



UHI Research Database pdf download summary

Improving eDNA-based protist diversity assessments using networks of amplicon sequence variants

Forster, Dominik; Lentendu, Guillaume; Filker, Sabine; Dubois, Elyssa; Wilding, Thomas A.; Stoeck, Thorsten

Published in:
Environmental Microbiology

Publication date:
2019

The re-use license for this item is:
CC BY-NC

The Document Version you have downloaded here is:
Peer reviewed version

The final published version is available direct from the publisher website at:
[10.1111/1462-2920.14764](https://doi.org/10.1111/1462-2920.14764)

[Link to author version on UHI Research Database](#)

Citation for published version (APA):

Forster, D., Lentendu, G., Filker, S., Dubois, E., Wilding, T. A., & Stoeck, T. (2019). Improving eDNA-based protist diversity assessments using networks of amplicon sequence variants. *Environmental Microbiology*. Advance online publication. <https://doi.org/10.1111/1462-2920.14764>

General rights

Copyright and moral rights for the publications made accessible in the UHI Research Database are retained by the authors and/or other copyright owners and it is a condition of accessing publications that users recognise and abide by the legal requirements associated with these rights:

- 1) Users may download and print one copy of any publication from the UHI Research Database for the purpose of private study or research.
- 2) You may not further distribute the material or use it for any profit-making activity or commercial gain
- 3) You may freely distribute the URL identifying the publication in the UHI Research Database

Take down policy

If you believe that this document breaches copyright please contact us at RO@uhi.ac.uk providing details; we will remove access to the work immediately and investigate your claim.

1 **Improving eDNA-based protist diversity assessments using networks of amplicon sequence**
2 **variants**

3

4 Dominik Forster¹, Guillaume Lentendu¹, Sabine Filker², Elyssa Dubois¹, Thomas A. Wilding³, Thorsten
5 Stoeck^{1*}

6

7 ¹Department of Ecology, University of Kaiserslautern, Germany

8 ²Department of Molecular Ecology, University of Kaiserslautern, Germany

9 ³Scottish Association for Marine Science, Scottish Marine Institute, Oban, Scotland

10

11 **Running title:** Network analyses enhance sequence grouping algorithms

12

13 *corresponding author

14 email: stoeck@rhrk.uni-kl.de

15 phone: +49-631-2052502

16 fax: +49-631-205-2496

17

18 **Summary**

19 Effective and precise grouping of highly similar sequences remains a major bottleneck in the
20 evaluation of high-throughput sequencing (HTS) datasets. Amplicon sequence variants (ASVs) offer a
21 promising alternative that may supersede the widely used operational taxonomic units (OTUs) in
22 environmental sequencing studies. We compared the performance of a recently developed pipeline
23 based on the algorithm DADA2 for obtaining ASVs against a pipeline based on the algorithm SWARM
24 for obtaining OTUs. Illumina-sequencing of 29 individual ciliate species resulted in up to 11 ASVs per
25 species, while SWARM produced up to 19 OTUs per species. To improve the congruency between
26 species diversity and molecular diversity, we applied sequence similarity networks (SSNs) for second-
27 level sequence grouping into network sequence clusters (NSCs). At 100% sequence similarity in
28 SWARM-SSNs, NSC numbers decreased from 7.9-fold overestimation without abundance filter, to
29 4.5-fold overestimation when an abundance filter was applied. For the DADA2-SSN approach, NSC
30 numbers decreased from 3.5-fold to 3-fold overestimation. Rand index cluster analyses predicted
31 best binning results between 97% and 94% sequence similarity for both DADA2-SSNs and SWARM-
32 SSNs. Depending on the ecological questions addressed in an environmental sequencing study with
33 protists we recommend ASVs as replacement for OTUs, best in combination with SSNs.

34

35 **Introduction**

36 Ever since high-throughput-sequencing (HTS) has been introduced in molecular ecology,
37 researchers have been looking for effective tools to evaluate the resulting sequence data in the
38 context of diversity measures. The standard way of addressing this issue is to group sequencing reads
39 obtained for example from environmental samples, into operational taxonomic units (OTUs; *e.g.* de
40 Vargas *et al.* 2015; Stoeck *et al.* 2010). This grouping can be achieved either by relying on global
41 clustering scores (Schloss *et al.*, 2009; Edgar, 2010; Fu *et al.*, 2012; Rognes *et al.*, 2016), or on local
42 clustering scores (Mahé *et al.*, 2015). While traditional clustering relies on fixed global clustering
43 thresholds expressed *e.g.* as sequence similarity between aligned sequences, local clustering allows
44 for a more fine-tuned evaluation by comparing all local pairs of nucleotides between the sequences
45 and iteratively grouping them into OTUs. It is well known, though, that every kind of OTU is at best a
46 bioinformatical approximation of a species (Schloss and Westcott, 2011; Tikhonov *et al.*, 2015) and
47 that we are far away from the one OTU – one species ideal. Currently available sequence grouping
48 methods tend to allocate sequencing reads of the same species into multiple OTUs. Even though
49 sequencing errors (Kunin *et al.*, 2010) and intraspecific genetic heterogeneity may also contribute to
50 diversity inflation in environmental HTS datasets, imprecise sequence grouping is the main cause of
51 severe biodiversity overestimations (Flynn *et al.*, 2015; Clare *et al.*, 2016). Likewise, imprecise
52 sequence grouping may also lead to biodiversity underestimations, when sequences of different
53 species are grouped into the same OTU (Bachy *et al.*, 2013; Grattepanche *et al.*, 2014). This
54 emphasizes the need for a more accurate grouping of sequences for providing more realistic
55 estimates of species richness and diversity within a sample.

56 Recently, Callahan and colleagues (Callahan *et al.*, 2017) presented an alternative approach
57 by replacing OTUs with amplicon sequence variants (ASVs). The authors developed the open-source
58 software package DADA2 (Callahan *et al.*, 2016) to model and correct Illumina-sequenced amplicon
59 errors. Bioinformatic tools that follow a similar denoising approach had already been introduced for
60 454 pyrosequencing datasets (Quince *et al.*, 2009; Reeder and Knight, 2010). Several terminologies

61 for the description of these tools can be found. We refer to them as “first-level sequence grouping
62 algorithms”, though “pre-clustering algorithms” is also commonly found in the literature (*e.g.*
63 Schloss *et al.*, 2011). They rely on statistical model-based evaluation of HTS data to infer which
64 sequence base-call differences represent true biological variants and which represent sequencing
65 artifacts (Callahan *et al.*, 2017; Knight *et al.*, 2018). The quality filtering models allow resolutions
66 down to single nucleotide differences between sequences and may thus affiliate error-prone
67 sequences with an existing ASV. Thus, even though designed as a denoising tool, DADA2 is at the
68 same time an elegant way of sequence grouping. Assigning sequences to existing ASVs is of further
69 importance when samples of different studies are analyzed in the same context. Since ASVs are
70 supposed to be consistent biological entities, they provide high levels of reproducibility and
71 comparability across independent studies (Callahan *et al.*, 2016, 2017; Amir *et al.*, 2017). Initial
72 studies comparing ASVs against OTUs obtained by global clustering approaches from the same
73 samples support these and other advantages of ASVs (Callahan *et al.*, 2016; Utter *et al.*, 2016; Allali
74 *et al.*, 2017; Needham *et al.*, 2017; Nearing *et al.*, 2018; Zoqratt *et al.*, 2018). One congruent finding
75 of all these studies was that distinctively fewer ASVs than OTUs were produced from the same
76 samples, regardless of the sampled habitat. Furthermore, when comparing OTU clustering methods
77 with ASV approaches, the latter could much more accurately reproduce a known diversity from mock
78 communities (Callahan *et al.*, 2016; Nearing *et al.*, 2018; Xue *et al.*, 2018).

79 Thus far, ASVs have never been directly compared with OTUs obtained from local clustering
80 approaches. One of the currently most widely used local clustering algorithms is SWARM (Mahé *et*
81 *al.*, 2014, 2015). In contrast to heuristic global clustering algorithms, SWARM relies on single-linkage
82 clustering and is input-order-independent, which results in much more robust OTU calling.
83 Furthermore, its high clustering stringency allows separation of even highly similar sequences (Mahé
84 *et al.*, 2014). This stringency is reflected by a small local clustering threshold that is by default set to
85 one nucleotide difference between two aligned sequences. SWARM-OTUs are created one after the
86 other by adding sequences in an iterative process. As long as sequences with equal or less than the

87 set nucleotide difference to any sequence already grouped into the OTU remain in the dataset, these
88 sequences will be added and the OTU will not be closed. That is, SWARM does not only avoid the
89 disadvantages of greedy heuristic global clustering algorithms, but also enables a very fine-grained
90 grouping of sequences. This generates distinct OTUs that differ in as little as two nucleotides from
91 one another. These advantages were recognized by several important studies on microbial
92 communities, which have further demonstrated that SWARM scales exceptionally well even to the
93 largest HTS datasets available to date (de Vargas *et al.*, 2015; Mahé *et al.*, 2017).

94 To further improve the congruency between species diversity and molecular diversity within
95 a sample, several authors have proposed a two-level sequence grouping approach, employing either
96 ASVs (Anslan *et al.*, 2018; Jusino *et al.*, 2018; Palmer *et al.*, 2018) or OTUs (Forster *et al.*, 2016) as a
97 first-level of sequence grouping. For instance, single-linkage first-level sequence grouping followed
98 by a subsequent second round of sequence grouping emerged as the most accurate strategy for
99 defining OTUs in a study by Huse *et al.* (2010). The goal of this combined two-level sequence
100 grouping approach is to fine-tune the obtained ASVs or OTUs to further improve biodiversity
101 estimates for accurate species richness predictions within a sample. A very promising example for a
102 second-level sequence grouping includes sequence similarity networks (SSNs; Forster *et al.*, 2015).
103 Sequence grouping in SSNs is achieved via pairwise sequence similarity scores. Two sequences are
104 connected by an edge, if their similarity passes a defined value. A group of sequences connected in
105 such a manner forms one connected component within the network, which is further interpreted as
106 a network sequence cluster (NSC), the result of second-level sequence grouping in SSNs. Since the
107 network approach is based on concepts from graph theory, there exists a full mathematical toolkit to
108 evaluate central ecological and evolutionary theories (*e.g.* Forster *et al.*, 2015; Corel *et al.*, 2016; Lord
109 *et al.*, 2016). But it was not until recently that the strength of SSNs could be exploited for large HTS
110 datasets. The *all-versus-all* pairwise sequence alignments on which the approach relies are time-
111 intensive and computationally demanding (Bik *et al.*, 2012; Sun *et al.*, 2012), and the computational
112 power as well as the tools became available for routine analyses only in the past years. After applying

113 a first level of sequence grouping (OTUs or ASVs), a dataset is represented by fewer sequences and
114 computational demands will exponentially decrease. Using this groundwork, the trade-off between
115 computational demands and scientific benefits becomes less disadvantageous when SSNs are
116 deployed as a second level of sequence grouping. Therefore, this strategy could enable SSN analyses
117 of even the largest datasets while allowing a more in-depth analyses of OTUs than possible with
118 heuristic clustering algorithms.

