

1 **Decoding rRNA sequences for improved metagenomics in sylvatic mosquito**
2 **species**

3 Cassandra Koh¹, Lionel Frangeul¹, Hervé Blanc¹, Carine Ngoagouni², Sébastien Boyer³, Philippe
4 Dussart⁴, Nina Grau⁵, Romain Girod⁵, Jean-Bernard Duchemin⁶ and Maria-Carla Saleh^{1*}

5 ¹ Viruses and RNA Interference Unit, Institut Pasteur, Université Paris Cité, CNRS UMR3569, Paris
6 F-75015, France

7 ² Medical Entomology Laboratory, Institut Pasteur de Bangui, Bangui PO Box 923, Central African
8 Republic

9 ³ Medical and Veterinary Entomology Unit, Institut Pasteur du Cambodge, Phnom Penh 12201,
10 Cambodia

11 ⁴ Virology Unit, Institut Pasteur du Cambodge, Phnom Penh 12201, Cambodia

12 ⁵ Medical Entomology Unit, Institut Pasteur de Madagascar, Antananarivo 101, Madagascar

13 ⁶ Vectopôle Amazonien Emile Abonnenc, Institut Pasteur de la Guyane, Cayenne 97306, French
14 Guiana

15 * To whom correspondence should be addressed.

16

17 Present address: Philippe Dussart, Institut Pasteur de Madagascar, Antananarivo 101, Madagascar;

18 Nina Grau, Sciences Economiques et Sociales de la Santé et Traitement de l'Information Médicale,

19 Faculté de Médecine, Marseille 13005, France

20

21 Email addresses:

22 Cassandra Koh, cassandra.koh@pasteur.fr; Lionel Frangeul, lionel.frangeul@pasteur.fr; Hervé Blanc,

23 herve.blanc@pasteur.fr; Carine Ngoagouni, carine.ngoagouni@pasteur-bangui.cf; Sébastien Boyer,

24 sboyer@pasteur-kh.org; Philippe Dussart, pdussart@pasteur.mg; Nina Grau, ngrau@pasteur.mg;

25 Romain Girod, rgirod@pasteur.mg; Jean-Bernard Duchemin, jbduchemin@pasteur-cayenne.fr; Maria-

26 Carla Saleh, carla.saleh@pasteur.fr.

27

28

29

30

31 **ABSTRACT**

32 **Background:** RNA-seq metagenomics on mosquitoes for surveillance and microbiome or pathogen
33 discovery allows us to understand the disease ecology of arboviruses. A major hurdle in these studies
34 is the depletion of overabundant ribosomal RNA (rRNA), commonly achieved using oligo-based
35 protocols. The lack of publicly available complete reference rRNA sequences for many mosquito
36 vector species narrows the range of such studies, causing a knowledge bias in mosquito vector
37 biology. Here we describe a strategy to assemble full-length 28S and 18S rRNA sequences of 29
38 sylvatic and peri-urban mosquito species sampled from Cambodia, the Central African Republic,
39 Madagascar, and French Guiana.

40 **Results:** Our score-based strategy successfully parses rRNA reads into insect and non-insect
41 sources, leading to the assembly of complete rRNA sequences for all specimens in the study. We
42 then evaluated the functionality of rRNA sequences as barcodes for taxonomy and phylogeny relative
43 to the mitochondrial *cytochrome c oxidase I* (COI) gene marker system. rRNA- and COI-based
44 phylogenetic inferences share little congruity. However, the former allowed for molecular species
45 identification when COI sequences were ambiguous or unavailable and revealed better supported
46 intergeneric evolutionary histories concordant with contemporary mosquito systematics.

47 **Conclusions:** The presented assembly strategy and the expansion of the rRNA reference library in
48 public databases by 234 novel complete 28S and 18S rRNA sequences provide a new tool in the form
49 of an RNA marker system to improve mosquito RNA-seq metagenomics. More holistic insights on
50 mosquito vector ecology will benefit the design of public health measures against arboviral diseases.

51 **Keywords:** ribosomal RNA, rRNA, depletion, RNA-seq, metagenomics, mosquito, phylogeny,
52 taxonomy

53 BACKGROUND

54 Mosquitoes top the list of vectors for arthropod-borne diseases and are implicated in the transmission
55 of many human pathogens responsible for arboviral diseases, malaria, and lymphatic filariasis (1).
56 Metagenomic studies on field-captured mosquito specimens for the purposes of surveillance and
57 microbiome or pathogen discovery are becoming increasingly important under the One Health
58 concept, which emphasises the importance of considering the role of biotic and abiotic elements
59 within the same ecosystem in contributing to zoonotic disease transmission (2). With next-generation
60 sequencing technologies becoming more accessible, these studies are increasing in frequency,
61 providing an unprecedented understanding of the interfaces among mosquitoes, their environment,
62 and their animal and human hosts. Currently, there is a strong focus on only a handful of species from
63 three genera of mosquitoes (*Aedes*, *Culex*, and *Anopheles*) due to their medical importance and
64 anthropophilic behaviour despite evidence that other species are also competent for the pathogens in
65 question. This narrows our knowledge of mosquito vector ecology to urban-dwelling species when
66 forest-dwelling mosquitoes are the ones responsible for maintaining the sylvatic transmission of
67 arboviruses among their reservoir hosts, which precedes autochthonous transmission in human
68 populations (3).

69 Ribosomal RNAs (rRNA) are non-coding RNA molecules that make up the ribosomal complexes
70 involved in translation of messenger RNA into proteins. In eukaryotes, 28S and 18S rRNA molecules
71 typically span lengths of four and two kilobases, respectively (4). They comprise at least 80% of the
72 total cellular RNA population. In RNA-seq experiments, their depletion is a necessary step during
73 library preparation where it is not possible to selectively enrich target signals (5). To achieve this, the
74 most routinely used depletion protocols require knowledge of rRNA sequence of the species of
75 interest. These protocols involve hybridizing antisense oligos (probes or primers) to rRNA molecules
76 followed by digestion by ribonucleases (5,6) or removal by bead capture (7).

77 For well-studied mosquito species, reference rRNA sequences are readily available on public
78 sequence databases such as GenBank or SILVA. As it is conventionally accepted that rDNA coding
79 regions are highly conserved, it may seem conceivable to use oligo-based depletion protocols
80 designed for one mosquito species on another. However, we found that within the family *Culicidae*
81 this is not always true. There is enough sequence divergence such that *Ae. aegypti*-based probes
82 produced poor depletion in *Culex* and *Anopheles* mosquitoes. In addition, full-length rRNA sequences

83 are much less represented compared to other molecular markers such as the *cytochrome c oxidase*
84 *subunit I* (COI) gene, which is the most widely used marker for molecular taxonomy and forms the
85 basis of the Barcode of Life Data System (BOLD) (8,9). The lack of reliable rRNA depletion methods
86 could deter mosquito metagenomic studies from expanding their sampling diversity. The inclusion of
87 lesser studied yet ecologically relevant species is imperative.

88 To address this, we sought to determine the 28S and 18S rRNA sequences of a diverse set of
89 sylvatic and peri-urban mosquito species across Cambodia, the Central African Republic,
90 Madagascar, and French Guiana. We employed a unique score-based read filtration strategy to
91 remove interfering non-mosquito rRNA reads to ensure accurate *de novo* assembly and generated
92 122 complete 28S and 114 complete 18S sequences from 29 mosquito species. This strategy would
93 facilitate the assembly of more rRNA sequences to expand the rRNA reference library. In parallel, we
94 obtained COI sequences to confirm morphology-based species identification and to compare
95 phylogenetic relationships inferred from the DNA and RNA markers, leading us to propose the use of
96 28S and 18S rRNA sequences as “rRNA barcodes”. Our sequence dataset enables rRNA-based
97 streamlined molecular species identification during RNA-seq and allows for the design of species-
98 specific oligos for cost-effective rRNA depletion for a broader range of mosquito species.

99

100 **RESULTS**

101

102 **Poor rRNA depletion using non-specific depletion methods**

103 During library preparations of mosquito samples for RNA-seq, routinely used methods for depleting
104 rRNA are commercial kits optimised for human or mice samples (10–15) or probe hybridisation
105 followed by ribonuclease digestion where the probes are 80-100 base pair antisense oligos. In cases
106 where the reference rRNA sequence of the target species is not known, oligos would be designed
107 based on the rRNA sequence of the closest related species (25, this study). These methods should,
108 in theory, be able to produce acceptable rRNA depletion efficiencies assuming that rRNA sequences
109 have high degrees of homology across species. However, in our hands we found that using probes
110 designed for the *Ae. aegypti* rRNA sequence followed by RNase H digestion according to the protocol
111 published by Morlan *et al.* (17) produced poor depletion in *Ae. albopictus*, and worse still in Culicine
112 and Anopheline species (Figure 1A). Additionally, the lack of reference rRNA sequences

113 compromises the clean-up of remaining rRNA reads from sequencing data, as reads belonging to
 114 more divergent regions do not map to a reference sequence from a different species. To solve this
 115 and to enable RNA-seq metagenomics on a wider range of mosquito species, we performed RNA-seq
 116 to obtain reference rRNA sequences for 29 mosquito species across nine genera from Cambodia, the
 117 Central African Republic, Madagascar, and French Guiana. Most of these species are associated with
 118 vector activity for various pathogens in their respective ecologies (Table 1).

119

120 **Table 1.** List of mosquito species represented in this study and their vector status.

Mosquito taxonomy*	Origin**	Collection site (ecosystem type)	Vector for***	Reference
<i>Aedes (Fredardsius) vittatus</i>	CF	rural (village)	DENV, ZIKV, CHIKV, YFV	(18)
<i>Aedes (Ochlerotatus) scapularis</i>	GF	rural (village)	YFV	(19)
<i>Aedes (Ochlerotatus) serratus</i>	GF	rural (village)	YFV, OROV	(20,21)
<i>Aedes (Stegomyia) aegypti</i>	CF	urban	DENV, ZIKV, CHIKV, YFV	(22)
<i>Aedes (Stegomyia) albopictus</i>	CF, KH	rural (village, nature reserve)	DENV, ZIKV, CHIKV, YFV, JEV	(22,23)
<i>Aedes (Stegomyia) simpsoni</i>	CF	rural (village)	YFV	(24)
<i>Anopheles (Anopheles) baezai</i>	KH	rural (nature reserve)	unreported	–
<i>Anopheles (Anopheles) coustani</i>	MG, CF	rural (village)	RVFV, malaria	(25–27)
<i>Anopheles (Cellia) funestus</i>	MG, CF	rural (village)	ONNV, malaria	(28)
<i>Anopheles (Cellia) gambiae</i>	MG, CF	rural (village)	ONNV, malaria	(29)
<i>Anopheles (Cellia) squamosus</i>	MG	rural (village)	RVFV, malaria	(27,30)
<i>Coquillettidia (Rhynchotaenia) venezuelensis</i>	GF	rural (village)	OROV	(21)
<i>Culex (Culex) antennatus</i>	MG	rural (village)	RVFV	(26,27)
<i>Culex (Culex) duttoni</i>	CF	rural (village)	unreported	–
<i>Culex (Culex) neavei</i>	MG	rural (village)	USUV	(31)
<i>Culex (Culex) orientalis</i>	KH	rural (nature reserve)	JEV	(32)
<i>Culex (Culex) perexiguus</i>	MG	rural (village)	WNV	(33)

<i>Culex (Culex) pseudovishnui</i>	KH	rural (nature reserve)	JEV	(23,34)
<i>Culex (Culex) quinquefasciatus</i>	MG, CF, KH	rural (village, nature reserve)	ZIKV, JEV, WNV, DENV, SLEV, RVFV, <i>Wuchereria bancrofti</i>	(34–36)
<i>Culex (Culex) tritaeniorhynchus</i>	MG, KH	rural (village, nature reserve)	JEV, WNV, RVFV	(23,34)
<i>Culex (Melanoconion) spissipes</i>	GF	rural (village)	VEEV	(37)
<i>Culex (Melanoconion) portesi</i>	GF	rural (village)	VEEV, TONV	(37,38)
<i>Culex (Melanoconion) pedroi</i>	GF	rural (village)	EEEV, VEEV, MADV	(38,39)
<i>Culex (Oculeomyia) bitaeniorhynchus</i>	MG, KH	rural (village, nature reserve)	JEV	(23,34)
<i>Culex (Oculeomyia) poecilipes</i>	MG	rural (village)	RVFV	(35)
<i>Eretmapodites intermedius</i>	CF	rural (village)	unreported	–
<i>Limatus durhamii</i>	GF	rural (village)	ZIKV	(40)
<i>Mansonia (Mansonia) titillans</i>	GF	rural (village)	VEEV, SLEV	(41,42)
<i>Mansonia (Mansonioides) indiana</i>	KH	rural (nature reserve)	JEV	(43)
<i>Mansonia (Mansonioides) uniformis</i>	MG, CF, KH	rural (village, nature reserve)	WNV, RVFV, <i>Wuchereria bancrofti</i>	(28,34,44)
<i>Mimomyia (Etorleptomyia) mediolineata</i>	MG	rural (village)	unreported	–
<i>Psorophora (Janthinosoma) ferox</i>	GF	rural (village)	ROCV	(45)
<i>Uranotaenia (Uranotaenia) geometrica</i>	GF	rural (village)	unreported	–

121 * () indicates subgenus

122 ** Origin countries are listed as their ISO alpha-2 codes: Central African Republic, CF; Cambodia, KH;
123 Madagascar, MG; French Guiana, GF.

