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

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

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 Guerrieri[1], 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), each of which ideally representing a

19  homogeneous taxonomic entity, e.g. a species or a genus. However, the choice of the

20  clustering threshold is rarely justified, and its impact on MOTU over-splitting or over-

21  merging even less tested. Here, we evaluated clustering threshold values for several

22  metabarcoding markers under different criteria: limitation of MOTU over-merging, limitation

23  of MOTU over-splitting, and trade-off between over-merging and over-splitting. We extracted

24  sequences from a public database for nine markers, ranging from generalist markers targeting

25  Bacteria or Eukaryota, to more specific markers targeting a class or a subclass (e.g. Insecta,

26  Oligochaeta). Based on the distributions of pairwise sequence similarities within species and

27  within genera, and on the rates of over-splitting and over-merging across different clustering

28  thresholds, we were able to propose threshold values minimizing the risk of over-splitting,

29  that of over-merging, or offering a trade-off between the two risks. For generalist markers,

30  high similarity thresholds (0.96-0.99) are generally appropriate, while more specific markers

31  require lower values (0.85-0.96). These results do not support the use of a fixed clustering

32  threshold. Instead, we advocate a careful examination of the most appropriate threshold based

33  on the research objectives, the potential costs of over-splitting and over-merging, and the

34  features of the studied markers.

35

36

37 **Keywords**

38 metabarcoding marker; sequence variant; COI; MOTU over-splitting, MOTU over-merging;

39 alpha diversity

40

**Introduction**

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 the 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

52 al., 2022), f) the pipeline selected for bioinformatics analyses (Boyer et al., 2016, Calderón-

53 Sanou et al., 2020, Capo et al., 2021, Couton et al., 2021, Macher et al., 2021, Mächler et al.,

54 2021), and g) the statistical approach used to translate metabarcoding data into ecological

55 information (Paliy & Shankar 2016, Chen & Ficetola 2020). Each of these methodological

56 choices can heavily influence the reliability and interpretation of results (Alberdi et al., 2018,

57 Zinger et al., 2019), and there is thus a critical need for 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

64 including genuine intraspecific diversity of the selected markers and errors occurring during

65    the 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, when performing

73    clustering based on sequence similarity, two key parameters have to be determined *a priori*.

74    The first 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). The choice of this threshold

80    is 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 level

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 of sequence similarity ranging from 0.87 to 0.99 to analyze multiple

87    microbial communities, and obtained community structures highly coherent across thresholds.

88    Nevertheless, levels of alpha diversity can be heavily impacted by the threshold selection.

89   Ideally, the threshold used for clustering would depend on a trade-off between MOTU over-

90   splitting and MOTU over-merging. A growing number of markers are currently being used in

91   metabarcoding studies (Taberlet et al., 2018), with some allowing broad-scale biodiversity

92   assessment but having limited taxonomic resolution (e.g. 18S rDNA primers amplifying all

93   eukaryotes; Guardiola et al., 2015) and others being highly specific to one single class or even

94   family (e.g. Baamrane et al., 2012, Ficetola et al., 2021). Biodiversity surveys generally aim

95   to generate a set of MOTUs that are each associated with a unique taxon, all taxa being

96   ideally situated at the same level in the taxonomic tree, in order to facilitate comparisons. In

97   these conditions, optimal clustering thresholds probably differ strongly across markers. One

98   can for example expect high values for highly conserved markers, and lower values for

99   markers showing high variability (Kunin et al., 2010, Brown et al., 2015). However, there is

100  limited quantitative assessment of how optimal clustering thresholds vary across markers (but

101  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 nine

104  metabarcoding markers (Table 1), ranging from generalist markers (i.e. targeting Bacteria or

105  Eukaryota) to more specific markers (e.g. targeting Oligochaeta [earthworms], Insecta

106  [insects] or Collembola [springtails]), and amplifying fragments situated either in protein

107  coding (e.g. cytochrome c oxidase subunit 1 mitochondrial gene) or non-protein coding (e.g.

108  rDNA genes) genomic regions. We evaluated how clustering thresholds can change for each

109  marker and taxonomic group, depending on the criterion adopted to set the threshold. We

110  used two alternative strategies to identify thresholds, each time with different objectives in

111  mind. First, following a procedure similar to the one adopted in barcoding studies (Machida et

112  al., 2009; Meyer & Paulay 2005), we compared the distribution probabilities of sequence

113    similarities among different individuals of the same species and among different species of

114    the same genus to identify values: *i*) minimizing the risk that different sequences of the same

115    species are split in different MOTUs (i.e. risk of over-splitting); *ii*) minimizing the risk that

116    distinct but related species are clustered in the same MOTU (i.e. risk of over-merging); *iii*)

117    balancing the risk of over-splitting and over-merging (Figure 1A). Second, we calculated the

118    over-splitting and over-merging rates of the studied markers for a range of clustering

119    thresholds, to identify values that minimize the two error rates (Figure 1B). We expect that, if

120    researchers want to minimize over-splitting, they should select lower clustering thresholds

121    than if they want to minimize over-merging. Furthermore, we expect higher clustering

122    thresholds for generalist markers compared to markers targeting one class or more restricted

123    taxonomic groups, because of the lower taxonomic resolution and slower evolutionary rate of

124    the former.

125

126    **Methods**

127

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

129    We focused on a set of nine DNA metabarcoding markers (Bact02, Euka02, Fung02, Sper01,

130    Arth02, COI-BF1/BR2, Coll01, Inse01, Olig01) targeting different taxonomic groups and

131    different genomic regions (Table 1). Four of these markers can be considered as generalist,

132    i.e. targeting entire superkingdoms or kingdoms: Bact02 targeting Bacteria, Euka02 targeting

133    Eukaryota, Fung02 targeting Fungi, and Sper01 targeting Spermatophyta (vascular plants).

134    Two markers were intermediate (Arth02 and COI-BF1/BR2, both targeting arthropods, i.e.

135    the most species-rich phylum on Earth). Finally, three markers were more specific, i.e.

136    targeting groups from classes to subclasses: Coll01 targeting Collembola (springtails), Inse01

137   targeting Insecta, and Olig01 targeting Oligochaeta (earthworms). Eight of these markers are

138   situated in non-protein coding genes (Bact02, Arth02, Coll01, Inse01 and Olig01: 16 rDNA

139   gene; Euka02: 18S rDNA gene; Fung02: ITS1 nuclear rDNA gene; Sper01: P6 loop of the

140   intron of the chloroplastic *trnL* gene). The last marker, COI-BF1/BR2, is situated in the

141   cytochrome c oxidase subunit 1 (COI) mitochondrial gene (Table 1).

