1

2    DR. LUCIE  ZINGER (Orcid ID : 0000-0002-3400-5825)

3    DR. HOLLY  BIK (Orcid ID : 0000-0002-4356-3837)

4    DR. ANTHONY A CHARITON (Orcid ID : 0000-0002-5809-3372)

5    DR. BRUCE E DEAGLE (Orcid ID : 0000-0001-7651-3687)

6    DR. IAN A. DICKIE (Orcid ID : 0000-0002-2740-2128)

7    DR. ALEX J. DUMBRELL (Orcid ID : 0000-0001-6282-3043)

8    DR. GENTILE FRANCESCO FRANCESCO FICETOLA (Orcid ID : 0000-0003-3414-5155)

9    DR. LUCA  FUMAGALLI (Orcid ID : 0000-0002-6648-2570)

10   DR. SIMON N JARMAN (Orcid ID : 0000-0002-0792-9686)

11   DR. LEHO  TEDERSOO (Orcid ID : 0000-0002-1635-1249)

12

13

14   Article type      : Editorial

15

16

17   *Corresponding author mail id: lucie@zinger.fr*

18   *Editorial note for Molecular Ecology*

19

20   **DNA metabarcoding - need for robust experimental designs to draw sound ecological**
21   **conclusions**

22

23   **Authors**

24   Lucie Zinger[1], Aurélie Bonin[2], Inger G. Alsos[3], Miklós Bálint[4,5], Holly Bik[6], Frédéric Boyer[2],
25   Anthony A. Chariton[7], Simon Creer[8], Eric Coissac[2], Bruce E. Deagle[9], Marta De Barba[2], Ian A.
26   Dickie[10], Alex J. Dumbrell[11], Gentile Francesco Ficetola[2,12], Noah Fierer[13,14], Luca Fumagalli[15], M.
27   Thomas P. Gilbert[16,17], Simon Jarman[18,19], Ari Jumpponen[20], Håvard Kauserud[21], Ludovic Orlando[22,23],
28   Johan Pansu[7,24,25], Jan Pawlowski[26,27], Leho Tedersoo[28], Philip Francis Thomsen[29], Eske
29   Willerslev[23,30,31], Pierre Taberlet[2, 3]

30

31   **Affiliations**

32 [1] Ecole Normale Supérieure, PSL Research University, Institut de Biologie de l'Ecole Normale
33 Supérieure (IBENS), CNRS UMR 8197, INSERM U1024, 46 rue d'Ulm, F-75005 Paris, France.

34 [2] Université Grenoble Alpes, CNRS, Laboratoire d'Ecologie Alpine (LECA), F-38000 Grenoble,
35 France.

36 [3] UiT – The Arctic University of Norway, Tromsø Museum, NO-9037 Tromsø, Norway.

37 [4] Senckenberg Biodiversity and Climate Research Centre, Senckenberganlage 25, 60325 Frankfurt am
38 Main, Germany.

39 [5] LOEWE Centre for Translational Biodiversity Genomics (LOEWE-TBG), Senckenberganlage 25,
40 60325 Frankfurt, Germany.

41 [6] Department of Nematology, University of California, Riverside, Riverside, CA, USA

42 [7] Department of Biological Sciences, Macquarie University, NSW 2113 Australia

43 [8] School of Natural Sciences, Bangor University, Gwynedd, LL57 2UW, UK

44 [9] Australian Antarctic Division, Kingston, Tasmania, Australia.

45 [10] BioProtection Research Centre, School of Biological Sciences, University of Canterbury,
46 Christchurch, New Zealand.

47 [11] School of Biological Sciences, University of Essex, Wivenhoe Park, Colchester, Essex, CO4 3SQ,
48 UK

49 [12] Department of Environmental Science and Policy, Università degli Studi di Milano. Via Celoria 26,
50 20133 Milano, Italy.

51 [13] Department of Ecology and Evolutionary Biology, University of Colorado, Boulder, CO, USA.

52 [14] Cooperative Institute for Research in Environmental Sciences, University of Colorado, Boulder,
53 CO, USA.

54 [15] Laboratory for Conservation Biology, Department of Ecology and Evolution, University of
55 Lausanne, Biophore building, CH-1015 Lausanne (Switzerland).

