

1 **Mini-barcodes are more suitable for large-scale species discovery in Metazoa than full-**  
2 **length barcodes**

3

4 **Darren Yeo<sup>1</sup>, Amrita Srivathsan<sup>1</sup>, Rudolf Meier<sup>1\*</sup>**

5

6

7 <sup>1</sup>Department of Biological Sciences, National University of Singapore, 14 Science 8 Drive 4, Singapore  
8 117543

9

10

11

12

13

14

15

16

17

18

19

20

21

22

23

24

25

26

27

28

29 \*Corresponding author: [meier@nus.edu.sg](mailto:meier@nus.edu.sg)

30 Keywords: DNA barcoding, mini-barcodes, species discovery, metabarcoding

31 Running Title: mini-barcodes for species discovery

32

### 33 **Abstract**

34 New techniques for the species-level sorting of millions of specimens have to be developed in  
35 order to answer the question of how many species live on earth. These methods should be  
36 reliable, scalable, and cost-effective as well as largely insensitive to the low-quality genomic  
37 DNA commonly obtained from museum specimens. Mini-barcodes seem to satisfy these  
38 criteria, but it is unclear whether they are sufficiently informative for species-level sorting. This  
39 is here tested based on 20 datasets covering ca. 30,000 specimens of 5,500 species. All  
40 specimens were first sorted based on morphology before being barcoded with full-length *cox1*  
41 barcodes. Mini-barcodes of different lengths and positions were then obtained *in silico* from  
42 the full-length barcodes using a sliding window approach (3 windows: 100-bp, 200-bp, 300-  
43 bp) as well as nine published mini-barcode primers (length: 94 – 407-bp). Afterwards, we  
44 determined whether barcode length and/or position reduces congruence between  
45 morphospecies and molecular Operational Taxonomic Units (mOTUs) that were obtained  
46 using three different species delimitation techniques (ABGD, PTP, objective clustering). We  
47 find that there is no significant difference in performance between full-length and mini-  
48 barcodes as long as they are of moderate length (>200-bp). Only very short mini-barcodes  
49 (<200-bp) perform poorly, especially when they are located near the 5' end of the Folmer  
50 region. Overall, congruence between morphospecies and mOTUs is ca. 80% for barcodes  
51 that are >200-bp. The congruent mOTUs contain ca. 75% of the specimens and we estimate  
52 that most of the conflict is caused by ca. 10% of the specimens that should be targeted for re-  
53 examination. Overall, barcode length (>200-bp) and species delimitation methods have minor  
54 effects on congruence. Our study suggests that large-scale species discovery and  
55 metabarcoding can utilize mini-barcodes without significant loss of information when  
56 compared to full-length barcodes. This is good news given that mini-barcodes can be obtained  
57 via cost-effective tagged amplicon sequencing using short-read sequencing platforms  
58 (Illumina: "NGS barcodes").

## 59 **Introduction**

60 The question of how many species live on earth has intrigued biologists for a very long time,  
61 but we are nowhere close to having a robust answer. We do know that fewer than 2 million of  
62 the estimated 10-100 million multicellular species have been described and that many are  
63 currently being extirpated by the “sixth mass extinction” (Ceballos et al., 2015; Sánchez-Bayo  
64 & Wyckhuys, 2019) with potentially catastrophic consequences for the environment (Cafaro,  
65 2015). Monitoring, halting, and perhaps even reversing this process is hampered by the  
66 “taxonomic impediment”. This impediment is particularly severe for “invertebrates” that  
67 collectively contribute much of the animal biomass (e.g., arthropods, annelids, nematodes:  
68 Stork et al., 2015; Bar-On et al., 2018). Most biologists thus agree that there is a pressing  
69 need for accelerating species discovery and description. This is very likely to require the  
70 development of new molecular methods. We would argue that they should not only be  
71 accurate, but also (1) rapid, (2) cost-effective, and (3) and largely insensitive to DNA quality.  
72 These criteria are important because tackling the earth’s biodiversity would likely require the  
73 processing of >500 million specimens even under the very conservative assumption that there  
74 are only 10 million species and a new species is discovered with every 50 specimens  
75 processed. Cost-effectiveness is similarly important because millions of species are found in  
76 countries with only basic research facilities and the large-scale international transfer of  
77 specimens for molecular work is becoming increasingly difficult under the Nagoya protocol.  
78 Fortunately, many species are already represented in museum holdings, but such specimens  
79 often yield degraded DNA (Cooper, 1994). Therefore, methods that require DNA of high-  
80 quality and quantity are not likely to be suitable for large-scale species discovery in  
81 invertebrates.

82 Conceptually, species discovery and description can be broken up into three steps. The first  
83 is obtaining specimens, the second, species-level sorting, and the third, species identification  
84 or description. Fortunately, centuries of collecting have generated many of the specimens that  
85 are needed for large-scale species discovery. Indeed, for many invertebrate groups it is likely

86 that the museum collections contain more specimens of undescribed than described species;  
87 i.e., this unsorted collection material represents vast and still underutilized source for species  
88 discovery (Lister & Climate Change Research Group, 2011; Kemp, 2015; Yeates et al., 2016).  
89 The second step in species discovery/description is species-level sorting, which often involves  
90 various levels of sorting according to the taxonomic expertise available. This step is in dire  
91 need for acceleration. Traditionally, it starts with the sorting of unsorted material into major  
92 taxa (e.g., order-level in insects). This task can be accomplished by parataxonomists but may  
93 in the future be taken over by machine sorting utilizing neural networks (Valan et al., 2019). In  
94 contrast, the subsequent species-level sorting is usually time-limiting because the specimens  
95 for many invertebrate taxa have to be prepared by highly-skilled specialists (e.g., dissected;  
96 slide-mounted) before the material can be sorted into putative species. This means that the  
97 traditional techniques are neither rapid nor cost-effective for many invertebrate groups. This  
98 impediment is likely to be largely responsible for why certain taxa that are known to be  
99 abundant and species-rich are particularly poorly studied (Bickel, 1999).

100 An alternative way to sort specimens to species-level would be with DNA sequences. This  
101 approach is particularly promising for metazoan species because most multicellular animal  
102 species can be distinguished based on cytochrome c oxidase subunit I (*cox1*) barcode  
103 sequences (Hebert et al., 2003). However, such sorting requires that each specimen is  
104 barcoded. This creates cost- and scalability problems when the barcodes are obtained with  
105 Sanger sequencing (see Taylor and Harris, 2012). Such sequencing is currently still the  
106 standard in barcoding studies because the animal barcode was defined as a 658-bp long  
107 fragment of *cox1* ("Folmer region": Folmer et al., 1994), although sequences >500-bp with <1%  
108 ambiguous bases are also considered BOLD-compliant (BOLDsystems.org). The 658-bp  
109 barcode was optimized for ABI capillary sequencers but has become a burden because it is  
110 not suitable for most new sequencing technologies, while Sanger sequencing remains  
111 expensive and only scalable when expensive liquid-handling robots are used. This approach  
112 is hence unlikely to become widely available in those countries that harbour most of the

113 species diversity. Due to these constraints, very few studies have utilized DNA barcodes to  
114 sort entire samples into putative species (but see Fagan-Jeffries et al., 2018). Instead, most  
115 studies use a mixed approach where species-level sorting is carried out based on morphology  
116 before a select few specimens per morphospecies are barcoded (e.g., Riedel et al., 2010).  
117 This two-step process requires considerable amounts of skilled manpower and time.

118 Scalability and cost-effectiveness are hallmark features of the new short-read high throughput  
119 sequencing technologies. In addition, these technologies are particularly suitable for  
120 sequencing the kind of degraded DNA that is typical for museum specimens. Indeed, anchored  
121 hybrid enrichment (AHE) has already been optimized for the use with old museum specimens  
122 (Bi et al., 2013; Guschanski et al., 2013; Blaimer et al., 2016) and is likely to play a major role  
123 for the integration of rare species into taxonomic and systematic projects. It will be difficult,  
124 however, to apply AHE to millions of specimens because it requires time-consuming and  
125 expensive molecular protocols (e.g., specimen-specific libraries). Fortunately, for most  
126 species it is likely that species-level sorting does not require a very large number of markers.

127 We would thus argue that the initial species-level sorting can be achieved using barcodes that  
128 are obtained via “tagged amplicon sequencing” on next-generation sequencing platforms  
129 (“NGS barcodes”; Wang et al., 2018; Yeo et al., 2018). Until recently, obtaining full-length NGS  
130 barcodes via tagged amplicon sequencing was difficult because the reads of most next-  
131 generation-sequencing platforms were too short for sequencing the full-length barcode. This  
132 has now changed with the arrival of third-generation platforms (ONT: MinION: Srivathsan et  
133 al., 2018; PacBio: Sequel: Hebert et al., 2018). These platforms, however, come with  
134 drawbacks; *viz* elevated sequencing error rates and higher cost. These problems are likely to  
135 be overcome in the future (e.g., Yang et al., 2018), but such solutions will not solve the main  
136 challenge posed by full-length barcodes; i.e., reliably obtaining amplicons from museum  
137 specimens with degraded DNA (e.g., Hajibabaei et al., 2006). We thus submit that one should  
138 optimize the barcode length based on empirical evidence; it should be only as long as needed  
139 for accurate pre-sorting of specimens into putative species.

140 Barcodes that are shorter than the full-length barcode are called mini-barcodes. They are  
141 obtained with primers that amplify shorter subsets of the original barcode region and have  
142 several advantages. Firstly, such amplicons are easier to obtain when the DNA in the sample  
143 is degraded (Hajibabaei & McKenna, 2012). Secondly, mini-barcodes can be sequenced at  
144 low cost using tagged amplicon sequencing on short-read sequencing platforms (e.g.,  
145 Illumina). Thirdly, mini-barcode primers are available for a large number of species-rich  
146 metazoan clades (Hajibabaei et al., 2006; Meusnier et al., 2008; Hebert et al., 2013; Little,  
147 2014) as well as for specific taxa such as fruit flies, catfish and sharks (Fan et al., 2009;  
148 Bhattacharjee & Ghosh, 2014; Fields et al., 2015). It is thus not surprising that short barcodes  
149 are already the barcodes of choice when the template DNA is degraded. This is often the case  
150 for museum specimens (Zuccon et al., 2012; Hebert et al., 2013) or for environmental DNA  
151 which is usually analysed via metabarcoding (e.g.: processed food: Armani et al., 2015;  
152 Shokralla et al., 2015; water, soil, fecal matter: Epp et al., 2012; Srivathsan et al., 2015; Lim  
153 et al., 2016). Mini-barcodes were initially obtained via Sanger sequencing, but they can now  
154 be sequenced much more efficiently via tagged amplicon sequencing on short-read platforms  
155 (“NGS barcoding”: Wang et al., 2018: sequencing cost < 4 cents). Wang et al. (2018) could  
156 thus implement a “reverse workflow” based on sequencing all specimens without any species-  
157 level pre-sorting based on morphology. Four-thousand specimens of ants were barcoded with  
158 a 313-bp mini-barcode. The specimens were then pre-sorted into 89 molecular operational  
159 taxonomic units (mOTUs) that were largely congruent with morphospecies (86 species).  
160 However, it remained unclear whether full-length DNA barcodes would have further improved  
161 congruence, whether the results from this study can be generalized, and which mini-barcode  
162 is optimal for large-scale pre-sorting of specimens into putative species.

