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

Aurélie Bonin¹, Alessia Guerrieri², and Francesco Ficetola²

¹Argaly ²University of Milan

October 19, 2021

Abstract

Clustering approaches are pivotal to handle the many sequence variants obtained in DNA metabarcoding datasets, therefore they have become a key step of metabarcoding analysis pipelines. Clustering often relies on a sequence similarity threshold to gather sequences in Molecular Operational Taxonomic Units (MOTUs) that ideally each represent a homogeneous taxonomic entity, e.g. a species or a genus. However, the choice of the clustering threshold is rarely justified, and its impact on MOTU oversplitting or over-merging even less tested. Here, we evaluated clustering threshold values for several metabarcoding markers under different criteria: limitation of MOTU over-merging, limitation of MOTU over-splitting, and trade-off between overmerging and over-splitting. We extracted sequences from a public database for eight markers, ranging from generalist markers targeting Bacteria or Eukaryota, to more specific markers targeting a class or a subclass (e.g. Insecta, Oligochaeta). Based on the distributions of pairwise sequence similarities within species and within genera and on the rates of over-splitting and overmerging, or offering a trade-off between the two risks. For generalist markers, high similarity thresholds (0.96-0.99) are generally appropriate, while more specific markers require lower values (0.85-0.96). These results do not support the use of a fixed clustering threshold (e.g. 0.97). Instead, we advocate a careful examination of the most appropriate threshold based on the research objectives, the potential costs of over-splitting and over-merging, and the features of the studied markers.

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

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

- 1. Department of Environmental Science and Policy, University of Milan. Via Celoria 10, 20126 Milano Italy
- 2. Argaly, Bâtiment CleanSpace, 354 Voie Magellan, 73800 Sainte-Hélène-du-Lac, France
- Univ. Grenoble Alpes, CNRS, Univ. Savoie Mont Blanc, LECA, Laboratoire d'Ecologie Alpine, F-38000 Grenoble, France

*Corresponding author: aurelie.bonin@argaly.com

Abstract

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

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

Keywords

metabarcoding marker; sequence variant; analysis parameter; MOTU over-splitting, MOTU over-merging; alpha diversity

Introduction

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

When analyzing metabarcoding data, bioinformatic pipelines generally produce a list of detected sequences, that can be assigned to a given taxon with a more or less precise taxonomic resolution. However, the number of unique sequences obtained after bioinformatic treatment is generally much higher than the number of taxa actually present in the sample (Calderon-Sanou et al. 2020, Machler et al. 2021). This stems from multiple reasons including genuine intraspecific diversity of the selected markers and errors occurring during the amplification or sequencing steps. Consequently, sequence clustering approaches are often used to collapse very similar sequences into one single Molecular Operational Taxonomic Unit (MOTU), which does not necessarily correspond to a species in the traditional sense (Kopylova et al. 2016, Froslev et al. 2017, Bhat et al. 2019, Antich et al. 2021). Sequence clustering can be performed using similarity thresholds, Bayesian approaches, or through single-linkage (Antich et al. 2021). Approaches based on similarity thresholds can have excellent performance and they display several advantages such as flexibility and easy implementation (Kopylova et al. 2016, Wei et al. 2021). However, two key parameters have to be determined a priori when performing clustering based on sequence similarity. The first one is the sequence to be selected as representative of the cluster. In the case of metabarcoding studies, keeping the most abundant sequence of the cluster as the cluster representative is a convenient way of merging sequence variants generated during the PCR or sequencing steps with the original sequence they derive from (Mercier et al. 2013). The second parameter is the similarity threshold (clustering threshold) used to build MOTUS (Clare et al. 2016, Calderon-Sanou et al. 2020, Wei et al. 2021). Choosing this threshold is delicate without prior knowledge of the maker and its intrinsic level of diversity. A too low threshold can collapse different taxa into the same MOTU (overmerging), while a too high threshold can create too many MOTUs (over-splitting) compared to the actual

diversity levels (Clare et al. 2016, Roy et al. 2019, Schloss 2021).