56 [16] Section for Evolutionary Genomics, Biological Institute, University of Copenhagen, Øster
57 Farimagsgade 5, 1353 Copenhagen, Denmark

58 [17] Norwegian University of Science and Technology, University Museum, Trondheim, Norway.

59 [18] Trace and Environmental DNA Laboratory, Department of Environment and Agriculture, Curtin
60 University, Perth, WA, Australia.

61 [19] Environomics Future Science Platform, CSIRO National Collections and Marine Infrastructure,
62 Crawley, WA, 6009, Australia.

63 [20] Division of Biology, Kansas State University, Manhattan, KS66506, USA.

64  [21] Section for Genetics and Evolutionary Biology (EVOGENE), University of Oslo, Blindernveien 31,

65  0316 Oslo, Norway.

66  [22] Laboratoire d'Anthropobiologie Moléculaire et d'Imagerie de Synthèse, CNRS UMR 5288,

67  Université de Toulouse, Université Paul Sabatier, 31000 Toulouse, France.

68  [23] Lundbeck Foundation GeoGenetics Center, University of Copenhagen, Øster Voldgade 5-7, 1350K

69  Copenhagen, Denmark.

70  [24] Station Biologique de Roscoff, UMR 7144 CNRS-Sorbonne Université, 29688 Roscoff, France.

71  [25] CSIRO Ocean & Atmosphere, Lucas Heights, NSW 2234, Australia.

72  [26] ID-Gene ecodiagnostics, Campus Biotech, Avenue Sécheron 15, 1202 Geneva, Switzerland.

73  [27] University of Geneva, Department of Genetics and Evolution, 1211 Geneva, Switzerland.

74  [28] Institute of Ecology and Earth Sciences, University of Tartu, 14a Ravila, 50411 Tartu, Estonia.

75  [29] Department of Bioscience, University of Aarhus, Ny Munkegade 116, DK-8000 Aarhus C,

76  Denmark

77  [30] Department of Zoology, University of Cambridge, Downing Street, Cambridge CB2 3EJ, UK.

78  [31] Wellcome Trust Sanger Institute, Wellcome Genome Campus, Hinxton, Cambridge CB10 1SA,

79  UK.

80

81  **ORCID**

82  Lucie Zinger : 0000-0002-3400-5825

83  Aurélie Bonin: 0000-0001-7800-8609

84  Inger G. Alsos: 0000-0002-8610-1085

85  Miklós Bálint: 0000-0003-0499-8536

86  Holly Bik: 0000-0002-4356-3837

87  Frédéric Boyer: 0000-0003-0021-9590

88  Anthony A. Chariton: 0000-0002-5809-3372

89  Simon Creer: 0000-0003-3124-3550

90  Eric Coissac: 0000-0001-7507-6729

91  Bruce E. Deagle: 0000-0001-7651-3687

92  Marta De Barba: 0000-0002-2979-3716

93  Ian A. Dickie: 0000-0002-2740-2128

94  Alex J. Dumbrell: 0000-0001-6282-3043

95  Gentile Francesco Ficetola: 0000-0003-3414-5155

96      Noah Fierer: 0000-0002-6432-4261

97      Luca Fumagalli: 0000-0002-6648-2570

98      M. Thomas P. Gilbert: 0000-0002-5805-7195

99      Simon Jarman: 0000-0002-0792-9686

100     Ari Jumpponen: 0000-0002-6770-2563

101     Håvard Kauserud: 0000-0003-2780-6090

102     Ludovic Orlando: 0000-0003-3936-1850

103     Johan Pansu: 0000-0003-0256-0258

104     Jan Pawlowski: 0000-0003-2421-388X

105     Leho Tedersoo: 0000-0002-1635-1249

106     Philip Francis Thomsen : 0000-0002-9867-4366

107     Eske Willerslev: 0000-0002-7081-6748

108     Pierre Taberlet: 0000-0002-3554-5954

109

111           DNA metabarcoding, especially when coupled with high-throughput DNA sequencing, is

112     currently revolutionizing our capacity to assess biodiversity across a full range of taxa and habitats,

113     from soil microbes (e.g. Thompson et al., 2017) to large marine fish (e.g. Thomsen et al., 2016), and

114     from contemporary to tens of thousands year-old biological communities (e.g. Willerslev et al., 2003).

