**1D MinION sequencing for large-scale species discovery: 7000 scuttle flies (Diptera: Phoridae) from one site in Kibale National Park (Uganda) revealed to belong to >650 species**

**Amrita Srivathsan[1], Emily Hartop[2,5,6], Jayanthi Puniamoorthy[1], Wan Ting Lee[1], Sujatha Narayanan Kutty[1,4], Olavi Kurina[3], Rudolf Meier[1,4*]**

[1] Department of Biological Sciences, National University of Singapore, 14 Science Drive 4

[2] Zoology Department, Stockholms Universitet, Stockholm, Sweden

[3] Estonian University of Life Sciences, Kreutzwaldi 5D, Tartu, Estonia

[4] Tropical Marine Science Institute, National University of Singapore, Singapore

[5] Station Linné, Öland, Sweden

[6] Naturhistoriska Riksmuseet, Stockholm, Sweden

*Corresponding author: Rudolf Meier: meier@nus.edu.sg

Other email addresses:

Amrita Srivathsan: asrivathsan@gmail.com

Emily Hartop:  emhartop@gmail.com

Jayanthi Puniamoorthy: dbsjay@nus.edu.sg

Wan Ting Lee:  dbslwt@nus.edu.sg

Sujatha Narayanan Kutty:  tmssunk@nus.edu.sg

Olavi Kurina: Olavi.Kurina@emu.ee

**Keywords**: NGS barcoding, DNA barcoding, Nanopore sequencing, MinION, large-scale species discovery

1  **ABSTRACT**

2

3  Background: Most animal species remain to be discovered. We recently proposed to tackle this

4  problem using a 'reverse workflow' where all specimens are barcoded via tagged amplicon

5  sequencing and then sorted into putative species (mOTUs). We furthermore suggested that the

6  COI barcodes can be obtained with minimal laboratory equipment using MinION sequencing, but

7  our test with $1D^2$ reads only yielded 500 barcodes per flowcell.

8

9  Results: Here we show how MinION 1D sequencing can be used to obtain ~3500 COI barcodes

10  per flowcell. Based on 7062 MinION barcodes for a hyper-diverse family of flies (Diptera:

11  Phoridae) collected by one trap in Kibale National Park, Uganda, we discover ~650 species which

12  exceeds the number of phorid species described for the entire Afrotropical region. Our updated

13  MinION pipeline increases processing speed via parallelization, improves demultiplexing, and yet

14  yields reliable barcodes (99.99% accuracy) and similar mOTUs as Illumina sequencing (match

15  ratio: 0.989). Morphological examination of 100 mOTUs confirms good congruence (93% of

16  mOTUs; >99% of specimens). Nearly 90% of species and specimens belong to the megadiverse

17  genus *Megaselia* which is routinely neglected because its species diversity and abundance is too

18  overwhelming. We show that it can be tackled with the reverse workflow. We also illustrate how

19  the molecular data guides the description of a new species: *Megaselia sepsioides* sp. nov..

20

21  Conclusions: MinION is suitable for reliable, rapid, and large-scale species discovery in

22  hyperdiverse taxa. Approximately 3500 specimens can be sequenced using one MinION flowcell

23  at a barcode cost of <0.35 USD.

24

25

## INTRODUCTION

Most life on earth has yet to be discovered with an estimated 80% of extant species being still unknown to science [1]. The majority of these species belong to hyper-diverse and species-rich invertebrate clades. These taxa are ubiquitous, contain most of the multicellular animal species, and often occur in great abundance. However, new species are difficult to find and delimit because it requires the study of thousands of specimens. Typically, this process starts with sampling specimens with bulk trapping methods (e.g. Malaise trap, fogging, pitfall traps, flight intercept traps). It usually yields thousands of specimens per site that need to be sorted: first to higher-level taxonomic groups by parataxonomists and then to species-level by taxonomic experts. The latter work has to be carried out by taxonomic experts because species-level sorting by parataxonomists tends to yield unreliable results [2]. Once morpho-species have been obtained, they are then often tested via DNA barcodes (658 bp fragment of COI) by sequencing a few representative specimens for each morpho-species [3]. This traditional workflow works well for taxa with small numbers of species and specimens but is so time-consuming for hyperdiverse and abundant clades that they are neglected. This is partially responsible for our lack of baseline data for many insect taxa. The traditional workflow has the additional downside that morphologically cryptic species are overlooked. This is particularly likely to happen when expensive Sanger sequencing is used because only few specimens can be barcoded and the probability of detecting cryptic species is low [4].

An alternative approach to species discovery is the 'reverse workflow' where every specimen in a sample is individually sequenced (or barcoded) without the destruction of the specimen [4-6]. The specimens are then grouped to molecular Operational Taxonomic Units (mOTUs) based on DNA barcodes. The morphological check of the putative species delimited with DNA barcodes comes last. The taxonomic expert works on pre-sorted material and rectifies mis-sorted specimens, identifies known species, and describes new species. This would have been deemed

52  unrealistically expensive prior to the advent of the High Throughput Sequencing technologies.

53  However, sequencing platforms like Illumina and PacBio are sufficiently cost-effective and are

54  now replacing barcoding via expensive Sanger sequencing [4, 5, 7-10]. For example, sequencing

55  tens of thousands of specimens with Illumina HiSeq can cost as little as 0.17 USD per specimen

56  (including PCR cost, see discussion in Wang et al., 2018 [4]). However, Illumina and PacBio

57  sequencing have some downsides. They are only cost-effective if >10,000 specimens have to be

58  barcoded, sequencing usually has to be outsourced (i.e., the amplicons have to be shipped to a

59  sequencing facility), and it often takes weeks to obtain data. It would be desirable to have

60  alternatives that are fast, scalable, and yet cost-effective. This would be particularly useful if

61  barcoding has to be accomplished under field conditions or in countries with limited access to

62  Illumina and PacBio sequencing [4, 5, 11, 12]. This is frequently the case and we therefore

63  strongly believe that the democratization of large-scale DNA barcoding will be important for

64  upscaling species discovery across the globe and encouraging the use of the reverse workflow

65  across many labs.

66

67  Oxford Nanopore's MinION has the potential to help with achieving these goals. It is a low-cost

68  and portable real-time sequencing device. However, its use and reliability for large scale

69  specimen handling remains to be fully understood. We recently showed that 500 reliable DNA

70  barcodes can be obtained using $1D^2$ sequencing on one flowcell of MinION. To our knowledge

71  this was the highest number of products that were successfully multiplexed in a single MinION

72  flowcell (5X higher than a study which examined ~100 amplicons of rDNA [11]) but such low

73  throughput still meant that the cost per barcode remained high (ca. 2 USD: [13]). While this scale

74  is useful for many species' identification projects, it is unlikely to be effective for large-scale

75  species discovery where samples can contain thousands of specimens. Furthermore, the 500

76  specimens were barcoded using $1D^2$ sequencing which requires a complicated library

77  preparation, base-calling is computationally intensive, and application to amplicon sequencing is

78    still under development. Unfortunately, it remained untested whether the more straightforward,

79    but less accurate 1D sequencing can be used for large-scale species discovery. This is addressed

80    in this manuscript.

81

82    However, we are here not only developing ways to use of MinION 1D sequencing for barcoding.

83    Instead, we are also investigating the curious observation that only 466 species of phorid flies

84    have been recorded for the Afrotropical Region [14]. Phoridae is a hyper-diverse clade belonging

85    to the true flies (Diptera). Diptera is one of several hyper-diverse insect orders that also include

86    beetles (Coleoptera), bees, wasps, and ants (Hymenoptera), and moths and butterflies

87    (Lepidoptera). The species estimates for all of Insecta vary between 3 and 13 million (reviewed

88    by Stork, 2018 [15]) with only ca. 1,000,000 currently being described [16]. Historically, the

89    *inordinate fondness* of taxonomists for beetles has led to Coleoptera outpacing Diptera and

90    Hymenoptera in numbers of described species. However, several recent studies suggest that

91    Hymenoptera and Diptera are likely to be more species-rich. For example, Forbes et al. [17]

92    hypothesize that Hymenoptera contained more species than either Diptera or Coleoptera based

93    on parasite host ratios for Microhymenoptera, but this study showed an underappreciation for

94    both the great numbers of Dipteran parasitoids and the diversity of true flies in general. Indeed,

95    Diptera has recently been proven to be surprisingly rich in a number of large-scale biodiversity

96    studies. In a large barcoding study of Canadian insects, Hebert et al. [18] found that Hymenoptera

97    and Diptera together accounted for two thirds of the 46,937 BINS acquired and predicted that one

98    dipteran family (Cecidomyiidae) has 16,000 species in Canada. The authors then extrapolated to

99    the worldwide fauna which they estimated to be 10 million insect species, of which 1.8 million

100   were predicted to be cecidomyiids [18]; i.e., a single family of Diptera may surpass the number of

101   described species in all of Coleoptera. Other studies similarly hint at the extraordinary richness of

102   Diptera. The Zurqui All Diptera Biodiversity Inventory (ZADBI) of a single site in Costa Rica was

103   heavily reliant on specimens collected with two Malaise traps run for a one-year period [19]. Only

104    41,001 specimens (a small fraction of the hundreds of thousands collected) could be studied by

105    taxonomic experts [20]. The specimens that were examined revealed 4,332 species of Diptera,

106    of which 800 were from Cecidomyiidae and 404 were for Phoridae [20], the fly family of focus

107    here.