142        For each of these markers, a sequence database was built from EMBL release 140

143   (https://www.ebi.ac.uk/about/news/service-news/release-140-enas-assembledannotated-

144   sequences-now-available-0, also available from the authors) as follows. An *in silico* PCR was

145   first carried out by running the program *ecoPCR* (Ficetola et al., 2010) using the

146   corresponding primers (Table S1). Three mismatches per primer were allowed (-e option),

147   and amplicon length (without primers) was restricted (-l and -L options) to the expected

148   length interval (Table S1). The amplified sequences were further filtered by keeping only

149   those belonging to the target taxonomic group, showing a taxonomic assignment (i.e. taxid) at

150   the species and genus levels and having no ambiguous nucleotides. This allowed assembling a

151   working dataset, from which we extracted two sub-datasets. The "within-species" dataset was

152   built by keeping only species for which at least two sequences (identical or not) were

153   available; if >2 sequences were available for a given species, we randomly selected two

154   sequences for that species using the *obiselect* command of the OBITools. The "within-genus"

155   dataset was built by keeping only genera for which at least two sequences were available; if

156   >2 sequences were available for a given genus, we randomly selected two sequences for that

157   genus using the *obiselect* command. For some markers (Bact02, Euka02, Fung02, Inse01,

158   Sper01), the within-species dataset and sometimes the within-genus dataset still contained a

159   very large number of sequences (>10,000). To limit computation time for these markers, we

160   randomly selected a subset of 5000 different taxa, to reach a final number of sequences equal

161    to 10,000. An example of dataset preparation is provided in

162    Script1_Arth02_DatasetsPreparation.sh (Supplementary Material), and Table S2 summarizes

163    the number of sequences in the different datasets.

164

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

166    As a measure of sequence similarity, we computed the pairwise LCS (Longest Common

167    Subsequence) scores between pairs of sequences in the within-species and within-genus

168    datasets using the *sumatra* program (Mercier et al., 2013; see

169    Script2A_Arth02_PairwiseSimilarities_Sumatra.sh from the Supplementary Material).

170    Methodological comparisons showed that this algorithm provides an excellent balance

171    between performance and computation efficiency (Jackson et al., 2016, Kopylova et al., 2016,

172    Bhat et al., 2019). As *sumatra* provides pairwise scores for all possible pairs of sequences, the

173    similarity scores resulting from the within-species dataset were filtered in R (R Core Team

174    2020) to keep only those representing similarities between sequences of the same species.

175    Similarly, the scores resulting from the within-genus dataset were filtered to keep only those

176    representing similarities between different species of the same genus (see first part of

177    Script2B_Arth02_DensityPlots.Rmd from the Supplementary Material).

178

179    **Approach to identify clustering thresholds on the basis of within-species and within-**

180    **genus sequence similarities**

181    We first examined within-species and within-genus sequence similarities to evaluate four

182    different strategies (Figure 1A) and determine the similarity value that: *i*) avoids over-

183    splitting; *ii*) avoids over-merging; *iii*) finds a balance between over-splitting and over-

184    merging, with two distinct procedures based on the intersection (*iii*-a) or on modes (*iii*-b) of

185  the density probability distributions (see Script2B_Arth02_DensityPlots.Rmd from the

186  Supplementary Material). These strategies are analogous to those adopted in traditional

187  barcoding studies to set the limit between intra-specific and inter-specific diversity (Meyer &

188  Paulay 2005).

189  ### *i)*  **Avoid over-splitting**

190  In this case, the aim is to avoid distributing different sequences belonging to the same species

191  in different clusters, i.e. to limit the probability of generating additional spurious MOTUs. For

192  this purpose, we selected as clustering threshold the 10% quantile of the distribution of

193  similarities between sequences from the same species (within-species dataset). With this

194  approach, the sequences belonging to the same species according to EMBL are gathered in

195  the same cluster in 90% of the cases.

196  ### *ii)*  **Avoid over-merging**

197  In this case, the aim is to avoid gathering sequences attributed to different species of the same

198  genus in the same cluster, i.e. to limit the probability of merging related species in the same

199  MOTU. For this purpose, we selected as clustering threshold the 90% quantile of the

200  distribution of similarities between different species belonging to the same genus. With this

201  approach, the sequences attributed to different species belonging to the same genus are

202  assigned to different clusters in 90% of the cases.

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

204  In this case, the aim was to minimize both over-splitting and over-merging. We considered

205  two distinct approaches. First, we obtained the probability distribution of within-species and

206  within-genus sequence pairwise similarities using the *density* function from R, with biased

207  cross-validation (bw="bcv") as smoothing bandwidth selector and a Gaussian smoothing

208  kernel (kernel="gaussian"; Venables & Ripley 2002). We tested other possible smoothing

209   bandwidth selectors, but biased cross-validation was the approach best fitting the score

210   histograms for all markers and all datasets (Figures S1 to S9). The balance threshold *iii*-a was

211   then identified as the intersection between the probability distributions of the within-species

212   and within-genus similarities. As an alternative approach to balance over-merging and over-

213   splitting (*iii*-b), we calculated the midpoint between the modes of the within-species and

214   within-genus probability distributions.

215

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

217   For each marker, over-merging and over-splitting rates were evaluated at different clustering

218   thresholds using the within-species dataset described in the paragraph "Markers examined and

219   construction of sequences datasets". This dataset contains two sequences at random, identical

220   or not, for a number of species belonging to the taxonomic group of interest.

221       For each within-species dataset, clustering was performed using the *sumaclust*

222   program (Mercier et al., 2013, see Script3A_Arth02_Clustering.sh from the Supplementary

223   material) with the *-n* option (normalization by alignment length) based on the sequence

224   similarities first calculated using the *sumatra* program (see above; Mercier et al., 2013).

225   Threshold values (*-t* option) ranging from 0.90 to 1 at 0.01 steps were tested for all markers

226   except Coll01 and Olig01 for which wider ranges ([0.70 – 1] and [0.80 – 1], respectively)

227   were selected based on the within-genus and within-species sequence similarity probability

228   distributions determined previously (see Figure 2). Clustered datasets were then explored to

229   calculate five different variables at each clustering threshold (see

230   Script3B_Arth02_Oversplitting_Overmerging.Rmd from the Supplementary Material): 1) the

231   number of clusters; 2) the percentage of MOTUs containing one single species; 3) the

232   percentage of MOTUs containing one single genus; 4) the percentage of species gathered in

1

233    one single MOTU; 5) the percentage of genera gathered in one single MOTU among genera

234    represented by several sequences. Variables 2 and 3 are indicative of appropriate MOTU

235    merging of sequences at the species and genus levels, respectively, while variables 4 and 5 are

236    indicative of appropriate MOTU splitting at the species and genus levels, respectively.

237         These values were also used to calculate three measures of error. We defined the over-

238    merging rate as *1 - the percentage of MOTUs containing one single species*; and the over-

239    splitting rate as *1 - the percentage of species gathered in one single MOTU*. The summed

240    error rate was then calculated as the sum of the over-merging and over-splitting rates. For this

241    estimate, we assigned the same weight to over-splitting and over-merging.

242

243    **Results**

244

245    Our *in-silico* PCRs amplified between 101,955 (Arth02) and 3,202,507 (Bact02) sequences

246    per marker (Table S2). After data filtering, we retained between 510 (Coll01) and 707,874

247    (Bact02) sequences per marker. The within-species dataset comprised between 118 (Coll01)

248    and 10,000 (Bact02, Euka02, Fung02, Sper01, COI-BF1/BR2, Inse01) sequences, while the

249    within-genus dataset comprised between 74 (Coll01) and 10,000 (Euka02 and Sper01)

250    sequences per marker.

251

252    **Clustering thresholds determined from probability distributions of within-species and**

253    **within-genus sequence similarities**

254         The probability distributions of within-species and within-genus sequence similarities

255    showed very contrasting patterns between the generalist and the specific markers (Figure 2).

256    For Arth02 and most of the markers targeting broad taxonomic groups (Bact02, Euka02, and

257   Sper01), the distributions of within-species and within-genus similarities were rather similar,

258   both showing a mode at very high similarity values (Figure 2). Fung02 showed a slightly

259   different pattern, as the within-genus similarities had a very broad distribution. Conversely,

260   for COI-BF1/BR2 and the more specific markers (Coll01, Inse01, and Olig01), the

261   distributions of sequence similarities were very different, with two clearly distinct peaks.

