

1 **Uncovering plant microbiomes using long-read metagenomic sequencing**

2 Sachiko Masuda ^a, Pamela Gan ^a, Yuya Kiguchi ^{bcd}, Mizue Anda ^b, Kazuhiro Sasaki ^e,
3 Arisa Shibata ^a, Wataru Iwasaki ^b, Wataru Suda ^d and Ken Shirasu ^{a,e}

4
5 ^a RIKEN Center for Sustainable Resource Science, Yokohama, Kanagawa, 230-0045,
6 Japan

7 ^b Department of Integrated Biosciences, Graduate School of Frontier Sciences, The
8 University of Tokyo, 277-8562, Japan

9 ^c Cooperative Major in Advanced Health Science, Graduate School of Advanced
10 Science and Engineering, Waseda University, Tokyo 169-8555, Japan

11 ^d RIKEN Center for Integrative Medical Sciences, Yokohama, Kanagawa, 230-0045,
12 Japan.

13 ^e Institute for Sustainable Agro-ecosystem Services, Graduate School of Agricultural
14 and Life Sciences, The University of Tokyo, Nishitokyo, Tokyo 188-0001, Japan

15 Present address: Japan International Research Center for Agricultural Sciences,
16 Tsukuba, Ibaraki 305-8686, Japan

17 ^e Graduate School of Science, The University of Tokyo, Bunkyo, 113-0032, Japan

18

19 Corresponding Author: Ken Shirasu, ken.shirasu@riken.jp

20

21 **Abstract**

22 The microbiome of plants plays a pivotal role in their growth and health. Despite its
23 importance, many fundamental questions about the microbiome remain largely
24 unanswered, such as the identification of colonizing bacterial species, the genes they carry,
25 and the location of these genes on chromosomes or plasmids. To gain insights into the
26 genetic makeup of the rice leaf microbiome, we performed a metagenomic analysis using
27 long-read sequences, and developed a genomic DNA extraction method that provides
28 relatively intact DNA for long-read sequencing. 1.8 Gb reads were assembled into 26,067
29 contigs, including 136 circular sequences of less than 1 Mbp, as well as 172 large (≥ 1
30 Mbp) sequences, six of which were circularized. Within these contigs, 669 complete 16S
31 rRNA genes were clustered into 166 bacterial species, 130 of which showed low identity
32 to previously defined sequences, suggesting that they represent novel species. The large
33 circular contigs contain novel chromosomes and a megaplasmid, and most of the smaller
34 circular contigs (<1 Mbp) were defined as novel plasmids or bacteriophages. One circular
35 contig represents the complete chromosome of an uncultivated bacterium in the candidate
36 phylum *Candidatus Saccharibacteria*. Our findings demonstrate the efficacy of long-

37 read-based metagenomics for profiling microbial communities and discovering novel
38 sequences in plant-microbiome studies.

39 **Introduction**

40 Plants provide a wide range of ecological niches for microbes, including the surfaces of
41 organs such as leaves, flowers, fruits and roots, internal structures that are colonized by
42 endosymbionts, the zone of intimate interactions between plant roots and microbes known
43 as the rhizosphere, and into the soil beyond the rhizosphere due to the diffusion of plant
44 products away from the root zone. These diverse environments allow microbes to form
45 complex communities collectively known as the plant microbiome¹ or phytobiome² that
46 can promote plant growth by affecting nutrient uptake, by suppressing pathogens, and by
47 inducing disease resistance³. While individual members of the plant microbiome may
48 express beneficial traits, because it is a complex system, the overall effects of the
49 microbiome on plant health cannot be predicted from individual microbial taxa³. The
50 community composition of a plant microbiome is influenced by host plant genetics^{4,5},
51 specificity and activity of plant defenses⁶, microbe-microbe interactions⁷, and
52 environmental factors such as soil geochemistry and UV light intensity. Understanding
53 the mechanisms by which the plant microbiome impacts plant health is a crucial area of
54 research, with the potential to inform the development of new strategies for improving
55 plant growth, productivity, and sustainability.

56 Plasmids and bacteriophages and the mobile genetic elements they carry,
57 including IS elements and transposons, have a role in shaping microbial communities by
58 providing novel capabilities^{8,9}, such as the well-studied VirB/VirD4 system found in both
59 pathogenic and symbiotic bacteria in plant-microbiome interactions¹⁰⁻¹³. The VirB/VirD4
60 system, typically encoded on a plasmid, enables pathogenic bacteria to invade host plants
61 by transporting effector proteins and manipulating the host's immune system. On the
62 other hand, the VirB/VirD4 system in root-nodulating bacteria is involved in protein
63 translocation and can have a host-dependent effect on symbiosis¹⁴. Despite their
64 ecological importance, our understanding of plasmid and bacteriophage sequences from
65 isolated bacteria is limited.

66 To gain deeper insights into the ecological and biological features of plant
67 microbiomes, culture-independent metagenomics has increasingly been employed as a
68 tool to analyze the genetic makeup of complex microbial communities¹. This approach
69 provides a catalog of the microbial diversity and functional potential within a given
70 community¹. However, traditional short-read (<500 bp) sequencing, typically of one or
71 two of the variable regions of 16S rRNA genes or other single genes, poses serious
72 limitations for accurate identification and genome reconstruction^{15,16}. Primer sequence
73 bias, and a large proportion of short reads that cannot be mapped to a reference genome
74 results in the loss of potentially useful information¹⁷. In contrast, long-read metagenomics

75 has the potential to generate longer contigs, thus improving genome reconstruction,
76 taxonomic assignment, and revealing previously undiscovered sequences, including
77 circular genomes and extrachromosomal elements. For instance, long-read metagenomic
78 sequencing of the human gut microbiome has uncovered a higher number of plasmids
79 than previously reported¹⁸. We anticipate that long-read metagenomics will be a valuable
80 approach for exploring plant microbiome metagenomes.

81 In this study, we used long-read metagenomics to better understand the genetic
82 makeup of the rice (*Oryza sativa*) microbiome. We enriched microbes from the
83 phyllosphere and established a genomic DNA extraction method for long-read
84 sequencing. Then, using reads from the Pacbio Sequel II sequencer, we reconstructed
85 26,067 contigs, including novel circularized chromosomes, plasmids and bacteriophages.
86 Notably, we identified the complete chromosome of the candidate phylum *Candidatus*
87 *Saccharibacteria*. Our results demonstrate that long-read based metagenomics provides
88 a powerful tool for profiling plant-associated microbial communities.

89 **Materials and Methods**

90 **Sampling and bacterial cell enrichment**

91 Rice plants (*Oryza sativa* cultivar ‘Koshihikari’) were grown in an experimental paddy
92 field at the Institute for Sustainable Agro-ecosystem Services, Graduate School of
93 Agricultural and Life Sciences, The University of Tokyo (35°74’N, 139°54’E). Plants
94 were sampled before heading on August 6th, 2018, and their roots and aerial parts were
95 separated and stored at -80 °C. The aerial parts were ground using a Roche GM200
96 grinder (2,000 rpm 15 sec Hit mode, 8,000 rpm 30 sec Cut mode and 8,000 rpm 15 sec
97 Cut mode) with dry ice. Plant-associated microbes were enriched from the ground
98 samples using a bacterial cell enrichment method as previously described¹⁹.

99

100 **DNA extraction**

101 Genomic DNA was extracted from the enriched plant-associated microbes using
102 enzymatic lysis. The cells were lysed by the addition of 20 mg/ml Lysozyme (Sigma-
103 Aldrich), 10 µl of Lysostaphin (>3,000 unit/ml, Sigma-Aldrich) and 10 µl of Mutanolysin
104 (> 4,000 units/ml, Sigma-Aldrich), and incubated for 3 h at 37 °C. SDS (20%, Sigma-
105 Aldrich) and Proteinase K (10mg/ml, Sigma-Aldrich) were then added, and DNA was
106 purified with cetrimonium bromide and phenol-chloroform. DNA was then incubated
107 with RNase for 30 min at 37 °C (Nippongene) and dissolved in TE buffer at 4 °C. The
108 sequencing library was constructed and sequenced within a week of DNA extraction. We
109 also extracted genomic DNA using a Fast DNA spin kit (MP-Biomedicals) for
110 mechanical lysis from the enriched microbes. DNA fragmentation was assessed using
111 pulsed-field gel electrophoresis.

112

113 **Sequencing of the V4 regions of 16S rRNA genes**

114 The V4 regions of 16S rRNA genes were amplified using the primers 515F (5’-ACA CTC
115 TTT CCC TAC ACG ACG CTC TTC CGA TCT GTG CCA GCM GCC GCG GTA A-
116 3’) and 806R (5’-GTG ACT GGA GTT CAG ACG TGT GCT CTT CCG ATC TGG
117 ACT ACH VGG GTW TCT AAT-3’) and sequenced with an Illumina MiSeq v3 platform.
118 The first 20 bases of primer sequences were trimmed from all paired reads. For low-
119 quality sequences, bases after 240 bp and 160 bp of forward and reverse primer sequences,
120 respectively, were truncated using Qiime2 v2018.11.0²⁰. The processed reads were
121 aligned to the SILVA123 dataset²¹, and their taxonomy was provisionally determined.

122

123 **Sequencing, assembly, and gene annotation of the plant microbiome**

124 SMRTbell libraries for sequencing were constructed according to the manufacturer's
125 protocol (Part Number 101-693-800 Version 01) without shearing. The libraries were cut
126 off at 20 kbp using the BluePippin size selection system (Sage Sciences). Libraries were
127 sequenced on two SMRT Cells 8M (Pacific Biosciences). We removed contaminating
128 plant sequences, subreads showing more than 80% identity and 80% length coverage
129 according to minimap2 v 2.14²² to 'Nipponbare' as the reference rice genome²³, and
130 PacBio's internal control sequences from subreads. 'Nipponbare' was used as the
131 reference genome because the draft genome of 'Koshihikari' is highly fragmented, with
132 an average read length of 32 bp²⁴. The remaining subreads were assembled using Canu
133 (version 1.8)²⁵ with the parameters previously described²⁶. After assembly, we removed
134 contaminated contigs derived from internal controls and reference genome using the same
135 method, and artificial contigs with long stretches of G, C, A or T by calculation of GC
136 contents with seqkit v0.11.0²⁷. For quality assessment of the contigs¹⁸, we aligned the
137 error-corrected reads generated during assembly to the contigs with pbmm2 v 1.2.1
138 (Pacific Biosciences) with maximum best alignment 1 and minimum concordance
139 percentage 90 set as parameters, and extracted the contigs with a depth >5. Contig
140 circularity was determined as previously described²⁶. For confirmation of quality, error
141 corrected reads were aligned to contigs using pbmm2 with the same parameters as for
142 contig quality assessment, then gaps and coverage were assessed using IGV browser v
143 2.8.2²⁸. Quast v 5.0.2²⁹ was used to evaluate the quality of genome assemblies. Functional
144 annotation of bacterial genes was conducted using PROKKA v 1.14.6 in the metagenomic
145 mode³⁰, COG database (BLASTP with the *e* value lower than 1e-05), Interproscan v 5.46-
146 81.0 (with the *e* value lower than 1e-05) and kofamscan v1.3.0³¹. Augustus v. 3.4.0 was
147 used to annotate the genes of fungal genomes³².

148

149 **Estimation of microbial composition using 16S rRNA genes**

150 We obtained bacterial 16S rRNA gene sequences from NCBI BioProject 33175 (Bacteria)
151 and 33317 (Archaea), removed those $\leq 1,400$ bp in length, and clustered the remaining
152 sequences using CD-HIT version 4.8.1 ($\geq 98.7\%$ identity)³³. The resulting curated 16S
153 rRNA gene database contains 11,782 distinct sequences. In the long read-based assembly
154 data, 16S rRNA genes longer than 1,400 bp on contigs having an average read depth ≥ 5
155 were aligned to our curated 16S rRNA gene database to obtain the maximum number of
156 target hits. Alignments with $<95\%$ length coverage were removed. We used 16S rRNA
157 genes with $\geq 98.7\%$ identity as the top hits for approximating bacterial community
158 composition at the species level. We counted the depth of the contigs carrying 16S rRNA
159 genes to estimate their abundance.

160 Full-length 16S rRNA genes were amplified using the primers 27F (5'- AGR
161 GTT YGA TYM TGG CTC AG) and 1492R (5'- RGY TAC CTT GTT ACG ACT T),
162 and sequenced on a SMRT cell 1M v3. Circular consensus sequences (>3 paths) were
163 constructed and demultiplexed using SMRTLink v 8.0.0 with default parameters. Primers
164 and chimeric reads were removed from demultiplexed CCS reads using dada2 v 1.16³⁴,
165 and reads $\geq 1,400$ bp were extracted. Full-length 16S rRNA amplicons were aligned with
166 our curated database to assign taxa and to estimate bacterial community composition at
167 the species level ($\geq 98.7\%$ identity). 16S rRNA gene sequences in the metagenomic data
168 were aligned with MAFFT v7.475³⁵ using default parameters. A phylogenetic tree was
169 constructed using RAxML v8.2.12³⁶ and visualized using FigTree v1.4.4
170 (<http://tree.bio.ed.ac.uk/software/figtree/>).