163 The answers to these questions remain elusive, because mini-barcodes remain insufficiently  
164 tested despite their ubiquitous use in metabarcoding. Arguably, existing tests suffer from lack  
165 of scale (the largest study includes 6695 barcodes for 1587 species: Meusnier et al., 2008)  
166 and taxonomic scope (usually only 1-2 family-level taxa: e.g. Hajibabaei et al., 2006; Yu & You,

167 2010). Furthermore, the tests yielded conflicting results. Hajibabaei et al. (2006) found high  
168 congruence with the full-length barcode when species are delimited based on mini-barcodes  
169 and Meusnier et al. (2008) find similar BLAST identification rates for mini-barcodes and full-  
170 length barcodes in their *in silico* tests. However, Yu & You (2010) conceded that mini-barcodes  
171 may have worse accuracy despite having close structural concordance with the full-length  
172 barcode. In addition, Sultana et al. (2018) concluded that the ability to identify species is  
173 compromised when the barcodes are too short (<150-bp), but it remained unclear at which  
174 length and in which position mini-barcodes start performing well. Furthermore, published tests  
175 of mini-barcodes compare their performance to results obtained with full-length barcodes. All  
176 conflict is then implicitly considered evidence for the failure of mini-barcodes to yield the  
177 “correct” mOTUs. However, results obtained with longer barcodes should not automatically be  
178 assumed to be accurate given that the Folmer region varies in nucleotide variability (Roe &  
179 Sperling, 2007). Lastly, the existing tests of mini-barcodes do not include a sufficiently large  
180 number of different mini-barcodes in order to be able to detect positional and lengths effects  
181 across the 658-bp barcode region.

182 Here, we address the lack of scale by including 20 studies covering 5500 species represented  
183 by ca. 30,000 barcodes. We furthermore test a large number of different mini-barcodes by  
184 applying a sliding window approach to generate mini-barcodes of different sizes (100, 200,  
185 300-bp window sizes, 60-bp intervals) and compare the results to the performance of nine  
186 mini-barcodes with published primers (mini-barcode length: 94 – 407-bp). The taxonomic  
187 scope of our study is broad enough to include a wide variety of metazoans ranging from  
188 earthworms to butterflies and birds. Lastly, we do not assume that mOTUs based on full-length  
189 barcodes are automatically more accurate than those obtained with mini-barcodes. Instead,  
190 we use morphology as an external criterion for assessing whether mOTUs obtained with  
191 different-length barcodes have different levels of congruence with morphospecies. Note that  
192 this does not imply that morphology is more suitable for species delimitation than molecular  
193 data. Instead, we test whether shortening barcodes influences congruence with morphology;

194 i.e. morphology is treated as a constant while testing whether barcode length and/or position  
195 influences the number of morphospecies that are recovered. Given that morphology is a  
196 generally accepted type of data that can be used for species delimitation, mini-barcodes that  
197 significantly lower congruence with morphospecies are unlikely to be useful for accurate  
198 species-level sorting.

199 We also compare the performance of different species delimitation methods. There has been  
200 substantial interest in developing algorithms for mOTU estimation, leading to the emergence  
201 of various species delimitation algorithms over the past decade (e.g., objective clustering:  
202 Meier et al., 2006; BPP: Yang & Rannala, 2010; jmOTU: Jones et al., 2011; ABGD; Puillandre  
203 et al., 2012; BINs: Ratnasingham & Hebert, 2013; PTP: Zhang et al., 2013; etc.). For the  
204 purposes of this study, we selected three algorithms as representatives of distance and tree-  
205 based methods: objective clustering, Automatic Barcode Gap Discovery (ABGD) and Poisson  
206 Tree Process (PTP). Objective clustering utilizes an *a priori* distance threshold to group  
207 sequences into clusters, ABGD groups sequences into clusters based on an initial prior and  
208 recursively uses incremental priors to find stable partitions, while PTP utilizes the branch  
209 lengths on the input phylogeny to delimit species units. Arguably, barcode data may not be  
210 appropriate for the application of PTP because a single marker is not likely to yield reliable  
211 phylogenetic trees (including branch lengths), but PTP has been frequently applied to barcode  
212 data in the literature (e.g. Ermakov et al., 2015; Han et al., 2016; Hollatz et al., 2016) and is  
213 thus included here. There are numerous additional techniques for species delimitation, but  
214 most require multiple markers and/or are usually even more reliant on accurately  
215 reconstructed phylogenetic trees and may not be easily scalable to millions of specimens.  
216 They are therefore not included in this study.

217

## 218 **Materials & Methods**

### 219 *Dataset selection*



220 We surveyed the barcoding literature in order to identify publications that cited the original  
221 barcode paper by Hebert et al. (2003) and met the following criteria: 1) have pre-identified  
222 specimens where the barcoded specimens were pre-sorted/identified based on morphology  
223 and 2) the dataset had at least 500 specimens with *cox1* barcodes >656-bp. We identified 20  
224 most recent datasets starting from 2017 (Table S1); all had >500 barcoded specimens even  
225 after removing those that were not sorted to species level (e.g., only identified to genus or  
226 higher) or had short sequences <657-bp (the full-length barcode is technically 658-bp long,  
227 but a 1-bp concession was made to prevent the loss of too much data). The barcode  
228 sequences were downloaded from BOLDSystems or NCBI GenBank and aligned with MAFFT  
229 v7 (Katoh & Standley, 2013) with a gap opening penalty of 5.0.

230 Using a custom python script, we generated three sets of mini-barcodes along a “sliding  
231 window”. They were of 100-, 200- and 300-bp lengths. The first iteration begins with the first  
232 base pair of the 658-bp barcode and the shifting windows jump 60-bp at each iteration,  
233 generating ten 100-bp windows, eight 200-bp windows and six 300-bp windows. Additionally,  
234 we identified nine mini-barcodes with published primers within the *cox1* Folmer region (Fig. 1  
235 & Table S2). These mini-barcodes have been repeatedly used in the literature published after  
236 2003 and were used for a broad range of taxa. The primers for the various mini-barcodes were  
237 aligned to the homologous regions of each dataset with MAFFT v7 --addfragments (Katoh &  
238 Standley, 2013) in order to identify the precise position of the mini-barcodes within the full-  
239 length barcode. The mini-barcode subsets from each barcode were then identified after  
240 alignment to full-length barcodes. Note that most of the published primers are in the 5' prime  
241 region of the full-length barcode.

#### 242 *Species delimitation*

243 The mini-barcodes and the full-length barcodes were clustered into putative species using  
244 three species delimitation algorithms: objective clustering (Meier et al., 2006), ABGD  
245 (Puillandre et al., 2012) and PTP (Zhang et al., 2013). For objective clustering, the mOTUs  
246 were clustered at 2 – 4% uncorrected p-distance thresholds (Srivathsan & Meier, 2012) using

247 a python script which reimplements the objective clustering of Meier et al. 2006 and allows for  
248 batch processing. The p-distance thresholds selected are the typical distance thresholds used  
249 for species delimitation in the literature (Meier et al. 2006; Ratnasingham & Hebert, 2013;  
250 Meier et al. 2016). The same datasets were also clustered with ABGD (Puillandre et al., 2012)  
251 using the default range of priors and with uncorrected p-distances, but the minimum slope  
252 parameter (-X) was reduced in a stepwise manner (1.5, 1.0, 0.5, 0.1) if the algorithm could not  
253 find a partition. We then considered the ABGD clusters at priors  $P=0.001$ ,  $P=0.01$  and  $P=0.04$   
254 in this study. The priors (P) refer to the maximum intraspecific divergence and functions  
255 similarly to p-distance thresholds at the first iteration, before being recursively refined by  
256 recursive application of the ABGD algorithm. Lastly, in order to use PTP, the datasets were  
257 used to generate maximum likelihood (ML) trees in RAxML v.8 (Stamatakis, 2014) via rapid  
258 bootstrapping (-f a) and the GTRCAT model. The best tree generated for each dataset was  
259 then used for species delimitation with PTP (Zhang et al., 2013) under default parameters.

#### 260 *Performance assessment*

261 We assess the performance of mini-barcodes by using morphospecies as an external arbiter.  
262 Species-level congruence was quantified using match ratios between molecular and  
263 morphological groups (Ahrens et al., 2016). The ratio is defined as  $\frac{2 \times N_{match}}{N_1 + N_2} \times 100$ , where  
264  $N_{match}$  is the number of clusters identical across both mOTU delimitation methods/thresholds  
265 ( $N_1$  &  $N_2$ ). Incongruence between morphospecies and mOTUs is usually caused by a few  
266 specimens that are assigned to the “incorrect” mOTUs. Conflict at the specimen-level can thus  
267 be quantified as the number of specimens that are in mOTUs that cause conflict with  
268 morphospecies.

269 In order to test whether barcode length is a significant predictor of congruence, MANOVA tests  
270 were carried out in R (R Core Team, 2017) with “match ratio” (species-level congruence) as  
271 the response variable and “dataset” and “mini-barcode” as categorical explanatory variables.  
272 We found that most of the variance in our study was generated by the “dataset” variable ( $P <$

273 0.05 in MANOVA tests). Given that we were particularly interested in the effect of barcode  
274 length and position, “dataset” was subsequently treated as a random effect “mini-barcode” as  
275 the explanatory variable (categorical) in a linear mixed effects model (R package *lme4*: Bates,  
276 2010). The *emmeans* R package (Lenth, 2018) was then used to perform pairwise post-hoc  
277 Tukey tests between mini- and full-length barcodes so as to assess whether either barcode  
278 was performing significantly better/worse. To compare the differences in performance  
279 between objective clustering, ABGD and PTP, ANOVA tests were performed in R. After which,  
280 pairwise Tukey tests were used to determine which species delimitation method was  
281 responsible for significant differences. Lastly, in order to explore the reasons for positional  
282 effects, the proportion of conserved sites for each mini-barcode was obtained using MEGA6  
283 (Tamura et al., 2013).

284 Match ratios indicate congruence at the species level, but it is also important to determine how  
285 many specimens have been placed in congruent units. Species- and specimen-level  
286 congruence are only identical when all mOTUs are represented by the same number of  
287 specimens. However, specimen abundances are rarely equal across species and hence  
288 match ratio is insufficient at characterizing congruence between mOTUs and morphospecies.  
289 It is straightforward to determine the number of congruent specimens as follows:

290 (1) Congruence Class I specimens: If  $A = B$  then number of congruent specimens is  $N_{c1} =$   
291  $|A|$  OR  $|B|$ .

