

1 Amplicon Sequencing of Single-copy Protein-coding Genes Reveals Accurate 2 Diversity for Sequence-discrete Microbiome Populations

3
4 Chengfeng Yang,^{a,b} Qinzhi Su,^{a,b} Min Tang,^b Shiqi Luo,^b Hao Zheng,^a Xue Zhang,^b Xin Zhou^b

5
6 ^aCollege of Food Science and Nutritional Engineering, China Agricultural University, Beijing, China

7 ^bDepartment of Entomology, College of Plant Protection, China Agricultural University, Beijing, China

8 Address correspondence to Xue Zhang, zhangxue05@cau.edu.cn or Xin Zhou, xinzhou@cau.edu.cn.

9
10
11 **ABSTRACT** An in-depth understanding of microbial function and the division of ecological niches
12 requires accurate delineation and identification of microbes at a fine taxonomic resolution.
13 Microbial phylotypes are typically defined using a 97% small subunit (16S) rRNA threshold.
14 However, increasing evidence has demonstrated the ubiquitous presence of taxonomic units of
15 distinct functions within phylotypes. These so-called sequence-discrete populations (SDPs) have
16 used to be mainly delineated by disjunct sequence similarity at the whole-genome level. However,
17 gene markers that could accurately identify and quantify SDPs are lacking in microbial community
18 studies. Here we developed a pipeline to screen single-copy protein-coding genes that could
19 accurately characterize SDP diversity via amplicon sequencing of microbial communities. Fifteen
20 candidate marker genes were evaluated using three criteria (extent of sequence divergence,
21 phylogenetic accuracy, and conservation of primer regions) and the selected genes were subject to
22 test the efficiency in differentiating SDPs within *Gilliamella*, a core honeybee gut microbial
23 phylotype, as a proof-of-concept. The results showed that the 16S V4 region failed to report
24 accurate SDP diversities due to low taxonomic resolution and changing copy numbers. In contrast,
25 the single-copy genes recommended by our pipeline were able to successfully quantify *Gilliamella*
26 SDPs for both mock samples and honeybee guts, with results highly consistent with those of
27 metagenomics. The pipeline developed in this study is expected to identify single-copy protein
28 coding genes capable of accurately quantifying diverse bacterial communities at the SDP level.

29
30 **IMPORTANCE** Microbial communities can be distinguished by discrete genetic and ecological
31 characteristics. These sequence-discrete populations are foundational for investigating the
32 composition and functional structures of microbial communities at high resolution. In this study, we
33 screened for reliable single-copy protein-coding marker genes to identify sequence-discrete
34 populations through our pipeline. Using marker gene amplicon sequencing, we could accurately and
35 efficiently delineate the population diversity in microbial communities. These results suggest that
36 single copy protein-coding genes can be an accurate, quantitative and economical alternative for
37 characterizing population diversity. Moreover, the feasibility of a gene as marker for any bacterial
38 population identification can be quickly evaluated by the pipeline proposed here.

39
40 **KEYWORDS** microbiota, SDP, quantification, 16S, metagenomics, *Gilliamella*

41 42 INTRODUCTION

43 Accurate identification of distinct functional units in natural bacterial communities is crucial in
44 understanding their ecological roles, interactions within the network, as well as the fine-scale
45 composition and dynamic changes within the whole community. As a rule of thumb, a bacterial
46 phylotype is often defined by grouping strains that share a sequence identify greater than 97% for a

47 selected fragment of the small subunit (16S) rRNA gene [1]. However, increasing evidence has
48 indicated that a bacterial phylotype may contain multiple finer lineages, each showing distinct
49 biological traits. For example, closely related enterotoxigenic *Escherichia coli* (ETEC) isolates
50 form discrete lineages with consistently definable variations in virulence profiles [2]. Such intra-
51 phylotype lineages could be delineated based on divergence in genomic sequences and phylogenetic
52 inferences. These finer subdivisions of phylotypes are called sequence-discrete populations (SDPs),
53 which typified by genetic and genealogical discontinuity from the rest of the community, and are
54 delineated by overall sequence divergence at the whole-genome level [3-5]. A broad comparison of
55 90,000 bacterial genomic sequences, with a close examination of pairwise genomic similarities in
56 natural bacterial communities, has proved the pervasive discontinuity in genetic similarity below
57 and above SDPs [3]. Bacteria in the same SDP normally show less than ca. 5% variation in whole-
58 genome sequences. This genetic divergence is much less than those among strains of the same
59 phylotype (ca. 30%) [6]. With respect to habitats, specific SDPs are likely ubiquitous in various
60 environments, such as human and animal guts [5, 7, 8], freshwater [9], ocean [10] and soil [11].
61 Therefore, SDPs are probably better than phylotypes, as taxonomic units that represent functional
62 entities in bacterial communities, which are likely shaped by ecological pressure and evolutionary
63 selection. As such, SDPs are important units of microbial diversity and should be considered as
64 baseline information for investing crucial questions, such as how do bacterial populations interact
65 and evolve within communities [4].

66 Despite the essential nature of accurate SDP identification, a rapid and accurate method that can
67 trace SDP boundaries is still lacking, especially with regards to the selection of proper markers for
68 evaluating sequence divergence. It is obvious that genetic divergence among bacterial strains is
69 dependent on which genes are compared. We now understand that the commonly used 16S gene
70 cannot generally provide sufficient resolution to characterize SDP diversity [12, 13]. For example,
71 in cases where the SDPs show a ~5-10% genome-wide divergence, they varied mostly merely < 0.1%
72 in the 16S sequences [14]. Moreover, the copy number of the 16S gene may vary significantly
73 among phylotypes or even among strains of the same phylotype, making quantitative
74 characterization of bacterial community a challenging, if not impossible, task [15, 16]. The 16S was
75 selected for phylotype delineation years ago because it has conserved primer sites that flank
76 relatively variable regions that made it easy to sequence with Sanger technology. Currently, much
77 effort has been put into developing genes or gene segments that can be easily sequenced, and that
78 vary enough to serve as practical proxies for SDP delineation [17-19]. However, a systematic
79 evaluation of the validity and performance of such genes in SDP delineation, which includes the
80 rapidly increasing but heterogeneously sampled database, has not been carried out.

81 Fortunately, recent developments in microbial genomics show a promising solution to
82 complement the coverage of bacterial genomes. The number of sequenced genomes of various
83 bacterial lineages has been growing rapidly. For example, the Genomes OnLine Database (GOLD)
84 now contains 437,099 bacterial genomes, the majority of which (397,945) are uncultured,
85 representing host-associated, environmental and engineered ecosystems [20]. The ever-growing
86 bacterial genome dataset offers a great opportunity to screen phylogenetically informative genes
87 that show good performance in taxonomic delineation, including those capable of quantitatively
88 characterizing bacterial communities at the SDP level [21, 22]. For instance, Wu and colleagues
89 identified 114 PhyEco universal markers for all bacteria [23]. From these universal markers, 15
90 single-copy protein-coding genes were successfully applied in estimating species abundances using
91 shotgun metagenomic data [24]. On the other hand, growing numbers of genomes and
92 metagenomes produced for particular bacterial communities or taxonomic groups allow for

93 comprehensive characterization of SDP diversity within focal environments and bacterial groups.
94 Taking social bee gut microbiota as an example, diverse strains derived from major honeybee hosts
95 have been isolated and deep-sequenced [25], including well-covered SDPs of nearly all core gut
96 bacterial phylotypes [5, 26, 27]. Thus, the relatively complete genome dataset provides a genome-
97 wide-based gold standard for defining SDPs for the honeybee core bacteria.