124 ** dengue virus, DENV; Zika virus, ZIKV; chikungunya virus, CHIKV; Yellow Fever virus, YFV; Oropouche virus,
125 OROV; Japanese encephalitis virus, JEV; Rift Valley Fever virus, RVFV; O’Nyong Nyong virus, ONNV; Usutu
126 virus, USUV; West Nile virus, WNV; Saint Louis encephalitis virus, SLEV; Venezuelan equine encephalitis
127 virus, VEEV; Tonate virus, TONV; Eastern equine encephalitis virus, EEEV; Madariaga virus, MADV; Rocio
128 virus, ROCV.

129

130 **rRNA reads filtering and sequence assembly**

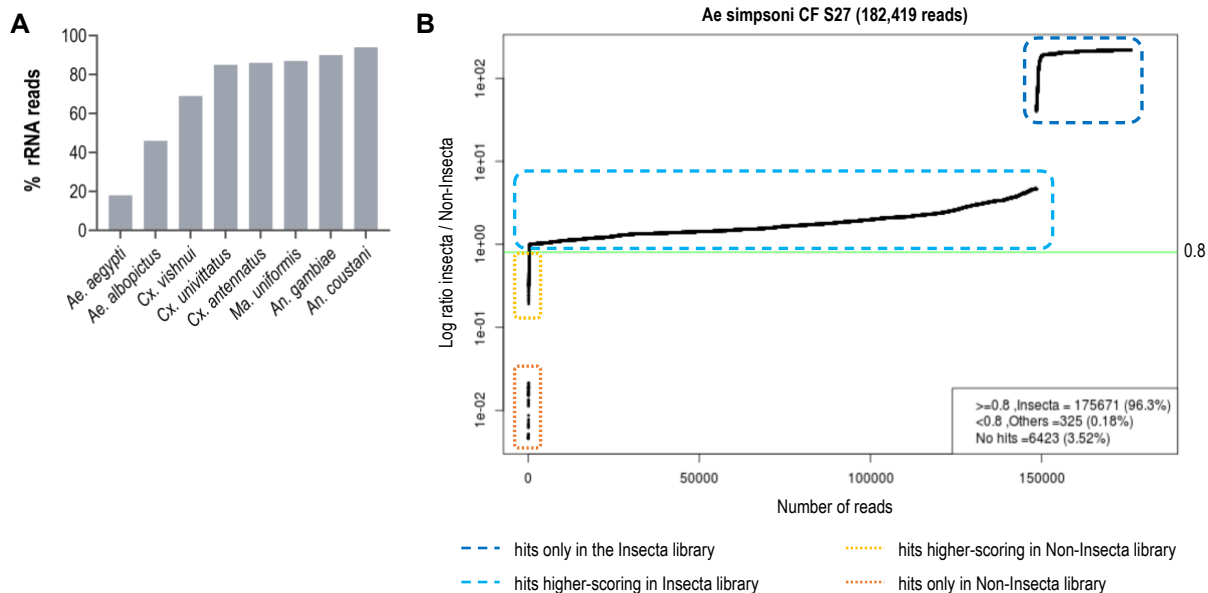
131 Assembling Illumina reads to reconstruct rRNA sequences from total mosquito RNA is not a
132 straightforward task. Apart from host rRNA, total RNA samples also contain rRNA from other
133 organisms associated with the host (microbiota, external parasites, or ingested diet). As all these
134 rRNA sequences contain highly conserved blocks, Illumina reads (150 bp) from these sequences can
135 interfere with and impede the contig assembly of host 28S and 18S rRNA.

136 Our score-based filtration strategy, described in detail in Methods, allowed us to bioinformatically
137 remove interfering rRNA reads and achieve successful de novo assembly of 28S and 18S rRNA
138 sequences for all our specimens. Briefly, for each Illumina read, we computed a ratio of BLAST
139 scores against an Insecta library over a Non-Insecta library. Reads were segregated into four
140 categories: (i) reads that map only to the Insecta library, (ii) reads that map better to the Insecta
141 relative to Non-Insecta library, (iii) reads that map better to the Non-Insecta relative to the Insecta
142 library, and finally (iv) reads that only map to the Non-Insecta library (Figure 1B, Figure S1). By
143 applying a conservative threshold of 0.8 to account for the non-exhaustiveness of reference libraries
144 used, we filtered out reads that likely do not originate from mosquito rRNA. Notably, 15 of our
145 specimens were engorged with vertebrate blood, a rich source of non-mosquito rRNA (Supplementary
146 Table 1). The successful assembly of complete 28S and 18S rRNA sequences demonstrates that this
147 strategy performs as expected even with high amounts of non-host rRNA reads. This is particularly
148 important in studies on field-captured mosquitoes as females are often sampled already having
149 imbibed a blood meal or captured using the human-landing catch technique.

150 We encountered challenges for three specimens morphologically identified as *Ma. africana*
151 (Specimen ID 33-35). COI amplification by PCR did not produce any product, hence COI barcoding
152 could not be used to confirm species identity. In addition, SPAdes was only able to assemble partial
153 length contigs, despite the high number of reads with high scores against the Insecta library. Among
154 other *Mansonia* specimens, the partial length contigs shared the highest similarity with contigs
155 obtained from “*Ma uniformis* CF S51”. We then performed a guided assembly using the 28S and 18S
156 sequences of this specimen as references, which successfully produced full-length contigs. In two of
157 these specimens (Specimen ID 34 and 35), our assembly initially produced two sets of 28S and 18S
158 rRNA sequences, one of which was similar to mosquito rRNA with low coverage and another with ten-
159 fold higher coverage and 95% nucleotide sequence similarity to a *Horreolanus* species of water mite

160 known to parasitize mosquitoes. Our filtration strategy allowed us to obtain rRNA sequences for the
161 mosquito as well as the unknown *Horreolanus* species. This shows that our strategy can be applied to
162 metabarcoding studies where the input material comprises multiple insect species, provided that
163 appropriate reference sequences of the target species or of a close relative are available.

164 Altogether, we were able to assemble 122 28S and 114 18S full-length mosquito rRNA sequences
165 for 29 mosquito species sourced from four countries across three continents. These sequences
166 represent, to our knowledge, the first records for seven mosquito genera: *Coquillettidia*, *Mansonia*,
167 *Limatus*, *Mimomyia*, *Uranotaenia*, *Psorophora*, and *Eretmapodites*. For *Culex*, *Aedes* and *Anopheles*
168 genera, where complete rRNA sequences were already available for a few species, this study
169 provides the first rRNA records for 18 species. The GenBank accession numbers for these sequences
170 and specimen information are listed in Table S1.



178 To verify the assembly accuracy of our rRNA sequences, we constructed a comprehensive

179 phylogenetic tree from the 28S rRNA sequences generated from our study alongside those publicly

180 available from NCBI databases (Figure 2). We applied a search criterion for NCBI sequences with at

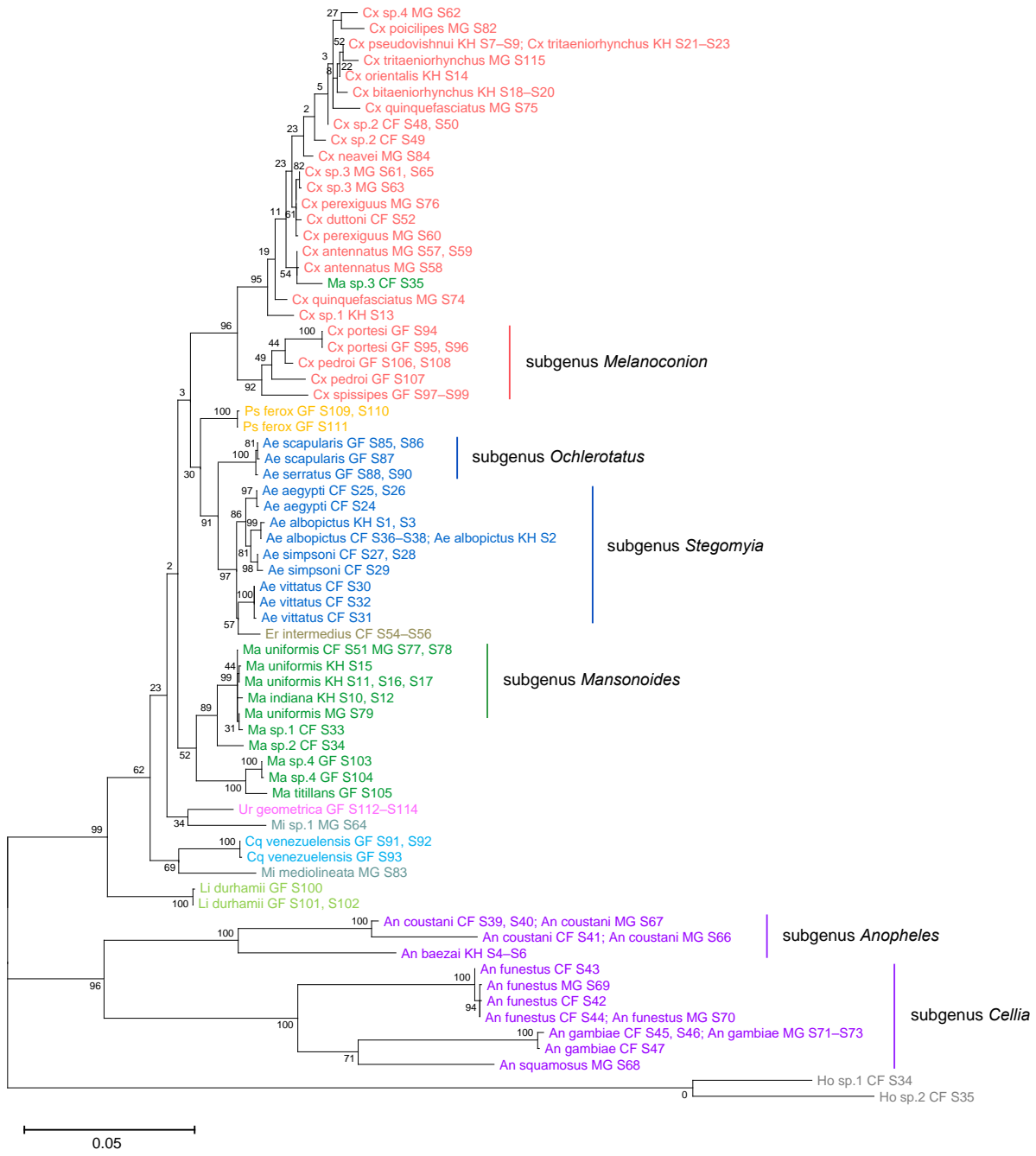
181 least 95% coverage of our sequence lengths (~4000 bp), aiming to have as many species as possible
182 represented. Although we rarely found NCBI entries for the same species represented in our study,
183 the resulting tree showed that our 28S sequences generally clustered according to their respective
184 species and subgenera, supported by moderate to good bootstrap values at terminal nodes. Species
185 taxa generally formed monophyletic clades, apart from *An. gambiae* and *Cx. quinquefasciatus*. *An.*
186 *gambiae* 28S rRNA sequences formed a clade with closely related sequences from *An. arabiensis*,
187 *An. merus*, and *An. coluzzii*, suggesting unusually high interspecies homology for Anophelines or
188 other members of subgenus *Cellia*. Meanwhile, *Cx. quinquefasciatus* 28S rRNA sequences formed a
189 taxon paraphyletic to sister species *Cx. pipiens*.



191 **Figure 2.** Phylogenetic tree based on 28S sequences generated from this study and from NCBI
192 databases (3900 bp) as inferred using maximum-likelihood method and constructed to scale in MEGA
193 X (46). Values at each node indicate bootstrap support (%) from 500 replications. For sequences from
194 this study, each specimen label contains information on its taxonomy, origin (as indicated in 2-letter
195 country codes), and specimen ID. Labels in bold indicate sequences derived from NCBI with
196 accession numbers shown. Label colours indicate genera: *Culex* in coral, *Anopheles* in purple, *Aedes*
197 in dark blue, *Mansonia* in dark green, *Culiseta* in maroon, *Limatus* in light green, *Coquillettidia* in light
198 blue, *Psorophora* in yellow, *Mimomyia* in teal, *Uranotaenia* in pink and *Eretmapodites* in brown. Scale
199 bar at 0.05 is shown.

200

201 28S sequence-based phylogenetic reconstructions (Figure 2, with NCBI sequences; Figure S3,
202 this study only) showed marked incongruence to that of 18S sequences (Figure 3). Although all rRNA
203 trees show clear segregation of genus *Anopheles* from tribes *Aedini* and *Culicini*, the phylogenetic
204 relationships of other genera in this study relative to the greater three are highly variable and weakly
205 supported, particularly in the 18S tree. The 18S tree also showed a number of taxonomic anomalies:
206 (i) the lack of definitive clustering by species within the *Culex* subgenus (ii) the inability to differentiate
207 between 18S sequences of *Cx. pseudovishnui* and *Cx. tritaeniorhynchus*; (iii) the placement of Ma sp.
208 3 CF S35 within a *Culex* clade; and (iv) the lack of a *Mimomyia* clade. The topology of the 18S tree
209 seem to suggest higher sequence divergence between the two *Cx. quinquefasciatus* taxa and
210 between the two *Mimomyia* taxa than in their 28S sequences. However, 28S and 18S rRNA
211 sequences are encoded by linked loci in rDNA clusters and should not be analysed separately.

