

MinION barcodes: biodiversity discovery and identification by everyone, for everyone

Amrita Srivathsan¹, Leshon Lee¹, Kazutaka Katoh^{2,3}, Emily Hartop^{4,5}, Sujatha Narayanan Kutty^{1,6}, Johnathan Wong¹, Darren Yeo¹, Rudolf Meier¹

¹ Department of Biological Sciences, National University of Singapore, Singapore

² Research Institute for Microbial Diseases, Osaka University, Japan

³ Artificial Intelligence Research Center, AIST, Tokyo, Japan

⁴ Zoology Department, Stockholms Universitet, Stockholm, Sweden

⁵ Station Linné, Öland, Sweden

⁶ Tropical Marine Science Institute, National University of Singapore, Singapore

1 **Abstract**

2 DNA barcodes are a useful tool for discovering, understanding, and monitoring biodiversity
3 which are critical at a time when biodiversity loss is a major problem for many countries.
4 However, widespread adoption of barcodes requires cost-effective and simple barcoding
5 methods. We here present a workflow that satisfies these conditions. It was developed via
6 “innovation through subtraction” and thus requires minimal lab equipment, can be learned
7 within days, reduces the barcode sequencing cost to <10 cents, and allows fast turnaround
8 from specimen to sequence by using the real-time sequencer MinION. We first describe
9 cost-effective and rapid procedures for obtaining tagged amplicons. We then demonstrate
10 how a portable MinION device can be used for real-time sequencing of tagged amplicons in
11 many settings (field stations, biodiversity labs, citizen science labs, schools). Small projects
12 can use the flow cell dongle (“Flongle”) while large projects can rely on MinION flow cells
13 that can be stopped and re-used after collecting sufficient data for a given project. We also
14 provide amplicon coverage recommendations that are based on several runs of MinION flow
15 cells (R10.3) which suggest that each run can generate >10,000 barcodes. Next, we present
16 a novel software, ONTbarcoder, which overcomes the bioinformatics challenges posed by
17 the sequencing errors of MinION reads. This software is compatible with Windows10,
18 Macintosh, and Linux, has a graphical user interface (GUI), and can generate thousands of
19 barcodes on a standard laptop within hours based on two input files (FASTQ, demultiplexing
20 file). We document that MinION barcodes are virtually identical to Sanger and Illumina
21 barcodes for the same specimens (>99.99%). Lastly, we demonstrate how rapidly MinION
22 data have improved by comparing the performance of sequential flow cell generations. We
23 overall assert that barcoding with MinION is the way forward for government agencies,
24 universities, museums, and schools because it combines low consumable and capital cost
25 with scalability. Biodiversity loss is threatening the planet and the use of MinION barcodes
26 will help with enabling an army of researchers and citizen scientists, which is necessary for
27 effective biodiversity discovery and monitoring.

29 **1. Background**

30 DNA sequences have been used for identification and taxonomic purposes for decades
31 (Hebert, Cywinska et al. 2003, Tautz, Arctander et al. 2003, Meier 2008), but for most of this
32 time been akin to mobile phones in the 1990s: of limited value due to sparse signal coverage
33 and high cost. Obtaining barcodes was problematic due largely to the complicated and
34 expensive procedures on which it relied. Some of these problems have since been
35 addressed by, for example, developing effective DNA extraction protocols and optimizing
36 Sanger sequencing procedures (Ivanova, Dewaard et al. 2006, Ivanova, Borisenko et al.
37 2009). These improvements enabled the establishment of a centralized barcoding facility in
38 2006. After 15 years and the investment of >200 million USD, ca. 8 million animal barcodes
39 are available for searches on BOLD Systems, but only ca. 6 million are in the public domain
40 (http://boldsystems.org/index.php/IDS_OpenIdEngine). Combined with barcodes from NCBI
41 GenBank, they are now a valuable resource to the global biodiversity community. However,
42 the cost of barcodes has remained high (<http://ccdb.ca/pricing/>) and the prevalent approach
43 for sizeable projects is sending specimens from all over the world to one center and then
44 only some back to the country of origin. This interferes with real-time biodiversity monitoring
45 and specimen accessibility. We therefore argue that access to barcodes has to be
46 democratized through decentralization. We here show that this achievable because the
47 application of a technique that is known as “innovation through subtraction” in engineering
48 readily yields simplified and cost-effective solutions for DNA barcoding and the amplification
49 and sequencing of a short mitochondrial COI fragment can be efficiently performed
50 anywhere and by biologists and citizen scientists alike.

51

52 A decentralized model for monitoring the world’s biodiversity is necessary given the scale,
53 urgency, and importance of the task at hand. For example, even if there were only 10 million
54 species of metazoan animals on the planet (Stork, McBroom et al. 2015) and a new species
55 is discovered with every 50th specimen that is processed, species discovery with barcodes
56 will require the sequencing of 500 million specimens (Yeo, Srivathsan et al. 2020). Yet,

57 species discovery is only a small part of the biodiversity challenge in the 21st century.
58 Biodiversity loss is now considered by the World Economic Forum as one of the top three
59 global risks based on likelihood and impact for the next 10 years (World Economic Forum
60 2020) and Swiss Re estimates that 20% of all countries face ecosystem collapse as
61 biodiversity declines (Swiss Re 2020). Biodiversity loss is no longer just an academic
62 concern; it is now a major threat to human communities and the health of the planet. This
63 also implies that biodiversity discovery and monitoring require completely different scales
64 than in the past. The old approaches thus need rethinking because all countries need real-
65 time distributional and abundance information to develop effective conservation strategies
66 and policies. In addition, they need information on how species interact with each other and
67 the environment (Abrego, Roslin et al. 2021). Many of these biodiversity monitoring and
68 environmental management activities have to focus on terrestrial invertebrates, whose
69 biomass surpasses that of all terrestrial vertebrates combined (Bar-On, Phillips et al. 2018)
70 and who occupy a broad range of ecological guilds. The main obstacles are high numbers of
71 specimens and species and the rapid decline of many of these taxa (Bell, Blumgart et al.
72 2020, Eisenhauer, Bonn et al. 2019, Hallman, Sorg et al. 2017, Hallman, Ssymank et al.
73 2021, Stepanian, Entekin et al. 2020, Wagner, Grames et al. 2021) which means that
74 monitoring should be locally conducted to allow for rapid turnaround. This requires simple
75 and cost-effective procedures that can be implemented anywhere by stakeholders with very
76 different scientific and skill backgrounds.

77

78 DNA barcoding was proposed at a time when biodiversity loss was not on the radar of
79 economists. Instead, barcodes were initially intended as an identification tool for biologists
80 (Hebert, Cywinska et al. 2003). Thus, most projects focused on taxa with a large following in
81 biology (e.g., birds, fish, butterflies) (Kwong, Srivathsan et al. 2012). However, this also
82 meant that these projects only covered a small proportion of the terrestrial animal biomass
83 (Bar-On, Phillips et al. 2018) and species-level diversity (Groombridge 1992). Yet, despite
84 targeting taxa with well-understood diversity, the projects struggled with covering >75% of

85 the described species in these groups (Kwong, Srivathsan et al. 2012). When the pilot
86 barcoding projects ran out of material from identified specimens, they started targeting
87 unidentified specimens; i.e., DNA barcoding morphed into a technique that was used for
88 biodiversity discovery (“dark taxa”: Page 2011, Kwong, Srivathsan et al. 2012). This shift
89 towards biodiversity discovery was gradual and incomplete because the projects used a
90 “hybrid approach” that started with subsampling or sorting specimens to “morphospecies”
91 before barcoding representatives of each morphospecies/sample (e.g., Barrett and Hebert
92 2005, Hendrich, Pons et al. 2010, Hebert, DeWaard et al. 2013, Ng'endo, Osiemo et al.
93 2013, Hebert, Ratnasingham et al. 2016, Thormann, Ahrens et al. 2016, Knox, Hogg et al.
94 2020). This is problematic, as morphospecies sorting is known to be labour-intensive and of
95 unpredictable quality because it is heavily dependent on the taxonomic expertise of the
96 sorters (Krell 2004, Stribling, Pavlik et al. 2008). Thus, such hybrid approaches are of limited
97 value for obtaining reliable quantitative data on biodiversity, but were adopted as a
98 compromise owing to the prohibitive cost of barcoding. The logical alternative is to barcode
99 all specimens and then group them into putative species based on sequence information.
100 Such a “reverse workflow” (Wang, Srivathsan et al. 2018), where every specimen is
101 barcoded as the initial pre-sorting step, yields quantitative data and corroborated species-
102 level units. However, the reverse workflow requires efficient and low-cost barcoding methods
103 that are also suitable for biodiverse countries with limited science funding.
104
105 Fortunately, such cost-effective barcoding methods are now becoming available. This is
106 partially due to the replacement of Sanger sequencing with second- and third-generation
107 sequencing technologies that have lowered sequencing costs dramatically (Shokralla, Spall
108 et al. 2012, Shokralla, Porter et al. 2015, Meier, Wong et al. 2016, Hebert, Braukmann et al.
109 2018, Krehenwinkel, Kennedy et al. 2018, Srivathsan, Baloglu et al. 2018, Wang, Srivathsan
110 et al. 2018, Srivathsan, Hartop et al. 2019, Yeo, Srivathsan et al. 2020). Such changes mean
111 that the reverse workflow is now available for tackling the species-level diversity of those
112 metazoan clades that are so specimen- and species-rich that they have been neglected in

113 the past (Ponder and Lunney 1999, Srivathsan, Hartop et al. 2019). Many of these clades
114 have high spatial species turnover, requiring many localities in each country to be sampled
115 and massive numbers of specimens to be processed (Yeo, Srivathsan et al. 2020). Such
116 intensive processing is best achieved close to the collecting locality to avoid the
117 unnecessary risks, delays and costs from shipping biodiversity samples across continents.
118 This is now feasible because biodiversity discovery can be readily pursued in decentralized
119 facilities at varied scales. Indeed, accelerated biodiversity discovery is a rare example of a
120 big science initiative that allows for meaningful engagement of students and citizen scientists
121 and can in turn significantly enhance biodiversity education and appreciation (Pomerantz,
122 Peñafiel et al. 2018, Watsa, Erkenwick et al. 2020). This is especially so when stakeholders
123 not only barcode, but also image specimens, determine species abundances, and map
124 distributions of newly discovered species. All of which can be based on specimens collected
125 in their own backyard.

126

127 But can such decentralized biodiversity discovery really be effective? Within the last five
128 years, the students and interns in the laboratory of the corresponding author at the National
129 University of Singapore barcoded >330,000 specimens. After analyzing the first >140,000
130 barcoded specimens for selected taxa representing different ecological guilds, the alpha and
131 beta diversity of Singapore's arthropod fauna was analyzed based on ~8,000 putative
132 species which revealed that some habitats were unexpectedly species-rich and harboured
133 very unique faunas (e.g., mangroves, freshwater swamp: Yeo, Srivathsan et al. 2020;
134 Baloğlu, Clews et al. 2018). Barcodes even helped with the conservation of charismatic taxa
135 when they were used to identify the larval habitats for more than half of Singapore's damsel-
136 and dragonfly species (Yeo, Puniamoorthy et al. 2018) and facilitated species interaction
137 research and biodiversity surveys based on eDNA (Lim, Tay et al. 2016, Srivathsan,
138 Nagarajan et al. 2019). Biodiversity appreciation by the public was fostered by featuring
139 newly discovered species and their species interactions on "Biodiversity of Singapore" (BOS
140 >15,000 species: <https://singapore.biodiversity.online/>) and dozens of new species have

141 been described and the descriptions of another 150 species are being finalized (Grootaert
142 2018, Tang, Grootaert et al. 2018, Tang, Yang et al. 2018, Wang, Yamada et al. 2018,
143 Wang, Yong et al. 2018, Grootaert 2019, Ismay and Ang 2019, Samoh, Satasook et al.
144 2019, Wang, Yamada et al. 2020).

145

146 **2. Methods for the democratization of DNA barcoding through simplification**

147 Barcoding a metazoan specimen requires the successful completion of three steps: (1)
148 obtaining DNA template, (2) amplifying *COI* via PCR, and (3) sequencing the *COI* amplicon.
149 Many biologists learn these techniques in university for a range of different genes – from
150 those that are easy to amplify (short fragments of ribosomal and mitochondrial genes with
151 well-established primers) to those are difficult (long, single-copy nuclear genes with few
152 known primers). Fortunately, amplification of short mitochondrial markers like *COI* does not
153 require the same level of care as nuclear markers. Learning how to barcode efficiently is
154 hence an exercise of unlearning by applying “innovation through subtraction”. Note that this
155 unlearning is of critical importance for the democratization of biodiversity discovery with DNA
156 barcodes and is particularly vital for boosting biodiversity research where it is most needed:
157 in biodiverse countries with limited science funding.

158

159 In this section, we first briefly summarize commonly used procedures for DNA extraction,
160 PCR, and sequencing. For each step we then describe how the procedures can be
161 simplified. In addition to the description, we provide videos for the described procedures
162 which are available from the YouTube channel “Integrative Biodiversity Discovery”
163 (<https://www.youtube.com/channel/UC1WowokomhQJRc71FmsUAcg>). Note that all
164 techniques have been extensively tested in our lab, primarily on invertebrates preserved in
165 ethanol. Regarding sequencing, we briefly introduce four methods, but our focus is on
166 MinION sequencing because this device is particularly suitable as the default sequencing
167 option for decentralized biodiversity discovery.

168

169 Methods for step 1: Obtaining DNA template

170 Most biologists learn that DNA extraction requires tissue digestion, DNA purification, and
171 DNA elution. This approach is slow and expensive because it frequently involves kits and
172 consumables that are designed for obtaining the kind of high-quality DNA that is needed for
173 amplifying “difficult” genes (e.g., long, single-copy nuclear markers). However *COI* is a
174 mitochondrial gene and thus naturally enriched. Indeed, the tiny mitochondrial genome (16
175 kbp) usually contributes 0.5-5% of the DNA in a genomic extraction (Arribas, Andújar et al.
176 2016, Crampton-Platt, Yu et al. 2016). Furthermore, barcoding requires only the
177 amplification of one short marker (<700 bp) so that not much DNA template is needed. This
178 allows for using the following simplified procedures that are designed for specimens
179 containing DNA template of reasonable quality (e.g., Malaise trap specimens collected within
180 the last 20 years).