119 Based on the available knowledge summarized above, we hypothesized that i) DADA2-
120 derived ASVs produce fewer molecular sequence clusters from individual protist species, and that ii)
121 sequence similarity networks as a second-level sequence grouping approach further improves the
122 congruency between species diversity and molecular diversity as revealed by HTS sequencing. Thus
123 far, these hypotheses have gone untested. If verified, DADA2 in combination with SSNs may allow for
124 better biodiversity estimates in molecular environmental diversity surveys and, consequently, for less
125 biased interpretations of diversity results obtained in such studies.

126 Our case study is based on a collection of 29 individual ciliate species, each of which was
127 independently processed from cultivation to Illumina sequencing of the V9 18S rDNA region and data
128 analysis. Ciliates were chosen as model organisms for unicellular eukaryotes, because they possess
129 morphological features that largely allow for clear species differentiation (but see also *e.g.* Kumar
130 and Foissner (2016) for cryptic ciliate species) and because bioinformatic delineation of ciliate species
131 has been thoroughly validated (*e.g.* Nebel *et al.*, 2011; Dunthorn *et al.*, 2014). Using the resulting HTS
132 datasets, we directly compared first-level sequence grouping by an ASV-producing pipeline based on
133 the algorithm DADA2 against first-level sequence grouping by an OTU-producing pipeline based on
134 the algorithm SWARM and assessed the degree to which both algorithms matched the expected
135 diversity (species richness of 29). For both algorithms we followed standard pipelines recommended
136 by the respective developers, available at <https://benjjneb.github.io/dada2/index.html> for DADA2
137 and at <https://github.com/frederic-mahe/swarm/wiki/Fred's-metabarcoding-pipeline> for SWARM.
138 We then applied sequence similarity networks as a second-level sequence grouping approach to

139 DADA2-derived ASV and Swarm-derived OTU results to further improve the degree to which the
140 expected species diversity can be obtained.

141

142 **Results**

143 *First-level sequence grouping*

144 Sequence grouping of the 29 species-specific ciliate datasets with the DADA2 pipeline
145 resulted in 101 different ASVs (Table 1). This represented a 3.5-fold overestimation of the species
146 diversity in the samples. For the species *Dexiotricha sp.*, *Metanophrys sp.*, *Trithigmostoma cucullulus*
147 and *Vorticella sp.*, DADA2 produced one ASV per species. For all other species, sequence grouping
148 with DADA2 resulted in two or more ASVs, eight of which yielded 5 or more ASVs. The most ASVs
149 were obtained for *Stentor coeruleus* (11). There was a moderate linear relationship between the
150 abundance of reads and the resulting number of ASVs in each species-specific dataset ($R^2=0.51$, p-
151 value < 0.001).

152 Sequence grouping of the 29 species-specific ciliate datasets with the SWARM pipeline
153 resulted in 229 different OTUs (Table 1), overestimating the species diversity of the samples by a
154 factor of 7.9. *Trithigmostoma cucullulus* was the only species for which SWARM produced one OTU
155 per species. Two OTUs were obtained for the species *Tetrahymena sp.* and *Tokophrya infusionum*. 22
156 of the species produced 5 or more OTUs, with most obtained for *Folliculina sp.* (18) and *Spathidium*
157 *ascendens* (19). There was a weak linear relationship between the abundance of reads and the
158 resulting number of OTUs in each species-specific dataset ($R^2=0.11$, p-value=0.045).

159 In direct comparison, DADA2 produced 2.3-times fewer ASVs than SWARM produced OTUs
160 (Table 1). We observed a weak linear relationship between the overestimation of DADA2 and the
161 overestimation of SWARM across all species ($R^2=0.31$, p-value=0.001). However, *Folliculina sp.* and
162 *Spathidium ascendens*, the two species with the largest overestimation in SWARM, resulted in only 7
163 and 4 ASVs in DADA2, respectively. By contrast, *Stentor coeruleus*, the species with the largest

164 overestimation in DADA2, also resulted in 12 SWARM-OTUs. For three species the number of ASVs
165 and OTUs was identical (*Tetrahymena sp.*, *Trithigmostoma cucullulus*, *Tokophrya infusionum*). For the
166 remaining 26 species fewer ASVs were produced than OTUs. There was not a single species for which
167 more OTUs than ASVs were produced.

168

169 *Quantitative evaluation of second-level sequence grouping*

170 Sequence clusters originating from second-level sequence grouping are generally designated
171 as network sequence clusters (NSCs; Figure 1). When resulting from DADA2-ASVs, these NSCs are
172 further described as ASV-NSCs, while NSCs resulting from SWARM-OTUs are further described as
173 OTU-NSCs. The sequence similarity network approach successfully reduced the amount of sequence
174 clusters from first-level sequence grouping in both DADA2 and SWARM. The reduction was higher at
175 lower sequence similarity binning levels, since more pairs of sequences could be connected in the
176 network, thus producing larger NSCs (*i.e.* comprising more first-level ASVs or OTUs, respectively; see
177 Figure 2 and Table S1). At a binning level of 99% similarity, SSNs could reduce the amount from 101
178 ASVs to 65 ASV-NSCs (35.6% reduction). While no reduction was possible for SWARM-derived OTUs
179 at the 99% binning level, the 229 SWARM-derived OTUs were reduced by more than half (53.3%
180 reduction) to 107 OTU-NSCs at the next-lower binning level of 98% similarity. DADA2-derived ASVs
181 were also reduced by more than half (51.5% reduction) at 98% similarity. At the lowest tested
182 binning level of 90%, SSNs could reduce the amount from 101 DADA2-derived ASVs to 16 ASV-NSCs
183 (84.2% reduction) and from 229 SWARM-derived OTUs to 15 OTU-NSCs (93.4% reduction).

184 Overall, the results of the DADA2-SSN combination were notably closer to the ideal of one
185 NSC per species than the results of the SWARM-SSN combination (Figure 2, Table S1). But for binning
186 levels lower than 95% similarity, the combination of SSNs with either DADA2 or SWARM produced
187 similar amounts of NSCs. Except for the SWARM-SSN results at 95% similarity, which matched the
188 ideally expected diversity of 29 species, all of these analyses underestimated the expected amount of

189 diversity. For most lower binning levels, SWARM-SSNs resulted in fewer NSCs than DADA2-SSNs, thus
190 diverging stronger from the expected amount of species than DADA2-SSNs (exceptions were 91% and
191 95% similarity). At all binning levels higher than 95% similarity, the SWARM-SSN combination
192 produced distinctively more NSCs than the DADA2-SSN combination. At 96% similarity, the DADA2-
193 SSN results (29 ASV-NSCs) matched the ideally expected diversity of 29 species. Up to similarity
194 binning levels of 98% similarity, the amounts of ASV-NSCs moderately increased (up to 49 NSCs,
195 meaning a 1.7-fold overestimation of the expected diversity). 99% similarity is the first level at which
196 more than twice as much ASV-NSCs were found as species were analyzed (65 ASV-NSCs, 2.2-fold
197 overestimation). The results at the 100% binning level were equivalent to the first-level sequence
198 grouping results of DADA2 without SSNs. By contrast, the amount of OTU-NSCs in the SWARM-SSN
199 approach increased much stronger at higher binning levels. When the binning level was set to 97%
200 similarity, which is a commonly used threshold for sequence grouping of ciliate data, the SWARM-
201 SSN approach produced 54 OTU-NSCs, which nearly doubles the expected diversity of the dataset
202 (1.9-fold overestimation). At 98% similarity, 107 OTU-NSCs were produced, which exceeded the
203 number of ASVs produced by DADA2 without a second-level grouping in SSNs (Figure 2).

204 The application of an abundance filter could further decrease the overestimation of diversity
205 by the second-level sequence grouping (Figure S1, Table S2). For each species-specific dataset, we
206 discarded every OTU that accounted for less than 0.01% of the total read abundance in that dataset.
207 Thereby, there was only little effect on the SSN results at binning levels lower than 95% similarity. At
208 higher binning levels, the effect was more apparent on the outcome of the SWARM-SSN approach.
209 The diversity overestimation at the 100% similarity binning level decreased from a 7.9-fold
210 overestimation without abundance filter, to a 4.5-fold overestimation when the abundance filter was
211 applied. For the DADA2-SSN approach, the diversity overestimation decreased at the same binning
212 level from 3.5-fold without abundance filter to 3-fold with abundance filter.

