1 Simultaneous Single-Cell Genome and Transcriptome Sequencing of Termite Hindgut Protists
2 Reveals Metabolic and Evolutionary Traits of Their Endosymbionts
3
4 Michael E. Stephens[1,2#], Jacquelynn Benjamino[1,3#], Joerg Graf[1], Daniel J. Gage[1*]
5
6 [1]Department of Molecular and Cell Biology, University of Connecticut, Storrs, CT USA
7 [2]Department of Entomology, Cornell University, Ithaca, NY USA
8 [3]The Jackson Laboratory for Genomic Medicine, Farmington, CT USA
9
10 #These two authors share first authorship. Both contributed equally to work described in the manuscript. Order was
11 determined by which did the majority of the manuscript writing.
12
13 * corresponding author
14
15

16 **Abstract**

17      Different protist species which colonize the hindguts of wood feeding *Reticulitermes*

18 termites are associated with endosymbiotic bacteria belonging to the genus *Endomicrobium*. In

19 this study, we focused on the endosymbionts of three protist species from *Reticulitermes*

20 *flavipes*, which included *Pyrsonympha vertens*, *Trichonympha agilis*, and *Dinenympha* species

21 II. Since these protist hosts represented members of difference taxa which colonize different

22 niches within the hindguts of their termite hosts, we investigated if these differences translated to

23 differential gene content and expression in their endosymbionts. Following assembly and

24 comparative genome and transcriptome analyses, we discovered that these endosymbionts

25 differed with respect to possible niche specific traits such carbon metabolism. Our analyses

26 supported that genes related to carbon metabolism were acquired by horizontal gene transfer

27 (HGT) from donor taxa which are present in termite's hindgut community. In addition, our

28 analyses supported that these endosymbionts have retained and expressed genes related to natural

29 transformation (competence) and recombination. Taken together, the presence of genes acquired

30 by HGT and a putative competence pathway supported that these endosymbionts are not cut-off

31 from gene flow and that competence may be a mechanism by which members of the

32 *Endomicrobium* can acquire new traits.

33

**Importance**

The composition and structure of wood, which contains cellulose, hemicellulose and lignin, prevents most organisms from using this common food source. Termites are a rare exception among animals, and they rely on a complex microbiome housed in their hindguts to use wood as a source of food. The lower termite *R. flavipes* houses a variety of protist and prokaryotes that are the key players the disassembly of lignocellulose. In this paper we describe the genomes and the gene expression profiles of five *Endomicrobium* endosymbionts living inside three different protist species from *R. flavipes*. Data from these genomes suggest that these *Endomicrobium* species have different mechanisms for using both carbon and nitrogen. In addition, they harbor genes that may be used to import free DNA from their environment. This process of DNA-uptake may contribute to the high levels of horizontal gene transfer often seen in the *Endomicrobium* species.

**Introduction**

Among the wood-feeding lower termites, symbiotic protists which reside in the hindgut are often colonized by endosymbionts (1–4). In *Reticulitermes* spp. termites both Oxymonadida (order) and Parabasilia (class) protists associate with endosymbiotic bacteria belonging to the genus *Endomicrobium* (phylum Elusimicrobia, class Endomicrobia) (2, 5–7). Members of *Endomicrobium* have been shown to comprise a significant portion of the core bacterial community in wood-feeding termites such as *R. flavipes* (8, 9). These endosymbiotic lineages are thought to have initiated their associations with hindgut protists approximately 70 - 40 million years ago (10) and arose from free-living relatives during multiple independent acquisition

56    events (11). Vertical passage from one protist cell to its progeny, has resulted in co-speciation as

57    inferred from congruent ribosomal RNA (rRNA) phylogenies (7, 10, 12, 13).

58         In addition to colonizing the cytoplasm of certain hindgut protist species, *Endomicrobium*

59    spp. are ectosymbionts of protists (14) can be free-living as well (15)(11, 16, 17). Because of

60    their distribution across these different niches, they provide an opportunity for studying bacterial

61    genome evolution across different association lifestyles: free-living, endosymbiotic, and

62    ectosymbiotic.

63         To determine differences two *Endomicrobium* species that are closely related but with

64    distinct lifestyles, a previous study compared genomes of a free-living *Endomicrobium*, *E.*

65    *proavitum* strain Rsa215 (16) isolated from *R. santonensis* (*R. flavipes*), and '*Candidatus*

66    Endomicrobium trichonymphae' strain Rs-D17 (3), an endosymbiont isolated from the

67    cytoplasm of a *Trichonympha* from the termite *R. speratus* (3, 18). The findings suggested that

68    the transition from the free-living state to an intracellular lifestyle involved genome reduction,

69    similar to that of endosymbionts of sap-feeding insects and many obligate intracellular

70    pathogens. However, the intracellular strain Rs-D17 also incorporated genes, possibly from other

71    termite gut inhabitants, by horizontal gene transfer (HGT) (18). For example, the genome of '*Ca.*

72    E. trichonymphae' Rs-D17 appeared to have acquired several pathways including those that

73    encode sugar and amino acid transporters and genes involved in amino acids biosynthesis (18).

74    These findings suggested that, unlike the endosymbionts of sap-feeding insects, *Endomicrobium*

75    species may not be completely cut-off from gene flow (18).

76         We expand upon these studies by presenting and comparing near-complete draft genomes

77    and transcriptomes of three different *Endomicrobium* organisms, that were assembled from

78    single protists cells of three different species that inhabit the hindgut of *R. flavipes*. One of these

79    protists species, *Pyrsonympha vertens* , lives attached to the oxic gut wall (19, 20), while the

80    other two *Trichonympha agilis* and *Dinenympha* species II, are found in the more anoxic hindgut

81    lumen. In addition, *P. vertens* and *D.* species II are both Oxymonads while *T. agilis* is a

82    Parabasalid.

83        The analyses indicate that these *Endomicrobium* have differences in gene content and

84    expression, related to carbon usage and metabolism. And as seen previously in '*Ca*. E.

85    trichonymphae' Rs-D17, they have likely acquired genes from putative donor taxa that are

86    commonly associated with termites. In addition, we describe data suggesting that these

87    *Endomicrobium* have retained competence genes which may allow them to import exogenous

88    DNA and that perhaps have contributed to HGT. Genes involved in this pathway are conserved

89    across several *Endomicrobium* species and were expressed in the endosymbionts examined in

90    this study.

91

92    **Methods**

93    **Termite collection and species identification**. *R. flavipes* termites were collected using

94    cardboard traps placed under logs for 2 to 4 weeks at the UConn Campus at Storrs, Connecticut

95    (Longitude -72.262216, Latitude 41.806543)  and their identity was verified as previously

96    described (7) by amplifying and sequencing the mitochondrial cytochrome oxidase II gene.

97    Termites were maintained in the lab with moistened sand and spruce wood that were initially

98    sterilized.

99

100   **Single protist cell isolation**. Termites from the worker caste were brought into an anaerobic

101   chamber and their hindguts were dissected with sterile forceps. Hindguts were ruptured in ice-

102   cold Trager's Solution U (TU) (21) and washed three times by centrifuging in 500 ul of TU at

103   3,000 rpm in an Eppendorf microcentrifuge for 90 seconds. This washed cell suspension was

104   then diluted 10-fold in TU buffer on ice. A 1 µl aliquot of the washed and diluted cell suspension

105   was added to a 9 µl droplet on a glass slide treated with RNase AWAY® Reagent (Life

106   Technologies) and UV light. Individual protist cells were isolated using a micromanipulator

107   (Eppendorf CellTram® Vario) equipped with a hand-drawn glass capillary. Individual cells were

108   washed three times in 10 µl droplets of TU via micromanipulation, transferring approximately

109   0.1 µl each time, and finally placed in 10µl molecular grade phosphate buffered saline (PBS),

110   flash frozen on dry ice, and immediately stored at -80°C. Meta-data regarding these protist cell

111   samples can be found as Supplementary Table 1 in S1 File.

112

113   **Whole genome and transcriptome amplification and sequencing.** The metagenome (DNA)

114   and metatranscriptome (cDNA) from individual protist cells and their associated bacteria were

115   simultaneously amplified 12-24 hours after isolation. Cell lysis and amplification was performed

116   using the Repli-g WGA/WTA kit (Qiagen). Cells were lysed using a Qiagen lysis buffer

117   followed immediately by incubation on ice. Two samples from each lysed cell were taken and

118   used in for whole genome amplification and whole transcriptome amplification. These were

119   carried using the manufacturer's standard protocol with exception that random hexamer primers

120   were used to amplify DNA and cDNA. DNA and cDNA were sheared using a Covaris M220

121   ultra-sonicator™ according to the manufacturer's protocol. WGA samples were sheared to a 550

122   bp insert size using 200 ng of DNA. WTA samples were sheared to a 350 bp insert size using

123   100 ng of cDNA. Sequencing libraries were prepared using the TruSeq Nano DNA Library Prep

124   kit from Illumina® according to the manufacturer's protocol. Each sample was prepared with a

125    forward and reverse barcode such that samples could be multiplexed on the same sequencing

126    run. The samples were sequenced using an Illumina® NextSeq 1x150 mid-output run and two

127    NextSeq 1x150 high-output runs. Meta-data regarding amplicon yields can be found in

128    Supplementary Table 1 in S1 File.

129

130    **Genomic read processing and assembly**. Reads were preprocessed before assembly using

131    BBmap (22). Reads were filtered for contaminating sequences by mapping reads to reference

132    genomes of potential contamination sources such as human DNA, human associated microbiota,

133    and organisms commonly used in our research laboratories. A list of references genomes used

134    for contamination filtering is provided in Supplementary Table 2 in S1 File.  Using BBmap

135    scripts, adaptor sequences were trimmed from reads and last base pair of 151 bp reads was

136    removed. Reads were then trimmed at both ends using a quality score cutoff of Q15.

137    Homopolymers were removed by setting an entropy cutoff of 0.2, a max G+C cutoff of 90%, and

138    by removing reads which possessed stretches of G's equal to or greater than 23 bases long. In

139    addition, reads which were below a minimum average quality of Q15 and/or 50 bases long were

140    removed. Genomic reads were then normalized to a minimum coverage of 2X and a maximum

141    coverage of 50X and then deduplicated using BBnorm. Genomic reads were assembled using the

142    A5 assembly pipeline (23) on the KBase web server (24). Meta-data regarding metagenome and

143    metatranscriptome reads numbers can be found in Supplementary Table 3 in S1 File.