Some works suggest that the ecological interpretation of metabarcoding data can be relatively robust to the threshold selected for sequence clustering. For instance, Botnen et al. (2018) used thresholds ranging from 0.87 to 0.99 of sequence similarity to analyze multiple microbial communities, and they obtained community structures highly coherent across thresholds. Nevertheless, levels of alpha diversity can be heavily impacted by the threshold selection. Ideally, the threshold used for clustering would depend on a trade-off between MOTU over-splitting and MOTU over-merging. A growing number of markers are currently being used in metabarcoding studies (Taberlet et al. 2018), with some allowing broad-scale biodiversity assessment but having limited taxonomic resolution (e.g. 18S rDNA primers amplifying all eukaryotes; Guardiola et al. 2015) and others being highly specific to one single class or even family (e.g. Baamrane et al. 2012, Ficetola et al. 2021). Biodiversity surveys generally aim to generate a set of MOTUs that are each associated with a unique taxon, and with all taxa situated at the same level in the taxonomic tree, to facilitate comparisons. In these conditions, optimal clustering thresholds probably strongly differ across markers. One can for example expect high similarity thresholds for highly conserved markers, and lower clustering thresholds for markers showing high intraspecific variability (Kunin et al. 2010, Brown et al. 2015). However, there is limited quantitative assessment of how optimal clustering thresholds vary across markers (but see Alberdi et al. 2018).

In this study, we analyzed sequences from a public database (EMBL) to identify clustering thresholds for different markers and under different criteria. We considered eight metabarcoding markers (Table 1), ranging from generalist ones (e.g. a 16S rDNA-based marker targeting Bacteria and a 18S rDNA-based marker targeting Eukaryota) to more specific markers (e.g. markers specific of earthworms, insects or springtails). We evaluated how clustering thresholds can change for each taxonomic group, depending on the criterion adopted to set the threshold. We used two alternative strategies to identify thresholds, each time with different objectives in mind. First, following a procedure similar to the one adopted in barcoding studies (Meyer and Paulay 2005), we compared the distribution probabilities of sequence similarities among different individuals of the same species and among different species of the same genus to identify thresholds: i) minimizing the risk that different sequences of the same species are split in different MOTUs (i.e. risk of over-splitting); ii) minimizing the risk that distinct but related species are clustered in the same MOTU (i.e. risk of over-merging); *iii*) balancing the risk of over-splitting and over-merging (Figure 1A). Second, we calculated the over-splitting and over-merging rates of the studied markers for a range of clustering thresholds, to identify values that minimize the two error rates (Figure 1B). We expect that, if researchers want to minimize over-splitting, they should select lower clustering thresholds than if they want to minimize over-merging. Furthermore, we expect higher clustering threshold values for generalist markers compared to markers targeting one class or more restricted taxonomic groups, because of the lower taxonomic resolution and slower evolutionary rate of the former.

Methods

Markers examined and construction of sequence datasets

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

For each of these markers, a sequence database was built from EMBL release 140 as 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 withingenus 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.

Calculation of sequence similarities and probability distributions

As a measure of sequence similarity, we computed the pairwise LCS (Longest Common Subsequence) scores between pairs of sequences in the within-species and within-genus datasets using the *sumatra* program (Mercier et al. 2013). Methodological comparisons showed that this algorithm provides an excellent balance between performance and computation efficiency (Jackson et al. 2016, Kopylova et al. 2016, Bhat et al. 2019). As *sumatra* provides pairwise scores for all possible pairs of sequences, the similarity scores resulting from the within-species dataset were filtered in R (R Core Team 2020) to keep only those representing similarities between sequences of the same species, while the scores resulting from the within-genus dataset were filtered to keep only those representing similarities between different species of the same genus.

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

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

Avoid over-splitting

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

Avoid over-merging

In this case, the aim is to avoid gathering sequences attributed to different species of the same genus in the same cluster (i.e. limiting the probability of merging related species in the same MOTU). For this approach, we selected as clustering threshold the 90% quantile of the distribution of similarities between different species belonging to the same genus. With this approach, the sequences attributed to different species belonging to the same genus are assigned to different clusters in 90% of the cases.

