

1 **Running head: Multilocus DNA barcoding**

2 **Title: Multilocus DNA barcoding – Species Identification with Multilocus Data**

3 **Authors:** Junning Liu^a, Jiamei Jiang^a, Shuli Song^a, Luke Tornabene^b, Ryan

4 Chabarria^c, Gavin J P Naylor^d, Chenhong Li^{a*}

5 **Affiliations:**

6 ^aKey Laboratory of Exploration and Utilization of Aquatic Genetic Resources,

7 Shanghai Ocean University, Ministry of Education, Shanghai 201306, China;

8 ^bSchool of Aquatic and Fisheries Sciences, University of Washington, Seattle, WA

9 98195, USA;

10 ^cCollege of Science & Engineering, Texas A&M University – Corpus Christi, Corpus

11 Christi, TX 78412-5806, USA;

12 ^dHollings Marine Laboratory, College of Charleston, Charleston, SC 29412, USA

13 *Correspondence to: Chenhong Li, chli@shou.edu.cn

14

14 **Abstract:** Species identification using DNA sequences, known as DNA barcoding has
15 been widely used in many applied fields. Current barcoding methods are usually
16 based on a single mitochondrial locus, such as cytochrome c oxidase subunit I (COI).
17 This type of barcoding is not always effective when applied to species separated by
18 short divergence times or that contain introgressed genes from closely related species.
19 Herein we introduce a more effective multi-locus barcoding framework that is based
20 on gene capture and “next-generation” sequencing and provide both empirical and
21 simulation tests of its efficacy. We examine genetic distinctness in two pairs of fishes
22 that are sister-species: *Siniperca chuatsi* vs. *S. kneri* and *Sicydium altum* vs. *S. adelum*,
23 where the COI barcoding approach failed species identification in both cases. Results
24 revealed that distinctness between *S. chuatsi* and *S. kneri* increased as more
25 independent loci were added. By contrast *S. altum* and *S. adelum* could not be
26 distinguished even with all loci. Analyses of population structure and gene flow
27 suggested that the two species of *Siniperca* diverged from each other a long time ago
28 but have unidirectional gene flow, whereas the two species of *Sicydium* are not
29 separated from each other and have high bidirectional gene flow. Simulations
30 demonstrate that under limited gene flow (< 0.00001 per gene per generation) and
31 enough separation time (> 100000 generation), we can correctly identify species using
32 more than 90 loci. Finally, we selected 500 independent nuclear markers for
33 ray-finned fishes and designed a three-step pipeline for multilocus DNA barcoding.

34 **Keywords:** DNA barcoding, incomplete lineage sorting, gene flow, nuclear gene
35 markers, gene capture, next-generation sequencing

36 DNA barcoding has been very successfully employed in many applied fields,
37 ranging from routine species identification (Sutou et al. 2011; Candek and Kuntner
38 2015; Hassold et al. 2016), to the discovery of cryptic species (Witt et al. 2006;
39 Kadarusman et al. 2012), tracking of invasive species (Saunders 2009;
40 Ghahramanzadeh et al. 2013; Marescaux and Van Doninck 2013), conservation, and
41 community ecology (Tanzler et al. 2012; Nevill et al. 2013; Hartvig et al. 2015;
42 Shapcott et al. 2015). The mitochondrial cytochrome c oxidase subunit I gene (COI)
43 has a good amount of variation and is easy to amplify using PCR based approaches in
44 most animal groups (Hebert et al. 2003; Smith et al. 2005; Vences et al. 2005; Ward et
45 al. 2005). It has become the most commonly used marker for animal DNA barcoding
46 since it was first proposed more than a decade ago (Hebert et al. 2003). In most cases,
47 single-locus (COI) DNA barcoding results in successful species identification. For
48 example, a success rate close to 100% were reported for Germany herpetofauna
49 (Hawlitschek et al. 2016), more than 90% for Chinese rodents (Li et al. 2015a), more
50 than 80% for freshwater fishes of the Congo basin (Collins et al. 2012; Decru et al.
51 2016), and 100% for mosquitoes (Chan et al. 2014). However, the success rate of
52 species identification was low for species complexes with gene flow (Hawlitschek et
53 al. 2016) or where species had only recently diverged (van Velzen et al. 2012).

54 In order to use barcoding for species identification, within species variation
55 must be less than between species variation. This generates a “break” in the
56 distribution of distances that is referred to as the “barcoding gap”. Indeed one of the
57 common causes of barcoding failure occurs when differences in demography

58 eliminate the barcoding gap, because intra-specific differences are greater than
59 inter-specific differences for the clades being compared. To an extreme, two
60 individuals could have the same COI sequence, while being distinctly different
61 species. Shared COI haplotypes have been reported in different species of spiders
62 (Spasojevic et al. 2016), birds (Aliabadian et al. 2013) and fishes (Mabragana et al.
63 2011). The single-locus barcoding is prone to misidentification when different species
64 share haplotypes.

65 Although haplotypes at a single locus, such as COI can be shared between two
66 species, it is unlikely that individuals of two species share alleles across multiple
67 independent genes. Accordingly, multilocus data should perform better for species
68 identification than any single locus could. Dowton et al. (2014) proposed
69 “next-generation” DNA barcoding based on multilocus data in which they
70 incorporated multispecies coalescent species delimitation. They analyzed *Sarcophaga*
71 flesh flies with two loci, mitochondrial COI and nuclear carbomoylphosphate
72 synthase (CAD), and found out that their coalescent-based *BEAST/BPP approach
73 was more successful than standard barcoding method (Dowton et al. 2014). However,
74 Collins and Cruickshank (2014) reanalyzed Dowton et al.’s data and showed that
75 standard single locus (COI) barcoding method could achieve the same accuracy as the
76 new multilocus framework did if an optimized distance threshold was applied (Brown
77 et al. 2012; Puillandre et al. 2012; Virgilio et al. 2012; Sonet et al. 2013).

78 The experiment of Dowton et al. (2014) seemed unsuccessful, but the likely
79 reason for this was that the data they used was not challenging enough for standard

80 single-locus barcoding methods, because there was only one unidentifiable individual
81 that was more divergent from its closest putative conspecific than the optimized
82 threshold (Collins and Cruickshank 2014). The other reason is that only a single
83 nuclear gene was used in the study of Dowton et al. (2014), thus providing little
84 additional information (Collins and Cruickshank 2014).

85 In the past it has been challenging to obtain sequences from sufficient
86 independent nuclear loci from a broad taxonomic group to make multilocus DNA
87 barcoding effective, but tools for finding thousands of nuclear gene markers (Bi et al.
88 2012; Li et al. 2012; Hedtke et al. 2013) and collecting their sequences through
89 cross-species gene capture and next-generation sequencing are now available (Li et al.
90 2013), providing an opportunity to rigorously test the power of multilocus DNA
91 barcoding. In this work, we tested the effectiveness of multilocus barcoding
92 employing hundreds of nuclear loci, to correctly identify species that were not
93 distinguishable based on only COI or a few nuclear loci.

94 MATERIALS AND METHODS

95 *Species discrimination using empirical data*

96 We have developed 4,434 single-copy nucleotide loci for ray-finned fishes,
97 and tested them in 83 species (29 families and 12 orders), covering major clades of
98 ray-finned fishes. Those markers have few missing data in the taxa tested, showing
99 promise for their deployment in phylogenetics and population genetic analyses. We
100 adopted those 4,434 loci as candidate barcoding markers in order to further optimize a
101 subset of universal markers for all ray-finned fishes. We choose loci that could be

102 readily captured and sequenced across taxa, and that were variable based on their
103 average p-distance values among taxa.

104 Some of the most challenging instances for DNA barcoding occur when taxa
105 are recently diverged or when gene flow exists between closely related species, or
106 both. In an effort to design a rigorous barcoding scheme, we picked empirical study
107 systems that would involve both challenges. The first involved sinipercid fishes, a
108 family of fishes containing two genera, 9 to 12 species depending on the authority
109 referenced (Zhou et al. 1988; Li 1991; Liu and Chen 1994; Nelson 2006). Among
110 them, two sister species, *Siniperca chuatsi* and *S. kneri* have distinct morphological
111 characters, such as number of pyloric caecum, ratio between eye length and head
112 length (Zhou et al. 1988), but they are not distinguishable using mitochondrial control
113 region sequences (Zhao et al. 2008). These two sister species are allopatric in most of
114 their distribution regions (Zhou et al. 1988; Li 1991), so the reason for unsuccessful
115 species identification in these sister species is likely due to their recency of speciation
116 (Zhao et al. 2008).

