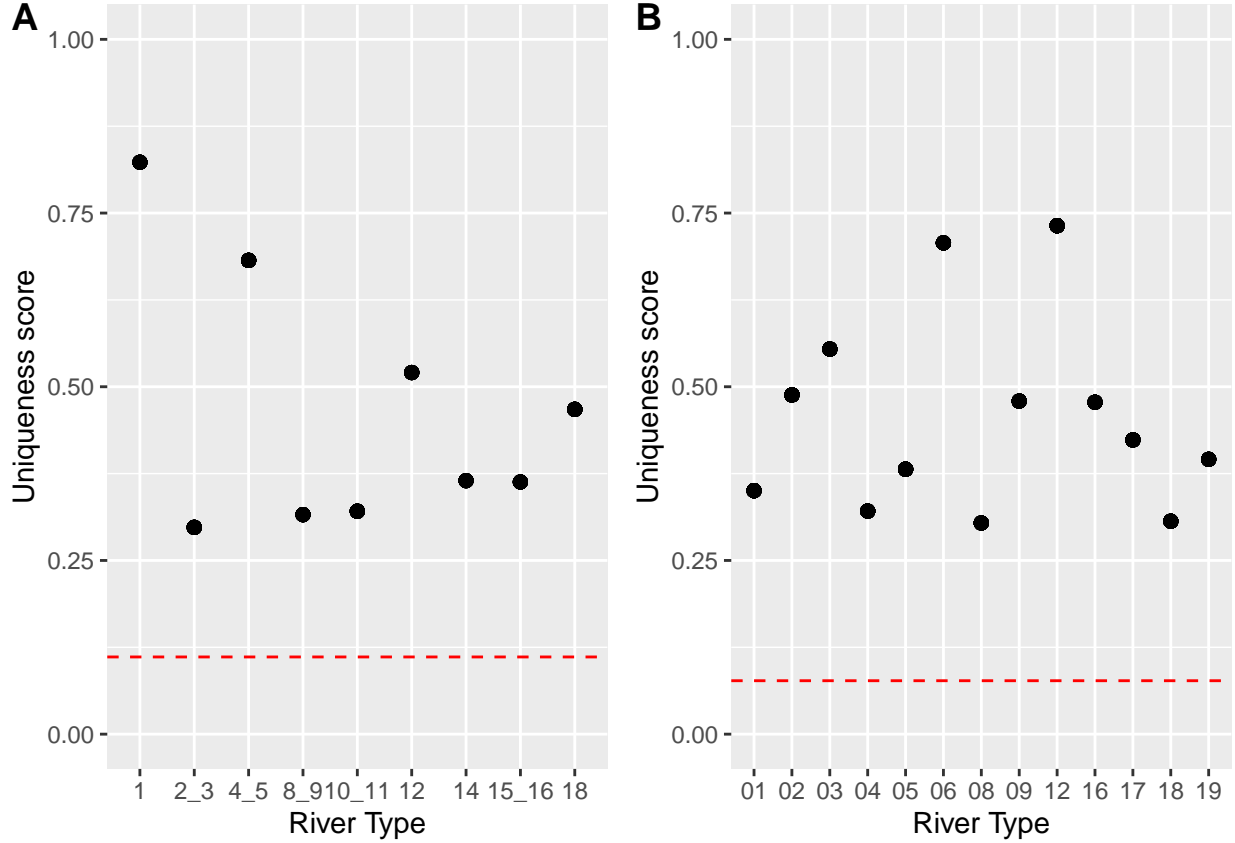


Figure 5: Uniqueness scores for typical assemblages of macro-invertebrates(A) and diatoms (B). The red dashed line indicates the lowest possible score.

We also used Principal Coordinate Analysis (PCoA, Gower (1966)) to visualize the similarity of TAs, based on Jaccard distance matrices (Figure 6).

Taxalists of the typical assemblages can be found [here](#).

8 Seasonal Typical Assemblages

In addition to the spatially defined TAs, we derived seasonal TAs (sTA) for a subset of river types. The four seasons were defined as follows: spring is March to May, Summer is June to August, Fall is September to November, and Winter is December to February.

