

# 1 An efficient and improved laboratory workflow and tetrapod database 2 for larger scale eDNA studies

3  
4 Jan Axtner<sup>1</sup>, Alex Crampton-Platt<sup>1</sup>, Lisa A. Hörig<sup>1</sup>, Charles C.Y. Xu<sup>2,3,4</sup>, Douglas W. Yu<sup>2,5</sup> and  
5 Andreas Wilting<sup>1</sup>

## 6 Affiliations:

7  
8 <sup>1</sup> Leibniz Institute for Zoo and Wildlife Research (*Leibniz-IZW*), Alfred-Kowalke-Str. 17,  
9 10315 Berlin, Germany

10 <sup>2</sup> State Key Laboratory of Genetic Resources and Evolution, Kunming Institute of Zoology,  
11 Chinese Academy of Sciences, 32 Jiaochang East Road, Kunming, Yunnan 650223, China

12 <sup>3</sup> Groningen Institute for Evolutionary Life Sciences, University of Groningen, P.O. Box  
13 11103, 9700 CC Groningen, The Netherlands

14 <sup>4</sup> Redpath Museum and Department of Biology, McGill University, Montreal, QC, Canada

15 <sup>5</sup> School of Biological Sciences, University of East Anglia, Norwich Research Park, Norwich,  
16 Norfolk NR47TJ, UK

- 17  
18     ▪ Abstract  
19     ▪ Introduction  
20     ▪ Methods  
21     ▪ Findings & Discussion  
22     ▪ Conclusion

## 23 24 Abstract

### 25 Background

26 The use of environmental DNA, ‘eDNA,’ for species detection via metabarcoding is growing  
27 rapidly and now, even terrestrial mammals can be monitored via ‘invertebrate-derived DNA’  
28 or ‘iDNA’ from hematophagous invertebrates. We present a co-designed lab workflow and  
29 bioinformatic pipeline to mitigate the two most important risks of e/iDNA: sample  
30 contamination and taxonomic mis-assignment. These risks arise from the need for  
31 amplification to detect the trace amounts of DNA and the necessity of using short target  
32 regions due to DNA degradation.

### 33 Findings

34 Here we present a high-throughput laboratory workflow that minimises these risks via a  
35 three-step strategy: (1) each sample is sequenced for two *PCR replicates* from each of two  
36 *extraction replicates*; (2) we use a ‘twin-tagging,’ two-step PCR protocol; (3) and a multi-

37 marker approach targeting three mitochondrial loci: *12S*, *16S* and *CytB*. As a test, 1532  
38 leeches were analysed from Sabah, Malaysian Borneo. Twin-tagging allowed us to detect  
39 and exclude chimeric sequences. The smallest DNA fragment (*16S*) amplified best for all  
40 samples but often at lower taxonomic resolution. We only accepted assignments that were  
41 found in both *extraction replicates*, totalling 174 assignments for 96 samples.

42 To avoid false taxonomic assignments, we also present an approach to create curated  
43 reference databases that can be used with the powerful taxonomic-assignment method  
44 *PROTAX*. For some taxonomic groups and some markers, curation resulted in over 50% of  
45 sequences being deleted from public reference databases, due mainly to: (1) limited overlap  
46 between our target amplicon and available reference sequences; (2) apparent mislabelling  
47 of reference sequences; (3) redundancy. A provided bioinformatics pipeline processes  
48 amplicons and conducts the *PROTAX* taxonomic assignment.

## 49 **Conclusions**

50 Our metabarcoding workflow should help research groups to increase the robustness of  
51 their results and therefore facilitate wider usage of e/iDNA, which is turning into a valuable  
52 source of ecological and conservation information on tetrapods.

53

## 54 **Introduction**

55 Monitoring, or even detecting, elusive or cryptic species in the wild can be challenging,  
56 particularly in dense vegetation or difficult terrain. In recent years there has been a rise in  
57 the availability of cost-effective DNA-based methods made possible by advances in high-  
58 throughput DNA sequencing (HTS). One such method is eDNA metabarcoding, which seeks  
59 to identify the species present in a habitat from traces of 'environmental DNA' (eDNA) in  
60 substrates such as water, soil, or faeces. A recent variation of eDNA metabarcoding, known  
61 as 'invertebrate-derived DNA' (iDNA) metabarcoding, targets the genetic material of prey or  
62 host species extracted from copro-, sarco- or haematophagous invertebrates. Examples  
63 include ticks [1], blow or carrion flies [2, 3, 4, 5], mosquitoes [6, 7, 8, 9] and leeches [10, 11,  
64 12,13]. Many of these parasites are ubiquitous, highly abundant, and easy to collect, making  
65 them an ideal source of biodiversity data, especially for terrestrial vertebrates that are  
66 otherwise difficult to detect [14, 15, 10]. In particular, the possibility for bulk collection and  
67 sequencing in order to screen large areas and minimise costs is attractive. However, most of  
68 the recent studies on iDNA studies focus on single-specimen DNA extracts and Sanger  
69 sequencing, and thus are not making use of the advances of HTS and a metabarcoding  
70 framework for carrying out larger scale biodiversity surveys.

71 That said, e/iDNA metabarcoding also poses several challenges, due to the low quality and  
72 low amounts of target DNA available, relative to non-target DNA (including the high-quality  
73 DNA of the live, invertebrate vector). In bulk iDNA samples comprised of many invertebrate  
74 specimens, this problem is further exacerbated by the variable time since each individual  
75 has fed, if at all, leading to differences in the relative amount and degradation of target DNA

76 per specimen. This makes e/iDNA studies similar to ancient DNA samples, which also pose  
77 the problem of low quality and low amounts of target DNA [16, 17]. The great disparity in  
78 the ratio of target to non-target DNA and the low overall amount of the former requires an  
79 enrichment step, which is achieved via the amplification of a short target sequence  
80 (amplicon) by polymerase chain reaction (PCR), to obtain enough target material for  
81 sequencing. However, this enrichment step can result in false-positive species detections,  
82 either through contamination or through volatile short PCR amplicons in the laboratory, and  
83 false negative results, through primer bias and low concentrations of template DNA.  
84 Although laboratory standards to prevent and control for such false results are well  
85 established in the field of ancient DNA, there are still no best-practice guidelines for e/iDNA  
86 studies, and thus few studies sufficiently account for such problems (but see [18]).

87 The problem is exacerbated by the use of ‘universal’ primers used for the PCR, which  
88 maximise the taxonomic diversity of the amplified sequences. This makes the method a  
89 powerful biodiversity assessment tool, even where little is known *a priori* about which  
90 species might be found. However, using such primers, in combination with low quality and  
91 quantity of target DNA, which often requires a high number of PCR cycles to generate  
92 enough amplicon products for sequencing, makes metabarcoding studies particularly  
93 vulnerable to false-results [13, 19; 20]. The high number of PCR cycles, combined with the  
94 high sequencing depth of HTS, also increase the likelihood that contaminants are amplified  
95 and detected, possibly to the same or greater extent as some true-positive trace DNA. As  
96 e/iDNA have been proposed as tools to detect very rare and priority conservation species  
97 such as the Saola, *Pseudoryx nghetinhensis* [10], false detection might result in misguided  
98 conservation activities worth several hundreds of thousands of US dollars e.g. [21].  
99 Therefore, similar to ancient DNA studies, great care must be taken to minimise the  
100 possibility for cross-contamination in the laboratory and to maximise the correct detection  
101 of species through proper experimental design. Replication in particular is an important tool  
102 for reducing the incidence of false negatives and detection of false positives but the trade-  
103 off is increased cost, workload, and analytical complexity [19].

104 A second source of false-positive species detections is the incorrect assignment of  
105 taxonomies to the millions of short HTS reads generated by metabarcoding. Although there  
106 has been a proliferation of tools focused on this step, most can be categorised into just  
107 three groups depending on whether the algorithm utilises sequence similarity searches,  
108 sequence composition models, or phylogenetic methods [22, 23, 24]. The one commonality  
109 among all methods is the need for a reliable reference database of correctly identified  
110 sequences, yet there are few curated databases currently appropriate for use in e/iDNA  
111 metabarcoding. Two exceptions are SILVA [25] for the nuclear markers SSU and LSU rRNA  
112 used in microbial ecology, and BOLD (Barcode of Life Database; citation) for the COI ‘DNA  
113 barcode’ region. For other loci, a non-curated database downloaded from the INSDC  
114 (International Nucleotide Sequence Database Collaboration, e.g. GenBank) is generally used.  
115 However, the INSDC places the burden for metadata accuracy, including taxonomy, on the

116 sequence submitters, with no restriction on sequence quality or veracity. For instance,  
117 specimen identification is often carried out by non-specialists, which increases error rates,  
118 and common laboratory contaminant species (e.g. human DNA sequences) are submitted in  
119 lieu of the sample itself. The rate of sequence mislabelling has not been assessed for  
120 GenBank, but for several curated microbial databases (Greengenes, LTP, RDP, SILVA),  
121 mislabelling rates have been estimated at between 0.2% and 2.5% [26]. It is likely that the  
122 true proportion of mislabelled samples in GenBank is higher than this given the lack of  
123 professional curation. Moreover, correctly identifying such errors is labour-intensive, so  
124 most metabarcoding studies simply base their taxonomic assignments on sequence-  
125 similarity searches of the whole INSDC database (e.g. with BLAST) [3, 10, 12] and thus can  
126 only detect errors if assignments are ecologically unlikely. Furthermore, reference  
127 sequences for the species that are likely to be sampled in iDNA studies are often  
128 underrepresented in or absent from these databases, which increases the possibility of  
129 incorrect assignment. For instance, fewer than 50% of species occurring in a tropical  
130 megadiverse rainforest are represented in Genbank (see findings below). When species-  
131 level matches are ambiguous, it might still be possible to assign a sequence to a higher  
132 taxonomic rank by using an appropriate algorithm such as MEGAN's Lowest Common  
133 Ancestor [27] or *PROTAX* [28].

134 We present here a complete laboratory workflow and complementary bioinformatics  
135 pipeline, starting from DNA extraction to taxonomic assignment of HTS reads using a  
136 curated reference database. The laboratory workflow allows for efficient screening of  
137 hundreds of e/iDNA samples: (1) two *extraction replicates* are separated during DNA  
138 extraction, and each is sequenced in two *PCR replicates* (Fig. 1); (2) a 'twin-tagged', two-step  
139 PCR protocol prevents cross-sample contamination as no unlabelled PCR products are  
140 produced (Fig. 2); (3) robustness of the taxonomic assignment is improved by using up to  
141 three mitochondrial markers. Our bioinformatics pipeline includes a standardized,  
142 automated, and replicable approach to create a curated database, which allows updating as  
143 new reference sequences become available, and to be expanded to other amplicons with  
144 minimal additional effort. We also provide scripts for processing the raw data to quality-  
145 controlled dereplicated reads and for taxonomic assignment of these reads using *PROTAX*  
146 [28], a probabilistic method that has been shown to be robust even when reference  
147 databases are incomplete [23, 4] (all scripts are available from URL  
148 <https://github.com/alexcrampton-platt/screenforbio-mbc>).