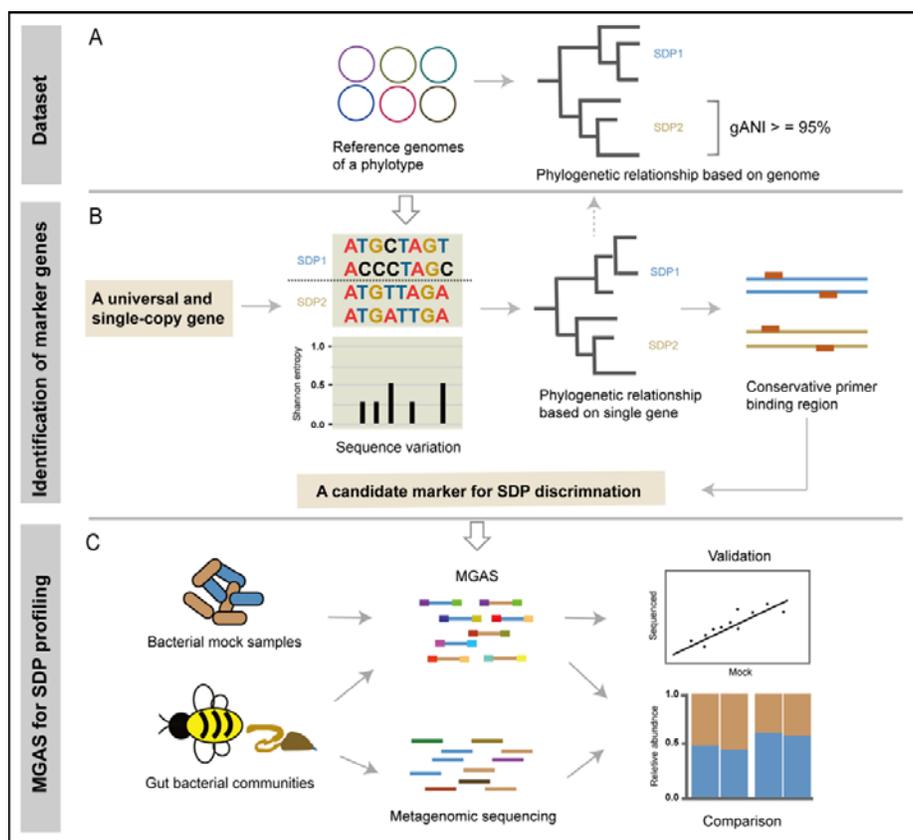
98 In the present study, we developed a pipeline to screen potential marker genes capable of
99 accurate identification and quantification of SDP diversity. We used the core bacterial phylotype
100 *Gilliamella* derived from the eastern honeybee *Apis cerana* as a proof of concept, and delineated
101 *Gilliamella* SDPs based on a set of comprehensive genome sequences. We further screened 15
102 single-copy protein-coding genes, which are present in all bacteria, to identify candidate marker
103 genes capable of differentiating the defined *Gilliamella* SDPs. Important characteristics such as the
104 level of sequence divergence, phylogenetic robustness, and the presence of conservative primer
105 regions, are considered in marker gene screening. Finally, we applied the candidate markers in
106 amplicon sequencing of both bacterial mock samples and real honeybee guts to verify their
107 efficiency in SDP profiling (Fig. 1). The markers we identified could accurately, consistently and
108 quantitatively capture SDP diversity.

109

110 RESULTS

111 **A comprehensive genome reference database for honeybee gut bacteria.** A comprehensive
112 genome reference database was constructed for honeybee gut bacteria (Table S1). A total of 242
113 genomes were included, covering 103 isolates from *A. cerana* and 139 from *A. mellifera*. SDPs
114 were identified for the core gut bacterial phylotypes using these reference genomes. SDPs differed
115 between honeybee species, which is consistent with previous studies [27, 28]. Within *A. cerana*
116 phylotypes, 5 SDPs were identified for *Gilliamella* (Gillia, n=65), 2 for *Bifidobacterium* (Bifido,
117 n=9), 1 for *Lactobacillus* Firm5 (Firm5, n=6), 1 for *Apibacter* (Apib, n=16) and 2 for *Snodgrassella*
118 (Snod, n=7). Within *A. mellifera* phylotypes, 6 SDPs were identified for Gillia (n=65), 9 for Bifido
119 (n=19), 2 for *Lactobacillus* Firm4 (Firm4, n=2), 6 for Firm5 (n=18) and 2 for Snod (n=35) (Table
120 S1). These SDPs delineated by genomes were used as references for subsequent taxonomic
121 assignments for the 16S, marker gene, or metagenome-based SDP identifications.

122



123
124
125
126
127
128
129
130
131

FIG 1 Screening marker genes suitable for SDP discrimination and quantification. (A) SDPs are identified for gut bacterial phylotypes based on phylogenetic relationships and genome-wide pairwise average nucleotide identities (gANI). (B) A candidate marker gene for SDP discrimination is selected from a set of universal and single-copy genes based on sequence variation, phylogenetic relationship and well-conserved regions for primer design. (C) The performance of marker gene amplicon sequencing (MGAS) on SDP identification and quantification is validated and compared as characterized using the mock samples and gut gut communities.

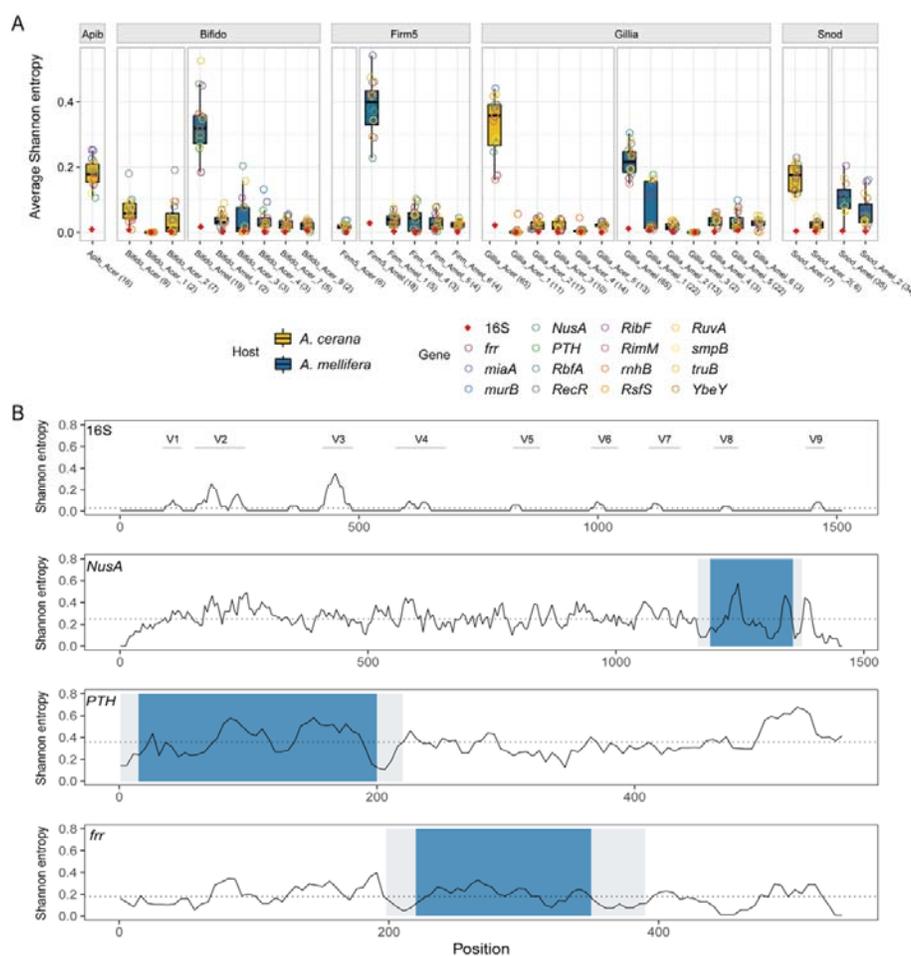
132 **Single-copy marker genes showed higher sequence variations at the SDP level than the**
133 **16S gene.** Sufficient sequence variation is crucial for high resolution discrimination of bacterial
134 SDPs. Here we compared the average Shannon entropy (ASE) between the whole-16S and the 15
135 single-copy marker genes. Our results clearly showed that the marker genes had much higher ASEs
136 at both phylotype and SDP levels compared to those of the 16S (Fig. 2A). The regional difference
137 in the variation levels between 16S and selected marker genes was also compared along the full
138 gene length. A slide-window (20 bp) ASE analysis showed that although several spikes of variable
139 regions were identified along the 16S gene, with the highest variable region corresponded to part of
140 the classic V3 region, its regional ASEs were generally lower compared to marker genes, e.g., *NusA*,
141 *PTH* and *frr* (Fig. 2B; Fig. S1).