292 Incongruence is caused by morphospecies that are split, lumped, or split and lumped in the  
293 mOTUs. However, any one mis-sorted specimen placed into a large-sized mOTU leads to all  
294 specimens in two mOTUs to be considered “incongruent” according to the criterion outlined  
295 above. Yet, most specimens are congruent and full congruence could be restored by re-  
296 examining the mis-sorted specimen. It is therefore also desirable to determine the number of  
297 specimens that require re-examination or, conversely, the number of specimens that would  
298 be congruent if one were to remove a few incongruent specimens. This number of specimens  
299 can be estimated by counting congruent specimens as follows:

300 (2) Congruence Class II specimens: Specimens that are in split or lumped mOTUs relative to  
301 morphospecies. Here, the largest subset of congruently placed specimens can be determined  
302 as follows. If  $A_1 \cup A_2 \cup \dots \cup A_x = B : Nc_{2=} \max(|A_1|, |A_2| \dots |A_x|)$

303 (3) Congruence Class III specimens: This covers specimens in sets of clusters that are both  
304 split and lumped relative to morphospecies. Here, only those specimens are considered  
305 potentially congruent that (1) are in one mOTU and one morphospecies and (2) combined  
306 exceed the number of the other specimens in the set of clusters. In detail, if  $A_1 \cup A_2 \cup \dots \cup$   
307  $A_x = B_1 \cup B_2 \cup \dots \cup B_y : Nc_3 = \max(|A_1 \cap B_1|, |A_2 \cap B_1| \dots |A_x \cap B_y|)$  only if  $\max(|A_1 \cap$   
308  $B_1|, |A_2 \cap B_1| \dots |A_x \cap B_y|) > \frac{1}{2} (|A_1 \cup A_2 \dots A_x|)$ .

309

## 310 Results

311 For species delimitation with objective clustering, we found that the 2% p-distance threshold  
312 yielded the highest congruence across the datasets. It was hence used as the upper-bound  
313 estimator for species- and specimen-level congruence. The corresponding results for the 3  
314 and 4% p-distance clusters are reported in the supplementary materials. For ABGD it was the  
315  $P=0.001$  prior that yielded the highest average match ratio and hence the clusters generated  
316 by this prior were used in the main analysis (see supplementary material for results under  
317  $P=0.01$  and  $P=0.04$ ). PTP does not require parameter choices post the input tree.

318 The MANOVA tests performed on all treatments (species delimitation method and distance  
319 threshold/prior) indicated that the test variable “dataset” was responsible for much more of the  
320 observed variance in “match ratio”. The choice of mini-barcode or mOTU algorithm that was  
321 used to generate the mOTUs was of secondary importance (Table S3). After accounting for  
322 “dataset”, we find that only mini-barcodes <200-bp perform significantly worse than full-length  
323 barcodes (Fig. 2); for all other mini-barcodes (>200-bp) the congruence with morphospecies  
324 does not differ significantly and is occasionally superior to what is observed for the full-length  
325 barcode. This is evident in the large number of significant differences ( $p < 0.05$  &  $p < 0.001$ )

326 in pairwise post-hoc Tukey tests applied to 100-bp mini- and 657-bp full-length barcodes. Only  
327 short <100-bp barcodes have a mean performance that is worse (<0 match ratio deviation)  
328 than the full-length barcode. Conversely, there is no significant difference between the 200  
329 and 300-bp mini-barcode and the full-length barcode when objective clustering or PTP are  
330 used to estimate mOTUs. Under ABGD, the mini-barcode outperform the full-length barcodes.  
331 For all mOTU delimitation methods, the variance across datasets appears to decline as the  
332 mini-barcode increases in length (Fig. 2). The results obtained for *in silico* mini-barcode are  
333 consistent with the performance of mini-barcode with published primers: the mini-barcode  
334 of 94-bp, 130-bp, and 145-bp length tend to perform worse than the longer mini-barcode (Fig.  
335 3). The results are also similar for specimen-level congruence (Table 1 & Fig. S8). However,  
336 there are some exceptions including the performance improvements of short mini-barcode,  
337 for example, for European marine fish and Northwest Pacific molluscs when grouped with  
338 objective clustering.

339 When the performance of the three different clustering methods was compared, significant  
340 differences ( $p < 0.05$  in ANOVA test) were found only for the 100-bp mini-barcode set (Fig. 4).  
341 Here, pairwise post-hoc Tukey tests find that objective clustering performs significantly better  
342 than the other delimitation methods ( $p < 0.001$ ) while ABGD and PTP do not differ significantly  
343 ( $p = 0.88$ ) but behave erratically for short mini-barcode (Fig. 2).

344 Mini-barcode situated at the 5' end of the full-length barcode appear to perform somewhat  
345 worse than those situated at the middle or at the 3' end (Fig. 2). For example, the 100-bp mini-  
346 barcode at the 5' end perform poorly for objective clustering (mini-barcode midpoints at 50,  
347 110 & 170-bp), ABGD (mini-barcode midpoints at 110 & 170-bp) and PTP (mini-barcode  
348 midpoint at 110-bp). This effect is, however, only statistically significant when the mini-  
349 barcode are very short (100-bp). This positional effect is present across all species  
350 delimitation techniques. Note that the 5' end of the full-length barcode appears to contain a  
351 large proportion of conserved sites, particularly around the 170-bp and 230-bp midpoint of the

352 100-bp mini-barcode (Fig. 5). This positional effect averages out as the mini-barcodes  
353 increase in length.

354 With regard to specimen-based congruence, we evaluated to the mini-barcodes with  
355 published primers and here report the results for those that a barcode length >200-bp.  
356 Approximately three quarters of all specimens are in the “Congruence Class I” (Tables 1 &  
357 S8); i.e., their placement is congruent between mOTUs and morphospecies (Average/Median:  
358 OC at 2%: 75/75%; ABGD P=0.001: 71/72%; PTP: 75/75%). The remaining specimens are  
359 placed in mOTUs that are split, lumped, or split and lumped. The number of specimens that  
360 are predominantly responsible for the splitting and lumping are here classified as Congruence  
361 Class II and III specimens. Overall, fewer than 10% of the specimens fall into these categories  
362 (Table 1: see Class II specimens across species delimitation methods). These are the  
363 specimens that should be studied when addressing conflict between morphospecies and  
364 mOTUs.

365

## 366 **Discussion**

367 Accelerating species discovery and description is arguably one of the foremost challenges in  
368 modern systematics. Material for many undescribed species is already in world’s natural  
369 history museums, but the specimens need to be sorted to species-level before they become  
370 available for species identification/description and can be used for large-scale analyses of  
371 biodiversity patterns. Pre-sorting specimens with DNA barcodes is a potentially promising  
372 solution because it is scalable, can be applied to millions of specimens, and much of the  
373 specimen handling can be automated. However, in order for this approach to be suitable, a  
374 sufficiently large proportion of the pre-sorted units need to accurately reflect species  
375 boundaries and the methods for obtaining the sequences need to be suitable for the  
376 processing of large numbers of specimens whose DNA is degraded.

377 *The main source of variance in congruence: datasets*

378 We here find that the average congruence between mOTUs and morphospecies is 80% for all  
379 barcodes >200-bp (median: 83%), with the median being higher (83%) because of outlier  
380 datasets with congruence <65% (OC at 2%; ABGD  $P=0.001$ , PTP). These outlier datasets are  
381 also likely to be responsible for the observation that much of the variance in congruence  
382 throughout our study is caused explained by the variable “dataset”. Despite the outliers, 72-  
383 75% (median) of the ca. 30,000 specimens are assigned to species that are supported by  
384 molecular and morphological data. Overall, this is a very high proportion when compared to  
385 species-level sorting by parataxonomists (Krell, 2004). Unfortunately, this specimen-based  
386 perspective on congruence is often underappreciated when mOTUs and morphospecies are  
387 compared. However, specimen-level congruence is an important criterion for evaluating the  
388 suitability of species-level sorting with barcodes. After all, the basic units in a museum  
389 collections or an ecological survey are specimens and not species. The correct placement of  
390 specimens into species is thus important for systematists and biodiversity researchers alike  
391 given that the former would like to see most of the specimens in a collection correctly placed  
392 and the latter often need abundance and biomass information at species-level resolution.

393 The remaining ca. 25% of specimens are placed in mOTUs whose boundaries do not agree  
394 with morphospecies. One may initially consider this an unacceptably high proportion, but it is  
395 important to keep in mind that the misplacement of one specimen (e.g., due to a contamination  
396 of a PCR) will render two mOTUs incongruent; i.e., all specimens in these mOTUs will be  
397 considered incongruent and included in the 25%. Arguably, one should instead estimate how  
398 many specimens are causing the conflict. These are the specimens that should be targeted in  
399 reconciliation studies. The proportion across the 20 datasets in our study is fairly low and  
400 ranges from 10-12% (median) depending on which mOTU delimitation technique is used.

401 Conflict between mOTUs and morphospecies can be caused by technical error or biology. A  
402 typical technical factor would be accidental misplacement of specimens due to lab  
403 contamination or error during morphospecies sorting. Indeed, the literature is replete with  
404 cases where mOTUs that were initially in conflict with morphospecies became congruent once

405 the study of additional morphological characters let to the revision of morphospecies  
406 boundaries (e.g., Smith et al., 2008; Tan et al., 2010; Baldwin et al., 2011; Ang et al., 2017).  
407 But there are also numerous biological reasons for why one should not expect perfect  
408 congruence between mOTUs and species. Lineage sorting, fast speciation, large amounts of  
409 intraspecific variability, and introgression are known to negatively affect the accuracy of DNA  
410 barcodes (Will & Rubinoff, 2004; Rubinoff et al., 2006; Meier, 2008). It is thus somewhat  
411 surprising that regardless of these issues, the final levels of congruence between  
412 morphospecies and DNA sequences are often quite high in animals (Ball et al., 2005; Cywinska  
413 et al., 2006; Renaud et al., 2012; Landi et al., 2014; Wang et al., 2018). This implies that the  
414 pre-sorting specimens to species-level units based on mini-barcodes is worth pursuing for  
415 many metazoan clades. High levels of congruence are, however, not a universal observation  
416 across all of life. This approach to specimen sorting is unlikely to be useful in groups with  
417 widespread barcode sharing between species. This phenomenon occurs within Metazoa (e.g.,  
418 Anthozoa: Huang et al., 2008) and is likely to be the default outside of Metazoa (e.g., Chase  
419 & Fay, 2009; Hollingsworth et al., 2011).

#### 420 *Barcode length and species delimitation methods*

421 We here tested the widespread assumption that mOTUs based on full-length barcodes are  
422 more reliable than those based on mini-barcodes (Burns et al., 2007; Min & Hickey, 2007). If  
423 this assumption was confirmed, then the use of mini-barcodes would have to be discouraged  
424 despite higher amplification success rates, improved suitability for degraded starting material,  
425 and the availability of cost-effective sequencing on short-read high-throughput platforms.  
426 However, we find that the performance of *cox1* mini-barcodes with a length >200-bp do not  
427 differ significantly from the performance of full-length barcodes. Indeed, compared to the  
428 dataset effect, the choice of barcode length is largely secondary. This conclusion is robust  
429 across 20 diverse datasets and holds across different clustering algorithms.