117 The other group of fishes that we used to test the multilocus DNA barcoding
118 method is *Sicydium*. *Sicydium* is a group of diadromous gobies native to fast-flowing
119 streams and rivers of the Americas (Central America, Mexico, Cocos Island, the
120 Caribbean, Colombia, Ecuador and Venezuela) and Africa. There are two syntopic
121 species, *S. altum* and *S. adelum* that could be separated according to distinct dental
122 papillae and other morphological characters (Bussing 1996), but they are
123 indistinguishable using mitochondrial or nuclear genes (Chabbarria 2015). Because

124 these two closely related species are frequently found together (Bussing 1996), it is
125 possible that they have been subject to interspecific gene flow which would account
126 for the high degree of genetic similarity between them. These two pairs of
127 sister-species were used as test cases to evaluate how gene flow and shallow
128 divergence times might affect species discrimination and identification based on
129 multilocus barcoding.

130 *Taxa sampling, target gene enrichment, sequencing and reads assembly*

131 We collected sequence data from 16,943 loci of nine siniperoid species,
132 including five *S. chuatsi* and five *S. kneri* (Song et al., accepted). Here, we reused this
133 dataset for testing multilocus barcoding. Sequences of the 4,434 loci of the siniperoids
134 were retrieved from Song et al.'s data. The samples included five *Coreoperca*
135 *whiteheadi*, one *S. scherzeri*, five *S. obscura*, two *S. undulata*, three *S. roulei*, five *S.*
136 *chuatsi* and five *S. kneri*.

137 For the goby study, nine *S. altum* and seven *S. adelum* were collected from
138 Costa Rica. Total genomic DNA was extracted from fin clips using a Tissue DNA kit
139 (Omega Bio-tek, Norcross, GA, USA) and the concentration of DNA was quantified
140 using NanoDrop 3300 Fluorospectrometer (Thermo Fisher Scientific, Wilmington,
141 DE, USA). The goby samples were enriched and sequenced for the same 4,434 loci.
142 The amount of DNA used for library preparation was 1 µg for each sample. The DNA
143 sample was first sheared to 250 bp using a Covaris M220 Focused-ultrasonicator™
144 (Covaris, Inc. Massachusetts, USA). A MYbaits kit containing baits for the 4,434 loci
145 was synthesized at MYcroarray (Ann Arbor, Michigan, USA). The baits were

146 designed on sequences of *Oreochromis niloticus* with 3 × tiling. Blunt-end repair,
147 adapter ligation, fill-in, pre-hybridization PCR and target gene enrichment steps
148 followed the protocol of cross-species gene capture (Li et al. 2013). The enriched
149 libraries were amplified with indexed primers, pooled equimolarly and sequenced on
150 a lane of Illumina HiSeq 2500 platform with other samples. The raw reads were
151 parsed to separate file for each species according to the indices on the adapter. Reads
152 assembling followed the pipeline of Yuan et al. (2016). Mitochondrial COI gene of
153 both the siniperoids and the gobies was also amplified and sequenced using Sanger
154 sequencing to compare COI barcoding with multilocus DNA barcoding using two
155 pairs of primers (siniF: AACCAGCGAGCATCCATCTA and siniR:
156 CAGTGGACGAAAGCAGCAAC for the siniperoids; sicyF:
157 GGTTGTGTTGAGGTTTCGGT and sicyR: TCCGAGCCGAACTAAGTCAA for
158 *Sicydium*).

159 *Effect of increasing number of loci on species discrimination*

160 Our assumption was that individuals of recently diverged species should be
161 more discernible using many loci than using fewer loci. Thus, we calculated
162 p-distance among 10 individuals of *Siniperca*, including five *S. chuatsi* and five *S.*
163 *kneri*, using different number of loci to test this hypothesis. Loci with no missing data
164 in all 10 individuals of *Siniperca* were picked using a custom Perl scripts
165 (picktaxagene.pl, Supplementary Materials). The obtained 2,612 loci were then sorted
166 by their average p-distance (distoutlier.pl, Supplementary Materials), so outlier loci
167 with extreme large p-distance could be checked by eye to spot bad data or bad

168 alignment. After removing the bad data, a different number (1, 3, 10, 30, 50, 70, 90,
169 100, 200, 300, 400, 500, 600, 700, 800, 900, 1000, 2000) of loci were randomly
170 picked and concatenated (samplegene.pl, Supplementary Materials) for calculating
171 p-distance among individuals (gapdis.pl, Supplementary Materials). The sampling at
172 each level of different number of loci was repeated two hundred times. The p-distance
173 among individuals vs. the number of loci used was drawn with GraphPad Prism 5
174 (San Diego, California). To check the effect of increasing number of loci on species
175 discrimination, the “all species barcodes” criteria was applied, that is queries was
176 considered successfully identified when they were followed by all conspecifics
177 according to their barcode (Meier et al. 2006). Custom Perl script was used to
178 calculate the rate of successful identification for 200 replicates at each level of
179 number of loci used (ID_correct_rate.pl, Supplementary Materials). Among individual
180 p-distance and rate of successful identification also were calculated for the *Sicydium*.
181 Sequences of COI gene also were used to calculate p-distance between individuals
182 from the same species and from different species to compare with the results of
183 nuclear genes. Spider (Brown et al. 2012) was used to optimize barcoding distance
184 threshold and to identify species using COI sequences as suggested by Collins and
185 Cruickshank (2014). The final number of loci recommended for DNA barcoding was
186 chosen based on the effect of increasing number of loci on the success rate of species
187 discrimination.

188 *Estimating species divergence and gene flow in the empirical data*

189 Gene flow and differentiation time of *S. chuatsi* and *S. kneri* was estimated

190 using IMA2 program with 200 loci (Hey 2010). The MCMC was run for 10 million
191 generations with sample recorded every hundred generations. The number of chains
192 was set to 20. The running parameters were set as -q2, -m1, -t3, -b 10000000, d100,
193 -hn20 and -s123. An additional run was performed with the same parameter but
194 different seeds -s111. These two run showed decent mixing, and similar results, so we
195 combined results from the two runs. Similar runs were done for the two species of
196 *Sicydium*. The genetic differentiation between the two species of *Siniperca* and the
197 two species of *Sicydium* also was estimated using Structure 2.3.4 (Pritchard et al.
198 2000). Three iterations for 100,000 generations (using a 100,000 burnin) were run for
199 each value of K (number of population clusters) ranging from 1 to 3. To identify the
200 number of population clusters that captures the major structure in the data, Structure
201 Harvester (Earl and vonHoldt 2012) was used to calculate the peak value for delta K
202 (Evanno et al. 2005).

203 *Simulating sister species sequence data with different divergence times and gene flow*

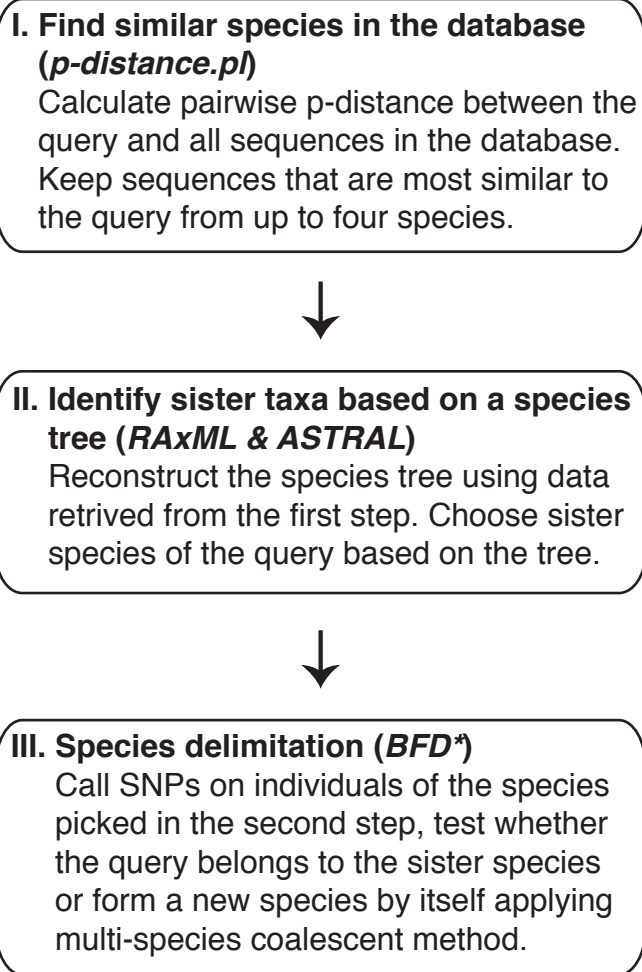
204 We simulated two diverging species with various splitting time and migration
205 rates to explore the effect of changing these two factors on species discrimination
206 over a broader range of parameter space. According to the IMA2 results of the
207 empirical data, the splitting time was set as 1000, 10000, 100000, and 700000
208 generations. The migration rate was set as 0, 0.000001, 0.00001, and 0.0001 per
209 generation. The simulation with 1000 generations splitting time was combined with
210 only 0 migration rate, because the two simulated species were already
211 indistinguishable under 1000 generations splitting time even when there was no gene

212 flow in the simulation. The simulations with 10000, 100000, and 700000 generation
213 splitting time were combined with all four migration rates. Fastsimcoal2 (Excoffier
214 and Foll 2011; Excoffier et al. 2013) was used to generate the simulated data. Twenty
215 thousand replicates were simulated for each scenario. The effective population size
216 used for simulation was 20000 in the ancestor species and the two descendant species.
217 Five sequences were sampled from each simulated species. The simulated data were
218 used to calculate p-distance among individuals of the same and different species.
219 Species identification success rate applying “all species barcodes” criteria was
220 calculated as described above. Identification success rate using different number (1, 3,
221 5, 10, 30, 50, 70, 90, 100, 200, 300, 400, 500, 600, 700, 800, 900, 1000) of simulated
222 loci was plotted against species splitting time and migration rate using R
223 (R_Core_Team 2015).