115     The breadth of potential applications is immense and spans surveys on the diversity or diet of species

116     native to specific ecosystems to bioindication (Pawlowski et al., 2018). The approach is also

117     especially cost-effective and easy to implement, which makes DNA metabarcoding one of the tools of

118     choice of the $21^{st}$ century for fundamental research and the future of large-scale biodiversity

119     monitoring programs (reviewed in Bohan et al., 2017; Creer et al., 2016; Taberlet, Bonin, Zinger, &

120     Coissac, 2018; Thomsen & Willerslev, 2015). However, as is often the case with any emerging

121     technology, we feel that the rise of DNA metabarcoding is occurring at a pace and in a manner that

122     often loses sight of the challenges in producing high-quality and reproducible data (Baker, 2016).

123     DNA metabarcoding is by essence a multidisciplinary approach building upon many complementary

124     expertises, including field and theoretical knowledge, taxonomic expertise, molecular biology,

125     bioinformatics, and computational statistics. Combining all these within single studies is necessary,

126     not so much for producing and analyzing the data *per se*, but rather for minimizing and controlling the

127     possible biases that can be introduced at any step of the experimental workflow - i.e. from the

128     sampling to data analysis - and that can lead to spurious ecological conclusions (reviewed in Bálint et

129     al., 2016; Nilsson et al., 2019; Dickie et al., 2018; Taberlet et al., 2018).

130    Whether the starting material consists of DNA from bulk samples (community DNA) and/or
131    from environmental DNA (eDNA), all DNA metabarcoding studies rely on a deceptively simple
132    succession of core experimental steps: (i) sampling and preservation of the starting material, (ii),
133    DNA extraction, (iii) PCR amplification of a taxonomically-informative genomic region, (iv) high-
134    throughput DNA sequencing of the amplicons, and (v) sequence analysis using bioinformatic
135    pipelines. Despite this apparent simplicity, each step can potentially introduce its own sources of
136    artifacts and biases (Figure 1). For example, the sampling design might not be effective for capturing
137    the full taxonomic diversity or the ecological processes under study, an undesired bias for studies
138    based on species detection. The availability of DNA in the samples is governed by its production rate,
139    transport and persistence, processes which are all largely dependent on the targeted organisms, their
140    biomass, and the ecosystem considered. A correct assessment of an ecological phenomenon based on
141    DNA metabarcoding require not only implementation of standardized standardized, randomized and
142    repeatable sampling designs and procedures (Dickie et al., 2018), but also consideration of DNA
143    dynamics in the underlying matrix (i.e. in gut, feces, water or soil matrices from tropical or boreal
144    organisms/ecosystems; Barnes & Turner, 2016). Likewise, the community under study can be
145    enriched - on purpose or not - with specific taxa depending on how the sample is collected (e.g. filter
146    size for water samples, removal of roots or not for soils), how it is transported/preserved, and how
147    DNA is extracted (differential extraction efficiencies). PCR amplification is also well known to be an
148    important source of biases, that are now fully revealed with high-throughput DNA sequencing
149    techniques. The preferential amplification of certain taxa over other ones due to inappropriate primers
150    provides one such example of potential bias (Clarke, Soubrier, Weyrich, & Cooper, 2014; Deagle,
151    Jarman, Coissac, Pompanon, & Taberlet, 2014). Primer biases can both skew abundance profiles and
152    lead to false negatives. PCR amplification can produce false negatives too through the presence of e.g.
153    PCR inhibitors, but also many false positives through the introduction of replication errors by the
154    DNA polymerase or the formation of chimeric fragments (reviewed in Taberlet et al., 2018). False
155    positives can also be introduced at any step of the experimental workflow through the presence of
156    reagent contaminants (Salter et al., 2014), or through samples, extractions or PCR cross-
157    contaminations. An even more insidious source of false positives pertains to the occurrence of "tag
158    jumps", sometimes referred to as "mistagging", "tag-switching", or "cross-talks" (Carlsen et al., 2012;
159    Edgar, 2018; Esling, Lejzerowicz, & Pawlowski, 2015; Schnell, Bohmann, & Gilbert, 2015). PCR
160    amplicons are indeed often tagged with unique short nucleotide sequences added on the 5'-end of the
161    primers (i.e. "tags"), which allow pooling all PCRs within a single sequencing run and reducing
162    sequencing costs. Each sequence obtained resulting in apparent cross-contaminations can then be
163    bioinformatically assigned back to its sample of origin on the basis of its tag (Schnell et al., 2015).
164    However, the procedures underlying the preparation of DNA libraries and/or the sequencing can
165    introduce these "tag jumps", when the tag assigned to one particular sample is in fact recombined to
166    the sequences belonging to another sample (Taberlet et al., 2018). This introduces additional,  non-

