

1 **Forensics and DNA Barcodes – Do Identification Errors Arise in the Lab**
2 **or in the Sequence Libraries?**

3

4

5 Mikko Pentinsaari¹

6 Sujeevan Ratnasingham¹

7 Scott E. Miller²

8 *Paul D. N. Hebert¹

9

10 ¹Centre for Biodiversity Genomics

11 University of Guelph

12 Guelph, ON N1G 2W1

13 Canada

14

15 ²National Museum of Natural History

16 Washington, DC

17 USA

18

19 *Corresponding author phebert@uoguelph.ca

20

21

22

23

24

25

26

27

28

29

30

31

32

33

34

35

36 **Abstract**

37 Forensic studies often require the determination of biological materials to a species level. As such, DNA-
38 based approaches to identification, particularly DNA barcoding, are attracting increased interest. The
39 capacity of DNA barcodes to assign newly encountered specimens to a species relies upon access to
40 informatics platforms, such as BOLD and GenBank, which host libraries of reference sequences and
41 support the comparison of new sequences to them. As parameterization of these libraries expands, DNA
42 barcoding has the potential to make valuable contributions in diverse forensic contexts. However, a recent
43 publication called for caution after finding that both platforms performed poorly in identifying specimens
44 of 17 common insect species. This study follows up on this concern by asking if the misidentifications
45 reflected problems in the reference libraries or in the query sequences used to test them. Because this
46 reanalysis revealed that missteps in acquiring and analyzing the query sequences were responsible for the
47 misidentifications, a workflow is described to minimize such errors in future investigations. The present
48 study also revealed the limitations imposed by the lack of a polished species-level taxonomy for many
49 groups. In such cases, forensic applications can be strengthened by mapping the geographic distributions
50 of sequence-based species proxies rather than waiting for the maturation of formal taxonomic systems
51 based on morphology.

52

53

54 **Introduction**

55 Species identifications play an important role in forensic analyses in contexts ranging from the
56 interception of trade in CITES-listed species [1] to ascertaining the post mortem interval [2]. There are
57 also expanding opportunities to track the movement of objects and organisms linked to their associated
58 DNA. Although species identifications can play an important role in these contexts, the lack of taxonomic
59 specialists often impedes analysis, a factor which has provoked interest in DNA-based approaches to

60 species identification. Past studies have established that DNA barcodes can often assign specimens to
61 their source species, but have also revealed differences in success among the kingdoms of eukaryotes.
62 For example, the three barcode regions (rbcl, matK, ITS2) for plants deliver lower success than the single
63 gene region (cytochrome *c* oxidase I, COI) used for animals [3]. Because COI generally has high accuracy
64 in species assignment [4–9], the conclusions from a recent study by Meiklejohn et al. [10] were surprising.
65 They assessed the capacity of reference sequences in BOLD, the Barcode of Life Data System [11], and
66 GenBank [12] to generate species-level identifications. Their analysis revealed that both platforms
67 performed similarly in identifying plants and macrofungi, but fared poorly in identifying insect species
68 with BOLD showing lower success than GenBank (35% vs. 53%). By evaluating the factors underpinning
69 the incorrect assignments, the present study revealed that errors in sequence acquisition and
70 interpretation accounted for most, if not all, of the misidentifications. To avoid similar issues in future
71 studies, there is a need to adopt more rigorous procedures for data acquisition and analysis, and to reduce
72 the current reliance on immature taxonomic systems.

73

74 **Material and Methods**

75 Meiklejohn et al. [10] analyzed 17 insects including representatives from 12 insect orders – Coleoptera
76 (1), Dermaptera (1), Diptera (5), Ephemeroptera (1), Hymenoptera (1), Lepidoptera (2), Mecoptera (1),
77 Neuroptera (1), Odonata (1), Orthoptera (1), Pthiraptera (1), and Siphonaptera (1). The specimens were
78 obtained from the Smithsonian’s National Museum of Natural History; most were collected 20+ years ago
79 (e.g. *Pediculus humanus* – 1955). Following DNA extraction, the barcode region of COI was PCR amplified
80 and then Sanger sequenced. Reflecting the DNA degradation typical of museum specimens, the sequences
81 recovered were often incomplete (e.g. 254 bp for *Hexagenia limbata*). The resultant sequences were
82 injected into the ID engine on BOLD [11] and into the BLAST function on GenBank [12]. This analysis
83 delivered correct species identifications for six specimens (35%) on BOLD and for nine (53%) on GenBank.