## 149 **Methods**

150 iDNA samples

151 We used 242 collections of haematophagous terrestrial leeches stored in *RNALater* (Sigma-  
152 Aldrich, Munich -Germany) from Deramakot Forest Reserve in Sabah, Malaysian Borneo as  
153 samples. Each sample consisted of one to 77 leech specimens (median 4). In total, 1532  
154 leeches were collected, exported under the permit (JKM/MBS.1000-2/3 JLD.2 (8) issued by  
155 the Sabah Biodiversity Council), and analysed at the laboratories of the Leibniz-IZW.

156 Laboratory workflow

157 The laboratory workflow is designed to both minimize the risk of sample cross-  
158 contamination and to aid identification of any instances that do occur. All laboratory steps  
159 (extraction, pre and post PCR steps, sequencing) took place in separate laboratories and no  
160 samples or materials were allowed to re-enter upstream laboratories at any point in the  
161 workflow. All sample handling was carried out under specific hoods that were wiped with  
162 bleach, sterilized, and UV irradiated for 30 minutes after each use. All labs are further UV  
163 irradiated for four hours each night.

164 *DNA extraction*

165 DNA was extracted from each sample in bulk. Leeches were cut into small pieces with a  
166 fresh scalpel blade and incubated in lysate buffer (proteinase K and ATL buffer at a ratio of  
167 1:10; 0.2 ml per leech) overnight at 55 °C (12 hours minimum) in an appropriately sized  
168 vessel for the number of leeches (2 or 5 ml reaction tube). For samples with more than 35  
169 leeches, the reaction volume was split in two and recombined after lysis.

170 Each lysate was split into two *extraction replicates* (A and B; maximum volume 600 µl) and  
171 all further steps were applied to these independently. We followed the DNeasy 96 Blood &  
172 Tissue protocol for animal tissues (Qiagen, Hilden -Germany) on 96 plates for clean-up. DNA  
173 was eluted twice with 100 µl TE buffer. DNA concentration was measured with PicoGreen  
174 dsDNA Assay Kit (Quant-iT, ThermoFisherScientific, Waltham -USA) in 384-well plate format  
175 using an appropriate plate reader (200 PRO NanoQuant, Tecan Trading AG, Männedorf -  
176 Switzerland). Finally, all samples were diluted to a maximum concentration of 10 ng/µl.

177 *Shot-gun sequencing to quantify mammalian DNA content*

178 To estimate the proportion of mammalian DNA in the leech samples, we ran a 75-cycle  
179 paired-end, shot-gun sequencing on an Illumina MiSeq on a subset of 58 samples. We used  
180 BLAST to compare the reads to GenBank and used MEGAN to find the lowest common  
181 ancestor for each read.

182 *PCR*

183 *Two-round PCR protocol.* – We amplified three mitochondrial markers – a short 93 bp  
184 fragment of 16S rRNA (16S), a 389 bp fragment of 12S rRNA (12S), and a 302 bp fragment of  
185 cytochrome b (*CytB*). For each marker, we ran a two-round PCR protocol (Figs. 1, 2). The  
186 first round amplified the target gene. The second round added the Illumina adapters for  
187 sequencing.

188 *Primer design.* – We used ‘twin-tagged’ PCR primers, meaning that *both* the forward and  
189 reverse primers were given the *same* sample-identifying sequence (i.e. ‘tags’) added as  
190 primer extensions (Fig. 2). This ensured that unlabelled PCR products were never produced  
191 and allowed us later to detect and delete tag jumping events [29] (Fig. 2). Primer sequences  
192 are in Table 1 [30, 31].

193 In the first PCR round, we used 25 different 5-bp *sample-identifying tags* (*tag 1*), with a  
194 minimum pairwise distance of three (Faircloth et al, 2012; Supplement Table 1). These

195 primers also contained different forward and reverse sequences (*Read 1 & Read 2 sequence*  
196 *primers*) (Supplement table 1) to act priming sites for the second PCR round (Fig. 2).

197 In the second PCR round, we used 20 different 5-bp *plate*-identifying tags (*tag 2*), with a  
198 minimum pairwise distance of three [32]. These primers also contained the Illumina P5 and  
199 P7 adapter sequences (Fig. 2). The product of the second PCR round could thus be cleaned  
200 up, quantified, pooled, and sequenced without needing to carry out a separate library  
201 preparation step (e.g. Nextera, TruSeq).

202 *Cycle number considerations.* – Because we know that our target DNA is at low  
203 concentration in the samples, we are faced with a trade-off between (1) using fewer PCR  
204 cycles (e.g. 30 cycles) to minimise amplification bias (caused by some target DNA binding  
205 better to the primer sequences and thus outcompeting during PCR other target sequences  
206 that bind less well, [33]) and (2) using more PCR cycles (e.g. 40 cycles) to ensure that low-  
207 concentration target DNA is sufficiently amplified in the first place. Rather than choose  
208 between these two extremes, we ran both low- and a high-cycle protocols and sequenced  
209 both sets of amplicons.

210 Thus, each of the two *extraction replicates* A and B was split and amplified using different  
211 cycle numbers (*PCR replicates* 1 and 2) for a total of four (= 2 *extraction replicates* X 2 *PCR*  
212 *replicates* -> A1/A2 and B1/B2 ) replicates per sample per marker (Fig. 1). For *PCR replicates*  
213 A1/B1, we used 30 cycles in the first PCR round to minimize the effect of amplification bias.  
214 For *PCR replicates* A2/B2, we used 40 cycles in the first PCR round to increase the likelihood  
215 of detecting species with very low input DNA (Fig. 1).

216 *PCR protocol.* – The first-round PCR reaction volume was 20 µl, including 0.1 µM primer mix,  
217 0.2 mM dNTPs, 1.5 mM MgCl<sub>2</sub>, 1x PCR buffer, 0.5 U AmpliTaq Gold™ (Invitrogen, Karlsruhe -  
218 Germany), and 2 µl of template DNA. Initial denaturation was 5 minutes at 95°C, followed  
219 by repeated cycles of 30 seconds at 95°C, 30 seconds at 54°C, and 45 seconds at 72°C. Final  
220 elongation was 5 minutes at 72°C. Samples were amplified in batches of 24 plus a negative  
221 (water) and a positive control (bank vole, *Myodes glareolus* DNA). All three markers were  
222 amplified simultaneously for each batch of samples in a single PCR plate. Non-target by-  
223 products were removed as required from some 12S PCRs by purification with magnetic  
224 Agencourt AMPure beads (Beckman Coulter, Krefeld -Germany).

225 In the second-round PCR, we used the same PCR protocol as above with 2 µl of the product  
226 of the first-round PCR and 10 PCR cycles.

227 *Quality control and sequencing*

228 Amplification was visually verified after the second-round PCR by gel electrophoresis on  
229 1.5% agarose gels. Controls were additionally checked with a TapeStation 2200 (D1000  
230 ScreenTape assay, Agilent, Waldbronn -Germany). All samples were purified with AMPure  
231 beads, using a beads-to-template ratio of 0.7:1 for 12S and *CytB* products, and a ratio of 1:1  
232 for 16S products. DNA concentration was measured with PicoGreen dsDNA as described  
233 above. Sequencing libraries were made for each PCR plate by equimolar pooling of all



234 positive samples; final concentrations were between 2 and 4 nmol. Generally, *12S* and *CytB*  
235 products were combined in a single library, whereas *16S* products were always separate,  
236 because of the difference in amplicon length. Up to 11 libraries were sequenced on each run  
237 of Illumina MiSeq following standard protocols. Libraries were sequenced with MiSeq  
238 Reagent Kit V3 (600 cycles, 300 bp paired-end reads) and had a final concentration of 11 pM  
239 spiked with 20 to 30% of PhiX control.

240 Establishment of the tetrapod reference database

241 *Reference database*

242 A custom bash script was written to generate a tetrapod reference database for each of the  
243 three markers, and additionally for a 250 bp mitochondrial cytochrome *c* oxidase subunit I  
244 amplicon (*COI*), which has previously been used in iDNA studies [2]. An important time-  
245 saving step was the use of the FASTA-formatted MIDORI mitochondrial databases [34]. The  
246 script updated the FASTA files for a subset of target species, removed errors and  
247 redundancy, and output FASTA files with species names and GenBank accessions in the  
248 headers. The script accepts four data inputs, two of which are optional. The required inputs  
249 are: (i) the MIDORI sequences (December 2015 'UNIQUE', downloaded from  
250 <http://www.reference-midori.info/download.php#>) for the relevant genes and (ii) an initial  
251 reference taxonomy. This taxonomy is needed to find or generate a full taxonomic  
252 classification for each sequence. Here we used the Integrated Taxonomic Information  
253 System (ITIS) classification for Tetrapoda, obtained with the R package *taxize* version 0.9.0  
254 [35], functions *downstream* and *classification*). The optional inputs are: (iii) supplementary  
255 FASTA files of reference sequences that should be added to the database, and (iv) a list of  
256 target species to be queried on GenBank to capture any sequences published since the  
257 MIDORI set was generated. For this study, 72 recently published [36] and 7 unpublished  
258 partial mitochondrial mammal genomes (ACCESSION No XXX) were added as input (iii). A list  
259 of 103 mammal species known to be present in the sampling area was added as input (iv).

260 With the above inputs, the seven curation steps are: 1) remove sequences not identified to  
261 species; 2) add any extra sequences from optional inputs (iii) and (iv) above; 3) select the  
262 target amplicon; 4) remove sequences with ambiguities; 5) compare species labels to the  
263 reference taxonomy from input (ii) and create a consensus taxonomy including any species  
264 known only from sequence data if genus already exists in reference; 6) identify and remove  
265 putatively mislabelled sequences; 7) discard redundant sequences, retaining one  
266 representative per haplotype per species.

267 The script is split into four modules, allowing optional manual curation at three key steps.  
268 The steps covered by each of the four modules are summarized in Table 2. The main  
269 programs used are highlighted and cited in the text where relevant, but many intermediate  
270 steps used common UNIX tools and unpublished lightweight utilities freely available from  
271 GitHub (Table 3).

