

# Derivation of Typical Assemblages

Jonathan Jupke

November 20, 2020

## Contents

1	Introduction	2
2	Harmonizing taxa names	2
3	What is the optimal taxonomic level?	3
4	Can we represent the stream types with our data?	4
5	What is a typical assemblage?	5
6	Redundancy between typical assemblages	7
7	Characteristics of typical assemblages	8
8	Patterns and overlap in seasonal assemblages	14
	References	16

# 1 Introduction

In GetReal, we investigate whether selected freshwater and terrestrial assemblages vary spatially and seasonally in their sensitivity towards chemicals and draw conclusions for risk assessment. After we agreed on using the Pan-European river typology proposed by Lyche Solheim *et al.* (2019) to delineate the receiving freshwater systems, the next step was to derive their typical assemblages (TA). Here we describe the methods we used to derive TAs of macroinvertebrates and diatoms for selected river types across Europe. We also describe and briefly discuss the results.

## 2 Harmonizing taxa names

International diatom occurrence datasets require extensive harmonization because of the taxonomic resolution differing between datasets, different working groups using different nomenclatures, identification errors, and ongoing changes to the accepted nomenclature (Kahlert *et al.* 2020). Though harmonizing occurrence datasets can reduce their taxonomic resolution, it also facilitates the detection of large-scale spatio-temporal patterns (Lee *et al.* 2019). We compared all our datasets against six databases that contain accepted names, synonyms with links to the respective accepted names, and suggestions for grouping contentious taxa in larger complexes. If we found a taxon name in one of the databases we either accepted it, changed it into the accepted name in case of a synonym, or included it in a complex. We did not query taxa we found in earlier databases against later ones, but in case the name changed from the original one, we queried the new one against earlier databases. Lastly, we controlled the results visually for consistency. We used the following databases 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)

We harmonized the macroinvertebrate data with gbif through the taxize R package (Chamberlain & Szöcs 2013).

### 3 What is the optimal taxonomic level?

One result of the second progress review for GetReal (29.04.2020) was that we should determine an optimal taxonomic level for each taxon separately instead of using one common level for all data. *Oligochaetes*, for example, are usually only determined to subclass level. In a setting with one common taxonomic level (e.g. genus) they would have to be omitted if this level would be higher than subclass. By using taxon-specific levels that take low-resolved taxa into account, we can thus use more of the data. However this poses the challenge of finding an optimal level for each taxon given a dataset.

The following refers exclusively to the macroinvertebrate data, as the taxonomic resolution was not an issue with diatom datasets. We describe the procedure for diatom data at the end of the section. The optimal level was established with a hierarchical approach. First, we removed all observations from Phyla and Classes that were not present in all datasets. We assumed that these represented differences in sampling rather than in communities. The classes Clitellata (Annelida), Insecta, Malacostraca (Arthropoda), Bivalvia and Gastropoda (Mollusca) remained. 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 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, we hold a taxon to be 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 salient change occurs between 85% and 86%. It occurs because for  $X > 86\%$  *Chironomidae* are represented at the family level. 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.

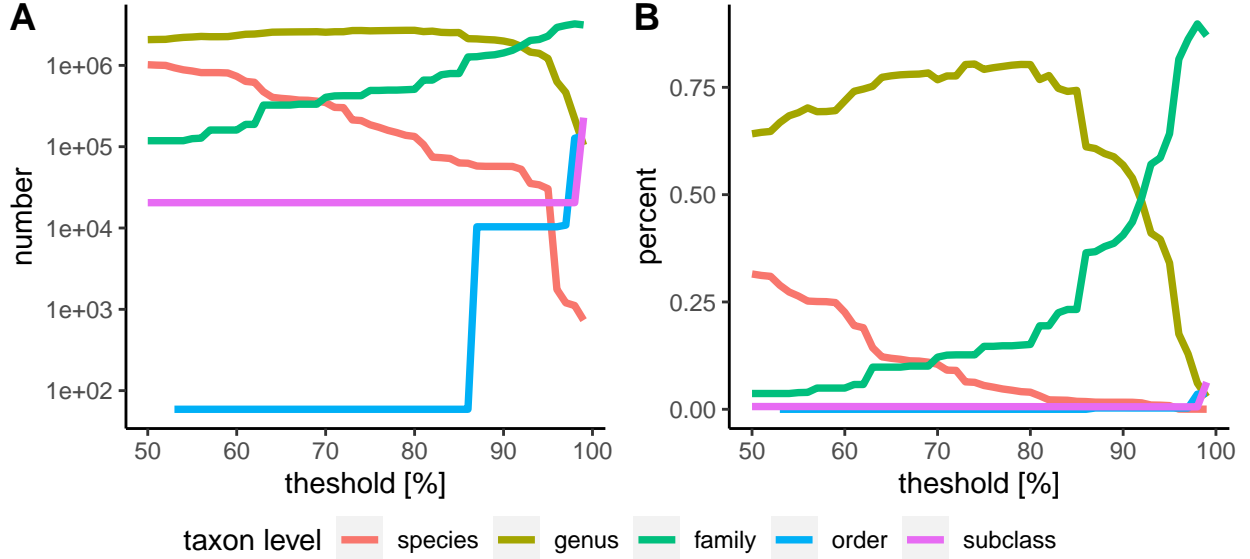


Figure 1: Effect of the threshold value on (A) total observations of each taxonomic level and (B) percentage of each taxonomic level.