144

145    **Genomic binning, draft genome assessment and annotation**. Metagenomic assemblies from

146    single protist host cells and their bacterial symbionts were binned using either 4mer or 6mer

147    frequencies with VizBin (25) and scaffolds at least 1Kb in size. Clustered scaffolds in genomic

148    bins of interest (low GC content) were selected in VizBin. Each scaffold from these bins were

149    used in a blastn (26) search against previously sequenced Elusimicrobia genomes

150    (Supplementary Table 4 in S1 File). Scaffolds which had a positive hit to other Elusimicrobia (at

151    least 70% identity over a 1kb alignment) were retained in the draft genomes and scaffolds which

152    did not have a significant hit to other Elusimicrobia genomes were used in a second blastn search

153    against the non-redundant (NR) database. Scaffolds which had positive hits to other

154    Elusimicrobia in the NR database were retained in the draft genomes. Draft genomes were

155    iteratively polished with the program Pilon (27). These draft genomes were then assessed for

156    contamination and completeness using CheckM which uses lineage specific marker genes to

157    perform analyses (28). The resulting near-complete draft genomes were then annotated on the

158    RAST Server  using a customized RASTtk workflow with options selected to call insertion

159    sequences and prophages (29, 30). Metabolic pathways pertaining to carbon metabolism, amino

160    acid biosynthesis, vitamin biosynthesis, and peptidoglycan biosynthesis were reconstructed from

161    the annotated genomes using pathways in the Kyoto Encyclopedia of Genes and Genomes

162    (KEGG) (31).

163

164    **Analysis of ribosomal gene phylogeny and average nucleotide identities.** Ribosomal 16S

165    genes from each of the *Endomicrobium* spp. draft genomes were trimmed and aligned to

166    references using MUSCLE (32), evolutionary models were tested and a Maximum likelihood

167    (ML) phylogenetic tree was made using IQ-TREE (33). JSpeciesWS (34) was used for

168    determining the genomic average nucleotide identities based on BLAST+ searches (ANIb)

169    between the *Endomicrobium* spp. draft genomes and the genome of '*Ca*. Endomicrobium

170    trichonymphae' Rs-D17, which is a close relative (3).

171    Assembled 18S rRNA genes were retrieved from metagenome assemblies by performing

172    a BLAST+ search using previously published 18S rRNA reference sequences for each protist

173    species as queries (7). When possible, protist 18S rRNA genes were amplified using leftover

174    DNA from WGA samples using universal primers 18SFU; 5'-

175    ATGCTTGTCTCAAAGGRYTAAGCCATGC-3' and 18SRU; 5'-

176    CWGGTTCACCWACGGAAACCTTGTTACG-3' (35) as previously described (7) and

177    sequenced by Sanger sequencing. This confirmation PCR was done on samples TA21, TA26,

178    and DS12. Assembled 18S rRNA genes were aligned to references using MUSCLE and a

179    Maximum likelihood (ML) phylogenetic tree was generated using IQ-TREE with model testing.

180

181    **Detection of horizontally acquired genes**. Genes that may have been acquired by horizontal

182    gene transfer were identified using phylogenetic methods. Initially, protein sequences of genes of

183    interest that were not shared across our draft genomes were aligned to references that spanned

184    eight different bacterial phyla (including Acidobacteria, Bacteroidetes, Nitrospiraceae,

185    Spirochaetes, Firmicutes, Actinobacteria, Proteobacteria, and group PVC) using MUSCLE and

186    phylogenetic trees were generated using IQ-TREE with model testing. Gene trees were then

187    compared to the 16S rRNA gene tree phylogeny (Supplemental Figure 3) to determine

188    evolutionary incongruence.

189

190    **Analysis of genes involved in competence and recombination**. Genes known to be involved in

191    DNA uptake, competence, and recombination were identified in each *Endomicrobium* spp. draft

192    genome based on their RAST annotations and homology to reference sequences. The distribution

193    of these genes was then compared across draft genomes and references which included free-

194     living relatives and other endosymbionts. To asses if these genes were complete and if the

195     encoded proteins likely retained their putative functions, homologs of each gene were obtained

196     from genomes of bacteria belonging to the phylum Elusimicrobia, aligned with MUSCLE, and

197     phylogenetic trees were generated using IQ-TREE (33) with model testing, and support values

198     generated using the "–abayes" and "–bb 1000" commands. The resulting phylogenetic trees were

199     used along with the MUSCLE alignments to perform a dN/dS analysis using the program

200     Codeml which is a part of the PAML and PAMLX packages (36, 37).

201

202     **Mapping transcriptome reads to draft genomes**. RNA-seq metatranscriptome reads were

203     quality trimmed and filtered as described above and error corrected in Geneious R11 (38) using

204     BBNorm with default settings. To remove rRNA reads before mapping, rRNA sequences were

205     identified from each metagenome assembly using RNAmmer (39) and reads were mapped to

206     these, as well as, rRNA references from refseq (40), SILVA (41) , and DictDb (42) databases

207     using BBmap (22). Remaining metatranscriptome reads were then mapped to their respective

208     *Endomicrobium* spp. draft genome in Geneious R11 using Bowtie 2 (43) with alignment type set

209     to "End to End" and using the "Medium Sensitivity" preset. Expression levels were then

210     calculated in Geneious R11, ambiguously mapped reads were excluded from the calculations.

211     RPKM values for genes in each genome are given in Supplementary Tables 5 and 6 in S1 File

212     **Verification of *comEC* expression by RT-PCR**. Primers were designed to amplify *comEC* from

213     '*Ca*. Endomicrobium agilae' in Geneious R11 using Primer3 (44) (Primers: endo_comec_F: 5'-

214     ATTTGCCTGTGTTTGAGAGT-3' and endo_comec_R: 5'-CCTGTTCCTGTGCTTTCAG-3').

215     Twenty termites were used to prepare RNA and cDNA samples for RT-PCR analysis. Termite

216     hindguts were dissected and ruptured in TU on ice in an anaerobic chamber. Hindgut contents

217     were washed with ice-cold TU three times at 3,000 rpm in an Eppendorf microcentrifuge for 90

218     seconds and then lysed in 1mL of TRIzol™ Reagent (Thermo Fisher Scientific). RNA was

219     isolated per the manufacturer's protocol and treated with TURBO™ DNase (Thermo Fisher

220     Scientific) following the manufacturer's protocol for 50μl reactions using ~5.6 mg of total RNA

221     in a 20 ul volume. 20 ul of the DNase-treated RNA was then used as template for cDNA

222     synthesis using SuperScript IV Reverse Transcriptase (Thermo Fisher Scientific) following the

223     manufacturer's protocol for first strand synthesis primed with random hexamers. The resulting

224     cDNA was treated with *E. coli* RNaseH for 20 minutes at 37°C.

225             RT-PCR reactions were performed using the *Endomicrobium comEC* primers with

226     RNaseH-treated cDNA serving as template and the no-RT control consisting of DNase-treated

227     RNA that did not undergo cDNA synthesis. RT-RCR was performed using Phusion Polymerase

228     (Thermo Fisher Scientific) with HF buffer and DMSO. Cycling conditions were: initial

229     denaturation at 94°C for 3 minutes, followed by 35 cycles of 94°C for 45 seconds, annealing at

230     59°C for 30 seconds, and extension at 72°C for 45 seconds. Final extension was done 72°C for 10

231     minutes. Hindgut DNA (washed protists cell fractions from five hindguts in molecular grade Tris

232     EDTA buffer) was used as a positive PCR control. RT-PCR products were visualized using a 1%

233     agarose gel with ethidium bromide. Products were purified using the Monarch DNA gel

234     purification kit (New England Biolabs), and Sanger sequenced.

235

236     **Data availability**. Raw reads and assemblies will be submitted to NCBI GenBank under

237     BioProject PRJNA644342.

238

239     **Results**

240     **Phylogeny of protist hosts**. Protist 18S rRNA genes were retrieved from metagenome

241     assemblies and confirmed (when possible) independently by PCR and Sanger sequencing. A

242     Maximum likelihood (ML) phylogenetic tree was made that indicated species of the protist cells

243     used in this study were *Trichonympha agilis* (cells TA21 and TA26), *Pyrsonympha vertens* (cells

244     PV1 and PV7), and *Dinenympha* species II (cell DS12) (Figure 1). These protist species have

245     been previously confirmed to live associated with *R. flavipes*, the termite species used in this

246     study.

247

248     ***Endomicrobium* genome statistics**. Five near-complete *Endomicrobium* genomes were obtained

249     from single protist cell metagenomic assemblies. The five genomes ranged from 1.12 - 1.37 Mb

250     in size, 35.3 – 36.6 % G+C, and 93.3 - 96.6% completeness (Figure 2A). To determine if these

251     genomes were from the same or different *Endomicrobium* species, we calculated pairwise

252     genomic s using an ANI score of 95% or greater as a marker for species-level cutoff (Figure 2B)

253     (45).

254        From *T. agilis* samples, we assembled two Endomicrobium draft genomes which had and

255     ANI score of greater than 96% to one another but less than 90% to '*Ca*. E. trichonymphae' Rs-

256     D17 indicating that they are likely different species. Based on this analysis, we refer to the draft

257     genomes as coming from '*Candidatus* Endomicrobium agilae' TA21 and '*Candidatus*

258     Endomicrobium agilae' TA26. We also assembled two *Endomicrobium* genomes from *P. vertens*

259     samples which had an ANI score greater than 97% identity to each another (Figure 2B) and

260     whose 16S rRNA genes were greater than 98% identical to a previously described species,

261     '*Candidatus* Endomicrobium pyrsonymphae' (6), which is the Candidatus species designation

262     that we use for PV1 and PV7. One additional *Endomicrobium* genome was assembled from

263   *Dinenympha* species II. This genome did not share an ANI score greater than 95% to other

264   *Endomicrobium* genomes and was thus given a new Candidatus species designation '*Candidatus*

265   Endomicrobium dinenymphae' DS12 (Figure 2B).

266        Individually, these *Endomicrobium* genomes contained between 1005 – 1230 orthologous

267   gene clusters (OGCs), of which 717 were found in all five genomes (Figure 2C). Additionally,

268   409 OCGs were unique to '*Ca*. E. agilae' TA21 and TA26 and another 183 OGCs were unique

269   to '*Ca*. E. pyrsonymphae' PV1 and PV7 (Figure 2C). Although the genome of '*Ca*. E.