272 *Module 1* - The first step is to select the tetrapod sequences from the MIDORI database for  
273 each of the four selected loci (input (i) above). This, and the subsequent step to discard  
274 sequences without strict binomial species names and reduce subspecies identifications to  
275 species-level, are made possible by the inclusion of the full NCBI taxonomic classification of  
276 each sequence in the FASTA header by the MIDORI pipeline. The headers of the retained  
277 sequences are then reformatted to include just the species name and GenBank accession  
278 separated by underscores. If desired, additional sequences from local FASTA files are then  
279 added to the MIDORI set (input (iii)). The headers of these FASTA files are required to be in  
280 the same format. Next, optional queries are made to the NCBI GenBank and RefSeq  
281 databases for each species in a provided list (input (iv)) for each of the four target loci, using  
282 NCBI's Entrez Direct [37]. Matching sequences are downloaded in FASTA format, sequences  
283 prefixed as "UNVERIFIED" are discarded, the headers are simplified as previously, and those  
284 sequences not already in the MIDORI set are added. The region of each sequence matching  
285 to the relevant target marker was extracted with a two-step process in which *usearch* (-  
286 *search\_pcr*) was used to select sequences where both primers were present, and these  
287 were in turn used as a reference to select partially matching sequences with *blastn* [38, 39].  
288 Sequences with a hit length of at least 90% of the expected marker length were retained by  
289 extracting the relevant subsequence based on the BLAST hit co-ordinates. Sequences with  
290 ambiguous bases were discarded at this stage. In the final step in module 1 a multiple-  
291 sequence alignment was generated with MAFFT [40, 41] for each partially curated amplicon  
292 dataset. The script then breaks to allow the user to check for any obviously problematic  
293 sequences that should be discarded before continuing.

294 *Module 2* - The species labels of the edited alignments are compared with the reference  
295 taxonomy (input (ii)). Any species not found is queried against the Catalogue of Life  
296 database (CoL) via *taxize* in case these are known synonyms, and their correct species label  
297 and classification is added to the reference taxonomy. The original species label is retained  
298 as a key to facilitate sequence renaming, and a note is added to indicate its status as a  
299 synonym. Finally, the genus name of any species not found in the CoL is searched against  
300 the consensus taxonomy, and if found, the novel species is added by taking the higher  
301 classification levels from of the other species in the genus. Orphan species labels are printed  
302 to a text file, and the script breaks to allow the user to check this list and manually create  
303 classifications for some or all if appropriate.

304 *Module 3* - This module begins by checking for any manually generated classification files  
305 (from the end of Module 2) and merging them with the reference taxonomy from Module 2.  
306 Any remaining sequences with unverifiable classifications are removed at this step. The next  
307 steps convert the sequences and taxonomy file to the correct formats for SATIVA [26].  
308 Sequence headers in the edited MAFFT alignments are reformatted to include only the  
309 GenBank accession and a taxonomy key file is generated with the correct classification listed  
310 for each accession number. In cases where the original species label was found to be a  
311 synonym, the corrected label is used. Putatively mislabelled sequences in each amplicon are



312 then detected with SATIVA, and the script breaks to allow inspection of the results. The user  
313 may choose to make appropriate edits to the taxonomy key file or list of putative mislabels  
314 at this point.

315 *Module 4* - Any sequences that are still flagged as mislabelled at the start of the fourth  
316 module are deleted from the SATIVA input alignments, and all remaining sequences are  
317 relabelled with the correct species name and accession. A final consensus taxonomy file is  
318 generated in the format required by *PROTAX*. Alignments are subsequently unaligned prior  
319 to species-by-species selection of a single representative per unique haplotype. Sequences  
320 that are the only representative of a species are automatically added to the final database.  
321 Otherwise, all sequences for each species are extracted in turn, aligned with MAFFT, and  
322 collapsed to unique haplotypes with *collapsetypes\_4.6.pl* (zero differences allowed; [42]).  
323 Representative sequences are then unaligned and added to the final database.

324 Bioinformatics workflow

325 *Read processing*

326 Although the curation of the reference databases is our main focus, it is just one part of the  
327 bioinformatics workflow for e/iDNA metabarcoding. A custom bash script was used to  
328 process raw basecall files to demultiplexed, cleaned, and dereplicated reads in FASTQ  
329 format on a run-by-run basis. All runs and amplicons were processed with the same settings  
330 unless otherwise indicated. *bcl2fastq* (Illumina) was used to convert basecall files to  
331 demultiplexed, paired-end FASTQ files for each library, allowing up to 1 mismatch in each  
332 *tag 2*. Each library was further demultiplexed into samples via unique *tag 1* pairs with  
333 *AdapterRemoval* (Schubert, Lindgreen and Orlando 2016), again allowing up to 1 mismatch  
334 in each tag. These steps allowed reads to be assigned to the correct samples via their four  
335 tags e.g. ABBA, ADDA, BDDB.

336 In all cases, amplicons were short enough to expect paired reads to overlap. Pairs were  
337 merged with *usearch (-fastq\_mergepairs; [43; 44])*, and only successfully merged pairs were  
338 retained. Primer sequences were trimmed with *cutadapt* [45], and only successfully  
339 trimmed reads at least 90% of expected amplicon length were passed to a quality filtering  
340 step with *usearch (-fastq\_filter)*. Lastly, reads were dereplicated with *usearch (-*  
341 *derep\_fulllength)* to retain only unique sequences, and singletons were discarded. The  
342 number of replicates that each unique sequence represented was also added to the read  
343 header at this step (option *-sizeout*).

344 *Taxonomic assignment*

345 The curated reference sequences and associated taxonomy were used for taxonomic  
346 classification of dereplicated reads using *PROTAX*, a recently published probabilistic method  
347 [28, 24]. *PROTAX* gives unbiased estimates of placement probability for each read at each  
348 taxonomic rank, allowing some assignments to be made to a higher rank even when there is  
349 a high degree of uncertainty at the species level. In other words, and unlike other taxonomic  
350 assignment methods, *PROTAX* can estimate the probability that a sequence belongs to a

351 taxon that is not included in the reference database. This was considered an important  
352 feature due to the expected incompleteness of the reference databases for tetrapods in the  
353 sampled location. As other studies have compared *PROTAX* with more established methods,  
354 e.g. MEGAN [27] (see [28, 4]), it was beyond the scope of this study to evaluate the  
355 performance of *PROTAX*.

356 Classification with *PROTAX* is a two-step process. Firstly, *PROTAX* selected a subset of the  
357 reference database that was used as training data to parameterise a *PROTAX* model for  
358 each marker, and secondly, the fitted models were used to assign four taxonomic ranks  
359 (species, genus, family, order) to each of the dereplicated reads, along with a probability  
360 estimate at each level. We also included the best similarity score of the assigned species or  
361 genus, mined from the LAST results (see below) for each read. This was helpful for flagging  
362 problematic assignments for downstream manual inspection, i.e. high probability  
363 assignments based on low similarity scores (implying that there are no better matches  
364 available) and low probability assignments based on high similarity scores (indicates  
365 conflicting database signal from several species with highly similar sequences).

366 Fitting the *PROTAX* model followed Somervuo et al. [24] except that 5000 training  
367 sequences were randomly selected for each target marker due to the large size of the  
368 reference database. In each case, 4500 training sequences represented a mix of known  
369 species with reference sequences (conspecific sequences retained in the database) and  
370 known species without reference sequences (conspecific sequences omitted, simulating  
371 species missing from the database), and 500 sequences represented previously unknown  
372 lineages distributed evenly across the four taxonomic levels (i.e. mimicked a mix of  
373 completely novel species, genera, families and orders). Pairwise sequence similarities of  
374 queries and references were calculated with LAST [46] following the approach of Somervuo  
375 et al. [24]. The models were weighted towards the Bornean mammals expected in the  
376 sampled area by assigning a prior probability of 90% to these 103 species and a 10%  
377 probability to all others ([24]; Supplement table 2). In cases of missing interspecific variation  
378 this helped to avoid unlikely assignments, especially in case of the very short 93 bp fragment  
379 of *16S*. Maximum *a posteriori* (MAP) parameter estimates were obtained following the  
380 approach of Somervuo et al. [28], but the models were parameterised for each of the four  
381 taxonomic levels independently, with a total of five parameters at each level (four  
382 regression coefficients and the probability of mislabelling).

383 Dereplicated reads for each sample were then classified using a custom bash script on a run-  
384 by-run basis. For each sample, reads in FASTQ format were converted to FASTA, and  
385 pairwise similarities were calculated against the full reference sequence database for the  
386 applicable marker with LAST. Assignments of each read to a taxonomic node based on these  
387 sequence similarities were made using a Perl script and the trained model for that level. The  
388 taxonomy of each node assignment was added with a second Perl script for a final table  
389 including the node assignment, probability, taxonomic level, and taxonomic path for each  
390 read. Read count information was included directly in the classification output via the size

391 annotation added to the read headers during dereplication. All Perl scripts to convert input  
392 files into the formats expected by *PROTAX*, *R* code for training the model following  
393 Somervuo et al. [24], and Perl scripts for taxonomic assignment were provided by P.  
394 Somervuo (personal communication).

#### 395 *Acceptance criteria*

396 In total we had twelve PCR reactions per sample: two *extraction replicates A* and *B* X two  
397 *PCR replicates 1* and *2* per extraction replication X the three markers (Fig. 1). We only  
398 accepted taxonomic assignments that were positively detected in both *extraction replicates*  
399 (*A* & *B*, Figure 3). The reason for conservatively omitting assignments that appeared in only  
400 one extraction replicate was to rule out sample cross-contamination during DNA extraction.  
401 In addition, we only accepted assignments with ten or more reads per marker, if only one  
402 marker was sequenced. If a species was assigned in more than one marker (e.g. *12S* and  
403 *16S*), we accepted the assignment even if in one sequencing run the number of reads was  
404 below ten.

405 Due to the imperfect PCR amplification of markers (the small *16S* fragment amplified better  
406 than the longer *CytB* fragment) and missing reference sequences in the database or shared  
407 sequence motifs between species, reads sometimes were assigned to species level for one  
408 marker but only to genus level for another marker. Thus, the final identification of species  
409 could not be automated and manual inspection and curation was needed. For each  
410 assignment, three parameters were taken into consideration: number of sequencing reads,  
411 the mean probability estimate derived from *PROTAX*, and the mean sequence similarity to  
412 the reference sequences based on LAST.

## 413 Findings & Discussion

### 414 *Database curation*

415 The MIDORI UNIQUE database (December 2015 version) contains 1,019,391 sequences  
416 across the four mitochondrial loci of interest (*12S*: 66,937; *16S*: 146,164; *CytB*: 223,247; *COI*:  
417 583,043), covering all Metazoa. Of these, 258,225 (25.3%) derive from the four tetrapod  
418 classes (Amphibia: 55,254; Aves: 51,096; Mammalia: 101,106; Reptilia: 50,769). The  
419 distribution of these sequences between classes and loci, and the losses at each curation  
420 step are shown in Figure 4. In three of the four classes, there is a clear bias towards *CytB*  
421 sequences, with over 50% of sequences derived from this locus. In both Aves and  
422 Mammalia, the *16S* and *12S* loci are severely underrepresented at less than 10% each, while  
423 for Reptilia, *COI* is the least sequenced locus in the database.