To avoid strong spatial signals in the sTA only those river types (RT) were considered in which samples were evenly distributed between seasons. In most cases, an even spatio-temporal distribution could only be achieved by omitting parts of the data (e.g. certain seasons or data sets). The maps for all RT with all available seasons as well as the respective subsets that were used in the further analyses can be found in the [here](#) for macroinvertebrates and [here](#)

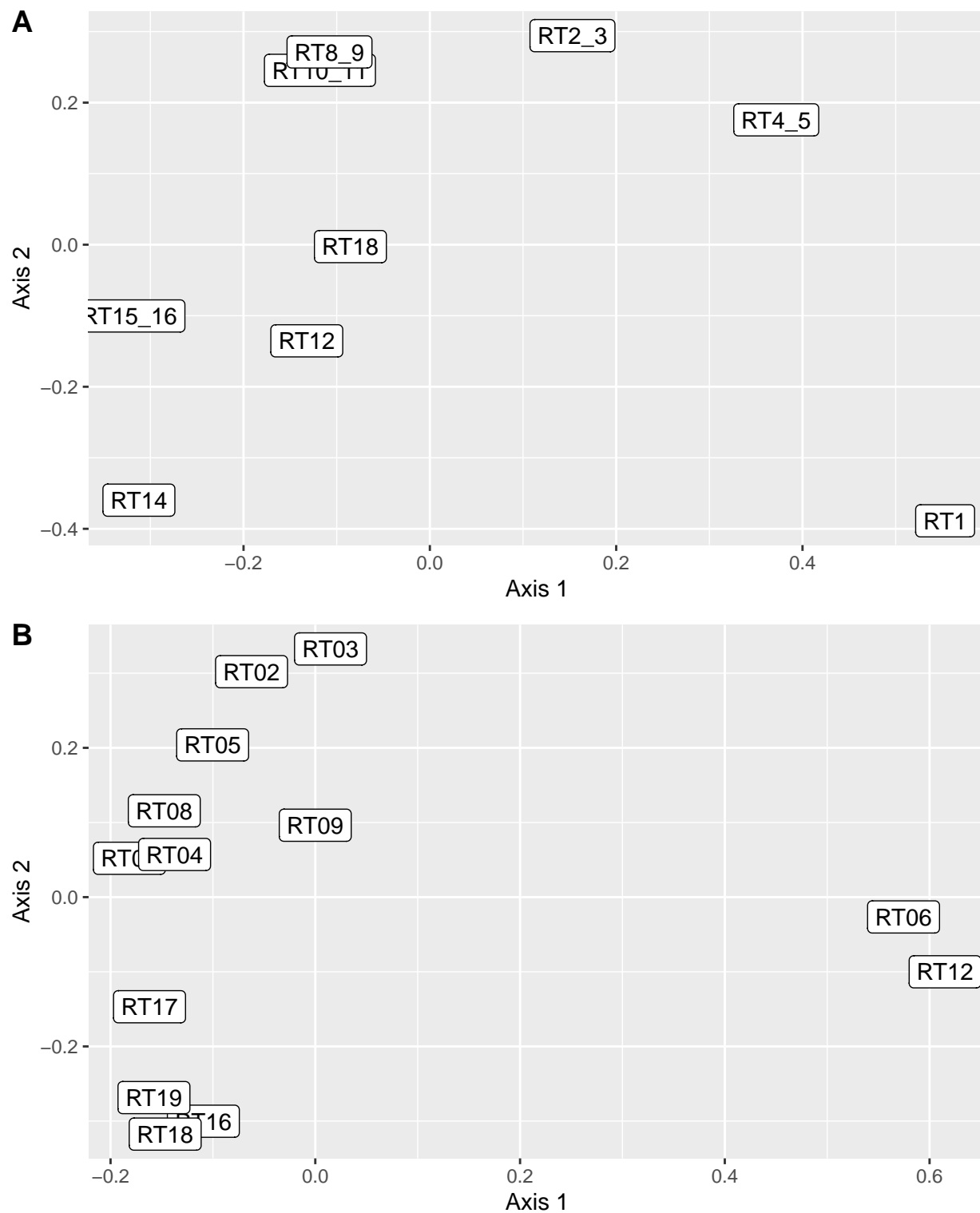


Figure 6: Principal Coordinate Analysis ordinations of typical assemblages based on Jaccard distance matrices. A shows the typical assemblages of macroinvertebrate and B those of diatoms.