270   dinenymphae' DS12 only had 24 unique OGCs, it shared 153 with '*Ca*. E. pyrsonymphae' PV1

271   and PV7 (Figure 2C) which may reflect similar selective pressures for gene retention in their

272   Oxymonad hosts *Dinenympha* and *Pyrsonympha*, or their more recent shared history, compared

273   with the *Endomicrobium* (TA21 and TA26) that associated with the Parabasalid *T. agilis*.

274

275   **Biosynthesis of amino acids, vitamins and peptidoglycan**. The presence of genes for the

276   various function discussed below were detected by tblastn using the queries listed in

277   Supplementary File 1 (Tables 8 and 9). When genes were not detected in this manner, read-

278   mapping using Geneious and read-mapping using Megan (63) were done as well to identify reads

279   from genes that may not have assembled into contigs. In general, each of the five

280   *Endomicrobium* genomes assembled in this study had similar gene content for processes

281   involved in the biosynthesis of amino acids (Supplementary Figure 1A), vitamins

282   (Supplementary Figure 1B), and peptidoglycan (Supplementary Figure 1D). Each genome

283   possessed complete pathways for alanine (from cysteine) aspartate, arginine, glutamine,

284   glutamate, glycine (from imported serine), isoleucine, leucine, valine, lysine, tyrosine,

285   phenylalanine and tryptophan biosynthesis (Supplementary Figure 1A).  Interestingly, the

286    *Endomicrobium* symbionts of Oxymonad protists (PV1, PV7 and DS12) lacked at least one gene

287    in the biosynthesis pathway for histidine (*hisG*) (Supplementary Figures 1A and 2). The

288    histidine biosynthetic pathway was complete in the genomes of '*Ca.* E. agilae' TA21 and TA26

289    (Supplementary Figure 1A). Conversely, it is likely that the *Endomicrobium* symbionts of

290    Oxymonad protists (PV1, PV7, and DS12) can make proline, while the symbionts represented by

291    genomes TA21, TA26 and RsD17 cannot (Supplementary Figure 1A). The five genomes

292    encoded incomplete pathways for the synthesis of cysteine and methionine. The three genomes

293    isolated from Oxymonad protists encoded a methionine transporter (MetT) and all contained a

294    gene encoding aB12-dependent methionine synthase system comprised of MetH and an

295    activation protein MetH2 (Supplementary Figure 1A). Also incomplete in all five genomes were

296    pathways for the synthesis of serine and asparagine. Each genome encoded a serine transporter

297    SdaC, a proline transporter (ProT) and PV1, PV7, and DS12 each encoded a glutamate

298    transporter, GltP (Supplementary Figure 1A).

299        The five *Endomicrobium* genomes also had similar gene content for processes involved

300    in the biosynthesis of vitamins and co-factors, with the pathways to pantothenate, CoA, NAD

301    and NADP being complete and other pathways being incomplete (Supplementary Figure 1B).

302    Interestingly, the biotin biosynthesis pathways in the five genomes are missing just a single gene

303    (*bioW*) needed to convert pimelate to pimelate-CoA suggesting that pimelate-CoA may be

304    synthesized by another enzyme or imported (Supplementary Figure 1B). Several genes in the

305    thiamine biosynthesis pathway were also missing in each of these genomes (Supplementary

306    Figure 1B). As noted previously for *E. proavitum* and "*Candidatus* Endomicrobium

307    trichonymphae" strain Rs-D17 the five genomes described here were also missing the steps in

308    the folate pathway needed to make 4-aminobenzoate, which may be transported into the cells

309    (18). The pathways for pyridoxine (B6) and vitamin B12 were also incomplete, though each of

310    the five *Endomicrobium* genomes appeared to encode ABC transport systems for vitamin B12

311    and heme.

312         Regarding peptidoglycan synthesis, each *Endomicrobium* genome was missing an

313    enzyme (BacA) that typically dephosphorylates undecaprenyl pyrophosphate. Since these

314    different *Endomicrobium* species, including the free-living *E. proavitum*, are missing the same

315    gene it may be that these bacteria utilize an alternate phosphatase to carry out the same function

316    as BacA.

317

318    **Differences in Carbon Metabolism**. Some of the more interesting differences between these

319    *Endomicrobium* genomes pertained to their carbon metabolisms. Each of these five

320    *Endomicrobium* genomes encoded relatively simple pathways for importing and using different

321    wood-derived carbon sources. Each had a complete phosphotransferase system (PTS) for

322    importing sugars. Present were two EIIA genes encoding sugar specific phosphorylation proteins

323    most closely related to those of the mannose and fructose type EIIA proteins (Supplementary

324    Figure 1C). Zheng et al. reported that *E. proavitum,* which contains a very similar PTS pathway,

325    did not grow on mannose or fructose, but did grow on glucose, suggesting that glucose may be

326    the carbohydrate transported by the PTS in that *Endomicrobium* species, and perhaps in the one

327    described here as well (16, 18).

328         Based on the gene content in the five genomes analyzed here, carbon sources capable of

329    being catabolized by endosymbiotic *Endomicrobium* species may often differ from each other

330    and from their free-living relatives. For example, '*Ca.* E. agilae' TA21 and TA26 encoded all the

331    genes necessary to import and use both glucuronate and glucose-6-phosphate (Figure 3A & 3B).

332      The closely-related '*Ca*. E. trichonymphae' Rs-D17 (3) also contained these genes. Interestingly,

333      genome analyses suggest that these two carbon sources cannot be used by the other

334      *Endomicrobium* species studied here which lack the glucuronate transporter ExuT, the

335      glucuronate isomerase UxaA and the glucose-6-phosphate transporter UhcP. The other

336      *Endomicrobium* genomes encoded either arabinose ('*Ca*. E. pyrsonymphae' PV1 and PV7) or

337      xylose ('*Ca*. E. dinenymphae' DS12) import and catabolism proteins that were not encoded in

338      the TA21, TA26, *E. proavitum* or '*Ca*. E. trichonymphae' Rs-D17 genomes. (Figures 4A & 4B;

339      Figures 5A & 5B, respectively). Transcriptome data indicated that each of the genes involved in

340      these carbon usage pathways were expressed in the respective *Endomicrobium* while they

341      resided in their protist hosts (Figures 3C, 4C, & 5C). Metabolites from these carbon sources are

342      typically fed into both the non-oxidative pentose phosphate pathway and glycolysis, both of

343      which are complete in the five genomes described here (Figures 3B, 4B, 5B, & Supplementary

344      Figure 1C).

345         Other likely differences in carbon metabolism of these *Endomicrobium* species, included

346      the production of fermentation end products (Supplementary Figure 1C). Analysis of the five

347      *Endomicrobium* genomes suggested that following glycolysis, pyruvate can be fermented to

348      acetate, however only the genomes of '*Ca.* E. agilae' and '*Ca*. E. dinenymphae', encoded AdhE

349      which can convert acetate to ethanol. (Figures 3B, 4B, & 5B). In addition, genes encoding lactate

350      dehyrodenase (LdH) were in the genomes of both '*Ca*. E. dinenymphae', and '*Ca*. E.

351      pyrsonymphae', but not '*Ca.* E. agilae' or '*Ca*. E. trichonymphae' Rs-D17 (Figures 3B, 4B, &

352      5B). Differences in these fermentation pathways between free-living *E. proavitum* and '*Ca*. E.

353      trichonymphae' Rs-D17 were described earlier by Zheng et al. (18).

354    Previous studies identified genes acquired by horizontal gene transfer (HGT) in other

355    *Endomicrobium* species (18), therefore, we tested whether HGT could, at least in part, explain

356    the differences seen in carbon metabolism across the genomes presented in this study.

357    Phylogenetic trees were made for each of the transport and isomerase proteins in the glucuronate,

358    arabinose and xylose degradation pathways (Figures 3D, 4D, & 5D) and the phylogenies

359    compared to the *Endomicrobium* 16S rRNA gene phylogeny to determine if they were congruent

360    (Supplementary Figure 3). In each case, these phylogenies were not congruent, suggesting that

361    these genes were acquired by HGT (Figures 3D, 4D, & 5D). Likely donor taxa include

362    Bacteroidetes, Actinobacteria, and Firmicutes (Figures 3D, 4D, & 5D), which are all part of the

363    hindgut community of *R. flavipes* (8). Similar data supporting HGT in '*Ca*. E. trichonymphae'

364    Rs-D17 have been reported and suggest that *Endomicrobium symbionts* are not cut off from gene

365    flow and HGT (18). This is in contrast to the older endosymbionts of sap-feeding insects, which

366    are traditionally thought to experience little to no gene flow, however recent analyses suggested

367    that HGT may occur more frequently than previously thought in these symbionts (46).

368

369    **Natural transformation and competence as a possible mechanism for acquiring genes**

370    Analyses of sequenced genomes of endosymbiotic *Endomicrobium* lineages indicate that

371    acquisition of genes by HGT is relatively common. Thus, their genomes could reveal insights

372    into the mechanisms by such genes were acquired. Interestingly, compared to other

373    endosymbionts, the *Endomicrobium* genomes were enriched in genes related to the uptake of

374    exogenous DNA and recombination (natural transformation/competence) (Supplementary

375    Figures 4A and 4C). Of special interest are the *Endomicrobium* genes *comEC*, *comEB*, *comF*,

376     *comM*, *ssb*, *drpA*, and *recA*  which are all involved in natural transformation in bacteria such as

377     *Vibrio cholerae* (47).

378         The dN/dS analyses of these genes supported the hypothesis that selection was acting to

379     maintain the amino acid sequences of their corresponding gene products (dN/dS < 1.0) with the

380     exception of *ssb* from TA21 (Figure 6A). In addition, transcriptome analysis indicated that these

381     genes were expressed (Figure 6B). Expression of *comEC*, which encodes a transporter that

382     imports single stranded DNA across the inner-membrane and into the cytoplasm of Gram-

383     negative bacteria (47, 48), was verified by RT-PCR and sequencing of *Ca*. Endomicrobium

384     agilae using *comEC* specific primers on a protist cell fraction sample prepared from 20 worker

385     termite hindguts (Figure 6C). Together these data support the hypothesis that genes involved in

386     this competence pathway are both conserved and expressed in these *Endomicrobium* symbionts

387     of hindgut protists of *R. flavipes*.