430 We also find that the choice of species delimitation algorithm matters little for mini-  
431 barcodes >200-bp (Fig. 4). This is fortunate as objective clustering and ABGD algorithms are



432 less computationally demanding than PTP, which necessitates the reconstruction of a ML  
433 trees. However, there are some exceptions. Firstly, when the mini-barcodes are extremely  
434 short (~100-bp), objective clustering tends to outperform ABGD and PTP. PTP's poor  
435 performance for the 100-bp mini-barcodes is not surprising given that it relies on tree  
436 topologies which cannot be estimated with confidence based on so little data. ABGD's poor  
437 performance is mostly observed for certain priors (e.g.,  $P=0.04$ : Fig. S5 & S6). Under these  
438 priors, ABGD tends to lump most of the 100-bp barcodes into one or few large clusters. Prior-  
439 choice also affects ABGD's performance for full-length barcodes. ABGD does not perform well  
440 with very low priors ( $P = 0.001$ : Fig. 2 & 3 vs.  $P = 0.01$ ;  $P = 0.04$ : Fig. S5). Overall, we conclude  
441 that the selection of the best priors and/or clustering thresholds remains a significant challenge  
442 for the study of largely unknown faunas that lack morphological information as an *a posteriori*  
443 method for selecting priors/thresholds. Overall, we recommend the use of multiple methods  
444 and thresholds in order to distinguish robust from labile mOTUs that are heavily dependent on  
445 threshold- or prior-choice.

#### 446 *Positional effects*

447 We find that in general, mini-barcodes at the 3' end of the Folmer region outperform mini-  
448 barcodes at the 5' end. This is consistent across all three species delimitation methods and  
449 was also reported by Shokralla et al. (2015) who concluded that mini-barcodes at the 5' end  
450 have worse species resolution for fish species. This positional effect is apparent when match  
451 ratios are compared across a "sliding window" (Fig. 2). The lowest congruence with  
452 morphology is observed for 100-bp mini-barcodes with midpoints at the 50, 110 and 170-bp  
453 marks. However, this positional effect is only significant when the barcode lengths are very  
454 short (<200-bp). Once the mini-barcodes are sufficiently long (>200-bp), there seems to be no  
455 appreciable difference in performance, which is not surprising because sampling more  
456 nucleotides helps with buffering against regional changes in nucleotide variability across the  
457 Folmer region. These changes may be related to the conformation of the Cox1 protein in the  
458 mitochondrion membrane. The Folmer region of Cox1 contains six transmembrane  $\alpha$ - helices

459 and connected by five loops (Tsukihara et al. 1996; Pentinsaari et al. 2016). Pentinsaari et al.  
460 (2016) compared 292 Cox1 sequences across 26 animal phyla and found high amino acid  
461 variability in helix I and the loop connecting helix I and helix II (corresponding to position 1-  
462 102 of *cox1*), as well as end of helix IV and loop connecting helix IV and V (corresponding to  
463 positions ~448-498). These regions of high variability are distant from the active sites and thus  
464 less likely to affect Cox1 function (Pentinsaari et al. 2016). This may lead to lower selection  
465 pressure and high variability in these areas which could impact the performance of mini-  
466 barcodes for species delimitation.

#### 467 *Accelerating biodiversity discovery and description*

468 We had earlier argued that species discovery and description can be broken up into three  
469 steps (1) obtaining specimens, (2) species-level sorting and (3) species identification or  
470 description. We here only address species-level sorting. This means that the impediments  
471 caused by slow species identification and description remain apparently unresolved. However,  
472 this is only partially correct. Firstly, some mOTUs delimited via barcodes can be identified via  
473 barcode databases. The proportion of successful identification differs depending on how well  
474 a particular fauna has been studied. This is illustrated by our recent work on dragon- and  
475 damselflies (Odonata), ants (Formicidae), and non-biting midges (Chironomidae) in Singapore  
476 (Wang et al., 2018, Yeo et al., 2018; Baloğlu et al., 2018). For odonates, BLAST-searches  
477 identified more than half of the 95 mOTUs and >75% of the specimens to species. The  
478 corresponding numbers for ants and midges were ca. 20% and 10% at mOTU-level, and 9%  
479 and 40% at the specimen-level. Secondly, mOTUs discovered via barcodes can be readily  
480 compared across studies and borders (Ratnasingham, et al. 2013). In contrast, species newly  
481 discovered based on morphological evidence usually remain unavailable to the scientific  
482 community until they are published. This is a very significant differences because a large  
483 amount of downstream biodiversity analysis can be carried out based on mOTUs instead of  
484 identified/described species. This includes studying species richness and abundance over  
485 time which is a task that is becoming increasingly important in the 21<sup>st</sup> century. This means

486 that only moderate harm is done if species identification or description are only completed at  
487 a later time.

488 The analyses of biodiversity patterns will be impacted by incorrectly delimited mOTUs. In our  
489 study, we find that ca. 80% of the mOTUs are congruent with morphospecies. This is prior to  
490 a reconciliation stage where the morphology of specimens with a conflicting assignment is  
491 revisited in order to rule out that the morphological evidence was misinterpreted and/or an  
492 insufficient number of characters was studied; i.e., we would predict that the overall  
493 congruence levels after reconciliation will be higher. Ideally, we would like to know which  
494 proportion of mOTUs that are in conflict with morphospecies will eventually be rejected, but  
495 unfortunately we still know fairly little about the congruence levels between morphology and  
496 barcodes after reconciliation. This is because rigorous studies would have to be based on  
497 datasets with dense taxon and geographic sampling where morphological and DNA sequence  
498 information is obtained for all specimens and all cases of conflict are re-studied. Unfortunately,  
499 there are very few datasets that satisfy these criteria. This is presumably because the high  
500 cost of full-length barcodes has prevented biologists from sequencing all specimens.

### 501 *Conclusions*

502 We here illustrate that mini-barcodes can be used for pre-sorting specimens into putative  
503 species and that they are arguably the preferred choice because (1) they are obtained more  
504 readily for specimens that only yield degraded DNA (Hajibabaei et al., 2006) and (2) are much  
505 cheaper. In particular, we recommend the use of mini-barcodes >200-bp at the 3' end of the  
506 Folmer region. It is encouraging that such mini-barcodes perform well across a large range of  
507 metazoan taxa. These conclusions are based on three species delimitation algorithms  
508 (objective clustering, ABGD and PTP) which, overall, have no appreciable differences in  
509 performance for such mini-barcodes. If the DNA of the specimens is so degraded that very  
510 short mini-barcodes have to be obtained, we advise against the use of PTP and ABGD  
511 (especially with high priors) in order to reduce the likelihood that species are lumped.

512

## 513 **Acknowledgements**

514 We would like to acknowledge support from a Ministry of Education grant on biodiversity  
515 discovery (R-154-000-A22-112). We would also like to thank Athira Adom for data processing  
516 and Emily Hartop for proofreading.

517

## 518 **References**

519 Ahrens, D., Fujisawa, T., Krammer, H. J., Eberle, J., Fabrizi, S., & Vogler, A. P. (2016). Rarity  
520 and incomplete sampling in DNA-based species delimitation. *Systematic Biology*, 65(3), 478-  
521 494.

522 Ang, Y., Rajaratnam, G., Su, K. F., & Meier, R. (2017). Hidden in the urban parks of New York  
523 City: *Themira lohmanus*, a new species of Sepsidae described based on morphology, DNA  
524 sequences, mating behavior, and reproductive isolation (Sepsidae, Diptera). *ZooKeys*, (698),  
525 95.

526 Armani, A., Guardone, L., Castigliero, L., D'Amico, P., Messina, A., Malandra, R., ... & Guidi,  
527 A. (2015). DNA and Mini-DNA barcoding for the identification of Porgies species (family  
528 Sparidae) of commercial interest on the international market. *Food Control*, 50, 589-596.

529 Baldwin, C. C., Castillo, C. I., & Weigt, L. A. (2011). Seven new species within western Atlantic  
530 *Starksia atlantica*, *S. lepicoelia*, and *S. sluiteri* (Teleostei, Labrisomidae), with comments on  
531 congruence of DNA barcodes and species. *ZooKeys*, (79), 21.

532 Ball, S. L., Hebert, P. D., Burian, S. K., & Webb, J. M. (2005). Biological identifications of  
533 mayflies (Ephemeroptera) using DNA barcodes. *Journal of the North American Benthological*  
534 *Society*, 24(3), 508-524.

- 535 Baloğlu, B., Clews, E., & Meier, R. (2018). NGS barcoding reveals high resistance of a  
536 hyperdiverse chironomid (Diptera) swamp fauna against invasion from adjacent freshwater  
537 reservoirs. *Frontiers in zoology*, 15(1), 31.
- 538 Bar-On, Y. M., Phillips, R., & Milo, R. (2018). The biomass distribution on Earth. *Proceedings*  
539 *of the National Academy of Sciences*, 115(25), 6506-6511.
- 540 Bates, D. M. (2010). *lme4*: Mixed-effects modeling with R.
- 541 Bhattacharjee, M. J., & Ghosh, S. K. (2014). Design of Mini-barcode for Catfishes for  
542 assessment of archival biodiversity. *Molecular Ecology Resources*, 14(3), 469-477.
- 543 Bi, K., Linderoth, T., Vanderpool, D., Good, J. M., Nielsen, R., & Moritz, C. (2013). Unlocking  
544 the vault: next-generation museum population genomics. *Molecular ecology*, 22(24), 6018-  
545 6032.
- 546 Bickel, D. J. (1999). What museum collections reveal about species accumulation, richness,  
547 and rarity: an example from the Diptera. *The other 99%: the conservation and biodiversity of*  
548 *invertebrates*, 174-181.
- 549 Blaimer, B. B., Lloyd, M. W., Guillory, W. X., & Brady, S. G. (2016). Sequence capture and  
550 phylogenetic utility of genomic ultraconserved elements obtained from pinned insect  
551 specimens. *PLoS One*, 11(8), e0161531.
- 552 Burns, J. M., Janzen, D. H., Hajibabaei, M., Hallwachs, W., & Hebert, P. D. (2007). DNA  
553 barcodes of closely related (but morphologically and ecologically distinct) species of skipper  
554 butterflies (Hesperiidae) can differ by only one to three nucleotides. *Journal of the*  
555 *Lepidopterists Society*, 61(3), 138-153.
- 556 Cafaro, P. (2015). Three ways to think about the sixth mass extinction. *Biological*  
557 *Conservation*, 192, 387-393.