424 The numbers of sequences and rates of loss due to our curation steps varied among  
425 taxonomic classes and the four loci, although losses were observed between steps in almost  
426 all instances. The most significant losses followed amplicon selection and removal of non-  
427 unique sequences. Amplicon selection led to especially high losses in Amphibia and *16S*,  
428 indicating that data published on GenBank for this class and marker do not generally overlap  
429 with the primer sets used here. Meanwhile, the high level of redundancy in public databases

430 was highlighted by the significant reduction in the number of sequences during the final  
431 step of removing redundant sequences – in all cases over 10% of sequences were discarded,  
432 but some losses exceeded 50% (Mammalia: *COI*, *CytB*, *16S*; Amphibia: *16S*).

433 Data loss due to apparent mislabelling ranged between 1.9% and 7.4% and was thus  
434 generally higher than similar estimates for curated microbial databases [26]. SATIVA flags  
435 potential mislabels and suggests an alternative label supported by the phylogenetic  
436 placement of the sequences, allowing the user to make an appropriate decision on a case by  
437 case basis. The pipeline pauses after this step to allow such manual inspection to take place.  
438 However, for the current database, the number of sequences flagged was large (4378 in  
439 total), and the required taxonomic expertise was lacking, so all flagged sequences from non-  
440 target species were discarded to be conservative. The majority of mislabels were identified  
441 at species level (3053), but there were also significant numbers at genus (788), family (364)  
442 and order (102) level. Two to three sequences from Bornean mammal species were  
443 unflagged in each amplicon to retain the sequences in the database. This was important as  
444 in each case these were the only reference sequences available for the species. Additionally,  
445 *Muntiacus vaginalis* sequences that were automatically synonymised to *M. muntjak* based  
446 on the available information in the Catalogue of Life were revised back to their original  
447 identifications to reflect current taxonomic knowledge.

#### 448 *Database composition*

449 The final database was skewed even more strongly towards *CytB* than was the raw  
450 database. It was the most abundant locus for each class and representing over 60% of  
451 sequences for both Mammalia and Reptilia. In all classes, *16S* made up less than 10% of the  
452 final database, with Reptilia *COI* also at less than 10%.

453 Figure 5 (frequency distributions) shows that most species represented in the curated  
454 database for any locus have just one unique haplotype against which HTS reads can be  
455 compared, while a few species have many haplotypes. The prevalence of species with 20 or  
456 more haplotypes is particularly notable in *CytB* where the four classes have between 25  
457 (Aves) and 265 (Mammalia) species in this category. Figure 5 (coloured circles in each plot)  
458 also shows, that the species in the taxonomy are incompletely sampled across all loci, but  
459 coverage varies significantly between categories. In spite of global initiatives to generate  
460 *COI* sequences [47], this marker does not offer the best species-level coverage in any class  
461 and is a poor choice for Amphibia and Reptilia (<15% of species included). Even the best  
462 performing marker, *CytB*, is not a universally appropriate choice, as Amphibia is better  
463 covered by *12S*. These differences in underlying database composition will impact the  
464 likelihood of obtaining accurate taxonomic assignment for any one species from any single  
465 marker. Further barcoding campaigns are clearly needed to fill gaps in all markers and all  
466 classes to increase the power of future e/iDNA studies. As the costs of HTS decrease, we  
467 expect that such gap-filling will increasingly shift towards whole mitochondrial genomes  
468 [36], reducing the effect of marker choice on detection likelihood. In the meantime,  
469 however, the total number of species covered by the database can be increased by

470 combining multiple loci (here, up to four) and thus the impacts of database gaps on  
471 correctly detecting species can be minimized ([48]; Fig. 6).

472 In the present study, the primary target for iDNA sampling was the mammal fauna of  
473 Malaysian Borneo, and the 103 species expected in the sampling area represent an  
474 informative case study highlighting the deficiencies in existing databases (Fig. 6). Nine  
475 species are completely unrepresented while only slightly over half (554 species) have at  
476 least one sequence for all of the loci. Individually, each marker covers over half of the target  
477 species, but none achieves more than 85% coverage (*12S*: 75 species; *16S*: 68; *CytB*: 88; *COI*:  
478 66). Equally striking is the lack of within-species diversity, as most of the incorporated  
479 species are represented by only a single haplotype per locus. Some of the species have large  
480 distribution ranges, so it is likely that in some cases the populations on Borneo differ  
481 genetically from the available reference sequences, possibly limiting assignment success.  
482 Only a few expected species have been sequenced extensively, and most are of economic  
483 importance to humans (e.g. *Bos taurus*, *Bubalus bubalis*, *Macaca* spp, *Paradoxurus*  
484 *hermaphroditus*, *Rattus* spp, *Sus scrofa*), with as many as 100 haplotypes available (*Canis*  
485 *lupus*). Other well-represented species ( $\geq 20$  haplotypes) present in the sampling area  
486 include several Muridae (*Chiropodomys gliroides*, *Leopoldamys sabanus*, *Maxomys surifer*,  
487 *Maxomys whiteheadi*) and leopard cat (*Prionailurus bengalensis*).

#### 488 *Laboratory workflow*

489 Shotgun sequencing of a subset of our samples revealed that the median mammalian DNA  
490 content was only 0.9%, ranging from 0% to 98%. These estimates are approximate, but with  
491 more than 75% of the samples being below 5%, this shows clearly the scarcity of target DNA  
492 in bulk iDNA samples. The generally low DNA content and the fact that the target DNA is  
493 often degraded make enrichment of the target barcoding loci necessary. We used PCR with  
494 high cycle numbers to obtain enough DNA for sequencing. However, this second step  
495 increases the risk of PCR error: artificial sequence variation, non-target amplification, and/or  
496 raising contaminations up to a detectable level.

497 We addressed these problems by running two *extraction replicates*, two *PCR replicates*, and  
498 a multi-marker approach. The need for *PCR replicates* has been acknowledged and  
499 addressed extensively in ancient DNA studies [16] and has also been highlighted for  
500 metabarcoding studies [18, 19, 20, 49]. Despite this, many e/iDNA studies do not carry out  
501 multiple *PCR replicates* to detect and omit potential false sequences. In addition, *extraction*  
502 *replicates* are seldom applied, despite the evidence that cross-sample DNA contamination  
503 can occur during DNA extraction [50, 51, 52]. Here we only accepted sequences that  
504 appeared in a minimum of two independent PCRs, one from each *extraction replicate A* and  
505 *B* (Fig. 1).

506 We also used three different loci to correct for potential PCR-amplification biases. We were,  
507 however, unable to quantify this bias in this study due to the high degradation of the target  
508 mammalian DNA, which resulted in much higher overall amplification rates for *16S*, the



509 shortest of our PCR amplicons. For *16S*, 85% of the samples amplified, whereas for *CytB* and  
510 *12S*, only 57% and 44% amplified, respectively. Despite the greater taxonomic resolution of  
511 the longer *12S* and *CytB* fragments, our poorer amplification results for these longer  
512 fragments emphasize that e/iDNA studies should generally focus on short PCR fragments to  
513 increase the likelihood of positive amplifications of the degraded target DNA. In the case of  
514 mammal-focussed e/iDNA studies, a shorter (100 bp) *CytB* fragment will likely be very  
515 useful.

516 Our second major precaution was the use of twin-tagging for both PCRs (Fig. 2). This ensures  
517 that unlabelled PCR products are never produced and allows us to multiplex a large number  
518 of samples on a single run of Illumina MiSeq run. Just 24 sample *tags 1* and 20 plate *tags 2*  
519 allow the differentiation of up to 480 samples. This greatly reduced sequencing and primer  
520 purchase costs while also largely eliminating sample-misassignment via tag jumping,  
521 because tag jump sequences have non-matching forward and reverse *tag 1* sequences [29].  
522 For our sequenced PCR plates, the rate of correct matching *tag 2* tags was 96%. We  
523 estimated the rate of tag jumps producing chimeric *tag 1* sequences to be of 1 to 5 % and  
524 these were removed from the dataset (Table 4). Twin-tagging increases costs because of the  
525 need to purchase a larger number of primer pairs. However, the risk of reporting false  
526 positives should compensate this, especially when it comes to rare or threatened species.

527 For the second PCR round, we used the same tag pair *tag 2* for all 24 samples of a PCR plate.  
528 In order to reduce cost we tested pooling these 24 samples prior to the second PCR round,  
529 but we detected a very high tag jumping rate of over 40% (Table 4), which ultimately would  
530 increase cost through reduced sequencing efficiency.

531 Tagging primers in the first PCR reduces the risk of cross-contamination via aerosolised PCR  
532 products. Previous studies have shown that unlabelled volatile PCR products pose a great  
533 risk of false detections [53], a risk that is greatly increased if a high number of samples are  
534 analysed in the laboratories [13]. Also, in laboratories where other research projects are  
535 conducted, this approach allows the detection of cross-experiment contamination.

536 Therefore, we see a clear advantage of our approach over ligation techniques when it  
537 comes to producing sequencing libraries, as the Illumina tags are only added after the first  
538 PCR, and thus the risk of cross contamination with unlabelled PCR amplicons is very low.

#### 539 *Assignment results*

540 A robust assignment of species is an important factor in metabarcoding as an incorrect  
541 identification might result incorrect management interventions. The reliability of taxonomic  
542 assignments is expected to vary with respect to both marker choice and database  
543 completeness, and this is reflected in the probability estimates provided by *PROTAX*. In a  
544 recent study, less than 10% of the mammal assignments made at species level against a  
545 worldwide reference database were considered reliable with the short 16S amplicon, but  
546 this increased to 46% with full-length 16S sequences [24]. In contrast, in the same study  
547 over 80% of insect assignments at species level were considered reliable with a more  
548 complete, geographically restricted database of full-length COI barcodes. A similar pattern



549 was observed in our data during manual curation of the assignment results – there was  
550 more ambiguity in the results for the short *16S* amplicon than for other markers. However,  
551 due to the limited amount of often degraded target DNA in e/iDNA samples, short  
552 amplicons amplify much better. In our case, this had the drawback that some species lacked  
553 any interspecific variation, and thus sequencing reads shared 99%-100% identity for several  
554 species. For example, our only *16S* reference of *Sus barbatus* was 100% identical to *S.*  
555 *scrofa*. But as latter species does not occur in the studied area we could assign all reads  
556 manually to *S. barbatus*. In several cases we were able to confirm *S. barbatus* by additional  
557 *CytB* results, highlighting the advantage of using multiple markers. Another important  
558 advantage of multiple markers is the opportunity to fill gaps in the reference database. For  
559 example, we lacked *16S* reference sequences for *Hystrix brachyura*, and reads were  
560 assigned by *PROTAX* only to the genus level: *Hystrix sp.*. In one sample, however, almost  
561 5000 *CytB* reads were assigned to *Hystrix brachyura* and thus we used the *Hystrix sp.* *16S*  
562 sequences in the same sample to build a consensus *16S* reference sequence for *Hystrix*  
563 *brachyura* for future analyses. We also inferred that PCR and sequencing errors resulted in  
564 reads being assigned to sister taxa. We observed that a high number of reads of a true  
565 sequence were assigned to a species and a lower number of noise sequences were assigned  
566 to a sister taxa. Such a pattern was observed for ungulates, especially deer that showed  
567 little variance in *16S*. It is hard to identify and control for such pattern automatically, and it  
568 highlights the importance of visual inspection of the results.

