

1 An endosymbiont harvest: Phylogenomic analysis
2 of *Wolbachia* genomes from the Darwin Tree of Life
3 biodiversity genomics project

4 Emmelien Vancaester^{1*}, Mark Blaxter¹

5 ¹ Tree of Life, Wellcome Sanger Institute, Hinxton, CB10 1SA, UK

6 * Corresponding author:

7

8 Emmelien Vancaester ORCID <https://orcid.org/0000-0002-9177-8808> ev3@sanger.ac.uk

9 Mark Blaxter ORCID <https://orcid.org/0000-0003-2861-949X> mb35@sanger.ac.uk

10 Keywords: *Wolbachia* alphaproteobacteria, endosymbiosis, host-symbiont coevolution,
11 *Wolbachia* phage, cytoplasmic incompatibility

12 Abstract

13 The Darwin Tree of Life project aims to sequence all described terrestrial and aquatic
14 eukaryotic species found in Britain and Ireland. Reference genome sequences are generated
15 from single individuals for each target species. In addition to the target genome, sequenced
16 samples often contain genetic material from microbiomes, endosymbionts, parasites and other
17 cobionts. *Wolbachia* endosymbiotic bacteria are found in a diversity of terrestrial arthropods
18 and nematodes, with supergroups A and B the most common in insects. We identified and
19 assembled 110 complete *Wolbachia* genomes from 93 host species spanning 92 families by
20 filtering data from 368 insect species generated by the Darwin Tree of Life project. From 15
21 infected species we assembled more than one *Wolbachia* genome, including cases where
22 individuals carried simultaneous supergroup A and B infections. Different insect orders had
23 distinct patterns of infection, with Lepidopteran hosts mostly infected with supergroup B, while
24 infections in Diptera and Hymenoptera were dominated by A-type *Wolbachia*. Other than these
25 large-scale order-level associations, host and *Wolbachia* phylogenies revealed no (or very
26 limited) cophylogeny. This points to the occurrence of frequent host switching events, including
27 between insect orders, in the evolutionary history of the *Wolbachia* pandemic. While
28 supergroup A and B genomes had distinct GC% and GC skew, and B genomes had a larger
29 core gene set and tended to be longer, it was the abundance of active and pseudogenised
30 copies of bacteriophage WO who was a strong determinant of *Wolbachia* genome size. Mining
31 raw genome data generated for reference genome assemblies is a robust way of identifying
32 and analysing cobiont genomes and giving greater ecological context for their hosts.

33 Introduction

34 The natural world is a complex web of interactions between living species. These interactions
35 can be mutualistic, commensal, pathogenic, parasitic, predatory or inconsequential, but each
36 individual lives alongside a rich diversity of cobionts. Most eukaryotes associate intimately with
37 a specific microbiota, and are commonly infected by a range of microbial and other pathogens.
38 For some microbial associates the distinction between mutualism and pathogenicity or
39 parasitism is fuzzy. For example, *Wolbachia* (Proteobacteria; Alphaproteobacteria;
40 Rickettsiales; Anaplasmataceae; Wolbachieae) are found living intracellularly in a range of
41 terrestrial arthropods and nematodes. No free-living *Wolbachia* are known: the association is
42 essential for their survival. In contrast, infection with *Wolbachia* can be beneficial to hosts, but
43 is not usually essential.

44 *Wolbachia* were first identified as mosquito endobacteria that were maternally transmitted,
45 through the oocyte, and that induced a range of reproductive manipulations on their hosts^{1,2}.
46 The most common manipulation by *Wolbachia* is to induce cytoplasmic incompatibility (CI).
47 Under CI, infected females are able to mate productively with all males, but uninfected females
48 are only able to mate with uninfected males (as mating with CI-inducing *Wolbachia*-infected
49 males results in zygotic death). This asymmetry in fitness can drive spread of the CI-inducing
50 *Wolbachia*. Other reproductive manipulations include feminisation of genetic males³, male
51 killing⁴ and induction of parthenogenesis in females⁵. All these manipulations promote the
52 transmission of infected oocytes to the next host generation, and thus boost the spread of
53 *Wolbachia*. In most species that can be infected, populations are a mix of infected and
54 infection-free individuals, and hosts can evolve to resist infection^{6,7}. While *Wolbachia* are often
55 described as reproductive parasites, association with *Wolbachia* can sometimes have
56 beneficial effects, providing nutritional supplementation to phloem-feeding hemiptera⁸, and
57 enhancing host immunity to viruses and *Plasmodium* parasites⁹. Indeed the host immunity-
58 boosting phenotype may explain the initial spread of *Wolbachia* in previously uninfected
59 populations. In nematodes, elimination of *Wolbachia* induces host sterility, and antibiotic
60 treatment is an effective addition to pharmacological treatment of human-infecting, *Wolbachia*-
61 positive filarial nematodes¹⁰.

62 *Wolbachia* infection of terrestrial arthropods is very common, with nearly half of all insect
63 species predicted to be infected¹¹. *Wolbachia* can be classified using molecular phylogenetic
64 analyses into a series of supergroups^{12,13}. Supergroups C, D and J are found only in filarial
65 nematodes, supergroups E and F are found in both nematodes and insects, and supergroups
66 A, B and S (and others for which full genome data are not available) are found only in
67 arthropods. Supergroups A and B are the most common *Wolbachia* found in terrestrial insects.

68 Analysis of *Wolbachia* biology has been expanded by the determination of genome sequences
69 for many isolates. The genome sequences for *Wolbachia* from over 90 host species are
70 publicly available, and mining of host genomic raw sequence data identified a large number
71 of additional partial genomes¹⁴. This understanding, that cobiont genomes can be assembled
72 from the “contamination” present in the data generated for a target host, has been especially
73 useful for the unculturable *Wolbachia*. We now have the opportunity to survey for the presence
74 of *Wolbachia* genomes at an unprecedented scale, as the Darwin Tree of Life (DTOL) project
75 aims to sequence all described terrestrial and aquatic eukaryotic species found in Britain and
76 Ireland¹⁵. This project is using high-accuracy long read and chromatin conformation long range

77 sequencing to generate and release publicly available chromosomal genome assemblies,
78 meeting exact standards of contiguity and completeness, for thousands of protists, fungi,
79 plants and animals. Several hundred terrestrial arthropod assemblies are already available
80 (<https://portal.darwintreeoflife.org>).

81 The DToL project sequences genomes from individual, wild-caught specimens of target
82 species, and thus will also generate data for the cobiont present in each specimen at the
83 time of sampling. Where possible, DToL processing usually avoids body parts or tissues that
84 are expected to have a high relative mass of cobionts. In smaller-bodied species, where the
85 whole organism is extracted, and in cases where *Wolbachia* disseminates widely within an
86 organism it is inevitable that cobiont genomes will be sequenced alongside the host genome.

87 Using k-mer classification tools, it is possible to efficiently and correctly separate out cobiont
88 data from that of the host, and to deliver clean host assemblies^{16–18}. The cobiont data are then
89 available for independent assembly and analysis. Here we present a survey of the first 368
90 terrestrial arthropod genome datasets produced in DToL for the presence of *Wolbachia*, and
91 assemble over 100 new *Wolbachia* genomes. We use these to explore patterns and processes
92 in bacterial genome evolution and coevolution of *Wolbachia* with its hosts and with its own
93 bacteriophage parasites. Lepidopteran hosts were mostly infected with supergroup B, while
94 infections in Diptera and Hymenoptera were mainly caused by A-type *Wolbachia*. However,
95 host and *Wolbachia* phylogenies revealed no (or very limited) cophylogeny. We show that
96 while B genomes tended to be longer compared to supergroup A, genome size in *Wolbachia*
97 is correlated with the level of integration of its double-stranded bacteriophage WO.

98 Results

99 Screening a diverse set of insect genome data for *Wolbachia* 100 infections

101 We screened raw genomic sequence data and primary assemblies for 368 insect species (204
102 Lepidoptera, 61 Diptera, 52 Hymenoptera, 24 Coleoptera, 9 Hemiptera, 5 Trichoptera, 4
103 Orthoptera, 3 Ephemeroptera, 3 Plecoptera, 2 Odonata and 1 Neuroptera) generated by DToL
104 for the presence of *Wolbachia* ([Table S1](#)) using the small subunit ribosomal RNA (SSU rRNA)
105 as a marker gene. *Wolbachia* SSU sequences were detected in 111 (30%) of the species.
106 This degree of infection is similar to previous estimates (ranging from 22%^{19,20} to 40%¹¹ of all
107 arthropods). While the DToL project aims to sequence eukaryotes from across Britain and
108 Ireland, 82% of the samples screened were sampled from the Wytham Woods Ecological
109 Observatory, Oxfordshire (<https://www.wythamwoods.ox.ac.uk/>)²¹. No correlation between
110 sampling location and incidence level was detected, with 29% of all samples collected in
111 Wytham Woods being *Wolbachia* positive, reflective of the overall incidence level (Figure S1).

112 *Wolbachia* prevalence and infection intensity varies between species and between
113 populations within a species^{22,23}. As only one individual was analysed for each taxon screened,
114 the true level of infection within the insect biota surveyed by DToL is likely much higher.
115 Incidence was lower in Coleoptera (4/24, 17%) compared to Lepidoptera (55/204, 27%),
116 Diptera (21/61, 34%) and Hymenoptera (23/52, 44%) (Figure 1A). We observed an equal
117 prevalence of infection in samples identified as female (39/138, 28%) and male (45/153, 29%)
118 (Figure 1B).

119 The DToL species were sequenced using PacBio Sequel II HiFi highly accurate long read
120 platform, generating consensus raw reads of 10-20 kb with base level accuracy of >99%
121 (~Q30-40). These long, accurate reads are ideal for assembly, particularly for bacterial
122 genomes where the information content per base is higher than in repeat-rich eukaryotes. The
123 average sequence length of HiFi reads identified as being derived from *Wolbachia* was 12 kb,
124 indistinguishable from host HiFi reads. We separated and assembled all *Wolbachia* reads in
125 each positive sample and screened these assemblies to identify complete genomes. We
126 generated 110 complete genomes, from 93 species, of which 77 were circular ([Table S2](#)). The
127 average completeness of these genomes, assessed using BUSCO, was 99.3%, with a mean
128 duplication level of 0.37%. The mean genome size of the new genomes was 1.47 Mb, which
129 is significantly larger than the average genome size of public *Wolbachia* genomes (1.32 Mb;
130 Wilcoxon rank sum test, p-value = 4.576×10^{-9}) (Figure S2). This is likely because it is possible
131 to assemble across repeated loci (such as integrated *Wolbachia* phage) with the long,
132 accurate HiFi reads. The mean number of contigs generated for the 33 genomes that could
133 not be circularised was 2.12 (ranging from 1 to 6).

