

1 A rapid and accurate MinION-based workflow for 2 tracking species biodiversity in the field

3 Simone Maestri^{1#}, Emanuela Cosentino^{1#}, Marta Paterno^{1,6}, Hendrik Freitag^{2,6,7}, Jhoana M. Garces²,
4 Luca Marcolungo¹, Massimiliano Alfano¹, Iva Njunji^{3,4,6}, Menno Schilthuizen^{3,4,5,6}, Ferry Slik⁶,
5 Michele Menegon⁸, Marzia Rossato¹ and Massimo Delledonne^{1,6*}

6 ¹ Department of Biotechnology, University of Verona, Strada Le Grazie 15, 37134, Verona, Italy;

7 ² Department of Biology, Ateneo de Manila University, School of Science & Engineering, Loyola Heights,
8 Quezon City 1101, Philippines;

9 ³ Taxon Expeditions B.V., Rembrandtstraat 20, 2311 VW Leiden, the Netherlands;

10 ⁴ Naturalis Biodiversity Center, Darwinweg 2, 2333 CR Leiden, the Netherlands;

11 ⁵ Institute for Biology Leiden, Leiden University, Sylviusweg 72, 2333 BE Leiden, the Netherlands;

12 ⁶ Faculty of Science, University of Brunei Darussalam, Jalan Tungku Link, Gadong BE1410, Brunei
13 Darussalam;

14 ⁷ Museum für Naturkunde, Leibniz Institute for Evolution and Biodiversity Science, Invalidenstraße 43,
15 Berlin, 10115, Germany

16 ⁸ Division of Biology & Conservation Ecology, Manchester Metropolitan University, UK; PAMS Foundation,
17 P.O. Box 16556, Arusha, Tanzania;

18 # S.M. and E.C. contributed equally

19 * Correspondence: massimo.delledonne@univr.it; Tel.: +39 045 8027962

20

21 **Abstract:** Genetic markers (DNA barcodes) are often used to support and confirm species
22 identification. Barcode sequences can be generated in the field using portable systems based on the
23 Oxford Nanopore Technologies (ONT) MinION platform. However, to achieve a broader application,
24 current proof-of-principle workflows for on-site barcoding analysis must be standardized to ensure
25 reliable and robust performance under suboptimal field conditions without increasing costs. Here we
26 demonstrate the implementation of a new on-site workflow for DNA extraction, PCR-based
27 barcoding and the generation of consensus sequences. The portable laboratory features inexpensive
28 instruments that can be carried as hand luggage and uses standard molecular biology protocols and
29 reagents that tolerate adverse environmental conditions. Barcodes are sequenced using MinION
30 technology and analyzed with ONTrack, an original *de novo* assembly pipeline that requires as few
31 as 500 reads per sample. ONTrack-derived consensus barcodes have high accuracy, ranging from
32 99,8% to 100%, despite the presence of homopolymer runs. The ONTrack pipeline has a user-friendly
33 interface and returns consensus sequences in minutes. The remarkable accuracy and low
34 computational demand of the ONTrack pipeline, together with the inexpensive equipment and
35 simple protocols, make the proposed workflow particularly suitable for tracking species under field
36 conditions.

37 **Keywords:** nanopore sequencing; long reads; field ecology; barcoding; portable lab; biodiversity

38

39 1. Introduction

40 Recent advances in molecular biology allow the use of genetic markers (DNA barcodes) to
41 support and confirm morphological evidence for species identification and to quantify interspecific
42 differences in order to compare species in terms of evolutionary distance. Most barcodes are still
43 generated using the Sanger sequencing method, which requires access to a well-equipped molecular
44 biology laboratory. Second-generation sequencing technologies are also used for barcoding, but they
45 depend on expensive equipment and the reads are often too short to distinguish species reliably. The
46 third-generation sequencer Oxford Nanopore Technologies (ONT) MinION based on nanopores has
47 proven successful for sequencing under extreme field conditions such as the tropical rainforests of
48 Tanzania, Ecuador and Brazil [1-3], the hot savannah of West Africa [4], and the ice floes of Antarctica

49 [5]. Bringing the laboratory to the field avoids the transport of samples to sequencing facilities, thus
50 greatly reducing the analysis time and the need to export genetic material from collection sites.

51 Although several groups have reported successful on-site barcoding, it remains difficult to
52 perform molecular biology procedures in sub-optimal and extreme environments. In our first
53 expeditions, the quality of sequences generated in the field was consistently lower than achieved in
54 the laboratory, suggesting that reagents and flow cells were affected by the unstable shipping and/or
55 environmental conditions [1]. Furthermore, a recent on-site MinION run produced a low output
56 consisting primarily of adapter sequences, probably reflecting the deterioration of the ligation
57 enzyme and flow cells during suboptimal storage [2]. Some groups used lyophilized reagents to
58 overcome adverse environments [1]. However, also equipment can be affected by extreme conditions,
59 as we found on two different expeditions to Borneo during which one of the two models of portable
60 PCR machine we brought with us lost temperature calibration resulting in the overheating and
61 consequent failure in barcode amplification. The identification of robust protocols and equipment
62 that tolerates suboptimal transport and operating conditions (but remains simple, inexpensive and
63 portable) is therefore highly desirable in order to exploit the full potential of barcode sequencing in
64 the field.

