

1 **Coming of age for COI metabarcoding of whole organism community DNA:**
2 **towards bioinformatic harmonisation**

3 Creedy T.J.¹, Andújar C.², Meramveliotakis E.³, Nogueras V.^{2,3}, Overcast I.⁴, Papadopoulou A.³,
4 Morlon H.⁴, Vogler A.P.^{1,5}, Emerson B.C.² & Arribas P.²

5 ¹Department of Life Sciences, Natural History Museum, Cromwell Road, London SW5 7BD

6 ²Instituto de Productos Naturales y Agrobiología (IPNA-CSIC), 38206, S.C. La Laguna, Spain

7 ³Department of Biological Sciences, University of Cyprus, Nicosia, Cyprus

8 ⁴Institut de Biologie de l'ENS (IBENS), Département de biologie, École normale supérieure,
9 CNRS, INSERM, Université PSL, Paris, France

10 ⁵Department of Life Sciences, Imperial College London Silwood Park Campus, Ascot

11 **Abstract**

12 Metabarcoding of DNA extracted from community samples of whole organisms (whole organism
13 community DNA, wocDNA) is increasingly being applied to terrestrial, marine and freshwater
14 metazoan communities to provide rapid, accurate and high resolution data for novel molecular
15 ecology research. The growth of this field has been accompanied by considerable development
16 that builds on microbial metabarcoding methods to develop appropriate and efficient sampling
17 and laboratory protocols for whole organism metazoan communities. However, considerably less
18 attention has focused on ensuring bioinformatic methods are adapted and applied
19 comprehensively in wocDNA metabarcoding. In this study we examined over 600 papers and
20 identified 111 studies that performed COI metabarcoding of wocDNA. We then systematically
21 reviewed the bioinformatic methods employed by these papers to identify the state-of-the-art.
22 Our results show that the increasing use of wocDNA COI metabarcoding for metazoan diversity
23 is characterised by a clear absence of bioinformatic harmonisation, and the temporal trends show
24 little change in this situation. The reviewed literature showed (i) high heterogeneity across
25 pipelines, tasks and tools used, (ii) limited or no adaptation of bioinformatic procedures to the
26 nature of the COI fragment, and (iii) a worrying underreporting of tasks, software and
27 parameters. Based upon these findings we propose a set of recommendations that we think the
28 wocDNA metabarcoding community should consider to ensure that bioinformatic methods are
29 appropriate, comprehensive and comparable. We believe that adhering to these recommendations
30 will improve the long-term integrative potential of wocDNA COI metabarcoding for biodiversity
31 science.

32 **Keywords:** *metabarcoding, COI barcode, animal communities, high-throughput sequencing,*
33 *bioinformatics, community ecology*

34 **Introduction**

35 Metabarcoding of DNA extracted from community samples of whole organisms (whole
36 organism community DNA, wocDNA) is a reliable and cost-efficient tool to study the
37 biodiversity of metazoan communities (Bush et al., 2019; Ji et al., 2013; Porter & Hajibabaei,
38 2018). This approach, which has also been referred to as community DNA (e.g. Andújar et al.,
39 2018b; Deiner et al., 2017) or bulk sample DNA (e.g. Braukmann et al., 2019; Yu et al., 2012)
40 metabarcoding, primarily differs from other approaches such as eDNA (environmental or extra-
41 organismal DNA; Taberlet et al., 2012) or iDNA (vertebrate DNA ingested by invertebrates;
42 Schnell et al., 2012) in that the source material is a community of whole organisms collected
43 through direct trapping or collection (e.g. malaise traps Ji et al., 2013, canopy fogging Creedy et
44 al., 2019) or separated from an environmental sample (e.g. from soil Arribas et al. 2016 or water
45 Suter et al., 2020). As a consequence, compared with eDNA and iDNA, wocDNA samples are
46 characterised by (i) a comparatively low level of DNA degradation in the target species, (ii) a
47 low proportion of non-target species, and (iii) the possibility for complementing, refining and/or
48 validating metabarcoding-derived community data against other conventional morphological and
49 molecular methods.

50 Metabarcoding of wocDNA samples is increasingly employed in community ecology,
51 evolutionary ecology, biogeography, conservation biology, environmental management, and
52 policy and decision-making (e.g. Bush et al., 2020; deWaard et al., 2019; Leese et al., 2018).
53 Metazoan wocDNA metabarcoding has been adapted from pioneering approaches developed to
54 inventory and characterise microbial diversity (e.g. Gilbert et al., 2010; Sogin et al., 2006). The
55 majority of these adaptations have focused on sampling, and molecular laboratory steps,
56 including adapted protocols to (i) sample, separate, enrich and/or clean animal wocDNA samples
57 (Creedy et al., 2019; Fonseca et al., 2010, 2011), (ii) perform wocDNA extractions (Marquina et
58 al., 2019; Nielsen et al., 2019), (iii) design and evaluate primers (Braukmann et al., 2019;

59 Elbrecht & Leese, 2017, Elbrecht et al. 2019), optimise amplification (Krehenwinkel et al.,
60 2017) and prepare libraries (Yang et al., 2020). There is a growing consensus on the use of the
61 mitochondrial cytochrome oxidase subunit I (COI) barcode, rather than other markers widely
62 used for metabarcoding of non-metazoan communities, as the standard for wocDNA
63 metabarcoding due to the range of COI primers with demonstrated efficiency (Braukmann et al.,
64 2019; Elbrecht & Leese, 2017), and the potential of COI to improve the utility, resolution and
65 reliability of wocDNA metabarcoding data (Andújar et al., 2018a; Turon et al., 2020).

66 However, in contrast to these advances in sampling and molecular processing, there has
67 been limited effort to review and evaluate how bioinformatic processing has been adapted to
68 metazoan wocDNA samples and the COI barcode, nor to examine consistency in bioinformatic
69 approaches across the field. Broadly, bioinformatic tasks involve the computational cleaning,
70 filtering and analysis of raw sequence data to produce biodiversity data comprising taxonomic
71 units and their incidence across samples, implemented in a particular order (a ‘pipeline’). There
72 are a wide array of software tools available for performing different bioinformatic tasks, from
73 standalone tools to catch-all software packages (e.g. *OBItools* Boyer et al., 2016; *QIIME*
74 Caporaso et al., 2010; *USEARCH/UPARSE* Edgar, 2013; and its open-source derivative
75 *VSEARCH* Rognes et al., 2016). This software has been largely developed for metabarcode loci
76 other than the COI region, with very few tools explicitly developed for protein coding
77 metabarcodes (although see Andújar et al., 2021; Nugent et al., 2020; Ramirez-Gonzalez et al.,
78 2013). To fully capitalise on the COI barcode for metabarcoding, bioinformatics should be
79 specifically tailored to its evolutionary properties, such as the ability to interrogate the amino
80 acid translation, and accounting for established patterns of sequence variation in protein coding
81 genes for strict filtering. Additionally, metabarcoding employs a number of key bioinformatic
82 tasks for which multiple alternative algorithms have been developed (e.g. denoising algorithms),
83 with considerable variation in outcomes depending on parameters and thresholds applied.

84 The structure of a bioinformatic metabarcoding pipeline will depend strongly on the
85 research aim, amplification and sequencing protocols, target locus, and target biodiversity
86 fraction. The diversity of bioinformatic tasks and the software approaches to implement them is
87 of course beneficial for designing appropriate pipelines, but such heterogeneity may also restrict
88 integrated, standardised and synergistic growth in the field. As metazoan wocDNA
89 metabarcoding becomes more accessible to researchers from a range of fields and backgrounds,
90 harmonisation of bioinformatic approaches is important to ensure (i) high-quality, reproducible
91 data amenable to qualitative or quantitative reviews and meta-analysis across studies, and (ii) a
92 reliable, consistent methodology for wider implementation, development and expansion of
93 wocDNA metabarcoding. We consider harmonisation not to mean strict prescription of the tasks
94 and software to use, nor their order. Instead a harmonised field would recognise the diversity of
95 approaches available, while recording key steps and establishing the effects of parameter choice
96 on the outcome of metabarcoding studies. This approach could be enabled by the adoption of
97 universal aligned standards for data generation and processing, while allowing for flexibility in
98 implementation to adapt to varying research goals and take advantage of novel methodological
99 development.

100 Harmonisation requires comprehensive examination of current practice to understand the
101 aims and approaches of prior work, and a synthesis of the successes and failures in past
102 implementations for the purposes of elaborating a framework to guide future research. Therefore
103 it is our aim to summarise the state of the art for bioinformatic processing of metazoan wocDNA
104 COI metabarcoding, and in doing so assess the potential for harmonisation. To this end, we
105 performed a systematic review of peer-reviewed studies, collating information on the different
106 bioinformatic pipelines, tasks and tools used in wocDNA COI metabarcoding in >100 recent
107 studies (2011-2020). We use this data to (i) describe the diversity, heterogeneity and
108 reproducibility of the bioinformatic procedures followed, (ii) identify the extent to which these

109 procedures are compatible with the evolutionary properties of the COI marker, and (iii) identify
110 the key bioinformatic tasks, provide a framework for successful metabarcoding bioinformatics
111 and make recommendations towards harmonised bioinformatic procedures for metazoan
112 wocDNA COI metabarcoding.

113 **Materials and Methods**

114 *Bibliographic search and screening*

115 We focused this work on studies using whole organism community DNA (wocDNA)
116 metabarcoding. In general, we define wocDNA samples as those where the target organisms
117 were: (i) likely alive at the time of sampling, (ii) present as a largely complete specimen, and (iii)
118 potentially identifiable using classical methods of morphological analysis. We exclude eDNA
119 and iDNA metabarcoding due to the potentially different bioinformatic processing needs
120 associated with these samples. In particular, eDNA and iDNA bioinformatic methods need to
121 accommodate degraded DNA and a potentially high proportion of non-target reads. Furthermore,
122 in many cases wocDNA metabarcoding is directly comparable to direct observation of
123 specimens and conventional methods of taxonomic assignment not available for eDNA
124 metabarcoding (Ji et al. 2013, Aylagas et al. 2016). This allows for more objective stringency
125 thresholds in bioinformatic filtering and delimitation of operational taxonomic units (OTU).

126 We conducted a systematic search of peer-reviewed studies in the Web of Science (WOS)
127 Core Collection (Science Citation Index Expanded, 1900-present) on 3rd November 2020, using
128 the search “TS = (metabarcoding) NOT TS = (*micro* OR *bacteria* OR *myco* OR
129 *archaea* OR fungi OR plant OR eDNA OR environmental DNA)”. These search parameters
130 were selected in order to obtain a comprehensive set of wocDNA metabarcoding studies limited
131 to Metazoa.