142 Because phylogenetic placement of the query sequence is a critical step in our SDP
143 identification method, each marker gene will need to first produce a “correct” phylogeny for the
144 phylotype in question. Therefore, we further examined whether each of the 15 marker genes could
145 produce the same SDP phylogeny as inferred from whole-genome sequences of *Gilliamella*. Here,
146 the tree based on all 65 *A. cerana Gilliamella* genomes was used as the gold standard. The results
147 showed that all 15 marker genes but *rnhB* reconstructed the SDP phylogeny, with all strains
148 assigned to corresponding SDPs (Fig. S2). On the *rnhB* gene tree, two *Gilliamella* genomes were

149 misplaced from SDP Acer_Gillia_4 to Acer_Gillia_2, which was likely due to a higher sequence
 150 similarity between these two SDPs at a value of $90.93\% \pm 0.18$ SD comparing to that between other
 151 SDPs ($79.98\% \pm 1.89$ SD). Therefore, *rnhB* was subsequently excluded from further screening.

152 For the 14 remaining marker genes, we further explored for regions that were suitable for
 153 amplicon sequencing, based on the presence of conserved primer regions flanking the hyper-
 154 variable region. The *RimM* gene lacked hyper variable regions across the full gene length (Fig. S1),
 155 while some other genes (*murB*, *RecR*, *miaA*, *RbfA*, *RibF*, *RuvA*, *RsfS* and *YebY*) did not demonstrate
 156 promising conserved regions for primer design. These genes were then excluded from the candidate
 157 gene pool. The 5 remaining candidates (*frr*, *NusA*, *PTH*, *truB* and *smpB*) all had a hyper-variable
 158 region of ~200-550 bp that was flanked by conservative primer regions. Among them, *frr*, *NusA*
 159 and *PTH* produced an amplicon of ~200 bp (Fig. 2B), which could be thoroughly sequenced with
 160 most current shotgun sequencing methods (e.g., PE100 or PE150). These 3 genes were then chosen
 161 for the final test for their performance in SDP discrimination in both identity and quantity, using
 162 *Gilliamella* mock samples and real honeybee guts.

163



164
 165 **FIG 2** Marker genes are highly variable among SDPs. (A) Average Shannon entropy of the 15 marker genes and
 166 the 16S gene at both phylotype and SDP levels of honey bee gut bacteria. Numbers in brackets for each of the
 167 SDP groups indicate the number of strains examined for that specific group. (B) The Shannon entropy across 16S
 168 and candidate marker genes of all *A. cerana Gilliamella*. The Shannon entropy value is subsequently averaged by
 169 a 20-bp slide-window at a 5-bp step. Gray shadows depict conserved regions optimal for primer-binding sites and
 170 blue shadows are considered as hypervariable regions in this study. Dash lines represent the mean Shannon

171 entropy values cross all sequences. Gray lines depict the classic variable regions of the 16S gene. Apib: *Apibacter*;
172 Bifido: *Bifidobacterium*; Firm5: *Lactobacillus* Firm5; Gillia: *Gilliamella*; Snod: *Snodgrassella alvi*.

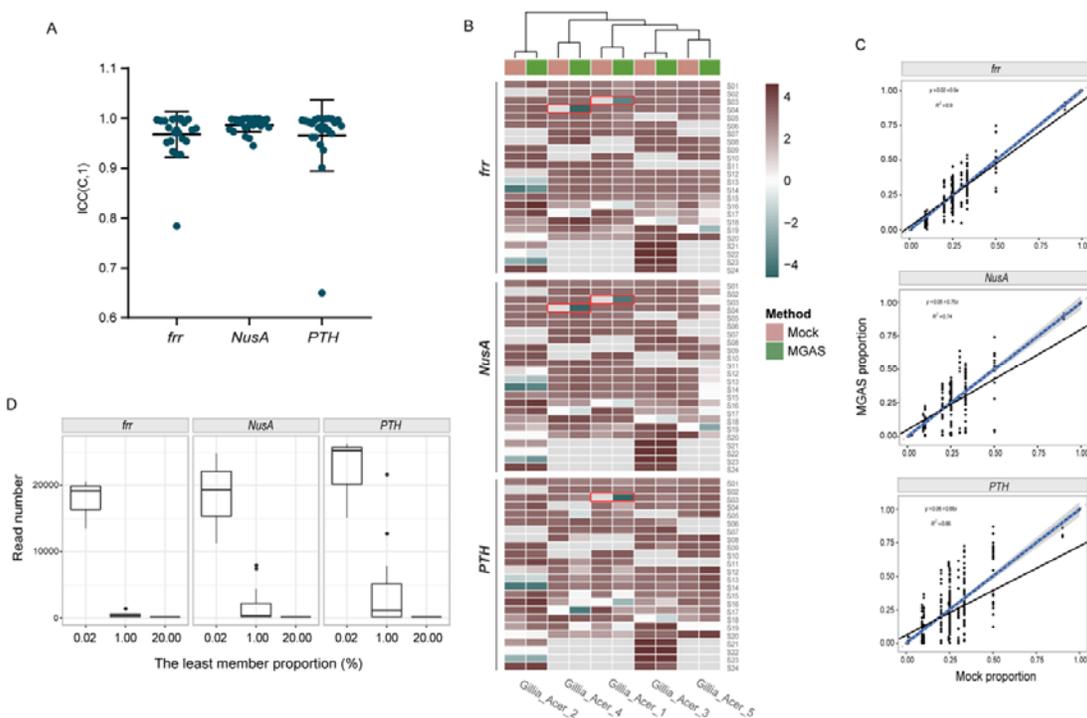
173

174 **Marker gene amplicon sequencing (MGAS) showed high accuracy, sensitivity and**
175 **repeatability in SDP profiling of mock samples.** Mock samples contained varied proportions of
176 the representative strain cultures of the 5 *Gilliamella* SDPs. These samples were extracted for DNA
177 and amplified for the hyper-variable regions of the 3 candidate marker genes (*frr*, *NusA* and *PTH*).
178 Twenty-four barcoded amplicons were pooled and shotgun sequenced for ca. 1 Gb data (ca. 2.5
179 million reads). Each mock sample was sequenced three times. An average of 73,462, 86,467 and
180 113,498 reads per sample was generated for *frr*, *NusA* and *PTH*, respectively.

181 The results of MGAS showed a high level of repeatability across the three replicates, where the
182 average ICC(C,1) > 0.9, except for *PTH*, which had an ICC(C,1) of 0.752 among samples with
183 equal proportion of bacterial DNA (Fig. 3A; Fig. S4C). With regards to detection accuracy, MGAS
184 correctly detected all bacterial members present in 22/24 samples, while two samples (S03 and S04)
185 showed false positive results, which was probably derived from sample contamination or
186 sequencing error (Fig. 3B). Because the sensitivity of amplicon sequencing was affected by
187 sequencing depth, we calculated the minimum read numbers required to detect members at low
188 abundances, using rarefaction curves (Fig. S5). The results suggested that strains with a relative
189 abundance of 1% could be detected by a minimum of ca. 1,123, 2,953 and 5,034 reads for *frr*, *NusA*
190 and *PTH* (equivalent to 0.49, 1.29 and 2.44 Mb data per sample), respectively. Accordingly, lower
191 abundance would require deeper sequencing. At a relative abundance of 0.02%, approximately
192 17,778, 18,518 and 22,222 reads (7.75, 8.07 and 10.76 Mb data) were required for *frr*, *NusA* and
193 *PTH*, respectively (Fig. 3D; Fig. S5). The sequencing depth was generally sufficient for SDP
194 detection in our study. Among the 216 sequenced samples, only two samples were sequenced with
195 only 963 (*frr*) and 2,348 (*PTH*) reads, respectively, and failed in identifying corresponding SDP
196 members at the lowest proportions (1% and 0.1%, respectively) due to insufficient sequencing
197 depth.

