1 **The first draft genomes of the ant *Formica exsecta,* and its *Wolbachia***

2 **endosymbiont reveal extensive gene transfer from endosymbiont to host**

3 Authors: Kishor Dhaygude[1], Abhilash Nair[1], Helena Johansson[1], Yannick Wurm[2],

4 and Liselotte Sundström[1, 3]

5 1) Organismal and Evolutionary Biology Research Programme, Faculty of Biological

6 and environmental sciences P.O. Box 65, FI-00014 University of Helsinki, Finland

7 2) Organismal Biology Department, School of Biological and Chemical Sciences,

8 Queen Mary University of London, Mile End Road, London E1 4NS, United

9 Kingdom

10 3) Tvärminne Zoological Station, University of Helsinki, J.A. Palménin tie 260, FI-

11 10900 Hanko, Finland

12 **Abstract:**

13 The wood ant *Formica exsecta* (Formicidae; Hymenoptera), is a common ant species

14 throughout the Palearctic region. The species is a well established model for studies

15 of ecological characteristics and evolutionary conflict. In this study, we sequenced

16 and assembled draft genomes for *Formica exsecta* and its endosymbiont *Wolbachia*.

17 The draft *F. exsecta* genome is 277.7 Mb long; we identify 13,767 protein coding

18 genes for which we provide gene ontology, and protein domain annotations. This is

19 also the first report of a *Wolbachia* genome from ants, and provides insights into the

20 phylogenetic position of this endosymbiont. We also identified multiple horizontal

21 gene transfer events (HGTs) from *Wolbachia* to *F. exsecta*. Some of these HGTs have

22 also occurred in parallel in multiple other insect genomes, highlighting the extent of

23 HGTs in eukaryotes. We expect that the *F. exsecta* genome will be valuable resource

24 in further exploration of the molecular basis of the evolution of social organization.

25

26 **Key words:** *Formica exsecta*, genome, endosymbionts*,* transposons, horizontal gene

27 transfer, *Wolbachia*

28

29 **Introduction**

30

31  Adapting to changes in the environment is the foundation of species survival, and is

32  usually thought to be a gradual process. Genomic changes, such as single nucleotide

33  substitutions play key roles in adaptive evolution, although few mutations are

34  beneficial. Besides nucleotide substitutions, other structural and regulatory units, such

35  as transposable elements (TEs) and epigenetic modifications, can also act as drivers in

36  adaptation (González et al., 2010; Rostant, Wedell & Hosken, 2012; Casacuberta &

37  González, 2013). Genetic material can also be acquired from other organisms by

38  means of horizontal gene transfer (HGTs), and this can also lead to novel adaptive

39  traits (Schönknecht, Weber & Lercher, 2014; Wybouw et al., 2016). Both mutations

40  and HGTs can drive rapid genome evolution (Dunning Hotopp, 2011; Boto, 2014).

41  Horizontal gene transfers have been reported in many taxa, most commonly from

42  bacteria to eukaryotes (Dunning Hotopp, 2011), plants (Yue et al., 2012; Matveeva &

43  Lutova, 2014), fungi (Rolland et al., 2009; Fitzpatrick, 2012; Bruto et al., 2014), but

44  the underlying mechanisms that underpin horizontal gene transfer events, and mode

45  by which bacterial genetic material is integrated into the eukaryote genome are not

46  well understood.

47

48  Many cases of horizontal gene transfer from bacteria to eukaryotes involve

49  intracellular endosymbionts, which are maternally transmitted through oocytes

50  (Werren, 1997; Ferree et al., 2005). The most common examples of endosymbiont to

51  host horizontal gene transfers involve the bacterium *Wolbachia*, a well described

52  intracellular, maternally inherited gram-negative bacterium known to infect over 40%

53  of the investigated insect species (Werren, 1997; Werren, Baldo & Clark, 2008).

54  *Wolbachia* infection is also prevalent in filarial nematodes, crustaceans, and arachnids

55  (Cordaux, Michel-Salzat & Bouchon, 2001; Fenn et al., 2006; Goodacre et al., 2006).

56  *Wolbachia-* host interactions can be mutualistic or pathogenic (Moya et al., 2008). A

57  number of ecdysozoan genomes have been reported to contain chromosomal

58  insertions originating from *Wolbachia*, including the mosquito *Aedes aegypti*

59  (Klasson et al., 2009a; Woolfit et al., 2009), the longhorn beetle *Monochamus*

60  *alternatus* (Aikawa et al., 2009), filarial nematodes of the genera *Onchocerca,*

61  *Brugia*, and *Dirofilaria* (Fenn et al., 2006; Hotopp et al., 2007), parasitoid wasps of

62  the genus *Nasonia*, the fruit fly *Drosophila ananassae*, the pea aphid *Acythosiphon*

63  *pisum* (Nikoh & Nakabachi, 2009; Nikoh et al., 2010), and the bean

64  beetle *Callosobruchus chinensis* (Kondo et al., 2002). Although most of the

65    transferred DNA is probably nonfunctional in the host genome (Kondo et al., 2002;

66    Hotopp et al., 2007; Nikoh et al., 2008), some of the transferred genes are functional

67    (Klasson et al., 2009a). These genes are expressed in specific tissues, are subject to

68    purifying selection, and are involved in processes such as protein synthesis inhibition,

69    membrane transport and metabolism (Hotopp et al., 2007; Woolfit et al., 2009;

70    McNulty et al., 2013).

71

72    Infection with *Wolbachia* is widespread in Hymenoptera. Most hymenopteran

73    *Wolbachia* infections have the cytoplasmic incompatibility phenotype (Werren &

74    Windsor, 2000), which leads to reproductive incompatibility between infected sperm

75    and uninfected eggs. Wenseleers et al. (1998) showed that 25 out of 50 species of ants

76    in Java and Sumatra screened positive for one strain of *Wolbachia*. By contrast, a

77    study on a single Swiss population of the ant *Formica exsecta*, found that all the ants

78    tested were infected with four or five different strains of *Wolbachia* (Keller et al.,

79    2001; Reuter & Keller, 2003).

80

81    The aims of this study are to test whether horizontally transferred genetic elements

82    exist in the genome of the ant *Formica exsecta*, and to describe the genomic

83    organization of any such elements. The genus *Formica* is listed by the Global Ant

84    Genome Alliance (GAGA) as one of the high-priority ant taxons to be sequenced

85    (Boomsma et al., 2017; http://antgenomics.dk/), owing to its key taxonomic position,

86    and the ecological and behavioral data that are available for the species. To date, no

87    genome sequence is available for this genus.

88

89    Our study population of *F. exsecta*, located on the Hanko peninsula, Southwestern

90    Finland, has been monitored since 1994, and data on demography, genetic structure,

91    and ecology are available (Sundström, Chapuisat & Keller, 1996; Sundström, Keller

92    & Chapuisat, 2003; Haag-Liautard et al., 2009; Vitikainen, Haag-Liautard &

93    Sundström, 2015). Based on genetic data on colony kin structure most (97%) of the

94    approximately 200 colonies are known to have a single reproductive queen, mated to

95    one or more (usually two) males (Sundström, Chapuisat & Keller, 1996; Sundström,

96    Keller & Chapuisat, 2003; Haag-Liautard et al., 2009; Vitikainen, Haag-Liautard &

97    Sundström, 2015). We report the whole genome sequencing of this species, and the

98    draft genome sequence of its associated cytoplasmic *Wolbachia* endosymbiont

99    (wFex). We further report the presence of multiple extensive insertions of *Wolbachia*

100    genetic material in the host genome, and compare the HGTs insertions discovered in

101    the assembled draft genome to other genomes, to understand the pattern of HGT

102    events between endosymbiont and host. We analyze in detail the genomic features of

103    *F. exsecta* along with its endosymbiont *Wolbachia,* and discuss our findings in the

104    light of genome evolution in *Wolbachia* and its host.

105

## Materials and Methods

107

### Sample collection and genome sequencing.

109    We selected one single-queen colony from our study population on the island

110    Furuskär (F162), and collected 200 adult males from this colony. We used males

111    because in Hymenoptera these arise through arrhenotoky (Normark, 2003) and are

112    haploid (Crozier, 1975), meaning that a pool of males together are representative of

113    the diploid genome of their mother. DNA extraction was done from testis, which

114    contains sperm cells and organ tissue, to avoid contamination by gut microbiota. We

115    used a Qiagen Genomic-tip 20/G extraction kit according to the manufacturer's

116    protocol. For Illumina sequencing we constructed three small insert paired-end

117    libraries (insert sizes of 200 bp, 500 bp, 800 bp), and four mate pair (large insert

118    paired-end) libraries (insert sizes of 2 kb, 5 kb, 10 kb and 20 kb), each containing

119    DNA from 15-50 pooled males. Libraries were prepared using protocols

120    recommended by the manufacturers. Sequencing was done at the Beijing Genomics

121    Institute (BGI) using HiSeq2000, which produced a total of 99.97 GB of raw data

122    (Table 1).

123

### Genome assembly

125    We assembled the *F. exsecta* genome using SOAPdenovo2 version 2.04 (Xie et al.,

126    2014) in three main steps. First, a de Bruijn graph was constructed using short length

127    insert library reads with default parameters (k-mer value of 45), to construct the

128    contigs. The initial contig assembly contained 104,190 contigs with an N50 size of

129    22,328 bp, and total length of 276.23 Mb of sequence, at an average depth of

130    coverage of 47.37X. Second, all individual reads were realigned onto the contigs.

131    Because reads are paired, they can aid with scaffolding: The number of reads

132    supporting the adjacency of each pair of contigs was calculated and weighted by the

133    ratio between consistent and conflicting paired ends. Scaffolds were constructed in a

134    stepwise manner using libraries of increasing sizes from 500bp insert size paired-end

135    reads up to mate-pair of 5 kb insert size. 80,473 contigs could not be placed in

136    scaffolds. These are highly similar repetitive sequences, since the cd-hit-est tool

137    (Huang et al., 2010) showed that 43% of these contigs clustered together at 80% of

138    the sequence length. Third, sequencing gaps in the scaffolds were closed with the two

139    mate-pair libraries (Insert size 10 kb and 20 kb). Overall, these steps produced an

140    initial assembly with an N50 scaffold length of 949,634 bp, and a total length of

141    289,843,734 bp with each scaffold longer than 200 bp.

142

143    We used blobology v1.0 (Kumar et al., 2013) to generate taxon-annotated GC-

144    coverage (TAGC) plots of scaffolds in the genome assembly, which can help to

145    identify bacterial contamination (Supplementary Figure S1). The scaffolds for the

146    TAGC plot were successfully annotated to the taxonomic order based on the best blast

147    match to the NCBI nt database (O'Leary et al., 2016). This analysis revealed that 74

148    scaffolds matched the endosymbiotic bacterium *Wolbachia*. Sixty-nine of these

149    scaffolds were removed as we concluded that they are part of the *Wolbachia* genome

150    (see analysis below), but five contigs were retained in the final assembly for *F.*

151    *exsecta* as they contained both *Wolbachia* and ant sequences. Following this curation,

152    the final draft genome assembly was 277.7 Mb long with an N50 value of 997,654 bp

153    and 36% Guanine-cytosine (GC) content (Table 2).

154

155    **Genome assembly of *Wolbachia***

156    All 25 published *Wolbachia* genomes were obtained from the NCBI database

157    (O'Leary et al., 2016). We aligned the 74 scaffolds from the initial *F. exsecta*

158    assembly that matched with *Wolbachia* against these genomes using MUMmer 3.23

159    (Kurtz et al., 2004), and inspected the alignments manually. Sixty-nine of the 74

160    scaffolds matched completely to *Wolbachia* genomic regions. These 69 scaffolds

161    represented 3.09 Mb total, with a N50 value of 104,167 bp, henceforth we refer to this

162    group of scaffolds as "the *Wolbachia* endosymbiont genome of *F. exsecta*" (wFex).