65 MinION-based sequencing is advantageous because it is portable, but it has a higher error rate
66 than other methods and thus appropriate analysis workflows are therefore needed to generate high-
67 quality barcode sequences [1,6]. High accuracy is particularly important in DNA-based taxonomy, as
68 the threshold for intra- versus interspecific divergence of the COI gene is usually at about 2% [7] and
69 in evolutionary 'young' species even lower [8]. We have previously attempted to reduce the high
70 error rate of MinION by using more accurate 2D reads derived from the consensus of the forward
71 and reverse strands. However, 2D sequencing kits are no longer available and have been replaced by
72 1D² kits, which have yet to be optimized for amplicon sequencing. Even so, new ONT chemistries
73 and software updates have greatly improved the throughput and 1D-read accuracy of nanopore
74 sequencing in the last 2 years [8, 9]. Based on this reduced error rate (10–15%, R9.4 chemistry), several
75 groups developed their own data analysis pipelines for barcoding, but none of the methods has yet
76 achieved the status of 'the gold standard' [1,2,6,9].

77 Two main strategies are used to generate high-quality barcode sequences: reference-based and
78 *de novo* pipelines. During the early development of nanopore sequencing, the high error rate in
79 homopolymer runs made reference-based methods the better approach [1,2]. In a typical workflow,
80 sequence reads are mapped to a reference sequence selected according to *a priori* knowledge, and the
81 consensus sequence is ultimately determined based on the majority rule. Reference-based pipelines
82 are useful when matching a target sequence to similar existing ones, but they struggle to reconstruct
83 an accurate barcode if the organism of interest has not been sequenced before. Notably, if the target
84 species carries an insertion compared to the reference species, the additional nucleotides are not
85 included in the final consensus sequence [2]. Unlike the reference-based approach, *de novo* assembly
86 pipelines rely only on the newly-generated reads. Therefore, they suffer more sequencing errors,
87 especially if they are distributed in a nonrandom manner, and *ad hoc* error correction methods are
88 needed to generate the barcodes using *de novo* assembly [2].

89 Recently, hybrid methods incorporating aspects of both approaches have been described [1,6].
90 One example is our *ONtoBAR* pipeline [1]. This creates a draft consensus sequence by assembling
91 MinION reads *de novo* and uses the draft to retrieve the most similar sequence from the NCBI nt
92 database, allowing the final consensus to be generated. Given the assumption that closely-related
93 species differ mainly due to the accumulation of single-nucleotide polymorphisms (SNPs) rather than
94 insertion/deletion polymorphisms (INDELs) that can generate frameshifts, the pipeline uses the
95 reference sequence as a scaffold, allowing the correction of mismatches derived from MinION errors.
96 Another hybrid method known as the *aacorrection* pipeline [6] is based on similar principles, in that a
97 draft consensus sequence is used to recover matching sequences from the NCBI nt database. These
98 are used to determine the correct reading frame, and generic bases (N) are introduced into the
99 MinION-derived consensus in order to preserve amino acid assignments. A recent study compared
100 reference-based and *de novo* approaches, finding that the *de novo* approach was more accurate because

101 the reference-based approach can introduce bias by missing INDELs [2]. However, the filtering step
102 in the proposed pipeline relied on quality scores (Q-scores) that are often recalibrated after basecaller
103 updates, making the results strongly dependent on the sequencing chemistry and the basecaller
104 version.

105 To fully exploit the potential of barcoding in the field, the proof-of-principle workflows reported
106 thus far must be translated into standardized systems allowing on-site sequencing by professional
107 users. Our involvement in conservation projects has motivated us not only to continuously improve
108 the analytical precision of the pipeline in order to track biodiversity at the species level more
109 accurately, but also to identify simple, rapid and inexpensive protocols. Here we demonstrate the
110 results achieved using an updated barcoding workflow that features improvements both to the
111 molecular biology field laboratory components and the subsequent data analysis.

112 2. Materials and Methods

113 2.1 Portable genomics laboratory

114 The portable genomics laboratory included the following equipment: three micropipettes (P1000,
115 P200 and P20, Eppendorf), a mini-microcentrifuge (Labnet Prism Mini Centrifuge, Labnet), a thermal
116 cycler (MiniOne PCR System, MiniOne), an electrophoresis system (MiniOne Electrophoresis System,
117 MiniOne), a fluorometer (Qubit 2.0, Thermo Fisher Scientific), the nanopore sequencer (MinION, ONT)
118 and an ASUS laptop (i7 processor, 16 GB RAM, 500 GB SSD) (**Figure 1**). The equipment was wrapped
119 in air-bubble packaging, transported in a single Peli case (55×45×20 cm) (**Figure 1**) and checked as
120 standard hold baggage in domestic and international flights (except the laptop, which was carried in
121 the cabin). Standard molecular biology reagents were selected and used as described below. Reagents
122 that required storage at 4 °C or –20 °C were transported in a foam box containing ice packs, and MinION
123 flow cells were stored in a thermal bag in the same box. PCR primers were transported lyophilized and
124 subsequently resuspended in 10 mM Tris-HCl (pH 8.0) supplemented with 1 mM EDTA and kept at
125 room temperature.

126

127 2.2 Sample collection, DNA extraction and barcode amplification

128 Sample collection, tissue dissection, total DNA extraction, barcode amplification, MinION library
129 preparation and sequencing were conducted in the field at the Ulu Temburong National Park (Brunei,
130 Borneo) in October 2018, during a Taxon Expedition (<https://taxonexpeditions.com/>). We analyzed
131 seven samples: two snails (Snail1 and Jap1) and five beetles (H36, H37, H42, H43 and Colen1). Two of
132 them (H42, H43) were collected in an emergence trap [10] in which the specimens were exposed to a
133 preserving agent consisting of ethanol (~65%) glycerol (~30%), water (~5%) and a little amount of dish-
134 washing detergent for several days.

135 Total genomic DNA was isolated using the DNeasy Blood and Tissue Kit (Qiagen) from a 1×1 mm
136 biopsy of snail tissue or from the whole beetle after cutting the thorax and abdomen. Samples were
137 incubated in ATL lysis buffer for 2 h at 56 °C and overnight at room temperature before DNA was
138 extracted according to the manufacturer's instructions and eluted in Tris-EDTA buffer (10 mM Tris, 1
139 mM EDTA, pH 8.0).

