A flexible pipeline combining clustering and correction tools for prokaryotic and

eukaryotic metabarcoding

Short title:

A flexible metabarcoding pipeline based on read correction

Miriam I. Brandt[1], Blandine Trouche[2], Laure Quintric[3], Patrick Wincker[4,5], Julie Poulain[4,5], and

Sophie Arnaud-Haond[1]


[1]MARBEC, Ifremer, Univ. Montpellier, IRD, CNRS, Sète, France

[2] Univ. Brest, CNRS, Ifremer, Laboratoire de Microbiologie des Environnements

Extrêmes, Plouzané, France

[3]Ifremer, Cellule Bioinformatique, Brest, France

[4] Génomique Métabolique, Génoscope, Institut François Jacob, CEA, CNRS, Univ. Evry,

Université Paris-Saclay, 91057 Evry, France

[5] Research Federation for the study of Global Ocean Systems Ecology and Evolution,

FR2022/ Tara


Corresponding author: sarnaud@ifremer.fr, miriam.isabelle.brandt@gmail.com,

## ABSTRACT

1    Environmental metabarcoding is an increasingly popular tool for studying biodiversity in

2    marine and terrestrial biomes. With sequencing costs decreasing, multiple-marker metabarcoding,

3    spanning several branches of the tree of life, is becoming more accessible. However, bioinformatic

4    approaches need to adjust to the diversity of taxonomic compartments targeted as well as to each

5    barcode gene specificities. We built and tested a pipeline based on Illumina read correction with

6    DADA2 allowing analysing metabarcoding data from prokaryotic (16S) and eukaryotic (18S, COI)

7    life compartments. We implemented the option to cluster Amplicon Sequence Variants (ASVs)

8    into Operational Taxonomic Units (OTUs) with swarm v2, a network-based clustering algorithm,

9    and to further curate the ASVs/OTUs based on sequence similarity and co-occurrence rates using

10   a recently developed algorithm, LULU. Finally, flexible taxonomic assignment was implemented

11   *via* Ribosomal Database Project (RDP) Bayesian classifier and BLAST. We validate this pipeline

12   with ribosomal and mitochondrial markers using eukaryotic mock communities and 42 deep-sea

13   sediment samples. The results show that ASVs, reflecting genetic diversity, may not be appropriate

14   for alpha diversity estimation of organisms fitting the biological species concept. The results

15   underline the advantages of clustering and LULU-curation for producing more reliable metazoan

16   biodiversity inventories, and show that LULU is an effective tool for filtering metazoan molecular

17   clusters, although the minimum identity threshold applied to co-occurring OTUs has to be

18   increased for 18S. The comparison of BLAST and the RDP Classifier underlined the potential of

19   the latter to deliver very good assignments, but highlighted the need for a concerted effort to build

20   comprehensive, ecosystem-specific, databases adapted to the studied communities.

21

22

23      Key words: Biodiversity, bioinformatics, environmental DNA, metabarcoding, mock

24    communities, eukaryotes (18S and COI), prokaryotes (16S)

25

26 **INTRODUCTION**

27      High-throughput sequencing (HTS) technologies are revolutionizing the way we assess

28 biodiversity. By producing millions of DNA sequences per sample, HTS allows broad taxonomic

29 biodiversity surveys through metabarcoding of bulk DNA from complex communities or from

30 environmental DNA (eDNA) directly extracted from soil, water, and air samples. First developed

31 to unravel cryptic and uncultured prokaryotic diversity, metabarcoding methods have been

32 extended to eukaryotes as powerful, non-invasive tools, allowing detection of a wide range of taxa

33 in a rapid, cost-effective way using a variety of sample types (Valentini et al. 2009; Taberlet et al.

34 2012; Creer et al. 2016; Stat et al. 2017). In the last decade, these tools have been used to describe

35 past and present biodiversity in terrestrial (Ji et al. 2013; Yoccoz et al. 2012; Yu et al. 2012; Slon

36 et al. 2017; Pansu et al. 2015), freshwater (Valentini et al. 2016; Deiner et al. 2016; Bista et al.

37 2015; Dejean et al. 2011; Evans et al. 2016), and marine (Fonseca et al. 2010; Sinniger et al. 2016;

38 Pawlowski et al. 2011; Massana et al. 2015; De Vargas et al. 2015; Salazar et al. 2016; Boussarie

39 et al. 2018; Bik et al. 2012) environments.

40      As every new technique brings on new challenges, a number of studies have put

41 considerable effort into delineating critical aspects of metabarcoding protocols to ensure robust and

42 reproducible results (see Fig.1 in Fonseca et al, 2018). Recent studies have addressed many issues

43 regarding sampling methods (Dickie et al. 2018), contamination risks (Goldberg et al. 2016), DNA

44 extraction protocols (Brannock and Halanych 2015; Deiner et al. 2015; Zinger et al. 2016),

45 amplification biases and required PCR replication levels (Nichols et al. 2018; Alberdi et al. 2017;

46 Ficetola et al. 2015). Similarly, computational pipelines, through which molecular data are

47 transformed into ecological inventories of putative taxa, have also been in constant improvement.

48 PCR-generated errors and sequencing errors are major bioinformatic challenges for metabarcoding

49 pipelines, as they can strongly bias biodiversity estimates (Coissac et al. 2012; Bokulich et al.

50   2013). A variety of tools have thus been developed for quality-filtering amplicon data to remove

51   erroneous reads and improve the reliability of Illumina-sequenced metabarcoding inventories

52   (Bokulich et al. 2013; Eren et al. 2013; Minoche et al. 2011). Studies that evaluated bioinformatic

53   processing steps have generally found that sequence quality-filtering parameters and clustering

54   thresholds most strongly affect molecular biodiversity inventories, resulting in considerable

55   variation during data analysis(Brannock and Halanych 2015; Clare et al. 2016; Brown et al. 2015;

56   Xiong and Zhan 2018).

57        There were historically two main reasons for clustering sequences into Operational

58   Taxonomic Units (OTUs). The first was to limit the bias due to PCR and sequencing errors (and

59   to some extent intra-individual variability linked to the existence of pseudogenes) by clustering

60   erroneous sequences with error-free target sequences. The second was to delineate OTUs as

61   clusters of homologous sequences (by grouping the alleles/haplotype at the same locus) that

62   would best fit a "species level", i.e. the Operational Taxonomic Units defined using a classical

63   phenetic *proxy* (Sokal and Crovello 1970). Recent bioinformatic algorithms alleviate the

64   influence of errors and intraspecific variability in metabarcoding datasets. First, amplicon-

65   specific error correction methods, commonly used to correct sequences produced by

66   pyrosequencing (Coissac et al. 2012), have now become available for Illumina-sequenced data.

67   Introduced in 2016, DADA2 effectively corrects Illumina sequencing errors and has quickly

68   become a widely used tool, particularly in the microbial world, producing more accurate

69   biodiversity inventories and resolving fine-scale genetic variation by defining Amplicon

70   Sequence Variants (ASVs) (Callahan et al. 2016; Nearing et al. 2018). Second, LULU is a

71   recently developed curation algorithm designed to filter out spurious clusters, originating from

72   PCR and sequencing errors, or intra-individual variability (pseudogenes, heteroplasmy), based on

73   their similarity and co-occurrence rate with more abundant clusters, allowing obtaining curated

5

74    datasets while avoiding arbitrary abundance filters (Frøslev et al. 2017). The authors validated

75    their approach on metabarcoding of plants using ITS2 (nuclear ribosomal internal transcribed

76    spacer region 2) and evaluated it on several pipelines. Their results show that ASV definition

77    with DADA2, subsequent clustering to address intraspecific variation, and final curation with

78    LULU is the safest pathway for producing reliable and accurate metabarcoding data. The authors

79    concluded that their validation on plants is relevant to other organism groups and other markers,

80    while recommending future validation of LULU on mock communities as LULU's minimum

81    match parameter may need to be adjusted to less variable marker genes.

82    The impact of errors being strongly decreased by correction algorithms such as DADA2

83    and LULU, the relevance of clustering sequences into OTUs is now being debated. Indeed, after

84    presenting their new algorithm on prokaryotic communities, the authors of DADA2 proposed that

85    the reproducibility and comparability of ASVs across studies challenge the need for clustering

86    sequences, as OTUs have the disadvantage of being study-specific and defined using arbitrary

87    thresholds (Callahan et al. 2017). However, clustering sequences may still be necessary in

88    metazoan datasets, where very distinct levels of intraspecific polymorphism can exist in the same

89    gene region among taxa due to both evolutionary and biological specificity (Bucklin et al. 2011;

90    Phillips et al. 2019). ASV-based inventories will thus be biased in favour of taxa with high levels

91    of intraspecific diversity, even though the latter are not necessarily the most abundant ones (Bazin

92    et al. 2006). Such bias in biodiversity inventories based on ASVs is likely to be magnified in

93    presence-absence metabarcoding datasets, commonly used for metazoan communities (Ji et al.

94    2013). Similarly, imposing a "universal" clustering threshold on metabarcoding datasets is also

95    introducing bias, penalizing groups with lower interspecific divergence, and overestimating species

96    diversity in groups with higher interspecific divergence. However, this can be alleviated with tools

97    such as swarm v2, a single-linkage clustering algorithm (Mahe et al. 2015). Based on network

98    theory, swarm v2 aggregates sequences iteratively and locally around seed sequences and

99    determines coherent groups of sequences, independent of amplicon input order, allowing highly

100    scalable and fine-scale clustering. Finally, it is widely recognized that homogeneous entities

101    sharing a set of evolutionary and ecological properties, i.e. *species* (Mayr 1942; de Queiroz 2005),

102    sometimes referred to "ecotypes" for prokaryotes (Cohan 2001; Gevers et al. 2005), represent a

103    fundamental category of biological organization that is the cornerstone of most ecological and

104    evolutionary theories and empirical studies. Maintaining ASV information for feeding databases

105    and cross-comparing studies is not incompatible with their clustering into OTUs, and this choice

106    depends on the purpose of the study, i.e. providing a census of the extent and distribution of genetic

107    polymorphism for a given gene, or a census of biodiversity to be used and manipulated in

108    ecological or evolutionary studies.

109    Here we evaluate DADA2 and LULU, using them alone and in combination with swarm

110    v2, to assess the performance of these new tools for metabarcoding of metazoan communities.

111    Using both mitochondrial COI (Leray et al. 2013) and the V1-V2 region of 18S ribosomal RNA

112    (rRNA) (Sinniger et al. 2016), we evaluated the need for clustering  and the effectiveness of LULU

113    curation to select pipeline parameters delivering the most accurate resolution of two deep-sea mock

114    communities. We then test the different bioinformatic tools on a deep-sea sediment dataset in order

115    to select an optimal trade-off between inflating biodiversity estimates and loosing rare biodiversity.

116    As a baseline for comparison, and in the perspective of the joint study of metazoan and microbial

117    taxa, we also analysed the 16S V4-V5 rRNA barcode on these natural samples (Parada et al. 2016).

118    Our objectives were to (1) discuss the use of ASV *vs* OTU-centred datasets depending on

119    taxonomic compartment and study objectives, and (2) determine the most adequate swarm-

120    clustering and LULU curation thresholds that avoid inflating biodiversity estimates while retaining

121    rare biodiversity.

7

122

## 1  MATERIALS AND METHODS

### 1.1  Preparation of samples

*Mock communities*

Genomic-DNA mass-balanced metazoan mock communities (5 ng/µL) were prepared using standardized 10 ng/µL DNA extracts of ten deep-sea specimens belonging to five taxonomic groups (Polychaeta, Crustacea, Anthozoa, Bivalvia, Gastropoda; Table S1). Specimen DNA was extracted using a CTAB extraction protocol, from muscle tissue or from whole polyps in the case of cnidarians. The mock communities differed in terms of ratios of total genomic DNA from each species, with increased dominance of three species and secondary species DNA input decreasing from 3% to 0.7%. We individually barcoded the species present in the mock communities: PCRs of both target genes were performed using the same primers as the ones used in metabarcoding (see below). The PCR reactions (25 µL final volume) contained 2 µL DNA template with 0.5 µM concentration of each primer, 1X *Phusion* Master Mix, and an additional 1 mM $MgCl_2$ for COI. PCR amplifications (98 °C for 30 s; 40 cycles of 10 s at 98 °C, 45 s at 48 °C (COI) or 57 °C (18S), 30 s at 72 °C; and 72 °C for 5 min) were cleaned up with ExoSAP (Thermo Fisher Scientific, Waltham, MA, USA) and sent to Eurofins (Eurofins Scientific, Luxembourg) for Sanger sequencing. The barcode sequences obtained for all mock specimens were added to the databases used for taxonomic assignments of metabarcoding datasets, and were submitted on Genbank under accession numbers MN826120-MN826130 and MN844176-MN844185.

142

*Environmental DNA*

Sediment cores were collected from thirteen deep-sea sites ranging from the Arctic to the Mediterranean during various cruises (Table S2). Sampling was carried out with a multicorer or

8

146 with a remotely operated vehicle. Three tube cores were taken at each sampling station (GPS

147 coordinates in Table S2). The latter were sliced into depth layers that were transferred into zip-lock

148 bags, homogenised, and frozen at −80°C on board before being shipped on dry ice to the laboratory.

149 The first layer (0-1 cm) was used in the present study. DNA extractions were performed using

150 approximately 10 g of sediment with the PowerMax Soil DNA Isolation Kit (Qiagen, Hilden,

151 Germany). To increase the DNA yield, the elution buffer was left on the spin filter membrane for

152 10 min at room temperature before centrifugation. The ~5 mL extract was then split into three parts,

153 one of which was kept in screw-cap tubes for archiving purposes and stored at -80°C. For the four

154 field controls, the first solution of the kit was poured into the control zip-lock bag, before following

155 the usual extraction steps. For the two negative extraction controls, a blank extraction (adding

156 nothing to the bead tube) was performed alongside sample extractions.

157

158 **1.2    Amplicon library construction and high-throughput sequencing**

159      Two primer pairs were used to amplify the mitochondrial COI and the 18S V1-V2 rRNA

160 barcode genes specifically targeting metazoans, and one pair of primer was used to amplify the

161 prokaryote 16S V4-V5 region. PCR amplifications, library preparation, and sequencing were

162 carried out at Genoscope (Evry, France) as part of the eDNAbyss project.