163

164    The remaining five scaffolds each contained several interspersed fragments with

165    similarity to *Wolbachia* genomes, whereas other parts of these scaffolds had high

166    similarity to genomes of ants. Furthermore, the sequencing coverage of these

167   scaffolds was similar to the *F. exsecta* scaffolds, rather than to the *Wolbachia*

168   scaffolds. Finally, detailed inspection of these scaffolds in a genome browser showed

169   no change in sequencing depth where we identify the interspersed fragments with

170   similarity to *Wolbachia*, which would be expected for erroneous chimeric assembly

171   (Lasken & Stockwell, 2007). These data thus suggest that fragments of *Wolbachia*

172   were horizontally transferred to the *F. exsecta* genome. To corroborate these results

173   with independent approaches, we re-assembled the raw sequencing data with two

174   additional independent algorithms that we expect would make different types of

175   assembly errors than SOAPdenovo. The first software, Velvet version 1.2.09 (Zerbino

176   & Birney, 2008), is also based on a de Bruijn graph; the second, SGA version 0.10.5

177   (Simpson & Durbin, 2012) is based on a string graph. Both resulting assemblies

178   confirmed the patterns we had seen, and validate the idea that the five SOAPdenovo

179   scaffolds containing sequence with similarity to both ants, and *Wolbachia* represent

180   horizontal gene transfers from *Wolbachia* to *F. exsecta*.

181

182   We further compared the sequences of the horizontally transferred fragments in the

183   five SOAPdenovo scaffolds against the NCBI (nr/nt) database (O'Leary et al., 2016),

184   using blast 2.2.27 (Altschul et al., 1990) to determine whether these fragments may

185   have also undergone horizontal gene transfer in other arthropod genomes. We

186   performed analogous searches on ant genomes present in the NCBI, and the

187   Fourmidable databases (Wurm et al., 2009). When a positive match with any other ant

188   or arthropod genomes was found, the exact location of the insertion was determined,

189   and compared with that of *F. exsecta*. Finally, the five scaffolds were also compared

190   to the *F. exsecta* transcriptome (Dhaygude et al., 2017), using blastn 2.2.27, to assess

191   similarity with expressed sequences.

192

193   **Quantitative assessment of genome assemblies**

194   The quality of the genome assembly is crucial, as it defines the quality of all

195   subsequent analyses that are based on the genome sequences. We explored multiple

196   assembly options (data not shown), and used two methods to assess assembly quality

197   and robustness in order to select the highest quality assembly. First, we evaluated

198   genome contiguity (number and length of contigs) using Quast 3.2 (Gurevich et al.,

199   2013) to assess whether our newly assembled draft genome is comparable to

200   published ant genomes (Favreau et al., 2018) based on assembly statistics (N50,N90).

201    Second, we used core gene content-based quality assessment using CEGMA 2.4

202    (Parra et al., 2007) to ascertain that the 248 most highly conserved eukaryotic proteins

203    are present in our genome assembly. We also compared genes present in our genome

204    assembly to single-copy orthologs across four lineage-specific sets (Eukaryota (303

205    genes), Insecta (1,658 genes), Arthropoda (2,675 genes), and Hymenoptera (4,415

206    genes)) using  the BUSCO 1.1(Simão et al., 2015). In addition, we compared the *F.*

207    *exsecta* genome with 13 other ant genomes, *Camponotus floridanus, Atta cephalotes,*

208    *Acromyrmex echinatior, Cardiocondyla obscurior, Cerapachys biroi, Lasius niger,*

209    *Linepithema humile, Monomorium pharaonis, Pogonomyrmex barbatus, Vollenhovia*

210    *emeryi, Wasmannia auropunctata, Harpegnathos saltator,* and *Solenopsis invicta*

211    (Wurm et al., 2009), using BUSCO. We report BUSCO quality metrics for the *F.*

212    *exsecta* genome. (Table 3).

213

214    The quality of the *Wolbachia* endosymbiont genome was quantified with a similar

215    approach, where we used BUSCO to examine the presence of Universal Single-Copy

216    Orthologs of the Bacteria (148 genes), and the Proteobacteria (221 genes) lineages

217    (Table 3). We also used BUSCO to compare the wFex genome with four other

218    *Wolbachia* genomes, including the *Wolbachia* endosymbionts of *Drosophila simulans*

219    *(wRi), Culex quinquefasciatus (wPip), Drosophila melanogaster (wMel),* and

220    *Drosophila simulans (wNo).*

221

222    **Gene prediction**

223    We combined several publicly available data sets and computational gene prediction

224    tools to establish an Official Gene Set (OGS) for the *F. exsecta* genome. First, we

225    used the MAKER version 2.28 pipeline (Cantarel et al., 2008; Holt & Yandell, 2011),

226    to derive consensus gene models from Augustus version 3.1.0 (Stanke &

227    Morgenstern, 2005), SNAP version 2016-07-28 (Korf, 2004), and Exonerate version

228    2.2.0 (Slater & Birney, 2005). For this MAKER prediction we used as input datasets

229    the *F. exsecta* transcriptome (ESTs) (Bioproject ID: PRJNA213662, (Dhaygude et al.,

230    2017)), and the proteomes of all available ant species (Uniprot download on 20-04-

231    2015). The longest protein at each genomic locus was retained, resulting in a set of

232    23,517 gene models. Because samples may have different sets of transcripts, owing to

233    different biological conditions or developmental stages (Dhaygude et al., 2017), we

234    additionally made a separate transcript-spliced assembly using RNA sequences

235    generated from separate libraries for different life stages (Dhaygude et al., 2017),

236    using the Tophat version 2.1.0 (Trapnell, Pachter & Salzberg, 2009), and Cufflinks

237    version 2.2.1 (Trapnell et al., 2010). The assemblies from the different samples were

238    then merged using cuffmerge (Trapnell et al., 2010). We further obtained separate

239    Augustus version 3.1.0 (Stanke & Morgenstern, 2005), and Glimmer version 3.02

240    (Salzberg et al., 1998) gene models with default settings (Augustus: --species=fly --

241    genemodel=partial, --strand=both, Glimmer: +f, +s, -g 60). The gene sets and gene

242    models from MAKER and from other programs were then merged. Redundancy was

243    removed by favoring for each transcript the longest prediction starting with a

244    methionine. If several transcripts had the same length we retained the one which had

245    the best support from the cufflinks transcript assembly. This redundancy removal

246    resulted in a final set of 13,637 protein coding gene models (final OGS), which

247    contained 33,121 transcripts.

248

249    **Genome Annotation**

250    We analyzed the complete official gene sets (OGS) of *F. exsecta* to identify sequence

251    and functional similarity by comparing with different sequence databases using blast.

252    By using a ribosomal database, we were able to annotate both the large (LSU), and

253    the small (SSU) subunit ribosomal RNAs. The remaining gene sequences were used

254    for retrieving functional information from other databases (SwissProt, Pfam,

255    PROSITE, and COG). Gene sequences were considered to be coding if they had a

256    strong unique hit to the SwissProt protein database (Magrane & Consortium, 2011;

257    The Uniprot Consortium, 2017), or appeared to be orthologs of known predicted

258    protein-coding genes from ant species based on TrEMBL (Translation of EMBL

259    nucleotide sequence database). We also assigned putative metabolic pathways,

260    functional classes, enzyme classes, GeneOntology terms, and locus names with the

261    AutoFact tool (Koski et al., 2005). To further improve annotation, and for assigning

262    biological function (e.g. gene expression, metabolic pathways), we also did

263    orthologous searches by comparing with other Hymenoptera sequences (Wurm et al.,

264    2009). To quantify variation in the numbers of protein family members, we performed

265    Pfam (version 24.0) (Bateman et al., 2004) and PROSITE profile (Sigrist et al., 2010)

266    analyses on proteins obtained from the *F. exsecta* gene set. Our final annotation

267    included gene sequences with retrieved protein-related names, functional domains,

268    and expression in other organisms along with enzyme commission (EC) numbers,

269 pathway information, Cluster of Orthologous Groups (COG), functional classes, and

270 Gene Ontology terms.

271

272 **Orthology and evolutionary rates**

273 Comparative genome-wide analysis of orthologous genes was performed with

274 OrthoVenn (Wang et al., 2015) to compare the predicted *F. exsecta* protein sequences

275 with those of four other ant species, *Camponotus floridanus*, *Lasius niger*, *Solenopsis*

276 *invicta,* and *Cerapachys biroi,* all of which were downloaded from their respective

277 public NCBI repositories. The predicted proteins of *F. exsecta* and the other four

278 species were uploaded into the OrthoVenn web server for identification and

279 comparison of orthologous clusters (Wang et al., 2015). Following clustering,

280 orthAgogue was used for the identification of putative orthology and inparalogy

281 relationships. To deduce the putative function of each ortholog, the first protein

282 sequence from each cluster was searched against the non-redundant protein database

283 UniProt using blastp 2.2.27. Pairwise sequence similarities among protein sequences

284 were determined for all species with a blastp 2.2.27 (E-value cut-off of $10^{-5}$, and an

285 inflation value of 1.5 for MCL). Finally, an interactive Venn diagram, summary

286 counts, and functional summaries of clusters shared between species were visualized

287 using OrthoVenn.

288

289 To identify genes under positive or relaxed purifying selection in *F. exsecta*, we

290 estimated the rates of non-synonymous to synonymous changes for core orthologous

291 genes (3,156) from five ant species (*F. exsecta*, *Camponotus floridanus*, *Lasius niger*,

292 *Solenopsis invicta,* and *Cerapachys biroi*). For this we only included orthologous

293 groups with one ortholog for each species (no paralogous genes were included) in the

294 analysis. We extracted coding and protein sequences for 3,156 orthologous groups

295 from the respective public NCBI repositories for the species included. We then

296 aligned all protein sequences using Clustal Omega (Sievers & Higgins, 2014), and

297 then converted them to nucleotide sequences with PAL2NAL version 14 (Yang,

298 1997). We then ran CODEML version 4.9e (Yang, 1997), using the branch site model

299 with *F. exsecta* as foreground branch, and the other five ant species as background

300 lineages. The Bayes empirical method (Yang et al. 2005) was used to estimate the

301 posterior probabilities, which were then used to identify sites under selection. We

302 additionally estimated pairwise dN/dS ratios for orthologous genes (5,148 genes)

303 between *Camponotus floridanus* and *F. exsecta* in CODEML.

304

305 We also ran an orthology analysis between the proteins from three *Wolbachia* species

306 published previously (wRi, wDac, wNo; (Klasson et al., 2009b; Ellegaard et al., 2013;

307 Ramirez-Puebla et al., 2016)), to find similarity with the predicted protein sets of the

308 newly assembled wFex genome. Orthologs were identified using OrthoVenn (E-value

309 cut-off of $10^{-5}$ and inflation value 1.5). In addition, we analyzed the paralogous genes

310 within the wFex genome, to help understand the increased genome size in comparison

311 to other *Wolbachia* genomes.

312

313 **Discovery and annotation of transposable elements**

314 We used RepeatMasker version 4.0.7 (Smit. et al., 2015), and the TransposonPSI

315 version 08-22-2010 (Brian J. Haas, 2011) to detect repetitive elements in the genome.

316 To retrieve and mask repetitive elements, we downloaded files from the Repbase and

317 Dfam databases, and aligned each of them with the *F. exsecta* genome sequences as

318 query sequences. Positive alignments were regarded as repetitive regions and

319 extracted for further analysis. To identify genome sequence region homology to

320 proteins encoded by different families of transposable elements, we used the

321 TransposonPSI analysis tool. This tool uses PSI-blast, with a collection of retro-

322 transposon ORF homology profiles to identify statistically significant alignments.

323

324 *Wolbachia* **phylogeny**

325 We analysed the phylogeny of *Wolbachia* in MrBayes v3.2.6 x64 (Ronquist &

326 Huelsenbeck, 2003), using a concatenated sequences of 35 genes. For this analysis,

327 each gene was considered as a different partition, and the most fitting nucleotide

328 substitution model was chosen for each gene, using the bayesian information criterion

329 (BIC) in the program jMODELTEST (Posada, 2008). The partitioned dataset was run

330 for 200,000 generations, sampling at every 100th generation with each partition

331 unlinked for the substitution parameters. Convergence of the runs was confirmed by

332 checking that the potential scale reduction factor was ~1.0 for all model parameters,

333 and by ensuring that an average split frequency of standard deviations < 0.01 was

334 reached (Ronquist & Huelsenbeck, 2003). The first 25% of the trees were discarded

335 as burn-in, and the remaining trees were used to create a 50% majority-rule consensus

336   tree, and to estimate the posterior probabilities. To check for consistency of the

337   phylogeny, Markov chain Monte Carlo (MCMC) runs were repeated to get a similar

338   50% majority-rule consensus tree with high posterior probabilities. The phylogenetic

339   tree generated was visualized using Figtree v1.4.2 (Rambaut, 2012).