132 The systematic search resulted in 692 records, which were screened to to select only
133 those studies that: (i) amplified some portion of the standard COI barcode “Folmer” region
134 (Folmer et al., 1994), (ii) fit our definition of wocDNA samples, comprising mixtures of
135 organisms extracted from the substrate, and (iii) provided a characterisation of metazoan
136 communities. Studies targeting extra-organismal DNA (i.e. eDNA, iDNA) were excluded. We
137 included studies of experimental mock communities composed of mixtures of DNA extracted
138 from individual specimens or mixtures of specimens, and we also included studies where the
139 target organisms remained partially or completely within an environmental substrate upon which
140 DNA extraction was performed (e.g. parasites within a host, arthropods within soil), if the
141 principal target was the whole organism community DNA. After reviewing the final set of
142 filtered papers, 24 additional papers fitting the selection criteria but not present in the systematic
143 WOS search were also included. A total of 111 articles constituted the set of core papers for
144 subsequent assessment (see Table S2 for a complete list).

145 ***The core papers***

146 All papers were systematically processed to record (i) the research aim and type of samples
147 analysed; (ii) the bioinformatic tasks and pipelines implemented; and (iii) the software tools used
148 and the reproducibility of the bioinformatic procedures employed. A detailed description of this
149 process is provided in the Supplementary Methods and is summarised as follows.

150 The research aim was categorised according to whether the focus was (i) the comparison of
151 molecular and/or bioinformatic procedures for metabarcoding, (ii) a proof-of-concept or
152 feasibility study into the success of metabarcoding for uncovering accurate community data in
153 the taxon/community/biome studied, or (iii) principally the study of ecological patterns and
154 processes. We recorded whether the metabarcoded communities were sampled from marine,

155 freshwater, terrestrial biomes or from a host species, and finally if the targets were invertebrates
156 or vertebrates.

157 Subsequently, the bioinformatic procedures for each paper were systematically parsed to
158 identify the different tasks implemented, i.e. specific bioinformatic actions with a clearly-defined
159 purpose and performed by a single tool. A total of 30 distinct bioinformatic tasks were identified
160 starting from initial procedures on raw sequencing files through to the generation of community
161 tables (see Table 1 for a description of each task). We focused solely on bioinformatic tasks that
162 were presented as necessary for the generation of information about the occurrence or incidence
163 of taxonomic units in the sampled communities (i.e. community data), and the taxonomic
164 identity of these units. For example, we did not record any steps performing phylogenetics with a
165 final OTU set, although we recorded steps where phylogeny-based methods were used as part of
166 OTU delimitation and filtering. Similarly, we recorded tasks that performed filtering of
167 community data for the purposes of removing OTUs or OTU records arising from erroneous
168 sequences or from cross-talk/contamination (Edgar, 2018), but we did not record tasks that
169 filtered community data for the purposes of statistical correction, such as normalisation or
170 rarefaction.

171 Once the different tasks implemented by each article were identified, the pipeline used (i.e.
172 the specific sequence of tasks in a particular order), was also recorded based on the order in
173 which the different tasks were mentioned in the text, figures, supplementary material and/or
174 cited papers. Where multiple mutually exclusive tasks were employed for the purposes of
175 comparison of pipelines, we recorded that pipeline that the authors concluded to be empirically
176 superior, or from which the authors used the output data for subsequent analysis.

177 For each of the bioinformatic tasks identified across the papers, we calculated (i) the number
178 of papers implementing the task, (ii) the task's relative position within the pipelines, (iii) the

179 information reported on the software, version and parameters used, and (iv) the homogeneity in
180 the software tools used to implement the task. We assessed homogeneity by calculating two
181 indices, the *software homogeneity rate* and the *software dominance rate* (see Fig 5). Finally, we
182 also summarised temporal trends in both the reporting and software heterogeneity of each task.

183 **Results and Discussion**

184 *Diversity of bioinformatic methods*

185 The 111 selected papers were published in 36 different journals with a broad focus on ecology
186 and molecular ecology. There has been a steady increase in the number of papers published in
187 this domain over time (Fig. 1). The earliest year of publication was 2011, but 77% of all papers
188 were published in the last four years (2017-2020, n = 86, Fig. 1). Almost all papers studied
189 invertebrate communities (n=108). Forty-five papers were focussed on terrestrial communities,
190 31 on freshwater, 30 on marine and five on parasite communities collected from a host vertebrate
191 (see Table S2 for all the details on the core papers set).

192 Despite a clear trend for increased use of wocDNA COI metabarcoding, the field remains in
193 a relatively early stage of implementation, reflected in the fact that in half of all papers (n=56,
194 n=38 in the last four years) metabarcoding was undertaken as a proof-of-concept and the authors
195 primarily discussed the feasibility of this method for the studied ecological system. Only 25
196 papers considered the sample sizes and metabarcoding procedures sufficiently rigorous to
197 answer ecological questions. Thirty papers were primarily methodological, assessing the
198 influences of primer choice, lab protocols and/or sequencing methods. However, within the
199 methodological category, no paper solely studied the effect of bioinformatic pipeline choices.
200 Indeed, only eight out of the 111 papers clearly stated that they compared different bioinformatic
201 tools for the same task, despite the use of 116 discrete pieces of software or functions in our final
202 count. These results illustrate the timely nature of this review, highlighting the inconsistent
203 implementation of bioinformatic methods, in contrast to the relative maturity and harmonisation
204 of field and laboratory methodologies.

205 ***High heterogeneity in tasks and pipelines***

206 The variety of bioinformatics pipelines reported across the 111 papers employed 108 unique
207 pipelines, i.e. sets of bioinformatics tasks carried out in a specific sequence. Three pipelines
208 were used twice; in two of these cases, a group of authors replicated their pipeline exactly, in the
209 other case the pipeline as reported consisted solely of a single step of searching raw reads against
210 a reference set. Although some of these pipelines were similar, with minor modifications to the
211 order, or the addition/removal of a few tasks, the heterogeneity of pipelines is remarkable. There
212 was also high heterogeneity in the number of tasks implemented within each pipeline, ranging
213 between 1 and 18 tasks, with half of the articles reporting fewer than 9 distinct bioinformatic
214 tasks (Fig. 2a). There was no particular trend in the number of tasks implemented over time
215 (Fig. 2b). The order in which these tasks were implemented also differed greatly (Fig. 2c),
216 although there was a tendency for certain tasks to be performed within similar general stages
217 within pipelines, that is, read preparation-based tasks tend to be implemented at the initial steps
218 of the pipelines, followed by filtering-based tasks and data generation tasks (Fig. 3).

219 Heterogeneity in the sequence of tasks may reflect the careful design and adaptation of
220 bioinformatic procedures within each study to the type and structure of sample and sequence
221 data and/or the specific research question, rather than the simple duplication of previously
222 published pipelines. However, high heterogeneity may equally result from the omission of
223 important tasks or their inappropriate implementation within the pipelines, and so result in low
224 comparability, integration and replication across studies. One clear example of this is associated
225 with the *Filtering* tasks of removal of erroneous sequence reads. Denoising (i.e. the removal of
226 sequencing errors based on models of error frequency parameterised by between-sequence
227 similarity, error sensitivity and/or relative frequency), was employed in just 18 studies and its
228 relative position within the pipelines was highly variable (see Table 1 and Fig. 3). While some
229 sequencing errors will be disregarded during OTU clustering, failure to incorporate denoising

230 can lead to false OTUs and thus OTU inflation (Shum & Palumbi, 2021) Furthermore, the trend
231 towards examining haplotypic variation in metazoan wocDNA metabarcoding through use of
232 amplicon sequence variants (ASVs, Callahan et al., 2017) requires minimising the number of
233 spurious sequences, relying on stringent filtering such as denoising. Similarly, filtering to
234 remove sequences with low copy number (that are often considered highly likely to be
235 erroneous) was reported in only half (n=57) of the studies, despite being generally recommended
236 (Calderón-Sanou et al., 2020; Ficaretola et al., 2017) and a critical step for reducing spurious
237 sequences surviving denoising including nuclear mitochondrial (NUMT, Lopez et al., 1994)
238 copies (Andújar et al., 2021). It should be noted that while many task absences are cases of
239 under-implementation, some may also be underreporting (see below).

240 ***Infrequent adaptation of pipelines to COI***

241 The COI locus differs from many other metabarcoding loci (e.g. 18S, 16S, 12S, ITS) in that it is
242 a protein coding gene, imparting strict expectations of amplicon sequence read properties that
243 can be exploited in metabarcoding bioinformatics (Andújar et al., 2018b). However, the
244 adaptation of pipelines to this fragment are in general rarely implemented in the papers of the
245 core set. For example, only 22 papers (20%) used amino acid translations to identify erroneous
246 sequences (“translation filtering”), using 11 different software tools for the task. The reason for
247 low implementation of translation filtering is likely twofold; first, none of the major
248 metabarcoding software packages include functions for translation filtering, and second, there is
249 no standard straightforward command line software for undertaking this task. Those papers that
250 carry out translation filtering do so by using one of three main approaches: (i) sequences are
251 viewed and translated in a GUI application such as Geneious (<https://www.geneious.com>) or
252 MEGA (Kumar et al., 2018), and those with stop codons manually removed, (ii) sequences are

253 processed through a custom script, some of which are available on github but none of which are
254 used by research groups separate from the author, and (iii) sequences are aligned against
255 references using MACSE (Ranwez et al., 2011) and those containing indels or stop codons are
256 removed. The first option is time consuming and prone to human error, and custom scripts are
257 challenging to document and maintain for a wider number of users. While MACSE is the most
258 frequent single approach, it is computationally efficient only for small datasets. There may be
259 some potential in the recent *coil* R package (Nugent et al., 2020) that uses Hidden Markov
260 Models to identify and filter translation-based errors and appears to scale well to large datasets,
261 although the R implementation presents a slight barrier to efficient inclusion in pipelines.
262 Furthermore, the majority of translation filtering approaches are based solely on removing stop
263 codons, while there may be other potential avenues for filtering based on amino acid translation.
264 The extent to which expectations for protein structural properties can be applied to
265 metabarcoding sequences for filtering other non-synonymous errors has been underexplored (but
266 see Turon et al., 2020).

267 In addition to the potential of amino acid translation, the protein coding nature of COI leads
268 to relatively stricter expectations of amplicon length. However, only half (n=54) of papers
269 reported using length filtering, despite this being a relatively trivial procedure and with functions
270 available in all metabarcoding software packages and as options in many more software tools.
271 There may be some underreporting here; given the implementation of a length filtering
272 parameter in many software tools that have a different primary purpose, authors may not have
273 explicitly reported that length thresholds had been applied as part of a different procedure (note
274 that we recorded when a single tool was reported to have fulfilled multiple tasks). Despite length
275 filtering being widely available, and the relative algorithmic simplicity of implementation, there
276 are no length filtering tools that allow for specification of thresholds outside of a simple
277 minimum-maximum range, despite the internal barcode region of protein coding genes generally

278 being expected to vary in length only by multiples of 3 bases. While trivial to implement this
279 programmatically for an experienced bioinformatician, this lack of straightforward user-friendly
280 availability presents a barrier to appropriate threshold implementation by those with less
281 experience.