- 558 Ceballos, G., Ehrlich, P. R., Barnosky, A. D., García, A., Pringle, R. M., & Palmer, T. M. (2015).  
559 Accelerated modern human-induced species losses: Entering the sixth mass  
560 extinction. *Science advances*, 1(5), e1400253.
- 561 Chase, M. W., & Fay, M. F. (2009). Barcoding of plants and fungi. *Science*, 325(5941), 682-  
562 683.
- 563 Cooper, A. (1994). DNA from museum specimens. In *Ancient DNA* (pp. 149-165). Springer,  
564 New York, NY.
- 565 Cywinska, A., Hunter, F. F., & Hebert, P. D. (2006). Identifying Canadian mosquito species  
566 through DNA barcodes. *Medical and veterinary entomology*, 20(4), 413-424.
- 567 Epp, L. S., Boessenkool, S., Bellemain, E. P., Haile, J., Esposito, A., Riaz, T., ... & Stenøien,  
568 H. K. (2012). New environmental metabarcodes for analysing soil DNA: potential for studying  
569 past and present ecosystems. *Molecular Ecology*, 21(8), 1821-1833.
- 570 Ermakov, O. A., Simonov, E., Surin, V. L., Titov, S. V., Brandler, O. V., Ivanova, N. V., &  
571 Borisenko, A. V. (2015). Implications of hybridization, NUMTs, and overlooked diversity for  
572 DNA barcoding of Eurasian ground squirrels. *PLoS One*, 10(1), e0117201.
- 573 Fagan-Jeffries, E. P., Cooper, S. J., Bertozzi, T., Bradford, T. M., & Austin, A. D. (2018). DNA  
574 barcoding of microgastrine parasitoid wasps (Hymenoptera: Braconidae) using high-  
575 throughput methods more than doubles the number of species known for Australia. *Molecular*  
576 *ecology resources*, 18(5), 1132-1143.
- 577 Fan, J. A., Gu, H., Chen, S., Mo, B., Wen, Y., He, W., ... & Zeng, X. (2009). Species  
578 identification of 36 kinds of fruit flies based on minimalist-barcode. *Chinese Journal of Applied*  
579 *& Environmental Biology*, 2, 215-219.
- 580 Fields, A. T., Abercrombie, D. L., Eng, R., Feldheim, K., & Chapman, D. D. (2015). A novel  
581 mini-DNA barcoding assay to identify processed fins from internationally protected shark  
582 species. *PloS One*, 10(2), e0114844.

- 583 Folmer, O., Black, M., Hoeh, W., Lutz, R. & Vrijenhoek, R. (1994). DNA primers for  
584 amplification of mitochondrial cytochrome c oxidase subunit I from diverse metazoan  
585 invertebrates. *Molecular Marine Biology and Biotechnology*, 3(5), 294-299.
- 586 Guschanski, K., Krause, J., Sawyer, S., Valente, L. M., Bailey, S., Finstermeier, K., ... &  
587 Lenglet, G. (2013). Next-generation museomics disentangles one of the largest primate  
588 radiations. *Systematic biology*, 62(4), 539-554.
- 589 Hajibabaei, M., Smith, M. A., Janzen, D. H., Rodriguez, J. J., Whitfield, J. B., & Hebert, P. D.  
590 (2006). A minimalist barcode can identify a specimen whose DNA is degraded. *Molecular*  
591 *Ecology Notes*, 6(4), 959-964.
- 592 Hajibabaei, M., & McKenna, C. (2012). DNA mini-barcodes. In *DNA barcodes* (pp. 339-353).  
593 Humana Press, Totowa, NJ.
- 594 Han, T., Lee, W., Lee, S., Park, I. G., & Park, H. (2016). Reassessment of species diversity of  
595 the subfamily Denticollinae (Coleoptera: Elateridae) through DNA Barcoding. *PloS one*, 11(2),  
596 e0148602.
- 597 Hebert, P. D., Cywinska, A., Ball, S. L., & Dewaard, J. R. (2003). Biological identifications  
598 through DNA barcodes. *Proceedings of the Royal Society of London. Series B: Biological*  
599 *Sciences*, 270(1512), 313-321.
- 600 Hebert, P. D., Zakharov, E. V., Prosser, S. W., Sones, J. E., McKeown, J. T., Mantle, B., & La  
601 Salle, J. (2013). A DNA 'Barcode Blitz': Rapid digitization and sequencing of a natural history  
602 collection. *PLoS One*, 8(7), e68535.
- 603 Hebert, P. D., Braukmann, T. W., Prosser, S. W., Ratnasingham, S., Ivanova, N. V., Janzen,  
604 D. H., ... & Zakharov, E. V. (2018). A Sequel to Sanger: amplicon sequencing that scales.  
605 *BMC genomics*, 19(1), 219.

- 606 Hollatz, C., Leite, B. R., Lobo, J., Froufe, H., Egas, C., & Costa, F. O. (2016). Priming of a  
607 DNA metabarcoding approach for species identification and inventory in marine macrobenthic  
608 communities. *Genome*, 60(3), 260-271.
- 609 Hollingsworth, P. M., Graham, S. W., & Little, D. P. (2011). Choosing and using a plant DNA  
610 barcode. *PloS one*, 6(5), e19254.
- 611 Huang, D., Meier, R., Todd, P. A., & Chou, L. M. (2008). Slow mitochondrial COI sequence  
612 evolution at the base of the metazoan tree and its implications for DNA barcoding. *Journal of*  
613 *Molecular Evolution*, 66(2), 167-174.
- 614 Jones, M., Ghoorah, A., & Blaxter, M. (2011). jmOTU and taxonator: turning DNA barcode  
615 sequences into annotated operational taxonomic units. *PLoS one*, 6(4), e19259.
- 616 Katoh, K., & Standley, D. M. (2013). MAFFT multiple sequence alignment software version 7:  
617 improvements in performance and usability. *Molecular biology and evolution*, 30(4), 772-780.
- 618 Kemp, C. (2015). Museums: The endangered dead. *Nature News*, 518(7539), 292.
- 619 Krell, F. T. (2004). Parataxonomy vs. taxonomy in biodiversity studies—pitfalls and applicability  
620 of 'morphospecies' sorting. *Biodiversity & Conservation*, 13(4), 795-812.
- 621 Landi, M., Dimech, M., Arculeo, M., Biondo, G., Martins, R., Carneiro, M., ... & Costa, F. O.  
622 (2014). DNA barcoding for species assignment: the case of Mediterranean marine  
623 fishes. *PLoS One*, 9(9), e106135.
- 624 Lenth, R. (2018). *Emmeans*: Estimated marginal means, aka least-squares means. *R*  
625 *Package Version*, 1(2).
- 626 Lim, N. K., Tay, Y. C., Srivathsan, A., Tan, J. W., Kwik, J. T., Baloğlu, B., ... & Yeo, D. C.  
627 (2016). Next-generation freshwater bioassessment: eDNA metabarcoding with a conserved  
628 metazoan primer reveals species-rich and reservoir-specific communities. *Royal Society*  
629 *Open Science*, 3(11), 160635.



- 630 Lister, A. M., & Climate Change Research Group. (2011). Natural history collections as  
631 sources of long-term datasets. *Trends in ecology & evolution*, 26(4), 153-154.
- 632 Little, D. P. (2014). A DNA mini-barcode for land plants. *Molecular Ecology Resources*, 14(3),  
633 437-446.
- 634 Meier, R., Shiyang, K., Vaidya, G., & Ng, P. K. (2006). DNA barcoding and taxonomy in Diptera:  
635 a tale of high intraspecific variability and low identification success. *Systematic biology*, 55(5),  
636 715-728.
- 637 Meier, R. (2008). DNA sequences in taxonomy: opportunities and challenges. *The New*  
638 *Taxonomy* (ed. Wheeler QD), 7, 95–127. CRC Press, New York.
- 639 Meier, R., Wong, W., Srivathsan, A., & Foo, M. (2016). \$1 DNA barcodes for reconstructing  
640 complex phenomes and finding rare species in specimen-rich samples. *Cladistics*, 32(1), 100-  
641 110.
- 642 Meusnier, I., Singer, G. A., Landry, J. F., Hickey, D. A., Hebert, P. D., & Hajibabaei, M. (2008).  
643 A universal DNA mini-barcode for biodiversity analysis. *BMC Genomics*, 9(1), 214.
- 644 Min, X. J., & Hickey, D. A. (2007). Assessing the effect of varying sequence length on DNA  
645 barcoding of fungi. *Molecular Ecology Resources*, 7(3), 365-373.
- 646 Pentinsaari, M., Salmela, H., Mutanen, M., & Roslin, T. (2016). Molecular evolution of a widely-  
647 adopted taxonomic marker (COI) across the animal tree of life. *Scientific Reports*, 6, 35275.
- 648 Puillandre, N., Lambert, A., Brouillet, S., & Achaz, G. (2012). ABGD, Automatic Barcode Gap  
649 Discovery for primary species delimitation. *Molecular Ecology*, 21(8), 1864-1877.
- 650 R Core Team (2017). R: A language and environment for statistical computing. *R Foundation*  
651 *for Statistical Computing*, Vienna, Austria. URL <http://www.R-project.org/>.
- 652 Ratnasingham, S., & Hebert, P. D. (2013). A DNA-based registry for all animal species: the  
653 Barcode Index Number (BIN) system. *PloS One*, 8(7), e66213.

- 654 Renaud, A. K., Savage, J., & Adamowicz, S. J. (2012). DNA barcoding of Northern Nearctic  
655 Muscidae (Diptera) reveals high correspondence between morphological and molecular  
656 species limits. *BMC ecology*, 12(1), 24.
- 657 Riedel, A., Daawia, D., & Balke, M. (2010). Deep *cox1* divergence and hyperdiversity of  
658 *Trigonopterus* weevils in a New Guinea mountain range (Coleoptera,  
659 Curculionidae). *Zoologica Scripta*, 39(1), 63-74.
- 660 Roe, A.D. and Sperling, F.A. (2007). Patterns of evolution of mitochondrial cytochrome c  
661 oxidase I and II DNA and implications for DNA barcoding. *Molecular Phylogenetics and*  
662 *Evolution*, 44(1), 325-345.
- 663 Rubinoff, D., Cameron, S., & Will, K. (2006). A genomic perspective on the shortcomings of  
664 mitochondrial DNA for “barcoding” identification. *Journal of heredity*, 97(6), 581-594.
- 665 Sánchez-Bayo, F., & Wyckhuys, K. A. (2019). Worldwide decline of the entomofauna: A review  
666 of its drivers. *Biological Conservation*, 232, 8-27.
- 667 Shokralla, S., Hellberg, R. S., Handy, S. M., King, I., & Hajibabaei, M. (2015). A DNA mini-  
668 barcoding system for authentication of processed fish products. *Scientific Reports*, 5, 15894.
- 669 Smith, M. A., Rodriguez, J. J., Whitfield, J. B., Deans, A. R., Janzen, D. H., Hallwachs, W., &  
670 Hebert, P. D. (2008). Extreme diversity of tropical parasitoid wasps exposed by iterative  
671 integration of natural history, DNA barcoding, morphology, and collections. *Proceedings of the*  
672 *National Academy of Sciences*, 105(34), 12359-12364.
- 673 Srivathsan, A., & Meier, R. (2012). On the inappropriate use of Kimura-2-parameter (K2P)  
674 divergences in the DNA-barcoding literature. *Cladistics*, 28(2), 190-194.
- 675 Srivathsan, A., Sha, J., Vogler, A. P., & Meier, R. (2015). Comparing the effectiveness of  
676 metagenomics and metabarcoding for diet analysis of a leaf-feeding monkey (*Pygathrix*  
677 *nemaeus*). *Molecular Ecology Resources*, 15(2), 250-261.

