

1

2 **Large scale analysis of over a thousand *Wolbachia*** 3 **genomes sheds new light on its evolution.**

4

5 Matthias Scholz^{1,2}, Davide Albanese¹, Kieran Tuohy¹, Claudio Donati¹,
6 Nicola Segata², Omar Rota-Stabelli¹

7

8 ¹ Research and Innovation Centre, Fondazione Edmund Mach (FEM), San Michele all'Adige, Italy

9 ² Department CIBIO, University of Trento, Trento, Italy

10 Correspondence and requests for materials should be addressed to NS (email:
11 nicola.segata@unimi.it) and O R-S (email:omar.rota@fmach.it).

12

13 *Supplementary information*

14

15 **Supplementary Notes**

16

17 **1. Detailed description of novel strains and affinities in the** 18 **phylogenetic trees of Fig. 2**

19

20 **General phylogenetic considerations**

21

22 Although they support very similar topologies, the presence-absence tree (Fig 2b) is generally less
23 supported than the alignment tree (Fig. 2a). This may be caused, at least partially, by some
24 homoplasies (independently shared genes) in the presence/absence tree and by the sequence
25 alignment being two orders of magnitude larger, therefore more robust to bootstrapping.

26

27 All genomes assembled in our analyses and depicted in Fig 2 belong to group A, B, C, D, F; none
28 are from other less known groups such as E, M, or L. The most likely explanation is that the latter
29 super-groups lack a high quality reference genomes: our pipeline uses a stringent annotation
30 criteria to obtain high quality assemblies, which does not allow to identify *Wolbachia* genomes very

31 distant from known high quality reference ones. We foresee that our method will further increase
32 *Wolbachia* sample diversity as soon as new high quality genomes will be available from less
33 represented supergroups.

34

35 **Supergroup specific phylogenetic considerations**

36

37 **Supergroup A**

38 Even though the large inner branch of **supergroup A** is mainly composed of *Wolbachia* present in
39 Diptera (flies) and Hymenoptera (bees, wasps, ants), we also found a related strain in Hemiptera
40 *Megacopta cribraria* (bug) centrally located in the same branch, and closer to the root located
41 strains from the Coleoptera *Diabrotica virgifera virgifera* (beetle) and the Hymenoptera
42 *Camponotus obliquus* (ant) closely related to the reference wDacA (Hemiptera). We show enlarged
43 diversity of genomes in various *Drosophila* species, including >500 *D. melanogaster*, >500 *D.*
44 *simulans*, and 22 *D. ananassae*. We also assembled two new genomes from *D. yakuba* and one
45 from *D. santomea*, located at the root of the *D. melanogaster* branch. Particularly interesting is an
46 additional new genome from the fly *Holcocephala fusca* located close to the root of the *D.*
47 *melanogaster* cluster. We also increased the diversity of *Diachasma alloeum* (wasp) by 7 genomes
48 and assembled new *Wolbachia* genomes from the flies *Megaselia abdita* and *Sphyracephala*
49 *brevicornis*.

50

51 **Supergroup B**

52 We expanded **supergroup B** by providing the first representative *Wolbachia* genomes of several
53 hosts. The deepest branch in the supergroup B contains *Wolbachia* strains from Lepidoptera
54 (butterflies, moths) which were not present in other supergroups. We found a new strain in
55 *Polygonia c-album* (butterfly) closely related to the reference wOb from *Operophtera brumata*
56 (moth). And we assembled a new strain from *Pararge aegeria* (butterfly) which is closely related to
57 the references wPip (mosquito) and wBol1 (butterfly). We also found a very distinct strain in
58 *Callosobruchus chinensis* (beetle) closely related to a new strain from an unidentified insect.
59 Finally, the position closely located to the root of this Lepidoptera dominated branch is occupied by
60 a new strain from *Tetranychus urticae* (mite). The central part of supergroup B is mainly composed
61 of Hemiptera. We enlarged the diversity of *Diaphorina citri* (bug) by 9 genomes and we found a
62 new strain in *Homalodisca vitripennis* (leafhopper).

63

64 **Supergroups C, D, F**

65 For the insect related **supergroup F** embedded in the nematode branch, we increased the
66 diversity of *Cimex Lectularius* (bed bug) and, we assembled a new *Wolbachia* genome from
67 *Melophagus ovinus* (fly) which shows to be closely related to *Wolbachia* wOc from *Osmia*
68 *caerulescens* (bee). In the nematode **supergroup C**, we found a novel *Wolbachia* strain from

69 *Onchocerca gutturosa* which is placed at the root of the largely expanded *Onchocerca* cluster
70 consisting of 39 *Onchocerca volvulus* and 7 *Onchocerca ochengi* genomes. In **supergroup D**, we
71 increased the diversity by 22 *Wolbachia* genomes from *Brugia malayi* and 4 from *Wuchereria*
72 *bancrofti*, and we provide 2 new *Wolbachia* genomes from *Brugia pahangi* that form a distinct
73 branch closely related to *B. malayi*.

74

75

76

77

2. Integration analyses in *Callosobruchus chinensis* and *Drosophila ananassae*.

For *Callosobruchus chinensis* and *Drosophila ananassae*, we manually inspected our core-genome data (the one used for our main tree of figure 2). According to Choi et al. (2015, GBE) integrated *Wolbachia* are characterised by an excess of non synonymous mutations as well as stop codons and frameshifts because of relaxed selection on the integrated *Wolbachia* compared to non-integrated one. We first realigned the core genome for these two hosts in order to recover the codon frame of genes (the core genome is mainly composed of conserved *Wolbachia* coding genes), we then look for intergenic stop codons, frameshifts, as well as for region poorly aligned. We did not find any internal stop codon nor disruptive insertions/deletions, except for few poorly aligned fragments, which we have been blasted: they all look genuinely as *Wolbachia*:

- In one of the two *Callosobruchus Wolbachia* (SRR949786), we found a poorly aligned sequence at position 22673 of the core genome alignment (position 4588 of *Callosobruchus* alignment) 1638 nucleotide long. This regions however does blast with 99.5% identity (and zero gaps) against well annotated *Wolbachia* genomes such as wPip and wAlbB.
- In *D. ananassae* in SRS2127163 we found a 1757 nt fragment at position 19245 of the core genome) which however blast with 100% similarity to *Drosophila ananassae* strain W2.1 chromosome (wAna, cp042094.1) and with wRi (CP001391.1).
- In *D. ananassae* SRS2127151-2127152-2127153 we found a 489 nt fragment at position 6132 which blast 100% with *Drosophila ananassae* strain W2.1 chromosome (wAna, cp042094.1) and wRi (CP001391.1).
- In *D. ananassae* SRS2126857-2126916-21235641-2127154 at position 36570 a fragment of 783 nt blasting 100% with *Drosophila ananassae* strain W2.1 chromosome (wAna, cp042094.1) and wRi (CP001391.1).
- In *D. ananassae* SRS2126857-2126916-21235641-2127154-2135644-2135642 at position 181210 of 993 nt blasting 100% with *Drosophila ananassae* strain W2.1 chromosome (wAna, cp042094.1) and wMel_ZH26 (CP042445.1)

The fact that all these fragments blast with perfect or nearly perfect similarity with well-curated reference genomes such as wRi, wMel and wPip, and that we did not found any present insertion/deletion when compared to them reassured us that they are genuinely *Wolbachia* fragments and not host integrated fragments.