282 ***Severe underreporting and increasing heterogeneity in the tools used for bioinformatic tasks***

283 Of the 30 bioinformatic tasks identified (see Table 1 for a description of the tasks), only 11 were
284 implemented in more than half of the papers (n<55) (Fig. 3). Quality filtering (n=92) and OTU
285 delimitation (n=89) were the tasks most reported. Some of the less reported tasks were those
286 associated with uncommon bioinformatic requirements of metabarcoding data, such as assembly
287 or degapping; others have become redundant with modern computational power, such as
288 preclustering. Low reporting of such tasks is likely an accurate reflection of rare implementation;
289 however, there are many other tasks that are fundamental in metabarcoding bioinformatics but
290 are poorly reported. For example, primer trimming was only reported by just over half of the
291 papers (n=67), yet is a completely necessary step. Similarly, adapter trimming was underreported
292 (n=21); while it is likely that in the majority of cases this is implemented by sequencing facilities
293 prior to the authors receiving data, its reporting, including parameters and tools used, is
294 fundamental to verify stringency of the read preparation procedures. The mapping of by-sample
295 reads to OTUs was reported by only one third (n=30) of the papers that employed OTU
296 delimitation, despite this being a necessary step for the production of ecological data for
297 downstream analysis. Furthermore, OTU mapping is not a trivial step; the level of
298 filtering/processing performed on the reads used for mapping (as opposed to filtering/processing
299 performed on the sequences used for OTU delimitation), and the similarity threshold and tie-
300 breaking algorithm employed to assign reads to OTU clusters could all substantially affect the

301 community data generated. The accurate reporting of this step is important to assess the validity
302 of a pipeline, its comparability across studies, and/or its ability to be reproduced.

303 In addition to the clear underreporting of tasks within the pipelines as discussed above, the
304 reporting of the bioinformatic tools and parameters used for those tasks cited in the papers was
305 also very poor (Table 1). Only 21 of the 111 papers reported software name, version and
306 parameters used for all of the bioinformatic tasks implemented, and 25 failed on all three counts
307 (Fig. 4a). When considering the degree of underreporting by task (Fig. 4b), the most
308 underreported software were used for some of the most perfunctory tasks (e.g. frequency
309 filtering, length filtering, dereplication) that can be easily reproduced using many equivalent
310 tools. Nonetheless, there remains relatively widespread underreporting, and this has remained
311 unchanged over time (Fig. 5b).

312 Within the reported software, we identified 93 software tools used in metabarcoding
313 bioinformatic pipelines (Table S3), of which 27% (25) were software ‘packages’. When taking
314 into account distinct functions within packages, a total of 169 unique tools were recorded,
315 however, this is likely an inaccurate picture given low reporting rates of functions used within
316 software packages across all steps. There is a clear increase in the number of different software
317 and software functions employed across all papers over time (Fig. 5a). Examining the diversity
318 of software used within tasks over time, controlling for the number of papers published, there is
319 limited improvement in homogeneity and a decrease in dominance of software (Figs 5c and 5d).
320 Given that the number of metabarcoding publications is increasing year-on-year, there is thus a
321 concomitant increase in the diversity of software used for a given task, and previously
322 commonly used software are being used less (Figs 5c and 5d). These trends reflect that while
323 new software tools are constantly being made available for metabarcoding, uptake is not
324 consistent across the field and while some researchers use more recent tools, many researchers
325 continue to use older methods, diversifying the field.

326 ***Toward a bioinformatic harmonisation of COI metabarcoding for metazoan wocDNA samples***

327 Our results show that the increasing use of wocDNA COI metabarcoding for metazoan diversity
328 is characterised by a clear absence of bioinformatic harmonisation, and the temporal trends show
329 little change in this situation. The reviewed literature showed (i) high heterogeneity across
330 pipelines, tasks and tools used, (ii) limited or no adaptation of bioinformatic procedures to the
331 nature of the COI fragment, and (iii) a worrying underreporting of tasks, software and
332 parameters.

333 The development of metabarcoding as a method for community ecology began with
334 microbial studies over a decade ago, which have revealed the extensive diversity of bacteria and
335 archaea on our planet and demonstrated the potential of metabarcoding for global biodiversity
336 syntheses (Bates et al., 2013; Thompson et al., 2017). Although the integration and meta-analysis
337 of microbial community data from independent studies is still challenging (e.g. Ramirez-
338 Gonzalez et al., 2013), the success of international consortia such as the Earth Microbiome
339 Project (EMP, Gilbert et al., 2010, 2014) has promoted the development of a harmonised
340 framework for data generation and analyses within microbial eDNA research (see e.g. Tedersoo
341 et al., 2015).

342 Through the adaptation of the microbial metabarcoding method to wocDNA samples,
343 specific protocols to sample, sort and enrich community samples for wocDNA metabarcoding
344 have been developed, targeting different taxonomic fractions and types of samples (e.g., Andújar
345 et al., 2018a; Arribas et al., 2016; Creedy et al., 2019; Elbrecht & Leese, 2017; Fonseca et al.,
346 2010; Yu et al., 2012). Additionally, recent efforts to adapt and optimise existing methods are
347 increasing efficiency and versatility, for example through non-destructive DNA extraction
348 techniques that retain specimens for morphological vouchering (Marquina et al., 2019; Nielsen

349 et al., 2019), or library preparation techniques tailored to metazoan samples (Yang et al., 2020).
350 Although wocDNA COI metabarcoding remains in an expansive phase of development,
351 standardisation in field and laboratory methods are emerging. This is in part boosted by
352 collaborative initiatives such as the BIOSCAN initiative and its regional extensions (e.g.
353 BIOALPHA), the Kruger Malaise Program, SITE-100, the Insect Biome Atlas Project,
354 LIFEPLAN, and iBioGen (Arribas et al., 2021).

355 In contrast, there has been little advance in the development and validation of best practices
356 associated with the bioinformatics processing of wocDNA COI metabarcoding data (but see
357 Yang et al., 2020 for error reduction). Outside of taxonomic assignment, discussion of
358 customising or parameterising tools for the purposes of working with wocDNA COI
359 metabarcoding is very rare, with most papers simply reporting using tools with default settings.
360 Our review has revealed heterogeneity in the number of tasks, the order of these within
361 pipelines, and the tools used to implement them, along with a lack of even basic adaptations to
362 the COI metabarcode for most of the papers. The majority of available software and resources
363 for metabarcoding bioinformatics are still those that have been developed around the 16S rRNA
364 gene (the primary target for microbiome metabarcoding), including the most popular software
365 packages (e.g. USEARCH) and sets of wrapper scripts (e.g. QIIME, OBITools). While in many
366 cases these methods may carry over to COI without issue, we observe very few studies that
367 report consideration or analysis that assesses or validates the suitability of software choices for
368 COI. These issues suggest that the expansion of wocDNA COI metabarcoding is proceeding at a
369 pace and manner that could lose sight of or simply ignore the challenges inherent in producing
370 high-quality data and reproducible methods (Baker et al., 2016; Zinger et al., 2019), and lose out
371 on the potential for exploiting the benefits of the COI marker for wocDNA metabarcoding of
372 Metazoa.

373 DNA metabarcoding has broad multidisciplinary potential, as demonstrated by the
374 expansion in use of metazoan wocDNA COI metabarcoding among users from very diverse
375 backgrounds. The diversity of applications of metabarcoding requires the concomitant
376 bioinformatic techniques to be flexible and adaptable, and the field remains under active
377 development. Thus it would not be productive to attempt to prescribe pipelines, tasks or even
378 software tools in the name of standardisation, as there is no one-size-fits-all approach in
379 metabarcoding. However, some degree of harmonisation is required to ensure quality,
380 reproducibility and potential integration in metastudies (Tedersoo et al., 2015). Additionally, the
381 absence of a harmonised framework of bioinformatic processing can act as a barrier for potential
382 new users (Liu et al., 2020), hampering the growth of the field. To these ends, we thus propose a
383 set of recommendations that we believe all researchers in the field should consider when
384 designing and reporting their wocDNA COI metabarcoding bioinformatics pipeline, with the
385 hope that they will catalyse harmonised implementation.

386 **Fully report all tasks, software, software versions and parameters used, even if just the**
387 **defaults.** Our results show that underreporting is a recurrent problem. Comprehensive reporting
388 of the tasks, pipelines and software used is essential for further integrating results in future
389 reviews or meta-analyses (Tedersoo et al., 2015). Furthermore, care should be taken to report not
390 just the name of the software package, but also the exact function, and if wrapper scripts are used
391 then the underlying functions should be reported. Considering the trade-off with current
392 constrictions for manuscript length, this could be achieved by the inclusion of a supporting table
393 following the STAR-METHODs philosophy (Marcus, 2016), where task reference, order within
394 the pipeline and software used are included. Note that the task lexicon and software lists
395 compiled in this review (see Table 1) are a very useful resource for this purpose. This reporting
396 effort for all the wocDNA COI metabarcoding will promote rigour and robustness with an

397 intuitive, consistent framework that makes reporting easier for the author and replication easier
398 for the reader.

399 **Implement filtering tasks such that spurious sequences are sufficiently removed to**
400 **meet the assumptions of the research question.** The quality of metabarcoding results is likely
401 to depend most on the appropriate inclusion of filtering into a pipeline (Calderón-Sanou et al.,
402 2020; Elbrecht et al., 2018; Zinger et al., 2019), so proper implementation of filtering tasks are
403 critical for robust and harmonised use of COI metabarcoding. In metabarcoding, real amplicon
404 sequence variants (ASVs, Callahan et al., 2017) amplified from target genes are inherently
405 accompanied by spurious sequences, arising from multiple sources. Indeed, taxonomic inflation
406 is a recurring issue demonstrated in communities with known haplotype composition (Creedy et
407 al., 2020; Elbrecht et al., 2018). This can be exacerbated for mitochondrial markers like COI,
408 due to the co-amplification of NUMTs and other non-authentic ASVs that are missed by
409 denoising and require stringent, optimised filtering based on read abundances such as that
410 implemented by the *metaMATE* software (Andújar et al., 2021). To ensure quality and
411 reproducibility, metabarcoding studies should consider implementing the six most common
412 filtering approaches, i.e Quality, Length, Chimera, Translation, and Frequency filtering, plus
413 Denoising. For each of these tasks, appropriate thresholds should be considered, implemented
414 and fully reported to a level that ensures reproducibility. Given the demonstrated importance of
415 these tasks for most wocDNA metabarcoding studies, if any are not employed by a study the
416 omission should be explained.

417 **Adapt pipelines to the COI fragment.** Suitable adaptations include read processing and
418 filtering steps that leverage evolutionary properties of the protein coding nature of this fragment,
419 or determining appropriate parameters for tools originally designed for other DNA regions.
420 Some recent advances have been made in filtering tasks (*metaMATE*, Andújar et al., 2021; *coil*,
421 Nugent et al., 2020; entropy-based denoising, Turon et al., 2020) but further development in

422 these promising areas is essential to fully exploit the potential of the COI gene for
423 metabarcoding. As mentioned previously, there are no tools that enable simple length filtering
424 variation that accounts for codon-level insertion or deletion. To our knowledge there is limited
425 work exploring the extent to which protein structure inference might allow identification of
426 erroneous sequences: for example the *SOAPbarcode* pipeline (Liu et al., 2013) includes a script
427 that filters sequences based on translation hydrophilicity, but this is not comprehensively
428 documented or discussed in the associated publications. Computation of protein structural
429 properties is relatively trivial to perform, and seems like a fertile ground for novel development
430 of filtering tools for protein coding markers.