262   Within-species similarities remained very high (mostly above 0.95), while within-genus

263   similarities generally showed lower values (mode around 0.88-0.90 for COI-BF1/BR2 and

264   Inse01, and below 0.80 for Olig01 and Coll01).

265       For all markers, criterion *i* (avoid over-splitting) yielded the lowest thresholds (Table

266   2), with very low values for Coll01 and Olig01. Conversely, criterion *ii* (avoid over-merging)

267   yielded extremely high values, except for Coll01. For all generalist markers and Arth02,

268   limiting over-merging would require setting clustering thresholds at 0.99 or higher. The same

269   objective would entail a slightly lower threshold for COI-BF1/BR2 and Inse01 (0.98) and

270   down to 0.94 for Olig01. For Coll01, criterion *ii* resulted in a very low threshold (0.77),

271   because many within-genus comparisons showed very low similarity values.

272       Criteria *iii*-a and *iii*-b searching a balance between over-merging and over-splitting

273   yielded somehow contrasting results across markers. For COI-BF1/BR2 and the three specific

274   markers (Coll01, Inse01, and Olig01), the within-genus and within-species similarities

275   showed clearly distinct peaks (Figure 2). As a consequence, the intersection between the two

276   curves could effectively represent the point minimizing both over-merging and over-splitting

277   (see discussion), and the midpoint between the modes also identified rather similar threshold

278   values. On the contrary, for the generalist markers and Arth02, the within-species and within-

279   genus similarities showed very high overlap and similar modes, and the density distributions

280   actually intersected at values lower than both modes. The midpoint between the modes

1

281 continued to identify threshold values intermediate between the peaks of within-species and

282 within-genus similarities.

283

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

285 For all markers, irrespective of the clustering threshold examined (values ≥ 0.70 for Coll01, ≥

286 0.80 for Olig01 and ≥ 0.90 for the other markers), the percentage of MOTUs containing one

287 single species was higher than 50%, and that of MOTUs containing one single genus was

288 higher or close to 70% (Figure 3). Overall, for the generalist and intermediate markers, these

289 two percentages showed a regular increase with the clustering threshold. For the specific

290 markers as well as Fung02 and COI-BF1/BR2, they reached values close to 100% for high

291 thresholds. Unsurprisingly, the two percentages tended to be lower for the generalist markers

292 than for the specific markers at a given threshold, indicating that the former are more sensitive

293 to over-merging. Fung02 was a notable exception, since about 87% and 97% of MOTUs

294 contained one single species and one single genus, respectively, at the 0.97 threshold, which

295 is frequently adopted as clustering threshold for fungal ITS sequences. These values were

296 comparable to those observed for COI-BF1/BR2 and the specific markers, for which > 85%

297 and > 98% of MOTUs contained one single species or one single genus, respectively, for

298 thresholds ≥ 0.95.

299 The percentages of species and genera gathered in one single MOTU decrease at a

300 similar rate with the clustering threshold, with generally a sharp drop at high thresholds (≥

301 0.98; Figure 3). However, the pattern of MOTU splitting was less characteristic of generalist

302 vs. specific markers. For some markers (Euka02, Sper01, Arth02, Inse01), the percentage of

303 species or genera gathered in a single MOTU remained higher or close to 50% up to high

304 thresholds (0.98). On the contrary, for Bact02, Fung02, COI-BF1/BR2, Coll01 and Olig01,

305     these percentages dropped quickly when the clustering threshold increased, indicating that

306     these markers are susceptible to over-splitting.

307     For all markers, the number of clusters generally increased regularly with the clustering

308     threshold up to 0.97-0.98 (Figure 3), followed by a sharp rise up to 1 (which was however less

309     obvious for Euka02 and Olig01). For example, for Bact02, the number of clusters more than

310     doubled between 0.97 (2862 clusters) and 1 (6461 clusters).

311     Our results showed clear patterns for over-merging and over-splitting rates, with over-

312     splitting quickly increasing and over-merging quickly decreasing at high clustering thresholds

313     (Figure 4). For several markers, the summed error showed a relatively clear minimum at

314     specific clustering thresholds (Figure 4): 0.96-0.99 for Bact02, 0.97-0.99 for Euka02 and

315     Arth02, 0.96-0.98 for Sper01, 0.93-0.96 for COI-BF1/BR2, and 0.94-0.97 for Inse01. The

316     minimum was much less evident for Fung02, Coll01 and Oligo01, these markers showing

317     relatively similar summed error rates over a broad range of clustering thresholds (Fung02:

318     0.91-0.98; Coll01: 0.89-0.97, with multiple minima; Oligo01: 0.84-0.96, with multiple

319     minima).

320

321     **DISCUSSION**

322

323     Sequence clustering approaches are routinely used for the identification of MOTUs in

324     metabarcoding studies, and they often resort to methods based on similarity values. Still,

325     selecting a clustering threshold for a given marker more than often relies on common

326     practices and rules of thumb rather than on proper scientific argument. By analyzing extensive

327     sequence data deposited in public databases for a range of generalist and specialist markers,

328     we showed that different thresholds can be selected depending on the marker and on the

1

329     criterion favored by researchers. All studied markers but one (COI-BF1/BR2) are situated in

330     non-protein coding genes (Table 1), and this has an influence on levels of sequence diversity.

331     More variability might be expected in protein-coding genes due to the redundancy of the

332     genetic code. Yet, for all markers including COI-BF1-BR2, the 10% quantile of the within-

333     species similarity probability distribution was almost always lower than the 0.97 clustering

334     threshold traditionally used in barcoding for markers targeting protein-coding genes like COI

335     (Hebert et al., 2003), or for microbial MOTU delimitation (Bálint et al., 2016). This indicates

336     indicating that some level of over-splitting can occur when using this threshold.

337         COI-BF1/BR2 is the only marker amplifying a fragment of a protein-coding gene, and

338     it would have been logical to observe singular patterns for this marker. However, this was not

339     the case, and COI-BF1/BR2, although designed to target arthropods (Elbrecht & Leese 2017)

340     like Arth02, actually showed a behavior very similar to the more specific Inse01 targeting

341     insects. The similarity between COI-BF1/BR2 and the more specific markers might be related

342     to their high resolution, which allows the successful distinction of closely related species even

343     on the basis of relatively short sequences (Elbrecht & Leese, 2017; Ficetola et al., 2021).

344     Furthermore, at 0.94, which is a suitable clustering threshold for COI-BF1/BR2, about 88% of

345     the MOTUs contain a single species, and about 88% of the species are gathered in a single

346     MOTU (Figure 3), indicating that MOTU richness at this threshold is a reasonably good

347     proxy for the number of species detected with this marker. This is corroborated by the number

348     of clusters observed at this threshold (5659), which is comparable to the expected number of

349     species (5000, Table S2) in the within-species dataset used to obtain Figure 3. Several COI

350     markers are routinely used in metabarcoding, and COI-BF1/BR2 shows a large overlap with

351     many of them (Elbrecht & Leese, 2017). We can thus expect that optimal clustering

352    thresholds for COI-BF1/BR2 can also be rightfully applied to markers targeting a slightly

353    different COI region.

354        Although the within-genus similarity values were generally lower than the within-

355    species similarities for all the markers, the overlap between the two distributions was

356    dependent on the generalist vs. specific nature of the marker. For some specific markers (e.g.