171

172 **Taxonomy assignment of large contigs**

173 16S rRNA gene similarity and average nucleotide identity (ANI) were used to assign
174 contigs larger than 1Mbp (n = 172, including 6 circular contigs) to taxa with GTDB-tk v
175 1.1.1 using default parameters³⁷. To provisionally assign taxa to contigs that were not
176 assignable using either of the above methods, the annotated genes on the contigs were
177 aligned to the nt database in NCBI using BLASTN with $\geq 80\%$ identity and $\geq 80\%$ length
178 coverage. Comparisons of circular contig sequences to reference genomes were plotted
179 with mummerplot, and contig completeness was calculated using checkM³⁸. Circular
180 contigs were classified as chromosomes or megaplastids according to the criteria of
181 diCenzo and Finan³⁹. Interproscan was used to search for genes encoding replication
182 proteins, such as DnaA and RepA.

183

184 **Classification of circular contigs less than 1 Mbp**

185 The sequences of the circular contigs were compared to known sequences obtained from
186 prokaryotic reference or representative genomes in the NCBI database and Plsdb version
187 2020_06_29⁴⁰ using nucmer. Interproscan was used to search for plasmid-enriched or
188 virus-related genes, and to classify those contigs as plasmid or phage. Kofamscan was
189 used to annotate VirB/VirD4 systems on plasmids. Plasflow and Mob-typer via MOB-
190 suite⁴¹ were also used for plasmid prediction, and for predictions of the plasmid host.
191 Virsorter2 (a score cut off >0.8)⁴² was used for predicting whether or not each contig
192 originated in a bacteriophage. CheckV⁴³ was used for assessing contig quality of viral
193 sequences. To provisionally assign the taxonomy of each contig, we aligned all genes on
194 the contig to the nt database in NCBI. The taxonomy of a contig was assigned as follows:
195 if more than one-fourth of the genes on a contig showed $> 80\%$ identity and $> 80\%$

196 coverage to the corresponding genes of the nt database in NCBI, we assigned the
197 taxonomy of the contig accordingly. If the genes on a contig were aligned to a strain of a
198 genus, but to a different species, the taxonomy of the contig was estimated at the genus
199 rank. Reliable taxonomy assignment was limited to cases where the number of the genes
200 aligned to known sequences was greater than one-fourth of the contig, and the genes on
201 the contig derived from one phylum. In all other cases, we concluded that the taxonomy
202 of those contigs was unassigned. However, in some cases Mob-typer could assign the
203 taxonomy of contigs which were not assigned using nt database. We tentatively assigned
204 taxonomy using the Mob-typer result for these contigs. Gene maps were drawn with
205 ‘ggplot2’ in R (<https://www.R-project.org/>).

206

207 **Taxonomic assignment and predicted gene function.**

208 We aligned all predicted genes to the COG database with an e value $< 1e-05$ to predict
209 gene function¹⁸. A similarity search of the annotated genes was conducted against the nt
210 database in NCBI using BLASTN with $\geq 95\%$ identity and $\geq 90\%$ coverage. From these
211 genes, we extracted those which were identified by species, and counted both the number
212 of genes and contigs that carried these genes.

213

214 **Genomic features of *Candidatus Saccharibacteria* (TM7)**

215 We obtained the 16S rRNA genes of *Candidatus Saccharibacteria* (formerly known as
216 TM7) from the NCBI database, removed sequences $\leq 1,400$ bp, and clustered the
217 remaining using CD-HIT at 98.7% identity. The genomic sequences of RAAC3
218 (GenBank: CP006915.1), *Candidatus Saccharimonas aalborgensis* (S_aal, GenBank:
219 CP005957.1), GWC2 (GenBank: CP011211.1), TM7x (GenBank: CP007496.1) and
220 YM_S32 (GenBank: CP025011.1) were obtained for genomic comparisons. Kofamscan
221 and BlastKOALA⁴⁴ were used to predict metabolic features. Average amino acid identity
222 was calculated using the Kostas lab AAI calculator with default parameters (<http://enve-omics.ce.gatech.edu/aai/>).

223

224 **Data availability**

225 Metagenomic data has been deposited in NCBI under accession number SAMN32580422
226 (BioSample), SRR23280466 and SRR23280465 (SRA) and PRJNA929667 (BioProject).

227

228 **Results**

229

230 **DNA preparation for long-read metagenomics**

231 To study rice-associated microbial communities in an agricultural environment, we
232 harvested 8-week old rice plants from a paddy field and extracted their leaf-associated
233 microbes using a density gradient centrifugation method to remove rice tissues¹⁹
234 (Extended Data Fig. 1A). We then compared enzymatic and mechanical cell lysis
235 methods for their ability to extract intact genomic DNA with minimal damage using
236 pulsed-field gel electrophoresis. This analysis showed that enzymatic lysis yielded intact
237 chromosomes, whereas mechanically lysed cells yielded fragmented DNA (Extended
238 Data Fig. 1B). Comparisons of bacterial community composition using the 16S rRNA
239 gene V4 region indicated that the two methods gave similar results (Pearson's correlation
240 coefficient = 0.84, Extended Data Fig. 2), although some phyla were more highly
241 represented in one method than in the other, possibly due to differences associated with
242 cell lysis of different taxa. Given the importance of intact genomic DNA for long-read
243 sequencing, we chose to use enzymatic lysis for further analyses.

244

245 **Long-read metagenomic sequencing of leaf-associated microbes.**

246 We sequenced gDNA from the leaf-associated microbes using two Sequel II 8M cells,
247 yielding 140 Gbp of data with a mean read length of 17 kbp and a mean library size of 15
248 kbp (Supplementary Table1), indicating that our DNA extraction method was suitable for
249 PacBio long-read sequencing. We obtained 26,067 contigs in total after assembly, with
250 an N₅₀ of 127 kbp, including 136 circular contigs smaller than 1 Mbp and 172 ≥ 1 Mbp,
251 6 of which were circularized (Table 1). A previous study reported that PacBio contigs
252 with ≥ 1 and 5 read depths had ≥ 98.5% and 99.4% identity, respectively, when aligned
253 to short-read contigs¹⁸. We thus were able to define 13,050 contigs with a depth of more
254 than 5 as high quality contigs. These contigs represented approximately half of the total,
255 and represented about 80% of the total reads (Table 1). Additionally, more than 90% of
256 the total nucleotides were found in the set of contigs with a length of ≥ 50 kbp. Importantly,
257 all large size contigs (≥ 1Mbp) were of high quality (Table 1). These data suggest that
258 nucleotide sequences obtained from the rice phyllosphere microbiome are reliable for
259 estimating bacterial community composition and their functions within the community.

260

261 **Estimation of microbial composition using 16S rRNA genes in long-read** 262 **metagenomics.**

263 We extracted 16S rRNA gene sequences $\geq 1,400$ bp in length and $\geq 95\%$ coverage of the
264 top hits from high quality contigs from the metagenomic data. A total of 669 16S rRNA
265 genes were identified on 561 contigs, representing 4.4% of high quality contigs
266 (Supplementary Table 2). A phylogenetic tree was used to summarize the taxonomy of
267 the detected 16S rRNA genes in the metagenome (Fig. 1). Many of 16S rRNA genes were
268 clustered with top-hit sequences at various taxonomic ranks, such as the species and genus
269 levels, but some were independently clustered. Of the 669 16S rRNA genes, 194 were
270 clustered into the *Methylobacterium* genus. The 669 16S rRNA genes clustered into 166
271 bacterial species, 130 of which had $\leq 98.7\%$ identity with any organism in the database,
272 suggesting that they represent novel species (Supplementary Table 3). Clustering the 16S
273 rRNA genes using a threshold for bacterial taxonomy⁴⁵ showed that 290 of the 16S rRNA
274 genes on 231 contigs were $\geq 98.7\%$ identical to sequences attributable to known taxa
275 (Table 2), but 378 were $\leq 98.7\%$ identical to known taxa (Table 2). Among the latter, 16
276 were ≤ 82 and, 1 was $< 78\%$, suggesting that they potentially represent a novel order and
277 class, respectively (Table 2). Ten of the 16S rRNA genes represent a putative novel order
278 within *Planctomycetes*.

279 Taxonomic analysis of the high-confidence identity 16S rRNA sequences (290
280 sequences with $\geq 98.7\%$ identity) identified 40 bacterial species (Supplementary Table
281 4). For example, *Curtobacterium pusillum* was the most abundant species; 57 of 16S
282 rRNA genes were identified on 46 contigs with a relative abundance of 15.2%. Similarly,
283 8 species of *Methylobacterium* were identified in the high-confidence identity sequences,
284 with three species (*M. indicum*, *M. radiotolerans* and *M. komagae*) having more than 10
285 16S rRNA genes (Supplementary Tables 2 and 4). The number of contigs containing 16S
286 rRNA genes was mostly lower than the number of 16S rRNA genes, indicating that these
287 bacteria carry multiple 16S rRNA gene copies (Supplementary Tables 2 and 4). In
288 particular, nine of the 16S rRNA genes from *Exiguobacterium acetylicum* were detected
289 on a single 1.7 Mbp contig (Contig ID: RRA86345 in Supplementary Table 2).

290 To verify bacterial community composition based on the long-read
291 metagenome, full-length 16S rRNA genes were amplified using universal bacterial
292 primers and sequenced to compare the taxonomic profiles. Among the 6,678 reads
293 obtained, 2,958 reads (44%) have $\geq 98.7\%$ identity to taxonomically classified 16S
294 rRNAs, a result similar to the metagenomic data (Extended Data Fig. 3). We identified
295 64 taxa (Supplementary Table 4) to the species level, with *C. pusillum* being the most
296 abundant, also consistent with the long-read metagenomic data. There were 16 species
297 with relative abundance greater than 1%, 13 of which were also detected in the long-read
298 metagenome (Supplementary Table 4). The combined relative abundance of the 13

299 species was 42.0% in the metagenome and 33.3% in the 16S rRNA nearly full-length
300 amplicon data sets.

301 We also compared the relative abundance of the rice bacterial community
302 determined by long-read metagenome sequencing versus short-read 16S rRNA
303 sequencing (Fig. 2), and found that the relative abundance of *Actinobacteria* was about
304 twice as high in the long-read metagenome (28.6%) than in either the nearly full-length
305 16S rRNA amplicon dataset (15.8%) or the V4 region dataset (15.2%). Comparing the
306 relative abundance of *Actinobacteria* showed that the relative abundance of
307 *Micrococcales* in the metagenome was about twice that of PCR-based amplicon
308 sequencing. The difference may be because certain classes of *Actinobacteria* were
309 difficult to detect using universal primers, even though the target sequences are identical⁴⁶.
310 Previous studies also showed that the V4 region is less reliable for classifying
311 *Actinobacteria* sequences⁴⁷. Therefore, our results suggest that long-read metagenome
312 sequencing provide more accurate identification about dominant bacterial communities
313 in the aerial parts of rice, particularly for *Actinobacteria*.

314

315 **Taxonomic assignment of predicted genes.**

316 We identified a total of 2,046,382 predicted genes in the metagenome (Table 1, Extended
317 Data Fig. 4). Of these putative genes, 364,262 had an e-value of 1.0e-05 or less and were
318 annotated using the COG database. Approximately 20% of the genes were categorized as
319 poorly characterized group 'R' (11.8%, general function prediction only) or 'S' (9.6%,
320 function unknown). 8% of the genes were annotated as amino acid metabolism (E), 6.5%
321 as carbohydrate transport and metabolism (G), and 6.2% as energy production and
322 conversion (C). Among these five categories, 50 - 70% of the genes were derived from
323 *Alphaproteobacteria*, particularly *Methylobacterium* (12.9 - 17.3 %). The putative genes
324 from Planctomycetes were the second most abundant (9.6 - 18.0 %, Extended Data Fig.
325 4). These results showed that *Methylobacterium* is the dominant genus in the rice
326 phyllosphere, supporting the bacterial community composition predicted by 16S rRNA
327 genes (Fig. 1).

328

329 **Taxonomic assignment of large contigs.**

330 16S rRNA gene sequences and ANI (GTDB-tk) were used to assign the taxonomy of 172
331 high quality contigs, which ranged from 1 Mbp to 8.5 Mbp (Fig. 3, Supplementary Table
332 5). 16S rRNA genes were found in 98 contigs (97 chromosomal and one in a
333 megaplasmid), and ANI was able to assign taxonomy of 147 of the contigs (Fig. 3,
334 Supplementary Table 5). 92 contigs were assigned by both the 16S rRNA gene and ANI.