181

182 *Simplified DNA “extraction”*: Obtaining template for DNA barcoding need not take more than
183 20 minutes, does not require DNA purification, and costs essentially nothing. The cheapest,
184 but not necessarily fastest, method is “directPCR”; i.e., deliberately “contaminating” a PCR
185 reaction with the DNA of the target organism by adding the entire specimen or a tissue
186 sample into the PCR reagent mix (Wong, Tay et al. 2014). This method is very fast and
187 effective for small specimens lacking thick cuticle or skin (Wong, Tay et al. 2014) and works
188 particularly well for many abundant aquatic invertebrates such as chironomid midges and
189 larvae. Larger specimens require the use of body parts [leg or antenna: Wong, Tay et al.
190 (2014)]. Such dissections tend to be labour-intensive if large numbers of specimens must be
191 processed, but it is a good method for small numbers of samples or in barcoding
192 experiments that are carried out in poorly equipped labs. Note that the whole body or body
193 part that is used for directPCR can be recovered after amplification, although soft-bodied
194 animals may become transparent.

195

196 An alternative to directPCR is buffer-based DNA extraction. This method is also essentially
197 cost-free because it involves alkaline buffers that are inexpensive, usually available in
198 molecular labs (e.g., PBS), or can be prepared easily (HotSHOT) (Truett, Heeger et al. 2000,
199 Thongjued, Chotigeat et al. 2019)). Our preferred method is extraction with HotSHOT, which
200 we have used for barcoding >50,000 arthropods (Yeo, Srivathsan et al. 2020). We use 10-15
201 μ L HotSHOT per specimen. Small specimens are submerged within the well of a microplate
202 while larger specimens are placed head-first into the well. The tissue need not be entirely
203 submerged in HotSHOT. DNA is obtained within 20 minutes in a thermocycler via two
204 heating steps (Truett, Heeger et al. 2000). After neutralization, >20 μ l of template is available
205 for amplifying *COI* and the voucher can be recovered. Note that HotSHOT extraction leaves
206 most of the DNA in the specimen untouched and more high quality DNA can subsequently
207 be extracted from the same specimen. An alternative to obtaining DNA via lab buffers is the
208 use of commercial DNA extraction buffers (Kranzfelder, Ekrem et al. 2016). These buffers
209 have a longer shelf life, and are good alternatives for users who only occasionally barcode
210 moderate numbers of specimens. In the past, we have used QuickExtract (Srivathsan,
211 Hartop et al. 2019) and found that 10 μ l is sufficient for obtaining DNA template from most
212 insect specimens. In summary, obtaining DNA templates for barcoding is fast and
213 straightforward and most published barcoding studies greatly overcomplicate this step. It
214 should be noted however, that all DNA extraction methods require the removal of excess
215 ethanol from specimens prior to extraction (e.g., by placing the specimen on tissue paper or
216 replacing ethanol with water prior to specimen processing) and that the DNA extracts
217 obtained with such methods should be stored at -20°C and be used within days.

218

219 Methods for step 2: amplifying *COI* via PCR. Most PCR recipes and reagents are optimized
220 to work for a wide variety of genes and not just for a gene like the *COI* barcode that is
221 naturally enriched, has a large number of known primers, and is fairly short. Standard PCR
222 recipes can therefore be simplified. However, the use of “modern” sequencing technologies
223 such as Illumina, PacBio, or Oxford Nanopore Technologies introduces one complication:

224 The amplicons have to be “tagged” (or “indexed”/“barcoded”). This is necessary because
225 pools of amplicons are sequenced simultaneously instead of processing one amplicon at a
226 time (as in Sanger sequencing). Tags are specimen identifiers consisting of short DNA
227 sequences at the 5’ ends of the amplicons. They allow for the assignment of each sequence
228 read obtained during sequencing to a specific specimen in the “demultiplexing”
229 bioinformatics step. Numerous tagging techniques have been described in the literature, but
230 most are too complicated for efficient DNA barcoding.

231

232 *Simplified techniques for obtaining tagged amplicons*

233 Published protocols tend to have five issues that increase workload and/or inflate cost, while
234 a fifth issue only affects amplicon tagging:

- 235 • *Issue 1: expensive polymerases or master mixes.*

236 These often utilize high-fidelity polymerases that are designed for amplifying low copy-
237 number nuclear genes based on low-concentration template but rarely make a difference
238 when amplifying *COI*. Indeed, even home-made polymerases can be used for barcoding.
239 This is important because high import taxes for consumables interfere with biodiversity
240 discovery in many biodiverse countries.

- 241 • *Issue 2: indiscriminate use of single-use consumables.*

242 Disposable products increase costs and damage the environment. Most biodiversity
243 samples are obtained under “unclean conditions” that create numerous opportunities for
244 cross-specimen contamination long before specimens reach the lab (e.g., thousands of
245 specimens rubbing against each other in sample containers and in the same
246 preservation fluid). Yet numerous studies have shown that the DNA from specimens
247 exposed to such conditions will usually outcompete contaminant DNA that is likely to
248 occur at much lower concentrations. Similarly, the probability that a washed/flushed and
249 autoclaved microplates or pipette tips retain enough viable contaminant DNA to
250 successfully outcompete the template DNA is extremely low. Indeed, we have repeatedly
251 tried and failed to amplify *COI* using reused plastic consumables and water as template.

252 That it is safe to reuse some consumables is again good news for biodiversity discovery
253 under severe financial constraints. Note, however, that we do not recommend the re-use
254 of consumables for handling stock chemicals such as primers and sequencing reagents.

255 • *Issue 3: large PCR volumes (25-50 μ l).*

256 Pools of tagged amplicons comprise hundreds or thousands of products and there is
257 typically more than enough DNA for preparing a library. Accordingly, even small PCR
258 volumes of 10-15 μ l are sufficient, thereby reducing consumable costs for PCR to nearly
259 half when compared to standard volumes of 25-50 μ l.

260 • *Issue 4: using gel electrophoresis for checking amplification success of each PCR*
261 *product.*

262 This time-consuming step is only justified when Sanger sequencing is used or when
263 high-priority specimens are barcoded. It is not necessary when barcoding large numbers
264 of specimens with modern sequencing technologies, because failed amplicons do not
265 add to the sequencing cost. Furthermore, specimens that failed to yield barcodes during
266 the first sequencing run can be re-sequenced or re-amplified and then added to
267 subsequent sequencing runs (Srivathsan, Hartop et al. 2019). We thus only use gel
268 electrophoresis to check a small number of reactions per microplate (N=8-12, including
269 the negative control) in order to make sure that there was no plate-wide failure.

270

271 The fifth issue requires more elaboration and concerns how to efficiently tag amplicons. We
272 tag via a single PCR reaction (Meier, Wong et al. 2016) using primers including the tag at
273 the 5' end because it is simpler than the dual-PCR tagging strategy dominating the literature.
274 The latter has numerous disadvantages when applied to one gene: it doubles the cost by
275 requiring two rounds of PCR, is more labour intensive, increases the risk for PCR errors by
276 requiring more cycles, and requires clean-up of every PCR product after the first round of
277 amplification. In contrast, tagging via a single PCR is simple and costs the same as any
278 gene amplification. It is here described for a microplate with 96 templates, but the protocol
279 can be adapted to the use of strip tubes or half-plates. What is needed is a 96-well primer

280 plate where each well contains a differently tagged reverse primer. This “primer plate” can
281 yield 96 unique combinations of primers once the 96 reverse primers are combined with the
282 one forward primer (f-primer x 96 differently tagged r-primers = 96 unique combinations).
283 This also means that if one purchases 105 differently tagged forward primers, one can
284 individually tag 10,800 specimens ($105 \times 96 = 10,800$ amplicons). This is the number of
285 amplicons that we consider appropriate for a MinION flow cell (R10.3; see below).

286

287 Assigning tag combinations is also straightforward. For each plate with 96 PCR reactions,
288 add the same f-primer to a tube with the PCR master mix (Taq DNA polymerase, buffer and
289 dNTPs) for the plate. Then dispense the “f-primed” master mix into the 96-wells. Afterwards,
290 use a multichannel pipette to add the DNA template and the tagged r-primers from the r-
291 primer plate into the PCR plate. All 96 samples in the plate now have a unique combination
292 of tagged primers because they only share the same tagged forward primer. This makes the
293 tracking of tag combinations simple because each PCR plate has its own tagged f-primer,
294 while the r-primer is consistently tied to well position. Each plate has a negative control to
295 ensure that no widespread contamination has occurred. The tagging information for each
296 plate is recorded in the demultiplexing file that is later used to demultiplex the reads obtained
297 during sequencing.

298

299 Some users may worry that the purchase of many primers is expensive, but one must keep
300 in mind that the amount of primer used in a PCR reaction is constant. Therefore, single
301 PCR-tagging only means a greater upfront investment. Ordering all primers at once,
302 however, does mean that one must be much more careful about avoiding primer
303 degeneration and contamination as the stock will last longer. Primer stock should be stored
304 at -80°C and the number of freeze-thaw cycles should be kept low (<10). This means that
305 upon receipt of the primer stock, it should be immediately aliquoted into plates/tubes holding
306 only enough primer for rapid use. For fieldwork, one should only bring enough dissolved
307 primer for the necessary experiments, or rely on lyophilised reagents.

308

309 The choice of tag length is determined by three factors. Longer tags reduce PCR success
310 rates (Srivathsan, Hartop et al. 2019) while they increase the proportion of reads that can be
311 assigned to a specific specimen (demultiplexing rate). Designing tags is not straightforward
312 because they must remain sufficiently distinct (>4bp from each other including
313 insertions/deletions) while avoiding homopolymers. We include a list of 13 bp tags that are
314 suitable in supplementary materials.

315

316 Methods for step 3: Amplicon sequencing. The use of the PCR techniques described so far
317 should keep the cost for a tagged barcode amplicon to 0.05-0.10 USD as long as the user
318 buys cost-effective consumables. What comes next is the purification of the amplicons via
319 the removal of unused PCR reagents, the adjustment of DNA concentration, and
320 sequencing. Sequencing can be done with Sanger sequencing, Oxford Nanopore
321 Technologies (ONT) (e.g., MinION: Srivathsan, Hartop et al. 2019), Illumina (Wang,
322 Srivathsan et al. 2018), or PacBio (e.g., Sequel: Hebert, Braukmann et al. 2018). Users
323 select the sequencing option that best suit their needs based on five major criteria: (1)
324 Scaling; i.e., ability to accommodate projects of different sizes, (2) turnaround times, (3) cost,
325 (4) amplicon length and (5) sequencing error rate. Sanger sequencing has fast turnaround
326 times but high sequencing cost per amplicon (\$3-4 USD). This is the only method where cost
327 scales linearly with the number of amplicons that need sequencing, while the other
328 sequencing techniques are fundamentally different in that each run has two fixed costs that
329 stay the same regardless of whether only a few or the maximum number of amplicons for the
330 respective flow cells are sequenced. The fixed costs are library preparation (preparing
331 amplicons for sequencing) and flow cell.

332

333 The MinION Flongle has the lowest fixed costs (library and flow cell: ca. \$120 USD for a
334 routine run) and we show here that it has sufficient capacity for ca. 250 barcodes. The
335 adapter for flongle can be purchased as part of a package that costs \$1460 USD, which also

336 includes 12 flow cells thus increasing the cost of the first 12 runs to ~\$170 USD. The
337 turnaround time is fast, so the MinION Flongle is arguably the best sequencing option for
338 small barcoding projects with > 50 barcodes. Full MinION flow cells also have fast
339 turnaround times, but the minimum run cost is ca. 1000 USD, so this option only becomes
340 more cost-effective than Flongle when >1800 amplicons are sequenced. As shown later, one
341 regular MinION flow cell can comfortably sequence 10,000 amplicons. This is a similar
342 volume to what has been described for PacBio (Sequel) (Hebert, Braukmann et al. 2018),
343 but the high instrument cost for PacBio means that sequencing usually has to be
344 outsourced, leading to longer wait times. By far the most cost-effective sequencing method
345 for barcodes is Illumina's NovaSeq sequencing. The fixed costs for library and lanes are high
346 (3000-4000 USD), but each flow cell yields 800 million reads which can comfortably
347 sequence 800,000 barcodes at a cost of < \$0.01 USD per barcode. This high capacity
348 means that the 6 million publicly available barcodes in BOLD Systems could have been
349 sequenced on just <8 NovaSeq flow cells for ~50,000 USD. However, Illumina sequencing
350 can only be used for mini-barcodes of up to 420 bp length (using 250bp PE sequencing
351 using SP flow cell). "Full-length" COI barcode (658 bp) can only be obtained by sequencing
352 two amplicons. Note that while Illumina barcodes are shorter than "full-length" barcodes,
353 there is no evidence that mini-barcodes have a negative impact on species delimitation or
354 identification as long as the mini-barcode is >250 bp in length (Yeo, Srivathsan et al. 2020).

355

356 *Simplified techniques for sequencing tagged amplicons:* Modern sequencing technologies
357 are used to sequence amplicon pools instead of individual amplicons. To obtain such a pool,
358 it is sufficient to combine only 1 µl per PCR product. The pool can be cleaned using several
359 PCR clean-up methods. We generally use SPRI bead-based clean-up, with Ampure
360 (Beckman Coulter) beads but Kapa beads (Roche) or the more cost-effective Sera-Mag
361 beads (GE Healthcare Life Sciences) in PEG (Rohland and Reich 2012) are also viable
362 options (Srivathsan, Hartop et al. 2019). We recommend the use of a 0.5X ratio for Ampure
363 beads for barcodes longer than 300 bp since it removes a larger proportion of primers and

364 primer dimers. However, this ratio is only suitable if yield is not a concern (e.g., pools
365 consisting of many and/or high concentration amplicons). Increasing the ratio to 0.7-1X will
366 improve yield but render the clean-up less effective. Amplicon pools containing large
367 numbers of amplicons usually require multiple rounds of clean-up, but only a small subset of
368 the entire pool needs to be purified because most library preparation kits require only small
369 amounts of DNA. Note that the success of the clean-up procedures should be verified with
370 gel electrophoresis, which should yield only one strong band of expected length. After the
371 clean-up, the pooled DNA concentration is measured in order to use an appropriate amount
372 of DNA for library preparation. Most laboratories use a Qubit, but less precise techniques are
373 probably also suitable.