431 **For each task, consider all software available and try to select the most appropriate**
432 **tool(s).** This can only be approached with sufficient information about available software, and to
433 this end we include a list of all software used for each task within Table S1, and Table S3
434 includes links to documentation and publications. The selection of the most appropriate tool is
435 not always straightforward, but we suggest considering (i) the extent to which the tool was
436 designed for the intended barcode region, purpose or dataset, (ii) the detail of available
437 documentation and explanation to ensure a tool performs as expected, (iii) the availability and
438 flexibility of options to appropriately apply the tool, (iv) the frequency of use of a tool in other
439 studies with similar research aims, and (v) all else being equal, the simplest approach. Ideally,
440 where multiple approaches exist, reasonable comparison between key methods should take place
441 to fully understand the potential variation in conclusions that might arise from different
442 bioinformatic choices, and the results of these comparisons should be reported. This is
443 particularly the case when considering alternative, conceptually distinct algorithms for more
444 bioinformatically complex tasks, such as denoising and OTU delimitation. The development of
445 software packages and open access platforms integrating a catalogue of common bioinformatic
446 tools, such mBRAVE (<http://www.mbrave.net/>), may play a fundamental role towards a proper

447 selection and harmonisation of the software used. However, software choices should be made on
448 the basis of appropriateness and usefulness, rather than simply ease of availability and
449 implementation due to inclusion in these packages/platforms.

450 **Verify the compatibility of the tasks within a pipeline, especially with respect to task**
451 **order.** It is important to ensure that the assumptions of one task have not been violated by
452 upstream processing; for example, UNOISE denoising employs a model of error rates in
453 Illumina sequencing, and if errors have been removed by prior length or frequency filtering this
454 model may not accurately fit to the data. Further, linked processes should be compatible: for
455 instance, if OTU delimitation is based on a linkage algorithm such as swarm (Mahé et al., 2015),
456 it is inappropriate to employ a simple similarity-based mapping method to assign reads to the
457 resultant OTUs.

458 Aside from these recommendations, we also urge researchers to make data publicly
459 available, both raw reads and final ASV and/or OTU sequences. Raw read datasets will become
460 an invaluable resource for future work integrating many wocDNA metabarcoding studies across
461 spatial and temporal scales, with continuing development and improvement of bioinformatic
462 pipelines allowing for forward-compatibility of the data as analytical tools continue to evolve.
463 Uploading ASV and/or OTU sequences, even with incomplete taxonomy, improves the
464 capability of methods for taxonomic assignment that draw on these resources and provides fertile
465 datasets for future development of bioinformatic methods.

466 **Conclusions**

467 The past decade has seen rapid growth in the development, testing and use of wocDNA COI
468 metabarcoding. Much effort has been expended in the development of laboratory, sequencing
469 and bioinformatic methodologies for wocDNA COI metabarcoding and for metabarcoding as a

470 whole. However, while much progress has been made towards harmonisation of lab and
471 sequencing methods, bioinformatic processes have remained a tangle of varying software,
472 pipelines and theoretical approaches, often suffering from underreported detail. This diversity
473 allows for versatility, especially for those who are well-informed and experienced in
474 bioinformatics and able to pick and choose the appropriate approach. However, choosing from
475 the range of approaches could easily hinder new applications of metabarcoding for researchers
476 coming from a limited bioinformatic background, and high heterogeneity can stymie the
477 potential for future reviews and meta-analyses. Our review, which is the first evaluating the state
478 of the art on this topic, highlights that this danger is clearly present in the field of metazoan
479 wocDNA COI metabarcoding. The results of our assessment and the recommendations derived
480 from it may help to improve bioinformatic harmonisation and thus the long-term integrative
481 potential of wocDNA COI metabarcoding for biodiversity science.

482 **Acknowledgements**

483 This research was supported by the iBioGen project, which has received funding from
484 the European Union's Horizon 2020 research and innovation programme under grant agreement
485 No 810729.

486 **Conflict of interest**

487 A.P.V. is a co-founder and scientific advisor of NatureMetrics, a private company providing
488 commercial services in DNA-based monitoring. The authors declare that they have no other
489 conflicts of interest.

490 **Author contributions**

491 T.J.C. and P.A, conceived the study. T.J.C. and P.A. assessed the initial paper set for
492 inclusion, T.J.C. evaluated the methods of the core paper set and analysed the data. T.J.C. and
493 P.A. wrote the initial draft and all co authors contributed to the final manuscript.

494 **Data accessibility**

495 Supporting Materials (methods, figures and tables) give the full details and methodological
496 evaluation of the 111 publications making up the core papers.

497 **References**

- 498 Andújar, C., Arribas, P., Gray, C., Bruce, C., Woodward, G., Yu, D. W., & Vogler, A. P. (2018a).
499 Metabarcoding of freshwater invertebrates to detect the effects of a pesticide spill. *Molecular*
500 *Ecology*, 27(1), 146–166. <https://doi.org/10.1111/mec.14410>
- 501 Andújar, C., Arribas, P., Yu, D. W., Vogler, A. P., & Emerson, B. C. (2018b). Why the COI barcode should
502 be the community DNA metabarcode for the metazoa. *Molecular Ecology*, 27(20), 3968–3975.
503 <https://doi.org/10.1111/mec.14844>
- 504 Andújar, C., Creedy, T. J., Arribas, P., López, H., Salces-Castellano, A., Pérez-Delgado, A. J., Vogler, A.
505 P., & Emerson, B. C. (2021). Validated removal of nuclear pseudogenes and sequencing artefacts
506 from mitochondrial metabarcode data. *Molecular Ecology Resources*.
507 <https://doi.org/10.1111/1755-0998.13337>
- 508 Arribas, P., Andújar, C., Bidartondo, M. I., Bohmann, K., Coissac, É., Creer, S., deWaard, J. R., Elbrecht,
509 V., Ficetola, G. F., Goberna, M., Kennedy, S., Krehenwinkel, H., Leese, F., Novotny, V.,
510 Ronquist, F., Yu, D. W., Zinger, L., Creedy, T. J., Meramveliotakis, E., ... Emerson, B. C. (2021).
511 Connecting high-throughput biodiversity inventories: Opportunities for a site-based genomic
512 framework for global integration and synthesis. *Molecular Ecology*, 30(5), 1120–1135.
513 <https://doi.org/10.1111/mec.15797>
- 514 Arribas, P., Andújar, C., Hopkins, K., Shepherd, M., & Vogler, A. P. (2016). Metabarcoding and
515 mitochondrial metagenomics of endogean arthropods to unveil the mesofauna of the soil.
516 *Methods in Ecology and Evolution*, 7(9), 1071–1081. <https://doi.org/10.1111/2041-210X.12557>
- 517 Arribas, P., Andújar, C., Salces-Castellano, A., Emerson, B. C., & Vogler, A. P. (2021). The limited spatial
518 scale of dispersal in soil arthropods revealed with whole-community haplotype-level
519 metabarcoding. *Molecular Ecology*, 30(1), 48–61. <https://doi.org/10.1111/mec.15591>
- 520 Baker, C. C. M., Bittleston, L. S., Sanders, J. G., & Pierce, N. E. (2016). Dissecting host-associated
521 communities with DNA barcodes. *Philosophical Transactions of the Royal Society B: Biological*

522 *Sciences*, 371(1702), 20150328. <https://doi.org/10.1098/rstb.2015.0328>

523 Bates, S. T., Clemente, J. C., Flores, G. E., Walters, W. A., Parfrey, L. W., Knight, R., & Fierer, N. (2013).
524 Global biogeography of highly diverse protistan communities in soil. *The ISME Journal*, 7(3),
525 652–659. <https://doi.org/10.1038/ismej.2012.147>

526 Boyer, F., Mercier, C., Bonin, A., Bras, Y. L., Taberlet, P., & Coissac, E. (2016). obitools: A unix-inspired
527 software package for DNA metabarcoding. *Molecular Ecology Resources*, 16(1), 176–182.
528 <https://doi.org/10.1111/1755-0998.12428>

529 Braukmann, T. W. A., Ivanova, N. V., Prosser, S. W. J., Elbrecht, V., Steinke, D., Ratnasingham, S.,
530 Waard, J. R. de, Sones, J. E., Zakharov, E. V., & Hebert, P. D. N. (2019). Metabarcoding a diverse
531 arthropod mock community. *Molecular Ecology Resources*, 19(3), 711–727.
532 <https://doi.org/10.1111/1755-0998.13008>

533 Bush, A., Compson, Z. G., Monk, W. A., Porter, T. M., Steeves, R., Emilson, E., Gagne, N., Hajibabaei,
534 M., Roy, M., & Baird, D. J. (2019). Studying Ecosystems With DNA Metabarcoding: Lessons
535 From Biomonitoring of Aquatic Macroinvertebrates. *Frontiers in Ecology and Evolution*, 7.
536 <https://doi.org/10.3389/fevo.2019.00434>

537 Bush, A., Monk, W. A., Compson, Z. G., Peters, D. L., Porter, T. M., Shokralla, S., Wright, M. T. G.,
538 Hajibabaei, M., & Baird, D. J. (2020). DNA metabarcoding reveals metacommunity dynamics in
539 a threatened boreal wetland wilderness. *Proceedings of the National Academy of Sciences*,
540 117(15), 8539–8545. <https://doi.org/10.1073/pnas.1918741117>

541 Calderón-Sanou, I., Münkemüller, T., Boyer, F., Zinger, L., & Thuiller, W. (2020). From environmental
542 DNA sequences to ecological conclusions: How strong is the influence of methodological
543 choices? *Journal of Biogeography*, 47(1), 193–206. <https://doi.org/10.1111/jbi.13681>

544 Callahan, B. J., McMurdie, P. J., & Holmes, S. P. (2017). Exact sequence variants should replace
545 operational taxonomic units in marker-gene data analysis. *The ISME Journal*, 11(12), 2639–2643.
546 <https://doi.org/10.1038/ismej.2017.119>

547 Caporaso, J. G., Kuczynski, J., Stombaugh, J., Bittinger, K., Bushman, F. D., Costello, E. K., Fierer, N.,
548 Peña, A. G., Goodrich, J. K., Gordon, J. I., Huttley, G. A., Kelley, S. T., Knights, D., Koenig, J.
549 E., Ley, R. E., Lozupone, C. A., McDonald, D., Muegge, B. D., Pirrung, M., ... Knight, R.
550 (2010). QIIME allows analysis of high-throughput community sequencing data. *Nature Methods*,
551 7(5), 335–336. <https://doi.org/10.1038/nmeth.f.303>

552 Creedy, T. J., Ng, W. S., & Vogler, A. P. (2019). Toward accurate species-level metabarcoding of
553 arthropod communities from the tropical forest canopy. *Ecology and Evolution*, 9(6), 3105–3116.
554 <https://doi.org/10.1002/ece3.4839>