108

109    Phoridae, or scuttle flies, are a worldwide family of true flies with approximately 4300 described

110    species [14]. Over 1600 of these are in the giant genus *Megaselia* Rondani, which has been

111    described as "one of the largest, most biologically diverse and taxonomically difficult genera in

112    the entire animal kingdom" [21]. In groups like *Megaselia*, the species discovery problem appears

113    insurmountable. Extremely large numbers of specimens are routinely collected, and they can

114    belong to very large numbers of species. Even in urban and suburban habitats, the diversity of

115    the family can be surprisingly high. Henry Disney, a world expert on the family, has recorded 75

116    species of phorids (48 of *Megaselia*) in his modest suburban garden [22]. Similarly, the BioSCAN

117    project in Los Angeles found that backyards in the city supported as many as 82 species [22]. In

118    natural areas, the diversity and abundance tend to be much higher and sorting such samples into

119    species-level units using traditional workflows is very labor-intensive. Rare and new species are

120    often hidden among very large numbers of common and described species. The study of groups

121    like *Megaselia* requires examination of thousands of specimens for which prodigious notes have

122    to be taken as specimens are compared. Many detailed drawings are prepared (for *Megaselia*

123    drawings of male genitalia are essential) – often based on dissections and slide mounts – because

124    many known/common species cannot be identified without detailed inspection. The process can

125    take hours, days, or longer for a single species discovery. This traditional workflow thus often

126    discourages all but the most tenacious taxonomists from taking up the study of hyper-diverse

127    genera within insects.

128

129      Here, we test whether 1D MinION sequencing may offer a rapid and revolutionary approach to

130      exploring phorid diversity more comprehensively. MinION sequencing is here applied to ca. 30%

131      of the phorid specimens that were collected in a single Malaise trap in Kibale National Park,

132      Uganda. We here describe how we processed ~8700 specimens, obtained ~7000 accurate

133      barcodes, and found 650 species (of which almost 90% are *Megaselia*). All this could be

134      accomplished using a workflow that requires less than a month.

135

136      **RESULTS**

137      **MinION based DNA barcoding**

138      The experiment was designed to obtain full-length COI barcodes for two sets of specimens via

139      tagged amplicon sequencing. A total of 8699 specimens were processed (Set 1: 4275; Set 2:

140      4519; 95 specimens were shared between the sets) (Fig. 1). In order to assess amplification

141      success rates, a subset of PCR products for each of the ninety-two 96-well plates were assessed

142      with agarose gels. The extrapolated amplification success rates were 86% and 74% for the two

143      sets of specimens (80.7% overall); i.e., we estimated that >3600 and >3300 DNA barcodes should

144      be obtainable via MinION sequencing given that gels tend to underestimate amplification success

145      rates (Table 1). The PCR products for each set were pooled and sequenced using MinION (set

146      1: 7,035,075; set 2: 7,179,121 1D nanopore reads). Both sets were sequenced in two sequencing

147      runs. The first run for each set was based on the pooled PCR products for all specimens in the

148      sets. It generated 3,069,048 and 4,853,363 reads, respectively. The results of the first run were

149      used to estimate coverage for each PCR product. Products with weak coverage (<=50x) were re-

150      pooled and re-sequenced (set 1: 2172 amplicons; set 2: 2211 amplicons). This added 3,966,027

151      and 2,325,758 reads to each set and improved the coverage of many low-coverage barcodes

152      (Fig. 2). The combined data were processed using an improved version of the bioinformatics

153      pipeline in Srivathsan et al. [13]. The improvements led to a higher demultiplexing rate (14%

154  increase for set 1: 898,979 vs. 787,239 reads; 9% increase for set 2: 647,152 vs. 593,131 reads)

155  and faster demultiplexing (10X using 4 cores: demultiplexing in 9 min vs 87 min for one of the

156  datasets).

157  Demultiplexing of all data and preliminary barcode calling revealed 3,797 and 3,476 MAFFT

158  barcodes with >=5X coverage and <1% ambiguous bases. These barcodes were subject to

159  correction using RACON [23] which yielded the same number of barcodes. When the barcodes

160  for the two sets of samples were combined, we overall obtained 7,220 MAFFT and RACON

161  barcodes. These preliminary barcodes still contain indel and substitution errors that were

162  addressed with an amino-acid correction pipeline that was first implemented in Srivathsan et al.

163  [13]. It yielded 7,178 MAFFT+AA and 7,194 RACON+AA barcodes. Some barcodes were not

164  retained because this pipeline rejects barcodes that have five or more consecutive indel errors.

165  Finally, the two sets of corrected barcodes were consolidated. This yielded a set of 7,115

166  consolidated barcodes. We rejected barcodes where the alignment of MAFFT+AA and

167  RACON+AA barcodes required the insertion of indels, as the +AA barcodes are expected to be

168  indel-free. Such indels in the alignments of MAFFT+AA and RACON+AA barcodes also indicate

169  discrepancies between MAFFT and RACON barcode estimates and there is no objective reason

170  to prefer one barcode over the other. The overall barcoding success rate was thus 81.9% (7,115

171  barcodes for 8,699 specimens). This was close to the expected 80.7% success rate based on gel

172  electrophoresis; i.e., MinION sequencing consistently produced sequence data for successfully

173  amplified products. A subsequent contamination check via BLAST revealed that of the 7,115

174  barcodes, 53 barcodes were unlikely to pertain to phorid flies (<1%) and we thus retained 7,062

175  barcodes for species richness estimation. Lastly, we inspected the reads obtained for the 92

176  negative controls (1 per microplate). Five negatives yielded MAFFT barcodes. Four of these had

177  a >97% match to non-phorids (two humans, one fish, one mollusc) and were eliminated. One low

178  coverage (13X) negative survived all filters and matched phorid COI. It was removed after

179   ascertaining that it did not impact the accuracy of the barcodes in the plate. This could be tested

180   by comparing the MinION barcodes for the plate with Illumina barcodes obtained from different

181   PCR products for the same DNA extraction plate (see below).

182   **Accuracy and selection of barcode sets**

183   To find the best strategy for obtaining accurate barcodes, we compared 5 sets of barcodes

184   (MAFFT, RACON, MAFFT+AA, RACON+AA, and consolidated barcodes) with the corresponding

185   barcodes based on Illumina sequencing. Illumina barcodes were obtained for 6,373 specimens

186   for the same specimens using different primers that amplified a 313 bp subset of the full-length

187   barcodes. The comparisons showed that the uncorrected MAFFT and RACON barcodes had an

188   accuracy of 99.61% and 99.51% (Table 2). Correction of these barcodes using the amino-acid

189   correction pipeline improved the accuracy considerably (>99.9% in all cases). The barcodes were

190   corrected after testing several "namino" parameters. Overall, namino=2 was found to yield the

191   most accurate barcodes, and minimized the number of inaccurate barcodes. We found that

192   MAFFT+AA barcodes were more accurate than RACON+AA barcodes, but MAFFT+AA barcodes

193   contained a much higher number of ambiguous nucleotides (Fig. 3). When RACON+AA and

194   MAFFT+AA barcodes were consolidated, the resulting "consolidated barcodes" were found to be

195   highly accurate (99.99%) and contained few ambiguous bases (median = 0.3 %).

196   We furthermore compared the mOTU richness estimated by the different barcode sets. mOTU

197   richness was very similar across MinION and Illumina barcodes for the consolidated and

198   uncorrected MAFFT/RACON barcodes. MAFFT+AA barcodes performed well in this comparison,

199   but yielded fewer mOTUs than Illumina barcodes: this may be due to a higher proportion of

200   ambiguous nucleotides (Fig. 3). However, comparison of mOTU richness alone does not imply

201   the same specimens were grouped into mOTUs across MinION and Illumina barcode sets. We

202   thus also calculated the match ratio for the datasets (3% clustering threshold). We found that all

203   five barcode sets (MAFFT, RACON, MAFFT+AA, RACON+AA, and consolidated barcodes,

204  namino=2) also had high match ratios (>0.97) with the consolidated barcodes and RACON

205  barcodes performing best with match ratios of >0.98 (consolidated barcodes: 0.989, RACON:

206  0.989, RACON+AA: 0.982). However, upon closer inspection only the multiple sequence

207  alignment (MSA) of the consolidated barcodes was indel-free while the other MSAs contained

208  indels. The largest number of indels was found in the MSA of uncorrected RACON barcodes

209  which indicates that the RACON barcodes retain a fair number of indels and may not be of

210  sufficient quality for submission to sequence databases. Overall, we would thus recommend the

211  usage of consolidated barcodes as the final barcode set. Because they maximize the per-base

212  accuracy, estimated mOTU diversity, match ratios, and yield high-quality alignments.

213  **Species richness estimation**

214  We thus proceeded to characterize the diversity of the phorid flies collected from the Malaise traps

215  based on the consolidated barcodes (namino=2). We overall obtained a mean of 683 mOTUs

216  (2%: 728, 3%: 685, 4%: 636) when the thresholds were varied from 2-4%. Species accumulation

217  and Chao 1 curves for mOTUs at 3% were not found to have reached a plateau, but the shape of

218  the curves suggest an estimated diversity is >1000 species in a single site collected by one

219  Malaise trap (Fig. 4).

220  **Congruence with morphology**

221  mOTUs are expected to be affected by mistakes caused by lab contamination and species

222  delimitation errors due to the biological properties of barcodes. We find that 6 of the 100 clusters

223  contained a single misplaced specimen and that there was one small cluster of four specimens

224  that appeared to consist of a mixture three morpho-species. This implies that 9 of the >1500

225  examined barcoded specimens were misplaced due to lab contamination. mOTUs based on

226  barcodes are expected to underestimate species for those that recently speciated and

227  overestimate species with deep splits [24]. This means that taxonomists working with mOTUs

228  should check for signs of lumping and splitting for closely related taxa. Our initial morphological

229  check of 100 randomly selected clusters (>1500 specimens) took ca. 30 hours. This preliminary

230  morphological screening covered all specimens while future studies could concentrate on

231  complex clusters, as any differences that may have occurred within 5% clusters were not

232  discernible without further preparations (dissection and slide mounting). Identifying these areas

233  of potential ambiguity is an advantage of this pipeline. It allows taxonomic experts to focus time

234  and energy on these complex clusters.