For the diatoms, we employed 75% as a threshold, because for *Gomphonema*, which is the fourth most common genus in our dataset, 81.43% of the observations were at the species level. The taxonomic resolution was higher than in the macroinvertebrate data and the lowest resolution is the genus level. The equivalent of Figure 1 for diatoms can be found [here](#).

## 4 Can we represent the stream types with our data?

We determined visually whether our dataset contained enough sampling sites in a given river type to derive meaningful TAs. The degree of representation for river types 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 that 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 macroinvertebrates and diatoms are shown in table 1.

For each river type we provide maps with the associated sampling sites for [macroinvertebrates](#) 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.*

Table 1: The ratings for all river types for macroinvertebrates and diatoms

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

(2019). We have fewer sampling sites for diatoms than for macroinvertebrates which entails that the representation of stream types is mostly lower.

## 5 What is a typical assemblage?

We derived typical assemblages using the commonness of each taxon in a respective river type. Here, the commonness of a taxon refers to the fraction of sampling sites in a river type where that taxon is found. This is conceptually equal to the  $B$  parameter used in the Species Indicator Value (Dufrêne & Legendre 1997; Cáceres & Legendre 2009). We also considered the second parameter of that statistic, the specificity  $A$  which refers to the fraction of occurrences of a taxon that fall within one river type. For example,  $A = 1$  indicates that all observations of taxon  $x$  occurred in river type  $y$  and hence that this taxon is highly specific to the given river type. A commonness of 1 shows that a taxon is present in every sample taken within a river type and therefore that it is very typical for that river type. A taxon belongs to the TA of a river type if it is more common than some threshold, which depends on the taxonomic level of the taxon, or if it is more specific than another threshold, which also depends on the taxonomic level. The B-threshold for species is 0.25, for genera 0.35 and for family or lower resolutions 0.55. In addition to the typical assemblages we also derived lists of specific taxa. To be included in these the specificity had to be above 0.9 for a species, 0.8 for genus or 0.7 for family or lower taxonomic levels.

We did not systematically optimize these thresholds. Such procedures would require optimization criteria, but we are not aware of any criterion that would work in this context. As TAs can be very similar in composition or harbor strongly differing numbers of taxa, neither criterion would optimize what we would consider a TA. 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 varying these parameters would alter the results. We derived TAs for 50 values of **A** and **B** ranging from 0.01 to 1 and computed taxon richness and uniqueness scores (see Figure 5 and text preceding the Figure)

of each TA. Please note that results are only shown and discussed for the non-redundant TAs (see section 6).