374

375 Obtaining a cleaned amplicon pool according to the outlined protocol is not time consuming.
376 However, many studies retain “old Sanger sequencing habits”. For example, they use gel
377 electrophoresis for each PCR reaction to test whether an amplicon has been obtained and
378 then clean and measure all amplicons one at a time for normalization (often with very
379 expensive techniques: Ampure beads: (Maestri, Cosentino et al. 2019); TapeStation,
380 BioAnalyzer, Qubit: (Seah, Lim et al. 2020)). This is presumably done to obtain a pool of
381 amplicons where each has equal representation. However, reads are cheap while individual
382 clean-ups and measurements for each PCR product are expensive. Furthermore, weak
383 products that failed to yield a barcode can be re-sequenced (Srivathsan, Hartop et al 2019).
384 Another strategy is normalization based on the band strength of a few PCR products per
385 plate as determined by gel electrophoresis. We used this strategy in the current study, but
386 have since determined that pooling without such normalization yields nearly identical
387 success rates based on the same number of reads (>9600 amplicons: 76.8% vs. 77.2%).

388

389 **3. MinION barcoding with new flow cells (R10.3, Flongle) and high-accuracy**

390 **basecalling**

391 Oxford Nanopore Technologies (ONT) instruments sequence DNA by passing single-
392 stranded DNA through a nanopore. This creates current fluctuations which can be measured
393 and translated into a DNA sequence via basecalling (Wick 2019). The sequencing devices
394 are small and inexpensive, but the read accuracy is only moderate (85% - 95%) (Wick 2019,
395 Silvestre-Ryan and Holmes 2021). This means that data analysis requires specialized
396 bioinformatics pipelines. The nanopores used for sequencing are arranged on flow cells, with
397 new flow cell chemistries and basecalling softwares regularly released. Recently, three
398 significant changes occurred. Firstly, ONT released a cheap flow cell (Flongle) that only has
399 126 pores (126 channels) instead of the customary 2048 pores (512 channels) of a full
400 MinION flow cell. We were interested in Flongle because it looked promising for small
401 barcoding projects that needed quick turnaround times. Secondly, ONT released a new flow
402 cell chemistry for full flowcells (R10.3) where the nanopores have a dual instead of a single
403 reader-head. Dual reading has altered the read error profile by giving better resolution to
404 homopolymers and improving consensus accuracy (Chang, Ip et al. 2020, Vereecke, Bokma
405 et al. 2020). Lastly, ONT released high accuracy (HAC) basecalling. We thus obtained
406 amplicons using techniques described in Section 2 and processed them further as described
407 below.

408

409 Library preparation. Library preparation was based on 200 ng of DNA for the full MinION flow
410 cells and 100 ng for the Flongle and used ligation-based kits(see Table 1 for details). We
411 generally followed kit instructions, but excluded the FFPE DNA repair mix in the end-repair
412 reaction, as this is mostly needed for formalin-fixed, paraffin-embedded samples. The
413 reaction volumes for the R10.3 flow cell libraries consisted of 45 µl of DNA, 7 µl of Ultra II
414 End-prep reaction buffer (New England Biolabs), 3 µl of Ultra II End Prep enzyme mix (New
415 England Biolabs) and 5 µl of molecular grade water. For the Flongle, only half of the
416 reagents were used to obtain a total volume of 30 µl. We further modified the Ampure ratio
417 to 1x for all steps as DNA barcodes are short whereas the recommended ratio in the manual
418 is for longer DNA fragments. The libraries were loaded and sequenced with a MinION Mk

419 1B. Data capture involved a MinIT or a Macintosh computer that meets the IT specifications
420 recommended by ONT. The bases were called using Guppy (versions provided in Table 2),
421 under the high-accuracy model in MinIT taking advantage of its GPU.

422

423 Sequencing. Six amplicon pools were sequenced (Table 1). For two of the pools, *Mixed*
424 *Diptera* (N=511) and *Afrotropical* (N=4,275) *Phoridae*, we had comparison barcodes that
425 were obtained with Sanger and Illumina sequencing and the same amplicon pools were
426 previously sequenced with earlier versions of MinION flow cells (Srivathsan, Baloglu et al.
427 2018, Srivathsan, Hartop et al. 2019) (Table 1). These two pools were used to assess the
428 accuracy of barcodes generated using R10.3 flow cells. Two additional datasets tested the
429 capacity of R10.3 flowcells for mini- and full length barcodes for the same specimens
430 (*Palearctic Phoridae*, 658 and 313 bp for ca. 9,930 specimens). The *Mixed Diptera*
431 *Subsample* and *Chironomidae* datasets test the performance of the Flongle. The *Mixed*
432 *Diptera Subsample* (N=257) is a subset of the *Mixed Diptera* (N=511) amplicon pool for
433 which we have Sanger barcodes for comparison. The *Chironomidae* dataset contains
434 sequences for 313 bp mini-barcodes for 191 specimens of Chironomidae that were newly
435 amplified for this study.

436

437 **Table 1.** Datasets used in the study and the corresponding experimental details.

Dataset Name	Number of specimens	Fragment size, primer information	Extraction/PCR setup	PCR cleanup	ONT Library Preparation kit/Flow cell used
R10.3 Datasets					
<u>Mixed Diptera</u> (see Srivathsan et al., 2018) - Sanger barcodes available	511 (257 mixed Diptera, 254 Dolichopodidae) 17 negatives	658 bp HCO2198, LCO1490 (Folmer et al., 1994)	Extraction Method: QuickExtract PCR Mix: Total volume: 20 µl 10x buffer: 2 µl dNTPs (2.5 mM): 1.5 µl Taq polymerase: 0.2 µl BSA (1 mg/ml): 2 µl Primer (5 µM): 2 µl each DNA: 2 µl	Ampure beads (Beckman Coulter)	SQK-LSK110/FLO-MIN111
<u>Afrotropical Phoridae</u> (see Srivathsan et al., 2019) - Illumina mini-barcode available	4275 (Phoridae) 45 negatives	658 bp HCO2198, LCO1490 (Folmer et al., 1994)	Extraction Method: QuickExtract PCR Mix: Total volume: 15.16 µl Mastermix (CWBio): 10 µl 25mM MgCl ₂ : 0.16 µl BSA (1 mg/ml): 2 µl Primer (10µM): 1 µl each DNA: 1 µl	Sera-Mag beads (GE Healthcare Life Sciences) in PEG	SQK-LSK109/FLO-MIN111
<u>Palaeartic Phoridae (658)</u>	9,929 (Phoridae) 105 negatives	658 bp jgHCO2198, LCO1490 (Folmer et al., 1994, Geller et al. 2013)	Extraction Method: HotSHOT PCR Mix: Total volume: 16 µl Mastermix (CWBio): 7 µl BSA (1mg/ml): 1 µl Primer (10µM): 1 µl each DNA: 6 µl	Ampure beads (Beckman Coulter)	SQK-LSK110/FLO-MIN111
<u>Palaeartic Phoridae (313)</u>	9,932 (Phoridae) 106 negatives	313 bp m1COLintF, jgHCO2198 (Leray et al. 2013, Geller et al. 2013)	Extraction Method: HotSHOT PCR Mix: Total volume: 14 µl Mastermix (CWBio): 7 µl BSA (1mg/ml): 1 µl Primer (10µM): 1 µl each DNA: 4 µl	Ampure beads (Beckman Coulter)	SQK-LSK110/FLO-MIN111
Flongle Datasets					
<u>Mixed Diptera subsample</u> (see Srivathsan et al., 2018) - Sanger barcodes available	257 7 negatives	See "Mixed Diptera" entry for R10.3	See "Mixed Diptera" entry for R10.3	Ampure beads (Beckman Coulter)	SQK-LSK109/Flongle
<u>Chironomidae</u>	191 (Chironomidae) 1 negative	313 bp m1COLintF, jgHCO2198 (Leray et al. 2013, Geller et al. 2013)	Extraction Method: HotSHOT PCR Mix: Total volume: 14 µl Mastermix (CWBio): 7 µl BSA (1 mg/ml): 1 µl Primer (10µM): 1 µl each DNA: 4 µl	Ampure beads (Beckman Coulter)	SQK-LSK109/Flongle

438

439

440

441 Bioinformatics

442 One of the most significant barriers to widespread barcoding with MinION is the high error
443 rates of ONT reads. In 2018, we developed a bioinformatics pipeline for error correction that
444 was too complex for the average user (Srivathsan, Baloglu et al. 2018, Srivathsan, Hartop et
445 al. 2019). After obtaining data with several R10.3 and new R9.4 flow cells, we initially applied
446 this pipeline (Srivathsan et al. 2019), but we noticed major improvements in terms of MinION
447 read quality and the total number of raw and demultiplexed reads produced by each flow
448 cell. This led to the development of a new user-friendly pipeline after considering alternative,
449 published pipelines which faced one or several of the following problems: they required high
450 read coverage, relied on external sequences, were complex, and/or needed several
451 command line steps, and included external dependencies that limit cross platform
452 compatibility (Menegon, Cantaloni et al. 2017, Maestri, Cosetino et al. 2019, Seah, Lim et al.
453 2020, Sahlin, Lim et al. 2021). We here present “ONTbarcoder”, which has a graphical user
454 interface (GUI) and is suitable for all major operating systems (Linux, Mac OS, Windows10).
455 Both are requirements for the democratization of barcoding with MinION. In addition, we
456 prepared a simple video tutorial
457 (<https://www.youtube.com/channel/UC1WowokomhQJRc71FmsUAcg>).

458

459 ONTbarcoder. ONTbarcoder (available at: <https://github.com/asrivathsan/ONTbarcoder>) has
460 three modules. (a) The first is a demultiplexing module which assigns reads to specimen-
461 specific bins. (b) The second is a barcode calling module which reconstructs the barcodes
462 based on the reads in each specimen bin. (c) The third is a barcode comparison module that
463 allows for comparing barcodes obtained via different software and software settings.

464

465 a. Demultiplexing. The user provides three pieces of information and two files: (1) primer
466 sequence, (2) expected fragment length, and (3) demultiplexing information (=tag
467 combination for each specimen). The latter is summarized in a demultiplexing file (see
468 supplementary information for format). The only other required file is the FASTQ file

469 obtained from MinKNOW/Guppy after basecalling. Demultiplexing by ONTbarcode starts by
470 analyzing the read length distribution in the FASTQ file. Only those reads that meet the read
471 length threshold are demultiplexed (default= 658 bp corresponding to metazoan COI
472 barcode). Technically, the threshold should be the amplicon length plus the length of both
473 tagged primers, but ONT reads have indel errors such that they are occasionally too short
474 and we therefore advise to specify the amplicon length as threshold. Reads that are twice
475 the expected fragment length are split into two parts. Splitting is based on the user given
476 fragment size, primer and tag lengths, and a window size to account for indel errors
477 (default=100 bp).

478

479 Once all reads suitable for demultiplexing have been identified, ONTbarcode finds the
480 primers via sequence alignment of the primer sequence to the reads (using python library
481 *edlib*). Up to 10 deviations from the primer sequence are allowed because this step is only
482 needed for determining the primer location and orientation within the read. For
483 demultiplexing, the flanking region of the primer sequence is retrieved whereby the number
484 of retrieved bases is equal to the user-specified tag length. The flanking sequences are then
485 matched against the tags from the user-provided tag combinations (demultiplexing file). In
486 order to account for sequencing errors, not only exact matches are accepted, but also
487 matches to “tag variants” that differ by up to 2 bps from the original tag
488 (substitutions/insertions/deletions). Note that accepting tag variants does not lead to
489 demultiplexing error because all tags differ by >4 bp. All reads thus identified as belonging to
490 the same specimen are pooled into the same bin. To increase efficiency, demultiplexing is
491 parallelized and the search space for primers and tags are restricted to user-specified parts
492 of each read.

493

494 **b. Barcode calling:** Barcode calling uses the reads within each specimen-specific bin to
495 reconstruct the barcode sequence. The reads are aligned to each other and a consensus
496 sequence is called. Barcode calling is done in three phases: “Consensus by Length”,

497 “Consensus by Similarity” and “Consensus by barcode comparison”. The user can opt to
498 only use some of these methods.

499

500 “Consensus by Length” is the main barcode calling mode. Alignment must be efficient in
501 order to obtain high-quality barcodes at reasonable speed for thousands of amplicons.
502 ONTbarcoder delivers speed by using an iterative approach that gradually increases the
503 number of reads (“coverage”) that is used during alignment. However, reconstructing
504 barcodes based on few reads could lead to errors which are here weeded out by using four
505 Quality Control (QC) criteria. The first three QC criteria are applied immediately after the
506 consensus sequence has been called: (1) the barcode must be translatable, (2) it has to
507 match the user-specified barcode length, and (3) the barcode has to be free of ambiguous
508 bases (“N”). To increase the chance of finding a barcode that meets all three criteria, we
509 subsample the reads in each bin by read length (thus the name “Consensus by Length”);
510 i.e., initially only those reads closest to the expected length of the barcode are used. For
511 example, if the user specified coverage=25x for a 658 bp barcode, ONTbarcoder would only
512 use the 25 reads that have the closest match to 658 bp. The fourth QC measure is only
513 applied to barcodes that have already met the first three QC criteria. A multiple sequence
514 alignment (MSA) is built for the barcodes obtained from the amplicon pool, and any barcode
515 that causes the insertion of gaps in the MSA is rejected. Note that if the user suspects that
516 barcodes of different length are in the amplicon pool, the initial analysis should use the
517 dominant barcode length. The remaining barcodes can then be recovered by re-analyzing all
518 data or only the failed read bins (“remaining”, see below) and bins that yielded barcodes that
519 had to be “fixed”. These bins can be reanalyzed using a different pre-set barcode length.

520

521 “Consensus by Similarity”. The barcodes that failed the QC during the “Consensus by
522 Length” stage are often close to the expected length and have few ambiguous bases, and/or
523 cause few gaps in the MSA. These “preliminary barcodes” can be improved through
524 “Consensus by Similarity”. This method eliminates outlier reads from the read alignments.

525 Such reads often differ considerably from the signal of the consensus barcode and
526 ONTbarcoder identifies them by sorting all reads by similarity to the preliminary barcode.
527 Only the top 100 reads (this default can be changed) that differ by <10% from the
528 preliminary barcode are retained and used for calling the barcodes again using the same
529 techniques described previously (including the same QC criteria). This distance threshold
530 accounts for errors generated by MinION but excludes highly erroneous or contaminating
531 reads. This improvement step converts many preliminary barcodes found during “Consensus
532 by Length” into barcodes that pass all four QC criteria by filling/removing indels or resolving
533 an ambiguous base.

534

535 “Consensus by barcode comparison”. The remaining preliminary barcodes that still failed to
536 convert into QC-compliant barcodes tend to be based on read bins with low coverage, but
537 some can yield good barcodes after subjecting them to a further improvement step that fixes
538 the remaining errors. ONTbarcoder identifies these errors by finding the 20 most similar QC-
539 compliant barcodes that have already been reconstructed for the other amplicons. The 21
540 sequences are aligned and ONTbarcoder finds the errors because they cause insertions and
541 deletions in the MSA. Insertions are deleted, gaps are filled with ambiguous bases (“N”), but
542 mismatches are retained. The number and kinds of “fixes” are recorded and added to the
543 FASTA header of the barcode.