140 Barcoding PCR was conducted by amplifying the mitochondrial gene encoding cytochrome
141 oxidase I (COI) using a MiniONE portable PCR device (MiniOne), lyophilized oligonucleotides and
142 PCR reagents previously kept at room temperature. We used the universal primers LCO1490 and
143 HC02198 [11] tailed with adaptors to allow indexing prior to MinION library preparation: 5'-TTT CTG
144 TTG GTG CTG ATA TTG CGG TCA ACA AAT CAT AAA GAT ATT GG-3' and 5'-ACT TGC CTG
145 TCG CTC TAT CTT CTA AAC TTC AGG GTG ACC AAA AAA TCA-3'. Each PCR (total volume 25
146 µl) comprised 2 µl of the DNA template, 0.25 µM of each primer, 0.25 mM of each dNTP, 1× Herculase
147 II reaction buffer, and 0.25 µl (20 U/µl) of Herculase II fusion DNA polymerase (Agilent Technologies).
148 The amplification profile consisted of an initial denaturation step (3 min at 95 °C) followed by 35 cycles
149 of 30 s at 95 °C, 30 s at 52 °C and 60 s at 72 °C, and a final extension for 5 min at 72 °C. PCR products
150 were verified by electrophoretic analysis (MiniOne Electrophoresis System, MiniOne) for the presence
151 of unique bands at the expected size (~700bp). The amplification of H37 and Colen1 was not successful,
152 so these samples were amplified using primers LepF1 (5'-TTT CTG TTG GTG CTG ATA TTG CAT
153 TCA ACC AAT CAT AAA GAT ATT GG-3') and LepR1 (5'-ACT TGC CTG TCG CTC TAT CTT CTA
154 AAC TTC TGG ATG TCC AAA AAA TCA-3') [12] using the reagents described above. The
155 amplification profile consisted of an initial denaturation step (1 min at 95 °C) followed by six cycles of
156 1 min at 95 °C, 90 s at 45°C and 75 s at 72 °C, then 36 cycles of 1 min at 95 °C, 90 s at 51°C and 75 s at 72
157 °C and a final extension for 5 min at 72 °C. PCR products were purified using 1.5X AMPureXP beads
158 (Beckman Coulter) and quantified using a Qubit 2.0 fluorimeter and the Qubit dsDNA BR assay kit
159 (Thermo Fisher Scientific).

160 To incorporate index sequences and allow the sequencing of multiple samples in each MinION
161 flow cell, a second round of PCR was carried out using 48 µl of the purified COI-PCR amplicons from
162 the first round (0.5 nM), 2 µl of indexed primers provided in the EXP-PBC001 kit (ONT), 0.25 mM of
163 each dNTP, 1× Herculase II reaction buffer, and 1 µl (20 U/µl) of Herculase II fusion DNA polymerase.

164 The amplification profile consisted of an initial denaturation step (3 min at 95 °C) followed by 15 cycles
165 of 15 s at 95 °C, 15 s at 62 °C and 30 s at 72 °C, and a final extension for 3 min at 72°C. Indexed PCR
166 products were purified using 0.8X AMPureXP beads (Beckman Coulter), quantified as described above
167 and pooled in equimolar concentrations.

168

169 **2.3 MinION library preparation and sequencing**

170 We used 1 µg of pooled amplicons to prepare sequencing libraries with the SQK-LSK108 DNA
171 Sequencing kit (ONT) according to the manufacturer's instructions (but omitting the DNA
172 fragmentation step). The library was loaded on a FLO-MIN106 flow cell (R9.4 sequencing chemistry).
173 Sequencing was carried out for 7 h in the field using MinKNOW v1.6.11 (ONT) on a portable laptop.

174

175 **2.4 Sanger sequencing**

176 Sanger sequencing was performed on COI PCR products prepared as described above and
177 purified using 1X AMPureXP beads. Sequencing was carried out at the BMR Genomics facilities in
178 Padova (Italy) or at the Museum für Naturkunde of Berlin (Germany), following our return from the
179 field expedition. Forward and reverse Sanger reads were assembled into a consensus sequence using
180 Geneious Prime v2019.0.4 (<http://www.geneious.com/>).

181

182 **2.5 Bioinformatic analysis of MinION reads**

183 After MinION sequencing, raw fast5 reads were basecalled and demultiplexed using Guppy
184 v2.3.7+e041753. To reduce the number of misassignments, a second round of demultiplexing was
185 performed requiring tags at both ends of reads using Porechop v0.2.3_seqan2.1.1
186 (<https://github.com/rrwick/Porechop>). Tags and adapters were trimmed using Porechop and reads of
187 abnormal length were filtered out using a custom script.