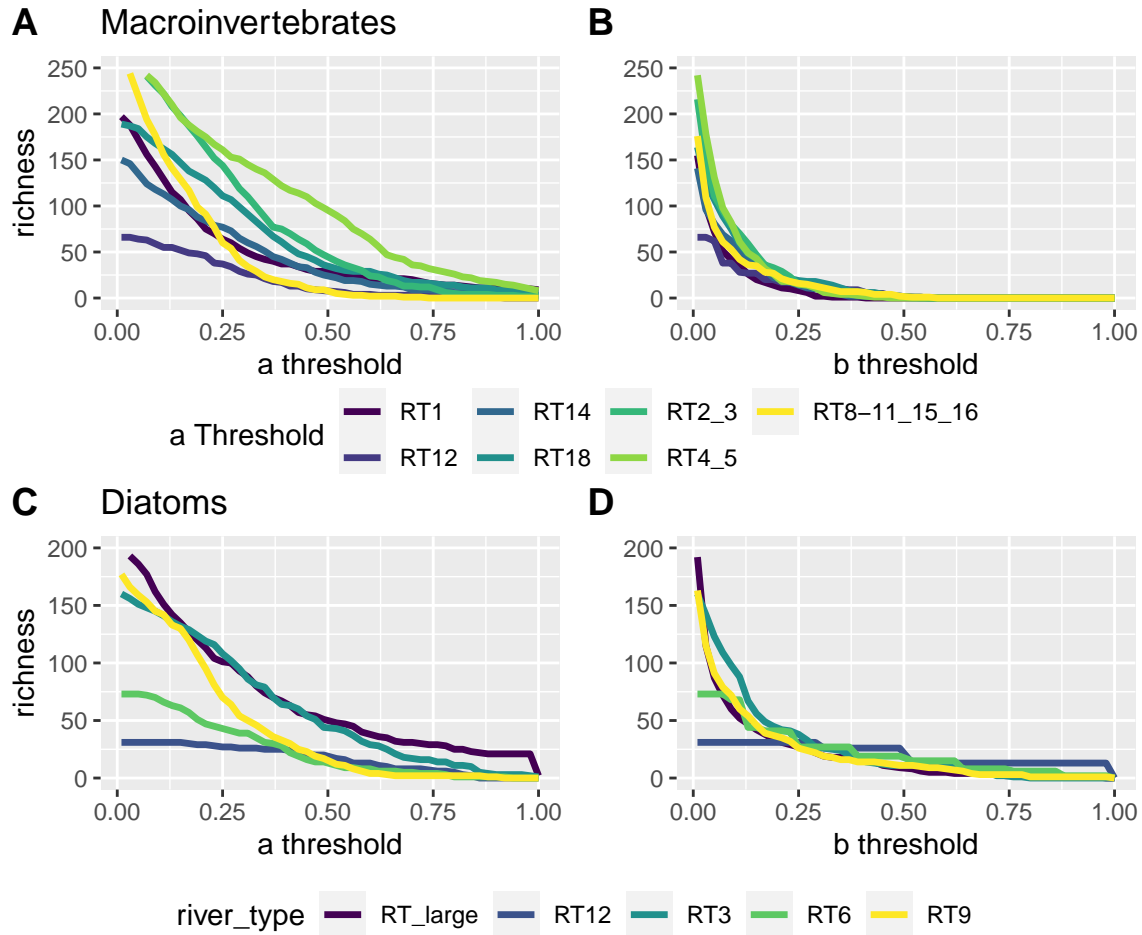


Figure 2: Changes in richness along gradients of A and B thresholds in macroinvertebrates (A, B) and diatoms (C, D). The line color indicates the river type.

Richness decreased with increasing **A** and **B** threshold in macroinvertebrates and diatoms (Figure 2). Uniqueness scores increased with **A** thresholds up to a certain level and then decreased (Figure 3). There is considerable variation between river types as to where this inflection occurs. Along the B-gradient, uniqueness score fluctuate erratically, driven by the relative rate at which other river types loose taxa and the identity of these. At some point no taxa are common enough to surpass the set B-threshold and the TAs are empty. For macroinvertebrates this takes place for B values between 0.41 (RT1) and 0.58 (RT8-11\_15\_16). In the diatom assemblages this occurs later (at B = 0.8 in RT3) and not at all in RT9, RT6, RT12. This highlights that there must be several widely distributed diatom taxa, which is likely to some degree also the result of our extensive harmonization efforts.

Plots for each taxon level separately are available for [macroinvertebrates](#) and [diatoms](#). How-

ever, the general patterns visible in Figure 2 and 3, hold for them as well.

