

MOLECULAR ECOLOGY RESOURCES

Optimal sequence similarity thresholds for clustering of molecular operational taxonomic units in DNA metabarcoding studies

Journal:	<i>Molecular Ecology Resources</i>
Manuscript ID	Draft
Manuscript Type:	Resource Article
Date Submitted by the Author:	n/a
Complete List of Authors:	Bonin, Aurélie; Argaly Guerrieri, Alessia; University of Milan, Environmental Science and Policy Ficaretola, Francesco; University of Milan, Department of Environmental Science and Policy
Keywords:	DNA metabarcoding marker, sequence variant, analysis parameter, MOTU over-splitting, MOTU over-merging, alpha diversity

1 **Optimal sequence similarity thresholds for clustering of molecular operational**
2 **taxonomic units in DNA metabarcoding studies**

3

4

5 **Aurélie Bonin^{1,2*}, Alessia Guerrieri¹, G. Francesco Ficetola^{1,3}**

6

7 1) Department of Environmental Science and Policy, University of Milan. Via Celoria 10,
8 20126 Milano Italy

9 2) Argaly, Bâtiment CleanSpace, 354 Voie Magellan, 73800 Sainte-Hélène-du-Lac, France

10 3) Univ. Grenoble Alpes, CNRS, Univ. Savoie Mont Blanc, LECA, Laboratoire d'Ecologie
11 Alpine, F-38000 Grenoble, France

12 *Corresponding author: aurelie.bonin@argaly.com

13 **Abstract**

14

15 Clustering approaches are pivotal to handle the many sequence variants obtained in DNA
16 metabarcoding datasets, therefore they have become a key step of metabarcoding analysis
17 pipelines. Clustering often relies on a sequence similarity threshold to gather sequences in
18 Molecular Operational Taxonomic Units (MOTUs) that ideally each represent a homogeneous
19 taxonomic entity, e.g. a species or a genus. However, the choice of the clustering threshold is
20 rarely justified, and its impact on MOTU over-splitting or over-merging even less tested.

21 Here, we evaluated clustering threshold values for several metabarcoding markers under
22 different criteria: limitation of MOTU over-merging, limitation of MOTU over-splitting, and
23 trade-off between over-merging and over-splitting. We extracted sequences from a public
24 database for eight markers, ranging from generalist markers targeting Bacteria or Eukaryota,
25 to more specific markers targeting a class or a subclass (e.g. Insecta, Oligochaeta). Based on
26 the distributions of pairwise sequence similarities within species and within genera and on the
27 rates of over-splitting and over-merging across different clustering thresholds, we were able
28 to propose threshold values minimizing the risk of over-splitting, that of over-merging, or
29 offering a trade-off between the two risks. For generalist markers, high similarity thresholds
30 (0.96-0.99) are generally appropriate, while more specific markers require lower values (0.85-
31 0.96). These results do not support the use of a fixed clustering threshold (e.g. 0.97). Instead,
32 we advocate a careful examination of the most appropriate threshold based on the research
33 objectives, the potential costs of over-splitting and over-merging, and the features of the
34 studied markers.

35

36

37 **Keywords**

38 metabarcoding marker; sequence variant; analysis parameter; MOTU over-splitting, MOTU

39 over-merging; alpha diversity

40

For Review Only

41 **Introduction**

42

43 DNA metabarcoding studies are typically based on a succession of experimental steps
44 governed by important methodological choices (Zinger et al. 2019). These include a) the
45 definition of sampling design and selection of sampling sites (Dickie et al. 2018), b) the
46 approach used for the preservation of the starting material (Tatangelo et al. 2014, Guerrieri et
47 al. 2021), c) the protocol used for DNA extraction (Taberlet et al. 2012, Eichmiller et al.
48 2016, Zinger et al. 2016, Lear et al. 2018, Capo et al. 2021), d) the selection of appropriate
49 primers to amplify a taxonomically-informative genomic region (Elbrecht et al. 2016, Fahner
50 et al. 2016, Ficetola et al. 2021), e) the strategy adopted for DNA amplification and high-
51 throughput sequencing of amplicons (Nichols et al. 2018, Taberlet et al. 2018, Bohmann et al.
52 2022), f) the pipeline selected for bioinformatics analyses (Boyer et al. 2016, Calderón-Sanou
53 et al. 2020, Capo et al. 2021, Couton et al. 2021, Macher et al. 2021, Mächler et al. 2021), and
54 g) the statistical approach used to translate metabarcoding data into ecological information
55 (Paliy and Shankar 2016, Chen and Ficetola 2020). Each of these methodological choices can
56 heavily influence the reliability and interpretation of results (Alberdi et al. 2018, Zinger et al.
57 2019), and there is thus a critical need for the development, proper assessment and
58 optimization of methods specially dedicated to DNA metabarcoding.

59 When analyzing metabarcoding data, bioinformatic pipelines generally produce a list
60 of detected sequences, that can be assigned to a given taxon with a more or less precise
61 taxonomic resolution. However, the number of unique sequences obtained after bioinformatic
62 treatment is generally much higher than the number of taxa actually present in the sample
63 (Calderón-Sanou et al. 2020, Mächler et al. 2021). This stems from multiple reasons including
64 genuine intraspecific diversity of the selected markers and errors occurring during the

65 amplification or sequencing steps. Consequently, sequence clustering approaches are often
66 used to collapse very similar sequences into one single Molecular Operational Taxonomic
67 Unit (MOTU), which does not necessarily correspond to a species in the traditional sense
68 (Kopylova et al. 2016, Froslev et al. 2017, Bhat et al. 2019, Antich et al. 2021). Sequence
69 clustering can be performed using similarity thresholds, Bayesian approaches, or through
70 single-linkage (Antich et al. 2021). Approaches based on similarity thresholds can have
71 excellent performance and they display several advantages such as flexibility and easy
72 implementation (Kopylova et al. 2016, Wei et al. 2021). However, two key parameters have
73 to be determined *a priori* when performing clustering based on sequence similarity. The first
74 one is the sequence to be selected as representative of the cluster. In the case of
75 metabarcoding studies, keeping the most abundant sequence of the cluster as the cluster
76 representative is a convenient way of merging sequence variants generated during the PCR or
77 sequencing steps with the original sequence they derive from (Mercier et al. 2013). The
78 second parameter is the similarity threshold (clustering threshold) used to build MOTUS
79 (Clare et al. 2016, Calderón-Sanou et al. 2020, Wei et al. 2021). Choosing this threshold is
80 delicate without prior knowledge of the maker and its intrinsic level of diversity. A too low
81 threshold can collapse different taxa into the same MOTU (over-merging), while a too high
82 threshold can create too many MOTUs (over-splitting) compared to the actual diversity levels
83 (Clare et al. 2016, Roy et al. 2019, Schloss 2021).

84 Some works suggest that the ecological interpretation of metabarcoding data can be
85 relatively robust to the threshold selected for sequence clustering. For instance, Botnen et al.
86 (2018) used thresholds ranging from 0.87 to 0.99 of sequence similarity to analyze multiple
87 microbial communities, and they obtained community structures highly coherent across
88 thresholds. Nevertheless, levels of alpha diversity can be heavily impacted by the threshold

89 selection. Ideally, the threshold used for clustering would depend on a trade-off between
90 MOTU over-splitting and MOTU over-merging. A growing number of markers are currently
91 being used in metabarcoding studies (Taberlet et al. 2018), with some allowing broad-scale
92 biodiversity assessment but having limited taxonomic resolution (e.g. 18S rDNA primers
93 amplifying all eukaryotes; Guardiola et al. 2015) and others being highly specific to one
94 single class or even family (e.g. Baamrane et al. 2012, Ficetola et al. 2021). Biodiversity
95 surveys generally aim to generate a set of MOTUs that are each associated with a unique
96 taxon, and with all taxa situated at the same level in the taxonomic tree, to facilitate
97 comparisons. In these conditions, optimal clustering thresholds probably strongly differ across
98 markers. One can for example expect high similarity thresholds for highly conserved markers,
99 and lower clustering thresholds for markers showing high intraspecific variability (Kunin et
100 al. 2010, Brown et al. 2015). However, there is limited quantitative assessment of how
101 optimal clustering thresholds vary across markers (but see Alberdi et al. 2018).