569 In total, we accepted 174 vertebrate detections (i.e. having positive detections in both  
570 *extraction replicates A and B*) within 96 bulk samples. 48% of these assignments were  
571 present in all four *A1, A2, B1 and B2*. 35% were present in at least three of replicates (e.g.  
572 *A1, A2, B1*). Although the true occurrence of species within our leeches was unknown, by  
573 accepting only positive *AB* assignment results, we increase the confidence of species  
574 detection, even if the total number of reads for that species was low. In almost all cases,  
575 however, the number of reads was high (median= 52,386; mean= 300,996; SD= 326,883).  
576 Keeping this in mind we do not believe that raw read numbers are the most reliable  
577 indicators of tetrapod DNA quantity in iDNA samples. PCR stochasticity, primer biases,  
578 multiple species in individual samples, and pooling of samples exert too many uncertainties  
579 that could bias the sequencing results. Replication of detection is inherently more reliable.  
580 In contrast to our expectation that higher cycle number might be necessary to amplify even  
581 the lowest amounts of target DNA, our data does not support this hypothesis. Although we  
582 observed an increase in positive PCRs for *A2/B2* (the 40-cycle PCR replicates), the total  
583 number of accepted assignments in *A1/B1* and *A2/B2* samples did not differ. This indicates  
584 first that high PCR cycle numbers mainly increased the risk of false positives and second that  
585 our multiple precautions successfully minimized the acceptance of false detections.

## 586 Conclusion

587 Metabarcoding of e/iDNA samples will certainly become a very valuable tool in assessing  
588 biodiversity, as it allows to detect species non-invasively without the need to capture and

589 handle the animals [54]. However, the technical and analytical challenges linked to sample  
590 types (low quantity and quality DNA) and poor reference databases have so far been  
591 insufficiently recognized. In contrast to ancient DNA studies where standardized laboratory  
592 procedures and specialized bioinformatics pipelines have been established and are followed  
593 in most cases, there is limited methodological consensus in e/iDNA studies, which reduces  
594 rigour. In this study, we present a robust metabarcoding workflow for e/iDNA studies. We  
595 hope that the provided scripts and protocols facilitate further development of rigour in this  
596 field. The use of e/iDNA metabarcoding to study the rarest and most endangered species  
597 such as the saola is exciting, but geneticists bear the heavy responsibility of providing  
598 correct answers to conservationists.

## 599 References

- 600 [1] Garipey TD, Lindsay R, Odgen N, Greory TR. Identifying the last supper: utility of the  
601 DNA barcode library for bloodmeal identification in ticks. *Mol Ecol Res.* 2012; 12: 646-  
602 52; doi: 10.1111/j.1755-0998.2012.03140.x
- 603 [2] Lee P-S, Gan HM, Clements GR, Wilson J-J. Field calibration of blowfly-derived DNA  
604 against traditional methods for assessing mammal diversity in tropical forests.  
605 *Genome* 2016; 59: 1008-22; doi:10.1139/gen-2015-0193
- 606 [3] Calvignac-Spencer S, Merkel K, Kutzner N, et al.. Carrion fly-derived DNA as a tool for  
607 comprehensive and cost-effective assessment of mammalian biodiversity. *Mol Ecol.*  
608 2013; 22: 915-24; doi:10.1111/mec.12183
- 609 [4] Rodgers, TW, Xu CCY, Giacalone J, et al.. Carrion fly-derived DNA metabarcoding is an  
610 effective tool for mammal surveys: Evidence from a known tropical mammal  
611 community. *Mol Ecol Res.* 2017; 17(6): 1-13; doi:10.1111/1755-0998.12701
- 612 [5] Hoffmann C, Merkel K, Sachse A, et al.. Blow flies as urban wildlife sensors. *Mol Ecol*  
613 *Res* 2018; 18(3): 502-10; doi: 10.1111/1755-0998.12754
- 614 [6] Schönberger AC, Wagner S, Tuten HC, et al.. Host preferences in host-seeking and  
615 blood-fed mosquitoes in Switzerland. *Med Vet Entomol.* 2015; 30(1): 39-52.
- 616 [7] Taylor L, Cummings RF, Velten R, et al.. Host (Avian) Biting Preference of Southern  
617 California *Culex* Mosquitoes (Diptera: Culicidae). *J Med Entomol.* 2012; 49(3): 687-96.
- 618 [8] Townzen JS, Brower AVZ, Judd DD. Identification of mosquito bloodmeals using  
619 mitochondrial cytochrome oxidase subunit I and cytochrome b gene sequences. *Med*  
620 *Vet Entomol.* 2008; 22. 386-93.
- 621 [9] Kocher A, Thoisy B, Catzeflies F, et al.. iDNA screening: Disease vectors as vertebrate  
622 samplers. *Mol Ecol.* 2017; 26(22): 6478-86.
- 623 [10] Schnell IB, Thomsen PF, Wilkinson N, et al.. Screening mammal biodiversity using DNA  
624 from leeches. *Curr Biol.* 2012, 22(8): R262—3.
- 625 [11] Tessler M, Weiskopf SR, Berniker L, et al.. Bloodlines: mammals, leeches, and  
626 conservation in southern Asia. *Syst Biodivers.* 2018; 1-9.
- 627 [12] Weiskopf SR, McCarthy KP, Tessler M, et al.. Using terrestrial haematophagous  
628 leeches to enhance tropical biodiversity monitoring programmes in Bangladesh. *J Appl*  
629 *Ecol.* 2018: 1-11.
- 630 [13] Schnell IB, Bohmann K, Schultze SE, et al.. Debugging diversity - a pan-continental  
631 exploration of the potential of terrestrial blood-feeding leeches as a vertebrate  
632 monitoring tool. *Mol Ecol Res.*2018; doi: 10.1111/1755-0998.12912
- 633 [14] Calvignac-Spencer S, Leendertz FH, Gilbert MT, Schubert G. An invertebrate stomach's  
634 view on vertebrate ecology: certain invertebrates could be used as "vertebrate

- 635 samplers" and deliver DNA-based information on many aspects of vertebrate ecology.  
636 BioEssays. 2013; 35(11): 1004-13.
- 637 [15] Schnell IB, Sollmann R, Calvignac-Spencer S, et al.. iDNA from terrestrial  
638 haematophagous leeches as a wildlife surveying and monitoring tool – prospects,  
639 pitfalls and avenues to be developed. Front Zool. 2015; 12:24.
- 640 [16] Pääbo S, Poinar H, Serre D, et al.. Genetic analyses from ancient DNA. Annu Rev  
641 Genet. 2004; 38: 645-79.
- 642 [17] Hofreiter M, Paijmans JL, Goodchild H, et al. The future of ancient DNA: Technical  
643 advances and conceptual shifts. BioEssays. 2015; 37(3): 284-93.
- 644 [18] Bonin A, Taberlet P, Zinger L, Coissac E. Environmental DNA: For Biodiversity Research  
645 and Monitoring. 1<sup>st</sup> ed. Oxford University Press; 2018.
- 646 [19] Ficetola GF, Pansu J, Bonin A, et al.. Replication levels, false presences and the  
647 estimation of the presence/absence from eDNA metabarcoding data. Mol Ecol Res.  
648 2014; 15(3): 543-56.
- 649 [20] Ficetola GF, Taberlet P., Coissac E. How to limit false positives in environmental DNA  
650 and metabarcoding? Mol Ecol Res. 2016; 16(3): 604-7.
- 651 [21] Dalton R. Still looking for that woodpecker. Nature. 2010; 463: 718-9.
- 652 [22] Bazinet AL, Cummings MP. A comparative evaluation of sequence classification  
653 programs. BMC bioinformatics. 2012; 13(1): 92.
- 654 [23] Richardson RT, Bengtsson-Palme J, Johnson RM. Evaluating and optimizing the  
655 performance of software commonly used for the taxonomic classification of DNA  
656 metabarcoding sequence data. Mol Ecol Res. 2017; 17(4): 760-9.
- 657 [24] Somervuo P, Yu DW, Xu CC, Ji Y, et al.. Quantifying uncertainty of taxonomic  
658 placement in DNA barcoding and metabarcoding. Methods Ecol Evol. 2017; 8(4): 398-  
659 407.
- 660 [25] Quast C, Gerken J, Schweer T, et al. SILVA Databases. In: Nelson KE. Encyclopedia of  
661 Metagenomics. 1<sup>st</sup> ed. Springer US; 2015. p. 626-635.
- 662 [26] Kozlov AM, Zhang J, Yilmaz P, Glöckner FO, Stamatakis A. (2016). Phylogeny-aware  
663 identification and correction of taxonomically mislabeled sequences. Nucleic Acids  
664 Res. 2016; 44(11): 5022-33.
- 665 [27] Huson DH, Auch AF, Qi J, Schuster SC. MEGAN analysis of metagenomic data. Genome  
666 Res. 2007; 17(3): 377-86.
- 667 [28] Somervuo P, Koskela S, Pennanen J, Henrik Nilsson R, Ovaskainen O. Unbiased  
668 probabilistic taxonomic classification for DNA barcoding. Bioinformatics. 2016; 32(19):  
669 2920-7.