84 The present study was initiated by downloading the 17 sequences from GenBank. They were then
85 resubmitted to the BOLD ID engine and to GenBank BLAST with self matches excluded. Because some of
86 the resultant identifications deviated from those reported in [10], the factors responsible for this
87 discordance were examined.

88 **Results and Discussion**

89 ***ID Results from BOLD and GenBank:*** Table S1A compares the ID results for the 17 specimens between
90 [10] and those obtained in the present study. The IDs from BLAST matched those reported by [10] as did
91 ten of the IDs from BOLD. The other seven IDs from BOLD corresponded to those from GenBank, but not
92 with the results in [10]. There was a simple explanation for this discordance. Meiklejohn et al. [10] had
93 submitted the reverse complement rather than the coding sequence into the ID engine on BOLD, an
94 approach which generated distant matches. Avoiding this misstep, the number of “correct” identifications
95 generated by BOLD and GenBank was similar (12/17 at the genus level, 9/17 at the species level).

96

97 ***Factors Responsible for Four ‘Errors’ in Generic Assignment:*** Both BOLD and GenBank delivered generic
98 identifications deemed incorrect for four specimens. In each case, the query sequence showed close
99 similarity (95–100% in three cases, 90% in one) to taxa belonging to a different order than that analyzed
100 (Supplementary Files 1 & 2). These discordances could either reflect errors in the reference libraries or in
101 the query sequences. The cause for one misidentification was certain; it arose through internal cross-
102 contamination as the sequence for *Hexagenia limbata* was a truncated version of that for *Glossina palpalis*
103 (identical at all 250 bp that overlapped). The other three mismatches involved taxa (springtail, gall midge,
104 strepsipteran) unrepresented among the 17 tested species ruling out internal contamination. Moreover,
105 because of their striking morphological differences to the test taxa (house fly, dragonfly, flea),
106 misidentification can be excluded as a cause. This leaves two possible explanations – contamination in the
107 reference sequence libraries or in the query sequences. Because each query sequence was embedded

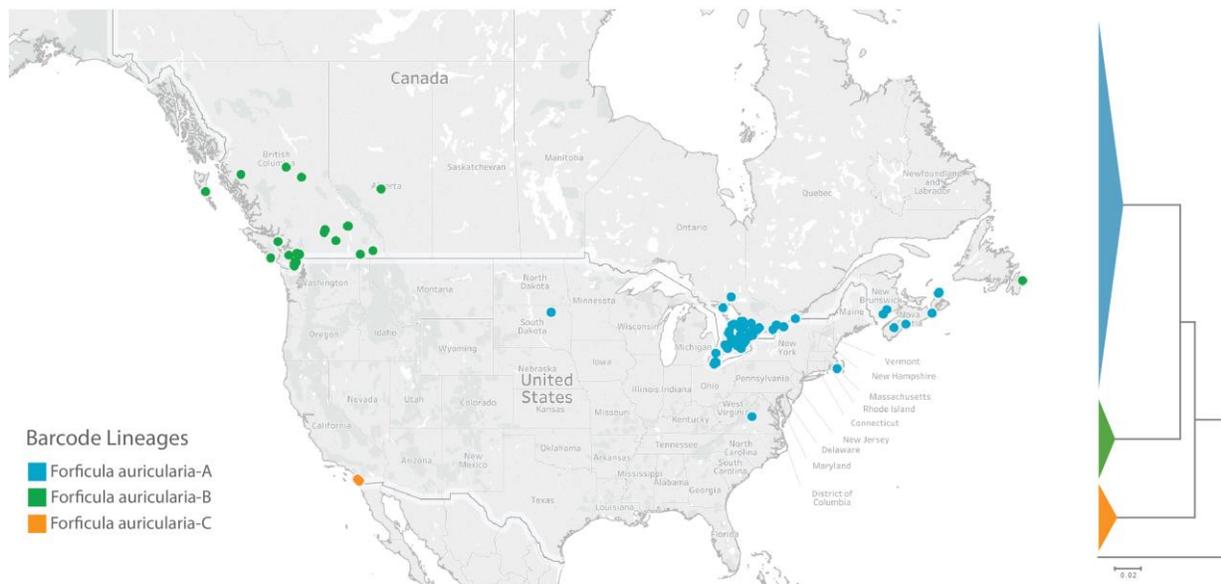
108 within many independently generated reference sequences from another order, these cases of
109 misidentification clearly arose from contamination of the query sequences. Cross-contamination is a well-
110 recognized risk when working with museum specimens so it is standard practice to check for its
111 occurrence [13,14], but Meiklejohn et al. [10] make no mention of exercising precautions in this regard.
112 After excluding these four cases, the number of correct identifications for BOLD and GenBank (12/13 for
113 genus, 9/13 for species) was identical.