555 Creedy, T. J., Norman, H., Tang, C. Q., Chin, K. Q., Andujar, C., Arribas, P., O'Connor, R. S., Carvell, C.,
556 Notton, D. G., & Vogler, A. P. (2020). A validated workflow for rapid taxonomic assignment and
557 monitoring of a national fauna of bees (Apiformes) using high throughput DNA barcoding.
558 *Molecular Ecology Resources*, 20(1), 40–53. <https://doi.org/10.1111/1755-0998.13056>

559 Deiner, K., Bik, H. M., Mächler, E., Seymour, M., Lacoursière-Roussel, A., Altermatt, F., Creer, S., Bista,
560 I., Lodge, D. M., Vere, N. de, Pfrender, M. E., & Bernatchez, L. (2017). Environmental DNA
561 metabarcoding: Transforming how we survey animal and plant communities. *Molecular Ecology*,
562 26(21), 5872–5895. <https://doi.org/10.1111/mec.14350>

563 deWaard, J. R., Levesque-Beaudin, V., deWaard, S. L., Ivanova, N. V., McKeown, J. T. A., Miskie, R.,
564 Naik, S., Perez, K. H. J., Ratnasingham, S., Sobel, C. N., Sones, J. E., Steinke, C., Telfer, A. C.,
565 Young, A. D., Young, M. R., Zakharov, E. V., & Hebert, P. D. N. (2019). Expedited assessment of
566 terrestrial arthropod diversity by coupling Malaise traps with DNA barcoding. *Genome*, 62(3),
567 85–95. <https://doi.org/10.1139/gen-2018-0093>

568 Edgar, R. C. (2013). UPARSE: Highly accurate OTU sequences from microbial amplicon reads. *Nature*
569 *Methods*, 10(10), 996–998. <https://doi.org/10.1038/nmeth.2604>

570 Edgar, R. C. (2018). UNCROSS2: Identification of cross-talk in 16S rRNA OTU tables. *BioRxiv*, 400762.
571 <https://doi.org/10.1101/400762>

572 Elbrecht, V., & Leese, F. (2017). Validation and Development of COI Metabarcoding Primers for
573 Freshwater Macroinvertebrate Bioassessment. *Frontiers in Environmental Science*, 5.
574 <https://doi.org/10.3389/fenvs.2017.00011>

575 Elbrecht, V., Vamos, E. E., Steinke, D., & Leese, F. (2018). Estimating intraspecific genetic diversity from
576 community DNA metabarcoding data. *PeerJ*, 6, e4644. <https://doi.org/10.7717/peerj.4644>

577 Elbrecht, V., Braukmann, T. W. A., Ivanova, N. V., Prosser, S. W. J., Hajibabaei, M., Wright, M.,
578 Zakharov, E. V., Hebert, P. D. N., & Steinke, D. (2019). Validation of COI metabarcoding primers
579 for terrestrial arthropods. *PeerJ*, 7, e7745. <https://doi.org/10.7717/peerj.7745>

580 Ficetola, G. F., Mazel, F., & Thuiller, W. (2017). Global determinants of zoogeographical boundaries.
581 *Nature Ecology & Evolution*, 1(4), 1–7. <https://doi.org/10.1038/s41559-017-0089>

582 Folmer, O., Black, M., Hoeh, W., Lutz, R., & Vrijenhoek, R. (1994). DNA primers for amplification of
583 mitochondrial cytochrome c oxidase subunit I from diverse metazoan invertebrates. *Molecular*
584 *Marine Biology and Biotechnology*, 3(5), 294–299.

585 Fonseca, V. G., Carvalho, G. R., Sung, W., Johnson, H. F., Power, D. M., Neill, S. P., Packer, M., Blaxter,
586 M. L., Lamshead, P. J. D., Thomas, W. K., & Creer, S. (2010). Second-generation environmental
587 sequencing unmasks marine metazoan biodiversity. *Nature Communications*, 1(1), 98.
588 <https://doi.org/10.1038/ncomms1095>

589 Fonseca, V. G., Packer, M., Carvalho, G. R., Power, D. M., Lamshead, P. J. D., & Creer, S. (2011).
590 Isolation of marine meiofauna from sandy sediments: From decanting to DNA extraction.
591 *Protocol Exchange*. <https://doi.org/10.1038/nprot.2010.157>

592 Gilbert, J. A., Jansson, J. K., & Knight, R. (2014). The Earth Microbiome project: Successes and
593 aspirations. *BMC Biology*, 12(1), 69. <https://doi.org/10.1186/s12915-014-0069-1>

594 Gilbert, J. A., Meyer, F., Jansson, J., Gordon, J., Pace, N., Tiedje, J., Ley, R., Fierer, N., Field, D.,
595 Kyripides, N., Glöckner, F.-O., Klenk, H.-P., Wommack, K. E., Glass, E., Docherty, K., Gallery,

- 596 R., Stevens, R., & Knight, R. (2010). The Earth Microbiome Project: Meeting report of the “1st
597 EMP meeting on sample selection and acquisition” at Argonne National Laboratory October 6th
598 2010. *Standards in Genomic Sciences*, 3(3), 249. <https://doi.org/10.4056/aigs.1443528>
- 599 Hajibabaei, M., Spall, J. L., Shokralla, S., & van Konynenburg, S. (2012). Assessing biodiversity of a
600 freshwater benthic macroinvertebrate community through non-destructive environmental
601 barcoding of DNA from preservative ethanol. *BMC Ecology*, 12(1), 28.
602 <https://doi.org/10.1186/1472-6785-12-28>
- 603 Ji, Y., Ashton, L., Pedley, S. M., Edwards, D. P., Tang, Y., Nakamura, A., Kitching, R., Dolman, P. M.,
604 Woodcock, P., Edwards, F. A., Larsen, T. H., Hsu, W. W., Benedick, S., Hamer, K. C., Wilcove,
605 D. S., Bruce, C., Wang, X., Levi, T., Lott, M., ... Yu, D. W. (2013). Reliable, verifiable and
606 efficient monitoring of biodiversity via metabarcoding. *Ecology Letters*, 16(10), 1245–1257.
607 <https://doi.org/10.1111/ele.12162>
- 608 Krehenwinkel, H., Wolf, M., Lim, J. Y., Rominger, A. J., Simison, W. B., & Gillespie, R. G. (2017).
609 Estimating and mitigating amplification bias in qualitative and quantitative arthropod
610 metabarcoding. *Scientific Reports*, 7(1), 17668. <https://doi.org/10.1038/s41598-017-17333-x>
- 611 Kumar, S., Stecher, G., Li, M., Knyaz, C., & Tamura, K. (2018). MEGA X: Molecular Evolutionary
612 Genetics Analysis across Computing Platforms. *Molecular Biology and Evolution*, 35(6), 1547–
613 1549. <https://doi.org/10.1093/molbev/msy096>
- 614 Leese, F., Bouchez, A., Abarenkov, K., Altermatt, F., Borja, Á., Bruce, K., Ekrem, T., Čiampor, F.,
615 Čiamporová-Zaťovičová, Z., Costa, F. O., Duarte, S., Elbrecht, V., Fontaneto, D., Franc, A.,
616 Geiger, M. F., Hering, D., Kahlert, M., Kalamujić Stroil, B., Kelly, M., ... Weigand, A. M.
617 (2018). Chapter Two - Why We Need Sustainable Networks Bridging Countries, Disciplines,
618 Cultures and Generations for Aquatic Biomonitoring 2.0: A Perspective Derived From the
619 DNAqua-Net COST Action. In D. A. Bohan, A. J. Dumbrell, G. Woodward, & M. Jackson (Eds.),
620 *Advances in Ecological Research* (Vol. 58, pp. 63–99). Academic Press.
621 <https://doi.org/10.1016/bs.aecr.2018.01.001>
- 622 Liu, M., Clarke, L. J., Baker, S. C., Jordan, G. J., & Burridge, C. P. (2020). A practical guide to DNA
623 metabarcoding for entomological ecologists. *Ecological Entomology*, 45(3), 373–385.
624 <https://doi.org/10.1111/een.12831>
- 625 Liu, S., Li, Y., Lu, J., Su, X., Tang, M., Zhang, R., Zhou, L., Zhou, C., Yang, Q., Ji, Y., Yu, D. W., & Zhou,
626 X. (2013). SOAPBarcode: Revealing arthropod biodiversity through assembly of Illumina
627 shotgun sequences of PCR amplicons. *Methods in Ecology and Evolution*, 4(12), 1142–1150.
628 <https://doi.org/10.1111/2041-210X.12120>
- 629 Lopez, J. V., Yuhki, N., Masuda, R., Modi, W., & O'Brien, S. J. (1994). Numt, a recent transfer and
630 tandem amplification of mitochondrial DNA to the nuclear genome of the domestic cat. *Journal*
631 *of Molecular Evolution*, 39(2), 174–190. <https://doi.org/10.1007/BF00163806>
- 632 Mahé, F., Rognes, T., Quince, C., Vargas, C. de, & Dunthorn, M. (2015). Swarm v2: Highly-scalable and

633 high-resolution amplicon clustering. *PeerJ*, 3, e1420. <https://doi.org/10.7717/peerj.1420>

634 Marcus, E. (2016). A STAR Is Born. *Cell*, 166(5), 1059–1060. <https://doi.org/10.1016/j.cell.2016.08.021>

635 Marquina, D., Esparza-Salas, R., Roslin, T., & Ronquist, F. (2019). Establishing arthropod community
636 composition using metabarcoding: Surprising inconsistencies between soil samples and
637 preservative ethanol and homogenate from Malaise trap catches. *Molecular Ecology Resources*,
638 19(6), 1516–1530. <https://doi.org/10.1111/1755-0998.13071>

639 Nielsen, M., Gilbert, M. T. P., Pape, T., & Bohmann, K. (2019). A simplified DNA extraction protocol for
640 unsorted bulk arthropod samples that maintains exoskeletal integrity. *Environmental DNA*, 1(2),
641 144–154. <https://doi.org/10.1002/edn3.16>

642 Nugent, C. M., Elliott, T. A., Ratnasingham, S., & Adamowicz, S. J. (2020). coil: An R package for
643 cytochrome *c* oxidase I (COI) DNA barcode data cleaning, translation, and error evaluation.
644 *Genome*, 63(6), 291–305. <https://doi.org/10.1139/gen-2019-0206>

645 Porter, T. M., & Hajibabaei, M. (2018). Scaling up: A guide to high-throughput genomic approaches for
646 biodiversity analysis. *Molecular Ecology*, 27(2), 313–338. <https://doi.org/10.1111/mec.14478>

647 Ramirez-Gonzalez, R., Yu, D. W., Bruce, C., Heavens, D., Caccamo, M., & Emerson, B. C. (2013).
648 PyroClean: Denoising Pyrosequences from Protein-Coding Amplicons for the Recovery of
649 Interspecific and Intraspecific Genetic Variation. *PLOS ONE*, 8(3), e57615.
650 <https://doi.org/10.1371/journal.pone.0057615>