negligible levels of sample cross-contaminations, which primarily involve the most abundant taxa and can have a disproportionate impact on samples with low DNA concentrations (Esling et al., 2015; Murray, Coghlan, & Bunce, 2015; Schnell et al., 2015). Similarly, the Illumina index located on the P5 sequencing adaptor can be subjected to "index jumps", resulting in apparent cross-contaminations (Taberlet et al., 2018). This bias happens when several individual Illumina sequencing libraries are pooled and loaded on the same sequencing lane (Kircher, Sawyer, & Meyer, 2012) Finally, high-throughput DNA sequencing instruments have their own error rates (Schirmer et al., 2015). The above list of problems is clearly not exhaustive, and the interested reader will find more complete reviews elsewhere (e.g. Bálint et al., 2016; Nilsson et al., 2019; Taberlet et al., 2018). Still, it illustrates that any potential bias must be considered carefully when designing an experimental protocol and when interpreting the results. This is crucial to limit their impact on downstream analyses, and to ensure that the conclusion drawn from such data are authentic.

There is now an increasingly diverse range of field, laboratory (e.g. Caporaso et al., 2011; Taberlet et al., 2018; Valentini et al., 2009) and bioinformatics (e.g. Boyer et al., 2016; Caporaso et al., 2010; Dumbrell, Ferguson, & Clark, 2016) procedures aiming at reducing the amount of both false negatives (i.e. due to partial sampling, extraction, amplification or sequencing bias) and false positives (i.e. due to contaminations, "tag/index jumps", or PCR and sequencing errors) in DNA metabarcoding experiments. However, using these protocols does not necessarily guarantee that the problem of false positives or negatives is completely under control. These protocols must continuously be reconsidered, especially alongside the emergence of novel DNA sequencing technologies that provide new opportunities, but also new challenges. Additionally, each individual study and each genomic marker comes with its own specificities, and this often requires customization of the above protocols. The sequence clustering threshold to be used to form molecular taxonomic units relevant to the question addressed (e.g. removing intraspecific marker variability when the species level is desired) provides such an example, and will critically depend on both the marker specificities and PCR/sequencing error rates. Bioinformatics tools can further fail to exclude molecular artifacts when the filtering thresholds are relaxed, which inflates sample diversity estimates. Likewise, they can also generate false negatives, for example when a genuine metabarcode is falsely flagged as an error or chimera, or when it is assigned to an incorrect taxon due to incomplete or inappropriate reference databases (Alsos et al., 2018; Coissac, Riaz, & Puillandre, 2012). This can be especially problematic when the question investigated strongly relies on species detection. It is therefore crucial to include several types of experimental controls so as to facilitate the exclusion of spurious signal and support the reliability of the biological conclusions (Figure 1). Amongst these controls, conducting pilot experiments is particularly helpful to assess how appropriate the sampling design is (Dickie et al., 2018). We also recommend that both biological replicates (i.e. multiple independent samples) and technical replicates (i.e. multiple extractions/PCR of the same sample and/or extract) are included in