114 ***Need for Taxonomic Validation of Museum Specimens:*** The four remaining ‘incorrect’ identifications all
115 involved cases where BOLD and GenBank assigned the query sequence to a species closely related to the
116 taxon analyzed by Meiklejohn et al. [10]. As such, the evidence for misidentification rests on the
117 presumption that their specimens were correctly identified. While the National Museum of Natural
118 History is considered one of the better curated of North American insect collections, the quality of
119 identification of individual specimens depends on the expertise available and the time elapsed since they
120 were assigned to a species. [15]. As such, specimens may be misidentified, mirroring the situation
121 reported in other studies. For example, Meier & Dikow [16] found that 12% of all species-level
122 identifications for a genus of asilid flies from various collections were wrong. Similarly, Muona [17] found
123 that from 1–25% of beetles belonging to two easily discriminated species pairs and one species tetrad
124 were incorrectly identified in a major collection. Similarly, efforts to build a DNA barcode reference library
125 for North American Lepidoptera exposed many misidentified specimens and overlooked cryptic species in
126 major collections [18]. Importantly, all four cases of apparent misidentification reported by Meiklejohn
127 et al. [10] involve species whose recognition is not straightforward. The sole case of generic
128 misidentification involved a presumptive specimen of the cat flea, *Ctenocephalides felis*, whose sequence
129 matched those for the human flea, *Pulex irritans*, on BOLD and GenBank. Because the latter species often
130 uses cats as a host and is morphologically similar to *C. felis*, there is a risk of misidentification. BOLD holds
131 nearly 1,200 records, contributed by 15 institutions, representing four species of *Ctenocephalides* and

132 each possesses a divergent array of barcode sequences. Although the taxonomy of these species is not
133 fully resolved [19], the barcode results support the monophyly of all species in the genus while *P. irritans*
134 forms a sister taxon. Because of the large number of records in the reference library and their derivation
135 from multiple laboratories, the supposed specimen of *C. felis* analyzed by Meiklejohn et al. [10] is almost
136 certainly *P. irritans*. The three remaining cases of presumptive species-level misidentifications involved
137 genera (*Gryllus*, *Glossina*, *Phaenaeus*) with complex taxonomy. One of the three species, *Gryllus*
138 *assimilans*, was formerly thought to be widely distributed in the New World, but it is now recognized to
139 be a complex of 8+ species, several of which can only be reliably distinguished by their call or life history
140 [20]. Similarly, the query species of tsetse fly (*G. palpalis*) is known to be a complex that includes *G.*
141 *brevipalpis* [21–24], the species identified by BOLD and GenBank. The third species, *Phanaeus vindex*, is
142 also a complex of at least two species [25], but it is likely more diverse as records for it on BOLD belong to
143 four distinct COI sequence clusters. Because of these taxonomic uncertainties, the four cases of
144 presumptive species- or genus-level misidentifications are best viewed as unconfirmed.