163

164 *Eukaryotic 18S V1-V2 rRNA gene amplicon generation*

165      Amplifications were performed with the *Phusion* High Fidelity PCR Master Mix with GC

166 buffer (Thermo Fisher Scientific, Waltham, MA, USA) and the SSUF04 (5'-

167 GCTTGTCTCAAAGATTAAGCC-3') and SSUR22*mod* (5'- CCTGCTGCCTTCCTTRGA-3')

168 primers (Sinniger et al. 2016), preferentially targeting metazoans, the primary focus of this study.

169 The PCR reactions (25 µL final volume) contained 2.5 ng or less of DNA template with 0.4 µM

9

170 concentration of each primer, 3% of DMSO, and 1X *Phusion* Master Mix. PCR amplifications

171 (98 °C for 30 s; 25 cycles of 10 s at 98 °C, 30 s at 45 °C, 30 s at 72 °C; and 72 °C for 10 min) of all

172 samples were carried out in triplicate in order to smooth the intra-sample variance while obtaining

173 sufficient amounts of amplicons for Illumina sequencing.

174

175 *Eukaryotic COI gene amplicon generation*

176 Metazoan COI barcodes were generated using the mlCOIintF (5'-

177 GGWACWGGWTGAACWGTWTAYCCYCC-3') and jgHCO2198 (5'-

178 TAIACYTCIGGRTGICCRAARAAYCA-3') primers (Leray et al. 2013). Triplicate PCR

179 reactions (20 µl final volume) contained 2.5 ng or less of total DNA template with 0.5 µM final

180 concentration of each primer, 3% of DMSO, 0.175 mM final concentration of dNTPs, and 1X

181 Advantage 2 Polymerase Mix (Takara Bio, Kusatsu, Japan). Cycling conditions included a 10 min

182 denaturation step followed by 16 cycles of 95 °C for 10 s, 30s at 62°C (−1°C per cycle), 68 °C for

183 60 s, followed by 15 cycles of 95 °C for 10 s, 30s at 46°C, 68 °C for 60 s and a final extension of

184 68 °C for 7 min.

185 *Prokaryotic 16S rRNA gene amplicon generation*

186 Prokaryotic barcodes were generated using 515F-Y (5'- GTGYCAGCMGCCGCGGTAA-

187 3') and 926R (5'- CCGYCAATTYMTTTRAGTTT-3') 16S-V4V5 primers (Parada et al. 2016).

188 Triplicate PCR mixtures were prepared as described above for 18S-V1V2, but cycling conditions

189 included a 30 s denaturation step followed by 25 cycles of 98 °C for 10 s, 53 °C for 30 s, 72 °C for

190 30 s, and a final extension of 72 °C for 10 min.

191

192

10

193    *Amplicon library preparation*

194          PCR triplicates were pooled and PCR products purified using 1X AMPure XP beads

195    (Beckman Coulter, Brea, CA, USA) clean up. Aliquots of purified amplicons were run on an

196    Agilent Bioanalyzer using the DNA High Sensitivity LabChip kit (Agilent Technologies, Santa

197    Clara, CA, USA) to check their lengths and quantified with a Qubit fluorimeter (Invitrogen,

198    Carlsbad, CA, USA). One hundred nanograms of pooled amplicon triplicates were directly end-

199    repaired, A-tailed and ligated to Illumina adapters on a Biomek FX Laboratory Automation

200    Workstation (Beckman Coulter, Brea, CA, USA). Library amplification was performed using a

201    Kapa Hifi HotStart NGS library Amplification kit (Kapa Biosystems, Wilmington, MA, USA) with

202    the same cycling conditions applied for all metagenomic libraries and purified using 1X AMPure

203    XP beads.

204

205    *Sequencing library quality control*

206          Amplicon libraries were quantified by Quant-iT dsDNA HS assay kits using a Fluoroskan

207    Ascent microplate fluorometer (Thermo Fisher Scientific, Waltham, MA, USA) and then by qPCR

208    with the KAPA Library Quantification Kit for Illumina Libraries (Kapa Biosystems, Wilmington,

209    MA, USA) on an MxPro instrument (Agilent Technologies, Santa Clara, CA, USA). Library

210    profiles were assessed using a high-throughput microfluidic capillary electrophoresis system

211    (LabChip GX, Perkin Elmer, Waltham, MA, USA).

212

213    *Sequencing procedures*

214          Library concentrations were normalized to 10 nM by addition of 10 mM Tris-Cl (pH 8.5)

215    and applied to cluster generation according to the Illumina Cbot User Guide (Part # 15006165).

216    Amplicon libraries are characterized by low diversity sequences at the beginning of the reads due

11

217    to the presence of the primer sequence. Low-diversity libraries can interfere in correct cluster

218    identification, resulting in a drastic loss of data output. Therefore, loading concentrations of

219    libraries were decreased (8–9 pM instead of 12–14 pM for standard libraries) and PhiX DNA spike-

220    in was increased (20% instead of 1%) in order to minimize the impacts on the run quality.

221    Libraries were sequenced on HiSeq2500 (System User Guide Part # 15035786) instruments

222    (Illumina, San Diego, CA, USA) in a 250 bp paired-end mode.

223

224    **1.3    Bioinformatic analyses**

225        All bioinformatic analyses were performed using a Unix shell script on a home-based

226    cluster (DATARMOR, Ifremer), available on Gitlab (https://gitlab.ifremer.fr/abyss-project/). The

227    mock communities were analysed alongside the natural samples, and used to validate the

228    metabarcoding pipeline in terms of detection of correct species and presence of false-positives. The

229    details of the pipeline, along with specific parameters used for all three metabarcoding markers are

230    listed in Table S3.

231

232    *Reads preprocessing*

233        Our multiplexing strategy relies on ligation of adapters to amplicon pools, meaning that

234    contrary to libraries produced by double PCR, the reads in each paired sequencing run can be

235    forward or reverse. DADA2 correction is based on error distribution differing between R1 and R2

236    reads. We thus developed a custom script (*abyss-preprocessing* in abyss-pipeline) allowing

237    separating forward and reverse reads in each paired run and reformatting the outputs to be

238    compatible with DADA2. Briefly, the script uses cutadapt v1.18 to detect and remove primers,

239    while separating forward and reverse reads in each paired sequence file to produce two pairs of

240    sequence files per sample named R1F/R2R and R2F/R1R. Cutadapt parameters (Table S3) were

12

241     set to require an overlap over the full length of the primer (default: 3 nt), with 2-4 nt mismatches

242     allowed for ribosomal loci, and 7 nt mismatches allowed for COI (default: 10%). Each identified

243     forward and reverse read is then renamed which the correct extension (/1 and /2 respectively),

244     which is a requirement for DADA2 to recognize the pairs of reads. Each pair of renamed sequence

245     files is then re-paired with BBMAP Repair v38.22 in order to remove singleton reads (non-paired

246     reads). Optionally, sequence file names can also be renamed if necessary using a CSV

247     correspondence file.

248

249     *Read correction, amplicon cluster generation and taxonomic assignment*

250     Pairs of Illumina reads were corrected with DADA2 v.1.10 (Callahan et al. 2016) following

251     the online tutorial for paired-end HiSeq data

252     (https://benjjneb.github.io/dada2/bigdata_paired.html). Reads were filtered and trimmed with the

253     *filterAndTrim* function and all reads containing ambiguous bases removed. The parameters were

254     set based on tutorial recommendations and trimming lengths were adjusted based on sequence

255     quality profiles, so that Q-scores remained above 30 (truncLen at 220 for 18S and 16S, 200 for

256     COI, maxEE at 2, truncQ at 11, maxN at 0).

257     The error model was calculated for forward and reverse reads (R1F/R2R pairs and then

258     R2F/R1R pairs) with *learnErrors* based on 100 million randomly chosen bases (default), and reads

259     were dereplicated using *derepFastq*. After read correction with the *dada* function, forward and

260     reverse reads were merged with a minimum overlap of 12 nucleotides, allowing no mismatches

261     (default). The amplicons were then filtered by size. The size range was set to 330-390 bp for the

262     18S SSU rRNA marker gene, 300-326 bp for the COI marker gene, and 350-390 bp for the 16S

263     rRNA marker gene.

13

264    Chimeras were removed with *removeBimeraDenovo* and ASVs were taxonomically

265    assigned via the RDP naïve Bayesian classifier method, the default assignment method

266    implemented in DADA2. A second taxonomic assignment method was optionally implemented in

267    the pipeline, allowing assigning ASVs using BLAST+ (Basic Local Alignment Search Tool v2.6.0)

268    based on minimum similarity and minimum coverage (-perc_identity 70 and –qcov_hsp 80). An

269    initial test implementing BLASTn+ to assign taxonomy only to the COI dataset using a 96%

270    percent identity threshold led to the exclusion of the majority of the clusters. Given observed inter-

271    specific mitochondrial DNA divergence levels of up to 30% within a same polychaete genus (Zanol

272    et al. 2010) or among some closely related deep-sea shrimp species (Shank et al. 1999), and

273    considering our interest in the identities of multiple, largely unknown taxa in poorly characterized

274    communities, more stringent BLAST thresholds were not implemented at this stage. The Silva132

275    reference database was used for the 16S and 18S SSU rRNA marker genes (Quast et al. 2012), and

276    MIDORI-UNIQUE (Machida et al. 2017) was used for COI. The databases were downloaded from

277    the DADA2 website (https://benjjneb.github.io/dada2/training.html) and from the FROGS website

278    (http://genoweb.toulouse.inra.fr/frogs_databanks/assignation/). Finally, to evaluate the effect of

279    clustering, ASV tables produced by DADA2 were clustered with swarm v2 (Mahe et al. 2015) at

280    *d=1,3,4,5 and 11* for 18S and 16S, and *d=1,5,6,7, and 13* for COI in FROGS

281    (http://frogs.toulouse.inra.fr/) (Escudié et al. 2018). Resulting OTUs were taxonomically assigned

282    via RDP and BLAST+ using the databases stated above.

283    Molecular clusters were refined in R v.3.5.1 (R Core Team 2018). A blank correction was

284    made using the *decontam* package v.1.2.1 (Davis et al. 2018), removing all clusters that were

285    prevalent (more frequent) in negative control samples. ASV/OTU tables were refined

286    taxonomically based on their RDP or BLAST taxonomy. For both assignment methods, unassigned

287    clusters were removed. Non-target 18S and COI clusters (bacterial, non-metazoan) as well as all

14

288    clusters with a terrestrial assignment (taxonomic groups known to be terrestrial-only, such as

289    Insecta, Arachnida, Diplopoda, Amphibia, terrestrial mammals, Stylommatophora, Aves,

290    Onychophora, Succineidae, Cyclophoridae, Diplommatinidae, Megalomastomatidae, Pupinidae,

291    Veronicellidae) were removed. Samples were checked to ensure that a minimum of 10,000

292    metazoan reads were left after refining. Finally, as tag-switching is always to be expected in

293    multiplexed metabarcoding analyses (Schnell et al. 2015), an abundance renormalization was

294    performed to remove spurious positive results due to reads assigned to the wrong sample

295    (Wangensteen    and    Turon    2016,    script    from

296    https://github.com/metabarpark/R_scripts_metabarpark).

297        To test LULU curation (Frøslev et al. 2017), refined 18S and COI ASVs/OTUs were

298    curated with LULU v.0.1 following the online tutorial (https://github.com/tobiasgf/lulu). The

299    LULU algorithm detects erroneous clusters by comparing their sequence similarities and co-

300    occurrence rate with more abundant ("parent") clusters. LULU was tested with a minimum relative

301    co-occurrence of 0.90, using a minimum similarity threshold (*minimum match*) at 84% (default)

302    and slightly higher at 90%, following recommendations of the authors for less variable loci than

303    ITS.

304        The vast majority of prokaryotes usually show low levels (< 1% divergence) of intra

305    genomic variability for the 16S SSU rRNA gene (Acinas et al. 2004; Pei et al. 2010). These low

306    intragenomic divergence levels can be efficiently removed with swarm clustering at d=1. Although

307    LULU curation may still be useful to merge redundant phylotypes in specific cases such as

308    haplotype network analyses, this was not tested in this study. Indeed, parallelization not being

309    currently available for LULU curation, the richness of prokaryote communities implied an

310    unrealistic calculation time, even on a powerful cluster (e.g. LULU curation was at 20-40% after 4

311    days of calculation on our cluster).

312

### 1.4  Statistical analyses

314   Sequence tables were analysed using R with the packages phyloseq v1.22.3 (McMurdie and

315   Holmes 2013) following guidelines on online tutorials (http://joey711.github.io/phyloseq/tutorials-

316   index.html), and vegan v2.5.2 (Oksanen et al. 2018). The datasets were normalized by rarefaction

317   to their common minimum sequencing depth, before analysis of mock communities and natural

318   samples.

319   To evaluate the functionality of the pipeline with the mock communities, taxonomically

320   assigned metazoan clusters were considered as derived from one of the ten species used for the

321   mock communities when the assignment delivered the corresponding species, genus, family, or

322   class. Clusters not fitting the expected taxa were labelled as 'Others'. Apart from PCR errors, these

323   non-target clusters may also originate from contamination by external DNA from associated

324   microfauna, or gut content in the case of whole polyps used for cnidarians.

325   Alpha diversity detected using each pipeline in the natural samples was evaluated with the

326   number of observed target-taxa in the rarefied datasets via analyses of variance (ANOVA) on

327   generalized linear models based on quasipoisson distribution models. Homogeneity of multivariate

328   dispersions were verified with the *betapart* package v.1.5.1 (Baselga and Orme 2012). Beta-

329   diversity patterns were visualised via Principal Coordinates Analyses (PCoA), using Jaccard

330   dissimilarities for metazoans and Bray-Curtis dissimilarities for prokaryotes. The effect of site and

331   LULU curation on community composition was tested by means of PERMANOVA, using the

332   function *adoni*s2 (vegan), with the same dissimilarities as in PCoAs, and permuting 999

333   times.Finally, BLAST and RDP taxonomic assignments of the mock samples and the global dataset

334   were compared at the most adequate pipeline settings for each locus. BLAST-refined (minimum

335   identity at 70%) and RDP-refined (minimum phylum bootstrap at 80%) datasets were compared

16

336      on ASV-level for prokaryotes, and OTU-level for metazoans (swarm *d=3*, LULU at 84% for COI

337      and 90% for 18S). As trials on MIDORI-UNIQUE resulted in very poor performance of RDP for

338      COI (assignments belonging mostly to Insecta), the comparison was performed with MIDORI-

339      UNIQUE subsampled to marine taxa only.

340

341      **2     RESULTS**

342      **2.1     Alpha diversity in mock communities**

343          A number of 2 million (18S) and 1.5 million (COI) raw reads were obtained from the two

344      mock communities (Table S4).After refining, these numbers were decreased to 1.3 million for 18S

345      and 0.7 million for COI.