134 The dataset includes the first *Wolbachia* genomes assembled from two insect orders, Odonata
135 (dragonflies and damselflies) and Orthoptera (grasshoppers and crickets). Both species of
136 dragonfly surveyed (Odonata) harboured *Wolbachia* (Figure 1A). The largest circular
137 *Wolbachia* genome generated, 2.19 Mb, was isolated from the blue-tailed damselfly. This is
138 the longest complete *Wolbachia* genome yet reported (Figure 5A). Although in most samples
139 infection by only a single *Wolbachia* strain was detected, 15 of 93 specimens (16%) were

140 infected with at least two *Wolbachia* genomes. Within *Phalera bucephala* (Lepidoptera) and
141 *Lasioglossum morio* (Hymenoptera) three genomes were assembled, while all other co-
142 infections involved two strains.

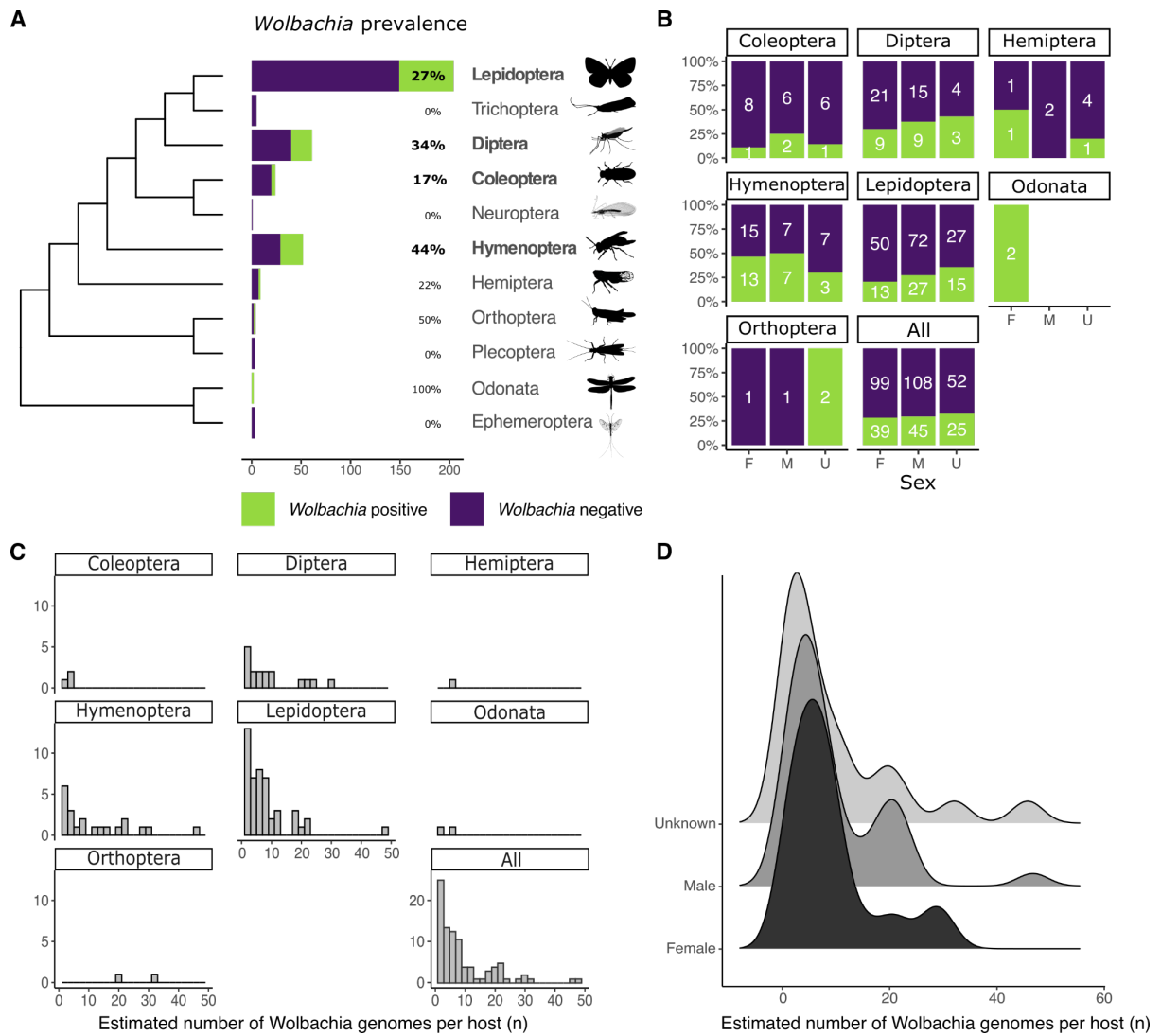
143 Having chromosomally-complete insect host genomes, as well as complete *Wolbachia* allows
144 for the estimation of the relative numbers of *Wolbachia* genomes per host genome. Most
145 infected hosts tightly control *Wolbachia* proliferation and have a relative abundance below ten
146 *Wolbachia* genomes per host nuclear genome. Particularly high abundances were observed
147 in *Thymelicus sylvestris* (48 *Wolbachia* per host) and *Athalia cordata* (47) ([Table S2](#)) (Figure
148 1C). The mean relative abundance in different taxonomic orders lay between 3 and 12, except
149 for the two crickets (Orthoptera), *Chorthippus brunneus* and *Chorthippus parallelus*, which
150 have a 33 and 20 *Wolbachia* genome copies per host genome, respectively (Figure 1C). No
151 significant difference was observed between relative *Wolbachia* abundance and sex of the
152 host (Figure 1D), with both male and female having a mean between nine and ten copies.

153 **Figure 1: Prevalence and relative abundance of *Wolbachia* in**
 154 **DTOL insect genomes.**

155 **A, B** Prevalence of *Wolbachia* in insect hosts, split by taxonomic order (A) and by sex (B). The
 156 cladogram of insect ordinal relationships is based on Misof et al²⁴. Orders with more than 10
 157 analysed species are shown in bold. Silhouettes are from PhyloPic (<http://phylopic.org/>). Sex
 158 of insects was classified as F (female), M (male) or U (unknown, where not recorded on
 159 collection).

160 **C, D** The estimated number of *Wolbachia* genomes per copy of the host nuclear genome split
 161 by taxonomic order (C) and by sex (D).

162



163

164 *Wolbachia* phylogeny suggests frequent host switching events

165 We selected 93 high-contiguity and high-completeness *Wolbachia* genomes from the public
166 INSDC databases, including genomes from *Wolbachia* infecting Nematoda (13 genomes),
167 Arachnida (4), Isopoda (1) and several orders of Hexapoda (75) ([Table S3](#)). Adding the 110
168 newly assembled genomes yielded a dataset of over 200 high-quality assemblies. We
169 annotated all protein-coding genes in those genomes using Prodigal²⁵, and clustered the
170 predicted protein sets into orthologous groups using OrthoFinder2²⁶. The resulting 634 near-
171 single copy genes were used to infer a phylogeny of *Wolbachia* (Figure 2A, Figure S3). From
172 this phylogeny we assigned each genome to the previously defined *Wolbachia*
173 supergroups^{12,13}. All newly assembled *Wolbachia* genomes belonged to either supergroup A
174 or B. While Lepidoptera were predominantly infected with supergroup B *Wolbachia* (42/53,
175 80%), *Wolbachia* supergroup A was most frequent in all other insect classes (46/57, 81%). It
176 has been previously observed that supergroup B is the most common *Wolbachia* type in
177 Lepidoptera^{19,27–29}. Of the 15 species where co-infections occurred, *Endotricha flammealis*,
178 *Phalera bucephala*, *Philonthus cognatus*, *Protocalliphora azurea* and *Sphaerophoria taeniata*
179 were co-infected with strains from both A and B supergroups, and the other ten co-infections
180 were of distinct strains within the same supergroup ([Table S2](#)).

181 *Wolbachia* generally do not show strict cophylogeny with their hosts^{7,23}. This pattern was also
182 observed when comparing host and *Wolbachia* phylogenies for the supergroup A and B
183 genomes (Figure 2B). Closely related insect species may be infected by dissimilar *Wolbachia*
184 strains and conversely, closely related *Wolbachia* can infect a diverse set of insects. For
185 example, the *Wolbachia* strains infecting the hoverfly *Eupeodes latifasciatus* and four
186 Lepidoptera (*Pararge aegeria*, *Celastrina argiolus*, *Hylaea fasciata*, *Watsonella binaria*)
187 (Figure 2C) share over 99% nucleotide identity. Because most of our new samples came from
188 a single site (Wytham Woods Genomic Observatory) we were also able to explore the
189 horizontal transfer of *Wolbachia* between hosts in a local context. Wytham Woods-derived
190 *Wolbachia* were no more likely to be related than any other *Wolbachia* subset (Figure S4).