Find a balance between over-splitting and over-merging

In this case, the aim was to minimize both over-splitting and over-merging. We considered two distinct approaches. First, we obtained the probability distribution of within-species and within-genus sequence pairwise similarities using the *density* function from R, with biased cross-validation (bw="bcv") as smoothing bandwidth selector and a Gaussian smoothing kernel (kernel="gaussian"; Venables and Ripley 2002). Other possible smoothing bandwidth selectors were tested, but biased cross-validation was the approach best fitting the score histograms for all markers and all datasets (data not shown). The balance threshold *iii-* a was then

identified as the intersection between the probability distributions of the within-species and within-genus similarities. As an alternative approach to balance over-merging and over-splitting (*iii*- b), we calculated the midpoint between the modes of the within-species and within-genus probability distributions.

Rates of over-merging and over-splitting

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

For each within-species dataset, clustering was performed using the sumaclust program (Mercier et al. 2013) with the -n option (normalization by alignment length) based on the sequence similarities first calculated using the sumatra program (see above; Mercier et al. 2013). Threshold values (-t option) ranging from 0.90 to 1 at 0.01 steps were tested for all markers except Coll01 and Olig01 for which wider ranges ([0.70 - 1] and [0.80 - 1], respectively) were selected based on the within-genus and within-species sequence similarity probability distributions determined previously (see Figure 2). Clustered datasets were then explored to calculate five different variables at each clustering threshold: 1) the number of clusters; 2) the percentage of MOTUs containing one single species; 3) the percentage of MOTUs containing one single genus; 4) the percentage of species gathered in one single MOTU; 5) the percentage of genera gathered in one single MOTU. Variables 2 and 3 are indicative of appropriate MOTU merging of sequences at the species and genus levels, respectively, while variables 4 and 5 are indicative of appropriate MOTU splitting at the species and genus levels, respectively.

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

Results

Our *in-silico* PCRs amplified between 17,000 (Coll01) and 3,2000,000 (Bact02) sequences per marker (Table S2). After data filtering, we retained between 510 (Coll01) and 708,000 (Bact02) sequences per marker. The within-species dataset comprised between 118 (Coll01) and 10,000 (Bact02, Euka02, Fung02, Sper01, Inse01) sequences, while the within-genus dataset comprised between 74 (Coll01) and 10,000 (Euka02 and Sper01) sequences per marker.

Clustering thresholds determined from probability distributions of within-species and withingenus sequence similarities

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

For all markers, criterion i (avoid over-splitting) yielded the lowest thresholds (Figure 3, Table S3), with very low levels for Coll01 and Olig01. Conversely, criterion ii (avoid over-merging) yielded extremely high values, except for Coll01. For all generalist markers, avoiding over-merging would require setting clustering thresholds at 0.99 or higher. For Coll01, criterion ii resulted in a rather low threshold (0.765), because many within-genus comparisons showed very low similarity values. Criteria *iii* -a and *iii* -b searching a balance between over-merging and over-splitting yielded somehow contrasting results across markers. For the three specific markers (Coll01, Inse01, and Olig01), the within-genus and within-species similarities showed clearly distinct peaks (Figure 2). As a consequence, the intersection between the two curves could effectively represent the point minimizing both over-merging and over-splitting (see discussion), and the midpoint between the modes also identified rather similar threshold values. On the contrary, for the generalist markers, the within-species and within-genus similarities showed very high overlap and similar modes, and the density distributions actually intersected at values lower than both modes. The midpoint between the modes continued to identify threshold values intermediate between the peaks of within-species and within-genus similarities.

Rates of over-splitting and over-merging

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

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

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

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

DISCUSSION

Sequence clustering approaches are routinely used for the identification of MOTUs in metabarcoding studies, and they often resort to methods based on similarity values. Still, selecting a clustering threshold for a given marker more than often relies on common practices and rules of thumb rather than on proper scientific argument. By analyzing extensive sequence data deposited in public databases for a range of generalist and specialist markers, we showed that different threshold values can be selected depending on the marker and on the criterion favored by researchers. All the markers we examined are situated in non-protein coding genes (Table S1), and this has an influence on levels of sequence intraspecific diversity. The 10% quantile of the within-species similarity probability distribution was almost always lower than the 0.97 clustering threshold

traditionally used in barcoding for markers targeting protein-coding genes like COI (Hebert et al. 2003), or for microbial MOTU delimitation (Balint et al. 2016), indicating that some level of over-splitting can occur at this threshold.