224 *A three-step multilocus DNA barcoding pipeline*

225 It is straightforward to use distance based methods to reveal divergence of two sister
226 species in the empirical and simulated data. But for more than two species, distance
227 based species identification becomes more complicated. Firstly, an ad hoc barcoding
228 threshold is needed to judge whether the query is one of the species represented in the
229 database or is a new and distinct species, but sometimes no barcoding gap exists for
230 establishing such thresholds. Secondly, the shortest distance does not guarantee a
231 sister species relationship either, because sister species with long branches might be
232 less similar to the query species than a non-sister species with a short branch. To

233 avoid these risks, we propose a three-step DNA barcoding method (Fig. 1).



234

235 **Figure 1.** A three-step multilocus DNA barcoding pipeline.

236 In the first step, p-distances between the query and all sequences in the
237 database are calculated. The sequences that are similar to the query are kept for
238 subsequent analyses (*p-distance.pl*, Supplementary Materials). This is a fast screening
239 process to retrieve all sequences from potential conspecifics or sister species. Because
240 the closest sequence might not be from a conspecifics or sister species, sequences
241 from up to four species are kept. In the second step, a species tree is reconstructed
242 using the sequences from the first step to identify potential conspecifics or sister
243 species of the query using RAxML version 8 (Stamatakis 2014) and ASTRAL 4.10.6

244 (Mirarab et al. 2014; Mirarab and Warnow 2015; Mirarab et al. 2016). Individual gene
245 trees are inferred using RAxML with the GTRGAMMA model, and then a species
246 tree is recovered from those gene trees using ASTRAL. The potential conspecifics or
247 sister species to the query are then chosen based on the phylogenetic relationship
248 depicted in the species tree. In the third step, species delimitation is done using a
249 Bayes factor delimitation approach, BFD* (Leaché et al. 2014). Single nucleotides
250 polymorphism (SNP) data are retrieved from the sequencing reads of the species
251 chosen in step two and used for the BFD* analysis. A path sampling with 48 steps was
252 conducted to estimate the marginal likelihood with a Markov chain Monte Carlo
253 (MCMC) chain length of 200,000 and a pre-burnin of 50,000 following the
254 recommended settings in BFD* (Leaché et al. 2014). The strength of support for
255 compared hypotheses was evaluated from Bayes factor scale, $2\ln(\text{BF})$ using the
256 framework of Kass and Raftery (1995). The BF scale is as follows: $0 < 2\ln(\text{BF}) < 2$ is
257 not worth more than a bare mention, $2 < 2\ln(\text{BF}) > 6$ means positive evidence, $6 <$
258 $2\ln(\text{BF}) < 10$ represents strong support, and $2\ln(\text{BF}) > 10$ represents decisive support.
259 If the result of BFD* analysis does not support two separate species, the query will be
260 assigned to the “sister species”; otherwise, the query will be considered as a new
261 species with its sequences add to the database and further study on its species status
262 will be recommended.

263 The final set of selected markers was used for testing the above-described
264 three-step multilocus DNA barcoding in the siniperoids, including 26 individuals of
265 seven species. An individual of *S. kneri* or *S. chuatsi* was randomly chosen as

266 unknown query that needs to be identified. The sequences of the unknown specimens
267 and all other sequences in the database were aligned using Clustal Omega v1.1.1
268 (Sievers et al. 2011). Custom Perl scripts, `concatnexus.pl` and `gapdis.pl` were used to
269 concatenate the sequences of individual loci, to calculate their p-distance between the
270 query and the sample in the database, and sorted them by the p-distance to find all
271 individuals that are close to the query sample.

272 *Testing effect of missing data in the database or in the query on the success rate of*
273 *species identification*

274 To test if our method could identify new species when the sequences of
275 conspecifics are not in the database, all samples of *S. kneri* were removed from the
276 database except that one random selected *S. kneri* individual was left as query. To
277 access the effect of missing data in the query sample, one *S. kneri* was selected as an
278 unknown sample, and 20 percent, 30 percent, and 50 percent of its loci were excluded,
279 then the data were used for multilocus DNA barcoding analysis.

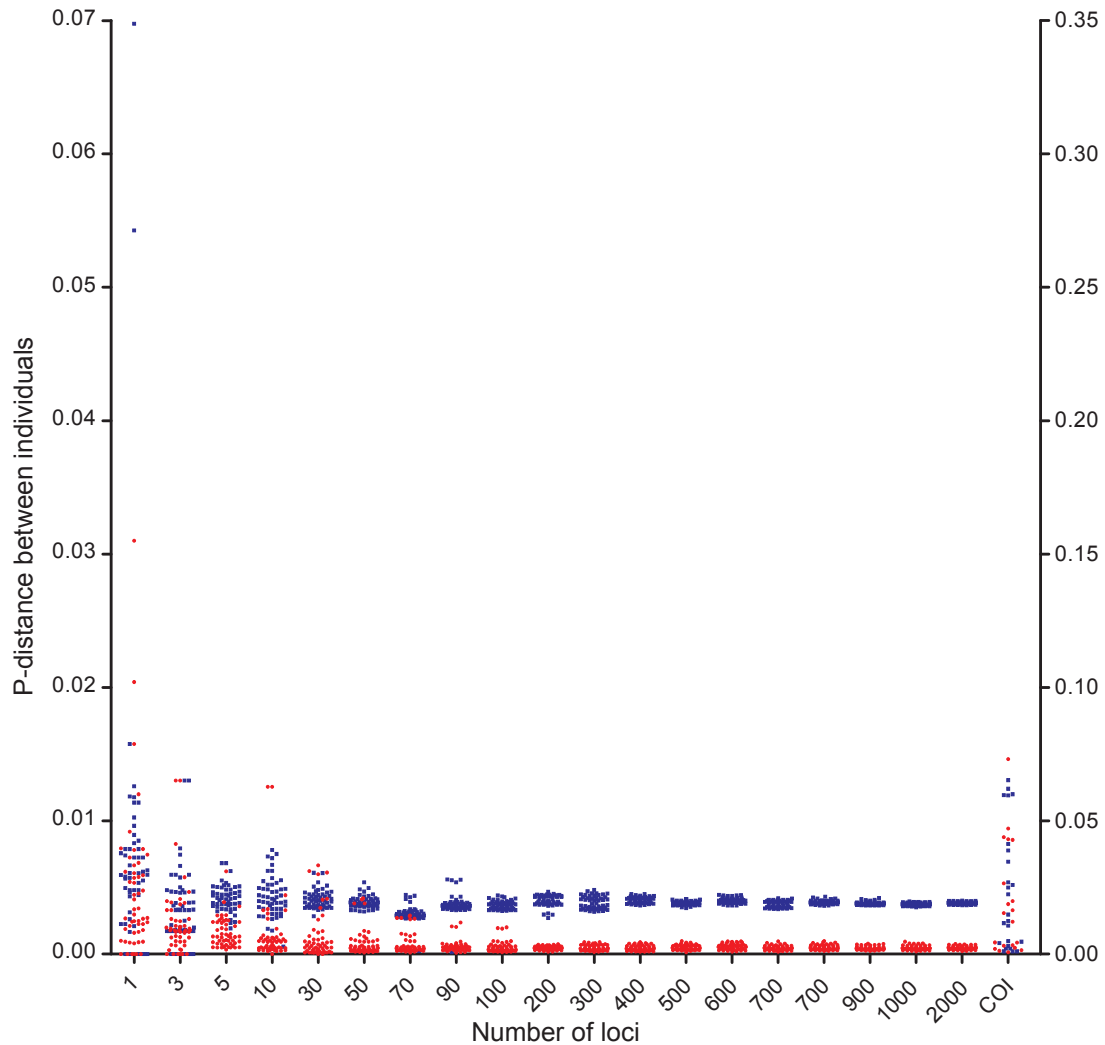
280 RESULTS

281 We investigated effect of increasing number of loci on species discrimination
282 and identification using empirical data (between *Siniperca chuatsi* and *S. kneri* and
283 between *Sicydium altum* and *S. adelum*). We subsequently estimated the population
284 parameters, gene flow and divergence time for both pairs of species. Guided by the
285 patterns seen in the empirical data, we simulated sequences with different splitting
286 times and migration rates, and explored the effect of divergence time and gene flow
287 on the success rate of species identification over a broader range of the relevant

288 parameter space. Finally, we selected 500 nuclear markers for ray-finned fishes and
289 designed a three-step pipeline for multilocus DNA barcoding.