335 We also assigned the taxonomy of one contig using a similarity search of its genes. In
336 total, the taxonomy of 157 contigs was assigned by these methods (Fig. 3, Supplementary
337 Table 5), but 19 contigs could not be assigned using either method. Among these, the
338 genes of 4 contigs showed a high identity ($\geq 95\%$) and length coverage ($\geq 90\%$) to the
339 *Moesziomyces antarcticus* (Fig. 3, Supplementary Table 5), suggesting that these four
340 contigs originated from a yeast.

341 We also attempted to discriminate between chromosomal and plasmid contigs
342 using ANI and the presence or absence of DNA replication initiators such as DnaA (for
343 bacterial chromosomes) and RepA (for some plasmids). We classified only one contig as
344 a megaplasmid (RRA6539; Fig. 3, orange color in the contig size column), as it showed
345 high similarity to a plasmid (NZ_CP049244.1) of *Rhizobium pseudoryzae*, which also
346 carries 16S rRNA genes on both its chromosome and plasmid. In addition, four contigs
347 were derived from the yeast *M. antarcticus* using blast search and three of four contigs
348 were carried minichromosome maintenance proteins (MCM2, 3, 6 and 10), suggesting
349 that those three contigs may be a minichromosome of *M. antarcticus* (Fig. 3,
350 Supplementary Table 5). In total, 156 of the large high quality contigs were classified as
351 bacterial chromosome, one was a megaplasmid, and four were fungal sequences. The
352 other 11 large contigs were not classified as either chromosomal or plasmids using these
353 methods.

354 Of the 156 bacterial chromosomal contigs, five were circularized, suggesting
355 that they are complete chromosomes (Black star in Fig. 3, Supplementary Table 5). The
356 taxonomy of these genomes can be tentatively assumed by 16S rRNA gene sequence
357 similarity and/or ANI, though the nucleotide sequences of some contigs do not match
358 those of sequenced strains (Extended Data Fig. 5). For example, four contigs (RRA2267,
359 RRA3045, RRA85519 and RRA944769) carry 16S rRNA genes with 87.4% - 99%
360 identity to the top hit (Fig. 3 and Supplementary Table 5). In particular, RRA944769
361 could be classified as a complete genome of a novel family of *Oligoflexales* based on 16S
362 rRNA gene identity and ANI (Fig. 3 and Supplementary Table 5). The other three contigs
363 (RRA2267, RRA3045 and RRA85519) and one additional contig (RRA944769) were
364 placed at the genus and order rank by ANI, respectively, all of which were consistent with
365 the 16S rRNA-based assessment (Fig. 3 and Supplementary Table 5). Curiously, no 16S
366 rRNA gene was detected in RRA2326, but it was assigned to the genus *Aureimonas* by
367 ANI (Fig. 3 and Supplementary Table 5). This confirms a previous report showing that
368 *Aureimonas* sp. AU20, isolated from the rice phyllosphere, has its 16S rRNA gene on a
369 small plasmid, but not on the chromosome⁴⁸. Additionally, the 11 unclassified contigs

370 were shown to have >90% completeness and <5% contamination by CheckM, suggesting
371 that those 11 contigs were nearly-complete chromosomes³⁸.

372

373 **Classification of circular contigs smaller than 1Mbp**

374 Of the 136 circular contigs ranging in size from 8.5 kbp to 832 kbp, with the GC content
375 from 36.8% to 75.2% (Fig 4, Supplementary Table 6), the sequences of 134 did not align
376 to any known sequences, suggesting that these are novel sequences. The remaining two
377 contigs were aligned with high similarity to a plasmid of *Methylobacterium*
378 *phyllosphaerae* strain CBMB27 (NZ_CP015369.1) (red color in contig size column in
379 Fig 4, Extended Data Fig. 6). 61 genes on 41 of the 136 contigs were annotated as *repC*
380 (Fig 4, Supplementary Table 6), suggesting that these contigs are novel *repABC* plasmids.
381 Plasmid hosts were identified using a similarity search from construction of a
382 phylogenetic tree using the RepC protein sequences and mob-typer. Sixteen contigs were
383 associated with *Alphaproteobacteria*, while the remaining 23 contigs could not be
384 assigned using these methods (Fig 4, Extended Data Fig. 7, Supplementary Table 6).
385 These results show that the likely origin of nearly two-thirds of the *repABC* plasmids
386 detected in this study is bacteria that have not been reported to carry *repABC* plasmids.

387 An additional 29 contigs were classified as dsDNA bacteriophages with a high
388 score (> 0.8) from virsorter2 with a CheckV completeness ranging from 14.3 to 100%
389 (Fig 4, Supplementary Table 6). Twenty-one of these contigs carried putative phage-
390 related genes, such as for capsid proteins. Thirteen carried putative partitioning protein
391 genes, seven contigs carried VirB/VirD4 component genes, and one had a relaxosome
392 protein TraY gene (Fig 4, Supplementary Table 6). Although we identified presumptive
393 novel bacteriophages in this study, we were unable to assign their taxonomy.

394 We identified 59 contigs that carried presumptive VirB/VirD4 component,
395 relaxosome protein, or type II toxin-antitoxin system genes, but were not classified as
396 either plasmids or bacteriophages by mob-typer and virsorter2. These genes are more
397 commonly plasmid-borne than chromosomal⁵, suggesting that these contigs may
398 represent incomplete plasmids (Fig 4, Supplementary Table 6). Similarity searches of the
399 genes on these contigs suggested that 21 out of the 59 contigs were from
400 *Alphaproteobacteria*, *Gammaproteobacteria* (*Pseudomonas*), or *Actinobacteria* (Fig 4,
401 Supplementary Table 6).

402 Pathogenic bacteria inject their Type 4 Secretion System (T4SS) effector
403 molecules directly into host cells, thereby altering host cell functions. The 11 gene
404 products of the *Agrobacterium tumefaciens* *virB* operon, together with the VirD4 protein,
405 are thought to form a membrane complex which facilitates the transfer of T-DNA to plant

406 cells. VirB/VirD4 T4SS components were present on 11 contigs (Fig. 5), nine of which
407 were *repABC* type plasmids, and the remaining two were possibly plasmids. The likely
408 origins of seven of these contigs were *Proteobacteria*, including *Rhizobiaceae*, and
409 *Rhodobacterales* in the *Alphaproteobacteria* group. The remaining four could not be
410 assigned to a taxonomic group. A comparison of the gene arrangements on the 11 contigs
411 with those of *A. tumefaciens* showed that all, or almost all components of VirB/VirD4
412 T4SS were present, although in a different order than in *A. tumefaciens*, with some gene
413 duplication and missing components. An additional 52 contigs carried at least one
414 component gene of the VirB/VirD4 T4SS.

415 We identified five small circular contigs carrying a presumptive *dnaA* gene,
416 which is typically found in bacterial chromosomes as part of the DNA replication
417 machinery (Fig. 4 and Supplementary Table 6). A similarity search revealed that the
418 origin of two of those contigs was *Rickettsia*, which are obligate intracellular α -
419 *proteobacteria* associated with various eukaryotic hosts. Notably, approximately half of
420 the 26 validated *Rickettsia* species have plasmids, some of which carry a *dnaA*-like gene
421 and range from in size from 12 kb to 83 kb⁴⁹. We also detected *dnaA* genes on contigs
422 that were classified as bacteriophage, potentially plasmid, or from the *Candidatus*
423 *Saccharibacteria* chromosome. Two other contigs originated from *Methylobacterium*,
424 but we were unable to classify these contigs as plasmid or bacteriophage based on the
425 available gene information. Four contigs could not be classified as chromosome, plasmid,
426 or bacteriophage due to a lack of similarity to any known bacterial or bacteriophage-
427 derived genes in public databases. In total, our analysis presumptively identified one
428 chromosome, 100 plasmids (41 *repABC*-type plasmids and 59 potentially plasmid-
429 associated contigs), 29 bacteriophages, and six unclassified contigs, demonstrating that
430 long-read metagenomic sequencing can effectively be used to identify a large number of
431 plasmids from a complex microbial community, most of which are novel.

432

433 **Complete genome of a bacterium in the *Candidatus Saccharibacteria* as-yet** 434 **uncultured phylum**

435 One of the key benefits of long-read metagenomic analysis is the potential to obtain
436 complete genome sequences of uncultivable microorganisms. Here, we obtained the
437 whole chromosome sequence of a member of the uncultured *Candidatus*
438 *Saccharibacteria* phylum as a circular contig (RRA8490). Phylogenetic analysis based
439 on 16S rRNA genes indicated that RRA8490 clusters with isolates found in the human
440 oral microflora (Extended Data Fig. 8A). A comparison of whole genome sequences and
441 amino acid identities between RRA8490 and previously determined strains in the

442 *Saccharibacteria*⁵⁰⁻⁵⁴ showed that the genomes of these strains and RRA8490 were
443 distinct, with average amino acid identities ranging from 52.2 to 54.2% (Extended Data
444 Fig. 8B and C). Using Kofamscan and Interproscan, we searched for metabolism-related
445 genes in RRA8490 and found that it encoded all the presumptive genes necessary for the
446 biosynthesis of peptidoglycan (MurABCDEFG, MraY and MtgA), suggesting that its cell
447 wall is of the gram-negative type (Fig. 6). Unlike other strains, RRA8490 did not encode
448 amino acid or fatty acids synthesis genes⁵⁵, but it did presumptively encode enzymes that
449 metabolize glucose to ribose and glycerate-3-phosphate, as well as phosphoenol-pyruvate
450 to malate, suggesting that these pathways may be used to generate ATP. RRA8490 also
451 encoded four regions of type IV pili (*pilM₁N₁O₁B₁TC₁D*, *pilB₂C₂M₂N₂O₂*, *pilB₃*, *pilB₄*),
452 similar to a previously reported *Saccharibacteria* (TM7) genome that carries type IV pili
453 for host cell attachment⁵³ (Fig. 6). In addition to these characteristics, RRA8490 also
454 encoded the cytochrome oxidase complex CyoABCDE, as previously reported.

455

456 **Discussion**

457 This study demonstrated that enzymatic genomic DNA extraction combined with long-
458 read metagenomic sequencing is an effective tool for profiling plant microbiota and their
459 genomes, as well as defining complete chromosomes, plasmids and bacteriophages from
460 long, high quality contigs. Importantly, comparisons of the community-profiling datasets
461 obtained using short-read 16S rRNA amplicon sequencing confirms that the enzymatic
462 DNA extraction method is largely unbiased, with the notable exception of the inclusion
463 of fungal DNA from the *Moesziomyces* genus. This is not surprising, as *Moesziomyces*
464 spp. are commonly detected in plants^{56,57}, but not amplified by bacterial 16S primers.
465 Therefore, the community profile of rice leaves obtained using long-read metagenomics
466 is consistent with previously reported datasets. For example, our data indicate that *C.*
467 *pusillum* is the dominant species in rice leaves (Supplementary Table 2 and 4), which is
468 consistent with the fact that *Curtobacterium* spp. have been isolated from the leaves of
469 many different plants⁵⁸⁻⁶¹ and are known to be abundant in a leaf litter communities⁶².
470 Similarly, the *Methylobacteriaceae* family is a dominant presence (Fig. 1), as seen in the
471 aerial parts of many plants⁶³.

472 A potential 130 novel species were identified from 669 16S rRNA genes. The
473 relative abundance of 16S rRNA genes detected in the metagenome ranged from 0.02%
474 to 1.8%, demonstrating the depth of coverage provided by long-read metagenomics.
475 However, a comparison between bacterial species detected in the metagenome and nearly
476 full-length 16S rRNA amplicon sequencing revealed that some species were only
477 detected using one method or the other, despite having relative abundances greater than

478 0.02% (Supplementary Table 4). This highlights the importance of using a combination
479 of long-read sequences, such as metagenomes, and 16S rRNA amplicon sequencing for
480 more comprehensive taxonomic assignments. Overall, long-read metagenomics is a
481 powerful tool for accurately identifying bacteria in the rice phyllosphere, with the
482 potential for greater discrimination between organisms as new analysis methods become
483 available.

484 The use of long-read metagenomics allows for the estimation of bacterial species
485 abundance currently possible, as it is not influenced by PCR amplification bias and
486 variations in the number of 16S rRNA genes present in each genome. Indeed, the
487 prevalence of *Micrococcales* (*Actinobacteria*) identified from metagenomic data was
488 about two-fold higher than from amplicon data (Fig. 2). The number and location of
489 rRNA operons (*rrn*) can vary significantly among bacteria, and some bacteria have
490 multiple copies of *rrn* on different chromosomes, as seen in *Brucella*⁶⁴ and *Vibrio*⁶⁵. Our
491 data demonstrated that the relative abundance of *E. acetylicum* is often overestimated due
492 to the presence of nine copies of 16S rRNA genes on four contigs in the draft genome of
493 *E. acetylicum*⁶⁶. In contrast, long-read metagenomics-based identification of the precise
494 number of 16S rRNA genes allows for accurate determination of bacteria abundance
495 within a community. When 16S rRNA genes were not present on contigs, we were able
496 to use ANI to tentatively assign taxonomy for nearly 30% of the contigs (Fig. 3 and
497 Supplementary Table 5), because long-read metagenomics allows the reconstruction of
498 very large contigs.