Although for all the markers the within-genus similarity values were generally lower than the within-species similarities, the overlap between the two distributions was dependent on the generalist vs. specific nature of the marker. For some specific markers (e.g. Coll01 and Olig01), distinct peaks were visible for the two similarity metrics (Figure 2). Within-species similarities generally were >0.90, while within-genus values were lower, frequently below 0.80. Such a pattern is not unexpected for markers with an excellent taxonomic resolution and designed to identify taxa at the species level. Conversely for the generalist markers, within-species and within-genus similarity probability distributions largely overlapped and the differences between the peaks were minimal. Nevertheless, even for these markers, the density of within-species similarity was consistently higher than that of within-genus similarity value is higher within species than within genera. In other words, at high clustering thresholds, a MOTU is more likely to represent a species than a genus. This result is confirmed by the fact that the percentage of MOTUs containing a single species is always higher than 50%, whatever the clustering threshold or the marker considered (Figure 4).

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

In this work, we came up with a global measure of the error associated with a given clustering threshold, that we called the "summed error". We calculate it by summing over-splitting and over-merging rates, assuming both have the same cost for biodiversity studies. However, it is possible to assign a differential weight to over-splitting and over-merging. For instance, if the aim is to reach conservative estimated of alpha diversity (i.e. avoid over-splitting), more weight can be assigned to over-splitting rate. Conversely, if the aim is to tease apart closely related species, that differ in their sensitivity to environmental stressors or in threat levels, one may prefer to avoid over-merging, particularly when extensive reference databases are available (Roy et al. 2019, Lopes et al. 2021).

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

However, for a few markers, the threshold values minimizing summed error yielded somewhat less clear patterns. For Fung02, the summed error rate was rather constant (36-37%) at all the thresholds between 0.91 and 0.98, while it quickly increased for higher clustering thresholds. For Coll01 and Oligo01, the summed

error rate showed multiple minima, some of which at very low clustering thresholds (Figure 5). In principle, increasing the threshold value should determine a monotone decrease of over-merging, and a monotone increase of over-splitting (Figure 1B). However, at low similarity values this was not always the case (Figure 5). This probably occurs because a very large number of sequences have pairwise similarities of 0.80-0.85 for these markers (Figure 2), and this might affect the identification of clusters, with some sequences clustering together e.g. at 0.85 but not at 0.86 similarity values. We also note that these similarity values match the ones corresponding to the intersection between the within-genus and within-species similarities for these markers (Figure 3). It is also possible that, at this level of sequence similarity, there is strong uncertainty between MOTUs representing different hierarchical levels of taxonomy.

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

The selection of clustering thresholds can have strong effect in the estimates of MOTUs richness (Figure 4), still it is important to remember that it often does not have a tremendous effect on the ecological message conveyed by metabarcoding data. For instance, Clare et al. (2016) examined different clustering thresholds to analyze dietary overlap between skinks and shrews in Mauritius. Although high clustering thresholds yielded a larger number of MOTUs, ecological conclusions remained rather consistent overall. Therefore, provided that appropriate parameters are considered (e.g. alpha diversity measured using Hill's numbers with q > 0 instead of richness, beta diversity estimates), the interpretation of data can be relatively robust (Clare et al. 2016, Roy et al. 2019, Calderón-Sanou et al. 2020, Mächler et al. 2021). Nevertheless, we discourage the blind application of one single clustering threshold like the classical 0.97, as it can have very different meaning across markers, and can inflate MOTU richness for fast-evolving markers. Instead, we advocate the ad-hoc definition of the most appropriate thresholds, on the basis of research aims, on the potential costs of over-splitting and over-merging, and on the features of the studied markers.