188 Starting from pre-processed MinION reads, the *ONTrack* pipeline consisted of the following steps.
189 First, VSEARCH v2.4.4_linux_x86_64 [13] was used to cluster reads at 70% identity and only reads in
190 the most abundant cluster were retained for subsequent analysis. Next, 200 reads were randomly
191 sampled using Seqtk sample v1.3-r106 (<https://github.com/lh3/seqtk>) and aligned using MAFFT v7.407
192 with parameters --localpair --maxiterate 1000, specific for iterative refinement, incorporating local
193 pairwise alignment information [14]. EMBOSS cons v6.6.6.0 ([http://emboss.open-](http://emboss.open-bio.org/rel/dev/apps/cons.html)
194 [bio.org/rel/dev/apps/cons.html](http://emboss.open-bio.org/rel/dev/apps/cons.html)) was then used to retrieve a draft consensus sequence starting from the
195 MAFFT alignment. The EMBOSS cons plurality parameter was set to the value obtained by multiplying
196 the number of aligned reads by 0.15, in order to include a base in the draft consensus sequence if at least
197 15% of the aligned reads carried that base. If less than 15% of the aligned reads carried the same base in
198 a specific position, and a generic base (N) was included in the consensus sequence, the generic base was
199 removed using a custom script. To polish the obtained consensus sequence, 200 reads were randomly
200 sampled using Seqtk sample, with a different seed to the one used before, and mapped to the draft
201 consensus sequence using Minimap2 v2.1.1-r341 [15]. The alignment file was filtered, sorted and
202 compressed to the *bam* format using Samtools v1.7 [16]. Nanopolish v0.11.0
203 (<https://github.com/jts/nanopolish>) was used to obtain a polished consensus sequence. When the
204 *ONTrack* pipeline was run multiple times, the polished consensus sequences produced during each
205 round were aligned with MAFFT, after setting the gap penalty to 0. The final consensus was retrieved
206 using EMBOSS cons based on the majority rule, namely including a base in the final consensus if it was
207 included in at least 50% of the iterations. PCR primers were trimmed from both sides of the consensus
208 sequence using Seqtk trimfq. As a final step, the consensus sequences were aligned using Blast v2.2.28+
209 against the NCBI nt database, which was downloaded locally. Seeds for subsampling reads in the three
210 iterations reported in the results were 1, 3 and 5 in the draft consensus step, and 2, 4 and 6 for the
211 polishing step, respectively. The accuracy of MinION consensus sequences was evaluated by aligning
212 the *ONTrack* consensus sequence to the corresponding Sanger-derived reference sequence using Blast
213 v2.2.28+ [17]. The accuracy of MinION reads was evaluated by aligning them to the corresponding
214 Sanger reference sequence using Minimap2 and running Samtools stats on the generated *bam* file.

215 All scripts were run within an Oracle Virtualbox v5.1.26 virtual machine emulating an Ubuntu
216 operating system on a Windows laptop without using any internet connection, and are available at
217 <https://github.com/MaestSi/ONTrack.git>. MinION-based consensus sequences and Sanger consensus
218 sequences are available as Supplementary Materials.

219
220 Sanger, MinION and consensus sequences are available at GenBank under the BioProject
221 PRJNA539982.
222

223 3. Results

224

225 3.1 COI barcode sequencing

226 To perform barcode sequencing in the field, the portable genomics laboratory we previously
227 described [1] was optimized further to include equipment and reagents with greater stability and
228 better performance in tropical environments (up to 35°C and 90% humidity) after transport on
229 standard domestic and international flights. Currently, the laboratory comprises seven portable
230 devices that can be fitted in one standard luggage item with dimensions of 55×45×20 cm (**Figure 1**).

231 After collecting two snails and five insects during a workshop held by Taxon Expeditions
232 (<https://taxonexpeditions.com/>) at the Ulu Temburong National Park (Borneo, Brunei) in October
233 2018, we dissected the tissue and extracted DNA. PCR products obtained by amplifying ~710 bp of
234 the COI gene were sequenced in the field using the MinION device with R9.4 sequencing chemistry.
235 The MinION flow cell showed 995 active pores during the pre-run quality control (starting from 1005
236 on delivery by the manufacturer) and produced 600,000 reads in 3.5 h. Raw fast5 reads were
237 basecalled, demultiplexed and trimmed offline, resulting in 9,000–77,000 reads per sample (**Table 1**).
238 When we returned to Europe, the same genomic fragments were amplified and sequenced from the
239 same DNA extracts using the Sanger method to evaluate the accuracy of the MinION-based
240 barcoding pipeline.

241 3.2 Barcode analysis using the *ONTrack* pipeline

242 The MinION reads were processed using *ONTrack*, a barcoding pipeline that we developed
243 using several samples collected over the last few years (**Figure 2**). The first step of the pipeline
244 involved clustering the reads to remove non-specific PCR products and nuclear mitochondrial DNA
245 segments (NUMTs), which can cause barcoding issues particularly when processing insect samples
246 [18,19]. We then randomly sampled 200 of the filtered reads and aligned them to produce a draft
247 consensus sequence. Starting from the draft consensus sequence, a polishing step was performed
248 using another set of 200 randomly sampled reads.

249 Despite the errors characterizing MinION reads (**Table 1**), the barcodes reconstructed using the
250 *ONTrack* pipeline had an average accuracy of 99.94% compared to the Sanger reference sequence. No
251 consistent differences were observed between the two distinct types of COI amplicons we analyzed
252 or the type of starting samples (**Table 2**).

253 The generated consensus sequences were finally used as BLAST queries against the NCBI nt
254 database, and the top hits for each sample were saved to a text file for operator analysis. Because the
255 database was downloaded locally, the whole pipeline from sequencing to the generation of consensus
256 sequences and the identification of BLAST top-hits could be completed without an internet
257 connection, which was in any case unavailable in the field on our expedition.

258 We found that, when running the *ONTrack* pipeline three times for the same sample, the results
259 differed slightly each time with an average accuracy ranging from 99.91% to 99.95%, depending on
260 the read group subsampled in each analysis (**Table 3**). The pipeline was therefore run iteratively by
261 aligning the consensus sequences generated during each round and extracting the ultimate consensus
262 sequence. This slightly increased the accuracy of our barcoding pipeline, removing errors present in
263 only one of the three iterations and thus achieving an average accuracy of 99.95%. The residual errors
264 were only present in homopolymer runs of at least 6 nt, although some homopolymer runs of 7 nt
265 were correctly reconstructed (**Figure 3**). The computational running time scaled linearly with the
266 number of iterations, making it feasible to perform three iterations in a reasonable amount of time
267 (~30 min per sample) on a standard laptop.