Furthermore, we also inspected our MUMmer alignment of *D. ananassae Wolbachia*, and did not find any evidence of wAna^{INT}: because according to Choi et al. 2015 (table 2 and table 3 therein) wAna^{INT} accumulate 20 times more mutations than wAna^{NF}, we would expect an excess of

113 mutations (therefore long branches in a phylogeny) in samples contaminated by wAnaINT
114 fragments: the branch length of all samples is instead homogenous (similar to each other) very
115 similar according to a RAxML analysis of the dataset.

116

117 3. Host verification

118 We verified the host species by reconstructing the 18S ribosomal RNA gene sequence using the
119 tool RiboTagger. In all cases of successfully 18S reconstruction, we could not find any clear
120 mislabelling of the host species. The host in the majority of our samples could be directly confirmed
121 at the species level. Remaining samples could not be distinguished at species levels, but were
122 confirmed at higher taxonomic levels, such as genus, and subgroup because annotated 18S are
123 missing from the reference database. In few samples we could not extract the rRNA sequence, as
124 the host sequences might have been removed and only the metagenome was made public (*D.*
125 *recens*, the source data for wRec is one of these cases). For two species, *Diabortica virgifera* and
126 *Drosophila ananassae*, RiboTagger could not confirm the taxonomic assignment of their SRA
127 sample, respectively finding likely contaminant fungi or the fruit fly *S. lebanonensis*. This may be
128 due to an issue in the RiboTagger pipeline or available references. We manually inspected COI
129 and could confirm that the best blast hits of those samples were reciprocally with *D. virgifera* and
130 *D. ananassae* (and no hit were found with *S. lebanonensis* in these samples).

131

132 We performed extra validation for two peculiar cases: *Holcocephala fusca* and *Caenorhabditis*
133 *remanei*. At the genus level, we clearly could confirm the host *Holcocephala fusca* (SRR1738186)
134 to be similar to the 18S reference of *Holcocephala abdominalis* 18S reference (the 18S of *H. fusca*
135 is not yet deposited in GenBank). We further check for signs of contamination by blasting the
136 SRR1738186 using COI of *Drosophila melanogaster* as query and found various reads with 95%
137 similarity: blasting of these reads against the nucleotide collection returned as best hits other
138 dipterans of the Syrphidae family (hoverflies such as *Sphaerophoria* sp., *Eristalis* sp, *Dasyhlea* sp).
139 It is therefore possible that the wMel-like *Wolbachia* is from hoverfly prey (but not from a
140 *Drosophila* as the COI found in SRR1738186 does not have highest hits to them) of the robber fly.
141 For *Caenorhabditis remanei* (SRR275642), we could not reconstruct any 18S sequence using
142 RiboTagger due to low sequencing depth, but this sample is considered to be of lower quality,
143 present only in the gene-content tree (Fig. 2b), and not part of our genome set in the core
144 sequence tree (Fig. 2a). We manually blasted SRR275642 using *Caenorhabditis remanei*
145 Cytochrome Oxidase subunit I (COI) and found reads covering the whole gene with 99-100%
146 identity, confirming the exact source of this sample. However, in order to exclude contaminants we
147 further blasted using *Drosophila mauritiana* COI (AF200831.1) because the putative new genome
148 for *Caenorhabditis remanei* has 98% identity to wNo of *Drosophila mauritiana* in our Fig. 2b.
149 Indeed, we found reads covering a portion (not all) of the gene with 98-100% identity, indicating a
150 contamination from a *Drosophila mauritiana* or from another closely related *Drosophila*.

151

152

4. Host "tags" used to search the NCBI SRA repository

154 insects, insect, bug, bugs, worm, roundworm, roundworms, silkworm, armyworm, termite, termites,
 155 ant, ants, mite, mites, ticks, tick, springtail, springtails, bees, bee, wasp, wasps, flea, fleas, moth,
 156 moths, Beetle, Beetles, spider, spiders, Wolbachia, Arthropods, Arthropod, Nematodes, Nematode,
 157 Anopheles, Aedes, Folsomia, Culex, Ctenocephalides, Onchocerca, Brugia, Osmia, Drosophila,
 158 Armadillidium, Diaphorina, Dirofilaria, crickets, Caenorhabditis, Aedes, Agelenopsis, Anastrepha,
 159 Anax, Anoplolepis, Argynnis, Autographa, Azteca, Bemisia, Brugia, Caleta, Camponotus, Cimex,
 160 Colias, Coptotermes, Crematogaster, Culex, Danaus, Diaphorina, Dirofilaria, Encarsia, Ephestia,
 161 Erebia, Euphydryas, Eurema, Euwallacea, Formica, Glossina, Gryllus, Hypolimnas, Incisitermes,
 162 Ischnura, Junonia, Kerria, Laodelphax, Leptopilina, Limenitis, Lycaeides, Melanitis, Melitaea,
 163 Minois, Monomorium, Myrmecorhynchus, Myrmica, Nasonia, Neptis, Nilaparvata, Notoncus,
 164 Odontomachus, Odontotermes, Onchocerca, thetrum, Ostrinia, Pantala, Papilio, Parantica,
 165 Parasite mite, Parnassius, Pheidole, Plebejus, Polistes, Polyergus, Polygonia, Polyommatus,
 166 Polyrhachis, Pseudomyrmex, Rhagoletis, Rhytidoponera, Sitona, Sogatella, Solenopsis,
 167 Spodoptera, Stenamma, Steriphus, Teleogryllus, Tetranychus, Tetrastichus, Tribolium,
 168 Trichogramma, Wasmannia, Xyleborus, Ypthima, Zootermes, cricket, mite, mites, Hexapoda,
 169 Collembola, Hexapoda, Isoptera, Zootermes, angusticollis, planthopper, leafhopper, butterfly,
 170 butterflies, crustacean, weevil, weevils, Ixodes, scorpion, scorpions, Acari, Acarina, cockroach,
 171 Cryptocercus, Cryptotermes, Termitidae, Alyscotermes, Macrotermitinae, Macrotermes,
 172 Microtermes, Odontotermes, Nasutitermitinae, Nasutitermes, Trinervitermes, Cornitermes,
 173 Syntermes, Termitinae, Amitermes, Amitermes, Globitermes, Microcerotermes, Cubitermes,
 174 Cubitermes, Ophiotermes, Neocapritermes, Promirotermes, Termes, "water bears", tardigrade,
 175 harvestman, symphylans, millipede, amphipod, isopod, oniscus, woodlice, carpenter, slater, simon,
 176 coneheads, Proturans, Diplurans, bristletails, silverfish, Mayflies, Dragonflies, damselflies,
 177 Grasshoppers, crickets, Earwigs, Stoneflies, Mantids, Cockroaches, cicadas, leafhoppers, aphids,
 178 psyllids, thrips, Booklice, Lice, Lacewings, Caddisflies, Caddisfly, Scorpionflies, Tardigrada,
 179 Onychophora, Chilopoda, Nematomorpha, Diplopoda, Symphyla, Arachnida, Amblypygi, Araneae,
 180 Opiliones, Pseudoscorpiones, Ricinulei, Scorpiones, Solifugae, Thelyphonida, Homoptera,
 181 Branchiopoda, Phyllophora, Sarsostraca, Remipedia, Maxillophora, Thecostraca, Branchiura,
 182 Copepoda, Ostracoda, Myodocopa, Malacostraca, Hoplocarida, Eumalacostraca, Protura, Diplura,
 183 Microcoryphia, Thysanura, Ephemeroptera, Odonata, thoptera, Phasmatodea, Grylloblattodea,
 184 Mantophasmatodea, Dermaptera, Plecoptera, Embiidina, Zoraptera, Mantodea, Blattodea,
 185 Hemiptera, Heteroptera, Sternorrhyncha, Thysanoptera, Psocoptera, Phthiraptera, Amblycera,
 186 Anoplura, Coleoptera, Neuroptera, Hymenoptera, Trichoptera, Lepidoptera, Siphonaptera,
 187 Mecoptera, Strepsiptera, Diptera, crayfish, crab, crabs, "water bears", tardigrade, harvestman,
 188 symphylans, millipede, amphipod, isopod, oniscus, woodlice, carpenter, slater, simon, coneheads,
 189 Proturans, Diplurans, bristletails, silverfish, Mayflies, Dragonflies, damselflies, Grasshoppers,