235

236  **New Species Description**

237

238  A primary aim of the reverse workflow is finding rare new species for description in bulk samples.

239  With specimens pre-sorted into mOTUs, morphologists have easy access to interesting and

240  potentially rare species that would otherwise remain undiscovered. While examining the 100

241  mOTUs, eight specimens were found to belong to a distinctive new species of *Megaselia*. This

242  species provided a good opportunity to demonstrate how the reverse workflow aids in species

243  discovery because a mOTU-specific haplotype network informed on which specimens should be

244  studied with morphology. The species is here described, and the description incorporates the

245  molecular data.

246

247  *Megaselia sepsioides* Hartop sp. n.

248  urn:lsid:zoobank.org:pub:ED268DF2-A886-4C31-A4FB-6271C382DECE

249

250  **Description**

251  See Fig. 5, 6.

252   In an effort to continue reducing redundancy and ambiguity in species descriptions, the

253   description of this species has excluded the character table from the method previously

254   established for *Megaselia* [25-27] and uses a barcode and photographic description. Photographs

255   are a key element in descriptions for large, diverse groups [28], where verbose descriptions can

256   be both overwhelmingly time consuming and insufficiently diagnostic. Most characters that would

257   have been in table form are clearly visible in the photographs provided; the few that are not were

258   deemed irrelevant to either the description or diagnosis of the species.

259

260   **Diagnosis**

261   This species is unmistakable even within the gargantuan and taxonomically difficult

262   genus *Megaselia*. The semi-circular expansion with modified peg-like setae on the forefemur is

263   unique among described members of the genus (Fig. 5, b). Similarly, the severe constriction of

264   the hind tibia basally is diagnostic (Fig. 5, d and e). The narrow and elongate form of the abdomen

265   is notable. Fig. 6 shows variations in setation between haplotypes. Only single specimens of the

266   two distinct haplotypes are available; more specimens will be necessary to determine if these are

267   eventually removed as distinct species or fall within a continuum of intraspecific variation.

268   DNA barcode for UGC0005996:

269   actttatattttattttggagcttgagctggaatagtaggtacttccttaagaatcataattcgtgctgaattaggacacccaggagcacttat

270   tggtgatgaccaaatttataatgtgattgttactgcacatgctttattataattttttttatagtaatacctattataataggaggttttggtaattg

271   acttgtacctttaatattaggagccccagatatggcattccctcgaatgaataatataagttttgaatattacctccttctttaactcttttatta

272   gccagaagtatagtagaaaatggagctggaactggttgaacagtttatcctcctttatcttctagaatcgctcatagtggagcttctgttgat

273   ttagcaattttctctcttcatttagctggaatttcatctattttaggagctgtaaattttattacaacaattattaatatacgatcatcaggtattaca

274   tttgaccgaatacctctatttgtttgatctgtaggtattacagctttattgctactcttatcacttcctgtttttagctggtgctattacaatactattaa

275   cagaccgaaattttaatacttcattttttgacccagcaggaggaggagatccaattttataccaacatttattc

276

277   **Material Examined**

278

279 Holotype:

280 Scientific Name: *Megaselia sepsioides* Hartop 2019; country: Uganda; state Province:

281 Kamwenge; locality: Kibale National Park: 1530 m; Coordinates: 00°33'54.2"N 30°21'31.3"E;

282 sampling protocol: Malaise trap; event date: iii-xii, 2010; individual count: 1; sex: male; life stage:

283 adult: UGC0005996; identified by: Emily Hartop: 2019; institution code: LKCNHM; collection code:

284 UGC; basis of record: preserved specimen.

285

286 Paratypes:

287 Scientific name: *Megaselia sepsioides* Hartop 2019; country: Uganda; state province:

288 Kamwenge; locality: Kibale National Park; elevation: 1530 m; coordinates: 00°33'54.2"N

289 30°21'31.3"E; sampling protocol: Malaise trap; event date: iii-xii, 2010; individual count: 7; sex:

290 male; life stage: adult; catalog number: UGC0012899, UGC0012244, UGC0012568,

291 UGC0003003, UGC0005864, UGC0012937, UGC0012971; identified by: Emily Hartop; date

292 identified: 2019; institution code: LKCNHM; collection code: UGC; basis of record: preserved

293 specimen

294

295 **Distribution**

296 Known from a single site in Kibale National Park, Uganda.

297

298 **Biology**

299 Unknown.

300

301 **Etymology**

302 Named by Yuchen Ang for the sepsid-like (Diptera: Sepsidae) foreleg modification.

303

**DISCUSSION**

*Large scale species discovery using MinION*

Our results suggest that MinION's 1D sequencing yields data of sufficient quality for producing high-quality DNA barcodes that can be used for large-scale species discovery. Through the development of new primer-tags, sequencing strategies, and improved bioinformatics procedures, we here increase the barcoding capacity of a MinION flowcell from 500 specimens (Srivathsan et al., 2018: $1D^2$ sequencing) to ~3500 specimens. This is achieved without a drop in accuracy because the new error correction pipeline is effective at eliminating most of the errors in the 1D reads (ca. 10%). Indeed, even the initial MinION barcodes (MAFFT & RACON) have very high accuracy (>99.5%) when compared to Illumina data. Note that this accuracy is even higher than what was obtained with $1D^2$ sequencing in Srivathsan et al. (2018: 99.2%). We suspect that this partially due to improvements in MinION sequencing chemistry and base-calling, but our upgraded bioinformatics pipeline also helps because it increases coverage for the amplicons. These findings are welcome news because 1D library preparations are much simpler than the library preps for $1D^2$. In addition, $1D^2$ reads are currently less suitable for amplicon sequencing (https://store.nanoporetech.com/kits-250/1d-sequencing-kit.html) [29].

We tested a range of different techniques for obtaining barcodes from 1D reads. Some techniques are computationally more expensive than others which raises the question of which bioinformatics pipeline should be used for obtaining accurate barcodes. Based on the current performance of MinION, we would recommend the usage of the "consolidated barcodes". They contain no indel errors and fewer substitution errors (99.99% accuracy) when compared to MAFFT+AA and RACON+AA barcodes. The mOTUs delimited with consolidated barcodes are virtually identical with the ones obtained via Illumina sequencing (Number of 3% mOTUs for Illumina: 661; MinION:

330  658; match-ratio: 0.989). In addition, consolidated barcodes were also the most accurate

331  barcodes based on $1D^2$ reads [13]. Nevertheless, we found that all barcode sets gave reliable

332  mOTU estimates. For corrected barcodes (MAFFT+AA and RACON+AA) this was consistent with

333  >99.9% per base accuracy, but the higher error rates of the uncorrected barcodes rarely affected

334  mOTU estimates. This is because indels were treated as missing data and mOTU estimation only

335  requires accurate estimates of distances for closely related taxa. This explains why the 99.5%

336  accurate RACON barcodes can group specimens into mOTUs in a very similar manner to what

337  is obtained with Illumina barcodes (match ratio of >0.98). The results imply that accurate mOTU

338  estimates can be obtained even based on preliminary barcodes.

339

340  Based on the procedures described here, MinION barcodes can be generated rapidly and at a

341  low sequencing cost of <0.35 USD per barcode. These properties make MinION a valuable tool

342  for species discovery whenever a few thousand specimens have to be sorted to species (<5000).

343  Even larger-scale barcoding is probably still best tackled with Illumina short-read or PacBio's

344  Sequel sequencing [4, 5, 7] because the cost is lower and the quality of the reads is higher.

345  However, both require access to expensive sequencers, the sequencing has to be outsourced,

346  and the users usually have to wait for several weeks in order to obtain the data. This is not the

347  case for barcoding with MinION where most of the data are collected within 10 hours of

348  sequencing. Our proposed MinION pipeline has the additional advantage that it only requires

349  basic molecular lab equipment including thermocyclers, magnetic rack, Qubit and potentially a

350  server computer. The latter is only needed for base-calling and should be replaceable by a new

351  portable, custom-built computational device for base-calling ONT data ("MinIT") while

352  demultiplexing and barcode-calling only requires a regular laptop computer. Overall, these

353  requirements are minimal which means that fully functional barcoding labs can be established for

354  <USD 5,000. Arguably, the biggest operational issue may be access to a sufficiently large number

355     of thermocyclers given that a study of the scale described here involved amplifying PCR products

356     in 92 microplates (=92 PCR runs).

357

358     Our new workflow for large-scale species discovery is based on sequencing the amplicons in two

359     sequencing runs. The second sequencing run can re-use the flowcell that was used for the first

360     run. Two runs are desirable because it improves overall barcoding success rates. The first run is

361     used to identify those PCR products with "weak" signal (=low coverage). These weak products

362     can then be re-sequenced in the second run. This dual-run strategy is designed to overcome the

363     challenges related to sequencing large numbers of PCR products: the quality and quantity of DNA

364     extracts are poorly controlled and PCR efficiency varies considerably. Pooling of products ideally

365     requires normalization, but this is not practical when thousands of samples are handled. Instead,

366     one can use the real-time sequencing provided by MinION to determine coverage and then boost

367     the coverage of low-coverage products by preparing and re-sequencing a separate library that

368     contains only the low coverage samples. Given that library preparations only require <200 ng of

369     DNA, even a set of weak amplicons will yield a sufficient amount of DNA.

370

371     *MinION sequencing and the "reverse workflow"*

372

373     Based on our data, we would argue that MinION is now suitable for implementing the "reverse-

374     workflow" where all specimens are sequenced first before mOTUs are assessed for morphological

375     consistency. This differs from the traditional workflow in that the latter requires sorting based on

376     morphology with only some morpho-species being subsequently tested via limited barcoding. We

377     would argue that the reverse-workflow is more suitable for handling species- and specimen-rich

378     clades because it requires less time than high-quality sorting based on morphology which often

379     involves genitalia preparations and slide-mounts. We would argue that expert sorting and

380  identification of material based on morphology is only very efficient for taxa where species-specific

381  morphological characters are easily accessible and the number of specimens is small. However,

382  most undiscovered and undescribed species are in poorly explored groups and regions where

383  specimen sorting and identification using the traditional techniques will be much slower. For

384  example, even if we assume that an expert can sort and identify 50 specimens of unknown phorids

385  per day, the reverse workflow pipeline would increase the species-level sorting rate by >10 times

386  (based on the extraction and PCR of six microplates per day). In addition, the sorting would be

387  carried out by lab personnel trained in amplicon sequencing while accurate morpho-species

388  sorting requires highly specialized taxonomic experts. Even such highly trained taxonomic experts

389  are often not able to match males and females of the same species (often one sex is ignored in

390  morphological sorting) while the matching of sexes and immatures is an automatic and desirable

391  byproduct of applying the reverse workflow [6].