340

341   **Results & Discussion**

342   **Assembly of the *Formica exsecta* genome**

343   We created Illumina sequencing libraries from DNA extracted from testes of males of

344   a *F. exsecta* colony to obtain >99 gigabases of Illumina sequence data. The final *F.*

345   *execta* genome resulting from assembly of this data was 277.7 megabases (Mb) long,

346   encompassing 14,617 scaffolds (Figure 1) with a N50 scaffold length of 997.7 kb

347   (Table 2). The number of scaffolds is higher than the number of chromosomes

348   reported for *F. exsecta* (n=26; Agosti & Hauschteck-Jungen, 1987; Rosengren,

349   Rosengren & Söderlund, 2009). Similarly, the *F. exsecta* genome assembly is

350   somewhat shorter than genome size estimates obtained by flow cytometry for species

351   in the subfamily Formicinae (range: 296-385 Mb; Tsutsui et al., 2008). These

352   discrepancies are unsurprising given the difficulty of assembling highly repetitive

353   gene content from short sequencing reads (Henson, Tischler & Ning, 2012). In line

354   with this, the genome assembly length metrics are similar to those of the 23 ant

355   genomes that have been published. The raw data, gapped scaffolds, and annotations

356   underpinning this assembly are deposited into public databases under BioProject

357   PRJNA393850 (accession NPMM00000000).

358

359   **Quantitative assessment of genome assembly**

360   Based on scaffold N50 and N75 statistics, contig size, and GC content, the *F. exsecta*

361   genome assembly is comparable in quality and completeness to other sequenced ant

362   genomes (Supplementary Table S1). All the 248 CEGMA eukaryotic core genes were

363   found, and 241 of these genes were complete in length. Similarly, 98.5% of 1634

364   BUSCO Insecta genes were complete in the genome (Table 3). These results held

365   with other BUSCO analysis levels including Eukaryota, Arthropoda, and

366   Hymenoptera, with low duplication levels (2.2% to 5.3%), and few missing genes

367   (0.6% to 1.27%; all details in Table 3). Such discrepancies can be due to technical

368   artifacts such as sequencing biases or assembly difficulties, as well as to true

369   differences between our *F. exsecta* sample and the BUSCO and CEGMA datasets. To

370   further evaluate genome completeness, we compared the independently generated *F.*

371   *exsecta* transcriptome (Dhaygude et al., 2017) to the genome reported here. More than

372   98.75 % of the 10.999 assembled ESTs mapped unambiguously to the genome (blastn

373   $E < 10^{-50}$). Together, these analyses show that the genome assembly has high

374   completeness.

375

376   **Gene Content in the *Formica exsecta* genome**

377   We identified 13,637 protein coding genes by combining *ab initio*, EST-based, and

378   sequence similarity based gene predictions methods. The GC content was higher in

379   exons (41.6%) than in introns (30.6%), a pattern similar to that reported in the honey

380   bee, *Apis mellifera*, and the fire ant, *Solenopsis invicta* (Weinstock et al., 2006; Wurm

381   et al., 2011). Despite this, as in other ant genomes (Schrader et al., 2014; Boomsma et

382   al., 2017), overall GC content in genes (35.1%) was similar to the rest of the genome

383   (36.0%).

384   We used blast and orthology analyses to characterize *F. exsecta* genes. The vast

385   majority (88%; 12,050) of these had the highest blastp similarity to genes in other

386   ants. A further 0.4% had the highest similarity to Apidae, and 0.6% to Braconidae,

387   Amniota, and *Wolbachia* (the latter probably due to HGT; see below and Figure 2).

388   The remaining 3.09% belong to other taxa not included in Figure 2 because they had

389   fewer than 20 hits. The remaining genes (7.91%, n= 1,080) lacked clear sequence

390   similarity [cutoff for blastx $E < 10^{-3}$] to known protein sequences or protein domains.

391   Some of these may represent erroneous gene predictions (Drăgan et al., 2016),

392   however 994 of them are ≥1000 bp and include an open reading frame >300 amino

393   acids long, which is unlikely to occur by chance. Importantly, although only a single

394   pooled transcriptome library, prepared from different developmental life stage

395   samples, was available for *F. exsecta*, 235 of the genes are expressed (FPKM ≥ 1;

396   Dhaygude et al., 2017). It is thus likely that a high proportion of the 1,080 genes are

397   taxonomically restricted genes unique to the *F. exsecta* lineage.

398   The total genes of *F. exsecta* (n=13,637) were grouped into 7,727 orthologous clusters

399   (Figure 3). Comparative analysis of the *F. exsecta* genes with the closely related

400   species *C. floridanus* and *L. niger*, and the more distantly related *S. invicta* and *C.*

401   *biroi* revealed, that 4,685 orthologous clusters out of 7,727 are shared between all five

402   species. In addition, we found 102 gene clusters that were exclusive to three

403   Formicinae genomes (*F. exsecta*, *C. floridanus* and *L. niger;* Supplementary Table

404   S2). Such genes are important candidates that could be involved in the evolution of

405   this subfamily. Many of the genes in these clusters had no detectable relation to

406   existing genes outside the Formicinae; those that did included GO annotations such as

407   glycerate kinase, transferase activity, deoxyribonucleoside diphosphate metabolic

408   process.

409   Interestingly, 633 of the *F. exsecta*-specific genes could be grouped into 197 ortholog

410   clusters of 2 or more genes (Supplementary Table S3), suggesting not only newly

411   evolved genes, but also potential gene duplication and subfunctionalisation. Previous

412   comparative genome studies have indicated that 10-20% of genes lack recognizable

413   homologs in other species in every taxonomic group so far studied (Wilson et al.,

414   2007; Khalturin et al., 2009; Johnson & Tsutsui, 2011; Tautz & Domazet-Lošo,

415   2011). Our lower percentage of orphan genes could be due to our hierarchical

416   approach to annotation, the wide range of databases used, and the large amounts of

417   ant genomic data generated over the past years (Favreau et al., 2018).

418   **Genes with signatures of evolution under positive selection**

419   We performed analyses to detect genes with signatures of positive selection in *F.*

420   *exsecta*. First, selection analysis (dN/dS ratio estimations) on 3,157 single-copy genes

421   shared between the five core ant species (without paralogous genes), revealed that 500

422   genes have signatures of positive selection in the lineage leading to *F. exsecta*. These

423   include genes involved in fatty acid metabolism, lipid catabolism, and chitin

424   metabolism (Supplementary Table S4). Interestingly, previous studies on ants, bees,

425   and flies also provide evidence for positive selection on genes in similar functional

426   categories as in our study (Roux et al., 2014). For example, genes involved in

427   biological functions such as carbohydrate metabolic processes, lipid metabolic

428   processes, cytoskeleton organization, cell surface receptor signaling pathways, and

429   RNA processing were overrepresented in the enrichment analysis, and such genes

430   were also previously reported as positively selected genes in ants, bees, and flies

431   (Viljakainen et al., 2009; Roux et al., 2014).

432

433    To perform a similar analysis on a larger number of genes, we used a second

434    approach based on pairwise comparisons between *F. exsecta* and *C. floridanus*. Out of

435    5,148 one-to-one- orthologs, 29 showed dN/dS > 1 (P < 0.005; Supplementary Table

436    S5). Although some of these putative genes could be artefactual or non-coding, they

437    all include an open reading frame of > 100 amino acids. Five (17%) out of 29 genes

438    are likely linked to transposon activity as they are transposase-like or have EpsG

439    domains. Among the other genes, only a few are annotated: the Icarapin-like protein

440    is a venom gene, and such genes have been shown to be under positive selection in

441    wasps (Werren et al., 2010). Perhaps more surprisingly we found high dN/dS for the

442    Homeobox protein gene orthopedia which is involved in early embryonic

443    development (Mackenzie et al., 1991).

444

445    **Repetitive elements**

446    Repetitive elements comprised 15.88% (44.10 Mb) of the *F. exsecta* assembly. This

447    proportion is similar to that found in other ants (16.5-31.5% (Schrader et al., 2014).

448    This is probably an underestimate because (i) genomic regions that cannot be

449    assembled are enriched with such repeats, (ii) multiple copies of a repetitive element

450    are often collapsed into a single copy during genome assembly, and (iii) only a

451    portion of repetitive elements in *F. exsecta* will have similarity to sequences in

452    standard repeat databases. Overall, 3.18% (8.8 Mb) of the assembly was composed of

453    simple repeats, whereas 12.73% (35.34 Mb) comprised interspersed repeats, most of

454    which (53.73%) could not be classified. Among those that could be classified, 10,542

455    retro element fragments represented 2.74% of the genome, and 53,438 DNA

456    transposons represented 4.23% of the genome. The *F. exsecta* genome contains copies

457    of the piggyBac transposon (23 in total, and 7 within intact ORFs). Higher numbers

458    (234) of piggyBac transposons have been found in *C. floridanus,* yet only 6 of these

459    were found within ORFs (Bonasio et al., 2010).

460

461    **The *Wolbachia* endosymbiont genome of *Formica exsecta***

462    The assembly of the *Wolbachia* endosymbiont, wFex, was 3.09 Mb long,

463    encompassing 69 scaffolds with a N50 scaffold length of 104,167 nt, and a GC

464    content of 35.13% (Table 2; GenBank: RCIU00000000, Bioproject: PRJNA436771).

465    This assembly of wFex shows extensive nucleotide similarity with the complete

466    genome of the *Wolbachia* endosymbiont of *Drosophila simulans,* wNo (GenBank ID:

467    NC_021084), and covers approximately 84% of its length (Supplementary Figure S2).

468    We determined that 549 genes are present as a single copy in the *Wolbachia* genomes

469    most closely related to wFex ((Lindsey et al., 2016) see below); 537 (99.6%) out of

470    these 539 core genes are present in the wFex genome, suggesting high completeness.

471

472    However, the wFex genome is considerably larger (3.09 Mb) than the *Wolbachia*

473    genomes reported previously (range: 0.95 to 1.66 Mb; Sun et al., 2001), and includes

474    a greater number of open reading frames (1,796 ORFs) than other published

475    *Wolbachia* genomes [range: 644 to 1,275 genes]. *Formica exsecta* is known to harbor

476    more than one *Wolbachia* strain (Reuter & Keller, 2003), thus these patterns could be

477    due to the presence of multiple endosymbiont strains. Two additional lines of

478    evidence support this idea. First, 212 genes (11.80 %), that are present as single-copy

479    genes in the wMel, wRi and wDac genomes (Klasson et al., 2009b; Ellegaard et al.,

480    2013; Ramirez-Puebla et al., 2016), are duplicated in our assembly (Supplementary

481    Table S6). Furthermore, 92 (12%) of the 775 genes present as a single copy in wFex,

482    included genetic variation within our sample, including in the cytochrome oxidase

483    subunit I; no such variation is normally expected. Despite extensive attempts, we

484    were unable to disentangle the two or more *Wolbachia* strains – this is likely because

485    differences in synteny between the strains cannot be resolved using short-read

486    sequence data. Similar assembly artifacts, due to multiple *Wolbachia* strains, have

487    also been reported by other studies (Ramírez-Puebla et al., 2016).