268

269 **3.3. Figures, Tables and Schemes**

270 **Tables**

271 **Table 1. Sequencing statistics.** For each sample, we show the COI primers used for PCR amplification, the
272 number of sequenced reads, the mean and the standard deviation of read length in base pairs, and the average
273 accuracy of MinION reads.

<i>Sample ID</i>	<i>Sample name</i>	<i>COI amplicon primers</i>	<i>Reads</i>	<i>Mean read length (sd)</i>	<i>Average read accuracy</i>
BC01	<i>Snail1</i>	LCO1490-HC02198	26,240	682 (16)	88.94%
BC02	<i>Jap1</i>	LCO1490-HC02198	68,822	681 (15)	87.95%
BC03	<i>H36</i>	LCO1490-HC02198	21,378	680 (17)	88.31%
BC04	<i>H37</i>	<i>LepF1 - LepR1</i>	21,115	564 (210)	86.74%
BC05	<i>H42</i>	LCO1490-HC02198	55,334	681 (15)	88.02%
BC06	<i>H43</i>	LCO1490-HC02198	76,680	683 (19)	87.13%
BC07	<i>Colen1</i>	<i>LepF1 - LepR1</i>	8,880	477 (231)	88.01%

274

275 **Table 2. Accuracy of consensus sequences generated by the ONTrack pipeline.** For each sample, we show
276 the mean percentage accuracy of the consensus sequences obtained.

<i>Sample ID</i>	<i>Consensus accuracy</i>
BC01	99.90%
BC02	100%
BC03	99.90%
BC04	100%
BC05	99.95%
BC06	99.89%
BC07	99.94%

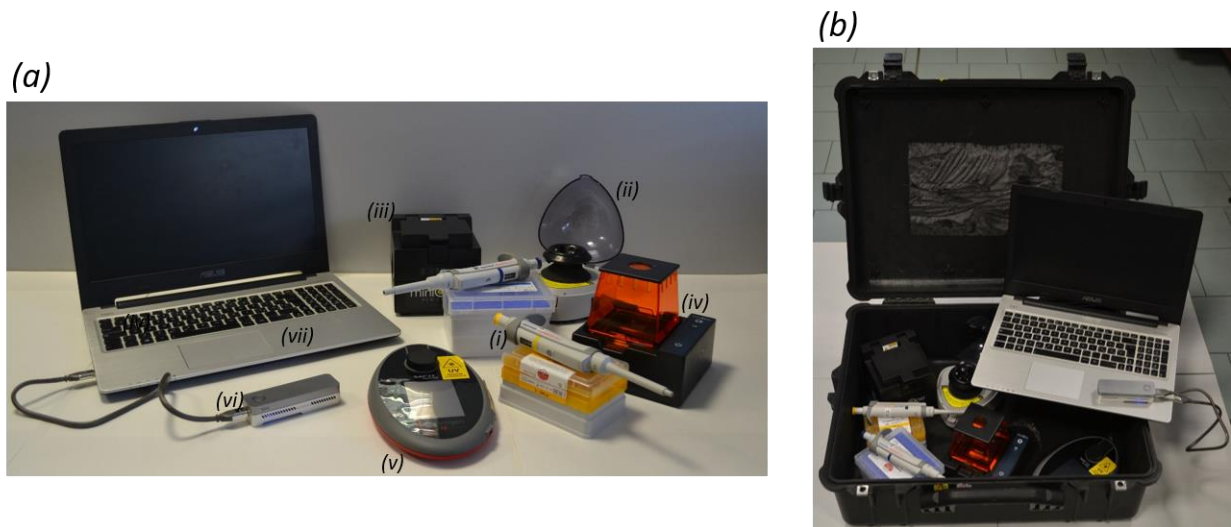
287 **Table 3. Accuracy of consensus sequences generated by combining three iterations of the ONTrack**
288 **pipeline.** For each sample, we show the number of properly reconstructed positions divided by the alignment
289 length and (in parentheses) the percentage accuracy of the consensus sequences for each of the three iterations,
290 the final consensus accuracy and the number of iterations supporting it.

<i>Sample ID</i>	<i>Consensus accuracy read set 1</i>	<i>Consensus accuracy read set 2</i>	<i>Consensus accuracy read set 3</i>	<i>Final consensus accuracy</i>	<i>Iterations supporting the final consensus</i>
BC01	650/651 (99.85%)	651/651 (100%)	650/651 (99.85%)	650/651 (99.85%)	2/3
BC02	656/656 (100%)	657/657 (100%)	657/657 (100%)	657/657 (100%)	3/3
BC03	647/64 (100%)	646/647 (99.85%)	646/647 (99.85%)	647/647 (100%)	1/3
BC04	606/606 (100%)	606/606 (100%)	606/606 (100%)	606/606 (100%)	3/3
BC05	656/656 (100%)	656/656 (100%)	657/658 (99.85%)	656/656 (100%)	2/3
BC06	576/576 (100%)	575/576 (99.83%)	574/575 (99.83%)	575/576 (99.83%)	2/3
BC07	535/536 (99.81%)	536/536 (100%)	536/536 (100%)	536/536 (100%)	2/3