102 In this study, we analyzed sequences from a public database (EMBL) to identify
103 clustering thresholds for different markers and under different criteria. We considered eight
104 metabarcoding markers (Table 1), ranging from generalist ones (e.g. a 16S rDNA-based
105 marker targeting Bacteria and a 18S rDNA-based marker targeting Eukaryota) to more
106 specific markers (e.g. markers specific of earthworms, insects or springtails). We evaluated
107 how clustering thresholds can change for each taxonomic group, depending on the criterion
108 adopted to set the threshold. We used two alternative strategies to identify thresholds, each
109 time with different objectives in mind. First, following a procedure similar to the one adopted
110 in barcoding studies (Meyer and Paulay 2005), we compared the distribution probabilities of
111 sequence similarities among different individuals of the same species and among different
112 species of the same genus to identify thresholds: *i*) minimizing the risk that different

113 sequences of the same species are split in different MOTUs (i.e. risk of over-splitting); *ii*)
114 minimizing the risk that distinct but related species are clustered in the same MOTU (i.e. risk
115 of over-merging); *iii*) balancing the risk of over-splitting and over-merging (Figure 1A).
116 Second, we calculated the over-splitting and over-merging rates of the studied markers for a
117 range of clustering thresholds, to identify values that minimize the two error rates (Figure
118 1B). We expect that, if researchers want to minimize over-splitting, they should select lower
119 clustering thresholds than if they want to minimize over-merging. Furthermore, we expect
120 higher clustering threshold values for generalist markers compared to markers targeting one
121 class or more restricted taxonomic groups, because of the lower taxonomic resolution and
122 slower evolutionary rate of the former.

123

124 **Methods**

125

126 **Markers examined and construction of sequence datasets**

127 We focused on a set of eight DNA metabarcoding markers (Bact02, Euka02, Fung02, Sper01,
128 Arth02, Coll01, Inse01, Olig01) targeting different taxonomic groups (Table 1). Four of these
129 markers can be considered as generalist, i.e. targeting entire superkingdoms or kingdoms:
130 Bact02 targeting Bacteria; Euka02 targeting Eukaryota; Fung02 targeting Fungi; Sper01
131 targeting Spermatophyta (vascular plants). One marker was intermediate (Arth02; targeting
132 arthropods, i.e. the most species-rich phylum on Earth). Finally, three were more specific, i.e.
133 targeting groups from classes to subclasses: Coll01 targeting Collembola (springtails); Inse01
134 targeting Insecta; Olig01 targeting Oligochaeta (earthworms).

135 For each of these markers, a sequence database was built from EMBL release 140 as
136 follows. An *in silico* PCR was first carried out by running the program *ecoPCR* (Ficetola et al.

2010) using the corresponding primers (Table S1). Three mismatches per primer were allowed (-e option), and the amplified amplicon length without primers was restricted (-l and -L options) to the expected length interval (Table S1). The amplified sequences were further filtered by keeping only those belonging to the target taxonomic group, showing a taxonomic assignment (i.e. taxid) at the species and genus levels and having no ambiguous nucleotides. This allowed assembling a working dataset, from which we extracted two sub-datasets. The “within-species” dataset was built by keeping only species for which at least two sequences (identical or not) were available; if >2 sequences were available for a given species, we randomly selected two sequences for that species. The “within-genus” dataset was built by keeping only genera for which at least two sequences were available; if >2 sequences were available for a given genus, we randomly selected two sequences for that genus. For some markers (Bact02, Euka02, Fung02, Inse01, Sper01), the within-species dataset and sometimes the within-genus dataset still contained a very large number of sequences (>10,000). To limit computation time for these markers, we randomly selected a subset of 5000 different taxa, to reach a final number of sequences equal to 10,000. Table S2 summarizes the number of sequences in the different datasets.

153

154 **Calculation of sequence similarities and probability distributions**

155 As a measure of sequence similarity, we computed the pairwise LCS (Longest Common
156 Subsequence) scores between pairs of sequences in the within-species and within-genus
157 datasets using the *sumatra* program (Mercier et al. 2013). Methodological comparisons
158 showed that this algorithm provides an excellent balance between performance and
159 computation efficiency (Jackson et al. 2016, Kopylova et al. 2016, Bhat et al. 2019).As
160 *sumatra* provides pairwise scores for all possible pairs of sequences, the similarity scores

161 resulting from the within-species dataset were filtered in R (R Core Team 2020) to keep only
162 those representing similarities between sequences of the same species, while the scores
163 resulting from the within-genus dataset were filtered to keep only those representing
164 similarities between different species of the same genus.

165

166 **Approaches to identify clustering thresholds on the basis of within-species and within-** 167 **genus sequence similarities**

168 We first examined within-species and within-genus sequence similarities to evaluate four
169 different strategies and determine the corresponding appropriate clustering threshold (Figure
170 1A) that: *i*) avoid over-splitting; *ii*) avoid over-merging; *iii*) find a balance between over-
171 splitting and over-merging, with two distinct procedures based on the intersection (*iii*-a) or on
172 modes (*iii*-b) of the density probability distributions. These strategies are analogous to those
173 adopted in traditional barcoding studies to set the limit between intra-specific and inter-
174 specific diversity (Meyer and Paulay 2005).

175 ***i*) Avoid over-splitting**

176 In this case, the aim is to avoid distributing different sequences belonging to the same species
177 in different clusters (i.e. limiting the probability of generating additional spurious MOTUs).
178 For this approach, we selected as clustering threshold the 10% quantile of the distribution of
179 similarities between sequences from the same species (within-species dataset). With this
180 approach, the sequences belonging to the same species according to EMBL are gathered in
181 the same cluster in 90% of the cases.

182 ***ii*) Avoid over-merging**

183 In this case, the aim is to avoid gathering sequences attributed to different species of the same
184 genus in the same cluster (i.e. limiting the probability of merging related species in the same

185 MOTU). For this approach, we selected as clustering threshold the 90% quantile of the
186 distribution of similarities between different species belonging to the same genus. With this
187 approach, the sequences attributed to different species belonging to the same genus are
188 assigned to different clusters in 90% of the cases.

189 **iii) Find a balance between over-splitting and over-merging**

190 In this case, the aim was to minimize both over-splitting and over-merging. We considered
191 two distinct approaches. First, we obtained the probability distribution of within-species and
192 within-genus sequence pairwise similarities using the *density* function from R, with biased
193 cross-validation (bw="bcv") as smoothing bandwidth selector and a Gaussian smoothing
194 kernel (kernel="gaussian"; Venables and Ripley 2002). Other possible smoothing bandwidth
195 selectors were tested, but biased cross-validation was the approach best fitting the score
196 histograms for all markers and all datasets (data not shown). The balance threshold *iii-a* was
197 then identified as the intersection between the probability distributions of the within-species
198 and within-genus similarities. As an alternative approach to balance over-merging and over-
199 splitting (*iii-b*), we calculated the midpoint between the modes of the within-species and
200 within-genus probability distributions.

201

202 **Rates of over-merging and over-splitting**

203 For each marker, over-merging and over-splitting rates were evaluated at different clustering
204 thresholds using the within-species dataset described in the paragraph "Markers examined and
205 construction of sequences datasets". This dataset contains two sequences at random, identical
206 or not, for a number of species belonging to the taxonomic group of interest.

207 For each within-species dataset, clustering was performed using the *sumaclust*
208 program (Mercier et al. 2013) with the *-n* option (normalization by alignment length) based

209 on the sequence similarities first calculated using the *sumatra* program (see above; Mercier et
210 al. 2013). Threshold values (*-t* option) ranging from 0.90 to 1 at 0.01 steps were tested for all
211 markers except Coll01 and Olig01 for which wider ranges ([0.70 – 1] and [0.80 – 1],
212 respectively) were selected based on the within-genus and within-species sequence similarity
213 probability distributions determined previously (see Figure 2). Clustered datasets were then
214 explored to calculate five different variables at each clustering threshold: 1) the number of
215 clusters; 2) the percentage of MOTUs containing one single species; 3) the percentage of
216 MOTUs containing one single genus; 4) the percentage of species gathered in one single
217 MOTU; 5) the percentage of genera gathered in one single MOTU. Variables 2 and 3 are
218 indicative of appropriate MOTU merging of sequences at the species and genus levels,
219 respectively, while variables 4 and 5 are indicative of appropriate MOTU splitting at the
220 species and genus levels, respectively.

221 These values were also used to calculate three measures of error. We defined the over-
222 merging rate as 1 - the percentage of MOTUs containing one single species; and the over-
223 splitting rate as 1 - the percentage of species gathered in one single MOTU. The summed
224 error rate was then calculated as the sum of the over-merging and over-splitting rates. It
225 should be noted that for this estimate, we assigned the same weight to over-splitting and over-
226 merging.

227

228 **Results**

229

230 Our *in-silico* PCRs amplified between 17,000 (Coll01) and 3,200,000 (Bact02) sequences
231 per marker (Table S2). After data filtering, we retained between 510 (Coll01) and 708,000
232 (Bact02) sequences per marker. The within-species dataset comprised between 118 (Coll01)

233 and 10,000 (Bact02, Euka02, Fung02, Sper01, Inse01) sequences, while the within-genus
234 dataset comprised between 74 (Coll01) and 10,000 (Euka02 and Sper01) sequences per
235 marker.