357    Coll01 and Olig01), distinct peaks were visible for the two similarity metrics (Figure 2).

358    Within-species similarities generally were >0.90, while within-genus values were <0.80. Such

359    a pattern is expected for markers with an excellent taxonomic resolution and designed to

360    identify taxa at the species level. Conversely for the generalist markers, within-species and

361    within-genus similarity probability distributions largely overlapped and the differences

362    between the peaks were minimal. Nevertheless, even for these markers, the density of within-

363    species similarity distribution was consistently higher than that of within-genus similarity

364    distribution at high similarity values. This suggests that the probability of observing the

365    corresponding sequence similarity is higher within species than within genera. In other words,

366    at high sequence similarities, a MOTU is more likely to represent a species than a genus. This

367    result is confirmed by the fact that the percentage of MOTUs containing a single species is

368    always higher than 50%, whatever the clustering threshold or the marker considered (Figure

369    3).

370        The sequences used as a primary source of information in this study were downloaded

371    from the EMBL public database, therefore our results are probably highly dependent on the

372    quality of the data deposited. Even though broad-scale analyses suggest that sequence data

373    from public database are generally reliable (Leray et al., 2019), errors in the sequence itself

374    (e.g. wrong nucleotide, or more complex errors like insertions, deletions, inversions,

375    duplications or pseudogene sequences) and taxonomic mislabeling can occur. Organisms that

376    are difficult to identify based on morphology are particular susceptible to wrong taxonomic

377    information (Bridge et al., 2003, Bidartondo 2008, Valkiūnas et al., 2008, Mioduchowska et

378    al., 2018). While errors in the sequence will affect within-species sequence similarity

379    negatively, the effect of taxonomic mislabeling is more diffuse. For example, in a group like

380    springtails where species delimitation is tricky (Porco et al., 2012), the existence of cryptic

381    species will decrease within-species sequence similarity while increasing over-splitting rates.

382    In a group like Bacteria, type strains are sometimes entered at the species level in the NCBI

383    (EMBL) taxonomy (Federhen 2015), leading to an inflation of within-genus similarity and

384    over-merging rates. In any case, database errors will make within-species and within-genus

385    similarities distributions more difficult to distinguish and clustering thresholds trickier to

386    identify, thus the over-splitting or over-merging rates reported here could be artificially higher

387    than in reality.

388        In this work, we came up with a global measure of the error associated with a given

389    clustering threshold, that we called the "summed error". We calculated it by summing over-

390    splitting and over-merging rates, assuming both have the same cost for biodiversity studies.

391    However, it is possible to assign a differential weight to over-splitting and over-merging. For

392    instance, if the aim is to reach conservative estimated of alpha diversity (i.e. avoid over-

393    splitting), more weight can be assigned to over-splitting rate. Conversely, if the aim is to tease

394    apart closely related species, that differ in their sensitivity to environmental stressors or in

395    threat levels, one may prefer to avoid over-merging, particularly when extensive reference

396    databases are available (Roy et al., 2019, Lopes et al., 2021).

397        For most of the markers we examined, the summed error approach provided relatively

398    clear results and identified a range of threshold values that minimized the summed error. For

399    instance, for Euka02, the summed error was relatively low at thresholds between 0.96 and

400 0.99 (Figure 4), indicating a good trade-off between over-merging and over-splitting.

401 Interestingly, this range of values was also highlighted by the analysis of probability

402 distributions (Table 2). Indeed, 0.96 is the threshold minimizing over-splitting for Euka02

403 while 0.99 is the balance (midpoint) threshold. The consistency of values obtained with very

404 different approaches supports the robustness of our conclusions.

405  However, for a few markers, the threshold values minimizing summed error yielded

406 somewhat less clear patterns. For Fung02, the summed error rate was rather constant (36-

407 37%) at all the thresholds between 0.91 and 0.98, while it quickly increased for higher

408 clustering thresholds. For Coll01 and Oligo01, the summed error rate showed multiple

409 minima, some of which at very low clustering thresholds (Figure 4). In principle, increasing

410 the threshold value should determine a monotone decrease of over-merging, and a monotone

411 increase of over-splitting (Figure 1B). However, at low similarity values this was not always

412 the case (Figure 4). This probably occurs because, for these markers a large proportion of

413 sequences have pairwise similarities of 0.80-0.85 (Figure 2), and this might affect the

414 identification of clusters, with some sequences clustering together e.g. at 0.85 but not at 0.86

415 similarity values. We also note that these similarity values match the ones corresponding to

416 the intersection between the within-genus and within-species similarities for these markers

417 (Table 2). It is also possible that, at this level of sequence similarity, there is strong

418 uncertainty between MOTUs representing different hierarchical levels of taxonomy.

419  Our results provide quantitative data that can help researchers set their optimal

420 clustering thresholds and understand the consequences of choosing low or high threshold

421 values. If a clear minimum exists for the summed error rate, it probably represents an

422 excellent trade-off between over-merging and over-splitting. In this sense, a threshold value

423 ranging from 0.97 to 0.99 is probably appropriate for both Bact02 and Euka02, while Arth02

424    should accommodate a slightly higher range (0.98-0.99) and a threshold of 0.97 seems to be

425    more suitable for Sper01. For Inse01 and COI-BF1/BR2, lower threshold values (0.94-0.97

426    and 0.93-0.96, respectively) are more judicious. All these values match with those obtained

427    on the basis of within-species and within-genus similarities (Table 2). However, for Coll01,

428    Oligo01 and Fung02, the summed error rate does not provide clear indications, and within-

429    species and within-genus similarity distributions (e.g. midpoint between modes) might be

430    more informative to set the clustering threshold (Figure 2 and Table 2).

431        The selection of clustering thresholds can have strong effect in the estimates of

432    MOTUs richness (Figure 3), still it is important to remember that it often does not have a

433    tremendous effect on the ecological message conveyed by metabarcoding data. For instance,

434    Clare et al. (2016) examined different clustering thresholds to analyze dietary overlap

435    between skinks and shrews in Mauritius. Although high clustering thresholds yielded a larger

436    number of MOTUs, ecological conclusions remained rather consistent overall. Therefore,

437    provided that appropriate parameters are considered (e.g. alpha diversity measured using

438    Hill's numbers with $q > 0$ instead of richness, beta diversity estimates), the interpretation of

439    data can be relatively robust (Clare et al., 2016, Roy et al., 2019, Calderón-Sanou et al., 2020,

440    Mächler et al., 2021). Nevertheless, we discourage the blind application of one single

441    clustering threshold like the classical 0.97, as it can have very different meaning across

442    markers, and can inflate MOTU richness for fast-evolving markers. Instead, we advocate the

443    ad-hoc definition of the most appropriate thresholds, depending on the research aims, the

444    potential costs of over-splitting and over-merging, and the features of the studied markers.

