

**A peer-reviewed version of this preprint was published in PeerJ on 9 April 2018.**

[View the peer-reviewed version](https://peerj.com/articles/4644) (peerj.com/articles/4644), which is the preferred citable publication unless you specifically need to cite this preprint.

Elbrecht V, Vamos EE, Steinke D, Leese F. 2018. Estimating intraspecific genetic diversity from community DNA metabarcoding data. PeerJ 6:e4644 <https://doi.org/10.7717/peerj.4644>

1 **Title: Assessing intraspecific genetic diversity from community DNA metabarcoding data**

2  
3 **Running Title (45 char max):** Extracting haplotypes from metabarcoding data

4 **Authors:** Vasco Elbrecht<sup>1,2\*</sup>, Ecaterina Edith Vamos<sup>1</sup>, Dirk Steinke<sup>2</sup>, Florian Leese<sup>1,3</sup>

5  
6 Affiliations:

7 1) Aquatic Ecosystem Research, Faculty of Biology, University of Duisburg-Essen, Universitätsstraße 5, 45141 Essen,  
8 Germany

9 2) Centre for Biodiversity Genomics, University of Guelph, 50 Stone Road East, Guelph, Ontario, N1G 2W1, Canada

10 3) Centre for Water and Environmental Research (ZWU) Essen, University of Duisburg-Essen, Universitätsstraße 2, 45141  
11 Essen, Germany

12 **\*\*Corresponding author:** Vasco Elbrecht (vasco.elbrecht@uni-due.de),

13  
14 **Abstract:**

15 **Background.** DNA metabarcoding is used to generate species composition data for entire communities. However,  
16 sequencing errors in high throughput sequencing instruments are fairly common, usually requiring reads to be clustered into  
17 operational taxonomic units (OTU), losing information on intraspecific diversity in the process.

18 **Methods.** This study combines sequence denoising strategies, normally applied in microbial research, with additional  
19 abundance based filtering to extract haplotypes from freshwater macroinvertebrate metabarcoding data sets. This novel  
20 approach is implemented in the R package "JAMP" and can be applied to Cytochrome c oxidase subunit I (COI) amplicon  
21 datasets. We tested our haplotyping method by sequencing i) a single-species mock community composed of 31 individuals  
22 with different haplotypes spanning three orders of magnitude in biomass and ii) 18 monitoring samples each amplified with  
23 four different primer sets and two PCR replicates.

24 **Results.** We detected all 15 haplotypes of the single specimens in the mock community with relaxed filtering and denoising  
25 settings. However, up to 480 additional unexpected haplotypes remained in both replicates. Rigorous filtering removes  
26 most unexpected haplotypes, but also can discard expected haplotypes mainly from the small specimens. In the monitoring  
27 samples, the different primer sets detected 177 - 200 OTUs, each containing an average of 2.40 to 3.30 haplotypes per OTU.  
28 Population structures were consistent between replicates, and similar between primer pairs, depending on the primer length.  
29 A closer look at abundant taxa in the data set revealed various population genetic patterns, e.g. *Taeniopteryx nebulosa* and

30 *Hydropsyche pellucidula* with a difference in north-south haplotype distribution, while *Oulimnius tuberculatus* and *Asellus*  
31 *aquaticus* display no clear population pattern but differ in genetic diversity.

32 **Discussion.** We developed a strategy to infer intraspecific genetic diversity from bulk invertebrate samples using  
33 metabarcoding data. It needs to be stressed that at this point metabarcoding-informed haplotyping is not capable to capture  
34 the full diversity present in bulk samples, due to variation in specimen size, primer bias and loss of sequence variants with  
35 low abundance. Nevertheless, for a high number of species intraspecific diversity is recovered, identifying potentially  
36 isolated populations and potential taxa for further more detailed phylogeographic investigation. While we are currently  
37 lacking large-scale metabarcoding data sets to fully take advantage our new approach, metabarcoding-informed haplotyping  
38 holds great promise for biomonitoring efforts that not only seek information about biological diversity but also underlying  
39 genetic diversity.

40  
41 **Keywords:** metabarcoding, high-throughput sequencing, haplotyping, population genetics, ecosystem assessment

## 42 43 **Introduction**

44 High-throughput analysis of DNA barcodes retrieved from environmental samples, i.e. DNA metabarcoding, allows for the  
45 rapid and standardized assessment of community composition without the need for morpho-taxonomy (Taberlet et al.,  
46 2012a; Creer et al., 2016). This new surge of data enables biodiversity surveys at speeds and scales that were previously  
47 inconceivable in ecological and evolutionary studies. While the approach has major strengths and is generally regarded as a  
48 game changer for ecological research (Creer et al., 2016), it still has limitations such as the fact that sequences are typically  
49 clustered into operational taxonomic units (OTUs, Fig. S1) thereby ignoring any intraspecific variation (Callahan,  
50 McMurdie & Holmes, 2017). However, clustering is a crucial step to reduce the influence of PCR and sequencing errors  
51 that can otherwise generate false sequence variation (Edgar, 2013). This inability to detect intraspecific variation hampers  
52 our ability to detect impacts at the population level. Simultaneous assessment of inter- and intraspecific diversity, however,  
53 would be a milestone forward in ecological research and management because haplotype data are direct proxies to register  
54 spatio-temporal dynamics of populations and both parameters can differ substantially (Taberlet et al., 2012b). Especially  
55 assessing fragmentation (e.g. Weiss & Leese 2016) or changes in population size in response to environmental impacts are  
56 key area in basic and applied ecological research (e.g. Sutherland et al. 2012). Also for management this parameter is  
57 important because because genetic variation is typically lost long before complete species or OTUs (Bálint et al., 2011).  
58 Unfortunately, methods to extract haplotype information from metabarcoding data sets are generally not widely available