203  the experimental workflow to disentangle the effect of both the biological and technical variances
204  (Ficetola et al., 2015). These replications are necessary because both sampling and PCR can introduce
205  biases in a stochastic manner, especially when the concentration of the target DNA is low. It is also
206  essential to analyze a sufficient number of negative controls at the field sampling, DNA extraction,
207  PCR, and sequencing steps, as well as positive controls consisting of mock communities, known DNA
208  samples, or even synthetic sequences reflecting the attributes of the targeted products (Figure 1). All
209  these controls must be sequenced along the biological samples, as they facilitate the detection of
210  sporadic contaminations and tag or index jumps while helping adjusting filtering and clustering
211  thresholds. Ultimately, they will be a token of the reliability of the whole data curation process (De
212  Barba et al., 2014). We also encourage careful consideration of the bioinformatics workflow itself,
213  since the filtering steps necessary to curate the data will critically depend on the experimental design
214  and the ecological question under study. Typically, sequences of low abundance in a given sample
215  may be genuine or artifacts deriving from PCR/sequencing errors or tag/index jumps. The retained
216  filtering threshold for taxon presence is thus dependent on the underlying rates of artifacts, as well as
217  on the sequencing depth. As the different experimental controls provide direct measurements of these
218  artifacts, they will therefore allow better tuning of the filtering thresholds. All of these technical
219  considerations should be precisely reported within publications together with relevant illustrations and
220  statistics characterizing the workflow, as they are necessary to assess the relevance and quality of the
221  data underpinning specific conclusions. A last, a most obvious example of control consists in
222  assessing the plausibility of the taxonomic composition based on *a priori* knowledge of the system or
223  taxa studied. Such knowledge can be derived from data obtained with complementary sensing
224  approaches such as visual observations. In this case, building exhaustive local reference databases of
225  the genomic marker used from local specimens will secure the taxonomic assignment step (e.g. Alsos
226  et al., 2018). When local information is unavailable, typically when studying microorganisms, it
227  remains possible to assess whether the community is composed of clades that are expected to occur in
228  the system surveyed or not, as e.g. soils, sediments, and gut environments harbour highly different
229  bacterial phyla (e.g. Thompson et al., 2017).

230  As users, readers, referees or editors, we realize that the above-mentioned issues remain too
231  often overlooked. This problematic stance can lead to unsubstantiated claims and undermine scientific
232  advances if not resolved. Inappropriate practices such as estimating species richness from fingerprint
233  profiles (Bent, Pierson, & Forney, 2007), the absence of biological replicates (Prosser, 2010), or that
234  of contaminant controls (Perez-Muñoz, Arrieta, Ramer-Tait, & Walter, 2017) have been repeatedly
235  criticized in the field of microbial ecology, and in the latter case, they contribute to the rising debate
236  about the existence or not of a womb microbiota. Ancient DNA research has also developed rigorous
237  standards to tackle issues related to contamination, sequencing errors, and data reproducibility (Poinar
238  and Cooper (2000). We believe that the community of DNA metabarcoding users has now come of

239 age and learnt from its past errors. At a time when more and more exhaustive guides of best practices

240 on the subject are emerging (Knight et al., 2018; Pollock, Glendinning, Wisedchanwet, & Watson,

241 2018; Taberlet et al., 2018), and where DNA sequencing costs are rapidly decreasing, we should be

242 always mindful of the adage "better safe than sorry". This note does not mean to imply that the

243 systematic use of the highest technical and analytical standards is reasonable nor the universal remedy

244 for all the challenges associated with DNA metabarcoding. Rather, we strongly encourage researchers

245 and end-users to adopt reflective decision-making when designing their experiment and to critically

246 appraise their results, with the ultimate aim to prove the robustness and reproducibility of their

247 conclusions.

248 **References**

249 Alsos, I. G., Lammers, Y., Yoccoz, N. G., Jørgensen, T., Sjögren, P., Gielly, L., & Edwards, M. E.

250     (2018). Plant DNA metabarcoding of lake sediments: How does it represent the contemporary

251     vegetation. *PLoS One*, *13*(4), e0195403.

252 Baker, M. (2016). 1,500 scientists lift the lid on reproducibility. *Nature News*, *533*(7604), 452.

253 Bálint, M., Bahram, M., Eren, A. M., Faust, K., Fuhrman, J. A., Lindahl, B., … Tedersoo, L. (2016).

254     Millions of reads, thousands of taxa: microbial community structure and associations analyzed

255     via marker genes. *FEMS Microbiology Reviews*, *40*(5), 686–700.

256 Barnes, M. A., & Turner, C. R. (2016). The ecology of environmental DNA and implications for

257     conservation genetics. *Conservation Genetics*, *17*(1), 1–17.

258 Bent, S. J., Pierson, J. D., & Forney, L. J. (2007). Measuring species richness based on microbial

259     community fingerprints: the emperor has no clothes. *Applied and Environmental Microbiology*,

260     *73*(7), 2399.