198 In addition to accurately identify *Gilliamella* SDPs, all three marker genes performed well in
199 quantifying relative abundances for mock samples. The relative abundances revealed by amplicon
200 reads were highly congruent with corresponding mock proportions in bacterial mock samples, with
201 the average R² values of 0.91, 0.74 and 0.66 for *frr*, *NusA* and *PTH*, respectively ($p < 2.2e-16$, Fig.
202 3C). The DNA mock samples yielded similar results, with the average R² values of 0.99, 0.91 and
203 0.99, for *frr*, *NusA* and *PTH*, respectively (Fig. S4B). Taken together, the MGAS method showed
204 high levels of accuracy, sensitivity and repeatability in characterizing SDP compositions, in both
205 taxonomic identity and relative abundance.

206

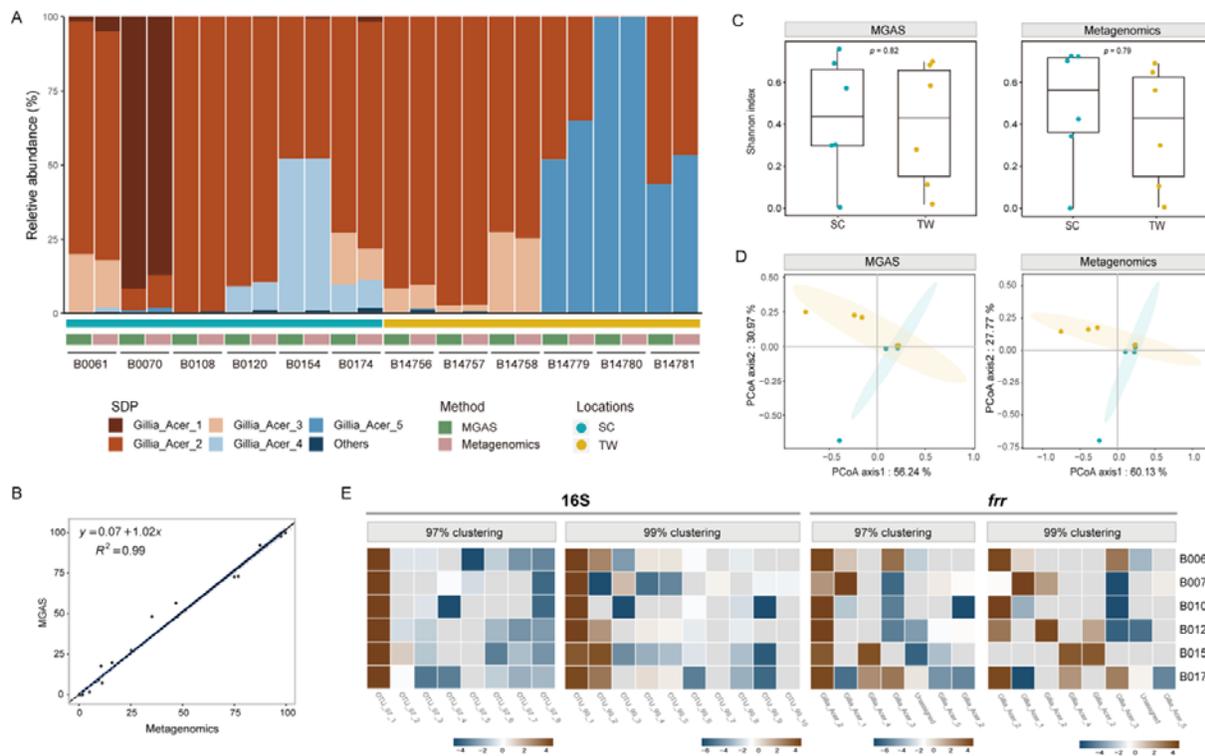


207
208
209
210
211
212
213
214
215
216
217
218
219
220

FIG 3 MGAS accurately identifies *A. cerana Gilliamella* SDPs. (A) Intra-class correlation coefficient (ICC) of relative abundance among the three replicates of MGAS samples. The ICC is calculated using the two-way mixed effects model with consistency (C) as the relationship among replicates, and single (1) result as the unit of measurement, i.e., ICC(C, 1). (B) Relative SDP abundances in mock samples revealed by marker gene sequencing. The results shown in the heatmap are the logarithms of the relative abundances of the five representative strains of the five SDPs of *A. cerana Gilliamella*. Grey box indicates a relative abundance at zero. False positive results are framed in red. (C) Spearman correlation of SDP abundances in *A. cerana Gilliamella* communities revealed by sequencing against mock samples. $p < 2.2e-16$. The black line presents the linear regression of the MGAS results against SDP abundances in mock samples. The blue solid and gray dashed lines represent a 1: 1 line and the fitted exponential regression (with 95 % confidence interval shown in gray shade), respectively. (D) Minimum read numbers required for detecting members at low abundances.

221 **MGAS performed equally well as metagenomics in characterizing honeybee gut SDP**
222 **diversity.** To examine the performance of the MGAS method in characterizing honeybee gut
223 microbiota, we used *frr* (Fig. 4) and *PTH* (Fig. S6) genes to calculate *Gilliamella* SDP diversities
224 for the 12 *A. cerana* workers from Sichuan and Taiwan, China. The MGAS was able to assign
225 strains to the correct SDP at accurate abundance for real gut samples, with results were highly
226 congruent with those from metagenomic sequencing (with $R^2 = 0.99$ for *frr* and 0.97 for *PTH*, $p <$
227 $2.2e-16$, Fig. 4B; Fig. S6B). Both results revealed that most individual bees were dominated by two
228 or three *Gilliamella* SDPs, yet with significant variations in dominant members and compositions
229 among individuals and across geographical locations (Fig. 4A). *Gillia_Acer_2* was the dominant
230 SDP in most of the sequenced bees, which was found in 11 out of the 12 samples, with 10 bearing
231 relative abundances of 48.06 - 98.37% (Fig. 4A). Both methods showed congruent results in alpha
232 diversity ($p = 0.82$ and 0.79 for MGAS and metagenomics sequencing, respectively, Wilcoxon
233 rank-sum test, Fig. 4C). At the beta diversity level, the principal coordinate analysis (PCoA) based
234 on Bray-Curtis dissimilarity revealed that the gut bacterial communities from bees of Sichuan and
235 Taiwan formed two distinct clusters, which separated along the first axis (Fig. 4D). This result was
236 again consistent between the MGAS and metagenomic methods (Adonis PERMANOVA, $R^2 =$

237 0.056, $p = 0.204$ for MGAS and $R^2 = 0.096$, $p = 0.134$ for metagenomics). Thus, the performance of
 238 SDP profiling using MGAS was parallel to the metagenomic gold standard in microbial community
 239 studies.
 240



241
 242
 243 **FIG 4** MGAS shows high congruence to metagenomic sequencing at SDP-level analysis. (A) Relative
 244 abundances of *Gilliamella* SDPs revealed by MGAS (*frf*) and metagenomics sequencing of *A. cerana* gut
 245 communities. (B) Spearman correlation coefficient between MGAS and metagenomics results, with $R^2 = 0.99$, $p <$
 246 $2.2e-16$. The black line presents the linear regression of the MGAS results in SDP abundances against those of
 247 metagenomics. The blue solid and gray dashed lines represent a 1: 1 line and the fitted exponential regression
 248 (with 95 % confidence interval shown in gray shade), respectively. (C) Shannon diversity index of SDP
 249 frequencies for bee guts from two locations calculated by MGAS (left panel) and metagenomic sequencing (right
 250 panel). The two methods showed no significant difference, with the p -value of 0.70 and 0.82 in SC and TW,
 251 respectively, by Wilcoxon rank-sum test. (D) Principal coordinate analysis (PCoA) based on Bray-Curtis
 252 dissimilarity of SDP compositions of honey bee workers from Sichuan and Taiwan using MGAS (left panel,
 253 Adonis PERMANOVA, $R^2 = 0.056$, $p = 0.204$) and metagenomic sequencing (right panel, Adonis
 254 PERMANOVA, $R^2 = 0.096$, $p = 0.134$). Each point represents the value for an individual bee and the color
 255 represent the location (Sichuan or Taiwan) of each bee. The shaded ellipses represent 95% confidence intervals on
 256 the ordination. (E) Relative abundances of *Gilliamella* OTUs in the gut microbiota of *A. cerana* assigned by
 257 clustering at 97% or 99% thresholds for 16S V4 and *frf*. The result shown in the heatmap are the logarithms of the
 258 relative abundances of the OTUs or five SDPs. Individual bees are marked to right of each row. Grey box
 259 indicates a relative abundance at zero.
 260