145 **Resolving Taxonomic Uncertainty:** As the preceding section reveals, efforts to assess the resolution of
146 DNA barcodes is often constrained by poor taxonomy. It is certain that some records on BOLD and
147 GenBank derive from misidentified specimens, but there is no easy path to correct them. This fact was
148 powerfully demonstrated by Mutanen et al. [26] study of DNA barcode variation in 4,977 species of
149 European Lepidoptera which revealed that 60% of the cases initially thought to indicate compromised
150 species resolution or DNA barcode sharing actually arose as a result of misidentifications, databasing
151 errors, or flawed taxonomy. As the taxonomic system for European Lepidoptera is very advanced, similar
152 issues will be a greater impediment in most other groups. Databases like BOLD and GenBank record these
153 divergences in taxonomic opinion, but they cannot resolve them, providing strong motivation for
154 approaches that sidestep this barrier. The Barcode Index Number (BIN) system is a good candidate as it
155 makes it possible to objectively register genetically diversified lineages [27]. One of the ‘species’ in the

156 current study, *Forficula auricularia*, provides a good example of the enhanced geographic resolution
157 offered by BINs that could be useful in forensic contexts. This taxon has been known to include two
158 lineages with differing distributions and life histories for >20 years, but it still remains a single recognized
159 species [28,29]. Barcode results indicate that North American populations actually include three divergent
160 lineages with allopatric distributions (Figure 1). As such, BIN assignments provide information on the
161 geographic distributions of the component lineages of this species complex that could be important in
162 certain forensic contexts, but that would be overlooked by a species-based assignment. Because most
163 species of multicellular organisms await description, it is certain that there are many other cases where
164 BIN-based analysis will enhance geographic resolution.



165
166 **Figure 1:** Geographic distributions and sequence clustering of the three barcode lineages of *Forficula*
167 *auricularia* in North America.

168
169 **Distinction Between BOLD and GenBank:** It is not surprising that BOLD and GenBank demonstrated similar
170 performance in identification, once operational issues were resolved, as many records appear in both
171 platforms. Sequences of COI submitted independently to GenBank are mined and entered into BOLD

172 periodically while records from BOLD are submitted to GenBank when they are published. At present,
173 11% of all COI barcode records on BOLD originate from GenBank, while 75% of the COI barcodes on
174 GenBank derive from BOLD. Although many records are shared, the two platforms diverge in collateral
175 data. For example, for the 17 species of insects analyzed in [10], 65% of the records originating from BOLD
176 possess GPS coordinates, 60% have trace electropherograms, and 40% have specimen images, while only
177 26% of those originating from GenBank had GPS coordinates and all lacked images and
178 electropherograms. In addition, BOLD employs BINs to integrate records that lack a genus or species
179 designation with those that possess them. These extended data elements and functionality are a
180 valuable, often essential, component in the evaluation of identification results.

181 **Conclusions and Path Forward:** Six of the 17 species examined by Meiklejohn et al. [10] escaped
182 operational errors, but the other 11 did not (Table 1), explaining the low identification success they
183 reported. Even after correcting for the use of reverse complements, the effectiveness of DNA barcoding
184 could not be evaluated for eight species, those impacted by sequence contamination or taxonomic
185 uncertainty. Importantly, DNA barcode records in BOLD and GenBank did deliver a correct species
186 assignment for the other nine species. While the outcome for these species is reassuring, the lack of an
187 outcome for other taxa reveals the need for improved protocols. Clearly, two conditions need to be
188 satisfied to ensure a correct identification – the query sequences must be legitimate and the reference
189 libraries must be well-validated. As a start, any study that aims to employ DNA barcodes for species
190 identification should include steps to ensure the sequences recovered are valid by including positive and
191 negative controls, by assessing sequence quality, and by checking for contaminants (Figure 2). Presuming
192 the query sequences pass these quality checks, the generation of a reliable identification requires a
193 comprehensive, well-validated reference library. Because BOLD is a workbench for the DNA barcode
194 research community, it will always contain sequences from specimens whose identifications are being
195 refined. The establishment of a Barcode REF library, based upon a small number of carefully validated

196 records for each species, would represent an important step towards improving its capacity to generate
197 reliable identifications. Under ideal circumstances, the reference sequence for each species would derive
198 from its holotype. However, because 90% of all multicellular organisms await description, and the status
199 of many described species groups is uncertain, these efforts will need to be reinforced by a BIN-based
200 approach.