346          Seven out of ten mock species were recovered in the 18S dataset and all species were

347      detected in the COI dataset (Table 1), even with minimum relative DNA abundance levels as low

348      as 0.7% (Mock 5). Taxonomically unresolved species were correctly assigned up to their common

349      family or class level. Dominant species generally produced more reads in both the clustered and

350      non-clustered datasets (Table S6).

351          When ASVs were clustered with swarm v2, this generally led to a slight loss of taxonomic

352      resolution: *Chorocaris* sp. was not detected in Mock 5 for 18S at d > *1,* and the two bivalves *P.*

353      *kilmeri* and *C. regab* were taxonomically misidentified for COI at d ≥ 1.

354          Clustering sequences with swarm v2 reduced the number of clusters produced per species,

355      but some species still produced multiple OTUs even at *d* values as high as *d=11* for 18S (*A.*

356      *arbuscula*, *Munidopsis* sp., and *E. norvegica*) and *d=13* for COI *D. dianthus, A. muricola,*

357      *Chorocaris* sp*.,* and *Paralepetopsis* sp.). Curating with LULU allowed reducing the number of

358      clusters produced per species to nearly one for both loci, but the best results were obtained in

359      datasets clustered at d > 1 for 18S and d ≥ 1 for COI. Moreover, LULU curation tended to decrease

17

360    the number of non-target clusters ("Others") (Table 1). In the clustered COI dataset, curating with

361    LULU at 84% *minimum match* resulted in the most accurate detection of community composition,

362    and this for all *d* values tested. However, curating with LULU the 18S data (ASVs or OTUs) led

363    to the loss of one shrimp species (*Chorocaris* sp) when the *minimum match* parameter was at 90%

364    and an additional species was lost (the limpet *Paralepetopsis* sp.) when this parameter was at 84%.

365    LULU consistently merged the shrimp species *Chorocaris* sp with another shrimp species as the

366    latter were always co-occurring in our mock samples.

367

18

Table 1. Number of ASVs/OTUs detected per species in the mock communities using different bioinformatic pipelines. White cells indicate an exact match with the number of OTUs expected, grey cells indicate a number of OTUs differing by ±3 from the number expected, and dark grey cells indicate a number of OTUs >3 from the one expected.

| 18S | DADA2 | DADA2+LULU 90% | DADA2+LULU 84% | | DADA2+swarm d1/d3/d4/d5/d11 | DADA2+swarm d1/d3/d4/d5/d11 + LULU 90% | DADA2+swarm d1/d3/d4/d5/d11 + LULU 84% |
|---|---|---|---|---|---|---|---|
| **Mock 3** | | | | | | | |
| Alcyonacea;*A.arbuscula* | 64 | 1 | 1 | Alcyonacea;*A.arbuscula* | 29/11/9/7/6 | 1/1/1/1/1 | 1/1/1/1/1 |
| Caryophylliidae;*D.dianthus* | 2 | 1 | 1 | Caryophylliidae;*D.dianthus* | 2/2/1/1/1 | 1/1/1/1/1 | 1/1/1/1/1 |
| *Alvinocaris muricola* | 2 | 1 | 1 | *Alvinocaris muricola* | 2/1/1/1/1 | 1/1/1/1/1 | 1/1/1/1/1 |
| *Chorocaris* sp. | 1 | 0 | 0 | *Chorocaris* sp. | 2/1/1/1/1 | 0/0/0/0/0 | 0/0/0/0/0 |
| *Munidopsis* sp. | 6 | 1 | 1 | *Munidopsis* sp. | 5/4/3/3/2 | 1/1/1/1/1 | 1/1/1/1/1 |
| Gastropoda;*Paralepetopsis* sp. | 1 | 1 | 0 | Gastropoda;*Paralepetopsis* sp. | 1/1/1/1/1 | 1/1/1/1/1 | 0/0/0/0/0 |
| Vesicomyidae;*P. kilmeri/C. regab/V. gigas* | 8 | 1 | 1 | Bivalvia;*P. kilmeri/C. regab/V. gigas* | 5/4/4/4/2 | 1/2/2/2/1 | 1/1/1/1/1 |
| Polychaeta;*E.norvegica* | 8 | 3 | 2 | Polychaeta;*E.norvegica* | 5/4/4/4/3 | 3/2/2/2/2 | 2/1/2/2/2 |
| Others | 3 | 3 | 2 | Others | 4/4/4/4/4 | 2/2/2/2/3 | 2/2/2/2/2 |
| **Mock 5** | | | | | | | |
| Alcyonacea;*A.arbuscula* | 54 | 1 | 1 | Alcyonacea;*A.arbuscula* | 28/11/9/7/6 | 1/1/1/1/1 | 1/1/1/1/1 |
| Caryophylliidae;*D.dianthus* | 1 | 1 | 1 | Caryophylliidae;*D.dianthus* | 1/1/1/1/1 | 1/1/1/1/1 | 1/1/1/1/1 |
| *Alvinocaris muricola* | 1 | 1 | 1 | *Alvinocaris muricola* | 1/1/1/1/1 | 1/1/1/1/1 | 1/1/1/1/1 |
| *Chorocaris* sp. | 1 | 0 | 0 | *Chorocaris* sp. | 1/0/0/0/0 | 0/0/0/0/0 | 0/0/0/0/0 |
| *Munidopsis* sp. | 4 | 1 | 1 | *Munidopsis* sp. | 4/3/3/3/2 | 1/1/1/1/1 | 1/1/1/1/1 |
| Gastropoda;*Paralepetopsis* sp. | 1 | 1 | 0 | Gastropoda;*Paralepetopsis* sp. | 1/1/1/1/1 | 1/1/1/1/1 | 0/0/0/0/0 |
| Vesicomyidae;*P. kilmeri/C. regab/V. gigas* | 5 | 1 | 1 | Bivalvia;*P. kilmeri/C. regab/V. gigas* | 5/3/3/3/2 | 1/1/1/1/1 | 1/1/1/1/1 |
| Polychaeta;*E.norvegica* | 11 | 3 | 2 | Polychaeta;*E.norvegica* | 5/4/4/4/3 | 3/2/2/2/1 | 2/1/2/2/2 |
| Others | 4 | 3 | 2 | Others | 3/4/4/4/2 | 4/2/2/2/1 | 4/2/2/2/3 |

| COI | DADA2 | DADA2+LULU 90% | DADA2+LULU 84% | | DADA2+swarm d1/d5/d6/d7/d13 | DADA2+swarm d1/d5/d6/d7/d13 + LULU 90% | DADA2+swarm d1/d5/d6/d7/d13 + LULU 84% |
|---|---|---|---|---|---|---|---|
| **Mock 3** | | | | | | | |
| *Acanella arbuscula* | 1 | 1 | 1 | *Acanella arbuscula* | 1/1/1/1/1 | 1/1/1/1/1 | 1/1/1/1/1 |
| Hexacorallia;*D.dianthus* | 3 | 3 | 3 | Hexacorallia;*D.dianthus* | 3/4/4/4/3 | 3/3/3/3/3 | 3/3/3/3/3 |
| *Alvinocaris ;A. muricola* | 26 | 2 | 2 | Alvinocaris;*A. muricola* | 21/12/10/10/5 | 1/1/1/1/1 | 1/1/1/1/1 |
| *Chorocaris* sp. | 2 | 1 | 1 | *Chorocaris* sp. | 3/3/3/3/3 | 1/1/1/1/1 | 1/1/1/1/1 |
| *Munidopsis* sp. | 2 | 1 | 1 | *Munidopsis* sp. | 3/2/1/1/1 | 2/1/1/1/1 | 1/1/1/1/1 |
| Gastropoda;*Paralepetopsis* sp. | 8 | 2 | 3 | Gastropoda;*Paralepetopsi s* sp. | 3/3/3/3/2 | 2/2/2/2/2 | 2/2/2/2/2 |
| *Phreagena kilmeri* | 2 | 1 | 1 | Bivalvia;*P. kilmeri* | 2/3/3/3/3 | 2/2/2/2/2 | 2/2/2/2/2 |
| Bivalvia;*C. regab* | 2 | 1 | 1 | Bivalvia;*C. regab* | | | |
| *Vesicomya gigas* | 1 | 1 | 1 | *Vesicomya gigas* | 1/1/1/1/1 | 1/1/1/1/1 | 1/1/1/1/1 |
| Polychaeta;*E.norvegica* | 3 | 2 | 1 | *Eunice norvegica* | 2/1/1/1/1 | 2/1/1/1/1 | 1/1/1/1/1 |
| Others | 7 | 6 | 6 | Others | 3/3/3/3/4 | 4/5/5/5/5 | 5/5/5/5/5 |
| **Mock 5** | | | | | | | |
| *Acanella arbuscula* | 1 | 1 | 1 | *Acanella arbuscula* | 1/1/1/1/1 | 1/1/1/1/1 | 1/1/1/1/1 |
| Hexacorallia;*D.dianthus* | 3 | 3 | 3 | Hexacorallia;*D.dianthus* | 3/3/3/3/3 | 3/3/3/3/3 | 3/3/3/3/3 |
| *Alvinocaris ;A. muricola* | 26 | 2 | 2 | Alvinocaris;*A. muricola* | 21/12/10/10/5 | 1/1/1/1/1 | 1/1/1/1/1 |
| *Chorocaris* sp. | 1 | 1 | 1 | *Chorocaris* sp. | 2/2/2/2/2 | 1/1/1/1/1 | 1/1/1/1/1 |
| *Munidopsis* sp. | 2 | 1 | 1 | *Munidopsis* sp. | 2/2/1/1/1 | 1/1/1/1/1 | 1/1/1/1/1 |
| Gastropoda;*Paralepetopsis* sp. | 5 | 2 | 2 | Gastropoda;*Paralepetopsis* sp. | 3/2/2/2/2 | 2/2/2/2/2 | 2/2/2/2/2 |
| *Phreagena kilmeri* | 1 | 1 | 1 | Bivalvia;*P. kilmeri* | 2/2/2/2/2 | 2/2/2/2/2 | 2/2/2/2/2 |
| Bivalvia;*C. regab* | 2 | 1 | 1 | Bivalvia;*C. regab* | | | |
| *Vesicomya gigas* | 1 | 1 | 1 | *Vesicomya gigas* | 1/1/1/1/1 | 1/1/1/1/1 | 1/1/1/1/1 |
| Polychaeta;*E.norvegica* | 3 | 2 | 1 | *Eunice norvegica* | 2/2/2/2/2 | 1/1/1/1/1 | 1/1/1/1/1 |
| Others | 6 | 5 | 4 | Others | 2/2/2/2/2 | 1/2/2/2/2 | 1/1/1/1/1 |

368

369

370          .

**2.2   Alpha-diversity patterns in natural samples**

*High-throughput sequencing results*

A number of 44 million (18S), 33 million (COI) and 16 million (16S) reads were obtained

from 42 sediment samples, 4 field controls, 2 extraction blanks, and 4-10 PCR blanks (Table S4).

Two sediment samples failed amplification for the COI marker gene (PCT_FA_CT2_0_1 and

CHR_CT1_0_1). For metazoans, less reads were retained after bioinformatic processing in

negative controls (36% for 18S, 47% for COI) compared to true samples (~60% for 18S, ~70%

for COI), while the opposite was observed for 16S (74% of reads retained in control samples

against 53% in true samples). Negative control samples (field, extraction, and PCR controls)

contained 2,186,230 (~8%) 18S reads, 1,015,700 (~4%) COI reads, and 2,618,729 (28%) 16S

reads. These reads were mostly originating from the field controls for metazoans (48% for 18S,

55% for COI) and extractions controls for 16S (50%).

After blank correction, data refining, and abundance renormalization, rarefaction curves

showed that a plateau was achieved for all samples in both clustered and non-clustered datasets,

suggesting an overall sequencing depth adequate to capture the diversity present (Fig. S1). The

final 18S datasets (with and without clustering at selected $d$ values) contained 8.9-9.6 million

marine metazoan reads in 42 sediment samples (Table S4), and comprised 57,661 ASVs and

19,504-44,948 OTUs (Table S6). The final COI datasets contained 4.5-6.9 million marine

metazoan reads in 40 sediment samples, and comprised 78,785 ASVs and 44,684-64,669 OTUs.

The 16S datasets contained from 6.6 to 6.7 million prokaryotic reads in 42 sediment samples,

producing 56,577 ASVs and 41,746-14,631 OTUs.