261 The 16S V4 region was also used to determine the *Gilliamella* SDP compositions for the 6 bee
 262 gut samples from Sichuan. We applied operational taxonomic unit (OTU) clustering based on
 263 sequence similarity at 97% and 99% identity thresholds, which are commonly adopted for
 264 surveying phylotype and intra-phylotype microbial diversities, respectively [12, 29], to assess the
 265 efficacy of 16S in SDP profiling. 16S amplicon sequencing resulted in 8 and 10 OTUs at 97% and
 266 99% thresholds, respectively, with a frequency cut off at > 100 . The identified OTU numbers

267 differed from those of the MGAS results at the same sequence similarity thresholds (Fig. 4E).
268 Alarming, 16S amplicons failed to assign OTUs to the correct SDPs via blast. And the relative
269 OTU proportions revealed by 16S disagreed with those from MGAS, where the numbers of
270 dominant OTUs (> 1%) revealed by MGAS were more congruent to those from metagenomics. The
271 improved performance with the MGAS method in characterizing SDP diversity is likely due to
272 greater sequence divergence of the marker genes. For instance, the average pairwise inter-SDPs
273 sequence similarity in the *frr* hyper-variable region was significantly lower ($90.92\% \pm 3.18$, $n = 65$)
274 than that of the 16S rRNA gene V4 region ($99.95\% \pm 0.65$, $n = 44$) (Wilcoxon rank-sum test, $p <$
275 $2e-16$).

276

277 SUMMARY AND DISCUSSION

278 We developed a pipeline to identify reliable marker genes for accurate identification and
279 quantification of SDPs from bacterial communities. Three important criteria were applied in the
280 assessment: the extent of sequence divergence, phylogenetic accuracy, and the presence of flanking
281 conservative primer regions. Single-copy protein-coding genes identified by our pipeline were
282 applied as marker genes in SDP quantification of honeybee gut microbiota, successfully producing
283 results consistent with those from metagenomics, which were used as the gold standard. Conversely,
284 we showed that the widely used 16S contained limited sequence divergence within phylotypes,
285 failing to provide sufficient resolution in differentiating SDPs. As a result, 16S V4 amplicon
286 sequencing cannot reflect fine scale bacterial diversity for the community. Consequently, dominant
287 OTUs delineated by 16S at 97% or 99% thresholds significantly differed from the defined SDPs.
288 On the other hand, the OTUs of single-copy protein-coding genes screened out by our pipeline were
289 successfully assigned to the correct SDPs, and the numbers of dominant OTUs showed more
290 congruent results to those from metagenomics.

291 Compared with whole-genome shotgun sequencing, amplicon sequencing of single-copy
292 protein-coding genes provides an alternative solution to characterize SDP diversity in an accurate,
293 quantitative and economical way. We address that not every single copy protein-coding gene is
294 efficacious in SDP quantification. The candidate gene must meet all three criteria integrated in our
295 pipeline to be a good marker gene. For a phylotype that is well represented by genomes of various
296 lineages, all single-copy genes, including protein-coding genes, can be evaluated by our pipeline. In
297 this case, we expect dozens to hundreds of proper marker genes to be filtered out. On the other hand,
298 a small set of core single-copy protein-copy genes that are determined to be universally present
299 among known bacteria, such as the 15 marker genes tested in this study, will likely provide
300 candidate genes suitable for accurate characterization of SDP diversity for less known bacterial taxa.

301 Accurate identification of the SDP composition will also facilitate the prediction of the
302 functional capacity of microbial communities. Functional attributes of a given bacterial lineage are
303 strongly correlated to its phylogenetic position [30]. Therefore, various approaches, e.g., PICRUTs
304 [31], have been developed to predict potential functions of a given microbial community based on
305 phylogenetic profiles of bacterial members. However, 16S sequences are employed in most current
306 programs for phylogenetic reconstruction. As demonstrated in this study, single-copy protein-
307 coding genes identified by our pipeline show better fidelity in revealing phylogenetic relationships
308 for the focal phylotype. Therefore, we anticipate that function prediction for microbial communities
309 will be further improved by integrating single-copy protein-coding genes and the screening pipeline
310 described here.

311

312 MATERIALS AND METHODS

313 **Genome references of core gut bacteria of honeybees.** A total of 242 bacterial genomes
314 associated with *A. mellifera* and *A. cerana* were downloaded from the NCBI genome database
315 (Table S1). These 242 genomes were used as the reference database of honeybee gut bacteria,
316 which comprised the 6 major phylotypes: *Apibacter* (n=16), *Bifidobacterium* (n=28), *Lactobacillus*
317 Firm4 (n=2), *Lactobacillus* Firm5 (n=24), *Gilliamella* (n=130) and *Snodgrassella* (n=42).

318 **SDP delineation for honeybee core phylotypes.** Protein-coding genes of all sequenced
319 genomes were annotated using Prokka (<https://github.com/tseemann/prokka>) [32]. Core genes,
320 which were defined as being shared by > 99% strains of a given phylotype, were identified using
321 Roary (version 3.13.0) [33] with the parameter -blastp 75. Multiple sequence alignments were
322 carried out using MAFFT (version v7.467, [https://github.com/The-Bioinformatics-](https://github.com/The-Bioinformatics-Group/Albiorix/wiki/mafft)
323 [Group/Albiorix/wiki/mafft](https://github.com/The-Bioinformatics-Group/Albiorix/wiki/mafft)) [34]. Phylogenetic trees were constructed using core single-copy genes
324 of each phylotype by RAxML (version 8.2.12, -x 12345 -N 1000 -p 12345 -f a -m GTRGAMMA)
325 [35]. Phylogenies were visualized in R (version 3.6.0) using the package ggtree_v2.4.1 [36] or
326 iTOL (version 6.1.1) [37]. Pairwise genome-wide average nucleotide identity (gANI) values were
327 calculated using pyani (version 0.2.10; <https://github.com/widdowquinn/pyani>) [38]. A clade with a
328 gANI \geq 95% from its closest clade was defined as an SDP.

329 **Screening for candidate marker genes capable of discriminating *Gilliamella* SDPs.** The
330 fifteen universal single-copy maker genes (*frr*, *NusA*, *PTH*, *RbfA*, *RecR*, *rnhB*, *RibF*, *RimM*, *RsfS*,
331 *RuvA*, *smpB*, *truB*, *miaA*, *murB* and *YebY*, listed in Table S2) [24] were evaluated as candidate
332 genes. The sequences of candidate marker genes were retrieved by MIDAS (version 1.3.2) [24],
333 whereas the 16S genes were retrieved from the reference genomes using an in-house script. The
334 average Shannon entropy (ASE) of the full gene length was used to assess sequence variation
335 between strains of inter- and intra-SDPs for all phylotypes, where the Shannon entropy for each
336 nucleotide site across genomes in comparison was calculated using oligotyping (version 2.1) [39].

337 The phylotype *Gilliamella*, which contains the most genomes available for this study, was used
338 as a proof of concept to examine the efficacy of marker genes in SDP differentiation. For each
339 SDPs in phylotype *Gilliamella*, the Shannon entropy values were subsequently averaged for each
340 20-bp slide-window with a 5-bp step to evaluate the regional genetic divergence along the full
341 length of the marker genes. Pairwise sequence similarities were determined by Clustal Omega [40].

342 From the candidate genes, potential marker genes that may efficiently distinguish all known
343 SDPs of the *Gilliamella* phylotype were screened. The following criteria were followed: 1) the
344 marker genes should contain conservative regions flanking the hyper-variable region for designing
345 primers enabling recovery target phylotype; 2) the amplicon length is between ~150-550 bps; 3) the
346 amplified region is sufficiently variable to allow the discrimination of SDPs; and 4) the primers are
347 specific to the focal phylotype to avoid off-target amplifications. The aforementioned 15 marker
348 genes were subject to these criteria, and 5 of them (*frr*, *NusA*, *PTH*, *truB* and *smpB*) were selected as
349 potential markers for identifying SDPs of *A. cerana* *Gilliamella*. Among these, three genes (*frr*,
350 *NusA* and *PTH*) were subjected to further testing as a proof of concept, because their amplicon
351 lengths were 206, 206 and 230 bp, respectively, which were ideal for current shotgun sequencing
352 platforms. To increase the throughput and cost efficiency, 24 amplicons were pooled for one
353 sequencing run. The 5' end of both forward and reverse primers were tagged with 6-bp unique
354 barcode sequences (see Table S3) to distinguish positive and negative DNA strains, and to
355 differentiate samples.