392

393  One key element of the reverse workflow is that vials with specimens that have haplotype

394  distances <5% are physically kept together. This helps when taxonomic experts scrutinize

395  mOTUs for congruence with morpho-species. Indeed, graphical representation of haplotype

396  relationships (e.g., haplotype networks) can be used to guide morphological re-examination as

397  illustrated in our description *of Megaselia sepsioides* (Fig. 7). The eight specimens belonged to

398  seven haplotypes, the most distant haplotypes were dissected in order to test whether the data

399  are consistent with one or two species. Variations in setation were observed (Fig. 6) that were

400  deemed likely to be consistent with intraspecific variation. Note that the morphological

401  examination of clusters is straightforward because the use of QuickExtract for DNA extractions

402  ensures morphologically intact specimens. We predict that taxa that have been historically

403  ignored or understudied due to extreme abundance of common species and high species diversity

404  will now become more accessible because taxonomic experts can work on pre-sorted material

405  instead of thousands of unsorted specimens.

406

407 But how good is the correspondence between mOTUs and morpho-species. We checked 100

408 randomly chosen mOTUs for congruence with morphospecies. We find that 93% of clusters are

409 congruent, and over 99% of specimens (six of the seven cases of incongruence were single

410 specimens). This is in line with congruence levels that we observed for ants and odonates [4, 6].

411 This means that MinION barcodes are not only reliable but also yield mOTUs that are largely

412 congruent with morphospecies. This also means they are suitable for characterizing even

413 complex and species-rich samples of hyper-diverse invertebrate clades.

414

415 *Remarkably high diversity of Phoridae in Kibale National Park*

416 It is astonishing that the barcodes obtained from a single site in Kibale National Park (Uganda)

417 revealed ca. 650 mOTUs of phorid flies. This diversity, obtained from one Malaise trap, contains

418 150% of the number of described phorid species (466) known from the entire Afrotropical region;

419 i.e., numerically at least 184 species must be new to the region or to science [14]. Note that the

420 barcoded specimens only represent 8 one-week samples between March 2010 and February

421 2011; i.e., forty-four samples obtained from the same Malaise trap remain un-sequenced. We

422 thus expect the diversity from this single site to eventually well exceed 1000 species (Fig. 4).

423

424 The unexpectedly high species richness found in this study inspired us to muse about the species

425 diversity of phorids in the Afrotropical region. This is what Terry Erwin did when he famously

426 estimated a tropical arthropod fauna of 30 million species based on his explorations of beetle

427 diversity in Panama [30]. Such speculation is useful because it raises new questions, inspires

428 follow-up research, and may be needed given that even with extensive sampling, the species

429 richness of diverse taxa can be remarkably difficult to estimate [31]. The Afrotropical region

430 comprises roughly 2000 squares of 100 km$^2$ size. In our study we only sampled a tiny area within

431    one of these squares and observed 650 species which are likely to represent a species

432    assemblage that exceeds >1000 species. This will only be a subset of the phorid fauna in the

433    area because many specialist species (e.g., termite inquilines) are not collected in Malaise traps.

434    Of course, the 1000 species are also only a subset of the species occurring in the remaining

435    habitats in the same 100 km$^2$ square which is likely to be home to several thousand species of

436    phorids. But let's only assume that on average each of the two-thousand 100 km$^2$ squares have

437    100 endemic phorid species. If so, the "endemic" phorids alone would contribute 200,000 species

438    of phorids to the Afrotropical fauna without even considering the contributions to species diversity

439    by the "normal" species turnover of the remaining species. These considerations make us believe

440    that it would be somewhat surprising if the Afrotropical region were to have fewer than half a

441    million species of phorids! Based on our sample from Kibale National Park, 90% of species and

442    specimens would belong to *Megaselia* as currently circumscribed. Could it be that there are

443    450,000 species of Afrotropical *Megaselia*? These guestimates would only be much lower if the

444    vast majority of phorid species had very wide distributions which we consider somewhat unlikely

445    given that the Afrotropical region covers a wide variety of climates and habitats.

446

447    As documented by our study, the bulk of phorid species are *Megaselia*. Based on the traditional

448    workflow, almost all these specimens would be relegated to unsorted Malaise trap residues for

449    decades or centuries. Indeed, there are thousands of vials labeled as "Phoridae" decorating the

450    shelves of all major museums worldwide. The traditional workflow is not able to keep pace with

451    the species numbers and abundance. This makes rapid species-level sorting with "NGS

452    barcodes" [4] so important. Biologists will finally be able to work on taxa that are so specimen-

453    and species-rich that they were considered unworkable with the traditional techniques. MinION

454    barcodes will be one of the techniques that can be used to tackle these clades. We predict that

455    these barcodes will become particularly important for setting up mobile laboratories that can

456 operate under field conditions, but MinION barcodes can also be generated by highschool

457 students and citizen scientists.

458

459 **METHODS**

460 **1. Sampling**

461 Samples were collected from a single Townes-type Malaise trap [32], in the Kibale National Park,

462 close to Kanyawara Biological Station in the evergreen primeval forest at an altitude of 1513 m

463 (00°33'54.2"N 30°21'31.3"E) (Fig. 4). Kibale National Park is characterized as a fragment of

464 submontane equatorial woodland being home to 216 tree species [33]. Temperatures in Kibale

465 range from 16°C to 23°C (annual mean daily minimum and maximum, respectively) [34]. As

466 described, the Malaise trap was checked every week when the collecting bottle with the material

467 was replaced by a resident parataxonomist ([35]: Mr. Swaibu Katusabe). The material in Kibale

468 National Park in Uganda was collected and transferred in accordance with approvals from the

469 Uganda Wildlife Authority (UWA/FOD/33/02) and Uganda National Council for Science and

470 Technology (NS 290/ September 8, 2011), respectively. The material was thereafter sorted to

471 higher-level taxa. Target Diptera taxa were sorted to family and we here used Phoridae. The

472 sampling was done over several months between 2010 and 2011. For the study carried out here,

473 we only barcoded ca. 30% of the phorid specimens. The flies were stored in ethyl alcohol at -20–

474 25°C until extraction.

475 **2. DNA extraction**

476 DNA was extracted using 10 ul of QuickExtract in a 96 well plate format and the whole fly was

477 used to extract DNA. The solution with the fly was incubated at 65°C for 15 min followed by 98°C

478 for 2 min. No homogenization was carried to ensure that intact specimen was available for

479 morphological examination.

480

**3. MinION based DNA barcoding**

*I.       Polymerase Chain Reactions (PCRs)*

Each plate with 96 QuickExtract extracts (95 specimens and 1 control, with exception of one plate

with no negative and one partial plate) was subjected to PCR in order to amplify the 658 bp

fragment of COI using LCO1490 5' GGTCAACAAATCATAAAGATATTGG 3' and HCO2198 5'

TAAACTTCAGGGTGACCAAAAAATCA 3' [36]. Each PCR product was amplified using primers

that included a 13 bp tag. For this study, 96 thirteen-bp tags were newly generated in order to

allow for upscaling of barcoding; these tags allow for multiplexing >9200 products in a single

flowcell of MinION through unique tag combinations (96x96 combinations).  To obtain these 96

tags, we first generated one thousand tags that differed by at least 6 bp using with

BarcodeGenerator      (http://comailab.genomecenter.ucdavis.edu/index.php/Barcode_generator)

[37] However, tag distances of >6 bp are not sufficient because they do not take into account

MinION's propensity for generating homopolymer and other indel errors. We thus excluded tags

with homopolymeric stretches that were >2 bp long. We next used a custom script that identified

tags that differed from each other by indel errors and eliminated tags recursively to ensure that

the final sets of tags differed from each other by >=3bp errors of any type (any combination of

insertions/deletions/substitutions). This corresponded with our bioinformatic strategy of using

<=2bp mismatch in tag identification. Lastly, we excluded tags that ended with "GG" because

LCO1490 starts with this motif. The PCR conditions were as follows, reaction mix: 10 µl Mastermix

from CWBio, 0.16 µl of 25mM MgCl$_2$, 2 µl of 1 mg/ml BSA, 1 µl each of 10 µM primers, and 1ul of

DNA. The PCR conditions were 5 min initial denaturation at 94°C followed by 35 cycles of

denaturation at 94°C (30 sec), 45°C (1 min), 72°C (1 min), followed by final extension of 72°C

(5min). For each plate, a subset of 7-12 products were run on a 2% agarose gel to ensure that

PCRs were successful. Of the 96 plates studied, 4 plates were excluded from further analyses as

506 they had <50% amplification success and one plate was duplicated across the two runs by

507 accident.