392

20

393      *Number of clusters among pipelines*

394      The number of metazoan clusters detected in the deep-sea sediment samples varied

395      significantly between bioinformatic pipelines chosen (, and also varied significantly among sites

396      (Table 2). However, the pipeline effect was consistent across sites although mean cluster numbers

397      detected per sample spanned a wide range in all loci (100-800 for 18S, 150-1,500 for COI datasets,

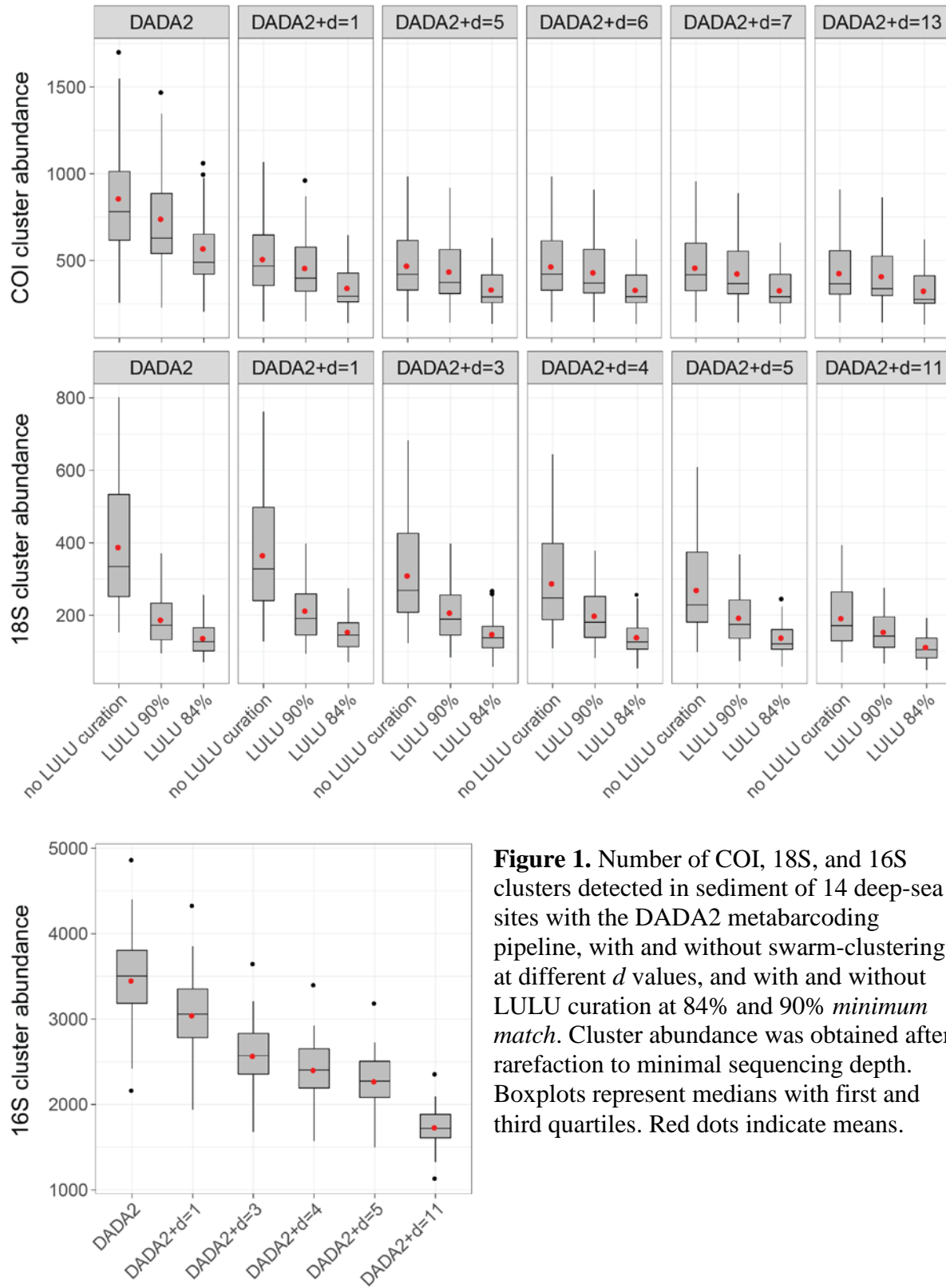398      and 1,500-5,000 for 16S, Fig. 1).

399      Expectedly, clustering significantly reduced the number of detected clusters per sample for

400      all loci. Consistent to results observed in mock communities, clustering at *d=1-13* resulted in

401      comparable OTU numbers for COI, while significantly higher OTU numbers were obtained at *d=1*

402      than with *d >1* for ribosomal loci (Fig. 1, Table 2). DADA2 detected on average 863 (SE=61)

403      metazoan COI ASVs per sample, and clustering reduced this number to around 500, regardless the

404      *d*-value. For ribosomal loci, clustering at *d=3*-5 reduced OTU numbers of around 25-30%

405      compared to without clustering, while at *d=11*, cluster numbers were halved.

406

Table 2. Effect of pipeline and site on the number of metazoan and prokaryote clusters. Results of the analysis of variance (ANOVA) of the rarefied cluster richness for the three genes studied. Pairwise comparisons were performed with Tukey's HSD tests. DS: Dada2+swarm; DSL: Dada2+swarm+LULU; d: swarm *d-value*. Significance codes: ***: p<0.001; **: p<0.01; *: p<0.05.

| LOCUS | F-value | p-value | Significant pairwise comparisons |
|---|---|---|---|
| **COI** | | | |
| Pipeline | 123.13 | p<0.001 | Dada2 > DS***;  DS(d1) > DS(d13)***; |
| Site | 356.37 | p<0.001 | Dada2 > DL***; DS > DSL 84% ***; D(S)L 90% > D(S)L 84% *** |
| Pipeline x Site | 0.16 | p>0.05 | DL > DSL***; DL 90% > DS*** |
| **18S V1-V2** | | | |
| Pipeline | 129.16 | p<0.001 | Dada2 > DS(d>1)***;  DS(d1) > DS(d>1)***; DS(d11) < DS(d1-5)***; |
| Site | 154.52 | p<0.001 | Dada2 > DL***; DS > DSL 84% ***; D(S)L 90% > D(S)L 84% ***; |
| Pipeline x Site | 0.49 | p>0.05 | DL 84% < DS*** |
| **16S V4-V5** | | | |
| Pipeline | 179.19 | p<0.001 | Dada2 > DS***; |
| Site | 18.46 | p<0.001 | DS(d1) > DS(d>1)***; DS(d11) < DS(d1-5)*** |
| Pipeline x Site | 0.06 | p>0.05 | |

407

21

**Figure 1.** Number of COI, 18S, and 16S clusters detected in sediment of 14 deep-sea sites with the DADA2 metabarcoding pipeline, with and without swarm-clustering at different *d* values, and with and without LULU curation at 84% and 90% *minimum match*. Cluster abundance was obtained after rarefaction to minimal sequencing depth. Boxplots represent medians with first and third quartiles. Red dots indicate means.

408

409       LULU curation of metazoan ASVs significantly decreased the number of clusters detected

410       at both tested *minimum match* values (Table 2). For OTU datasets, the decrease was significant

411       only when the *minimum match* parameter was at 84%. The effect of LULU curation was stronger

412       at a lower *minimum match* value for both loci, as LULU curation at 90% of ASVs or OTUs resulted

413       in significantly more clusters than when the minimum match was at 84% (Table 2). The effect of

414       LULU curation of was also more pronounced for the 18S locus: LULU decreased by 31-65% the

415       number of 18S ASVs/OTUs, compared to 7-33% for COI. LULU curation of ASVs or OTUs

416       resulted in comparable cluster numbers in the 18S datasets, regardless the *d*-value used for

417       clustering. For example, at 84% *minimum match*, LULU curation produced on average $137 \pm 7$ and

418       $140 \pm 8$ clusters per sample after application on ASVs and OTUs (*d=4*) respectively. At 90%, these

419       numbers were at $189 \pm 11$ and $200 \pm 12$ (Fig. 1). This was not the case for COI, where LULU

420       curation of ASVs resulted in significantly more clusters ($574 \pm 38$ at 84% and $742 \pm 53$ at 90%)

421       than LULU curation of OTUs ($334 \pm 21$ and $433 \pm 31$ for *d=6*).

422       Looking at mean ASV and OTU numbers detected per phylum with each pipeline showed

423       consistent effects of swarm clustering and LULU curation, but highlighted strong differences in

424       the amount of intragenomic variation between taxonomic groups. For all loci investigated, some

425       taxa displayed high ASV to OTU ratios, while others were hardly affected by clustering or LULU

426       curation in terms of numbers of clusters detected (Fig S2).

427

428    **2.3    Patterns of beta-diversity between pipelines**

429       Community differences were visualized using PCoA ordinations (Jaccard and Bray-Curtis

430       dissimilarities for metazoans and prokaryotes respectively) in clustered and non-clustered datasets

431       (Fig. 2, Fig. S3). Expectedly, PERMANOVAs confirmed that sites differed significantly in terms

432       of community structure, accounting from 45% to 89% of variation in data. Evaluating the effect of

433    LULU curation (at 84% and 90%) for metazoans showed that LULU-curated data resolved similar

434    ecological patterns than non-curated data, accounting from 0.5% (COI) to 1.3% (18S) of variation

435    in data (Fig. 2).

436         Although ASV and OTU datasets detected similar levels of variation due to sites in

437    PERMANOVAs, clustering levels affected the ecological patterns resolved by ordinations in rRNA

438    loci (Fig 2). At low $d$ values ($d=1-3$), ecological patterns were consistent to patterns observed in

439    the ASV datasets, with samples segregating by site and depth. Increasing $d$ values produced

440    stronger segregation among sites, thus resulting in differentiation among ocean basins rather than

441    depth. This change in resolution occurred with $d$ values as low as $d=4$ for 18S but was strongest at

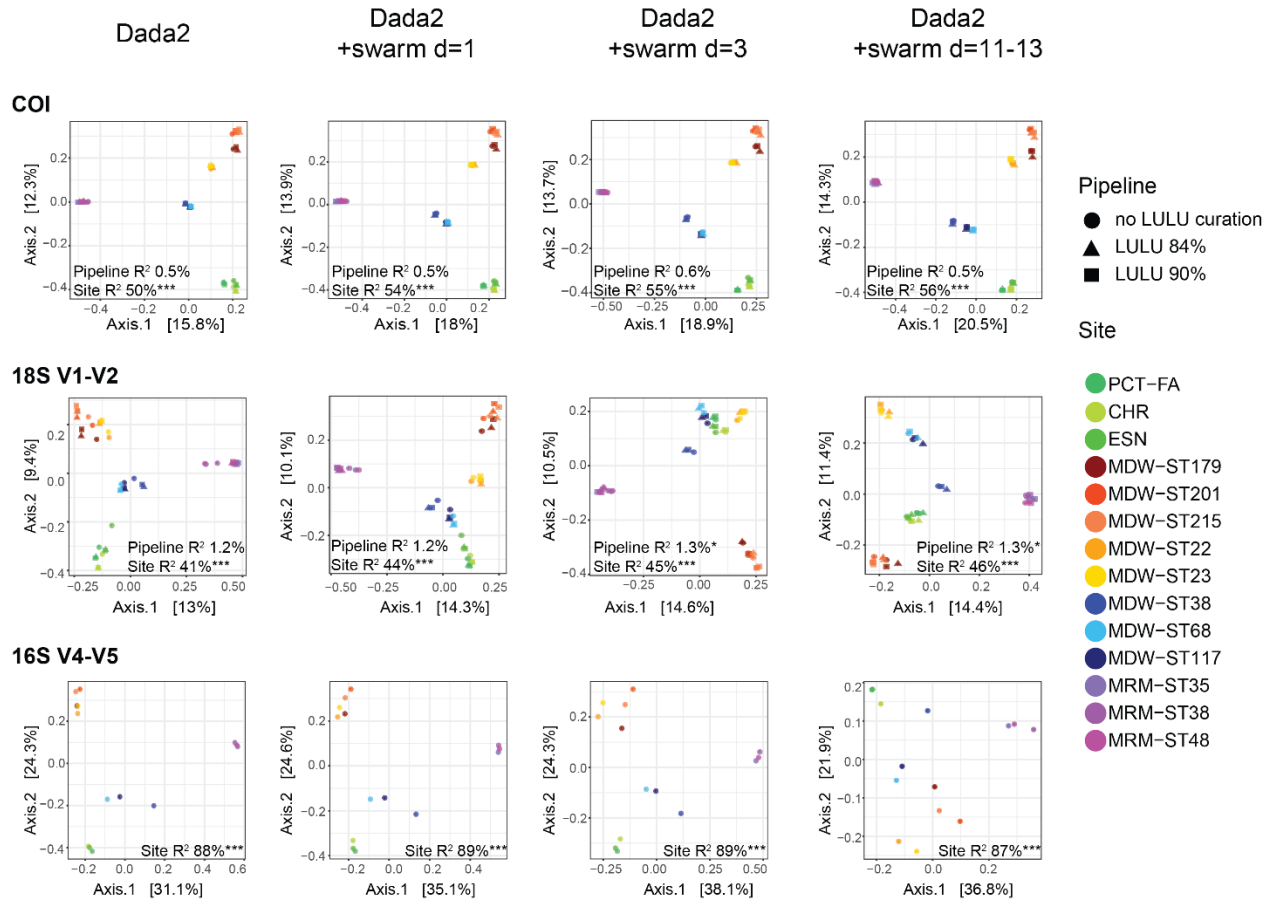442    $d=11$ for both rRNA loci (Fig. S3, Fig. 2).

443

Figure 2. Beta-diversity patterns in ASV and OTU-centred datasets. PCoA ordinations showing community differentiation observed between sites and LULU *vs* not LULU curated samples, for the DADA2 metabarcoding pipeline with and without clustering. Metazoan datasets were clustered at *d=1-13* (COI) *d=1-11* (18S) and curated with LULU at two minimum match values. The prokaryote 16S dataset was clustered at *d=1-11*. $R^2$ values and associated p-values obtained in PERMANOVAs are shown in the ordination plots. Significance codes: ***: p<0.001; **: p<0.01; *: p<0.05. Colour codes: Green: Mediterranean < 1,000 m; Red-yellow: Mediterranean-Atlantic transition zone 300-1,000 m; Blue: North Atlantic < 1,000 m; Purple: Arctic < 1,000 m.

444

445

446

25

447     **2.4    Taxonomic assignment quality**

448          BLAST and RDP Bayesian Classifier assignments were compared in the mock

449     communities and natural samples, on data clustered at *d=3* and curated with LULU at 84% for COI

450     and 90% for 18S. For prokaryotes, assignment methods were compared on the ASV-level. BLAST

451     and RDP assigned similar amounts of OTUs in the prokaryote dataset, but BLAST assigned 20-

452     70% less OTUs in the metazoan datasets (Table S7). Assigning with BLAST at a minimum of 70%

453     hit identity resulted in comparable results as described above. Eight of the ten species were

454     recovered with COI and six species were recovered with 18S, while the vesicomyid bivalves were

455     taxonomically unresolved with both loci (Fig. S4). Although most species produced one single

456     OTU, between one and three species still resulted in 2-3 OTUs in each mock sample. Assigning

457     the 18S dataset with RDP resulted in comparable taxonomic resolutions, although more species

458     produced more than one OTU. Assigning the COI dataset with RDP using the MIDORI-UNIQUE

459     database resulted in assignments of the mock samples that did not match the expected taxa and

460     were mostly belonging to arthropods, a problem not observed with BLAST (data not shown). When

461     the database was reduced to marine-only taxa, all 10 species were detected, and this at expected

462     OTU abundances, once data was filtered for phylum bootstrap levels $\geq$ 80% (Fig S4). However,

463     applying a phylum bootstrap minimum of 80% resulted in a strong decrease in the number of final

464     target OTUs, particularly for COI where only 226 OTUs remained after filtering (Table S7). This

465     reduced recovery with RDP after applying a minimum phylum bootstrap level was not observed in

466     prokaryotes, where 51,000-55,000 ASVs were left after filtering with both assignment methods

467     (Table S7).