388         The competence genes discussed above are involved in the translocation of single-

389     stranded DNA across the inner membrane of Gram-negative bacteria and subsequent

390     recombination. Also present in the genomes of all five *Endomicrobium* species analyzed in this

391     study are genes which encode proteins that are similar to Type IV pilins. The TA21 and TA26

392     genomes contained a large chromosomal region devoted to Type IV Tad-like pilus synthesis as

393     does E. *proavitum* and *E. minutum*. The bacterium *Ca*. E. trichonymphae' Rs-D17 has a similar

394     region, but it appears that many of the genes have become pseudogenes. The *P. vertens* and

395     *Dinenympha* species II *Endomicrobium* symbionts had genes encoding pilins similar to the Type

396     II PulG pilins. Some pilins from classes Type-IV and Type-II can bind and import double-

397     stranded DNA across the outer membrane and have been shown to work in conjunction with

398     ComEC-type proteins (47, 49). In addition, all five genomes possessed a pre-pilin peptidase

399     (PilD). A graphical summary of these findings along with a model of how competence may work

400     in these organisms is provided as Supplementary Figure 6. A list of genes and their putative

401     function can be found in Supplementary Table 7 in S1 File.

402

403     **Transcriptome analysis of *Endomicrobium* populations inside single protist cells.**

404     Transcriptome analyses of the *Endomicrobium* populations inside single protist cells (from

405     which the five genomes were derived) revealed similar gene expression profiles with a few

406     notable exceptions (Figure 7). While our sample size for this work was necessarily small and

407     while we were unable to do time-resolved sampling of the hindgut community from single

408     termites, some general trends did appear in the transcriptomic data. Among COG categories,

409     which are quite broad, the expression by each endosymbiont population was relatively similar

410     with one exception being that there was higher expression of genes related to carbohydrate

411     transport and metabolism in *Trichonympha* hosts (TA) compared to the Oxymonad hosts (PV

412     and DS12) (Figure 7A).

413            Analyses which focused on narrower categories such genes in related biosynthetic

414     pathways, carbohydrate transport and break down, peptidoglycan synthesis and DNA uptake and

415     repair revealed further differences not only between the endosymbionts of different protist

416     species but even between the populations of endosymbionts of individual protist cells of the

417     same type (Figure 7B). This is demonstrated by the differences related to the expression of genes

418     in the glutamine and glutamate biosynthesis pathway (Figure 7B). Overall this pathway is more

419     highly expressed by the endosymbionts of *Pyrsonympha* hosts compared to those in other protist

420     species but there was variation in the expression of this pathway between individual protist host

421     cells. For example, in host cell PV1 this pathway represented 4.2% of the total transcriptome

422    reads, mostly from the gene *glnN* encoding a glutamine synthase, whereas in host cell PV7 it

423    comprised only 0.2% (Figure 7B). Similar variation in this pathway was seen in the TA21 and

424    TA26 transcriptomes. This data suggests that even populations of the bacterium residing in

425    different host cells may not be expressing the same functions at any given point in time.

426         Core genes, which were shared between all five of the *Endomicrobium* species,

427    represented an average of 30% – 36% of the transcription of each endosymbiont population

428    (Figure 7C). The expression of genes that were specific to each *Endomicrobium* species ranged

429    from 11% in 'Ca. E. pyrsonymphae' to 22% in 'Ca. E. agilae' indicating that there is likely

430    differential gene content and gene expression endosymbionts of different protist host species

431    (Figure 7C).  Transcriptomic data are in Supplementary Tables 5 and 6 in S1 File

432

433    **Discussion**

434         Single-cell protist metagenomics has enabled the assembly of genomes from several

435    protist-associated bacterial symbionts from termites (1, 3, 4, 50, 51). In this study, we present

436    near-complete draft genomes and transcriptomes of five endosymbiotic *Endomicrobium* samples,

437    from three different protist species. These endosymbionts displayed differences with regards to

438    their gene content and expression. For example, these organisms possessed different carbon

439    usage pathways. One hypothesis that may explain such differences in carbon utilization is that

440    different carbon sources may be provided to the *Endomicrobium* endosymbionts as by-products

441    of the protist hosts' hydrolysis and fermentation of the polysaccharides that originated in wood.

442    For example, glucuronate may be present in the cytoplasm of *Trichonympha* spp. because they

443    possessed the enzymes needed to cleave those monomers from polysaccharides found in wood,

444    whereas the other protists, *P. vertens* and *D.* species II, may not be able to generate such

445    monomers, or they may use them for other purposes (see below).  If true, this suggests that there

446    may be specialization among the protists with regards to polysaccharide hydrolysis in the

447    hindgut of wood-feeding termites. A recent study demonstrated a division of labor among

448    symbiotic protist species in a different termite, *Coptotermes formosanus*, where certain protist

449    species produce different hydrolytic enzymes to degrade polysaccharides found in wood (52).

450    These differences in protist functions may also explain why there was higher expression of genes

451    related to carbohydrate transport and metabolism in the endosymbionts of *Trichonympha* hosts

452    compared to the other protist species (Figure 7).

453        However, an alternative hypothesis is that metabolites can be partitioned within the host

454    and some are specifically provided to certain symbionts. This, if true, may allow the host to

455    control endosymbiont population densities through selective carbon source provision. Such host

456    control of carbon provisioning is thought to operate in nitrogen-fixing root nodule symbioses

457    ensuring that bacterial symbionts continue to provide fixed nitrogen in return for plant-provided

458    carbon. In support of the second hypothesis, the membrane-embedded symbiont '*Ca*.

459    Desulfovibrio trichonymphae' which co-colonizes the same *Trichonympha* host as '*Ca*. E.

460    trichonymphae' Rs-D17, uses malate and citrate as carbon sources whereas its co-inhabitant

461    likely uses glucuronate and glucose-6-phosphate (51).

462        Evidence suggests that *Endomicrobium* species have acquired genes by HGT and some of

463    the donor taxa may include termite-associated bacteria. Endosymbiotic lineages of

464    *Endomicrobium*, and their free-living relatives, possess many genes involved in DNA uptake,

465    repair, and recombination (Supplementary Figure 4). Our analyses showed that the genes

466    *comEC*, *comEB*, *comF*, *comM*, *ssb*, *drpA*, and *recA* are usually conserved within the

467    Elusimicrobia phylum and were expressed in the endosymbiotic *Endomicrobium* species

468    characterized in this study (Figure 6). Collectively these genes have been shown to be involved

469    in the translocation of single stranded DNA across the inner membrane of other Gram-negative

470    bacteria and in homologous recombination (47). The gene *comEC*, in particular, has an important

471    function in this process as it encodes an essential part of the DNA transporter (47, 48). Using

472    both transcriptome data, RT-PCR, and sequencing we were able to show that '*Ca*.

473    Endomicrobium agilae' *comEC* is expressed in these endosymbionts (Figure 6B and 6C).

474    Collectively the genes involved in the competence pathway comprised between 0.4% to 1.3% of

475    the total transcriptome reads of each endosymbiont population (Figure 7). These data suggest

476    that *Endomicrobium* species may have the ability to become competent which may allow them to

477    acquired DNA from the wider termite gut community and could result in HGT.

478            It is not clear how these organisms transport DNA across their outer membranes. None

479    of these *Endomicrobium* species possessed all the components of a Type IV pili-based DNA-

480    translocation system, but '*Ca*. E. agilae' TA21 and TA26 contained a near-complete Type IV *tad*

481    system which may allow DNA uptake (Supplemental Figures 4 and 5). It is also puzzling why

482    each of the five genomes have retained *pilD,* a prepilin peptidase as well as genes encoding Type

483    II and Type IV pilins. These may carry out some function in DNA-uptake or they may be non-

484    functional and are in the process of being lost. If competence is a common trait among the

485    *Endomicrobium*, this could explain why these organisms have many genes acquired through

486    HGT and this capability may allow for rapid adaptation to new and diverse niches. Because

487    hindgut protists phagocytize wood, wood-associated bacteria and perhaps free-living bacteria

488    (53), this may be a route through which endosymbiotic *Endomicrobium* could be exposed to

489    exogenous DNA.

490    However, it is worth noting that competence is not the only plausible avenue for DNA

491    acquisition in these endosymbionts. HGT could also occur by bacteriophage transduction,

492    conjugation, or other routes. Several lines of evidence indicate that these endosymbionts are

493    susceptible to molecular parasites, such as bacteriophages and plasmids. Previous studies have

494    reported that *Endomicrobium* species possessed several intact defense mechanisms to combat

495    molecular parasites such as CRISPR-Cas and restriction-modification systems (54, 55). The

496    *Endomicrobium* species sequenced in this study also contained those defense systems. The

497    complete genome sequence of a bacteriophage of an endosymbiont ('*Candidatus* Azobacteroides

498    pseudotrichonymphae') of a termite hindgut protist has previously been published, indicating

499    that phage infection is not limited to *Endomicrobium* endosymbionts and may be common in

500    termite hindguts (56).

501    Our analysis of *Endomicrobium* genomes and transcriptomes obtained from single protist

502    cell metagenomes, highlighted several important differences across protist hosts which have led

503    to hypotheses that warrant further investigation. In each case, the major hurdle of testing these is

504    the current inability to culture the protists hosts which restricts their experimental tractability.

505    However, the use of additional and new -omics approaches should further our understanding of

506    these symbioses by focusing on the hosts' and their symbionts' genes, mRNA, metabolic and

507    protein contents (57–62).

508

512

## References

1. Hongoh Y, Sharma VK, Prakash T, Noda S, Toh H, Taylor TD, Kudo T, Sakaki Y, Toyoda A, Hattori M, Ohkuma M. 2008. Genome of an endosymbiont coupling N2 fixation to cellulolysis within protist cells in termite gut. Science 322:1108–1109.

2. Ohkuma M. 2008. Symbioses of flagellates and prokaryotes in the gut of lower termites. Trends Microbiol 16:345–352.

3. Hongoh Y, Sharma VK, Prakash T, Noda S, Taylor TD, Kudo T, Sakaki Y, Toyoda A, Hattori M, Ohkuma M. 2008. Complete genome of the uncultured Termite Group 1 bacteria in a single host protist cell. PNAS 105:5555–5560.

4. Strassert JFH, Mikaelyan A, Woyke T, Brune A. 2016. Genome analysis of ' Candidatus Ancillula trichonymphae'', first representative of a deep-branching clade of Bifidobacteriales , strengthens evidence for convergent evolution in flagellate endosymbionts.' Environ Microbiol Rep 8:865–873.