Acknowledgments

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

References

Alberdi, A., Aizpurua, O., Gilbert, M. T. P., & Bohmann, K. (2018). Scrutinizing key steps for reliable metabarcoding of environmental samples. *Methods in Ecology and Evolution*, 9, 134-147.

Antich, A., Palacin, C., Wangensteen, O. S., & Turon, X. (2021). To denoise or to cluster, that is not the question: optimizing pipelines for COI metabarcoding and metaphylogeography. *BMC Bioinformatics*, 22, 177.

Baamrane, M. A. A., Shehzad, W., Ouhammou, A., Abbad, A., Naimi, M., Coissac, E., ... Znari, M. (2012). Assessment of the food habits of the Moroccan dorcas gazelle in M'Sabih Talaa, West Central Morocco, using the *trnL* approach. *PLoS ONE*, 7, e35643.

Bálint, M., Bahram, M., Eren, A. M., Faust, K., Fuhrman, J. A., Lindahl, B., ... Tedersoo, L. (2016). Millions of reads, thousands of taxa: microbial community structure and associations analyzed via marker genes. *FEMS Microbiology Reviews*, 40, 686-700. Bhat, A. H., Prabhu, P., & Balakrishnan, K. (2019). A critical analysis of state-of-the-art metagenomics OTU clustering algorithms. *Journal of Biosciences*, 44, 9.

Bidartondo, M. I. (2008). Preserving accuracy in GenBank. Science, 319, 1616.

Bienert, F., De Danieli, S., Miquel, C., Coissac, E., Poillot, C., Brun, J. J., & Taberlet, P. (2012). Tracking earthworm communities from soil DNA. *Molecular Ecology*, 21, 2017-2030.

Bohmann, K., Elbrecht ,V., Carøe, C., Bista, L., Leese, F., Bunce, M., Yu, D. W., ... Creer, S. (in press). Strategies for sample labelling and library preparation in DNA metabarcoding studies. *Molecular Ecology Resources*.

Botnen, S. S., Davey, M. L., Halvorsen, R., & Kauserud, H. (2018). Sequence clustering threshold has little effect on the recovery of microbial community structure. *Molecular Ecology Resources*, 18, 1064-1076.

Boyer, F., Mercier, C., Bonin, A., Le Bras, Y., Taberlet, P., & Coissac, E. (2016). OBITOOLS: a Unix-inspired software package for DNA metabarcoding. *Molecular Ecology Resources*, 16, 176-182.

Bridge, P. D., Roberts, P. J., Spooner, B. M., & Panchal, G. (2003). On the unreliability of published DNA sequences. *New Phytologist*, 160, 43-48.

Brown, E. A., Chain, F. J. J., Crease, T. J., MacIsaac, H. J., & Cristescu, M. E. (2015). Divergence thresholds and divergent biodiversity estimates: can metabarcoding reliably describe zooplankton communities? *Ecology* and *Evolution*, 5, 2234-2251.

Calderón-Sanou, I., Münkemüller, T., Boyer, F., Zinger, L., & Thuiller, W. (2020). From environmental DNA sequences to ecological conclusions: How strong is the influence of methodological choices? *Journal of Biogeography*, 47, 193–206.

Capo, E., Giguet-Covex, C., Rouillard, A., Nota, K., Heintzman, P., Vuillemin, A. ... Parducci, L. (2021). Lake sedimentary DNA research on past terrestrial and aquatic biodiversity: Overview and recommendations. *Quaternary* 4, 6.

Chen, W., & Ficetola, G. F. (2020). Statistical and numerical methods for Sedimentary-ancient-DNA-based study on past biodiversity and ecosystem functioning. *Environmental DNA*, 2, 115–129.

Clare, E. L., Chain, F. J. J., Littlefair, J. E., & Cristescu, M. E. (2016). The effects of parameter choice on defining molecular operational taxonomic units and resulting ecological analyses of metabarcoding data. *Genome*, 59, 981-990.

Couton, M., Baud, A., Daguin-Thiébaut, C., Corre, E., Comtet, T., & Viard, F. (2021). High-throughput sequencing on preservative ethanol is effective at jointly examining infraspecific and taxonomic diversity, although bioinformatics pipelines do not perform equally. *Ecology and Evolution*, 11, 5533-5546.