468          BLAST hit identities of the overall datasets varied strongly depending on phyla and

469     marker gene (Fig. 3). For 18S, most clusters had hit identities $\geq$ 90%. Poorly assigned clusters

470     (hit identity < 90%) represented less than 20% of the dataset and were mostly assigned to

26

471    Nematoda, Cnidaria, Tardigrada, Porifera, and Xenacoelomorpha. For COI, nearly all clusters

472    had similarities to sequences in databases lower than 90%. Overall, arthropods and echinoderms

473    were detected at similar levels by both markers. The 18S barcode marker performed better in the

474    detection of nematodes, annelids, platyhelminths, and xenacoelomorphs while COI mostly

475    detected cnidarians, molluscs, and poriferans (Fig. 3), highlighting the complementarity of these

476    two loci. BLAST hit identity was much higher for prokaryotes, with most clusters assigned with

477    more than 90% similarity to sequences in databases. When datasets were filtered for RDP

478    phylum bootstrap levels ≥ 80%, most assignments also had high genus bootstrap values for

479    ribosomal loci. However, for COI, a considerable number of OTUs assigned to arthropods,

480    cnidarians, molluscs, vertebrates, and poriferans still had genus bootstraps < 60%.
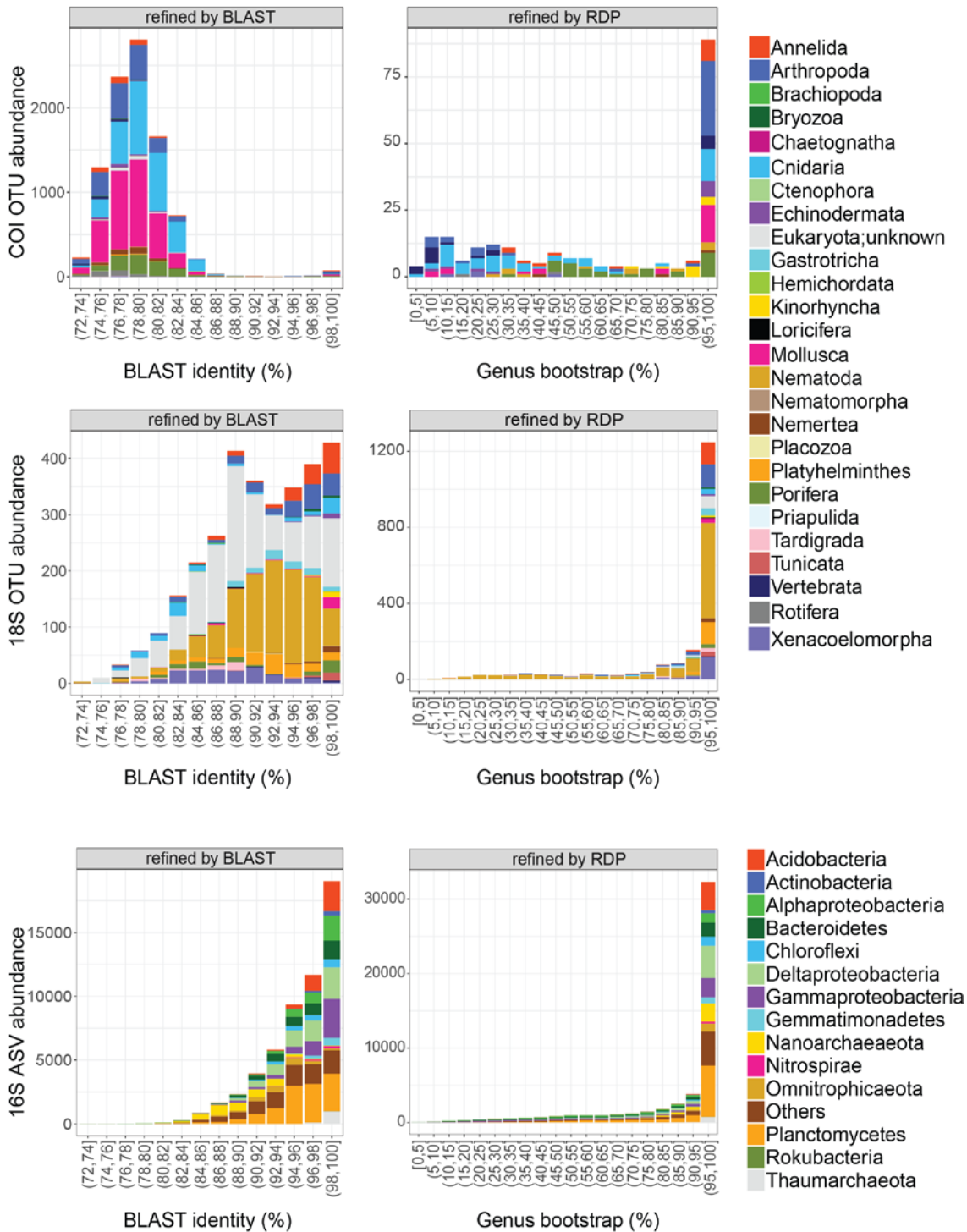
481

482

483



Figure 3. Taxonomic assignment quality of BLAST and RDP methods on metazoan and prokaryote metabarcoding datasets of 14 deep-sea sites. BLAST hit identity of all target clusters detected is given at hit identities > 70%. RDP-assigned data was filtered for phylum bootstraps ≥ 80%, and associated genus bootstraps are displayed. Taxonomic assignments were performed on the Silva132 database for 18S and 16S, and on the MIDORI-UNIQUE database, subsampled to marine taxa for COI.

486   **3    DISCUSSION**

487   **3.1    ASVs and OTUs for genetic *vs* species diversity**

488          The rise of HTS and the subsequent use of DNA metabarcoding have revolutionized

489   microbiology by unlocking the access to uncultivable microorganisms, which represent by far the

490   great majority of prokaryotes (Klappenbach et al. 2001). The development and improvement of

491   molecular and bioinformatic methods to perform inventories were historically primarily developed

492   for 16S rRNA barcode loci, before being transferred to the eukaryotic kingdom based on the use

493   of barcode markers such as 28S and 18S rRNA, ITS, or mitochondrial markers such as COI

494   (Valentini et al. 2009; Bellemain et al. 2010). Thus, most bioinformatics pipelines were initially

495   developed accounting for intrinsic properties of prokaryotes and concepts inherent to microbiology

496   (Caporaso et al. 2010; Schloss et al. 2009; Boyer et al. 2016), before being transferred to eukaryotes

497   in general, or metazoans in particular. Such application transfers require adaptations to account for

498   differences in both concepts and basic biological features. One example is the question of the

499   relevance of using ASVs, advocated to replace OTUs "*as the standard unit of marker-gene analysis*

500   *and reporting*" (Callahan et al., 2017): an advice for microbiologists that may not apply when

501   working on metazoans.

502          First, metazoans are well known to exhibit variable and sometimes very high intraspecific

503   polymorphism in 18S-V1 and above all in COI. Second, the results on the mock samples showed

504   that single individuals produced very different numbers of ASVs, indicating that ASV-centred

505   datasets do not reflect actual species composition in metazoans. As this "demultiplication" will be

506   highly variable across taxa (as seen in Fig. S2, and references such as Plouviez et al. 2009 and

507   Teixeira et al. 2013), the taxonomic compositions of samples based on ASVs will reflect genetic

508   rather than species diversity.

509      Clustering ASVs into OTUs and/or curating with LULU alleviated the numerical inflation,

510      but some species still produced more than one OTU, even at *d*-values as high as *d=11-13*. While

511      clustering and LULU curation improved numerical results in the mock communities, they were

512      associated with a decrease in taxonomic resolution, especially for 18S where some closely related

513      species were merged with increasing clustering/filtering thresholds (i.e. the vesicomyid bivalves,

514      the gastropod, and the shrimp species; Table 1). When studying natural habitats, very likely to

515      harbour closely related co-occurring species, both LULU curation and clustering are thus likely to

516      lead to the loss of true species diversity, particularly for low-resolution markers such as 18S.

517      Optimal results in the mock samples, i.e. delivering the best balance between the limitation of

518      spurious clusters and the loss of true species diversity, were obtained with LULU curation at 90%

519      for 18S and 84% for COI, highlighting the importance of adjusting bioinformatic correction tools

520      to each barcode marker, a step for which mock communities are most adequate.

521

522   **3.2   ASVs *vs* OTUs in natural communities: adapting pipeline parameters to marker**

523         **properties**

524       Life histories of organisms, together with intrinsic properties of marker genes, determine

525      the level of intragenomic and intraspecific diversity. Intraspecific variation is a recognised problem

526      in metabarcoding, known to generate spurious clusters (Brown et al. 2015), especially in the COI

527      barcode marker. Indeed, this gene region has increased intragenomic variation due to its high

528      evolutionary rate but also due to heteroplasmy and the abundance of pseudogenes, such as NUMTs,

529      playing an important part of the supernumerary OTU richness in COI-metabarcoding (Bensasson

530      et al. 2001; Song et al. 2008). Together with clustering, LULU curation at 84% proved effective in

531      limiting the number of multiple clusters produced by single individuals, confirming its efficiency

532      to correct for intragenomic diversity (Table 1).

533  The mock communities we used here did not contain several haplotypes of the same species

534 (intraspecific variation), as is most often the case in environmental samples. This prevents us from

535 generalizing the comparable results obtained after LULU curation of ASVs and OTUs, and the

536 apparently limited effect of clustering in the mock samples to communities that are more complex.

537 However, LULU curation of ASVs is not suited to account for natural haplotype diversity: not all

538 haplotypes co-occur and when they do so, they may vary in proportion and dominance

539 relationships, making clustering more suited to account for natural haplotypic diversity. Thus,

540 clustering ASVs will still be necessary to produce inventories of metazoan communities that reflect

541 species rather than gene diversity.

542  As expected, evaluation of clustering and LULU curation on natural samples showed

543 distinct results for 18S and COI. Indeed, concerted evolution, a common feature of SSU rRNA

544 markers such as 16S (Hashimoto et al. 2003; Klappenbach et al. 2001) and 18S (Carranza et al.

545 1996), limits the amount of intragenomic polymorphism. In metazoans, a lower level of diversity

546 is expected for the slower evolving 18S gene (Carranza et al. 1996), than for COI which exhibits

547 faster evolutionary rates (Machida and Knowlton 2012; Machida et al. 2012). This is reflected in

548 the lower ASV (DADA2) to OTU (DADA2+swarm) ratios observed here for 18S (1.0-2.2.)

549 compared to COI (2.0-2.7) data at clustering $d$-values comprised between one and seven (Table

550 S6), underlining the different influence –and importance– of clustering on these loci, and the need

551 for a versatile, marker by marker choice for clustering and curation parameters. When applying

552 LULU to ASVs (DADA2) *versus* OTUs (DADA2+swarm) on 18S, similar cluster numbers were

553 obtained (Fig. 1), suggesting a limited added effect of clustering for this marker once DADA2 and

554 LULU are applied. This is in line with its slow evolutionary rate (Carranza et al. 1996) leading to

555 a limited number of haplotypes per species compared to COI. In contrast, for COI, LULU curation

556 of the ASV dataset led to nearly twice the number of clusters ($574 \pm 38$ at 84% and $742 \pm 53$ at

31

557  90%) compared to LULU curation of OTUs (at $d=6$: 334 ± 21 for 84% and 433 ± 31 for 90%).

558  This confirms the higher intraspecific diversity of COI, and the need to combine clustering with

559  LULU curation to account for intraspecific diversity in natural samples, especially with highly

560  polymorphic markers such as COI.

561      Finally, the reproductive mode and pace of selection in microbial populations may lead to

562  locally lower levels of intraspecific variation than the one expected for metazoans. Prokaryotic

563  alpha diversity was however also affected by the clustering of ASVs (Fig. 1), supporting the

564  estimation of a 2.5-fold greater number of 16S rRNA variants than the actual number of bacterial

565  "species" (Acinas et al. 2004). The significant decrease in the number of OTUs after clustering at

566  $d=1$ (Table 2, Fig. 1, decrease of ~25%) suggests the occurrence of very closely related 16S rRNA

567  sequences, possibly belonging to the same ecotype/species. Such entities may still be important to

568  delineate in studies aiming for example at identifying species associations (i.e. symbiotic

569  relationships) across large distances and ecosystems, where drift or selection can lead to slightly

570  different ASVs in space and time, with their function and association remaining stable.

571

572  **3.3    Influence on beta diversity**

573      After focusing on alpha diversity estimates, i.e. on the numerical accuracy of inventories,

574  the analysis of community structures showed that the LULU-curated datasets resolved similar

575  ecological patterns as datasets not curated with LULU. However, clustering affected resolution of

576  ecological patterns in ribosomal loci when $d$ values were high, and this was not the case for COI,

577  where similar patterns were resolved in all datasets (Fig. 2). This is in accordance with other studies

578  reporting severe impacts of bioinformatic parameters on alpha diversity while comparable patterns

579  of beta diversity are observed in ASV and OTU datasets, at least down to a minimum level of

580  clustering stringency (Xiong and Zhan 2018; Bokulich et al. 2013).

32

581        Clustering and LULU curation mainly led to the decrease of the number of clusters assigned

582    to particular taxa in both loci, such as annelids, arthropods, nematodes, or platyhelminthes for 18S,

583    and chordates, cnidarians, echinoderms, or poriferans for COI (Fig. S2). The strong decrease in

584    cluster numbers observed in these phyla suggests that the latter have greater intraspecific

585    polymorphism, although the decrease could also be due to the merging of closely related species,

586    as both markers have lower taxonomic resolution in particular taxa. This has been acknowledged

587    for 18S in general, but in nematodes in particular (Derycke et al. 2010), and reported in cnidarians

588    with COI (Hebert et al. 2003).