291

292

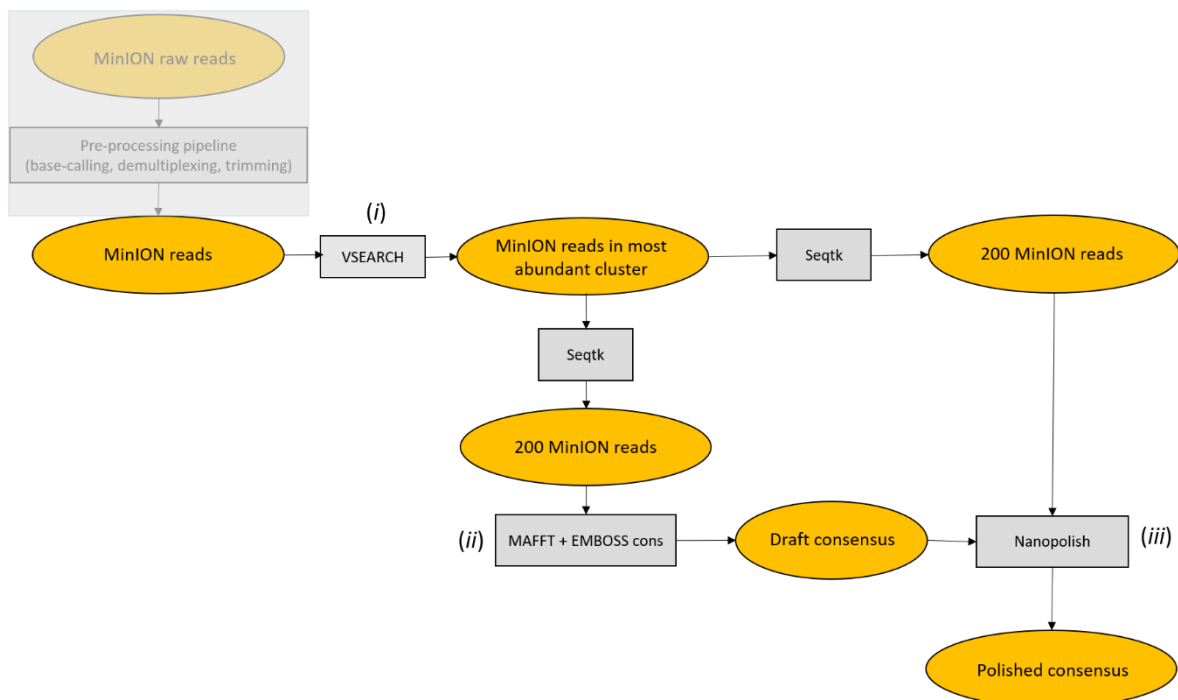
293 **Figures**



294

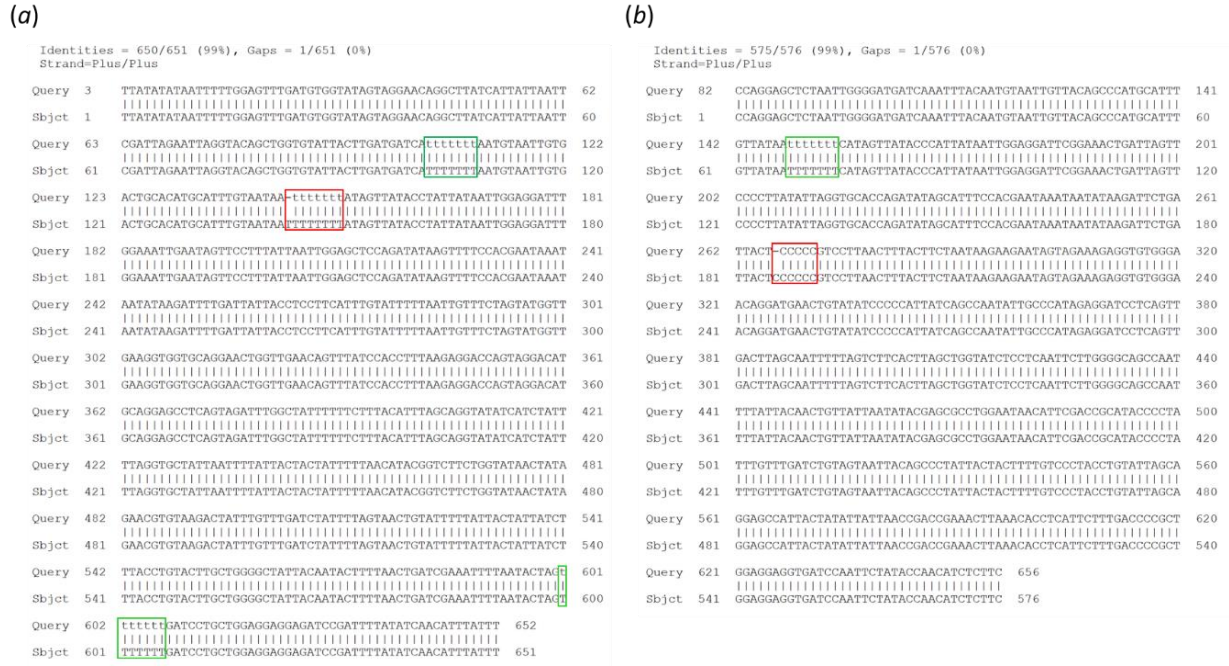
295 **Figure 1. The portable genomics laboratory.** Panel (a) shows the equipment comprising the portable
296 genomics laboratory, namely (i) micropipettes, (ii) a mini-microcentrifuge, (iii) a thermal cycler, (iv) an
297 electrophoresis system, (v) a fluorometer, (vi) the nanopore sequencer MinION, and (vii) a laptop. Panel (b)
298 shows how the laboratory is transported.

299



300

301 **Figure 2: ONTrack pipeline flowchart.** (i) MinION reads are clustered at 70% identity using VSEARCH
302 and only reads in the most abundant cluster are retained for subsequent analysis. (ii) Next, 200 reads are then
303 subsampled by Seqtk, aligned with MAFFT and a draft consensus is extracted with EMBOSS cons. (iii) The
304 draft consensus sequence is then polished using Nanopolish, based on a second set of 200 randomly sampled
305 reads.