5. Ohkuma M, Sato T, Noda S, Ui S, Kudo T, Hongoh Y. 2007. The candidate phylum "Termite Group 1" of bacteria: Phylogenetic diversity, distribution, and endosymbiont members of various gut flagellated protists. FEMS Microbiol Ecol 60:467–476.

6. Stingl U, Radek R, Yang H, Brune A. 2005. "Endomicrobia": Cytoplasmic symbionts of termite gut protozoa form a separate phylum of prokaryotes. Appl Environ Microbiol 71:1473–1479.

7. Stephens ME, Gage DJ. 2020. Single-cell amplicon sequencing reveals community structures and transmission trends of protist-associated bacteria in a termite host. PLoS One 15:1–19.

8. Benjamino J, Graf J. 2016. Characterization of the Core and Caste-Specific Microbiota in

536      the Termite, Reticulitermes flavipes. Front Microbiol 7:171.

537  9.    Boucias DG, Cai Y, Sun Y, Lietze VU, Sen R, Raychoudhury R, Scharf ME. 2013. The

538      hindgut lumen prokaryotic microbiota of the termite Reticulitermes flavipes and its

539      responses to dietary lignocellulose composition. Mol Ecol 22:1836–1853.

540  10.    Ikeda-Ohtsubo W, Brune A. 2009. Cospeciation of termite gut flagellates and their

541      bacterial endosymbionts: Trichonympha species and "Candidatus Endomicrobium

542      trichonymphae." Mol Ecol 18:332–342.

543  11.    Mikaelyan A, Thompson CL, Meuser K, Zheng H, Rani P, Plarre R, Brune A. 2017. High-

544      resolution phylogenetic analysis of Endomicrobia reveals multiple acquisitions of

545      endosymbiotic lineages by termite gut flagellates. Environ Microbiol Rep 9:477–483.

546  12.    Ikeda-Ohtsubo W, Desai M, Stingl U, Brune A. 2007. Phylogenetic diversity of

547      "Endomicrobia" and their specific affiliation with termite gut flagellates. Microbiology

548      153:3458–3465.

549  13.    Zheng H, Dietrich C, L. Thompson C, Meuser K, Brune A. 2015. Population Structure of

550      Endomicrobia in Single Host Cells of Termite Gut Flagellates (*Trichonympha* spp.).

551      Microbes Environ 30:92–98.

552  14.    Izawa K, Kuwahara H, Sugaya K, Lo N, Ohkuma M, Hongoh Y. 2017. Discovery of

553      ectosymbiotic Endomicrobium lineages associated with protists in the gut of stolotermitid

554      termites. Environ Microbiol Rep 9:411–418.

555  15.    Herlemann DPR, Geissinger O, Ikeda-Ohtsubo W, Kunin V, Sun H, Lapidus A,

556      Hugenholtz P, Brune A. 2009. Genomic analysis of "Elusimicrobium minutum" the first

557      cultivated representative of the phylum "Elusimicrobia" (formerly termite group 1). Appl

558      Environ Microbiol 75:2841–2849.

559   16.   Zheng H, Dietrich C, Radek R, Brune A. 2016. Endomicrobium proavitum , the first

560         isolate of Endomicrobia class . nov . ( phylum Elusimicrobia ) – an ultramicrobacterium

561         with an unusual cell cycle that fixes nitrogen with a Group IV nitrogenase. Environ

562         Microbiol 18:191–204.

563   17.   Ikeda-Ohtsubo W, Faivre N, Brune A. 2010. Putatively free-living 'Endomicrobia'-

564         ancestors of the intracellular symbionts of termite gut flagellates? Environ Microbiol Rep

565         2:554–9.

566   18.   Zheng H, Dietrich C, Brune A. 2017. Genome analysis of Endomicrobium proavitum

567         suggests loss and gain of relevant functions during the evolution of intracellular

568         symbionts. Appl Environ Microbiol 83:1–14.

569   19.   Brune A, Emerson D, Breznak JA. 1995. The termite gut microflora as an oxygen sink:

570         Microelectrode determination of oxygen and pH gradients in guts of lower and higher

571         termites. Appl Environ Microbiol 61:2681–2687.

572   20.   Yang H, Schmitt-Wagner D, Stingl U, Brune A. 2005. Niche heterogeneity determines

573         bacterial community structure in the termite gut (Reticulitermes santonensis). Environ

574         Microbiol 7:916–932.

575   21.   Trager W. 1934. The Cultivation of a Cellulose-Digesting Flagellate, Trichomonas

576         Termopsidis, and of Certain other Termite Protozoa. Biol Bull 66:182–190.

577   22.   Bushnell B. 2014. BBMap: A Fast, Accurate, Splice-Aware Aligner.

578   23.   Coil D, Jospin G, Darling AE. 2015. A5-miseq: An updated pipeline to assemble

579         microbial genomes from Illumina MiSeq data. Bioinformatics 31:587–589.

580   24.   Arkin AP, Cottingham RW, Henry CS, Harris NL, Stevens RL, Maslov S, Dehal P, Ware

581         D, Perez F, Canon S, Sneddon MW, Henderson ML, Riehl WJ, Murphy-Olson D, Chan

582    SY, Kamimura RT, Kumari S, Drake MM, Brettin TS, Glass EM, Chivian D, Gunter D,

583    Weston DJ, Allen BH, Baumohl J, Best AA, Bowen B, Brenner SE, Bun CC, Chandonia

584    J-M, Chia J-M, Colasanti R, Conrad N, Davis JJ, Davison BH, DeJongh M, Devoid S,

585    Dietrich E, Dubchak I, Edirisinghe JN, Fang G, Faria JP, Frybarger PM, Gerlach W,

586    Gerstein M, Greiner A, Gurtowski J, Haun HL, He F, Jain R, Joachimiak MP, Keegan KP,

587    Kondo S, Kumar V, Land ML, Meyer F, Mills M, Novichkov PS, Oh T, Olsen GJ, Olson

588    R, Parrello B, Pasternak S, Pearson E, Poon SS, Price GA, Ramakrishnan S, Ranjan P,

589    Ronald PC, Schatz MC, Seaver SMD, Shukla M, Sutormin RA, Syed MH, Thomason J,

590    Tintle NL, Wang D, Xia F, Yoo H, Yoo S, Yu D. 2018. KBase: The United States

591    Department of Energy Systems Biology Knowledgebase. Nat Biotechnol 36:566–569.

592  25.  Laczny CC, Sternal T, Plugaru V, Gawron P, Atashpendar A, Margossian HH, Coronado

593    S, der Maaten L V., Vlassis N, Wilmes P. 2015. VizBin - An application for reference-

594    independent visualization and human-augmented binning of metagenomic data.

595    Microbiome 3:1–7.

596  26.  Camacho C, Coulouris G, Avagyan V, Ma N, Papadopoulos J, Bealer K, Madden TL.

597    2009. BLAST+: architecture and applications. BMC Bioinformatics 10:421.

598  27.  Walker BJ, Abeel T, Shea T, Priest M, Abouelliel A, Sakthikumar S, Cuomo CA, Zeng Q,

599    Wortman J, Young SK, Earl AM. 2014. Pilon: An integrated tool for comprehensive

600    microbial variant detection and genome assembly improvement. PLoS One 9.

601  28.  Parks DH, Imelfort M, Skennerton CT, Hugenholtz P, Tyson GW. 2015. CheckM:

602    assessing the quality of microbial genomes recovered from. Cold Spring Harb Lab Press

603    Method 1:1–31.

604  29.  Aziz RK, Bartels D, Best A, DeJongh M, Disz T, Edwards RA, Formsma K, Gerdes S,

605    Glass EM, Kubal M, Meyer F, Olsen GJ, Olson R, Osterman AL, Overbeek RA, McNeil

606    LK, Paarmann D, Paczian T, Parrello B, Pusch GD, Reich C, Stevens R, Vassieva O,

607    Vonstein V, Wilke A, Zagnitko O. 2008. The RAST Server: Rapid annotations using

608    subsystems technology. BMC Genomics 9:1–15.

609    30.    Brettin T, Davis JJ, Disz T, Edwards RA, Gerdes S, Olsen GJ, Olson R, Overbeek R,

610    Parrello B, Pusch GD, Shukla M, Thomason JA, Stevens R, Vonstein V, Wattam AR, Xia

611    F. 2015. RASTtk: A modular and extensible implementation of the RAST algorithm for

612    building custom annotation pipelines and annotating batches of genomes. Sci Rep 5.

613    31.    Kanehisa M, Goto S. 2000. KEGG: Kyoto Encyclopedia of Genes and Genomes. Nucleic

614    Acids Res 28:27–30.

615    32.    Edgar RC. 2004. MUSCLE: Multiple sequence alignment with high accuracy and high

616    throughput. Nucleic Acids Res 32:1792–1797.

617    33.    Nguyen LT, Schmidt HA, Von Haeseler A, Minh BQ. 2015. IQ-TREE: A fast and

618    effective stochastic algorithm for estimating maximum-likelihood phylogenies. Mol Biol

619    Evol 32:268–274.

620    34.    Richter M, Rosselló-Móra R, Oliver Glöckner F, Peplies J. 2016. JSpeciesWS: a web

621    server for prokaryotic species circumscription based on pairwise genome comparison.

622    Bioinformatics 32:929–931.

623    35.    Tikhonenkov D V., Janouškovec J, Keeling PJ, Mylnikov AP. 2016. The Morphology,

624    Ultrastructure and SSU rRNA Gene Sequence of a New Freshwater Flagellate, Neobodo

625    borokensis n. sp. (Kinetoplastea, Excavata). J Eukaryot Microbiol 63:220–232.

626    36.    Yang Z. 1997. PAML: a program package for phylogenetic analysis by maximum

627    likelihood. Bioinformatics 13:555–556.

628    37.    Xu B, Yang Z. 2013. PamlX: A graphical user interface for PAML. Mol Biol Evol

629            30:2723–2724.

630    38.    Kearse M, Moir R, Wilson A, Stones-Havas S, Cheung M, Sturrock S, Buxton S, Cooper

631            A, Markowitz S, Duran C, Thierer T, Ashton B, Meintjes P, Drummond A. 2012.

632            Geneious Basic: An integrated and extendable desktop software platform for the

633            organization and analysis of sequence data. Bioinformatics 28:1647–1649.