59 and thus most studies are based on single-specimen analyses. Some of those are based on denoising algorithms capable of  
60 distinguishing between true haplotypes and sequencing noise (e.g. Tikhonov, Leach & Wingreen, 2015; Eren et al., 2015;  
61 Edgar, 2016; Callahan et al., 2016; Amir et al., 2017) and have been tested for microbial samples (e.g. Eren et al., 2015;  
62 Callahan et al., 2016; Needham, Sachdeva & Fuhrman, 2017). Wares & Pappalardo (2016) suggested that haplotype  
63 information in metazoan datasets can be used to, for instance, improve taxa abundance estimates, which was successfully  
64 demonstrated with bat diet samples (Corse et al., 2017). Recent studies were also able to infer haplotypes with  
65 metabarcoding for single specimens (Shokralla et al., 2014), arthropod bulk samples (Elbrecht & Leese, 2015; Pedro et al.,  
66 2017) and environmental water samples (Sigsgaard et al., 2016), all highlighting the possibility to extract sequence variant  
67 information within OTUs when targeting metazoan taxa.

68 We here further explore bioinformatics strategies in order to unlock the potential of metabarcoding based haplotyping of  
69 entire and complex metazoan communities. Therefore, we combined stringent quality filtering of reads with the recently  
70 developed *unoise3* denoising strategy (Edgar, 2016) and calibrated this approach using a previously characterized single-  
71 species mock sample composed of specimens with known haplotypes (Elbrecht & Leese, 2015; Vamos, Elbrecht & Leese,  
72 2017). Subsequently, we multi-species metabarcoding data collected from 18 sample sites as part of a governmental  
73 freshwater macroinvertebrates biomonitoring program (Elbrecht et al., 2017). These were denoised with the developed  
74 strategy and we tested the potential to detect intraspecific variation across the broad geographic gradient across multiple  
75 taxa.

76

## 77 **Materials & Methods**

78 We tested our haplotyping strategy on two available DNA metabarcoding datasets, 1) a single-species mock sample  
79 containing 31 specimens with known haplotypes (Vamos, Elbrecht & Leese, 2017) and 2) a multi-species macroinvertebrate  
80 community dataset from the Finnish governmental stream monitoring program (Elbrecht et al., 2017). The samples were  
81 sequenced for a region nested within the classical Folmer COI region (Folmer et al., 1994) with two replicates each. Hereby,  
82 the single-species sample (1) was sequenced using a short primer set amplifying 178 bp, while the multi-species monitoring  
83 samples were amplified using four different primer sets targeting a region of up to 421 bp (Elbrecht & Leese, 2017). Paired-  
84 end sequencing (250 bp) was performed on Illumina MiSeq and HiSeq systems with high sequencing depth (on average  
85 1.53 million reads per sample, SD = 0.29).

86 To extract individual haplotypes from the metabarcoding datasets, we used strict quality filtering followed by denoising  
87 (*unoise3* (Edgar, 2016), with additional threshold-based filtering steps (see Fig. 1B). The full metabarcoding and

88 haplotyping pipelines are available as R package (<https://github.com/VascoElbrecht/JAMP>), which requires Usearch  
89 v10.0.240 (Edgar, 2013). All used pipeline commands are also available as supporting information (Fig. S2, Scripts S1,  
90 JAMP v0.28). In short, pre-processing of reads involved sample demultiplexing, paired-end merging, primer trimming,  
91 generation of reverse complements where needed (to align all reads in the forward direction), max ee filtering = 0.5 (Edgar  
92 & Flyvbjerg, 2015), only keeping reads of exact length targeted by the respective primer set, subsampling to 1 and 0.4  
93 million reads, respectively, to generate the same sequencing depth for the single species and monitoring bulk samples. To  
94 further reduce the amount of sequences affected by sequencing errors we applied read denoising with unoise3 as  
95 implemented in Usearch (Edgar, 2016) to all samples of a dataset using only reads with  $\geq 10$  abundance in each sample  
96 after dereplication. Different expected error cutoffs and alpha values were tested, with ee = 0.5 and alpha = 5 being used for  
97 the final analysis of the 18 bulk samples.

98 For the single-species mock sample, the denoised and quality filtered reads (prior to denoising) were mapped against the  
99 expected 15 haplotype sequences using Vsearch (v2.4.3) (Rognes et al., 2016). The unoise3 implementation into the JAMP  
100 package adds additional threshold-based filtering after the denoising step, which we used for the Finnish multi-species  
101 monitoring samples in order to discard low abundant haplotypes and OTUs "Denoise(..., minhaplosize = 0.01, OTUmin =  
102 0.1)". Additionally, within each OTU and sample site, only haplotypes with at least 5% abundance per sample were  
103 considered for generating haplotype maps and networks, in order to exclude low abundance OTUs which can be difficult to  
104 separate from PCR artifacts and sequencing errors.

105

## 106 Results

107 Our approach starts with denoising of quality filtered reads using unoise3 (Edgar, 2016) followed by an additional  
108 threshold-based filtering step which includes OTU clustering of denoised reads (Edgar, 2013) and the removal of low  
109 abundant OTUs / haplotypes (see Fig. 1B). We validated this approach by using a single species mock community of known  
110 haplotype composition (Elbrecht & Leese, 2015), in which we found 943 unexpected haplotypes above 0.003% abundance  
111 with no expected error filtering applied (Fig. 1A). Filtering the raw sequence data with different quality thresholds (max ee,  
112 Edgar & Flyvbjerg, 2015) reduced the number of unexpected haplotypes by only up to 10.22% (Fig S3). The consistency  
113 between the two independent sequencing replicates indicates that a major fraction of the detected haplotypes represent in  
114 fact real biological signal (e.g. somatic mutations, numts or heteroplasmy, Bensasson et al., 2001; Shokralla et al., 2014),  
115 which is difficult to differentiate from PCR and sequencing errors. Even after using different alpha values for the unoise3  
116 algorithm some unexpected sequence variants remained (Fig S4). An error filtering of max ee = 0.5 in combination with an

117 alpha of 5 was chosen for subsequent analysis (Fig. 1C), as it offers the best trade-off between expected and unexpected  
118 haplotypes (9 of 15 expected, 6 unexpected with low abundance), while retaining 67.08% (SD = 17.69%) of the original  
119 sequence data after quality filtering and before denoising.