651 Ranwez, V., Harispe, S., Delsuc, F., & Douzery, E. J. P. (2011). MACSE: Multiple Alignment of Coding
652 SEquences Accounting for Frameshifts and Stop Codons. *PLOS ONE*, 6(9), e22594.
653 <https://doi.org/10.1371/journal.pone.0022594>

654 Rognes, T., Flouri, T., Nichols, B., Quince, C., & Mahé, F. (2016). VSEARCH: A versatile open source
655 tool for metagenomics. *PeerJ*, 4, e2584. <https://doi.org/10.7717/peerj.2584>

656 Shum, P., & Palumbi, S. R. (2021). Testing small-scale ecological gradients and intraspecific
657 differentiation for hundreds of kelp forest species using haplotypes from metabarcoding.
658 *Molecular Ecology*. <https://doi.org/10.1111/mec.15851>

659 Schnell, I. B., Thomsen, P. F., Wilkinson, N., Rasmussen, M., Jensen, L. R. D., Willerslev, E., Bertelsen,
660 M. F., & Gilbert, M. T. P. (2012). Screening mammal biodiversity using DNA from leeches.
661 *Current Biology*, 22(8), R262–R263. <https://doi.org/10.1016/j.cub.2012.02.058>

662 Sogin, M. L., Morrison, H. G., Huber, J. A., Welch, D. M., Huse, S. M., Neal, P. R., Arrieta, J. M., &
663 Herndl, G. J. (2006). Microbial diversity in the deep sea and the underexplored “rare biosphere”.
664 *Proceedings of the National Academy of Sciences*, 103(32), 12115–12120.
665 <https://doi.org/10.1073/pnas.0605127103>

666 Sogin, M. L., Morrison, H. G., Huber, J. A., Welch, D. M., Huse, S. M., Neal, P. R., Arrieta, J. M., &
667 Herndl, G. J. (2006). Microbial diversity in the deep sea and the underexplored “rare biosphere”.
668 *Proceedings of the National Academy of Sciences*, 103(32), 12115–12120.
669 <https://doi.org/10.1073/pnas.0605127103>

670 Taberlet, P., Coissac, E., Pompanon, F., Brochmann, C., & Willerslev, E. (2012). Towards next-generation
671 biodiversity assessment using DNA metabarcoding. *Molecular Ecology*, 21(8), 2045–2050.
672 <https://doi.org/10.1111/j.1365-294X.2012.05470.x>

673 Tedersoo, L., Anslan, S., Bahram, M., Põlme, S., Riit, T., Liiv, I., Kõljalg, U., Kisand, V., Nilsson, H.,
674 Hildebrand, F., Bork, P., & Abarenkov, K. (2015). Shotgun metagenomes and multiple primer
675 pair-barcode combinations of amplicons reveal biases in metabarcoding analyses of fungi.
676 *MycoKeys*, 10, 1–43. <https://doi.org/10.3897/mycokeys.10.4852>

677 Thompson, L. R., Sanders, J. G., McDonald, D., Amir, A., Ladau, J., Locey, K. J., Prill, R. J., Tripathi, A.,
678 Gibbons, S. M., Ackermann, G., Navas-Molina, J. A., Janssen, S., Kopylova, E., Vázquez-Baeza,
679 Y., González, A., Morton, J. T., Mirarab, S., Zech Xu, Z., Jiang, L., ... Knight, R. (2017). A
680 communal catalogue reveals Earth’s multiscale microbial diversity. *Nature*, 551(7681), 457–463.
681 <https://doi.org/10.1038/nature24621>

682 Turon, X., Antich, A., Palacín, C., Præbel, K., & Wangenstein, O. S. (2020). From metabarcoding to
683 metaphylogeography: Separating the wheat from the chaff. *Ecological Applications*, 30(2),
684 e02036. <https://doi.org/10.1002/eap.2036>

685 Yang, C., Bohmann, K., Wang, X., Cai, W., Wales, N., Ding, Z., Gopalakrishnan, S., & Yu, D. W. (2020).
686 Biodiversity Soup II: A bulk-sample metabarcoding pipeline emphasizing error reduction.
687 *BioRxiv*, 2020.07.07.187666. <https://doi.org/10.1101/2020.07.07.187666>

688 Yu, D. W., Ji, Y., Emerson, B. C., Wang, X., Ye, C., Yang, C., & Ding, Z. (2012). Biodiversity soup:
689 Metabarcoding of arthropods for rapid biodiversity assessment and biomonitoring. *Methods in*
690 *Ecology and Evolution*, 3(4), 613–623. <https://doi.org/10.1111/j.2041-210X.2012.00198.x>

691 Zinger, L., Bonin, A., Alsos, I. G., Bálint, M., Bik, H., Boyer, F., Chariton, A. A., Creer, S., Coissac, E.,
692 Deagle, B. E., Barba, M. D., Dickie, I. A., Dumbrell, A. J., Ficetola, G. F., Fierer, N., Fumagalli,
693 L., Gilbert, M. T. P., Jarman, S., Jumpponen, A., ... Taberlet, P. (2019). DNA metabarcoding—
694 Need for robust experimental designs to draw sound ecological conclusions. *Molecular Ecology*,
695 28(8), 1857–1862. <https://doi.org/10.1111/mec.15060>

696 **Tables**

697 **Table 1:** Table of all bioinformatic tasks performed across the core papers set. Tasks are
698 grouped into four groups by broad purposes, and a detailed definition of each task is given along
699 with summary statistics of the implementation of each task across the 111 papers. For a list of
700 the software used for each task, Table S1 is an expanded version of this table.

701 **Figures**

702 **Figure 1:** Year of publication of the articles in the core papers set. Bar fills and numbers
703 refer to the number of articles within each research aim category. Note that only articles indexed
704 by Web of Science by 3rd November 2020 were included.

705 **Figure 2:** Bioinformatic pipelines implemented by the core papers set. A) Frequency
706 distribution of the number of tasks by study, B) Number of tasks by study against the year of
707 publication, with best fit regression line in blue with shaded 95% confidence intervals around the
708 line. Slight horizontal jitter added to points to better show density. C) Network diagram of tasks
709 and different pipeline routes through these tasks. All pipelines start and end on the respective
710 orange nodes. All other nodes are coloured according to the four main categories of
711 bioinformatic tasks; red for read preparation tasks, blue for sequence processing, green for
712 filtering and purple for data generation tasks. Arrows link tasks performed consecutively, with
713 direction of arrow showing order of tasks. Thickness of arrows shows relative frequency of pairs
714 of consecutive tasks. Arrows coloured orange are the top 10% of consecutive task pairs by
715 relative frequency. Note that while this illustrates a possible complete pipeline from Start to End,
716 this “average” pipeline is not in fact performed by any of the papers assessed by this review.

717 **Figure 3:** Violin plot of standardised task position within pipelines. Increasing x-axis
718 position denotes later placement of task within pipelines, vertical dashed lines denote 25%, 50%
719 and 75% of the way through the pipeline respectively. Tasks are separated into task groups and
720 ordered within task group by mean standardised pipeline position. Points denote task positions
721 where tasks occurred too infrequently to compute density profile for violin plots. Values report
722 the total number of papers implementing each task.

723 **Figure 4:** Plots summarising the reporting of three key aspects of bioinformatic tools
724 (software name, version and parameters) by the core papers. A). Venn diagram shows the number
725 of papers fully reporting each detail, i.e. giving the software used for every task reported, and
726 giving the parameters and version for each task where software is given; 86 papers reported at

727 least one of the three details for all steps, 25 further papers failed to fully report all three details
728 in all steps. B) Bar chart details the proportion of papers employing a specific task that failed to
729 report the software used for that task, with longer bars denoting a greater proportion of papers
730 not reporting software for that specific task

731 **Figure 5:** Consistency in software reporting and use over time. A) The total number of
732 unique software functions reported across all papers for each year of publication. B) For each
733 paper, the proportion of the total number of bioinformatic tasks for which the software used for a
734 task was not reported. C) The software homogeneity rate is one minus the number of different
735 software tools used for a given task in a given year, divided by the number of papers employing
736 the task in that year, calculated only when more than one paper reported a task in a given year. A
737 value of 1 means all papers used the same tool for a given task in a given year. D) The software
738 dominance rate is the proportion of papers that use the most common software tool for a given
739 task in a given year, calculated only when more than one paper reported a task in a given year. A
740 value of 1 means all papers used the same tool for a given task in a given year. B-D) Best fit
741 regression lines are shown in blue with shaded 95% confidence intervals around the line.
742 Horizontal jitter added to points to illustrate density within years; C & D) colours denote
743 different tasks, see Figure S1.

Task Group	Task	Description	Number papers reporting task	Number papers not reporting software	Total number of software tools	Total number of software functions	Number of papers performing manually
Read preparation	<i>quality control</i>	<i>Generating a report of sequence quality information from a sample or set of samples - no modification is done to data</i>	19	0	4	4	0
	<i>adapter trimming</i>	<i>Trimming of sequencing adapters</i>	9	1	6	6	0
	<i>demultiplexing</i>	<i>Separation of sequences from a mixed pool into separate pools based on the occurrence of a unique set of bases (index or tag)</i>	55	17	16	19	0
	<i>pair merging</i>	<i>The assembly of mate pair reads into a single contig</i>	63	1	10	18	0
	<i>quality trimming</i>	<i>The removal of bases from either or both ends of sequences in a pool based on quality scores</i>	20	1	8	10	0
	<i>mate pairing</i>	<i>The identification and synchronisation of mate pair reads between two samples, often involving arranging reads in identical orders and/or removal of reads without a mate pair</i>	3	0	3	3	0
	<i>primer trimming</i>	<i>Trimming of PCR primers</i>	66	8	15	17	0
	<i>reverse complementation</i>	<i>Reverse complementing the sequences in a pool</i>	7	3	2	2	0
	<i>sequence conversion</i>	<i>Converting sequences from fastq to fasta</i>	3	0	2	3	0
	<i>length trimming</i>	<i>The removal of bases from either or both ends of sequences in a pool, either the removal of a fixed number of bases or the removal of a variable number of bases to reduce sequences to a standard length</i>	10	3	6	7	0
	<i>pair concatenation</i>	<i>Concatenating mate pair reads into a single contig (where reads don't overlap)</i>	8	4	4	4	0
	<i>assembly</i>	<i>The assembly of reads into contigs, applied when more than one pair of overlapping fragments have been metabarcoded</i>	6	0	4	4	0
	<i>degapping</i>	<i>Removal of gaps from sequences</i>	1	0	1	1	0
Sequence processing	<i>dereplication</i>	<i>The removal of duplicate reads to retain only unique sequences in a pool; often the total number of copies of a sequence is recorded in the header of the retained sequence</i>	58	10	11	19	0
	<i>size sorting</i>	<i>The sorting of a fasta file according to a size annotation in the header</i>	10	2	3	4	0
Filtering	<i>quality filtering</i>	<i>Removal and/or trimming of sequences from a pool based on quality information. Also often converts from fastq to fasta.</i>	81	11	20	27	0
	<i>similarity filtering</i>	<i>Removal of sequences based on similarity to an alignment, either based on sequence identity or alignment position</i>	9	1	4	4	0