190 crickets, Earwigs, Stoneflies, Mantids, Cockroaches, cicadas, leafhoppers, aphids, psyllids, thrips,
191 Booklice, Lice, Lacewings, Caddisflies, Caddisfly, Scorpionflies, Tardigrada, Onychophora,
192 Chilopoda, Nematomorpha, Diplopoda, Symphyla, Arachnida, OR Amblypygi, Araneae, Opiliones,
193 Pseudoscorpiones, Ricinulei, Scorpiones, Solifugae, Thelyphonida, Homoptera, Branchiopoda,
194 Phyllopoda, Sarsostraca, Remipedia, Maxillopoda, Thecostraca, Branchiura, Copepoda,
195 Ostracoda, Myodocopa, Malacostraca, Hoplocarida, Eumalacostraca, Protura, Diplura,
196 Microcoryphia, Thysanura, Ephemeroptera, Odonata, thoptera, Phasmatodea, Grylloblattodea,
197 Mantophasmatodea, Dermaptera, Plecoptera, Embiidina, Zoraptera, Mantodea, Blattodea,
198 Hemiptera, Heteroptera, Sternorrhyncha, Thysanoptera, Psocoptera, Phthiraptera, Amblycera,
199 Anoplura, Coleoptera, Neuroptera, Hymenoptera, Trichoptera, Lepidoptera, Siphonaptera,
200 Mecoptera, Strepsiptera, Diptera

201

202 **5. Keywords searched but not present in NCBI**

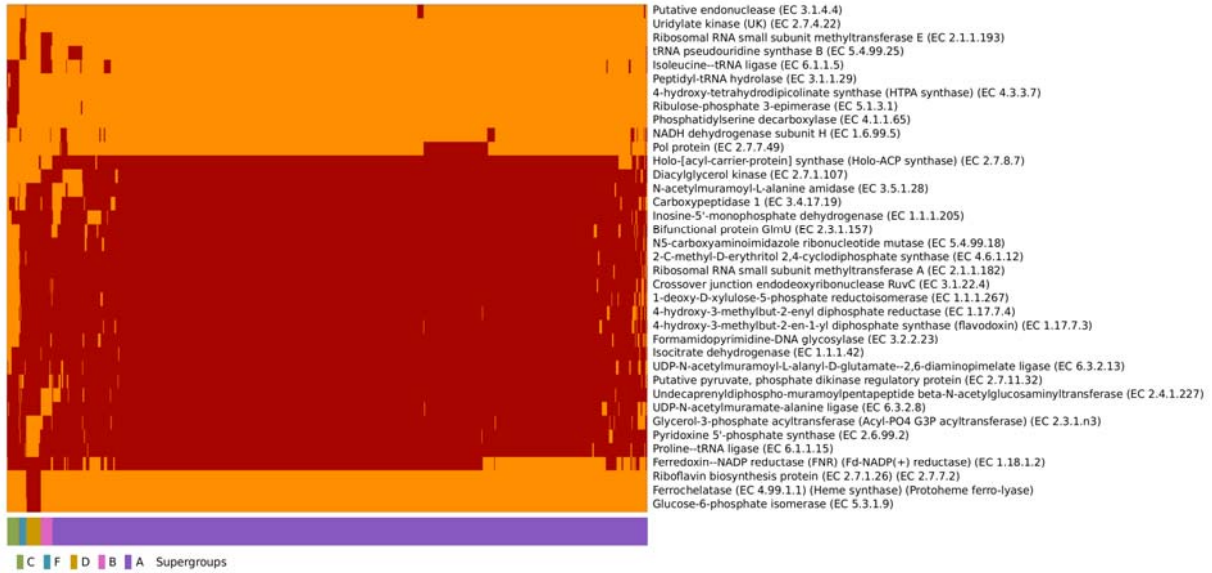
203 Aciagrion, Acisoma, Acraea, Acromis, Aenictus, Aganaspis, Amblyptilia, Anthene, Apanteles,
204 Aphantopus, Apoica, Araschnia, Ariadne, Aricia, Asobara, Azanus, Barronopsis, Brachythemis,
205 Brachythems, Brangas, Cabera, Carcharodus, Carpomya, Carterocephalus, Castalius, Catopsilia,
206 Celastrina, Cepora, Ceriagrion, Ceutorhynchus, Chelonus, Chelymorpha, Clossiana,
207 Coenonympha, Colotis, Corcyra, Cordylochernes, Crocothemis, Cupido, Delias, Diaphorencyrtus,
208 Dictyophara, Diplacodes, Diplazon, Doryctobracon, Dorymyrmex, Epophthalmia, Eretmocerus,
209 Eriborus, Eumedonia, Evagetes, Everes, Fabriciana, Gambrus, Geometra, Glaucopsyche, Heodes,
210 Heteropterus, Hipparchia, Hirtodrosophila, Horaga, Hylyphantes, Hyponephela, Hyposoter,
211 Ictinogomphus, Iraota, Ixias, Jalmenus, Jamides, Lasiommata, Leptidea, Leptogenys,
212 Leptomymex, Leptosia, Leucophenga, Libythea, Lissorhoptrus, Lophomyrmex, Lopinga, Lycaena,
213 Macrosteles, Maculinea, Mansonia, Melanargia, Metapone, Microgaster, Muscidifurax Nacaduba,
214 Neurothemis, Nymphalis, Ochetellus, Ochloides, Ocymymex, Odontosema, Opistograptis,
215 Opisthophthalmus, Opius, eina, ius, nipholidotos, Pareronia, Petrobia, Pleuroptia, Polybia, Pontia,
216 Protocalliphora Pseudozizeeria, Pyrgus, Satyrium, Slavum, Spalgis, Sueus, Suillia, Surendra,
217 Sycoscapter, Syrphophilus, Tarucus, Technomyrmex, Telicada, Teractrocera, Thecla,
218 Thersamonia, Tongeia, Trithemis, Tymmophorus, Udaspes, Walkerella, Xylosandrus, Zizeeria,
219 Demodex, Acherongia, Acherontides, Acherontiella, Acheroxenylla, Austrogastrura, Barbagastrura,
220 Biscoia, Bonetogastrura, Celegastrura, Ceratophysella, Choreutinula, Cosberella, Denigastrura,
221 Ecuadogastrura, Gnathogastrura, Gomphiocephalus, Hypogastrura, Jacutogastrura,
222 Mesachorutes, Mesogastrura, Microgastrura, Neobeckerella, Octoacanthella,
223 Ongulogastrura, ogastrura, Parawillemia, Paraxenylla, Pseudacherontides, Schaefferia,
224 Stenogastrura, Tafallia, Taurogastrura, Thibaudylla, Triacanthella, OR Typhlogastrura,
225 Willemgastura, Willemia, Xenylla, Xenyllogastrura, Hypogastruridae