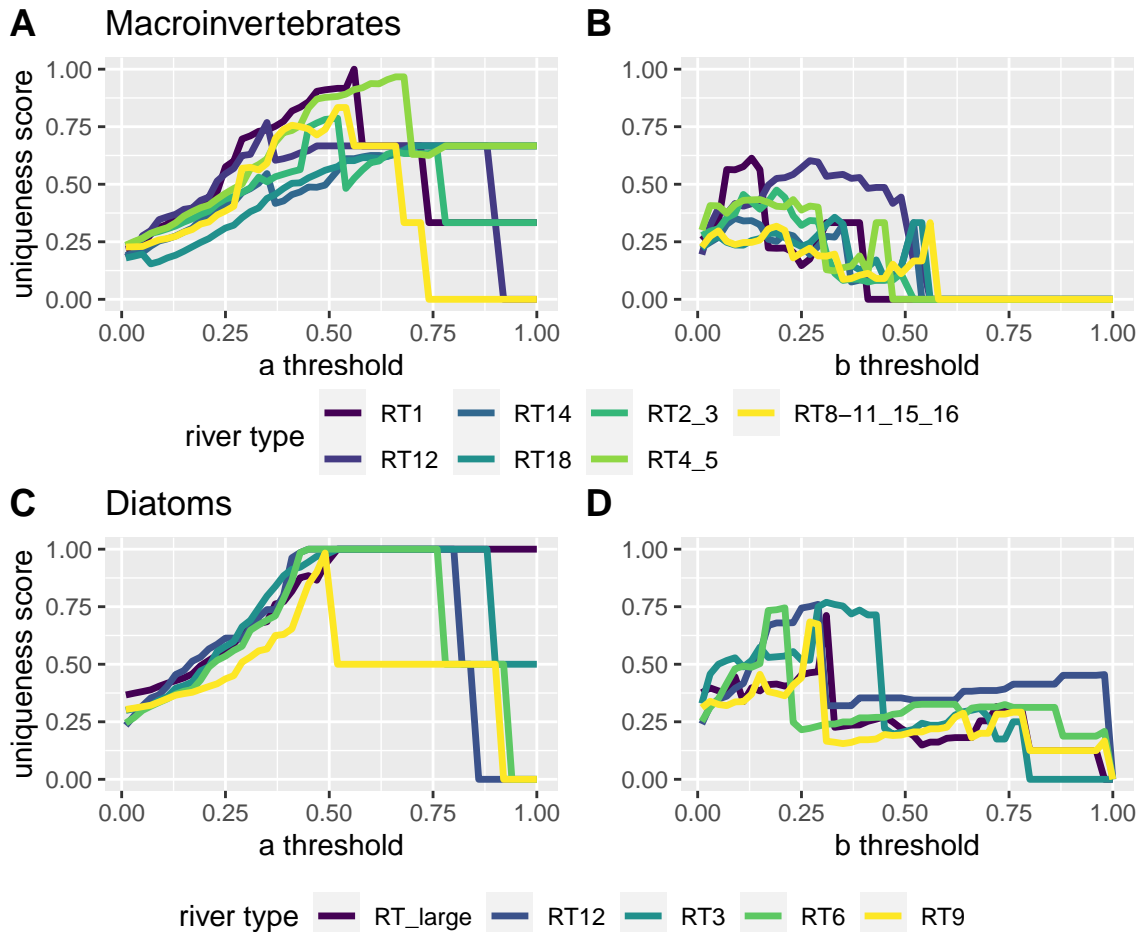


Figure 3: Changes in uniqueness scores along gradients of A and B thresholds in macroinvertebrates (A, B) and diatoms (C, D). The line color indicates the river type.

## 6 Redundancy between typical assemblages

We assessed the degree to which the different TAs overlap. 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. After assessing the redundancy we joined the two river types with the highest redundancy, if that redundancy was above 0.8, and computed the new TAs. This procedure was iterated until no redundancies were above 0.8. In macroinvertebrates, we combined RT2 and 3, RT 4 and 5 and RT 8,9,10,11,15 and 16. RT2 and 3 are both lowland, siliceous rivers of medium to large (RT2) or very small to small (RT3) size. Their combination is plausible resulting in

very small to large lowland siliceous rivers. Note that very large rivers (catchment area  $> 10,000 \text{ km}^2$ ) are a separate type (RT1). RT4 and 5 follow the same logic as RT2 and 3. They are lowland calcareous or mixed rivers of of medium to large (RT4) or very small to small (RT5) size. The larger cluster of RT 8,9,10,11,15, and 16 is driven by altitude. The river types 8 to 11 are mid-altitude (200 - 800 m.a.s.l.) rivers differing in size and geology. There are three mid-altitude river types that did not fall into this cluster: RT12 which represents siliceous rivers with histosol soils in more than 20% of the catchment area, RT13 which is exceedingly rare and was omitted from the analysis, and RT18 which delineates Mediterranean mid-altitude rivers. The two river types RT15 and RT16 are two of three types of high altitude ( $> 800 \text{ m.a.s.l.}$ ) rivers. Both occur mainly in southern Europe with differentiates them from the northern high altitude rivers in RT14.

In summary this large cluster represents mid- to high-altitude rivers of various sizes and geologies.

In diatoms, the river types 1, 2, 4, 8, 9, 17, 18, and 19 are all combined into one large cluster. This cluster comprises large (RT2 and RT4) to very large (RT1) lowland rivers, all three types perennial Mediterranean rivers (RT 17, 18 and 19) as well as some mid-altitude rivers (RT8 and RT9). This collection seems eclectic and might indicate a dominance of widespread generalist species. Hence it is more informative to focus on what is not comprised in it.

Small lowland rivers (RT3 and RT5) have distinct typical TAs from their larger counterparts and from each other. The two river types that are characterized by histosol soils (RT6 and RT12), are most similar to one another (55.2% and 51.6%) but their overlap is the lowest we observed in diatoms. Previous studies indicated, that there is strong turnover in diatom communities along organic matter gradients (Kelly *et al.* 1995; Coring 1999; Hering *et al.* 2006). Lastly, high altitude rivers (RT16) seem to be different from the large cluster despite a considerable spatial proximity to many sites in the latter. Wang *et al.* (2019) recently showed that altitude affects assembly processes and community composition in diatoms and Göthe *et al.* (2013) found similar patterns.

## 7 Characteristics of typical assemblages

In all macroinvertebrate TAs, except RT14, genus is the prevalent taxonomic level, followed by families or lower taxonomic levels and lastly species (Figure 4A). The mean number of species was 1.33, mean number of genera 15.57, and the mean number of families or lower 6. RT1 had the least taxa with 16 taxa and RT2\_3 and RT18 were the most taxa rich assemblages with 30 taxa. For diatoms, species is the prevalent taxonomic level in all TAs (Figure 4B). The mean number of species per TA is 29.5 and the mean number of genera 3.33. RT3 has the most taxa rich TA with 41 taxa and RT16 has the least taxa in its TA with 26. Note that RT3 or a combined river type including it (RT2\_3) had the most taxa rich TA for both taxa and that diatom TAs encompass more taxa than macroinvertebrate



TAs which supports the trends from the sensitivity analysis (Figure 3).