- 678 Srivathsan, A., Baloglu, B., Wang, W., Tan, W. X., Bertrand, D., Ng, A. H., ... & Meier, R.  
679 (2018). A Min ION™-based pipeline for fast and cost-effective DNA barcoding. *Molecular*  
680 *Ecology Resources*, 0, 1-15.
- 681 Stamatakis, A. (2014). RAxML version 8: a tool for phylogenetic analysis and post-analysis of  
682 large phylogenies. *Bioinformatics*, 30(9), 1312-1313.
- 683 Stork, N. E., McBroom, J., Gely, C., & Hamilton, A. J. (2015). New approaches narrow global  
684 species estimates for beetles, insects, and terrestrial arthropods. *Proceedings of the National*  
685 *Academy of Sciences*, 112(24), 7519-7523.
- 686 Sultana, S., Ali, M. E., Hossain, M. M., Naquiah, N., & Zaidul, I. S. M. (2018). Universal mini  
687 COI barcode for the identification of fish species in processed products. *Food Research*  
688 *International*, 105, 19-28.
- 689 Tamura, K., Stecher, G., Peterson, D., Filipiński, A., & Kumar, S. (2013). MEGA6: molecular  
690 evolutionary genetics analysis version 6.0. *Molecular biology and evolution*, 30(12), 2725-  
691 2729.
- 692 Tan, D. S., Ang, Y., Lim, G. S., Ismail, M. R. B., & Meier, R. (2010). From 'cryptic species' to  
693 integrative taxonomy: an iterative process involving DNA sequences, morphology, and  
694 behaviour leads to the resurrection of *Sepsis pyrrhosoma* (Sepsidae: Diptera). *Zoologica*  
695 *Scripta*, 39(1), 51-61.
- 696 Tamura, K., Stecher, G., Peterson, D., Filipiński, A., & Kumar, S. (2013). MEGA6: molecular  
697 evolutionary genetics analysis version 6.0. *Molecular Biology and Evolution*, 30(12), 2725-  
698 2729.
- 699 Taylor, H. R., & Harris, W. E. (2012). An emergent science on the brink of irrelevance: a review  
700 of the past 8 years of DNA barcoding. *Molecular Ecology Resources*, 12(3), 377-388.

- 701 Tsukihara, T., Aoyama, H., Yamashita, E., Tomizaki, T., Yamaguchi, H., Shinzawa-Itoh, K., ...  
702 & Yoshikawa, S. (1996). The whole structure of the 13-subunit oxidized cytochrome c oxidase  
703 at 2.8 Å. *Science*, 272(5265), 1136-1144.
- 704 Valan, M., Makonyi, K., Maki, A., Vondráček, D., & Ronquist, F. (2019). Automated Taxonomic  
705 Identification of Insects with Expert-Level Accuracy Using Effective Feature Transfer from  
706 Convolutional Networks. *Systematic biology*, syz014, <https://doi.org/10.1093/sysbio/syz014>.
- 707 Wang, W. Y., Srivathsan, A., Foo, M., Yamane, S. K., & Meier, R. (2018). Sorting specimen-  
708 rich invertebrate samples with cost-effective NGS barcodes: Validating a reverse workflow for  
709 specimen processing. *Molecular Ecology Resources*, 18(3), 490-501.
- 710 Yang, Z., & Rannala, B. (2010). Bayesian species delimitation using multilocus sequence  
711 data. *Proceedings of the National Academy of Sciences*, 107(20), 9264-9269.
- 712 Yang, C., Tan, S., Meng, G., Bourne, D. G., O'brien, P. A., Xu, J., ... & Liu, S. (2018). Access  
713 COI barcode efficiently using high throughput Single End 400 bp sequencing. *BioRxiv*, 498618.  
714 doi: <http://dx.doi.org/10.1101/498618>.
- 715 Yeates, D. K., Zwick, A., & Mikheyev, A. S. (2016). Museums are biobanks: unlocking the  
716 genetic potential of the three billion specimens in the world's biological collections. *Current*  
717 *opinion in insect science*, 18, 83-88.
- 718 Yeo, D., Puniamoorthy, J., Ngiam, R. W. J., & Meier, R. (2018). Towards holomorphology in  
719 entomology: rapid and cost-effective adult–larva matching using NGS barcodes. *Systematic*  
720 *entomology*, 43(4), 678-691.
- 721 Yu, H. J., & You, Z. H. (2010). Comparison of DNA truncated barcodes and full-barcodes for  
722 species identification. In *Advanced Intelligent Computing Theories and Applications. With*  
723 *Aspects of Artificial Intelligence* (pp. 108-114). Springer, Berlin, Heidelberg.
- 724 Zhang, J., Kapli, P., Pavlidis, P., & Stamatakis, A. (2013). A general species delimitation  
725 method with applications to phylogenetic placements. *Bioinformatics*, 29(22), 2869-2876.

726 Zuccon, D., Brisset, J., Corbari, L., Puillandre, N., Utge, J., & Samadi, S. (2012). An optimised  
727 protocol for barcoding museum collections of decapod crustaceans: a case-study for a 10–  
728 40-years-old collection. *Invertebrate Systematics*, 26(6), 592-600.

729 **Table 1.** Proportion of specimens congruent between morphospecies and mOTU clusters  
 730 under the three stringency classes. Values in brackets represent the estimated number of  
 731 specimens causing conflict.

| <b>Objective clustering, 2% p-distance</b> |               |               |               |               |               |                |               |
|--|---------------|---------------|---------------|---------------|---------------|----------------|---------------|
| <b>Length</b>                              | <b>295</b>    | <b>307</b>    | <b>313</b>    | <b>407</b>    | <b>657</b>    | <b>Average</b> | <b>Median</b> |
| <b>Midpoint</b>                            | <b>405</b>    | <b>154</b>    | <b>502</b>    | <b>455</b>    | <b>329</b>    |                |               |
| <b>Class I</b>                             | 74%<br>(7475) | 74%<br>(7511) | 75%<br>(7275) | 75%<br>(7246) | 76%<br>(7118) | 75%<br>(7325)  | 75%<br>(7372) |
| <b>Class II</b>                            | 89%<br>(2049) | 89%<br>(2214) | 89%<br>(2079) | 89%<br>(2030) | 90%<br>(1967) | 89%<br>(2068)  | 89%<br>(2079) |
| <b>Class III</b>                           | 90%<br>(975)  | 90%<br>(743)  | 90%<br>(746)  | 90%<br>(786)  | 91%<br>(795)  | 90%<br>(809)   | 90%<br>(802)  |
| <b>ABGD, P=0.001</b>                       |               |               |               |               |               |                |               |
| <b>Length</b>                              | <b>295</b>    | <b>307</b>    | <b>313</b>    | <b>407</b>    | <b>657</b>    | <b>Average</b> | <b>Median</b> |
| <b>Midpoint</b>                            | <b>405</b>    | <b>154</b>    | <b>502</b>    | <b>455</b>    | <b>329</b>    |                |               |
| <b>Class I</b>                             | 73%<br>(7916) | 74%<br>(7685) | 72%<br>(8169) | 71%<br>(8583) | 66%<br>(9833) | 71%<br>(8437)  | 72%<br>(8126) |
| <b>Class II</b>                            | 89%<br>(2178) | 88%<br>(2014) | 89%<br>(2138) | 88%<br>(2129) | 87%<br>(2604) | 88%<br>(2213)  | 88%<br>(2057) |
| <b>Class III</b>                           | 90%<br>(806)  | 90%<br>(895)  | 90%<br>(849)  | 90%<br>(915)  | 88%<br>(881)  | 89%<br>(869)   | 90%<br>(862)  |
| <b>PTP</b>                                 |               |               |               |               |               |                |               |
| <b>Length</b>                              | <b>295</b>    | <b>307</b>    | <b>313</b>    | <b>407</b>    | <b>657</b>    | <b>Average</b> | <b>Median</b> |
| <b>Midpoint</b>                            | <b>405</b>    | <b>154</b>    | <b>502</b>    | <b>455</b>    | <b>329</b>    |                |               |
| <b>Class I</b>                             | 74%<br>(7643) | 74%<br>(7744) | 75%<br>(7369) | 76%<br>(7130) | 75%<br>(7207) | 75%<br>(7419)  | 75%<br>(7318) |
| <b>Class II</b>                            | 89%<br>(2362) | 88%<br>(2433) | 89%<br>(2227) | 89%<br>(2095) | 89%<br>(1965) | 89%<br>(2216)  | 89%<br>(2179) |
| <b>Class III</b>                           | 89%<br>(732)  | 89%<br>(904)  | 90%<br>(702)  | 90%<br>(721)  | 90%<br>(897)  | 90%<br>(791)   | 90%<br>(733)  |