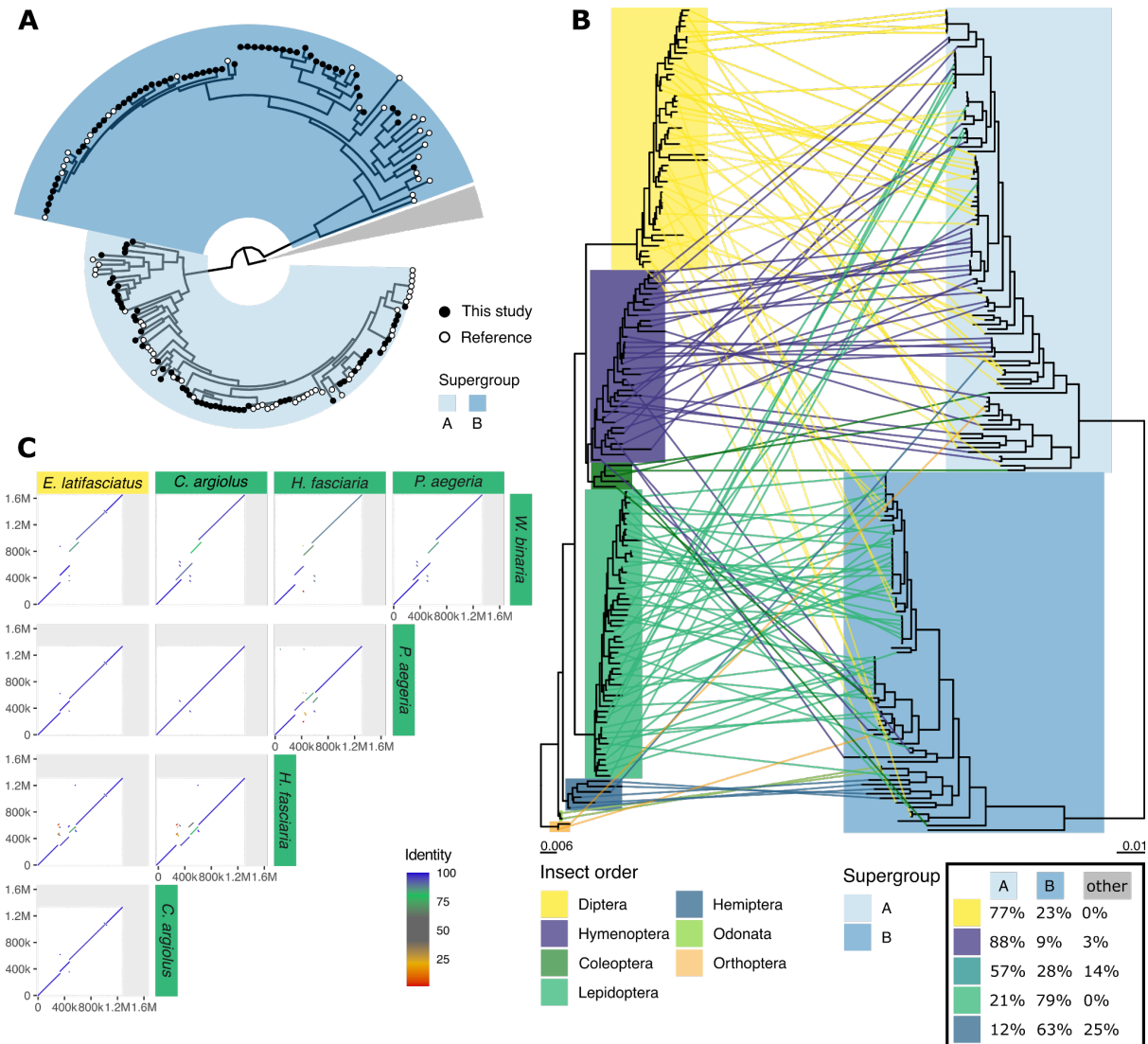
191 **Figure 2: *Wolbachia* DTOL genomes expand known phylogeny.**

192 **A** Circular phylogeny of supergroup A and B *Wolbachia*, visualised with the root placed
 193 between the A and B supergroups and the remaining supergroups (C,D,E,F, J, S; nodes
 194 collapsed as grey wedge), highlighting newly sequenced genomes (black tip labels) and
 195 genomes from public databases (white).

196 **B** Incongruence between host topology (left) and supergroup A and B *Wolbachia* topology
 197 (right) is shown as a tanglegram. Overview of the supergroups infecting diverse insect orders
 198 is given in a table (inset, bottom right).

199 **C** Example of a host switching event, where the *Wolbachia* of the hoverfly *Eupeodes*
 200 *latifasciatus* has high nuclear sequence identity and genome colinearity to four *Wolbachia*
 201 genomes assembled from Lepidoptera.

202



203

204 Intrinsic properties of *Wolbachia* distinguish supergroups

205 The completeness of the new genomes, and in particular the circular assemblies achieved for
206 77 of them, permits analyses of genome properties that are not possible with fragmented and
207 partial genomes. All circularised genomes, including those from public databases, were
208 rotated to start at the presumed origin of replication. The average pairwise whole genome
209 nucleotide identity between all *Wolbachia* genomes ranged between 77.3% and 100.0%, with
210 at least 92.8% and 93.5% identity within supergroup A and B, respectively (Figure 3A). The
211 number of breakpoints interrupting pairwise whole-genome alignments was counted,
212 normalised for the total alignable length, and compared to average nucleotide identity (ANI) of
213 the compared genomes (Figure 3A). A significant correlation was observed between
214 nucleotide divergence and the number of breakpoints in supergroups A (0.90, $p < 2.2e^{-16}$,
215 Spearman correlation) and B (0.69, $p < 2.2e^{-16}$, Spearman correlation) (Figure 3A).

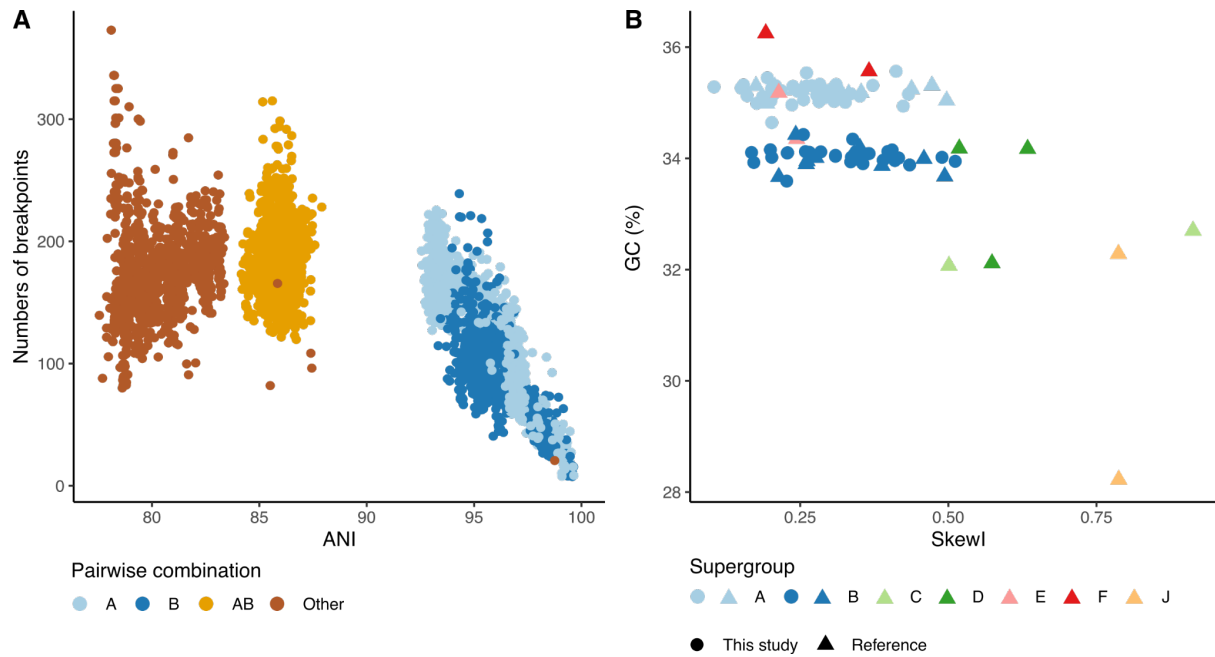
216 GC skew accumulates in stable bacterial genomes through differential mutation pressures on
217 leading versus lagging strands. Genomes that have undergone frequent rearrangement are
218 expected to have lower overall GC skew, which can be summarised across the genome as a
219 single metric, Skewl³⁰. Genomes from supergroups A and B had distinct GC contents (Figure
220 3B), with supergroup A having a higher mean GC (35.2%, standard deviation 0.15%),
221 compared to B (34.0%, standard deviation 0.16%) (two-sample t-test p-value $< 2.2e^{-16}$).
222 Genomes from other supergroups had distinct GC content, often very different from A and B
223 genomes, but as so few examples have been sequenced, general patterns are not discernible.
224 In both A and B supergroups Skewl values were relatively low, but genomes from *Wolbachia*
225 from nematode hosts (C, D, J) had higher Skewl values (Figure 3B). A high degree of GC
226 skew had already been observed in the *Wolbachia* strain infecting *Dirofilaria immitis*
227 (supergroup C)³¹. This suggests that nematode-associated *Wolbachia* have retained
228 chromosome stability across a long timeframe, not observed in supergroups A and B, also
229 evident by their large number of re-arrangements (Figure 3A).

230 **Figure 3: Comparative genomics of *Wolbachia*.**

231 **A** Average nucleotide identity (ANI) plotted against the number of breakpoints in comparisons
232 within A supergroup genomes, within B, between A and B and between other supergroup
233 *Wolbachia*.

234 **B** Index of skewness compared to GC content for all circularised *Wolbachia* genomes.

235



236

237 Conservation and diversity in gene content of *Wolbachia*

238 *Wolbachia*, because they are sheltered within the cells of their hosts, may be relatively isolated
239 from other bacteria, and thus have somewhat closed pan-genomes. One route to acquisition
240 and sharing of new genes is through the *Wolbachia* phage (WO phage), which alongside the
241 essential phage particle structural genes carry a cargo of genes that have been implicated in
242 host manipulation. We re-annotated all 203 *Wolbachia* with the same, standard gene finding
243 toolkit, Prodigal, to normalise annotations. While this may have lost careful manual revision in
244 previously determined gene sets, it avoids issues of data incompatibility. Gene number
245 correlated with genome size, and the average gene number in the newly assembled set of
246 supergroup A and B *Wolbachia* was larger than in A and B genomes from the public databases
247 (Figure S5). Comparing all genomes, the mean number of predicted genes was larger in
248 supergroup B (1467) compared to A (1385).