213

214 *Qualitative evaluation of second-level sequence grouping*

215 Species-specific datasets were used for first-level sequence grouping to preclude DADA2-
216 derived ASVs or SWARM-derived OTUs that contained sequences from more than one dataset (Figure
217 1). The preclusion did not apply for SSN analyses since representative sequences of either ASVs or
218 OTUs from all species were pooled for *all-versus-all* pairwise sequence analyses during second-level
219 sequence grouping. Therefore, we distinguished between three different types of NSCs (Figure 1): i)
220 closed NSCs, which contained all DADA2-derived ASVs or SWARM-derived OTUs of one species and
221 none from any other species; ii) open NSCs, which contained some, but not all, DADA2-derived ASVs
222 or SWARM-derived OTUs of one species and none from any other species; iii) hybrid NSCs, which
223 contained DADA2-derived ASVs or SWARM-derived OTUs of at least two different species. The goal is
224 to determine which approach maximizes the number of closed NSCs and minimizes the number of
225 hybrid NSCs for each species-specific dataset. The distinction between hybrid, open and closed NSCs
226 should not be confounded with the terminology used for defining reference databases in some OTU
227 clustering methods (for more information on the latter see *e.g.* Bik *et al.*, 2012). Figures 3A and B
228 illustrate the results at each binning level from 90 to 100% sequence similarity (see also Table S3).
229 Independent of the first-level sequence grouping algorithm, most closed NSCs were produced at
230 intermediate binning levels. Except for 91% similarity, the DADA2-SSN approach produced at each
231 binning level more closed NSCs than the SWARM-SSN approach. The maxima of 24 closed ASV-NSCs
232 and 20 closed OTU-NSCs were detected at a binning level of 94% similarity. The least number of
233 closed NSCs was in both cases detected at the 100% similarity level (4 closed ASV-NSCs, 1 closed
234 OTU-NSCs). Most hybrid NSCs were produced at low similarity binning levels. At 90% similarity we
235 detected maxima of 4 hybrid ASV-NSCs and OTU-NSCs. In general, both approaches produced similar
236 numbers of hybrid NSCs at each binning level, with decreasing numbers of hybrid NSCs at increasing
237 similarity binning levels. For the DADA2-SSN approach one hybrid NSC was consistently observed
238 even at binning levels up to 99% similarity; for the SWARM-SSN approach no more hybrid NSCs were
239 observed at binning levels higher than 97% similarity. By contrast, increasing similarity binning levels
240 led to increasing numbers of open NSC for both approaches. Open NSCs were observed in neither
241 approach up to a binning level of 94%. Starting from 2 open ASV-NSCs and 7 open OTU-NSCs at 95%

242 similarity, the numbers of open NSCs steadily increased until reaching 97 open ASV-NSCs and 228
243 open OTU-NSCs at 100% similarity. The maximal number of open ASV-NSCs was surpassed in the
244 SWARM-SSN approach already at a binning level of 98% similarity (amounting to 99 OTU-NSCs).

245 A detailed evaluation of the DADA2-SSN results (Figure 3A) revealed that at binning levels
246 from 94-99% sequence similarity, one hybrid ASV-NSC was repeatedly produced. This always
247 comprised *Hypotrachida sp.* and *Oxytricha granulifera* (as well as *Urospinula succisa* at binning levels
248 lower than 97%). Another hybrid ASV-NSC, comprising *Spathidium ascendens* and *Spathidium*
249 *foissneri*, was continuously detected until a binning level of 97%, after which the two species from
250 the same genus no longer formed a hybrid ASV-NSC. The hybrid NSC patterns of the SWARM-SSN
251 approach were similar to the ones observed in the DADA2-SSN approach. However, hybrid OTU-NSCs
252 comprising *Hypotrachida sp.*, *Oxytricha granulifera* and *Urospinula succisa* were only detected up to
253 95% similarity. Of these three, only *Hypotrachida sp.* and *Urospinula succisa* also aggregated into
254 hybrid OTU-NSCs at 96% and 96.5% similarity. Hybrid OTU-NSCs comprising *Spathidium ascendens*
255 and *Spathidium foissneri* were detected from 93-97% similarity. Starting from the 96% similarity
256 binning level, the amount of open OTU-NSCs was always higher than the amount of closed OTU-
257 NSCs, with increasing numbers of open OTU-NSCs and decreasing numbers of closed OTU-NSCs
258 towards higher similarity binning levels.

259 The cluster distribution of NSCs was compared against an artificial perfect NSC cluster
260 distribution, in which only closed NSCs existed. The results of this comparison are expressed as rand
261 index (RI) and adjusted rand index (ARI) values in Figure 4. For both indices, the DADA2-SSN
262 approach resulted in higher values than the SWARM-SSN approach. A notable exception are the
263 values for the binning level of 96% similarity at which the maxima for the SWARM-SSN approach
264 were observed. The maximum RI and ARI values for the DADA2-SSN approach were observed at the
265 94% sequence similarity binning level. The progression of RI values with a distinct plateau phase at
266 intermediate binning levels is similar for DADA2-SSNs and SWARM-SSNs. For DADA2-SSNs the phase
267 starts at 93% similarity (RI = 0.9871) and reaches until 97% similarity (RI = 0.9907). For SWARM-SSNs,

268 the plateau phase is shorter and reaches from 95% similarity (RI = 0.9873) to 97% similarity (RI =
269 0.985). At binning levels higher than 97%, RI values slowly start to decrease, while ARI values
270 drastically decrease. For both approaches, RI values were lowest at the 90% similarity and ARI values
271 were lowest at the 100% similarity binning level.

272 Applying an abundance-filter affected the quantitative output of the second-level sequence
273 grouping stronger than the qualitative output. The removal of low abundant DADA2-derived ASVs or
274 SWARM-derived OTUs resulted in first place in a general decrease of open NSCs, which was coupled
275 to a slight increase of closed NSCs at similarity binning levels of 97% and higher (Figures S2A and B).
276 But this did not lead to different trends comparing RI values calculated from abundance-filtered and
277 non-abundance-filtered data (Figure S3). There was, however, a trend towards higher ARI values for
278 the DADA-SSN approach and lower ARI values for the SWARM-SSN approach compared between
279 abundance-filtered and non-abundance-filtered data.

280 We also tested the complete sequence grouping workflow for DADA2 and SWARM
281 when merging all species-specific datasets into one ciliate dataset before first-level sequence
282 grouping. However, the results were not substantially different from those obtained when species-
283 specific datasets were used for first-level sequence grouping (Table S5, Figures S4 and S5). A stronger
284 effect could be observed when all species-specific datasets were merged before first-level sequence
285 grouping and, in addition, an abundance-filter was applied before second-level sequence grouping in
286 SSNs (Table S5). Since the initial idea of this study was to analyze distinct ciliate species in the same
287 way and each species-specific dataset indeed represented a sample of a distinct ciliate species, we
288 decided to focus on the results obtained from species-specific datasets.

289

290 **Discussion**

291 *First-level sequence grouping results overestimate species richness*

292 A major goal of most environmental HTS studies is to provide realistic estimates of
293 biodiversity within a sample. However, it is not a trivial task to place the tremendous amount of
294 resulting data into an ecologically meaningful context. One of the main difficulties is the delineation
295 of species based on HTS datasets. Ciliates have been used as model organisms for addressing this
296 problem in the past, but so far all comparisons of morphospecies richness and OTU richness within
297 the same samples have been incongruent (Bachy *et al.*, 2013; Grattepanche *et al.*, 2014; Stoeck *et al.*,
298 2014). While some studies relying on clone library and Sanger sequencing technologies reported an
299 underestimation of diversity (Bachy *et al.*, 2013; Grattepanche *et al.*, 2014), more recent studies
300 relying on short HTS reads reported a diversity overestimation (Stoeck *et al.*, 2014; Flynn *et al.*, 2015;
301 Clare *et al.*, 2016). The same trend can be inferred from the results of our study: the DADA2 as well
302 as the SWARM algorithm delineated highly similar sets of sequences obtained from species-specific
303 Illumina-sequencing datasets and assigned to the same taxonomic hit, into multiple ASVs or OTUs.
304 However, of the two tested first-level sequence grouping approaches, DADA2-derived ASVs came
305 much closer to the known number of species than SWARM-derived OTUs.

306 The lower level of diversity overestimation in DADA2 is a consequence of the algorithm's
307 error-model based approach (Callahan *et al.*, 2017; Knight *et al.*, 2018). Not only does DADA2 exclude
308 a large fraction of reads based on statistical models from the analyses before assigning ASVs, it also
309 performs a sequencing artifact correction. Together, these steps aim at removing all spurious reads
310 and retaining only those variants that are not a product of erroneous sequencing. Conclusions about
311 whether or not an ASV represents a true organismal variant are still hard to draw, since ciliates and
312 other microbial species are known for intraspecific and intra-individual sequence polymorphism (*e.g.*
313 Miao *et al.*, 2004; Coleman, 2005; Gong *et al.*, 2013; Wang *et al.*, 2017). There may be cases in which
314 DADA2 retains artificial sequences, just as well as there are some cases in which DADA2 discards true
315 organismal variants. For instance, previous work of Gong and colleagues (2007) reported an
316 intraspecific polymorphism with a mean sequence divergence of 1.6% for the SSU rRNA gene within
317 the ciliate genus *Gaestrostyla*. Given that their results were based on sequences of three individuals,

318 this may only represent a fraction of the complete intraspecific polymorphism for these organisms. In
319 our analyses, we detected only two ASVs for *Gaestostyla steinii* while we detected considerably more
320 ASVs of ciliates (e.g. 11 ASVs for *Stentor coeruleus*) for which no high variation rates of intraspecific
321 polymorphism have been reported (Kusch, 1998; Zhang *et al.*, 2012). Thus, it is possible that some of
322 the more divergent true organismal variants of *G.steinii* may have been error-corrected by DADA2.

323 In contrast to the DADA2-pipeline, the SWARM algorithm does not perform any denoising
324 but is purely designed for grouping sequences. As such, SWARM will use any sequence provided in
325 the input dataset without deciding whether or not the sequence may be a true organismal variant or
326 a sequencing artifact. Denoising is a very sensible and important step that has to be performed by
327 other bioinformatic tools, preferentially on results obtained after sequence grouping in SWARM
328 (Mahé *et al.*, 2014). SWARM's focus is on finding smallest differences between sequences and
329 enabling a very fine-scaled resolution of genetic diversity in a sample. As reflected by our results, this
330 is counterproductive when looking at alpha diversity or species richness within a sample. The
331 *Spathidium ascendens* dataset of 19 SWARM-derived OTUs was the most severe example of over-
332 splitting in our data. Since the dataset of *S. ascendens* was subjected to the same bioinformatic
333 treatment as the other datasets with standard SWARM parameters, it is unlikely that this high
334 SWARM-OTU number is an artifact of the SWARM algorithm. The inflation of OTUs may be caused by
335 intraspecific sequence polymorphism in *Spathidium ascendens*, but this feature has not been studied
336 for this species before. For other species which yielded high numbers of OTUs, intra-specific and
337 intra-individual sequence polymorphism is documented. Among the order Heterotrichida,
338 polymorphic sites were found within the V9 region of the 18S rDNA (Wang *et al.*, 2017), which could
339 explain our finding of 18 SWARM-OTUs for *Folliculina sp.*. Sequence polymorphism is also a
340 widespread feature in the genus *Paramecium* (Coleman, 2005). At least some part of the SWARM-
341 OTUs from *Paramecium bursaria* (10 OTUs) and *Paramecium tetraurelia* (9 OTUs) may reflect this
342 true intraspecific genetic diversity. Using dataset replicates, as proposed by Prosser (2010), might
343 help for deciding which sequence is artificial and which is a true organismal variant. But even this

344 strategy cannot be the *ultima ratio* for species-specific datasets from single cell sequencing, since
345 intra-individual genetic diversity may not contain all variants of intraspecific genetic diversity.