290 *Species discrimination using empirical data*

291 After all loci with missing taxa were excluded, 2586 loci were retained for
292 *Siniperca*. The intra- and interspecific p-distances between five individuals of *S.*
293 *chuatsi* and five *S. kneri* using different numbers of nuclear loci or COI are shown in
294 Figure 2. The intraspecific p-distance (red) calculated using one locus or a small
295 number of loci overlap with interspecific p-distance (blue). There is no barcoding gap
296 separating the intra- and interspecific distances. Intraspecific distances did not
297 become distinguishable from interspecific distances until more than 90 loci were used.
298 The gap separating the intra- and interspecific distance increased as more loci were
299 added, but had little effect after 400 loci were used. The variance of the intra- and
300 interspecific p-distance decreased when more loci were included in calculating the
301 p-distance. The p-distance calculated on COI sequences also had mixed intra- and
302 interspecific values but they were an order of magnitude higher than those calculated
303 using nuclear loci (Fig. 2 right y-axis).

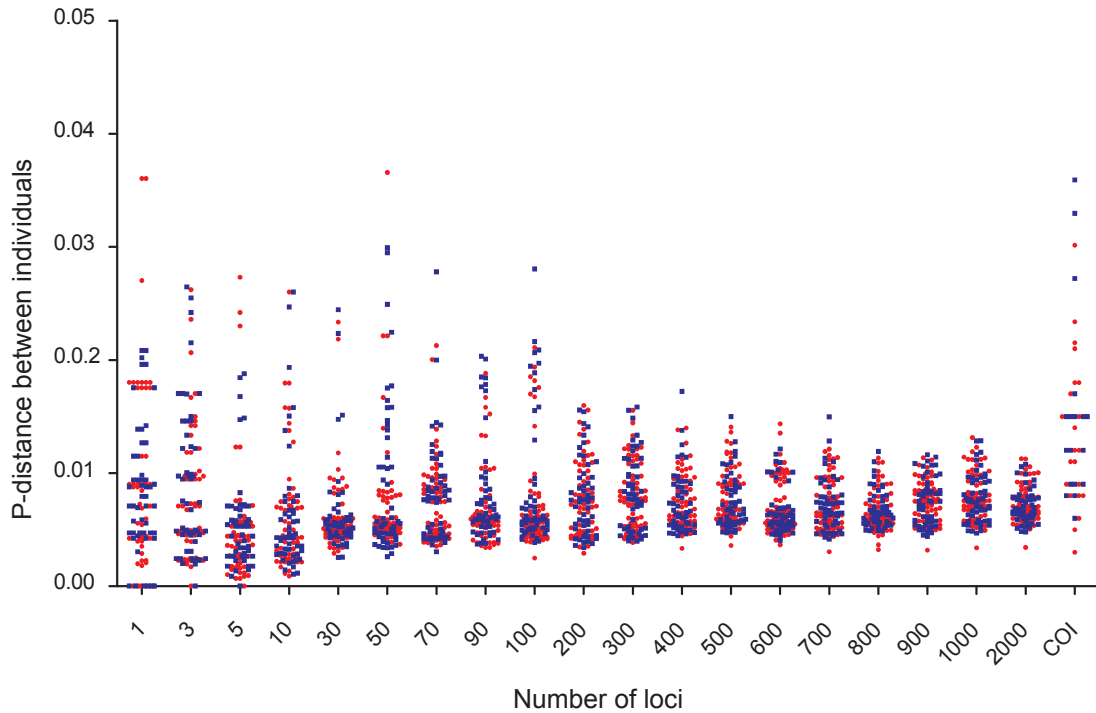


304

305 **Figure 2.** Intra- (red) and interspecific (blue) p-distance of *Siniperca chuatsi* and *S.*
306 *kneri* calculated using different number of nuclear loci or COI gene. The distances
307 based on COI (right y-axis) are one order of magnitude higher than distances
308 calculated using nuclear loci (left y-axis).

309 Similar p-distance calculations on *S. altum* and *S. adelum* resulted in a
310 different pattern than the one observed for *Siniperca*. The intra- (red) and interspecific
311 (blue) p-distances in *Sicydium* were always mixed together, no matter how many loci
312 were included in the analysis. The variance of intra- and interspecific p-distance
313 decreased when more loci were included. The intra- and interspecific p-distances

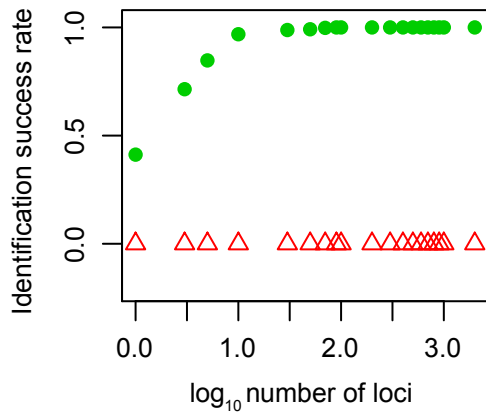
314 calculated using COI also were indistinguishable (Fig. 3).



315

316 **Figure 3.** Intra- (red) and interspecific (blue) p-distance of *Sicydium altum* and *S.*
317 *adelum* calculated using different numbers of nuclear loci or COI gene.

318 The success rate of identification was low (0.412) in *Siniperca* when only one
319 locus was used based on “all species barcodes” criterion with two hundred trials, but
320 it rose up quickly and reached 1.0 after more than 90 loci were added to the dataset
321 (green dots, Fig. 4; Table S1). The identification success rate was zero in *Sicydium*
322 according to the “all species barcodes” criterion, no matter how many loci were
323 included in the analysis (red triangles, Fig. 4). We also applied the COI barcoding
324 approach with an optimized threshold. The success rate of species identification using
325 COI was zero in both *Siniperca* and *Sicydium*.



326

327 **Figure 4.** The relationship between number of loci used and success rate of
328 identification between *Siniperca chuatsi* and *S. kneri* (green dots), and between
329 *Sicydium altum* and *S. adelum* (red triangles).

330 *Population parameters inferred for the two species pairs*

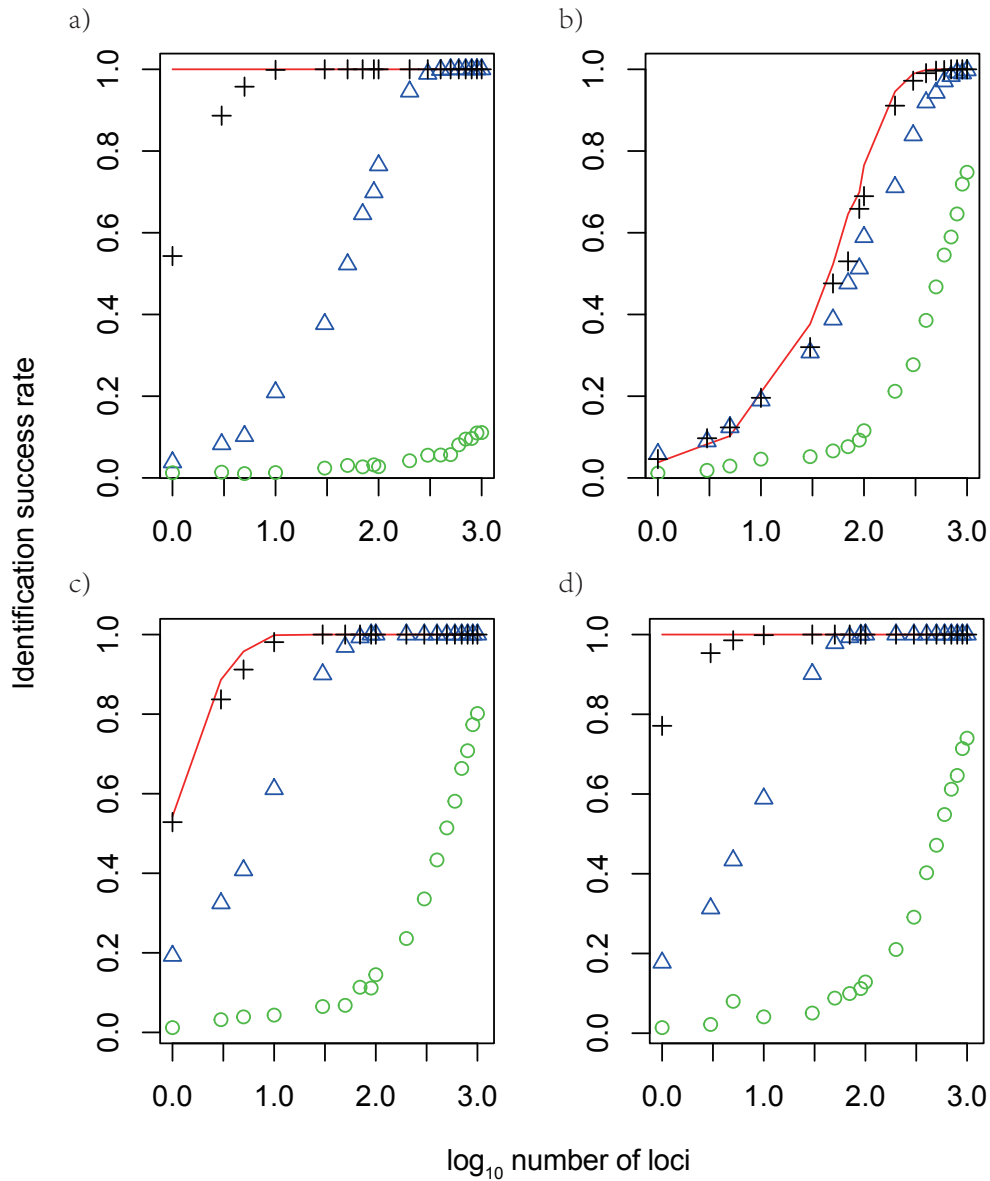
331 To investigate the difference seen in the results of the *Siniperca* and *Sicydium*
332 analyses, we explored some of the population attributes associated with each of these
333 two groups. Structure analysis showed that K equaled to 2 had the highest probability
334 when analyzing the two *Siniperca* species (Fig. S1 a), but the two *Sicydium* species
335 were indistinguishable (Fig. S1 b). The divergence time between *S. chuatsi* and *S.*
336 *kneri* was estimated as $t_0 = 1.754$, which would be equal to ~800,000 generations if
337 we assume an average locus size of 300 bp, a generation time of 2 - 3 years for
338 *Siniperca* and a substitution rate of 2.22×10^{-9} per site per year (Kumar and
339 Subramanian 2002). Gene flow from *S. chuatsi* to *S. kneri* was 0.157 (not significant
340 by LLRtest), but gene flow from *S. kneri* to *S. chuatsi* was highly significant, 0.640 (p
341 < 0.001). The divergence time between *S. altum* and *S. adelum* was estimated as $t_0 =$
342 0.003195, which was not significantly different from zero (HPD95_{L0}=0). Gene flow