Table 3: Overlap between seasonal typical assemblages (sTA) of diatoms in river type 11 expressed in percent of taxa in row sTA also present in column sTA. N is the number of taxa in the respective sTA.

	summer	autumn	winter	N
summer	100.0	46.4	35.7	28
autumn	61.9	100.0	38.1	21
winter	45.5	36.4	100.0	22

for diatoms. As an example, the map of macroinvertebrate samples for the combined RT 10+11 is shown in Figure ??.

To visualize differences between the seasons we used Nonmetric multidimensional scaling (NMDS) on Jaccard dissimilarity matrices. The resulting NMDS plots are available in the for macroinvertebrates and diatoms. Figure ?? shows the NMDS plot for invertebrate samples in RT10+11. Further, we evaluated whether the Jaccard dissimilarity between sites would be better explained by spatial distance or by season. To this end, we employed generalized dissimilarity modeling (GDM, Ferrier *et al.* (2007)). In GDMs, the response variable is the ecological dissimilarity between two sites (expressed in some *a priori* chosen dissimilarity metric, here Jaccard). Smooth functions are fitted to the environmental data and the differences between the values of these functions at the two sites of interest are used as explanatory variables. By using a generalized modeling framework we can account for the bounded nature of dissimilarity metrics (between 0-1) and the smooth functions allow for variation in the rate of compositional turnover along gradients. The plots comparing the effect of spatial distance to that of season for all GDMs can be found in the for macroinvertebrates and diatoms. The plot for invertebrates in RT10 + 11 is shown in Figure ??

Based on the results of NMDS and GDMs, we selected RT 10 + 11 and RT 15 + 16 for invertebrates and RT10 and RT15 for diatoms. For these four river types sTA were derived in the same way as the non-seasonal TAs.

9 Patterns and overlap in seasonal assemblages

In river type 11, the number of diatom taxa in the sTAs did not vary strongly between the seasons (Table 3). The summer and autumn sTAs were more similar to each other than either of them was to the winter sTA. The latter was most similar to the summer sTA, as they share some Gomphonema species (*Gomphonema olivaceum olivaceoides* and *Gomphonema parvulum* Complex) which are absent from the autumn sTA with exception of *Gomphonema pumilum* Complex.

For diatom in RT 15, the winter sTA is considerably larger than the summer and autumn sTAs (Table @ref{tab:tbl-sta-overlap-dia-15}).

Table 4: Overlap between seasonal typical assemblages (sTA) of diatoms in river type 15 expressed in percent of taxa in row sTA also present in column sTA. N is the number of taxa in the respective sTA.

	summer	autumn	winter	N
summer	100.0	57.9	63.2	19
autumn	68.8	100.0	81.2	16
winter	41.4	44.8	100.0	29

Table 5: Overlap between seasonal typical assemblages (sTA) of macroinvertebrates in the combined river type 10+11 expressed in percent of taxa in row sTA also present in column sTA. N is the number of taxa in the respective sTA.

	spring	summer	autumn	winter	N
spring	100.0	0.0	50.0	100.0	2
summer	0.0	100.0	71.4	28.6	7
autumn	5.9	29.4	100.0	17.6	17
winter	33.3	33.3	50.0	100.0	6

Both, the summer and the autumn sTAs, overlap 81.2% with the winter sTA. Therefore, they cross the threshold of 75% overlap we used to delineate redundant TAs. The overlap between the winter sTA and either summer or autumn sTA is of a similar size (41.4% and 44.8%). In general, the overlaps in river type 15 are larger than those in river type 11 which might indicate a weaker seasonal turnover in these ecosystems.