346

347 *Second-level sequence grouping with sequence similarity networks improves diversity estimates*

348 The shortcomings for species diversity estimations of both first-level sequence groupings
349 (DADA2-ASVs and SWARM-OTUs) were distinctively alleviated when sequence similarity networks
350 were used as second-level sequence grouping. Even though species richness was still overestimated,
351 the extent of overestimation decreased through the application of SSNs. The greater effect of SSNs,
352 in terms of reducing the amount of first-level sequence grouping results, could be observed for
353 SWARM. These findings corroborate predictions of an earlier study, in which SWARM and SSNs had
354 not been used in combination, but as two different means of sequence grouping (Forster *et al.*,
355 2016). Interestingly, though, no reduction of SWARM's first-level sequence grouping results could be
356 achieved at sequence similarity binning levels of 99% and higher. This result can be attributed to our
357 application of SWARM on short V9 gene regions of ciliates, which reach an average length of 120
358 nucleotides (Dunthorn *et al.*, 2012). Because a single nucleotide difference was used on these
359 sequences in SWARM's first-level sequence grouping step, SWARM-OTUs already contained all pairs
360 of sequences with sequences similarities higher than 99% to each other. Thus, a further reduction by
361 SSNs in second-level sequence grouping was not possible at this binning level. To rely on local
362 nucleotide differences and avoid global sequence similarity values (typically 97%) is a central aspect
363 of SWARM (Mahé *et al.*, 2014, 2015). Although this increases the resolution of genetic diversity, the
364 higher resolution comes at the cost of generating numerous OTUs characterized by similar
365 sequences. Sequence similarity networks are well suited for reducing the apparent number of OTUs,
366 since they offer a straightforward way of grouping SWARM-derived OTUs while treating each of them
367 as an equal entity and allowing further downstream comparisons of these entities (*e.g.* as in Forster
368 *et al.*, 2015, 2016). The same also applies to DADA2-derived ASVs. As outlined in the discussion of the
369 DADA2 approach, first-level sequence grouping in DADA2 resulted in a much smaller overestimation

370 of species richness than SWARM. Whilst there was inherently less room for further reduction of the
371 diversity overestimation by SSNs, the DADA2-SSN combination did produce diversity estimates that
372 better mirrored the species richness than those obtained from the SWARM-SSN combination.

373 Previous studies (*e.g.* Huse *et al.*, 2010; Bonder *et al.*, 2012) predicted that two-levels of
374 sequence grouping will have positive effects on the accuracy with which taxonomic units are defined
375 from HTS datasets, but without using sequence similarity networks. Huse and colleagues conducted
376 an initial single-linkage sequence grouping followed by a second-level sequence grouping, similar to
377 the combination of SWARM and SSNs. Although our results confirm their predictions about the
378 positive effect of the two-level sequence grouping strategy, the combination of SWARM (for single-
379 linkage first-level grouping) with SSNs (for second-level grouping) was not as accurate for defining
380 species-specific groups as the DADA2-SSN combination. The latter combination is somewhat similar
381 to the strategy used by Bonder and colleagues (2012): for a mock dataset of 15 species, they reduced
382 the number of OTUs by 93.4% when a denoising step was employed before final sequence grouping.
383 Likewise, our combination of DADA2 as a denoising step and SSNs for sequence grouping led to a
384 reduction of OTUs from first- to second-level grouping by 74.3% at the binning level with the highest
385 RI and ARI score (94% similarity). Further reduction of NSCs would be possible, but would not lead to
386 a qualitatively better output. The lower NSC reduction rate in the current study compared to the
387 study of Bonder and colleagues (2012) is merely a product of the more effective denoising by DADA2
388 (Callahan *et al.*, 2016; Knight *et al.*, 2018).

389 DADA2 has recently been introduced as an effective first-level sequence grouping and
390 denoising tool and combined with different sequence grouping algorithms (Frøslev *et al.*, 2017;
391 Anslan *et al.*, 2018; Jusino *et al.*, 2018; Palmer *et al.*, 2018), but none of the studies used SSNs.
392 Frøslev *et al.* (2017) decided to combine DADA2 with hierarchical sequence grouping in VSEARCH,
393 followed by a post-clustering treatment based on ecological patterns to remove erroneous groups of
394 sequences. In contrast to this, we decided to combine DADA2 with hierarchical sequence grouping in
395 VSEARCH and analyzed the resulting pairwise sequence similarities in SSNs. Without post-clustering

396 treatment, both approaches result in quite similar patterns of diversity overestimation, though to a
397 lesser extent by the DADA2-SSN approach on the ciliate dataset. When post-clustering is taken into
398 account, the approach of Frøslev and colleagues led to an underestimation of diversity by one third
399 at their suggested level of 97% sequence similarity. At the same level of similarity, the DADA2-SSN
400 combination overestimated the diversity by one fifth (6 ASV-NSCs more than species expected). In
401 addition, our approach led to more accurate diversity estimates and an underestimation of diversity
402 by only three ASV-NSCs at the binning level with the highest RI and ARI scores (94% similarity). The
403 marginal differences between RI and ARI scores at intermediate similarity binning levels (*e.g.* RI
404 difference of 0.0002 and ARI difference of 0.0026 between 94% and 95% similarity), indicate that
405 there is not one universally applicable level for optimal sequence grouping. Instead, there exists a
406 range of sequence similarities, at which qualitatively highly similar and equally precise sequence
407 grouping can be achieved. To identify the best binning level, we advise to create SSNs on such a
408 range of similarities and evaluate the outcome for best sequence grouping results. Our data suggests
409 that the use of more conservative and lower similarity binning levels has a positive effect towards
410 more accurate diversity estimates of ciliates while also allowing for more precise species delineation.
411 Similar conclusions can be drawn from the SWARM-SSN approach, which yielded the best results at
412 96% similarity and thus, below the 97% similarity threshold widely used for clustering ciliate
413 sequence data. While our two-level sequence grouping strategy can easily be adapted to datasets of
414 other taxonomic groups or other barcode gene regions, it is important to note that the similarity
415 binning level, which produced the best output for our ciliate V9 18S rDNA dataset, is not
416 generalizable to other datasets without prior tests. The extent of genetic diversity varies among
417 different taxonomic groups (*e.g.* Brown *et al.*, 2015) and even within ciliates, different sequence
418 similarity thresholds have been shown to be more effective for delineating species when working
419 with datasets of different hypervariable regions of the 18S rDNA (Dunthorn *et al.*, 2012).

420 *All-versus-all* pairwise sequence alignments in hierarchical sequence grouping are a
421 prerequisite for attaining the advantages of sequence similarity networks (Forster *et al.*, 2015, 2016;

422 Corel *et al.*, 2016). One of the main benefits for using hierarchical instead of heuristic sequence
423 grouping is that an over-splitting of diversity is avoided (Mahé *et al.*, 2014; Flynn *et al.*, 2015). This
424 benefit displays in the avoidance of diversity over-splitting in our study as well. By contrast, an over-
425 splitting was observed in the study of Frøslev *et al.* (2017) when no additional post-clustering of the
426 sequence grouping results was conducted. A further benefit of using networks for second-level
427 grouping is that they allow for detailed evaluation of the information provided within each NSC. For
428 instance, the persistent grouping of *Spathidium ascendens* and *Spathidium foissneri* in hybrid NSCs is
429 not only explained by them belonging to the same genus, but can be further related to rapid
430 radiation events and incomplete lineage sorting in the order Spathidiida (Vďačný *et al.*, 2014). By
431 considering the internal structure of NSCs, one can thus draw additional ecological and evolutionary
432 conclusions about the organisms under study.

433 Although SSNs emerged as a powerful tool for improving first-level sequence grouping
434 results, we have to state that it remained impossible to perfectly reproduce the diversity in the
435 samples. This conclusion is not unexpected, because evolutionary processes can only be
436 approximated, but not predicted by bioinformatic algorithms and molecular proxies. Different
437 lineages may evolve at different rates (Brown *et al.*, 2015) or gene transfer may occur (Baptiste and
438 Boucher, 2008), all of which complicate the evaluation of sequencing data and the estimation of
439 diversity from an environmental community dataset. There is, however, still room for improvement
440 and drawing even more accurate pictures when working with SSNs as second-level sequence
441 grouping. Our results indicated that SWARM-derived OTUs are strongly affected by subsequent
442 denoising steps. For SWARM sequence grouping and denoising, we followed the same strategy used
443 in several benchmarks studies (*e.g.* de Vargas *et al.*, 2015; Mahé *et al.*, 2017), but still observed a
444 distinct overestimation of species richness. To limit the overestimation, we propose additional
445 filtering steps relying on sequence abundances. Other studies have shown that abundance filter can
446 efficiently remove noise from HTS datasets (Quince *et al.*, 2009, 2011; Reeder and Knight, 2010;
447 Bokulich *et al.*, 2013; Auer *et al.*, 2017). Likewise, an abundance-filter step is also implemented in

448 DADA2 (Callahan *et al.*, 2016). Applying an abundance filter had a positive effect on the outcome of
449 the DADA2-SSN and especially on the outcome of the SWARM-SSN strategy (see Figures S1, S2, S3
450 and Table S2). The effect was further enhanced when species-specific datasets were merged before
451 first-level sequence grouping and an abundance-filter was applied before second-level sequence
452 grouping in SSNs. But even without merging the datasets, the ASV-NSCs and OTU-NSCs reflected the
453 real species diversity quantitatively closer after the application of an abundance filter. The filter had
454 little effect, however, on the formation of hybrid and closed NSCs. Most of all, the formation of open
455 NSCs was distinctively reduced. Since open NSCs result from ASV or OTU variants, which cannot be
456 affiliated with ASV or OTU variants from the same species, this implies that many open NSCs were
457 actually low abundant variants of a given species. They possibly emerge from either sequencing
458 errors or real but rare genetic polymorphisms. In how far the application of the abundance filter and
459 the removal of these variants can be justified with regard to an environmental study and become a
460 standard step of the second-level sequence grouping by SSN remains to be tested in the future.