343 from *S. altum* to *S. adelum* was 0.494, and gene flow from *S. adelum* to *S. altum* was
344 0.502.

345 *Simulation results*

346 To explore the effect of divergence time and gene flow on the success rate of
347 species identification, we conducted a series of simulations using twenty thousand
348 loci for two species with a range of splitting times and migration profiles. Five
349 sequences from each species were sampled to calculate species identification success
350 rate. Different number of simulated loci were randomly picked and used to identify
351 species. The identification success rate rose with increasing number of loci included
352 in the analyses in all scenarios (Fig. 5). When there was no migration between the two
353 simulated species, the identification success rate increased with splitting time (Fig.
354 5a). The simulation with a splitting time of 1000 generations had the worst
355 identification success rate, only 0.111 even with 1000 loci used (green circle, Fig. 5 a;
356 Table S2). The samples with a splitting time of 10000 generations had low success
357 rates with a small number of loci used, but rose to 1 when more than 400 loci were
358 added to the analyses (blue triangles, Fig. 5a). The samples with a splitting time of
359 100000 generations had a success rate of 1 when more than 10 loci were used (black
360 crosses, Fig. 5 a). Samples with a splitting time of 700000 had success rate of 1 for all

361 analyses (red line, Fig. 5a).



362

363 **Figure 5.** Identification success rate using simulated sequences under different

364 scenarios. a) migration rate equals zero and divergence time equals 700000

365 generations (red line), 100000 generations (black crosses), 10000 generations (blue

366 triangles), and 1000 generations (green circles); b) divergence time equals 10000 and

367 migration rate equals 0 (red line), 0.000001 (black crosses), 0.00001 (blue triangles)

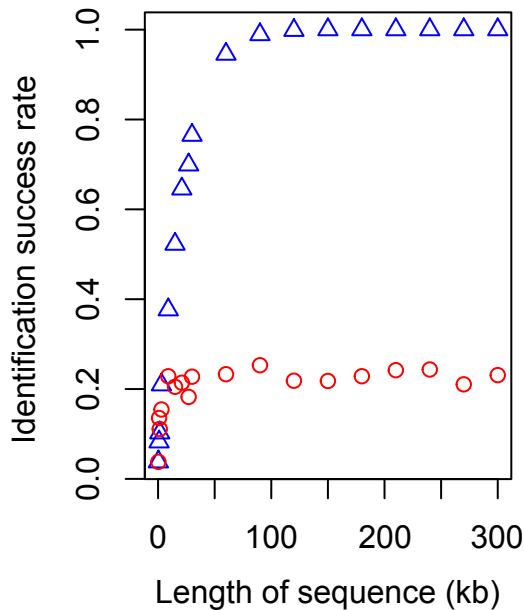
368 and 0.0001 (green circles); c) divergence time equals 100000 and migration rate

369 equals 0 (red line), 0.000001 (black crosses), 0.00001 (blue triangles) and 0.0001
370 (green circles); d) divergence time equals 700000 and migration rate equals 0 (red
371 line), 0.000001 (black crosses), 0.00001 (blue triangles) and 0.0001 (green circles).

372 When gene flow was considered, high gene flow worked in concert with
373 shallow divergence time to reduce the identification success rate (Fig. 5 b-d). A
374 migration rate of 0.0001 (per gene per generation) always led to the worst success rate,
375 and failed to reach a success rate of 1.0 even when all 1000 loci were used in analyses
376 (green circles, Fig. 5 b-d; Table S3-S5). At a migration rate of 0.00001 (blue triangles,
377 Fig. 5 b-d) or a migration rate of 0.000001 (black crosses, Fig. 5 b-d), the
378 identification success rate improved quickly with increasing number of loci (Fig. 5
379 b-d; Table S3-S5). When the divergence time was greater than 100000 generations
380 and gene flow was lower than 0.00001, the identification success rate reached 1.0
381 when more than 90 loci were added to the analysis (Fig. 5 c-d; Table S4-S5).

382 To test whether the length of sequence or the number of loci was the key for
383 success in species identification, we simulated a single locus with increasing size
384 matching the total length of multiple loci. We found that increasing the length of a
385 single locus from 300 bp to 9000 bp improved the success rate slightly, but the
386 success rate did not change when longer sequences were used (Fig. 6 red circles;
387 Table S6). In contrast, concatenating more independent loci with the same total length
388 as the single locus continuously improved the identification success rate, until it

389 reached one (Fig. 6 blue triangles; Table S6).



390

391 **Figure 6.** Comparison between success rates of species identification based on a) a
392 single locus and b) multiple loci. The length of the single locus equals the total length
393 of multiple loci (300 bp each).

394

Multilocus DNA barcoding using empirical data

395

396

397

398

399

400

401

402

403

Based on the results from p-distance and species identification analyses of simulated and empirical data, we decided to pick 500 loci for multilocus DNA barcoding. First, we filtered the 4,434 markers developed for all ray-finned fishes and kept 750 loci with the lowest number of missing taxa. Next, we sorted the 750 loci by their average p-distance and picked from them 500 independent loci with large p-distances. This design was implemented both to minimize missing data when applying to ray-finned fishes and to ensure that loci would be variable enough for multilocus DNA barcoding. Information describing the 500 loci is listed in Supplementary Materials (Table S7).

404 Three individuals, 839_3 (*S. kneri*), 839_6 (*S. kneri*), and 938_1 (*S. chuatsi*)
405 were randomly selected. Each of the randomly picked individuals was used to
406 simulate “a putatively unknown” query for identification. Firstly, the p-distance
407 between the unknown query and the other sinipercids in the database was calculated
408 (Table S8). Secondly, based on the sorted list of p-distances, we selected five closely
409 related taxa, including the query. For example, for 839_3, we used sequence data of
410 839_3, *S. kneri*, *S. chuatsi*, *S. undulata* and *S. obscura* to reconstruct a species tree, in
411 which 839_3 was found to be sister to *S. kneri* (Fig. S2). We then ran a BFD* test to
412 delimitate the unknown query (839_3) with *S. kneri* using *S. chuatsi* as outgroup. The
413 BFD* analyses correctly grouped 839_3 (*S. kneri*) with *S. kneri* (Table 1). The two
414 other randomly picked samples, 839_6 (*S. kneri*), and 938_1 (*S. chuatsi*) were also
415 correctly identified (Table S8 and S9).

416 DNA barcoding using only COI data was unsuccessful. In many cases, the
417 closest taxa of the unknown samples were not their conspecifics either in the tree or
418 measured by p-distances (Fig. S3, Table S10).