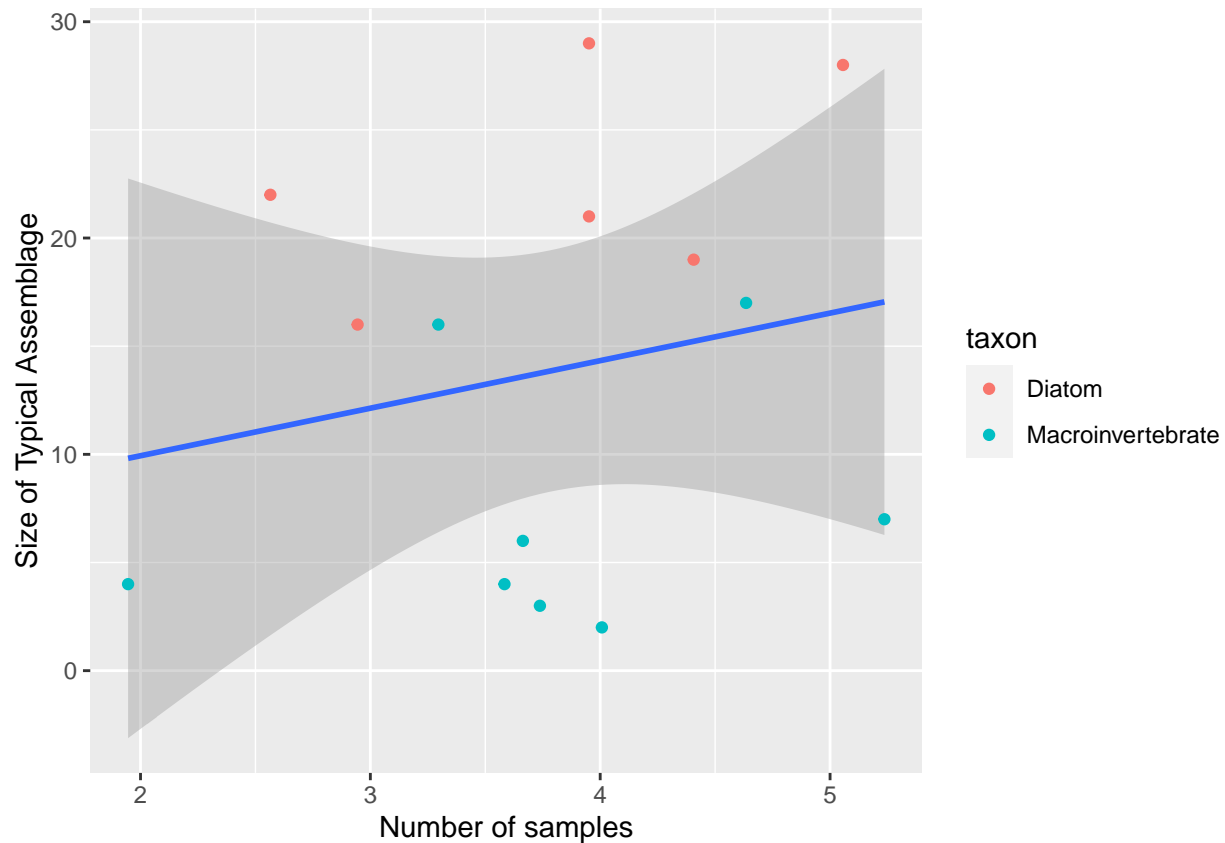
For the macroinvertebrates, the number of taxa in the sTAs is lower than for diatoms. In both river types, the number of taxa in the autumn sTA is also markedly higher than for all other macroinvertebrate sTAs. In the combined river type 10 + 11, the spring sTA was nested in the winter sTA and had no overlap with the summer sTA (Table ??). The summer sTA was most similar to the autumn TA (71.4% overlap) and vice versa (29.4% overlap). Half of the taxa in the winter TA are also part of the autumn TA which is the highest overlap for the winter TA.

Table 6: Overlap between seasonal typical assemblages (sTA) of diatoms in river type 15 expressed in percent of taxa in row sTA also present in column sTA. N is the number of taxa in the respective sTA.

	spring	summer	autumn	winter	N
spring	100.0	75.0	50.0	50.0	4
summer	100.0	100.0	66.7	66.7	3
autumn	12.5	12.5	100.0	18.8	16
winter	50.0	50.0	75.0	100.0	4

In the other combined river type considered here, RT 15 + 16, the summer sTA is nested within the spring sTA and the winter sTA is almost nested within the autumn sTA (Table ??). *Limnoidae* is the only taxon that occurs in the winter sTA but not the autumn sTA. Across this divide, the sTAs only share the two taxa which are common to all four: *Baetis* and *Chironomidae*.