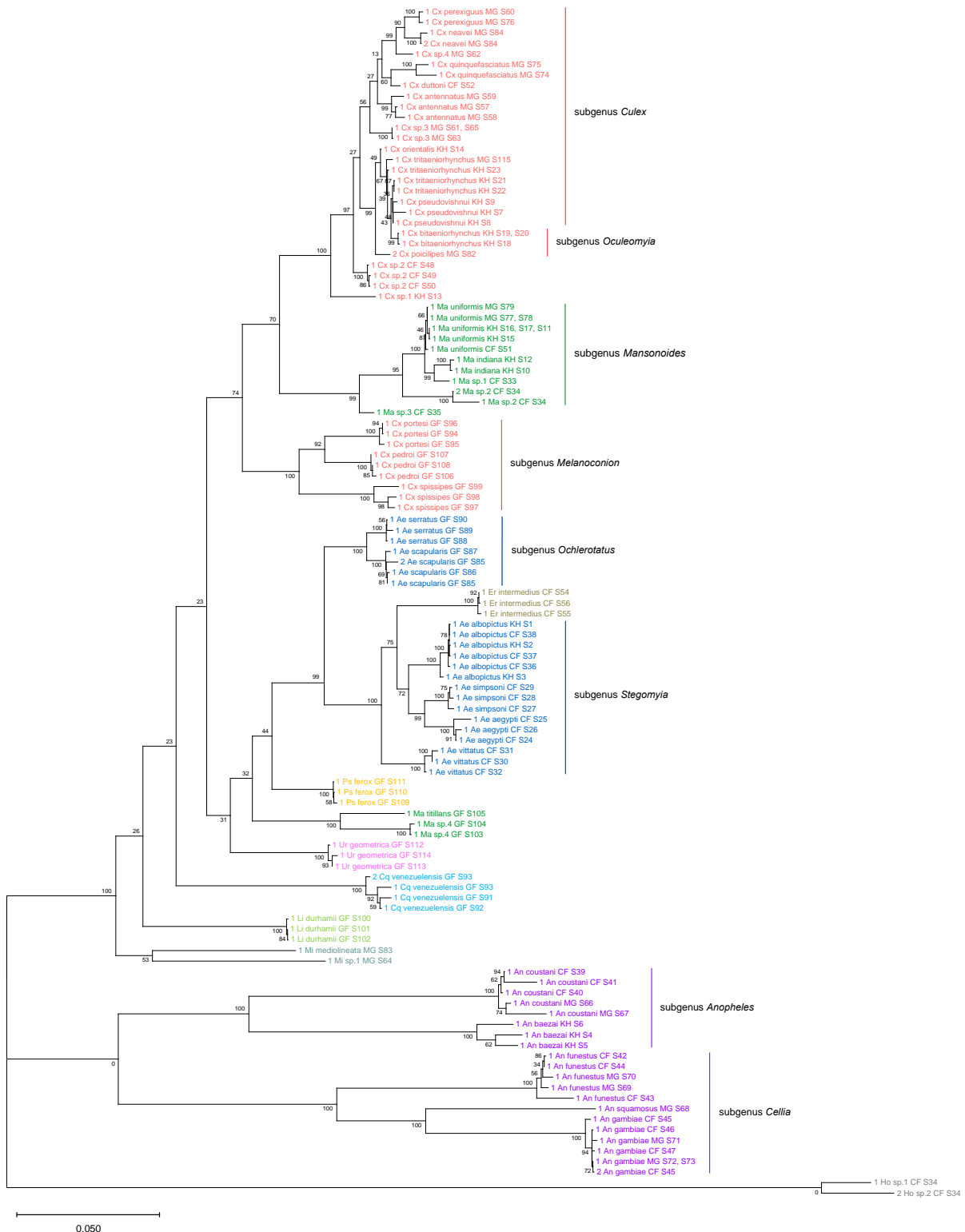


212

213 **Figure 3.** Phylogenetic tree based on 18S sequences (1900 bp) as inferred using maximum-likelihood
 214 method and constructed to scale in MEGA X (46). Values at each node indicate bootstrap support (%)
 215 from 500 replications. Each specimen label contains information on its taxonomy, origin (as indicated
 216 in 2-letter country codes), and specimen ID. Label colours indicate genera: *Culex* in coral, *Anopheles*
 217 in purple, *Aedes* in dark blue, *Mansonia* in dark green, *Limatus* in light green, *Coquillettidia* in light
 218 blue, *Psorophora* in yellow, *Mimomyia* in teal, *Uranotaenia* in pink and *Eretmapodites* in brown. Scale
 219 bar at 0.05 is shown.

220 Indeed, when concatenated 28S+18S rRNA sequences were generated from the same specimens
221 (Figure 4), the phylogenetic tree resulting from these sequences more closely resembles the 28S tree
222 (Figure 2) with regard to the basal position of the *Mimomyia* clade within the *Culicinae* subfamily with
223 good bootstrap support in either tree (84% in 28S tree, 100% in concatenated 28S+18S tree). For
224 internal nodes, bootstrap support values were higher in the concatenated tree compared to the 28S
225 tree. Interestingly, the 28S+18S tree formed an *Aedini* tribe-clade encompassing taxa from genera
226 *Psorophora*, *Aedes*, and *Eretmapodites*, possibly driven by the inclusion of 18S sequences.
227 Concatenating the 28S and 18S sequences also resolved the anomalies found in the 18S tree and
228 added clarity to the close relationship between *Culex* and *Mansonia* taxa. Of note, the *Culex* and
229 *Mansonia* genera are no longer monophyletic in the concatenated 28S+18S tree. Genus *Culex* is
230 paraphyletic with respect to subgenus *Mansonoides* of genus *Mansonia* (Figure 2). *Ma. titillans* and
231 *Ma* sp. 4, which we suspect to be *Ma. pseudotitillans*, always formed a distinct branch in 28S or 18S
232 phylogenies, thus possibly representing a clade of subgenus *Mansonia*.

233 The concatenated 28S+18S tree recapitulates what is classically known about the systematics of
234 our specimens, namely (i) the early divergence of genus *Anopheles* from other *Culicidae* genera, (ii)
235 the division of genus *Anopheles* into two subgenera, *Anopheles* and *Cellia*, (iii) the division of genus
236 *Aedes* into subgenera *Stegomyia* and *Ochlerotatus*, (iv) the divergence of monophyletic subgenus
237 *Melanoconion* within the *Culex* genus (47,48).



238

239 **Figure 4.** Phylogenetic tree based on concatenated 28S and 18S sequences generated from this
 240 study (3900+1900 bp) as inferred using maximum-likelihood method and constructed to scale in
 241 MEGA X (46). Values at each node indicate bootstrap support (%) from 500 replications. For
 242 sequences from this study, each specimen label contains information on its taxonomy, origin (as
 243 indicated in 2-letter country codes), and specimen ID. Label colours indicate genera: *Culex* in coral,

244 *Anopheles* in purple, *Aedes* in dark blue, *Mansonia* in dark green, *Culiseta* in maroon, *Limatus* in light
245 green, *Coquillettidia* in light blue, *Psorophora* in yellow, *Mimomyia* in teal, *Uranotaenia* in pink and
246 *Eretmapodites* in brown. Scale bar at 0.05 is shown.

247

248 **rRNA as a molecular marker for taxonomy and phylogeny**

249 We sequenced a 621 bp region of the COI gene not only to confirm morphological identification of our
250 specimens but also to compare the functionality of rRNA and COI sequences as molecular markers
251 for taxonomic and phylogenetic investigations. COI sequences were able to unequivocally determine
252 the species in most specimens except for the following cases. *An. coustani* COI sequences from our
253 study regardless of specimen origin shared remarkably high nucleotide similarity (>98%) with several
254 other *Anopheles* species such as *An. rhodesiensis*, *An. rufipes*, *An. ziemanni*, *An. tenebrosus*,
255 although *An. coustani* remained the most frequent and closest match. In the case of *Ae. simpsoni*,
256 three specimens were morphologically identified as *Ae. opok* although their COI sequences showed
257 97-100% similarity to that of *Ae. simpsoni*. As NCBI held no records of *Ae. opok* COI, we instead
258 aligned the putative *Ae. simpsoni* COI sequences against *Ae. luteocephalus* and *Ae. africanus*, sister
259 species of *Ae. opok* and found they shared only 90% and 89% similarity, respectively. Given this
260 significant divergence, we concluded these specimens to be *Ae. simpsoni*. Ambiguous results were
261 especially frequent among *Culex* specimens belonging to the *Cx. pipiens* or *Cx. vishnui* species
262 groups, where the query sequence differed with either of the top two hits by a single nucleotide. For
263 example, between *Cx. quinquefasciatus* and *Cx. pipiens* of the *Cx. pipiens* species group, and
264 between *Cx. vishnui* and *Cx. tritaeniorhynchus* of the *Cx. vishnui* species group.

265 Among our three specimens of *Ma. titillans*, two appeared to belong to a single species that is
266 different but closely related to *Ma. titillans*. We surmised that these specimens could instead be *Ma.*
267 *pseudotitillans* based on morphological similarity but were not able to verify this by molecular means
268 as no COI reference sequence is available for this species. These specimens are hence putatively
269 labelled as “Ma sp.4”.

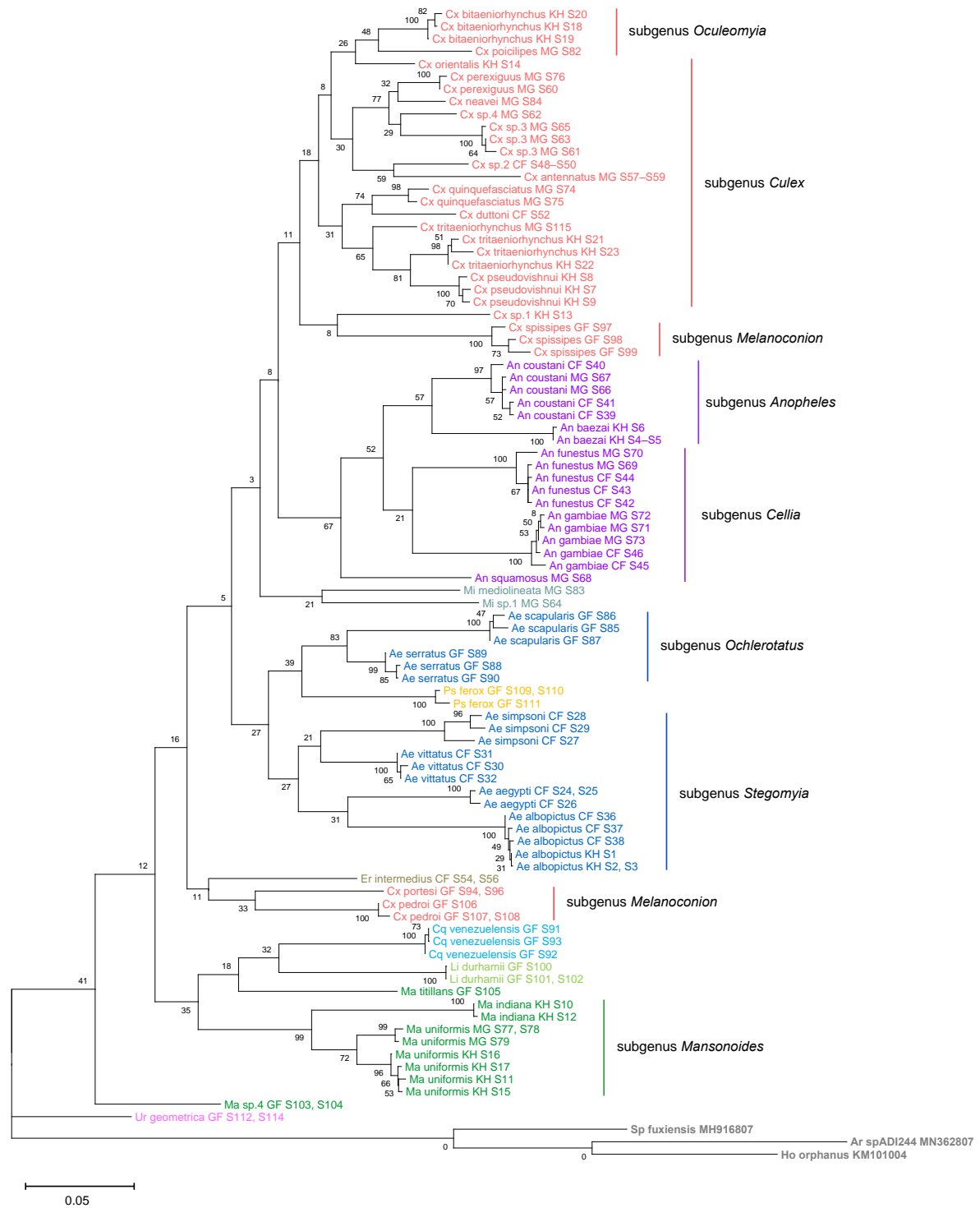
270 Phylogenetic reconstruction based on the COI sequences showed clustering of all species taxa
271 into distinct clades, underlining the utility of the COI gene in molecular taxonomy (Figure 5). However,
272 species delineation among members of *Culex* species groups were not as clear cut, although sister

273 species were correctly placed as sister taxa. This is comparable to the 28S+18S tree (Figure 4) and is
274 indicative of lower intraspecies distances relative to interspecies distances.

275 To evaluate the utility of 28S and 18S rRNA sequences for taxonomy-based species identification,
276 we used 28S+18S rRNA phylogenetic inference to discern the identity of six specimens for which COI
277 barcoding could not be performed. These specimens include three unknown *Mansonia* species
278 (Specimen ID 33–35), a *Ma. uniformis* (Specimen ID 51), an *An. gambiae* (Specimen ID 47), and a
279 *Ur. geometrica* (Specimen ID 113). Their positions in the 28S rRNA tree relative to adjacent taxa
280 confirms the morphological identification of all six specimens to the genus level and, for three of them,
281 to the species level (Figure 4).