length filtering	<i>The removal of sequences from a pool that are less than, more than, or fall within or outside of a specified length threshold or thresholds</i>	54	21	17	23	0	
preclustering	<i>Reduction of sequence variation in a dataset prior to further processing - a form of denoising</i>	12	1	3	6	0	
denoising	<i>The removal of reads containing putative PCR or sequencing errors based on statistical assessment</i>	18	1	8	8	0	
normalisation	<i>A process by which the number of sequences for each of a set of samples is reduced where necessary such that the output set of samples all have the same number of sequences while maintaining the relative frequencies of OTUs</i>	2	0	1	1	1	
chimera filtering	<i>The filtering of putative chimeric assemblies from a pool of mate paired reads</i>	63	4	6	16	1	
translation filtering	<i>Removal of sequences from a set of sequence based on their translation, usually removing sequences with inframe stop codons or frameshifts due to erroneous indels or substitutions caused by sequencing errors</i>	22	3	11	12	0	
frequency filtering	<i>Removal of sequences based on their frequency in a pool</i>	51	37	11	15	1	
taxonomy filtering	<i>Removal of sequences based on an assigned taxonomy or a taxonomic classification</i>	9	5	1	1	1	
mistag filtering	<i>Removal of sequences based on putative tagging errors</i>	3	1	1	1	0	
Data generation	OTU delimitation	<i>The grouping of a set of sequences into OTUs by some method</i>	84	5	12	22	0
	OTU mapping	<i>The mapping of sequences to OTUs to provide read counts for each OTU</i>	30	3	7	11	0
	uncurated taxonomic assignment	<i>The assignment (identification or classification) of taxonomy to OTUs using a global uncurated reference database (e.g. GenBank, BOLD)</i>	55	2	11	13	0
	reference taxonomic assignment	<i>The assignment (identification or classification) of taxonomy to OTUs using a purpose-built and/or specially curated reference set of sequences</i>	60	9	18	23	1

Table 1: Table of all bioinformatic tasks performed across the core papers set. Tasks are grouped into four groups by broad purposes, and a detailed definition of each task is given along with summary statistics of the implementation of each task across the 111 papers. For a list of the software used for each task, Table S1 is an expanded version of this table.

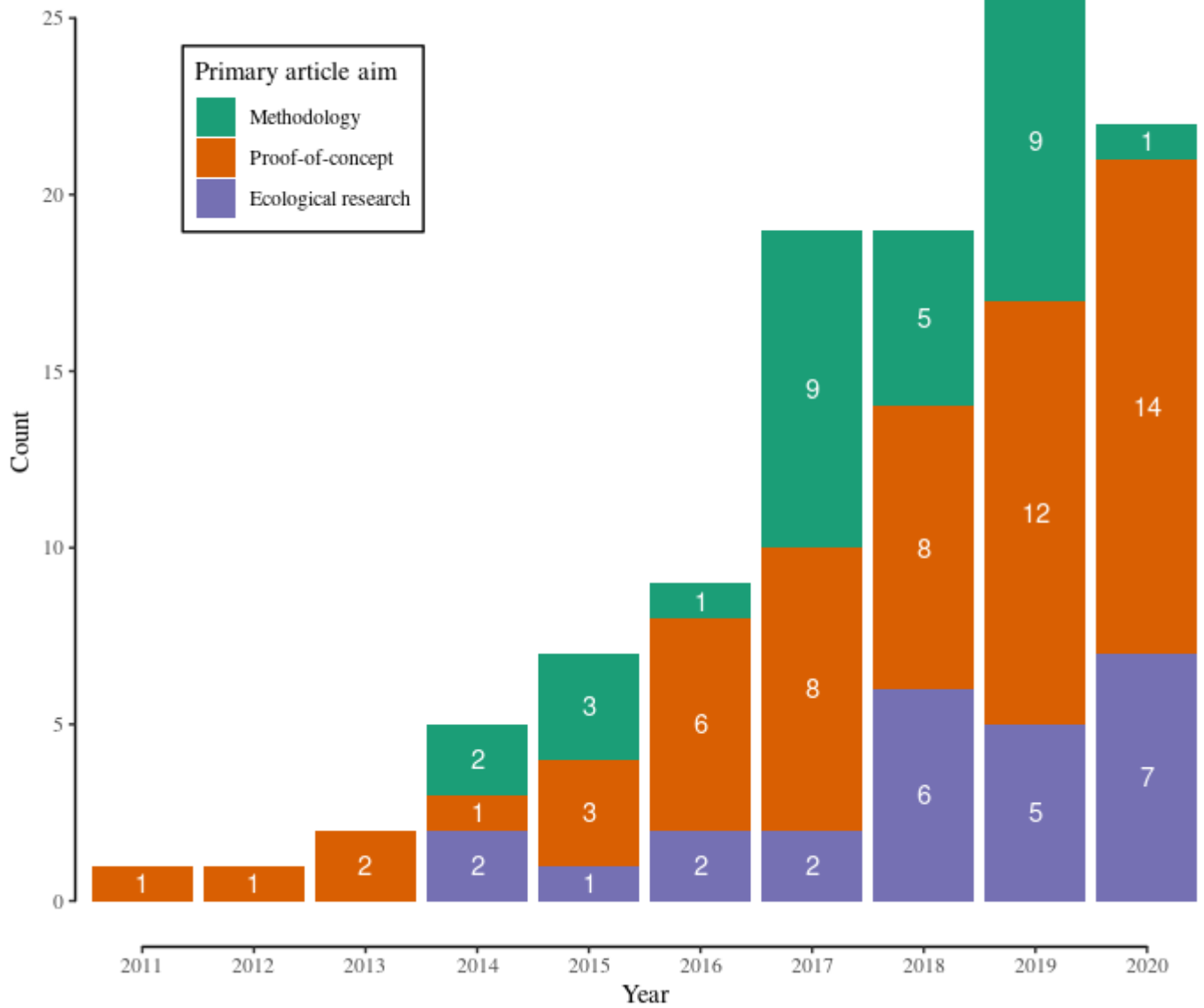


Figure 1: Year of publication of the articles in the core papers set. Bar fills and numbers refer to the number of articles within each research aim category. Note that only articles indexed by Web of Science by 3rd November 2020 were included.

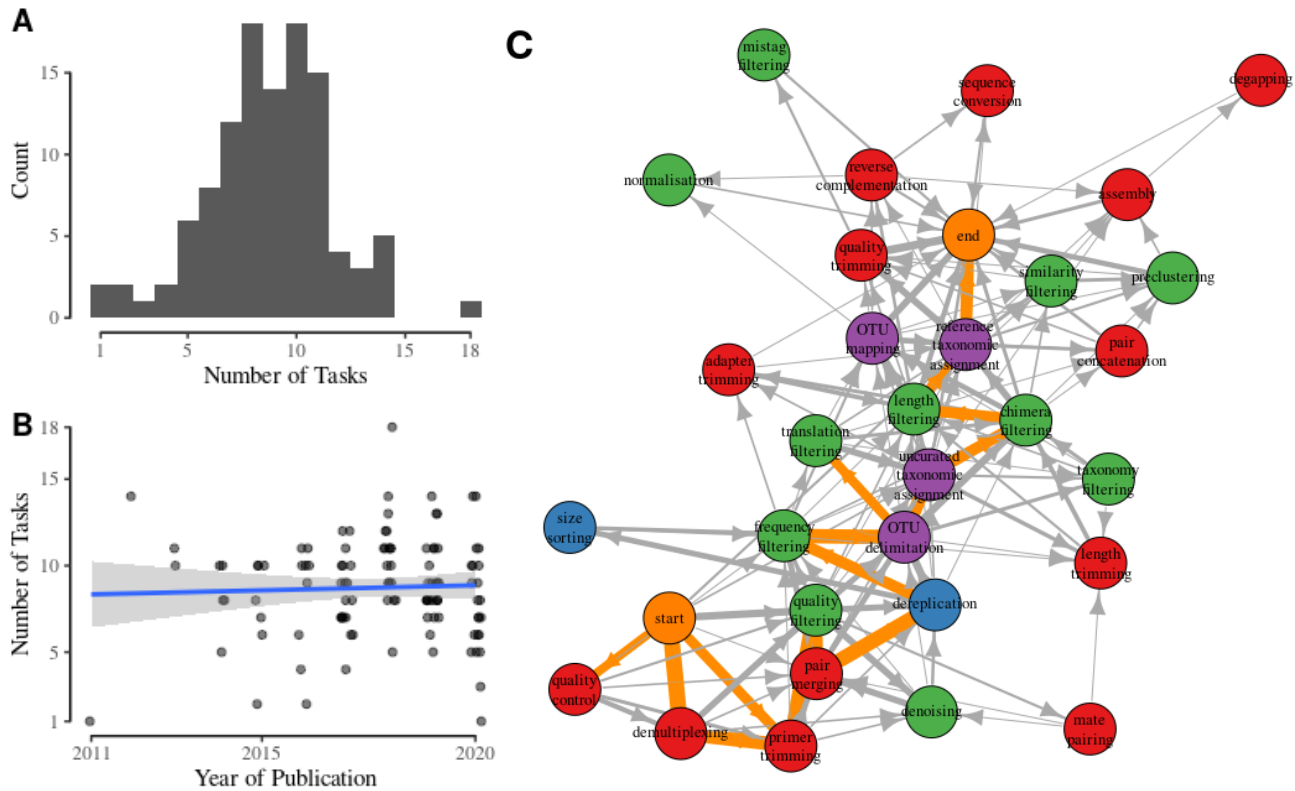


Figure 2: Bioinformatic pipelines implemented by the core papers set. Left: A) Frequency distribution of the number of tasks by study, B) Number of tasks by study against the year of publication, with best fit regression line in blue with shaded 95% confidence intervals around the line. Slight horizontal jitter added to points to better show density. Right: C) Network diagram of tasks and different pipeline routes through these tasks. All pipelines start and end on the respective orange nodes. All other nodes are coloured according to the four main categories of bioinformatic tasks; red for read preparation tasks, blue for sequence processing, green for filtering and purple for data generation tasks. Arrows link tasks performed consecutively, with direction of arrow showing order of tasks. Thickness of arrows shows relative frequency of pairs of consecutive tasks. Arrows coloured orange are the top 10% of consecutive task pairs by relative frequency. Note that while this illustrates a possible complete pipeline from Start to End, this “average” pipeline is not in fact performed by any of the papers assessed by this review.