- 670 [29] Schnell IB, Bohmann K, Gilbert MTP. (2015). Tag jumps illuminated—reducing  
671 sequence-to-sample misidentifications in metabarcoding studies. *Mol Ecol Res.* 2015;  
672 15(6): 1289-1303.
- 673 [30] Kocher TD, Thomas WK, Meyer A, et al.. Dynamics of mitochondrial DNA evolution in  
674 animals: amplification and sequencing with conserved primers. *Proc. Natl. Acad. Sci.*  
675 *U.S.A.*. 1989; 86(16): 6196-6200.
- 676 [31] Taylor PG. Reproducibility of ancient DNA sequences from extinct Pleistocene fauna.  
677 *Mol Biol Evol.* 2996; 13(1): 283-5.
- 678 [32] Faircloth BC, Glenn TC. Not all sequence tags are created equal: designing and  
679 validating sequence identification tags robust to indels. *PLoS One.* 2012; 7(8): e42543
- 680 [33] Murray DC, Coghlan ML, Bunce M. From benchtop to desktop: important  
681 considerations when designing amplicon sequencing workflows. *PLoS One.* 2015;  
682 10(4): e0124671.
- 683 [34] Machida RJ, Leray M, Ho SL, Knowlton N. Metazoan mitochondrial gene sequence  
684 reference datasets for taxonomic assignment of environmental samples. *Sci Data.*  
685 2017; 4: 170027.
- 686 [35] Chamberlain SA, Szöcs E. taxize: taxonomic search and retrieval in R. Version 2.  
687 *F1000Res.* 2013; 2: 191.
- 688 [36] Salleh FM, Ramos-Madriral J, Peñalosa F, et al.. An expanded mammal mitogenome  
689 dataset from Southeast Asia. *GigaScience.* 2017; 6(8): 1-8
- 690 [37] Kans, Jonathan. Entrez Direct: E-utilities on the UNIX Command Line. In: Entrez  
691 Programming Utilities Help [Internet]. Bethesda (MD): National Center for  
692 Biotechnology Information (US). 2010.
- 693 [38] Altschul SF, Gish W, Miller W, Myers EW, Lipman DJ. Basic local alignment search tool.  
694 *Journal of molecular biology.* 1990; 215(3):, 403-10.
- 695 [39] Camacho C, Coulouris G, Avagyan V, et al.. BLAST+: architecture and applications. *BMC*  
696 *bioinformatics.* 2009; 10(1): 421.
- 697 [40] Katoh K, Standley DM. (2013). MAFFT multiple sequence alignment software version  
698 7: improvements in performance and usability. *Mol Biol Evol.* 2013; 30(4): 772-80.
- 699 [41] Katoh K, Misawa K, Kuma KI, Miyata T. MAFFT: a novel method for rapid multiple  
700 sequence alignment based on fast Fourier transform. *Nucleic Acids Res.* 2002; 30(14):  
701 3059-66.
- 702 [42] Chesters D. (2013) *collapsetypes.pl* [computer software available at  
703 <http://sourceforge.net/projects/collapsetypes/>]
- 704 [43] Edgar RC. Search and clustering orders of magnitude faster than BLAST.  
705 *Bioinformatics.* 2010; 26(19): 2460-2461.

- 706 [44] Edgar RC, Flyvbjerg H. Error filtering, pair assembly and error correction for next-  
707 generation sequencing reads. *Bioinformatics*. 2015; 31(21): 3476-82.
- 708 [45] Martin M. Cutadapt removes adapter sequences from high-throughput sequencing  
709 reads. *EMBnet. Jjournal*. 2011; 17(1): 10-12.
- 710 [46] Kiełbasa SM, Wan R, Sato K, Horton P, Frith MC. Adaptive seeds tame genomic  
711 sequence comparison. *Genome Res*. 2011; 21(3): 487-493.
- 712 [47] Ratnasingham S, Hebert PDN. BOLD: The Barcode of Life Data System  
713 ([www.barcodinglife.org](http://www.barcodinglife.org)). *Mol Ecol Notes*. 2007; 3: 355-64.
- 714 [48] Evans NT, Li Y, Renshaw MA, et al. Fish community assessment with eDNA  
715 metabarcoding: effects of sampling design and bioinformatic filtering. *Can J Fish Aquat  
716 Sci*. 2017; 74(9):, 1362-74.
- 717 [49] Zepeda-Mendoza ML, Bohmann K, Baez AC, Gilbert MTP. DAME: a toolkit for the initial  
718 processing of datasets with PCR replicates of double-tagged amplicons for DNA  
719 metabarcoding analyses. *BMC Res Notes*. 2016; 9(1): 255.
- 720 [50] Racimo F, Renaud G, Slatkin M. Joint estimation of contamination, error and  
721 demography for nuclear DNA from ancient humans. *PLoS Genet*. 2016; 12(4):  
722 e1005972.
- 723 [51] Orlando L, Gilbert MTP, Willerslev E. Reconstructing ancient genomes and  
724 epigenomes. *Nat Rev Genet* 2015; 16(7): 395
- 725 [52] Laurin-Lemay S, Brinkmann H, Philippe H. Origin of land plants revisited in the light of  
726 sequence contamination and missing data. *Current Biology*. 2012; 22(15): R593-4.
- 727 [53] Kwok S, Higuchi R. Avoiding false positives with PCR. *Nature*. 1989; 339: 237-8.
- 728 [54] Bush A, Sollmann R, Wilting A, et al.. Connecting Earth observation to high-throughput  
729 biodiversity data. *Nat Ecol Evol* 2017; 1(7): 0176.



730 **Table 1:** Sequence motifs that compose the 25 different target primers for the first and the  
 731 second PCR. First PCR primers consist of target specific primer followed by an overhang out  
 732 of sample specific *tag 1* and *read 1* and *read 2* sequencing primer, respectively. The second  
 733 PCR primers consist of the *read 1* or the *read 2* sequencing primer followed by an plate  
 734 specific *tag 2* and the P5 and P7 adapters, respectively (see also Fig. 2).  
 735

Name	Sequence	Reference
tag A	TGCAT	Faircloth & and Glenn 2012
tag B	TCAGC	Faircloth & and Glenn 2012
tag C	AAGCG	Faircloth & and Glenn 2012
tag D	ACAAG	Faircloth & and Glenn 2012
tag E	AGTGG	Faircloth & and Glenn 2012
tag F	TTGAC	Faircloth & and Glenn 2012
tag G	CCTAT	Faircloth & and Glenn 2012
tag H	GGATG	Faircloth & and Glenn 2012
tag I	CTAGG	Faircloth & and Glenn 2012
tag K	CACCT	Faircloth & and Glenn 2012
tag L	GTCAA	Faircloth & and Glenn 2012
tag M	GAAGT	Faircloth & and Glenn 2012
tag N	CGGTT	Faircloth & and Glenn 2012
tag O	ACCGA	Faircloth & and Glenn 2012
tag P	ACGTC	Faircloth & and Glenn 2012
tag Q	AGACT	Faircloth & and Glenn 2012
tag R	AGGAA	Faircloth & and Glenn 2012
tag S	ATCC	Faircloth & and Glenn 2012
tag T	CAATC	Faircloth & and Glenn 2012
tag V	CATGA	Faircloth & and Glenn 2012
tag W	CCACA	Faircloth & and Glenn 2012
tag X	GCTTA	Faircloth & and Glenn 2012
tag Y	GGTAC	Faircloth & and Glenn 2012
tag Z	AACAC	Faircloth & and Glenn 2012
Tag Control	ATCTG	Faircloth & and Glenn 2012
<i>CytB</i> -fw	AAAAAGCTTCCATCCAACATCTCAGCATGATGAAA	Kocher et al. 1989
<i>CytB</i> -rv	AAACTGCAGCCCCTCAGAATGATATTTGTCCTCA	Kocher et al. 1989
<i>16S</i> -fw	CGGTTGGGGTGACCTCGGA	Taylor 1996
<i>16S</i> -rv	GCTGTTATCCCTAGGGTAACT	Taylor 1996
<i>12S</i> -fw	AAAAAGCTTCAAACCTGGGATTAGATACCCCACTAT	Kocher et al. 1989
<i>12S</i> -rv	TGACTGCAGAGGTGACGGCGGTGTGT	Kocher et al. 1989
Read 1 sequence primer	ACACTCTTTCCCTACACGACGCTCTTCCGATCT	Illumina Document # 1000000002694 v03
Read 2 sequence primer	GTGACTGGAGTTCAGACGTGTGCTCTTCCGATCT	Illumina Document # 1000000002694 v03
P5 adapter	AATGATACGGCGACCACCGAGATCTACAC	Illumina Document # 1000000002694 v03
P7 adapter	CAAGCAGAAGACGGCATACGAGAT	Illumina Document # 1000000002694 v03

736

737 **Table 2:** Main steps undertaken by each module of the database curation script.

<b>MODULE</b>	<b>STEPS</b>
Module 1	<p>Extract subset of raw MIDORI database for query taxon and loci.</p> <p>Remove sequences with non-binomial species names, reduce subspecies to species labels</p> <p>Add local sequences (optional)</p> <p>Check for relevant new sequences for list of query species on NCBI (GenBank and RefSeq) (optional)</p> <p>Select amplicon region and remove primers</p> <p>Remove sequences with ambiguous bases</p> <p>Align</p> <p>End of module: Optional check of alignments</p>
Module 2	<p>Compare sequence species labels with taxonomy</p> <p>Non-matching labels queried against Catalogue of Life to check for known synonyms</p> <p>Remaining mismatches kept if genus already exists in taxonomy, otherwise flagged for removal</p> <p>End of module: Optional check of flagged species labels</p>
Module 3	<p>Discard flagged sequences</p> <p>Update taxonomy key file for sequences found to be incorrectly labelled in Module 2</p> <p>Run SATIVA</p> <p>End of module: Optional check of putatively mislabelled sequences</p>
Module 4	<p>Discard flagged sequences</p> <p>Finalise consensus taxonomy and relabel sequences with correct species label and accession number</p> <p>Select one representative sequence per haplotype per species</p>

738 **Table 3:** GNU core utilities and other lightweight tools used for manipulation of text and  
739 sequence files

<b>TOOL</b>	<b>FUNCTION</b>	<b>SOURCE</b>
awk, cut, grep, join, sed, sort, tr	Processing text files	GNU core utilities
seqbuddy	Processing FASTA/Q files	<a href="https://github.com/biologyguy/BuddySuite">https://github.com/biologyguy/BuddySuite</a>
seqkit	Processing FASTA/Q files	<a href="https://github.com/shenwei356/seqkit">https://github.com/shenwei356/seqkit</a>
seqtk	Processing FASTA/Q files	<a href="https://github.com/lh3/seqtk">https://github.com/lh3/seqtk</a>
tabtk	Processing tab-delimited text files	<a href="https://github.com/lh3/tabtk">https://github.com/lh3/tabtk</a>