544

545 Output. ONTbarcoder extensively documents the barcoding results so that users can check
546 the output and potentially modify the barcode calling parameters. For example, it produces a
547 summary table (Outputtable.csv) and FASTA files that contain the different classes of
548 barcodes. Each barcode header contains information on coverage used for barcode calling,
549 coverage of the specimen bin, length of the barcode, number of ambiguities and number of
550 indels fixed. Five sets of barcodes are provided, here discussed in the order of barcode
551 quality: (1) “QC_compliant”: The barcodes in this set satisfy all four QC criteria without
552 correction and are the highest quality barcodes. (2) “Filtered_barcodes”: this file contains the

553 barcodes that are translatable, have <1% ambiguities and have up to 5 indels fixed during
554 the last step of the bioinformatics pipeline. These filtering thresholds were calibrated based
555 on two datasets for which we have Sanger/Illumina barcodes and the resulting barcodes are
556 found to be highly accurate. Note that the file with filtered barcodes also includes the
557 QC_compliant barcodes and that all results discussed in this manuscript are based on
558 filtered barcodes given that they are of much higher quality than the average barcode in
559 BOLDSystems (assessment in Srivathsan, Baloglu et al 2018).

560

561 The remaining files include barcodes of lesser and/or suspect quality. (3)

562 “Fixed_barcodes_XtoY”: these files contain barcodes that had indel errors fixed and are
563 grouped by the number of errors fixed. Only the barcodes with 1-5 errors overlap with
564 Filtered barcodes file, if they have <1% ambiguities. (4) “Allbarcodes”: this file contains all
565 barcodes in sets (1)-(3). (5) “Remaining”: these are barcodes that fail to either translate or
566 are not of predicted length. Note that all barcodes should be checked via BLAST against
567 comprehensive databases in order to detect contamination. There are several online tools
568 available for this and we recommend the use of GBIF sequence ID tool
569 (<https://www.gbif.org/tools/sequence-id>) which gives straightforward output including a
570 taxonomic summary.

571

572 The output folder also includes the FASTA files that were used for alignment and barcode
573 calling. The raw read bins are in the “demultiplexed” folder, while the resampled bins (by
574 length, coverage, and similarity) are in their respective subfolders named after the search
575 step. Note that the raw reads are encoded to contain information on the orientation of the
576 sequence and thus cannot be directly used in other software without modifications (see
577 ONTbarcoder manual on Github). Lastly, for each barcode FASTA file (1-5), there are
578 folders with the files that were used to call the barcodes. This means that the user can, for
579 example, reanalyze those bins that yielded barcodes with high numbers of ambiguous

580 bases. Lastly a “runsummary.xlsx” document allows the user to explore the details of the
581 barcodes obtained at every step of the pipeline.

582 Algorithms. ONTbarcoder uses the following published algorithms. All alignments utilize
583 MAFFTv7 (Kato and Standley 2013). The MSAs that use MinION reads to form a
584 consensus barcode are constructed in an approach similar to lamassemble (Frith,
585 Mitsuhashi et al. 2020), using parameters optimized for nanopore data by “last-train”
586 (Hamada, Ono et al. 2017) which accounts for strand specific error biases. The MAFFT
587 parameters can be modified in the “parfile” supplied with the software which will help with
588 adjusting the values given the rapidly changing nanopore technology. All remaining MSAs in
589 the pipeline (e.g., of preliminary barcodes) use MAFFT’s default settings. All read and
590 sequence similarities are determined with the *edlib* python library under the Needle-Wunsch
591 (“NW”) setting, while primer search is using the infix options (“HW”). All consensus
592 sequences are called from within the software. This is initially done based on a minimum
593 frequency of 0.3 for each position. This threshold was empirically determined based on
594 datasets where MinION barcodes can be compared to Sanger/Illumina barcodes. The
595 threshold is applied as follows. All sites where >70% of the reads have a gap are deleted.
596 For the remaining sites, ONTbarcoder accepts those consensus bases that are found in at
597 least >30% of the reads. If no base/multiple bases reach this threshold, an “N” is inserted.
598 To avoid reliance on a single threshold, ONTbarcoder allows the user to change the
599 consensus calling threshold from 0.2 to 0.5 for all barcodes that fail the QC criteria at 0.3
600 frequency. However, barcodes called at different frequencies are only accepted if they pass
601 the first three QC criteria and are identical. If no such barcode is found, the 0.3 frequency
602 consensus barcode is used for further processing.

603

604 c. Barcode comparison. Many users may want to call their barcodes under different settings
605 and then compare barcode sets. The ONTbarcoder GUI simplifies such comparisons. A set
606 of barcodes is dragged into the window and the user can select a barcode set as the
607 reference. The barcode comparisons are conducted using *edlib* library. The barcodes in the

608 sets are compared and classified into three categories: “identical” where sequences are a
609 perfect match and lack ambiguities, “compatible” where the sequences only differ by
610 ambiguities, and “incorrect” where the sequences differ by at least one base pair. Several
611 output files are provided. A summary sheet, a FASTA file each for “identical”, “compatible”,
612 and the sequences only found in one dataset. Lastly, there is a folder with FASTA files
613 containing the different barcodes for each incompatible set of sequences. This module can
614 be used for either comparing set(s) of barcodes to reference sequences, or for comparing
615 barcode sets against each other. It furthermore allows for pairwise comparisons and
616 comparisons of multiple sets in an all-vs-all manner. This module was used here to get the
617 final accuracy values presented in Table 3.

618

619 **4. Performance of flow cells (R10.3, Flongle) and high-accuracy basecalling**

620 The pools used to test the new ONT products contained amplicons for 191 - 9,932
621 specimens and were run for 15-49 hours (Table 2). The fast5 files were basecalled using
622 Guppy in MinIT under the high accuracy (HAC) model. Basecalling large datasets under
623 HAC is currently still very slow and took 12 days in MinIT for the *Palaeartic Phoridae* (658
624 bp) dataset (Table 2) but the reads yielded high demultiplexing rates for three of the four
625 R10.3 MinION datasets (= 30-49%). The exception was the *Palaeartic Phoridae* (313 bp)
626 dataset (15.5%). Flongle datasets showed overall also lower demultiplexing rates (17-21%).
627

628 **Table 2.** Datasets generated in this study and the results of barcoding using ONTbarcoder at
 629 200X coverage (Consensus by Length) and 100X coverage (Consensus by Similarity).
 630

<u>Dataset Name</u>	<u>Flow cell details</u> <u>Run time/Guppy</u> <u>version</u>	<u>Raw reads/reads passing</u> <u>length threshold/reads of</u> <u>suitable length/</u> <u>demultiplexed</u>	<u>Demultiplexing</u> <u>rate/#</u> <u>QC compliant</u> <u>barcodes /#</u> <u>Filtered barcodes</u> <u>with 1N/# Filtered</u> <u>barcodes with >1N</u> <u>/# Unreliable</u> <u>barcodes</u>
MinION R10.3 Datasets			
Mixed Diptera (658 bp, N=511)	R10.3: reused flow cell: 71 pores according to QC, but 500+ active during run Runtime: 27.5 hrs Guppy: 4.2.3+f90bd04	3,864,000/3,425,357/3,560 ,389/1,544,758	43.39%/495/2/5/8 Total success rate= 502/511 (98.2%)
Afrotropical Phoridae (658 bp, N=4,275)	R10.3: new flow cell: QC: 1,101 pores Runtime: 49.5 hrs Guppy: 4.0.11+f1071ce	6,838,903/5,465,164/5,474 ,306/2,681,029	48.97%/3,725/121/5 9/247 Total success rate= 3905/4275 (91.3%)
Palaeartic Phoridae (658 bp, N=9,932)	R10.3: new flow cell: QC: 1,239 pores Runtime: 47.5 hrs Guppy: 4.2.3+f90bd04	16,595,984/15,658,174/16, 100,505/5,012,489	31.13%/8,026/108/2 31/780 Total success rate= 8,365/9,932 (84.2%)
Palaeartic Phoridae (313 bp, N=9,929)	R10.3: new flow cell: QC: 1,297 pores Runtime: 37 hrs Guppy: 4.2.3+f90bd04	13,690,869/13,221,764/10, 366,455/12,983,260/2,015, 135	15.52%/8,705/118/1 12/899 Total success rate= 8,935/9,929 (90%)
Flongle Datasets			
Mixed Diptera Subsample (658 bp, N=257)	Flongle: new QC: 81 pores Runtime: 24 hrs Guppy: v 4.0.11+f1071ce	294,896/222,189/190,952/ 33,270	17.42%/185/35/20/9 Total success rate= 240/257 (93.4%)
Chironomidae (313 bp, N=191)	Flongle: new QC: 74 pores Runtime: 15 hrs Guppy: 4.2.3+f90bd04	560,062/525,087/504,621/ 108,574	21.52%/178/1/2/6 Total success rate= 181/191 (94.8%)

631

632 We used ONTbarcoder to analyze the MinION data for all six datasets by analyzing all
 633 specimen-specific read bins at different coverages (5-200x in steps of 5x). This means that
 634 the barcodes for a bin with 27 reads were called five times at 5x, 10x, 15x, 20x, and 25x
 635 coverages while bins with >200x were analyzed 40 times at 5x increments. Instead of using
 636 conventional rarefaction via random subsampling reads, we used the first reads provided by

637 the flow cell because this accurately reflects how the data accumulated during the
638 sequencing run and how many barcodes would have been obtained if the run had been
639 stopped early. This rarefaction approach also allowed for mapping the barcode success
640 rates against either coverage or time.

641

642 In order to obtain a “best” estimate for how many barcodes can be obtained, we also carried
643 out one analysis at 200x coverage with the maximum number of “Comparison by Similarity”
644 reads set to 100. This means that ONTbarcoder selected up to 200 reads from the
645 specimen-specific read bin that had the closest match to the length of the target barcode
646 (i.e., 313 or 658 bp), then produced an MSA and consensus barcode using MAFFT. If the
647 resulting consensus barcode did not satisfy all four QC criteria, ONTbarcoder would select
648 up to 100 reads that had at least a 90% match to the preliminary barcode. These reads
649 would then be used to call another barcode with MAFFT. Only if this also failed to produce a
650 QC-compliant barcode, ONTbarcoder would “fix” the preliminary barcode using its 20 closest
651 matches in the dataset. All analyses produced a “filtered” set of barcodes (barcodes with
652 <1% Ns and up to 5 fixes) that were used for assessing the accuracy and quality via
653 comparison with Sanger and Illumina barcodes for *Mixed Diptera (MinION R10.3)*,
654 *Afrotropical Phoridae (MinION R10.3)*, and *Mixed Diptera Subsample (Flongle R9.4)*. For the
655 comparisons of the barcode sets obtained at the various coverages, we used MAFFT and
656 the `assess_corrected_barcode.py` script in miniBarcoder (Srivathsan et al., 2019).

657

658 We investigated barcode accuracy (Figure 1) by directly aligning the MinION barcodes with
659 the corresponding Sanger and Illumina barcodes. We find that MinION barcodes are virtually
660 identical to Sanger and Illumina barcodes (>99.99% identity, Table 3). We then established
661 that the number of ambiguous bases (“N”) is also very low for barcodes obtained with R10.3
662 (<0.01%). Indeed, more than 90% of all barcodes are entirely free of ambiguous bases. In
663 comparison, Flongle barcodes have a slightly higher proportion of ambiguous bases
664 (<0.06%). They are concentrated in ~20% of all sequences so that 80% of all barcodes

665 again lack Ns. This means that MinION barcodes more than just match the Consortium for
666 the Barcode of Life (CBOL) criteria for “barcode” designation with regard to length, accuracy,
667 and ambiguity.

668

669 Rarefaction at different read coverage levels reveals that 80-90% of high-quality barcodes
670 are obtained within a few hours of sequencing. In addition, the number of barcodes
671 generated by MinION exceeded or was comparable to what could be obtained with Sanger
672 or Illumina sequencing (Figure 1). We then determined the coverage needed for obtaining
673 reliable barcodes. For this purpose, we plotted the number of barcodes obtained against
674 coverage (Figure 2). This revealed that the vast majority of specimen bins yield high-quality
675 barcodes at coverages between 25x and 50x when R10.3 reads are used. Increasing
676 coverage beyond 50x leads to only modest improvements of quality and few additional
677 specimen amplicons yield new barcodes. The coverage needed for obtaining Flongle
678 barcodes is somewhat higher, but the main difference between the R9.4 technology of the
679 Flongle flow cell and R10.3 is that more barcodes retain ambiguous bases even at high
680 coverage for data from R9.4 flow cells. The differences in read quality between R9.4 and
681 R10.3 become even more obvious when the read bins for the “Mixed Diptera Subsample”
682 are analyzed based an identical numbers of R10.3 and R9.4 reads. The barcodes based on
683 Flongle and R10.3 data are compatible, but the R10.3 barcodes are ambiguity-free while
684 some of the corresponding Flongle barcodes retain 1-2 ambiguous bases.

685

686 Overall, these results imply that 100x raw read coverage is sufficient for obtaining barcodes
687 with either R10.3 or R9.4 flow cells. Given that most MinION flow cells yield >10 million
688 reads of an appropriate length, this means that one could, in principle, obtain 100,000
689 barcodes in one flow cell. However, this would require that all amplicons are represented by
690 similar numbers of copies and that all reads could be correctly demultiplexed. In reality, only
691 30-50% of the reads can be demultiplexed and the number of reads per amplicon fluctuates
692 widely (Figure 3). Very-low coverage bins tend to yield no barcodes or barcodes of lower

693 quality (errors or Ns). These low-coverage barcodes can be improved by collecting more
694 data, but this comes at a high cost and increased risk of a small number of contaminant
695 reads yielding barcodes. For example, we observed that some “negative” PCR controls
696 yielded low-quality barcodes for 4 of 106 negatives in the Palaeartic Phoridae (313 bp) and
697 1 of 105 negatives in the Palaeartic Phoridae (658 bp) datasets.