488

489    To determine how wFex is related to other *Wolbachia,* we used Bayesian

490    phylogenetic analysis based on 35 conserved genes (Supplementary Table S7) from

491    the 25 available *Wolbachia* genomes from the NCBI database. The analysis revealed

492    three distinct monophyletic clades, all with posterior probabilities >0.9. Each of these

493    clades represent one super group of *Wolbachia* (Figure 4). Of these three supergroups,

494    two have been found only in arthropods (super groups A and B), and the third super

495    group is found only in filarial nematodes (super group C; Werren, Baldo & Clark,

496    2008). In the phylogenetic analysis, wFex clustered with the *Wolbachia* strains within

497    super group A, and most closely matched the strain that infects the scale insect,

498    *Dactylopius coccus,* (wDacA). This is consistent with earlier studies on *Wolbachia* in

499    ants, which also found supergroup A in the majority of the infected ants (Werren &

500    Windsor, 2000).

501    Given that wFex affiliates with the supergroup A in our phylogenetic analysis, we

502    investigated the extent to which its gene content aligned with that of other *Wolbachia*

503    genomes in the same supergroup. We found that 525 genes were shared across all

504    strains in this supergroup, including wFex (Figure 5). About 20% of these genes had

505    no match to known proteins, whereas the remaining genes matched a wide range of

506    predicted functions (Ellegaard et al., 2013; Lindsey et al., 2016). We also found

507    strain-specific genes (wFex - 50 genes, wMel - 4 genes, wRi - 3 genes, wDac - 9

508    genes). The wFex-specific genes included inferred annotations including Ankyrin

509    repeat protein, ATP synthase, and chromosome partition protein (Supplementary

510    Table S8). These strain-specific genes can provide an interesting snapshot of the

511    evolutionary dynamics of a species. For example, ankyrin repeat proteins are involved

512    in numerous functional processes, and have been suggested to play an important role

513    in host-symbiont interactions (Li, Mahajan & Tsai, 2006). Comparative analyses

514    suggest that they may be involved in host communication and reproductive

515    phenotypes (Voronin & Kiseleva, 2008).

516

517    To explore differences in gene content between CI-inducing and mutualist strains of

518    *Wolbachia*, homologous genes in six CI-inducing strains, and three mutualist strains

519    were aligned and compared (Lindsey et al., 2016). The mutualist *Wolbachia* strains

520    (range: 644-805 genes) had fewer genes than the CI-inducing ones (range: 911-1,275

521    genes). The CI-inducing strains shared 84 genes not found in the mutualist strains. We

522    found 80 (95.23%) of these 84 genes in wFex (Supplementary Figure S3), suggesting

523    that wFex may be CI-inducing.

524

525    **Horizontal gene transfers, and functional novelty**

526    Intracellular symbionts can contribute new genes or fragments of genes to the host

527    genome via horizontal gene transfer (Keeling & Palmer, 2008; Werren, Baldo &

528    Clark, 2008; Dunning Hotopp, 2011). We found evidence for ancestral horizontal

529    transfer from *Wolbachia* to the host *F. exsecta* in five scaffolds (scaffold83,

530    scaffold233, scaffold574, scaffold707, scaffold741). The four largest transfers are 13

531    to 47 kb long, and include 83 putative functional protein coding genes, whereas the

532    fifth and smallest insertion (475 bp) lacks protein coding genes other than a

533    degenerate *Wolbachia* transposase. This transposase is present in 7 out of 29

534    published *Wolbachia* genomes. Our analysis shows that similar transfer events of this

535    homologous fragment apparently also have occurred from *Wolbachia* to the genomes

536    of the ants *Vollenhovia emeryi* (gene: LOC105557741), and *Cardiocondyla obscurior*

537    (scaffolds scf7180001101632 and scf7180001108526), as well as the microfilarial

538    nematode *Brugia pahangi*, the Arizona spittle bug *Clastoptera arizona,* and the

539    parasitoid wasp *Diachasma alloeum*.

540    One-third of invertebrate genomes are thought to contain recent *Wolbachia* gene

541    insertions, ranging in size from short segments (<600 bp), to nearly the entire genome

542    (Hotopp et al., 2007; Werren, Baldo & Clark, 2008). Most of these transferred

543    fragments contained transposable elements, as well as some other functional genes

544    from the *Wolbachia* genome. The HGT events from *Wolbachia* to *F. exsecta* are

545    located in or near regions with transposases. Our blast results suggest that four of the

546    insert regions had *Wolbachia* transposases, whereas one insert region has a

547    transposase of ant origin. Whether the presence of such transposases close to HGT

548    sites facilitates insertions is unknown. Interestingly, the putative functional protein-

549    coding genes of *Wolbachia* inserted in the *F. exsecta* genome are similar to the genes

550    reported in similar HGTs events in other insect genomes (eg: ABC transporter,

551    Ankyrin repeat containing protein (Table 4) (Brelsfoard et al., 2014; International

552    Glossina Genome Initiative, 2014). This could indicate that some HGT events are

553    either more likely to occur or to be retained for reasons that could be neutral or

554    adaptive to the host or to the endosymbiont. The transcriptome of *F. exsecta* shows

555    that at least 6 out of the 83 genes from the *Wolbachia* HGT regions are transcribed

556    but with a low FPKM values (range 0.04 to 1.6). These low level transcription trait

557    often observed in bacteria-eukaryote HGTs (Hotopp et al., 2007; Nikoh et al., 2008;

558    Dunning Hotopp, 2011).

559    **Conclusions**

560    Here we present the first draft genome of the ant *F. exsecta,* and its *Wolbachia*

561    endosymbiont. This is the first report of a *Wolbachia* genome from ants, and provides

562    insights into its phylogenetic position. We further identified multiple HGT events

563    from *Wolbachia* to *F. exsecta*. Some of these have also occurred in parallel in several

564    other insect genomes, highlighting the extent of HGTs in eukaryotes. We expect that

565    the *F. exsecta* genome will be a valuable resource in understanding the molecular

566 basis of the evolution of social organization in ants: Recent genomic comparisons

567 between *Formica selysi* and *S. invicta* have shown convergent evolution of a social

568 chromosome, that underpins social organisation in these ants (Purcell et al., 2014).

569 Additional comparison of these genomic regions with *F. exsecta* could provide

570 valuable insights on the evolution of genomic architectures underlying social

571 organization.

572

## Acknowledgements

582

## References

584
585 Agosti D., Hauschteck-Jungen E. 1987. Polymorphism of males in *Formica exsecta*
586 Nyl. (Hym.: Formicidae). Insectes Sociaux 34:280–290. DOI: 10.1007/BF02224360.
587
588 Aikawa et.al., 2009. Longicorn beetle that vectors pinewood nematode carries many
589 *Wolbachia* genes on an autosome. Proceedings of the Royal Society B: Biological
590 Sciences 276:3791–3798. DOI: 10.1098/rspb.2009.1022.
591
592 Altschul et.al., 1990. Basic local alignment search tool. Journal of Molecular Biology
593 215:403–410. DOI: 10.1016/S0022-2836(05)80360-2.
594
595 Bateman et.al., 2004. The Pfam protein families database. Nucleic Acids Research
596 32:D138-41. DOI: 10.1093/nar/gkh121.
597
598 Bonasio et.al 2010. Genomic comparison of the ants *Camponotus floridanus* and
599 *Harpegnathos saltator*. Science 329:1068–1071. DOI: 10.1126/science.1192428.
600
601 Boomsma et.al., 2017. The Global Ant Genomics Alliance (GAGA). Myrmecological
602 News 25:61–66.
603

604    Boto L. 2014. Horizontal gene transfer in the acquisition of novel traits by metazoans.
605    Proceedings of the Royal Society B: Biological Sciences 281:20132450. DOI:
606    10.1098/rspb.2013.2450.
607
608
609    Brelsfoard et.al., 2014. Presence of extensive Wolbachia symbiont insertions
610    discovered in the genome of its host *Glossina morsitans morsitans*. PLoS Neglected
611    Tropical Diseases 8:e2728. DOI: 10.1371/journal.pntd.0002728.
612
613    Brian J. Haas. 2011. TransposonPSI. http://transposonpsi.sourceforge.net
614
615    Bruto et.al., 2014. Frequent, independent transfers of a catabolic gene from bacteria to
616    contrasted filamentous eukaryotes. Proceedings of the Royal Society of London B:
617    Biological Sciences. 281(1789): 20140848. DOI:  10.1098/rspb.2014.0848
618
619    Cantarel et.al., 2008. MAKER: an easy-to-use annotation pipeline designed for
620    emerging model organism genomes. Genome Research 18:188–96. DOI:
621    10.1101/gr.6743907.
622
623
624    Casacuberta E., González J. 2013. The impact of transposable elements in
625    environmental adaptation. Molecular Ecology 22:1503–1517. DOI:
626    10.1111/mec.12170.
627
628    Cordaux R., Michel-Salzat A., Bouchon D. 2001. Wolbachia infection in crustaceans:
629    novel hosts and potential routes for horizontal transmission. Journal of Evolutionary
630    Biology 14:237–243. DOI: 10.1046/j.1420-9101.2001.00279.x.
631
632    Crozier RH. 1975. Hymenoptera. Animal Cytogenetics:95. ISBN:9783443260040
633
634    Dhaygude et.al. 2017. Transcriptome sequencing reveals high isoform diversity in the
635    ant *Formica exsecta*. PeerJ 5:e3998. DOI: 10.7717/peerj.3998.
636
637    Drăgan et.al., 2016. GeneValidator: identify problems with protein-coding gene
638    predictions. Bioinformatics 32:1559–1561. DOI: 10.1093/bioinformatics/btw015.
639
640    Dunning Hotopp JC. 2011. Horizontal gene transfer between bacteria and animals.
641    Trends in Genetics 27:157–163. DOI: 10.1016/j.tig.2011.01.005.
642
643    Ellegaard et.al., 2013. Comparative genomics of Wolbachia and the bacterial species
644    concept. PLoS Genetics 9:e1003381. DOI: 10.1371/journal.pgen.1003381.
645
646    Favreau et.al., 2018. Genes and genomic processes underpinning the social lives of
647    ants. Current Opinion in Insect Science 25:83–90. DOI: 10.1016/J.COIS.2017.12.001.
648
649    Fenn et.al., 2006. Phylogenetic relationships of the Wolbachia of nematodes and
650    arthropods. PLoS Pathogens 2:e94. DOI: 10.1371/journal.ppat.0020094.
651

652    Ferree et al. 2005. Wolbachia utilizes host microtubules and dynein for anterior
653    localization in the *Drosophila* oocyte. PLoS Pathogens 1:e14. DOI:
654    10.1371/journal.ppat.0010014.
655
656    Fitzpatrick DA. 2012. Horizontal gene transfer in fungi. FEMS Microbiology Letters
657    329:1–8. DOI: 10.1111/j.1574-6968.2011.02465.x.
658
659    González et.al., 2010. Genome-wide patterns of adaptation to temperate environments
660    associated with transposable elements in Drosophila. PLoS Genetics 6:e1000905.
661    DOI: 10.1371/journal.pgen.1000905.
662
663    Goodacre et.al. 2006. Wolbachia and other endosymbiont infections in spiders.
664    Molecular Ecology 15:517–527. DOI: 10.1111/j.1365-294X.2005.02802.x.
665
666    Gurevich A., Saveliev V., Vyahhi N., Tesler G. 2013. QUAST: quality assessment
667    tool for genome assemblies. Bioinformatics 29:1072–1075. DOI:
668    10.1093/bioinformatics/btt086.
669
670    Haag-Liautard et al., 2009. Fitness and the level of homozygosity in a social insect.
671    Journal of Evolutionary Biology 22:134–42. DOI: 10.1111/j.1420-9101.2008.01635.x
672
673    Henson J., Tischler G., Ning Z. 2012. Next-generation sequencing and large genome
674    assemblies. Pharmacogenomics 13:901–15. DOI: 10.2217/pgs.12.72.
675
676    Holt C., Yandell M. 2011. MAKER2: an annotation pipeline and genome-database
677    management tool for second-generation genome projects. BMC Bioinformatics
678    12:491. DOI: 10.1186/1471-2105-12-491.
679
680    Hotopp et.al., 2007. Widespread lateral gene transfer from intracellular bacteria to
681    multicellular eukaryotes. Science 317:1753–1756. DOI: 10.1126/science.1142490.
682
683    Huang Y., Niu B., Gao Y., Fu L., Li W. 2010. CD-HIT Suite: a web server for
684    clustering and comparing biological sequences. Bioinformatics 26:680–682. DOI:
685    10.1093/bioinformatics/btq003.
686
687    International Glossina Genome Initiative. 2014. Genome sequence of the tsetse fly
688    (Glossina morsitans): vector of African trypanosomiasis. Science (New York, N.Y.)
689    344:380–6. DOI: 10.1126/science.1249656.
690
691    Johnson BR., Tsutsui ND. 2011. Taxonomically restricted genes are associated with
692    the evolution of sociality in the honey bee. BMC Genomics 12:164. DOI:
693    10.1186/1471-2164-12-164.
694
695    Keeling PJ., Palmer JD. 2008. Horizontal gene transfer in eukaryotic evolution.
696    Nature Reviews Genetics 9:605–618. DOI: 10.1038/nrg2386.
697
698    Keller et.al., 2001. Sex ratio and Wolbachia infection in the ant *Formica exsecta*.
699    Heredity 87:227–33.DOI:10.1046/j.1365-2540.2001.00918.x
700