282 The phylogenetic relationships indicated by the COI tree compared to the 28S+18S tree present
283 only a few points of congruence. COI-based phylogenetic inference indeed showed clustering of
284 generic taxa into clades albeit with very weak bootstrap support, except for genera *Culex* and
285 *Mansonia* (Figure 5). Contrary to the 28S+18S tree (Figure 4), *Culex* subgenus *Melanoconion* was
286 depicted as a polyphyletic taxon with *Cx. spissipes* being a part of the greater *Culicini* clade with
287 members from subgenera *Oculeomyia* and *Culex* while *Cx. pedroi* and *Cx. portesi* formed a distantly
288 related clade. Among the *Mansonia* specimens, the two unknown *Ma* sp.4 specimens were not
289 positioned as the nearest neighbours of *Ma. titillans* and instead appeared to have diverged earlier
290 from most of the other taxa from the *Culicidae* family. Notably, the COI sequences of genus
291 *Anopheles* is not basal to the other members of *Culicidae* and is instead shown to be sister to *Culex*
292 COI sequences (8% bootstrap support). This is a direct contrast to what is suggested by the rRNA
293 phylogenies (Figures 2–4), which suggests *Culex* rRNA sequences to be among the most recently
294 diverged. Bootstrap values for the more internal nodes of the COI trees are remarkably low compared
295 to those of rRNA-based trees.

296 In all rRNA trees, it is clear that the interspecific and intersubgeneric evolutionary distances within
297 the genus *Anopheles* are high relative to any other genera, indicating a greater degree of divergence.
298 This is evidenced by the longer branch lengths connecting Anopheline species-clades to the node of
299 the most recent common ancestor for subgenera *Anopheles* and *Cellia* (Figures 2-4, Supplementary
300 Figure 3). This feature is not evident in the COI tree, where the Anopheline interspecies distances are
301 comparable to those within the *Culex*, *Aedes*, and *Mansonia* taxa.



302

303 **Figure 5.** Phylogenetic tree based on COI sequences (621–699 bp) as inferred using maximum-
 304 likelihood method and constructed to scale in MEGA X (46). Values at each node indicate bootstrap
 305 support (%). Each specimen label contains information on its taxonomy, origin
 306 (as indicated in 2-letter country codes), and specimen ID. Label colours indicate genera: *Culex* in
 307 coral, *Anopheles* in purple, *Aedes* in dark blue, *Mansonia* in dark green, *Limatus* in light green,

308 *Coquillettidia* in light blue, *Psorophora* in yellow, *Mimomyia* in teal, *Uranotaenia* in pink and
309 *Eretmapodites* in brown. Scale bar at 0.05 is shown.

310

311 **On *Culex* species groups**

312 *Culex* (subgenus *Culex*) specimens of this study comprise several closely related sister species
313 belonging to the *Cx. vishnui* and *Cx. univittatus* species groups, which are notoriously difficult to
314 differentiate based on morphology. Accordingly, in the 28S+18S rRNA (Figure 4) and COI (Figure 5)
315 trees these species and their known sister species were clustered together within the *Culex*
316 (subgenus *Culex*) clade: *Cx. tritaeniorhynchus* with *Cx. pseudovishnui* (*Cx. vishnui* species group);
317 *Cx. perexiguus* with *Cx. neavei* (*Cx. univittatus* species group).

318 The use of COI barcoding to discern between members of the *Culex* species groups was limited.
319 For example, for the two *Cx. quinquefasciatus* samples in our dataset (Specimen ID 74 and 75),
320 BLAST analyses of their COI sequences revealed they are a single nucleotide away from *Cx. pipiens*
321 or *Cx. quinquefasciatus* COI sequences (Table S2). In the 28S rRNA tree with NCBI sequences
322 (Figure 2), two NCBI sequences of *Cx. pipiens* sequences formed a clade sister to another containing
323 three *Cx. quinquefasciatus* NCBI sequences and the “*Cx quinquefasciatus* MG S74” sequence with
324 78% bootstrap support. This is in accordance with other studies examining mitochondrial sequences
325 (49) and morphological attributes (50). This shows that the 28S rRNA sequence can distinguish the
326 two species and confirms that “*Cx quinquefasciatus* MG S74” is indeed a *Cx. quinquefasciatus*
327 specimen. However, “*Cx quinquefasciatus* MG S75” is shown to be basal from other sequences within
328 this *Cx. pipiens* species group-clade with 100% bootstrap support. Given that *Cx. quinquefasciatus*
329 and *Cx. pipiens* are known to interbreed, it is plausible that this individual is a hybrid of the two
330 species (51).

331

332 **DISCUSSION**

333 In metagenomics or surveillance studies on field-captured mosquitoes, the lack of reference rRNA
334 sequences hinders good oligo-based depletion and efficient clean-up of RNA-seq data. Additionally,
335 *de novo* assembly of rRNA sequences is complicated due to the high sequence conservation across
336 all distantly related organisms that could be present in a single specimen, i.e., microbiota, parasites,

337 or vertebrate blood meal. Hence, we sought out to establish a method to bioinformatically filter out
338 non-host rRNA reads for the accurate assembly of novel 28S and 18S rRNA reference sequences.

339 We found that phylogenetic reconstructions based on 28S sequences or concatenated 28S+18S
340 rRNA sequences were able to cluster mosquito taxa according to species correctly and corroborates
341 current mosquito classification. This demonstrates that our bioinformatics methodology reliably
342 generates bona fide 28S and 18S rRNA sequences, even in specimens parasitized by water mites or
343 engorged with vertebrate blood. Further, we were able to use 28S+18S rRNA taxonomy for molecular
344 species identification when COI sequences were unavailable or ambiguous, thus supporting the use
345 of 28S rRNA sequences as an rRNA barcode. rRNA barcodes would have the advantage of
346 circumventing the need to additionally isolate and sequence DNA from specimens, as RNA-seq reads
347 can be directly mapped against rRNA reference sequences. Even after depletion, there are sufficient
348 numbers of reads (5-10% of reads per sample) to assemble complete rRNA contigs (Frangeul L,
349 personal communication).

350 Phylogenetic inferences based on 28S or 18S rRNA sequences do not produce the same
351 interspecific relationships, suggesting a difference in mutation rates between the two gene regions.
352 Relative to 28S sequences, we observed more instances where multiple specimens have near-
353 identical 18S rRNA sequences. This can occur for specimens belonging to the same species, but also
354 for conspecifics sampled from different geographic locations, such as *An. coustani*, *An. gambiae*, or
355 *Ae. albopictus*. More rarely, specimens from the same species group, such as *Cx. pseudovishnui* and
356 *Cx. tritaeniorhynchus*, were also found to share 18S rRNA sequences. This was surprising given that
357 the 18S rRNA sequences in our dataset is 1900 bp long. Concatenation of 28S and 18S rRNA
358 sequences resolved this issue, enabling species delineation even among members of *Culex* species
359 groups.

360 Taking advantage of our multi-country sampling, we examined whether rRNA or COI phylogeny
361 can be used to discriminate conspecifics originating from different countries. Our dataset contains five
362 of such species: *An. coustani*, *An. funestus*, *An. gambiae*, *Ae. albopictus*, and *Ma. uniformis*. Among
363 the rRNA trees, the concatenated 28S+18S tree and 28S tree were able to discriminate between *Ma.*
364 *uniformis* specimens from Madagascar, Cambodia, and the Central African Republic, and between
365 *An. coustani* specimens from Madagascar and the Central African Republic (100% bootstrap value).
366 In the COI tree, only *Ma. uniformis* was resolved into geographical clades comprising specimens from

367 Madagascar and specimens from Cambodia (72% bootstrap value). No COI sequence was obtained
368 from one *Ma. uniformis* from the Central African Republic. The use of rRNA sequences seemingly
369 provides more accurate phylogeographic information than COI alone.

370 Morphological identification suffers in accuracy when dealing with *Culex* species groups. Aside
371 from sharing many morphological traits, sister species are often sympatric and show at least some
372 competence for a number of viral and filarial pathogens, such as Japanese encephalitis virus, St
373 Louis encephalitic virus, Usutu virus, and *Wuchereria bancrofti* (52). However, each of these species
374 have distinct ecologies and host preferences, thus the challenge of correctly identifying vector species
375 can affect epidemiological risk estimation for these diseases (51). In Asia, for example, cryptic
376 members of the *Cx. vishnui* species group confound tracking of Japanese encephalitis virus
377 transmission (53). The morphological differences between the *Culex* species *bitaeniorhynchus*,
378 *tritaeniorhynchus*, *vishnui*, and *pseudovishnui* are often elusive, the former three having been
379 morphologically identified in our study but later revealed by COI barcoding to be another species.

380 The *Cx. pipiens* species group is especially challenging as its member species are capable of
381 interbreeding, showing genetic introgression to varying extents depending on the geographical
382 population (54). The seven member species of this complex are practically indistinguishable
383 morphologically and require molecular methods to discern (51,55). However, the 621 bp COI
384 sequence amplified in our study did not contain enough nucleotide divergence to allow clear
385 identification, given that the COI sequence of *Cx. quinquefasciatus* specimens differed from that of
386 *Cx. pipiens* by a single nucleotide. Similarly, Batovska *et al* (56) found that even the Internal
387 Transcribed Spacer 2 (ITS2) rDNA region, another common molecular marker, could not differentiate
388 the two species. Other DNA molecular markers such as nuclear *Ace-2* or *CQ11* genes (55,57) or
389 *Wolbachia pipientis* infection status (54) are typically employed in tandem. In our study, we used on
390 28S rRNA sequence-based taxonomy (Figure 2) to validate the identity of specimen “Cx
391 quinquefasciatus MG S74” and suggests that specimen “Cx quinquefasciatus MG S75” might have
392 been a *pipiens-quinquefasciatus* hybrid. These examples demonstrate how 28S rRNA sequences,
393 concatenated with 18S or alone, contain enough resolution to differentiate between *Cx. pipiens* and
394 *Cx. quinquefasciatus*. rRNA barcode taxonomy thus allows for more accurate species identification
395 and ecological observations in the context of disease transmission. Additionally, tracing the genetic
396 flow across hybrid populations within the *Cx. pipiens* species group can inform estimates of vectorial

397 capacity for each species. Only one or two members from the *Cx. pipiens* and *Cx. vishnui* species
398 groups were represented in our dataset. An explicit investigation including all member species of
399 these species groups in greater sample numbers is warranted to further test the degree of accuracy
400 with which 28S and 18S rRNA sequences can delineate sister species.

401 Our study also included French Guianese *Culex* species *Cx. spissipes* (group Spissipes), *Cx.*
402 *pedroi* (group Pedroi), and *Cx. portesi* (group Vomerifer). These species belong to the New World
403 subgenus *Melanoconion*, section Spissipes, with well-documented distribution in North and South
404 Americas (58) and are vectors of encephalitic alphaviruses EEEV and VEEV among others (37–39).
405 Indeed, our rooted rRNA and COI trees show the divergence of the three *Melanoconion* species from
406 the major *Culex* clade comprising species broadly found across Africa and Asia (23,51,52,59). The
407 topology of the concatenated 28S+18S tree puts the *Cx. portesi* and *Cx. pedroi* species-clades as
408 sister groups (92% bootstrap support), with *Cx. spissipes* as a basal group within the *Melanoconion*
409 clade (100% bootstrap support). This corroborates the systematics elucidated by Navarro and
410 Weaver (60) using the ITS2 marker, and those by Sirivanakarn (58) and Sallum and Forattini (61)
411 based on morphological features. Curiously, in the COI tree *Cx. spissipes* sequences were clustered
412 with unknown species *Cx. sp1*, forming a clade sister to one containing other *Culex* (*Culex*) and *Culex*
413 (*Oculeomyia*) species, albeit with very low bootstrap support. Previous phylogenetic studies based on
414 the COI gene have consistently placed *Cx. spissipes* or the Spissipes group basal to other groups
415 within the *Melanoconion* subgenus (62,63). However, these studies contain only *Culex*
416 (*Melanoconion*) species in their dataset, apart from *Cx. quinquefasciatus* to act as an outgroup. This
417 clustering of *Cx. spissipes* with non-*Melanoconion* species in our COI phylogeny could be an artefact
418 of a much more diversified dataset rather than a true phylogenetic link.

419 The evolutionary histories inferred from rRNA-based and COI-based phylogenies in our study
420 hardly correspond. rRNA phylogenies suggest the world of *Anopheles* is seemingly immense
421 compared to any other genera with remarkably large evolutionary distances between one Anopheline
422 species to another. This is not apparent in the COI phylogeny, perhaps reflecting the higher
423 mutational rate of mitochondrial genomes relative to nuclear genomes (64). It would be interesting to
424 further compare rRNA and COI phylogenies among other Anopheline subgenera beyond the
425 subgenera *Anopheles* and *Cellia* represented in this study. Lamentably, we found during our search in

426 NCBI databases that many Anopheline rRNA records lack subgenus information, stressing the
427 importance of including detailed taxonomy of mosquito specimens when reporting sequence data.

428 In contrast to the *Anopheles* case, two specimens of an unknown *Mansonia* species, “Ma sp.4 GF
429 S103” and “Ma sp.4 GF S104”, provided an example where interspecies relatedness based on their
430 COI sequences is greater than that based on their rRNA sequences in relation to “Ma titillans GF
431 S105”. While all rRNA trees (Figure 2–4) placed “Ma titillans GF S105” as a sister taxon with 100%
432 bootstrap support, the COI tree places Ma sp.4 basal to all other species except *Ur. geometrica*. This
433 may hint at a historical selective sweep in mitochondrial genome, whether arising from mutations or
434 linkage disequilibrium with inherited symbionts (65), resulting in the drastically distinct mitochondrial
435 haplogroup found in Ma sp.4. To note, the COI sequences of “Ma sp.4 GF S103” and “Ma sp.4 GF
436 S104” share 87.12 and 87.39% nucleotide similarity, respectively, with that of “Ma titillans GF S105”.
437 Interestingly, the endosymbiont *Wolbachia pipientis* has been detected in *Ma. titillans* sampled from
438 Brazil (66), which may contribute to the divergence of “Ma titillans GF S105” COI sequence away from
439 those of other *Mansonia* species. The COI phylogeny of these *Mansonia* specimens highlights the
440 drawbacks of using a mitochondrial DNA marker in determining evolutionary relationships (65), 28S
441 and 18S rRNA sequences may be better able to illustrate evolutionary history than COI sequence
442 alone.