236

237 **Clustering thresholds determined from probability distributions of within-species and** 238 **within-genus sequence similarities**

239 The probability distributions of within-species and within-genus sequence similarities
240 showed very contrasting patterns between the generalist and the specific markers (Figure 2).
241 For the five markers targeting a phylum or broader taxonomic groups (Bact02, Euka02,
242 Fung02, Sper01, and Arth02), the distributions of within-species and within-genus similarities
243 were rather similar, both showing a mode at very high similarity values (Figure 2). Fung02
244 showed a slightly different pattern, as the within-genus similarities had a very broad
245 distribution. Conversely, for the more specific markers, the distributions of sequence
246 similarities were very different, with two clearly distinct peaks. Within-species similarities
247 remained very high (mostly above 0.95), while within-genus similarities generally showed
248 lower values (mode around 0.90 for Inse01, and below 0.80 for Olig01 and Coll01).

249 For all markers, criterion *i* (avoid over-splitting) yielded the lowest thresholds (Figure
250 3, Table S3), with very low levels for Coll01 and Olig01. Conversely, criterion *ii* (avoid over-
251 merging) yielded extremely high values, except for Coll01. For all generalist markers,
252 avoiding over-merging would require setting clustering thresholds at 0.99 or higher. For
253 Coll01, criterion *ii* resulted in a rather low threshold (0.765), because many within-genus
254 comparisons showed very low similarity values.

255 Criteria *iii-a* and *iii-b* searching a balance between over-merging and over-splitting
256 yielded somehow contrasting results across markers. For the three specific markers (Coll01,

257 Inse01, and Olig01), the within-genus and within-species similarities showed clearly distinct
258 peaks (Figure 2). As a consequence, the intersection between the two curves could effectively
259 represent the point minimizing both over-merging and over-splitting (see discussion), and the
260 midpoint between the modes also identified rather similar threshold values. On the contrary,
261 for the generalist markers, the within-species and within-genus similarities showed very high
262 overlap and similar modes, and the density distributions actually intersected at values lower
263 than both modes. The midpoint between the modes continued to identify threshold values
264 intermediate between the peaks of within-species and within-genus similarities.

265

266 **Rates of over-splitting and over-merging**

267 For all markers, whatever the clustering threshold examined (values ≥ 0.70 for Coll01, ≥ 0.80
268 for Olig01 and ≥ 0.90 for the other markers), the percentage of MOTUs containing one single
269 species was higher than 50%, and that of MOTUs containing one single genus was higher or
270 close to 70% (Figure 4). Overall, for the generalist and intermediate markers, these two
271 percentages showed a regular increase with the clustering threshold, and for the specific
272 markers, they tended to values close to 100% for high thresholds. Unsurprisingly, the two
273 percentages tended to be lower for the generalist markers than for the specific markers at a
274 given threshold, indicating that the former are more sensitive to over-merging. Fung02 was a
275 notable exception, since about 87% and 97% of MOTUs contained one single species and one
276 single genus, respectively, at the 0.97 threshold, which is a frequently adopted clustering
277 threshold for fungal ITS sequences. These values were comparable to those observed for the
278 specific markers, for which $> 85\%$ and $> 98\%$ of MOTUs contained one single species or one
279 single genus, respectively, for thresholds ≥ 0.95 .

280 For all markers, the percentages of species and genera gathered in one single MOTU
281 decrease both at a similar rate with the clustering threshold, with generally a sharp drop at
282 high thresholds (≥ 0.98 ; Figure 4). However, the pattern of MOTU splitting was less
283 characteristic of generalist vs. specific markers. For some markers (Euka02, Sper01, Arth02,
284 Inse01), the percentage of species or genera gathered in a single MOTU remained higher or
285 close to 50% up to high thresholds (0.98). On the contrary, for Bact02, Fung02, Coll01,
286 Olig01, these percentages dropped quickly when the clustering threshold increased, indicating
287 that these markers are susceptible to over-splitting.

288 For all markers, the number of clusters generally increased regularly with the clustering
289 threshold up to 0.97-0.98 (Figure 4), followed by a sharp rise up to 1 (which was however less
290 obvious for Euka02 and Olig01). For example, for Bact02, the number of clusters more than
291 doubled between 0.97 (2862 clusters) and 1 (6461 clusters).

292 Our results showed clear patterns for over-merging and over-splitting rates, with over-
293 splitting quickly increasing and over-merging quickly decreasing at high clustering thresholds
294 (Figure 5). For several markers, the summed error showed a relatively clear minimum at
295 specific clustering thresholds (Figure 5): 0.96-0.99 for Bact02 and Euka02, 0.97-0.99 for
296 Arth02, 0.94-0.96 for Inse01, and 0.96-0.98 for Sper01. The minimum was much less evident
297 for Fung02, Coll01 and Oligo01, these markers showing relatively similar summed error rates
298 over a broad range of clustering thresholds (Fung02: 0.91-0.98; Coll01: 0.82-0.96, with
299 multiple minima; Oligo01: 0.84-0.96, with multiple minima).

300

301

302 **DISCUSSION**

303

304 Sequence clustering approaches are routinely used for the identification of MOTUs in
305 metabarcoding studies, and they often resort to methods based on similarity values. Still,
306 selecting a clustering threshold for a given marker more than often relies on common
307 practices and rules of thumb rather than on proper scientific argument. By analyzing extensive
308 sequence data deposited in public databases for a range of generalist and specialist markers,
309 we showed that different threshold values can be selected depending on the marker and on the
310 criterion favored by researchers. All the markers we examined are situated in non-protein
311 coding genes (Table S1), and this has an influence on levels of sequence intraspecific
312 diversity. The 10% quantile of the within-species similarity probability distribution was
313 almost always lower than the 0.97 clustering threshold traditionally used in barcoding for
314 markers targeting protein-coding genes like COI (Hebert et al. 2003), or for microbial MOTU
315 delimitation (Bálint et al. 2016), indicating that some level of over-splitting can occur at this
316 threshold.

317 Although for all the markers the within-genus similarity values were generally lower
318 than the within-species similarities, the overlap between the two distributions was dependent
319 on the generalist vs. specific nature of the marker. For some specific markers (e.g. Coll01 and
320 Olig01), distinct peaks were visible for the two similarity metrics (Figure 2). Within-species
321 similarities generally were >0.90 , while within-genus values were lower, frequently below
322 0.80. Such a pattern is not unexpected for markers with an excellent taxonomic resolution and
323 designed to identify taxa at the species level. Conversely for the generalist markers, within-
324 species and within-genus similarity probability distributions largely overlapped and the
325 differences between the peaks were minimal. Nevertheless, even for these markers, the
326 density of within-species similarity was consistently higher than that of within-genus
327 similarity at high clustering thresholds, indicating that the probability of observing the

328 corresponding similarity value is higher within species than within genera. In other words, at
329 high clustering thresholds, a MOTU is more likely to represent a species than a genus. This
330 result is confirmed by the fact that the percentage of MOTUs containing a single species is
331 always higher than 50%, whatever the clustering threshold or the marker considered (Figure
332 4).

333 The sequences used as a primary source of information in this study were downloaded
334 from EMBL, and our results are thus highly dependent on the quality of the data deposited in
335 this public database. Even though broad-scale analyses suggest that these data are generally
336 reliable (Leray et al. 2019), errors in the sequence itself (e.g. wrong nucleotide, or more
337 complex errors like insertions, deletions, inversions, duplications or pseudogene sequences)
338 and taxonomic mislabeling can occur in public sequence databases, especially for organisms
339 which are difficult to identify based on morphology (Bridge et al. 2003, Bidartondo 2008,
340 Valkiūnas et al. 2008, Mioduchowska et al. 2018). While the first type of error will affect
341 within-species sequence similarity negatively, sometimes substantially, the effect of the
342 second type is more diffuse. For example, in a group like springtails where species
343 delimitation is tricky (Porco et al. 2012), the existence of cryptic species will decrease within-
344 species sequence similarity while increasing over-splitting rates. In a group like Bacteria, type
345 strains are sometimes entered at the species level in the NCBI (EMBL) taxonomy (Federhen
346 2015), leading to an inflation of within-genus similarity and over-merging rates. In every case
347 though, database errors will make within-species and within-genus similarities distributions
348 more difficult to distinguish and clustering thresholds trickier to identify, thus the over-
349 splitting or over-merging rates reported here could be artificially higher than in reality.