499 The use of long reads allowed us to identify the genes contained within six
500 large circular contigs. Two of the contigs (RRA2326 and RRA3045) appeared to be
501 derived from known species, but the others had little similarity to previously whole-
502 genome sequenced strains (Fig. 3 and Supplementary Table 5), suggesting that the four
503 remaining circular contigs represent the chromosomes and a megaplasmid of novel
504 species. Notably, the complete chromosome of *Rhizobium giardini* (RRA85519) was
505 sequenced and identified for the first time, which can serve as a valuable reference for
506 this species. A novel strain in the genus *Oligoflexia* (RRA944769) was also identified,
507 which differed in genomic composition and presumptive energy metabolism pathways
508 from previously isolated and sequenced strains^{67,68}. For instance, RRA944769 encodes
509 putative nitrate- and nitrite reductases (NO₃ to NO), whereas *Oligoflexia tunisience*
510 Shr3^T has genes converting from NO₂ to N₂. *O. tunisience* Shr3^T also has *aa3*- and *cbb3*-
511 type cytochrome c oxidases, whereas RRA944769 has cytochrome *b*₆ in addition to *aa3*-
512 and *cbb3*. Our study also confirmed the natural occurrence of a chromosome lacking
513 16S rRNA genes in an *Aureimonas* sp. genome (RRA3045), resolving previous

514 uncertainty around chromosomal rearrangements during cultivation⁴⁸. These examples
515 demonstrate the power of long-read metagenomics in accurately identifying and
516 characterizing the genetic makeups of a complex sample.

517 Long-read metagenomics has a major advantage over short-read metagenomics
518 in that it can be used to define circular mobile elements such as plasmids and temperate
519 bacteriophages. These elements play a significant role in microbiome interactions and
520 horizontal gene transfer^{8,69}. Short-read methods have largely been unable to fully
521 assemble complete plasmids and bacteriophages, meaning that many of the genes present
522 on these elements have not been identified. In this study, among 136 small circular contigs,
523 only two were found to align well with known plasmids of *M. phyllosphaerae*, suggesting
524 that the majority of these sequences represent novel plasmids or bacteriophages. Contigs
525 that encode functions often found on plasmids, such as toxin-antitoxin systems and T4SS,
526 are likely to be plasmids¹⁸. Those functions are often important for modulating
527 interactions with plant cells and other microorganisms⁷⁰. Approximately 22% of these
528 contigs putatively encode RepC, a protein involved in *repABC* plasmid replication, which
529 are common in *Alphaproteobacteria*⁷¹. In fact, almost half of the RepC-encoding genes
530 identified in the novel plasmids appear to be clustered with *Rhodobacteraceae*,
531 *Rhizobiales* and *Hyphomicrobiales*. Other RepC-encoding genes could not be assigned to
532 any taxon (Fig. 4 and Extended Data Fig. 7), suggesting that they originated from
533 unidentified bacteria. However, predicting a plasmid host, particularly for broad host
534 range plasmids, is a challenging task in itself, let alone for metagenomic studies⁹. While
535 we have made suggestions about the host of origin for some of the plasmids, the origin of
536 others remains unclear. New technologies, such as droplet microfluidics for isolation of
537 single bacterial cells combined with plasmid-specific markers may help to address this
538 deficiency in the future⁷².

539 One of the best studied plasmid-specific functions is the VirB/VirD4 system of
540 the Ti plasmid from *Agrobacterium tumefaciens*. The *virB* operon (*B1-11*) together with
541 *virD4* encode a putative T4SS. T4SS are highly diverse⁷³⁻⁷⁷, and the exact number of
542 genes and their role in T4SS assembly or function is unknown in many classes of T4SS⁷⁸.
543 The VirB/VirD4 system identified in RRA9653 carried *virB1-11* and *virD4* gene
544 homologs (Fig. 5). Four additional putative plasmids (RRA13979, RRA16697,
545 RRA22037 and RRA86817) were conserved for *virB2-B11* of the eleven *virB* and *virD4*
546 genes, which are essential for construction of the VirB/VirD4 system in *A. tumefaciens*
547 (Fig. 5), suggesting that these five plasmids are the VirB/VirD4 system of *A. tumefaciens*
548 type. Interestingly, Ti plasmids belong to the *repABC* family, which is widely distributed
549 among many species of *Alphaproteobacteria*. However, no *repABC* genes were found on

550 two of the plasmids (RRA9653 and RRA22037), suggesting that these plasmids may have
551 been horizontally transferred from other bacterial taxa.

552 Long-read metagenomics provide genomic information for poorly studied or as-
553 yet uncultivated bacteria. We analyzed the genes identified in the metagenomic data that
554 occurred with high confidence identity and coverage to specific bacterial species, and
555 counted the number of contigs carrying these genes. Interestingly, the number of genes
556 from *Planctomyces bacterium* was much lower than expected, at 1,497 genes, despite
557 being carried by the largest number (405) of contigs (Extended Data Fig. 9). Similarly,
558 the number of the genes from *Microbacterium testaceum* and *Phreatobacter*
559 *cathodiphilus* was also low, but these genes were carried by a large number of contigs.
560 These findings suggest that the genomic information of these three species is largely
561 unknown. For instance, our study found a high abundance of *Planctomyces*: 56 out of
562 669 16S rRNA genes and 42 out of 172 large size contigs were derived from novel
563 *Planctomyces* (Figs. 1 and 3, Supplementary Tables 2 and 5). Because *Planctomyces*
564 are difficult to culture⁷⁹, there is limited gene/genomic information available for this
565 group. To our knowledge, the complete genome sequence of *M. testaceum* has only been
566 deposited for one strain (3.98 Mbp)⁸⁰, but using long-read metagenomics, we were able
567 to obtain seven large contigs that represent complete or nearly complete chromosomes of
568 *Microbacterium* (Fig. 3). We were able to obtain 79 contigs from *P. cathodiphilus*, of
569 which only one strain has been sequenced⁸¹ (Extended Data Fig. 9). Thus, our
570 methodology can be used to parse the ecological and biological functions of fastidious
571 bacterial groups.

572 Our long-read metagenomic analysis was able to define the complete circular
573 genome (RRA8490) of an uncultured bacterium belonging to the *Candidatus*
574 *Saccharibacteria* phylum. Members of this phylum have been detected in numerous
575 natural environments such as soils, animals, and plants, but lack of cultured isolates has
576 limited our understanding of their biology. Consequently, only a few complete genomes
577 from this phylum have been reported⁵⁰⁻⁵⁴. Compared to the recently nearly completed
578 genome (1.45 Mb) of an oat-associated member of the *Candidatus Saccharibacteria*
579 phylum tentatively designated *Teamsevenus rhizospherense* strain YM_S32, RRA8490
580 is much smaller (0.83 Mb) and belongs to a different clade (Extended Data Fig. 8). Both
581 *T. rhizospherense* and RRA8490 apparently lack the ability to synthesize amino acids
582 from central metabolites, but RRA8490 is predicted to carry type IV pili and cytochrome
583 bo3, similar to others in *Candidatus Saccharibacteria* phylum. RRA8490 is predicted to
584 be able to assimilate and metabolize glucose and fructose, which are compounds found
585 in leaf exudates of the rice phyllosphere⁸². This suggest that RRA8490 may utilize these

586 compounds as carbon sources. Additionally, *Candidatus Saccharibacteria* are obligate
587 epibionts of *Actinobacteria*, which they lyse to obtain nutrients⁵³, suggesting that
588 RRA8490 may not rely solely on plant exudates for nutrient acquisition, but may also
589 degrade *Actinobacteria*. The CyoABCDE, cytochrome o oxidase complex, is used by
590 *Rhizobium etli* to adapt to anaerobic conditions⁸³, but Cyo appears to be produced only
591 under oxygen-rich growth conditions in *E. coli*⁸⁴. These results suggest that the ability to
592 function at a wide range of oxygen concentrations, as demonstrated by the presence of
593 Cyo in RRA8490, would be beneficial for this bacterium as it adapts to a variety of
594 oxygen conditions in its natural environment.

595 In conclusion, long-read metagenomics fueled by high quality DNA extraction
596 provides an efficient method for exploring uncharted organisms in the plant microbiome,
597 and the resulting data represents an emerging primary resource for a deeper understanding
598 of plant-associated microbial ecology.

599 **Acknowledgements**

600 We are grateful to the technical staffs of the Department of Technical Development in
601 ISAS of The University Tokyo for their maintenance of plant materials. We also thank
602 Prof. Dr. K. Minamisawa and Dr. S. Hara for sharing the cell-density centrifugation from
603 rice plants protocol. This work was supported by the JSPS KAKENHI grants
604 JP20H05909 and JP22H00364 (K.Sh.) and JP20H05592 (S.M.) and by JPNP18016
605 commissioned by the New Energy and Industrial Technology Development Organization
606 (NEDO).

607

608 **Contributions**

609 All authors substantially contributed to this work. All authors approved the submitted
610 version of this manuscript. Data analysis and interpretation of data were performed by all
611 authors. S.M., and K.Sh. contributed to the design of the work. S.M., P. G., K. Sa and K.
612 Sh carried out rice sampling from an experimental field plot.

613 Table 1. Summary of assembly results.

Assembly results	total contigs	High quality contigs
Number of the contigs	26,067	13,050
Total nucleotides (bp)	1,763,254,923	1,521,823,352
Number of nucleotides (bp, \geq 50 kbp)	1,390,656,635	1,370,599,048
Largest contig size (bp)	8,528,088	8,528,088
Contig N ₅₀	127,510	185,687
Predicted CDS	2,046,382	1,674,802
Number of contigs \geq 1 Mbp	172	172
Number of circular contigs	142	132

614

615 Table 2. 16S rRNA genes above the threshold for bacterial taxonomy detected in the
616 metagenome.

Identity (%)	Taxonomic rank	number of 16S rRNA genes	number of contigs
98.7	species	290	231
94.5	genus	234	198
86.5	family	127	122
82	order	16	16
78	class	1	1

617

618 **Figure legends**

619 Figure 1. Overview of the phylogeny of 16S rRNA genes detected in the metagenome.
620 Pink circles indicate the number of 16S rRNA genes detected in the metagenome for each
621 major branch: large pink circles with numbers represent branches with more than ten
622 genes, small pink circles each represent one gene. Blue circles represent the number of
623 top species/genus identities of metagenome-derived 16S rRNA genes.

624
625 Figure 2. Relative abundance of 16S rRNA genes detected in the metagenome, compared
626 to 16S rRNA gene full-length amplicon sequences and short reads. (A) Relative
627 abundance of bacterial phyla and (B) relative abundance of the class *Actinobacteria* for
628 each sequencing method.

629
630 Figure 3. Characteristics of large contigs (>1Mbp, n=172). The taxonomic assignment of
631 each contig was determined based on 16S rRNA genes, ANI, and gene similarity searches.
632 Contigs carrying 16S rRNA genes are shown in yellow blocks. Contigs classified by ANI
633 are shown in green blocks. Contigs carrying *dnaA*, *repA* and mini-chromosome
634 maintenance genes are shown in purple, pink and light blue blocks, respectively. The blue,
635 red and black bars in the contig size column represent chromosomes, megaplasmids and
636 unclassified (neither chromosome nor plasmid), respectively. The black star indicates a
637 circular contig.

638
639 Figure 4. Characteristics of small circular contigs (<1Mbp, n=136). Gene annotations are
640 indicated by color blocks. Dark purple and light purple represent complete/nearly
641 complete, and partial genes of VirB/VirD T4SS systems, respectively. Red bars in the
642 contig size column indicate known plasmid sequences. The contig ID shown in the contig
643 size column corresponds to contigs carrying complete/nearly complete VirB/VirD T4SS.

644
645 Figure 5. Gene arrangements, predicted host, and estimated type of plasmid for T4SS
646 genes discovered in the metagenome. Purple indicates hypothetical or non-T4SS
647 component genes.

648
649 Figure 6. Predicted metabolism of RRA8490, a potential new strain in the *Candidatus*
650 *Saccharibacteria* genus. Metabolic pathways were reconstructed using kofamscan,
651 Interproscan and similarity searches.