226

227 **Supplementary Figures**

228

229



230

231

232 **Supplementary Figure 1.** Heatmap of the comparative genomic analysis of 989 novel high
 233 quality *Wolbachia* assemblies associated to the 14 host groups in Fig. 4b. Shown are enzyme
 234 categories (EC) that are significantly different between host groups (Fisher test).

235

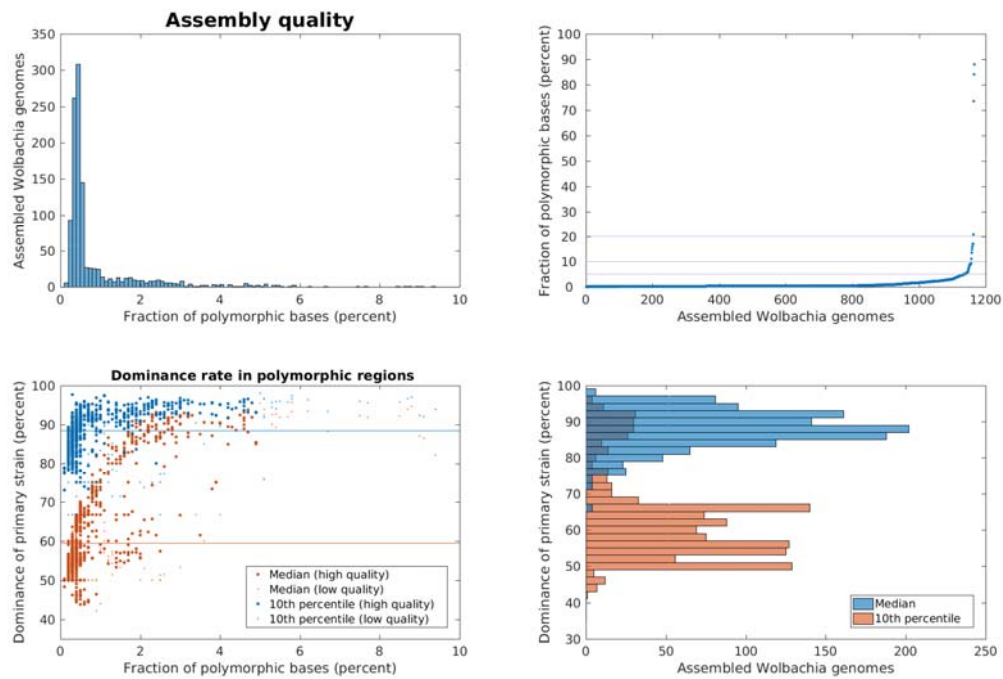
236

237

238

239

240



241

242

243 **Supplementary Figure 2.** Assembly quality control based on polymorphic site identification. (a,b)

244 Most assembled *Wolbachia* genomes show polymorphic pattern in less than one percent of the

245 total genome size. (c,d) In these polymorphic regions, the primary bases show a median

246 dominance of 88.3 percent. In the lower 10th percentile, the median dominance of the primary

247 sequence is 59.45 percent (horizontal lines). Based on four criteria including polymorphism,

248 assembled genomes are considered as of high quality (big dots), low quality (small dots).

249

250

251

252

253

254

255

256

257

258

259

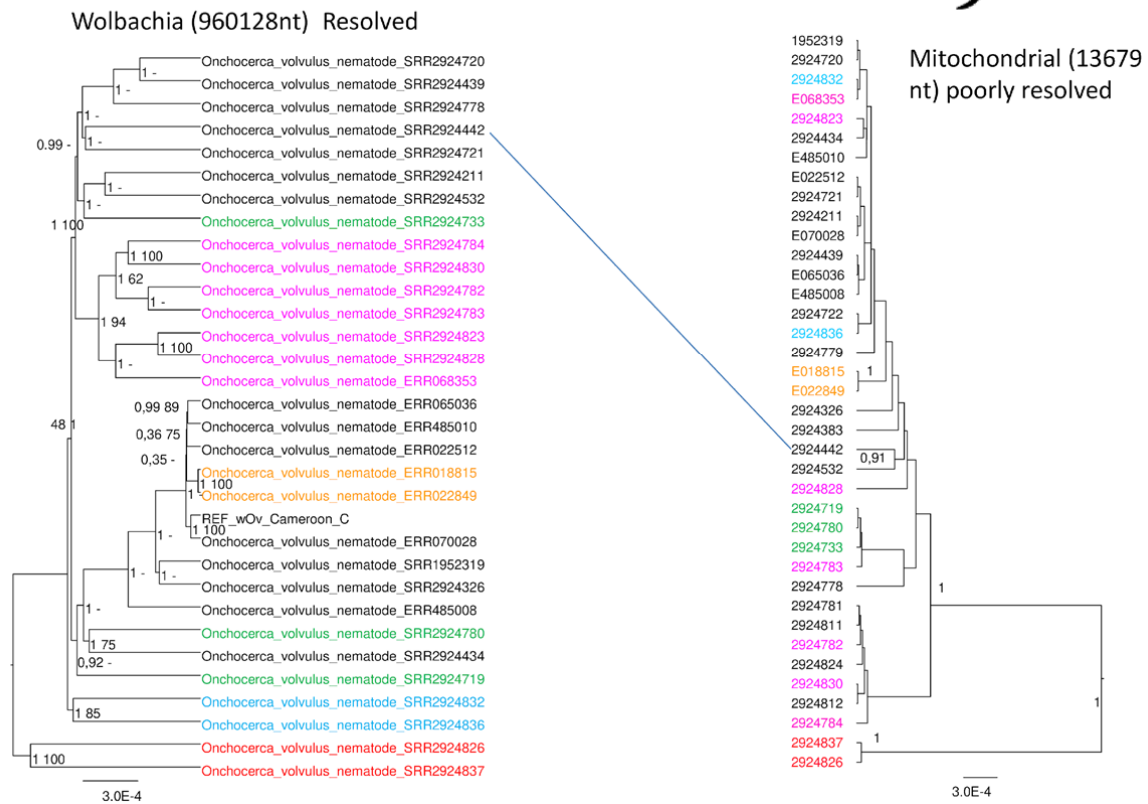
260

261

262 **Supplementary Figure 3:** Co-phylogenies of *Wolbachia* and mitochondria reconstructed from the
263 same Sequence Read Archive file for 11 different host species (in 15 different panels). Trees are
264 the posterior consensus and values at nodes are posterior probabilities from BEAST analysis
265 employing for all analyses GTR+G replacement model, strict uncalibrated clock, and constant
266 coalescent tree priors. Plus in some cases Bootstrap supports from a GTR+G RAXML analysis. For
267 *D. melanogaster* and *D. simulans* we further show a cladogram of the BEAST tree to ease
268 discriminating nodes and supports.