249 We used OrthoFinder with default settings to define clusters of orthologous proteins across all
250 *Wolbachia* genomes. Each genome contained between 0 and 184 novel, strain-specific genes
251 (average 19). These novel genes were shorter than all genes (average gene length overall
252 was 875 nucleotides or ~290 amino acids, while novel genes averaged 434 nucleotides or
253 ~145 amino acids). As expected, supergroups which were not well represented often
254 contained more strain-specific genes. For example, wCfeT from supergroup E (which infects
255 cat fleas, *Ctenocephalides felis*) uniquely encoded genes for pantothenate (panC-panG-panD-
256 panB)³² and thiamine (thiG-thiC) biosynthesis. Nonetheless, out of the ten genomes with most
257 strain-specific genes, seven belonged to either supergroup A or B. These novel genes were
258 not preferentially associated with WO phage regions (Figure S6) but the majority (78%) had
259 annotations that associated them with transposon and mobile element function. This suggests
260 that much of the novelty arose through mobile elements other than WO phage. Other than
261 clusters with one or two members, the most frequently observed cluster sizes were 203±2.
262 These clusters contained the single-copy (and near-single-copy) orthologs deployed in
263 phylogenetic analyses (Figure 4A). Overall, the majority of the proteins encoded in the
264 *Wolbachia* genomes were members of orthology clusters that were present in at least 95% of
265 all strains.

266 The abundant sampling of supergroup A and B genomes allowed us to address and compare
267 the sizes of the core- and pan-proteomes of these groups. The larger genome and proteome
268 size found in supergroup B was reflected in a larger core proteome (Figure 4B), but supergroup
269 A had a larger pan-proteome (Figure 4B). While the core proteomes differed, very few of the
270 protein families that were part of each supergroup's core proteome were unique to that
271 supergroup. One supergroup-restricted set of protein families was found to comprise the
272 operon for arginine transport (ArtM, ArtQ and ArtP and the repressor of arginine degradation
273 ArgR)³³, which was uniquely detected and conserved in supergroup A (present in 83/103 or
274 80% of all *Wolbachia* A genomes). Although the periplasmic arginine-specific binding protein
275 (ArtI or ArtJ) was not detected, the presence of this ATP-binding cassette-type (ABC)
276 transporter suggests that these *Wolbachia* are acquiring arginine from their hosts.

277 The operon producing biotin (vitamin B7)³⁴ was detected in seven of the 110 new genomes,
278 all belonging to supergroup A (Figure 4C). One derived from *Icerya purchasi* (Hemiptera) and
279 six were from Hymenoptera (two from *Lasioglossus malacharum*, which carried two strains,
280 and single strains from three *Andrena* and a *Nomada* species). The biotin synthesis cluster
281 has been described previously from a restricted but diverse set of supergroups, including two

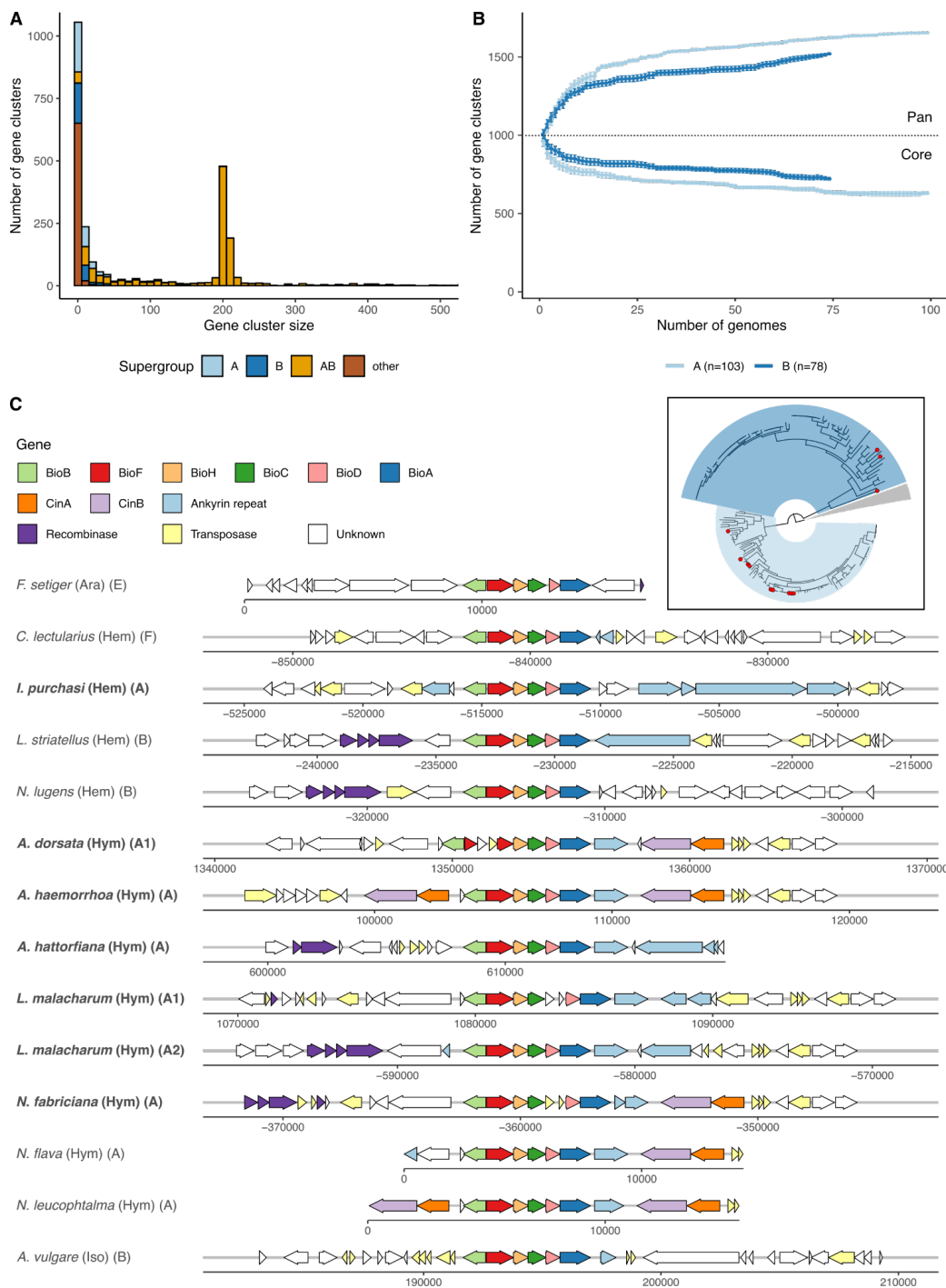
282 A genomes from additional *Nomada* bee hosts. This distribution suggests possible ecological
283 linkage³⁵, as *Andrena* bees are kleptoparasitized by *Nomada* cuckoo bees and phylogenetic
284 analyses of both the biotin gene clusters and the *Wolbachia* core proteomes show close
285 relationships between this cluster of genomes (Figure S7). The gene cluster is strongly
286 conserved in physical organisation of all six necessary genes (bioA-D,F,H). In the genomic
287 region immediately surrounding the operon we identified recombinase and transposase
288 genes, as well as ankyrin repeat containing genes and toxin-antitoxin cytoplasmic
289 incompatibility Cin gene pairs. In three genomes (from *Andrena dorsata*, *Nomada fabricium*
290 and one of the *L. malacharum* strains) the operon was independently disrupted by
291 transposases. The region containing the biotin operon thus has the hallmarks of a “virulence
292 island” that may be mobile between genomes, and may have accrued additional genes
293 (ankyrin, Cin) that hitchhike with the biotin operon.

294 **Figure 4: Exploration of *Wolbachia* protein-coding gene**
 295 **diversity.**

296 **A** Histogram of protein family size per supergroup.

297 **B** Rarefaction analysis of pan- and core proteomes of supergroups A and B, based on 500,000
 298 random addition-order permutations of co-occurring orthogroups excluding novel genes.

299 **C** Synteny of the biotin cluster shows conserved gene order and punctuated pattern of species
 300 presence (inset, species with biotin cluster present are highlighted with red circles).



301

302 WO prophage insertions expand genome size

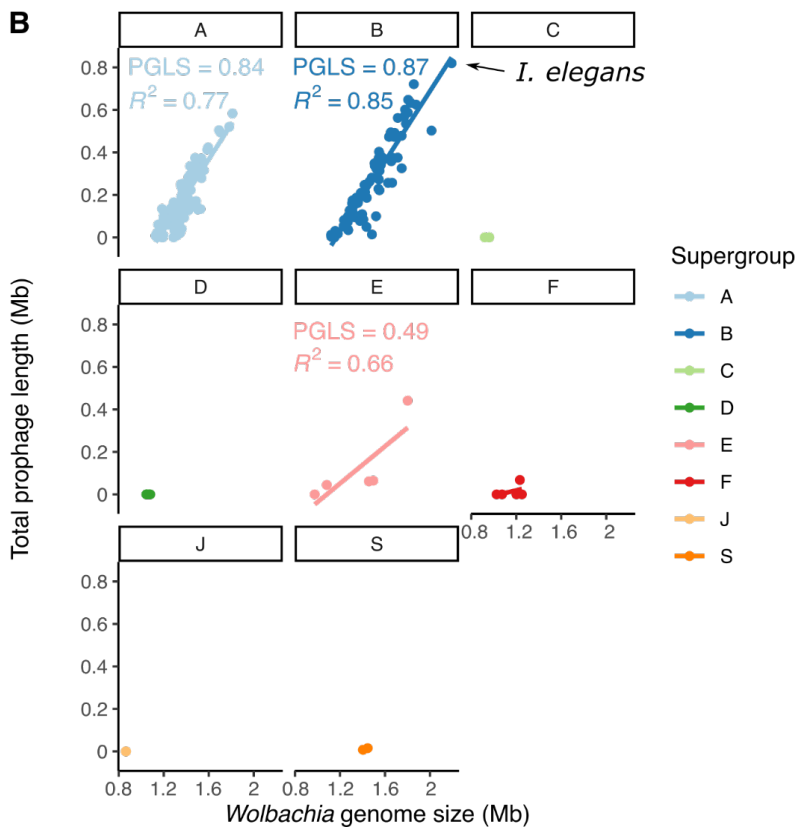
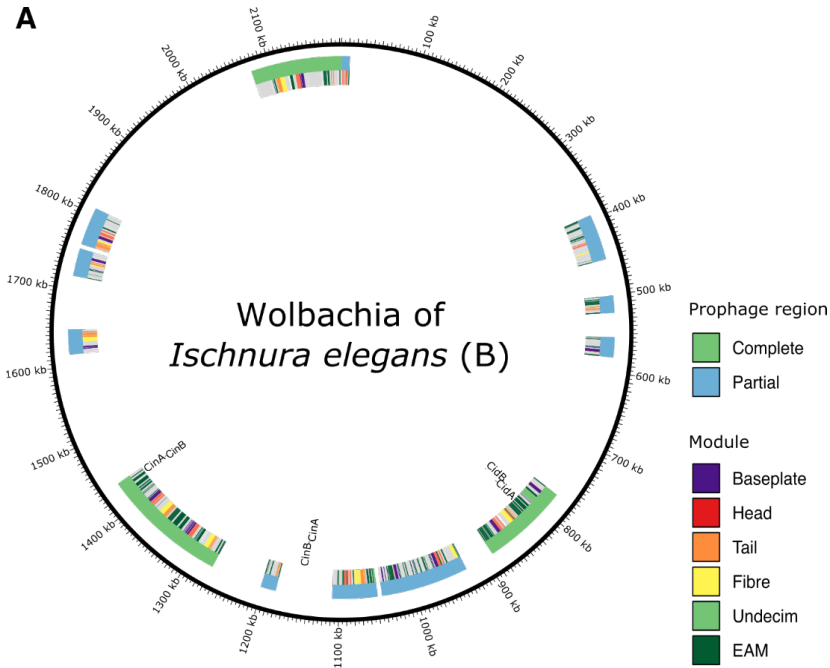
303 *Wolbachia* can itself be infected by double-stranded DNA temperate bacteriophages, WO
304 phage, which can integrate in the genome of its host as a prophage. Four modules are
305 necessary for construction and function of phage particles during the lytic stage: head,
306 baseplate, tail and fibre, and inserted and pseudogenised WO phage can be identified and
307 discriminated based on the presence and completeness of these components. Regions of a
308 *Wolbachia* genome flanked by WO phage modules are likely to form components that are
309 transduced by the phage during infection of new cells, “cargo” loci that form the Eukaryotic
310 Association Module (EAM)^{36,37}. All the *Wolbachia* genomes were screened for prophage
311 regions using essential module genes from previously annotated WO insertions ([Table S4](#)).
312 Prophage regions were deemed complete when all four modules were observed with at least
313 80% of genes of each module present. An abundance of putative intact and pseudogenised
314 WO phage were identified. For example the supergroup B *Wolbachia* from *Ischneura elegans*
315 (the bluetail damselfly; the largest *Wolbachia* genome assembled) contained three putative
316 intact prophage and nine WO phage fragments (Figure 5A) summing to 0.8 Mb of the genome.