419 *Effect of missing data on multilocus DNA barcoding*

420 When all conspecifics were excluded from the database, the unknown query,
421 893_3 (*S. kneri*) was found to be closely related to its sister species *S. chuatsi* (Fig.
422 S4). The p-distances also indicated that the unknown was related to *S. chuatsi* (Table
423 S11). A species delimitation analysis was run with the BFD* method to test whether
424 the unknown should be assigned to *S. chuatsi* or not. The result of BFD* strongly
425 support the unknown query as a separate species ($2\ln BF = 2255.7$; Table 1). In other

426 tests, we keep the database intact, but excluded 20%, 30% and 50% of the loci from
427 the unknown query (893_3 *S. kneri*). We still identified the unknown correctly using
428 the multilocus DNA barcoding approach (Table 1).

429 TABLE 1. Results for species delimitation on unknown sample 893_3 (*S. kneri*) using BFD*
430 based on all 500 nuclear loci, missing 20%, 30% and 50% of the 500 loci or missing
431 conspecific of *S. kneri* in the database.

Data treatment	Model	Marginal likelihood	2lnBF
Using all data	Lumping 839_3 and <i>S. kneri</i>	-1575.80	20.62
	Splitting 839_3 and <i>S. kneri</i>	-1586.11	
Excluding conspecifics of 839_3	Lumping 839_3 and <i>S. chuatsi</i>	-2350.77	
	Splitting 839_3 and <i>S. chuatsi</i>	-1222.90	2255.7
Excluding 20% loci of the 839_3	Lumping 839_3 and <i>S. kneri</i>	-1467.41	26.54
	Splitting 839_3 and <i>S. kneri</i>	-1480.68	
Excluding 30% loci of the 839_3	Lumping 839_3 and <i>S. kneri</i>	-1247.44	22.12
	Splitting 839_3 and <i>S. kneri</i>	-1258.50	
Excluding 50% loci of the 839_3	Lumping 839_3 and <i>S. kneri</i>	-914.60	22.40
	Splitting 839_3 and <i>S. kneri</i>	-925.80	

432

DISCUSSION

433 *Species are more distinguishable with more independent loci than with a few*

434 Our results demonstrated that the difference between species become more
435 distinct when more independent loci are used. The intra- (red) and interspecific (blue)
436 p-distance between individuals of *S. chuatsi* and *S. kneri* were largely overlapping
437 when only COI gene or a few randomly picked nuclear gene were used to calculate
438 the p-distance (Fig. 2). When more loci were added to the analyses, the intra- and
439 interspecific distance became better separated. At 90 loci, a “barcoding gap” between
440 the intra- and interspecific distance emerged. The variance of the intra- and
441 interspecific distances also decreased as the number of loci used in the analyses
442 increased. Based on these findings we conclude that the lack of an apparent barcoding
443 gap between *S. chuatsi* and *S. kneri* using COI or a few nuclear genes is due to
444 sampling error. Using more independent loci would likely improve the estimates of
445 population parameters (Lee and Edwards 2008). Similarly, more independent loci
446 should improve precision of both the estimated intra- and interspecific genetic
447 distance, resulting in increased discriminatory power (Fig. 4). The same patterns were
448 observed in all of our simulated analyses, namely that the species identification
449 success rate rose with increasing number of loci (Fig. 5). Interestingly, using longer
450 genes instead of more genes did not improve species identification (Fig. 6).

451

Divergence time vs. gene flow

452 Gene flow between sister species can cause problems that are similar to those

453 caused by a lack of divergence. *S. chuatsi* and *S. kneri* were estimated split at around
454 80 thousand generations ago, with uni-directional introgression flowing from *S. kneri*
455 to *S. chuatsi*, $m_{1>0} = 0.640$. Therefore, the lack of reciprocal monophyly or barcoding
456 gap between *S. chuatsi* and *S. kneri* using COI or a few nuclear loci could, in fact, be
457 caused by gene flow between the two species rather than the short divergence time
458 originally hypothesized.

459 *Sicydium altum* and *S. adelum* were estimated to have split very recently, $t_0 =$
460 0.003195. Bi-directional gene flow was estimated as 0.494 from *S. altum* to *S. adelum*,
461 and 0.502 from *S. adelum* to *S. altum*. All of our analyses could not differentiate
462 between *S. altum* and *S. adelum* genetically. Structure analysis (Fig. S1), species
463 identification and p-distance assessments (Fig. 3 and 4) all indicated that *S. altum* and
464 *S. adelum* are indistinguishable. Accordingly, we suggest that the taxonomic status of
465 *S. altum* and *S. adelum* be revisited by a more detailed morphological analysis.

466 It is difficult to tell whether gene flow or short species divergence time played
467 a more prominent role in obstructing DNA barcoding. It has been reported that a
468 considerable proportion of animal species do not form monophyletic groups (Funk
469 and Omland 2003; Ross 2014), but the causes for such patterns have not yet been
470 fully explored. From the results of our empirical and simulated analyses, we conclude
471 that when the splitting time between sister species was more than 100000 generations
472 old and the migration rate was lower than 0.00001, using multilocus DNA barcoding
473 (with more than 90 loci) we could correctly determine the species status of unknown
474 samples, whereas single-locus DNA barcoding suffered from lacking of power in

475 species discrimination.

476 *Markers for multilocus DNA barcoding*

477 A suite of universal gene markers that could be used on a whole group of
478 organisms is a prerequisite for multilocus DNA barcoding. Because of improvements
479 in sequencing technology and the increasing number of publicly accessible genome
480 data bases, more and more genome-scale markers have been developed for different
481 group of organisms, such as turtles (Shen et al. 2011), birds (McCormack et al. 2013),
482 tapeworms (Yuan et al. 2016), flower flies (Young et al. 2016), plants (Schmickl et al.
483 2016), echinoderms (Hugall et al. 2016), insects (Blaimer et al. 2016) and vertebrates
484 (Li et al. 2013). Some of these markers can be applied across broad groups of
485 organisms, whereas other have only been tested for restricted groups. We predict that
486 obtaining suitable sets of markers for multilocus DNA barcoding will not be a
487 limitation, but a lot of testing will need to be carried out across a broad range of taxa
488 before an agreed set of common markers can be established for each major group of
489 organisms.

490 Our pick of 500 markers for ray-finned fishes has been tested in all major
491 lineages of fishes. We chose markers that were found to be present in most groups of
492 fishes and that were variable across groups. We recommend using them as standard
493 multilocus DNA barcode markers for all ray-finned fishes. Our results indicate that
494 more than 90 loci should be enough for species identification, but we advocate using
495 the complete set of 500 loci, as there is almost no extra cost in capturing 500 rather
496 than 90 loci. Additionally, targeting more loci provides insurance against missing

497 data. We found that missing 20%, 30% and up to 50% loci in the unknown sample
498 had no effect in identification success.

499 Other alternatives to collecting large datasets for DNA barcoding include
500 genome skimming (Coissac et al. 2016) and whole-chloroplast genome sequencing
501 (Li et al. 2015b). Genome skimming employs low-coverage shotgun sequencing of
502 genomic DNA, which circumvents the need for PCR, avoiding the needs for universal
503 primers. Because genome skimming is unselective, it involves collecting a lot of data
504 that ultimately is not used, but requires data storage and analysis resources.

505 Low-coverage shotgun sequencing also yields a high proportion of missing data.

506 Sequencing genomes of chloroplasts or other organelles is focused on a single long
507 sequence, which tends to yield low success rate of species identification, as shown in
508 our simulation.

509 *Three-step data analysis pipeline of multilocus barcoding*

510 Dowton et al. (2014) proposed an pipeline integrating species tree
511 reconstruction and species delimitation. They used Beast* to build a species tree
512 (Heled and Drummond 2010), and took the species tree as guide tree for delimitating
513 species using BPP (Rannala and Yang 2003; Yang and Rannala 2010). Our method is
514 similar to the method of Dowton et al. (2014). We first screened the reference
515 database for individuals from closely related species based on p-distance between the
516 unknown query and sequences in the database. We only choose four closely related
517 species as potential conspecifics or sister species. We think the number of species
518 selected is enough for the current study, because our p-distance calculation was based

519 on many independent loci, which reduces random error. The small number of selected
520 species could also help to relieve the computational burden associated with
521 reconstructing the species tree in the second step. Using a combination of RAxML
522 and ASTRAL program, we could reconstruct a species tree of five taxa, four selected
523 species plus the query in minutes using 500 loci. In the last step, we included only
524 three taxa, one conspecific or sister species, one outgroup species and the unknown
525 query for species delimitation using BFD*, which also saved computation time. We
526 anticipate that the computational burden associated with multilocus DNA barcoding
527 will be further reduced as new algorithms are developed, to make multilocus
528 barcoding a real-time tool.

529 *Cost of multilocus DNA barcoding*

530 Finally, from a practical standpoint, multilocus barcoding through target gene
531 enrichment is efficient. We estimate around \$90 for the total cost of capturing and
532 sequencing 500 loci per sample, which is less than the cost of amplifying and
533 sequencing 10 loci using the traditional methods of PCR and Sanger sequencing. The
534 cost of target gene capture comprises: library prep, \$50; RNA baits, \$32; and
535 sequencing, \$8 per sample. The major costs are associated with the purchase of
536 commercial RNA bait kits and the library preparation step, which can be lowered by
537 purchasing kits in bulk and by using robots to automate library preparation.