589        Overall, based on alpha and beta diversity results observed in mock communities and

590    natural samples, applying LULU at 84% seems to efficiently curate metazoan COI datasets without

591    significant loss of species, but clustering is required, at least at $d=1$, in order to address high

592    intraspecific polymorphism. For 18S, LULU curation seems to require values above 84% (e.g.

593    90%) in order to avoid the loss of species, as seen in the mock communities. However, the low

594    taxonomic resolution obtained with this marker suggests that clustering should be performed at low

595    $d$-values ($d<4$) to address intraspecific polymorphism without affecting beta-diversity patterns. For

596    prokaryotes, clustering 16S ASVs at $d=1$ reduces the number of detected clusters by ~25%, which

597    may help addressing intragenomic variation when needed.

598

599    **3.4    Taxonomic resolution and assignment quality**

600        The COI locus allowed the detection of all ten species present in the mock samples,

601    compared to seven in the 18S dataset (Table 1). This locus also provided much more accurate

602    assignments, most of them correct at the genus (and species) level, confirming that COI uncovers

603    more metazoan species and offers a better taxonomic resolution than 18S (Tang et al. 2012; Clarke

604    et al. 2017; Andújar et al. 2018). Our results also support approaches combining nuclear and

33

605   mitochondrial markers to achieve more comprehensive biodiversity inventories (Cowart et al.

606   2015; Drummond et al. 2015; Zhan et al. 2014). Indeed, strong differences exist in amplification

607   success among taxa (Bhadury et al. 2006; Carugati et al. 2015), exemplified by nematodes, which

608   are well detected with 18S but not with COI (Bucklin et al. 2011). The high complementarity of

609   COI and 18S in terms of targeted taxa (highlighted in Fig. S2), also supports the approach taken

610   by Stefanni et al. (2018), as subsampling each gene dataset for its "best targeted phyla" and

611   subsequently combining both seems to be a very efficient way to produce comprehensive and non-

612   redundant biodiversity inventories.

613        Finally, compared to BLAST assignments, similar taxonomic resolution was observed

614   using the RDP Bayesian Classifier on the mock samples for 18S (Fig. S4) and for COI when using

615   the MIDORI-UNIQUE marine-only database. Poor performance of RDP using the full MIDORI

616   database is likely due to the size of the database, and to its low coverage of deep-sea species. The

617   problem of underrepresentation of deep-sea taxa is especially apparent with the BLAST

618   assignments, which generally displayed low identities to sequences in databases, especially for COI

619   (Fig. 3). Using minimum similarities of 80% for COI and 86% for 18S as cut-off values for

620   metazoans has been shown to improve the taxonomic quality of metazoan metabarcoding datasets

621   (Stefanni et al. 2018). However, phylogenies of marine invertebrates have found high levels of

622   species divergence (up to ~30%), even within genera (Zanol et al. 2010). Consequently, studies on

623   deep-sea taxa have found that some invertebrate species had COI sequences diverging more than

624   20% from any other species present in molecular databases (Shank et al. 1999; Herrera et al. 2015).

625   At present, it thus seems difficult to work at taxonomic levels beyond phylum-level with deep-sea

626   metabarcoding data when using large public databases. Small databases, taxonomically similar to

627   the targeted communities, and with sequences of the same length as the amplified fragment of

628   interest, are known to maximise accurate identification (Macheriotou et al. 2019). When using the

34

629    reduced marine-only COI database, RDP provided the most accurate assignments in the mock

630    samples when the phylum bootstrap level was $\geq 80$ (Fig. S 4), although this filtering threshold

631    drastically reduced the number of OTUs in the overall dataset (Table S7). The development of

632    custom-built marine RDP training sets, without overrepresentation of terrestrial species, is

633    therefore needed for this Bayesian assignment method to be effective on deep-sea datasets. With

634    reduced and more specific databases, removing clusters with phylum bootstrap-level $< 80$ should

635    be an efficient way to increase taxonomic quality of deep-sea metabarcoding datasets. At present,

636    if accurate taxonomic assignments are sought while using universal primers, we advocate assigning

637    taxonomy in two steps: first, using BLAST and a large database including all phyla amplifiable by

638    the primer set, extracting the clusters belonging to the groups of interest, then re-assigning

639    taxonomy to these target taxa using RDP and a smaller, taxon-specific database.

640

641    **CONCLUSIONS AND PERSPECTIVES**

642         Using mock communities and natural samples, we evaluate several recent algorithms and

643    assess their capacity to improve the quality of molecular biodiversity inventories of metazoans and

644    prokaryotes. Our results support the fact that ASV data should be produced and communicated for

645    reusability and reproducibility following the recommendations of Callahan et al. (2017). This is

646    especially useful in large projects spanning wide geographic zones and time scales, as different

647    ASV datasets can be easily merged *a posteriori,* and clustered if necessary afterwards.

648    Nevertheless, clustering ASVs into OTUs will be required to obtain accurate species-level

649    inventories, at least for metazoan communities, with a more severe influence of clustering observed

650    on alpha diversity estimates than beta-diversity patterns. Considering 16S polymorphism observed

651    in prokaryotic species (Acinas et al. 2004) and the possible geographic segregation of their

652    populations, clustering may also be required in prokaryotic datasets, for example in studies

653     screening for species associations (i.e. symbiotic relationships) as symbionts may be prone to

654     differential fixation through enhanced drift (Shapiro, Leducq, & Mallet, 2016).

655     Our results also demonstrated that LULU effectively curates metazoan biodiversity

656     inventories obtained through metabarcoding. They also underline the need to adapt parameters for

657     curation (e.g. LULU curation at 90% for 18S and 84% for COI) and clustering to each gene used

658     and taxonomic compartment targeted, in order to identify an optimal balance between the

659     correction for spurious clusters and the merging of closely related species.

660     Finally, our findings also showed that accurate taxonomic assignments of deep-sea species

661     can be obtained with the RDP Bayesian Classifier, but only with reduced databases containing

662     ecosystem-specific sequences.

663     The pipeline is publicly available on Gitlab (https://gitlab.ifremer.fr/abyss-project/), and

664     allows the use of sequence data obtained from libraries produced by double PCR or adaptor ligation

665     methods, as well as having built-in options for using six commonly used metabarcoding primers.

666

667

**668 ACKNOWLEDGEMENTS**

36

677     MEDWAVES (Covadonga Orejas) cruises. The MEDWAVES cruise was organised in the

678     framework of the ATLAS Project, supported by the European Union 2020 Research and

679     Innovation Programme under grant agreement number: 678760. We thank Stefaniya Kamenova,

680     Tiago Pereira, and an anonymous referee for their comments on a previous version of this

681     manuscript.

682

683     **REFERENCES**

684     Acinas, Silvia G, Luisa A Marcelino, Vanja Klepac-Ceraj, and Martin F Polz. 2004. 'Divergence

685         and Redundancy of 16S RRNA Sequences in Genomes with Multiple Rrn Operons'.

686         *Journal of Bacteriology* 186 (9): 2629–35. https://doi.org/10.1128/JB.186.9.2629-

687         2635.2004.

688     Alberdi, Antton, Ostaizka Aizpurua, M. Thomas P. Gilbert, and Kristine Bohmann. 2017.

689         'Scrutinizing Key Steps for Reliable Metabarcoding of Environmental Samples'. Edited by

690         Andrew Mahon. *Methods in Ecology and Evolution*, 2017. https://doi.org/10.1111/2041-

691         210X.12849.

692     Andújar, Carmelo, Paula Arribas, Douglas W. Yu, Alfried P. Vogler, and Brent C. Emerson.

693         2018. 'Why the COI Barcode Should Be the Community DNA Metabarcode for the

694         Metazoa'. *Molecular Ecology* 27 (20): 3968–75. https://doi.org/10.1111/mec.14844.

695     Baselga, Andrés, and C. David L. Orme. 2012. 'Betapart : An R Package for the Study of Beta

696         Diversity'. *Methods in Ecology and Evolution* 3 (5): 808–12. https://doi.org/10.1111/j.2041-

697         210X.2012.00224.x.

698     Bazin, Eric, Sylvain Glémin, and Nicolas Galtier. 2006. 'Population Size Does Not Influence

699        Mitochondrial Genetic Diversity in Animals'. *Science* 312 (5773): 570–72.

700        https://doi.org/10.1126/science.1122033.

701    Bellemain, Eva, Tor Carlsen, Christian Brochmann, Eric Coissac, Pierre Taberlet, and Håvard

702        Kauserud. 2010. 'ITS as an Environmental DNA Barcode for Fungi: An in Silico Approach

703        Reveals Potential PCR Biases'. *BMC Microbiology* 10 (July): 189.

704        https://doi.org/10.1186/1471-2180-10-189.

705    Bensasson, Douda, De Xing Zhang, Daniel L. Hartl, and Godfrey M. Hewitt. 2001.

706        'Mitochondrial Pseudogenes: Evolution's Misplaced Witnesses'. *Trends in Ecology and*

707        *Evolution*. https://doi.org/10.1016/S0169-5347(01)02151-6.

708    Bhadury, P, M C Austen, D T Bilton, P J D Lambshead, A D Rogers, and G R Smerdon. 2006.

709        'Molecular Detection of Marine Nematodes from Environmental Samples: Overcoming

710        Eukaryotic Interference'. *Aquatic Microbial Ecology* 44 (1): 97–103. https://doi.org/Doi

711        10.3354/Ame044097.

712    Bik, Holly M., Way Sung, Paul De Ley, James G Baldwin, Jyotsna Sharma, Axayácatl Rocha-

713        Olivares, and W Kelley Thomas. 2012. 'Metagenetic Community Analysis of Microbial

714        Eukaryotes Illuminates Biogeographic Patterns in Deep-Sea and Shallow Water Sediments.'

715        *Molecular Ecology* 21 (5): 1048–59. https://doi.org/10.1111/j.1365-294X.2011.05297.x.

716    Bista, Iliana, G Carvalho, K Walsh, M Christmas, Mehrdad Hajibabaei, P Kille, D Lallias, and

717        Simon Creer. 2015. 'Monitoring Lake Ecosystem Health Using Metabarcoding of

718        Environmental DNA: Temporal Persistence and Ecological Relevance'. *Genome* 58 (5):

719        197.

720    Bokulich, Nicholas A, Sathish Subramanian, Jeremiah J Faith, Dirk Gevers, Jeffrey I Gordon,

721        Rob Knight, David A Mills, and J Gregory Caporaso. 2013. 'Quality-Filtering Vastly

722        Improves Diversity Estimates from Illumina Amplicon Sequencing'. *Nature Methods* 10 (1):

723        57–59. https://doi.org/10.1038/nmeth.2276.

724    Boussarie, Germain, Judith Bakker, Owen S. Wangensteen, Stefano Mariani, Lucas Bonnin, Jean

725        Baptiste Juhel, Jeremy J. Kiszka, et al. 2018. 'Environmental DNA Illuminates the Dark

726        Diversity of Sharks'. *Science Advances* 4 (5): eaap9661.

727        https://doi.org/10.1126/sciadv.aap9661.

728    Boyer, F, C Mercier, A Bonin, Y Le Bras, Pierre Taberlet, and Eric Coissac. 2016. 'OBITOOLS:

729        A UNIX-Inspired Software Package for DNA Metabarcoding'. *Molecular Ecology*

730        *Resources* 16 (1): 176–82. https://doi.org/10.1111/1755-0998.12428.

731    Brannock, P M, and K M Halanych. 2015. 'Meiofaunal Community Analysis by High-

732        Throughput Sequencing: Comparison of Extraction, Quality Filtering, and Clustering

733        Methods'. *Marine Genomics* 23: 67–75. https://doi.org/10.1016/j.margen.2015.05.007.

734    Brown, E A, F J J Chain, T J Crease, H J MacIsaac, and M E Cristescu. 2015. 'Divergence

735        Thresholds and Divergent Biodiversity Estimates: Can Metabarcoding Reliably Describe

736        Zooplankton Communities?' *Ecology and Evolution* 5 (11): 2234–51.

737        https://doi.org/10.1002/ece3.1485.

738    Bucklin, Ann, Dirk Steinke, and Leocadio Blanco-Bercial. 2011. 'DNA Barcoding of Marine

739        Metazoa'. *Annual Review of Marine Science* 3 (1): 471–508.

740        https://doi.org/10.1146/annurev-marine-120308-080950.

741    Callahan, Benjamin J., Paul J. McMurdie, Michael J. Rosen, Andrew W. Han, Amy Jo A.

742        Johnson, and Susan P. Holmes. 2016. 'DADA2: High-Resolution Sample Inference from

743        Illumina Amplicon Data'. *Nature Methods* 13 (7): 581–83.

744        https://doi.org/10.1038/nmeth.3869.

745    Callahan, Benjamin J., Paul J McMurdie, and Susan P Holmes. 2017. 'Exact Sequence Variants

746        Should Replace Operational Taxonomic Units in Marker-Gene Data Analysis'. *ISME*

747　　　　*Journal* 11 (12): 2639–43. https://doi.org/10.1038/ismej.2017.119.

748　Caporaso, J Gregory, Justin Kuczynski, Jesse Stombaugh, Kyle Bittinger, Frederic D Bushman,

749　　　　Elizabeth K Costello, Noah Fierer, et al. 2010. 'QIIME Allows Analysis of High-

750　　　　Throughput Community Sequencing Data'. *Nature Methods* 7 (5): 335–36.

751　　　　https://doi.org/10.1038/nmeth.f.303.

752　Carranza, Salvador, Gonzalo Giribet, Carles Ribera, Jaume Baguñà, and Marta Riutort. 1996.

753　　　　'Evidence That Two Types of 18S RDNA Coexist in the Genome of Dugesia (Schmidtea)

754　　　　Mediterranea (Platyhelminthes, Turbellaria, Tricladida)'. *Molecular Biology and Evolution*

755　　　　13 (6): 824–32. https://doi.org/10.1093/oxfordjournals.molbev.a025643.

756　Carugati, Laura, Cinzia Corinaldesi, Antonio Dell'Anno, and Roberto Danovaro. 2015.

757　　　　'Metagenetic Tools for the Census of Marine Meiofaunal Biodiversity: An Overview'.

758　　　　*Marine Genomics* 24 (December): 11–20. https://doi.org/10.1016/j.margen.2015.04.010.