120 For the denoising of our multi-species environmental biomonitoring samples, additional and more conservative  
121 filtering steps were introduced to ensure only true sequence variants are included in the analysis (discarding low abundant  
122 OTUs and haplotypes below 0.1% and 0.01%, as well as haplotypes below 5% read abundance within each OTU of the  
123 respective sample, Fig. 1C green line). Denoising of metabarcoding data from 18 macroinvertebrate samples of the Finnish  
124 routine stream monitoring, recovered 177 - 200 OTUs containing 534 - 646 haplotypes (on average 2.40 - 3.30 haplotypes  
125 per OTU, SD = 2.13 - 3.26) for the different primer pairs (Table S1). Most OTUs were only present in a few sample  
126 locations, allowing for only limited population genetic analysis (Fig. S5, see also Fig. S7 in Elbrecht et al., 2017). Fig. 2  
127 depicts some examples of haplotype diversity and geographic distribution for more common and widely distributed taxa in  
128 this study. For *Taeniopteryx nebulosa* (Plecoptera) and *Hydropsyche pellucidula* (Trichoptera) we found distinct patterns of  
129 latitudinal variation in haplotype composition (Fig. 2A, B), while *Oulimnius tuberculatus* (Coleoptera) showed low genetic  
130 variation across all primer combinations (Fig. 2C, Fig. S3C). *Asellus aquaticus* (Isopoda) on the other hand showed very  
131 high genetic diversity for endemic haplotypes (Fig. 2D).  
132 Extracted haplotype patterns between replicates were highly reproducible ( $R^2 = 0.751$ , SD = 0.242), while at the same time  
133 recovering more sequence variants with longer amplicons (Fig. S6). Taxon occurrence for the four taxa analysed in detail  
134 matched morphology based identifications (Elbrecht et al., 2017) in most cases (only four false positive detections, Fig. 2).  
135 The few inconsistencies between replicates in haplotypes and taxa occurrence are mostly affecting low abundance reads. In  
136 the sequence alignments, all four primer sets shared most of the variable positions (Fig. S6).

137

## 138 Discussion

139 In this case study, we developed and demonstrated a bioinformatic strategy to process metabarcoding data first using a  
140 controlled single-species approach, in order to extract intraspecific genetic diversity information from complex multi-  
141 species metazoan environmental samples. While our multi-species dataset was limited to only 18 sampling sites, and many  
142 taxa were not widely distributed (Elbrecht et al., 2017), we could still infer potential population genetic patterns for some of  
143 the abundant and more widespread taxa. Where available, observed population genetic patterns were also consistent with  
144 previous studies, e.g. earlier work reported high genetic diversity for *A. aquaticus* (Sworobowicz et al., 2015). Other

145 published work, e.g. on *H. pellucidula* (Múrria et al., 2010) and *O. tuberculatus* (Čiampor & Kodada, 2010) was too limited  
146 in sampling size and region for proper comparison.

147 Deriving haplotypes from metabarcoding data does not require specialized field or laboratory protocols, as existing data is  
148 analyzed. And while our dataset is very limited with just 18 sample sites, there are efforts underway to implement DNA  
149 metabarcoding based monitoring of stream water quality in Europe, potentially generating HTS data for thousands of  
150 sample sites every year (Leese et al., 2016, Leese et al. in press). Such haplotype data, even though limited in resolution and  
151 based only on a single gene marker, could be used to formulate hypotheses about taxa dispersal at an unprecedented scale  
152 (Hughes, Schmidt & FINN, 2009), which would be highly beneficial for the renaturation and management of aquatic  
153 ecosystems.

154 While the detection of haplotypes from bulk samples was demonstrated in this and other studies (Sigsgaard et al., 2016;  
155 Corse et al., 2017; Pedro et al., 2017), the limitations of metabarcoding based haplotyping remain relatively unexplored.  
156 Metabarcoding data sets can be affected by primer bias (Elbrecht & Leese, 2015), tag switching (Esling, Lejzerowicz &  
157 Pawlowski, 2015; Schnell, Bohmann & Gilbert, 2015), as well as PCR and sequencing errors (Nakamura et al., 2011;  
158 Tremblay et al., 2015). Such issues can lead to artificial haplotypes, which are usually sufficiently different to distinguish  
159 them from actual haplotypes in the samples, especially if they are less abundant and thus likely influenced by stochastic  
160 effects (Leray & Knowlton, 2017). We applied very strict quality filtering in our pipeline, and cautiously discarded all  
161 haplotypes below 5% abundance within an OTU. This is necessary, as low abundant OTUs can not be separated from  
162 sequencing errors, somatic mutations (Shokralla et al., 2014) and other noise in the data, as we have shown for the single  
163 species mock samples. Strict filtering will remove rare and low abundant OTUs, but it is necessary to reduce the amount of  
164 false positive artificial OTUs that result from the currently rather high error rates of HTS instruments. Even with such strict  
165 filtering settings, we can not be fully confident that all false haplotypes were excluded e.g. as the result of undetected  
166 chimeric sequences (Edgar et al., 2011) or systematic sequencing errors (Nakamura et al., 2011; Schirmer et al., 2015;  
167 Schirmer, 2016) that likely persist across replicates. Approaches relying on the comparison of replicate samples could be an  
168 appropriate strategy in particular when working with unicellular organisms (Lange et al., 2015). However, as for our  
169 metazoan communities many variants occur across both replicates (Fig. 1). Macroinvertebrate communities can vary  
170 considerably in biomass, which means rare and small specimens will be underrepresented when extracting DNA from bulk  
171 samples (Elbrecht, Peinert & Leese, 2017). Thus, taxa in the sample are sequenced at different sequencing depth, which  
172 likely has an influence on the amount of false haplotypes detected within each OTU. Additionally, differences in specimen  
173 biomass can skew the detection of haplotypes, as only those of large specimens will be retained in bioinformatics analysis  
174 (haplotypes of small specimens are likely below 5% abundance). Such uncertainties need to be considered when doing