652

653 **Extended Data**

654 Extended Data Fig. 1. Preparation of genomic DNA from the rice-microbiome for long-
655 read metagenomic sequencing. (A) Rice plants were sampled from experimental field and
656 ground with dry ice. Bacterial cells were purified from aerial parts of rice plants using
657 cell density centrifugation. (B) Genomic DNA was extracted from the purified
658 microbiome using physical and enzymatic lysis. The presence of chromosomal DNA was
659 confirmed using Pulse-field gel electrophoresis.

660

661 Extended Data Fig. 2. Comparison of the relative abundance of 16S rRNA genes extracted
662 with mechanical lysis and enzymatic lysis.

663

664 Extended Data Fig. 3. Comparison of the relative abundance of 16S rRNA genes at
665 different taxonomical ranks in the metagenome and 16S rRNA full-length amplicon
666 sequences.

667

668 Extended Data Fig. 4. Gene categorization using the COG database. The number of genes
669 categorized in each group is indicated in the heatmap.

670

671 Extended Data Fig. 5. Alignment of whole genomic sequences between the six circular
672 contigs and the closest bacterial relative genome. The genome of the reference strain,
673 *Rhizobium giardinii*, was determined by whole genome sequencing (WGS). The bold
674 dotted line of the horizontal axis represents each contig of *R. giardinii*.

675

676 Extended Data Fig. 6. Alignment of the whole genomic sequences of RRA17620 and
677 RRA19473 to the plasmid of *Methylobacterium phyllosphaerae* strain CBMB27
678 (NZ_CP015369.1).

679

680 Extended Data Fig. 7. RAxML phylogenetic tree of RepC protein encoded in small
681 circular contigs (< 1Mbp). The 61 copies of RepC on 39 contigs were used to construct
682 the tree. The accession number indicates the representative RepC in each cluster. The
683 taxonomic assignment of RepC copies that are independently clustered with known RepC
684 are defined as unknown. *Klebsiella pneumoniae* RepC was used as the outgroup.

685

686 Extended Data Fig. 8. Comparison of the *Candidatus Saccharibacteria* strain located in
687 the long read metagenomic dataset with published genomes of other strains in the genus.
688 (A) Phylogenetic tree using 16S rRNA genes of strains in *Candidatus Saccharibacteria*.

689 Contig RRA8490 is indicated by the red box, adjacent to the human oral cavity group.
690 *Candidatus Gracilibacteria* and *Candidatus Absconditabacteria* were used as outgroups.
691 (B) Comparison of whole genomic sequences between RRA8490 and five *Candidatus*
692 *Saccharibacteria* (TM7) isolates. Whole genomic sequences were compared using
693 nucmer, showing that RRA8490 is not similar to the others. (C) Amino acid identity
694 between RRA8490 and five isolates of *Candidatus Saccharibacteria*. The average amino
695 acid identity was calculated using the AAI calculator developed by the Kostas lab with
696 default parameters (<http://enve-omics.ce.gatech.edu/aai/>).

697

698 Extended Data Fig. 9. The number of genes identified in the metagenome dataset with
699 high identity and coverage to specific bacterial species and of the contigs carrying these
700 genes. The contigs with a minimum of 50 predicted genes are shown.

701

702 **Supplementary tables**

703 Table 1. Summary of the sequencing results in this study.

704 Table 2. Summary of 16S rRNA genes detected in the metagenome dataset.

705 Table 3. Summary of the 16S rRNA genes clustered with $\geq 98.7\%$ identity in metagenome
706 dataset.

707 Table 4. All 16S rRNA genes (≥ 98.7 identity) detected in the metagenome, and full-
708 length 16S rRNA gene amplicon sequences.

709 Table 5. Summary of large contigs ($>1\text{Mbp}$, $n=172$)

710 Table 6. Summary of small circular contigs ($<1\text{Mbp}$, $n=136$)

711 REFERENCES

- 712 1 Vorholt, J. A., Vogel, C., Carlstrom, C. I. & Muller, D. B. Establishing Causality:
713 Opportunities of Synthetic Communities for Plant Microbiome Research. *Cell*
714 *Host Microbe* **22**, 142-155, doi:10.1016/j.chom.2017.07.004 (2017).
- 715 2 Beans, C. Core Concept: Probing the phytobiome to advance agriculture. *Proc.*
716 *Natl. Acad. Sci. U. S. A.* **114**, 8900-8902, doi:10.1073/pnas.1710176114 (2017).
- 717 3 Trivedi, P., Leach, J. E., Tringe, S. G., Sa, T. & Singh, B. K. Plant-microbiome
718 interactions: from community assembly to plant health. *Nat Rev Microbiol* **18**,
719 607-621, doi:10.1038/s41579-020-0412-1 (2020).
- 720 4 Edwards, J. *et al.* Structure, variation, and assembly of the root-associated
721 microbiomes of rice. *Proc. Natl. Acad. Sci. U. S. A.* **112**, E911-920,
722 doi:10.1073/pnas.1414592112 (2015).
- 723 5 Ofek, M., Voronov-Goldman, M., Hadar, Y. & Minz, D. Host signature effect on
724 plant root-associated microbiomes revealed through analyses of resident vs. active
725 communities. *Environ Microbiol* **16**, 2157-2167, doi:10.1111/1462-2920.12228
726 (2014).
- 727 6 Liu, H., Carvalhais, L. C., Schenk, P. M. & Dennis, P. G. Effects of jasmonic acid
728 signalling on the wheat microbiome differ between body sites. *Sci Rep* **7**, 41766,
729 doi:10.1038/srep41766 (2017).
- 730 7 Agler, M. T. *et al.* Microbial Hub Taxa Link Host and Abiotic Factors to Plant
731 Microbiome Variation. *PLoS Biol* **14**, e1002352,
732 doi:10.1371/journal.pbio.1002352 (2016).
- 733 8 Schierstaedt, J. *et al.* Role of Plasmids in Plant-Bacteria Interactions. *Curr Issues*
734 *Mol Biol* **30**, 17-38, doi:10.21775/cimb.030.017 (2019).
- 735 9 Koskella, B. & Taylor, T. B. Multifaceted Impacts of Bacteriophages in the Plant
736 Microbiome. *Annu Rev Phytopathol* **56**, 361-380, doi:10.1146/annurev-phyto-
737 080417-045858 (2018).
- 738 10 Gordon, J. E. & Christie, P. J. The Agrobacterium Ti Plasmids. *Microbiol Spectr*
739 **2**, doi:10.1128/microbiolspec.PLAS-0010-2013 (2014).
- 740 11 Sugawara, M. *et al.* Comparative genomics of the core and accessory genomes of
741 48 Sinorhizobium strains comprising five genospecies. *Genome Biol* **14**, R17,
742 doi:10.1186/gb-2013-14-2-r17 (2013).
- 743 12 Wasai-Hara, S. *et al.* Diversity of Bradyrhizobium in Non-Leguminous Sorghum
744 Plants: B. ottawaense Isolates Unique in Genes for N(2)O Reductase and Lack of
745 the Type VI Secretion System. *Microbes Environ* **35**,
746 doi:10.1264/jsme2.ME19102 (2020).

- 747 13 Kaneko, T. *et al.* Complete genomic structure of the cultivated rice endophyte
748 *Azospirillum* sp. B510. *DNA Res* **17**, 37-50, doi:10.1093/dnares/dsp026 (2010).
- 749 14 Hubber, A. M., Sullivan, J. T. & Ronson, C. W. Symbiosis-induced cascade
750 regulation of the *Mesorhizobium loti* R7A VirB/D4 type IV secretion system. *Mol*
751 *Plant Microbe Interact* **20**, 255-261, doi:10.1094/MPMI-20-3-0255 (2007).
- 752 15 Bai, Y. *et al.* Functional overlap of the Arabidopsis leaf and root microbiota.
753 *Nature* **528**, 364-369, doi:10.1038/nature16192 (2015).
- 754 16 Ikeda, S. *et al.* Low nitrogen fertilization adapts rice root microbiome to low
755 nutrient environment by changing biogeochemical functions. *Microbes Environ*
756 **29**, 50-59, doi:10.1264/jsme2.me13110 (2014).
- 757 17 Zhu, Z., Ren, J., Michail, S. & Sun, F. MicroPro: using metagenomic unmapped
758 reads to provide insights into human microbiota and disease associations. *Genome*
759 *Biol* **20**, 154, doi:10.1186/s13059-019-1773-5 (2019).
- 760 18 Suzuki, Y. *et al.* Long-read metagenomic exploration of extrachromosomal
761 mobile genetic elements in the human gut. *Microbiome* **7**, 119,
762 doi:10.1186/s40168-019-0737-z (2019).
- 763 19 Ikeda, S. *et al.* Development of a bacterial cell enrichment method and its
764 application to the community analysis in soybean stems. *Microb Ecol* **58**, 703-
765 714, doi:10.1007/s00248-009-9566-0 (2009).
- 766 20 Bolyen, E. *et al.* Reproducible, interactive, scalable and extensible microbiome
767 data science using QIIME 2. *Nat Biotechnol* **37**, 852-857, doi:10.1038/s41587-
768 019-0209-9 (2019).
- 769 21 Quast, C. *et al.* The SILVA ribosomal RNA gene database project: improved data
770 processing and web-based tools. *Nucleic Acids Res* **41**, D590-596,
771 doi:10.1093/nar/gks1219 (2013).
- 772 22 Li, H. Minimap2: pairwise alignment for nucleotide sequences. *Bioinformatics* **34**,
773 3094-3100, doi:10.1093/bioinformatics/bty191 (2018).
- 774 23 Kawahara, Y. *et al.* Improvement of the *Oryza sativa* Nipponbare reference
775 genome using next generation sequence and optical map data. *Rice (N Y)* **6**, 4,
776 doi:10.1186/1939-8433-6-4 (2013).
- 777 24 Kobayashi, A., Hori, K., Yamamoto, T. & Yano, M. Koshihikari: a premium
778 short-grain rice cultivar - its expansion and breeding in Japan. *Rice (N Y)* **11**, 15,
779 doi:10.1186/s12284-018-0207-4 (2018).
- 780 25 Koren, S. *et al.* Canu: scalable and accurate long-read assembly via adaptive k-
781 mer weighting and repeat separation. *Genome Res* **27**, 722-736,
782 doi:10.1101/gr.215087.116 (2017).

- 783 26 Kiguchi, Y., Nishijima, S., Kumar, N., Hattori, M. & Suda, W. Long-read
784 metagenomics of multiple displacement amplified DNA of low-biomass human
785 gut phageomes by SACRA pre-processing chimeric reads. *DNA Res* **28**,
786 doi:10.1093/dnares/dsab019 (2021).
- 787 27 Shen, W., Le, S., Li, Y. & Hu, F. SeqKit: A Cross-Platform and Ultrafast Toolkit
788 for FASTA/Q File Manipulation. *PLoS One* **11**, e0163962,
789 doi:10.1371/journal.pone.0163962 (2016).
- 790 28 Robinson, J. T. *et al.* Integrative genomics viewer. *Nat Biotechnol* **29**, 24-26,
791 doi:10.1038/nbt.1754 (2011).
- 792 29 Mikheenko, A., Prjibelski, A., Saveliev, V., Antipov, D. & Gurevich, A. Versatile
793 genome assembly evaluation with QUAST-LG. *Bioinformatics* **34**, i142-i150,
794 doi:10.1093/bioinformatics/bty266 (2018).
- 795 30 Seemann, T. Prokka: rapid prokaryotic genome annotation. *Bioinformatics* **30**,
796 2068-2069, doi:10.1093/bioinformatics/btu153 (2014).
- 797 31 Aramaki, T. *et al.* KofamKOALA: KEGG Ortholog assignment based on profile
798 HMM and adaptive score threshold. *Bioinformatics* **36**, 2251-2252,
799 doi:10.1093/bioinformatics/btz859 (2020).
- 800 32 Keller, O., Kollmar, M., Stanke, M. & Waack, S. A novel hybrid gene prediction
801 method employing protein multiple sequence alignments. *Bioinformatics* **27**, 757-
802 763, doi:10.1093/bioinformatics/btr010 (2011).
- 803 33 Li, W. & Godzik, A. Cd-hit: a fast program for clustering and comparing large
804 sets of protein or nucleotide sequences. *Bioinformatics* **22**, 1658-1659,
805 doi:10.1093/bioinformatics/btl158 (2006).
- 806 34 Callahan, B. J. *et al.* DADA2: High-resolution sample inference from Illumina
807 amplicon data. *Nat Methods* **13**, 581-583, doi:10.1038/nmeth.3869 (2016).
- 808 35 Katoh, K. & Standley, D. M. MAFFT multiple sequence alignment software
809 version 7: improvements in performance and usability. *Mol Biol Evol* **30**, 772-780,
810 doi:10.1093/molbev/mst010 (2013).
- 811 36 Kozlov, A. M., Darriba, D., Flouri, T., Morel, B. & Stamatakis, A. RAXML-NG:
812 a fast, scalable and user-friendly tool for maximum likelihood phylogenetic
813 inference. *Bioinformatics* **35**, 4453-4455, doi:10.1093/bioinformatics/btz305
814 (2019).
- 815 37 Chaumeil, P. A., Mussig, A. J., Hugenholtz, P. & Parks, D. H. GTDB-Tk: a toolkit
816 to classify genomes with the Genome Taxonomy Database. *Bioinformatics*,
817 doi:10.1093/bioinformatics/btz848 (2019).