306

307 **Figure 3. Analysis of residual errors in the ONTrack final consensus sequences.** Alignment of the
 308 MinION consensus sequence (Query) to the Sanger sequence (Sbjct) is shown for samples BC01 (a) and BC06
 309 (b). The residual errors, present in homopolymer runs of 6 and 8 nt, are highlighted in red. Properly
 310 reconstructed homopolymers of 7 nt are highlighted in green.
 311

312 4. Discussion

313 We have described the implementation of a new workflow for barcoding in the field, from DNA
314 extraction to the generation of consensus sequences. The selected protocols allowed the extraction of
315 DNA from tiny snail-tissue biopsies and from whole beetles after cutting the abdomen to release soft
316 tissues, as required to preserve the integrity of the specimens for detailed morphological evaluation.
317 PCR products were successfully obtained despite the transport of our equipment in a standard Peli
318 case and the storage of molecular biology reagents in local fridges and freezers powered for only 10 h
319 per day. The MinION flow cells, which were not adversely affected by the transportation and storage
320 conditions, retained most of their active pores and produced a good number of reads in a few hours.
321 These results indicate that the molecular biology field laboratory workflow was robust, allowing us
322 to barcode organisms at the collection site even under adverse environmental conditions (in this case
323 a rainforest characterized by high temperatures and humidity).

324 On the software side, the new bioinformatics pipeline allowed us to analyze MinION reads using
325 open-source and custom-developed scripts that run locally on a Linux Virtual Machine. The
326 sequencing and data analysis could therefore be combined on a standard Windows laptop with a
327 user-friendly interface. Most importantly, the improvements addressed some of the weaknesses of
328 earlier pipelines, such as their dependence on sequence databases and Q-score calibration. The
329 *ONTrack* pipeline works with as few as ~500 reads per sample and achieves high accuracy when
330 applied to MinION sequencing data obtained from COI barcode amplicons. Moreover, starting from
331 processed MinION reads, the *ONTrack* pipeline returns consensus sequences in a few minutes,
332 making it particularly suitable for work in the field.

333 The residual error rate in our consensus sequences never exceeded ~0.2%. The proposed
334 workflow can therefore be considered as a powerful tool for species identification given that most
335 species pairs show sequence divergence exceeding 2% [7]. Further improvements may be achieved
336 thanks to the software and chemistry enhancements regularly provided by ONT. A new flip-flop
337 basecalling algorithm (<https://github.com/nanoporetech/flappie>) was recently implemented in the
338 Guppy production basecaller and it should further reduce the error rate, albeit at the expense of
339 basecalling time. A new sequencing chemistry (R10) will be released soon, increasing the accuracy
340 especially in homopolymer runs and thus bringing on-site sequencing ever closer to the quality of
341 Sanger analysis.

342 Sequencing and basecalling currently remain the most time-consuming steps in the pipeline, but
343 both the hardware and software solutions provided by ONT are likely to become much more agile in
344 the near future. Indeed, ONT recently released MinIT, a rapid analysis and device-control accessory
345 for nanopore sequencing that connects to the MinION sequencer and performs GPU-accelerated and
346 real-time basecalling. Moreover, the Medaka tool (<https://github.com/nanoporetech/medaka>) is
347 expected to create polished consensus sequences faster than Nanopolish because it starts from
348 basecalled data rather than raw signals. Finally, new MinION flow cells (Flongle) were recently made
349 available and these are suitable for experiments that do not require a massive throughput, thus
350 substantially reducing sequencing costs for small datasets. Because the *ONTrack* pipeline provides
351 high-quality results with as few as ~500 reads per sample (0.35 Mbp), multiple samples could be
352 multiplexed in a single run and still fit Flongle specifications (1 Gbp) further reducing the cost.
353 Considering a multiplex of 12 samples in a Flongle run, currently the maximum supported by
354 standard ONT kits, we estimated a cost of about 30 € per sample to generate a barcode sequence with
355 the workflow described herein. This is not far from the costs of standard Sanger sequencing (~15 €
356 per sample when sequencing both strands, without considering the extra shipment costs).
357 Remarkably, the entire portable genomics laboratory described in this article can be acquired with a
358 modest budget of 6000 €, compared to ~80,000 € for a Sanger sequencer (ABI capillary). Dedicated,
359 expert personnel are required to run the latter instrument, whereas the MinION sequencer is very
360 simple and requires no special training. An additional significant advantage is that, unlike other
361 sequencing technologies, the real-time MinION device does not require the number of sequenced
362 reads to be set before the experiment begins. Therefore, the sequencing run can be stopped at any

363 time when the necessary number of reads has been generated, achieving further cost and time
364 savings.
365

366 **Author Contributions:** Conceptualization, M.D. and M.S.; methodology, S.M., E.C., M.M., M.R. and M.D.;
367 software, S.M.; validation, S.M. and E.C.; formal analysis, S.M. and E.C.; investigation, E.C., M.R., M.P., H.F.,
368 M.A., I.N., L.M., M.S., F.S., J.G.; writing—original draft preparation, S.M., M.R. and M.D.; writing—review and
369 editing, M.R. and M.D.; visualization, S.M.; supervision, M.R. and M.D.; project administration, M.R. and M.D.;
370 funding acquisition, M.D., M.S. and I.N.

371 **Funding:** This research received no external funding

372 **Acknowledgments:** We gratefully acknowledge the Ulu Temburong National Park (Brunei, Borneo) for
373 permission to conduct research in the field; export of biological materials was done under permit
374 BioRIC/HoB/TAD/51 from the Ministry of Primary Resources and Tourism, Brunei Darussalam. We thank
375 Davide Canevazzi for the support in bioinformatic analysis.