759　Clare, Elizabeth L., Frédéric J.J. Chain, Joanne E. Littlefair, and Melania E. Cristescu. 2016. 'The

760　　　　Effects of Parameter Choice on Defining Molecular Operational Taxonomic Units and

761　　　　Resulting Ecological Analyses of Metabarcoding Data'. Edited by Kristy Deiner. *Genome*

762　　　　59 (11): 981–90. https://doi.org/10.1139/gen-2015-0184.

763　Clarke, Laurence J., Jason M. Beard, Kerrie M. Swadling, and Bruce E. Deagle. 2017. 'Effect of

764　　　　Marker Choice and Thermal Cycling Protocol on Zooplankton DNA Metabarcoding

765　　　　Studies'. *Ecology and Evolution* 7 (3): 873–83. https://doi.org/10.1002/ece3.2667.

766　Cohan, Frederick M. 2001. 'Bacterial Species and Speciation'. Edited by M. Kane. *Systematic*

767　　　　*Biology* 50 (4): 513–24. https://doi.org/10.1080/10635150118398.

768　Coissac, Eric, Tiayyba Riaz, and Nicolas Puillandre. 2012. 'Bioinformatic Challenges for DNA

769　　　　Metabarcoding of Plants and Animals'. *Molecular Ecology* 21 (8): 1834–47.

770　　　　https://doi.org/10.1111/j.1365-294X.2012.05550.x.

771     Cowart, Dominique A, Miguel Pinheiro, Olivier Mouchel, Marion Maguer, Jacques Grall,

772         Jacques Miné, and Sophie Arnaud-Haond. 2015. 'Metabarcoding Is Powerful yet Still Blind:

773         A Comparative Analysis of Morphological and Molecular Surveys of Seagrass

774         Communities'. *PLoS One* 10 (2): e0117562. https://doi.org/10.1371/journal.pone.0117562.

775     Creer, Simon, Kristy Deiner, Serita Frey, Dorota Porazinska, Pierre Taberlet, W. Kelley Thomas,

776         Caitlin Potter, and Holly M. Bik. 2016. 'The Ecologist's Field Guide to Sequence-Based

777         Identification of Biodiversity'. Edited by Robert Freckleton. *Methods in Ecology and*

778         *Evolution* 7 (9): 1008–18. https://doi.org/10.1111/2041-210X.12574.

779     Davis, Nicole M., Diana M. Proctor, Susan P. Holmes, David A. Relman, and Benjamin J.

780         Callahan. 2018. 'Simple Statistical Identification and Removal of Contaminant Sequences in

781         Marker-Gene and Metagenomics Data'. *Microbiome* 6 (1): 226.

782         https://doi.org/10.1186/s40168-018-0605-2.

783     Deiner, Kristy, Emanuel A. Fronhofer, Elvira Mächler, Jean Claude Walser, and Florian

784         Altermatt. 2016. 'Environmental DNA Reveals That Rivers Are Conveyer Belts of

785         Biodiversity Information'. *Nature Communications* 7 (1): 12544.

786         https://doi.org/10.1038/ncomms12544.

787     Deiner, Kristy, Jean-Claude C Walser, E Machler, Florian Altermatt, Elvira Mächler, and Florian

788         Altermatt. 2015. 'Choice of Capture and Extraction Methods Affect Detection of Freshwater

789         Biodiversity from Environmental DNA'. *Biological Conservation* 183 (March): 53–63.

790         https://doi.org/10.1016/j.biocon.2014.11.018.

791     Dejean, Tony, Alice Valentini, Antoine Duparc, Stephanie Pellier-Cuit, Francois Pompanon,

792         Pierre Taberlet, and Claude Miaud. 2011. 'Persistence of Environmental DNA in Freshwater

793         Ecosystems'. *PLoS One* 6 (8). https://doi.org/10.1371/journal.pone.0023398.

794     Derycke, Sofie, Jan Vanaverbeke, Annelien Rigaux, Thierry Backeljau, and Tom Moens. 2010.

41

795    'Exploring the Use of Cytochrome Oxidase c Subunit 1 (COI) for DNA Barcoding of Free-

796    Living Marine Nematodes'. Edited by Peter Roopnarine. *PLoS ONE* 5 (10): e13716.

797    https://doi.org/10.1371/journal.pone.0013716.

798  Dickie, Ian A., Stephane Boyer, Hannah L. Buckley, Richard P. Duncan, Paul P. Gardner, Ian D.

799    Hogg, Robert J. Holdaway, et al. 2018. 'Towards Robust and Repeatable Sampling Methods

800    in EDNA-Based Studies'. *Molecular Ecology Resources*. Wiley/Blackwell (10.1111).

801    https://doi.org/10.1111/1755-0998.12907.

802  Drummond, A J, R D Newcomb, T R Buckley, D Xie, A Dopheide, B C M Potter, J Heled, et al.

803    2015. 'Evaluating a Multigene Environmental DNA Approach for Biodiversity

804    Assessment'. *Gigascience* 4. https://doi.org/ARTN 4610.1186/s13742-015-0086-1.

805  Eren, A Murat, Joseph H Vineis, Hilary G Morrison, and Mitchell L Sogin. 2013. 'A Filtering

806    Method to Generate High Quality Short Reads Using Illumina Paired-End Technology'.

807    *PLoS ONE* 8 (6): e66643. https://doi.org/10.1371/journal.pone.0066643.

808  Escudié, Frédéric, Lucas Auer, Maria Bernard, Mahendra Mariadassou, Laurent Cauquil, Katia

809    Vidal, Sarah Maman, et al. 2018. 'FROGS: Find, Rapidly, OTUs with Galaxy Solution'.

810    Edited by Bonnie Berger. *Bioinformatics* 34 (8): 1287–94.

811    https://doi.org/10.1093/bioinformatics/btx791.

812  Evans, N T, B P Olds, M A Renshaw, C R Turner, Y Y Li, C L Jerde, A R Mahon, M E Pfrender,

813    G A Lamberti, and D M Lodge. 2016. 'Quantification of Mesocosm Fish and Amphibian

814    Species Diversity via Environmental DNA Metabarcoding'. *Molecular Ecology Resources*

815    16 (1): 29–41. https://doi.org/10.1111/1755-0998.12433.

816  Ficetola, Gentile Francesco, Johan Pansu, Aurélie Bonin, Eric Coissac, Charline Giguet-Covex,

817    Marta De Barba, Ludovic Gielly, et al. 2015. 'Replication Levels, False Presences and the

818    Estimation of the Presence/Absence from EDNA Metabarcoding Data'. *Molecular Ecology*

819     *Resources* 15 (3): 543–56. https://doi.org/10.1111/1755-0998.12338.

820    Fonseca, Vera G. 2018. 'Pitfalls in Relative Abundance Estimation Using Edna Metabarcoding'.

821     *Molecular Ecology Resources* 18 (5): 923–26. https://doi.org/10.1111/1755-0998.12902.

822    Fonseca, Vera G., Gary R Carvalho, Way Sung, Harriet F Johnson, Deborah M Power, Simon P

823     Neill, Margaret Packer, et al. 2010. 'Second-Generation Environmental Sequencing

824     Unmasks Marine Metazoan Biodiversity'. *Nature Communications* 1.

825     https://doi.org/9810.1038/ncomms1095.

826    Frøslev, Tobias Guldberg, Rasmus Kjøller, Hans Henrik Bruun, Rasmus Ejrnæs, Ane Kirstine

827     Brunbjerg, Carlotta Pietroni, and Anders Johannes Hansen. 2017. 'Algorithm for Post-

828     Clustering Curation of DNA Amplicon Data Yields Reliable Biodiversity Estimates'.

829     *Nature Communications* 8 (1). https://doi.org/10.1038/s41467-017-01312-x.

830    Gevers, Dirk, Frederick M. Cohan, Jeffrey G. Lawrence, Brian G. Spratt, Tom Coenye, Edward

831     J. Feil, Erko Stackebrandt, et al. 2005. 'Re-Evaluating Prokaryotic Species'. *Nature Reviews*

832     *Microbiology* 3 (9): 733–39. https://doi.org/10.1038/nrmicro1236.

833    Goldberg, Caren S., Cameron R. Turner, Kristy Deiner, Katy E. Klymus, Philip Francis

834     Thomsen, Melanie A. Murphy, Stephen F. Spear, et al. 2016. 'Critical Considerations for the

835     Application of Environmental DNA Methods to Detect Aquatic Species'. Edited by M.

836     Gilbert. *Methods in Ecology and Evolution* 7 (11): 1299–1307.

837     https://doi.org/10.1111/2041-210X.12595.

838    Hashimoto, Joel G, Bradley S Stevenson, and Thomas M Schmidt. 2003. 'Rates and

839     Consequences of Recombination between RRNA Operons'. *Journal of Bacteriology* 185

840     (3): 966–72. https://doi.org/10.1128/JB.185.3.966-972.2003.

841    Hebert, Paul D.N., Sujeevan Ratnasingham, and Jeremy R. de Waard. 2003. 'Barcoding Animal

842     Life: Cytochrome c Oxidase Subunit 1 Divergences among Closely Related Species'.

843    *Proceedings of the Royal Society of London. Series B: Biological Sciences* 270 (suppl_1):

844    S96-9. https://doi.org/10.1098/rsbl.2003.0025.

845    Herrera, Santiago, Hiromi Watanabe, and Timothy M. Shank. 2015. 'Evolutionary and

846    Biogeographical Patterns of Barnacles from Deep-Sea Hydrothermal Vents'. *Molecular*

847    *Ecology* 24 (3): 673–89. https://doi.org/10.1111/mec.13054.

848    Ji, Yinqiu, Louise Ashton, Scott M Pedley, David P Edwards, Yong Tang, Akihiro Nakamura,

849    Roger Kitching, et al. 2013. 'Reliable, Verifiable and Efficient Monitoring of Biodiversity

850    via Metabarcoding'. *Ecology Letters* 16 (10): 1245–57. https://doi.org/10.1111/ele.12162.

851    Klappenbach, J A, Paul R. Saxman, Cole James R., and Thomas M. Schmidt. 2001. 'Rrndb: The

852    Ribosomal RNA Operon Copy Number Database'. *Nucleic Acids Research* 29 (1): 181–84.

853    https://doi.org/10.1093/nar/29.1.181.

854    Leray, Matthieu, J Y Yang, C P Meyer, S C Mills, N Agudelo, V Ranwez, J T Boehm, and Ryuji

855    J. Machida. 2013. 'A New Versatile Primer Set Targeting a Short Fragment of the

856    Mitochondrial COI Region for Metabarcoding Metazoan Diversity: Application for

857    Characterizing Coral Reef Fish Gut Contents'. *Front Zool* 10: 34.

858    https://doi.org/10.1186/1742-9994-10-34.

859    Macheriotou, Lara, Katja Guilini, Tania Nara Bezerra, Bjorn Tytgat, Dinh Tu Nguyen, Thi Xuan

860    Phuong Nguyen, Febe Noppe, et al. 2019. 'Metabarcoding Free-Living Marine Nematodes

861    Using Curated 18S and CO1 Reference Sequence Databases for Species-Level Taxonomic

862    Assignments'. *Ecology and Evolution* 9 (1): 1–16. https://doi.org/10.1002/ece3.4814.

863    Machida, Ryuji J., and Nancy Knowlton. 2012. 'PCR Primers for Metazoan Nuclear 18S and 28S

864    Ribosomal DNA Sequences'. Edited by Jack Anthony Gilbert. *PLoS ONE* 7 (9): e46180.

865    https://doi.org/10.1371/journal.pone.0046180.

866    Machida, Ryuji J., Matthew Kweskin, and Nancy Knowlton. 2012. 'PCR Primers for Metazoan

44

867  Mitochondrial 12S Ribosomal DNA Sequences'. *PLoS ONE* 7 (4).

868  https://doi.org/10.1371/journal.pone.0035887.

869  Machida, Ryuji J., Matthieu Leray, Shian Lei Ho, and Nancy Knowlton. 2017. 'Data Descriptor:

870  Metazoan Mitochondrial Gene Sequence Reference Datasets for Taxonomic Assignment of

871  Environmental Samples'. *Scientific Data* 4. https://doi.org/10.1038/sdata.2017.27.

872  Mahe, F, Torbjørn Rognes, C Quince, Colomban De Vargas, and M Dunthorn. 2015. 'Swarm v2:

873  Highly-Scalable and High-Resolution Amplicon Clustering'. *PeerJ* 3. https://doi.org/Artn

874  E142010.7717/Peerj.1420.

875  Massana, Ramón Ramon, Angélique Gobet, Stéphane Audic, David Bass, Lucie Bittner,

876  Christophe Boutte, Aurélie Chambouvet, et al. 2015. 'Marine Protist Diversity in European

877  Coastal Waters and Sediments as Revealed by High-Throughput Sequencing'.

878  *Environmental Microbiology* 17 (10): 4035–49. https://doi.org/10.1111/1462-2920.12955.

879  Mayr, Ernst. 1942. *Systematics and the Origin of Species, from the Viewpoint of a Zoologist.*

880  New York, NY: Columbia University Press.

881  http://www.hup.harvard.edu/catalog.php?isbn=9780674862500.

882  McMurdie, Paul J., and Susan Holmes. 2013. 'Phyloseq: An R Package for Reproducible

883  Interactive Analysis and Graphics of Microbiome Census Data'. Edited by Michael Watson.

884  *PLoS ONE* 8 (4): e61217. https://doi.org/10.1371/journal.pone.0061217.

885  Minoche, André E, Juliane C Dohm, and Heinz Himmelbauer. 2011. 'Evaluation of Genomic

886  High-Throughput Sequencing Data Generated on Illumina HiSeq and Genome Analyzer

887  Systems'. *Genome Biology* 12 (11): R112. https://doi.org/10.1186/gb-2011-12-11-r112.

888  Nearing, Jacob T., Gavin M. Douglas, André M. Comeau, and Morgan G.I. Langille. 2018.

889  'Denoising the Denoisers: An Independent Evaluation of Microbiome Sequence Error-

890  Correction Approaches'. *PeerJ* 6: e5364. https://doi.org/10.7717/peerj.5364.

45

891    Nichols, Ruth V., Christopher Vollmers, Lee A. Newsom, Yue Wang, Peter D. Heintzman,

892        McKenna Leighton, Richard E. Green, and Beth Shapiro. 2018. 'Minimizing Polymerase

893        Biases in Metabarcoding'. *Molecular Ecology Resources* 18 (5): 927–39.

894        https://doi.org/10.1111/1755-0998.12895.