201

202

Specimen #	ID	Reverse Complement	Contamination	Incorrect ID
1	<i>Phanaeus vindex</i>	—	—	Yes
2	<i>Forficula auricularia</i>	Yes	—	—
3	<i>Chrysomya rufifacies</i>	—	—	—
4	<i>Calliphora vicina</i>	—	—	—
5	<i>Aedes aegypti</i>	—	—	—
6	<i>Glossina palpalis</i>	—	—	Yes
7	<i>Musca domestica</i>	Yes	Yes	N.D.
8	<i>Hexagenia limbata</i>	—	Yes	N.D.
9	<i>Vespula squamosa</i>	—	—	—
10	<i>Callosamia promethea</i>	—	—	—
11	<i>Danaus plexippus</i>	—	—	—
12	<i>Merope tuber</i>	Yes	—	—
13	<i>Ululodes quadripunctatus</i>	Yes	—	—
14	<i>Gomphus exilis</i>	Yes	Yes	N.D.
15	<i>Gryllus assimilis</i>	—	—	Yes
16	<i>Ctenocephalic felis</i>	Yes	—	Yes
17	<i>Pediculus humanus capitis</i>	Yes	Yes	N.D.

203

204

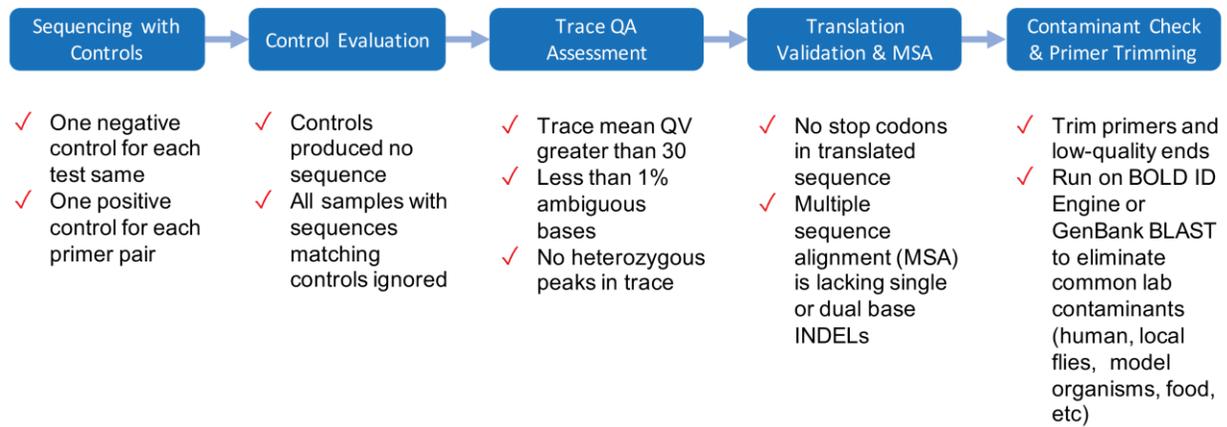
205

206 **Table 1:** Three categories of operational errors which compromised efforts by Meiklejohn et al. [8] to test
207 the effectiveness of the BOLD and GenBank reference libraries in identifying 17 insect species.

208

209

210



211

212 **Figure 2:** Five key workflow features to maximize the chance of recovering reliable sequence records

213

214 Acknowledgments

215 Funding from the Canada First Research Excellence Fund, the Ontario Ministry of Research and
216 Innovation, the Canada Foundation for Innovation, and NSERC support development of the BOLD
217 platform and its computational infrastructure.