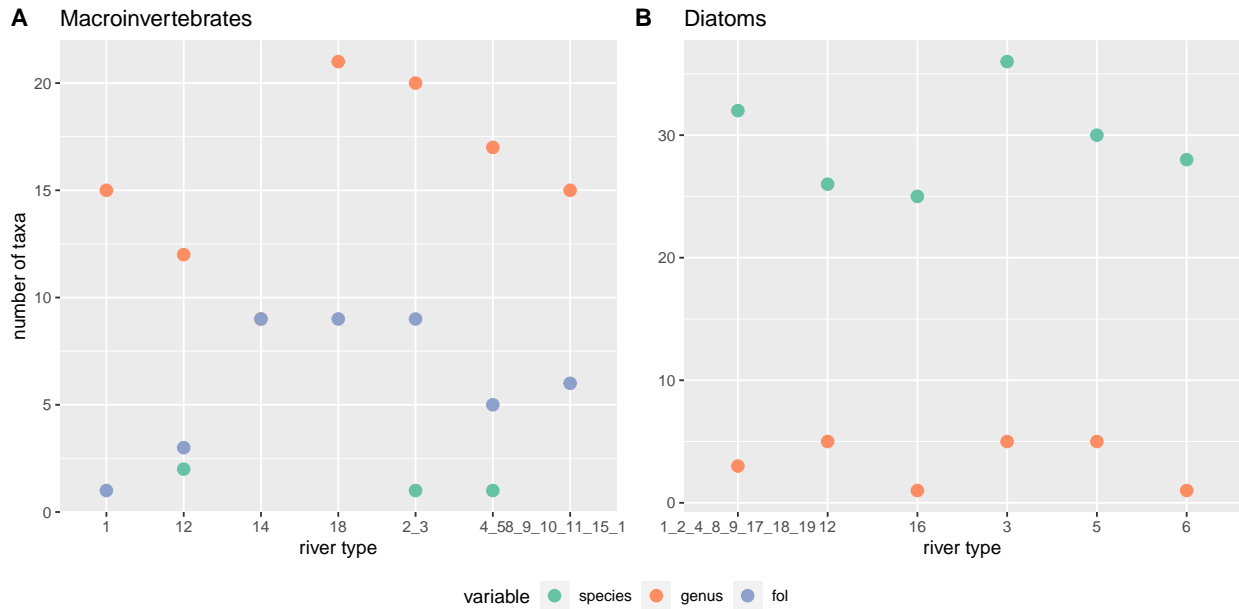


Figure 4: Numbers of taxa on each taxonomical level for all typical assemblages.

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 5. The dashed horizontal lines indicate the minimum scores.

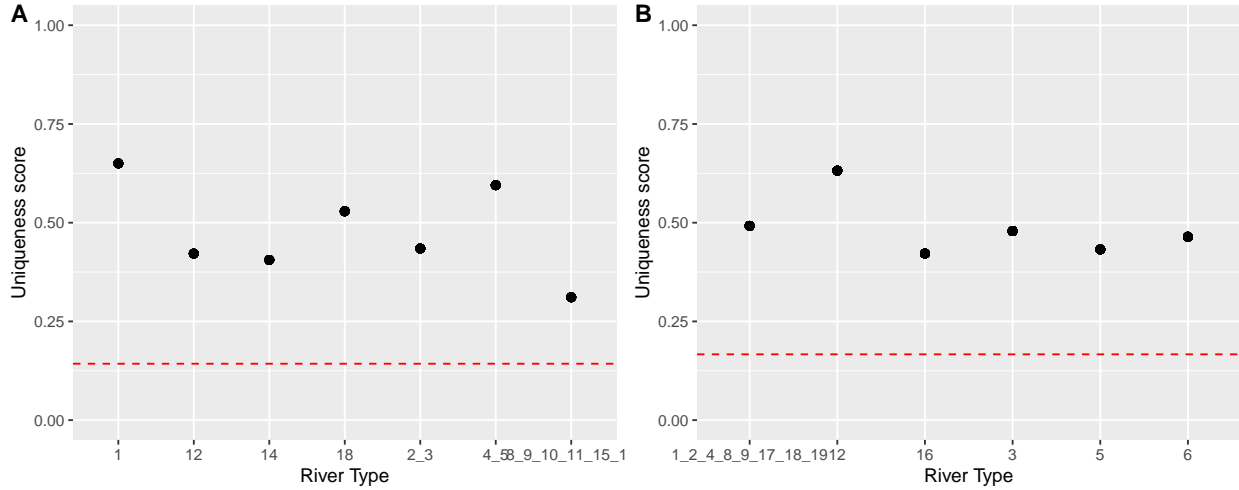


Figure 5: Uniqueness scores for typical assemblages of macroinvertebrates(A) and diatoms (B). The red dashed line indicates the lowest possible score.

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

The first axis of the macroinvertebrate PCoA separates the TAs into two groups: i) the large cluster 8-11\_15\_16, RT12, RT14, RT18, RT2\_3 and ii) RT4\_5 and RT1. The larger group represents all mid to high-altitude rivers and RT2\_3.

The two river types in the smaller group are more distinct than those in the larger but could be identified as low-land rivers in contrast to the first group. However the cumulative relative eigenvalues of the first two principal components is only a little higher than half, which means that there is considerable variation in the distance matrix that is not displayed in this PCoA