732

733

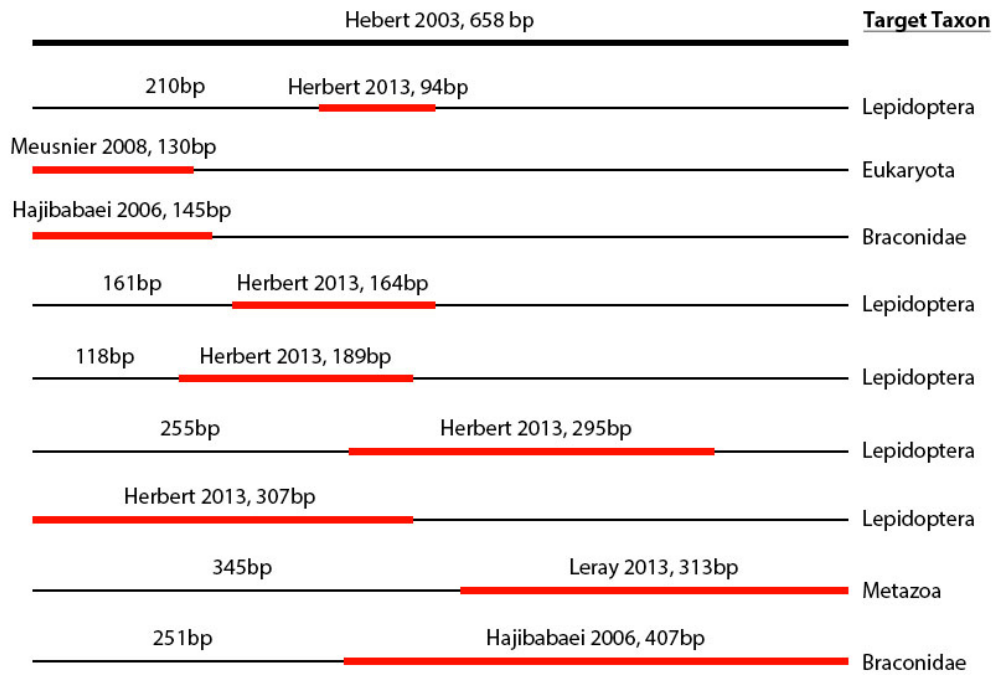
734

735

736

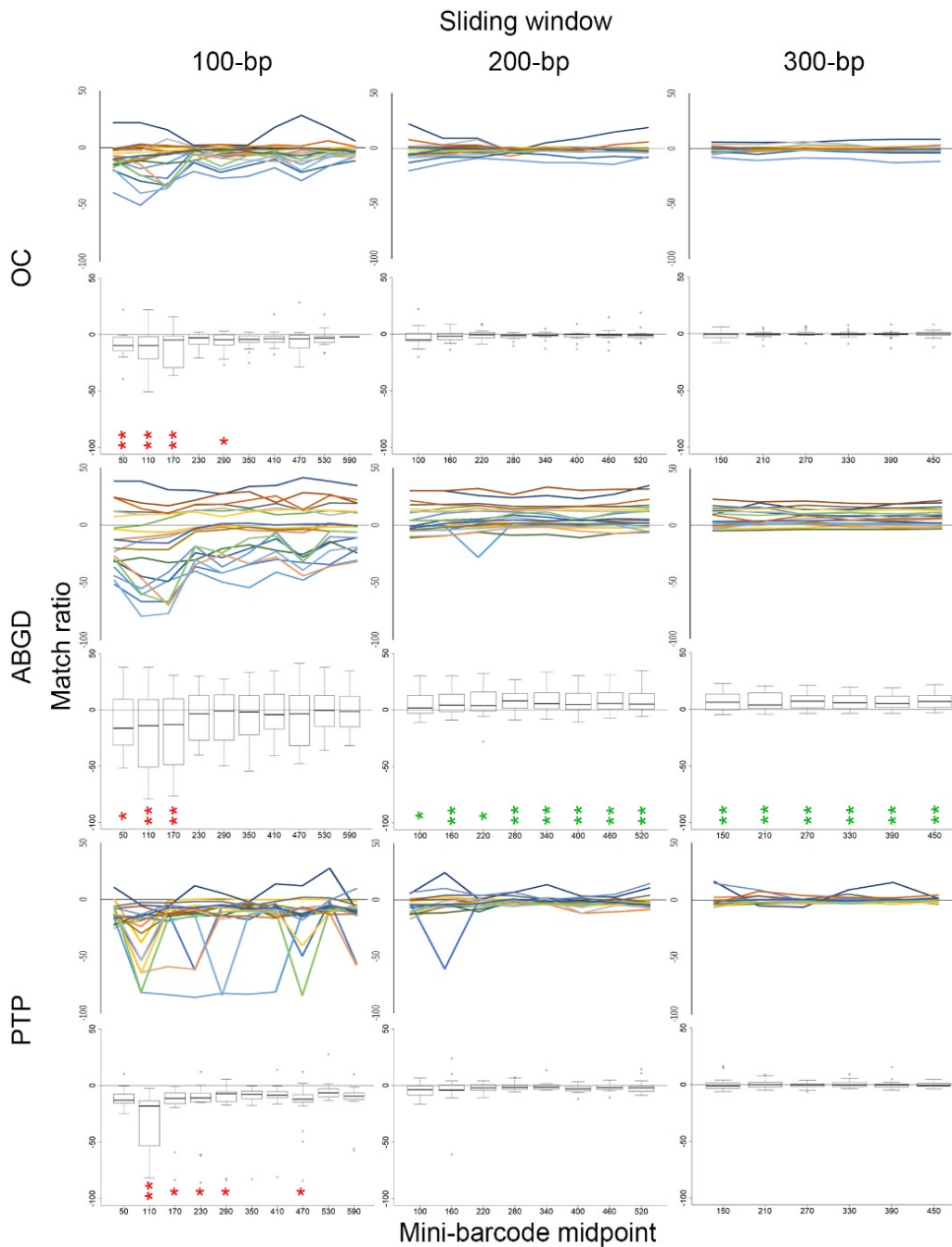
737

738



739

740 **Figure 1.** Position of the mini-barcode with established primers in this study.



741

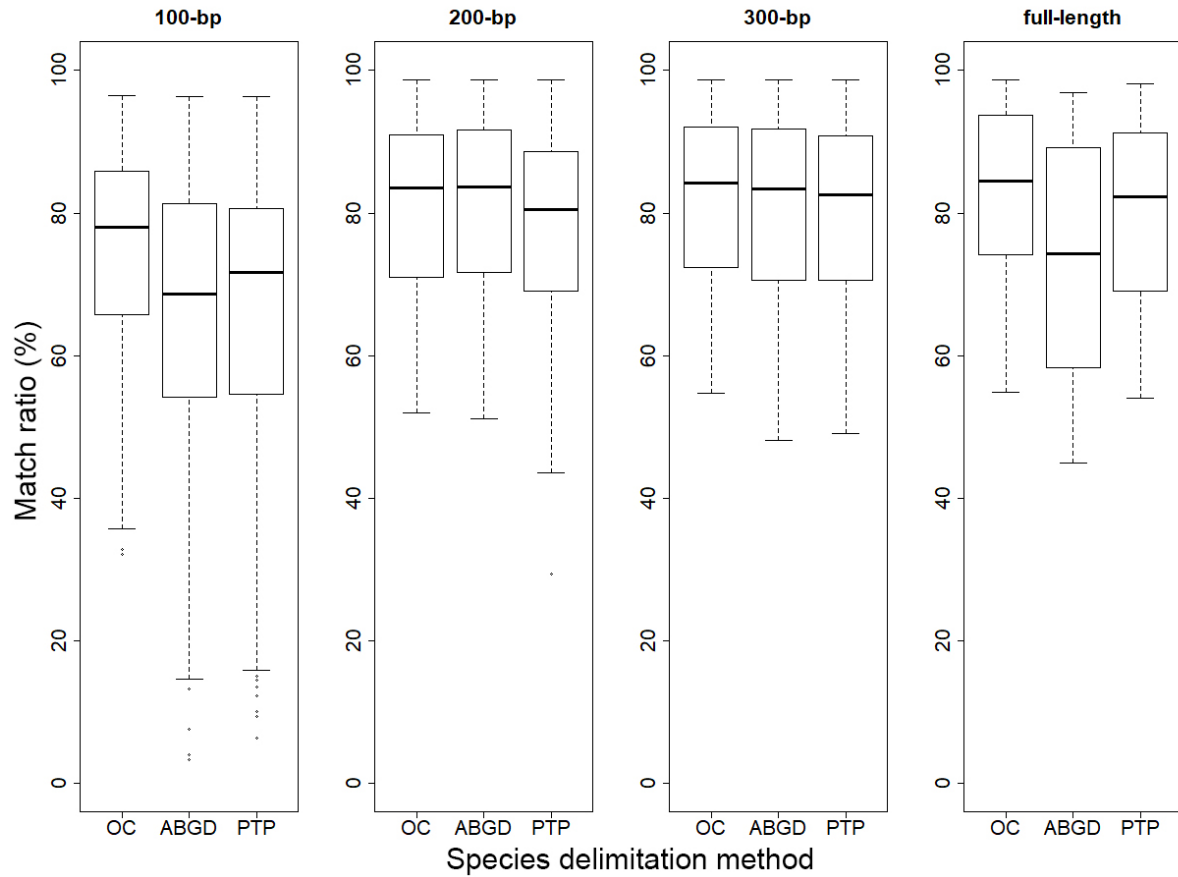
742 **Figure 2.** Performance of mini-barcodes along a sliding window (100, 200, 300-bp). Mini-  
743 barcode position is indicated on the x-axis and congruence with morphology on the y-axis.  
744 mOTUs were obtained with Objective Clustering (2%), ABGD P=0.001 prior), and PTP. Each  
745 line represents one data set while the boxplots summarise the values across datasets.  
746 Significant deviations from the results obtained with full-length barcodes are indicated with  
747 color-coded asterisks (\* =  $p < 0.05$ ; \*\* =  $p < 0.001$ ; red = poorer and green = higher  
748 congruence with morphology).