445

451

452    **References**

453    Alberdi, A., Aizpurua, O., Gilbert, M. T. P., & Bohmann, K. (2018). Scrutinizing key steps for
454        reliable metabarcoding of environmental samples. *Methods in Ecology and Evolution,*
455        *9,* 134-147.
456    Antich, A., Palacin, C., Wangensteen, O. S., & Turon, X. (2021). To denoise or to cluster, that
457        is not the question: optimizing pipelines for COI metabarcoding and
458        metaphylogeography. *BMC Bioinformatics, 22,* 177.
459    Baamrane, M. A. A., Shehzad, W., Ouhammou, A., Abbad, A., Naimi, M., Coissac, E., …
460        Znari, M. (2012). Assessment of the food habits of the Moroccan dorcas gazelle in
461        M'Sabih Talaa, West Central Morocco, using the *trnL* approach. *PLoS ONE, 7,* e35643.
462    Bálint, M., Bahram, M., Eren, A. M., Faust, K., Fuhrman, J. A., Lindahl, B., … Tedersoo, L.
463        (2016). Millions of reads, thousands of taxa: microbial community structure and
464        associations analyzed via marker genes. *FEMS Microbiology Reviews, 40,* 686-700.
465    Bhat, A. H., Prabhu, P., & Balakrishnan, K. (2019). A critical analysis of state-of-the-art
466        metagenomics OTU clustering algorithms. *Journal of Biosciences, 44,* 9.
467    Bidartondo, M. I. (2008). Preserving accuracy in GenBank. *Science, 319,* 1616.
468    Bienert, F., De Danieli, S., Miquel, C., Coissac, E., Poillot, C., Brun, J. J., & Taberlet, P. (2012).
469        Tracking earthworm communities from soil DNA. *Molecular Ecology, 21,* 2017-2030.
470    Bohmann, K., Elbrecht,V., Carøe, C., Bista, L., Leese, F., Bunce, M., Yu, D. W., … Creer, S.
471        (in press). Strategies for sample labelling and library preparation in DNA
472        metabarcoding studies. *Molecular Ecology Resources.*
473    Botnen, S. S., Davey, M. L., Halvorsen, R., & Kauserud, H. (2018). Sequence clustering
474        threshold has little effect on the recovery of microbial community structure. *Molecular*
475        *Ecology Resources, 18,* 1064-1076.
476    Boyer, F., Mercier, C., Bonin, A., Le Bras, Y., Taberlet, P., & Coissac, E. (2016). OBITOOLS:
477        a Unix-inspired software package for DNA metabarcoding. *Molecular Ecology*
478        *Resources, 16,* 176-182.
479    Bridge, P. D., Roberts, P. J., Spooner, B. M., & Panchal, G. (2003). On the unreliability of
480        published DNA sequences. *New Phytologist, 160,* 43-48.
481    Brown, E. A., Chain, F. J. J., Crease, T. J., MacIsaac, H. J., & Cristescu, M. E. (2015).
482        Divergence thresholds and divergent biodiversity estimates: can metabarcoding reliably
483        describe zooplankton communities? *Ecology and Evolution, 5,* 2234-2251.

484   Calderón-Sanou, I., Münkemüller, T., Boyer, F., Zinger, L., & Thuiller, W. (2020). From
485         environmental DNA sequences to ecological conclusions: How strong is the influence
486         of methodological choices? *Journal of Biogeography, 47*, 193–206.
487   Capo, E., Giguet-Covex, C., Rouillard, A., Nota, K., Heintzman, P., Vuillemin, A. …Parducci,
488         L. (2021). Lake sedimentary DNA research on past terrestrial and aquatic biodiversity:
489         Overview and recommendations. *Quaternary 4,* 6.
490   Chen, W., & Ficetola, G. F. (2020). Statistical and numerical methods for Sedimentary-ancient-
491         DNA-based study on past biodiversity and ecosystem functioning. *Environmental DNA,*
492         *2*, 115–129.
493   Clare, E. L., Chain, F. J. J., Littlefair, J. E., & Cristescu, M. E. (2016). The effects of parameter
494         choice on defining molecular operational taxonomic units and resulting ecological
495         analyses of metabarcoding data. *Genome, 59*, 981-990.
496   Couton, M., Baud, A., Daguin-Thiébaut, C., Corre, E., Comtet, T., & Viard, F. (2021). High-
497         throughput sequencing on preservative ethanol is effective at jointly examining
498         infraspecific and taxonomic diversity, although bioinformatics pipelines do not perform
499         equally. *Ecology and Evolution, 11*, 5533-5546.
500   Dickie, I. A., Boyer, S., Buckley, H. L., Duncan, R. P., Gardner, P. P., Hogg, I. D., … Weaver,
501         L. (2018). Towards robust and repeatable sampling methods in eDNA-based studies.
502         *Molecular Ecology Resources, 18*, 940-952.
503   Eichmiller, J. J., Miller L. M., & Sorensen, P.W. (2016). Optimizing techniques to capture and
504         extract environmental DNA for detection and quantification of fish. *Molecular Ecology*
505         *Resources, 16*, 56-68.
506   Elbrecht, V., & Leese, F. (2017). Validation and development of COI metabarcoding primers
507         for freshwater macroinvertebrate bioassessment. *Frontiers in Environmental Science,*
508         *5*, 11.
509   Elbrecht, V., Taberlet, P., Dejean, T., Valentini, A., Usseglio-Polatera, P., Beisel, J. N., …
510         Leese, F. (2016). Testing the potential of a ribosomal 16S marker for DNA
511         metabarcoding of insects. *PeerJ, 4*, 12.
512   Epp, L. S., Boessenkool, S., Bellemain, E. P., Haile, J., Esposito, A., Riaz, T., … Brochmann,
513         C. (2012). New environmental metabarcodes for analysing soil DNA: potential for
514         studying past and present ecosystems. *Molecular Ecology, 21,* 1821-1833.
515   Fahner, N. A., Shokralla, S., Baird, D. J., & Hajibabaei, M. (2016). Large-scale monitoring of
516         plants through environmental DNA metabarcoding of soil: Recovery, resolution, and
517         annotation of four DNA markers. *PLoS ONE, 11*, e0157505.
518   Federhen, S. (2015). Type material in the NCBI Taxonomy Database. *Nucleic Acids Research,*
519         *43*, D1086-D1098.
520   Ficetola, G. F., Boyer, F, Valentini, A. Bonin, Meyer, A., Dejean, T., … Taberlet, P. (2021).
521         Comparison of markers for the monitoring of freshwater benthic biodiversity through
522         DNA metabarcoding. *Molecular Ecology, 30*, 3189–3202.
523   Ficetola, G. F., Coissac, E., Zundel, S., Riaz, T., Shehzad, W., Bessière, J., … Pompanon, F.
524         (2010). An *in silico* approach for the evaluation of DNA barcodes. *BMC Genomics, 11*,
525         434.
526   Froslev, T. G., Kjoller, R., Bruun, H. H., Ejrnaes, R., Brunbjerg, A. K., Pietroni, C., & Hansen,
527         A. J. (2017). Algorithm for post-clustering curation of DNA amplicon data yields
528         reliable biodiversity estimates. *Nature Communications, 8*, 11.
529   Guardiola, M., Uriz, M. J., Taberlet, P., Coissac, E., Wangensteen, O. S., & Turon, X. (2015).
530         Deep-sea, deep-sequencing: metabarcoding extracellular DNA from sediments of
531         marine canyons. *PLoS ONE, 10*, e0139633.