634    39.    Lagesen K, Hallin P, Rødland EA, Stærfeldt HH, Rognes T, Ussery DW. 2007.

635            RNAmmer: Consistent and rapid annotation of ribosomal RNA genes. Nucleic Acids Res

636            35:3100–3108.

637    40.    Pruitt KD, Tatusova T, Maglott DR. 2007. NCBI reference sequences (RefSeq): A curated

638            non-redundant sequence database of genomes, transcripts and proteins. Nucleic Acids Res

639            35:61–65.

640    41.    Quast C, Pruesse E, Yilmaz P, Gerken J, Schweer T, Yarza P, Peplies J, Glöckner FO.

641            2013. The SILVA ribosomal RNA gene database project: Improved data processing and

642            web-based tools. Nucleic Acids Res 41:590–596.

643    42.    Mikaelyan A, Köhler T, Lampert N, Rohland J, Boga H, Meuser K, Brune A. 2015.

644            Classifying the bacterial gut microbiota of termites and cockroaches: a curated

645            phylogenetic reference database (DictDb). Syst Appl Microbiol 38:472–482.

646    43.    Langmead B, Salzberg SL. 2012. Fast gapped-read alignment with Bowtie 2. Nat Methods

647            9:357–359.

648    44.    Untergasser A, Cutcutache I, Koressaar T, Ye J, Faircloth BC, Remm M, Rozen SG. 2012.

649            Primer3-new capabilities and interfaces. Nucleic Acids Res 40:1–12.

650    45.    Jain C, Rodriguez-R LM, Phillippy AM, Konstantinidis KT, Aluru S. 2018. High

651         throughput ANI analysis of 90K prokaryotic genomes reveals clear species boundaries.

652         Nat Commun 9:1–8.

653    46.   López-Madrigal S, Gil R. 2017. Et tu, brute? Not even intracellular mutualistic symbionts

654         escape horizontal gene transfer. Genes (Basel) 8:1–16.

655    47.   Seitz P, Blokesch M. 2013. DNA-uptake machinery of naturally competent Vibrio

656         cholerae. Proc Natl Acad Sci 110:17987–17992.

657    48.   Pimentel ZT, Zhang Y. 2018. Evolution of the natural transformation protein, ComEC, in

658         Bacteria. Front Microbiol 9:1–10.

659    49.   Ellison CK, Dalia TN, Vidal Ceballos A, Wang JCY, Biais N, Brun Y V., Dalia AB.

660         2018. Retraction of DNA-bound type IV competence pili initiates DNA uptake during

661         natural transformation in Vibrio cholerae. Nat Microbiol 3:773–780.

662    50.   Utami YD, Kuwahara H, Igai K, Murakami T, Sugaya K, Morikawa T, Nagura Y, Yuki

663         M, Deevong P, Inoue T, Kihara K, Lo N, Yamada A, Ohkuma M, Hongoh Y. 2019.

664         Genome analyses of uncultured TG2/ZB3 bacteria in 'Margulisbacteria' specifically

665         attached to ectosymbiotic spirochetes of protists in the termite gut. ISME J 13:455–467.

666    51.   Sato T, Hongoh Y, Noda S, Hattori S, Ui S, Ohkuma M. 2009. Candidatus Desulfovibrio

667         trichonymphae, a novel intracellular symbiont of the flagellate Trichonympha agilis in

668         termite gut. Environ Microbiol 11:1007–15.

669    52.   Nishimura Y, Otagiri M, Yuki M, Shimizu M, Inoue J ichi, Moriya S, Ohkuma M. 2020.

670         Division of functional roles for termite gut protists revealed by single-cell transcriptomes.

671         ISME J.

672    53.   Brune A. 2014. Symbiotic digestion of lignocellulose in termite guts. Nat Rev Microbiol

673         12:168–80.

674    54.    Izawa K, Kuwahara H, Kihara K, Yuki M, Lo N, Itoh T, Ohkuma M, Hongoh Y. 2016.

675           Comparison of intracellular "ca. endomicrobium trichonymphae" genomovars illuminates

676           the requirement and decay of defense systems against foreign DNA. Genome Biol Evol

677           8:3099–3107.

678    55.    Zheng H, Dietrich C, Hongoh Y, Brune A. 2016. Restriction-modification systems as

679           mobile genetic elements in the evolution of an intracellular symbiont. Mol Biol Evol

680           33:721–725.

681    56.    Pramono AK, Kuwahara H, Itoh T, Toyoda A, Yamada A, Hongoh Y. 2017. Discovery

682           and Complete Genome Sequence of a Bacteriophage from an Obligate Intracellular

683           Symbiont of a Cellulolytic Protist in the Termite Gut. Microbes Environ Environ 32:112–

684           117.

685    57.    Bhattacharya D, Yoon HS, Price DC, Stepanauskas R, Rajah VD, Sieracki ME, Wilson

686           WH, Yang EC, Duffy S. 2011. Single-Cell Genomics Reveals Organismal Interactions in

687           Uncultivated Marine Protists. Science (80- ) 332:714–717.

688    58.    Karnkowska A, Vacek V, Zubáčová Z, Treitli SC, Petrželková R, Eme L, Novák L,

689           Žárský V, Barlow LD, Herman EK, Soukal P, Hroudová M, Doležal P, Stairs CW, Roger

690           AJ, Eliáš M, Dacks JB, Vlček Č, Hampl V. 2016. A eukaryote without a mitochondrial

691           organelle. Curr Biol 26:1274–1284.

692    59.    Mangot JF, Logares R, Sánchez P, Latorre F, Seeleuthner Y, Mondy S, Sieracki ME,

693           Jaillon O, Wincker P, Vargas C De, Massana R. 2017. Accessing the genomic information

694           of unculturable oceanic picoeukaryotes by combining multiple single cells. Sci Rep 7:1–

695           12.

696    60.    Vacek V, Novák LVF, Treitli SC, Táborský P, Cepicka I, Kolísko M, Keeling PJ, Hampl

697          V. 2018. Fe-S Cluster Assembly in Oxymonads and Related Protists. Mol Biol Evol

698          35:2712–2718.

699   61.    Kolisko M, Boscaro V, Burki F, Lynn DH, Keeling PJ. 2014. Single-cell transcriptomics

700          for microbial eukaryotes. Curr Biol 24:R1081–R1082.

701   62.    Hamann E, Tegetmeyer HE, Riedel Di, Littmann S, Ahmerkamp S, Chen J, Hach PF,

702          Strous M. 2017. Syntrophic linkage between predatory Carpediemonas and specific

703          prokaryotic populations. ISME J 11:1205–1217.

704   63.     Huson DH, Beier S, Flade I, Górska A, El-Hadidi M, et al. (2016) MEGAN Community

705   Edition - Interactive Exploration and Analysis of Large-Scale Microbiome Sequencing Data.

706   PLOS Computational Biology 12(6): e1004957. https://doi.org/10.1371/journal.pcbi.1004957

707

708

709

710     **Figures and Legends**



712     **Figure 1. Protist 18S rRNA gene phylogeny.** 18S rRNA genes were retrieved from single
713     protist cell metagenome assemblies, aligned to references, and a Maximum likelihood (ML)
714     phylogenetic tree was made using IQ-Tree using substitution model TIM2+G4. All 18S rRNA
715     gene sequences obtained in this study (denoted by *) are shown grouped with their respective
716     references. Branch support values represent the Bayesian posterior probability and Bootstrap
717     support values respectively.
718

719

**Figure 2. *Endomicrobium* draft genomes statistics, speciation, and shared gene content.** (A)
16S rRNA gene Maximum likelihood tree (unrooted) of the three *Endomicrobium* species,
genome sizes, percent G+C content, and estimated percent genome completeness. (B) Pairwise
genomic ANI scores of *Endomicrobium* genomes obtained by this study and a previously
sequenced relative Rs-D17. (C) UpSet graph of the number of orthologous gene clusters (OGCs)
of protein coding sequences within and across each of the *Endomicrobium* draft genomes.

726

**Figure 3. Carbon metabolism and HGT in 'Ca. Endomicrobium agilae'.** (A) Gene neighborhood of the genes involved in the metabolism of glucuronate in the 'Ca. Endomicrobium agilae' TA21 and TA26 genomes. (B) Diagram of a protist host and an *Endomicrobium* cell showing the inferred metabolic conversions of carbon sources based on gene content data. (C) Gene expression data of genes of interest (rows) pertaining to carbon metabolism in 'Ca. Endomicrobium agilae' TA21 and TA26 (columns). (D) Maximum likelihood phylogenetic trees of amino acid sequences of the transporter (ExuT, using substitution model LG+F+G4) and isomerase (UxaC, using substitution model LG+I+G4) in the glucuronate metabolism pathway. Support values represent the Bayesian posterior probability and Bootstrap support values respectively.
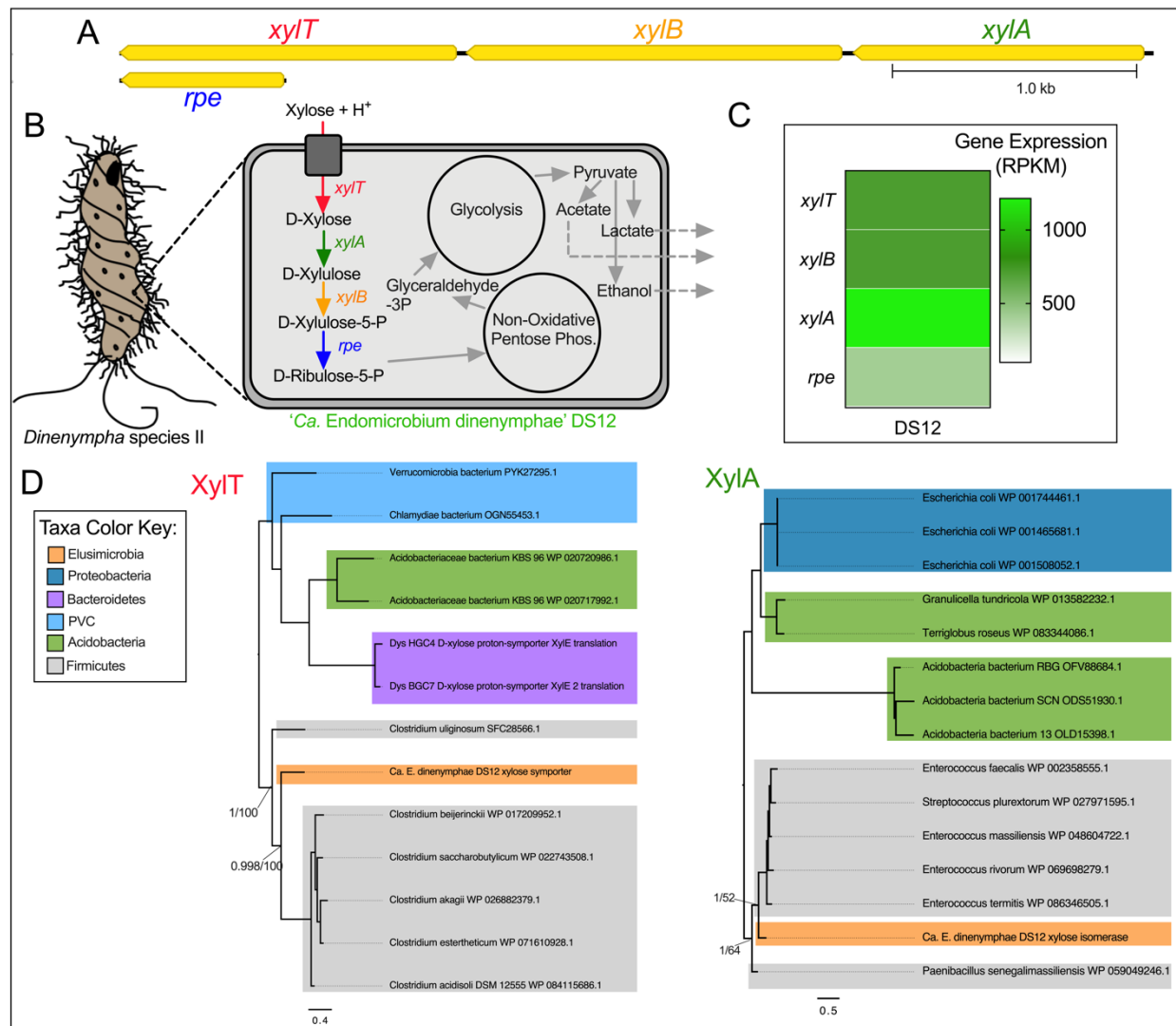
739