317 The fraction of total prophage region in each genome ranged from 0-38%. Nematode-
318 associated *Wolbachia* typically are not infected by WO phage³⁸ and no prophage regions were
319 detected in genomes of supergroup C, D, J and nematode-infecting F (Figure 5B). A significant
320 correlation was found between genome size and WO prophage span in supergroups A and B
321 (Figure 5B). This association was robust to correction for phylogenetic relatedness of the
322 genomes (model fit increased to 0.84 and 0.87 respectively with p-values <10⁻¹⁶).

323 **Figure 5: WO prophage in *Wolbachia***

324 **A** Annotation of the WO prophage integrated in the genome of the *Wolbachia* strain infecting *Ischnura elegans*.

326 **B** *Wolbachia* genome size is strongly correlated with integrated prophage span in supergroups with WO
 327 phage association. Phylogenetic Generalised Least Squares (PGLS) analyses were performed to
 328 assess the correlation between prophage length and genome size in a phylogenetically-aware manner.



329

330 Toxins are often associated with mobile elements

331 We identified several potential cargo genes within intact and fragmented prophage. These
332 included transposases and integrases associated with mobile elements, and other loci
333 previously associated with eukaryotic manipulation, such as cytoplasmic incompatibility loci
334 and ankyrin repeat containing genes, as expected from the EAM model^{36,37}.

335 *Wolbachia* produces a suite of toxins³⁹ which can have dramatic effects on their hosts, such
336 as cytoplasmic incompatibility (CI). The CI phenotype is caused by two adjacent genes, CifA
337 and CifB, which function as a toxin-antitoxin pair^{40,41}. Phylogenetic analysis classified most
338 *Wolbachia* Cif gene pairs into four types (I-IV)⁴². A fifth type (V) is much more variable in
339 structure. The toxin component can have nuclease activity (in which case the gene pair is
340 frequently referred to as CinA-CinB), deubiquitinase (CidA-CidB) or both (CndA,CndB)⁴³. All
341 type II, III and IV pairs have nuclease domains, while all type I have deubiquitinase and most
342 have nuclease⁴². Two hundred and sixty one full-length and likely functional Cif pairs were
343 detected in 133 of the 181 (73%) supergroup A and B genomes. One Cif pair was detected in
344 most genomes, but many had several, with seven copies in the *Wolbachia* strain infecting the
345 holly tortrix moth (*Rhopobota naevana*). Most of the gene pairs (93) contained a
346 deubiquitinase domain (type I, Cid), while the other four types occurred in roughly equal
347 proportions (II: 40, III: 43, IV:35 and V:50). Many pairs (177/261; 69%) were located in the
348 predicted EAM of the prophage.

349 Loci encoding additional toxins such as RelE/RelB, Spaid-like and latrotoxin were identified in
350 multiple *Wolbachia* genomes, frequently in prophage regions (71/586 [12%], 136/597 [23%]
351 and 227/382 [59%] genes, respectively). The Tc pore-forming toxin complex, which consists
352 of two genes TcA and TcB-C, was detected in a limited number of A and B supergroup
353 genomes and also showed a predisposition to occur within prophage (13/69 [19%] and 9/35
354 [26%], respectively). Additional toxin-encoding loci had limited presence in different subgroups
355 and were not associated with prophage regions. ParD/ParE only occurred in supergroups A,
356 B and E, and FIC only occurred in supergroups A, E, F and S. The type IV toxin-antitoxin gene
357 pair AbiEii/AbiGii-AbiEi, which protects against the spread of phage infection⁴⁴, was only
358 detected in two genomes in supergroup E. It is noteworthy that these two genomes had very
359 low levels of prophage-derived DNA (4.3% of their genome span).

360 Discussion

361 Isolation of cobiont genomes, and specifically *Wolbachia* genomes, from shotgun high-
362 throughput sequencing data has been established for many years⁴⁵. In the field of prokaryotic
363 and eukaryotic microbial metagenomics, metagenome-assembled genomes (MAGs) are likely
364 to be the only way to access many unculturable microbial genomes, even if the species they
365 derive from are hyperabundant^{46,47}. The abundance of raw sequencing data in the
366 International Nucleotide Sequence Database Collaboration (INSDC) databases has been an
367 attractive prospecting ground for microbial associates of eukaryotic target species. To date,
368 most raw data available for such searches have been short reads from Illumina and other
369 platforms. These reads are too short to partition efficiently into bins corresponding to putative
370 distinct genomes. Preliminary assembly of such datasets is more likely to be able to separate
371 cobionts from target genomes. These approaches have been applied to hunt for *Wolbachia*
372 with a recent *tour de force* generating nearly 1,200 *Wolbachia* MAGs from publicly available
373 data¹⁴. However, these MAGs suffer from the expected issues of low completeness (due to
374 low effective coverage), fragmentation (due to coverage and sequence repeat issues),
375 undetected contamination and inability to distinguish co-infecting strains. Moreover, the biased
376 nature of public data meant that these derived from only 37 different host species.

377 We generated 110 *Wolbachia* assemblies from 368 terrestrial arthropod HiFi datasets, and 77
378 of these were fully circular genome assemblies. The genomes were uniformly of high
379 completeness (Figure S2). Due to the high intrinsic base quality of HiFi reads (Q30 to Q40;
380 from one error in 1000 to one error in 10,000) we were able to distinguish insertions of
381 *Wolbachia* DNA into the host genome from true components of the *Wolbachia* genome, and
382 to independently assemble even closely related strains with confidence. As we were screening
383 raw data from a biodiversity genomics programme that aims to sample a wide phylogenetic
384 diversity of hosts, the new *Wolbachia* genomes presented here more than double the number
385 of different host species from which *Wolbachia* genomes have been assembled. The
386 assembled genomes include the first representatives isolated from Odonata (damselflies) and
387 Orthoptera (crickets). In 16 additional datasets we identified likely *Wolbachia* content but were
388 not able to produce credible genome assemblies (see Supplemental Data, Table S2). This
389 was usually because the *Wolbachia* sequence was present in very low effective coverage (~
390 threefold) but in some samples no credible assembly was generated despite high coverage.
391 These datasets may contain multiple recombining strains, or contain large insertions in the
392 host genome and deserve further exploration.

393 The distribution of *Wolbachia* in insect hosts is a function of the balance between co-speciation
394 (vertical transmission of *Wolbachia* among daughters of the host species) and horizontal
395 transmission where strains move between species. Transmission among insect hosts was the
396 dominant pattern underpinning *Wolbachia* distribution, but we identified two features of the
397 distribution, one local and one general, that are of note. Lepidoptera were more likely to be
398 infected with supergroup B *Wolbachia* than A, and Hymenoptera, Diptera and Coleoptera were
399 more likely to be infected with supergroup A strains. Multi locus sequence typing (MLST) has
400 previously shown that supergroup B is the most common *Wolbachia* type in Lepidoptera^{19,27-}
401 ²⁹. This suggests some non-exclusive specialisation of *Wolbachia* on their hosts, which may
402 be driven by *Wolbachia* genetics, host genetics or (less likely) a distinct set of ecological
403 transmission routes in each insect group. Many of our genomes derived from insects were
404 collected at one site, the Wytham Woods Genomic Observatory (Figure S1) but this subset

405 was no more closely related than other genomes from widely separated sites (Figure S5). It is
406 likely that the mobility of hosts, including through seasonal migration, means that sampling
407 from one geographical site is a valid approximation of more global sampling.

408 Close ecological association between host species may promote sharing of *Wolbachia*
409 isolates and localised genetic exchange, for example within predator-prey systems. The close
410 similarity of *Wolbachia* genomes from *Andrena* solitary bees and their *Nomada* cuckoo bee
411 kleptoparasites, and the shared occurrence of the biotin synthesis operon (Figure 4C) may be
412 a case of transmission within an ecological network. The presence of the biotin operon in
413 *Wolbachia* of insects that largely or solely feed on low-protein plant fluids (nectar or phloem)
414 suggests that the *Wolbachia* may be offering nutritional support to their hosts⁴⁸, and thus that
415 this cluster of genomes may have been positively selected for their mutualist tendencies. Other
416 genes whose distribution among isolates is driven by horizontal gene transfer, including by
417 mobile elements such as phage, might be expected to have a distribution that is not explained
418 by overall genome relatedness, and might reflect ecological association. We note that previous
419 work has suggested that horizontal transmission rather than cospeciation may also explain
420 closely related *Wolbachia* in closely related taxa. For example, genomic divergence between
421 closely related *Wolbachia* in sister *Drosophila* species was too low to be the product of
422 independent evolution since the last common ancestor of the flies^{49,50}.