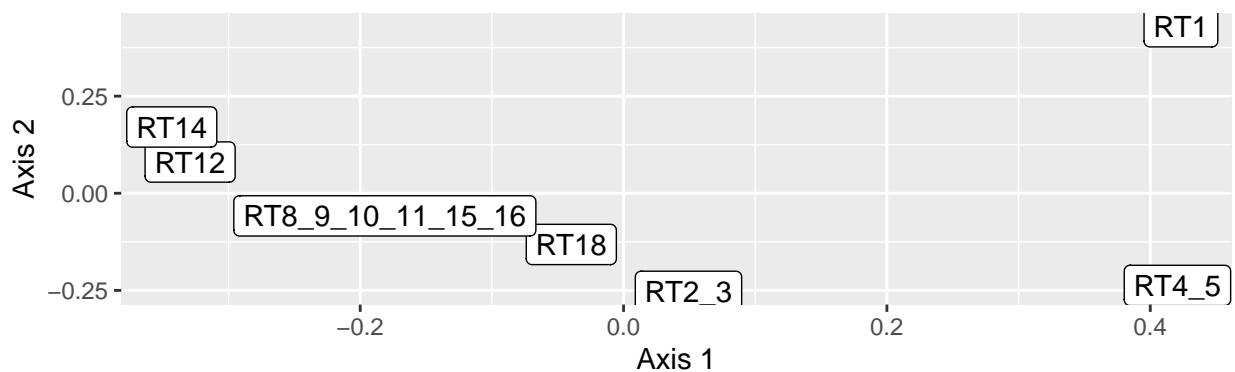
In the PCoA of diatom TAs, river types that are expected be similar to some included in the large cluster are close to the latter (i.e. RT9 and RT3). The two river types with strong influence of organic soils (RT6 and RT12) are distinct from each other and from the larger group.

Online, we provide the taxa lists for all [macroinvertebrate](#) and [diatom](#) TAs. # 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 sTAs, we only considered those river types (RT) in which samples were evenly distributed between seasons. In most cases, we had to omit parts of the data (e.g. certain seasons or datasets) to achieve an even spatio-temporal distribution. Online, we provide maps for all RT with all available seasons as well as the respective subsets

**A** PCoA of typical macroinvertebrate assemblages

cumulative eigenvalues of first two axes: 0.55



**B** PCoA of typical diatoms assemblages

cumulative eigenvalues of first two axes: 0.66

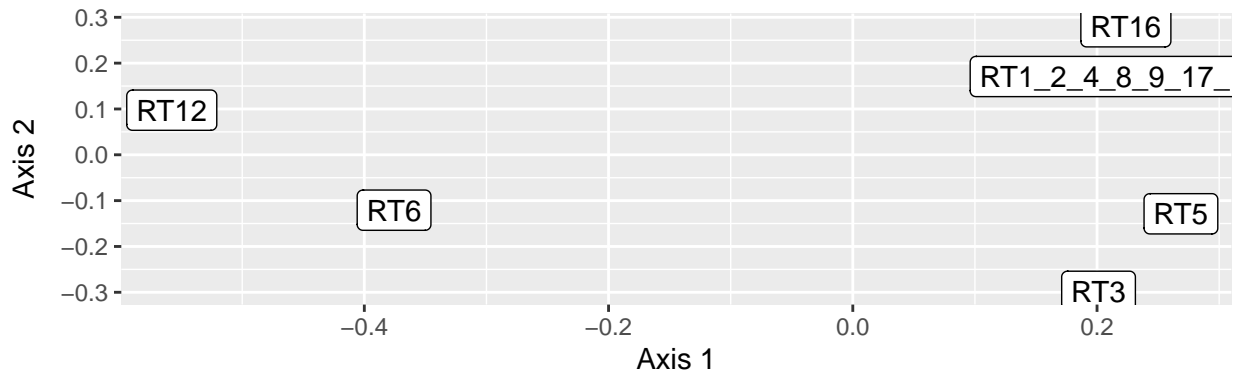


Figure 6: Principal Coordinate Analysis ordinations of typical assemblages based on Jaccard distance matrices. A shows the typical assemblages of macroinvertebrates and B those of diatoms.

that we used in the further analyses for [macroinvertebrates](#) and [diatoms](#).

As an example, the map of macroinvertebrate samples for the combined RT 4\_5 is shown in Figure 7.

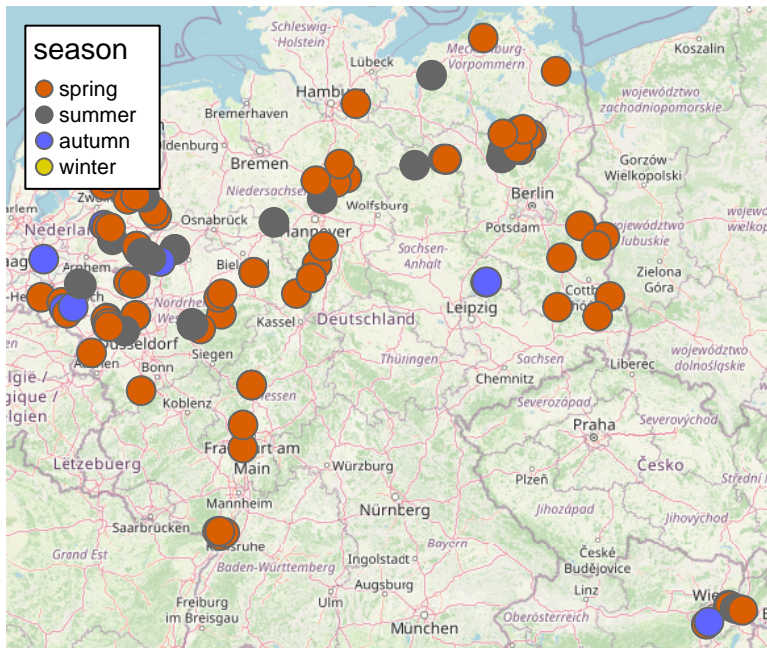


Figure 7: Map of sampling sites for the combined River Type 4\_5. The color of the points shows the season of sampling.