356 **Bacterial mock samples.** One representative strain from each of the five *Gilliamella* SDPs
357 associated with *A. cerana* was cultured at 35°C and 5% CO₂ for 48 h, on heart infusion agar (HIA)

358 medium containing 5% sheep's blood [41]. To screen potential contaminations, the full-length 16S
359 gene was amplified for each bacterial culture using universal primers 27F and 1492R [41] and was
360 subject to Sanger sequencing. 16S sequences were checked against those of the reference strains for
361 identification, before strains were mixed for mock samples. Each *Gilliamella* culture was adjusted
362 to OD₆₀₀ = 0.5. Twenty-four mock SDP communities were prepared by mixing up 2-5 of the
363 representative strains at varied proportions. The compositions of the mock samples were set as:
364 equal proportion of each of the five strains, equal proportion of four strains with the absence of one
365 strain at a time, equal proportion of three strains with the absence of two randomly selected strains,
366 and a series of varied compositions with relative abundances ranging from ca. 0.02% to 50%. DNA
367 of the bacterial mixtures were extracted using a CTAB-based DNA extraction protocol followed by
368 recovery in 10 mM Tris-EDTA buffer (1×TE, pH 7.4) and quantified using the Qubit® DNA Assay
369 Kit on a Qubit® 3.0 Fluorometer (Life Technologies, CA, USA). Alternatively, genomic DNA of
370 each of the five representative strain cultures was extracted separately and the mixed at varied
371 compositions and proportions (see Table S4).

372 **SDP identification and quantification for mock samples using amplicon sequencing of the**
373 **three marker genes.** PCR amplification was performed for *frr* (*frr*-F 5'
374 GCTGAAGATGCAAGAAC and *frr*-R 5' GCATCACGACGAATATT), *NusA* (*NusA*-F 5'
375 CTTGAAATTGAAGAACT and *NusA*-R 5' GTACCTTGTTTCAGCTAA), and *PTH* (*PTH*-F 5'
376 AAAGTTATTGTAGG and *PTH*-R 5' CCACTTAAATTCATAAA) for each mock sample with
377 three replicates. Triplicate 50-μl reactions were carried out with 25 μl of 2 × Phanta Max Master
378 Mix (Vazyme Biotech, Nanjing, China), 2 μl (each) of 10 μM primer, 19 μl of ddH₂O, and 2 μl of
379 template DNA. The thermocycling profile consisted of an initial 3-min denaturation at 95 °C, 35
380 cycles of 15 s at 95 °C, 15 s at 52 °C for *NusA* and *frr* or at 42 °C for *PTH*, and 20 s at 72 °C and a
381 final 10-min extension step at 72 °C. After being visualized on 2% agarose gels, DNA was purified
382 using a gel extraction kit (Qiagen, Germany) and quantified using the Qubit® DNA Assay Kit on a
383 Qubit® 3.0 Fluorometer. Barcoded amplicons of up to 24 mock samples were pooled together and
384 subject to Illumina sequencing using a NovaSeq 6000 platform (PCR-free library, 150 PE) at
385 Novogene (Beijing, China). Approximately 1 Gb of raw data were obtained from each pooled
386 library (Table S5).

387 The program fastq-multx (version 1.3.1. <https://github.com/brwnj/fastq-multx>) was employed
388 to demultiplex sequencing reads based on barcode sequences. The 6-bp barcodes in reverse
389 sequences were trimmed using Seqtk (<https://github.com/lh3/seqtk>). The demultiplexed paired-end
390 reads were then analyzed in QIIME2 (version 2020.2. <https://qiime2.org>) [42]. A plugin DATA2
391 [43] was used to denoise reads and to group sequences into amplicon sequence variants (ASVs).
392 Individual ASVs were then taxonomically classified using blast (classify-consensus-blast) at a 97%
393 identity threshold (Fig. S3) against the 3 marker genes (*frr*, *NusA* and *PTH*) derived from the
394 customized bee gut bacterial dataset. The relative abundance of each SDP (RA_{SDP}) was calculated
395 as: $RA_{SDP} = (NR_{SDP}) / (NR_{Gillia}) * 100$, where NR_{SDP} represents the number of reads mapped to the
396 focal SDP and NR_{Gillia} represents the number of reads mapped to all *Gilliamella* SDPs. These
397 estimated abundances were then compared to those of the mock samples. The performance of SDP
398 profiling of the 3 marker genes was evaluated on the basis of accuracy, sensitivity and repeatability.
399 Intraclass correlation coefficient (ICC) with a two way random/mixed (ICC(C,1)) model was used
400 to assess the repeatability of this method using SPSS (version 20.1) [44].

401 Rarefaction curves were plotted using identified SDP numbers against read numbers, which
402 were used to infer the minimum read number required to detect strains at varied proportions. For
403 each sample, ASVs with a depth <100 were filtered out. Rarefaction was performed using QIIME2

404 with the plugin alpha-rarefaction and a sampling depth of 40,000 reads per sample and default
405 parameters. Minimum read numbers for identifying SDPs with relative abundances of 0.02%, 1%
406 and 20% were chosen manually.

407 **SDP identification and quantification for *A. cerana* gut microbiota using 16S, marker**
408 **genes, and metagenome sequencing.** Adult worker bees collected in Sichuan were used to
409 quantify *Gilliamella* SDP diversity using three different methods (16S V4 region amplicon
410 sequencing, MGAS and metagenomic sequencing). Bees were first cooled at 4 °C for 10 min. Then
411 the entire guts were dissected from the abdomen using sterile forceps and DNA was extracted using
412 a CTAB bead-beating protocol described previously [45].

413 Firstly, the 16S V4 region was amplified for six bee guts from Sichuan and sequenced using an
414 Illumina Hiseq X Ten platform (250-300 bp insert size, 250 PE) at BGI-Shenzhen (Shenzhen,
415 China). Raw reads obtained for each sample were summarized in Table S6. Data quality control
416 was performed using fastp (version 0.13.1, -q 20 -u 10 -w 16) [46]. The demultiplexed sequences
417 were denoised and grouped into ASVs using an open reference method VSEARCH [47] embedded
418 in QIIME 2. The taxonomic identification for ASVs was subsequently performed using the naive-
419 Bayesian classifier trained on the BGM-Db, a curated 16S reference database for the classification
420 of honeybee and bumblebee gut bacteria [48]. A feature table and ASVs consisting of filtered 16S
421 reads pertaining to *Gilliamella* was constructed. OTU clustering was performed at both 97% and 99%
422 identity thresholds, respectively, using VSEARCH with cluster-features-de-novo method.
423 Additionally, low-abundant OTUs comprising of <100 reads were removed. Taxonomic
424 assignments for OTUs were performed using blast against the BGM-Db with SDP-level taxonomy.
425 OTU composition heatmaps were generated based on relative abundances and visualized in R.

426 Secondly, for each sample, the marker genes *frr* and *PTH*, which demonstrated the best and
427 worst performances in accuracy and sensitivity, respectively, among the 3 marker genes, were
428 applied following the same pipeline used in the mock samples. ASVs of the six sample from
429 Sichuan were clustered into OTUs and filtered following the abovementioned 16S V4 pipeline.
430 Taxonomic assignments for OTUs were performed by blast against *frr* sequences derived from the
431 customized bee gut bacterial genome sequence database.

432 Finally, metagenome sequencing of four bee (B0108, B0120, B0154 and B0174) guts was
433 performed using an Illumina Hiseq X Ten platform (300-400 bp insert size, 150 PE) at BGI-
434 Shenzhen. Additional metagenomes of eight worker bee guts (BioProject PRJNA705951) were
435 download from NCBI (Table S6). The metagenome sequencing was used as the gold standard for
436 *Gilliamella* diversity distributed in the honeybee guts. Shotgun reads mapped to the *A. cerana*
437 genome (GCF_001442555.1) using BWA aln (version 0.7.16a-r1181, -n 1) [49] were identified as
438 host reads and subsequently excluded. We used the 'run_midias.py species' script in MIDAS with
439 default parameters to estimate the relative abundances of SDPs for each sample. Finally, the results
440 from MGAS were compared to those from metagenome sequencing to assess the performance of
441 the marker genes.