740

**Figure 4. Carbon metabolism and HGT in '*Ca*. Endomicrobium pyrsonymphae'.** (A) Gene neighborhood of the genes involved in the metabolism of arabinose in the '*Ca*. Endomicrobium pyrsonymphae' PV1 and PV7 genomes. (B) Diagram of a protist host and an *Endomicrobium* cell showing the inferred metabolic conversions of carbon sources based on gene content data. (C) Gene expression data of genes of interest (rows) pertaining to carbon metabolism in '*Ca*. Endomicrobium pyrsonymphae' PV1 and PV7 (columns). (D) Maximum likelihood phylogenetic trees of amino acid sequences from the transporter (AraE, using substitution model LG+F+G4) and isomerase (AraA, using substitution model LG+I+G4) in the arabinose metabolism pathway. Support values represent the Bayesian posterior probability and Bootstrap support values respectively.
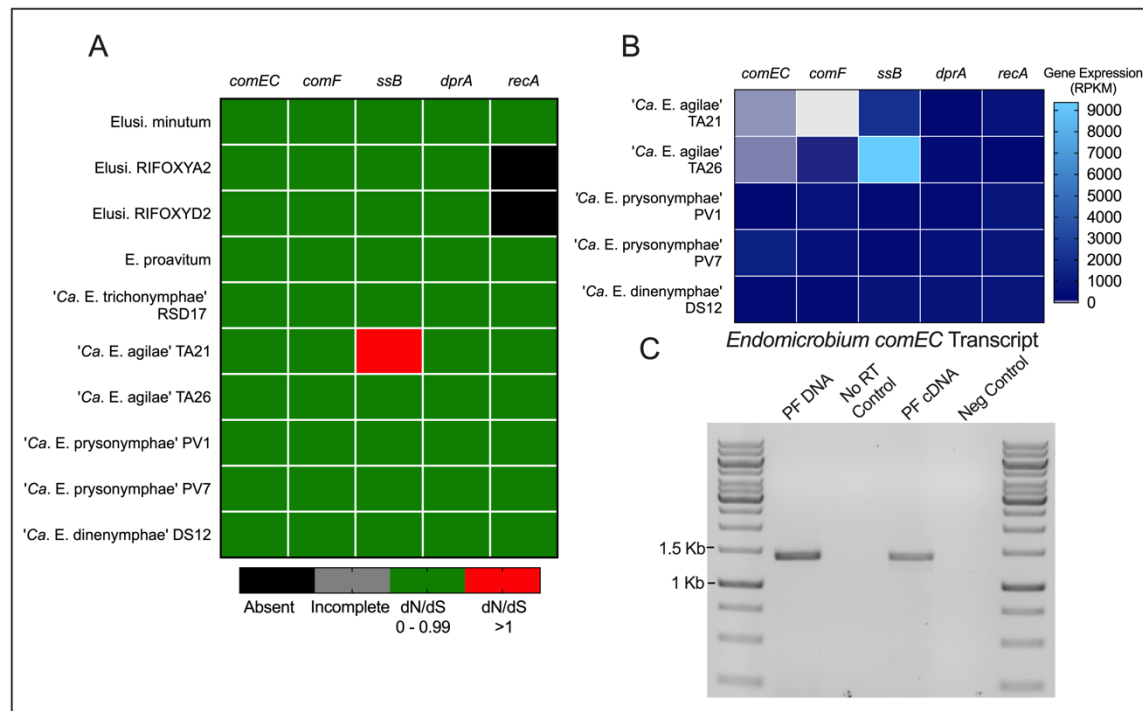
751

752

**Figure 5. Carbon metabolism and HGT in '*Ca*. Endomicrobium dinenymphae'.** (A) Gene neighborhood of the genes involved in the metabolism of xylose in the '*Ca*. Endomicrobium dinenymphae' DS12 genome. (B) Diagram of a protist host and an *Endomicrobium* cell showing the inferred metabolic conversions of carbon sources based on gene content data. (C) Gene expression data of genes of interest (rows) pertaining to carbon metabolism in '*Ca*. Endomicrobium dinenymphae' DS12 (column). (D) Maximum likelihood phylogenetic trees of amino acid sequences from the transporter (XylT, using substitution model LG+F+G4) and isomerase (XylA, using substitution model LG+G4) in the xylose metabolism pathway. Support values represent the Bayesian posterior probability and Bootstrap support values respectively.

762

763



764

**Figure 6. Analysis of genes involved in a putative competence pathway in *Endomicrobium* spp.** (A) Heatmap showing the results of dN/dS analyses of genes involved in competence and recombination (columns) from *Endomicrobium* spp. and *Elusimicrobium* relatives (rows). (B) Gene expression data of those genes (columns) in the *Endomicrobium* spp. (rows) presented in this study. (C) RT-PCR gel image of *Endomicrobium comEC* transcript. Samples consisted of protist fraction (PF) DNA (positive control), No RT Control, PF cDNA, and molecular grade water (Negative control). Accession numbers for reference genomes used can be found in Supplementary Table 2.

773

774



Percent of total, genomic, RPKM values

**Figure 7. Transcriptome analysis of *Endomicrobium* populations from individual protist cells.** (Top) Expression analysis of genes grouped into clusters of orthologous groups (COG) functional categories. (Middle) Expression analysis of genes in certain metabolic pathways pertaining to amino acid and cofactor biosynthesis, DNA processing, and peptidoglycan biosynthesis. (Bottom) Expression of genes grouped by their distributions among the different *Endomicrobium* species. RPKM values for individual cells, and averages of genera and all cells are in the various columns.
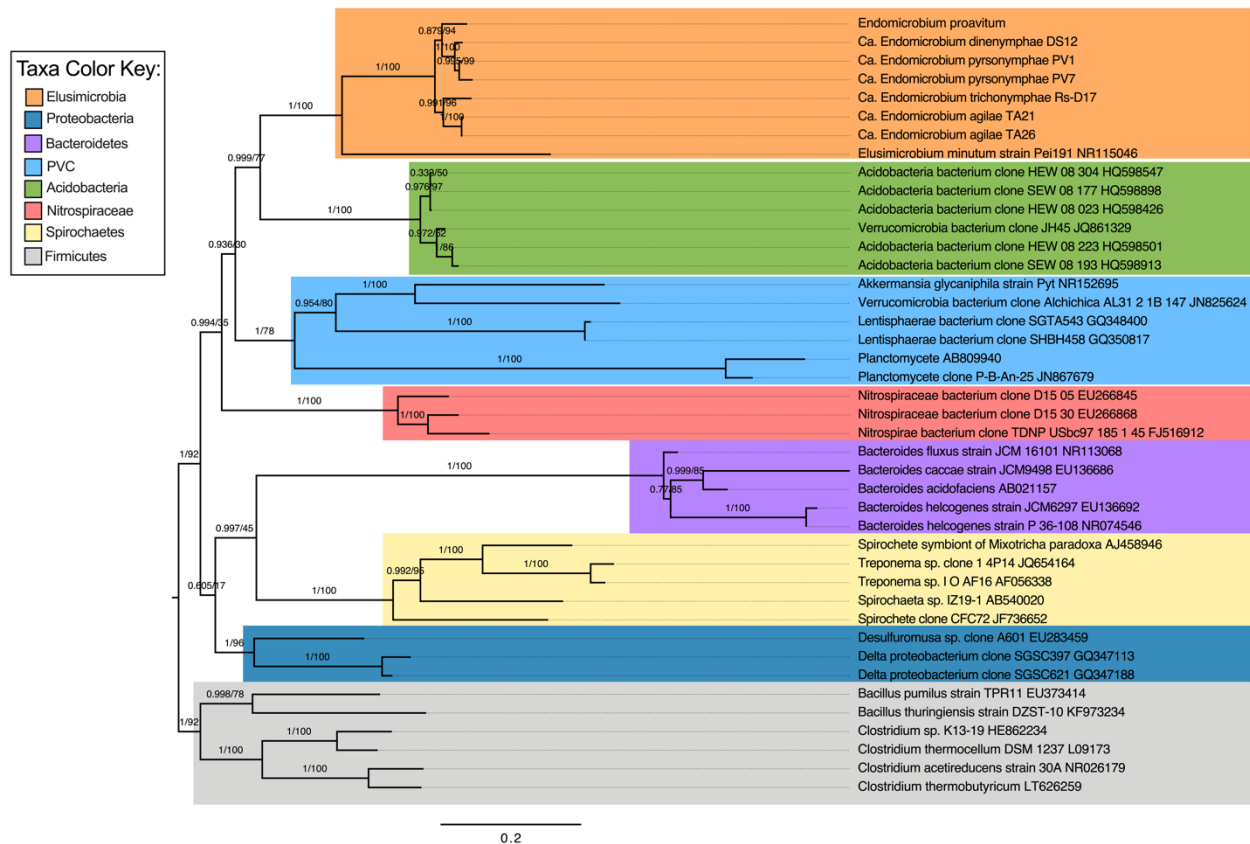
783 **Supplementary Figures and Legends**