532    Guerrieri, A., Bonin, A., Münkemüller, T., Gielly, L., Thuiller, W., & Ficetola, G. F. (2021).
533        Effects of soil preservation for biodiversity monitoring using environmental DNA.
534        *Molecular Ecology, 30*, 3313-3325.
535    Hebert, P. D. N., Ratnasingham, S., & deWaard, J. R. (2003). Barcoding animal life:
536        cytochrome c oxidase subunit 1 divergences among closely related species. *Proceedings*
537        *of the Royal Society B-Biological Sciences, 270,* S96-S99.
538    Jackson, M. A., Bell, J. T., Spector, T. D., & Steves, C. J. (2016). A heritability-based
539        comparison of methods used to cluster 16S rRNA gene sequences into operational
540        taxonomic units. *PeerJ, 4*, 19.
541    Janssen, P., Bec, S., Fuhr, M., Taberlet, P., Brun, J.-J., & Bouget, C. (2018). Present conditions
542        may mediate the legacy effect of past land-use changes on species richness and
543        composition of above- and below-ground assemblages. *Journal of Ecology, 106*, 306-
544        318.
545    Kopylova, E., Navas-Molina, J. A., Mercier, C., Xu, Z. Z., Mahe, F., He, Y., … Knight, R.
546        (2016). Open-source sequence clustering methods improve the state of the art.
547        *mSystems, 1*, 16.
548    Kunin, V., Engelbrektson, A., Ochman, H., & Hugenholtz, P. (2010). Wrinkles in the rare
549        biosphere: pyrosequencing errors can lead to artificial inflation of diversity estimates.
550        *Environmental Microbiology, 12*, 118-123.
551    Lear, G., Dickie, I., Banks, J., Boyer, S., Buckley, H. L., Buckley, T. R., … Holdaway, R.
552        (2018). Methods for the extraction, storage, amplification and sequencing of DNA from
553        environmental samples. *New Zealand Journal of Ecology, 42*, 10.
554    Leray, M., Knowlton, N., Ho, S. L., Nguyen, B. N., & Machida, R. J. (2019). GenBank is a
555        reliable resource for 21st century biodiversity research. *Proceedings of the National*
556        *Academy of Sciences of the United States of America, 116*, 22651-22656.
557    Lopes, C. M., Baêta, D., Valentini, A., Lyra, M. L., Sabbag, A. F., Gasparini, J. L., … Zamudio,
558        R. K. (2021). Lost and found: Frogs in a biodiversity hotspot rediscovered with
559        environmental DNA. *Molecular Ecology, 30*, 3289-3298.
560    Macher, T.-H., Beermann, A. J., & Leese, F. (2021). TaxonTableTools: A comprehensive,
561        platform-independent graphical user interface software to explore and visualise DNA
562        metabarcoding data. *Molecular Ecology Resources, 21*, 1705-1714.
563    Machida, R. J., Hashiguchi, Y., Nishida, M., & Nishida, S. (2009). Zooplankton diversity
564        analysis through single-gene sequencing of a community sample. *BMC Genomics, 10*,
565        438.
566    Mächler, E., Walser, J.-C., & Altermatt, F. (2021). Decision-making and best practices for
567        taxonomy-free environmental DNA metabarcoding in biomonitoring using Hill
568        numbers. *Molecular Ecology, 30*, 3326-3339.
569    Mercier, C., Boyer, F., Bonin, A., & Coissac, E. (2013). SUMATRA and SUMACLUST: fast
570        and exact comparison and clustering of sequences. *Programs and Abstracts of the*
571        *SeqBio 2013 Workshop*, 27-29.
572    Meyer, C. P., & Paulay, G. (2005). DNA barcoding: Error rates based on comprehensive
573        sampling. *PLoS Biology, 3*, 2229-2238.
574    Mioduchowska, M., Czyz, M. J., Goldyn, B., Kur, J., & Sell, J. (2018). Instances of erroneous
575        DNA barcoding of metazoan invertebrates: Are universal *cox1* gene primers too
576        "universal" ? *PLoS ONE, 13*, e0199609.
577    Nichols, R. V., Vollmers, C., Newsom, L. A., Wang, Y., Heintzman, P. D., Leighton, M., …
578        Shapiro, B. (2018). Minimizing polymerase biases in metabarcoding. *Molecular*
579        *Ecology Resources, 18*, 927-939.

580     Paliy, O., & Shankar, V. (2016). Application of multivariate statistical techniques in microbial
581         ecology. *Molecular Ecology, 25*, 1032-1057.
582     Porco, D., Bedos, A., Penelope, G., Janion, C., Skarżyński, D., Stevens, M. I., … Deharveng,
583         L. (2012). Challenging species delimitation in Collembola: cryptic diversity among
584         common springtails unveiled by DNA barcoding. *Invertebrate Systematics, 26*, 470-
585         477.
586     R Core Team. (2020). R: A language and environment for statistical computing. R Foundation
587         for Statistical Computing, Vienna.
588     Roy, J., Mazel, F., Sosa-Hernández, M. A., Dueñas, J. F., Hempel, S., Zinger, L., & Rillig, M.
589         C. (2019). The relative importance of ecological drivers of arbuscular mycorrhizal
590         fungal distribution varies with taxon phylogenetic resolution. *New Phytologist, 224,*
591         936-948.
592     Schloss, P. D. (2021). Amplicon sequence variants artificially split bacterial genomes into
593         separate clusters. *mSphere, 6*, e00191-00121.
594     Taberlet, P., Bonin, A., Zinger, L., & Coissac, E. (2018). Environmental DNA for biodiversity
595         research and monitoring. Oxford University Press, Oxford.
596     Taberlet, P., Coissac, E., Pompanon, F., Gielly, L., Miquel, C., Valentini, A., … Willerslev, E.
597         (2007). Power and limitations of the chloroplast *trnL* (UAA) intron for plant DNA
598         barcoding. *Nucleic Acids Research, 35*, e14.
599     Taberlet, P., Prud'homme, S. M., Campione, E., Roy, J., Miquel, C., Shehzad, W., … Coissac,
600         E. (2012). Soil sampling and isolation of extracellular DNA from large amount of
601         starting material suitable for metabarcoding studies. *Molecular Ecology, 21*, 1816-1820.
602     Tatangelo, V., Franzetti, A., Gandolfi, I., Bestetti, G., & Ambrosini, R. (2014). Effect of
603         preservation method on the assessment of bacterial community structure in soil and
604         water samples. *FEMS Microbiology Letters, 356*, 32-38.
605     Valkiūnas, G., Atkinson, C. T., Bensch, S., Sehgal, R. N., & Ricklefs, R. E. (2008). Parasite
606         misidentifications in GenBank: how to minimize their number? *Trends in Parasitology,*
607         *24*, 247-248.
608     Venables, W. N., & Ripley, B. D. (2002). Modern applied statistics with S. Fourth Edition.
609         Springer, New York.
610     Wei, Z.-G., Zhang, X.-D., Cao, M., Liu, F., Qian, Y., & Zhang, S.-W. (2021). Comparison of
611         methods for picking the operational taxonomic units from amplicon sequences.
612         *Frontiers in Microbiology, 12*, 644012.
613     Zinger, L., Bonin, A., Alsos, I., Bálint, M., Bik, H., Boyer, F., … Taberlet, P. (2019). DNA
614         metabarcoding - need for robust experimental designs to draw sound ecological
615         conclusions. *Molecular Ecology, 28*, 1857-1862.
616     Zinger, L., Chave, J., Coissac, E., Iribar, A., Louisanna, E., Manzi, S., … Taberlet, P. (2016).
617         Extracellular DNA extraction is a fast, cheap and reliable alternative for multi-taxa
618         surveys based on soil DNA. *Soil Biology & Biochemistry, 96*, 16-19.

619

620

621     **Data Accessibility**

622     Raw data obtained from EMBL r140 (*ecopcr* files) and example scripts to prepare the datasets

623     and perform the analyses are available on Dryad: https://doi.org/10.5061/dryad.crjdfn353.