- 818 38 Parks, D. H., Imelfort, M., Skennerton, C. T., Hugenholtz, P. & Tyson, G. W.
819 CheckM: assessing the quality of microbial genomes recovered from isolates,
820 single cells, and metagenomes. *Genome Res* **25**, 1043-1055,
821 doi:10.1101/gr.186072.114 (2015).
- 822 39 diCenzo, G. C. & Finan, T. M. The Divided Bacterial Genome: Structure,
823 Function, and Evolution. *Microbiol Mol Biol Rev* **81**, doi:10.1128/MMBR.00019-
824 17 (2017).
- 825 40 Galata, V., Fehlmann, T., Backes, C. & Keller, A. PLSDB: a resource of complete
826 bacterial plasmids. *Nucleic Acids Res* **47**, D195-D202, doi:10.1093/nar/gky1050
827 (2019).
- 828 41 Robertson, J. & Nash, J. H. E. MOB-suite: software tools for clustering,
829 reconstruction and typing of plasmids from draft assemblies. *Microb Genom* **4**,
830 doi:10.1099/mgen.0.000206 (2018).
- 831 42 Guo, J. *et al.* VirSorter2: a multi-classifier, expert-guided approach to detect
832 diverse DNA and RNA viruses. *Microbiome* **9**, 37, doi:10.1186/s40168-020-
833 00990-y (2021).
- 834 43 Nayfach, S. *et al.* CheckV assesses the quality and completeness of metagenome-
835 assembled viral genomes. *Nat Biotechnol* **39**, 578-585, doi:10.1038/s41587-020-
836 00774-7 (2021).
- 837 44 Kanehisa, M., Sato, Y. & Morishima, K. BlastKOALA and GhostKOALA:
838 KEGG Tools for Functional Characterization of Genome and Metagenome
839 Sequences. *J Mol Biol* **428**, 726-731, doi:10.1016/j.jmb.2015.11.006 (2016).
- 840 45 Yarza, P. *et al.* Uniting the classification of cultured and uncultured bacteria and
841 archaea using 16S rRNA gene sequences. *Nat Rev Microbiol* **12**, 635-645,
842 doi:10.1038/nrmicro3330 (2014).
- 843 46 Farris, M. H. & Olson, J. B. Detection of Actinobacteria cultivated from
844 environmental samples reveals bias in universal primers. *Lett Appl Microbiol* **45**,
845 376-381, doi:10.1111/j.1472-765X.2007.02198.x (2007).
- 846 47 Palkova, L. *et al.* Evaluation of 16S rRNA primer sets for characterisation of
847 microbiota in paediatric patients with autism spectrum disorder. *Sci Rep* **11**, 6781,
848 doi:10.1038/s41598-021-86378-w (2021).
- 849 48 Anda, M. *et al.* Bacterial clade with the ribosomal RNA operon on a small plasmid
850 rather than the chromosome. *Proc. Natl. Acad. Sci. U. S. A.* **112**, 14343-14347,
851 doi:10.1073/pnas.1514326112 (2015).

- 852 49 El Karkouri, K., Pontarotti, P., Raoult, D. & Fournier, P. E. Origin and Evolution
853 of Rickettsial Plasmids. *PLoS One* **11**, e0147492,
854 doi:10.1371/journal.pone.0147492 (2016).
- 855 50 Albertsen, M. *et al.* Genome sequences of rare, uncultured bacteria obtained by
856 differential coverage binning of multiple metagenomes. *Nat Biotechnol* **31**, 533-
857 538, doi:10.1038/nbt.2579 (2013).
- 858 51 Kantor, R. S. *et al.* Small genomes and sparse metabolisms of sediment-associated
859 bacteria from four candidate phyla. *mBio* **4**, e00708-00713,
860 doi:10.1128/mBio.00708-13 (2013).
- 861 52 He, X. *et al.* Cultivation of a human-associated TM7 phylotype reveals a reduced
862 genome and epibiotic parasitic lifestyle. *Proc. Natl. Acad. Sci. U. S. A.* **112**, 244-
863 249, doi:10.1073/pnas.1419038112 (2015).
- 864 53 Batinovic, S., Rose, J. J. A., Ratcliffe, J., Seviour, R. J. & Petrovski, S.
865 Cocultivation of an ultrasmall environmental parasitic bacterium with lytic ability
866 against bacteria associated with wastewater foams. *Nat Microbiol* **6**, 703-711,
867 doi:10.1038/s41564-021-00892-1 (2021).
- 868 54 Starr, E. P. *et al.* Stable isotope informed genome-resolved metagenomics reveals
869 that Saccharibacteria utilize microbially-processed plant-derived carbon.
870 *Microbiome* **6**, 122, doi:10.1186/s40168-018-0499-z (2018).
- 871 55 McLean, J. S. *et al.* Acquisition and Adaptation of Ultra-small Parasitic Reduced
872 Genome Bacteria to Mammalian Hosts. *Cell Rep* **32**, 107939,
873 doi:10.1016/j.celrep.2020.107939 (2020).
- 874 56 Toju, H., Okayasu, K. & Notaguchi, M. Leaf-associated microbiomes of grafted
875 tomato plants. *Sci Rep* **9**, 1787, doi:10.1038/s41598-018-38344-2 (2019).
- 876 57 Tanaka, E., Koitabashi, M. & Kitamoto, H. A teleomorph of the ustilaginalean
877 yeast *Moesziomyces antarcticus* on barnyardgrass in Japan provides bioresources
878 that degrade biodegradable plastics. *Antonie Van Leeuwenhoek* **112**, 599-614,
879 doi:10.1007/s10482-018-1190-x (2019).
- 880 58 Magnani, G. S. *et al.* Diversity of endophytic bacteria in Brazilian sugarcane.
881 *Genet Mol Res* **9**, 250-258, doi:10.4238/vol9-1gmr703 (2010).
- 882 59 West, E. R., Cother, E. J., Steel, C. C. & Ash, G. J. The characterization and
883 diversity of bacterial endophytes of grapevine. *Can J Microbiol* **56**, 209-216,
884 doi:10.1139/w10-004 (2010).
- 885 60 Pereira, S. I. & Castro, P. M. Diversity and characterization of culturable bacterial
886 endophytes from *Zea mays* and their potential as plant growth-promoting agents

- 887 in metal-degraded soils. *Environ Sci Pollut Res Int* **21**, 14110-14123,
888 doi:10.1007/s11356-014-3309-6 (2014).
- 889 61 Kooner, A. & Soby, S. Draft Genome Sequence of *Curtobacterium* sp. Strain
890 MWU13-2055, Isolated from a Wild Cranberry Fruit Surface in Massachusetts,
891 USA. *Microbiol Resour Announc* **11**, e0056522, doi:10.1128/mra.00565-22
892 (2022).
- 893 62 Matulich, K. L. *et al.* Temporal variation overshadows the response of leaf litter
894 microbial communities to simulated global change. *ISME J* **9**, 2477-2489,
895 doi:10.1038/ismej.2015.58 (2015).
- 896 63 Vorholt, J. A. Microbial life in the phyllosphere. *Nat Rev Microbiol* **10**, 828-840,
897 doi:10.1038/nrmicro2910 (2012).
- 898 64 Michaux, S. *et al.* Presence of two independent chromosomes in the *Brucella*
899 *melitensis* 16M genome. *J Bacteriol* **175**, 701-705, doi:10.1128/jb.175.3.701-
900 705.1993 (1993).
- 901 65 Yamaichi, Y., Iida, T., Park, K. S., Yamamoto, K. & Honda, T. Physical and
902 genetic map of the genome of *Vibrio parahaemolyticus*: presence of two
903 chromosomes in *Vibrio* species. *Mol Microbiol* **31**, 1513-1521,
904 doi:10.1046/j.1365-2958.1999.01296.x (1999).
- 905 66 Vishnivetskaya, T. A. *et al.* Draft genome sequences of 10 strains of the genus
906 *exiguobacterium*. *Genome Announc* **2**, doi:10.1128/genomeA.01058-14 (2014).
- 907 67 Nakai, R. *et al.* Genome sequence and overview of *Oligoflexus tunisiensis*
908 *Shr3(T)* in the eighth class *Oligoflexia* of the phylum *Proteobacteria*. *Stand*
909 *Genomic Sci* **11**, 90, doi:10.1186/s40793-016-0210-6 (2016).
- 910 68 Hahn, M. W. *et al.* *Silvanigrella aquatica* gen. nov., sp. nov., isolated from a
911 freshwater lake, description of *Silvanigrellaceae* fam. nov. and *Silvanigrellales*
912 ord. nov., reclassification of the order *Bdellovibrionales* in the class *Oligoflexia*,
913 reclassification of the families *Bacteriovoracaceae* and *Halobacteriovoraceae* in
914 the new order *Bacteriovoracales* ord. nov., and reclassification of the family
915 *Pseudobacteriovoracaceae* in the order *Oligoflexales*. *Int J Syst Evol Microbiol* **67**,
916 2555-2568, doi:10.1099/ijsem.0.001965 (2017).
- 917 69 Morella, N. M., Gomez, A. L., Wang, G., Leung, M. S. & Koskella, B. The impact
918 of bacteriophages on phyllosphere bacterial abundance and composition. *Mol*
919 *Ecol* **27**, 2025-2038, doi:10.1111/mec.14542 (2018).
- 920 70 Nelson, M. S., Chun, C. L. & Sadowsky, M. J. Type IV Effector Proteins Involved
921 in the *Medicago-Sinorhizobium* Symbiosis. *Mol Plant Microbe Interact* **30**, 28-
922 34, doi:10.1094/MPMI-10-16-0211-R (2017).

- 923 71 Cevallos, M. A., Cervantes-Rivera, R. & Gutierrez-Rios, R. M. The repABC
924 plasmid family. *Plasmid* **60**, 19-37, doi:10.1016/j.plasmid.2008.03.001 (2008).
- 925 72 Hosokawa, M., Nishikawa, Y., Kogawa, M. & Takeyama, H. Massively parallel
926 whole genome amplification for single-cell sequencing using droplet
927 microfluidics. *Sci Rep* **7**, 5199, doi:10.1038/s41598-017-05436-4 (2017).
- 928 73 Berger, B. R. & Christie, P. J. Genetic complementation analysis of the
929 *Agrobacterium tumefaciens* virB operon: virB2 through virB11 are essential
930 virulence genes. *J Bacteriol* **176**, 3646-3660, doi:10.1128/jb.176.12.3646-
931 3660.1994 (1994).
- 932 74 Grohmann, E., Keller, W. & Muth, G. Mechanisms of Conjugative Transfer and
933 Type IV Secretion-Mediated Effector Transport in Gram-Positive Bacteria. *Curr*
934 *Top Microbiol Immunol* **413**, 115-141, doi:10.1007/978-3-319-75241-9_5 (2017).
- 935 75 Christie, P. J. The Mosaic Type IV Secretion Systems. *EcoSal Plus* **7**,
936 doi:10.1128/ecosalplus.ESP-0020-2015 (2016).
- 937 76 Christie, P. J., Gomez Valero, L. & Buchrieser, C. Biological Diversity and
938 Evolution of Type IV Secretion Systems. *Curr Top Microbiol Immunol* **413**, 1-30,
939 doi:10.1007/978-3-319-75241-9_1 (2017).
- 940 77 Chetrit, D., Hu, B., Christie, P. J., Roy, C. R. & Liu, J. A unique cytoplasmic
941 ATPase complex defines the *Legionella pneumophila* type IV secretion channel.
942 *Nat Microbiol* **3**, 678-686, doi:10.1038/s41564-018-0165-z (2018).
- 943 78 Guglielmini, J. *et al.* Key components of the eight classes of type IV secretion
944 systems involved in bacterial conjugation or protein secretion. *Nucleic Acids Res*
945 **42**, 5715-5727, doi:10.1093/nar/gku194 (2014).
- 946 79 Wiegand, S., Jogler, M. & Jogler, C. On the maverick Planctomycetes. *FEMS*
947 *Microbiol Rev* **42**, 739-760, doi:10.1093/femsre/fuy029 (2018).
- 948 80 Morohoshi, T., Wang, W. Z., Someya, N. & Ikeda, T. Genome sequence of
949 *Microbacterium testaceum* StLB037, an N-acylhomoserine lactone-degrading
950 bacterium isolated from potato leaves. *J Bacteriol* **193**, 2072-2073,
951 doi:10.1128/JB.00180-11 (2011).
- 952 81 Baek, K. & Choi, A. Complete Genome Sequence of *Phreatobacter* sp. Strain
953 NMCR1094, a Formate-Utilizing Bacterium Isolated from a Freshwater Stream.
954 *Microbiol Resour Announc* **8**, doi:10.1128/MRA.00860-19 (2019).
- 955 82 Sanjenbam, P., Buddidathi, R., Venkatesan, R., Shivaprasad, P. V. & Agashe, D.
956 Phenotypic diversity of *Methylobacterium* associated with rice landraces in
957 North-East India. *PLoS One* **15**, e0228550, doi:10.1371/journal.pone.0228550
958 (2020).