350 In this work, we came up with a global measure of the error associated with a given
351 clustering threshold, that we called the “summed error”. We calculate it by summing over-

352 splitting and over-merging rates, assuming both have the same cost for biodiversity studies.
353 However, it is possible to assign a differential weight to over-splitting and over-merging. For
354 instance, if the aim is to reach conservative estimated of alpha diversity (i.e. avoid over-
355 splitting), more weight can be assigned to over-splitting rate. Conversely, if the aim is to tease
356 apart closely related species, that differ in their sensitivity to environmental stressors or in
357 threat levels, one may prefer to avoid over-merging, particularly when extensive reference
358 databases are available (Roy et al. 2019, Lopes et al. 2021).

359 For most of the markers we examined, the summed error approach provided relatively
360 clear results, and identified a range of threshold values that minimized the summed error. For
361 instance, for Euka02, the summed error was relatively low at thresholds between 0.96 and
362 0.99 (Figure 5), indicating a good trade-off between over-merging and over-splitting.
363 Interestingly, this range of values was also highlighted by the analysis of probability
364 distributions (Figure 3, Table S3). Indeed, 0.96 is the threshold minimizing over-splitting for
365 Euka02 while 0.99 is the balance (midpoint) threshold. The consistency of values obtained
366 with very different approaches supports the robustness of our conclusions.

367 However, for a few markers, the threshold values minimizing summed error yielded
368 somewhat less clear patterns. For Fung02, the summed error rate was rather constant (36-
369 37%) at all the thresholds between 0.91 and 0.98, while it quickly increased for higher
370 clustering thresholds. For Coll01 and Oligo01, the summed error rate showed multiple
371 minima, some of which at very low clustering thresholds (Figure 5). In principle, increasing
372 the threshold value should determine a monotone decrease of over-merging, and a monotone
373 increase of over-splitting (Figure 1B). However, at low similarity values this was not always
374 the case (Figure 5). This probably occurs because a very large number of sequences have
375 pairwise similarities of 0.80-0.85 for these markers (Figure 2), and this might affect the

376 identification of clusters, with some sequences clustering together e.g. at 0.85 but not at 0.86
377 similarity values. We also note that these similarity values match the ones corresponding to
378 the intersection between the within-genus and within-species similarities for these markers
379 (Figure 3). It is also possible that, at this level of sequence similarity, there is strong
380 uncertainty between MOTUs representing different hierarchical levels of taxonomy.

381 Our results provide quantitative data that can help researchers set their optimal
382 clustering thresholds, and understand the consequences of choosing low or high threshold
383 values. If a clear minimum exists for the summed error rate, it probably represents an
384 excellent trade-off between over-merging and over-splitting. In this sense, a threshold value
385 ranging from 0.96 to 0.99 is probably appropriate for both Bact02 and Euka02, while Arth02
386 should accommodate a slightly higher range (0.98-0.99) and a fixed threshold of 0.97 seems
387 to be more suitable for Sper01. For Inse01, lower threshold values (0.94-0.96) are more
388 judicious. All these values match with those obtained on the basis of within-species and
389 within-genus similarities (Figure 3). However, for Coll01, Oligo01 and Fung02, the summed
390 error rate does not provide clear indications, and within-species and within-genus similarity
391 distributions (e.g. midpoint between modes) might be more informative to set the threshold
392 value (Figures 2 and 3).

393 The selection of clustering thresholds can have strong effect in the estimates of
394 MOTUs richness (Figure 4), still it is important to remember that it often does not have a
395 tremendous effect on the ecological message conveyed by metabarcoding data. For instance,
396 Clare et al. (2016) examined different clustering thresholds to analyze dietary overlap
397 between skinks and shrews in Mauritius. Although high clustering thresholds yielded a larger
398 number of MOTUs, ecological conclusions remained rather consistent overall. Therefore,
399 provided that appropriate parameters are considered (e.g. alpha diversity measured using

400 Hill's numbers with $q > 0$ instead of richness, beta diversity estimates), the interpretation of
401 data can be relatively robust (Clare et al. 2016, Roy et al. 2019, Calderón-Sanou et al. 2020,
402 Mächler et al. 2021). Nevertheless, we discourage the blind application of one single
403 clustering threshold like the classical 0.97, as it can have very different meaning across
404 markers, and can inflate MOTU richness for fast-evolving markers. Instead, we advocate the
405 ad-hoc definition of the most appropriate thresholds, on the basis of research aims, on the
406 potential costs of over-splitting and over-merging, and on the features of the studied markers.

407

408 **Acknowledgments**

409 This study was supported by the European Research Council under the European Community's
410 Horizon 2020 Programme, Grant Agreement no. 772284 (IceCommunities).

411

412 **References**

- 413 Alberdi, A., Aizpurua, O., Gilbert, M. T. P., & Bohmann, K. (2018). Scrutinizing key steps for
414 reliable metabarcoding of environmental samples. *Methods in Ecology and Evolution*,
415 9, 134-147.
- 416 Antich, A., Palacin, C., Wangenstein, O. S., & Turon, X. (2021). To denoise or to cluster, that
417 is not the question: optimizing pipelines for COI metabarcoding and
418 metaphylogeography. *BMC Bioinformatics*, 22, 177.
- 419 Baamrane, M. A. A., Shehzad, W., Ouhammou, A., Abbad, A., Naimi, M., Coissac, E., ...
420 Znari, M. (2012). Assessment of the food habits of the Moroccan dorcas gazelle in
421 M'Sabih Talaa, West Central Morocco, using the *trnL* approach. *PLoS ONE*, 7, e35643.
- 422 Bálint, M., Bahram, M., Eren, A. M., Faust, K., Fuhrman, J. A., Lindahl, B., ... Tedersoo, L.
423 (2016). Millions of reads, thousands of taxa: microbial community structure and
424 associations analyzed via marker genes. *FEMS Microbiology Reviews*, 40, 686-700.
- 425 Bhat, A. H., Prabhu, P., & Balakrishnan, K. (2019). A critical analysis of state-of-the-art
426 metagenomics OTU clustering algorithms. *Journal of Biosciences*, 44, 9.
- 427 Bidartondo, M. I. (2008). Preserving accuracy in GenBank. *Science*, 319, 1616.
- 428 Bienert, F., De Danieli, S., Miquel, C., Coissac, E., Poillot, C., Brun, J. J., & Taberlet, P. (2012).
429 Tracking earthworm communities from soil DNA. *Molecular Ecology*, 21, 2017-2030.
- 430 Bohmann, K., Elbrecht, V., Carøe, C., Bista, L., Leese, F., Bunce, M., Yu, D. W., ... Creer, S.
431 (in press). Strategies for sample labelling and library preparation in DNA
432 metabarcoding studies. *Molecular Ecology Resources*.