Dickie, I. A., Boyer, S., Buckley, H. L., Duncan, R. P., Gardner, P. P., Hogg, I. D., ... Weaver, L. (2018). Towards robust and repeatable sampling methods in eDNA-based studies. *Molecular Ecology Resources*, 18, 940-952.

Eichmiller, J. J., Miller L. M., & Sorensen, P.W. (2016). Optimizing techniques to capture and extract environmental DNA for detection and quantification of fish. *Molecular Ecology Resources*, 16, 56-68.

Elbrecht, V., Taberlet, P., Dejean, T., Valentini, A., Usseglio-Polatera, P., Beisel, J. N., ... Leese, F. (2016). Testing the potential of a ribosomal 16S marker for DNA metabarcoding of insects. *PeerJ*, 4, 12.

Epp, L. S., Boessenkool, S., Bellemain, E. P., Haile, J., Esposito, A., Riaz, T., ... Brochmann, C. (2012). New environmental metabarcodes for analysing soil DNA: potential for studying past and present ecosystems. *Molecular Ecology*, 21, 1821-1833. Fahner, N. A., Shokralla, S., Baird, D. J., & Hajibabaei, M. (2016). Large-scale monitoring of plants through environmental DNA metabarcoding of soil: Recovery, resolution, and annotation of four DNA markers. *PLoS ONE, 11*, e0157505.

Federhen, S. (2015). Type material in the NCBI Taxonomy Database. *Nucleic Acids Research, 43*, D1086-D1098.

Ficetola, G. F., Boyer, F, Valentini, A. Bonin, Meyer, A., Dejean, T., ... Taberlet, P. (2021). Comparison of markers for the monitoring of freshwater benthic biodiversity through DNA metabarcoding. *Molecular Ecology*, 30, 3189–3202.

Ficetola, G. F., Coissac, E., Zundel, S., Riaz, T., Shehzad, W., Bessiere, J., ... Pompanon, F. (2010). An *in silico* approach for the evaluation of DNA barcodes. *BMC Genomics*, 11, 434.

Froslev, T. G., Kjoller, R., Bruun, H. H., Ejrnaes, R., Brunbjerg, A. K., Pietroni, C., & Hansen, A. J. (2017). Algorithm for post-clustering curation of DNA amplicon data yields reliable biodiversity estimates. *Nature Communications*, 8, 11.

Guardiola, M., Uriz, M. J., Taberlet, P., Coissac, E., Wangensteen, O. S., & Turon, X. (2015). Deep-sea, deep-sequencing: metabarcoding extracellular DNA from sediments of marine canyons. *PLoS ONE*, 10, e0139633.

Guerrieri, A., Bonin, A., Munkemuller, T., Gielly, L., Thuiller, W., & Ficetola, G. F. (2021). Effects of soil preservation for biodiversity monitoring using environmental DNA. *Molecular Ecology*, 30, 3313-3325.

Hebert, P. D. N., Ratnasingham, S., & deWaard, J. R. (2003). Barcoding animal life: cytochrome c oxidase subunit 1 divergences among closely related species. *Proceedings of the Royal Society B-Biological Sciences*, 270, S96-S99.

Jackson, M. A., Bell, J. T., Spector, T. D., & Steves, C. J. (2016). A heritability-based comparison of methods used to cluster 16S rRNA gene sequences into operational taxonomic units. *PeerJ*, 4, 19.

Janssen, P., Bec, S., Fuhr, M., Taberlet, P., Brun, J.-J., & Bouget, C. (2018). Present conditions may mediate the legacy effect of past land-use changes on species richness and composition of above- and below-ground assemblages. *Journal of Ecology*, 106, 306-318.

Kopylova, E., Navas-Molina, J. A., Mercier, C., Xu, Z. Z., Mahe, F., He, Y., ... Knight, R. (2016). Opensource sequence clustering methods improve the state of the art. *mSystems*, 1, 16.

Kunin, V., Engelbrektson, A., Ochman, H., & Hugenholtz, P. (2010). Wrinkles in the rare biosphere: pyrosequencing errors can lead to artificial inflation of diversity estimates. *Environmental Microbiology*, 12, 118-123.