624

**Authors Contribution**

625

626  All authors conceived the idea for the manuscript, AB and GFF designed the study, AB

627  performed the analyses, AB and GFF generated the figures and drafted the manuscript, and all

628  authors contributed with discussions and edits.

629 **Table 1. Characteristics of the nine studied markers.**

| Marker | Target gene | Target group | Taxonomic level | Taxonomic resolution * | | | | Reference(s) |
|---|---|---|---|---|---|---|---|---|
| | | | | Species level | Genus level | Family level | Order level | |
| Bact02 | V4 region of the 16S rDNA gene | Bacteria | Superkingdom | 19.6% | 55.7% | 55.1% | 60.2% | Taberlet et al., (2018) |
| Euka02 | V7 region of the 18S rDNA gene | Eukaryota | Superkingdom | 47.0% | 59.5% | 68.3% | 67.1% | Guardiola et al., (2015) |
| Fung02 | ITS1 nuclear rDNA gene | Fungi | Kingdom | 72.5% | 90.2% | 87.7% | 85.5% | Epp et al., (2012), Taberlet et al., (2018) |
| Sper01 | P6 loop of the intron of the chloroplastic *trnL* gene | Spermatophyta | Clade < kingdom | 21.5% | 36.9% | 77.4% | 89.6% | Taberlet et al., (2007) |
| Arth02 | 16S mitochondrial rDNA gene | Arthropoda | Phylum | 68.6% | 89.6% | 97.5% | 100.0% | Taberlet et al., (2018) |
| COI-BF1/BR2 | Cytochrome c oxidase subunit 1 mitochondrial gene | Arthropoda | Phylum | 85.6% | 97.0% | 95.1% | 93.5% | Elbrecht & Leese (2017) |
| Coll01 | 16S mitochondrial rDNA gene | Collembola | Class | 80.5% | 87.2% | 75.0% | NA | Janssen et al., (2018) |
| Inse01 | 16S mitochondrial rDNA gene | Insecta | Class | 87.8% | 96.8% | 95.4% | 79.3% | Taberlet et al., (2018) |
| Olig01 | 16S mitochondrial rDNA gene | Oligochaeta | Subclass | 89.3% | 95.7% | 100.0% | 100.0% | Bienert et al., (2012), Taberlet et al., (2018) |

630 *Percentage of discriminated taxa among taxa amplified *in silico*, as calculated by the *ecotaxspecificity* program from the OBITools.
631 Reported from Taberlet et al., (2018) for all markers, except for COI-BF1/BR2 for which these values were determined using the
632 sequences amplified *in silico* from EMBL r140.

633

634

635 **Table 2. Values of the different thresholds estimated for the nine studied markers on the basis of within-species and within-genus**
636 **sequence similarities.**

| Target | Bact02 | Euka02 | Fung02 | Sper01 | Arth02 | COI-BF1/BR2 | Coll01 | Inse01 | Olig01 |
|---|---|---|---|---|---|---|---|---|---|
| Criterion *i*: Avoid over-splitting (10% quantile of within-species probability distribution) | 0.961 | 0.962 | 0.885 | 0.967 | 0.986 | 0.937 | 0.739 | 0.944 | 0.855 |
| Criterion *ii*: Avoid over-merging (90% quantile of within-genus probability distribution) | 1.000 | 1.000 | 0.986 | 1.000 | 1.000 | 0.975 | 0.765 | 0.981 | 0.944 |
| Criterion *iii*-a: Balance-a (intersection of within-species and within-genus probability distributions) | 0.982 | 0.976 | 0.949 | 0.980 | 0.989 | 0,955 | 0.849 | 0.964 | 0.920 |
| Criterion *iii*-b: Balance-b (midpoint between modes) | 0.997 | 0.995 | 0.972 | 0.997 | 0.996 | 0.936 | 0.856 | 0.948 | 0.880 |

637

638

**Figure captions**


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


**Figure 2. Density probability distributions of sequence pairwise similarities within**
**species (blue lines) and within genera (red lines) for the nine studied markers.** For each
marker, vertical dotted lines represent the 10% quantile of the within-species probability
distribution (blue; threshold limiting over-splitting) and the 90% quantile of the within-genus
probability distribution (red; threshold limiting over-merging). Vertical full lines represent the
intersection of the within-species and within-genus probability distributions (yellow, balance-
a) and the midpoint between modes (grey, balance-b)

663    **Figure 3. Evolution of over-splitting and over-merging rates for a range of clustering**

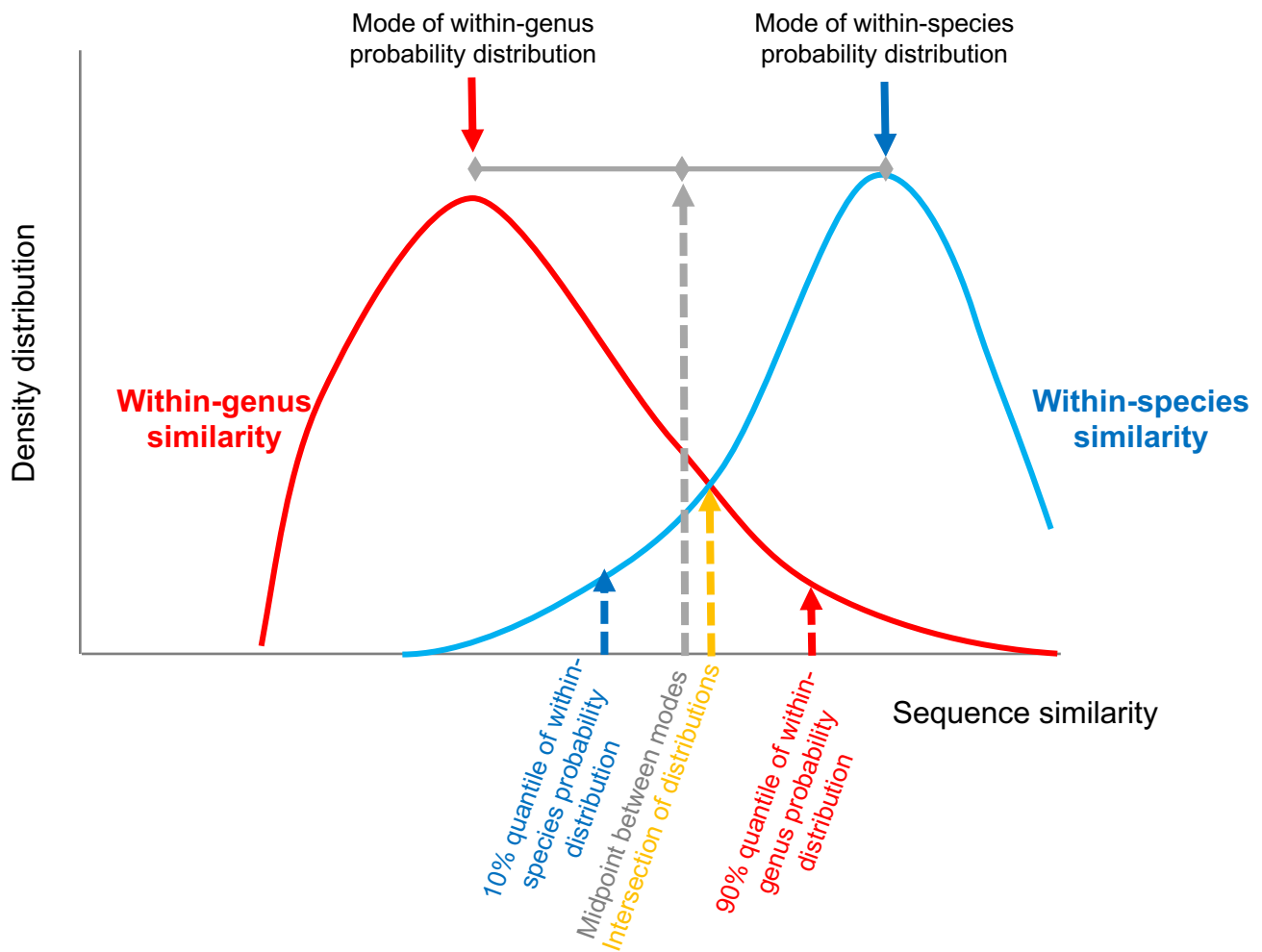664    **thresholds, for the nine studied markers.**