461

462 **Conclusion**

463 The consequences and relevance of notably different outputs using DADA2 and SWARM in
464 ecological studies remains to be tested in real case scenarios using field samples. While some specific
465 ecological questions (such as results of beta diversity analyses) may not suffer, others (such as results
466 of alpha diversity analyses or the identification of ecologically relevant key species) may be highly
467 compromised by different information obtained from DADA2 and SWARM. Beyond that, second-level
468 sequence grouping with sequence similarity networks clearly improved diversity estimates compared
469 to first-level sequence grouping. We expect that future studies will benefit from implementing this
470 strategy, especially by relying on the combination of DADA2 and sequence similarity networks.

471

472 **Experimental procedures**

473 *Ciliate specimen collection*

474 Single ciliate cells were hand-picked from either pure cell cultures or from environmental
475 culture samples, leading to a collection of 29 different ciliate species (Table 1, see Table S4 for
476 material sources). To allow for testing species delineation at higher taxonomic levels, six of the
477 known ciliate classes were covered, including Heterotrichea, Litostomatea, Oligohymenophorea,
478 Phyllopharyngea, Prostomatea and Spirotrichea. To test the species delineation on lower taxonomic
479 levels, the collection also included 2 species from each of the genera *Paramecium* and *Spathidium*.
480 Each of the 29 species was treated independently in our workflow to create one species-specific
481 high-throughput sequencing dataset per ciliate for downstream analyses (Figure 1).

482

483 *DNA extraction, amplification and high-throughput sequencing*

484 For each ciliate species, DNA was extracted from individually picked specimens from pure
485 cultures, enrichments or environmental samples using Qiagen's DNeasy Blood & Tissue Kit, followed
486 by a clean-up step with Qiagen's MinElute PCR Purification Kit. Both kits were used according to the
487 manufacturer's instructions (Qiagen; Hilden, Germany). Extracted DNA was then amplified in a semi-
488 nested PCR targeting the hypervariable V9 region of the 18S rDNA. The V9 is a routinely used
489 barcode gene region for eukaryotic community analyses (de Vargas *et al.*, 2015) and has been
490 specifically tested for delineation of ciliate species (Dunthorn *et al.*, 2012, 2014). In the first step of
491 the semi-nested PCR, we used the ciliate specific forward primer CilF (5'-
492 TGGTAGTGTATTGGACWACCA-3'; Lara *et al.*, 2007) in combination with the universal eukaryotic 18S
493 reverse primer EukB (5'-TGATCCTTCTGCAGGTTACCTAC-3'; Medlin *et al.*, 1988). This initial step was
494 especially important for ciliate species from environmental samples to avoid amplification of non-
495 target eukaryotic organisms. The CilF-EukB protocol consisted of an initial denaturation step of 30
496 seconds (s) at 98°C, followed by 35 cycles of 10 s at 98°C, 30 s at 62°C and 60 s at 72°C; the final
497 extension lasted for 300 s at 72°C. After the first step of the semi-nested protocol, resulting PCR
498 products were cleaned once more with Qiagen's MinElute PCR Purification Kit according to the
499 manufacturer's instructions. In the second step of the semi-nested PCR, the CilF-EukB products were

500 amplified with a universal eukaryotic V9 forward primer (5'-GTACACACCGCCCGTC-3'; Lane, 1991) and
501 EukB as reverse primer. Each primer pair was tagged with one of ten different barcodes, so that the
502 resulting species-specific barcoded sequences could easily be retraced in downstream analyses. The
503 V9-EukB protocol consisted of an initial denaturation step of 30 s at 98°C, followed by 25-35 cycles
504 (adjusted for each species, see Table S4) of 10 s at 98°C, 20 s at 64°C and 25 s at 72°C; the final
505 extension lasted for 300 s at 72°C. Once the barcoded V9 PCR products were successfully amplified,
506 samples were pooled into libraries that contained up to ten different barcodes (one for each ciliate
507 species). Libraries were paired-end sequenced on either an Illumina MiSeq or NextSeq platform (see
508 Table S4 for details), each generating 250-base pair reads. High-throughput sequencing was
509 conducted by SeqIT (Kaiserslautern, Germany).

510 After sequencing, libraries were split into single species datasets using CUTADAPT v1.18
511 (Martin, 2011) according to the barcodes initially applied during the semi-nested PCR. Filtering of
512 reads with an expected barcode and removal of those was followed by matching the primer
513 sequences at both 5' and 3'-ends in CUTADAPT, as well. Reads were kept if they exactly matched the
514 forward or reverse V9 primers at their 5'-end and, at the same time, if they exactly matched the
515 reverse or forward V9 primers near to the 3'-end ("linked adapter" approach in CUTADAPT). Primer
516 sequences and overhang at 5'-end were removed during this filtering process to keep only the
517 targeted V9 region. Reads were oriented in the same direction using VSEARCH v2.8.0 (Rognes *et al.*,
518 2016). Thus, the first paired libraries only contained reads oriented in the forward V9 primer
519 direction and the second paired libraries only contained reads oriented in the reverse V9 primer
520 direction. The barcode- and primer-filtering produced 29 single species paired libraries (for the
521 species *Epistylis plicatilis* two libraries existed, but the data were merged for subsequent steps)
522 which were used as input for the downstream first-level sequence grouping approaches.

523 The sequence data were deposited at the Sequence Read Archive of the National Center for
524 Biotechnology Information (NCBI) and are available under accession number PRJNA548847. All

525 bioinformatic procedures and commands for statistical analyses are provided in HTML format as
526 supplemental material (File S1, Table S6).

527

528 *Quality filtering and first-level sequence grouping with DADA2*

529 Reads were truncated to their first 100 bp and filtered with a maximum expected error of 0.2
530 using the DADA2 package v1.8 in R v3.5.1 (Callahan *et al.*, 2016; R Core Team, 2018). Then error rates
531 were learned for each sequencing run separately, reads were dereplicated per single species paired
532 libraries and the core DADA2 denoising algorithm was applied on each read of each single species'
533 paired library using their respective sequencing run error model. Denoised reads were paired for
534 each single species by requesting a minimum overlap of 50 bp and by allowing a maximum of 5
535 mismatches using the function *mergePairs* in DADA2. ASVs produced in this manner were merged
536 among all 29 single species libraries and chimeras were *de-novo* removed from this ASV table using
537 the consensus approach from the *removeBimeraDenovo* function.

538

539 *Quality filtering and first-level sequence grouping with SWARM*

540 Demultiplexed libraries, from which barcodes and primers were previously removed, were
541 used for first-level sequence grouping with SWARM. For each paired end, species-specific library, we
542 precisely followed the publicly available instructions at [https://github.com/frederic-](https://github.com/frederic-mahe/swarm/wiki/Fred's-metabarcoding-pipeline)
543 [mahe/swarm/wiki/Fred's-metabarcoding-pipeline](https://github.com/frederic-mahe/swarm/wiki/Fred's-metabarcoding-pipeline). In short, paired-end reads were merged and
544 subsequently dereplicated into amplicons using VSEARCH. First-level sequence grouping in SWARM
545 v2.0.5 (Mahé *et al.*, 2015) was performed on the amplicons with *-d 1* and the fastidious option *-f*. The
546 resulting SWARM-OTUs were subjected to a *de novo* chimera detection in UCHIME (Edgar *et al.*,
547 2011) and singletons were removed from the output. From each non-chimeric SWARM-OTU we
548 extracted the seed amplicon as representative sequence for downstream taxonomic assignment.

549

550 *Taxonomic assignment of first-level sequence grouping results*

551 ASVs and representative seed sequences (*i.e.*, the most abundant amplicon) of SWARM OTUs
552 were annotated using VSEARCH (Altschul *et al.*, 1990) against a modified version of the PR2-derived
553 database provided in de Vargas *et al.*, 2015. This reference database was specifically designed for
554 taxonomic assignment of V9 sequences via VSEARCH (*i.e.* by containing only references trimmed to
555 the V9 region). We manually trimmed and added ciliate reference sequences that were still missing
556 in the database. The added references represent NCBI GenBank entries deposited under accession
557 numbers AB558117, AF429900, AF508776, FJ998037, KC991098, KF301567, KF411460, KF733753,
558 KF733756, KF878932, KU525298, MG589318. Only ASVs and OTUs assigned with at least 90%
559 sequence similarity to the most abundant ciliate species in each species-specific library were kept for
560 further analyses, while sequences of *e.g.* prey, parasitic or mutualistic organisms present in the
561 samples were discarded.

562

563 *Second-level sequence grouping with sequence similarity networks*

564 Identical second-level sequence grouping steps were performed independently for the
565 DADA2 and SWARM first-level sequence grouping outputs (Figure 1). Representative sequences of all
566 species-specific target ASVs or OTUs were first pooled to create one dataset for either DADA2 and
567 SWARM. The initial step for sequence similarity network construction employed *all-versus-all*
568 pairwise sequence analyses of these datasets in VSEARCH (Rognes *et al.*, 2016) using the settings -
569 *allpairs_global, -iddef 1* and a similarity cutoff of 90%. This resulted in an edge table in which each
570 line represented a pair of sequences that shared a sequence similarity of at least 90% to each other.
571 From the edge table, unweighted and undirected SSNs were calculated in R with the package *igraph*
572 v1.2.2 (Csardi and Nepusz, 2006). To determine the sequence similarity which gave the maximal
573 congruence between the number of NSCs and the number of ciliate species, similarity binning levels
574 for SSN construction from 90%-100% were tested in single percentage steps, except between 96%

575 and 98% for which we tested 0.5 percentage steps. This was because previous studies suggested the
576 best cutoff level for species delineation at this range of sequence similarities (*e.g.* Worden, 2006;
577 Caron *et al.*, 2009). Every node in a SSN represented one DADA2-derived ASV or one SWARM-derived
578 OTU representative sequence. Every edge in a SSN represented a sequence similarity between two
579 ASVs or two representative sequences that exceeded the applied binning level (Figure 1, see also
580 Forster *et al.*, 2015). Thus, NSCs either represented connected components, *i.e.* a cluster of ASVs or
581 OTU representative sequences that could be further grouped based on the applied binning level; or a
582 single node, *i.e.* a single ASV or OTU representative sequence which sequence similarity to any other
583 ASV or OTU representative sequence in the dataset was lower than the applied binning level. For
584 instance, a binning level of 100% similarity reproduced the results from the first-level sequence
585 grouping, meaning that the SSNs consisted exclusively of single nodes that represented DADA2-
586 derived ASVs or SWARM-derived OTUs.