- 959 83 Lunak, Z. R. & Noel, K. D. A quinol oxidase, encoded by *cyoABCD*, is utilized
960 to adapt to lower O₂ concentrations in *Rhizobium etli* CFN42. *Microbiology*
961 (*Reading*) **161**, 203-212, doi:10.1099/mic.0.083386-0 (2015).
- 962 84 Cotter, P. A., Chepuri, V., Gennis, R. B. & Gunsalus, R. P. Cytochrome o
963 (*cyoABCDE*) and d (*cydAB*) oxidase gene expression in *Escherichia coli* is
964 regulated by oxygen, pH, and the *fnr* gene product. *J Bacteriol* **172**, 6333-6338,
965 doi:10.1128/jb.172.11.6333-6338.1990 (1990).

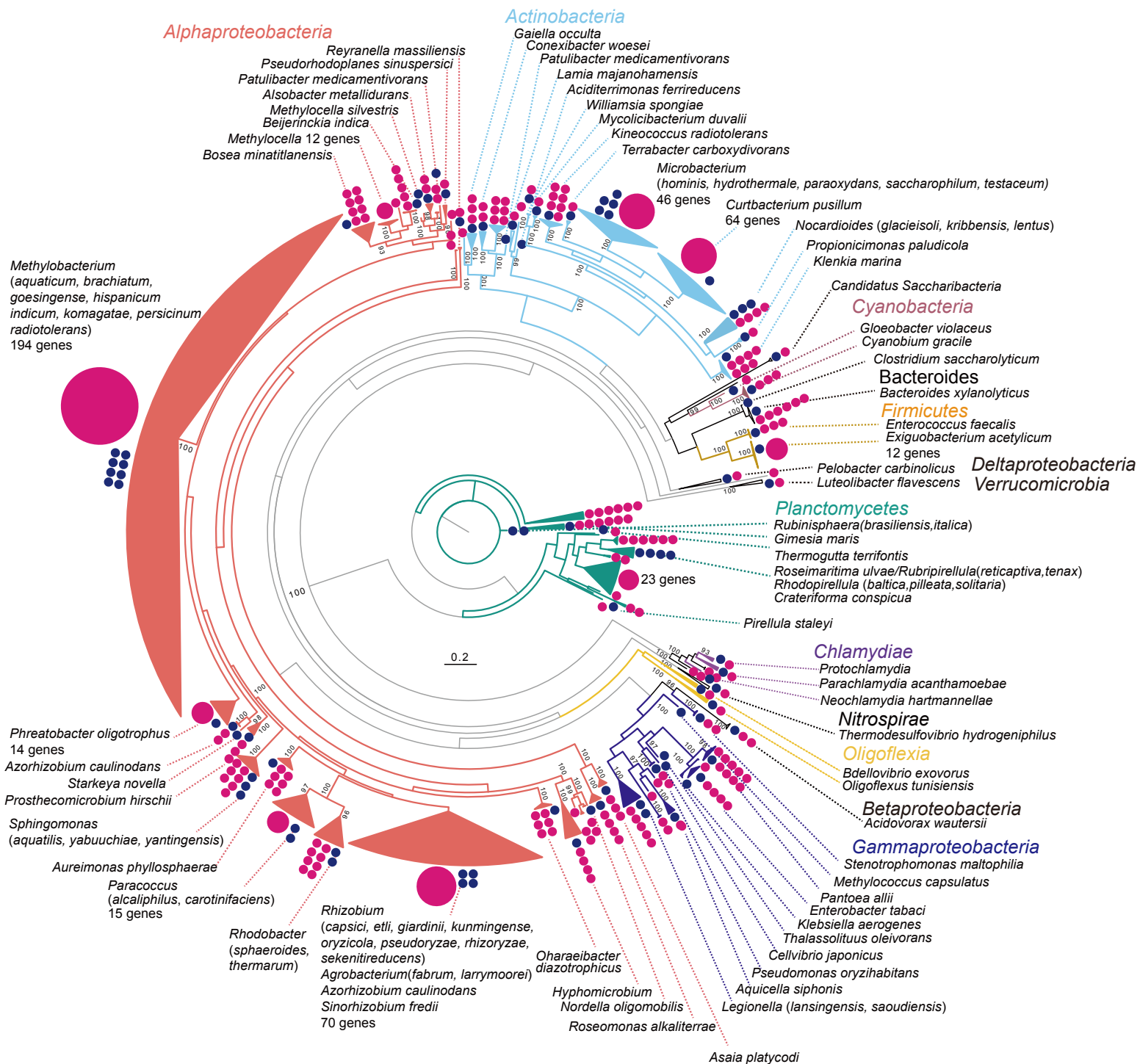


Fig. 1

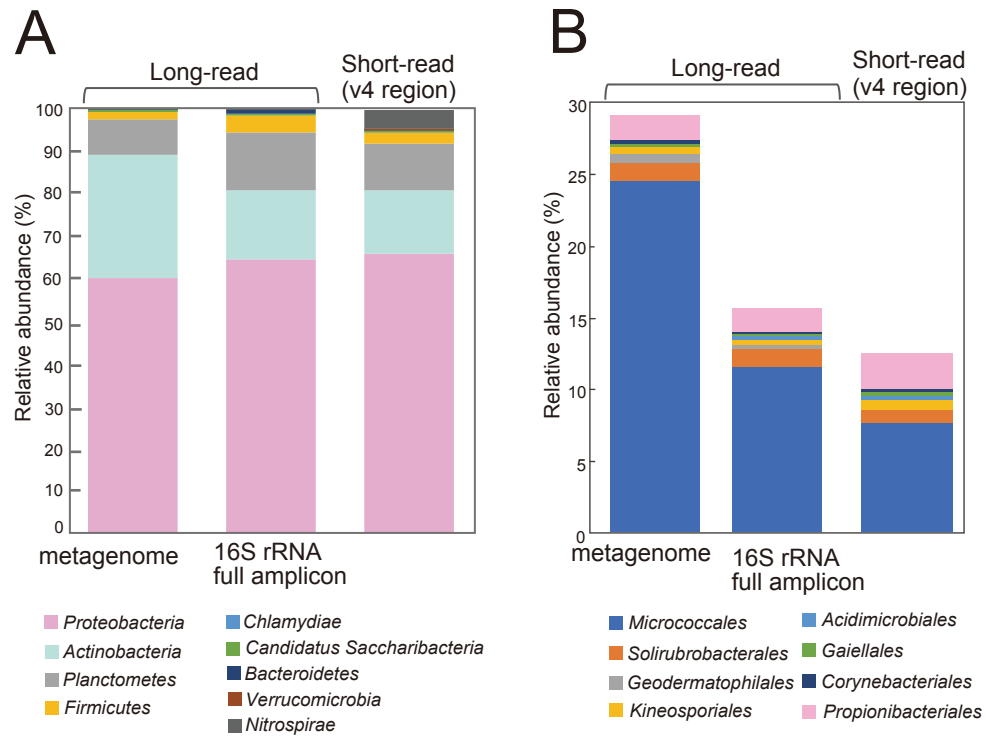


Fig. 2

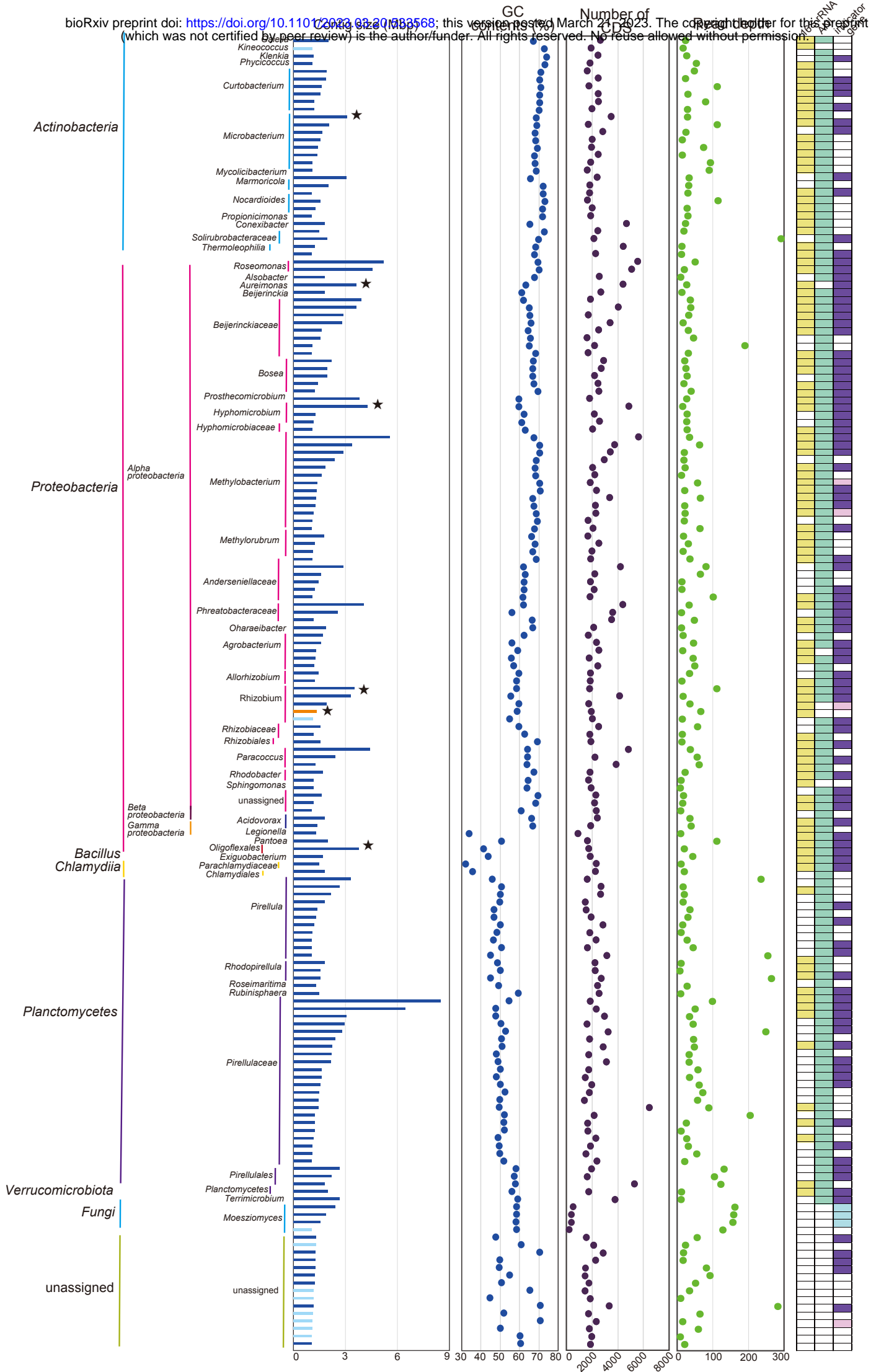