443

444 **Conclusions**

445 Surveillance and microbiome discovery studies in wild mosquitoes are paramount for the
446 establishment of public health measures to control arboviral diseases. Here we present a score-based
447 rRNA assembly strategy and 234 newly generated 28S and 18S mosquito rRNA sequences. Our
448 work has expanded the current rRNA reference library by presenting, to our knowledge, the first
449 records for many species not previously present in public databases and paves the way for the
450 assembly of many more. These novel rRNA sequences can improve mosquito RNA-seq
451 metagenomics by expanding reference sequence data for the optimization of species-specific oligo-
452 based depletion protocols, for streamlined species identification by rRNA barcoding and for improved
453 RNA-seq data clean-up. In addition, rRNA barcodes could serve as an additional tool for mosquito
454 taxonomy and phylogeny although further studies are necessary to reveal how they measure up
455 against other nuclear or mitochondrial DNA marker systems (9,56,67–69).

456 We showed that phylogenetic inferences from a tree based on 28S rRNA sequences alone or
457 concatenated 28S +18S rRNA sequences largely agree with contemporary mosquito classification
458 and can be used for species identification given a reference sequence. In analysing the same set of
459 specimen by COI or rRNA sequences, we found deep discrepancies in phylogenetic inferences. We
460 conclude that while COI-based phylogeny is fairly useful to study recent evolutionary events, rRNA
461 sequences may be better suited for investigations of more ancient evolutionary history.

462

463 **METHODS**

464

465 **Sample collection**

466 Mosquito specimens were sampled from 2019 to 2020 by medical entomology teams from the Institut
467 Pasteur International Network: Institut Pasteur de Bangui (Central African Republic, Africa; CF),
468 Institut Pasteur de Madagascar (Madagascar, Africa; MG), Institut Pasteur du Cambodge (Cambodia,
469 South East Asia; KH), and Institut Pasteur de la Guyane (French Guiana, South America; GF). Adult
470 mosquitoes were sampled using a combination of techniques including CDC light traps, BG sentinels,
471 and human-landing catches. Sampling sites are non-urban locations including rural settlements in the
472 Central African Republic, Madagascar, and French Guiana and national parks in Cambodia.
473 Mosquitoes were identified using morphological identification keys on cold tables before preservation
474 by flash freezing in liquid nitrogen and transportation in dry ice to Institut Pasteur Paris for analysis.
475 The full list of the 112 mosquito specimens used in this study and their related information are
476 provided in Supplementary Table 1.

477

478 **RNA and DNA isolation**

479 Nucleic acids were isolated from mosquito specimens using TRIzol reagent according to
480 manufacturer's protocol (Invitrogen, Thermo Fisher Scientific, Waltham, Massachusetts, USA, Cat.
481 No. 15596018). Single mosquitoes were homogenised into 200 μ L of TRIzol reagent and other of the
482 reagents within the protocol were volume-adjusted accordingly. Following phase separation, RNA
483 samples were isolated from the aqueous phase while DNA samples were isolated from the remaining
484 interphase and phenol-chloroform phase. From here, RNA is used to prepare cDNA libraries for next

485 generation sequencing while DNA is used for PCR amplification and Sanger sequencing of the
486 mitochondrial *cytochrome c oxidase subunit I* (COI) gene.

487

488 **Probe depletion of rRNA**

489 We tested a method for selective depletion of rRNA by Morlan *et al.* (17) on several mosquito species
490 from the *Aedes*, *Culex*, and *Anopheles* genera. We designed 77 tiled 80 bp DNA probes antisense to
491 the *Ae. aegypti* 28S, 18S, and 5.8S rRNA sequences. A pool of probes at a concentration of 0.04 μM
492 were prepared. To bind probes to rRNA, 1 μL of probes was added to rRNA samples along with 2 μL
493 of Hybridisation Buffer (100 mM Tris-HCl and 200 mM NaCl) to obtain a final volume of 20 μL and
494 subjected to a slow-cool incubation starting at 95 $^{\circ}\text{C}$ for 2 minutes, followed by cooling to 22 $^{\circ}\text{C}$ at a
495 rate of 0.1 $^{\circ}\text{C}$ per second, followed by an additional 5 minutes at 22 $^{\circ}\text{C}$. The resulting RNA:DNA
496 hybrids were treated with 2.5 μL Hybridase™ Thermostable RNase H (Epicentre, Illumina, Madison,
497 Wisconsin, USA, Cat No. H39500) and incubated at 37 $^{\circ}\text{C}$ for 30 minutes. To remove DNA probes,
498 the mix was treated with 1 μL DNase I (Invitrogen, Cat No. 18047019) and purified with Agencourt
499 RNAClean XP Beads (Beckman Coulter, Brea, California, USA, Cat No. A63987). The resulting RNA
500 is used for total RNA sequencing to check depletion efficiency.

501

502 **Total RNA sequencing**

503 To obtain rRNA sequences, RNA samples were quantified on a Qubit Fluorometer (Invitrogen) using
504 the Qubit RNA BR Assay kit (Invitrogen, Cat No. Q10211) for concentration adjustment. Non-depleted
505 total RNA was used for library preparation for next generation sequencing using the NEBNext Ultra II
506 RNA Library Preparation Kit for Illumina (New England Biolabs, Ipswich, Massachusetts, USA, Cat.
507 No. E7770L) and the NEBNext Multiplex Oligos for Illumina (Dual Index Primers Set 1) (New England
508 Biolabs, Cat. No. E7600S). Sequencing was performed on a NextSeq500 sequencing system
509 (Illumina, San Diego, California, USA). Quality control of fastq data and trimming of adapters were
510 performed with FastQC and cutadapt, respectively.

511

512 **28S and 18S rRNA assembly**

513 To obtain 28S and 18S rRNA contigs, we had to first clean our fastq library by separating the reads
514 representing mosquito rRNA from all other reads. To achieve this, we used the SILVA RNA sequence

515 database to create 2 libraries: one containing all rRNA sequences recorded under the "Insecta" node
516 of the taxonomic tree, the other containing the rRNA sequences of many others nodes distributed
517 throughout the taxonomic tree, hence named "Non-Insecta" (70). Each read was aligned using the
518 nucleotide Basic Local Alignment Search Tool (BLASTn, <https://blast.ncbi.nlm.nih.gov/>) of the
519 National Center for Biotechnology Information (NCBI) against each of the two libraries and the scores
520 of the best high-scoring pairs from the two BLASTns are subsequently used to calculate a ratio of
521 Insecta over Non-Insecta scores (71). Only reads with a ratio greater than 0.8 were used in the
522 assembly. The two libraries being non-exhaustive, we chose this threshold of 0.8 to eliminate only
523 reads that were clearly of a non-insect origin. Selected reads were assembled with the SPAdes
524 assembler using the "-rna" option, allowing more heterogeneous coverage of contigs and kmer
525 lengths of 31 , 51 and 71 bases (72). This method successfully assembled rRNA sequences for all
526 specimens, including a parasitic *Horreolanus* water mite (122 sequences for 28S and 114 sequences
527 for 18S).

528 Initially, our filtration technique had two weaknesses. First, there is a relatively small number of
529 complete rRNA sequences in the Insecta library from SILVA. To compensate for this, we carried out
530 several filtration cycles and in between cycles, added all the complete sequences produced in
531 previous cycles to the Insecta library. Second, when our mosquito specimens were parasitized by
532 other insects, it was not possible to bioinformatically filter out rRNA reads belonging to the parasite.
533 For these rare cases, we used the "--trusted-contigs" option of SPAdes, giving it access to the 28S
534 and 18S sequences of the mosquito closest in terms of taxonomy to the one we were assembling. By
535 doing this, the assembler was able to reconstruct the rRNA of the mosquito as well as the rRNA of the
536 parasitizing insect. All assembled rRNA sequences from this study have been deposited in GenBank
537 with accession numbers OM350214–OM350327 for 18S rRNA sequences and OM542339–
538 OM542460 for 28S rRNA sequences.

539

540 **COI amplicon sequencing**

541 The mitochondrial COI gene was amplified from DNA samples using the universal "Folmer" primer set
542 LCO1490 (5'-GGTCAACAAATCATAAAGATATTGG-3') and HCO2198 (5'-
543 TAAACTTCAGGGTGACCAAAAAATCA-3'), as is the standard for COI barcoding, producing a 658 bp
544 product (73). PCRs were performed using Phusion High-Fidelity DNA Polymerase (Thermo Fisher

545 Scientific, Cat. No. F530L). Every 50 μ L reaction contained 10 μ L of 5X High Fidelity buffer, 1 μ L of 10
546 mM dNTPs, 2.5 μ L each of 10 mM forward (LCO1490) and reverse (HCO2198) primer, 28.5 μ L of
547 water, 5 μ L of DNA sample, and 0.5 μ L of 2 U/ μ L Phusion DNA polymerase. A 3-step cycling
548 incubation protocol was used: 98 °C for 30 seconds; 35 cycles of 98 °C for 10 seconds, 60 °C for 30
549 seconds, and 72 °C for 15 seconds; 72 °C for 5 minutes followed by a 4 °C hold. PCR products were
550 size-verified using gel electrophoresis and then gel-purified using the QIAquick Gel Extraction Kit
551 (Qiagen, Hilden, Germany, Cat. No. 28706). Sanger sequencing of the COI amplicons were
552 performed by Eurofins Genomics, Ebersberg, Germany.

553

554 **COI sequence analysis**

555 Forward and reverse COI DNA sequences were end-trimmed to remove bases of poor quality (Q
556 score < 30). At the 5' ends, sequences were trimmed at the same positions such that all forward
557 sequences start with 5'- TTTTGG and all reverse sequences start with 5'- GGNTCT. Forward and
558 reverse sequences were aligned using BLAST to produce a 621 bp consensus sequence, the
559 minimum length of COI sequence obtained for all specimens. In cases where good quality sequences
560 extends beyond 621 bp, forward and reverse sequences were assembled using Pearl
561 (<https://www.gear-genomics.com/pearl/>) and manually checked for errors against trace files
562 (74,75). We successfully assembled a total of 106 COI sequences. All assembled COI sequences
563 from this study have been deposited in GenBank with accession numbers OM630610–OM630715.

564

565 **COI validation of morphology-based species identification**

566 We analysed assembled COI sequences with BLASTn against the nucleotide collection (nr/nt)
567 database to confirm morphology-based species identification. BLAST analyses revealed 32 cases
568 where top hits indicated a different species identity, taking <95% nucleotide sequence similarity as the
569 threshold to delineate distinct species (Supplementary Table 1). In these cases, the COI sequence of
570 the specimen was then BLAST-aligned against a GenBank record representing the morphological
571 species to verify that the revised identity is a closer match by a significant margin, i.e., more than 2%
572 nucleotide sequence similarity. All species names reported hereafter reflect identities determined by
573 COI barcoding except for cases where COI-based identities were ambiguous, in which case
574 morphology-based identities were retained. In cases where matches were found within a single genus

575 but of multiple species, specimens were indicated as an unknown member of their genus (e.g., *Culex*
576 sp.). Information of the highest-scoring references for all specimens, including details of ambiguous
577 BLASTn results, are recorded in Supplementary Table 2.

578 Within our COI sequences, we found six unidentified *Culex* species (including two that matched to
579 NCBI entries identified only to the genus level), four unidentified *Mansonia* species, and one
580 unidentified *Mimomyia* species. For *An. baezai*, no existing NCBI records were found at the time the
581 analysis was performed.

582

583 **Phylogenetic analysis**

584 Multiple sequence alignment (MSA) were performed on assembled COI and rRNA sequences using
585 the MUSCLE software (Additional Files 6–9) (76,77). As shown in Supplementary Figure 2 on the
586 conservation of identity along the alignment, the 28S sequences contain many blocks of highly
587 conserved nucleotides throughout the sequence, which makes the result of multiple alignment
588 particularly obvious. We therefore did not test other alignment programs. The multiple alignment of
589 the COI amplicon is even more evident since no gaps are necessary for this alignment.

590 Phylogenetic tree reconstructions were performed with the MEGA X software using the maximum-
591 likelihood method (46). The default parameters were used with the addition of bootstrapping with 500
592 replications in order to be able to quantify the level of confidence in the branches of the trees
593 obtained. For the 28S and 18S rRNA trees, two sequences belonging to an unknown species of
594 parasitic mite from the genus *Horreolanus* found in our specimens were included to serve as an
595 outgroup taxon. In addition, we created and analysed a separate dataset combining our 28S rRNA
596 sequences and full-length 28S rRNA sequences from the NCBI databases totalling 169 sequences
597 from 58 species (12 subgenera). To serve as outgroups for the COI tree, we included sequences
598 obtained from NCBI of three water mite species, *Horreolanus orphanus* (KM101004), *Sperchon*
599 *fluxiensis* (MH916807), and *Arrenurus* sp. (MN362807).