587 Previous studies indicated that the implementation of an abundance filter can be beneficial
588 for increasing the accuracy of sequence grouping (*e.g.* Quince *et al.*, 2009, 2011; Reeder and Knight,
589 2010; Bokulich *et al.*, 2013; Auer *et al.*, 2017). The rationale behind this is that sequences which
590 emerge from organisms that are actually occurring in a sample should be much more abundant in a
591 HTS dataset, than sequences which represent sequencing artifacts or elusive contaminations.
592 Following this idea, we tested in a separate approach if an abundance filter could improve the
593 outcome of our two-level sequence grouping strategy. For this test, we removed before SSN
594 construction from each species' dataset all DADA2-derived ASVs with an abundance of less than
595 0.01% with regard to all sequences of that dataset. For abundance filtering in SWARM we first set
596 SWARM's *-b* option to the species-specific 0.01% abundance threshold of each dataset during first-
597 level sequence grouping, then removed before SSN construction all SWARM-OTUs with an
598 abundance of less than 0.01% with regard to all sequences of that dataset. All other second-level
599 sequence grouping steps in this test were performed as outlined above.

600

601 *Statistical evaluation of sequence grouping*

602 All statistical tests were run in R v3.5.1. The focus of the statistical evaluation for second-
603 level sequence grouping was to identify the similarity binning level, which maximized the number of
604 closed NSCs and minimized the number of hybrid NSCs. To express this level mathematically, we
605 listed the NSC membership of each representative sequence and applied both the rand index (RI;
606 Rand, 1971) and Hubert's and Arabie's adjusted rand index (ARI; Hubert and Arabie, 1985). In short,
607 the RI compares the congruence between two cluster distributions, whereas the ARI is the corrected-
608 by-chance version of the RI. The values of the indices range from 0 to 1, with 0 describing completely
609 different cluster distributions and 1 describing perfectly matching cluster distributions between two
610 sets of data. In this study, we compared the sequence grouping of representative sequences within
611 NSCs for each binning level to the optimal distribution of representative sequences, in which case
612 DADA2-derived ASVs and SWARM-derived OTUs exclusively form species-specific, closed NSCs (that
613 is 29 NSCs, one for each species). RI and ARI were calculated with the R package *clues* v0.5.9 (Chang
614 *et al.*, 2010).

615

616 **Acknowledgements**

617 We thank Gianna Pitsch, Bettina Sonntag, Thomas Pröschold, Ewa Przyboś and the staff of CCAP
618 (SAMS, Oban, Scotland) for providing cultured ciliate species. Furthermore, we thank Hans-Werner
619 Breiner and Sara Neß for their help with laboratory works. We are grateful for the constructive
620 criticism of three anonymous reviewers. This study was supported by a research grant of the Carl
621 Zeiss Stiftung to DF and an ASSEMBLE PLUS project to TS under European Union's Horizon 2020
622 research and innovation program (grant agreement no. 730984).

623 The authors declare that there is no conflict of interest.

624 **References**

- 625
- 626 Allali, I., Arnold, J.W., Roach, J., Cadenas, M.B., Butz, N., Hassan, H.M., *et al.* (2017) A comparison of
627 sequencing platforms and bioinformatics pipelines for compositional analysis of the gut
628 microbiome. *BMC Microbiol.* **17**: 194.
- 629 Altschul, S.F., Gish, W., Miller, W., Myers, E.W., and Lipman, D.J. (1990) Basic local alignment search
630 tool. *J. Mol. Biol.* **215**: 403–410.
- 631 Amir, A., McDonald, D., Navas-Molina, J.A., Kopylova, E., Morton, J.T., Xu, Z.Z., *et al.* (2017) Deblur
632 rapidly resolves single-nucleotide community sequence patterns. *mSystems* **2**: e00191-16.
- 633 Anslan, S., Nilsson, R.H., Wurzbacher, C., Baldrian, P., Leho Tedersoo, and Bahram, M. (2018) Great
634 differences in performance and outcome of high-throughput sequencing data analysis
635 platforms for fungal metabarcoding. *MycoKeys* **29**–40.
- 636 Auer, L., Mariadassou, M., O’Donohue, M., Klopp, C., and Hernandez-Raquet, G. (2017) Analysis of
637 large 16S rRNA Illumina data sets: Impact of singleton read filtering on microbial community
638 description. *Mol. Ecol. Resour.* **17**: e122–e132.
- 639 Bachy, C., Dolan, J.R., López-García, P., Deschamps, P., and Moreira, D. (2013) Accuracy of protist
640 diversity assessments: morphology compared with cloning and direct pyrosequencing of 18S
641 rRNA genes and ITS regions using the conspicuous tintinnid ciliates as a case study. *ISME J.* **7**:
642 244–255.
- 643 Baptiste, E. and Boucher, Y. (2008) Lateral gene transfer challenges principles of microbial
644 systematics. *Trends Microbiol.* **16**: 200–207.
- 645 Benson, D.A., Cavanaugh, M., Clark, K., Karsch-Mizrachi, I., Lipman, D.J., Ostell, J., and Sayers, E.W.
646 (2013) GenBank. *Nucleic Acids Res.* **41**: D36–D42.
- 647 Bik, H.M., Porazinska, D.L., Creer, S., Caporaso, J.G., Knight, R., and Thomas, W.K. (2012) Sequencing
648 our way towards understanding global eukaryotic biodiversity. *Trends Ecol. Evol.* **27**: 233–
649 243.
- 650 Bokulich, N.A., Subramanian, S., Faith, J.J., Gevers, D., Gordon, J.I., Knight, R., *et al.* (2013) Quality-
651 filtering vastly improves diversity estimates from Illumina amplicon sequencing. *Nat.*
652 *Methods* **10**: 57–59.
- 653 Bonder, M.J., Abeln, S., Zaura, E., and Brandt, B.W. (2012) Comparing clustering and pre-processing in
654 taxonomy analysis. *Bioinformatics* **28**: 2891–2897.
- 655 Brown, E.A., Chain, F.J.J., Crease, T.J., Maclsaac, H.J., and Cristescu, M.E. (2015) Divergence
656 thresholds and divergent biodiversity estimates: can metabarcoding reliably describe
657 zooplankton communities? *Ecol. Evol.* **5**: 2234–2251.
- 658 Callahan, B.J., McMurdie, P.J., and Holmes, S.P. (2017) Exact sequence variants should replace
659 operational taxonomic units in marker-gene data analysis. *ISME J.* **11**: 2639–2643.
- 660 Callahan, B.J., McMurdie, P.J., Rosen, M.J., Han, A.W., Johnson, A.J.A., and Holmes, S.P. (2016)
661 DADA2: High-resolution sample inference from Illumina amplicon data. *Nat. Methods* **13**:
662 581–583.
- 663 Caron, D.A., Countway, P.D., Savai, P., Gast, R.J., Schnetzer, A., Moorthi, S.D., *et al.* (2009) Defining
664 DNA-based operational taxonomic units for microbial-eukaryote ecology. *Appl. Environ.*
665 *Microbiol.* **75**: 5797–5808.
- 666 Chang, F., Qiu, W., Zamar, R.H., Lazarus, R., and Wang, X. (2010) clues: An R package for
667 nonparametric clustering based on local shrinking. *J. Stat. Softw.* **033**.
- 668 Clare, E.L., Chain, F.J.J., Littlefair, J.E., and Cristescu, M.E. (2016) The effects of parameter choice on
669 defining molecular operational taxonomic units and resulting ecological analyses of
670 metabarcoding data. *Genome* **59**: 981–990.
- 671 Coleman, A.W. (2005) *Paramecium aurelia* revisited. *J. Eukaryot. Microbiol.* **52**: 68–77.
- 672 Corel, E., Lopez, P., Méheust, R., and Baptiste, E. (2016) Network-thinking: Graphs to analyze
673 microbial complexity and evolution. *Trends Microbiol.* **24**: 224–237.
- 674 Csardi, G. and Nepusz, T. (2006) The igraph software package for complex network research. *InterJ*
675 *Complex Systems* **1695**: 1–9.
- 676 Dunthorn, M., Klier, J., Bunge, J., and Stoeck, T. (2012) Comparing the hyper-variable V4 and V9

677 regions of the small subunit rDNA for assessment of ciliate environmental diversity. *J.*
678 *Eukaryot. Microbiol.* **59**: 185–187.

679 Dunthorn, M., Otto, J., Berger, S.A., Stamatakis, A., Mahé, F., Romac, S., *et al.* (2014) Placing
680 environmental next-generation sequencing amplicons from microbial eukaryotes into a
681 phylogenetic context. *Mol. Biol. Evol.* **31**: 993–1009.

682 Edgar, R.C. (2010) Search and clustering orders of magnitude faster than BLAST. *Bioinformatics* **26**:
683 2460–2461.

684 Flynn, J.M., Brown, E.A., Chain, F.J.J., Maclsaac, H.J., and Cristescu, M.E. (2015) Toward accurate
685 molecular identification of species in complex environmental samples: testing the
686 performance of sequence filtering and clustering methods. *Ecol. Evol.* **5**: 2252–2266.

687 Forster, D., Bittner, L., Karkar, S., Dunthorn, M., Romac, S., Audic, S., *et al.* (2015) Testing ecological
688 theories with sequence similarity networks: marine ciliates exhibit similar geographic
689 dispersal patterns as multicellular organisms. *BMC Biol.* **13**: 16.

690 Forster, D., Dunthorn, M., Stoeck, T., and Mahé, F. (2016) Comparison of three clustering approaches
691 for detecting novel environmental microbial diversity. *PeerJ* **4**: e1692.