895    Oksanen, Jari, Michael Blanchet, Guillaume F. Friendly, Roeland Kindt, Pierre Legendre, Dan

896        McGlinn, R. Peter Minchin, R.B. O'Hara, et al. 2018. 'Vegan: Community Ecology

897        Package'. https://cran.r-project.org/package=vegan.

898    Pansu, Johan, Charline Giguet-Covex, Francesco Ficetola, Ludovic Gielly, Frederic Boyer, Eric

899        Coissac, Isabelle Domaizon, Lucie Zinger, Jerome Poulenard, and Fabien Arnaud. 2015.

900        'Environmental DNA Metabarcoding to Investigate Historic Changes in Biodiversity'.

901        *Genome* 58 (5): 264.

902    Parada, A E, D M Needham, and J A Fuhrman. 2016. 'Every Base Matters: Assessing Small

903        Subunit RRNA Primers for Marine Microbiomes with Mock Communities, Time Series and

904        Global Field Samples'. *Environ Microbiol* 18 (5): 1403–14. https://doi.org/10.1111/1462-

905        2920.13023.

906    Pawlowski, Jan W., Richard Christen, Beatrice Lecroq, Dipankar Bachar, Hamid Reza

907        Shahbazkia, Linda Amaral-Zettler, and Laure Guillou. 2011. 'Eukaryotic Richness in the

908        Abyss: Insights from Pyrotag Sequencing'. *PLoS One* 6 (4).

909        https://doi.org/e1816910.1371/journal.pone.0018169.

910    Pei, Anna Y., William E Oberdorf, Carlos W Nossa, Ankush Agarwal, Pooja Chokshi, Erika A

911        Gerz, Zhida Jin, et al. 2010. 'Diversity of 16S RRNA Genes within Individual Prokaryotic

912        Genomes'. *Applied and Environmental Microbiology* 76 (12): 3886–97.

913        https://doi.org/10.1128/AEM.02953-09.

914    Phillips, Jarrett D., Daniel J. Gillis, and Robert H. Hanner. 2019. 'Incomplete Estimates of

915       Genetic Diversity within Species: Implications for DNA Barcoding'. *Ecology and*

916       *Evolution*. John Wiley & Sons, Ltd. https://doi.org/10.1002/ece3.4757.

917   Plouviez, S., T. M. Shank, B. Faure, C. Daguin-Thiebaut, F. Viard, F. H. Lallier, and D. Jollivet.

918       2009. 'Comparative Phylogeography among Hydrothermal Vent Species along the East

919       Pacific Rise Reveals Vicariant Processes and Population Expansion in the South'. *Molecular*

920       *Ecology* 18 (18): 3903–17. https://doi.org/10.1111/j.1365-294X.2009.04325.x.

921   Quast, Christian, Elmar Pruesse, Pelin Yilmaz, Jan Gerken, Timmy Schweer, Pablo Yarza, Jörg

922       Peplies, and Frank Oliver Glöckner. 2012. 'The SILVA Ribosomal RNA Gene Database

923       Project: Improved Data Processing and Web-Based Tools'. *Nucleic Acids Research* 41 (D1):

924       D590–96. https://doi.org/10.1093/nar/gks1219.

925   Queiroz, Kevin de. 2005. 'Ernst Mayr and the Modern Concept of Species'. *Proceedings of the*

926       *National Academy of Sciences* 102 (Supplement 1): 6600–6607.

927       https://doi.org/10.1073/pnas.0502030102.

928   R Core Team. 2018. 'R: A Language and Environment for Statistical Computing.' R Foundation

929       for Statistical Computing, Vienna, Austria.

930   Salazar, Guillem, Francisco M Cornejo-Castillo, Veronica Benitez-Barrios, Eugenio Fraile-Nuez,

931       X Anton Alvarez-Salgado, Carlos M Duarte, Josep M Gasol, and Silvia G Acinas. 2016.

932       'Global Diversity and Biogeography of Deep-Sea Pelagic Prokaryotes'. *Isme Journal* 10 (3):

933       596–608. https://doi.org/10.1038/ismej.2015.137.

934   Schloss, Patrick D., Sarah L. Westcott, Thomas Ryabin, Justine R. Hall, Martin Hartmann, Emily

935       B. Hollister, Ryan A. Lesniewski, et al. 2009. 'Introducing Mothur: Open-Source, Platform-

936       Independent, Community-Supported Software for Describing and Comparing Microbial

937       Communities'. *Applied and Environmental Microbiology* 75 (23): 7537–41.

938       https://doi.org/10.1128/AEM.01541-09.

939    Schnell, Ida Bærholm, Kristine Bohmann, and M. Thomas P. Gilbert. 2015. 'Tag Jumps

940        Illuminated - Reducing Sequence-to-Sample Misidentifications in Metabarcoding Studies'.

941        *Molecular Ecology Resources* 15 (6): 1289–1303. https://doi.org/10.1111/1755-0998.12402.

942    Shank, Timothy M., Michael B. Black, Kenneth M. Halanych, Richard A. Lutz, and Robert C.

943        Vrijenhoek. 1999. 'Miocene Radiation of Deep-Sea Hydrothermal Vent Shrimp (Caridea:

944        Bresiliidae): Evidence from Mitochondrial Cytochrome Oxidase Subunit I'. *Molecular*

945        *Phylogenetics and Evolution* 13 (2): 244–54. https://doi.org/10.1006/mpev.1999.0642.

946    Shapiro, B. Jesse, Jean Baptiste Leducq, and James Mallet. 2016. 'What Is Speciation?' Edited

947        by Ivan Matic. *PLoS Genetics* 12 (3): e1005860.

948        https://doi.org/10.1371/journal.pgen.1005860.

949    Sinniger, Frederic, Jan W. Pawlowski, Saki Harii, Andrew J. Gooday, Hiroyuki Yamamoto,

950        Pierre Chevaldonné, Tomas Cedhagen, Gary Carvalho, and Simon Creer. 2016. 'Worldwide

951        Analysis of Sedimentary DNA Reveals Major Gaps in Taxonomic Knowledge of Deep-Sea

952        Benthos'. *Frontiers in Marine Science* 3 (June): 92.

953        https://doi.org/10.3389/FMARS.2016.00092.

954    Slon, Viviane, Charlotte Hopfe, Clemens L Weiß, Fabrizio Mafessoni, Marco De La Rasilla,

955        Carles Lalueza-Fox, Antonio Rosas, et al. 2017. 'Neandertal and Denisovan DNA from

956        Pleistocene Sediments'. *Science* 356 (6338): 605–8.

957        https://doi.org/10.1126/science.aam9695.

958    Sokal, Robert R., and Theodore J. Crovello. 1970. 'The Biological Species Concept : A Critical

959        Evaluation'. *The American Naturalist* 104 (936): 127–53.

960    Song, Hojun, Jennifer E Buhay, Michael F Whiting, and Keith A Crandall. 2008. 'Many Species

961        in One: DNA Barcoding Overestimates the Number of Species When Nuclear

962        Mitochondrial Pseudogenes Are Coamplified'. *Proceedings of the National Academy of*

963      *Sciences of the United States of America* 105 (36): 13486–91.

964      https://doi.org/10.1073/pnas.0803076105.

965    Stat, Michael, Megan J. Huggett, Rachele Bernasconi, Joseph D. Dibattista, Tina E. Berry,

966      Stephen J. Newman, Euan S. Harvey, and Michael Bunce. 2017. 'Ecosystem Biomonitoring

967      with EDNA: Metabarcoding across the Tree of Life in a Tropical Marine Environment'.

968      *Scientific Reports* 7. https://doi.org/10.1038/s41598-017-12501-5.

969    Stefanni, Sergio, David Stanković, Diego Borme, Alessandra de Olazabal, Tea Juretić, Alberto

970      Pallavicini, and Valentina Tirelli. 2018. 'Multi-Marker Metabarcoding Approach to Study

971      Mesozooplankton at Basin Scale'. *Scientific Reports* 8 (1): 12085.

972      https://doi.org/10.1038/s41598-018-30157-7.

973    Taberlet, Pierre, Eric Coissac, Mehrdad Hajibabaei, and Loren H. Rieseberg. 2012.

974      'Environmental DNA'. *Molecular Ecology* 21 (8): 1789–93. https://doi.org/10.1111/j.1365-

975      294X.2012.05542.x.

976    Tang, Cuong Q., Francesca Leasi, Ulrike Obertegger, Alexander Kieneke, Timothy G

977      Barraclough, and Diego Fontaneto. 2012. 'The Widely Used Small Subunit 18S RDNA

978      Molecule Greatly Underestimates True Diversity in Biodiversity Surveys of the Meiofauna.'

979      *Proceedings of the National Academy of Sciences of the United States of America* 109 (40):

980      16208–12. https://doi.org/10.1073/pnas.1209160109.

981    Teixeira, Sara, Karine Olu, Carole Decker, Regina L Cunha, Sandra Fuchs, Stéphane Hourdez,

982      Ester A. Serrão, and Sophie Arnaud-Haond. 2013. 'High Connectivity across the

983      Fragmented Chemosynthetic Ecosystems of the Deep Atlantic Equatorial Belt: Efficient

984      Dispersal Mechanisms or Questionable Endemism?' *Molecular Ecology* 22 (18): 4663–80.

985      https://doi.org/10.1111/mec.12419.

986    Valentini, Alice, François Pompanon, and Pierre Taberlet. 2009. 'DNA Barcoding for

987        Ecologists'. *Trends in Ecology and Evolution*. Elsevier Current Trends.

988        https://doi.org/10.1016/j.tree.2008.09.011.

989    Valentini, Alice, Pierre Taberlet, Claude Miaud, Raphaël Raphael Civade, Jelger Herder, Philip

990        Francis Thomsen, Eva Bellemain, et al. 2016. 'Next-Generation Monitoring of Aquatic

991        Biodiversity Using Environmental DNA Metabarcoding'. *Molecular Ecology* 25 (4): 929–

992        42. https://doi.org/10.1111/mec.13428.

993    Vargas, Colomban De, Stéphane Audic, Nicolas Henry, Johan Decelle, Frédéric Mahé, Ramiro

994        Logares, Enrique Lara, et al. 2015. 'Eukaryotic Plankton Diversity in the Sunlit Ocean'.

995        *Science* 348 (6237). https://doi.org/10.1126/science.1261605.

996    Wangensteen, Owen S., and Xavier Turon. 2016. 'Metabarcoding Techniques for Assessing

997        Biodiversity of Marine Animal Forests'. In *Marine Animal Forests*, edited by S. Rossi, L.

998        Bramanti, A. Gori, and C. Orejas Saco del Valle, 1–29. Cham: Springer International

999        Publishing. https://doi.org/10.1007/978-3-319-17001-5_53-1.

1000   Xiong, Wei, and Aibin Zhan. 2018. 'Testing Clustering Strategies for Metabarcoding-Based

1001       Investigation of Community–Environment Interactions'. *Molecular Ecology Resources* 18

1002       (6): 1326–38. https://doi.org/10.1111/1755-0998.12922.

1003   Yoccoz, N G, K A Brathen, L Gielly, J Haile, M E Edwards, T Goslar, H von Stedingk, et al.

1004       2012. 'DNA from Soil Mirrors Plant Taxonomic and Growth Form Diversity'. *Molecular*

1005       *Ecology* 21 (15): 3647–55. https://doi.org/10.1111/j.1365-294X.2012.05545.x.

1006   Yu, Douglas W, Yinqiu Ji, Brent C Emerson, Xiaoyang Wang, Chengxi Ye, Chunyan Yang, and

1007       Zhaoli Ding. 2012. 'Biodiversity Soup: Metabarcoding of Arthropods for Rapid Biodiversity

1008       Assessment and Biomonitoring'. *Methods in Ecology and Evolution* 3 (4): 613–23.

1009       https://doi.org/10.1111/j.2041-210X.2012.00198.x.

1010   Zanol, Joana, Kenneth M. Halanych, Torsten H. Struck, and Kristian Fauchald. 2010. 'Phylogeny

1011      of the Bristle Worm Family Eunicidae (Eunicida, Annelida) and the Phylogenetic Utility of

1012      Noncongruent 16S, COI and 18S in Combined Analyses'. *Molecular Phylogenetics and*

1013      *Evolution* 55 (2): 660–76. https://doi.org/10.1016/j.ympev.2009.12.024.

1014 Zhan, Aibin, Sarah A. Bailey, Daniel D. Heath, and Hugh J. Macisaac. 2014. 'Performance

1015      Comparison of Genetic Markers for High-Throughput Sequencing-Based Biodiversity

1016      Assessment in Complex Communities'. *Molecular Ecology Resources* 14 (5): 1049–59.

1017      https://doi.org/10.1111/1755-0998.12254.

1018 Zinger, Lucie, Jérôme Chave, Eric Coissac, Amaia Iribar, Eliane Louisanna, Sophie Manzi,

1019      Vincent Schilling, Heidy Schimann, Guilhem Sommeria-Klein, and Pierre Taberlet. 2016.

1020      'Extracellular DNA Extraction Is a Fast, Cheap and Reliable Alternative for Multi-Taxa

1021      Surveys Based on Soil DNA'. *Soil Biology and Biochemistry* 96: 16–19.

1022      https://doi.org/10.1016/j.soilbio.2016.01.008.

1023

## DATA ACCESSIBILITY

1025      The data for this work can be accessed in the European Nucleotide Archive (ENA)

1026 database (Study accession number will be given upon manuscript acceptance). The data set,

1027 including sequences, databases, as well as raw and refined ASV/OTU tables, has been deposited

1028 on ftp://ftp.ifremer.fr/ifremer/dataref/bioinfo/merlin/abyss/BioinformaticPipelineComparisons/.

1029 Bioinformatic scripts, config files, and R scripts are available on Gitlab

1030 (https://gitlab.ifremer.fr/abyss-project/).

1031  **AUTHOR CONTRIBUTIONS**

1032        MIB and SAH designed the study, MIB and JP carried out the laboratory and molecular

1033  work; MIB and BT performed the bioinformatic and statistical analyses. LQ assisted in the

1034  bioinformatic development and participated in the study design. MIB and SAH wrote the

1035  manuscript. All authors contributed to the final manuscript.