698

699 To facilitate the planning of barcode projects, we illustrate the trade-offs between barcode
700 yield, time, and amount of raw data needed for six amplicon pools (Figure 4: 191-9,932
701 specimens). These standard curves can be used to roughly estimate the amount of raw
702 reads needed to achieve a specific goal for a barcoding project of a specific size (e.g.,
703 obtaining 80% of all barcodes for a project with 1000 amplicons). For each dataset, we
704 illustrate how many reads were needed to recover a certain proportion of barcodes. The
705 number of recoverable barcodes was set to the number of all error-free, filtered barcodes
706 obtained in an analysis of all data. We would argue that this is a realistic estimate of
707 recoverable barcodes given the saturation plots in Figure 1 that suggest that most barcodes
708 with significant amounts of data have been called at 200x coverage. Note, however, that
709 Figure 4 can only provide very rough guidance on how many reads are needed because, for
710 example, the demultiplexing rates differ between flow cells and different amplicon pools have
711 very different read abundance distributions (see Figure 3).

712

713 **Table 3.** Quality assessment of barcodes generated by ONTbarcoder at 200X read
714 coverage (Consensus by Length) and 100X coverage (Consensus by Similarity). The
715 accuracy of MinION barcodes is compared with the barcodes obtained for the same
716 specimens using Illumina/Sanger sequencing. Errors are defined as sum of substitution or
717 indel errors. All denominators for calculating percentages are the total number of nucleotides
718 assessed.
719

Dataset	No. of comparison barcodes	No. of barcodes with errors/No. of errors/% identity	# of Ns/%Ns
R10.3: Mixed Diptera: Sanger barcodes available	476	2/10/99.997%	19 (0.006%)
R10.3: Afrotropical Phoridae: Illumina barcodes available*	3316	23/48/99.995%	284 (0.011%)
Flongle-Mixed Diptera Subsample: Sanger barcodes available	231	5/8/99.994%	91 (0.058%)

720

721 *5 barcodes with very high distances from reference were excluded for R10.3: Afrotropical Phoridae dataset as
722 they likely represent lab contamination (see Srivathsan, Hartop et al. (2019).

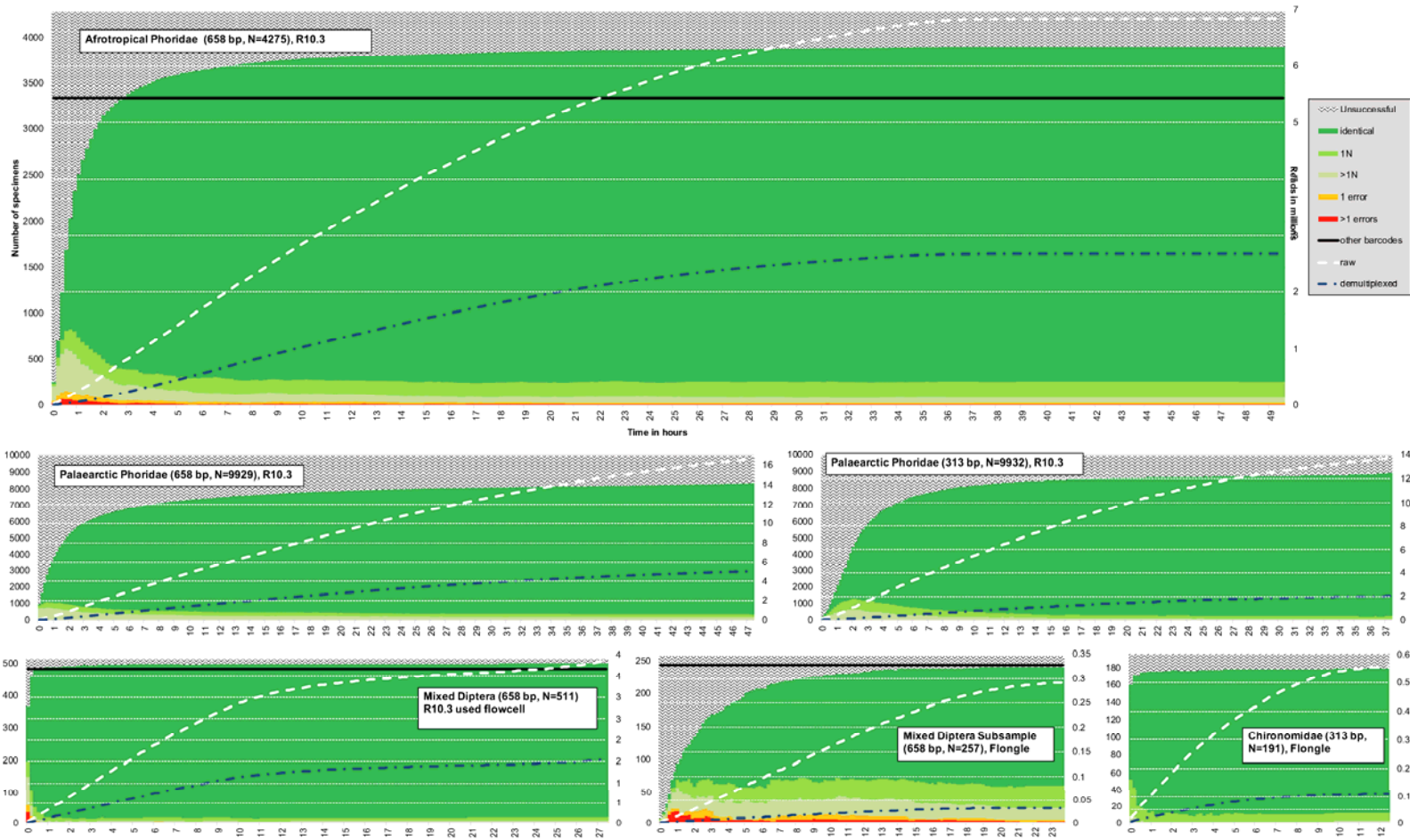


Figure 1. Rapid recovery of accurate MinION barcodes over time (in hours, x-axis) (filtered barcodes: dark green = barcodes passing all 4 QC criteria, light green = one ambiguous base; lighter green = more than 1N, no barcode = white with pattern, 1 mismatch = orange, >1 mismatch = red). The solid black line represents the number of barcodes available for comparison. White dotted line represents the amount of raw reads collected over time, blue represents number of demultiplexed reads over time (plotted against Z-axis)

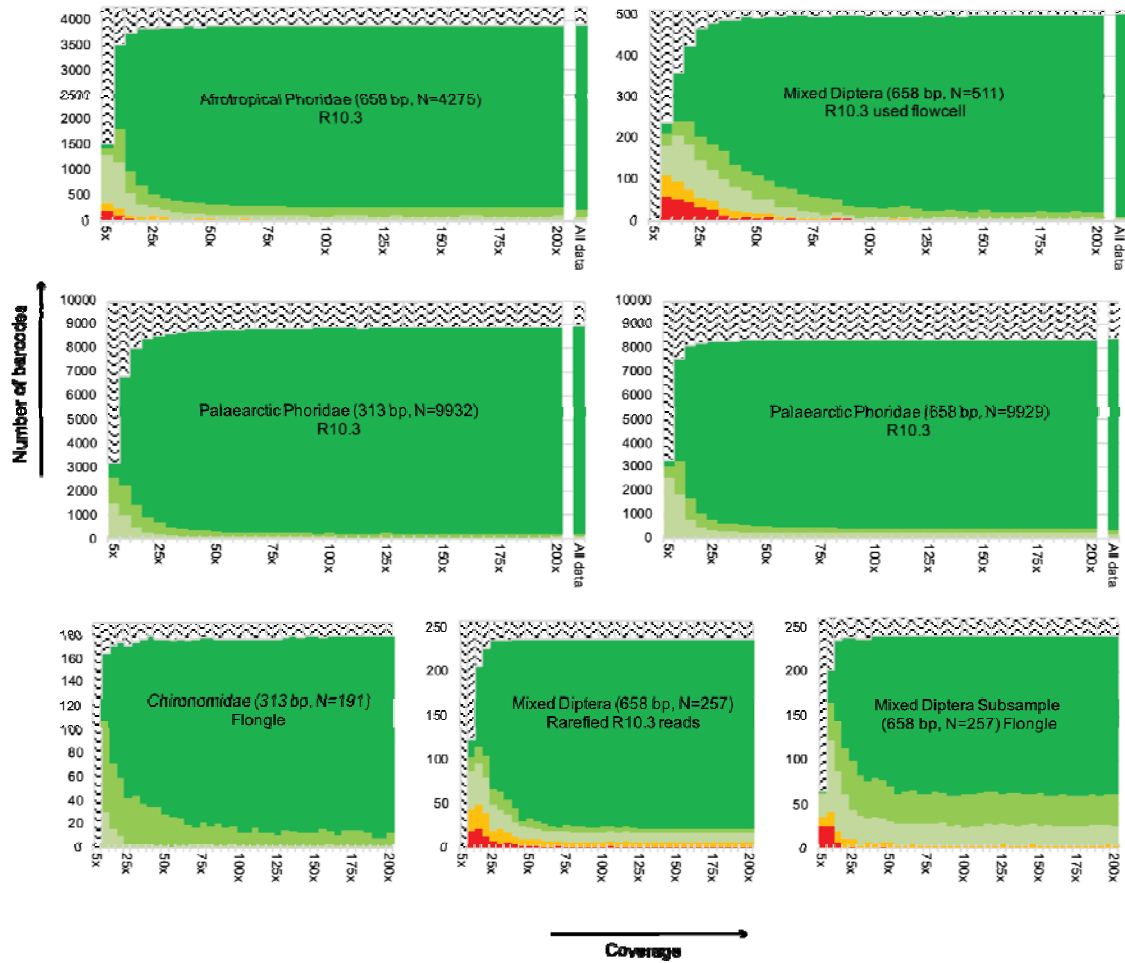


Figure 2. Relationship between barcode quality and coverage. Subsetting the data to 5-200X coverage shows that there are very minor gains to barcode quality after 25-50X coverage. (filtered barcodes: dark green = barcodes passing all 4 QC criteria, light green = one ambiguous base; lighter green = more than 1N, no barcode = white with pattern, 1 mismatch = orange, >1 mismatch = red).

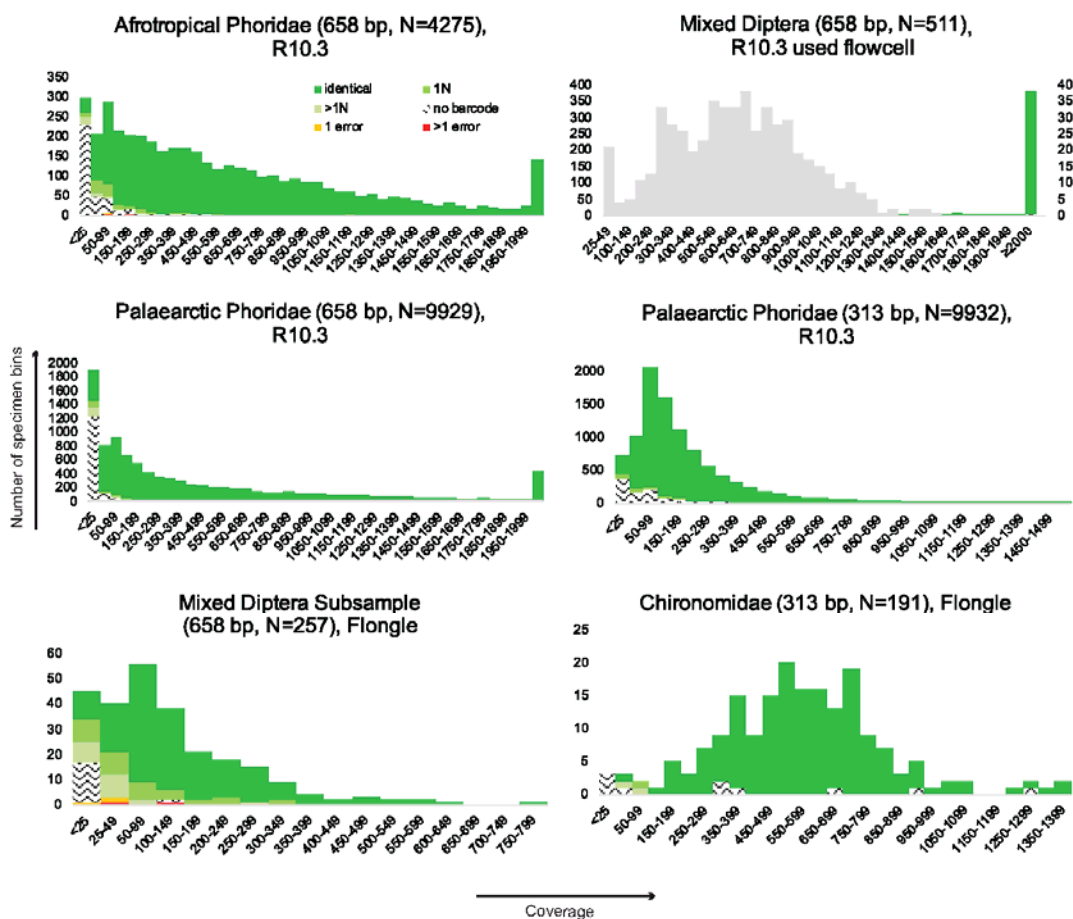


Figure 3. Read bin size distribution for six amplicon pools (color-coding as in Figs 1-2). Due to the very generous coverage for the “Mixed Diptera” dataset, we also use grey to show the bin size distribution after dividing the bin read totals by 5.