175 population genetic analysis, which is usually done at specimen level, with the exact number of specimens and haplotypes  
176 known for each sampling site. It has to be emphasized that at this point metabarcoding based haplotyping only provides  
177 very limited information of genetic diversity and phylogeography of a given taxon. However, interesting patterns emerging  
178 from such studies can be subsequently explored by collecting taxa of interest and using standard population genetic markers  
179 with a higher resolution (e.g. microsatellites, ddRAD (Peterson et al., 2012)). Our study demonstrates the feasibility and  
180 potential of metabarcoding data for the investigation of population genetic patterns of entire complex environmental  
181 communities. The shortcomings and the level of resolution of this novel approach need to be carefully tested. Additionally,  
182 more bioinformatics approaches suited for the analysis of metazoan bulk samples need to be developed, especially with  
183 respect to variation in specimen biomass (Elbrecht, Peinert & Leese, 2017). Furthermore, most software currently used in  
184 this field was developed for microbial samples and should therefore be further tested and benchmarked for its feasibility in  
185 studies involving eukaryotes. Despite the clear limitations of this haplotyping approach, we are confident that it will be  
186 useful in future large-scale studies of genetic diversity. While metabarcoding studies will remain affected by sequencing  
187 errors (potentially leading to false haplotypes), we expect that most of these issues can be mitigated by increasing the  
188 number of sampling sites to several hundred or even thousands. For large-scale efforts such as routine monitoring using  
189 metabarcoding (Baird & Hajibabaei, 2012; Gibson et al., 2015; Elbrecht et al., 2017), this might soon become a feasible  
190 option if not standard.

191

## 192 **Conclusions**

193 Our study demonstrates that haplotypes can be extracted from complex metazoan metabarcoding datasets. This proof of  
194 concept work already shows emerging population genetic patterns for a few species, but more large-scale validation studies  
195 are needed to explore the limitations and the potential of metabarcoding based haplotyping. While some shortcomings such  
196 as occasional false positive detections and loss of rare and small taxa are difficult to overcome per sample for such complex  
197 communities, these can be partly offset by studying comparative patterns of intraspecific variation across many taxa and  
198 sites. As metabarcoding becomes more accessible and larger DNA-based biodiversity assessment and monitoring initiatives  
199 emerge, sampling and extracting haplotypes from hundreds of sites might become a feasible path of future research.

200

201

202

203 **Data availability.** Unprocessed raw sequence data are available from previous studies on the NCBI SRA archive. Single



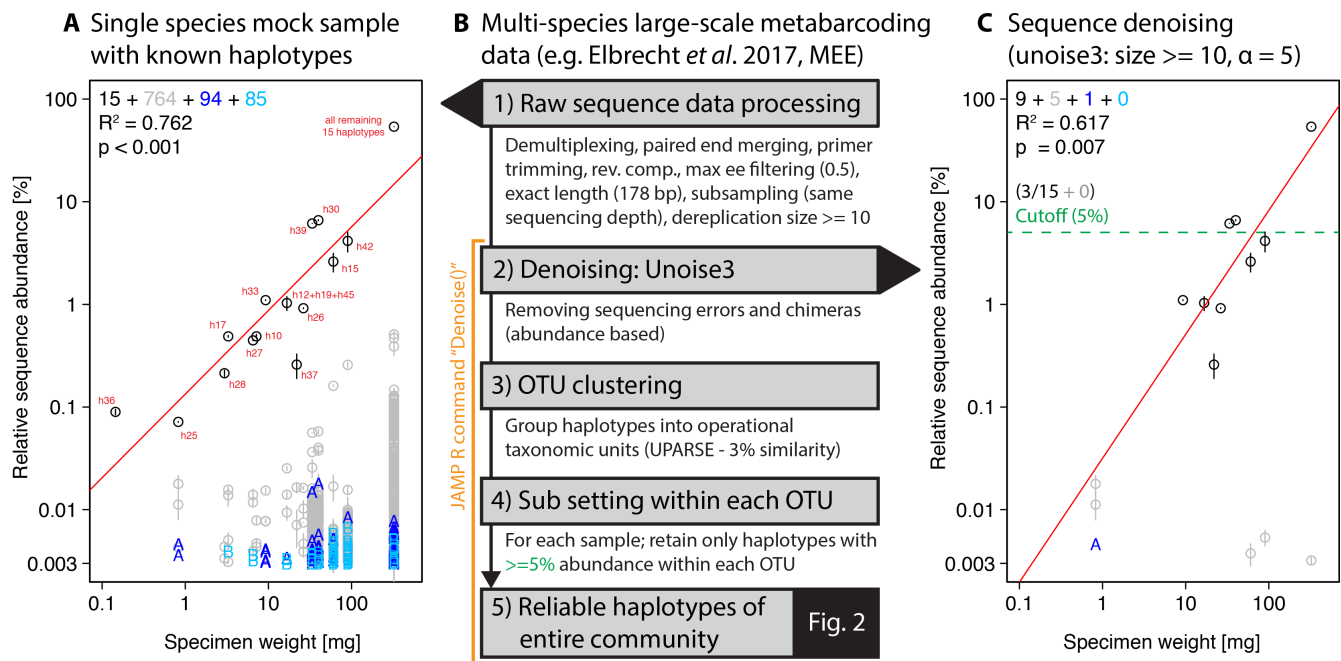
204 species mock sample: SRR5295658 and SRR5295659 (Vamos, Elbrecht & Leese, 2017), monitoring samples: SRR4112287  
205 (Elbrecht et al., 2017). The JAMP R package is available on GitHub ([github.com/VascoElbrecht/JAMP](https://github.com/VascoElbrecht/JAMP)) with the used R  
206 scripts (Script S1) and full haplotype tables (Table S1) available as supporting information.