701  Khalturin et.al. 2009. More than just orphans: are taxonomically-restricted genes
702  important in evolution? Trends in Genetics 25:404–413. DOI:
703  10.1016/j.tig.2009.07.006.
704
705  Klasson et.al., 2009a. Horizontal gene transfer between *Wolbachia* and the mosquito
706  *Aedes aegypti*. BMC Genomics 10:33. DOI: 10.1186/1471-2164-10-33.
707
708  Klasson et.al. 2009b. The mosaic genome structure of the *Wolbachia* wRi strain
709  infecting *Drosophila simulans*. Proceedings of the National Academy of Sciences
710  106:5725–5730. DOI: 10.1073/pnas.0810753106.
711
712  Kondo et.al. 2002. Genome fragment of *Wolbachia* endosymbiont transferred to X
713  chromosome of host insect. Proceedings of the National Academy of Sciences
714  99:14280–14285. DOI: 10.1073/pnas.222228199.
715
716  Korf I. 2004. Gene finding in novel genomes. BMC Bioinformatics 5:59. DOI:
717  10.1186/1471-2105-5-59.
718
719  Koski LB., Gray MW., Lang BF., Burger G. 2005. AutoFACT: an automatic
720  functional annotation and classification tool. BMC Bioinformatics 6:151. DOI:
721  10.1186/1471-2105-6-151.
722
723  Kumar et.al. 2013. Blobology: exploring raw genome data for contaminants,
724  symbionts and parasites using taxon-annotated GC-coverage plots. Frontiers in
725  Genetics 4:237. DOI: 10.3389/fgene.2013.00237.
726
727  Kurtz et.al. 2004. Versatile and open software for comparing large genomes. 5.
728  Genome Biology 20045:R12. DOI:10.1186/gb-2004-5-2-r12
729
730  Lasken RS., Stockwell TB. 2007. Mechanism of chimera formation during the
731  Multiple Displacement Amplification reaction. BMC Biotechnology 7:19. DOI:
732  10.1186/1472-6750-7-19.
733
734  Li J., Mahajan A., Tsai M-D. 2006. Ankyrin Repeat: A unique motif mediating
735  protein–protein interactions. Biochemistry 45:15168–15178. DOI:
736  10.1021/bi062188q.
737
738  Lindsey ARI., Werren JH., Richards S., Stouthamer R. 2016. Comparative genomics
739  of a parthenogenesis-inducing *Wolbachia* symbiont. G3 (Bethesda, Md.) 6:2113–23.
740  DOI: 10.1534/g3.116.028449.
741
742  Mackenzie A., Leeming GL., Jowett AK., Ferguson MW., Sharpe PT. 1991. The
743  homeobox gene Hox 7.1 has specific regional and temporal expression patterns
744  during early murine craniofacial embryogenesis, especially tooth development in vivo
745  and in vitro. Development (Cambridge, England) 111:269–85.
746
747  Magrane M., Consortium U. 2011. UniProt Knowledgebase: a hub of integrated
748  protein data. Database 2011:bar009-bar009. DOI: 10.1093/database/bar009.
749

750  Matveeva T V., Lutova LA. 2014. Horizontal gene transfer from Agrobacterium to
751  plants. Frontiers in Plant Science 5:326. DOI: 10.3389/fpls.2014.00326.
752
753  McNulty SN., Fischer K., Curtis KC., Weil GJ., Brattig NW., Fischer PU. 2013.
754  Localization of Wolbachia-like gene transcripts and peptides in adult Onchocerca
755  flexuosa worms indicates tissue specific expression. Parasites & Vectors 6:2. DOI:
756  10.1186/1756-3305-6-2.
757
758  Moya A., Peretó J., Gil R., Latorre A. 2008. Learning how to live together: genomic
759  insights into prokaryote–animal symbioses. Nature Reviews Genetics 9:218–229.
760  DOI: 10.1038/nrg2319.
761
762  Nikoh N., McCutcheon JP., Kudo T., Miyagishima S., Moran NA., Nakabachi A.
763  2010. Bacterial genes in the aphid genome: absence of functional gene transfer from
764  Buchnera to its host. PLoS Genetics 6:e1000827. DOI:
765  10.1371/journal.pgen.1000827.
766
767  Nikoh N., Nakabachi A. 2009. Aphids acquired symbiotic genes via lateral gene
768  transfer. BMC Biology 7:12. DOI: 10.1186/1741-7007-7-12.
769
770  Nikoh et.al. 2008. Wolbachia genome integrated in an insect chromosome: evolution
771  and fate of laterally transferred endosymbiont genes. Genome Research 18:272–80.
772  DOI: 10.1101/gr.7144908.
773
774  Normark BB. 2003. The evolution of alternative genetic systems in insects. Annual
775  Review of Entomology 48:397–423. DOI: 10.1146/annurev.ento.48.091801.112703.
776
777  O'Leary et.al. 2016. Reference sequence (RefSeq) database at NCBI: current status,
778  taxonomic expansion, and functional annotation. Nucleic Acids Research 44:D733–
779  D745. DOI: 10.1093/nar/gkv1189.
780
781  Parra G., Bradnam K., Korf I., Bateman A. 2007. CEGMA: a pipeline to accurately
782  annotate core genes in eukaryotic genomes. Bioinformatics 23:1061–1067. DOI:
783  10.1093/bioinformatics/btm071.
784
785  Purcell J., Brelsford A., Wurm Y., Perrin N., Chapuisat M. 2014. Convergent genetic
786  architecture underlies social organization in ants. Current Biology : CB 24:2728–32.
787  DOI: 10.1016/j.cub.2014.09.071.
788
789  Rambaut A. 2012. Figtree Tool (http://tree.bio.ed.ac.uk/software/figtree/)
790
791  Ramírez-Puebla et.al., 2016. Genomes of Candidatus Wolbachia bourtzisii wDacA
792  and Candidatus Wolbachia pipientis wDacB from the cochineal insect Dactylopius
793  coccus (Hemiptera: Dactylopiidae). G3 (Bethesda, Md.) 6:3343–3349. DOI:
794  10.1534/g3.116.031237.
795
796  Reuter M., Keller L. 2003. High levels of multiple Wolbachia infection and
797  recombination in the ant *Formica exsecta*. Molecular Biology and Evolution 20:748–
798  753. DOI: 10.1093/molbev/msg082.
799

800  Rolland T., Neuvéglise C., Sacerdot C., Dujon B. 2009. Insertion of horizontally
801  transferred genes within conserved syntenic regions of Yeast genomes. PLoS One
802  4:e6515. DOI: 10.1371/journal.pone.0006515.

804  Ronquist F., Huelsenbeck JP. 2003. MrBayes 3: Bayesian phylogenetic inference
805  under mixed models. Bioinformatics (Oxford, England) 19:1572–4.

807  Rosengren M., Rosengren R., Söderlund V. 2009. Chromosome numbers in the genus
808  Formica with special reference to the taxonomical position of *Formica uralensis*
809  *Ruzsk.* and *Formica truncorum* Fabr. Hereditas 92:321–325. DOI: 10.1111/j.1601-
810  5223.1980.tb01715.x.

812  Rostant WG., Wedell N., Hosken DJ. 2012. Transposable elements and insecticide
813  resistance. Advances in Genetics 78:169–201. DOI: 10.1016/B978-0-12-394394-
814  1.00002-X.

816  Roux J., Privman E., Moretti S., Daub JT., Robinson-Rechavi M., Keller L. 2014.
817  Patterns of positive selection in seven ant genomes. Molecular Biology and Evolution
818  31:1661–1685. DOI: 10.1093/molbev/msu141.

820  Salzberg SL., Delcher AL., Kasif S., White O. 1998. Microbial gene identification
821  using interpolated Markov models. Nucleic Acids Research 26:544–8.

823  Schönknecht G., Weber APM., Lercher MJ. 2014. Horizontal gene acquisitions by
824  eukaryotes as drivers of adaptive evolution. BioEssays 36:9–20. DOI:
825  10.1002/bies.201300095.

827  Schrader et.al. 2014. Transposable element islands facilitate adaptation to novel
828  environments in an invasive species. Nature Communications 5:5495. DOI:
829  10.1038/ncomms6495.

831  Sievers F., Higgins DG. 2014. Clustal Omega, accurate alignment of very large
832  numbers of sequences. In: Methods in Molecular Biology (Clifton, N.J.). 105–116.
833  DOI: 10.1007/978-1-62703-646-7_6.

835  Sigrist CJA., Cerutti L., de Castro E., Langendijk-Genevaux PS., Bulliard V., Bairoch
836  A., Hulo N. 2010. PROSITE, a protein domain database for functional
837  characterization and annotation. Nucleic Acids Research 38:D161-6. DOI:
838  10.1093/nar/gkp885.

840  Simão  et.al. 2015. BUSCO: assessing genome assembly and annotation completeness
841  with single-copy orthologs. Bioinformatics 31:3210–3212. DOI:
842  10.1093/bioinformatics/btv351.

844  Simpson JT., Durbin R. 2012. Efficient de novo assembly of large genomes using
845  compressed data structures. Genome Research 22:549–556. DOI:
846  10.1101/gr.126953.111.

848  Slater GSC., Birney E. 2005. Automated generation of heuristics for biological
849  sequence comparison. BMC Bioinformatics 6:31. DOI: 10.1186/1471-2105-6-31.

850

851   Smit., AFA., Hubley R., Green P. 2015. RepeatMasker Open-4.0.
852   http://www.repeatmasker.org/

853

854   Stanke M., Morgenstern B. 2005. AUGUSTUS: a web server for gene prediction in
855   eukaryotes that allows user-defined constraints. Nucleic Acids Research 33:W465-7.
856   DOI: 10.1093/nar/gki458.

857

858   Sun L V., Foster JM., Tzertzinis G., Ono M., Bandi C., Slatko BE., O'Neill SL. 2001.
859   Determination of Wolbachia genome size by pulsed-field gel electrophoresis. Journal
860   of Bacteriology 183:2219–25. DOI: 10.1128/JB.183.7.2219-2225.2001.

861

862   Sundström L. 1994. Sex ratio bias, relatedness asymmetry and queen mating
863   frequency in ants. Nature. DOI:10.1038/367266a0

864

865   Sundström L., Chapuisat M., Keller L. 1996. Conditional manipulation of sex ratios
866   by ant workers: a test of kin selection theory. Science 274:993–995. DOI:
867   10.1126/science.274.5289.993

868

869   Sundström L., Keller L., Chapuisat M. 2003. Inbreeding and sex-biased gene flow in
870   the ant Formica exsecta. Evolution; international journal of organic evolution
871   57:1552–61. DOI:10.1111/j.0014-3820.2003.tb00363.x

872

873   Tautz D., Domazet-Lošo T. 2011. The evolutionary origin of orphan genes. Nature
874   Reviews Genetics 12:692–702. DOI: 10.1038/nrg3053.

875

876   The Uniprot Consortium. 2017. UniProt: the universal protein knowledgebase.
877   Nucleic Acids Research 45:D158–D169. DOI: 10.1093/nar/gkw1099.