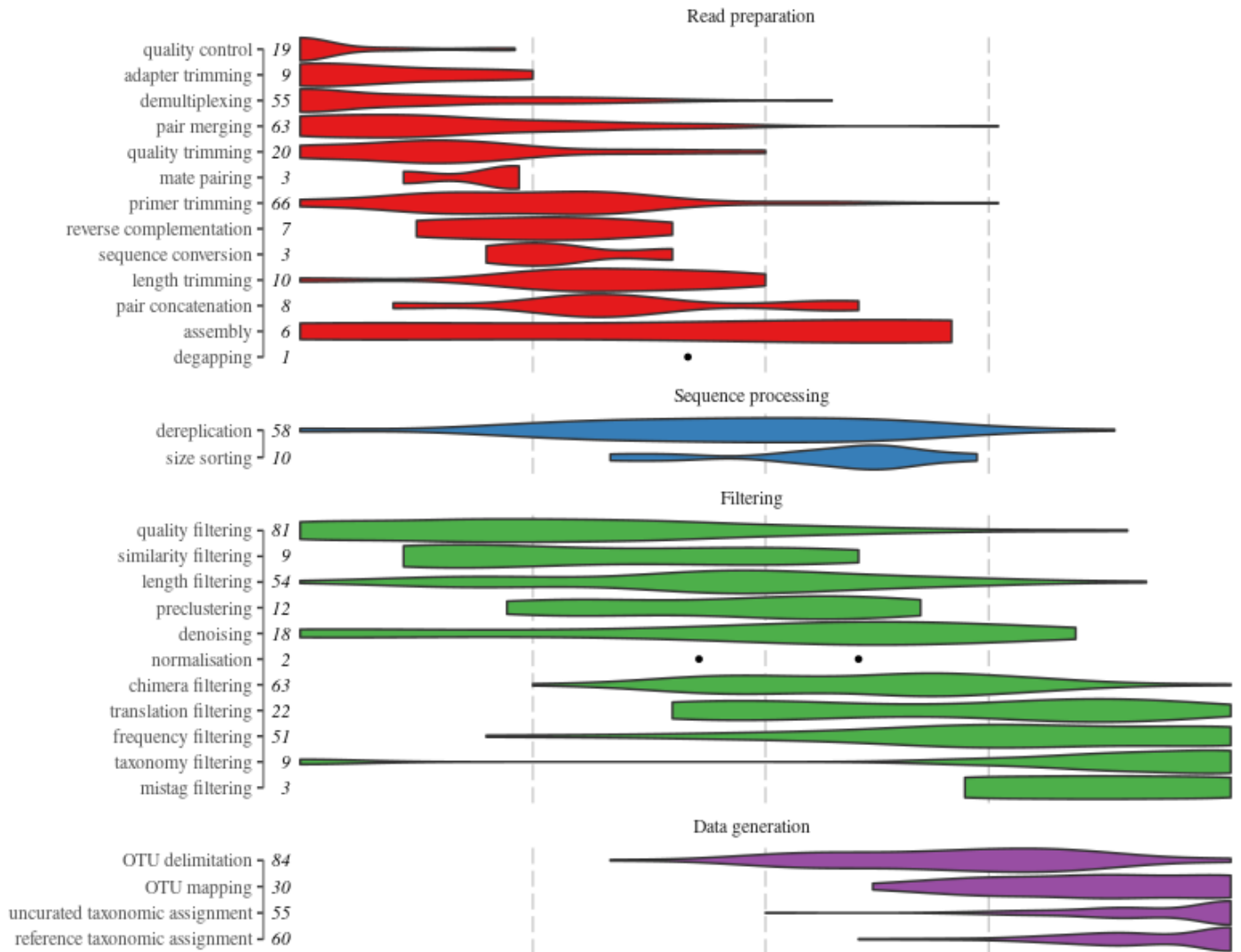


Figure 3: Violin plot of standardised task position within pipelines. Increasing x-axis position denotes later placement of task within pipelines, vertical dashed lines denote 25%, 50% and 75% of the way through the pipeline respectively. Tasks are separated into task groups and ordered within task group by mean standardised pipeline position. Points denote task positions where tasks occurred too infrequently to compute density profile for violin plots. Values report the total number of papers implementing each task.

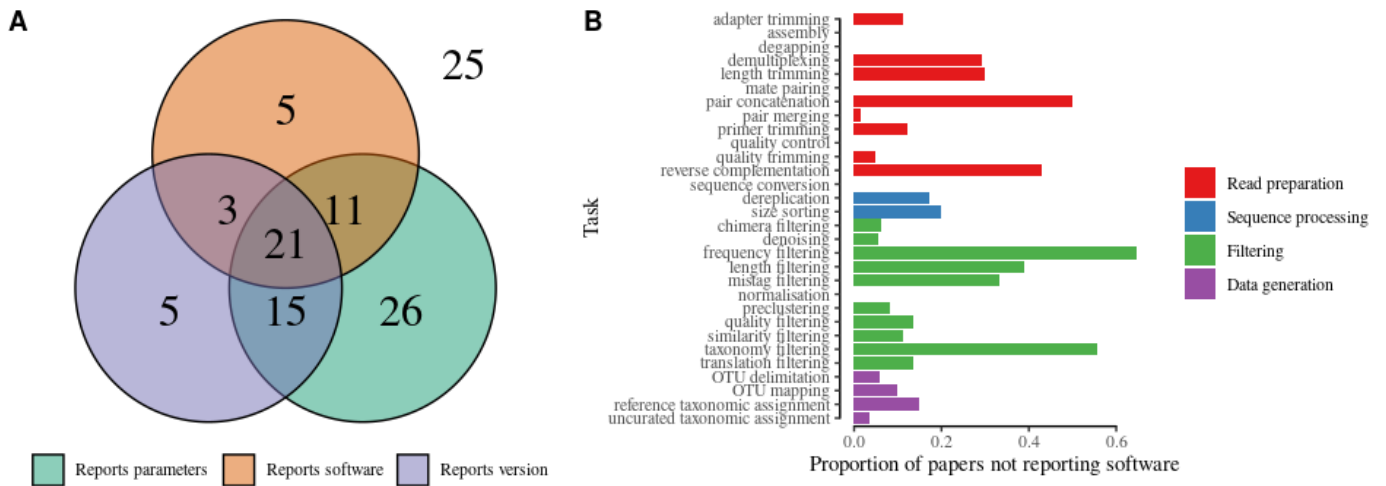


Figure 4: Plots summarising the reporting of 3 key methodological details by papers. A) Venn diagram shows the number of papers fully reporting each detail, i.e. giving the software used for every task reported, and giving the parameters and version for each task where software is given; 86 papers reported at least one of the three details for all steps, 25 further papers failed to fully report all three details in all steps. B) Bar chart details the proportion of papers employing a specific task that failed to report the software used for that task, with longer bars denoting a greater proportion of papers not reporting software for that specific task.

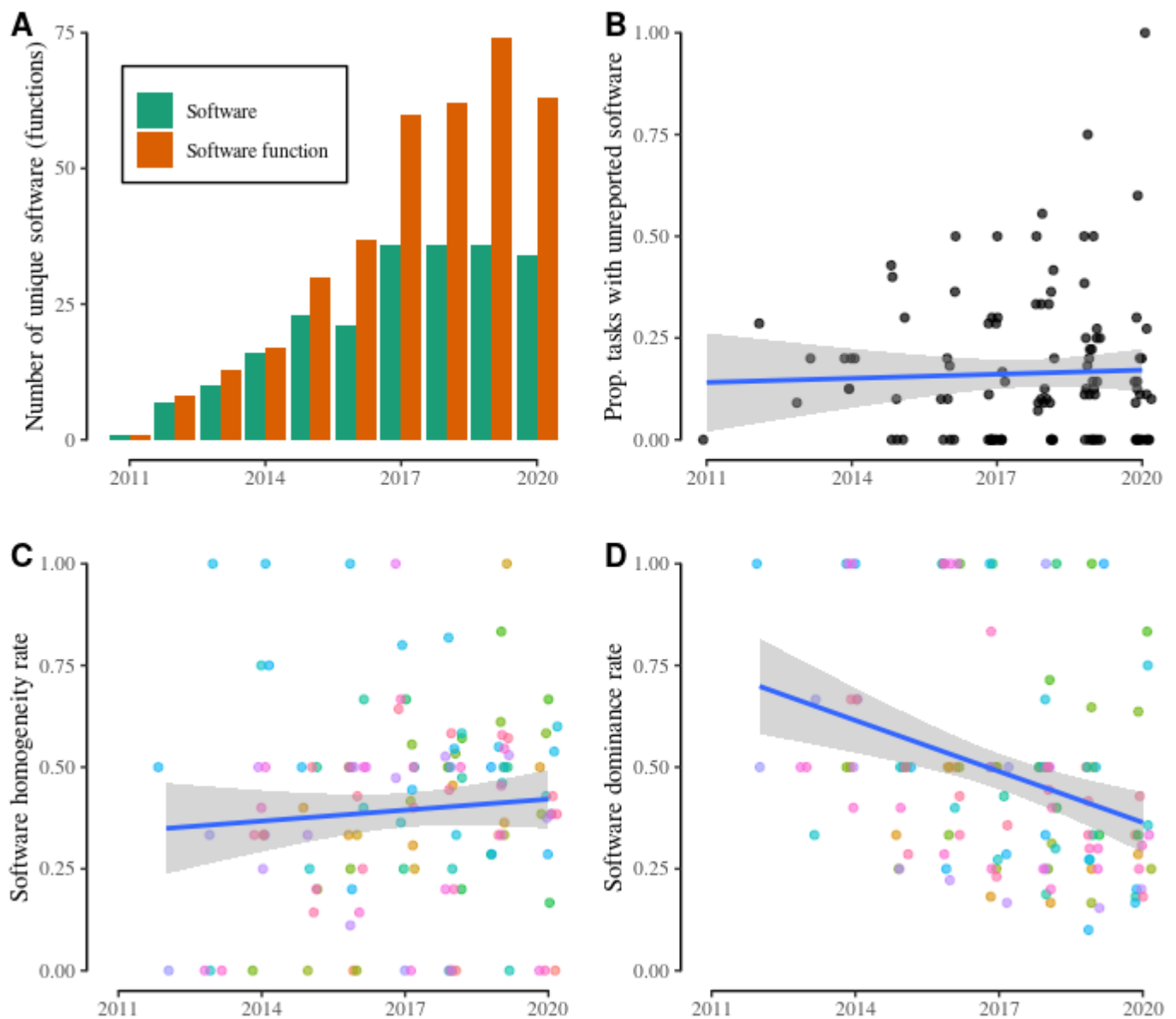


Figure 5: Consistency in software reporting and use over time. *A)* The total number of unique software functions reported across all papers for each year of publication. *B)* For each paper, the proportion of the total number of bioinformatic tasks for which the software used for a task was not reported. *C)* The software homogeneity rate is one minus the number of different software tools used for a given task in a given year, divided by the number of papers employing the task in that year, calculated only when more than one paper reported a task in a given year. A value of 1 means all papers used the same tool for a given task in a given year. *D)* The software dominance rate is the proportion of papers that use the most common software tool for a given task in a given year, calculated only when more than one paper reported a task in a given year. A value of 1 means all papers used the same tool for a given task in a given year. *B-D)* Best fit regression lines are shown in blue with shaded 95% confidence intervals around the line. Horizontal jitter added to points to illustrate density within years; *C & D)* colours denote different tasks, see Figure S1.