207

208

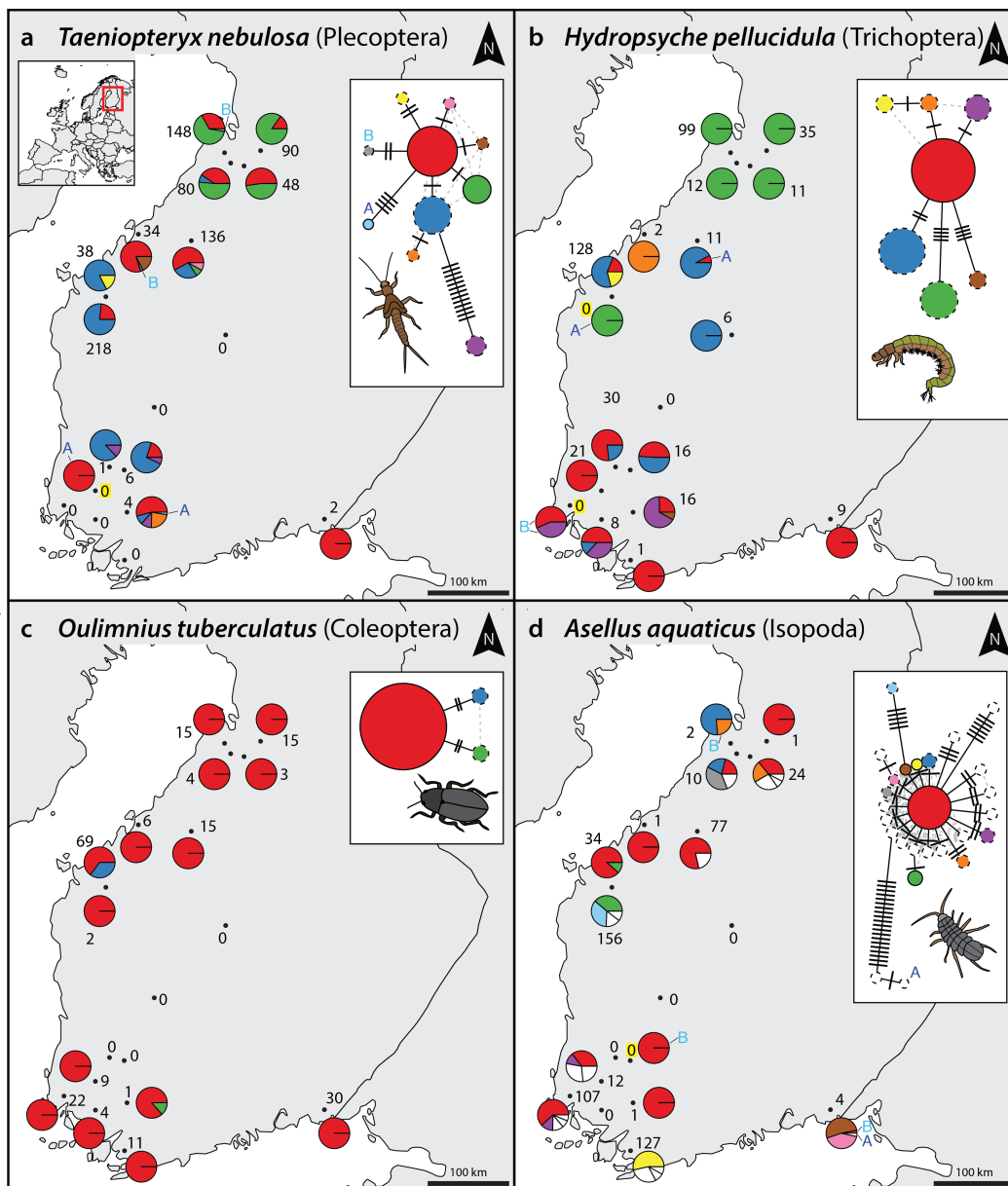
209

210



211

212 **Figure 1:** Overview of DNA metabarcoding data of a single-species mock sample containing specimens with 15 distinct  
 213 haplotypes (black circles). Detected haplotypes (unexpected ones shown in grey and blue) plotted against specimen biomass  
 214 for the processed data (A) and followed by read denoising using unoise3 (C). Denoising was applied to both replicates  
 215 individually, with a circle if the read was detected in both samples (error bar = SD) and A or B if the read was found in only  
 216 one replicate. For processing of large-scale samples (B, Fig. 2), all samples were pooled and jointly denoised, followed by  
 217 OTU clustering and read mapping then followed by discarding of haplotypes below a 5% threshold within each sample.