269

***Onchocerca volvulus*: poorly resolved mitochondria, but at least one incongruence**



270

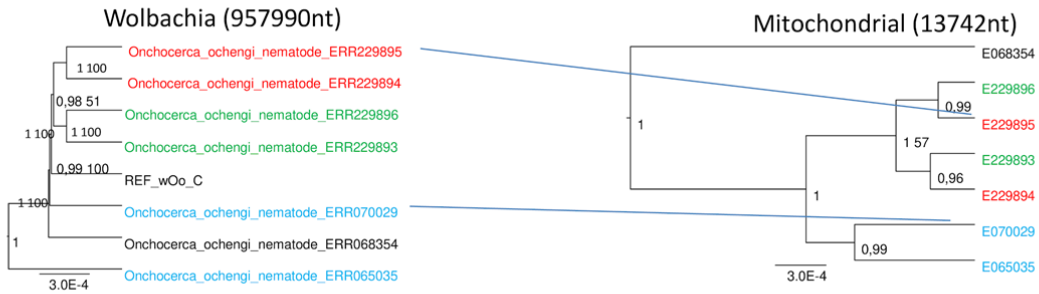
271

272

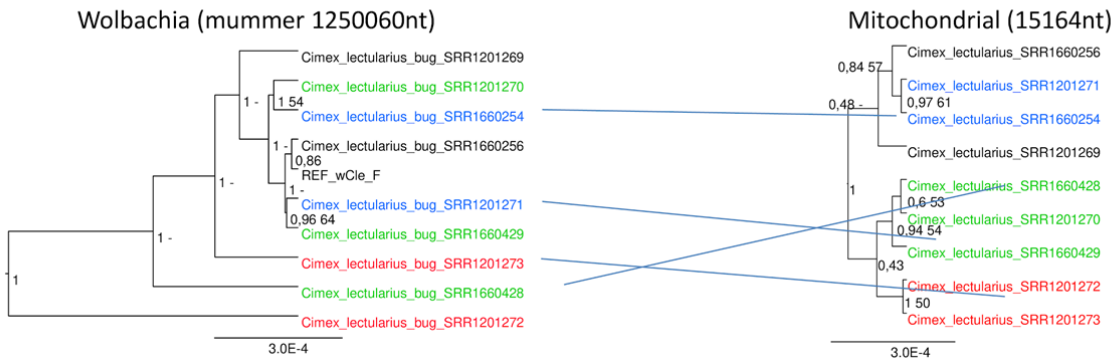
273

274

***Onchocerca ochengi*: well resolved topologies and at least 2 well supported incongruencies**

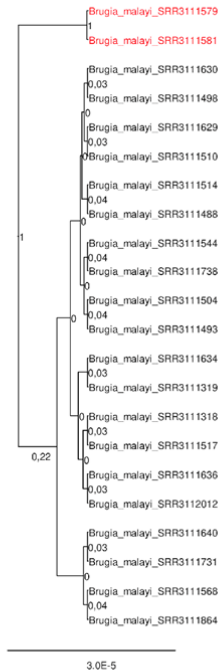


***Cimex lectularius*: partially resolved mitochondria, at least four incongruencies**

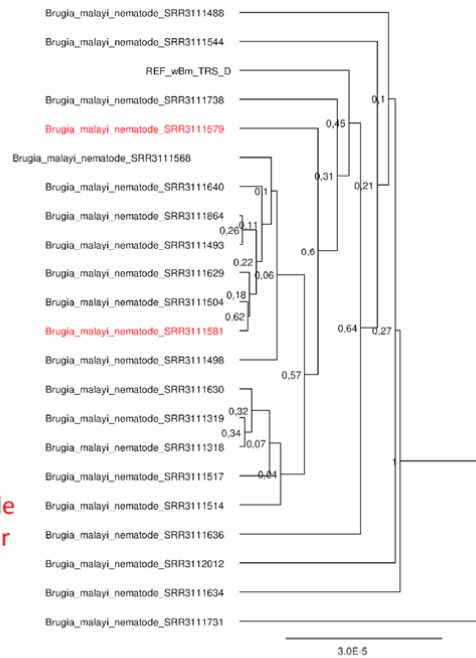


***Brugia malayi*: non significant incongruence due to unresolved mito phylogeny**

Wolbachia (1080084nt)

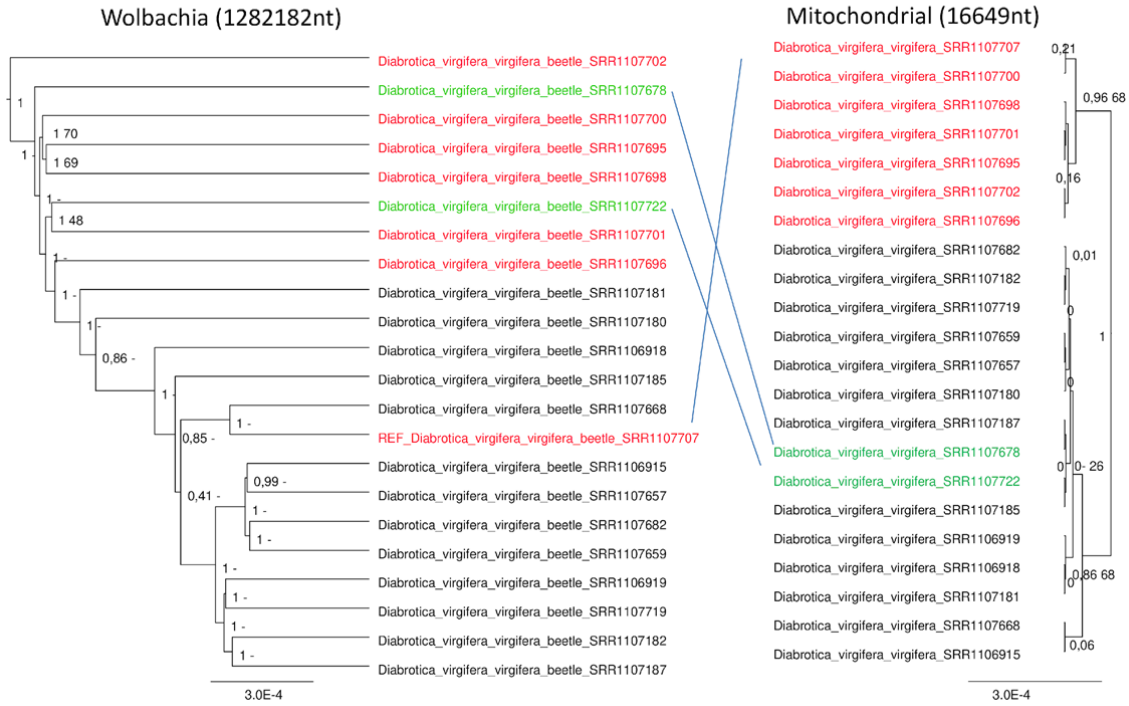


Mitochondrial (13657nt)