508

509 *II.      MinION sequencing*

510 The most cost-effective strategy for nanopore sequencing was optimized during the study. For

511 the initial experiment (set 1) we sequenced amplicons for 4275 phorid flies. The flowcell was used

512 for 48 hours and yielded barcodes for ~3200 products, but we noticed lack of data for products

513 for which amplification bands could be observed on the agarose gel. We thus re-pooled all

514 products with a sequencing depth <=50X (2119 specimens), prepared a new library and

515 sequenced them on a new flowcell. The experiment was successful. However, in order reduce

516 sequencing cost, we pursued a different strategy for the second set of specimens. This set

517 consisted of pooled amplicons for 4519 flies, but here we stopped the sequencing on the flowcell

518 after 24 hours. The flowcell was then washed using ONT's flowcell wash kit and prepared for

519 reuse. The results from the first 24 hours of sequencing were used to identify amplicons with weak

520 coverage. They were re-pooled, a second library was prepared, and sequenced on the pre-used

521 and washed flowcell.

522

523 Amplicon pooling strategy: For set 1, all plates were grouped by amplicon strength as judged by

524 the intensity of products on agarose gels (5 strong pools +2 weak pools). For set 2, each plate

525 was pooled, quantified, and cleaned using either 1X Ampure beads or 1.1X Sera-Mag beads in

526 PEG. For the re-sequencing of weak or "problematic" products (see below), we identified the latter

527 based on the results of the initial sequencing run. We located (1) specimens <=10X coverage (set

528 1: 1054, set 2: 1054) and (2) samples with coverage between 10X and 50X (set 1: 1118, set 2:

529 1065). Lastly, we also created a (3) third pool of specimens with problematic products that were

530 defined as those that were found to be of low accuracy during comparisons with Illumina barcodes

531 and those that had high levels of ambiguous bases (>1% ambiguous bases during preliminary

532    barcode calling). Very few amplicons belonged to this category (set 1: 68, set 2: 92) and it is thus

533    not included in the flowchart in Figure 1. In order to efficiently re-pool hundreds of specimens

534    across plates we wrote a script that generates visual maps of the all microplates that illustrate

535    where the weak products are found.

536

537    Library preparation and sequencing: We used the SQK-LSK109 ligation sequencing kit for library

538    preparation and sequencing. Our first experiment on set 1 used 1 ug of starting DNA while all

539    other libraries used 200 ng pooled product. Library preparation was carried out as per

540    manufacturer's instructions with one exception: the various clean-up procedures at the end-prep

541    and ligation stages used 1X Ampure beads instead of 0.4 X as suggested in the instructions

542    because the amplicons in our experiments were short (~735 bp with primers and tags). The

543    sequencing was carried out using MinION sequencer with varying MinKNOW versions between

544    August 2018 - January 2019. Fast5 files generated were uploaded onto a computer server and

545    base-calling was carried out using Guppy 2.3.5+53a111f. No quality filtering criteria were used.

546    Our initial work with Albacore suggested that quality filtering improved demultiplexing rate but

547    overall more reads could be demultiplexed without the filtering criterion.

548

549    *III.    Data analyses for MinION barcoding*

550    We analysed the data using miniBarcoder [13], which was improved in several ways for the

551    present study. Overall, the pipeline starts with a primer search with *glsearch36,* then flanking

552    nucleotide sequences are identified, and reads are demultiplexed based on tags. For the latter,

553    an error of up to 2 bp errors are allowed. The demultiplexed reads are aligned using MAFFT v7

554    (--op 0) (here v7) [38]. In order to improve speed, we used only a random subset of 100 reads

555    from each demultiplexed file for alignment. Based on these alignments, a majority rule consensus

556    is called to obtain what we call "MAFFT barcodes". A second preliminary barcode is generated

557    by mapping all reads back to the MAFFT barcode using Graphmap (here v0.5.2) [39] and calling

558  the consensus using Racon (here, v1.3.1) [40]. This yields what we call "RACON barcodes". Both

559  MAFFT and RACON barcodes are subject to further correction based on publicly available

560  barcodes in GenBank. These corrections are advisable in order to fix remaining indel errors. The

561  correction takes advantage of the fact that COI sequences are translatable; i.e., an amino-acid

562  based error correction pipeline can be used (details can be found in Srivathsan et al. (2018).

563  Applying this pipeline to MAFFT and RACON barcodes respectively yields MAFFT+AA and

564  RACON+AA barcodes. Lastly, these barcodes can be consolidated by comparing them to yield

565  "consolidated barcodes".

566

567  The version of the pipeline in Srivathsan et al. [13] was modified as follows:

568

569  a. *Tackling 1D reads for massively multiplexed data:* The large number of simultaneously

570  barcoded specimens presented many challenges related to varying coverage and quality. We

571  hence sought to develop ways to account for increased error rates of 1D sequencing and develop

572  objective ways for quality assessments based on the MinION data and publically available data

573  (GenBank): (1) The GraphMap max error was increased from 0.05 to 0.15 to account for error

574  rates of 1D reads. (2) We modified the approach for calculating consolidated barcodes. We here

575  use the strict consensus of MAFFT+AA and RACON+AA barcodes in order to resolve conflicts

576  between MAFFT and RACON barcodes if there are substitution conflicts. In Srivathsan et al.

577  (2018) we accepted MAFFT+AA barcodes in cases of conflict, but for the 1D data we found that

578  MAFFT+AA barcodes had more ambiguities than RACON+AA barcodes which could be resolved

579  via calculating a strict consensus. We also introduced a criterion for identifying "problematic"

580  barcodes based on observing too many differences between MAFFT+AA and RACON+AA

581  barcodes; i.e., if indels are found in the alignment of MAFFT+AA and RACON+AA barcode during

582  consolidation, the barcode is rejected. This criterion is based on the fact that +AA barcodes should

583  be indel-free. Either the MAFFT+AA or the RACON+AA barcode is likely to be incorrect, but there

584    is no objective way to identify the correct sequence. (3) We assessed how different window sizes

585    can impact the amino-acid correction pipeline by varying the "namino" parameter (number of AA

586    to be inspected in each side of an indel).

587

588    *b. Demultiplexing rate*: (1) We introduced a "homopolymer compression" of the putative tag

589    sequences in order to improve demultiplexing rates. After primer searches, the old pipeline used

590    to identify the flanking 13 bps that were likely to be tag sequence. Instead we now use a 20 bp

591    flanking sequence and then compress any homopolymer >3bp before searching for 13 bp tags.

592    (2) We now analyze reads that are likely the result of two products that ligated during library

593    preparation. This was based on our experience with one library where ligation of products was

594    prominent. Long reads are split into size classes of <1300, >1300 and <2000. These settings

595    were set based on 658 bp barcode of COI: the total product size including tags and primers is

596    735 bp, and hence, a sequence with two products ligated to each other is expected to be 1470

597    bp long. The sequences were split in a manner that ensured that the tags of first product are not

598    affecting the tag found in the second, i.e., primer search for second product was conducted after

599    the first 650 bp of the sequence. Currently, this option is only available for reads that consist of

600    two ligated products.

601

602    *c. Processing speed and memory requirements*: (1) For primer identification we now limit the

603    searches to the first and last 100 bp of each read which allowed for improving speed of the primer

604    search.  (2) We parallelized the process of demultiplexing and MAFFT barcode calling using

605    multiprocessing in python. This allows for fast demultiplexing whenever computational power is

606    limited. (3) We optimized the pipeline with regard to memory consumption and ensured that all

607    processes are applied to batches of <=20000 sequences. The latter modification is needed given

608    the rapid increase in MinION output and batch processing is scalable with increased output.

609

**4. Illumina Based NGS Barcoding for validation**

In order to validate MinION barcodes and optimizing error correction strategies, we used reference COI barcodes obtained via Illumina sequencing for the same QuickExtract DNA extractions. Illumina sequencing was carried out for a 313 bp fragment of the same COI barcoding region using m1COlintF: 5'-GGWACWGGWTGAACWGTWTAYCCYCC-3' [41] and modified jgHCO2198: 50-TANACYTCNGGRTGNCCRAARAA YCA-3 [42]. We recently conducted an extensive analyses to understand if a 313 bp minibarcode is able to provide similar identifications and species delimitations as the 658 bp barcodes and found this to be the case when examining >5000 species [43]. Tagged primers were used for the PCRs as specified in Wang et al. (2018) [44]. The PCR mix was as follows: 4 µl Mastermix from CWBio, 1µl of 1 mg/ml BSA, 1µl of 10 µM of each primer, and 1ul of DNA. PCR conditions: 5 min initial denaturation at 94°C followed by 35 cycles of denaturation at 94°C (1 min), 47°C (2 min), 72°C (1 min), followed by final extension of 72°C (5 min). The PCR products were pooled and sequenced along with thousands of other specimens in a lane of HiSeq 2500 (250 bp PE sequencing). The data processing followed Wang et al. (2018): paired end reads were merged using PEAR (v 0.9.6) [45], reads were demultiplexed by an inhouse pipeline that looks for perfect tags while allowing for 2 bp mismatches in primer sequence. For each sample, demultiplexed reads are merged to identify the most dominant sequence, and a barcode is accepted only if this sequence is 5X as common as the next most common sequence.

**5. Assessment of MinION barcodes and mOTU delimitations**

Both MinION and Illumina barcodes were subject to contamination check. For MinION barcodes we used preliminary MAFFT barcodes given that this is the largest barcode set. Barcodes were matched to GenBank using BLASTN and taxonomic classifications were assigned using readsidentifier [46]. Any barcode with >95% match to a non-phorid sequence was excluded from the dataset. Furthermore, if any barcode best matched to Bacteria, it was also excluded.

635

MinION barcodes were assessed by scripts provided in the miniBarcoder package (assess_corrected_barcodes.py and assess_uncorrected_barcodes.py). For uncorrected barcodes, this was done by aligning barcodes to reference Illumina barcodes using dnadiff [47]. For corrected barcodes (+AA), we used MAFFT to obtain pairwise alignments. This allowed us to compare per base accuracy. Here we excluded those barcodes that had >3% difference between MinION and Illumina barcodes. These are likely due to wet lab errors (<0.5% of specimens). Such error is not entirely surprising given that the MinION and Illumina barcodes were generated using different amplicons.