423 *Wolbachia* can promote reproductive success of their hosts^{1,2}, and thus their own Darwinian
424 fitness, through reproductive manipulations such as CI. The loci underpinning CI are a diverse
425 set of toxin-antitoxin gene pairs. Our survey of *Wolbachia* identified many additional CI gene
426 pairs, mainly of the I Cid type and mostly associated with WO phage. Many genomes had
427 more than one toxin-antitoxin pair, and some individual hosts were infected with multiple
428 *Wolbachia* strains carrying different CI gene pairs. These CI genes likely mediate conflict
429 between *Wolbachia* strains and the ecosystem of toxin deposition and rescue in individual
430 zygotes must be complex. Interestingly we identified CI gene pairs next to 5 of the 14 biotin
431 synthesis operons, suggesting that the mobile elements that transduce this presumably
432 mutualist physiology are also engaged in CI conflict.

433 One striking feature of the genomes assembled from the HiFi reads was that their average
434 span was ~10% greater than the average size of previously-assembled *Wolbachia* genomes.
435 As we also observed a correlation between content of WO phage in the genome and genome
436 size (Figure 5B), we speculate that the lower average size of previous assemblies may be
437 because the presence of near-identical segments of phage and other mobile elements led to
438 collapse of repeats and artificial underestimation of true genome size. This underestimation of
439 genome size may also have biased understanding of WO phage diversity and of the diversity
440 of genes that can be transduced by the phage. WO phage carry genes necessary for
441 production of phage particles and cargo genes that have been hypothesised to form an
442 Eukaryotic Association Module (EAM)^{36,37}. The increased genome size and increased
443 resolution of WO phage copies might also mean increased gene content and diversity, and an
444 increased set of common EAM loci. We estimated the pan-proteome of A and B supergroup
445 strains, and found that the A supergroup had a higher pan-proteome but a smaller core
446 proteome than supergroup B. Coupled with the observation of host-association bias between
447 these supergroups, and other major genomic features such as GC proportion, this suggests
448 that these divergent groups have followed very distinct evolutionary trajectories, despite
449 evidence for transduction of loci between supergroups, and perhaps have evolved distinct
450 physiologies and host-manipulation or -cooperation strategies. We note that the average

451 nucleotide identity (ANI) between A and B supergroup strains, and between strains from all
452 supergroups, is relatively low (within-supergroup identity >93%, between supergroup identity
453 <88%). This pattern of significant phylogenetic separation between supergroups suggests, as
454 others have noted, that these supergroups have the features expected of bacterial species³³.

455 The Darwin Tree of Life project¹⁵ is one of a growing constellation of biodiversity genomics
456 initiatives worldwide that, under the banner of the Earth BioGenome Project⁵¹, intend to
457 “sequence life for the future of life” (<https://www.earthbiogenome.org>). These projects, based
458 around ecological, regional or taxonomic lists of target species, will lay the foundations for
459 biological research, bioindustry and conservation for the next decades. While their focus is to
460 generate reference genomes for eukaryotic species, these projects will also yield critical
461 resources for the study of the microbial cobionts – mutualists, pathogens, parasites and
462 commensals – that live on and in eukaryotic organisms. Our understanding of *Wolbachia* and
463 other common endosymbionts will thrive on a rich harvest of cobiont genomes from the tens
464 to hundreds of thousands of host genomes that will be generated in the next decade. The
465 assembly of 110 high-quality *Wolbachia* genomes shows the power of the long read data now
466 being generated and the analytic approach that allowed these low complexity metagenomes
467 to be effectively separated into their constituent parts. Analysis of these genomes revealed a
468 propensity to infect different insect orders among supergroups, while simultaneously
469 pinpointing to several host switching events during the course of the *Wolbachia* pandemic.
470 Moreover, we observed that genome size in *Wolbachia* is correlated with the abundance of
471 active and pseudogenised copies of bacteriophage WO.

472 Methods

473 Detection and assembly of *Wolbachia* genomes from DToL 474 species data

475 DToL raw data are generated from whole or partial single specimens, and thus contain
476 sequence from any cobionts in or on the specimen at the time of sampling. We screened data
477 for 368 insect genomes generated by the Darwin Tree of Life project¹⁵ for the presence of the
478 intracellular endosymbiont *Wolbachia* ([Table S1](#)) using a marker gene scan approach by
479 searching for the small subunit rRNA locus. The prokaryotic 16S rRNA alignment from RFAM
480 (RF00177)⁵² was transformed into a HMMER profile and the profile was used to screen contigs
481 with nhmmscan⁵³. We defined a positive match as having an e-value $<10^{-150}$ or an aligned
482 length of >1000 nucleotides. Putative positive regions were extracted from the sequences,
483 and compared to the SILVA SSU database (version 138.1)⁵⁴ using sina⁵⁵. Matches were
484 filtered to retain only those with $>90\%$ identity. Taxonomic classification of each positive was
485 determined via a consensus rule of 80 percent of the top 20 best hits, using both the NCBI⁵⁶
486 and SILVA⁵⁷ taxonomies.

487 For *Wolbachia*-positive samples, all PacBio HiFi reads were analysed using kraken2⁵⁸ with a
488 custom database consisting of a genome from a species closely related to the host, all RefSeq
489 genomes of Anaplasmataceae and reference genomes of additionally detected cobionts
490 downloaded using NCBI datasets and masked using dustmasker⁵⁹. Horizontal transfer of
491 fragments of endosymbiont and organellar DNA to the nuclear genome is a common
492 phenomenon. To avoid inadvertently classifying nuclear *Wolbachia* insertions (NUWTs) as
493 deriving from an independent bacterial replicon, *Wolbachia* reads identified by kraken2 were
494 mapped to the insect genome assembly and only contigs fully covered by these reads were
495 retained. The *Wolbachia* reads were also independently re-assembled using several assembly
496 tools: flye (version 2.9) (flye --pacbio-hifi {reads} -o {dir} -t {threads} --asm-coverage 50 --
497 genome-size 1.6m --scaffold)⁶⁰, hifiasm (version 0.14) (hifiasm -o {prefix} -t {threads} {reads}
498 -D 10 -l 1 -s 0.999)⁶¹ and hifiasm-meta (version 0.1-r022) (hifiasm_meta -o {prefix} -t {threads}
499 {reads} -l 1)⁶². The several assemblies generated for each sample were ranked based on their
500 completeness using BUSCO version 5.2.2⁶³ and the Rickettsiales_odb10 dataset, alignment
501 to reference genomes using nucmer (version 4.0.0)⁶⁴, evenness of coverage and circularity.
502 The best (most complete, single-contig circular preferred) assembly per sample was chosen.
503 For samples where 10X Genomics Chromium data were available, polishing was performed
504 using FreeBayes-called variants⁶⁵ from 10X short reads aligned with LongRanger. The host
505 origin, span and completeness of all *Wolbachia* detected is presented in [Table S2](#).

506 Collation of *Wolbachia* genome dataset, gene prediction and 507 orthology inference

508 All available *Wolbachia* genomes were downloaded from NCBI GenBank on 01/02/2022, and
509 supplemented with assemblies generated from short-read insect datasets by Scholz et al.¹⁴.
510 This dataset contained replicate genomes for very closely related *Wolbachia* from the same
511 host, and many fragmented and partial assemblies. Only the most contiguous assembly per
512 host species was retained. These genomes were renamed using the schema

513 “R_Xyz_GenSpec_§”, where Xyz is the first three letters of the insect order of the host,
514 GenSpec is an abbreviation derived from the generic and specific epithets of the host, and §
515 indicates the supergroup. Retained assemblies were assessed for the presence of
516 contamination by performing a contig analysis by kraken2 using a database of only circular
517 *Wolbachia* genomes. A list of all removed contigs can be found in [Table S3](#). Furthermore, we
518 only included database-sourced *Wolbachia* genomes with at least 90% BUSCO
519 completeness⁶³ and at most 3% duplication with the Rickettsiales_odb10 dataset ([Table S3](#)).
520 The exception to this filtering was the inclusion of genomes belonging to the most divergent
521 supergroup S.

522 All of the publicly available and newly assembled genomes were annotated using Prodigal
523 (version 2.6.3)²⁵. Protein families were inferred using OrthoFinder (version 2.4.0)²⁶. We
524 identified 624 protein families which were single-copy in more than 95% of all *Wolbachia*
525 genomes. These were individually aligned using mafft in automatic mode (version 7.490)⁶⁶.
526 Individual maximum likelihood gene trees were calculated using iqtree (version 2.1.4) (iqtree
527 -s {alignment} -nt {threads})⁶⁷, and coalescence of these gene trees was determined using
528 ASTRAL (version 5.7.4)⁶⁸. The individual alignments were trimmed using trimAl (version 1.4)⁶⁹,
529 and concatenated to form a supermatrix. This was used to infer a maximum likelihood
530 phylogeny with iqtree using 1000 ultrafast bootstrap approximation iterations (version 2.1.4)
531 (iqtree -s {supermatrix} -m LG+G4 -bb 1000 -nt {threads})⁶⁷. The insect topology was
532 subsampled from Chesters and al⁷⁰. Incongruence in topology between the insect host and
533 *Wolbachia*, host phylogeny was determined with ggtree in R⁷¹.

534 Intrinsic genomic properties

535 All circular genomes were rotated to start with HemE (OG0000716) on the positive strand, as
536 this gene is located next to the origin of replication⁷². All pairwise alignments were calculated
537 using nucmer (version 4.0.0)⁶⁴, and breakpoints were inferred and adjusted for the aligned
538 coverage. Average nucleotide diversity was calculated using FastANI (version 1.33)⁷³. GC and
539 GC skew index values were calculated for all genomes using SkewIT³⁰.