218

219 References

- 220 1. Chang C-H, Dai W-Y, Chen T-Y, Lee A-H, Hou H-Y, Liu S-H, et al. DNA barcoding reveals CITES-listed
221 species among Taiwanese government seized chelonian specimens. *Genome*. 2018 61;615–624.
- 222 2. Koroiva R, de Souza MS, Roque FO, Pepinelli M. DNA barcodes for forensically important fly
223 species in Brazil. *J Med Entomol*. 2018;55: 1055–1061.
- 224 3. Coissac E, Hollingsworth PM, Lavergne S, Taberlet P. From barcodes to genomes: extending the
225 concept of DNA barcoding. *Mol Ecol*. 2016;25: 1423–1428. doi:10.1111/mec.13549
- 226 4. Pentinsaari M, Hebert PDN, Mutanen M. Barcoding beetles: A regional survey of 1872 species

- 227 reveals high identification success and unusually deep interspecific divergences. PLoS One. 2014;9:
228 e108651. doi:10.1371/journal.pone.0108651
- 229 5. Hajibabaei M, Janzen DH, Burns JM, Hallwachs W, Hebert PDN. DNA barcodes distinguish
230 species of tropical Lepidoptera. Proc Natl Acad Sci USA. 2006;103: 968–971.
231 doi:10.1073/pnas.0510466103
- 232 6. Hausmann A, Godfray HC, Huemer P, Mutanen M, Rougerie R, van Nieuwerkerken EJ, et al. Genetic
233 patterns in European geometrid moths revealed by the Barcode Index Number (BIN) system. PLoS
234 One. 2013;8: e84518. doi:10.1371/journal.pone.0084518
- 235 7. Hendrich L, Morinière J, Haszprunar G, Hebert PDN, Hausmann A, Köhler F, et al. A
236 comprehensive DNA barcode database for Central European beetles with a focus on Germany:
237 Adding more than 3,500 identified species to BOLD. Mol Ecol Resour. 2015;15: 795–818.
238 doi:10.1111/1755-0998.12354
- 239 8. Huemer P, Mutanen M, Sefc KM, Hebert PDN. Testing DNA barcode performance in 1000
240 species of European Lepidoptera: large geographic distances have small genetic impacts. PLoS One
241 2014;9: e115774. doi:10.1371/journal.pone.0115774
- 242 9. Kerr KCR, Stoeckle MY, Dove CJ, Weigt LA, Francis CM, Hebert PDN. Comprehensive DNA
243 barcode coverage of North American birds. Mol Ecol Notes. 2007;7: 535–543. doi:10.1111/j.1471-
244 8286.2007.01670.x
- 245 10. Meiklejohn KA, Damaso N, Robertson JM. Assessment of BOLD and GenBank – Their accuracy
246 and reliability for the identification of biological materials. PLoS One 2019;14: e0217084.
247 doi:10.1371/journal.pone.0217084
- 248 11. Ratnasingham S, Hebert PDN. BOLD: The Barcode of Life Data System
249 (<http://www.barcodinglife.org>). Mol Ecol Notes. 2007;7: 355–364. doi:10.1111/j.1471-

250 8286.2007.01678.x

251 12. Benson DA, Cavanaugh M, Clark K, Karsch-Mizrachi I, Lipman DJ, Ostell J, Sayers EW. GenBank.
252 Nucleic Acids Res. 2017 4;45(D1):D37-D42. doi: 10.1093/nar/gkw1070.

253 13. Siddall ME, Fontanella FM, Watson SC, Kvist S, Erséus C. Barcoding bamboozled by bacteria:
254 Convergence to metazoan mitochondrial primer targets by marine microbes. Syst Biol. 2009;58:
255 445–451. doi:10.1093/sysbio/syp033

256 14. Mioduchowska M, Czyż MJ, Gołdyn B, Kur J, Sell J. Instances of erroneous DNA barcoding of
257 metazoan invertebrates: Are universal cox1 gene primers too “universal”? PLoS One. 2018;13:
258 e0199609. doi:10.1371/journal.pone.0199609

259 15. McGinley R J. Where's the management in collections management? Planning for improved care,
260 greater use and growth of collections. International Symposium and First World Congress on the
261 Preservation and Conservation of Natural History Collections 1993;3: 309–333.

262 16. Meier R, Dikow T. Significance of specimen databases from taxonomic revisions for estimating
263 and mapping the global species diversity of invertebrates and repatriating reliable specimen data.
264 Conserv Biol. 2004;18: 478–488. doi:10.1111/j.1523-1739.2004.00233.x

265 17. Muona J. Huomioita eläinmuseon kuoriaskokoelmien virhemäärittämisestä [Some observations
266 concerning incorrectly determined beetles in public collections. (Coleoptera)]. Sahlbergia. 2001;6:
267 34–36.