261 Bohan, D. A., Vacher, C., Tamaddoni-Nezhad, A., Raybould, A., Dumbrell, A. J., & Woodward, G.

262     (2017). Next-generation global biomonitoring: large-scale, automated reconstruction of

263     ecological networks. *Trends in Ecology & Evolution*, *32*(7), 477–487.

264 Boyer, F., Mercier, C., Bonin, A., Le Bras, Y., Taberlet, P., & Coissac, E. (2016). obitools: a unix-

265     inspired software package for DNA metabarcoding. *Molecular Ecology Resources*, *16*(1), 176–

266     182.

267 Caporaso, J. G., Kuczynski, J., Stombaugh, J., Bittinger, K., Bushman, F. D., Costello, E. K., …

268     Knight, R. (2010). QIIME allows analysis of high-throughput community sequencing data.

269     *Nature Methods*, *7*(5), 335–336.

270 Caporaso, J. G., Lauber, C. L., Walters, W. A., Berg-Lyons, D., Lozupone, C. A., Turnbaugh, P. J., …

271     Knight, R. (2011). Global patterns of 16S rRNA diversity at a depth of millions of sequences per

272     sample. *Proceedings of the National Academy of Sciences of the United States of America*, *108

273     Suppl 1*, 4516–4522.

274 Carlsen, T., Aas, A. B., Lindner, D., Vrålstad, T., Schumacher, T., & Kauserud, H. (2012). Don't

275   make a mista(g)ke: is tag switching an overlooked source of error in amplicon pyrosequencing
276   studies? *Fungal Ecology*, *5*(6), 747–749.
277   Clarke, L. J., Soubrier, J., Weyrich, L. S., & Cooper, A. (2014). Environmental metabarcodes for
278   insects: *in silico* PCR reveals potential for taxonomic bias. *Molecular Ecology Resources*, *14*(6),
279   1160–1170.
280   Coissac, E., Riaz, T., & Puillandre, N. (2012). Bioinformatic challenges for DNA metabarcoding of
281   plants and animals. *Molecular Ecology*, *21*, 1834–1847.
282   Creer, S., Deiner, K., Frey, S., Porazinska, D., Taberlet, P., Thomas, W. K., … Bik, H. M. (2016). The
283   ecologist's field guide to sequence-based identification of biodiversity. *Methods in Ecology and*
284   *Evolution*, *7*(9), 1008–1018.
285   Deagle, B. E., Jarman, S. N., Coissac, E., Pompanon, F., & Taberlet, P. (2014). DNA metabarcoding
286   and the cytochrome c oxidase subunit I marker: not a perfect match. *Biology Letters*, *10*(9),
287   20140562.
288   De Barba, M., Miquel, C., Boyer, F., Rioux, D., Coissac, E., & Taberlet, P. (2014). DNA
289   metabarcoding multiplexing for omnivorous diet analysis and validation of data accuracy.
290   *Molecular Ecology Resources*, *14*(2), 306–323.
291   Dickie, I. A., Boyer, S., Buckley, H. L., Duncan, R. P., Gardner, P. P., Hogg, I. D., … Weaver, L.
292   (2018). Towards robust and repeatable sampling methods in eDNA-based studies. *Molecular*
293   *Ecology Resources*, *18*(5), 940–952.
294   Dumbrell, A. J., Ferguson, R. M. W., & Clark, D. R. (2016). Microbial community analysis by single-
295   amplicon high-throughput next generation sequencing: data analysis – from raw output to
296   ecology. In T. J. McGenity, K. N. Timmis, & B. Nogales (Eds.), *Hydrocarbon and lipid*
297   *microbiology protocols*. Humana Press.
298   Edgar, R. (2018). *UNCROSS2: identification of cross-talk in 16S rRNA OTU tables. bioRxiv.*
299   doi:10.1101/400762
300   Esling, P., Lejzerowicz, F., & Pawlowski, J. (2015). Accurate multiplexing and filtering for high-
301   throughput amplicon-sequencing. *Nucleic Acids Research*, *43*(5), 2513–2524.
302   Ficetola, G. F., Pansu, J., Bonin, A., Coissac, E., Giguet-Covex, C., De Barba, M., … Taberlet, P.
303   (2015). Replication levels, false presences and the estimation of the presence/absence from
304   eDNA metabarcoding data. *Molecular Ecology Resources*, *15*(3), 543–556.
305   Kircher, M., Sawyer, S., & Meyer, M. (2012). Double indexing overcomes inaccuracies in multiplex
306   sequencing on the Illumina platform. *Nucleic Acids Research*, *40*. doi:10.1093/nar/gkr771
307   Knight, R., Vrbanac, A., Taylor, B. C., Aksenov, A., Callewaert, C., Debelius, J., … Dorrestein, P. C.
308   (2018). Best practices for analysing microbiomes. *Nature Reviews. Microbiology*, *16*(7), 410–
309   422.
310   Murray, D. C., Coghlan, M. L., & Bunce, M. (2015). From benchtop to desktop: important
311   considerations when designing amplicon sequencing workflows. *PLoS One*, *10*(4), e0124671.