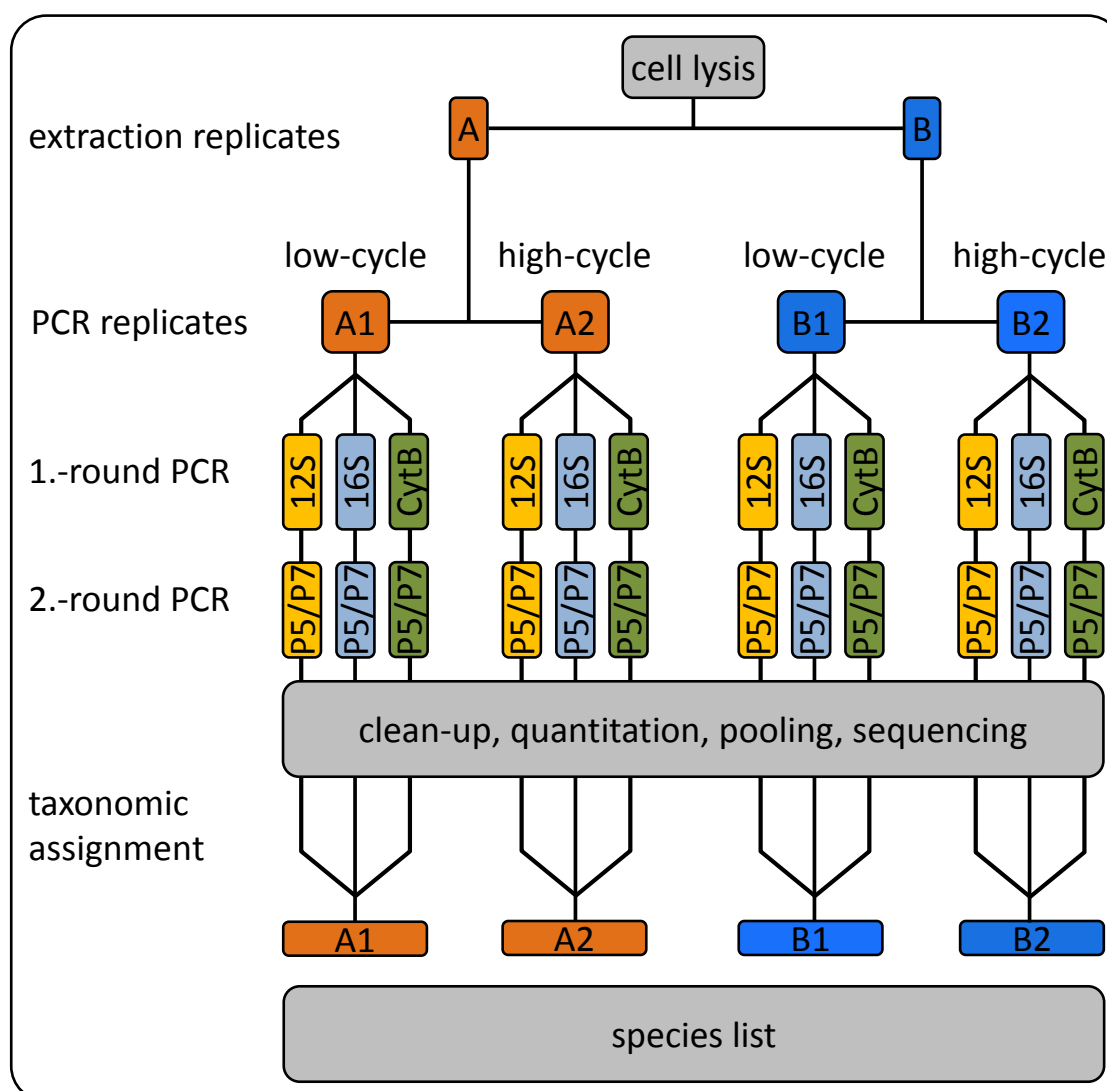
740

741 **Table 4:** Number of reads per sequencing run and the numbers of reads with matching, chimeric or unidentifiable tags.

	<b>total reads</b>	<b>matching tag 2 reads</b>	<b>chimeric tag 2 reads</b>	<b>%<sup>1</sup></b>	<b>matching tag 1 reads</b>	<b>chimeric tag 1 reads</b>	<b>%<sup>2</sup></b>	<b>erroneous tag 1 reads</b>	<b>%<sup>2</sup></b>
<b>SeqRun01</b>	18,438,517	18,102,702	282,419	1.5	17,514,515	451,028	2.5	137,159	0.8
<b>SeqRun02</b>	25,385,558	24,596,380	626,245	2.5	23,426,084	612,045	2.5	558,251	2.3
<b>SeqRun03</b>	14,875,796	14,393,884	343,528	2.3	13,766,187	426,181	3.0	201,516	1.4
<b>SeqRun04</b>	2,027,794	1,935,149	56,077	2.8	1,806,655	88,307	4.6	40,187	2.1
<b>SeqRun05</b>	18,221,504	17,500,366	421,588	2.3	16,793,851	482,365	2.8	161,458	0.9
<b>SeqRun06</b>	20,718,202	19,874,913	429,048	2.1	19,317,305	371,048	1.9	81,422	0.4
<b>SeqRun07</b>	24,604,610	23,746,938	663,730	2.7	22,446,187	497,366	2.1	803,385	3.4
<b>Total</b>	124,271,981	120,150,332	2,822,635	2.3	115,070,784	2,928,340	2.5	1,983,378	1.7
<b>IndexRun</b>	10,276,093	10,116,808	NA	NA	5,841,190	4,186,688	41.4	88,930	0.9

<sup>1</sup> refers to total reads

<sup>2</sup> refers to matching tag 2



742

743

744

745

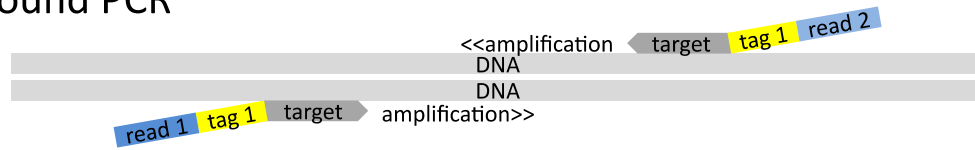
746

747

748

**Figure 1:** laboratory scheme; during DNA extraction the sample is split into two extraction replicates A & B. Our Protocol consists of two rounds of PCR that were the sample tags, the necessary sequencing primer and sequencing adapters are added to the the amplicons. For each extraction replicate we ran a low cycle PCR and a high cycle PCR for each marker that we have twelve independent PCR replicates per sample. All PCR products were sequenced and the obtained reads were taxonomically identified with PROTAX.

### 1.-round PCR



### 1.-round product:



### 2.-round PCR



### 2.-round product:



749

750

751

752

753

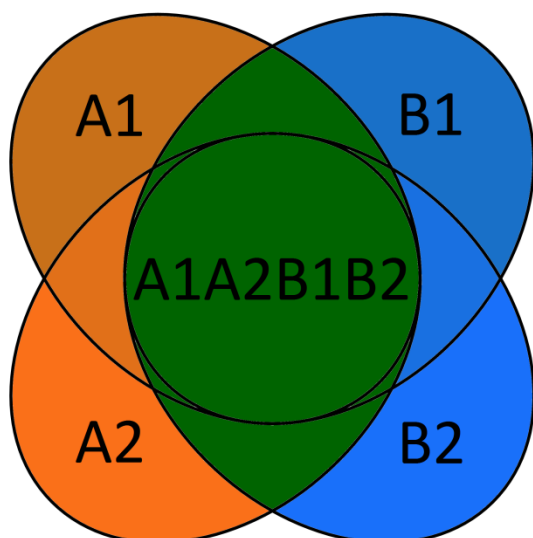
754

755

756

**Figure 2:** Scheme to build double ‘twin-tagged’ PCR libraries. The first round of PCR uses target-specific primers (12S, 16S, or CytB, dark grey) that have both been extended with the same (i.e. ‘twin’) sample-identifying *tag* sequences *tag 1* (yellow) and then with the different *read 1* (dark blue) and *read 2* (light blue) sequence primers. The second round of PCR uses the priming sites of the *read 1* and *read 2* sequencing primers to add twin plate-identifying *tag* sequences *tag 2* (orange) and the P5 (dark red) and P7 (light red) Illumina adapters.