878

879   Trapnell C., Pachter L., Salzberg SL. 2009. TopHat: discovering splice junctions with
880   RNA-Seq. Bioinformatics 25:1105–1111. DOI: 10.1093/bioinformatics/btp120.

881

882   Trapnell et.al. 2010. Transcript assembly and quantification by RNA-Seq reveals
883   unannotated transcripts and isoform switching during cell differentiation. Nature
884   Biotechnology 28:511–515. DOI: 10.1038/nbt.1621.

885

886   Tsutsui ND., Suarez A V., Spagna JC., Johnston JS. 2008. The evolution of genome
887   size in ants. BMC Evolutionary Biology 8:64. DOI: 10.1186/1471-2148-8-64.

888

889   Viljakainen L., Evans JD., Hasselmann M., Rueppell O., Tingek S., Pamilo P. 2009.
890   Rapid evolution of immune proteins in social insects. Molecular Biology and
891   Evolution 26:1791–1801. DOI: 10.1093/molbev/msp086.

892

893   Vitikainen E., Haag-Liautard C., Sundström L. 2011. Inbreeding and reproductive
894   investment in the ant *Formica exsecta*. Evolution 65. DOI: 10.1111/j.1558-
895   5646.2011.01273.x.

896

897   Vitikainen EIK., Haag-Liautard C., Sundström L. 2015. Natal dispersal, mating
898   patterns, and inbreeding in the ant *Formica exsecta*. The American naturalist
899   186:716–27. DOI: 10.1086/683799.

900
901    Voronin DA., Kiseleva E V. 2008. Functional role of proteins containing ankyrin
902    repeats. Cell and Tissue Biology 2:1–12. DOI: 10.1134/S1990519X0801001X.
903
904    Wang Y., Coleman-Derr D., Chen G., Gu YQ. 2015. OrthoVenn: a web server for
905    genome wide comparison and annotation of orthologous clusters across multiple
906    species. Nucleic Acids Research 43:W78-84. DOI: 10.1093/nar/gkv487.
907
908    Weinstock et.al., 2006. Insights into social insects from the genome of the honeybee
909    *Apis mellifera*. Nature 443:931–949. DOI: 10.1038/nature05260.
910
911    Wenseleers T., Ito F., Van Borm S., Huybrechts R., Volckaert F., Billen J. 1998.
912    Widespread occurrence of the microorganism Wolbachia in ants. Proceedings of the
913    Royal Society B: Biological Sciences 265:1447–1452. DOI: 10.1098/rspb.1998.0456.
914
915    Werren JH. 1997. Wolbachia run amok. Proceedings of the National Academy of
916    Sciences of the United States of America 94:11154–5.
917
918    Werren JH., Baldo L., Clark ME. 2008. Wolbachia: master manipulators of
919    invertebrate biology. Nature Reviews Microbiology 6:741–751. DOI:
920    10.1038/nrmicro1969.
921
922    Werren et.al. 2010. Functional and evolutionary insights from the genomes of three
923    parasitoid Nasonia species. Science 327:343–348. DOI: 10.1126/science.1178028.
924
925
926    Werren JH., Windsor DM. 2000. Wolbachia infection frequencies in insects: evidence
927    of a global equilibrium? Proceedings. Biological sciences 267:1277–85. DOI:
928    10.1098/rspb.2000.1139.
929
930    Wilson GA., Feil EJ., Lilley AK., Field D. 2007. Large-scale comparative genomic
931    ranking of taxonomically restricted genes (TRGs) in bacterial and archaeal genomes.
932    PLoS One 2:e324. DOI: 10.1371/journal.pone.0000324.
933
934    Woolfit M., Iturbe-Ormaetxe I., McGraw EA., O'Neill SL. 2009. An ancient
935    horizontal gene transfer between mosquito and the endosymbiotic bacterium
936    *Wolbachia pipientis*. Molecular Biology and Evolution 26:367–374. DOI:
937    10.1093/molbev/msn253.
938
939    Wurm et.al. 2009. Fourmidable: a database for ant genomics. BMC Genomics 10:5.
940    DOI: 10.1186/1471-2164-10-5.
941
942    Wurm et.al., 2011. The genome of the fire ant *Solenopsis invicta*. Proceedings of the
943    National Academy of Sciences 108:5679–5684. DOI: 10.1073/pnas.1009690108.
944
945    Wybouw N., Pauchet Y., Heckel DG., Van Leeuwen T. 2016. Horizontal gene
946    transfer contributes to the evolution of arthropod herbivory. Genome Biology and
947    Evolution 8:1785–801. DOI: 10.1093/gbe/evw119.
948

949 Xie et.al. 2014. SOAPdenovo-Trans: de novo transcriptome assembly with short
950 RNA-Seq reads. Bioinformatics 30:1660–1666. DOI: 10.1093/bioinformatics/btu077.
951
952 Yang Z. 1997. PAML: a program package for phylogenetic analysis by maximum
953 likelihood. Computer Applications in the Biosciences□: CABIOS 13:555–6.
954
955 Yue J., Hu X., Sun H., Yang Y., Huang J. 2012. Widespread impact of horizontal
956 gene transfer on plant colonization of land. Nature Communications 3:1152. DOI:
957 10.1038/ncomms2148.
958
959 Zerbino DR., Birney E. 2008. Velvet: Algorithms for de novo short read assembly
960 using de Bruijn graphs. Genome Research 18:821–829. DOI: 10.1101/gr.074492.107.
961
962

963 **Data Accessibility**

964

965 The raw Illumina sequences of paired-end and mate-pair libraries are deposited on the

966 National Center for Biotechnology Information (NCBI) under the bio-project number

967 PRJNA393850, with the accession numbers SAMN07344805-SAMN07344811. The

968 assembled genome sequence of *F. exsecta* is deposited on Genbank with the accession

969 number NPMM00000000. Similarly, the draft genome assembly of wFex is deposited

970 under the project number PRJNA436771.

971

972

973 **List of Tables:**

974

975 Table 1: Summary statistics for the raw sequencing data, before and after filtering

976 reads. "Coverage depth" was calculated based on the estimated assembled genome

977 size (300 Mb).

978

979 Table 2: Genome assembly statistics for *F. exsecta* and its *Wolbachia* endosymbiont.

980

981 Table 3: BUSCO quality metrics for the *F. exsecta* genome and the *Wolbachia*

982 endosymbiont *of F. exsecta (wFex)* genome assembly.

983

984 Table 4: HGT inserts from *Wolbachia* present in the genome of *F. exsecta* with details

985 of length and position in the *F. exsecta* genome. The presence of similar insert regions

986 in other eukaryote genomes is also shown.

987

**List of Figures:**

989

990 Figure 1. *De novo* genome assembly of *F. exsecta* genome, summarized by the

991 following metrics: a) Overall assembly length, b) Number of scaffolds/contigs, c)

992 Length of the longest scaffold/contig, d) Scaffold/contig N50 and N90, e) Percentage

993 GCs and percentage Ns , f) BUSCO completeness, g) Scaffold/contig length/count

994 distribution.

995

996 Figure 2. Taxonomic distribution of the best blastp hits of *F. exsecta* proteins to the

997 non-redundant (nr) protein database (E < $10^{-5}$).

998

999 Figure 3. Venn diagram showing the distribution of gene families (orthologous

1000 clusters) among five ant species including three closely related members of the

1001 subfamily Formicinae (*Formica exsecta, Camponotus floridanus, Lasius niger*)*,* and

1002 two distinctly related ants *(Solenopsis invicta* and *Cerapachys biroi).*

1003

1004 Figure 4: Phylogeny of the *Wolbachia* supergroups A, B, and C strains with the newly

1005 assembled wFex genome. The phylogenetic reconstructions are based on individual

1006 analyses of 35 core genes of 25 *Wolbachia* strains. The support values on the branch

1007 labels indicate Bayesian posterior probabilities. The letters A-C indicate the separate

1008 supergroups.

1009

1010 Figure 5. Venn diagram displaying the overlap in orthologous genes among four

1011 *Wolbachia* species including the newly assembled wFex strain and the wDac, wRi,

1012 wMel strains reported previously.

1013 **Supplementary Tables:**

1014 S1:  Comparison of assembly statistics of the *F. exsecta* genome and 13 other

1015 published ant genomes.

1016 S2: List of genes specific to the Formicinae as identified by OrthoVenn.

1017 S3: List of species-specific genes in *F. exsecta,* as identified by OrthoVenn.

1018    S4: List of *F. exsecta* genes under positive or relaxed purifying selection (dN/dS

1019    ratios > 1) in comparison to five other ant species (*Camponotus floridanus, Lasius*

1020    *niger, Solenopsis invicta* and *Cerapachys biroi*)

1021    S5: List of *F. exsecta* genes showing dN/dS ratios > 1 in pairwise comparison to

1022    *Camponotus floridanus.*

1023    S6: List of genes with paralogs in the wFex genome, which are present as single

1024    copies in the wMel, wRi, wDac genomes.

1025    S7: List of conserved *Wolbachia* genes used for phylogenetic analysis.

1026    S8: List of species-specific genes in wFEX genome, as identified by OrthoVenn.

1027    **Supplementary Figures:**

1028    S1. TAGC plot of *F. exsecta,* and its *Wolbachia* endosymbiont. The TAGC plots were

1029    taxonomically annotated, and the contigs with best similarity to Arthropoda and

1030    Proteobacteria are highlighted in color.

1031

1032    S2. Visualization of genome coverage of wFex against the *Wolbachia* endosymbiont

1033    of *Drosophila simulans* (wNo) genome, using the alignment software Mummer.

1034

1035    S3. Venn diagram displaying the overlap in orthologous genes across CI-inducing and

1036    mutualist *Wolbachia* species.