Fig. 3

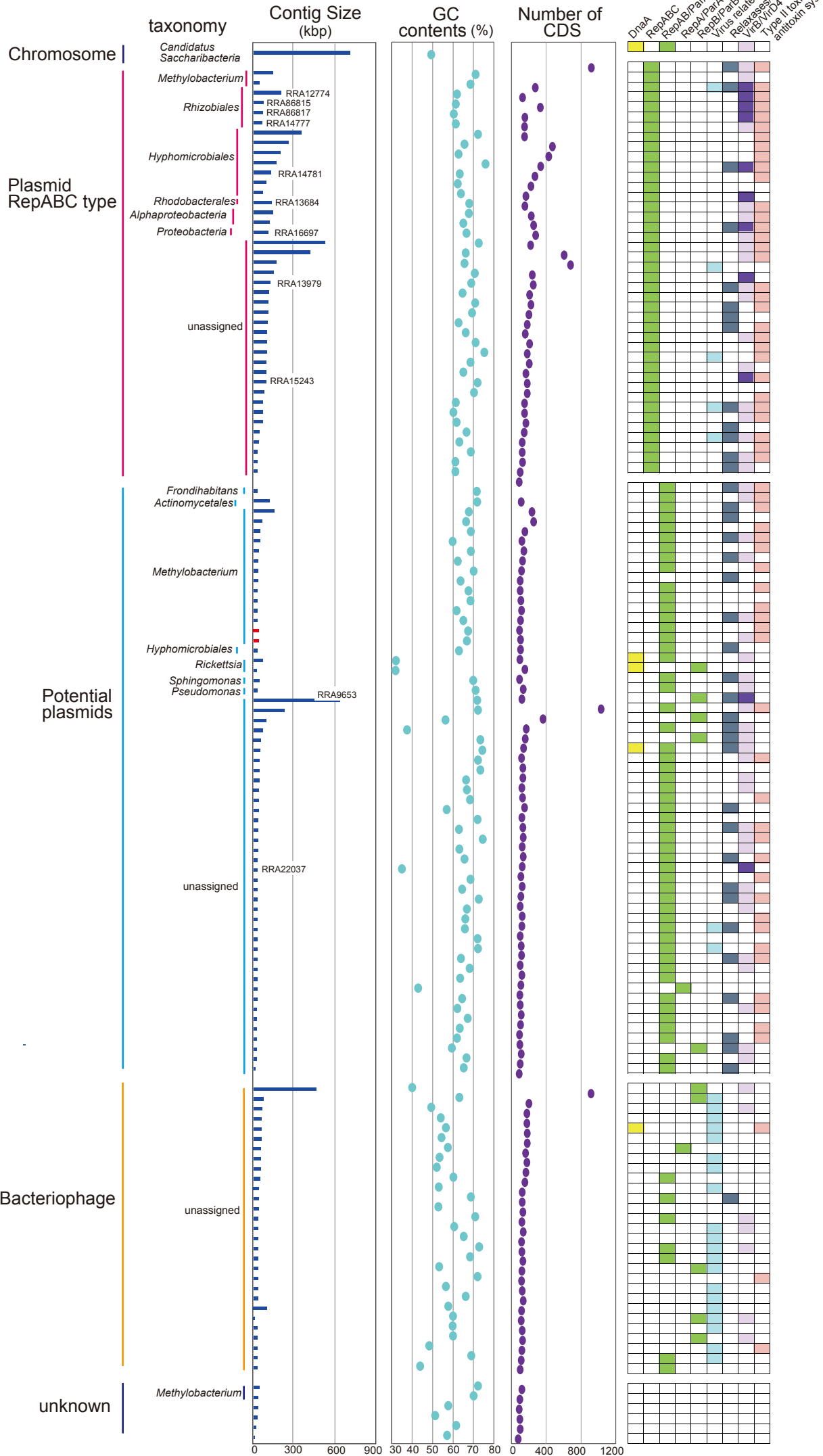


Fig. 4

Agrobacterium tumefaciens A6
pTiA6

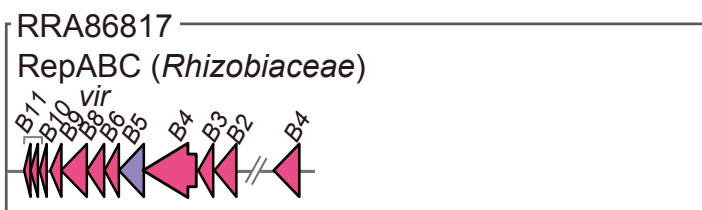
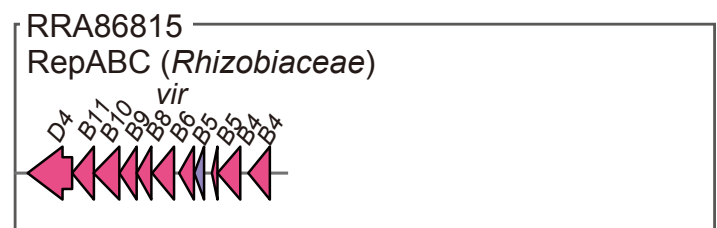
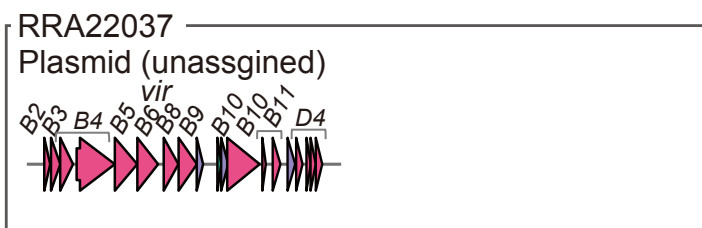
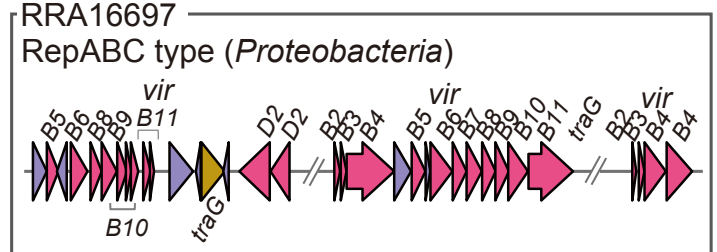
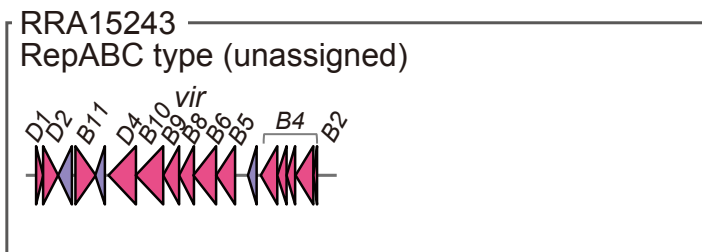
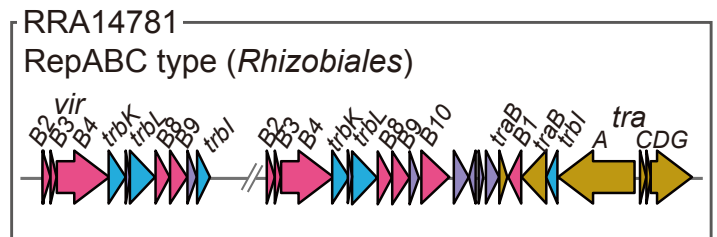
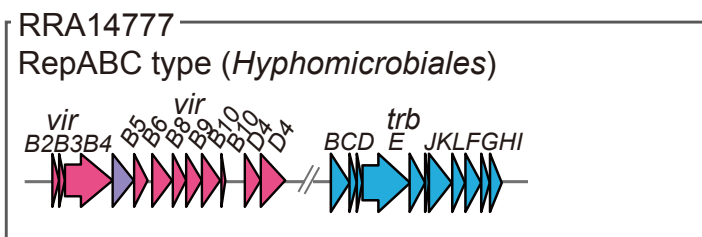
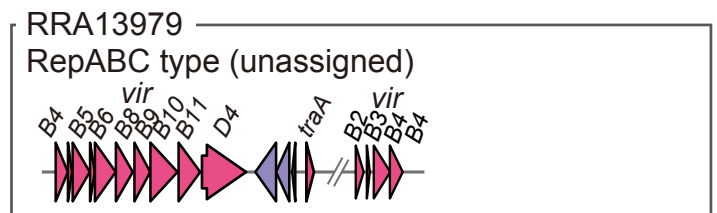
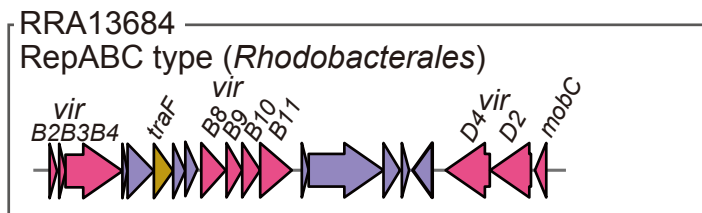
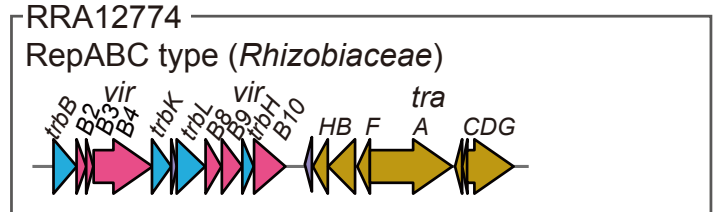
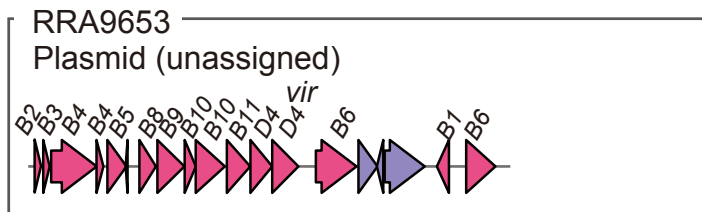
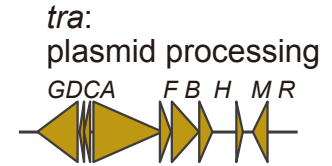
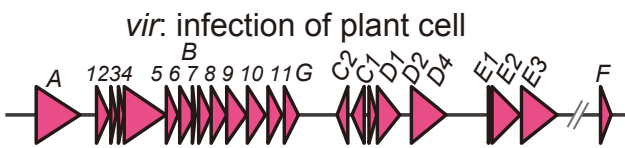


Fig. 5

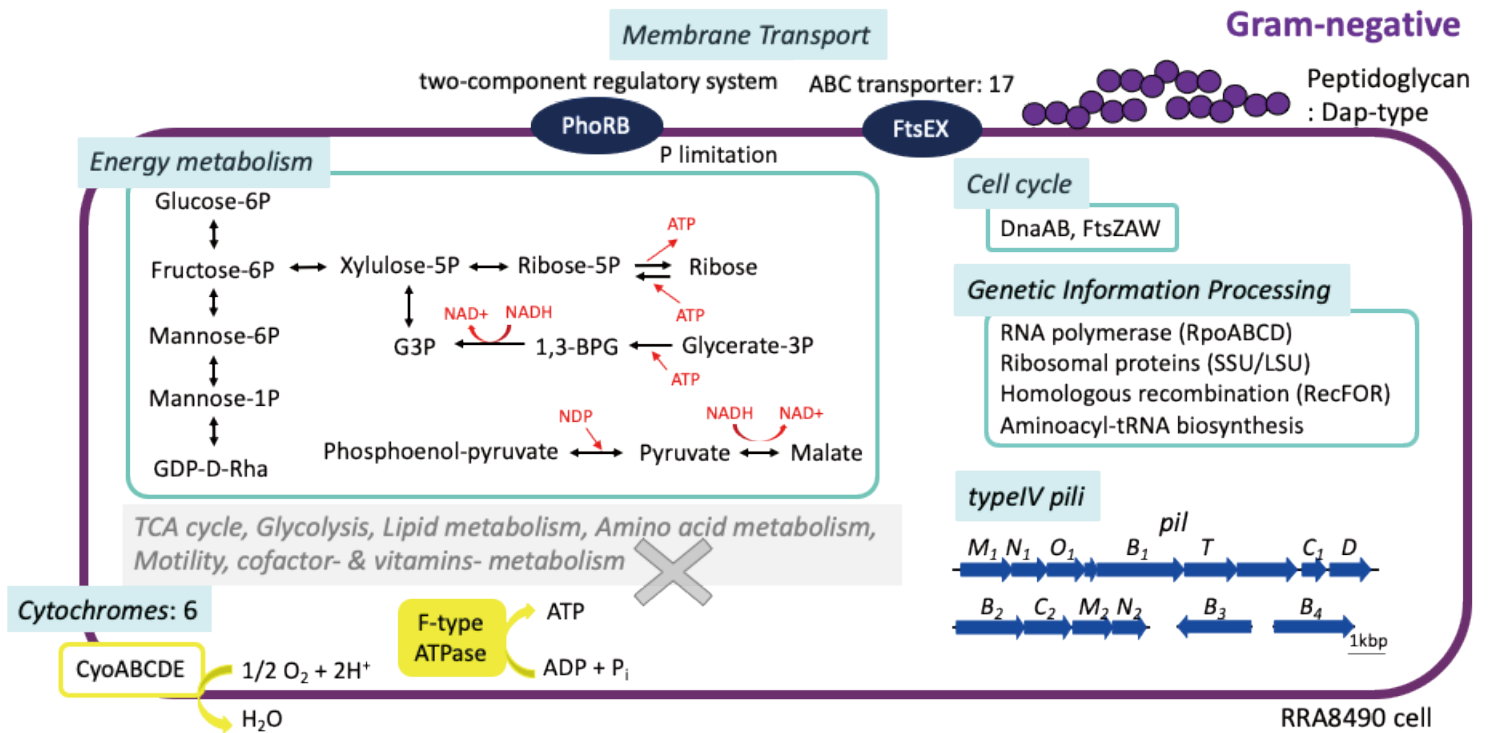


Fig.6