We were further interested in understanding how mOTU delimitation is affected by error correction. Barcodes were aligned using MAFFT and MinION barcodes were further trimmed to the 313 bp region of the Illumina barcode. mOTU delimitation was done at 2, 3,and 4% using SpeciesIdentifier (objective clustering) [48]. mOTU richness was estimated and we furthermore calculated match ratio between two sets of clusters [49]. Match ratio is given by $\frac{2N_{match}}{N_1+N_2}$.

**6. Morphological examination**

For morphological examination of the clustered specimens we used 100 randomly selected non-singleton mOTUs delimited at 5% but also kept track of sub-clusters within the large mOTUs that were splitting at 1-5%. This allowed for examination of closely related, but potentially distinct, species. We were mostly interested in understanding if MinION barcodes were placing specimens into mOTUs incorrectly and hence we examined if specimens were consistent morphologically in each of these 5% clusters. The choice of 5% cluster may seem initially inconsistent with the choice of 3% for mOTU delimitation for other analyses, but examination of all specimens within 5% clusters allows for comparing multiple closely related 3% (or 1-5%) mOTUs. This often requires genitalia preparations which will be carried out at a larger scale once more specimens have been

659   sequenced. In this study, the process is only illustrated for the newly described species for which

660   we illustrate how the haplotype network obtained with the median joining method in PopART (Fig.

661   7) [50] guides morphological examination of specimens for the new species, *Megaselia*

662   *sepsioides* sp. nov.. Specimens with conspicuously large haplotype differences were dissected in

663   order to rule out the presence of multiple closely related species. Differences in setation were

664   observed between the two distant haplotypes (UGC0012899 and UGC001224) and the main

665   cluster    (UGC0003003,    UGC0005864,    UGC0005996,    UGC0006224,    UGC0012568,

666   UGC0012937, and UGC0012971) and are detailed in Fig. 6. It is deemed likely that with further

667   collection of specimens from this haplotype cluster, these now rather distinct differences will be

668   absorbed into a continuum of intraspecific variation. Specimen examination was done with a Leica

669   m80 and Leica M205 C stereo microscopes and the images were obtained with a Dun Inc.

670   Microscope Macrophotography system (Canon 7D chassis with 10X Mitutoyo lens). Photography

671   stacking was done with Zerene stacker. Specimens were cleared with clove oil and mounted on

672   slides in Canada balsam following the protocol of Disney [51].

673   **7. Species Richness estimation**

674   The most accurate set of barcodes were used for assessing the overall diversity of the barcoded

675   phorids.  mOTU delimitation was based on SpeciesIdentifier [48]. The 3% mOTUs were used for

676   estimating species richness. Here we used EstimateS9 [52], and changed the diversity settings

677   to use the classical formula of Chao1 and Chao2. This is because coefficient of variation of the

678   abundance or incidence distribution was >0.5.

679

680   **AVAILABILITY OF SOURCE CODE AND REQUIREMENTS**

681   - Project name: miniBarcoder

682   - Project home page: https://github.com/asrivathsan/miniBarcoder

- Operating system(s): Linux, MacOSX (for initial barcode calling)

- Programming language: Python

- Other requirements: MAFFT BARCODE: MAFFT, glsearch36, RACON BARCODE: GraphMap, Racon, and amino acid correction: BioPython, MAFFT

- License: GNU GPL

- New scripts included: mb_parallel_consensus.py mb_parallel_demultiplex.py, repool_by_plate.py

- Updated scripts: consolidate.py, bug fix in aacorrection.py

**AVAILABILITY OF SUPPORTING DATA AND MATERIALS**

Besides the source code, the nanopore data and input files and barcode sets will be made available once FTP link is obtained.

**DECLARATIONS**

**Author's Contributions:** R.M. and A.S. conceived the workflow and the analytical approach. O.K. conducted/organized the sampling and Diptera sorting and commented on the manuscript. Molecular work was conducted by J.P., W.T.L., S.N.K., E.H..and A.S. Pipeline development and data analyses was conducted by A.S. Morphological examination, and species description was conducted by E.H. Figures were prepared by E.H., S.N.K., A.S. Manuscript was written by R.M., A.S. and E.H..

713

714 **REFERENCES**

715 1.     Wilson EO. Biodiversity research requires more boots on the ground. Nat Ecol Evol.

716        2017;1 11:1590-1. doi:10.1038/s41559-017-0360-y.

717 2.     Krell F-T. Parataxonomy vs. taxonomy in biodiversity studies – pitfalls and applicability of

718        'morphospecies' sorting. Biodiversity and Conservation. 2004;13 4:795-812.

719        doi:10.1023/b:Bioc.0000011727.53780.63.

720 3.     Hebert PD, Cywinska A, Ball SL and deWaard JR. Biological identifications through DNA

721        barcodes. Proceedings of Royal Society B. 2003;270 1512:313-21.

722 4.     Wang WY, Srivathsan A, Foo M, Yamane SK and Meier R. Sorting specimen-rich

723        invertebrate samples with cost-effective NGS barcodes: Validating a reverse workflow

724        for specimen processing. Molecular Ecology Resources. 2018;18 3:490-501.

725 5.     Meier R, Wong W, Srivathsan A and Foo M. $1 DNA barcodes for reconstructing

726        complex phenomes and finding rare species in specimen-rich samples. Cladistics.

727        2016;32 1:100-10.

728 6.     Yeo D, Puniamoorthy J, Ngiam RWJ and Meier R. Towards holomorphology in

729        entomology: rapid and cost-effective adult–larva matching using NGS barcodes.

730        Systematic Entomology. 2018;43 4:678-91.

731 7. Hebert PDN, Braukmann TWA, Prosser SWJ, Ratnasingham S, deWaard JR, Ivanova

732   NV, et al. A Sequel to Sanger: amplicon sequencing that scales. BMC Genomics.

733   2018;19:219.

734 8. Shokralla S, Porter TM, Gibson JF, Dobosz R, Janzen DH, Hallwachs W, et al.

735   Massively parallel multiplex DNA sequencing for specimen identification using an

736   Illumina MiSeq platform. Scientific Reports. 2015;5:9687.

737 9. Creedy TJ, Norman H, Tang CQ, Chim KQ, Andujar C, Arribas P, et al. A validated

738   workflow for rapid taxonomic assignment and monitoring of a national fauna of bees

739   (Apiformes) using high throughput barcoding. BioRxiv. 2019;  doi:10.1101/575308.

740 10. Krehenwinkel H, Kennedy SR, Rueda A, Lam A and Gillespie RG. Scaling up DNA

741   barcoding – Primer sets for simple and cost efficient arthropod systematics by multiplex

742   PCR and Illumina amplicon sequencing. 9. 2018;11 2181-2193.

743 11. Krehenwinkel H, Pomerantz A, Henderson JB, Kennedy SR, Lim JY, Swamy V, et al.

744   Nanopore sequencing of long ribosomal DNA amplicons enables portable and simple

745   biodiversity assessments with high phylogenetic resolution across broad taxonomic

746   scale. GigaScience. 2019;  doi:10.1093/gigascience/giz006.

747 12. Krehenwinkel H, Wolf M, Lim JY, Rominger AJ, Simison WB and Gillespie RG.

748   Estimating and mitigating amplification bias in qualitative and quantitative arthropod

749   metabarcoding. Scientific Reports. 2017;7 17668.

750 13. Srivathsan A, Baloglu B, Wang W, Tan WX, Bertrand D, Ng AHQ, et al. A MinION™-

751   based pipeline for fast and cost-effective DNA barcoding. Molecular Ecology Resources.

752   2018;18 5:1035-49.

753 14. Brown BV: Phorid Catalog. http://phorid.net/pcat/. Accessed 27 April 2019.

754 15. Stork NE. How Many Species of Insects and Other Terrestrial Arthropods Are There on

755   Earth? Annu Rev Entomol. 2018;63:31-45. doi:10.1146/annurev-ento-020117-043348.

756    16.    Zhang ZQ. Animal biodiversity: An introduction to higher-level classification and

757          taxnomic richness. Magnolia Press; 2011.

758    17.    Forbes AA, Bagley RK, Beer MA, Hippee AC and Widmayer HA. Quantifying the

759          unquantifiable: why Hymenoptera – not Coleoptera – is the most speciose animal order.

760          BMC Ecology. 2018;18 21 doi:10.1101/274431.

761    18.    Hebert PD, Ratnasingham S, Zakharov EV, Telfer AC, Levesque-Beaudin V, Milton MA,

762          et al. Counting animal species with DNA barcodes: Canadian insects. Philos Trans R

763          Soc Lond B Biol Sci. 2016;371 1702 doi:10.1098/rstb.2015.0333.

764    19.    Borkent ART, Brown BV, Adler PH, Amorim DDS, Barber K, Bickel D, et al. Remarkable

765          fly (Diptera) diversity in a patch of Costa Rican cloud forest: Why inventory is a vital

766          science. Zootaxa. 2018;4402 1 doi:10.11646/zootaxa.4402.1.3.

767    20.    Brown BV, Borkent A, Adler PH, De Souza Amorim D, Barber K, Bickel D, et al.

768          Comprehensive inventory of true flies (Diptera) at a tropical site. Communications

769          Biology. 2018;1 21:8. doi:10.1038/s42003-018-0022-x.

770    21.    Marshall SA. Flies: The Natural History and Diversity of Diptera. Buffalo, New York:

771          Firefly Books; 2012.

772    22.    Brown BV and Hartop EA. Big data from tiny flies: patterns revealed from over 42,000

773          phorid flies (Insecta: Diptera: Phoridae) collected over one year in Los Angeles,

774          California, USA. Urban Ecosystems. 2016;  doi:10.1007/s11252-016-0612-7.

775    23.    Vaser R, Sovic I, Nagarajan N and Sikic M. Fast and accurate de novo genome

776          assembly from long uncorrected reads. Genome Res. 2017;27 5:737-46.