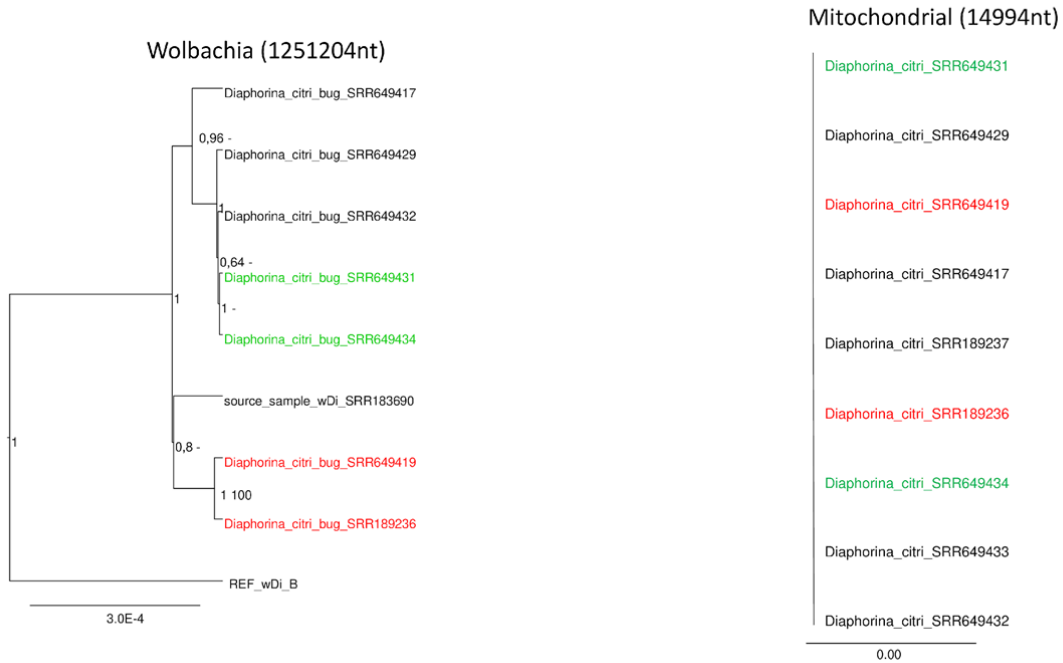
Note that in this figure genetic scale is one fold smaller than in other populations

***Diabrotica virgifera*: partially resolved mitochondria, at least three incongruences**



276

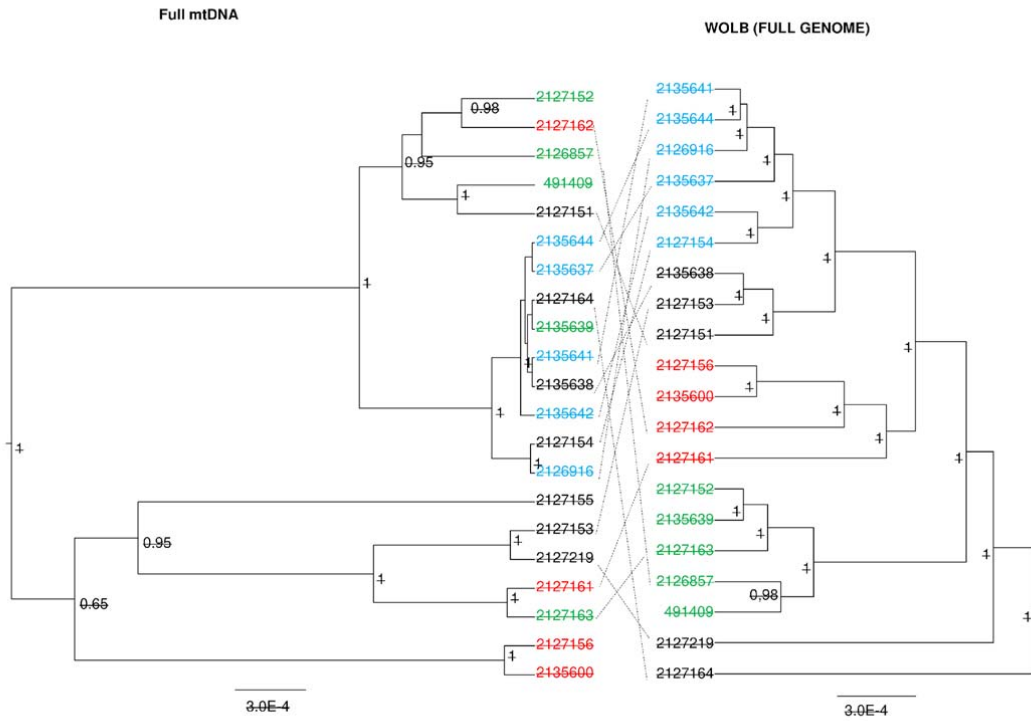
***Diaphorina citri*: unresolved (identical) mitochondria**



277

278

D. ananassae: MANY well supported incogruences.

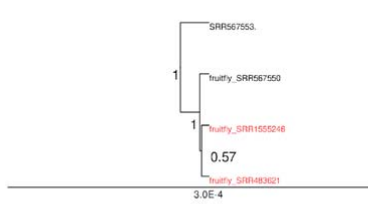


279

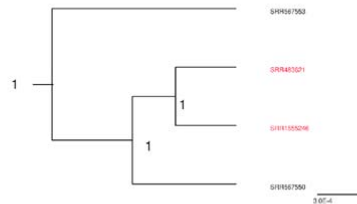
D. mauritiana: one (poorly supported) incogruences



Wolbachia (nt)



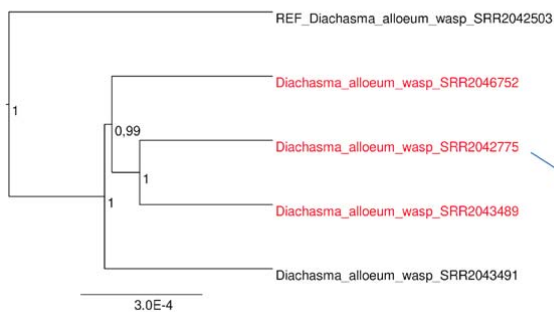
Mitochondrial (14960nt)



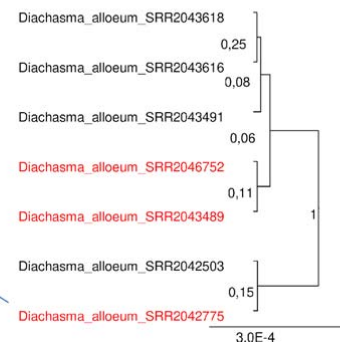
Diachasma alloeum: unresolved mitochondria, one incogruence due to well supported root