Several possible mechanisms that could explain the higher richness in diatom sTAs compared to macroinvertebrate sTAs as well as the higher richness in autumn sTAs observed for invertebrates are explored below. There might be a connection between the number of samples taken and the size of the TAs. The effect could plausibly be hypothesized to increase or decrease the size of the TA with an increasing number of sampling locations. A decrease in TA size could occur because the absolute number of occurrences required to be included in the TA is increased with the number of sampling sites. In this model, the high number of taxa would be caused by taxa that are only included because of noise. The number should decrease if more samples would be taken. An increase in TA size with an increasing number of sampling sites could occur because the total impact of atypical sites which might be overrepresented in the sample by happenstance would most certainly decrease. Here, atypical refers to the community composition, i.e. taxa that are rare in other sites occur and otherwise frequently occurring taxa are absent. Similar effects might also account for the difference between diatoms and macroinvertebrates. In a linear regression of the size of typical assemblages against the number of sampling locations (log-transformed) across taxa groups, no relationship was identified ($F = 0.569$, $df = 12$, $p = 0.47$, $R^2 = 0.05$, Figure ??).



10 Notes for Traits

CWM with B value then RR-VGLM (glm RDA)

References

- Cáceres, M.D. & Legendre, P. (2009). Associations between species and groups of sites: indices and statistical inference. *Ecology*, 90, 3566–3574.
- Dufrêne, M. & Legendre, P. (1997). Species Assemblages and Indicator Species: The need for a flexible asymmetrical Approach. *Ecological Monographs*, 67, 345–366.
- Ferrier, S., Manion, G., Elith, J. & Richardson, K. (2007). Using generalized dissimilarity modelling to analyse and predict patterns of beta diversity in regional biodiversity assessment. *Diversity and Distributions*, 13, 252–264.
- GBIF.org. (2020). *GBIF home page*.
- Gower, J.C. (1966). Some distance properties of latent root and vector methods used in multivariate analysis. *Biometrika*, 53, 325–338.

- Guiry, G.M., M. D. & Guiry. (2020). *AlgaeBase. World-wide electronic publication, national university of ireland, galway.*
- Kahlert, M., Rühland, K.M., Lavoie, I., Keck, F., Saulnier-Talbot, E. & Bogan, D. *et al.* (2020). Biodiversity patterns of Arctic diatom assemblages in lakes and streams: Current reference conditions and historical context for biomonitoring. *Freshwater Biology*, 1–25.
- Lecointe, C., Coste, M. & Prygiel, J. (1993). “Omnidia”: Software for taxonomy, calculation of diatom indices and inventories management. *Hydrobiologia*, 269, 509–513.
- Lee, S.S., Bishop, I.W., Spaulding, S.A., Mitchell, R.M. & Yuan, L.L. (2019). Taxonomic harmonization may reveal a stronger association between diatom assemblages and total phosphorus in large datasets. *Ecological Indicators*, 102, 166–174.
- Lyche Solheim, A., Austnes, K., Globevnik, L., Kristensen, P., Moe, J. & Persson, J. *et al.* (2019). A new broad typology for rivers and lakes in Europe: Development and application for large-scale environmental assessments. *Science of the Total Environment*, 697, 134043.
- Mauch, E., Schmedtje, U., Maetze, A. & Fischer, F. (2017). Taxaliste der Gewässerorganismen Deutschlands. *Informationsberichte des Bayerischen Landesamtes für Wasserwirtschaft*, 1.
- Rimet, F., Gusev, E., Kahlert, M., Kelly, M.G., Kulikovskiy, M. & Maltsev, Y. *et al.* (2019). Diat.barcode, an open-access curated barcode library for diatoms. *Scientific Reports*, 9, 1–12.
- Scott Chamberlain & Eduard Szocs. (2013). Taxize - taxonomic search and retrieval in r. *F1000Research*.