- 433 Botnen, S. S., Davey, M. L., Halvorsen, R., & Kauserud, H. (2018). Sequence clustering
434 threshold has little effect on the recovery of microbial community structure. *Molecular*
435 *Ecology Resources*, *18*, 1064-1076.
- 436 Boyer, F., Mercier, C., Bonin, A., Le Bras, Y., Taberlet, P., & Coissac, E. (2016). OBITOOLS:
437 a Unix-inspired software package for DNA metabarcoding. *Molecular Ecology*
438 *Resources*, *16*, 176-182.
- 439 Bridge, P. D., Roberts, P. J., Spooner, B. M., & Panchal, G. (2003). On the unreliability of
440 published DNA sequences. *New Phytologist*, *160*, 43-48.
- 441 Brown, E. A., Chain, F. J. J., Crease, T. J., MacIsaac, H. J., & Cristescu, M. E. (2015).
442 Divergence thresholds and divergent biodiversity estimates: can metabarcoding reliably
443 describe zooplankton communities? *Ecology and Evolution*, *5*, 2234-2251.
- 444 Calderón-Sanou, I., Münkemüller, T., Boyer, F., Zinger, L., & Thuiller, W. (2020). From
445 environmental DNA sequences to ecological conclusions: How strong is the influence
446 of methodological choices? *Journal of Biogeography*, *47*, 193–206.
- 447 Capo, E., Giguet-Covex, C., Rouillard, A., Nota, K., Heintzman, P., Vuillemin, A. ... Parducci,
448 L. (2021). Lake sedimentary DNA research on past terrestrial and aquatic biodiversity:
449 Overview and recommendations. *Quaternary* *4*, 6.
- 450 Chen, W., & Ficetola, G. F. (2020). Statistical and numerical methods for Sedimentary-ancient-
451 DNA-based study on past biodiversity and ecosystem functioning. *Environmental DNA*,
452 *2*, 115–129.
- 453 Clare, E. L., Chain, F. J. J., Littlefair, J. E., & Cristescu, M. E. (2016). The effects of parameter
454 choice on defining molecular operational taxonomic units and resulting ecological
455 analyses of metabarcoding data. *Genome*, *59*, 981-990.
- 456 Couton, M., Baud, A., Daguin-Thiébaud, C., Corre, E., Comtet, T., & Viard, F. (2021). High-
457 throughput sequencing on preservative ethanol is effective at jointly examining
458 infraspecific and taxonomic diversity, although bioinformatics pipelines do not perform
459 equally. *Ecology and Evolution*, *11*, 5533-5546.
- 460 Dickie, I. A., Boyer, S., Buckley, H. L., Duncan, R. P., Gardner, P. P., Hogg, I. D., ... Weaver,
461 L. (2018). Towards robust and repeatable sampling methods in eDNA-based studies.
462 *Molecular Ecology Resources*, *18*, 940-952.
- 463 Eichmiller, J. J., Miller L. M., & Sorensen, P.W. (2016). Optimizing techniques to capture and
464 extract environmental DNA for detection and quantification of fish. *Molecular Ecology*
465 *Resources*, *16*, 56-68.
- 466 Elbrecht, V., Taberlet, P., Dejean, T., Valentini, A., Usseglio-Polatera, P., Beisel, J. N., ...
467 Leese, F. (2016). Testing the potential of a ribosomal 16S marker for DNA
468 metabarcoding of insects. *PeerJ*, *4*, 12.
- 469 Epp, L. S., Boessenkool, S., Bellemain, E. P., Haile, J., Esposito, A., Riaz, T., ... Brochmann,
470 C. (2012). New environmental metabarcodes for analysing soil DNA: potential for
471 studying past and present ecosystems. *Molecular Ecology*, *21*, 1821-1833.
- 472 Fahner, N. A., Shokralla, S., Baird, D. J., & Hajibabaei, M. (2016). Large-scale monitoring of
473 plants through environmental DNA metabarcoding of soil: Recovery, resolution, and
474 annotation of four DNA markers. *PLoS ONE*, *11*, e0157505.
- 475 Federhen, S. (2015). Type material in the NCBI Taxonomy Database. *Nucleic Acids Research*,
476 *43*, D1086-D1098.
- 477 Ficetola, G. F., Boyer, F., Valentini, A., Bonin, Meyer, A., Dejean, T., ... Taberlet, P. (2021).
478 Comparison of markers for the monitoring of freshwater benthic biodiversity through
479 DNA metabarcoding. *Molecular Ecology*, *30*, 3189–3202.

- 480 Ficetola, G. F., Coissac, E., Zundel, S., Riaz, T., Shehzad, W., Bessière, J., ... Pompanon, F.
481 (2010). An *in silico* approach for the evaluation of DNA barcodes. *BMC Genomics*, *11*,
482 434.
- 483 Froslev, T. G., Kjoller, R., Bruun, H. H., Ejrnaes, R., Brunbjerg, A. K., Pietroni, C., & Hansen,
484 A. J. (2017). Algorithm for post-clustering curation of DNA amplicon data yields
485 reliable biodiversity estimates. *Nature Communications*, *8*, 11.
- 486 Guardiola, M., Uriz, M. J., Taberlet, P., Coissac, E., Wangenstein, O. S., & Turon, X. (2015).
487 Deep-sea, deep-sequencing: metabarcoding extracellular DNA from sediments of
488 marine canyons. *PLoS ONE*, *10*, e0139633.
- 489 Guerrieri, A., Bonin, A., Münkemüller, T., Gielly, L., Thuiller, W., & Ficetola, G. F. (2021).
490 Effects of soil preservation for biodiversity monitoring using environmental DNA.
491 *Molecular Ecology*, *30*, 3313-3325.
- 492 Hebert, P. D. N., Ratnasingham, S., & deWaard, J. R. (2003). Barcoding animal life:
493 cytochrome c oxidase subunit 1 divergences among closely related species. *Proceedings*
494 *of the Royal Society B-Biological Sciences*, *270*, S96-S99.
- 495 Jackson, M. A., Bell, J. T., Spector, T. D., & Steves, C. J. (2016). A heritability-based
496 comparison of methods used to cluster 16S rRNA gene sequences into operational
497 taxonomic units. *PeerJ*, *4*, 19.
- 498 Janssen, P., Bec, S., Fuhr, M., Taberlet, P., Brun, J.-J., & Bouget, C. (2018). Present conditions
499 may mediate the legacy effect of past land-use changes on species richness and
500 composition of above- and below-ground assemblages. *Journal of Ecology*, *106*, 306-
501 318.
- 502 Kopylova, E., Navas-Molina, J. A., Mercier, C., Xu, Z. Z., Mahe, F., He, Y., ... Knight, R.
503 (2016). Open-source sequence clustering methods improve the state of the art.
504 *mSystems*, *1*, 16.
- 505 Kunin, V., Engelbrekton, A., Ochman, H., & Hugenholtz, P. (2010). Wrinkles in the rare
506 biosphere: pyrosequencing errors can lead to artificial inflation of diversity estimates.
507 *Environmental Microbiology*, *12*, 118-123.
- 508 Lear, G., Dickie, I., Banks, J., Boyer, S., Buckley, H. L., Buckley, T. R., ... Holdaway, R.
509 (2018). Methods for the extraction, storage, amplification and sequencing of DNA from
510 environmental samples. *New Zealand Journal of Ecology*, *42*, 10.
- 511 Leray, M., Knowlton, N., Ho, S. L., Nguyen, B. N., & Machida, R. J. (2019). GenBank is a
512 reliable resource for 21st century biodiversity research. *Proceedings of the National*
513 *Academy of Sciences of the United States of America*, *116*, 22651-22656.
- 514 Lopes, C. M., Baêta, D., Valentini, A., Lyra, M. L., Sabbag, A. F., Gasparini, J. L., ... Zamudio,
515 R. K. (2021). Lost and found: Frogs in a biodiversity hotspot rediscovered with
516 environmental DNA. *Molecular Ecology*, *30*, 3289-3298.
- 517 Macher, T.-H., Beermann, A. J., & Leese, F. (2021). TaxonTableTools: A comprehensive,
518 platform-independent graphical user interface software to explore and visualise DNA
519 metabarcoding data. *Molecular Ecology Resources*, *21*, 1705-1714.
- 520 Mächler, E., Walser, J.-C., & Altermatt, F. (2021). Decision-making and best practices for
521 taxonomy-free environmental DNA metabarcoding in biomonitoring using Hill
522 numbers. *Molecular Ecology*, *30*, 3326-3339.
- 523 Mercier, C., Boyer, F., Bonin, A., & Coissac, E. (2013). SUMATRA and SUMACLUSt: fast
524 and exact comparison and clustering of sequences. *Programs and Abstracts of the*
525 *SeqBio 2013 Workshop*, 27-29.
- 526 Meyer, C. P., & Paulay, G. (2005). DNA barcoding: Error rates based on comprehensive
527 sampling. *PLoS Biology*, *3*, 2229-2238.