600

601 **DECLARATIONS**

602

603 **Ethics approval and consent to participate**

604 Not applicable

605

606 **Consent for publication**

607 Not applicable

608

609 **Availability of data and materials**

610 RNA-seq fastq sequence data are available from the corresponding author on reasonable request.

611 Multiple sequence alignment files are included in this article as additional files. All sequences

612 generated in this study have been deposited in GenBank under the accession numbers OM350214-

613 OM350327 for 18S rRNA sequences, OM542339-OM542460 for 28S rRNA sequences, and

614 OM630610-OM630715 for COI sequences (Table S1).

615

616 **Competing interests**

617 The authors declare that they have no competing interests.

618

619 **Funding**

620 This work was supported by the Defence Advanced Research Projects Agency PREEMPT program

621 managed by Dr. Rohit Chitale and Dr. Kerri Dugan [Cooperative Agreement HR001118S0017] (the

622 content of the information does not necessarily reflect the position or the policy of the U.S.

623 government, and no official endorsement should be inferred).

624

625 **Acknowledgements**

626 We thank members of the Saleh lab for valuable discussions and to Inès Partouche for laboratory

627 assistance. We especially thank all medical entomology staff of IP Bangui, IP Cambodge (Sony Yean,

628 Kimly Heng, Kalyan Chhuoy, Sreynik Nhek, Moeun Chhum, Kimhuor Sour and Pierre-Olivier

629 Maquart), IP Madagascar, and IP Guyane for assistance in field missions, laboratory work, and

630 logistics. We are also grateful to Dr Catherine Dauga for advice on phylogenetic analyses and to

631 Amandine Guidez for providing a French Guiana-specific COI reference library.

632

633 **ADDITIONAL FILES**

634 **Additional File 1: Figure S1.** Study workflow from specimens to sequences (PNG).

635 **Additional File 2: Table S1.** Taxonomic and sampling information on mosquito specimens and
636 associated accession numbers of their 28S, 18S, and COI sequences (XLSX).

637 **Additional File 3: Table S2.** COI sequence BLAST analyses summary (XLSX).

638 **Additional File 4: Figure S2.** Sequence conservation among 169 28S rRNA sequences obtained
639 from this study and from the NCBI databases (PDF).

640 **Additional File 5: Figure S3.** Phylogenetic tree based on 28S sequences generated from this study
641 (3900 bp) as inferred using maximum-likelihood method and constructed to scale in MEGA X (46).
642 Values at each node indicate bootstrap support (%) from 500 replication. Each specimen label
643 contains information on its taxonomy, origin (as indicated in 2-letter country codes), and specimen ID.
644 Label colours indicate genera: *Culex* in coral, *Anopheles* in purple, *Aedes* in dark blue, *Mansonia* in
645 dark green, *Limatus* in light green, *Coquillettidia* in light blue, *Psorophora* in yellow, *Mimomyia* in teal,
646 *Uranotaenia* in pink and *Eretmapodites* in brown. Scale bar at 0.05 is shown (PDF).

647 **Additional File 6:** Multiple sequence alignment of 169 28S rRNA sequences from this study and from
648 NCBI databases (FASTA).

649 **Additional File 7:** Multiple sequence alignment of 122 28S rRNA sequences, including two
650 sequences from *Horreolanus sp.* (FASTA).

651 **Additional File 8:** Multiple sequence alignment of 114 18S rRNA sequences, including two
652 sequences from *Horreolanus sp.* (FASTA).

653 **Additional File 9:** Multiple sequence alignment of 106 COI sequences (FASTA).

654

655 REFERENCES

- 656 1. WHO. Global vector control response 2017–2030. World Health Organization. 2017.
- 657 2. Webster JP, Gower CM, Knowles SCL, Molyneux DH, Fenton A. One health - an ecological
658 and evolutionary framework for tackling Neglected Zoonotic Diseases. *Evol Appl.* 2016 Feb
659 1;9(2):313–33.
- 660 3. Gould E, Pettersson J, Higgs S, Charrel R, de Lamballerie X. Emerging arboviruses: Why
661 today? *One Heal.* 2017 Dec 1;4:1–13.
- 662 4. GALE K, CRAMPTON J. The ribosomal genes of the mosquito, *Aedes aegypti*. *Eur J Biochem.*
663 1989;185(2):311–7.
- 664 5. Phelps WA, Carlson AE, Lee MT. Optimized design of antisense oligomers for targeted rRNA

- 665 depletion. *Nucleic Acids Res.* 2021;49(1):e5.
- 666 6. Fauver JR, Akter S, Morales AIO, Black WC, Rodriguez AD, Stenglein MD, et al. A reverse-
667 transcription/RNase H based protocol for depletion of mosquito ribosomal RNA facilitates viral
668 intrahost evolution analysis, transcriptomics and pathogen discovery. *Virology.* 2019;528:181–
669 97.
- 670 7. Kukutla P, Steritz M, Xu J. Depletion of ribosomal RNA for mosquito gut metagenomic RNA-
671 seq. *J Vis Exp.* 2013;(74):50093.
- 672 8. Hebert PDN, Cywinska A, Ball SL, DeWaard JR. Biological identifications through DNA
673 barcodes. *Proc R Soc B Biol Sci.* 2003;270(1512):313–21.
- 674 9. Ratnasingham S, Hebert PDN. BOLD: The Barcode of Life Data System: Barcoding. *Mol Ecol*
675 *Notes.* 2007;7(3):355–64.
- 676 10. Bishop-Lilly KA, Turell MJ, Willner KM, Butani A, Nolan NME, Lentz SM, et al. Arbovirus
677 detection in insect vectors by Rapid, high- throughput pyrosequencing. *PLoS Negl Trop Dis.*
678 2010;4(11):e878.
- 679 11. Kumar N, Creasy T, Sun Y, Flowers M, Tallon LJ, Dunning Hotopp JC. Efficient subtraction of
680 insect rRNA prior to transcriptome analysis of *Wolbachia*-*Drosophila* lateral gene transfer.
681 *BMC Res Notes.* 2012;5:230.
- 682 12. Weedall GD, Irving H, Hughes MA, Wondji CS. Molecular tools for studying the major malaria
683 vector *Anopheles funestus*: Improving the utility of the genome using a comparative poly(A)
684 and Ribo-Zero RNAseq analysis. *BMC Genomics.* 2015;16(1):931.
- 685 13. Zakrzewski M, Rašić G, Darbro J, Krause L, Poo YS, Filipović I, et al. Mapping the virome in
686 wild-caught *Aedes aegypti* from Cairns and Bangkok. *Sci Rep.* 2018;8(1):4690.
- 687 14. Belda E, Nanfack-Minkeu F, Eiglmeier K, Carissimo G, Holm I, Diallo M, et al. De novo
688 profiling of RNA viruses in *Anopheles malaria* vector mosquitoes from forest ecological zones
689 in Senegal and Cambodia. *BMC Genomics.* 2019;20(1):664.
- 690 15. Chandler JA, Liu RM, Bennett SN. RNA Shotgun Metagenomic Sequencing of Northern
691 California (USA) Mosquitoes Uncovers Viruses, Bacteria, and Fungi. *Front Microbiol.*
692 2015;6:185.
- 693 16. Thongsripong P, Chandler JA, Kittayapong P, Wilcox BA, Kapan DD, Bennett SN.
694 Metagenomic shotgun sequencing reveals host species as an important driver of virome

- 695 composition in mosquitoes. *Sci Rep.* 2021;11(1):8448.
- 696 17. Morlan JD, Qu K, Sinicropi D V. Selective depletion of rRNA enables whole transcriptome
697 profiling of archival fixed tissue. *PLoS One.* 2012;7(8):e42882.
- 698 18. Diallo D, Fall G, Diagne CT, Gaye A, Ba Y, Dia I, et al. Concurrent amplification of Zika,
699 chikungunya, and yellow fever virus in a sylvatic focus of arboviruses in Southeastern
700 Senegal, 2015. *BMC Microbiol.* 2020;20:181.
- 701 19. Vasconcelos PFC, Costa ZG, Travassos da Rosa ES, Luna E, Rodrigues SG, Barros VLRS, et
702 al. Epidemic of jungle yellow fever in Brazil, 2000: Implications of climatic alterations in disease
703 spread. *J Med Virol.* 2001;65(3):598–604.
- 704 20. Cardoso J da C, de Almeida MAB, dos Santos E, da Fonseca DF, Sallum MAM, Noll CA, et al.
705 Yellow fever virus in *Haemagogus leucocelaenus* and *Aedes serratus* mosquitoes, Southern
706 Brazil, 2008. *Emerg Infect Dis.* 2010;16(12):1918–24.
- 707 21. Gaillet M, Pichard C, Restrepo J, Lavergne A, Perez L, Enfissi A, et al. Outbreak of oropouche
708 virus in french guiana. *Emerg Infect Dis.* 2021;27(10):2711–4.
- 709 22. Kraemer MUG, Reiner RC, Brady OJ, Messina JP, Gilbert M, Pigott DM, et al. Past and future
710 spread of the arbovirus vectors *Aedes aegypti* and *Aedes albopictus*. *Nat Microbiol.*
711 2019;4(5):854–63.
- 712 23. Auerswald H, Maquart PO, Chevalier V, Boyer S. Mosquito vector competence for japanese
713 encephalitis virus. Vol. 13, *Viruses.* 2021. p. 1154.
- 714 24. Mukwaya LG, Kayondo JK, Crabtree MB, Savage HM, Biggerstaff BJ, Miller BR. Genetic
715 differentiation in the yellow fever virus vector, *Aedes simpsoni* complex, in Africa: Sequence
716 variation in the ribosomal DNA internal transcribed spacers of anthropophilic and non-
717 anthropophilic populations. *Insect Mol Biol.* 2000;9(1):85–91.
- 718 25. Mwangangi JM, Muturi EJ, Muriu SM, Nzovu J, Midega JT, Mbogo C. The role of *Anopheles*
719 *arabiensis* and *Anopheles coustani* in indoor and outdoor malaria transmission in Taveta
720 District, Kenya. *Parasites and Vectors.* 2013;6:114.
- 721 26. Nepomichene TNJJ, Raharimalala FN, Andriamandimby SF, Ravalohery JP, Failloux AB,
722 Heraud JM, et al. Vector competence of *Culex antennatus* and *Anopheles coustani*
723 mosquitoes for Rift Valley fever virus in Madagascar. *Med Vet Entomol.* 2018;32(2):259–62.
- 724 27. Ratovonjato J, Olive MM, Tantely LM, Andrianaivolambo L, Tata E, Razainirina J, et al.

- 725 Detection, isolation, and genetic characterization of Rift Valley fever virus from anopheles
726 (Anopheles) coustani, anopheles (Anopheles) squamosus, and culex (Culex) antennatus of
727 the haute matsiatra region, Madagascar. *Vector-Borne Zoonotic Dis.* 2011;11(6):753–9.
- 728 28. Lutomiah J, Bast J, Clark J, Richardson J, Yalwala S, Oullo D, et al. Abundance, diversity, and
729 distribution of mosquito vectors in selected ecological regions of Kenya: Public health
730 implications. *J Vector Ecol.* 2013 Jun 1;38(1):134–42.
- 731 29. Brault AC, Foy BD, Myles KM, Kelly CLH, Higgs S, Weaver SC, et al. Infection patterns of
732 o'nyong nyong virus in the malaria-transmitting mosquito, *Anopheles gambiae*. *Insect Mol Biol.*
733 2004;13(6):625–35.
- 734 30. Stevenson JC, Simubali L, Mbambara S, Musonda M, Mweetwa S, Mudenda T, et al.
735 Detection of plasmodium falciparum infection in anopheles squamosus (diptera: Culicidae) in
736 an area targeted for malaria elimination, Southern Zambia. *J Med Entomol.* 2016;53(6):1482–
737 7.
- 738 31. Nikolay B, Diallo M, Boye CSB, Sall AA. Usutu virus in Africa. Vol. 11, *Vector-Borne and*
739 *Zoonotic Diseases.* 2011. p. 1417–23.
- 740 32. Kim H, Cha GW, Jeong YE, Lee WG, Chang KS, Roh JY, et al. Detection of Japanese
741 encephalitis virus genotype V in *Culex orientalis* and *Culex pipiens* (Diptera: Culicidae) in
742 Korea. *PLoS One.* 2015;10(2):e0116547.
- 743 33. Vázquez González A, Ruiz S, Herrero L, Moreno J, Molero F, Magallanes A, et al. West Nile
744 and Usutu viruses in mosquitoes in Spain, 2008-2009. *Am J Trop Med Hyg.* 2011;85(1):178–
745 81.
- 746 34. Maquart PO, Sokha C, Boyer S. Mosquito diversity (Diptera: Culicidae) and medical
747 importance, in a bird sanctuary inside the flooded forest of Prek Toal, Cambodia. *J Asia Pac*
748 *Entomol.* 2021;24(4):1221–7.
- 749 35. Ndiaye EH, Fall G, Gaye A, Bob NS, Talla C, Diagne CT, et al. Vector competence of *Aedes*
750 *vexans* (Meigen), *Culex poicilipes* (Theobald) and *Cx. quinquefasciatus* Say from Senegal for
751 West and East African lineages of Rift Valley fever virus. *Parasites and Vectors.* 2016;9:94.
- 752 36. Bhattacharya S, Basu P, Sajal Bhattacharya C. The Southern House Mosquito, *Culex*
753 *quinquefasciatus*: profile of a smart vector. *J Entomol Zool Stud JEZS.* 2016;73(42):73–81.
- 754 37. Weaver SC, Ferro C, Barrera R, Boshell J, Navarro JC. Venezuelan Equine Encephalitis.