777          doi:10.1101/gr.214270.116.

778    24.    Kwong S, Srivathsan A, Vaidya G and Meier R. Is the COI barcoding gene involved in

779          speciation through intergenomic conflict?. Molecular Phylogenetics and Evolution.

780          2012;62 3:1009.

781  25.  Hartop EA and Brown BV. The tip of the iceberg: a distinctive new spotted-wing

782       *Megaselia* species (Diptera: Phoridae) from a tropical cloud forest survey and a new,

783       streamlined method for *Megaselia* descriptions. Biodiversity Data Journal. 2014;2:e4093.

784       doi:10.3897/BDJ.2.e4093.

785  26.  Hartop EA, Brown BV and Disney RHL. Opportunity in our ignorance: urban biodiversity

786       study reveals 30 new species and one new Nearctic record for *Megaselia* (Diptera:

787       Phoridae) in Los Angeles (California, USA). Zootaxa. 2015;3941:451-84.

788  27.  Hartop EA, Brown BV and Disney RHL. Flies from L.A., The Sequel: Twelve further new

789       species of *Megaselia* (Diptera: Phoridae) from the BioSCAN Project in Los Angeles

790       (California, USA). Biodiversity Data Journal. 2016, p. e7756.

791  28.  Riedel A, Sagata K, Surbakti S, Rene T and Michael B. One hundred and one new

792       species of *Trigonopterus* weevils from New Guinea. Zookeys. 2013; 280:1-150.

793       doi:10.3897/zookeys.280.3906.

794  29.  ONT: https://store.nanoporetech.com/kits-250/1d-sequencing-kit.html. Accessed 28 April

795       2019.

796  30.  Erwin TL. Tropical Forests: Their Richness in Coleoptera and Other Arthropod Species.

797       The Coleopterists Bulletin. 1982;36 1:74-5.

798  31.  Longino JT, Coddington J and Colwell RK. The ant fauna of a tropical rain forest:

799       estimating species richness three different ways. Ecology. 2002;83:689-702.

800  32.  Townes H. A light-weight Malaise trap. Entomological News. 1972;83:239-47.

801  33.  Howard PC. Nature conservation in Uganda's tropical forest reserves. The IUCN

802       Tropical Forest Programme. 1991.

803  34.  Chapman CA and Chapman LJ. Forest regeneration in logged and unlogged forests of

804       Kibale National Park, Uganda. Biotropica. 1997;29:396-412.

805  35.  Kurina O. Description of four new species of *Zygomyia* Winnertz from Ethiopia and

806       Uganda (Diptera: Mycetophilidae). African Invertebrates. 2012;53 1:205-20.

807   36.   Folmer O, Black M, Hoeh W, Lutz R and Vrijenhoek R. DNA primers for amplification of

808         mitochondrial cytochrome c oxidase subunit I from diverse metazoan invertebrates.

809         Molecular Marine Biology and Technology. 1994;3 5:294-9.

810   37.   Comai L and Howell T. Barcode Generator.

811         *http://comailabgenomecenterucdavisedu/indexphp/Barcode_generator*. 2012.

812   38.   Katoh K and Standley DM. MAFFT multiple sequence alignment software version 7:

813         improvements in performance and usability. Molecular Biology and Evolution. 2013;30

814         4:772-80.

815   39.   Sovic I, Sikic M, Wilm A, Fenlon SN, Chen S and Nagarajan N. Fast and sensitive

816         mapping of nanopore sequencing reads with GraphMap. Nature Communications.

817         2016;7:11307.

818   40.   Vaser R, Sovic I, Nagarajan N and Sikic M. Fast and accurate de novo genome

819         assembly from long uncorrected reads. Genome Research. 2017;27 5:737-46.

820   41.   Leray M, Yang JY, Meyer CP, Mills SC, Agudelo N, Ranwez V, et al. A new versatile

821         primer set targeting a short fragment of the mitochondrial COI region for metabarcoding

822         metazoan diversity: application for characterizing coral reef fish gut contents. Frontiers in

823         Zoology. 2013;10:34.

824   42.   Geller JM, C., Parker M and Hawk H. Redesign of PCR primers for mitochondrial

825         cytochrome c oxidase subunit I for marine invertebrates and application in all-taxa biotic

826         surveys. Molecular Ecology Resources. 2013;13 5:851-61.

827   43.   Yeo D, Srivathsan A and Meier R. Mini-barcodes are more suitable for large-scale

828         species discovery in Metazoa than full-length barcodes. bioRxiv. 2019.

829         doi:10.1101/594952.

830   44.   Wang WY, Srivathsan A, Foo M, Yamane SK and Meier R. Sorting specimen-rich

831         invertebrate samples with cost-effective NGS barcodes: Validating a reverse workflow

832        for specimen processing. Molecular ecology resources. 2018;18 3:490-501.

833        doi:10.1111/1755-0998.12751.

834   45.   Zhang J, Kobert K, Flouri T and Stamatakis A. PEAR: a fast and accurate Illumina

835        Paired-End reAd mergeR. Bioinformatics. 2014;30 5:614-20.

836   46.   Srivathsan A, Sha JC, Vogler AP and Meier R. Comparing the effectiveness of

837        metagenomics and metabarcoding for diet analysis of a leaf-feeding monkey (*Pygathrix*

838        *nemaeus*). Mol Ecol Resour. 2015;15 2:250-61. doi:10.1111/1755-0998.12302.

839   47.   Kurtz S, Phillippy A, Delcher AL, Smoot M, Shumway M, Antonescu C, et al. Versatile

840        and open software for comparing large genomes. Genome biology. 2004;5 2:R12.

841        doi:10.1186/gb-2004-5-2-r12.

842   48.   Meier R, Shiyang K, Vaidya G and Ng PKL. DNA Barcoding and Taxonomy in Diptera: A

843        Tale of High Intraspecific Variability and Low Identification Success. Systematic Biology.

844        2006;55 5:715–28.

845   49.   Ahrens D, Fujisawa T, Krammer HJ, Eberle J, Fabrizi S and Vogler AP. Rarity and

846        Incomplete Sampling in DNA-based Species Delimitation. Systematic Biology. 2016;65

847        3:478-94.

848   50.   Leigh JW and Bryant D. PopART: Full-feature software for haplotype network

849        construction. Methods Ecol Evol. 2015;6 9:1110-6.

850   51.   Disney RHL. Scuttle flies (Diptera: Phoridae) Part II: the genus *Megaselia*. Fauna of

851        Arabia. 2009;24:249-357.

852   52.   Colwell RK. EstimateS: Statistical estimation of species richness and shared species

853        from samples. Version 9 and earlier. User's Guide and application. 2013.

854

855

**Figure Legends**

**Figure 1**: Flowchart for generating MinION barcodes.

**Figure 2**: Effect of re-pooling on coverage of barcodes for both sets of specimens.

**Figure 3**: Ambiguities in MAFFT+AA (Purple), RACON+AA (Yellow) and Consolidated barcodes (Green) with varying namino parameters (1,2 and 3). One outlier value for Racon+3AA barcode corresponding to 13% ambiguities was excluded from the plot.

**Figure 4**: The Malaise trap that revealed the estimated >1000 mOTUs as shown by the species richness estimation curve. Green: Chao1 Mean, Pink: S (Mean), Orange: Singleton Mean, Purple: Doubleton mean.

**Figure 5**: Lateral habitus (a) and diagnostic features of *Megaselia sepsioides* spec. nov. (a, inset) terminalia, (b) posterior view of foreleg, (c) anterior view of midleg (d,e) anterior and postero-dorsal views of hindleg, (e) dorsal view of thorax and abdomen.

**Figure 6**: Haplotype variation of Megaselia sepsioides spec. nov. (a) UGC0005996, (b) UGC0012244, (c) UGC0012899. UGC numbers refer to specimen IDs.

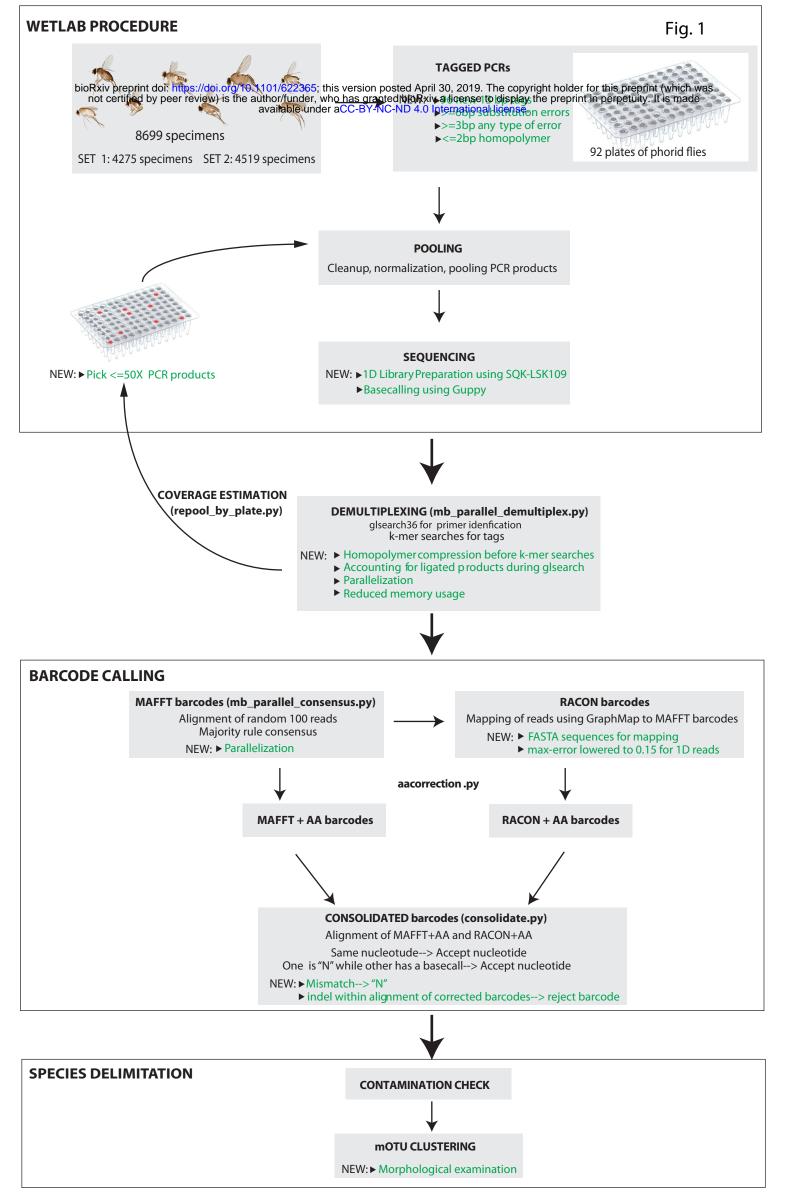**Figure 7**: Haplotype network for Megaselia sepsioides spec. nov. UGC numbers refer to specimen IDs.