538 SUPPLEMENTARY MATERIAL

539 Data available from the Dryad Digital Repository: <http://dx.doi.org/XXX>.

540

FUNDING

541

This work was supported by the Shanghai Pujiang Program, the Program for

542

Professor of Special Appointment (Eastern Scholar) at Shanghai Institutions of Higher

543

Learning to C. Li.

544

ACKNOWLEDGEMENTS

545

The authors would like to thank Shanghai Oceanus Supercomputing Center

546

(SOSC) for providing computational resources.

547

REFERENCES

548

Aliabadian M, Beentjes KK, Roselaar CS, van Brandwijk H, Nijman V, Vonk R. 2013. DNA barcoding of Dutch birds. *Zookeys*:25-48.

549

550

Bi K, Vanderpool D, Singhal S, Linderoth T, Moritz C, Good JM. 2012. Transcriptome-based exon

551

capture enables highly cost-effective comparative genomic data collection at moderate

552

evolutionary scales. *BMC Genomics*, 13:403.

553

Blaimer BB, Lloyd MW, Guillory WX, Brady SG. 2016. Sequence Capture and Phylogenetic Utility of

554

Genomic Ultraconserved Elements Obtained from Pinned Insect Specimens. *PLoS One*,

555

11:e0161531.

556

Brown SD, Collins RA, Boyer S, Lefort MC, Malumbres-Olarte J, Vink CJ, Cruickshank RH. 2012. Spider:

557

an R package for the analysis of species identity and evolution, with particular reference to DNA

558

barcoding. *Mol Ecol Resour*, 12:562-565.

559

Bussing WA. 1996. *Sicydium adelum*, a new species of gobiid fish (Pisces: Gobiidae) from Atlantic

560

slope streams of Costa Rica. *Rev. Biol. Trop.*, 44:819-825.

561

Candek K, Kuntner M. 2015. DNA barcoding gap: reliable species identification over morphological

562

and geographical scales. *Mol Ecol Resour*, 15:268-277.

563

Chabarría RE. 2015. Evolution of the genus *Sicydium* (Gobiidae: Sicydiinae). Department of Life

564

Sciences. Corpus Christi, TX, Texas A&M University – Corpus Christi.

565

Chan A, Chiang LP, Hapuarachchi HC, Tan CH, Pang SC, Lee R, Lee KS, Ng LC, Lam-Phua SG. 2014. DNA

566

barcoding: complementing morphological identification of mosquito species in Singapore.

567

Parasite Vector, 7:569.

568

Coissac E, Hollingsworth PM, Lavergne S, Taberlet P. 2016. From barcodes to genomes: extending the

569

concept of DNA barcoding. *Molecular ecology*, 25:1423-1428.

570

Collins RA, Armstrong KF, Meier R, Yi Y, Brown SD, Cruickshank RH, Keeling S, Johnston C. 2012.

571

Barcoding and border biosecurity: identifying cyprinid fishes in the aquarium trade. *PLoS One*,

572

7:e28381.

573

Collins RA, Cruickshank RH. 2014. Known knowns, known unknowns, unknown unknowns and

574

unknown knowns in DNA barcoding: a comment on Dowton et al. *Syst. Biol.*, 63:1005-1009.

- 575 Decru E, Moelants T, De Gelas K, Vreven E, Verheyen E, Snoeks J. 2016. Taxonomic challenges in
576 freshwater fishes: a mismatch between morphology and DNA barcoding in fish of the
577 north-eastern part of the Congo basin. *Mol Ecol Resour*, 16:342-352.
- 578 Downton M, Meiklejohn K, Cameron SL, Wallman J. 2014. A preliminary framework for DNA barcoding,
579 incorporating the multispecies coalescent. *Syst. Biol.*, 63:639-644.
- 580 Earl DA, vonHoldt BM. 2012. STRUCTURE HARVESTER: a website and program for visualizing
581 STRUCTURE output and implementing the Evanno method. *Conservation Genetics Resources*,
582 4:359-361.
- 583 Evanno G, Regnaut S, Goudet J. 2005. Detecting the number of clusters of individuals using the
584 software STRUCTURE: a simulation study. *Molecular ecology*, 14:2611-2620.
- 585 Excoffier L, Dupanloup I, Huerta-Sanchez E, Sousa VC, Foll M. 2013. Robust demographic inference
586 from genomic and SNP data. *PLoS Genet*, 9:e1003905.
- 587 Excoffier L, Foll M. 2011. fastsimcoal: a continuous-time coalescent simulator of genomic diversity
588 under arbitrarily complex evolutionary scenarios. *Bioinformatics*, 27:1332-1334.
- 589 Funk DJ, Omland KE. 2003. SPECIES-LEVEL PARAPHYLY AND POLYPHYLY: Frequency, Causes, and
590 Consequences, with Insights from Animal Mitochondrial DNA. *Annu. Rev. Ecol. Evol. Syst.*,
591 34:397-423.
- 592 Ghahramanzadeh R, Esselink G, Kodde LP, Duistermaat H, van Valkenburg JL, Marashi SH, Smulders
593 MJ, van de Wiel CC. 2013. Efficient distinction of invasive aquatic plant species from non-invasive
594 related species using DNA barcoding. *Mol Ecol Resour*, 13:21-31.
- 595 Hartvig I, Czako M, Kjaer ED, Nielsen LR, Theilade I. 2015. The Use of DNA Barcoding in Identification
596 and Conservation of Rosewood (*Dalbergia* spp.). *PLoS One*, 10:e0138231.
- 597 Hassold S, Lowry PP, 2nd, Bauert MR, Razafintsalama A, Ramamonjisoa L, Widmer A. 2016. DNA
598 Barcoding of Malagasy Rosewoods: Towards a Molecular Identification of CITES-Listed *Dalbergia*
599 Species. *PLoS One*, 11:e0157881.
- 600 Hawlitschek O, Moriniere J, Dunz A, Franzen M, Rodder D, Glaw F, Haszprunar G. 2016.
601 Comprehensive DNA barcoding of the herpetofauna of Germany. *Mol Ecol Resour*, 16:242-253.
- 602 Hebert PD, Cywinska A, Ball SL, deWaard JR. 2003. Biological identifications through DNA barcodes.
603 *Proc Biol Sci*, 270:313-321.
- 604 Hedtke SM, Morgan MJ, Cannatella DC, Hillis DM. 2013. Targeted enrichment: maximizing
605 orthologous gene comparisons across deep evolutionary time. *PLoS One*, 8:e67908.
- 606 Heled J, Drummond AJ. 2010. Bayesian inference of species trees from multilocus data. *Mol Biol Evol*,
607 27:570-580.
- 608 Hey J. 2010. Isolation with migration models for more than two populations. *Mol Biol Evol*,
609 27:905-920.
- 610 Hugall AF, O'Hara TD, Hunjan S, Nilsen R, Moussalli A. 2016. An Exon-Capture System for the Entire
611 Class Ophiuroidea. *Mol Biol Evol*, 33:281-294.
- 612 Kadarusman, Hubert N, Hadiaty RK, Sudarto, Paradis E, Pouyaud L. 2012. Cryptic diversity in
613 Indo-Australian rainbowfishes revealed by DNA barcoding: implications for conservation in a
614 biodiversity hotspot candidate. *PLoS One*, 7:e40627.
- 615 Kass RE, Raftery AE. 1995. Bayes factors. *Journal of the American Statistical Association* 90:773-795.
- 616 Kumar S, Subramanian S. 2002. Mutation rates in mammalian genomes. *Proceedings of the National
617 Academy of Sciences of the United States of America*, 99:803-808.
- 618 Leaché AD, Fujita MK, Minin VN, Bouckaert RR. 2014. Species delimitation using genome-wide SNP