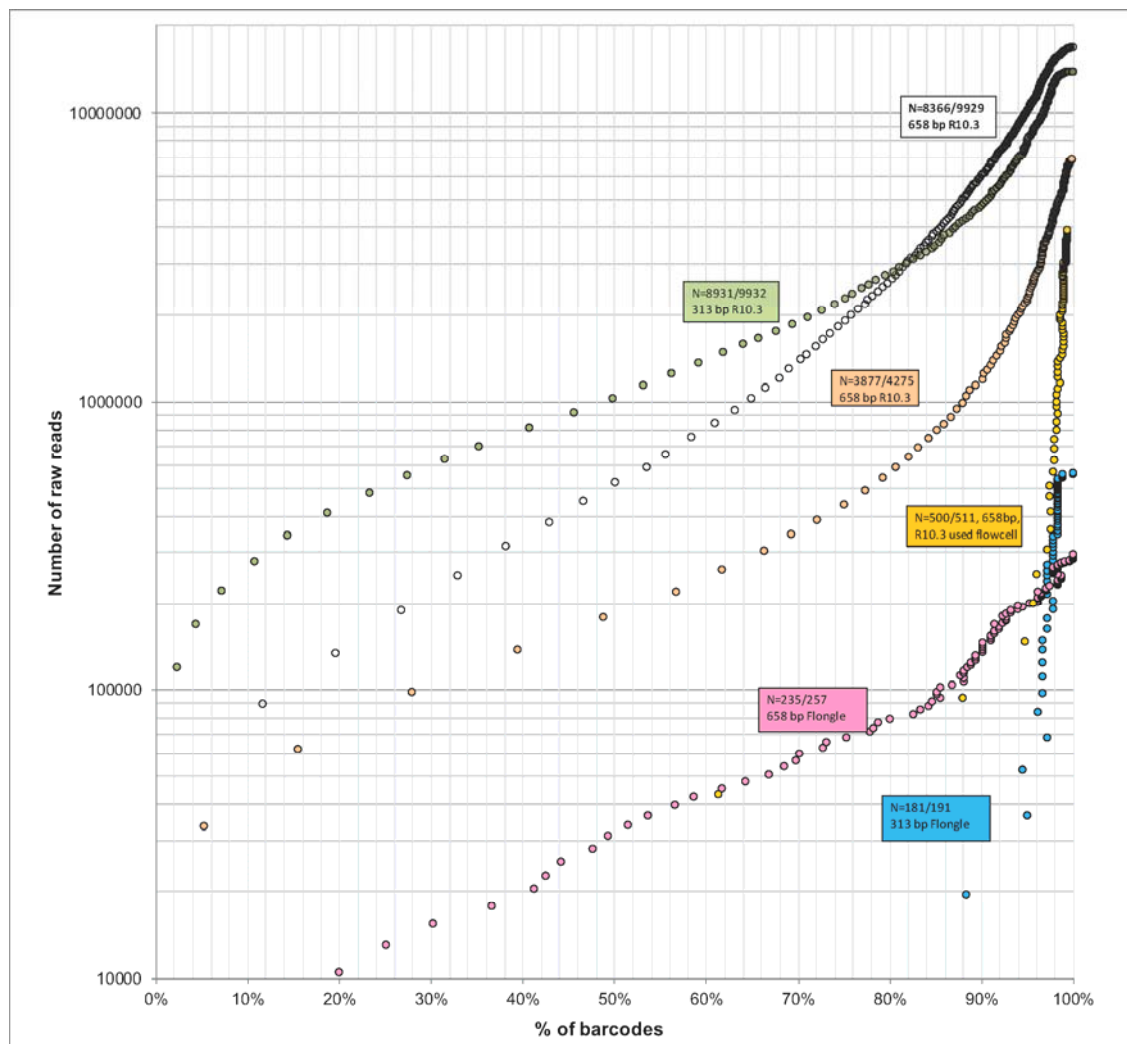


Figure 4. Relationship between barcoding success and number of raw reads for six amplicon pools (191-9932 specimens; barcoding success rates 84-97%). Percentage of barcodes recovered is relative to the final estimate based on all data.

722 **Discussion**

723 *Democratization of Barcoding*

724 Biodiversity research needs new scalable techniques for large-scale species discovery and
725 monitoring. This task is particularly urgent and challenging for invertebrates that collectively
726 make up most of the terrestrial animal biomass. We argued earlier that this is likely to be a
727 task that requires the processing of at least 500 million specimens from all over the world
728 with many tropical countries with limited research funding requiring much of the biodiversity
729 discovery work. Pre-sorting these specimens into putative species-level units with DNA
730 sequences is a promising solution as long as obtaining and analyzing the data are
731 sufficiently straightforward and cost-effective. We believe that the techniques described in
732 this manuscript will help with achieving these goals. Generating DNA barcodes involves
733 three steps. The first is obtaining a DNA template, and we have herein outlined some
734 simplified procedures that render this process essentially free-of-cost, although automation
735 and AI-based solutions will be useful for processing very large numbers of specimens in
736 countries with high manpower cost. The third step is the sequencing of the amplicon.
737 Fortunately, there are now several cost-effective solutions based on 2nd and 3rd generation
738 sequencing technologies so that barcodes can be sequenced for as little as a penny (USD).

739

740 We here argue that sequencing with MinION is particularly attractive although the cost is
741 higher (0.10 USD) than with Illumina sequencing. There are several reasons. MinION library
742 preparation can be learned within hours and an automated library preparation instrument is
743 in development that will eventually work for ligation-based libraries. Furthermore, MinION
744 flow cells can accommodate projects of varying scales. Flongle can be used for amplicon
745 pools with a few hundred products, while an R10.3 flow cell can accommodate projects with
746 up to 10,000 specimens. The collection of data on MinION flow cells can be stopped
747 whenever enough have been acquired. Flow cells can then be washed and re-used again
748 although the remaining capacity declines over time because some nanopores will become
749 unavailable. We have re-used flow cells up to four times. Traditionally, the main obstacles to

750 using MinION have been poor read quality and high cost. Both issues are fading into the
751 past. The quality of MinION reads has improved to such a degree that the laptop-version of
752 our new software “ONTbarcoder” can generate thousands of very high quality barcodes
753 within hours. There is no longer a need to polish reads or rely on external data or algorithms.
754 The greater ease with which MinION barcodes can be obtained is due to several factors.
755 Firstly, much larger numbers of reads can now be obtained with one MinION flow cell.
756 Secondly, R10.3 reads have a different error profile which allows for reconstructing higher-
757 quality barcodes. Thirdly, high accuracy basecalling has improved raw read quality and thus
758 demultiplexing rates. Lastly, we can now use parameter settings for MAFFT that are
759 designed for MinION reads. These changes mean that even low-coverage bins yield very
760 accurate barcodes; i.e., both barcode quality and quantity are greatly improved.

761

762 *Rapid progress in barcode quality and quantity*

763 We previously tested MinION for barcoding (2018, 2019) and here re-sequenced some of
764 the same amplicon pools. This allowed for a precise assessment of the improvements. In
765 2018, sequencing the 511 amplicons of the *Mixed Diptera* sample required one flow cell and
766 we obtained 488 barcodes of which only one lacked ambiguous bases. In 2021, we used the
767 remaining ~500 pores of a used R10.3 flow cell (1st use was for 49 hours). After washing, we
768 obtained 502 barcodes and >98% (496) of them were free of ambiguous bases. The results
769 obtained for the 2019 amplicon pools were also better. In 2019, one flow cell (R9.4) allowed
770 us to obtain 3,223 barcodes from a pool of amplicons obtained from 4,275 specimens of
771 *Afrotropical Phoridae*. Resequencing weak amplicons increased the total number of
772 barcodes by approximately 500 to 3,762 (Srivathsan, Hartop et al. 2019). Now, one R10.3
773 flow cell yielded 3,905 barcodes (+143) for the same amplicon pool, while retaining an
774 accuracy of >99.99% and reducing the ambiguities from 0.45% to 0.01%. If progress
775 continues at this pace, we predict that MinION will be the default barcoding tool for most
776 users. This, too, is because all barcoding steps can now be carried out in one laboratory with
777 a modest set of equipment (see Table 4). With MinION being readily available, there is no

778 longer the need to outsource sequencing and/or to wait until enough barcode amplicons
779 have been prepared for an Illumina or PacBio flow cell (Ho, Puniamoorthy et al. 2020). This
780 democratizes biodiversity discovery and allows many biologists, government agencies,
781 students, and citizen scientists from around the globe to get involved in these initiatives.
782 Biodiversity discovery with cost-effective barcodes will also facilitate biodiversity discovery in
783 countries with high biodiversity but limited science funding.

784

785

Table 4. Equipment required for MinION barcoding

Required	
1	MinION sequencer (preferably Mk1C for basecalling)
2	Thermocycler(s)
3	Gel Electrophoresis setup
4	Magnetic Separation Rack
5	Qubit for DNA quantification
6	<u>Standard equipment:</u> Vortex, Mini-centrifuge, pipettes, freezer, fridge
7	Standard laptop or PC
Optional but highly desirable	
1	Multichannel pipette(s)
2	Hula Mixer

786

787 This raises the question of how much it costs to sequence a barcode with MinION. There is
788 no straightforward answer because the cost depends on user targets. For example, a user
789 who wants to sequence a pool of 5000 barcodes may want a 80% success rate in order to
790 identify the dominant species in a sample. Based on Figure 4, only ca. 1.5 million raw
791 MinION reads would be needed. On average, MinION flow cells yield >10 million reads and
792 cost USD 475-900 depending on how many cells are purchased at the same time. Including

793 a library cost of ca. USD 100, the overall sequencing cost of a project that requires 1.5
794 million reads is USD 180-235. This experiment would be expected to yield 4000 barcodes for
795 the 5000 amplicons (4-6 cents/barcode). Given the low cost of 1 million MinION reads (\$50-
796 90), we predict that most users will opt for sequencing at a greater depth since this will likely
797 yield several hundred additional barcodes. However, this will then increase the sequencing
798 cost per barcode, because the first 1.5 million reads already recovered barcodes for all
799 strong amplicons. Additional reads will predominantly strengthen read coverage for these
800 amplicons and relatively few reads will be added to the read bins that were too weak to yield
801 barcodes at low coverage; i.e., there are diminishing returns for additional sequencing.

802

803 Overall, we thus predict that most users will, at most, try to multiplex 10,000 amplicons in the
804 same MinION flow cell so that the sequencing cost per specimen would be 0.06-0.10 USD
805 depending on the bulk purchase of flow cells. However, we also predict that large-scale
806 biodiversity projects will switch to sequencing with PromethION, a larger sequencing unit
807 that can accommodate up to 48 flow cells. This will lower the sequencing cost by more than
808 60%, as PromethION flow cells have 6 times the number of pores for twice the cost (capacity
809 per flow cell should be 60,000 barcodes). At the other end of the scale are those users who
810 occasionally need a few hundred barcodes. They can use Flongle flow cells, which are
811 comparatively expensive (0.50 USD) because each flow cell costs \$70 and requires a library
812 that is prepared with half the normal reagents (ca. \$50). A change of the flow cell chemistry
813 from that of R9.4 to R10.3 would, however, help with improving the quality of the barcodes
814 obtained from Flongle. Lastly the initial setup cost for MinION/Flongle, can be as low as
815 1000 USD, but we recommend purchase of Mk1C unit (currently 4900 USD) for easy access
816 to a GPU that is required for high accuracy basecalling. Note also, that obtaining flow cells at
817 low cost often requires collaboration between several labs because it allows for buying flow
818 cells in bulk.

819

820 *ONTbarcode for large-scale species discovery with MinION*

821 There are a number of studies that have used MinION for barcoding fungi, animals, and
822 plants (Menegon, Cantaloni et al. 2017, Pomerantz, Peñafiel et al. 2018, Wurzbacher,
823 Larsson et al. 2018, Krehenwinkel, Pomerantz et al. 2019, Maestri, Cosentino et al. 2019,
824 Chang, Ip et al. 2020, Chang, Ip et al. 2020, Knot, Zouganelis et al. 2020, Seah, Lim et al.
825 2020, Sahlin, Lim et al. 2021). There is one fundamental difference between these studies
826 and the vision presented here. These studies tended to show that MinION sequencing can
827 be done in the field. Thus only a very small number of specimens were analysed (<150 with
828 the exception of >500 in Chang, Ip et al 2020). The field use is an attractive feature for time-
829 sensitive samples that could degrade before reaching a lab. However, for the time being it is
830 unlikely to help substantially with tackling the challenges related to large-scale biodiversity
831 discovery and monitoring because obtaining few MinION barcodes per flow cell is too
832 expensive for most researchers in biodiverse countries. Additionally, the bioinformatic
833 pipelines that were developed for these small-scale projects were not suitable for large-
834 scale, decentralized barcoding in a large variety of facilities. For example, some of the
835 studies used ONT's commercial barcoding kit that only allows for multiplexing up to 96
836 samples in one flow cell (Maestri, Cosentino et al. 2019, Seah, Lim et al. 2020); i.e., each
837 amplicon had very high read coverage which influenced the corresponding bioinformatics
838 pipelines (e.g. ONTrack's recommendation is 1000x: Maestri, Cosentino et al. 2019). The
839 generation of such high coverage datasets also meant that the pipelines were only tested for
840 such a small number of samples (<60: Menegon, Cantaloni et al. 2017, Maestri, Cosentino et
841 al. 2019, Seah, Lim et al. 2020, Sahlin, Lim et al. 2021) that these tests were unlikely to
842 represent the complexities of large, multiplexed amplicon pools (e.g., nucleotide diversity,
843 uneven coverage).

844

845 ONTbarcoder evolved from miniBarcoder, whose barcodes have been assessed for
846 accuracy in four different studies covering >8000 barcodes (Chang, Ip et al. 2020, Chang, Ip
847 et al. 2020, Srivathsan, Baloglu et al. 2018, Srivathsan, Hartop et al. 2019). The new
848 software introduced here addresses two drawbacks of its precursor, miniBarcoder. Firstly,

849 we dropped the translation-based error correction that tended to increase the number of Ns.
850 This step used to be essential because indel errors were prevalent in consensus barcodes
851 obtained with older flow cell models. Secondly, ONTbarcoder can be installed by unzipping a
852 file and is easy to maintain on different operating systems. Until now, external dependencies
853 were a major drawback of all MinION bioinformatics pipelines. For example, the one
854 described by Sahlin et al. (2021) involved minibar/qcat and nanofilt, while NGSpeciesID
855 relies on isONclust SPOA, Parasail, and optionally, Medaka (Daily 2016, Krehenwinkel,
856 Pomerantz et al. 2019, Sahlin and Medvedev 2020). These dependencies and complexities
857 meant that Watsa et al. (2020) recommended bioinformatics training before MinION
858 barcoding could be used in schools (e.g., training in UNIX command-line) and additionally
859 required the installation of several software tools onto the teaching computers. Neither is
860 needed for ONTbarcoder, which runs on a regular laptop and has been extensively tested
861 (>4000 direct comparisons to Sanger and Illumina barcodes). In addition, ONTbarcoder is
862 designed in a way that thousands of barcodes can be obtained rapidly without impairing
863 accuracy; i.e., one can run a very fast analysis by using low read coverage. However, at very
864 low coverages, fewer barcodes would be recovered because many would not pass the 4 QC
865 criteria. Speed is also achieved through the parallelization of most steps on UNIX systems
866 (Mac and Linux; parallelization is restricted to demultiplexing in Windows).