- 755 Annu Rev Entomol. 2004;49:141–74.
- 756 38. Talaga S, Duchemin JB, Girod R, Dusfour I. The Culex Mosquitoes (Diptera: Culicidae) of
757 French Guiana: A Comprehensive Review With the Description of Three New Species. J Med
758 Entomol. 2021;58(1):182–221.
- 759 39. Turell MJ, O'guinn ML, Dohm D, Zyzak M, Watts D, Fernandez R, et al. Susceptibility of
760 Peruvian Mosquitoes to Eastern Equine Encephalitis Virus. J Med Entomol. 2008;45(4):720–5.
- 761 40. Barrio-Nuevo KM, Cunha MS, Luchs A, Fernandes A, Rocco IM, Mucci LF, et al. Detection of
762 Zika and dengue viruses in wild-caught mosquitoes collected during field surveillance in an
763 environmental protection area in São Paulo, Brazil. PLoS One. 2020;15(10):e0227239.
- 764 41. Turell MJ. Vector competence of three Venezuelan mosquitoes (Diptera: Culicidae) for an
765 epizootic IC strain of Venezuelan equine encephalitis virus. J Med Entomol. 1999;36(4):407–9.
- 766 42. Hoyos-López R, Soto SU, Rúa-Uribe G, Gallego-Gómez JC. Molecular identification of Saint
767 Louis encephalitis virus genotype IV in Colombia. Mem Inst Oswaldo Cruz [Internet]. 2015 Aug
768 21 [cited 2022 Apr 21];110(6):719–25. Available from:
769 <http://www.scielo.br/j/mioc/a/YbSNxJQWw6dFLx5d4TtwfG/?lang=en>
- 770 43. Arunachalam N, Philip Samuel P, Hiriyan J, Thenmozhi V, Gajanana A. Japanese encephalitis
771 in Kerala, South India: Can *Mansonia* (Diptera: Culicidae) play a supplemental role in
772 transmission? J Med Entomol. 2004;41(3):456–61.
- 773 44. Ughasi J, Bekard HE, Coulibaly M, Adabie-Gomez D, Gyapong J, Appawu M, et al. *Mansonia*
774 *africana* and *Mansonia uniformis* are Vectors in the transmission of *Wuchereria bancrofti*
775 lymphatic filariasis in Ghana. Parasites and Vectors. 2012;5(1):1–5.
- 776 45. Mitchell CJ, Forattini OP, Miller BR. Vector competence experiments with Rocio virus and
777 three mosquito species from the epidemic zone in Brazil. Rev Saude Publica. 1986;20(3):171–
778 7.
- 779 46. Kumar S, Stecher G, Li M, Knyaz C, Tamura K. MEGA X: Molecular evolutionary genetics
780 analysis across computing platforms. Mol Biol Evol. 2018;35(6):1547–9.
- 781 47. Harbach RE, Kitching IJ. The phylogeny of Anophelinae revisited: Inferences about the origin
782 and classification of Anopheles (Diptera: Culicidae). Zool Scr. 2016;45(1):34–47.
- 783 48. Harbach RE. The Culicidae (Diptera): A review of taxonomy, classification and phylogeny.
784 Zootaxa. 2007;1668(1):591–638.

- 785 49. Sun L, Li TJ, Fu WB, Yan ZT, Si FL, Zhang YJ, et al. The complete mt genomes of *Lutzia*
786 *halifaxia*, *Lt. fuscanus* and *Culex pallidothorax* (Diptera: Culicidae) and comparative analysis of
787 16 *Culex* and *Lutzia* mt genome sequences. *Parasites and Vectors*. 2019;12:368.
- 788 50. Harbach RE, Culverwell CL, Kitching IJ. Phylogeny of the nominotypical subgenus of *Culex*
789 (Diptera: Culicidae): insights from analyses of anatomical data into interspecific relationships
790 and species groups in an unresolved tree. *Syst Biodivers*. 2017;15(4):296–306.
- 791 51. Farajollahi A, Fonseca DM, Kramer LD, Marm Kilpatrick A. “Bird biting” mosquitoes and human
792 disease: A review of the role of *Culex pipiens* complex mosquitoes in epidemiology. *Infect*
793 *Genet Evol*. 2011;11(7):1577–85.
- 794 52. Nchoutpouen E, Talipouo A, Djiappi-Tchamen B, Djamouko-Djonkam L, Kopya E, Ngadjeu
795 CS, et al. *Culex* species diversity, susceptibility to insecticides and role as potential vector of
796 Lymphatic filariasis in the city of Yaoundé, Cameroon. *PLoS Negl Trop Dis*.
797 2019;13(4):e0007229.
- 798 53. Maquart PO, Boyer S. *Culex vishnui*. *Trends Parasitol*. 2022;(Vector of the Month).
- 799 54. Cornel AJ, Mcabee RD, Rasgon J, Stanich MA, Scott TW, Coetzee M. Differences in Extent of
800 Genetic Introgression between Sympatric *Culex pipiens* and *Culex quinquefasciatus* (Diptera:
801 Culicidae) in California and South Africa. *J Med Entomol*. 2003;40(1):36–51.
- 802 55. Zittra C, Flechl E, Kothmayer M, Vitecek S, Rossiter H, Zechmeister T, et al. Ecological
803 characterization and molecular differentiation of *Culex pipiens* complex taxa and *Culex*
804 *torrentium* in eastern Austria. *Parasites and Vectors*. 2016;9:197.
- 805 56. Batovska J, Cogan NOI, Lynch SE, Blacket MJ. Using next-generation sequencing for DNA
806 barcoding: Capturing allelic variation in ITS2. *G3 Genes, Genomes, Genet*. 2017;7(1):19–29.
- 807 57. Aspen S, Savage HM. Polymerase chain reaction assay identifies North American members of
808 the *Culex pipiens* complex based on nucleotide sequence differences in the
809 acetylcholinesterase gene *Ace.2*. *J Am Mosq Control Assoc*. 2003;19(4):323–8.
- 810 58. Sirivanakarn S. A review of the Systematics and a Proposed Scheme of Internal Classification
811 of the New World Subgenus *Melanoconion* of *Culex* (Diptera, Culicidae). *Mosq Syst*.
812 1982;14(4):265–333.
- 813 59. Takhampunya R, Kim HC, Tippayachai B, Kengluetcha A, Klein TA, Lee WJ, et al. Emergence
814 of Japanese encephalitis virus genotype v in the Republic of Korea. *Virology*. 2011;8:449.

- 815 60. Navarro JC, Weaver SC. Molecular phylogeny of the Vomerifer and Pedroi Groups in the
816 spissipes section of the subgenus Culex (Melanoconion). *J Med Entomol.* 2004;41(4):575–81.
- 817 61. Sallum MAM, Forattini OP. Revision of the Spissipes Section of Culex (Melanoconion)
818 (diptera: Culicidae). *J Am Mosq Control Assoc.* 1996;12(3):517–600.
- 819 62. Torres-Gutierrez C, Bergo ES, Emerson KJ, de Oliveira TMP, Greni S, Sallum MAM.
820 Mitochondrial COI gene as a tool in the taxonomy of mosquitoes Culex subgenus
821 Melanoconion. *Acta Trop.* 2016;164:137–49.
- 822 63. Torres-Gutierrez C, De Oliveira TMP, Emerson KJ, Bergo ES, Sallum MAM. Molecular
823 phylogeny of Culex subgenus Melanoconion (Diptera: Culicidae) based on nuclear and
824 mitochondrial protein-coding genes. *R Soc Open Sci.* 2018;5:171900.
- 825 64. Arctander P. Comparison of a mitochondrial gene and a corresponding nuclear pseudogene.
826 *Proc R Soc B Biol Sci.* 1995;262(1363):13–9.
- 827 65. Hurst GDD, Jiggins FM. Problems with mitochondrial DNA as a marker in population,
828 phylogeographic and phylogenetic studies: The effects of inherited symbionts. *Proc R Soc B*
829 *Biol Sci.* 2005;272:1525–34.
- 830 66. De Oliveira CD, Gonçalves DS, Baton LA, Shimabukuro PHF, Carvalho FD, Moreira LA.
831 Broader prevalence of Wolbachia in insects including potential human disease vectors. *Bull*
832 *Entomol Res.* 2015;105(3):305–15.
- 833 67. Beebe NW. DNA barcoding mosquitoes: Advice for potential prospectors. *Parasitology.*
834 2018;145(5):622–33.
- 835 68. Behura SK. Molecular marker systems in insects: current trends and future avenues. *Mol Ecol.*
836 2006;15(11):3087–113.
- 837 69. Vezenegho SB, Issaly J, Carinci R, Gaborit P, Girod R, Dusfour I, et al. Discrimination of 15
838 Amazonian Anopheline Mosquito Species by Polymerase Chain Reaction—Restriction
839 Fragment Length Polymorphism. *J Med Entomol.* 2022;tjac008.
- 840 70. Quast C, Pruesse E, Yilmaz P, Gerken J, Schweer T, Yarza P, et al. The SILVA ribosomal
841 RNA gene database project: Improved data processing and web-based tools. *Nucleic Acids*
842 *Res.* 2013;41(Database issue):D590–6.
- 843 71. Altschul SF, Gish W, Miller W, Myers EW, Lipman DJ. Basic local alignment search tool. *J Mol*
844 *Biol.* 1990;215(3):403–10.

- 845 72. Bankevich A, Nurk S, Antipov D, Gurevich AA, Dvorkin M, Kulikov AS, et al. SPAdes: A new
846 genome assembly algorithm and its applications to single-cell sequencing. *J Comput Biol.*
847 2012;19(5):455–77.
- 848 73. Folmer O, Black M, Hoeh W, Lutz R, Vrijenhoek R. DNA primers for amplification of
849 mitochondrial cytochrome c oxidase subunit I from diverse metazoan invertebrates. *Mol Mar*
850 *Biol Biotechnol.* 1994;3(5):294–9.
- 851 74. Rausch T, Fritz MHY, Untergasser A, Benes V. Tracy: Basecalling, alignment, assembly and
852 deconvolution of sanger chromatogram trace files. *BMC Genomics.* 2020;21(1):230.
- 853 75. Rausch T, Hsi-Yang Fritz M, Korb J, Benes V. Alfred: Interactive multi-sample BAM
854 alignment statistics, feature counting and feature annotation for long- and short-read
855 sequencing. *Bioinformatics.* 2019;35(14):2489–91.
- 856 76. Madeira F, Park YM, Lee J, Buso N, Gur T, Madhusoodanan N, et al. The EMBL-EBI search
857 and sequence analysis tools APIs in 2019. *Nucleic Acids Res.* 2019;47(W1):W636–41.
- 858 77. Edgar RC. MUSCLE: A multiple sequence alignment method with reduced time and space
859 complexity. *BMC Bioinformatics.* 2004;5:113.

860

861 TABLE AND FIGURES LEGENDS

862 **Table 1.** List of mosquito species represented in this study and their vector status. Origin countries
863 are listed as their ISO alpha-2 codes: Central African Republic, CF; Cambodia, KH; Madagascar, MG;
864 French Guiana, GF.

865 **Figure 1. (A)** Proportion of rRNA reads found in mosquito specimen pools depleted by probe
866 hybridisation followed by RNase H digestion. Probes were antisense to *Ae. aegypti* rRNA sequences.
867 **(B)** Read vs. score ratio plot of “*Ae simpsoni* CS S27”. Green line indicates a 0.8 cut-off where only
868 reads above this threshold are used in rRNA assembly.