442 **Data availability.** Raw data from MGAS, 16S V4 amplicon and metagenomic sequencing have
443 been submitted to NCBI under BioProject PRJNA772085.

444

445 **SUPPLEMENTAL MATERIAL**

446 **FIG S1** The Shannon entropy across the remain marker genes of all *A. cerana Gilliamella*. The
447 Shannon entropy value is subsequently averaged by a 20-bp slide-window at a 5-bp step. Dash lines
448 represent the mean Shannon entropy values cross all sequences.

449 **FIG S2** All but *mhB* of the 15 marker genes produce five SDPs for *A. cerana Gilliamella*
450 phylotype in concert with the whole-genome result.

451 **FIG S3** Histograms of average nucleotide identity values of the 3 marker genes from comparisons
452 between strains belonging to the same SDPs (green) or different SDPs (red). Vertical black line
453 indicates the threshold for bacterial SDPs taxonomy for the present method.

454 **FIG S4** MGAS accurately identifies the *A. cerana Gilliamella* SDPs in DNA mock samples. (A)
455 Relative SDP abundances in mock samples revealed by MGAS. The results shown in the heatmap
456 are the logarithms of the relative abundances percentage of the five representative strains of the five
457 SDPs of *A. cerana Gilliamella*. Grey box indicates a relative abundance at zero. (B) Spearman
458 correlation of SDP abundances in *A. cerana Gilliamella* communities revealed by sequencing against
459 mock samples, $p < 2.2e-16$. The black line presents the linear regression of the MGAS results
460 against SDP abundances in mock samples. The blue solid and gray dashed lines represent a 1: 1 line
461 and the fitted exponential regression (with 95 % confidence interval shown in gray shade),
462 respectively. (C) Repeatability of relative abundance between replicates of DNA mock samples. $n =$
463 6, ICC(C,1) is 0.936, 0.974 and 0.752 for *frr*, *NusA* and *PTH* genes, respectively.

464 **FIG S5** Rarefaction curves of detected bacterial SDPs in bacterial mock samples reach the
465 saturation stage with increasing read numbers.

466 **FIG S6** Amplicon sequencing with the *PTH* gene showed high congruence to metagenomic
467 sequencing at SDP-level analyses. (A) Relative abundances of *Gilliamella* SDPs revealed by
468 MGAS (*PTH* gene) and metagenomics sequencing of *A. cerana* gut communities. (B) Spearman
469 correlation coefficient between MGAS and metagenomics results, with $R^2 = 0.97$, $p < 2.2e-16$. The
470 black line presents the linear regression of the MGAS results in SDP abundances against those of
471 metagenomics. The blue solid and gray dashed lines represent a 1: 1 line and the fitted exponential
472 regression (with 95 % confidence interval shown in gray shade), respectively.

473 **TABLE S1** Information of the reference genomes.

474 **TABLE S2** Information of the marker genes.

475 **TABLE S3** List of barcode sequences.

476 **TABLE S4** Mixing ratio of mock samples.

477 **TABLE S5** Statistics of data outputs.

478 **TABLE S6** Summary of read processing and data obtained from marker gene, 16S V4 amplicon
479 and metagenomic sequencing of honey bee guts.

480

481 **ACKNOWLEDGEMENTS**

482 This work was supported by the Program of Ministry of Science and Technology of China
483 (2018FY100403), National Natural Science Foundation of China (No. 31772493) and National
484 Natural Science Foundation of China (No. 32000346).

485 The authors declare no competing financial interests.

486 Xin Z. and Xue Z. designed, organized and coordinated the study. C.Y. conducted the screening
487 pipeline development, marker gene and 16S V4 amplicon sequencing analysis. Q.S. retrieved the
488 sequences of the 16S and single-copy protein-coding genes, and conducted reference-based
489 metagenome mapping. M.T. assisted in sequence variation analysis. S.L. conducted sample
490 collection and SDP identification. Xin Z., Xue Z., C.Y. and Hao Z. wrote the first drafts and all
491 authors contributed to and proofed the manuscript.