540 Gene content analysis

541 To functionally annotate predicted genes, both Prokka (version 1.14.6)⁷⁴ and InterProScan
542 (version 5.54-87.0)⁷⁵ were run. The synteny plot of the biotin locus was created using
543 gggenes⁷⁶. All six genes that make up the biotin locus (BioA-D, BioF, BioH) were individually
544 aligned with mafft in automatic mode (version 7.490)⁶⁶ and transformed into a concatenated
545 nucleotide alignment. A phylogenetic tree was built using the model GTR+F+G4 in iqtree
546 (version 2.1.4)⁶⁷. Genes responsible for cytoplasmic incompatibility (CI) were identified by a
547 BLAST search⁷⁷ using the following genes as queries: CidA: WP_010962721.1,
548 WP_182158704.1; CidB: WP_010962722.1, WP_182158703.1 and CinA: CAQ54402.1;
549 CinB: CAQ54403.1. Only pairs of identified neighbouring genes (e-value $1e^{-30}$, coverage 80-
550 120%) were retained.

551 WO prophage analysis

552 A list of known prophage sequences was generated based on annotated regions described in
553 the literature^{37,40,78} ([Table S4](#)) for a set of genomes (R_Dip_DroSim_A, R_Hym_NasVit_A,

554 R_Dip_DroAna_A, R_Dip_Haelrr_A, R_Hym_CerSol_A and R_Hym_WiePum_A) and linked
555 to their respective gene families. Each *Wolbachia* genome was screened for continuous
556 stretches of linked prophage genes with at most five other genes in-between and these were
557 annotated as prophage regions if they contained at least one gene from one of the four core
558 phage modules (head, baseplate, tail, fibre). This permitted detection of novel prophage-
559 associated genes. Regions which contained at least 5 of 6 head, 7 of 8 baseplate, 5 of 6 fibre
560 and 5 of 6 tail module genes were deemed complete. Genomic maps of prophage integration
561 were created with circos⁷⁹. Phylogenetic Generalised Least Squares analyses were performed
562 to assess the correlation between prophage length and genome size using the ape R
563 package⁸⁰, using a Brownian model of evolution and the phylogenetic tree in Figure 2A. R
564 squared values were calculated using the package rr2⁸¹.

565 Acknowledgements

566 We thank our many colleagues in the Darwin Tree of Life project – from field collectors to data
567 curators – for the production of the raw data we analysed. We also thank Tree of Life
568 colleagues, especially Claudia Weber, Charlotte Wright and Ellen Cameron for fruitful
569 discussions and Andrew Varley, James Torrance and Shane McCarthy for help with sequence
570 deposition. This research was funded by the Wellcome Trust Grants 206194 and 218328. For
571 the purpose of Open Access, the author has applied a CC BY public copyright licence to any
572 Author Accepted Manuscript version arising from this submission.

573 Data availability

574 The raw data for each species analysed is available under BioProject PRJEB40665
575 (<https://www.ebi.ac.uk/ena/browser/view/PRJEB40665>). Darwin Tree of Life species and data
576 are collated in the project portal at <https://portal.darwintreeoflife.org>. The *Wolbachia* genome
577 assemblies are currently in progress for sequence deposition in INSDC, but can already be
578 accessed on Zenodo (<https://doi.org/10.5281/zenodo.7092419>).

579 References

- 580 1. Yen, J. H. & Barr, A. R. New Hypothesis of the Cause of Cytoplasmic Incompatibility in *Culex*
581 *pipiens* L. *Nature* **232**, 657–658 (1971).
- 582 2. Yen, J. H. & Barr, A. R. The etiological agent of cytoplasmic incompatibility in *Culex pipiens*.
583 *Journal of Invertebrate Pathology* **22**, 242–250 (1973).
- 584 3. Bordenstein, S. R., O'Hara, F. P. & Werren, J. H. Wolbachia-induced incompatibility precedes
585 other hybrid incompatibilities in *Nasonia*. *Nature* **409**, 707–710 (2001).
- 586 4. Hurst, G. D. D. *et al.* Male-killing Wolbachia in two species of insect. *Proc. R. Soc. Lond. B* **266**,
587 735–740 (1999).
- 588 5. Stouthamer, R., Breeuwer, J. A. J., Luck, R. F. & Werren, J. H. Molecular identification of
589 microorganisms associated with parthenogenesis. *Nature* **361**, 66–68 (1993).
- 590 6. Hornett, E. A. *et al.* Evolution of Male-Killer Suppression in a Natural Population. *PLoS Biol* **4**, e283
591 (2006).
- 592 7. Werren, J. H., Baldo, L. & Clark, M. E. Wolbachia: master manipulators of invertebrate biology. *Nat*
593 *Rev Microbiol* **6**, 741–751 (2008).
- 594 8. Nikoh, N. *et al.* Evolutionary origin of insect–*Wolbachia* nutritional mutualism. *Proc. Natl. Acad.*
595 *Sci. U.S.A.* **111**, 10257–10262 (2014).
- 596 9. Pan, X. *et al.* The bacterium Wolbachia exploits host innate immunity to establish a symbiotic
597 relationship with the dengue vector mosquito *Aedes aegypti*. *ISME J* **12**, 277–288 (2018).
- 598 10. Hoerauf, A., Mand, S., Adjei, O., Fleischer, B. & Büttner, D. W. Depletion of wolbachia
599 endobacteria in *Onchocerca volvulus* by doxycycline and microfilaridermia after ivermectin
600 treatment. *The Lancet* **357**, 1415–1416 (2001).
- 601 11. Zug, R. & Hammerstein, P. Still a Host of Hosts for Wolbachia: Analysis of Recent Data
602 Suggests That 40% of Terrestrial Arthropod Species Are Infected. *PLoS ONE* **7**, e38544 (2012).
- 603 12. Zhou, W., Rousset, F. & O'Neill, S. Phylogeny and PCR-based classification of Wolbachia
604 strains using *wsp* gene sequences. *Proc. R. Soc. Lond. B* **265**, 509–515 (1998).
- 605 13. Glowska, E., Dragun-Damian, A., Dabert, M. & Gerth, M. New Wolbachia supergroups
606 detected in quill mites (Acari: Syringophilidae). *Infection, Genetics and Evolution* **30**, 140–146
607 (2015).
- 608 14. Scholz, M. *et al.* Large scale genome reconstructions illuminate Wolbachia evolution. *Nat*
609 *Commun* **11**, 5235 (2020).
- 610 15. The Darwin Tree of Life Project Consortium. Sequence locally, think globally: The Darwin
611 Tree of Life Project. *Proceedings of the National Academy of Sciences* **119**, e2115642118 (2022).
- 612 16. Challis, R., Richards, E., Rajan, J., Cochrane, G. & Blaxter, M. BlobToolKit – Interactive
613 Quality Assessment of Genome Assemblies. *G3 Genes|Genomes|Genetics* **10**, 1361–1374 (2020).
- 614 17. Eren, A. M. *et al.* Anvi'o: an advanced analysis and visualization platform for 'omics data.
615 *PeerJ* **3**, e1319 (2015).
- 616 18. Regan, T. *et al.* Characterisation of the British honey bee metagenome. *Nat Commun* **9**, 4995
617 (2018).

- 618 19. West, S. A., Cook, J. M., Werren, J. H. & Godfray, H. C. J. *Wolbachia* in two insect host–
619 parasitoid communities. *Molecular Ecology* **7**, 1457–1465 (1998).
- 620 20. Duron, O. *et al.* The diversity of reproductive parasites among arthropods: *Wolbachia* do not
621 walk alone. *BMC Biol* **6**, 27 (2008).
- 622 21. Savill, P., Perrins, C., Kirby, K. & Fisher, N. *Wytham woods: Oxford's ecological laboratory*.
623 (Oxford Univ. Press, 2011).
- 624 22. Hilgenboecker, K., Hammerstein, P., Schlattmann, P., Telschow, A. & Werren, J. H. How
625 many species are infected with *Wolbachia*? – a statistical analysis of current data: *Wolbachia*
626 infection rates. *FEMS Microbiology Letters* **281**, 215–220 (2008).
- 627 23. Ahmed, M. Z., Breinholt, J. W. & Kawahara, A. Y. Evidence for common horizontal
628 transmission of *Wolbachia* among butterflies and moths. *BMC Evol Biol* **16**, 118 (2016).
- 629 24. Misof, B. *et al.* Phylogenomics resolves the timing and pattern of insect evolution. *Science*
630 **346**, 763–767 (2014).
- 631 25. Hyatt, D. *et al.* Prodigal: prokaryotic gene recognition and translation initiation site
632 identification. *BMC Bioinformatics* **11**, 119 (2010).
- 633 26. Emms, D. M. & Kelly, S. OrthoFinder: phylogenetic orthology inference for comparative
634 genomics. *Genome Biol* **20**, 238 (2019).
- 635 27. Russell, J. A. *et al.* Specialization and geographic isolation among *Wolbachia* symbionts from
636 ants and lycaenid butterflies. *Evolution* **63**, 624–640 (2009).
- 637 28. Werren, J. H. & Windsor, D. M. *Wolbachia* infection frequencies in insects: evidence of a
638 global equilibrium? *Proc. R. Soc. Lond. B* **267**, 1277–1285 (2000).
- 639 29. Tagami, Y. & Miura, K. Distribution and prevalence of *Wolbachia* in Japanese populations of
640 Lepidoptera: *Wolbachia* in Japanese Lepidoptera. *Insect Molecular Biology* **13**, 359–364 (2004).
- 641 30. Lu, J. & Salzberg, S. L. SkewIT: The Skew Index Test for large-scale GC Skew analysis of
642 bacterial genomes. *PLoS Comput Biol* **16**, e1008439 (2020).
- 643 31. Comandatore, F. *et al.* Supergroup C *Wolbachia*, mutualist symbionts of filarial nematodes,
644 have a distinct genome structure. *Open Biol.* **5**, 150099 (2015).
- 645 32. Mahmood, S., Nováková, E., Martinů, J., Sychra, O. & Hypša, V. *Extremely reduced*
646 *supergroup F Wolbachia : transition to obligate insect symbionts*.
647 <http://biorxiv.org/lookup/doi/10.1101/2021.10.15.464041> (2021) doi:10.1101/2021.10.15.464041.
- 648 33. Ellegaard, K. M., Klasson, L., Näslund, K., Bourtzis, K. & Andersson, S. G. E. Comparative
649 Genomics of *Wolbachia* and the Bacterial Species Concept. *PLoS Genet* **9**, e1003381 (2013).
- 650 34. Gerth, M. & Bleidorn, C. Comparative genomics provides a timeframe for *Wolbachia* evolution
651 and exposes a recent biotin synthesis operon transfer. *Nat Microbiol* **2**, 16241 (2017).
- 652 35. Gerth, M., Röthe, J. & Bleidorn, C. Tracing horizontal *Wolbachia* movements among bees
653 (*Anthophila*): a combined approach using multilocus sequence typing data and host phylogeny.
654 *Mol Ecol* **22**, 6149–6162 (2013).
- 655 36. Bordenstein, S. R. & Bordenstein, S. R. Eukaryotic association module in phage WO
656 genomes from *Wolbachia*. *Nat Commun* **7**, 13155 (2016).
- 657 37. Bordenstein, S. R. & Bordenstein, S. R. Widespread phages of endosymbionts: Phage WO