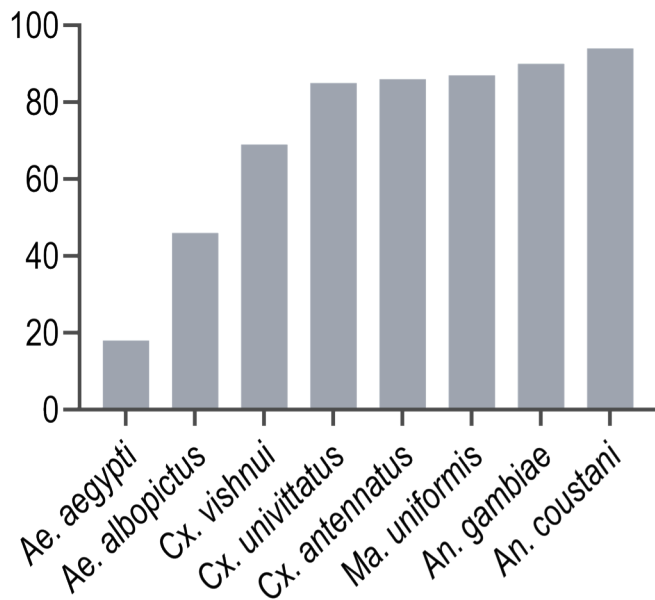
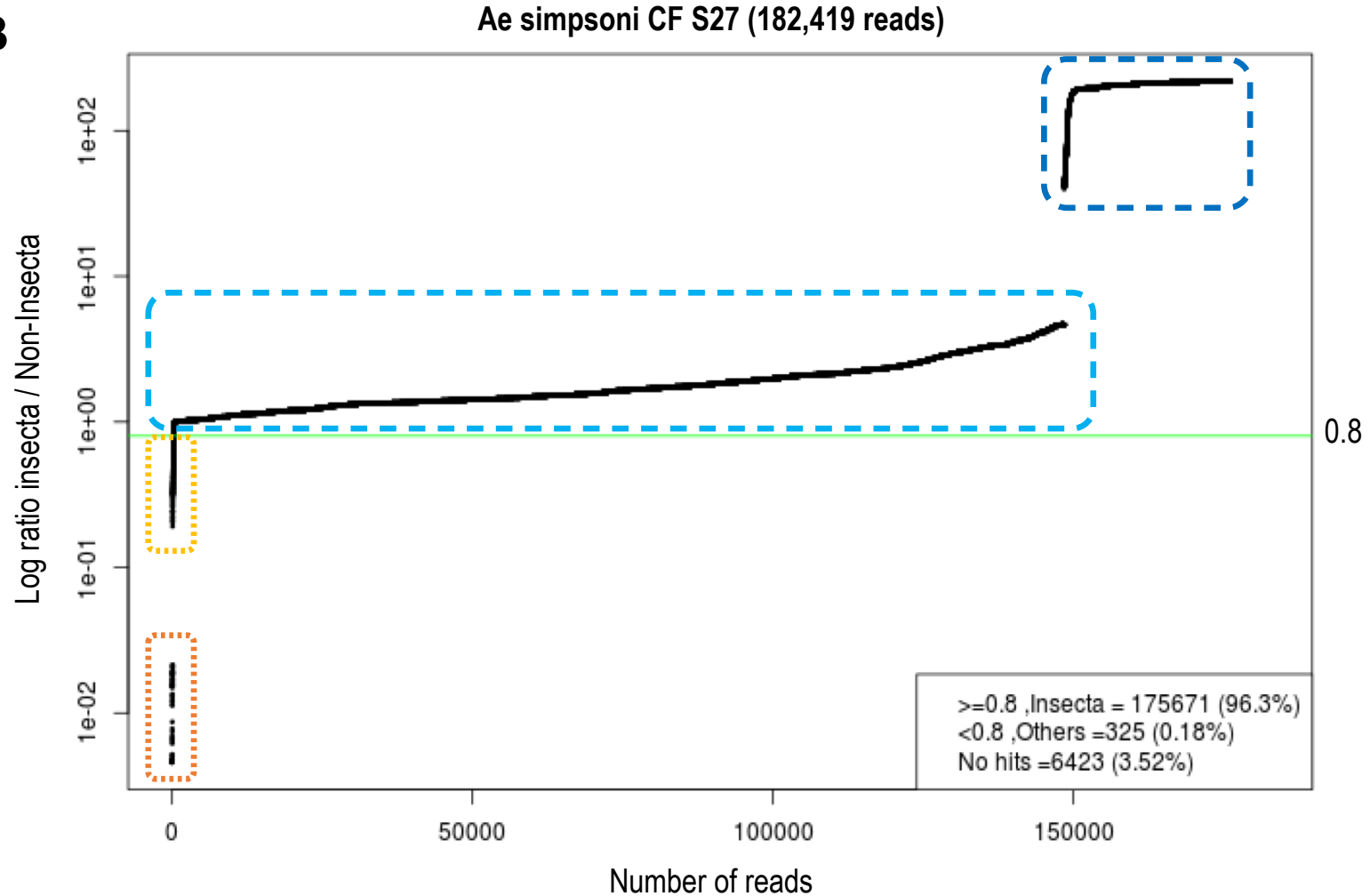
869 **Figure 2.** Phylogenetic tree based on concatenated 28S and 18S sequences generated from this
870 study (3900+1900 bp) as inferred using the maximum-likelihood method and constructed to scale in
871 MEGA X (46). Values at each node indicate bootstrap support (%) from 500 replications. For
872 sequences from this study, each specimen label contains information on its taxonomy, origin (as
873 indicated in 2-letter country codes), and specimen ID. Label colours indicate genera: *Culex* in coral,
874 *Anopheles* in purple, *Aedes* in dark blue, *Mansonia* in dark green, *Culiseta* in maroon, *Limatus* in light

875 green, *Coquillettidia* in light blue, *Psorophora* in yellow, *Mimomyia* in teal, *Uranotaenia* in pink and
876 *Eretmapodites* in brown. Scale bar at 0.05 is shown.

877 **Figure 3.** Phylogenetic tree based on 28S sequences generated from this study and from NCBI
878 databases combined (3900 bp) as inferred using the maximum-likelihood method and constructed to
879 scale in MEGA X (46). Values at each node indicate bootstrap support (%) from 500 replications. For
880 sequences from this study, each specimen label contains information on its taxonomy, origin (as
881 indicated in 2-letter country codes), and specimen ID. Labels in bold indicate sequences derived from
882 NCBI. Label colours indicate genera: *Culex* in coral, *Anopheles* in purple, *Aedes* in dark blue,
883 *Mansonia* in dark green, *Culiseta* in maroon, *Limatus* in light green, *Coquillettidia* in light blue,
884 *Psorophora* in yellow, *Mimomyia* in teal, *Uranotaenia* in pink and *Eretmapodites* in brown. Scale bar
885 at 0.05 is shown.

886 **Figure 4.** Phylogenetic tree based on concatenated 28S and 18S sequences generated from this
887 study (3900+1900 bp) as inferred using maximum-likelihood method and constructed to scale in
888 MEGA X (46). Values at each node indicate bootstrap support (%) from 500 replications. For
889 sequences from this study, each specimen label contains information on its taxonomy, origin (as
890 indicated in 2-letter country codes), and specimen ID. Label colours indicate genera: *Culex* in coral,
891 *Anopheles* in purple, *Aedes* in dark blue, *Mansonia* in dark green, *Culiseta* in maroon, *Limatus* in light
892 green, *Coquillettidia* in light blue, *Psorophora* in yellow, *Mimomyia* in teal, *Uranotaenia* in pink and
893 *Eretmapodites* in brown. Scale bar at 0.05 is shown.

894 **Figure 5.** Phylogenetic tree based on COI sequences (621–699 bp) as inferred using the maximum-
895 likelihood method and constructed to scale in MEGA X (46). Values at each node indicate bootstrap
896 support (%) from 500 replications. Each specimen label contains information on its taxonomy, origin
897 (as indicated in 2-letter country codes), and specimen ID. Label colours indicate genera: *Culex* in
898 coral, *Anopheles* in purple, *Aedes* in dark blue, *Mansonia* in dark green, *Limatus* in light green,
899 *Coquillettidia* in light blue, *Psorophora* in yellow, *Mimomyia* in teal, *Uranotaenia* in pink and
900 *Eretmapodites* in brown. Scale bar at 0.05 is shown.

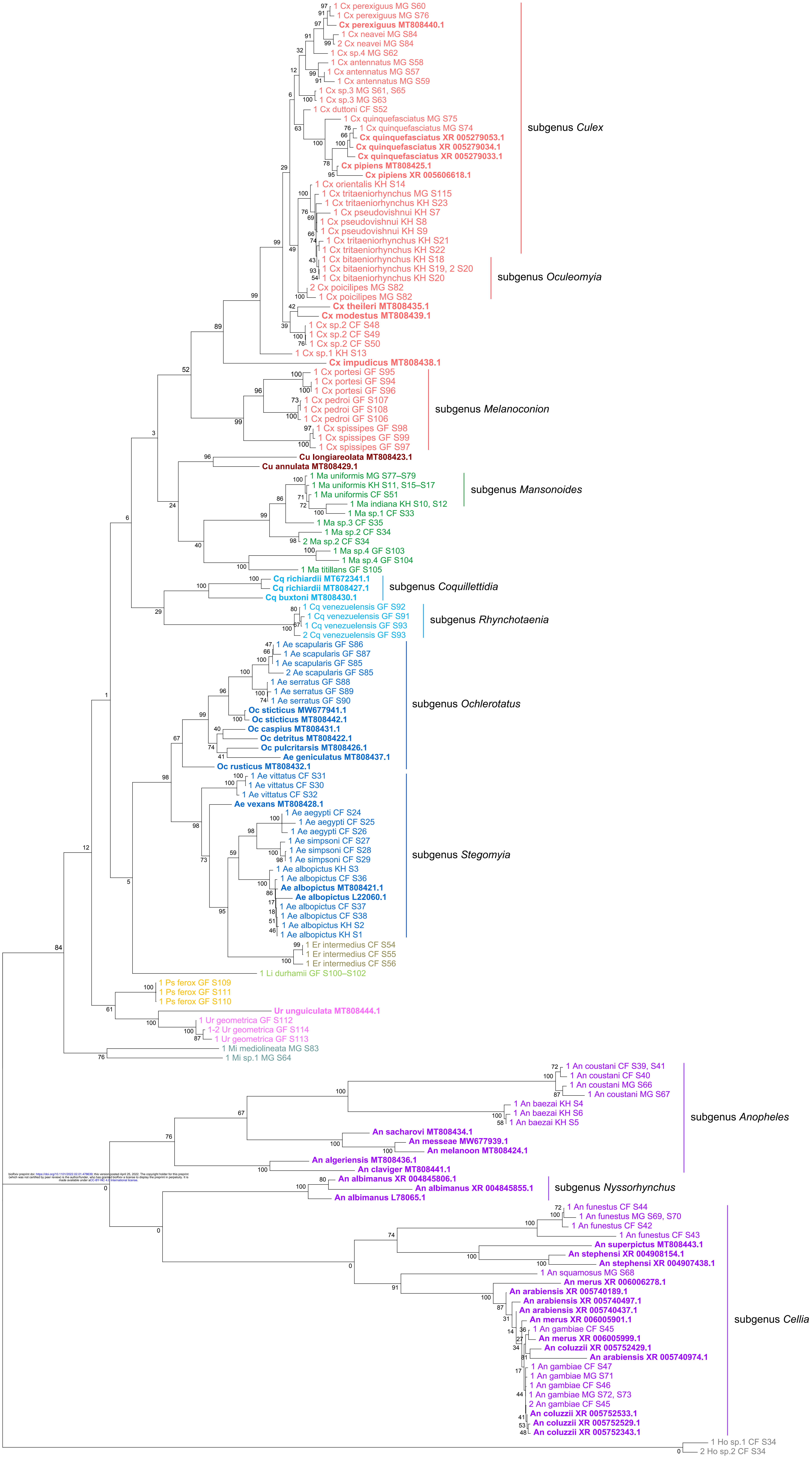
A**B**

--- hits only in the Insecta library

--- hits higher-scoring in Insecta library

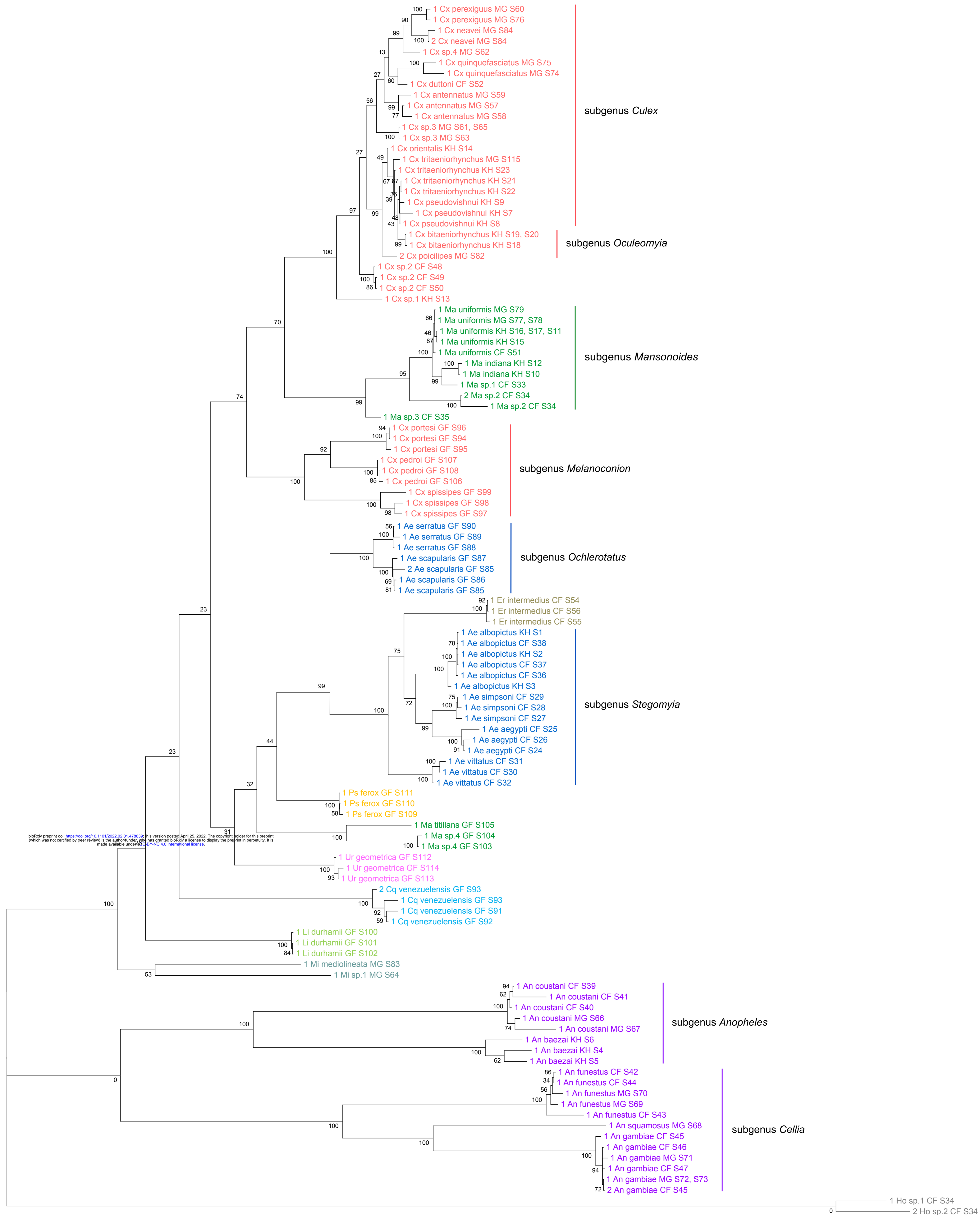
--- hits higher-scoring in Non-Insecta library

--- hits only in Non-Insecta library





0.05



bioRxiv preprint doi: <https://doi.org/10.1101/2022.02.01.478633>; this version posted April 25, 2022. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted bioRxiv a license to display the preprint in perpetuity. It is made available under aCC-BY-NC 4.0 International license.