867

868 ONTbarcoder also allows for updating the parameter file for alignment. This is advisable
869 because MinION continues to evolve quickly. We expect flow cell capacity to increase further
870 and basecalling to improve (see Xu, Mai et al. 2020). For example, a new basecaller
871 ("*bonito*") developed by ONT has shown promise by improving raw read accuracy
872 ([https://nanoporetech.com/about-us/news/new-research-algorithms-yield-accuracy-gains-](https://nanoporetech.com/about-us/news/new-research-algorithms-yield-accuracy-gains-nanopore-sequencing)
873 [nanopore-sequencing](https://nanoporetech.com/about-us/news/new-research-algorithms-yield-accuracy-gains-nanopore-sequencing)). This basecaller is currently suitable for research teams equipped
874 with GPU infrastructure and for advanced users familiar with Linux command lines.
875 However, our preliminary tests of *bonito* for barcoding (Flongle: *Mixed Diptera Subsample*,
876 *Chironomidae*; R10.3: *Palaeartic Phoridae*, 313 bp; bonito version=0.3.6) does not yet

877 significantly affect barcode quality or quantity (unpublished data). However, this may change
878 in the immediate future and readers are advised to watch out for developments. Fortunately,
879 these changes will only further improve MinION barcodes that are already highly accurate
880 and cost-effective.

881

882 *Biodiversity monitoring*

883 Some readers are likely to argue that large-scale biodiversity discovery and monitoring can
884 be more efficiently carried out via metabarcoding of whole samples consisting of hundreds or
885 thousands of specimens. This would question the need for large-scale, decentralized
886 barcoding of individual specimens. However, large-scale barcoding and metabarcoding will
887 more likely complement each other. For example, large-scale barcoding of individual
888 specimens remains essential for discovering and describing species as it preserves
889 individual voucher specimens associated with the barcode which can be used for further
890 research. Taxonomic research can be guided by examination of putative species units
891 (molecular Operational Taxonomic Units or mOTUs) using species delimitation algorithms
892 (either distance based clustering of sequences: Meier, Shiyang et al. 2006; Puillandre,
893 Brouillet et al. 2020) or tree based methods (Pons, Barraclough et al. 2006; Zhang, Kapli et
894 al. 2013). In this process, it is important to remember that *COI* lumps recently diverged
895 species and divides species with deep allopatric splits (Hickerson, Meyer et al. 2006),
896 making the ability to relate barcodes to individual specimens critical for barcode cluster
897 validation. High quality barcode databases are important for the analysis of metabarcoding
898 data because they facilitate the identification of numts, heteroplasmy, contaminants and
899 errors. Large-scale barcoding will also be needed in order to benefit from another new
900 technique that may become critical for biodiversity discovery and monitoring; i.e. AI-assisted
901 analysis of images (Valan, Makonyi et al. 2019). Large-scale barcoding generates identified
902 specimens that can be imaged and utilized for training neural networks. With increasing
903 advancements in imaging hardware, computational processing power and machine learning
904 systems, AI-assisted biodiversity monitoring could be the method of choice in the future

905 because it could quickly determine and count many common species and only specimens
906 from new/rare species would still require barcoding.

907

908 **Conclusions**

909 Many biologists would like to have ready access to barcodes without having to run large and
910 complex laboratories or send specimens halfway around the world. Many have been
911 impressed by MinION's low cost, portability, and ability to deliver real-time sequencing, but
912 they were worried about high cost and complicated bioinformatics pipelines. We here
913 demonstrate that these concerns are no longer justified. MinION barcodes obtained by
914 R10.3 flow cells are virtually identical to barcodes obtained with Sanger and Illumina
915 sequencing. Barcoding with MinION is now also cost-effective and the new "ONTbarcoder"
916 software makes it straightforward for researchers with little bioinformatics background to
917 analyze the data on a standard laptop. This will make biodiversity discovery scalable and
918 accessible to all.

919

920 **Acknowledgements**

921 We would like to thank John T. Longino and Michael Branstetter for providing valuable
922 comments on the manuscript. For the Palaearctic phorid samples, we would like to thank
923 Dave Karlsson, the Swedish Insect Inventory Project, and the crew at Station Linné that
924 sorted out the phorids. We would also like to thank Wan Ting Lee for help with molecular
925 work, and the numerous staff, students and interns who have contributed to the
926 establishment of the pipeline in the NUS laboratory. We would also like to acknowledge
927 Suphavitai Chayaporn and Niranjan Nagarajan from Genome Institute of Singapore for their
928 help with basecalling using bonito. This work was supported by a Ministry of Education grant
929 on biodiversity discovery (R-154-000-A22-112).

930

931 **Software and test dataset availability**

932 ONTbarcoder is available at <https://github.com/asrivathsan/ONTbarcoder>, which also
933 contains the link to download the raw data and demultiplexing files. The manual for the
934 software is included in the repository
935 https://github.com/asrivathsan/ONTbarcoder/blob/main/ONTBarcoder_manual.pdf. The
936 videos tutorials can be found in the YouTube channel Integrative Biodiversity Discovery:
937 <https://www.youtube.com/channel/UC1WowokomhQJRc71FmsUAcg>.

938 **Literature cited**

939 Abgrego, N., T. Roslin, T. Huotari, Y. Ji, N.M. Schmidt, J. Wang, D. Yu and O. Ovaskainen
940 (2021). "Accounting for species interactions is necessary for predicting how arctic
941 arthropod communities respond to climate change." *Ecography* doi:
942 10.1111/ecog.05547.

943 Arribas, P., C. Andújar, K. Hopkins, M. Shepherd and A. P. Vogler (2016). "Metabarcoding
944 and mitochondrial metagenomics of endogean arthropods to unveil the mesofauna of
945 the soil." *Methods in Ecology and Evolution* **7**(9): 1071-1081.

946 Baloğlu, B., E. Clews and R. Meier (2018). "NGS barcoding reveals high resistance of a
947 hyperdiverse chironomid (Diptera) swamp fauna against invasion from adjacent
948 freshwater reservoirs." *Frontiers in Zoology* **15**(1): 31.

949 Bar-On, Y. M., R. Phillips and R. Milo (2018). "The biomass distribution on Earth."
950 *Proceedings of the National Academy of Sciences* **115**(25): 6506-6511.

951 Barrett, R. D. H. and P. D. Hebert (2005). "Identifying spiders through DNA barcodes."
952 *Canadian Journal of Zoology* **83**: 481-491.

953 Bell, J.R., D. Blumgart and C.R. Shortall (2020). "Are insects declining and at what rate? An
954 analysis of standardised, systematic catches of aphid and moth abundances across
955 Great Britain". *Insect Conservation and Diversity* **13**(2): 115-126.

956 Chang, J. J. M., Y. C. A. Ip, A. G. Bauman and D. Huang (2020). "MinION-in-ARMS:
957 Nanopore sequencing to expedite barcoding of specimen-rich macrofaunal samples
958 from autonomous reef monitoring structures." *Frontiers in Marine Science* **7**: 448.

- 959 Chang, J. J. M., Y. C. A. Ip, C. S. L. Ng and D. Huang (2020). "Takeaways from mobile DNA
960 barcoding with BentoLab and MinION." Genes **11**: 1121.
- 961 Crampton-Platt, A., D. W. Yu, X. Zhou and A. P. Vogler (2016). "Mitochondrial
962 metagenomics: letting the genes out of the bottle." Gigascience **5**(1): s13742-13016-
963 10120-y.
- 964 Daily, J. (2016). "Parasail: SIMD C library for global, semi-global, and local pairwise
965 sequence alignments." BMC Bioinformatics **17**: 81.
- 966 Eisenhauer, N., A. Bonn and C.A. Guerra (2019). "Recognizing the quiet extinction of
967 invertebrates". Nature Communications **10**: 50.
- 968 Forum, W. E. (2020). "World Economic Forum. The Global Risks Report 2020.", from
969 <https://www.weforum.org/reports/the-global-risks-report-2020>.
- 970 Frith, M. C., S. Mitsuhashi and K. Katoh (2020). Iamassemble: Multiple Alignment and
971 Consensus Sequence of Long Reads. Multiple Sequence Alignment. K. Katoh. New
972 York, Humana: 135-145.
- 973 Groombridge, B., Ed. (1992). Global Biodiversity: Status of the Earth's Living Resources.
974 World Conservation Monitoring Centre. London, Chapman & Hall.
- 975 Grootaert, P. (2018). "Revision of the genus *Thinophihis* Wahlberg (Diptera: Dolichopodidae)
976 from Singapore and adjacent regions: A long term study with a prudent reconciliation
977 of a genetic to a classic morphological approach." Raffles Bulletin of Zoology **66**: 413-
978 473.
- 979 Grootaert, P. (2019). "Species turnover between the northern and southern part of the South
980 China Sea in the Elaphropeza Macquart mangrove fly communities of Hong Kong and
981 Singapore (Insecta: Diptera: Hybotidae)." European Journal of Taxonomy **554**: 1-27.
- 982 Hallman, C.A., M. Sorg, E. Jongejans, H. Siepel, N. Hofland, H. Schwan, W. Stenmans, A.
983 Müller, H. Sumser, T. Hörren, D. Goulson and H. de kroon (2017). "More than 75
984 percent decline over 27 years in total flying insect biomass in protected areas." PLoS
985 One **12**(10): e0185809.

- 986 Hallman, C.A., A. Ssymank, M. Sorg, H. de Kroon and E Jongejans (2021). "Insect biomass
987 decline scaled to species diversity: General patterns derived from a hoverfly
988 community." Proceedings of the National Academy of Sciences **118**(2): e2002554117.
- 989 Hamada, M., Y. Ono, K. Asai and M. C. J. B. Frith (2017). "Training alignment parameters
990 for arbitrary sequencers with LAST-TRAIN." Bioinformatics **33**(6): 926-928.
- 991 Hebert, P. D., T. W. A. Braukmann, S. W. J. Prosser, S. Ratnasingham, J. R. deWaard, N.
992 V. Ivanova, D. Janzen, W. Hallwachs, S. Naik, J. E. Sones and E. V. Zakharov (2018).
993 "A Sequel to Sanger: amplicon sequencing that scales." BMC Genomics **19**: 219.
- 994 Hebert, P. D., J. R. DeWaard, E. V. Zakharov, S. W. J. Prosser, J. E. Sones, J. T. A.
995 McKeown, B. Mantle and J. La Salle (2013). "A DNA 'Barcode Blitz': Rapid digitization
996 and sequencing of a Natural History collection." PLoS One **8**(7): e68535.
- 997 Hebert, P. D. N., A. Cywinska, S. L. Ball and J. R. deWaard (2003). "Biological identifications
998 through DNA barcodes." Proceedings of the Royal Society Biological Sciences Series
999 B **270**(1512): 313-321.
- 1000 Hebert, P. D. N., S. Ratnasingham, E. V. Zakharov, A. C. Telfer, V. Levesque-Beaudin, M.
1001 A. Milton, S. Pedersen, P. Jannetta and J. R. deWaard (2016). "Counting animal
1002 species with DNA barcodes: Canadian insects." Philosophical Transactions of the
1003 Royal Society B: Biological Sciences **371**: 20150333.
- 1004 Hendrich, L., J. Pons, I. Ribera and M. Balke (2010). "Mitochondrial Cox1 sequence data
1005 reliably uncover patterns of insect diversity but suffer from high lineage-idiosyncratic
1006 error rates." PLoS One **5**(12): e14448.
- 1007 Hickerson, M. J., C. P. Meyer and Moritz (2006). "DNA barcoding will often fail to discover
1008 new animal species over broad parameter space." Systematic Biology **55**(5): 729-739.
- 1009 Ho, J. K. I., J. Puniamoorthy, A. Srivathsan and R. Meier (2020). "MinION sequencing of
1010 seafood in Singapore reveals creatively labelled flatfishes, confused roe, pig DNA in
1011 squid balls, and phantom crustaceans." Food Control **112**: 107144.

- 1012 Ismay, B. and Y. Ang (2019). "First records of *Pseudogaurax* Malloch 1915 (Diptera:
1013 Chloropidae) from Singapore, with the description of two new species discovered with
1014 NGS barcodes." Raffles Bulletin of Zoology **67**: 412-420.
- 1015 Ivanova, N. V., A. V. Borisenko and P. D. N. Hebert (2009). "Express barcodes: racing from
1016 specimen to identification." Molecular Ecology Resources **9**: 35-41.
- 1017 Ivanova, N. V., J. R. Dewaard and P. D. N. Hebert (2006). "An inexpensive, automation-
1018 friendly protocol for recovering high-quality DNA." Molecular Ecology Notes **6**(4): 998-
1019 1002.
- 1020 Katoh, K. and D. M. Standley (2013). "MAFFT Multiple Sequence Alignment Software
1021 Version 7: Improvements in performance and usability." Molecular Biology and
1022 Evolution **30**(4): 772-780.
- 1023 Knot, I. E., G. D. Zouganelis, G. D. Weedall, S. A. Wich and R. Rae (2020). "DNA barcoding
1024 of nematodes using the MinION." Frontiers in Ecology and Evolution **8**: 100.
- 1025 Knox, M. A., I. D. Hogg, C. A. Pilditch, J. C. Garcia-R, P. D. N. Hebert and D. Steinke (2020).
1026 "Contrasting patterns of genetic differentiation for deep-sea amphipod taxa along New
1027 Zealand's continental margins." Deep Sea Research Part I: Oceanographic Research
1028 Papers **162**: 103323.
- 1029 Kranzfelder, P., T. Ekrem and E. Stur (2016). "Trace DNA from insect skins: a comparison of
1030 five extraction protocols and direct PCR on chironomid pupal exuviae." Molecular
1031 Ecology Resources **16**(1): 353-363.
- 1032 Krehenwinkel, H., S. R. Kennedy, A. Rueda, A. Lam and R. G. Gillespie (2018). "Scaling up
1033 DNA barcoding – Primer sets for simple and cost efficient arthropod systematics by
1034 multiplex PCR and Illumina amplicon sequencing." Methods in Ecology and Evolution
1035 **9**(11): 2181-2193.
- 1036 Krehenwinkel, H., A. Pomerantz, J. B. Henderson, S. R. Kennedy, J. Y. Lim, V. Swamy, J. D.
1037 Shoobridge, N. Graham, N. H. Patel, R. G. Gillespie and S. Prost (2019). "Nanopore
1038 sequencing of long ribosomal DNA amplicons enables portable and simple biodiversity

- 1039 assessments with high phylogenetic resolution across broad taxonomic scale."
1040 Gigascience **8**(5): giz006.
- 1041 Krell, F. T. (2004). "Parataxonomy vs. taxonomy in biodiversity studies - pitfalls and
1042 applicability of 'morphospecies' sorting." Biodiversity and Conservation **13**(4): 795-812.
- 1043 Kwong, S., A. Srivathsan and R. Meier (2012). "An update on DNA barcoding: low species
1044 coverage and numerous unidentified sequences." Cladistics **28**(6): 639-644.
- 1045 Lim, N. K. M., Y. C. Tay, A. Srivathsan, J. W. T. Tan, J. T. B. Kwik, B. Baloğlu, R. Meier and
1046 D. C. J. Yeo (2016). "Next-generation freshwater bioassessment: eDNA
1047 metabarcoding with a conserved metazoan primer reveals species-rich and reservoir-
1048 specific communities." Royal Society Open Science **3**: 160635.
- 1049 Maestri, S., E. Cosentino, M. Paterno, H. Freitag, J. M. Garces, L. Marcolungo, M. Alfano, I.
1050 Njunjić, M. Schilthuizen, F. Slik, M. Menegon, M. Rossato and M. Delledonne (2019).
1051 "A rapid and accurate MinION-based workflow for tracking species biodiversity in the
1052 field." Genes **10**(6): 468.
- 1053 Meier, R., K. Shiyang, G. Vaidya and P.K.L. Ng (2006). "DNA barcoding and taxonomy in
1054 Diptera: a tale of high intraspecific variability and low identification success."
1055 Systematic Biology **55**(5): 715-728.
- 1056 Meier, R. (2008). DNA sequences in taxonomy - Opportunities and challenges. New
1057 Taxonomy. Q. D. Wheeler. **76**: 95-127.
- 1058 Meier, R., W. H. Wong, A. Srivathsan and M. S. Foo (2016). "\$1 DNA barcodes for
1059 reconstructing complex phenomes and finding rare species in specimen-rich samples."
1060 Cladistics **32**(1): 100-110.
- 1061 Menegon, M., C. Cantaloni, A. Rodriguez-Prieto, C. Centomo, A. Abdelfattah, M. Rossato,
1062 M. Bernardi, L. Xumerle, S. Loader and M. Delledonne (2017). "On site DNA barcoding
1063 by nanopore sequencing." PLoS One **12**(10): e0184741.
- 1064 Ng'endo, R. N., Z. B. Osiemo and R. Brandl (2013). "DNA barcodes for species identification
1065 in the hyperdiverse ant genus *Pheidole* (Formicidae: Myrmicinae)." Journal of Insect
1066 Science **13**: 27.