To visualize differences between the seasons we used Non-metric multidimensional scaling (NMDS) on Jaccard dissimilarity matrices. The resulting NMDS plots are available for [macroinvertebrates](#) and [diatoms](#). Figure 8 shows the NMDS plot for invertebrate samples in RT4\_5. For diatoms and macroinvertebrates there are no or little discernible seasonal patterns in most river types. This also shows in high NMDS stress values (typically above 0.2).

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. Plots that show the effect of spatial distance and that of season for all GDMs are available for [macroinvertebrates](#) and [diatoms](#). The plot for

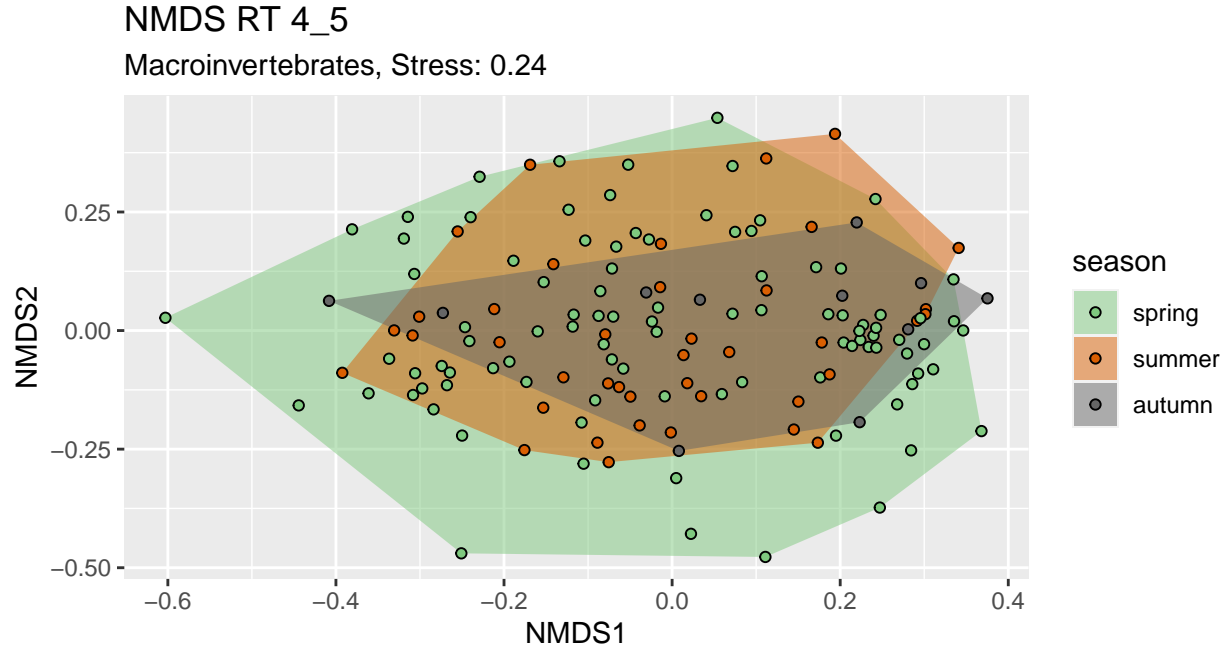


Figure 8: Nonmetric multidimensional scaling plot of Jaccard dissimilarity matrices for macroinvertebrate communities in RT4\_5. The color of the points shows the season. Convex hulls surround all sampling points from one season.

invertebrates in RT4\_5 is shown in Figure 9. The findings from the NMDS are confirmed in the GDMs - spatial distance explains more of the variation in jaccard distances than seasons.