692 Frøslev, T.G., Kjølner, R., Bruun, H.H., Ejrnæs, R., Brunbjerg, A.K., Pietroni, C., and Hansen, A.J. (2017)
693 Algorithm for post-clustering curation of DNA amplicon data yields reliable biodiversity
694 estimates. *Nat. Commun.* **8**: 1188.

695 Fu, L., Niu, B., Zhu, Z., Wu, S., and Li, W. (2012) CD-HIT: accelerated for clustering the next-generation
696 sequencing data. *Bioinformatics* **28**: 3150–3152.

697 Gong, J., Dong, J., Liu, X., and Massana, R. (2013) Extremely high copy numbers and polymorphisms
698 of the rDNA operon estimated from single cell analysis of oligotrich and peritrich ciliates.
699 *Protist* **164**: 369–379.

700 Gong, J., Kim, S.-J., Kim, S.-Y., Min, G.-S., Roberts, D.M., Warren, A., and Choi, J.-K. (2007) Taxonomic
701 redescrptions of two ciliates, *Protogastrostyla pulchra n. g., n. comb.* and *Hemigastrostyla*
702 *enigmatica* (Ciliophora: Spirotrichea, Stichotrichia), with phylogenetic analyses based on 18S
703 and 28S rRNA gene sequences. *J. Eukaryot. Microbiol.* **54**: 468–478.

704 Grattepanche, J.-D., Santoferrara, L.F., McManus, G.B., and Katz, L.A. (2014) Diversity of diversity:
705 conceptual and methodological differences in biodiversity estimates of eukaryotic microbes
706 as compared to bacteria. *Trends Microbiol.* **22**: 432–437.

707 Hoshina, R., Hayashi, S., and Imamura, N. (2006) Intraspecific genetic divergence of *Paramecium*
708 *bursaria* and re-construction of the paramecian phylogenetic tree. *Acta Protozool.* **45**: 377–
709 386.

710 Hubert, L. and Arabie, P. (1985) Comparing partitions. *J. Classification* **2**: 193–218.

711 Huse, S.M., Welch, D.M., Morrison, H.G., and Sogin, M.L. (2010) Ironing out the wrinkles in the rare
712 biosphere through improved OTU clustering. *Environ. Microbiol.* **12**: 1889–1898.

713 Jusino, M.A., Banik, M.T., Palmer, J.M., Wray, A.K., Xiao, L., Pelton, E., *et al.* (2018) An improved
714 method for utilizing high-throughput amplicon sequencing to determine the diets of
715 insectivorous animals. *Mol. Ecol. Resour.* **19**: 176–190.

716 Knight, R., Vrbanac, A., Taylor, B.C., Aksenov, A., Callewaert, C., Debelius, J., *et al.* (2018) Best
717 practices for analysing microbiomes. *Nat. Rev. Microbiol.* **16**: 410.

718 Kumar, S. and Foissner, W. (2016) High cryptic soil ciliate (Ciliophora, Hypotrichida) diversity in
719 Australia. *Eur. J. Protistol.* **53**: 61–95.

720 Kunin, V., Engelbrektson, A., Ochman, H., and Hugenholtz, P. (2010) Wrinkles in the rare biosphere:
721 pyrosequencing errors can lead to artificial inflation of diversity estimates. *Environ. Microbiol.*
722 **12**: 118–123.

723 Kusch, J. (1998) Local and temporal distribution of different genotypes of pond-dwelling *Stentor*
724 *coeruleus*. *Protist* **149**: 147–154.

725 Landis, W.G. (1986) The interplay among ecology, breeding systems, and genetics in the *Paramecium*
726 *aurelia* and *Paramecium bursaria* complexes. *Prog. Protistol.* **1**: 287–307.

727 Lane, D.J. (1991) 16S/23S rRNA sequencing. In, Stackebrandt, E. and Goodfellow, M. (eds), *Nucleic acid*
728 *techniques in bacterial systematics*. John Wiley and Sons: Chichester, UK, pp. 115–175.

729 Lara, E., Berney, C., Harms, H., and Chatzinotas, A. (2007) Cultivation-independent analysis reveals a

730 shift in ciliate 18S rRNA gene diversity in a polycyclic aromatic hydrocarbon-polluted soil.
731 *FEMS Microbiol. Ecol.* **62**: 365–373.

732 Lord, E., Cam, M.L., Baptiste, É., Méheust, R., Makarenkov, V., and Lapointe, F.-J. (2016) BRIDES: A
733 new fast algorithm and software for characterizing evolving similarity networks using
734 Breakthroughs, Roadblocks, Impasses, Detours, Equals and Shortcuts. *PLOS ONE* **11**:
735 e0161474.

736 Mahé, F., Rognes, T., Quince, C., Vargas, C. de, and Dunthorn, M. (2014) Swarm: robust and fast
737 clustering method for amplicon-based studies. *PeerJ* **2**: e593.

738 Mahé, F., Rognes, T., Quince, C., Vargas, C. de, and Dunthorn, M. (2015) Swarm v2: highly-scalable
739 and high-resolution amplicon clustering. *PeerJ* **3**: e1420.

740 Martin, M. (2011) Cutadapt removes adapter sequences from high-throughput sequencing reads.
741 *EMBnet.journal* **17**: 10–12.

742 Masella, A.P., Bartram, A.K., Truszkowski, J.M., Brown, D.G., and Neufeld, J.D. (2012) PANDAseq:
743 paired-end assembler for illumina sequences. *BMC Bioinformatics* **13**: 31.

744 Medlin, L., Elwood, H.J., Stickel, S., and Sogin, M.L. (1988) The characterization of enzymatically
745 amplified eukaryotic 16S-like rRNA-coding regions. *Gene* **71**: 491–499.

746 Miao, W., Fen, W.-S., Yu, Y.-H., Zhang, X.-Y., and Shen, Y.-F. (2004) Phylogenetic relationships of the
747 subclass Peritrichia (Oligohymenophorea, Ciliophora) inferred from small subunit rRNA gene
748 sequences. *J. Eukaryot. Microbiol.* **51**: 180–186.

749 Nearing, J.T., Douglas, G.M., Comeau, A.M., and Langille, M.G.I. (2018) Denoising the Denoisers: an
750 independent evaluation of microbiome sequence error-correction approaches. *PeerJ* **6**:
751 e5364.

752 Nebel, M., Pfabel, C., Stock, A., Dunthorn, M., and Stoeck, T. (2011) Delimiting operational taxonomic
753 units for assessing ciliate environmental diversity using small-subunit rRNA gene sequences.
754 *Environ. Microbiol. Rep.* **3**: 154–158.

755 Needham, D.M., Sachdeva, R., and Fuhrman, J.A. (2017) Ecological dynamics and co-occurrence
756 among marine phytoplankton, bacteria and myoviruses shows microdiversity matters. *ISME J.*
757 **11**: 1614–1629.

758 Palmer, J.M., Jusino, M.A., Banik, M.T., and Lindner, D.L. (2018) Non-biological synthetic spike-in
759 controls and the AMPtk software pipeline improve mycobiome data. *PeerJ* **6**: e4925.

760 Prosser, J.I. (2010) Replicate or lie. *Environ. Microbiol.* **12**: 1806-1810.

761 Quince, C., Lanzén, A., Curtis, T.P., Davenport, R.J., Hall, N., Head, I.M., *et al.* (2009) Accurate
762 determination of microbial diversity from 454 pyrosequencing data. *Nat. Methods* **6**: 639–
763 641.

764 Quince, C., Lanzen, A., Davenport, R.J., and Turnbaugh, P.J. (2011) Removing noise from
765 pyrosequenced amplicons. *BMC Bioinformatics* **12**: 38.

766 R Core Team (2018) R: A language and environment for statistical computing. R Foundation for
767 Statistical Computing, Vienna, Austria.

768 Rand, W.M. (1971) Objective criteria for the evaluation of clustering methods. *J. Am. Stat. Assoc.* **66**:
769 846–850.

770 Reeder, J. and Knight, R. (2010) Rapidly denoising pyrosequencing amplicon reads by exploiting rank-
771 abundance distributions. *Nat. Methods* **7**: 668–669.

772 Rognes, T., Flouri, T., Nichols, B., Quince, C., and Mahé, F. (2016) VSEARCH: a versatile open source
773 tool for metagenomics. *PeerJ* **4**: e2584.

774 Schloss, P.D., Gevers, D., and Westcott, S.L. (2011) Reducing the effects of PCR amplification and
775 sequencing artifacts on 16S rRNA-based studies. *PLOS ONE* **6**: e27310.

776 Schloss, P.D. and Westcott, S.L. (2011) Assessing and improving methods used in operational
777 taxonomic unit-based approaches for 16S rRNA gene sequence analysis. *Appl. Environ.*
778 *Microbiol.* **77**: 3219–3226.

779 Schloss, P.D., Westcott, S.L., Ryabin, T., Hall, J.R., Hartmann, M., Hollister, E.B., *et al.* (2009)
780 Introducing mothur: Open-source, platform-independent, community-supported software
781 for describing and comparing microbial communities. *Appl. Environ. Microbiol.* **75**: 7537–
782 7541.

783 Stoeck, T., Bass, D., Nebel, M., Christen, R., Jones, M.D.M., Breiner, H.-W., and Richards, T.A. (2010)
784 Multiple marker parallel tag environmental DNA sequencing reveals a highly complex
785 eukaryotic community in marine anoxic water. *Mol. Ecol.* **19**: 21–31.

786 Stoeck, T., Breiner, H.-W., Filker, S., Ostermaier, V., Kammerlander, B., and Sonntag, B. (2014) A
787 morphogenetic survey on ciliate plankton from a mountain lake pinpoints the necessity of
788 lineage-specific barcode markers in microbial ecology. *Environ. Microbiol.* **16**: 430–444.

789 Sun, Y., Cai, Y., Huse, S.M., Knight, R., Farmerie, W.G., Wang, X., and Mai, V. (2012) A large-scale
790 benchmark study of existing algorithms for taxonomy-independent microbial community
791 analysis. *Brief Bioinformatics* **13**: 107–121.

792 Tikhonov, M., Leach, R.W., and Wingreen, N.S. (2015) Interpreting 16S metagenomic data without
793 clustering to achieve sub-OTU resolution. *ISME J.* **9**: 68–80.