Table 1: Number of reads and barcodes generated via MinION sequencing.

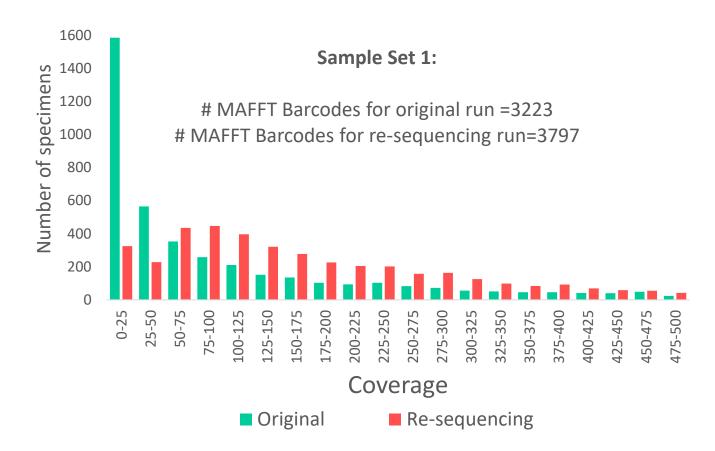| | Set 1: Two flowcells | Set 2: One flowcell | Combined (set 1 & 2)* |
|---|---|---|---|
| # Specimens | 4275 | 4519 | 8699 |
|     Resequencing (re-pooled) | 2172 | 2211 | |
| # reads/# reads >600 bp | 7,035,075/3,703,712 | 7,179,121/2,652,657 | |
|     Initial sequencing (all) | 3,069,048/1,942,212 | 4,853,363/2,250,591 | |
|     Resequencing (re-pooled) | 3,966,027/1,761,500 | 2,325,758/402,066 | |
| # demultiplexed reads | 898,979 (24.3%) | 647,152 (24.4%) | NA |
|     Initial sequencing (all) | 562,434 (29%) | 561,383 (24.9%) | |
|     Resequencing (re-pooled | 336,545 (19%) | 85,769 (21.3%) | |
| **Combined results of original and resequencing runs** | | | |
| # specimens with >=5X coverage | 4227 (98.9%) | 4287 (94.9%) | 8428 (96.9%) |
| # MAFFT barcodes <1% N's | 3797 (88.8%) | 3476 (76.9%) | 7220 (83%) |
| # MAFFT + AA barcodes | 3774 (88.3%) | 3464 (75.7%) | 7178 (82.5%) |
| # RACON barcodes | 3797 (88.8%) | 3476 (76.9%) | 7220 (83%) |
| # RACON +AA barcodes | 3790 (88.7%) | 3469 (76.7%) | 7194 (83%) |
| # Consolidated barcodes | 3740 (87.4%) | 3394 (75%) | 7115 (81.8%) |
| # Consolidated barcodes (non-phorids removed) | 3700 (86.7%) | 3369 (74.6%) | 7062 (81.2%) |
| **# mOTUS (2/3/4%)** | | | **728/685/636** |

* one plate was accidentally sequenced in both runs, wherever duplicates were present, second run was selected.
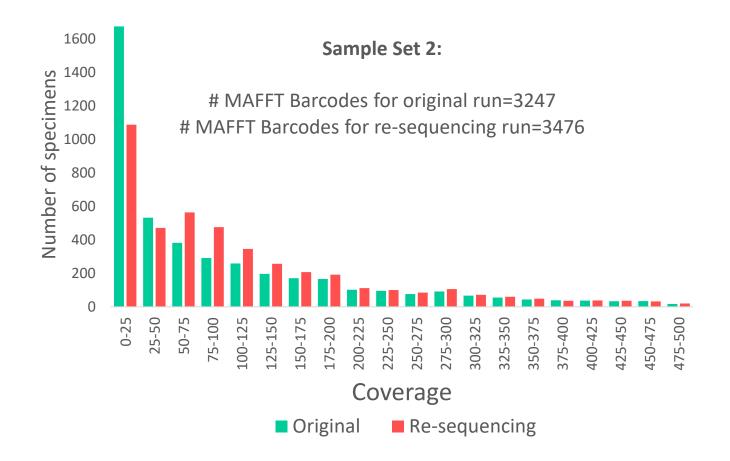
Table 2. Accuracy of MinION as assessed by Illumina barcodes. The overall optimal strategy is "Consolidated (namino=2)". Optimal congruence values are highlighted in green.

| Dataset | # compared with Illumina | Accuracy | # barcodes with errors/# >3% errors | mOTU richness deviation between MinION and Illumina barcodes | | |
|---|---|---|---|---|---|---|
| | | | | 2% | 3% | 4% |
| MAFFT | 6330 | 99.6091 | 4473/31 | -3 (-0.43%) | **-1 (-0.15)** | -8 (-1.3) |
| RACON | 6330 | 99.5075 | 4526/37 | 15 (2.08%) | 5 (0.75) | **-3 (-0.48)** |
| MAFFT + AA (namino=1) | 6146 | 99.9689 | 269/30 | -6 (-0.85%) | -4 (-0.61) | -15 (-2.7) |
| MAFFT + AA (namino=2) | 6146 | 99.9795 | 218/28 | -8 (-1.15%) | -6 (-0.91) | -13 (-2.2) |
| MAFFT + AA (namino=3) | 6286 | 99.967 | 350/28 | -7 (-0.1%) | -4 (-0.61) | -14 (-2.2) |
| RACON + AA (namino=1) | 6319 | 99.9617 | 419/27 | 5 (0.72%) | **1 (0.15)** | **-3 (-0.48)** |
| RACON + AA (namino=2) | 6319 | 99.9711 | 326/27 | 2 (0.29%) | -2 (-0.31) | -5 (-0.81) |
| RACON + AA (namino=3) | 6319 | 99.9664 | 383/29 | 4 (0.57%) | -2 (-0.3) | -8 (-1.28) |
| Consolidated (namino=1) | 6230 | 99.9855 | 191/25 | -1 (-0.14%) | -2 (-0.30) | **-3 (-0.48)** |
| **Consolidated (namino=2)** | 6255 | **99.9864** | **185/25** | **0 (0%)** | -3 (-0.45) | -4 (-0.64) |
| Consolidated (namino=3) | 6250 | 99.9762 | 329/26 | 1 (0.14%) | **-1 (-0.15)** | **-4 (-0.48)** |

Fig. 1

**WETLAB PROCEDURE**

8699 specimens

SET 1: 4275 specimens    SET 2: 4519 specimens

**TAGGED PCRs**

NEW: ▶ a clean 1D barcode
▶ <=6bp substitution errors
▶ >=3bp any type of error
▶ <=2bp homopolymer

92 plates of phorid flies

**POOLING**

Cleanup, normalization, pooling PCR products

NEW: ▶ Pick <=50X PCR products

**SEQUENCING**

NEW: ▶ 1D Library Preparation using SQK-LSK109
▶ Basecalling using Guppy

**COVERAGE ESTIMATION**
**(repool_by_plate.py)**

**DEMULTIPLEXING (mb_parallel_demultiplex.py)**

glsearch36 for primer idenfication
k-mer searches for tags

NEW: ▶ Homopolymer compression before k-mer searches
▶ Accounting for ligated products during glsearch
▶ Parallelization
▶ Reduced memory usage

**BARCODE CALLING**

**MAFFT barcodes (mb_parallel_consensus.py)**

Alignment of random 100 reads
Majority rule consensus

NEW: ▶ Parallelization

**RACON barcodes**

Mapping of reads using GraphMap to MAFFT barcodes

NEW: ▶ FASTA sequences for mapping
▶ max-error lowered to 0.15 for 1D reads

aacorrection .py

**MAFFT + AA barcodes**

**RACON + AA barcodes**

**CONSOLIDATED barcodes (consolidate.py)**

Alignment of MAFFT+AA and RACON+AA

Same nucleotide--> Accept nucleotide
One is "N" while other has a basecall--> Accept nucleotide

NEW: ▶ Mismatch--> "N"
▶ indel within alignment of corrected barcodes--> reject barcode

**SPECIES DELIMITATION**

**CONTAMINATION CHECK**

**mOTU CLUSTERING**

NEW: ▶ Morphological examination

Fig. 2

Fig. 3

Fig. 4

Fig. 5

**Fig. 6**

Fig. 7