- 528 Mioduchowska, M., Czyz, M. J., Goldyn, B., Kur, J., & Sell, J. (2018). Instances of erroneous
529 DNA barcoding of metazoan invertebrates: Are universal *cox1* gene primers too
530 "universal" ? *PLoS ONE*, *13*, e0199609.
- 531 Nichols, R. V., Vollmers, C., Newsom, L. A., Wang, Y., Heintzman, P. D., Leighton, M., ...
532 Shapiro, B. (2018). Minimizing polymerase biases in metabarcoding. *Molecular*
533 *Ecology Resources*, *18*, 927-939.
- 534 Paliy, O., & Shankar, V. (2016). Application of multivariate statistical techniques in microbial
535 ecology. *Molecular Ecology*, *25*, 1032-1057.
- 536 Porco, D., Bedos, A., Penelope, G., Janion, C., Skarżyński, D., Stevens, M. I., ... Deharveng,
537 L. (2012). Challenging species delimitation in Collembola: cryptic diversity among
538 common springtails unveiled by DNA barcoding. *Invertebrate Systematics*, *26*, 470-
539 477.
- 540 R Core Team. (2020). R: A language and environment for statistical computing. R Foundation
541 for Statistical Computing, Vienna.
- 542 Roy, J., Mazel, F., Sosa-Hernández, M. A., Dueñas, J. F., Hempel, S., Zinger, L., & Rillig, M.
543 C. (2019). The relative importance of ecological drivers of arbuscular mycorrhizal
544 fungal distribution varies with taxon phylogenetic resolution. *New Phytologist*, *224*,
545 936-948.
- 546 Schloss, P. D. (2021). Amplicon sequence variants artificially split bacterial genomes into
547 separate clusters. *mSphere*, *6*, e00191-00121.
- 548 Taberlet, P., Bonin, A., Zinger, L., & Coissac, E. (2018). Environmental DNA for biodiversity
549 research and monitoring. Oxford University Press, Oxford.
- 550 Taberlet, P., Coissac, E., Pompanon, F., Gielly, L., Miquel, C., Valentini, A., ... Willerslev, E.
551 (2007). Power and limitations of the chloroplast *trnL* (UAA) intron for plant DNA
552 barcoding. *Nucleic Acids Research*, *35*, e14.
- 553 Taberlet, P., Prud'homme, S. M., Campione, E., Roy, J., Miquel, C., Shehzad, W., ... Coissac,
554 E. (2012). Soil sampling and isolation of extracellular DNA from large amount of
555 starting material suitable for metabarcoding studies. *Molecular Ecology*, *21*, 1816-1820.
- 556 Tatangelo, V., Franzetti, A., Gandolfi, I., Bestetti, G., & Ambrosini, R. (2014). Effect of
557 preservation method on the assessment of bacterial community structure in soil and
558 water samples. *FEMS Microbiology Letters*, *356*, 32-38.
- 559 Valkiūnas, G., Atkinson, C. T., Bensch, S., Sehgal, R. N., & Ricklefs, R. E. (2008). Parasite
560 misidentifications in GenBank: how to minimize their number? *Trends in Parasitology*,
561 *24*, 247-248.
- 562 Venables, W. N., & Ripley, B. D. (2002). Modern applied statistics with S. Fourth Edition.
563 Springer, New York.
- 564 Wei, Z.-G., Zhang, X.-D., Cao, M., Liu, F., Qian, Y., & Zhang, S.-W. (2021). Comparison of
565 methods for picking the operational taxonomic units from amplicon sequences.
566 *Frontiers in Microbiology*, *12*, 644012.
- 567 Zinger, L., Bonin, A., Alsos, I., Bálint, M., Bik, H., Boyer, F., ... Taberlet, P. (2019). DNA
568 metabarcoding - need for robust experimental designs to draw sound ecological
569 conclusions. *Molecular Ecology*, *28*, 1857-1862.
- 570 Zinger, L., Chave, J., Coissac, E., Iribar, A., Louisanna, E., Manzi, S., ... Taberlet, P. (2016).
571 Extracellular DNA extraction is a fast, cheap and reliable alternative for multi-taxa
572 surveys based on soil DNA. *Soil Biology & Biochemistry*, *96*, 16-19.
- 573

574

575 **Data Accessibility**

576 Raw data obtained from EMBL r140 (*ecopcr* files) and example scripts run to perform the
577 analyses are available on Dryad: <https://doi.org/10.5061/dryad.crjdfn353>.

578

579 **Authors Contribution**

580 All authors conceived the idea for the manuscript, AB and GFF designed the study, AB
581 performed the analyses, AB and GFF generated the figures and drafted the manuscript, and all
582 authors contributed with discussions and edits.

583

For Review Only

584 **Table 1. Target groups and taxonomic resolution of the eight studied markers.**

Marker	Target group	Taxonomic level	Taxonomic resolution *				Reference(s)
			Species level	Genus level	Family level	Order level	
Bact02	Bacteria	Superkingdom	19.6%	55.7%	55.1%	60.2%	Taberlet et al. (2018)
Euka02	Eukaryota	Superkingdom	47.0%	59.5%	68.3%	67.1%	Guardiola et al. (2015)
Fung02	Fungi	Kingdom	72.5%	90.2%	87.7%	85.5%	Epp et al. (2012), Taberlet et al. (2018)
Sper01	Spermatophyta	Clade < kingdom	21.5%	36.9%	77.4%	89.6%	Taberlet et al. (2007)
Arth02	Arthropoda	Phylum	68.6%	89.6%	97.5%	100.0%	Taberlet et al. (2018)
Coll01	Collembola	Class	80.5%	87.2%	75.0%	NA	Janssen et al. (2018)
Inse01	Insecta	Class	87.8%	96.8%	95.4%	79.3%	Taberlet et al. (2018)
Olig01	Oligochaeta	Subclass	89.3%	95.7%	100.0%	100.0%	Bienert et al. (2012), Taberlet et al. (2018)

585

586 * Estimated as the percentage of discriminated taxa among amplified taxa; reported from

587 Taberlet et al. (2018).

588

589 **Figure captions**

590

591 **Figure 1.** Different approaches to identify the most appropriate clustering thresholds. A):
592 approaches based on similarities between sequences belonging to different individuals from
593 the same species (blue curve), and similarities between sequences belonging to different
594 species from the same genus (red curve). One can choose to minimize the risk that different
595 sequences from the same species are split in different MOTUs (over-splitting risk; e.g. 10%
596 quantile of the distribution of within-species similarities), the risk that sequences from
597 different species belonging to the same genus are clustered in the same MOTU (over-merging
598 risk; e.g. 90% quantile of within-genus similarities), or one can try to find a balance between
599 the risks of over-splitting and over-merging (e.g. with the intersection between probability
600 distributions, or the midpoint between the modes of both distributions). B) Approaches based
601 on rates of over-splitting and over-merging. One can compare the over-splitting (blue) and the
602 over-merging (red) rates, and/or one can identify the thresholds minimizing the sum of these
603 rates (violet).

604

605 **Figure 2.** Density probability distributions of sequence pairwise similarities within species
606 (blue lines) and within genera (red lines) for the eight studied markers. For each marker,
607 dotted lines represent the 10% quantile of the within-species probability distribution (blue;
608 threshold limiting over-splitting), the 90% quantile of the within-genus probability
609 distribution (red; threshold limiting over-merging), the intersection of the within-species and
610 within-genus probability distributions (green, balance-a) and the midpoint between modes
611 (black, balance-b)

612

613 **Figure 3.** Different possible clustering thresholds for the eight studied markers, depending on
614 the selected criterion.

615

616 **Figure 4.** Evolution of over-splitting and over-merging rates for a range of clustering
617 thresholds, for the eight studied markers. The left y-axes report percentage values; the right y-
618 axes indicate the number of obtained clusters.

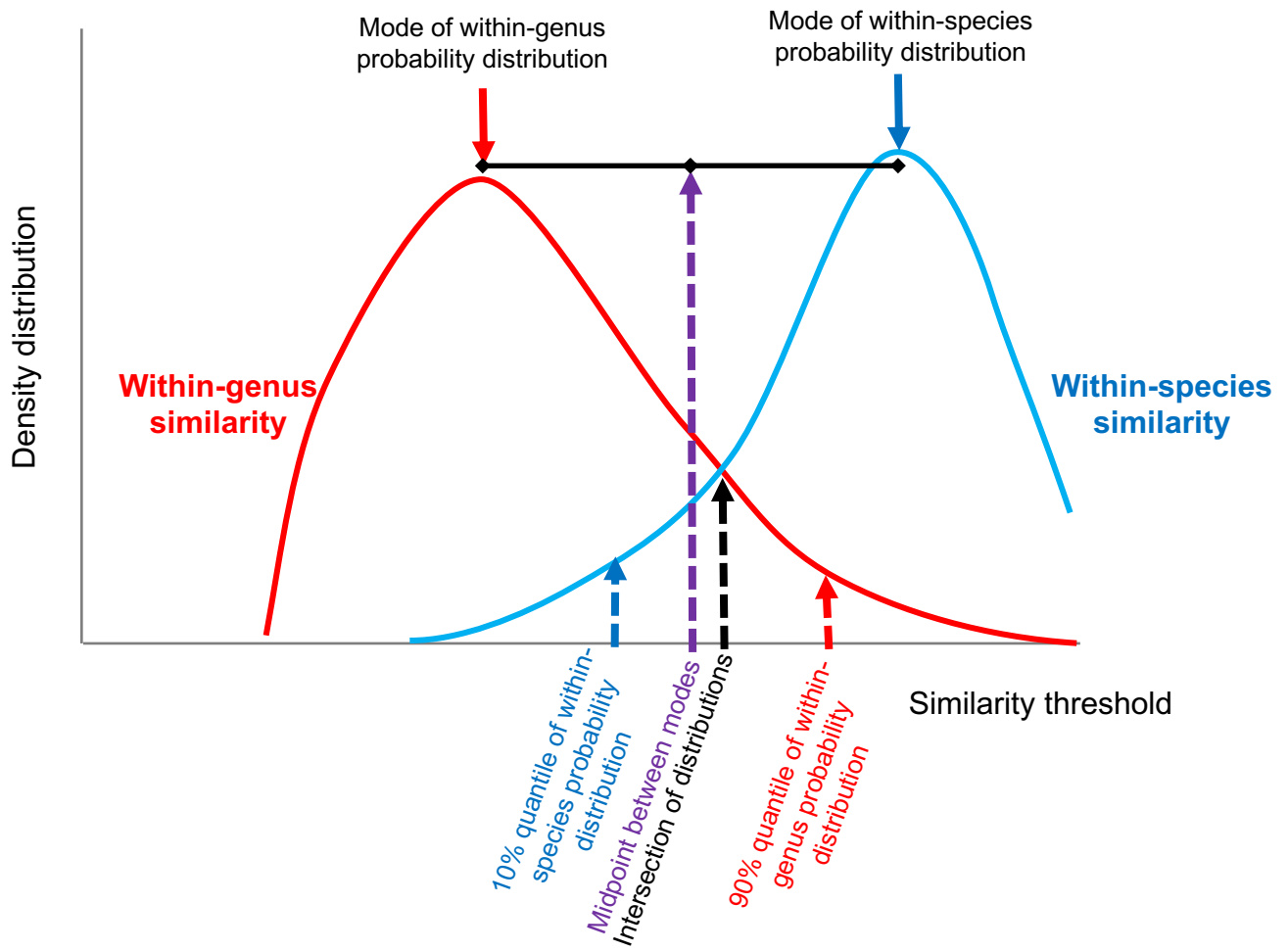
619

620 **Figure 5.** Over-splitting (blue) and over-merging (red) rates, as well as the summed error rate
621 (i.e. over-splitting rate + over-merging rate; violet), for the eight studied markers across a
622 range of clustering thresholds.