**Figure 1** The de novo genome assembly of the *F. exsecta* genome summarized by different metrics: a) Overall assembly length, b) Number of scaffolds/contigs, c) Length of longest scaffold/contig, d) Scaffold/contig N50 and N90, e) Percentage GC and percentage Ns, f) BUSCO completeness. g) Scaffold/contig length/count distribution.
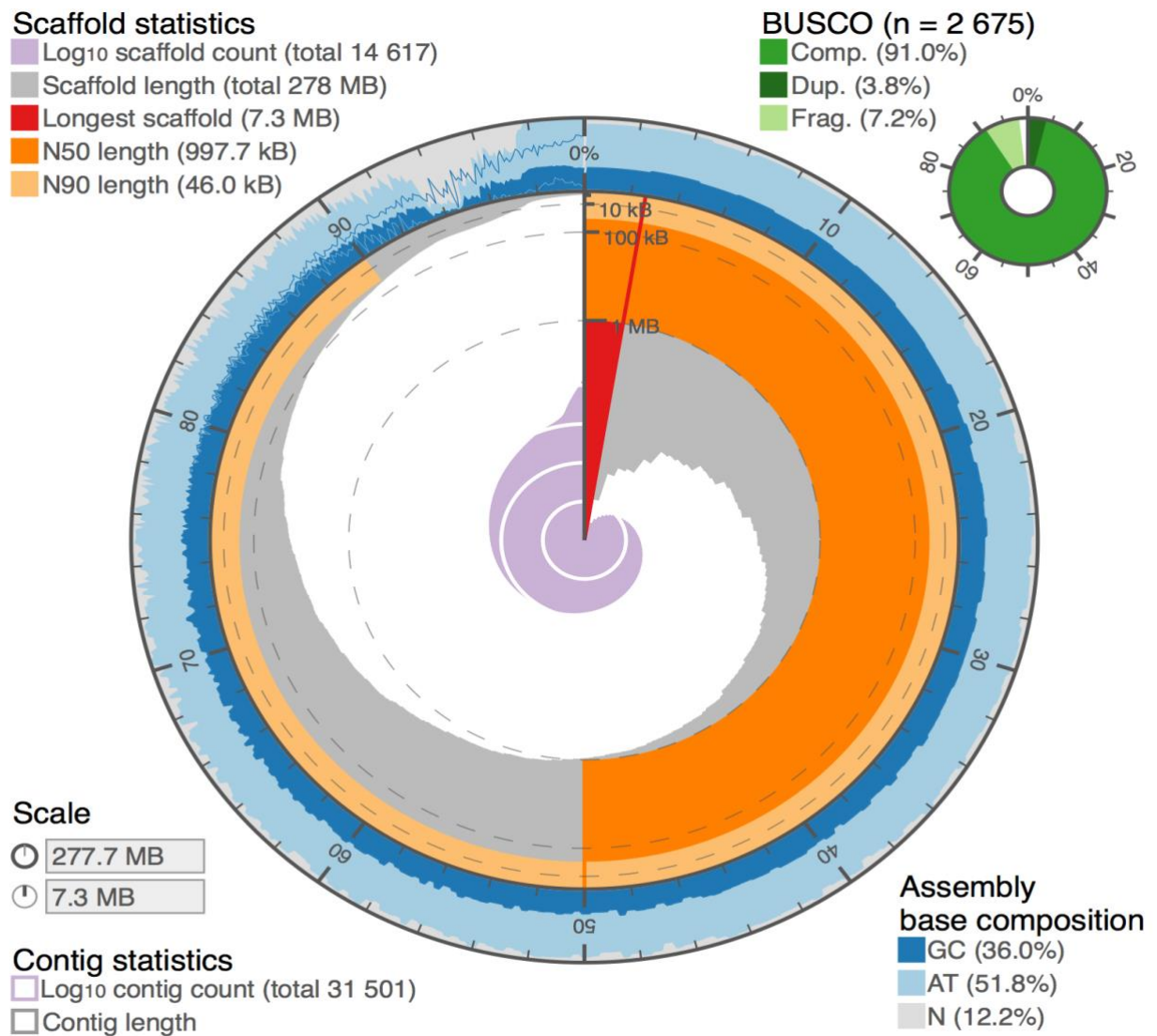
**Figure 2.** Taxonomic distribution of best blastP hits of *F. exsecta* proteins to the nonredundant (nr) protein database (E < 10−5).
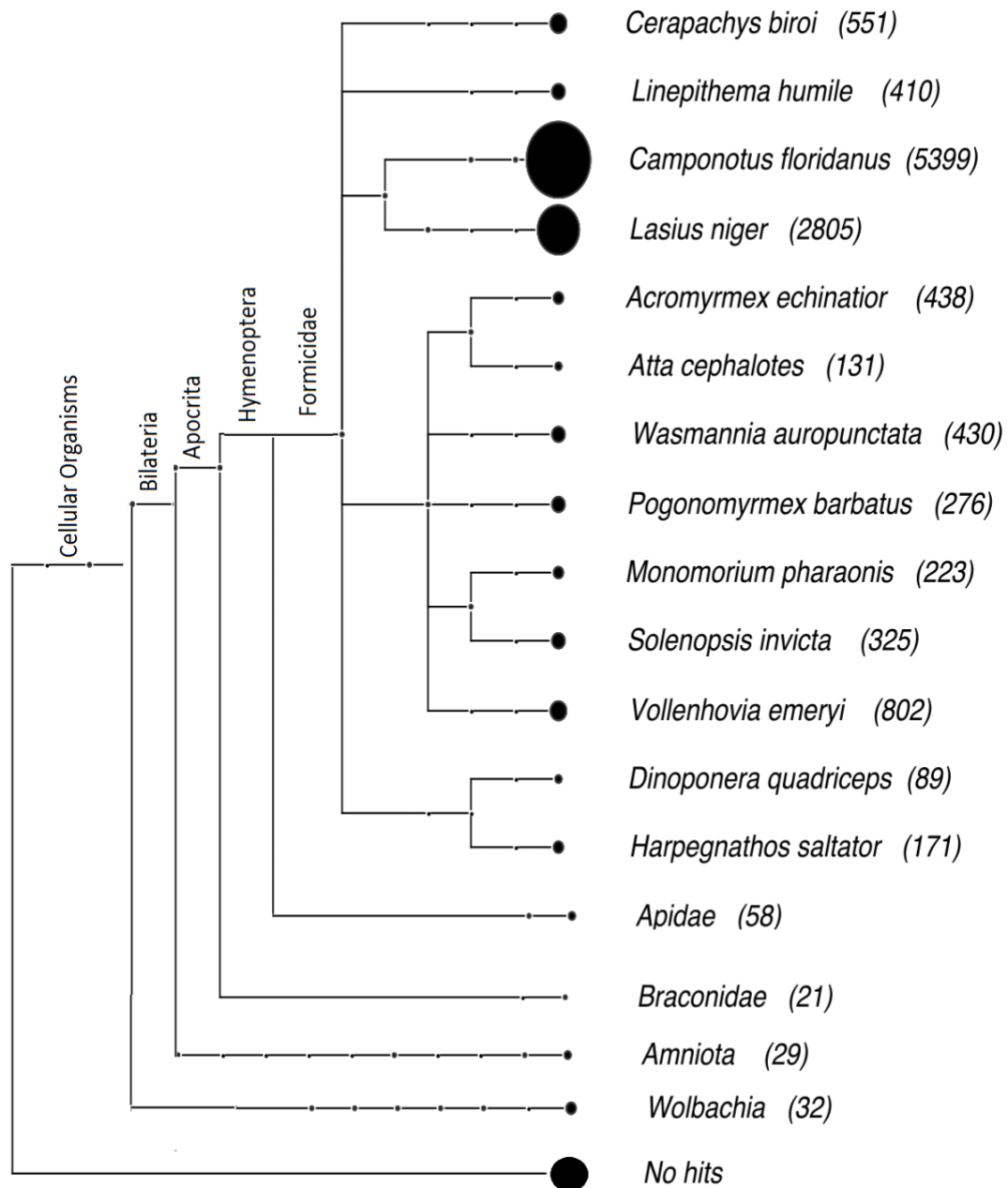
**Figure 3.** Venn diagram showing the distribution of gene families (orthologous clusters) among five ant species including closely related three members of subfamily Formicinae (*Formica exsecta, Camponotus floridanus, Lasius niger)* and other two distinctly related ants *(Solenopsis invicta* and *Cerapachys biroi).*
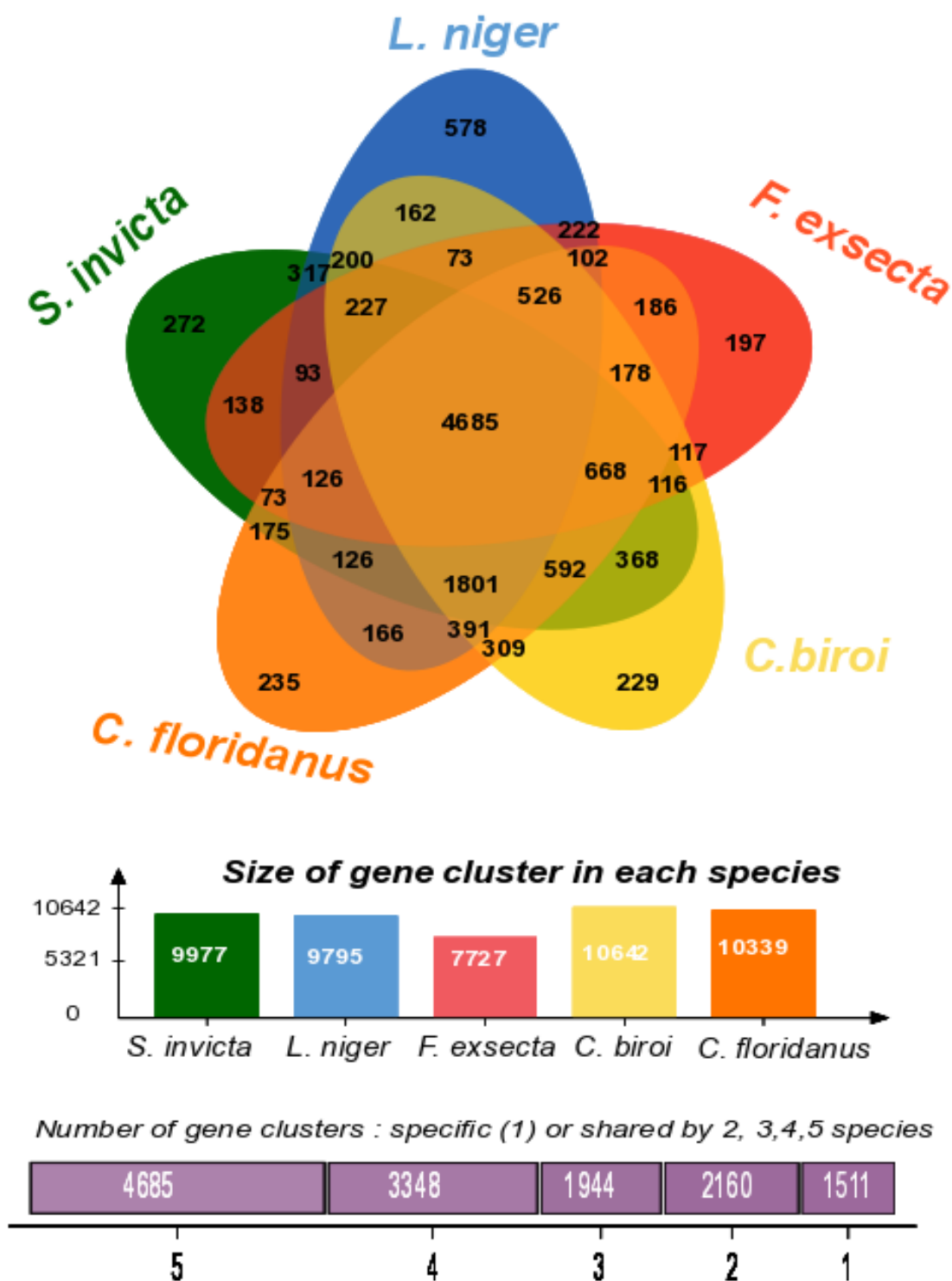
**Figure 4.** Phylogeny of the *Wolbachia* supergroups A, B, and C strains with the newly assembled wFex genome. Phylogenetic reconstructions based on individual analyses of 35 core gene of 25 *Wolbachia* strains. The numbers at the node indicate the posterior probabilities obtained from Bayesian phylogenetic analysis. Each Supergroup is labeled with letters A-C.
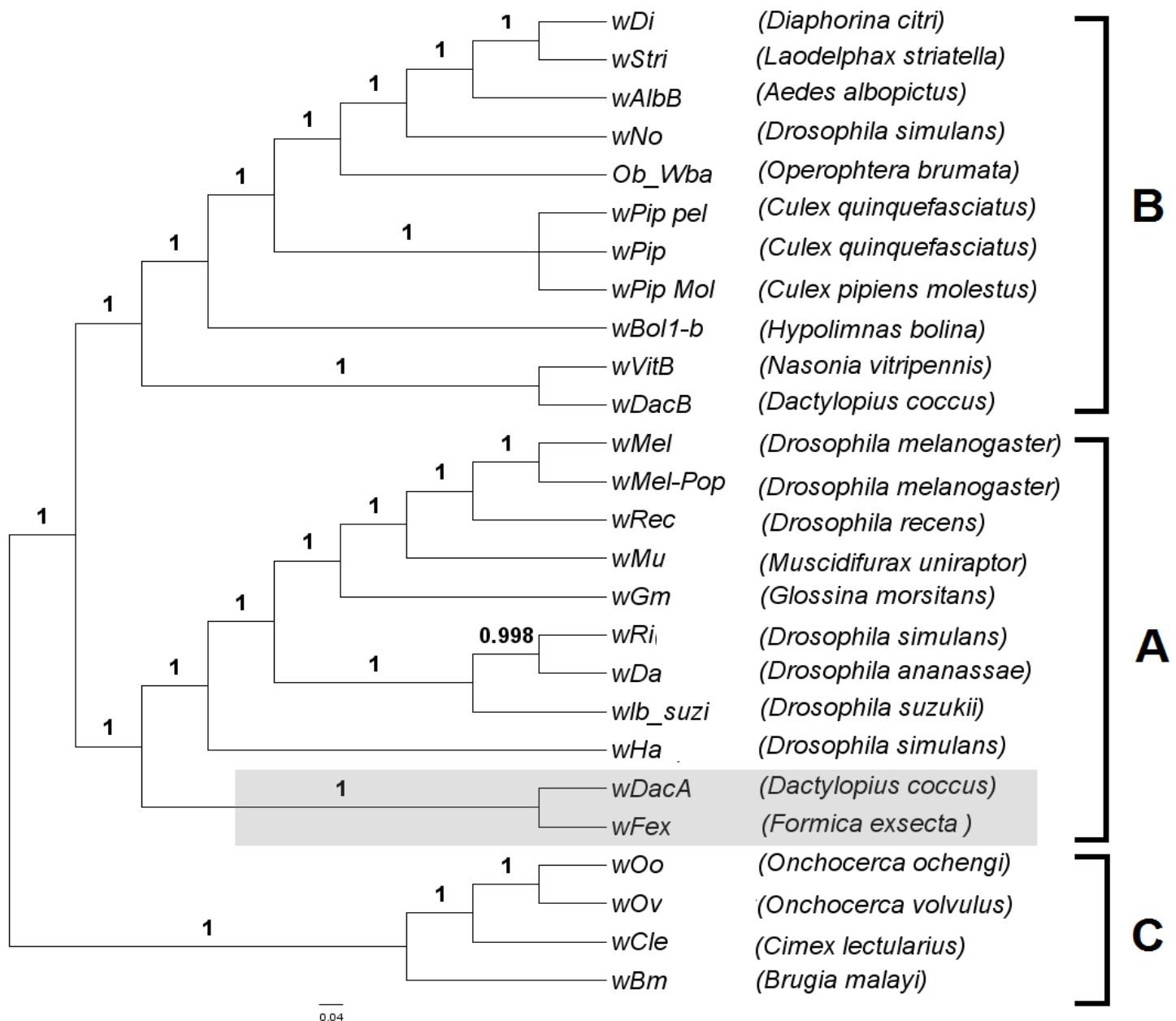
**Figure 5.** Venn diagram displaying overlap in orthologous gene among four