312 Nilsson, R. H., Anslan, S., Bahram, M., Wurzbacher, C., Baldrian, P., & Tedersoo, L. (2019).
313    Mycobiome diversity: high-throughput sequencing and identification of fungi. *Nature Reviews.*
314    *Microbiology*, *17*(2), 95–109.

315 Pawlowski, J., Kelly-Quinn, M., Altermatt, F., Apothéloz-Perret-Gentil, L., Beja, P., Boggero, A., …
316    Kahlert, M. (2018). The future of biotic indices in the ecogenomic era: Integrating (e)DNA
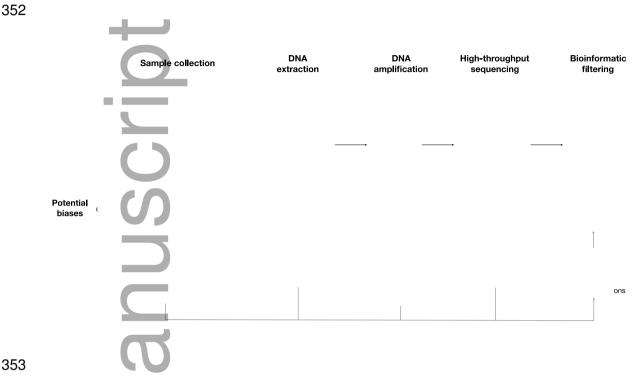317    metabarcoding in biological assessment of aquatic ecosystems. *The Science of the Total*
318    *Environment*, *637-638*, 1295–1310.

319 Perez-Muñoz, M. E., Arrieta, M.-C., Ramer-Tait, A. E., & Walter, J. (2017). A critical assessment of
320    the ``sterile womb'' and ``in utero colonization" hypotheses: implications for research on the
321    pioneer infant microbiome. *Microbiome*, *5*(1), 48.

322 Poinar, H. N., & Cooper, A. (2000). Ancient DNA: do it right or not at all. *Science*, *5482*, 1139.

323 Pollock, J., Glendinning, L., Wisedchanwet, T., & Watson, M. (2018). The madness of microbiome:
324    attempting to find consensus "best practice" for 16S microbiome studies. *Applied and*
325    *Environmental Microbiology*, *84*(7), e02627–17.

326 Prosser, J. I. (2010). Replicate or lie. *Environmental Microbiology*, *12*(7), 1806–1810.

327 Salter, S. J., Cox, M. J., Turek, E. M., Calus, S. T., Cookson, W. O., Moffatt, M. F., … Walker, A. W.
328    (2014). Reagent and laboratory contamination can critically impact sequence-based microbiome
329    analyses. *BMC Biology*, *12*, 87.

330 Schirmer, M., Ijaz, U. Z., D'Amore, R., Hall, N., Sloan, W. T., & Quince, C. (2015). Insight into
331    biases and sequencing errors for amplicon sequencing with the Illumina MiSeq platform. *Nucleic*
332    *Acids Research*, *43*(6), e37.

333 Schnell, I. B., Bohmann, K., & Gilbert, M. T. P. (2015). Tag jumps illuminated – reducing sequence-
334    to-sample misidentifications in metabarcoding studies. *Molecular Ecology Resources*, *15*(6),
335    1289–1303.