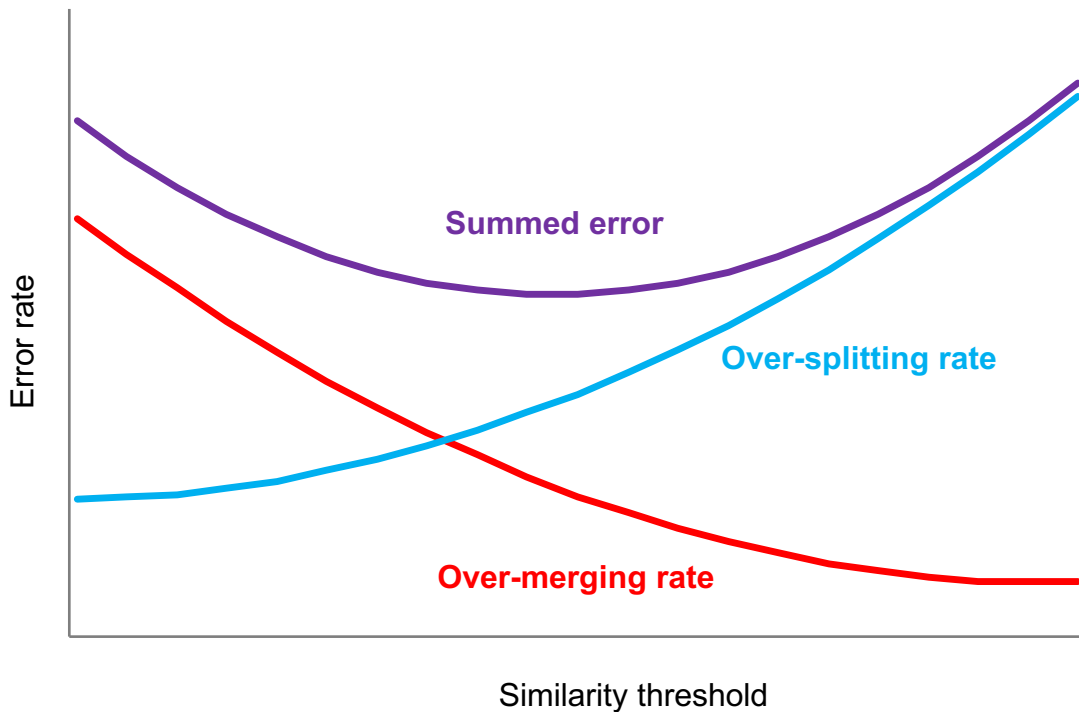
623

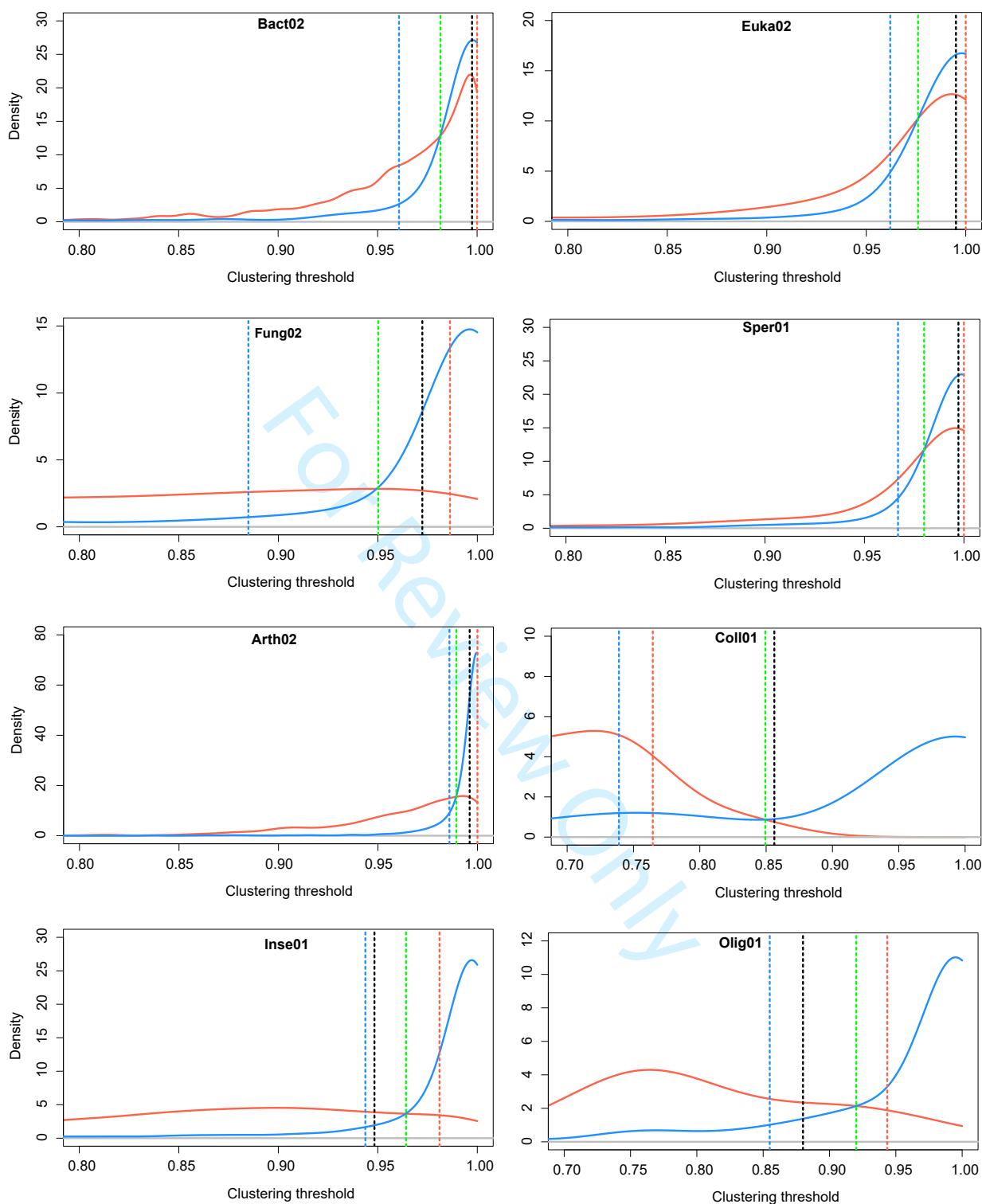
624

A) Approaches based on within-species and within-genus sequence similarities



B) Approaches based on over-splitting and over-merging rates



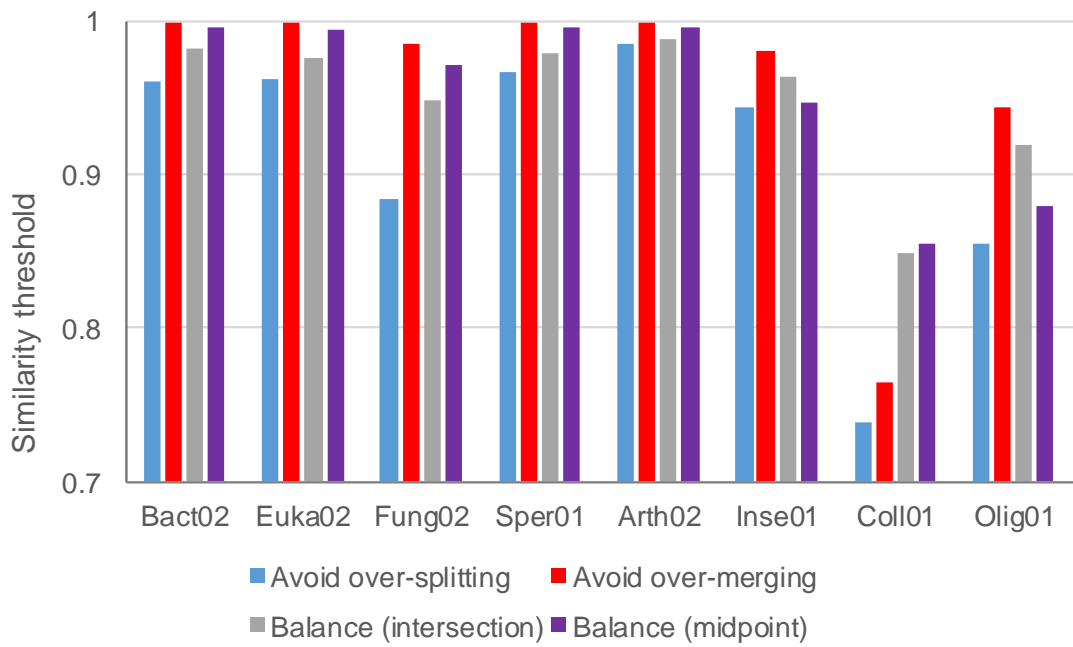


Probability distributions of

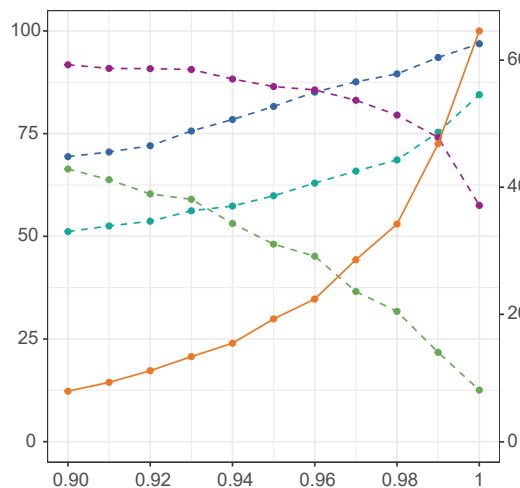
- sequence pairwise similarities within species
- sequence pairwise similarities within genera

Clustering thresholds

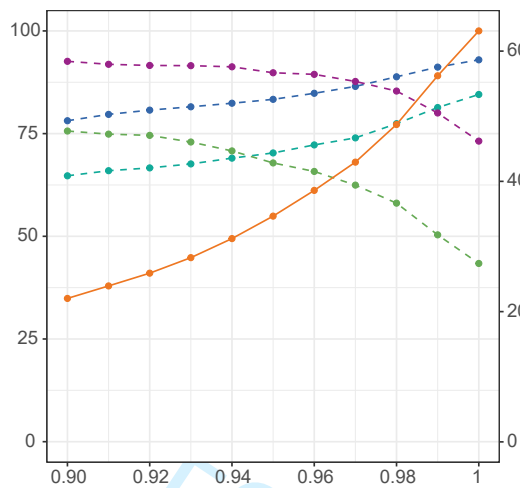
- Intersection
- Species 10% quantile
- Genus 90% quantile
- Midpoint between modes



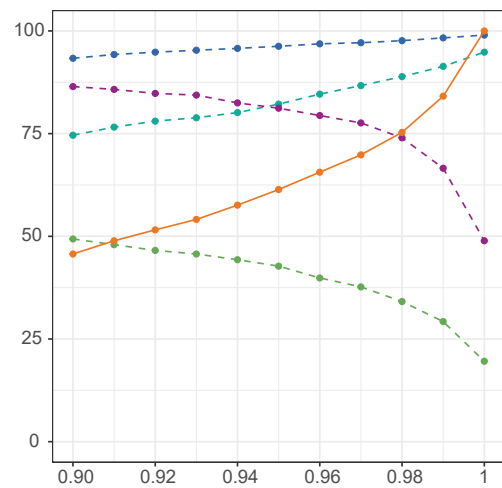