Wolbachia (1266700nt)

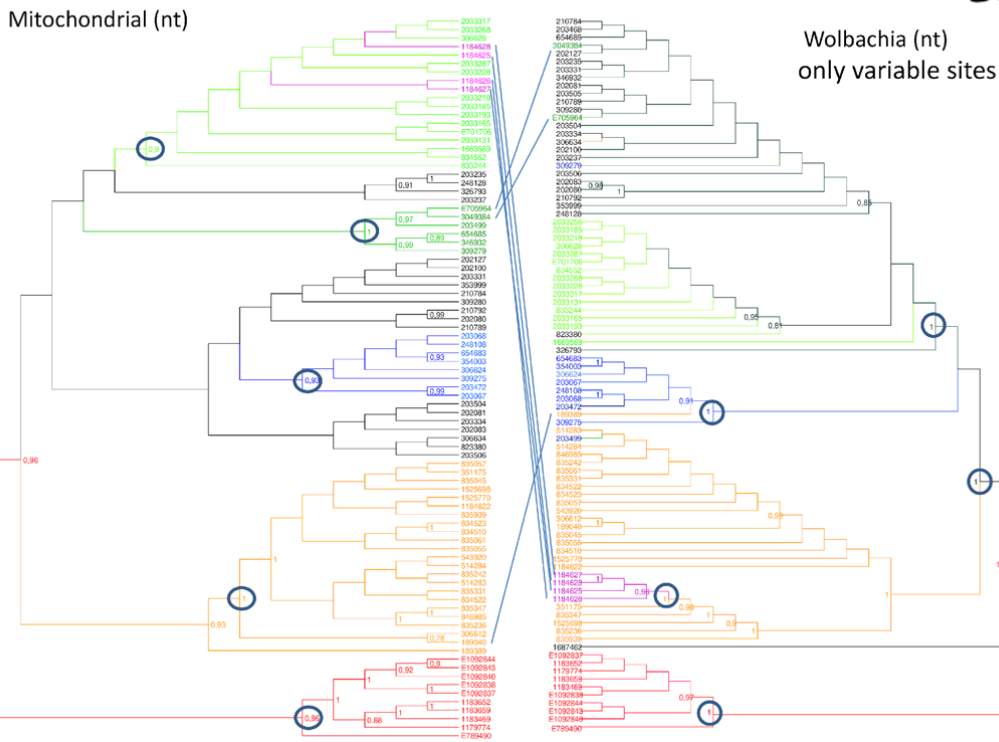


Mitochondrial (PARTIAL mtDNA)

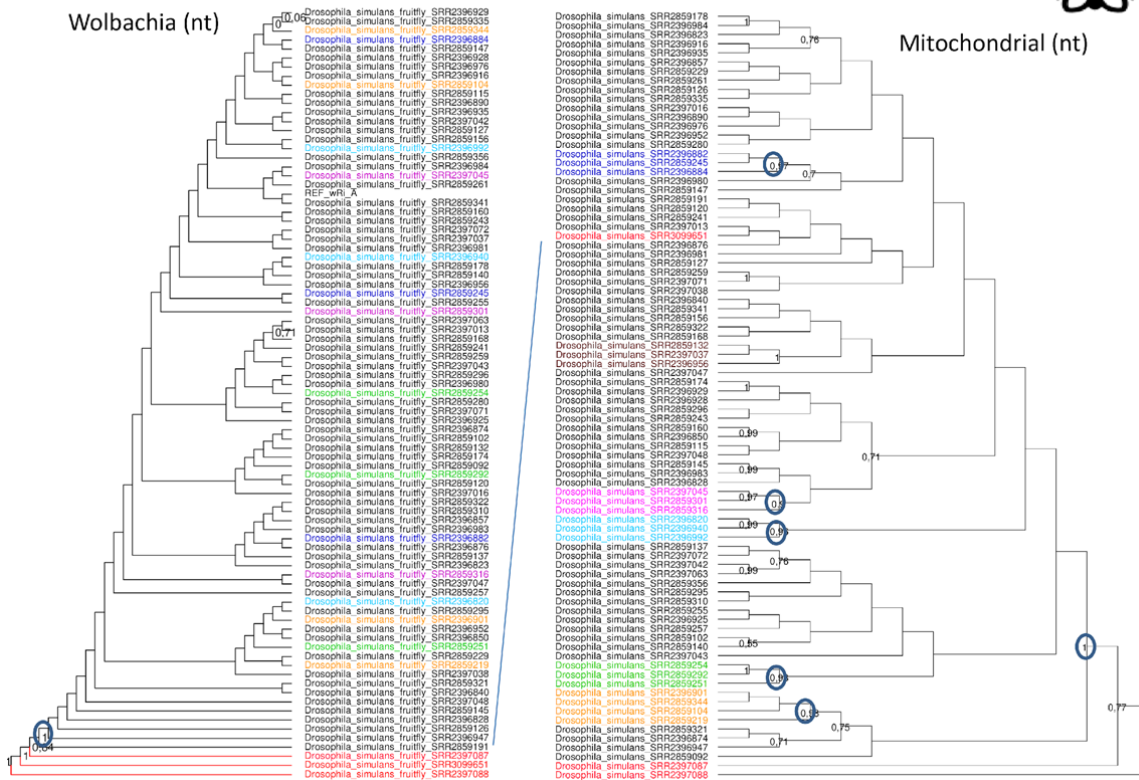


280

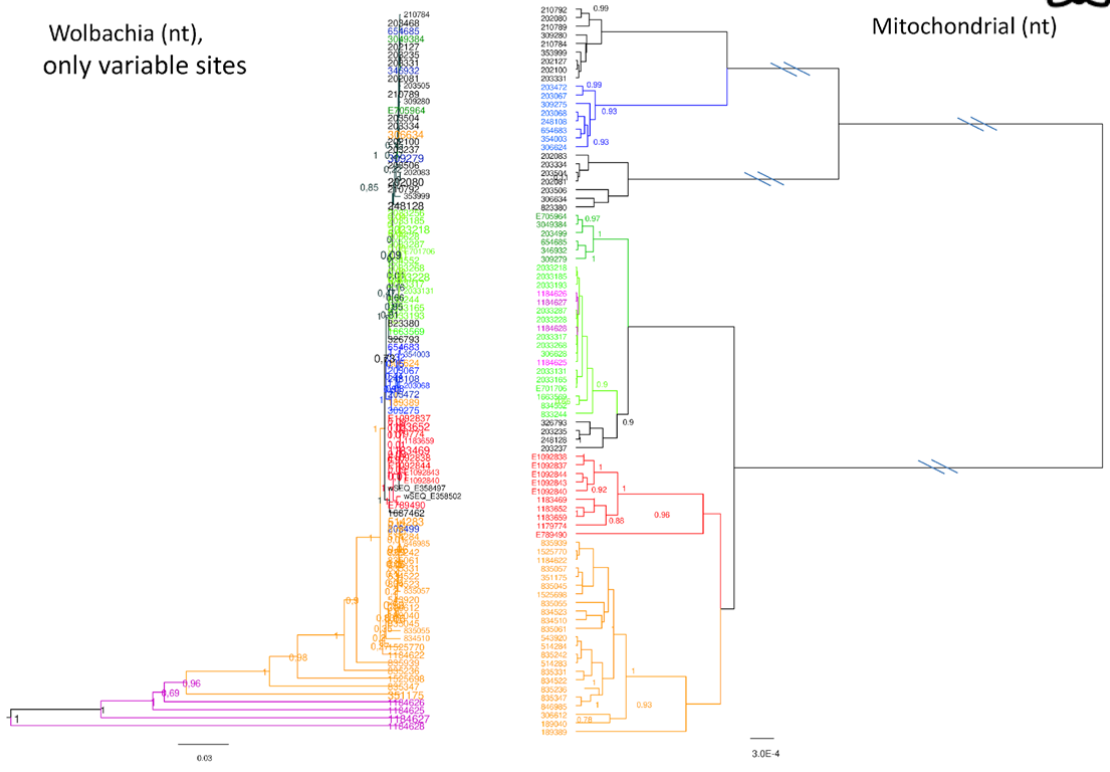
D. melanogaster: BEAST CLADOGRAM, partially resolved, but various incongruences