Lear, G., Dickie, I., Banks, J., Boyer, S., Buckley, H. L., Buckley, T. R., ... Holdaway, R. (2018). Methods for the extraction, storage, amplification and sequencing of DNA from environmental samples. *New Zealand Journal of Ecology*, 42, 10.

Leray, M., Knowlton, N., Ho, S. L., Nguyen, B. N., & Machida, R. J. (2019). GenBank is a reliable resource for 21st century biodiversity research. *Proceedings of the National Academy of Sciences of the United States of America*, 116, 22651-22656.

Lopes, C. M., Baeta, D., Valentini, A., Lyra, M. L., Sabbag, A. F., Gasparini, J. L., ... Zamudio, R. K. (2021). Lost and found: Frogs in a biodiversity hotspot rediscovered with environmental DNA. *Molecular Ecology*, 30, 3289-3298.

Macher, T.-H., Beermann, A. J., & Leese, F. (2021). TaxonTableTools: A comprehensive, platformindependent graphical user interface software to explore and visualise DNA metabarcoding data. *Molecular Ecology Resources, 21*, 1705-1714. Machler, E., Walser, J.-C., & Altermatt, F. (2021). Decision-making and best practices for taxonomy-free environmental DNA metabarcoding in biomonitoring using Hill numbers. *Molecular Ecology*, 30, 3326-3339.

Mercier, C., Boyer, F., Bonin, A., & Coissac, E. (2013). SUMATRA and SUMACLUST: fast and exact comparison and clustering of sequences. *Programs and Abstracts of the SeqBio 2013 Workshop*, 27-29.

Meyer, C. P., & Paulay, G. (2005). DNA barcoding: Error rates based on comprehensive sampling. *PLoS Biology*, 3, 2229-2238.

Mioduchowska, M., Czyz, M. J., Goldyn, B., Kur, J., & Sell, J. (2018). Instances of erroneous DNA barcoding of metazoan invertebrates: Are universal *cox1* gene primers too "universal"? *PLoS ONE*, 13, e0199609.

Nichols, R. V., Vollmers, C., Newsom, L. A., Wang, Y., Heintzman, P. D., Leighton, M., ... Shapiro, B. (2018). Minimizing polymerase biases in metabarcoding. *Molecular Ecology Resources*, 18, 927-939.

Paliy, O., & Shankar, V. (2016). Application of multivariate statistical techniques in microbial ecology. *Molecular Ecology*, 25, 1032-1057.

Porco, D., Bedos, A., Penelope, G., Janion, C., Skarżyński, D., Stevens, M. I., ... Deharveng, L. (2012). Challenging species delimitation in Collembola: cryptic diversity among common springtails unveiled by DNA barcoding. *Invertebrate Systematics*, 26, 470-477.

R Core Team. (2020). R: A language and environment for statistical computing. R Foundation for Statistical Computing, Vienna.

Roy, J., Mazel, F., Sosa-Hernández, M. A., Dueñas, J. F., Hempel, S., Zinger, L., & Rillig, M. C. (2019). The relative importance of ecological drivers of arbuscular mycorrhizal fungal distribution varies with taxon phylogenetic resolution. *New Phytologist*, 224,936-948.

Schloss, P. D. (2021). Amplicon sequence variants artificially split bacterial genomes into separate clusters. mSphere, 6, e00191-00121.

Taberlet, P., Bonin, A., Zinger, L., & Coissac, E. (2018). Environmental DNA for biodiversity research and monitoring. Oxford University Press, Oxford.

Taberlet, P., Coissac, E., Pompanon, F., Gielly, L., Miquel, C., Valentini, A., ... Willerslev, E. (2007). Power and limitations of the chloroplast *trnL* (UAA) intron for plant DNA barcoding. *Nucleic Acids Research*, 35, e14.

Taberlet, P., Prud'homme, S. M., Campione, E., Roy, J., Miquel, C., Shehzad, W., ... Coissac, E. (2012). Soil sampling and isolation of extracellular DNA from large amount of starting material suitable for meta-barcoding studies. *Molecular Ecology*, 21, 1816-1820.