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

338 Thompson, L. R., Sanders, J. G., McDonald, D., Amir, A., Ladau, J., Locey, K. J., … Earth
339    Microbiome Project Consortium. (2017). A communal catalogue reveals Earth's multiscale
340    microbial diversity. *Nature*, *551*(7681), 457–463.

341 Thomsen, P. F., Møller, P. R., Sigsgaard, E. E., Knudsen, S. W., Jørgensen, O. A., & Willerslev, E.
342    (2016). Environmental DNA from seawater samples correlate with trawl catches of subarctic,
343    deepwater fishes. *PLoS One*, *11*(11), e0165252.

344 Thomsen, P. F., & Willerslev, E. (2015). Environmental DNA as an emerging tool in conservation for
345    monitoring past and present biodiversity. *Biological Conservation*, *183*(C), 4–18.

346 Valentini, A., Miquel, C., Nawaz, M. A., Bellemain, E., Coissac, E., Pompanon, F., … Taberlet, P.
347    (2009). New perspectives in diet analysis based on DNA barcoding and parallel pyrosequencing:
348    the *trn*L approach. *Molecular Ecology Resources*, *9*, 51–60.

349  Willerslev, E., Hansen, A. J., Binladen, J., Brand, T. B., Gilbert, M. T. P., Shapiro, B., … Cooper, A.

350      (2003). Diverse plant and animal genetic records from Holocene and Pleistocene sediments.

351      *Science*, *300*, 791–795.

352



353

354  **Figure 1 Summarized workflow of DNA metabarcoding and biases in the data production**

355  **process, with the potential associated controls to assess data quality.** Expectations on the local

356  community, either from a priori knowledge on the site or organisms targeted, or obtained through e.g.

357  vizual census, specimen collection, or building of a local reference database, constitute a first

358  assessment of the DNA metabarcoding experiment success. Pilot experiments are essential for

359  optimizing the whole experimental design, from the sampling strategy (e.g. number of biological

360  replicates) to the entire technical approach. Field, extraction, PCR, and tagging-system negative and

361  positive controls should be sequenced along with biological samples. They all aim at identifying (i)

362  potential contaminants that could be introduced at any experimental step, and (ii) potential

363  experimental artifacts due to the DNA extraction, PCR, and sequencing steps. Field negative controls

364  consists of extracting DNA from storage/extraction buffers brought to the field or used to clean

365  sampling instruments. Tagging-system negative controls can only be implemented when amplicons

366  are identified by a unique combination of tags attached to the 5' end of each amplification primer, and

367  where one or several tag combinations remain unused in the experimental design. In such conditions,

368  tagging-system controls can be performed at the bioinformatics analysis step, by monitoring the

369  number of sequences harboring unexpected tag combinations. This number is actually a direct

370  measurement of the tag-jump rate. "Index jumps" are more difficult to evaluate, and can be controlled

371  either by indexing both library adapters (P5 and P7) or when the libraries sequenced together have

372   identifiable sequences that could indicate their origin. The positive controls (constructed using either
373   synthetic DNA with the primer target sequences on both sides, DNA extracted from a mock
374   community, or known environmental samples), as well as prior expectations on the taxa that should
375   occur in the system can be used to evaluate the effectiveness of the data production process, the
376   impact of contaminants on the retrieved ecological signal and the adequacy of bioinformatics filtering
377   procedures.

**Sample collection**

**DNA dynamics**
**(~season, system, organism)**

**DNA extraction**

**DNA amplification**

reagents/aerosols contaminants

**High-throughput sequencing**

**Bioinformatic filtering**

**Potential biases**

Undersampling
Contamination from past/neighbouring events
Experimental contamination

Undersampling
Taxon-specific inefficiency
Experimental contamination

+ Polymerase errors/chimeras
+ Inappropriate primers

+ Tag/index jump &
sequencing errors

Inappropriate filtering thresholds
Mis-classifications

**Potential controls**

- Expected target (or non-target) taxa
- Building of a local reference database
- Pilot experiment
- Biological replicates
- Field negative controls

- Technical replicates
- Extraction negative controls
- Positive controls

- Technical replicates
- PCR negative controls
- Positive controls
- Use of multiple primer set or *in silico* pre-evaluation of primers

- Tagging system negative controls

- Filtering/clustering criteria and threshold adjustments based on all controls and replicates
- Taxonomic congruence with *a priori* expectations