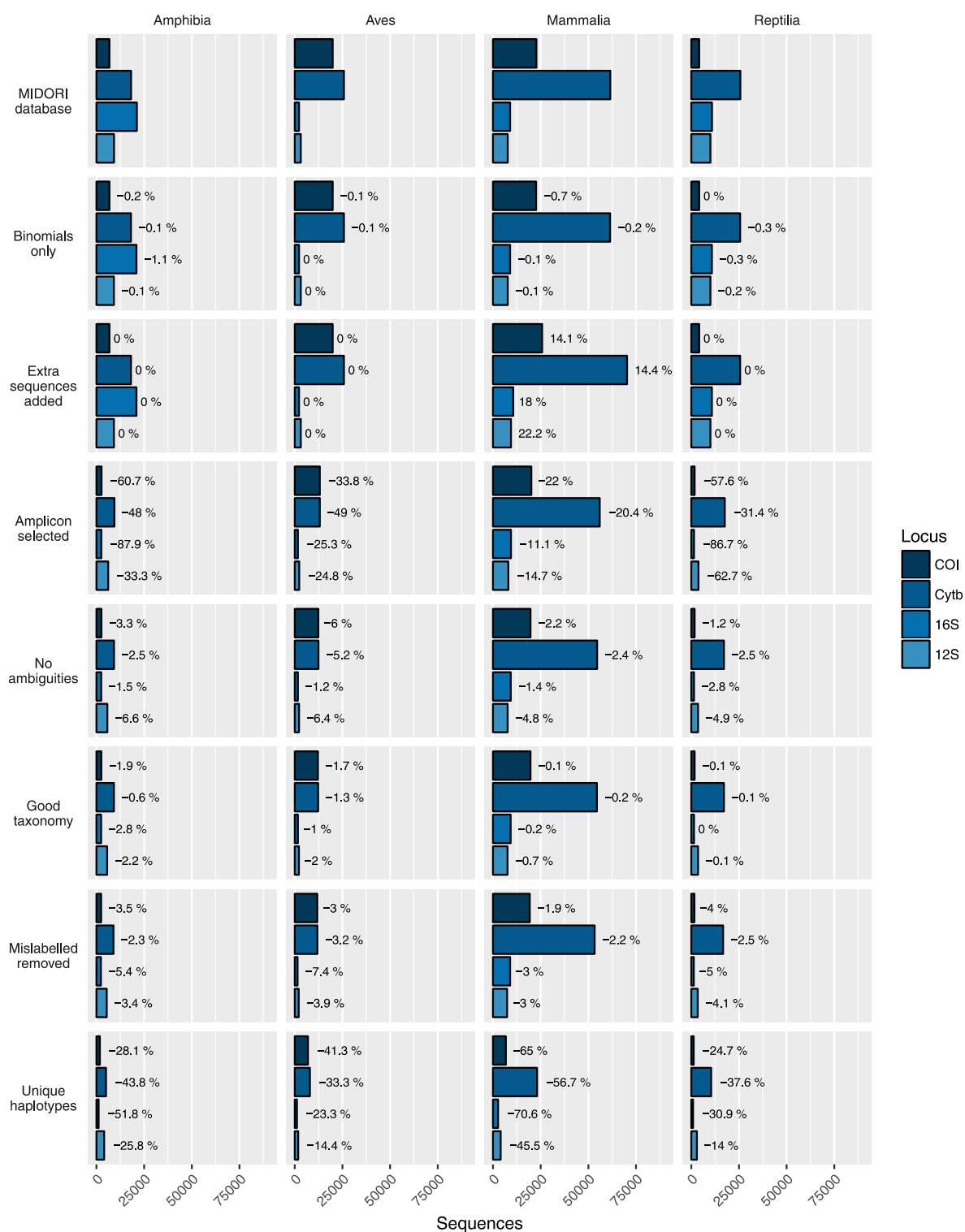


757

758

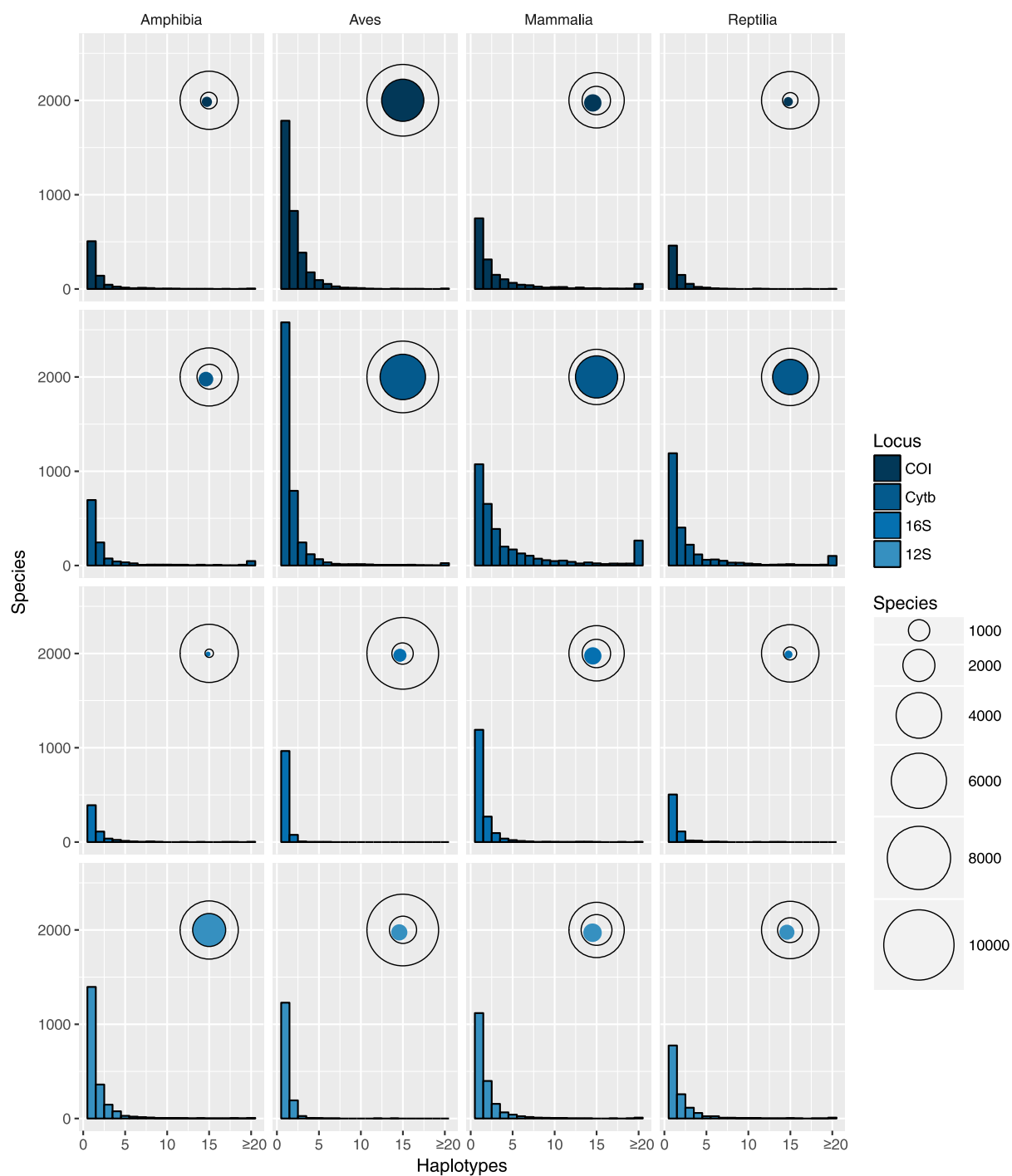
759

**Figure 3:** We only accepted taxonomic assignments that were positively detected in both *extraction replicates* A and B (green colour).



760

761 **Figure 4:** Data availability and percentage loss at each major step in the database curation procedure  
 762 for each target amplicon and class of Tetrapoda. The number of sequences decreases between steps  
 763 except “Extra sequences added” where additional target sequences are included for Mammalia and  
 764 there is no change for the other three classes.



765

766

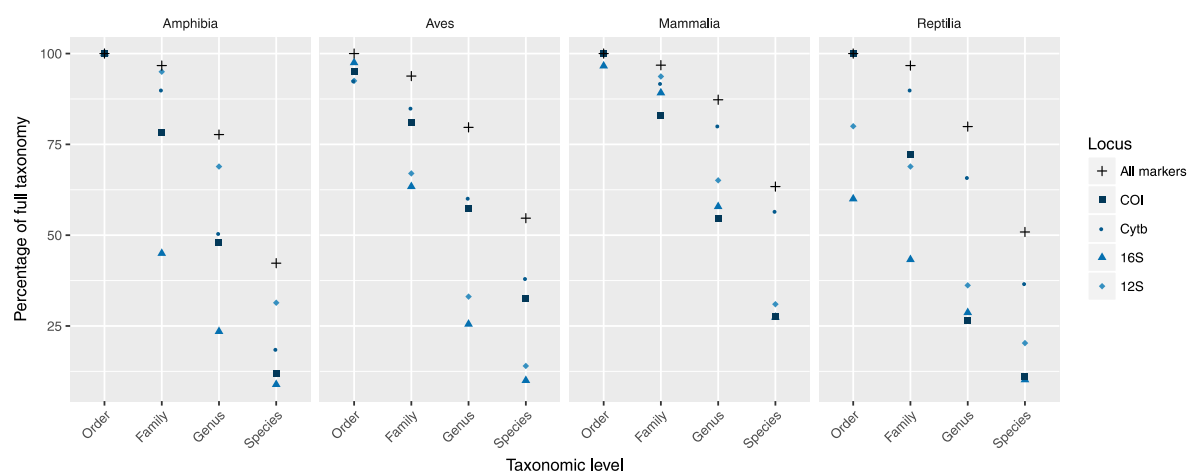
767

768

769

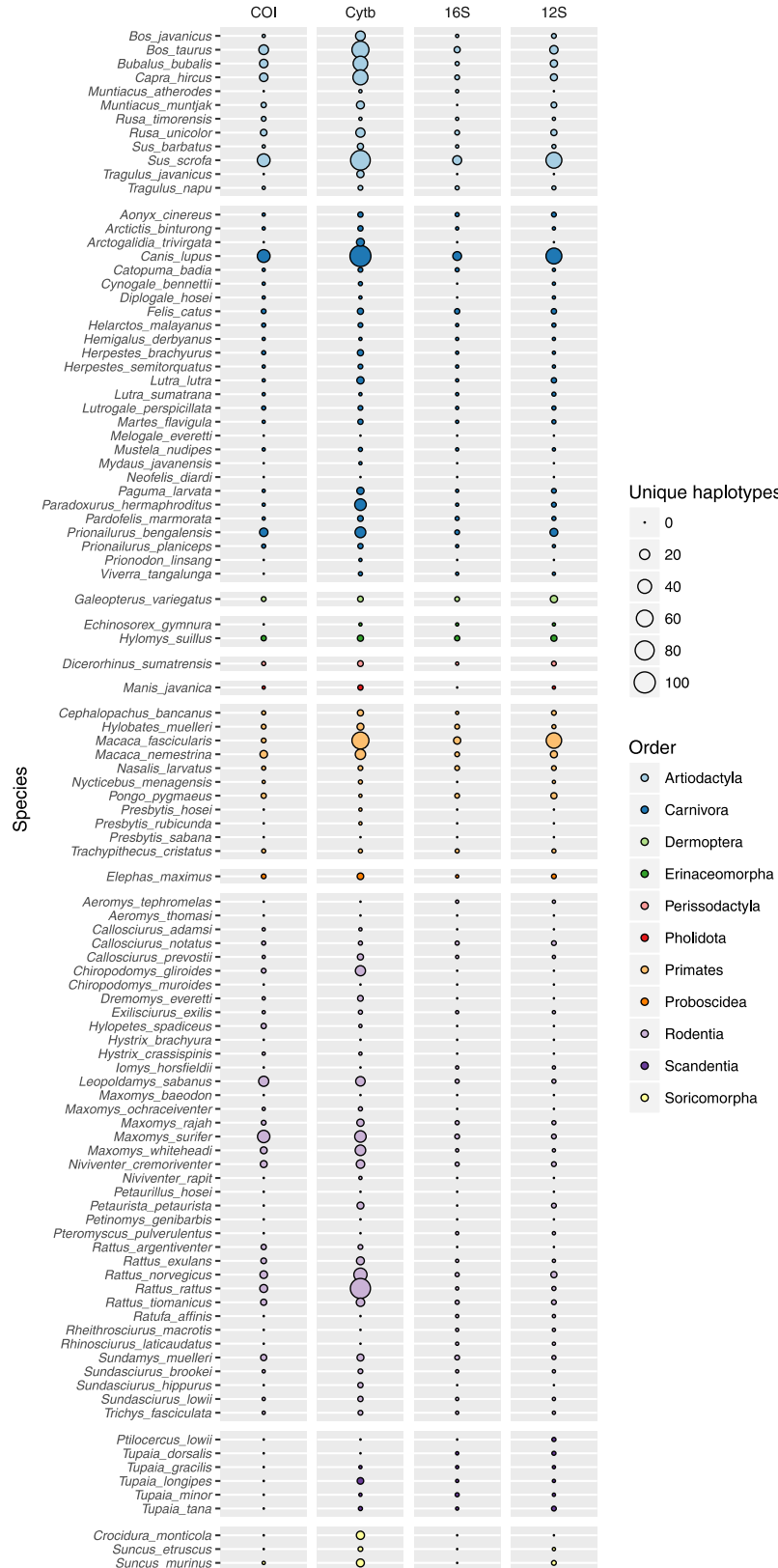
770

**Figure 5:** Haplotype number by species (frequency distribution) and the total number of species with at least one haplotype, shown relative to the total number of species in the taxonomy for that category (bubbles), shown for each marker and class of Tetrapoda. The proportion of species covered by the database varies between categories but in all cases a majority of recovered species are represented by a single unique haplotype.



771

772 **Figure 6:** The percentage of the full taxonomy covered by the final database at each taxonomic level  
 773 for each class of Tetrapoda. Includes the percentage of taxa represented by each marker and all  
 774 markers combined. In all cases taking all four markers together increases the proportion of species,  
 775 genera and families covered by the database but it remains incomplete when compared with the full  
 776 taxonomy.



777

778 **Figure 7:** The number of unique haplotypes per marker for each of the 103 mammal species  
 779 expected in the study area. Bubble size is proportional to the number of haplotypes and varies  
 780 between 0 and 100. Only 554 species have at least one sequence per marker and nine species are  
 781 completely unrepresented in the current database.