218

219 **Figure 2:** Haplotype maps and networks extracted from multi-species community metabarcoding datasets amplified with  
 220 the BF2+BR2 primer set for four abundant macroinvertebrate taxa (A = *Taeniopteryx nebulosa*, B = *Hydropsyche*  
 221 *pellucidula*, C = *Oulimnius tuberculatus*, D = *Asellus aquaticus*). Numbers next to each sampling site indicate sample size  
 222 of the respective taxa based on morphological identification in a sample (Elbrecht et al., 2017). Conflicts between DNA and  
 223 morphology based-detections are highlighted in yellow. Haplotype frequency composition per site is indicated by pie charts.  
 224 For *A. aquaticus* only the 10 most common haplotypes are visualised with different colours (remaining ones in white). Each  
 225 crossline in a network represents one base pair difference between the respective haplotypes. Dashed lines around a circle  
 226 indicate novel haplotypes that were not available in the BOLD reference database. An A or B next to a haplotype in the map  
 227 or network indicates the presence of this haplotype in only in one replicate.

228

## 229 Acknowledgements

230 We would like to thank members of the leeselab for helpful discussions. This study is part of the European Cooperation in  
231 Science and Technology (COST) Action DNAqua-Net (CA15219). D.S. was supported by the Canada First Research  
232 Excellence Fund for the Food from Thought initiative. E.E.V. was supported by a grant of the Bodnarescu Foundation.

233

## 234 Author contributions

235 V.E. developed the haplotyping concept, with contributions from E.E.V. and F.L., V.E. developed the bioinformatics and  
236 analysed the data, V.E., E.E.V., D.S., and F.L. wrote and revised the paper.

237

238

239

## 240 Supporting information

241 **Figure S1:** Schematic overview of errors affecting metabarcoding data and clustering / denoising strategies to reduce them.

242 **Figure S2:** Overview of the haplotyping strategy used here and their implementation in the JAMP R package.

243 **Figure S3:** Effect of different quality filtering (max ee) on reads of the single species mock sample.

244 **Figure S4:** Effect of different alpha values in read denoising of the single-species mock sample.

245 **Figure S5:** Bar plots of haplotype distribution within each OTU.

246 **Figure S6:** Detailed plots of four example taxa from the denoised multi-species monitoring samples, showing haplotype  
247 maps & networks, similarity between replicates and sequence alignment for all BF/BR primer sets.

248 **Table S1:** Finland haplotype table (for all four different primer combinations).

249 **Scripts S1:** Metabarcoding and denoising pipeline, and additional scripts used to produce the figures.

250

251

252

253 Amir A, McDonald D, Navas-Molina JA, Kopylova E, Morton JT, Zech Xu Z, Kightley EP, Thompson LR, Hyde ER,  
254 Gonzalez A, Knight R 2017. Deblur Rapidly Resolves Single-Nucleotide Community Sequence Patterns. *mSystems*  
255 2:e00191–16–7. DOI: 10.1128/mSystems.00191-16.

256 Baird DJ, Hajibabaei M 2012. Biomonitoring 2.0: a new paradigm in ecosystem assessment made possible by next-  
257 generation DNA sequencing. 21:2039–2044.

258 Bálint M, Domisch S, Engelhardt CHM, Haase P, Lehrian S, Sauer J, Theissinger K, Pauls SU, Nowak C 2011. Cryptic  
259 biodiversity loss linked to global climate change. *Nature Climate Change* 1:1–6. DOI: 10.1038/nclimate1191.

260 Bensasson D, Zhang DX, Hartl DL, Hewitt GM 2001. Mitochondrial pseudogenes: evolution's misplaced witnesses. *Trends*  
261 *in ecology & evolution (Personal edition)* 16:314–321.

262 Callahan BJ, McMurdie PJ, Holmes SP 2017. Exact sequence variants should replace operational taxonomic units in  
263 marker-gene data analysis. :1–5. DOI: 10.1038/ismej.2017.119.