| Objective clustering       |           |                   |     |           |          |     |           |     |     |     | ABGD                       |           |                   |     |           |          |     |           |     |     |     | PTP              |                            |                   |     |           |          |     |           |     |     |     |    |
|----------------------------|-----------|-------------------|-----|-----------|----------|-----|-----------|-----|-----|-----|----------------------------|-----------|-------------------|-----|-----------|----------|-----|-----------|-----|-----|-----|------------------|----------------------------|-------------------|-----|-----------|----------|-----|-----------|-----|-----|-----|----|
| Length                     | 94        | 130               | 145 | 164       | 189      | 295 | 307       | 313 | 407 | 657 | Length                     | 94        | 130               | 145 | 164       | 189      | 295 | 307       | 313 | 407 | 657 | Length           | 94                         | 130               | 145 | 164       | 189      | 295 | 307       | 313 | 407 | 657 |    |
| Midpoint                   | 257       | 65                | 73  | 243       | 213      | 405 | 154       | 502 | 455 | 329 | Midpoint                   | 257       | 65                | 73  | 243       | 213      | 405 | 154       | 502 | 455 | 329 | Midpoint         | 257                        | 65                | 73  | 243       | 213      | 405 | 154       | 502 | 455 | 329 |    |
| Great Barrier Reef Fish    | 94        | 91                | 92  | 95        | 95       | 97  | 94        | 97  | 97  | 97  | North Sea Molluscs         | 94        | 90                | 92  | 94        | 94       | 89  | 90        | 89  | 86  | 80  | 90               | Great Barrier Reef Fish    | 89                | 85  | 87        | 91       | 91  | 95        | 90  | 95  | 96  | 96 |
| South China Sea Fish       | 94        | 95                | 96  | 98        | 88       | 95  | 93        | 95  | 95  | 97  | South China Sea Fish       | 92        | 78                | 83  | 89        | 93       | 94  | 92        | 93  | 91  | 90  | 89               | Pakistan Lepidoptera       | 95                | 83  | 82        | 87       | 90  | 95        | 95  | 95  | 95  | 94 |
| North Sea Molluscs         | 91        | 86                | 87  | 94        | 91       | 94  | 94        | 93  | 93  | 95  | Great Barrier Reef Fish    | 79        | 75                | 79  | 90        | 89       | 97  | 92        | 95  | 95  | 94  | 88               | Ecuador Geometridae        | 14                | 89  | 91        | 94       | 97  | 98        | 97  | 96  | 99  | 98 |
| Canada Echinoderms         | 81        | 89                | 85  | 91        | 93       | 91  | 91        | 91  | 94  | 95  | Canada Echinoderms         | 87        | 78                | 76  | 91        | 91       | 87  | 91        | 95  | 92  | 89  | 88               | North Sea Molluscs         | 78                | 71  | 82        | 94       | 96  | 86        | 92  | 86  | 93  | 93 |
| Ecuador Geometridae        | 78        | 70                | 73  | 85        | 98       | 99  | 99        | 99  | 99  | 99  | Pakistan Lepidoptera       | 61        | 67                | 75  | 84        | 89       | 95  | 91        | 95  | 96  | 94  | 85               | South China Sea Fish       | 86                | 75  | 77        | 87       | 88  | 92        | 89  | 89  | 92  | 92 |
| <b>Set Mean:</b>           | <b>92</b> | <b>Set Range:</b> |     | <b>70</b> | <b>-</b> |     | <b>99</b> |     |     |     | <b>Set Mean:</b>           | <b>88</b> | <b>Set Range:</b> |     | <b>61</b> | <b>-</b> |     | <b>97</b> |     |     |     | <b>Set Mean:</b> | <b>89</b>                  | <b>Set Range:</b> |     | <b>14</b> | <b>-</b> |     | <b>99</b> |     |     |     |    |
| Germany Aranea & Opiliones | 83        | 89                | 88  | 90        | 84       | 84  | 85        | 84  | 85  | 85  | Germany EPT                | 85        | 84                | 85  | 86        | 87       | 83  | 85        | 84  | 78  | 67  | 82               | Canada Echinoderms         | 75                | 71  | 78        | 83       | 90  | 89        | 90  | 87  | 91  | 89 |
| European Marine Fish       | 87        | 88                | 89  | 90        | 81       | 83  | 85        | 83  | 83  | 84  | Ecuador Geometridae        | 48        | 58                | 57  | 96        | 70       | 99  | 99        | 99  | 99  | 97  | 82               | Germany EPT                | 82                | 80  | 81        | 81       | 83  | 79        | 83  | 82  | 79  | 80 |
| South America Butterflies  | 72        | 74                | 73  | 77        | 90       | 92  | 93        | 92  | 93  | 93  | South America Butterflies  | 60        | 70                | 77  | 84        | 84       | 89  | 91        | 91  | 91  | 81  | 82               | North Europe Tachinidae    | 77                | 70  | 73        | 82       | 80  | 84        | 83  | 83  | 82  | 85 |
| Germany EPT                | 78        | 85                | 84  | 85        | 85       | 85  | 82        | 84  | 85  | 85  | European Marine Fish       | 80        | 72                | 72  | 85        | 86       | 83  | 84        | 80  | 78  | 69  | 79               | Amazon Moths               | 32                | 88  | 38        | 89       | 88  | 91        | 89  | 91  | 90  | 90 |
| Northwest Pacific Molluscs | 88        | 87                | 90  | 91        | 83       | 78  | 77        | 77  | 78  | 78  | North America Birds        | 61        | 67                | 73  | 74        | 76       | 81  | 83        | 83  | 84  | 87  | 77               | Germany Aranea & Opiliones | 72                | 73  | 77        | 77       | 79  | 82        | 80  | 77  | 81  | 79 |
| <b>Set Mean:</b>           | <b>84</b> | <b>Set Range:</b> |     | <b>72</b> | <b>-</b> |     | <b>93</b> |     |     |     | <b>Set Mean:</b>           | <b>80</b> | <b>Set Range:</b> |     | <b>48</b> | <b>-</b> |     | <b>99</b> |     |     |     | <b>Set Mean:</b> | <b>80</b>                  | <b>Set Range:</b> |     | <b>32</b> | <b>-</b> |     | <b>91</b> |     |     |     |    |
| Amazon Moths               | 69        | 77                | 76  | 79        | 91       | 89  | 92        | 84  | 85  | 87  | Amazon Moths               | 72        | 43                | 42  | 88        | 75       | 87  | 90        | 90  | 91  | 89  | 77               | North America Birds        | 72                | 67  | 66        | 77       | 78  | 81        | 82  | 81  | 85  | 85 |
| North America Birds        | 84        | 77                | 77  | 82        | 77       | 83  | 83        | 83  | 83  | 84  | North Europe Tachinidae    | 71        | 60                | 66  | 77        | 80       | 77  | 84        | 82  | 82  | 84  | 76               | French Guiana Earthworms   | 71                | 83  | 69        | 85       | 80  | 81        | 62  | 78  | 75  | 69 |
| French Guiana Earthworms   | 77        | 85                | 84  | 90        | 78       | 74  | 75        | 78  | 77  | 84  | French Guiana Earthworms   | 78        | 89                | 93  | 81        | 81       | 73  | 80        | 70  | 59  | 55  | 76               | European Marine Fish       | 71                | 60  | 61        | 71       | 78  | 80        | 74  | 79  | 79  | 80 |
| Pakistan Lepidoptera       | 52        | 74                | 74  | 54        | 90       | 91  | 90        | 91  | 91  | 91  | Northwest Pacific Molluscs | 76        | 75                | 75  | 77        | 77       | 76  | 77        | 76  | 76  | 63  | 75               | Northwest Pacific Molluscs | 72                | 67  | 68        | 75       | 74  | 72        | 74  | 76  | 77  | 78 |
| Tanytarsus                 | 81        | 81                | 82  | 81        | 63       | 61  | 62        | 63  | 61  | 65  | Germany Aranea & Opiliones | 76        | 18                | 79  | 85        | 82       | 71  | 74        | 73  | 65  | 52  | 67               | Congo Fish                 | 57                | 44  | 50        | 67       | 69  | 68        | 64  | 68  | 71  | 69 |
| <b>Set Mean:</b>           | <b>79</b> | <b>Set Range:</b> |     | <b>52</b> | <b>-</b> |     | <b>92</b> |     |     |     | <b>Set Mean:</b>           | <b>74</b> | <b>Set Range:</b> |     | <b>18</b> | <b>-</b> |     | <b>93</b> |     |     |     | <b>Set Mean:</b> | <b>72</b>                  | <b>Set Range:</b> |     | <b>44</b> | <b>-</b> |     | <b>85</b> |     |     |     |    |
| North Europe Tachinidae    | 47        | 51                | 49  | 57        | 77       | 80  | 82        | 82  | 80  | 81  | Congo Fish                 | 68        | 44                | 55  | 67        | 71       | 70  | 63        | 67  | 71  | 65  | 64               | Tanytarsus                 | 51                | 58  | 64        | 61       | 58  | 68        | 63  | 64  | 66  | 63 |
| Congo Fish                 | 63        | 61                | 60  | 68        | 70       | 65  | 65        | 65  | 65  | 60  | Tanytarsus                 | 67        | 71                | 71  | 61        | 62       | 50  | 56        | 48  | 47  | 45  | 58               | North America Pyraustinae  | 45                | 58  | 54        | 56       | 49  | 58        | 71  | 80  | 76  | 58 |
| North America Pyraustinae  | 54        | 54                | 55  | 55        | 57       | 72  | 66        | 70  | 71  | 70  | North America Pyraustinae  | 25        | 29                | 41  | 67        | 76       | 70  | 74        | 67  | 58  | 59  | 57               | Ecuador Chrysomelidae      | 55                | 57  | 55        | 56       | 55  | 55        | 54  | 55  | 55  | 55 |
| Ecuador Chrysomelidae      | 61        | 56                | 64  | 69        | 63       | 58  | 59        | 55  | 56  | 57  | Ecuador Chrysomelidae      | 55        | 55                | 54  | 56        | 56       | 57  | 57        | 55  | 57  | 58  | 56               | Iberia Butterflies         | 32                | 47  | 48        | 47       | 48  | 59        | 58  | 61  | 58  | 54 |
| Iberia Butterflies         | 51        | 55                | 56  | 59        | 56       | 55  | 55        | 55  | 55  | 55  | Iberia Butterflies         | 26        | 36                | 42  | 53        | 57       | 60  | 65        | 59  | 55  | 52  | 50               | South America Butterflies  | 0                 | 0   | 0         | 0.5      | 0   | 89        | 90  | 93  | 91  | 91 |
| <b>Set Mean:</b>           | <b>62</b> | <b>Set Range:</b> |     | <b>47</b> | <b>-</b> |     | <b>82</b> |     |     |     | <b>Set Mean:</b>           | <b>57</b> | <b>Set Range:</b> |     | <b>25</b> | <b>-</b> |     | <b>76</b> |     |     |     | <b>Set Mean:</b> | <b>55</b>                  | <b>Set Range:</b> |     | <b>0</b>  | <b>-</b> |     | <b>93</b> |     |     |     |    |

749

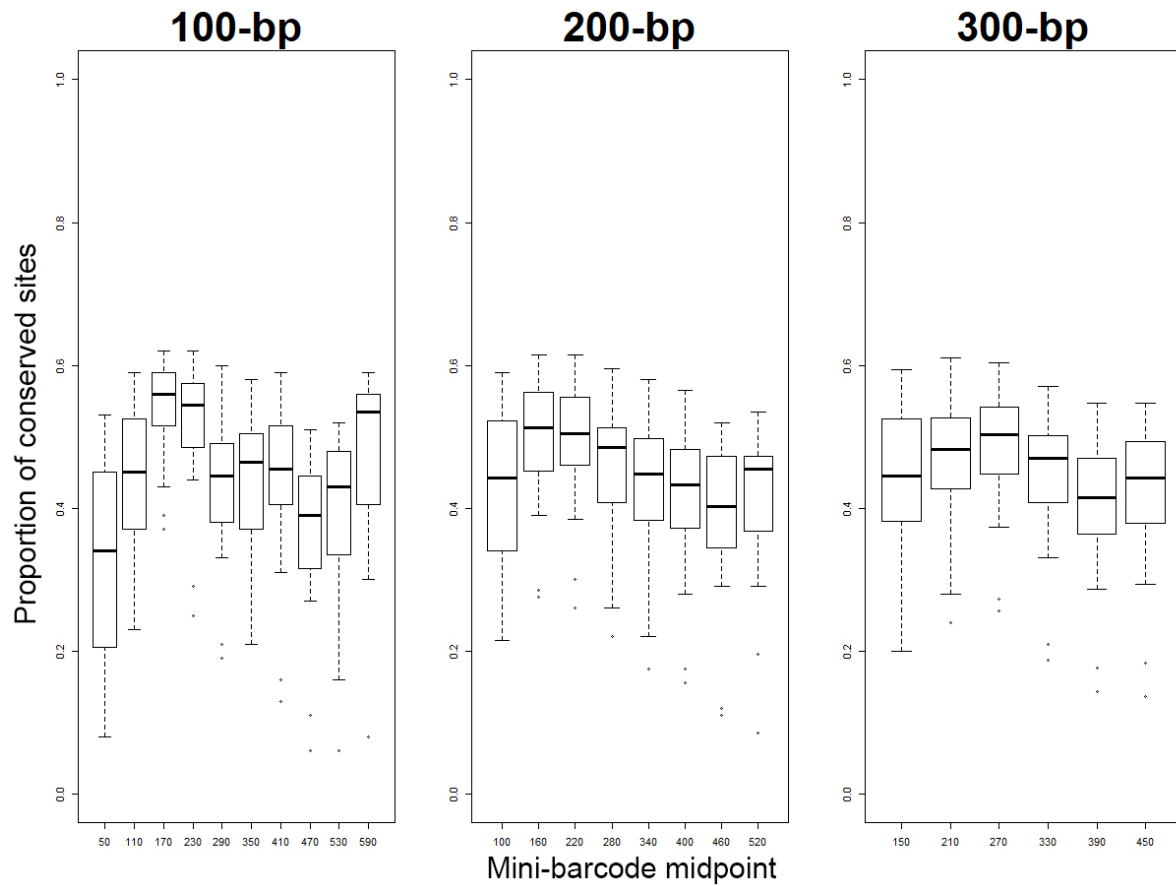
750 **Figure 3.** Match ratios across three different species delimitation methods. Mini-barcodes (columns) are sorted by primer length while the  
751 datasets (rows) are grouped into 4 classes according to average match ratio. Colours are applied separately to each class.



752

753 **Figure 4.** Comparison of species delimitation methods for full-length and mini-barcodes

754 generated by “sliding windows” (100-bp, 200-bp, 300-bp).



755

756 **Figure 5.** Proportion of conserved sites along the full-length barcode (sliding windows of  
757 100-bp, 200-bp, 300-bp).

758