376 **Conflicts of Interest:** The authors declare no conflict of interest

377 **References**

- 378 1. Menegon M; Cantaloni C; Rodriguez-Prieto A; Centomo C; Abdelfattah A; Rossato M; Bernardi M;
379 Xumerle L; Loader S; Delledonne M. On site DNA barcoding by nanopore sequencing. *PLoS ONE* **2017**,
380 12.
- 381 2. Pomerantz A; Peñafiel N; Arteaga A; Bustamante L; Pichardo F; Coloma LA; Barrio-Amorós CL;
382 Salazar-Valenzuela D; Prost S. Real-time DNA barcoding in a rainforest using nanopore sequencing:
383 opportunities for rapid biodiversity assessments and local capacity building. *Gigascience* **2018**, 7.
- 384 3. Faria NR; Quick J; Claro IM; Thézé J; de Jesus JG; Giovanetti M; Kraemer MUG; Hill SC; Black A; da
385 Costa AC, et al. Establishment and cryptic transmission of Zika virus in Brazil and the Americas. *Nature*
386 **2017**, 546, 406-410.
- 387 4. Quick J; Loman N; Duraffour S; Simpson JT; Severi E; Cowley L; Bore JA; Koundouno R; Dudas G;
388 Mikhail A, et al. Real-time, portable genome sequencing for Ebola surveillance. *Nature* **2016**, 530, 228-
389 232.
- 390 5. Edwards A; Debbonaire AR; Nicholls SM; Rassner SME; Sattler B; Cook JM; Davy T; Soares AR; Mur
391 LAJ; Hodson AJ. In-field metagenome and 16S rRNA gene amplicon nanopore sequencing robustly
392 characterize glacier microbiota. *bioRxiv* **2019**, doi: <https://doi.org/10.1101/073965>.
- 393 6. Srivathsan A; Baloğlu B; Wang W; Tan WX; Bertrand D; Ng AHQ; Boey EJH; Koh JY; Nagarajan N;
394 Meier R. A MinION™-based pipeline for fast and cost-effective DNA barcoding. *Mol Ecol Resour* **2018**.
- 395 7. Hebert PDN; Ratnasingham S; deWaard JR. Barcoding animal life: cytochrome c oxidase subunit 1
396 divergences among closely related species. *Proc Biol Sci* **2003**, 270, S96-S99.
- 397 8. Freitag H; Kodaka J. A taxonomic review of the genus *Ancyronix* Erichson, 1847 from Sulawesi (Insecta:
398 Coleoptera: Elmidae). *Journal of Natural History* **2017**, 51, 561-606.
- 399 9. Krehenwinkel H; Pomerantz A; Henderson JB; Kennedy SR; Lim JY; Swamy V; Shoobridge JD; Patel
400 NH; Gillespie RG; Prost S. Nanopore sequencing of long ribosomal DNA amplicons enables portable
401 and simple biodiversity assessments with high phylogenetic resolution across broad taxonomic scale.
402 *Gigascience* **2019**, doi: 10.1093/gigascience/giz006.
- 403 10. Freitag H. Adaptation of an Emergence Trap for Use in Tropical Streams. *International Review of*
404 *Hydrobiology* **2004**, 89, 363-374.

- 405 11. Folmer O; Black M; Hoeh W; Lutz R; Vrijenhoek R. DNA primers for amplification of mitochondrial
406 cytochrome c oxidase subunit I from diverse metazoan invertebrates. *Mol Mar Biol Biotechnol.* **1994**, *3*,
407 294-299.
- 408 12. Hebert PDN; Penton EH; Burns J; Janzen DH; Hallwachs W. Ten species in one: DNA barcoding reveals
409 cryptic species in the neotropical skipper butterfly, *Astrartes fulgerator*. *Proc Nat Acad Sci USA* **2004**,
410 *101*, 14812–14817.
- 411 13. Rognes T; Flouri T; Nichols B; Quince C; Mahé F. VSEARCH: a versatile open source tool for
412 metagenomics. *PeerJ* **2016**, *4*.
- 413 14. Katoh K; Misawa K; Kuma K; Miyata T. MAFFT: a novel method for rapid multiple sequence alignment
414 based on fast Fourier transform. *Nucleic Acids Res.* **2002**, *30*, 3059–3066.
- 415 15. Li H. Minimap2: pairwise alignment for nucleotide sequences. *Bioinformatics* **2018**, *191*.
- 416 16. Li H; Handsaker B; Wysoker A; Fennell T; Ruan J; Homer N; Marth G; Abecasis G; Durbin R; 1000
417 Genome Project Data Processing Subgroup. The Sequence alignment/map (SAM) format and
418 SAMtools. *Bioinformatics* **2009**, *25*, 2078-2079.
- 419 17. Altschul SF; Gish W; Miller W; Myers EW; Lipman DJ. Basic local alignment search tool. *J. Mol. Biol*
420 **1990**, *215*, 403-410.
- 421 18. Kang AR; Kim MJ; Park IA; Kim KY; Kim I. Extent and divergence of heteroplasmy of the DNA
422 barcoding region in *Anapodisma miramae* (Orthoptera: Acrididae). *Mitochondrial DNA A DNA Mapp*
423 *Seq Anal.* **2016**, *27*, 3405-3414.
- 424 19. Meza-Lázaro RN; Poteaux C; Bayona-Vásquez NJ; Branstetter MG; Zaldívar-Riverón A. Extensive
425 mitochondrial heteroplasmy in the neotropical ants of the *Ectatomma ruidum* complex (Formicidae:
426 Ectatomminae). *Mitochondrial DNA A DNA Mapp Seq Anal* **2018**, *29*, 1203-1214.
427