264 Callahan BJ, McMurdie PJ, Rosen MJ, Han AW, Johnson AJA, Holmes SP 2016. DADA2: High-resolution sample

- 265 inference from Illumina amplicon data. *Nature Methods* 13:581–583. DOI: 10.1038/nmeth.3869.
- 266 Corse E, MEGLÉCZ E, Archambaud G, Ardisson M, MARTIN J-F, Tougard C, Chappaz R, DUBUT V 2017. A from-  
267 benchtop-to-desktop workflow for validating HTS data and for taxonomic identification in diet metabarcoding studies.  
268 *Molecular ecology resources* 17:e146–e159. DOI: 10.1111/1755-0998.12703.
- 269 Creer S, Deiner K, Frey S, Porazinska D, Taberlet P, Thomas WK, Potter C, Bik HM 2016. The ecologist's field guide to  
270 sequence-based identification of biodiversity. *Methods in Ecology and Evolution* 7:1008–1018. DOI: 10.1111/2041-  
271 210X.12574.
- 272 Čiampor F Jr, Kodada J 2010. Taxonomy of the *Oulimnius tuberculatus* species group (Coleoptera: Elmidae) based on  
273 molecular and morphological data. *Zootaxa*.
- 274 Edgar RC 2013. UPARSE: highly accurate OTU sequences from microbial amplicon reads. *Nature Methods* 10:996–998.  
275 DOI: 10.1038/nmeth.2604.
- 276 Edgar RC 2016. UNOISE2: improved error-correction for Illumina 16S and ITS amplicon sequencing. *bioRxiv*. DOI:  
277 10.1101/081257.
- 278 Edgar RC, Flyvbjerg H 2015. Error filtering, pair assembly and error correction for next-generation sequencing reads.  
279 *Bioinformatics* 31:3476–3482. DOI: 10.1093/bioinformatics/btv401.
- 280 Edgar RC, Haas BJ, Clemente JC, Quince C, Knight R 2011. UCHIME improves sensitivity and speed of chimera detection.  
281 *Bioinformatics* 27:2194–2200. DOI: 10.1093/bioinformatics/btr381.
- 282 Elbrecht V, Leese F 2015. Can DNA-Based Ecosystem Assessments Quantify Species Abundance? Testing Primer Bias and  
283 Biomass—Sequence Relationships with an Innovative Metabarcoding Protocol. *PloS one* 10:e0130324–16. DOI:  
284 10.1371/journal.pone.0130324.
- 285 Elbrecht V, Leese F 2017. Validation and development of freshwater invertebrate metabarcoding COI primers for  
286 Environmental Impact Assessment. *Frontiers in Freshwater Science*. DOI: 10.3389/fenvs.2017.00011.
- 287 Elbrecht V, Peinert B, Leese F 2017. Sorting things out: Assessing effects of unequal specimen biomass on DNA  
288 metabarcoding. *Ecology and Evolution* 7:6918–6926. DOI: 10.1002/ece3.3192.
- 289 Elbrecht V, Vamos E, Meissner K, Aroviita J, Leese F 2017. Assessing strengths and weaknesses of DNA metabarcoding  
290 based macroinvertebrate identification for routine stream monitoring. *Methods in Ecology and Evolution*:1–21. DOI:  
291 10.7287/peerj.preprints.2759v2.
- 292 Eren AM, Morrison HG, Lescault PJ, Reveillaud J, Vineis JH, Sogin ML 2015. Minimum entropy decomposition:  
293 Unsupervised oligotyping for sensitive partitioning of high- throughput marker gene sequences. 9:968–979. DOI:  
294 10.1038/ismej.2014.195.
- 295 Esling P, Lejzerowicz F, Pawlowski J 2015. Accurate multiplexing and filtering for high-throughput amplicon-sequencing.  
296 *Nucleic acids research* 43:2513–2524. DOI: 10.1093/nar/gkv107.
- 297 Folmer O, Black M, Hoeh W, Lutz R, Vrijenhoek R 1994. DNA primers for amplification of mitochondrial cytochrome c  
298 oxidase subunit I from diverse metazoan invertebrates. *Molecular marine biology and biotechnology* 3:294–299.
- 299 Gibson JF, Shokralla S, Curry C, Baird DJ, Monk WA, King I, Hajibabaei M 2015. Large-Scale Biomonitoring of Remote  
300 and Threatened Ecosystems via High-Throughput Sequencing. *PloS one* 10:e0138432–15. DOI:  
301 10.1371/journal.pone.0138432.
- 302 Hughes JM, Schmidt DJ, FINN DS 2009. Genes in Streams: Using DNA to Understand the Movement of Freshwater Fauna  
303 and Their Riverine Habitat. *BioScience* 59:573–583. DOI: 10.1525/bio.2009.59.7.8.
- 304 Lange A, Jost S, Heider D, Bock C, Budeus B, Schilling E, Strittmatter A, Boenigk J, Hoffmann D 2015. AmpliconDuo: A  
305 Split-Sample Filtering Protocol for High-Throughput Amplicon Sequencing of Microbial Communities. *PloS one*  
306 10:e0141590–22. DOI: 10.1371/journal.pone.0141590.
- 307 Leese F, Altermatt F, Bouchez A, Ekrem T 2016. DNAqua-Net: Developing new genetic tools for bioassessment and  
308 monitoring of aquatic ecosystems in Europe. *Research Ideas and Outcomes*. DOI: 10.3897/rio.2.e11321.
- 309 Leese F, Bouchez A, Abarenkov K, Altermatt F, Borja A et al. (in press). Why We Need Sustainable Networks Bridging  
310 Countries, Disciplines, Cultures and Generations for Aquatic Biomonitoring 2.0: A Perspective Derived From the  
311 DNAqua-Net COST Action. *Advances in Ecological Research* DOI: 10.1016/bs.aecr.2018.01.001
- 312 Leray M, Knowlton N 2017. Random sampling causes the low reproducibility of rare eukaryotic OTUs in Illumina COI  
313 metabarcoding. *PeerJ* 5:e3006–27. DOI: 10.7717/peerj.3006.
- 314 Múrria C, Zamora-Muñoz C, Bonada N, Ribera C, Prat N 2010. Genetic and morphological approaches to the problematic  
315 presence of three Hydropsychespecies of the pellucidulagroup (Trichoptera: Hydropsychidae) in the westernmost  
316 Mediterranean Basin. *Aquatic Insects* 32:85–98. DOI: 10.1080/01650424.2010.482939.
- 317 Nakamura K, Oshima T, Morimoto T, Ikeda S, Yoshikawa H, Shiwa Y, Ishikawa S, Linak MC, Hirai A, Takahashi H,  
318 Altaf-UI-Amin M, Ogasawara N, Kanaya S 2011. Sequence-specific error profile of Illumina sequencers. *Nucleic acids*  
319 *research* 39:e90. DOI: 10.1093/nar/gkr344.
- 320 Needham DM, Sachdeva R, Fuhrman JA 2017. Ecological dynamics and co-occurrence among marine phytoplankton,  
321 bacteria and myoviruses shows microdiversity matters. *The ISME Journal*:1–16. DOI: 10.1038/ismej.2017.29.
- 322 Pedro PM, Piper R, Bazilli Neto P, Cullen L Jr., Dropa M, Lorencao R, Matté MH, Rech TC, Rufato MO Jr., Silva M,  
323 Turati DT 2017. Metabarcoding Analyses Enable Differentiation of Both Interspecific Assemblages and Intraspecific