*Wolbachia* species including newly assembled wFex strain and previously reported
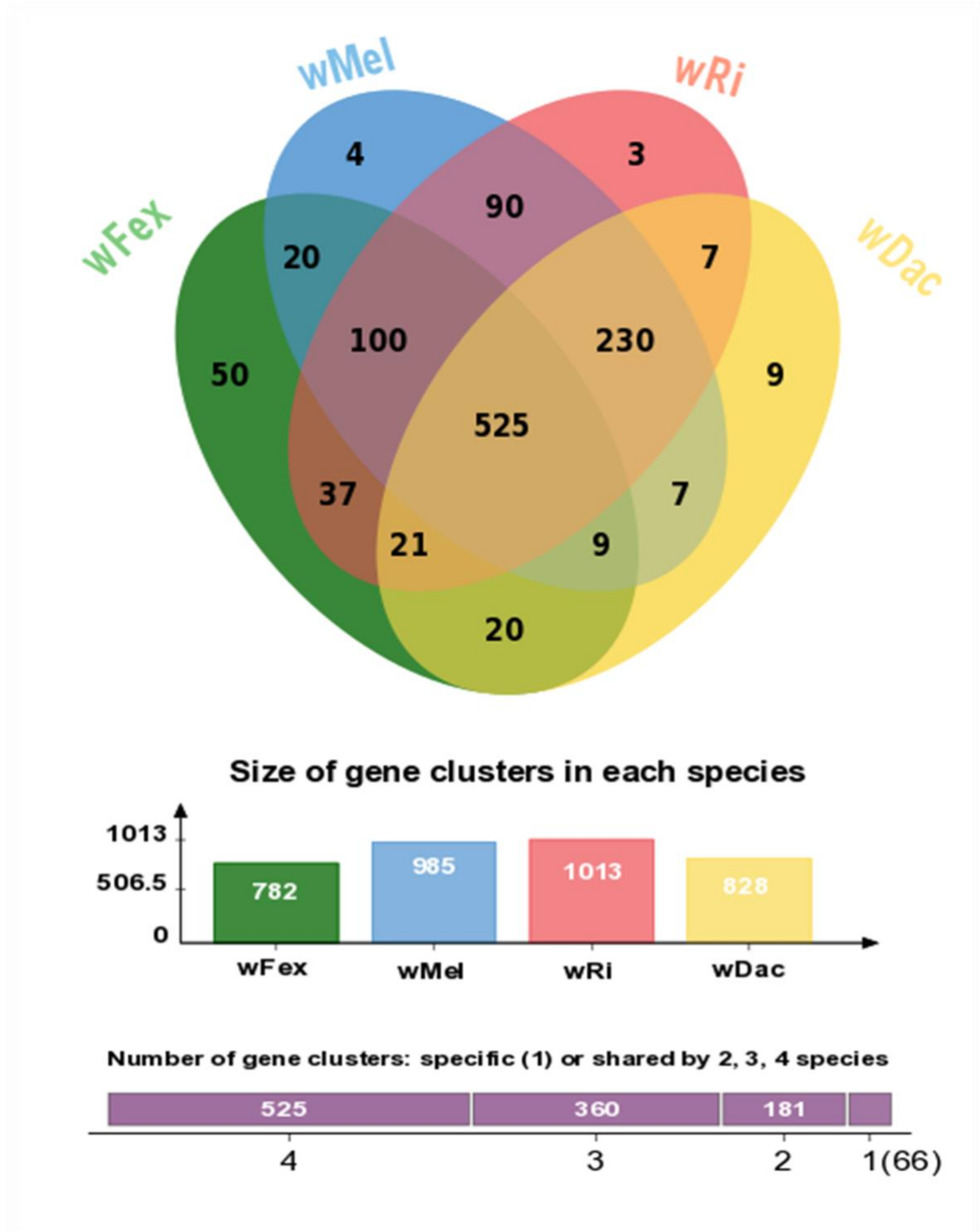
wDac, wRi, wMel strains.

**Table 1.** Summary statistics of raw sequencing data before and after filtering the reads which was further used for the genome assembly. "Coverage depth" was calculated based on the assembled genome size (500 Mb).

| Insert Size | Pair reads Length (bp) | Raw | | After Filter | |
|---|---|---|---|---|---|
| | | Total Data (G) | Sequence coverage (X) | Total Data (G) | Sequence coverage (X) |
| 170bp | 100 bp | 22.68 | 45.36 | 20.96 | 41.93 |
| 500bp | 100 bp | 8.54 | 17.08 | 7.34 | 14.69 |
| 800bp | 100 bp | 8.84 | 17.69 | 5.14 | 10.29 |
| 2kb | 100 bp | 13.23 | 26.46 | 7.05 | 14.10 |
| 5kb | 100 bp | 14.51 | 29.02 | 4.74 | 9.49 |
| 10kb | 100 bp | 11.77 | 23.53 | 5.51 | 11.02 |
| 20kb | 100 bp | 20.40 | 40.81 | 2.91 | 5.81 |
| Total | -- | 99.97 | 199.95 | 53.66 | 107.32 |

**Table 2.** Genome assembly statistics for *F. exsecta* and its *Wolbachia endosymbiont*.

| Genome Assembly Stats | *Formica exsecta* Genome | FE *Wolbachia* endosymbiont Genome |
|---|---|---|
| Total length | 277719392 (277 MB) | 3096460 (3.09 MB) |
| Total contigs | 14617 | 69 |
| Contigs (>= 1000 bp) | 3136 (98.24% genome) | 68(99.97% genome) |
| Contigs (>= 50000 bp) | 545 (89.59% genome) | 22(75.48% genome) |
| N50: | 997654 bp | 104167 bp |
| N75: | 318356 bp | 54296 bp |
| L50: | 73 | 11 |
| L75: | 185 | 22 |
| GC (%) | 36.00 | 35.13 |

**Table 3.** BUSCO quality metrics for the *F. exsecta* genome and the *Wolbachia* endosymbiont *of F. exsecta (wFex)* genome assembly.

| BUSCO metric | *F. exsecta* genome | | | | *FE Wolbachia endosymbiont Genome* (wFex) | |
| --- | --- | --- | --- | --- | --- | --- |
| | **Eukaryota** | **Insecta** | **Arthropoda** | **Hymenoptera** | **Bacteria** | **Proteobacteria** |
| Complete | 299 (98.7%) | 1634(98.5%) | 2549(95.29%) | 4249 (96.2%) | 107 (72.30%) | 158 (71.49%) |
| Complete and single copy | 283 (93.4%) | 1572 (94.8%) | 2446(91.44%) | 4151 (94.0%) | 35 (23.65%) | 55 (24.88) |
| Complete and duplicated | 16 (5.3%) | 62 (3.7%) | 103 (3.86%) | 98 (2.2%) | 72 (48.65%) | 103 (46.60%) |
| Fragmented | 1 (0.3%) | 15 (0.9%) | 195 (7.29%) | 123 (2.8%) | 9 (6.08%) | 11 (4.97%) |
| Missing | 3 (1.0%) | 9 (0.6%) | 34 (1.27%) | 43 (1.0%) | 32 (21.62%) | 52 (23.52%) |
| Total | 303 (100%) | 1658 (100%) | 2675 (100%) | 4415 (100%) | 148 (100%) | 221 (100%) |

**Table 4.** HGT inserts from *Wolbachia* present in the genome of *F. exsecta* with details of its length and position in the *F. exsecta* genome. The presence of similar insert regions in other eukaryote genomes is also shown.

| *Wolbachia* gene name | HGT region in *F. exsecta* | Length HGT (bp) | Transposon region near HGT | Transposon Name | Observed in other species | Other Host Species name with position of similar insertion |
|---|---|---|---|---|---|---|
| Transposase | scaffold83: 2271642-2272117 | 475 | scaffold83:2271642-2272117 | transposase | Complete | *Vollenhovia emeryi (LOC105557741), Cardiocondyla obscurior (genes: scf7180001101632 and scf7180001108526), Diachasma alloeum (LOC107035412),Brugia pahangi (BPAG_contig0001587),* |
| ABC transporter ATP-binding protein, porphobilinogen deaminase,D-alanine--D-alanine ligase, DNA processing protein DprA , triose-phosphate isomerase | scaffold233: 1712452-1725498 | 13046 | scaffold233:1714122-1714241 | transposase | Partial (few gene region) | *Vollenhovia emeryi(NW_011967015.1,NW_011967060.1),Wasmannia auropunctata (scf7180000683207,scf7180000730160),Rhagoletis zephyria(NW_016158779.1),Planococcus citri(KF021963.1) ,Ctenocephalides felis (KC177865.1)* |
| DNA repair protein RadC,transposase,DNA ligase,ABC transporter permease,ATP-dependent protease La | scaffold574: 102007-116197 | 14190 | scaffold574:105963-106483 | transposase | Partial (few gene region) | *Vollenhovia emeryi (LOC105557101,NW_011966940.1,NW_011966751.1 ), Monomorium pharaonis (scf7180001140281), Rhagoletis zephyria (LOC108377626),Parasteatoda tepidariorum(LOC107444616, LOC107450900)* |
| probable carboxypeptidase,type IV secretion system,conjugal transfer protein TrbL,lysyl-tRNA synthetase,UDP-N-acetylmuramoylalanine-D-glutamate ligase | scaffold707: 1-38814 | 38813 | scaffold707:35826-36154 | Mariner Mos1 transposase (Ant origin) | Partial (few gene region) | *Vollenhovia emeryi (NW_011966954.1,NW_011966496 ), Wasmannia auropunctata (scf7180000735528), Brugia pahangi (BPAG_contig0000608, BPAG_scaffold0000225)* |
| DNA methylase, Ankyrin repeat domain protein, regulatory protein RepA,site-specific recombinase, cytochrome b-like | scaffold741: 1-47265 | 47264 | scaffold741:54020-54482, scaffold741:52587-52910 | IS110 family transposase, Integrase | Partial (few gene region) | *Vollenhovia emeryi (LOC105557561,NW_011966954.1,NW_011967060.1,NW_011967015.1 ),Wasmannia auropunctata (LOC105460331 , scf7180000733651), Drosophila ananassae (WD_0580 gene)* |