Bact02



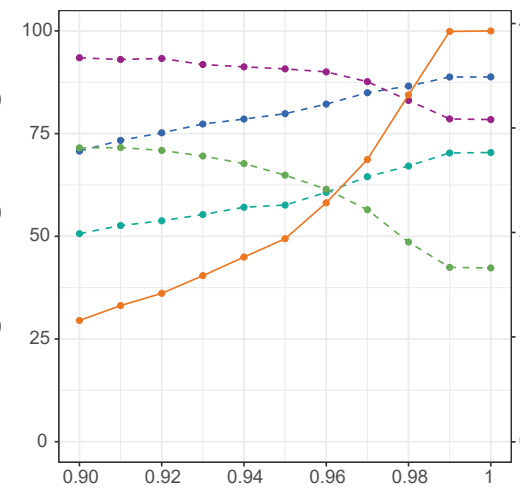
Euka02



Molecular Ecology Resources

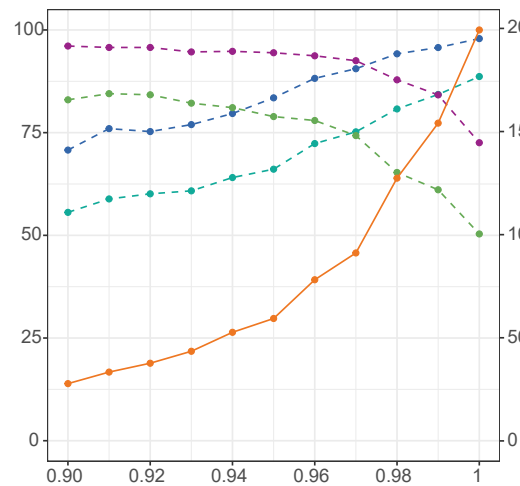


Sper01

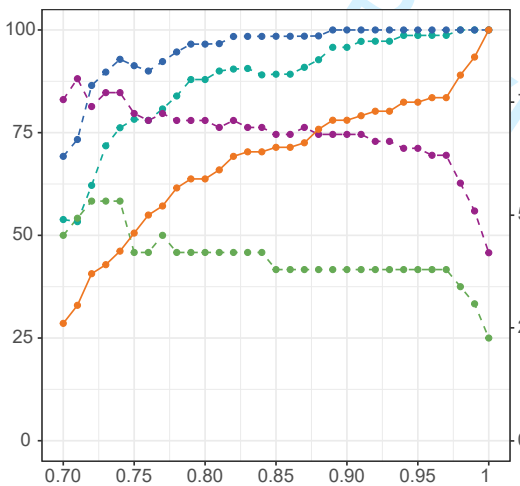


Page 30 of 31

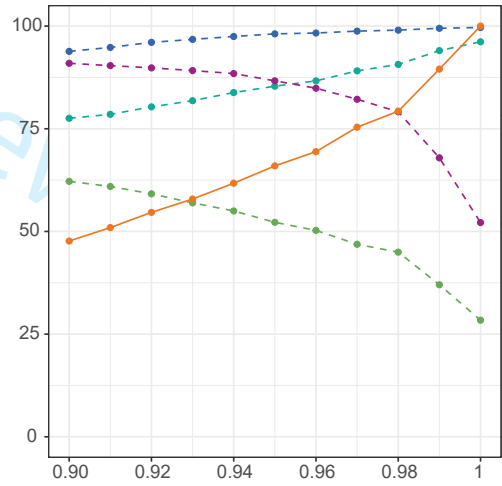
Arth02



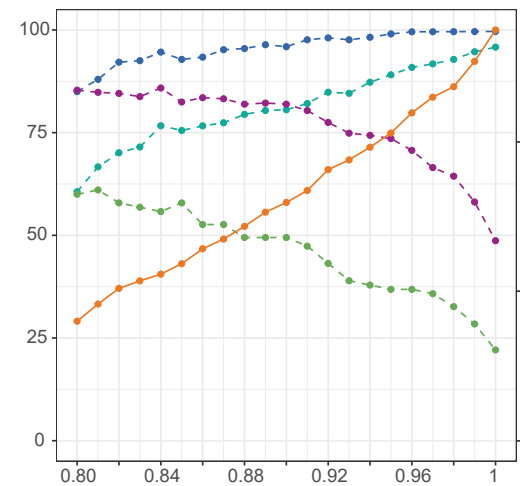
Coll01



Inse01



Olig01



—●— % of MOTUs containing a single genus

—●— % of genera gathered in a single MOTU among genera represented by several sequences

—●— Number of clusters

—●— % of MOTUs containing a single species

—●— % of species gathered in a single MOTU among species represented by several sequences