665

666    **Figure 4. Over-splitting (blue) and over-merging (red) rates, as well as the summed**

667    **error rate (i.e. over-splitting rate + over-merging rate; violet), for the nine studied**

668    **markers across a range of clustering thresholds.** Horizontal grey arrows indicate the range
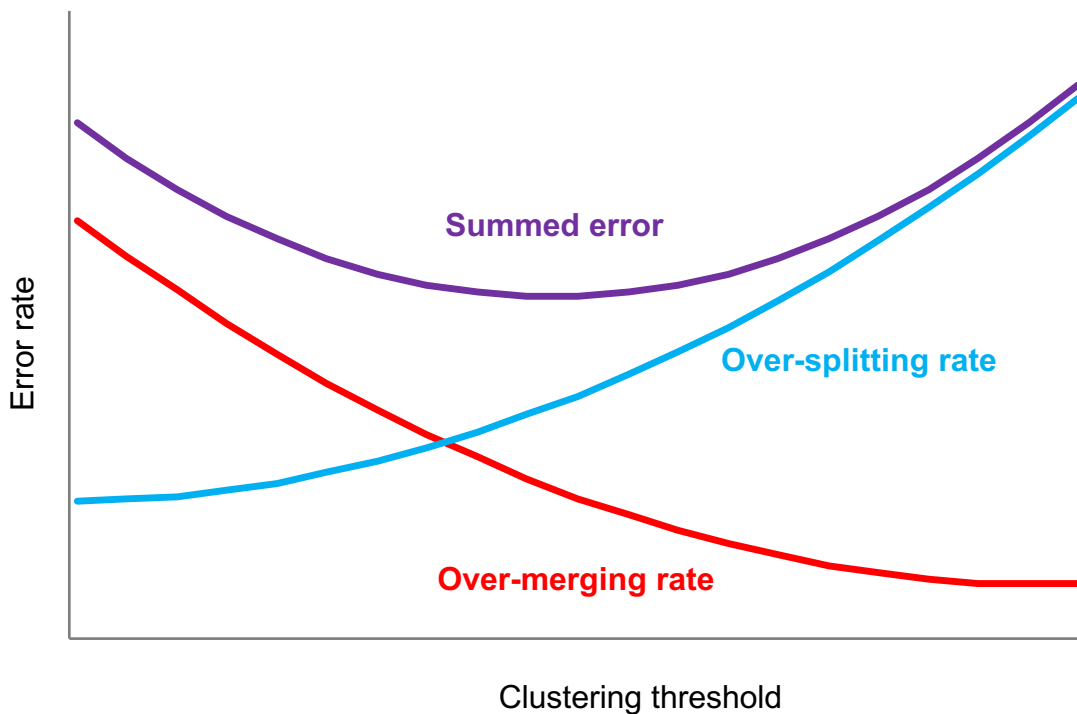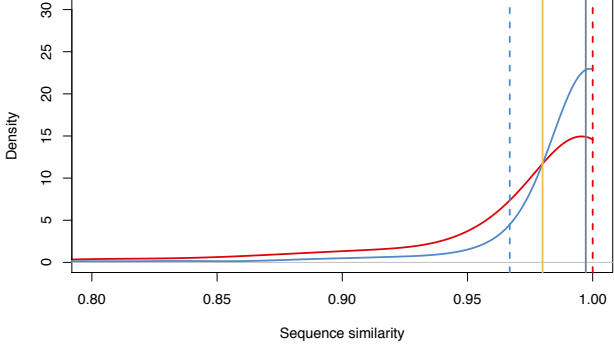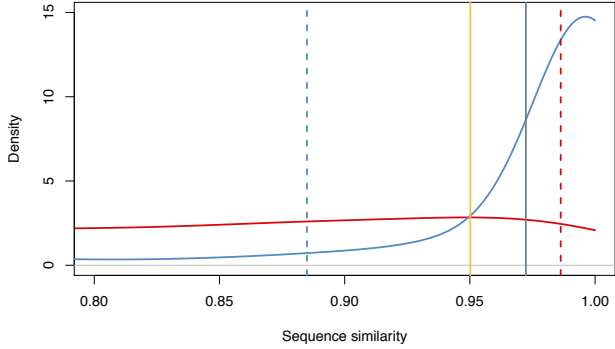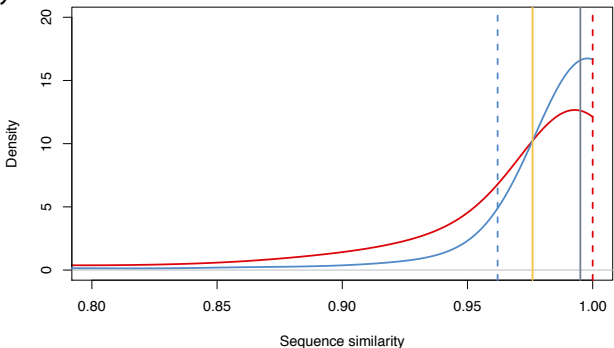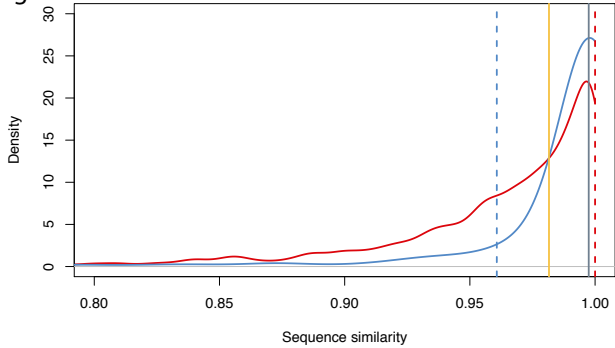
669    for which the summed error rate is minimal.

670

671

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



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

Molecular Ecology Resources

Bact02, Euka02, Fung02, Sper01, Arth02, COI–BF1/BR2, Coll01, Inse01, Olig01

Clustering threshold

Legend:
— Number of clusters
● % of MOTUs containing a single genus
▲ % of MOTUs containing a single species
● % of genera gathered in a single MOTU (among genera represented by several sequences)
▲ % of species gathered in a single MOTU

Molecular Ecology Resources



Bact02, Euka02, Fung02, Sper01, Arth02, COI–BF1/BR2, Coll01, Inse01, Oligo01

Clustering threshold

Over-splitting rate
Over-merging rate
Summed error
Range of optimal clustering thresholds