- 619 data. *Syst Biol*, 63:534-542.
- 620 Lee JY, Edwards SV. 2008. Divergence across Australia's Carpentarian barrier: statistical
621 phylogeography of the red-backed fairy wren (*Malurus melanocephalus*). *Evolution*,
622 62:3117-3134.
- 623 Li C, Hofreiter M, Straube N, Corrigan S, Naylor GJ. 2013. Capturing protein-coding genes across highly
624 divergent species. *Biotechniques*, 54:321-326.
- 625 Li C, Riethoven JJ, Naylor GJ. 2012. *EvolMarkers*: a database for mining exon and intron markers for
626 evolution, ecology and conservation studies. *Mol Ecol Resour*, 12:967-971.
- 627 Li J, Zheng X, Cai Y, Zhang X, Yang M, Yue B. 2015a. DNA barcoding of Murinae (Rodentia: Muridae)
628 and Arvicolinae (Rodentia: Cricetidae) distributed in China. *Mol Ecol Resour*, 15:153-167.
- 629 Li S. 1991. Geographical distribution of the Sinipercine fishes. *Chinese Journal of Zoology*, 26:40-44.
- 630 Li X, Yang Y, Henry RJ, Rossetto M, Wang Y, Chen S. 2015b. Plant DNA barcoding: from gene to
631 genome. *Biol Rev Camb Philos Soc*, 90:157-166.
- 632 Liu H, Chen Y. 1994. Phylogeny of the sinipercine fishes with some taxonomic notes. *Zoological*
633 *Research*, 15:1-12.
- 634 Mabragana E, Diaz de Astarloa JM, Hanner R, Zhang J, Gonzalez Castro M. 2011. DNA barcoding
635 identifies Argentine fishes from marine and brackish waters. *PLoS One*, 6:e28655.
- 636 Marescaux J, Van Doninck K. 2013. Using DNA barcoding to differentiate invasive *Dreissena* species
637 (Mollusca, Bivalvia). *Zookeys*:235-244.
- 638 McCormack JE, Harvey MG, Faircloth BC, Crawford NG, Glenn TC, Brumfield RT. 2013. A phylogeny of
639 birds based on over 1,500 loci collected by target enrichment and high-throughput sequencing.
640 *PLoS One*, 8:e54848.
- 641 Meier R, Shiyang K, Vaidya G, Ng PK. 2006. DNA barcoding and taxonomy in Diptera: a tale of high
642 intraspecific variability and low identification success. *Syst. Biol.*, 55:715-728.
- 643 Mirarab S, Bayzid MS, Warnow T. 2016. Evaluating Summary Methods for Multilocus Species Tree
644 Estimation in the Presence of Incomplete Lineage Sorting. *Syst. Biol.*, 65:366-380.
- 645 Mirarab S, Reaz R, Bayzid MS, Zimmermann T, Swenson MS, Warnow T. 2014. *ASTRAL*: genome-scale
646 coalescent-based species tree estimation. *Bioinformatics*, 30:i541-548.
- 647 Mirarab S, Warnow T. 2015. *ASTRAL-II*: coalescent-based species tree estimation with many hundreds
648 of taxa and thousands of genes. *Bioinformatics*, 31:i44-52.
- 649 Nelson JS. 2006. *Fishes of the world*. 4th ed. New York, John Wiley and Sons, Inc.
- 650 Nevill PG, Wallace MJ, Miller JT, Krauss SL. 2013. DNA barcoding for conservation, seed banking and
651 ecological restoration of *Acacia* in the Midwest of Western Australia. *Mol Ecol Resour*,
652 13:1033-1042.
- 653 Pritchard JK, Stephens M, Donnelly P. 2000. Inference of population structure using multilocus
654 genotype data. *Genetics*, 155:945-959.
- 655 Puillandre N, Lambert A, Brouillet S, Achaz G. 2012. *ABGD*, Automatic Barcode Gap Discovery for
656 primary species delimitation. *Molecular ecology*, 21:1864-1877.
- 657 R_Core_Team. 2015. *R: A language and environment for statistical computing*. Vienna, Austria., R
658 Foundation for Statistical Computing.
- 659 Rannala B, Yang Z. 2003. Bayes estimation of species divergence times and ancestral population sizes
660 using DNA sequences from multiple loci. *Genetics*, 164:1645-1656.
- 661 Ross HA. 2014. The incidence of species-level paraphyly in animals: a re-assessment. *Mol. Phylogenet.*
662 *Evol.*, 76:10-17.

- 663 Saunders GW. 2009. Routine DNA barcoding of Canadian Gracilariales (Rhodophyta) reveals the
664 invasive species *Gracilaria vermiculophylla* in British Columbia. *Mol Ecol Resour*, 9 Suppl
665 s1:140-150.
- 666 Schmickl R, Liston A, Zeisek V, Oberlander K, Weitemier K, Straub SC, Cronn RC, Dreyer LL, Suda J.
667 2016. Phylogenetic marker development for target enrichment from transcriptome and genome
668 skim data: the pipeline and its application in southern African *Oxalis* (Oxalidaceae). *Mol Ecol*
669 *Resour*, 16:1124-1135.
- 670 Shapcott A, Forster PI, Guymer GP, McDonald WJ, Faith DP, Erickson D, Kress WJ. 2015. Mapping
671 biodiversity and setting conservation priorities for SE Queensland's rainforests using DNA
672 barcoding. *PLoS One*, 10:e0122164.
- 673 Shen XX, Liang D, Wen JZ, Zhang P. 2011. Multiple genome alignments facilitate development of NPCL
674 markers: a case study of tetrapod phylogeny focusing on the position of turtles. *Mol Biol Evol*,
675 28:3237-3252.
- 676 Sievers F, Wilm A, Dineen D, Gibson TJ, Karplus K, Li W, Lopez R, McWilliam H, Remmert M, Soding J,
677 *et al.* 2011. Fast, scalable generation of high-quality protein multiple sequence alignments using
678 Clustal Omega. *Mol Syst Biol*, 7:539.
- 679 Smith MA, Fisher BL, Hebert PD. 2005. DNA barcoding for effective biodiversity assessment of a
680 hyperdiverse arthropod group: the ants of Madagascar. *Philos Trans R Soc Lond B Biol Sci*,
681 360:1825-1834.
- 682 Sonet G, Jordaens K, Nagy ZT, Breman FC, De Meyer M, Backeljau T, Virgilio M. 2013. Adhoc: an R
683 package to calculate ad hoc distance thresholds for DNA barcoding identification.
684 *Zookeys*:329-336.
- 685 Spasojevic T, Kropf C, Nentwig W, Lasut L. 2016. Combining morphology, DNA sequences, and
686 morphometrics: revising closely related species in the orb-weaving spider genus *Araniella*
687 (Araneae, Araneidae). *Zootaxa*, 4111:448-470.
- 688 Stamatakis A. 2014. RAxML version 8: a tool for phylogenetic analysis and post-analysis of large
689 phylogenies. *Bioinformatics*, 30:1312-1313.
- 690 Sutou M, Kato T, Ito M. 2011. Recent discoveries of armyworms in Japan and their species
691 identification using DNA barcoding. *Mol Ecol Resour*, 11:992-1001.
- 692 Tanzler R, Sagata K, Surbakti S, Balke M, Riedel A. 2012. DNA barcoding for community ecology--how
693 to tackle a hyperdiverse, mostly undescribed Melanesian fauna. *PLoS One*, 7:e28832.
- 694 van Velzen R, Weitschek E, Felici G, Bakker FT. 2012. DNA barcoding of recently diverged species:
695 relative performance of matching methods. *PLoS One*, 7:e30490.
- 696 Vences M, Thomas M, Bonett RM, Vieites DR. 2005. Deciphering amphibian diversity through DNA
697 barcoding: chances and challenges. *Philos Trans R Soc Lond B Biol Sci*, 360:1859-1868.
- 698 Virgilio M, Jordaens K, Breman FC, Backeljau T, De Meyer M. 2012. Identifying insects with incomplete
699 DNA barcode libraries, African fruit flies (Diptera: Tephritidae) as a test case. *PLoS One*, 7:e31581.
- 700 Ward RD, Zemlak TS, Innes BH, Last PR, Hebert PD. 2005. DNA barcoding Australia's fish species.
701 *Philos Trans R Soc Lond B Biol Sci*, 360:1847-1857.
- 702 Witt JD, Threlloff DL, Hebert PD. 2006. DNA barcoding reveals extraordinary cryptic diversity in an
703 amphipod genus: implications for desert spring conservation. *Molecular ecology*, 15:3073-3082.
- 704 Yang Z, Rannala B. 2010. Bayesian species delimitation using multilocus sequence data. *Proceedings*
705 *of the National Academy of Sciences of the United States of America*, 107:9264-9269.
- 706 Young AD, Lemmon AR, Skevington JH, Mengual X, Stahls G, Reemer M, Jordaens K, Kelso S, Lemmon

707 EM, Hauser M, *et al.* 2016. Anchored enrichment dataset for true flies (order Diptera) reveals
708 insights into the phylogeny of flower flies (family Syrphidae). *BMC Evol Biol*, 16:143.
709 Yuan H, Jiang J, Jimenez FA, Hoberg EP, Cook JA, Galbreath KE, Li C. 2016. Target gene enrichment in
710 the cyclophyllidean cestodes, the most diverse group of tapeworms. *Mol Ecol Resour.*
711 Zhao J, Li C, Zhao L, Wang W, Cao Y. 2008. Mitochondrial diversity and phylogeography of the Chinese
712 perch, *Siniperca chuatsi* (Perciformes: Sinipercidae). *Mol. Phylogenet. Evol.*, 49:399–404.
713 Zhou C, Yang Q, Cai D. 1988. On the classification and distribution of the sinipercinae fishes (family
714 Serranidae). *Zoological Research*, 9:113-125.
715
716
717