268 18. Levesque-Beaudin V, Rosati ME, Silversen N, Warne CP, Brown A, Telfer AC et al. Museum
269 harvesting in major natural history collections. Genome 2017; 60:962. doi: 10.1139/gen-2017-0178

270 19. Lawrence AL, Brown GK, Peters B, Spielman DS, Morin-Adeline V, Šlapeta J. High phylogenetic
271 diversity of the cat flea (*Ctenocephalides felis*) at two mitochondrial DNA markers. Med Vet Entomol.

- 272 2014;28: 330–336. doi:10.1111/mve.12051
- 273 20. Weissman DB, Gray DA, Pham HT, Tijssen P. Billions and billions sold: Pet-feeder crickets
274 (Orthoptera: Gryllidae), commercial cricket farms, and epizootic densovirus, and government
275 regulations make for a potential disaster. 2012; Zootaxa 3504:67-88.
- 276 21. Gooding RH, Krafur ES. TSETSE GENETICS: Contributions to biology, systematics, and control of
277 Tsetse flies. Annu Rev Entomol. 2005;50: 101–123. doi:10.1146/annurev.ento.50.071803.130443
- 278 22. Gooding RH, Solano P, Ravel S. X-chromosome mapping experiments suggest occurrence of
279 cryptic species in the tsetse fly *Glossina palpalis palpalis*. Can J Zool. 2004;82: 1902–1909.
280 doi:10.1139/z05-002
- 281 23. Dyer NA, Furtado A, Cano J, Ferreira F, Odete Afonso M, Ndong-Mabale N, et al. Evidence for a
282 discrete evolutionary lineage within Equatorial Guinea suggests that the tsetse fly *Glossina palpalis*
283 *palpalis* exists as a species complex. Mol Ecol. 2009;18: 3268–3282. doi:10.1111/j.1365-
284 294X.2009.04265.x
- 285 24. De Meeûs T, Bouyer J, Ravel S, Solano P. Ecotype evolution in *Glossina palpalis* subspecies,
286 major vectors of Sleeping Sickness. PLoS Negl Trop Dis. 2015;9: e0003497.
287 doi:10.1371/journal.pntd.0003497
- 288 25. Price DL. Phylogeny and biogeography of the dung beetle genus *Phanaeus* (Coleoptera:
289 Scarabaeidae). Systematic Entomology. 2009;34: 131–150.
- 290 26. Mutanen M, Kivelä SM, Vos RA, Doorenweerd C, Ratnasingham S, Hausmann A, et al. Species-
291 level para- and polyphyly in DNA barcode gene trees: Strong operational bias in European
292 Lepidoptera. Syst Biol. 2016;65(6):1024-1040. doi: 10.1093/sysbio/syw044
- 293 27. Ratnasingham S, Hebert PDN. A DNA-based registry for all animal species: the Barcode Index

- 294 Number (BIN) system. PLoS One. 2013;8:e66213. doi: 10.1371/journal.pone.0066213.
- 295 28. Guillet S, Guiller A, Deunff J, Vancassel M. Analysis of a contact zone in the *Forficula auricularia*
296 L. (Dermaptera: Forficulidae) species complex in the Pyrenean Mountains. Heredity. 2000;85: 444–
297 449. doi:10.1046/j.1365-2540.2000.00775.x
- 298 29. Guillet S, Josselin N, Vancassel M. Multiple Introductions of the *Forficula auricularia* species
299 complex (Dermaptera: Forficulidae) in eastern North America. Can Entomol. 2000;132: 49–57.
300 doi:10.4039/Ent13249-1

301

302

303

304

305

306

307

308 **Supplementary Data**

309 Table S1. Comparison of query results (top matches) for 17 insect species between Meiklejohn et al.
310 2019 (doi:10.1371/journal.pone.0217084) and the present study.

311 Supplementary file 1 (xlsx). Top 20 matches in GenBank BLAST queries for the four specimens deemed
312 cross-contaminations.

313 Supplementary file 2 (xlsx). Top 20 matches from queries to the BOLD ID engine for four specimens
314 whose COI sequences derive from cross-contamination