- 324 Divergence in Habitats With Differing Management Practices. *Environmental Entomology*:1–9. DOI:  
325 10.1093/ee/nvx166.
- 326 Peterson BK, Weber JN, Kay EH, Fisher HS, Hoekstra HE 2012. Double Digest RADseq: An Inexpensive Method for De  
327 Novo SNP Discovery and Genotyping in Model and Non-Model Species. *PLoS one* 7:e37135. DOI:  
328 10.1371/journal.pone.0037135.t001.
- 329 Rognes T, Flouri T, Nichols B, Quince C, Mahé F 2016. VSEARCH: a versatile open source tool for metagenomics. *PeerJ*  
330 4:e2584–22. DOI: 10.7717/peerj.2584.
- 331 Schirmer M 2016. Illumina Error Profiles: Resolving Fine-Scale Variation in Metagenomic Sequencing Data. *BMC*  
332 *bioinformatics*:1–15. DOI: 10.1186/s12859-016-0976-y.
- 333 Schirmer M, Ijaz UZ, D'Amore R, Hall N, Sloan WT, Quince C 2015. Insight into biases and sequencing errors for  
334 amplicon sequencing with the Illumina MiSeq platform. *Nucleic acids research*:1–16. DOI: 10.1093/nar/gku1341.
- 335 Schnell IB, Bohmann K, Gilbert MTP 2015. Tag jumps illuminated - reducing sequence-to-sample misidentifications in  
336 metabarcoding studies. *Molecular ecology resources* 15:1289–1303. DOI: 10.1111/1755-0998.12402.
- 337 Shokralla S, Gibson JF, Nikbakht H, Janzen DH, Hallwachs W, Hajibabaei M 2014. Next-generation DNA barcoding: using  
338 next-generation sequencing to enhance and accelerate DNA barcode capture from single specimens. *Molecular ecology*  
339 *resources*:n/a–n/a. DOI: 10.1111/1755-0998.12236.
- 340 Sigsgaard EE, Nielsen IB, Bach SS, Lorenzen ED, Robinson DP, Knudsen SW, Pedersen MW, Jaidah MA, Orlando L,  
341 Willerslev E, Møller PR, THOMSEN PF 2016. Population characteristics of a large whale shark aggregation inferred  
342 from seawater environmental DNA. *Nature Ecology & Evolution* 1:0004–5. DOI: 10.1038/s41559-016-0004.
- 343 Sutherland WJ, Freckleton RP, Godfray HCJ, Beissinger SR, Benton T, Cameron DD, et al. (2012). Identification of 100  
344 fundamental ecological questions. *Journal of Ecology*, 101(1), 58–67. DOI: 10.1111/1365-2745.12025
- 345 Sworobowicz L, Grabowski M, Mamos T, Burzyński A, Kilikowska A, Sell J, Wysocka A 2015. Revisiting the  
346 phylogeography of *Asellus aquaticus* in Europe: insights into cryptic diversity and spatiotemporal diversification.  
347 *Freshwater biology* 60:1824–1840. DOI: 10.1111/fwb.12613.
- 348 Taberlet P, Coissac E, Pompanon F, Brochmann C, Willerslev E 2012a. Towards next-generation biodiversity assessment  
349 using DNA metabarcoding. *Molecular Ecology* 21:2045–2050. DOI: 10.1111/j.1365-294X.2012.05470.x.
- 350 Taberlet P, Zimmermann NE, Englisch T, Tribsch A, Holderegger R, Alvarez N, Niklfeld H, Coldea G, Mirek Z, Moilanen  
351 A, Ahlmer W, Marsan PA, Bona E, Bovio M, Choler P, Cieślak E, Colli L, Cristea V, Dalmas J-P, Frajman B, Garraud  
352 L, Gaudeul M, Gielly L, Gutermann W, Jogan N, Kagalo AA, Korbecka G, Küpfer P, Lequette B, Letz DR, Manel S,  
353 Mansion G, Marhold K, Martini F, Negrini R, Niño F, Paun O, Pellecchia M, Perico G, Piękoś-Mirkowa H, Prosser F,  
354 Puşcaş M, Ronikier M, Scheuerer M, Schneeweiss GM, Schönswetter P, Schratt-Ehrendorfer L, Schüpfer F, Selvaggi  
355 A, Steinmann K, Thiel-Egenter C, van Loo M, Winkler M, Wohlgemuth T, Wraber T, Gugerli F, IntraBioDiv  
356 Consortium 2012b. Genetic diversity in widespread species is not congruent with species richness in alpine plant  
357 communities. *Ecology letters* 15:1439–1448. DOI: 10.1111/ele.12004.
- 358 Tikhonov M, Leach RW, Wingreen NS 2015. Interpreting 16S metagenomic data without clustering to achieve sub-OTU  
359 resolution. *The ISME Journal* 9:68–80. DOI: 10.1038/ismej.2014.117.
- 360 Tremblay J, Singh K, Fern A, Kirton ES, He S, Woyke T, Lee J, Chen F, Dangl JL, Tringe SG 2015. Primer and platform  
361 effects on 16S rRNA tag sequencing. *Frontiers in Microbiology* 6:8966–15. DOI: 10.3389/fmicb.2015.00771.
- 362 Vamos EE, Elbrecht V, Leese F 2017. Short COI markers for freshwater macroinvertebrate metabarcoding. *Metabarcoding*  
363 *and Metagenomics*. DOI: 10.7287/peerj.preprints.3037v1.
- 364 Wares J, Pappalardo P 2016. Can Theory Improve the Scope of Quantitative Metazoan Metabarcoding? *Diversity* 8:1–15.  
365 DOI: 10.3390/d8010001.
- 366 Weiss M & Leese F (2016). Widely distributed and regionally isolated! Drivers of genetic structure in *Gammarus fossarum*  
367 in a human-impacted landscape. *BMC Evolutionary Biology*, 1–14. DOI: 10.1186/s12862-016-0723-z