- 1067 Page, R. (2011). "Dark taxa: GenBank in a post-taxonomic world."
1068 <https://iphylo.blogspot.com/2011/04/dark-taxa-genbank-in-post-taxonomic.html>,
1069 Accessed February 2021.
- 1070 Pearson, W. R. (2017). "Finding protein and nucleotide similarities with FASTA." Current
1071 Protocols in Bioinformatics **53**: 3.9.1-3.9.25.
- 1072 Pomerantz, A., N. Peñafiel, A. Arteaga, L. Bustamante, F. Pichardo, L. A. Coloma, C. L.
1073 Barrio-Amorós, D. Salazar-Valenzuela and S. J. G. Prost (2018). "Real-time DNA
1074 barcoding in a rainforest using nanopore sequencing: opportunities for rapid
1075 biodiversity assessments and local capacity building." **7**(4): giy033.
- 1076 Ponder, W. and D. Lunney (1999). The Other 99% - the Conservation and Biodiversity of
1077 Invertebrates. Sydney, Transactions of the Royal Zoological Society of New South
1078 Wales.
- 1079 Pons, J., T.G. Barraclough, J. Gomez-Zurita, A. Cardoso, D.P. Duran, S. Hazell, S. Kamoun,
1080 W.D. Sumlin and A.P. Vogler (2006). Sequence-Based Species Delimitation for the
1081 DNA Taxonomy of Undescribed Insects. Systematic Biology **55**(4): 595-609.
- 1082 Puillandre, N., S. Brouillet, G. Achaz (2021). ASAP: assemble species by automatic
1083 partitioning. Molecular Ecology Resources **21**: 609-620.
- 1084 Stepanian, P.M., S.A. Entrekin, C.E. Wainwright, D.Mirkovic, J.L. Tank and J.F. Kelly (2020).
1085 "Declines in an abundant aquatic insect, the burrowing mayfly, across major North
1086 American waterways." Proceedings of the National Academy of Sciences **117**(6):
1087 2987-2992.
- 1088 Swiss Re. (2020). " Biodiversity and Ecosystem Services A business case for re/insurance."
1089 Zurich, Swiss Re Management Ltd.
- 1090 Rohland, N. and D. Reich (2012). "Cost-effective, high-throughput DNA sequencing libraries
1091 for multiplexed target capture." Genome research **22**: 939-946.
- 1092 Sahlin, K., M. C. W. Lim and S. Prost (2021). "NGSpeciesID: DNA barcode and amplicon
1093 consensus generation from long-read sequencing data." Ecology and Evolution **11**(3):
1094 1392-1398.

- 1095 Sahlin, K. and P. Medvedev (2020). "De novo clustering of long-read transcriptome data
1096 using a greedy, quality value-based algorithm." Journal of Computational Biology
1097 **27**(4): 472-484.
- 1098 Samoh, A., C. Satasook and P. Grootaert (2019). "NGS-barcodes, haplotype networks
1099 combined to external morphology help to identify new species in the mangrove genus
1100 *Ngirhaphium* Evenhuis & Grootaert, 2002 (Diptera: Dolichopodidae: Rhaphiinae) in
1101 Southeast Asia." Raffles Bulletin of Zoology **67**: 640-659.
- 1102 Seah, A., M. C. W. Lim, D. McAloose, S. Prost and T. A. Seimon (2020). "MinION-based
1103 DNA barcoding of preserved and non-Invasively collected wildlife samples." Genes
1104 **11**(4): 445.
- 1105 Shokralla, S., T. M. Porter, J. F. Gibson, R. Dobosz, D. Janzen, W. Hallwachs, G. B. Golding
1106 and M. Hajibabaei (2015). "Massively parallel multiplex DNA sequencing for specimen
1107 identification using an Illumina MiSeq platform." Scientific Reports **5**: 9687.
- 1108 Shokralla, S., J. L. Spall, J. F. Gibson and M. Hajibabaei (2012). "Next-generation
1109 sequencing technologies for environmental DNA research." Molecular Ecology **21**(8):
1110 1794-1805.
- 1111 Silvestre-Ryan, J. and I. Holmes (2021). "Pair consensus decoding improves accuracy of
1112 neural network basecallers for nanopore sequencing." Genome Biology **22**: 38.
- 1113 Sović, I., M. Šikić, A. Wilm, S. N. Fenlon, S. Chen and N. Nagarajan (2016). "Fast and
1114 sensitive mapping of nanopore sequencing reads with GraphMap." Nature
1115 Communications **7**: 11307.
- 1116 Srivathsan, A., B. Baloğlu, W. Wang, W. X. Tan, D. Bertrand, A. H. Q. Ng, E. J. H. Boey, J.
1117 J. Y. Koh, N. Nagarajan and R. Meier (2018). "A MinION-based pipeline for fast and
1118 cost-effective DNA barcoding." Molecular Ecology Resources **18**(5): 1035-1049.
- 1119 Srivathsan, A., E. Hartop, J. Puniamoorthy, W. T. Lee, S. N. Kutty, O. Kurina and R. Meier
1120 (2019). "Rapid, large-scale species discovery in hyperdiverse taxa using 1D MinION
1121 sequencing." BMC Biology **17**(1): 96.

- 1122 Srivathsan, A., N. Nagarajan and R. Meier (2019). "Boosting natural history research via
1123 metagenomic clean-up of crowdsourced feces." PLoS Biology **17**(11): e3000517.
- 1124 Stork, N. E., J. McBroom, C. Gely and A. J. Hamilton (2015). "New approaches narrow
1125 global species estimates for beetles, insects, and terrestrial arthropods." Proceedings
1126 of the National Academy of Sciences **112**(24): 7519-7523.
- 1127 Stribling, J. B., K. L. Pavlik, S. M. Holdsworth and E. W. Leppo (2008). "Data quality,
1128 performance, and uncertainty in taxonomic identification for biological assessments."
1129 Journal of the North American Benthological Society **27**(4): 906-919.
- 1130 Tang, C. F., P. Grootaert and D. Yang (2018). "*Protomedetera*, a new genus from the
1131 Oriental and Australasian realms (Diptera, Dolichopodidae, Medeterinae)." Zookeys
1132 **743**: 137-151.
- 1133 Tang, C. F., D. Yang and P. Grootaert (2018). "Revision of the genus *Lichtwardtia* Enderlein
1134 in Southeast Asia, a tale of highly diverse male terminalia (Diptera, Dolichopodidae)."
1135 Zookeys **798**: 63-107.
- 1136 Tautz, D., P. Arctander, A. Minelli, R. H. Thomas and A. P. Vogler (2003). "A plea for DNA
1137 taxonomy." Trends in Ecology & Evolution **18**(2): 70-74.
- 1138 Thongjued, K., W. Chotigeat, S. Bumrungsri, P. Thanakiatkrai and T. Kitpipit (2019). "A new
1139 cost-effective and fast direct PCR protocol for insects based on PBS buffer." Molecular
1140 Ecology Resources **19**(3): 691-701.
- 1141 Thormann, B., D. Ahrens, D. M. Armijos, M. K. Peters and T. Wagner (2016). "Exploring the
1142 leaf beetle fauna (Coleoptera: Chrysomelidae) of an Ecuadorian mountain forest using
1143 DNA barcoding." PLoS One **11**(2): e0148268.
- 1144 Truett, G., P. Heeger, R. Mynatt, A. Truett, J. Walker and M. J. B. Warman (2000).
1145 "Preparation of PCR-quality mouse genomic DNA with hot sodium hydroxide and tris
1146 (HotSHOT)." Biotechniques **29**(1): 52-54.
- 1147 Valan, M., K. Makonyi, A. Maki, D. Vondráček and F. Ronquist (2019). "Automated
1148 taxonomic identification of insects with expert-level accuracy using effective feature
1149 transfer from convolutional networks." Systematic Biology **68**(6): 876-895.

- 1150 Vaser, R., I. Sovic, N. Nagarajan and M. Sikic (2017). "Fast and accurate de novo genome
1151 assembly from long uncorrected reads." Genome Res **27**(5): 737-746.
- 1152 Vereecke, N., J. Bokma, F. Haesebrouck, H. Nauwynck, F. Boyen, B. Pardon and S. Theuns
1153 (2020). "High quality genome assemblies of *Mycoplasma bovis* using a taxon-specific
1154 Bonito basecaller for MinION and Flongle long-read nanopore sequencing." BMC
1155 Bioinformatics **21**: 517.
- 1156 Wagner, D.L. E.M. Grames, M.L. Forister, M.R. Berenbaum and D. Stopak (2021). "Insect
1157 decline in the Anthropocene: Death by a thousand cuts." Proceedings of the National
1158 Academy of Sciences **118**(2): e2023989118.
- 1159 Wang, W. Y., A. Srivathsan, M. Foo, S. K. Yamane and R. Meier (2018). "Sorting specimen-
1160 rich invertebrate samples with cost-effective NGS barcodes: Validating a reverse
1161 workflow for specimen processing." Molecular Ecology Resources **18**(3): 490-501.
- 1162 Wang, W. Y., A. Yamada and K. Eguchi (2018). "First discovery of the mangrove ant
1163 *Pheidole sexspinosa* Mayr, 1870 (Formicidae: Myrmicinae) from the Oriental region,
1164 with redescrptions of the worker, queen and male." Raffles Bulletin of Zoology **66**:
1165 652-663.
- 1166 Wang, W. Y., A. Yamada and S. Yamane (2020). "Maritime trap-jaw ants (Hymenoptera,
1167 Formicidae, Ponerinae) of the Indo-Australian region - redescription of *Odontomachus*
1168 *malignus* Smith and description of a related new species from Singapore, including
1169 first descriptions of males." Zookeys **915**: 137-174.
- 1170 Wang, W. Y., G. W. J. Yong and W. Jaitrong (2018). "The ant genus *Rhopalomastix*
1171 (Hymenoptera: Formicidae: Myrmicinae) in Southeast Asia, with descriptions of four
1172 new species from Singapore based on morphology and DNA barcoding." Zootaxa
1173 **4532**(3): 301-340.
- 1174 Watsa, M., G. A. Erkenwick, a. Pomerantz and S. Prost (2020). "Portable sequencing as a
1175 teaching tool in conservation and biodiversity research." PLoS Biology **18**(4):
1176 e3000667.

- 1177 Wick, R. R. (2019). "Performance of neural network basecalling tools for Oxford Nanopore
1178 sequencing." Genome Biology **20**: 129.
- 1179 Wong, W. H., Y. C. Tay, J. Puniamoorthy, M. Balke, P. S. Cranston and R. Meier (2014).
1180 "Direct PCR' optimization yields a rapid, cost-effective, nondestructive and efficient
1181 method for obtaining DNA barcodes without DNA extraction." Molecular Ecology
1182 Resources **14**(6): 1271-1280.
- 1183 Wurzbacher, C., E. Larsson, J. Bengtsson-Palme, S. V. den Wyngaert, S. Svantesson, E.
1184 Kristiansson, M. Kagami and R. H. Nilsson (2018). "Introducing ribosomal tandem
1185 repeat barcoding for fungi." Molecular Ecology Resources **19**(1): 118-127.
- 1186 Xu, Z., Y. Mai, D. Liu, W. He, X. Lin, C. Xu, L. Zhang, X. Meng, J. Mafofo, W. A. Zaher, Y. Li
1187 and N. Qiao (2020). "Fast-Bonito: A faster basecaller for nanopore sequencing."
1188 BioRxiv: doi:10.1101/2020.1110.1108.318535.
- 1189 Yeo, D., J. Puniamoorthy, R. W. J. Ngiam and R. Meier (2018). "Towards holomorphology in
1190 entomology: rapid and cost-effective adult-larva matching using NGS barcodes."
1191 Systematic Entomology **43**(4): 678-691.
- 1192 Yeo, D., A. Srivathsan and R. Meier (2020). "Longer is Not Always Better: Optimizing
1193 Barcode Length for Large-Scale Species Discovery and Identification." Systematic
1194 Biology **69**(5): 999-1015.
- 1195 Yeo, D., A. Srivathsan, J. Puniamoorthy, M. Foo, P. Grootaert, L. Chan, B. Guenard, C.
1196 Damken, R. A. Wahab and Y. J. b. Ang (2020). "Mangroves are an overlooked hotspot
1197 of insect diversity despite low plant diversity." BioRxiv:
1198 doi:10.1101/2020.12.17.423191.
- 1199 Zhang, J. P. Kapli, P. Pavlidis and A. Stamatakis (2013). "A general species delimitation
1200 method with applications to phylogenetic placements." Bioinformatics **29**(22): 2869-
1201 2876.