- 658 genomics and the proposed taxonomic classification of Symbioviridae. *PLoS Genet* **18**, e1010227
659 (2022).
- 660 38. Gavotte, L. *et al.* A Survey of the Bacteriophage WO in the Endosymbiotic Bacteria
661 *Wolbachia*. *Molecular Biology and Evolution* **24**, 427–435 (2006).
- 662 39. Massey, J. H. & Newton, I. L. G. Diversity and function of arthropod endosymbiont toxins.
663 *Trends in Microbiology* **30**, 185–198 (2022).
- 664 40. LePage, D. P. *et al.* Prophage WO genes recapitulate and enhance *Wolbachia*-induced
665 cytoplasmic incompatibility. *Nature* **543**, 243–247 (2017).
- 666 41. Beckmann, J. F., Ronau, J. A. & Hochstrasser, M. A *Wolbachia* deubiquitylating enzyme
667 induces cytoplasmic incompatibility. *Nat Microbiol* **2**, 17007 (2017).
- 668 42. Martinez, J., Klasson, L., Welch, J. J. & Jiggins, F. M. Life and Death of Selfish Genes:
669 Comparative Genomics Reveals the Dynamic Evolution of Cytoplasmic Incompatibility. *Molecular*
670 *Biology and Evolution* **38**, 2–15 (2021).
- 671 43. Beckmann, J. F. *et al.* The Toxin–Antidote Model of Cytoplasmic Incompatibility: Genetics and
672 Evolutionary Implications. *Trends in Genetics* **35**, 175–185 (2019).
- 673 44. Dy, R. L., Przybilski, R., Semeijn, K., Salmond, G. P. C. & Fineran, P. C. A widespread
674 bacteriophage abortive infection system functions through a Type IV toxin–antitoxin mechanism.
675 *Nucleic Acids Research* **42**, 4590–4605 (2014).
- 676 45. Kumar, S. & Blaxter, M. L. Simultaneous genome sequencing of symbionts and their hosts.
677 *Symbiosis* **55**, 119–126 (2011).
- 678 46. Almeida, A. *et al.* A unified catalog of 204,938 reference genomes from the human gut
679 microbiome. *Nat Biotechnol* **39**, 105–114 (2021).
- 680 47. Parks, D. H. *et al.* Recovery of nearly 8,000 metagenome-assembled genomes substantially
681 expands the tree of life. *Nat Microbiol* **2**, 1533–1542 (2017).
- 682 48. Ju, J.-F. *et al.* *Wolbachia* supplement biotin and riboflavin to enhance reproduction in
683 planthoppers. *ISME J* **14**, 676–687 (2020).
- 684 49. Conner, W. R. *et al.* Genome comparisons indicate recent transfer of *w* Ri-like *Wolbachia*
685 between sister species *Drosophila suzukii* and *D. subpulchrella*. *Ecol Evol* **7**, 9391–9404 (2017).
- 686 50. Turelli, M. *et al.* Rapid Global Spread of *w*Ri-like *Wolbachia* across Multiple *Drosophila*.
687 *Current Biology* **28**, 963–971.e8 (2018).
- 688 51. Lewin, H. A. *et al.* The Earth BioGenome Project 2020: Starting the clock. *Proc. Natl. Acad.*
689 *Sci. U.S.A.* **119**, e2115635118 (2022).
- 690 52. Kalvari, I. *et al.* Rfam 14: expanded coverage of metagenomic, viral and microRNA families.
691 *Nucleic Acids Research* **49**, D192–D200 (2021).
- 692 53. Eddy, S. R. Accelerated Profile HMM Searches. *PLOS Computational Biology* **7**, e1002195
693 (2011).
- 694 54. Quast, C. *et al.* The SILVA ribosomal RNA gene database project: improved data processing
695 and web-based tools. *Nucleic Acids Research* **41**, D590–D596 (2013).
- 696 55. Pruesse, E., Peplies, J. & Glöckner, F. O. SINA: Accurate high-throughput multiple sequence
697 alignment of ribosomal RNA genes. *Bioinformatics* **28**, 1823–1829 (2012).

- 698 56. Schoch, C. L. *et al.* NCBI Taxonomy: a comprehensive update on curation, resources and
699 tools. *Database* **2020**, baaa062 (2020).
- 700 57. Yilmaz, P. *et al.* The SILVA and “All-species Living Tree Project (LTP)” taxonomic
701 frameworks. *Nucl. Acids Res.* **42**, D643–D648 (2014).
- 702 58. Wood, D. E., Lu, J. & Langmead, B. Improved metagenomic analysis with Kraken 2. *Genome*
703 *Biology* **20**, 257 (2019).
- 704 59. Morgulis, A., Gertz, E. M., Schäffer, A. A. & Agarwala, R. A fast and symmetric DUST
705 implementation to mask low-complexity DNA sequences. *J Comput Biol* **13**, 1028–1040 (2006).
- 706 60. Kolmogorov, M., Yuan, J., Lin, Y. & Pevzner, P. A. Assembly of long, error-prone reads using
707 repeat graphs. *Nat Biotechnol* **37**, 540–546 (2019).
- 708 61. Cheng, H., Concepcion, G. T., Feng, X., Zhang, H. & Li, H. Haplotype-resolved de novo
709 assembly using phased assembly graphs with hifiasm. *Nat Methods* **18**, 170–175 (2021).
- 710 62. Feng, X., Cheng, H., Portik, D. & Li, H. Metagenome assembly of high-fidelity long reads with
711 hifiasm-meta. *Nat Methods* **19**, 671–674 (2022).
- 712 63. Manni, M., Berkeley, M. R., Seppey, M., Simão, F. A. & Zdobnov, E. M. BUSCO Update:
713 Novel and Streamlined Workflows along with Broader and Deeper Phylogenetic Coverage for
714 Scoring of Eukaryotic, Prokaryotic, and Viral Genomes. *Molecular Biology and Evolution* **38**, 4647–
715 4654 (2021).
- 716 64. Marçais, G. *et al.* MUMmer4: A fast and versatile genome alignment system. *PLOS*
717 *Computational Biology* **14**, e1005944 (2018).
- 718 65. Garrison, E. & Marth, G. *Haplotype-based variant detection from short-read sequencing*.
719 <http://arxiv.org/abs/1207.3907> (2012) doi:10.48550/arXiv.1207.3907.
- 720 66. Katoh, K. & Standley, D. M. MAFFT Multiple Sequence Alignment Software Version 7:
721 Improvements in Performance and Usability. *Molecular Biology and Evolution* **30**, 772–780 (2013).
- 722 67. Minh, B. Q. *et al.* IQ-TREE 2: New Models and Efficient Methods for Phylogenetic Inference in
723 the Genomic Era. *Molecular Biology and Evolution* **37**, 1530–1534 (2020).
- 724 68. Mirarab, S. *et al.* ASTRAL: genome-scale coalescent-based species tree estimation.
725 *Bioinformatics* **30**, i541–i548 (2014).
- 726 69. Capella-Gutiérrez, S., Silla-Martínez, J. M. & Gabaldón, T. trimAl: a tool for automated
727 alignment trimming in large-scale phylogenetic analyses. *Bioinformatics* **25**, 1972–1973 (2009).
- 728 70. Chesters, D. Construction of a Species-Level Tree of Life for the Insects and Utility in
729 Taxonomic Profiling. *Syst Biol* syw099 (2016) doi:10.1093/sysbio/syw099.
- 730 71. Yu, G., Lam, T. T.-Y., Zhu, H. & Guan, Y. Two Methods for Mapping and Visualizing
731 Associated Data on Phylogeny Using *Ggtree*. *Molecular Biology and Evolution* **35**, 3041–3043
732 (2018).
- 733 72. Ioannidis, P. *et al.* New criteria for selecting the origin of DNA replication in Wolbachia and
734 closely related bacteria. *BMC Genomics* **8**, 182 (2007).
- 735 73. Jain, C., Rodriguez-R, L. M., Phillippy, A. M., Konstantinidis, K. T. & Aluru, S. High throughput
736 ANI analysis of 90K prokaryotic genomes reveals clear species boundaries. *Nat Commun* **9**, 5114
737 (2018).

- 738 74. Seemann, T. Prokka: rapid prokaryotic genome annotation. *Bioinformatics* **30**, 2068–2069
739 (2014).
- 740 75. Quevillon, E. *et al.* InterProScan: protein domains identifier. *Nucleic Acids Research* **33**,
741 W116–W120 (2005).
- 742 76. Wilkins, D. gggenes. <https://github.com/wilkox/gggenes>.
- 743 77. Altschul, S. F., Gish, W., Miller, W., Myers, E. W. & Lipman, D. J. Basic local alignment
744 search tool. *Journal of Molecular Biology* **215**, 403–410 (1990).
- 745 78. Miao, Y., Xiao, J. & Huang, D. Distribution and Evolution of the Bacteriophage WO and Its
746 Antagonism With Wolbachia. *Front. Microbiol.* **11**, 595629 (2020).
- 747 79. Krzywinski, M. *et al.* Circos: an information aesthetic for comparative genomics. *Genome Res*
748 **19**, 1639–1645 (2009).
- 749 80. Paradis, E. & Schliep, K. ape 5.0: an environment for modern phylogenetics and evolutionary
750 analyses in R. *Bioinformatics* **35**, 526–528 (2019).
- 751 81. Ives, A. R. R2s for Correlated Data: Phylogenetic Models, LMMs, and GLMMs. *Systematic*
752 *Biology* **68**, 234–251 (2019).