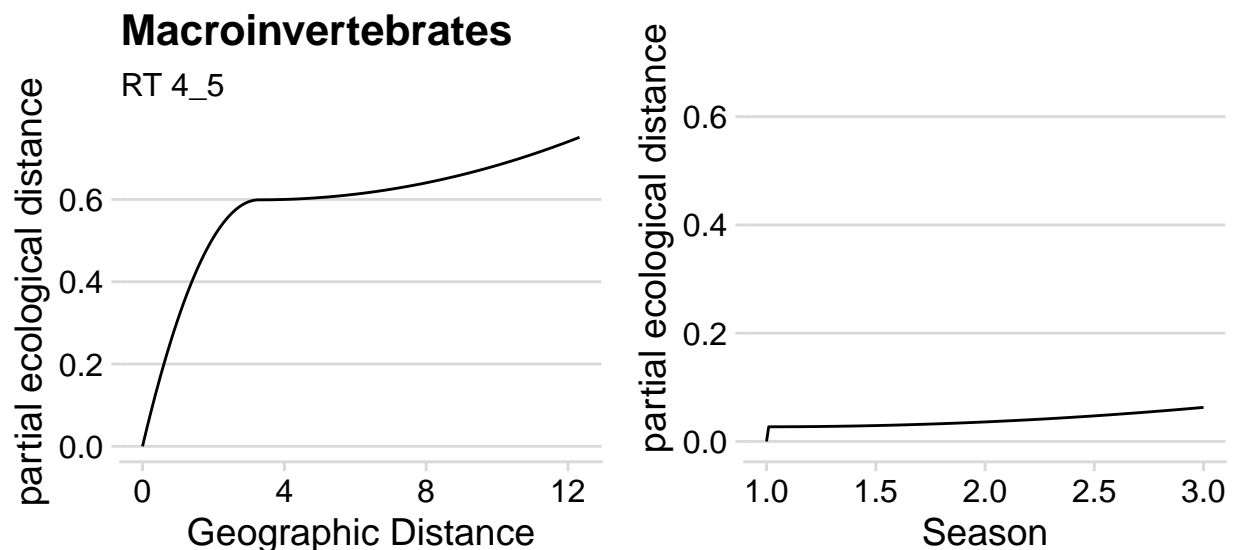


Figure 9: Partial ecological distance between sites with increasing geographic distance or changing season (1 = spring, 2 = summer, 3 = autumn, and 4 = winter) predicted with Generalized Dissimilarity Models

We selected RT2\_3 and 8-11\_15\_16 for invertebrates and RT11 and 15 for diatoms because they showed the strongest effect of season in GDM and NMDS. For these four river types sTA were derived in the same way as the non-seasonal TAs.

## 8 Patterns and overlap in seasonal assemblages

In RT11, the number of diatom taxa in the sTAs did not vary strongly between the seasons (Table 2). The summer and autumn sTAs were more similar to each other than either to the winter sTA. The latter was lightly more similar to the summer sTA (66.7%) than to the autumn sTA (63.6%).

Table 2: 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	87.9	66.7	33
autumn	82.9	100.0	60.0	35
winter	66.7	63.6	100.0	33

In RT15 similar patterns can be seen: The number of taxa does not vary strongly between seasons and summer and autumn sTAs are more similar to each other than to the winter

sTA. The difference between summer and winter sTAs in RT15 is the largest seasonal change we found in the two river types.

Table 3: 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	88.5	61.5	26
autumn	85.2	100.0	66.7	27
winter	55.2	62.1	100.0	29

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 higher than in other seasons. In the combined type RT2\_3, the summer sTA is nested in the autumn sTA (Table 4).

Table 4: Overlap between seasonal typical assemblages (sTA) of diatoms in river type 2+3 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	N
summer	100.0	100	5
autumn	27.8	100	18

In the other combined river type considered here, RT 8-11\_15\_16, the sTA sizes are more similar to each other than in RT2\_3 and the sTAs are not as similar (Table 5). Spring and summer have a higher overlap with autumn than with each other and autumn has a marginally higher overlap with summer (56.5%) than with spring (52.2%).

Table 5: Overlap between seasonal typical assemblages (sTA) of diatoms in river type 8-11+15+16 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	N
spring	100.0	71.4	85.7	14
summer	62.5	100.0	81.2	16
autumn	52.2	56.5	100.0	23

Online, we provide the complete taxa lists for [macroinvertebrate](#) and [diatom](#) sTAs.

## References

- Cáceres, M.D. & Legendre, P. (2009). Associations between species and groups of sites: indices and statistical inference. *Ecology*, 90, 3566–3574.
- Chamberlain, S. & Szöcs, E. (2013). Taxize - taxonomic search and retrieval in R. *F1000Research*.
- Coring, E. (1999). Situation and developments of algal (diatom)-based techniques for monitoring rivers in Germany. *Use of Algae for Monitoring Rivers III. Agence de l'Eau Artois-Picardie, Douai*, 122–127.
- 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.
- Göthe, E., Angeler, D.G., Gottschalk, S., Löfgren, S. & Sandin, L. (2013). The Influence of Environmental, Biotic and Spatial Factors on Diatom Metacommunity Structure in Swedish Headwater Streams. *PLoS ONE*, 8.
- Guiry, G.M., M. D. & Guiry. (2020). *AlgaeBase. World-wide electronic publication, national university of Ireland, Galway*.
- Hering, D., Johnson, R.K., Kramm, S., Schmutz, S., Szoszkiewicz, K. & Verdonschot, P.F.M. (2006). Assessment of European streams with diatoms, macrophytes, macroinvertebrates and fish: A comparative metric-based analysis of organism response to stress. *Freshwater Biology*, 51, 1757–1785.
- 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.
- Kelly, M.G., Penny, C.J. & Whitton, B.A. (1995). Comparative performance of benthic diatom indices used to assess river water quality. *Hydrobiologia*.
- 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.
- Wang, J., Hu, J., Tang, T., Heino, J., Jiang, X. & Li, Z. *et al.* (2019). Seasonal shifts in the assembly dynamics of benthic macroinvertebrate and diatom communities in a subtropical river. *Ecology and Evolution*, 1–13.