794 Utter, D.R., Mark Welch, J.L., and Borisy, G.G. (2016) Individuality, stability, and variability of the
795 plaque microbiome. *Front. Microbiol.* **7**.

796 de Vargas, C., Audic, S., Henry, N., Decelle, J., Mahé, F., Logares, R., *et al.* (2015) Eukaryotic plankton
797 diversity in the sunlit ocean. *Science* **348**: 1261605.

798 Vďačný, P., Breiner, H.-W., Yashchenko, V., Dunthorn, M., Stoeck, T., and Foissner, W. (2014) The
799 chaos prevails: Molecular phylogeny of the Haptoria (Ciliophora, Litostomatea). *Protist* **165**:
800 93–111.

801 Wang, C., Zhang, T., Wang, Y., Katz, L.A., Gao, F., and Song, W. (2017) Disentangling sources of
802 variation in SSU rDNA sequences from single cell analyses of ciliates: impact of copy number
803 variation and experimental error. *Proc. R. Soc. B* **284**: 20170425

804 Worden, A.Z. (2006) Picoeukaryote diversity in coastal waters of the Pacific Ocean. *Aquat. Microb.*
805 *Ecol.* **43**: 165–175.

806 Xue, Z., Kable, M.E., and Marco, M.L. (2018) Impact of DNA sequencing and analysis methods on 16S
807 rRNA gene bacterial community analysis of dairy products. *mSphere* **3**: e00410-18.

808 Zhang, W.-J., Lin, Y.-S., Cao, W.-Q., and Yang, J. (2012) Genetic diversity and variance of *Stentor*
809 *coeruleus* (Ciliophora: Heterotrichea) inferred from inter-simple sequence repeat (ISSR)
810 fingerprinting. *J. Eukaryot. Microbiol.* **59**: 157–162.

811 Zoqratt, M.Z.H.M., Eng, W.W.H., Thai, B.T., Austin, C.M., and Gan, H.M. (2018) Microbiome analysis
812 of Pacific white shrimp gut and rearing water from Malaysia and Vietnam: implications for
813 aquaculture research and management. *PeerJ* **6**: e5826.

814

815

816 **Table 1: First-level sequence grouping results of each ciliate species-specific dataset.** Results are

817 shown for first-level grouping with both DADA2 and SWARM.

818

Ciliate species	DADA2-ASVs	SWARM-OTUs
<i>Chilodonella uncinata</i>	5	9
<i>Coleps hirtus hirtus</i>	3	5
<i>Deviata rositae</i>	2	5
<i>Dexiotricha sp.</i>	1	6
<i>Epistylis plicatilis</i>	5	9
<i>Euplotes sp.</i>	5	6
<i>Folliculina sp.</i>	7	18
<i>Fuscheria uluruensis</i>	2	8
<i>Gastrostyla steinii</i>	2	4
<i>Hypotrichida sp.</i>	3	5
<i>Lagynophrya acuminata</i>	2	11
<i>Metanophrys sp.</i>	1	4
<i>Oxytricha granulifera</i>	3	7
<i>Paramecium bursaria</i>	7	10
<i>Paramecium tetraurelia</i>	3	9
<i>Pelagodileptus trachelioides</i>	3	7
<i>Platynematum salinarum</i>	3	15
<i>Schmidingerothrix salinarum</i>	7	11
<i>Spathidium ascendens</i>	4	19
<i>Spathidium foissneri</i>	3	4
<i>Spirostomum ambiguum</i>	2	3
<i>Stentor coeruleus</i>	11	12
<i>Tetrahymena sp.</i>	2	2
<i>Tokophrya infusionum</i>	2	2
<i>Trithigmostoma cucullulus</i>	1	1
<i>Uroleptus willii</i>	5	16
<i>Urospinula succisa</i>	4	8
<i>Usconophrys sp.</i>	2	7
<i>Vorticella sp.</i>	1	6

819

820

821 **Figure legends**

822 **Fig. 1: Schematic workflow of the two-level sequence grouping approach.** The workflow is shown
823 for four hypothetical species and includes sample processing, first-level sequence grouping in either
824 DADA2 or SWARM and second-level sequence grouping with SSNs. A unique color was chosen for
825 each of the species to highlight the separate handling of each dataset. Second-level sequence
826 grouping also displays in how far network sequence clusters (NSCs) can be used to distinguish
827 between closed, open and hybrid NSCs.

828

829 **Fig.2: Results of second-level sequence grouping at different sequence similarity binning levels.** The
830 bars show the gradual increase of NSCs with increasing similarity binning level. Please note that the
831 x-axis contains 0.5% steps between 96% and 98% similarity. DADA2-SSN results are displayed in blue,
832 SWARM-SSN results in orange. The dashed line indicates the ideally expected diversity richness of 29
833 NSCs. Raw values of the bars can also be found in Table S1.

834

835 **Fig.3: Qualitative evaluation of second-level sequence grouping with the DADA2-SSN approach (A)**
836 **and the SWARM-SSN approach (B).** The evaluation focused on the questions how many closed NSCs
837 (green), open NSCs (yellow) and hybrid NSCs (red) were found. Ideally, a binning level should be
838 chosen at which both the number of closed NSCs is maximal and the number of hybrid NSCs is
839 minimal. The bars show the gradual changes of NSC types with increasing similarity binning levels.

840

841 **Fig.4: Rand index (RI) and adjusted Rand index (ARI) results for second-level sequence grouping.**
842 The graphs show the values for comparing the observed cluster distributions of the DADA2-SSN (in
843 blue) and SWARM-SSN (in orange) approaches against a perfect cluster distribution of the same data
844 in which one cluster existed for each of the 29 species under study. The values were calculated for all

845 tested sequence similarity binning levels. RI results are displayed as a straight line, ARI results as a
846 dashed line.

847

848 **Supplementary files**

849 **Table S1: NSC numbers at different sequence similarity binning levels.** The values document the
850 gradual increase of ASV-NSCs and OTU-NSCs with increasing similarity binning level, which is also
851 displayed in Figure 2.

852

853 **Table S2: Difference between non-abundance filtered and abundance-filtered data before second-**
854 **level sequence grouping.** For testing the effect of an abundance filter, all ASVs or SWARM-OTUs
855 which amounted to less than 0.01% of the sequences of a species were filtered from the dataset.
856 That is, the numbers shown here are also equivalent to the number of NSCs at the 100% binning level
857 for second-level sequence grouping with and without abundance filter.

858

859 **Table S3: Numbers of NSC types at each similarity binning level.** The amounts of closed, open and
860 hybrid NSCs shown in this table reflect the bars shown in Figures 3A and B.

861

862 **Table S4: Sample information for each of the 29 ciliate species.** In addition to species name and
863 taxonomy, the table also shows the source of the sample, to which GenBank accession number the
864 best hit refers and with which Illumina platform it was sequenced.

865

866 **Table S5: NSC numbers at different sequence similarity binning levels when species-specific**
867 **datasets were merged before first-level sequence grouping.** The values document the gradual
868 increase of ASV-NSCs and OTU-NSCs with increasing similarity binning level. NSC numbers at the
869 similarity binning level of 100% are equivalent to the outcome of first-level sequence grouping after
870 merging species-specific ciliate datasets. The third and fourth columns show NSC numbers when

871 species-specific datasets were merged before first-level sequence grouping and an abundance-filter
872 was applied before second-level sequence grouping.

873

874 **Table S6: Library information on all used datasets.**

875

876 **Fig. S1: Results of second-level sequence grouping at different sequence similarity binning levels**
877 **after abundance-filtering.** The bars show the gradual increase of NSCs with increasing similarity
878 binning level, similar to the results shown without abundance filter in Figure 2. Please note that the
879 x-axis contains 0.5% steps between 96% and 98% similarity. DADA2-SSN results are displayed in blue,
880 SWARM-SSN results in orange. The dashed line indicates the ideally expected diversity richness of 29
881 NSCs.

882

883 **Fig. S2: Qualitative evaluation of second-level sequence grouping with the DADA2-SSN approach**
884 **(A) and the SWARM-SSN approach (B) after abundance-filtering.** Similar to Figures 3A and B, the
885 evaluation focused on the questions how many closed NSCs (green), open NSCs (yellow) and hybrid
886 NSCs (red) were found when working with abundance-filtered data. Ideally, a binning level should be
887 chosen at which both the number of closed NSCs is maximal and the number of hybrid NSCs is
888 minimal. The bars show the gradual changes of NSC types with increasing similarity binning levels.

889

890 **Fig. S3: Rand index (RI) and adjusted Rand index (ARI) results for second-level sequence grouping**
891 **after abundance-filtering.** The graphs show the values for comparing the observed cluster
892 distributions of the DADA2-SSN (in blue) and SWARM-SSN (in orange) approaches after an initial
893 abundance filtering step of the data, against a perfect cluster distribution of the same data in which
894 one cluster existed for each of the 29 species under study. The values were calculated for all tested

895 sequence similarity binning levels. RI results are displayed as a straight line, ARI results as a dashed
896 line.

897

898 **Fig. S4: Qualitative evaluation of second-level sequence grouping with the DADA2-SSN approach**

899 **(A) and the SWARM-SSN approach (B) when species-specific datasets were merged before first-**

900 **level sequence grouping.** Similar to figures 3A and B, as well as figures S2A and S2B, the evaluation

901 focused on the questions how many closed NSCs (green), open NSCs (yellow) and hybrid NSCs (red)

902 were found. Ideally, a binning level should be chosen at which both the number of closed NSCs is

903 maximal and the number of hybrid NSCs is minimal. The bars show the gradual changes of NSC types

904 with increasing similarity binning levels.

905

906 **Fig. S5: Rand index (RI) and adjusted Rand index (ARI) results for second-level sequence grouping,**

907 **when species-specific datasets had been merged before first-level sequence grouping.** The graphs

908 show the values for comparing the observed cluster distributions of the DADA2-SSN (in blue) and

909 SWARM-SSN (in orange) approaches, against a perfect cluster distribution of the same data in which

910 one cluster existed for each of the 29 species under study. The values were calculated for all tested

911 sequence similarity binning levels. RI results are displayed as a straight line, ARI results as a dashed

912 line.

913

914 **File S1: Supplementary codes in HTML format.**