D. simulans: BEAST cladogram, mostly unresolved at least one incongruence



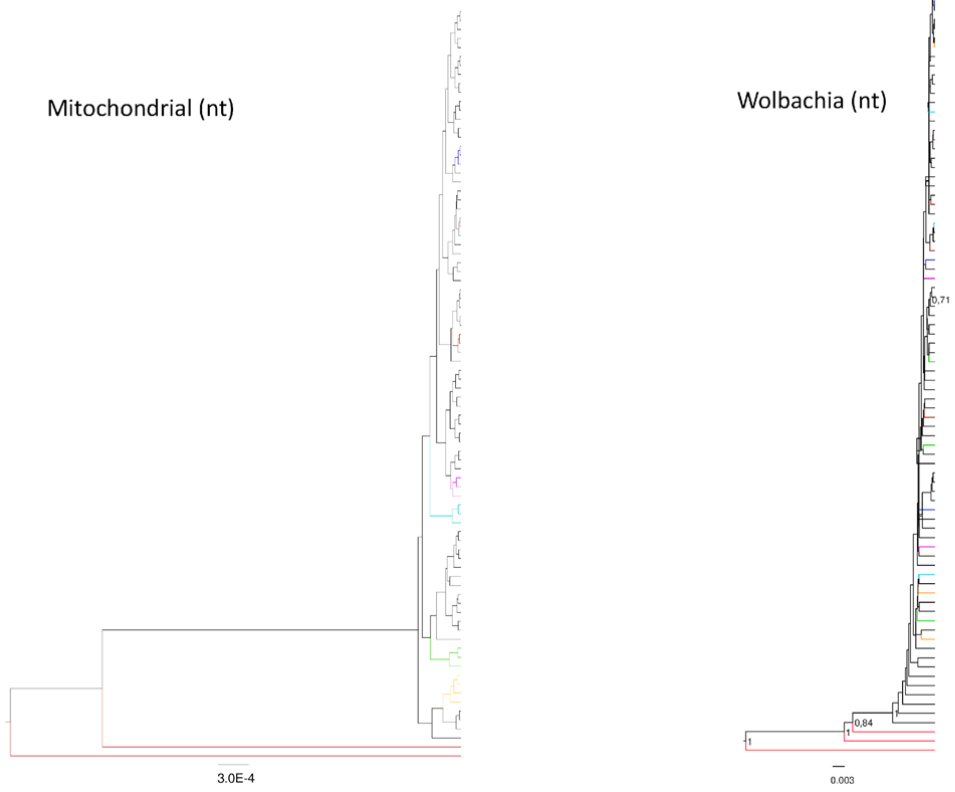
D. melanogaster: original BEAST trees with actual branch lengths.



284

285

D. simulans : original BEAST trees with actual branch lengths



286

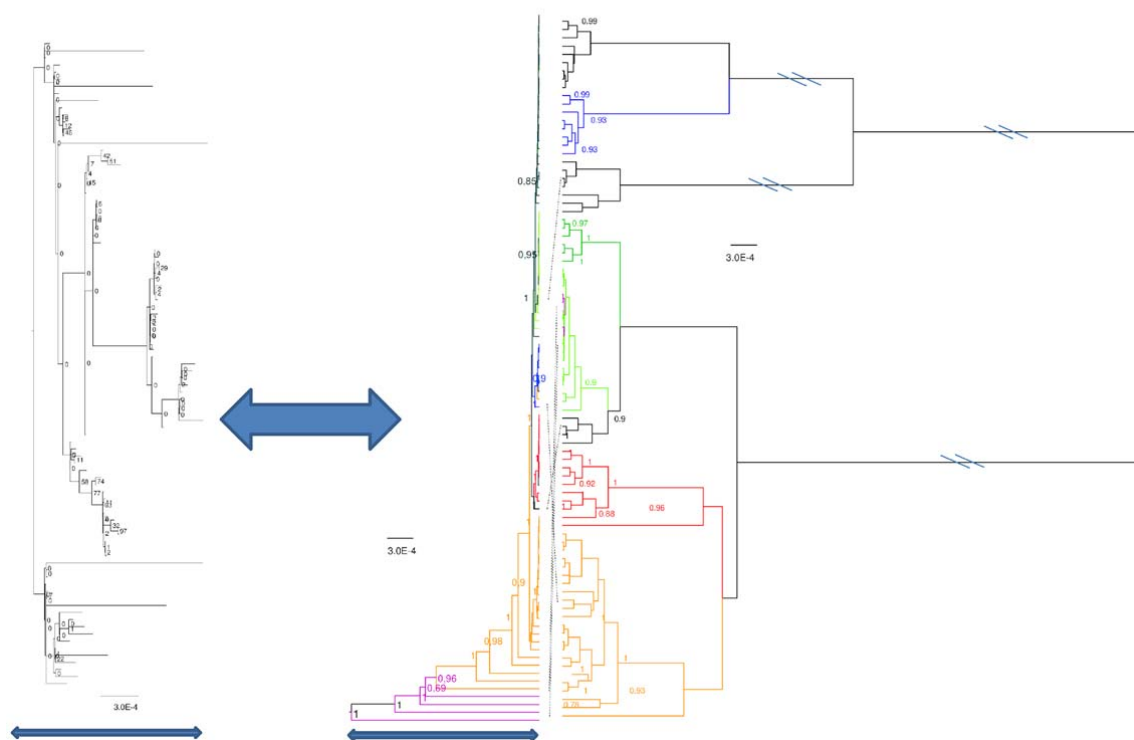
D. melanogaster, BEAST rescaling using Raxml tree length



Wolbachia RaxML tree

Wolbachia BEAST rescaled on ML tree

mtDNA BEAST tree



287

288

D. simulans, BEAST rescaling using Raxml tree length



Wolbachia RaxML tree

Wolbachia BEAST rescaled on ML tree

mtDNA BEAST tree



289