784 **Supplementary Figure 1. Gene content in *Endomicrobium* spp. genomes regarding**
785 **metabolic functions and cell wall biosynthesis.** (A) Gene content of amino acid biosynthesis
786 pathways. (B) Gene content of vitamins and co-factor biosynthesis pathways. (C) Genes
787 involved in central metabolism. (D) Gene content of peptidoglycan biosynthesis. Note that the
788 genes marked with a "*" or a "!" were not in the final assemblies, but their reads were detected
789 by either Megan (!) or by read mapping onto the same gene from another closely related
790 organisms (*) or both (!*) (see Supplementary Figure 2 as an example).
791

792



793

794 **Supplementary Figure 2. Mapping coverage of the histidine biosynthesis pathway in '*Ca.***
795 **Endomicrobium pyrsonymphae' PV1.** Metagenomic reads from sample PV1 were mapped to
796 the draft genome of '*Ca*. Endomicrobium pyrsonymphae PV7'. The resulting coverage of these
797 genes indicate that '*Ca*. Endomicrobium pyrsonymphae PV1' encoded the histidine biosynthesis
798 pathway and that the reason those genes are missing from the draft genome is likely due to an
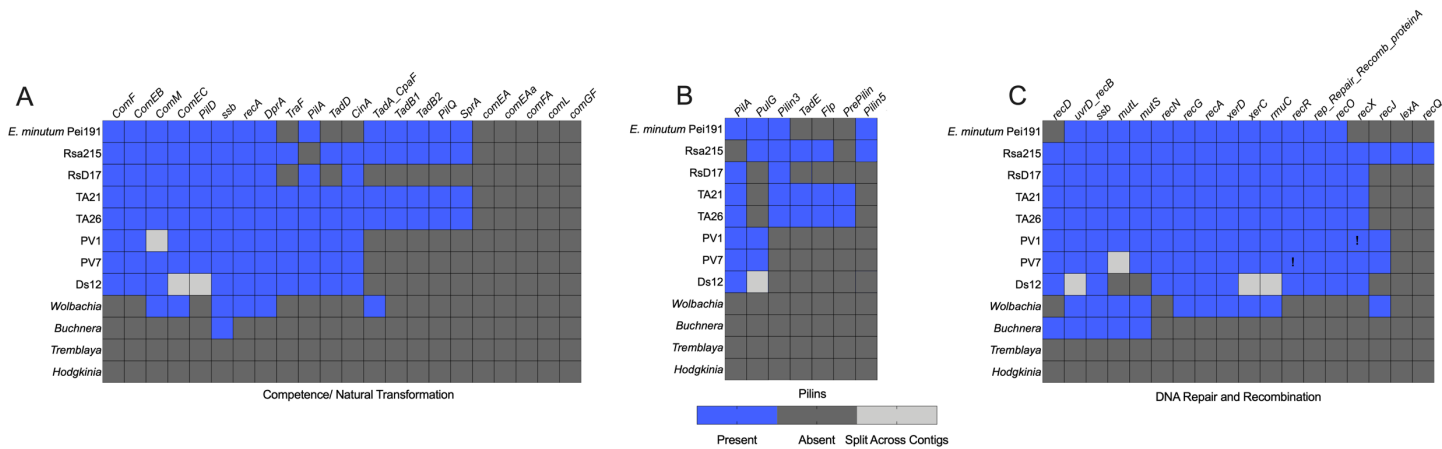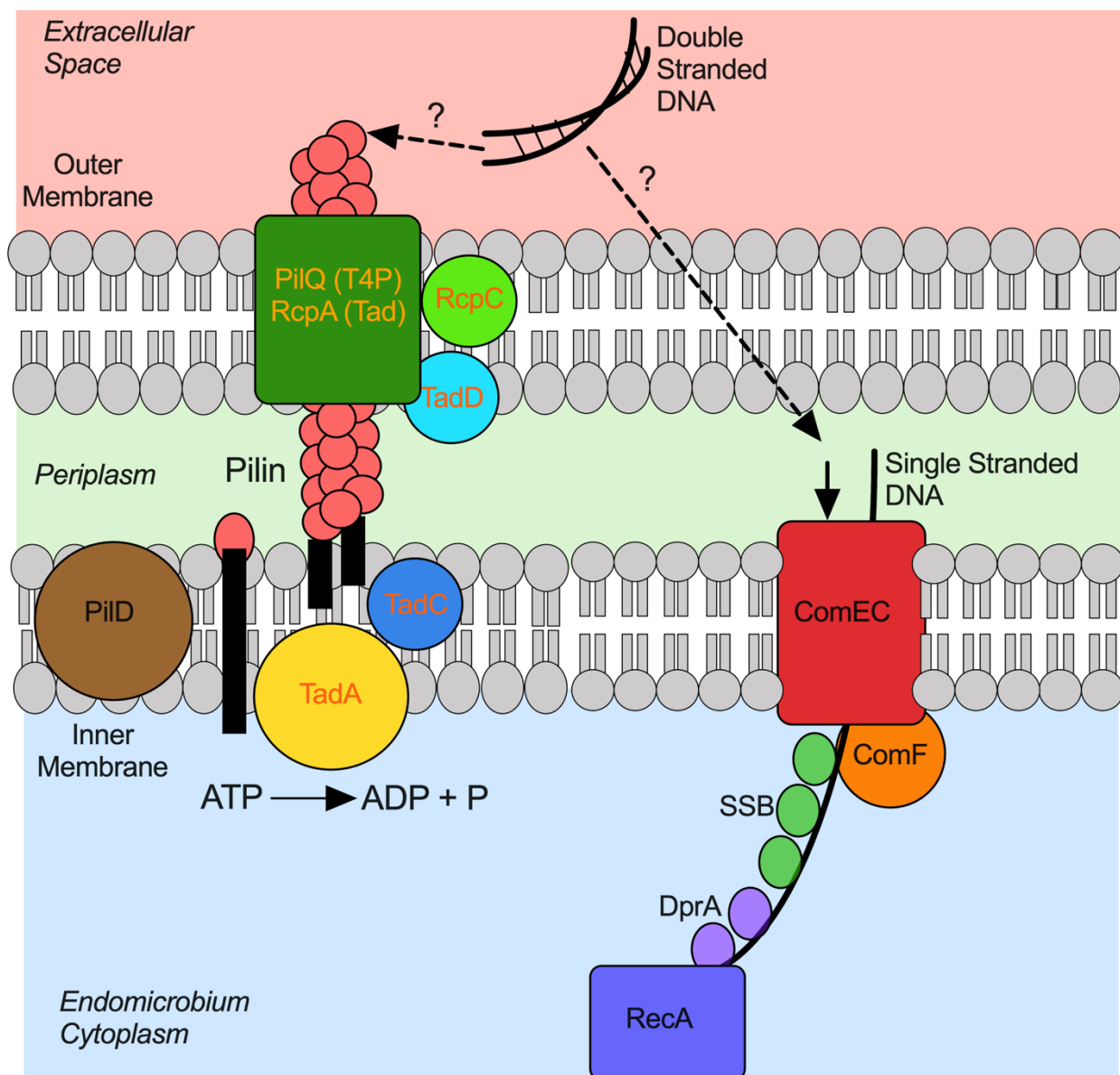799 artifact of assembly or binning.

800

801



802

**Supplemental Figure 3. 16S rRNA phylogeny of Elusimicrobia with respect to other phyla.** Maximum likelihood phylogenetic tree of 16S rRNA genes (using substitution model TPM3u + G4) which was used a marker-gene to establish an organismal phylogeny of the Elusimicrobia phylum. This phylogeny was used to determine incongruence between the 16S rRNA gene and other genes of interest that may have been acquired via HGT by the *Endomicrobium* species.

808

809



810 **Supplementary Figure 4. Gene content of regarding genes involved in natural**
811 **transformation, pilus assembly, and DNA recombination/repair.** (A) Presence and absence
812 matrices of genes involved in natural transformation, (B) pilus assembly, and (C) DNA repair
813 and recombination found the endosymbiotic *Endomicrobium* spp. genomes. The presence of
814 these genes was investigated in their free-living relatives (Rsa215 and *E. minutum* Pei191) and
815 other endosymbiotic bacteria. Accession numbers for references genomes used in this analysis
816 are provided in Supplementary Table 4. Note that the genes marked with a "*" or a "!" were not
817 in the final assemblies, but their reads were detected by either Megan (!) or by read mapping
818 onto the same gene from another closely related organisms (*) or both (!*) (see Supplementary
819 Figure 2 as an example).

820

821
822 **Supplemental Figure 5. Graphical summary of a putative competence pathway and**
823 **proteins involved in pilus assembly in *Endomicrobium* species.** Proteins shared by all
824 *Endomicrobium* species are in black font while those that are only retained in '*Ca*. E. agilae' are
825 colored in orange. All *Endomicrobium* possessed the pre-pilin peptidase (PilD) and one or more
826 genes that encode pilins. '*Ca*. E. agilae' possessed a near-complete *tad* locus/ T2SS as well the
827 T4P secretin (PilQ).
828
829

830     Supplementary File S1 spreadsheet contains the following:
831     Sup. Table 1 Meta-Data
832     Sup. Table 2 Filtering ref.
833     Sup. Table 3 Read Numbers
834     Sup. Table 4 Genome refs.
835     Sup. Table 5 RPKM values
836     Sup. Table 6 RPKM charts
837     Sup. Table 7 Type IV tad genes
838     Sup. Table 8 Biosynthetic genes
839     Sup. Table 9 Repair_Compet genes
840