Tatangelo, V., Franzetti, A., Gandolfi, I., Bestetti, G., & Ambrosini, R. (2014). Effect of preservation method on the assessment of bacterial community structure in soil and water samples. *FEMS Microbiology Letters*, 356, 32-38.

Valkiūnas, G., Atkinson, C. T., Bensch, S., Sehgal, R. N., & Ricklefs, R. E. (2008). Parasite misidentifications in GenBank: how to minimize their number? *Trends in Parasitology*, 24, 247-248.

Venables, W. N., & Ripley, B. D. (2002). Modern applied statistics with S. Fourth Edition. Springer, New York.

Wei, Z.-G., Zhang, X.-D., Cao, M., Liu, F., Qian, Y., & Zhang, S.-W. (2021). Comparison of methods for picking the operational taxonomic units from amplicon sequences. *Frontiers in Microbiology*, 12, 644012.

Zinger, L., Bonin, A., Alsos, I., Bálint, M., Bik, H., Boyer, F., ... Taberlet, P. (2019). DNA metabarcoding - need for robust experimental designs to draw sound ecological conclusions. *Molecular Ecology*, 28, 1857-1862.

Zinger, L., Chave, J., Coissac, E., Iribar, A., Louisanna, E., Manzi, S., ... Taberlet, P. (2016). Extracellular DNA extraction is a fast, cheap and reliable alternative for multi-taxa surveys based on soil DNA. *Soil Biology & Biochemistry*, 96, 16-19.

Data Accessibility

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

Authors Contribution

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

Marker	Target group	Taxonomic level	Taxonomic resolution \ast	Taxonomic resolution \ast	Taxonomic resolution
			Species level	Genus level	Family level
Bact02	Bacteria	Superkingdom	19.6%	55.7%	55.1%
Euka02	Eukaryota	Superkingdom	47.0%	59.5%	68.3%
Fung02	Fungi	Kingdom	72.5%	90.2%	87.7%
Sper01	Spermatophyta	Clade < kingdom	21.5%	36.9%	77.4%
Arth02	Arthropoda	Phylum	68.6%	89.6%	97.5%
Coll01	Collembola	Class	80.5%	87.2%	75.0%
Inse01	Insecta	Class	87.8%	96.8%	95.4%
Olig01	Oligochaeta	Subclass	89.3%	95.7%	100.0%

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

* Estimated as the percentage of discriminated taxa among amplified taxa; reported from Taberlet et al. (2018).

Figure captions

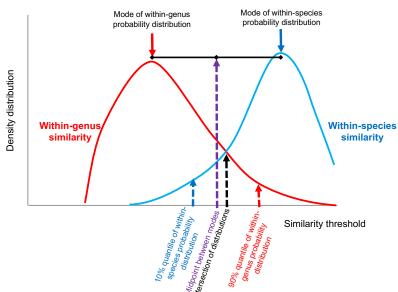
Figure 1. Different approaches to identify the most appropriate clustering thresholds. A): approaches based on similarities between sequences belonging to different individuals from the same species (blue curve), and similarities between sequences belonging to different species from the same genus (red curve). One can choose to minimize the risk that different sequences from the same species are split in different MOTUs (oversplitting 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) Approaches 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 eight studied markers. For each marker, dotted lines represent the 10% quantile of the within-species probability distribution (blue; threshold limiting over-splitting), the 90% quantile of the within-genus probability distribution (red; threshold limiting over-merging), the intersection of the within-species and within-genus probability distributions (green, balance-a) and the midpoint between modes (black, balance-b)

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

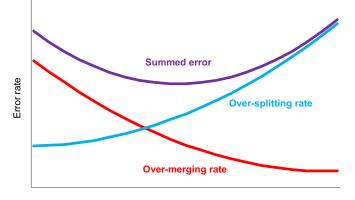
Figure 4. Evolution of over-splitting and over-merging rates for a range of clustering thresholds, for the eight studied markers. The left y-axes report percentage values; the right y-axes indicate the number of obtained clusters.

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



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

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



Similarity threshold