492

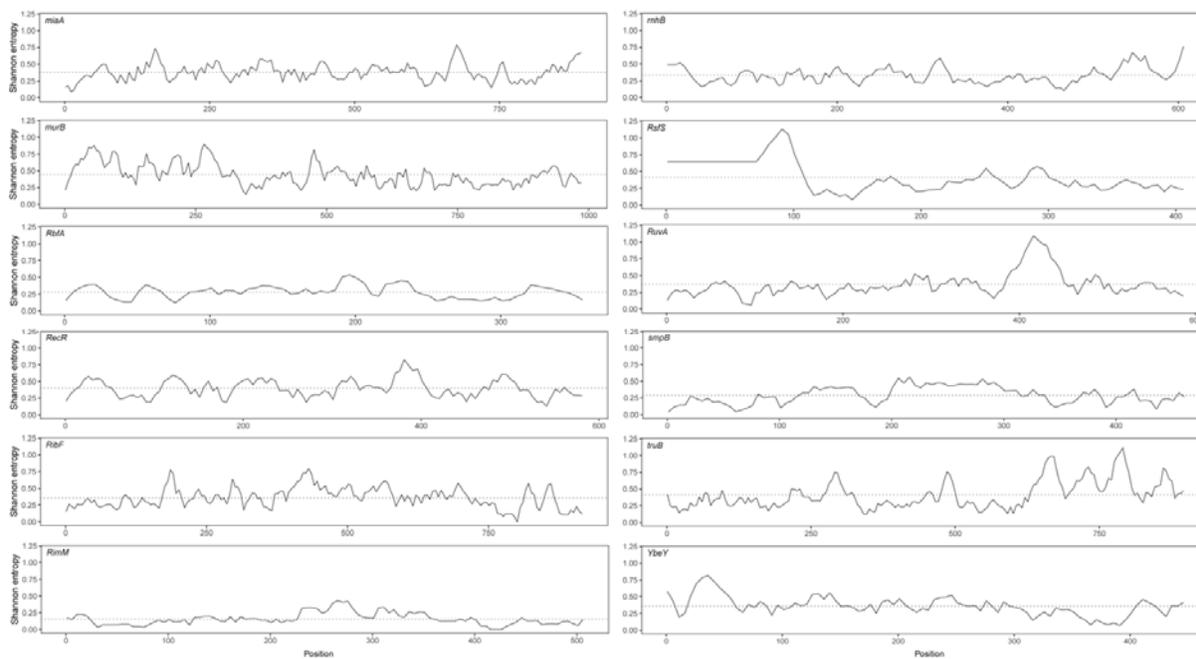
493 **REFERENCES**

- 494 1. Moreira D, López-García P. 2011. Phylotype, p 1254-1254. *In* Gargaud M, Amils R,
495 Quintanilla JC, Cleaves HJ, Irvine WM, Pinti DL, Viso M. Berlin, Heidelberg, Encyclopedia of
496 Astrobiology. Springer Berlin Heidelberg.
- 497 2. von Mentzer A, Connor TR, Wieler LH, Semmler T, Iguchi A, Thomson NR, Rasko DA, Joffre
498 E, Corander J, Pickard D, Wiklund G, Svennerholm A-M, Sjöling Å, Dougan G. 2014.
499 Identification of enterotoxigenic *Escherichia coli* (ETEC) clades with long-term global
500 distribution. *Nat Genet* 46:1321-1326.
- 501 3. Jain C, Rodriguez RL, Phillippy AM, Konstantinidis KT, Aluru S. 2018. High throughput ANI
502 analysis of 90K prokaryotic genomes reveals clear species boundaries. *Nat Commun* 9(1):5114.
- 503 4. Caro-Quintero A, Konstantinidis KT. 2012. Bacterial species may exist, metagenomics reveal.
504 *Environ Microbiol* 14:347-355.
- 505 5. Ellegaard KM, Engel P. 2019. Genomic diversity landscape of the honey bee gut microbiota.
506 *Nat Commun* 10:446.
- 507 6. Konstantinidis KT, Tiedje JM. 2005. Genomic insights that advance the species definition for
508 prokaryotes. *Proc Natl Acad Sci U S A* 102:2567-2572.
- 509 7. Fehlner-Peach H, Magnabosco C, Raghavan V, Scher JU, Tett A, Cox LM, Gottsegen C,
510 Watters A, Wiltshire-Gordon JD, Segata N, Bonneau R, Littman DR. 2019. Distinct
511 polysaccharide utilization profiles of human intestinal *Prevotella copri* isolates. *Cell Host*
512 *Microbe* 26:680-690 e685.
- 513 8. Tett A, Huang KD, Asnicar F, Fehlner-Peach H, Pasolli E, Karcher N, Armanini F, Manghi P,
514 Bonham K, Zolfo M, Filippis FD, Magnabosco C, Bonneau R, Lusingu J, Amuasi J, Reinhard
515 K, Rattei T, Boulund F, Engstrand L, Zink A, Collado MC, Littman DR, Eibach D, Ercolini D,
516 Rota-Stabelli O, Huttenhower C, Maixner F, and Segata N. 2019. The *Prevotella copri*
517 complex comprises four distinct clades underrepresented in westernized populations. *Cell Host*
518 *Microbe* 26:666-679 e667.
- 519 9. Oh S, Caro-Quintero A, Tsementzi D, DeLeon-Rodriguez N, Luo C, Poretsky R,
520 Konstantinidis KT. 2011. Metagenomic insights into the evolution, function, and complexity of
521 the planktonic microbial community of Lake Lanier, a temperate freshwater ecosystem. *Appl*
522 *Environ Microbiol* 77:6000-6011.
- 523 10. Konstantinidis KT, DeLong EF. 2008. Genomic patterns of recombination, clonal divergence
524 and environment in marine microbial populations. *ISME J* 2:1052-1065.
- 525 11. Olm MR, Crits-Christoph A, Diamond S, Lavy A, Matheus Carnevali PB, Banfield JF. 2020.
526 Consistent metagenome-derived metrics verify and delineate bacterial species boundaries.
527 *mSystems* 5:e00731-19.
- 528 12. Kim M, Oh HS, Park SC, Chun J. 2014. Towards a taxonomic coherence between average
529 nucleotide identity and 16S rRNA gene sequence similarity for species demarcation of
530 prokaryotes. *Int J Syst Evol Microbiol* 64:346-351.
- 531 13. Yarza P, Yilmaz P, Pruesse E, Glockner FO, Ludwig W, Schleifer KH, Whitman WB, Euzéby
532 J, Amann R, Rossello-Mora R. 2014. Uniting the classification of cultured and uncultured
533 bacteria and archaea using 16S rRNA gene sequences. *Nat Rev Microbiol* 12:635-645.
- 534 14. DeLong EF, Preston CM, Mincer T, Rich V, Hallam SJ, Frigaard N-U, Martinez A, Sullivan
535 MB, Edwards R, Brito BR, Chisholm SW, Karl DM. 2006. Community genomics among
536 stratified microbial assemblages in the ocean's interior. *Science* 311:486-503.
- 537 15. Sun DL, Jiang X, Wu QL, Zhou NY. 2013. Intragenomic heterogeneity of 16S rRNA genes
538 causes overestimation of prokaryotic diversity. *Appl Environ Microbiol* 79:5962-5969.

- 539 16. Mende DR, Sunagawa S, Zeller G, Bork P. 2013. Accurate and universal delineation of
540 prokaryotic species. *Nat Methods* 10:881-884.
- 541 17. Powell E, Ratnayeke N, Moran NA. 2016. Strain diversity and host specificity in a specialized
542 gut symbiont of honeybees and bumblebees. *Mol Ecol* 25:4461-4471.
- 543 18. Raymann K, Bobay LM, Moran NA. 2018. Antibiotics reduce genetic diversity of core species
544 in the honeybee gut microbiome. *Mol Ecol* 27:2057-2066.
- 545 19. Bobay LM, Wissel EF, Raymann K. 2020. Strain structure and dynamics revealed by targeted
546 deep sequencing of the honey bee gut microbiome. *mSphere* 5:e00694-20.
- 547 20. Mukherjee S, Stamatis D, Bertsch J, Ovchinnikova G, Sundaramurthi JC, Lee J, Kandimalla M,
548 Chen IA, Kyrpides NC, Reddy TBK. 2021. Genomes OnLine Database (GOLD) v.8: overview
549 and updates. *Nucleic Acids Res* 49:D723-D733.
- 550 21. Baldaufnl SL, Baldaufnl AJ, Wenk-Siefert I, Doolittle WF. 2000. A kingdom-level phylogeny
551 of Eukaryotes based on combined protein data. *Science* 290:972-977.
- 552 22. Konstantinidis KT, Tiedje JM. 2005. Towards a genome-based taxonomy for prokaryotes. *J*
553 *Bacteriol* 187:6258-6264.
- 554 23. Wu D, Jospin G, Eisen JA. 2013. Systematic identification of gene families for use as
555 "markers" for phylogenetic and phylogeny-driven ecological studies of bacteria and archaea
556 and their major subgroups. *PLoS One* 8:e77033.
- 557 24. Nayfach S, Rodriguez-Mueller B, Garud N, Pollard KS. 2016. An integrated metagenomics
558 pipeline for strain profiling reveals novel patterns of bacterial transmission and biogeography.
559 *Genome Res* 26:1612-1625.
- 560 25. Zheng H, Steele MI, Leonard SP, Motta EVS, Moran NA. 2018. Honey bees as models for gut
561 microbiota research. *Lab Anim (NY)* 47:317-325.
- 562 26. Voulgari-Kokota A, McFrederick QS, Steffan-Dewenter I, Keller A. 2019. Drivers, diversity,
563 and functions of the solitary-bee microbiota. *Trends Microbiol* 27:1034-1044.
- 564 27. Kwong WK, Medina LA, Koch H, Sing KW, Soh EJY, Ascher JS, Jaffé R, Moran NA. 2017.
565 Dynamic microbiome evolution in social bees. *Science Advances* 3:e1600513.
- 566 28. Ellegaard KM, Suenami S, Miyazaki R, Engel P. 2020. Vast differences in strain-level
567 diversity in the gut microbiota of two closely related honey bee species. *Curr Biol* 30:2520-
568 2531 e2527.
- 569 29. Edgar RC. 2018. Updating the 97% identity threshold for 16S ribosomal RNA OTUs.
570 *Bioinformatics* 34:2371-2375.
- 571 30. Martiny AC, Treseder K, Pusch G. 2013. Phylogenetic conservatism of functional traits in
572 microorganisms. *ISME J* 7:830-838.
- 573 31. Langille MG, Zaneveld J, Caporaso JG, McDonald D, Knights D, Reyes JA, Clemente JC,
574 Burkepille DE, Vega Thurber RL, Knight R, Beiko RG, Huttenhower C. 2013. Predictive
575 functional profiling of microbial communities using 16S rRNA marker gene sequences. *Nat*
576 *Biotechnol* 31:814-821.
- 577 32. Seemann T. 2014. Prokka: rapid prokaryotic genome annotation. *Bioinformatics* 30:2068-2069.
- 578 33. Page AJ, Cummins CA, Hunt M, Wong VK, Reuter S, Holden MT, Fookes M, Falush D,
579 Keane JA, Parkhill J. 2015. Roary: rapid large-scale prokaryote pan genome analysis.
580 *Bioinformatics* 31:3691-3693.
- 581 34. Katoh K, Standley DM. 2013. MAFFT multiple sequence alignment software version 7:
582 improvements in performance and usability. *Mol Biol Evol* 30:772-780.
- 583 35. Stamatakis A. 2014. RAxML version 8: a tool for phylogenetic analysis and post-analysis of
584 large phylogenies. *Bioinformatics* 30:1312-1313.

- 585 36. Yu G, Lam TT, Zhu H, Guan Y. 2018. Two Methods for mapping and visualizing associated
586 data on phylogeny using ggtree. *Mol Biol Evol* 35:3041-3043.
- 587 37. Letunic I, Bork P. 2011. Interactive Tree Of Life v2: online annotation and display of
588 phylogenetic trees made easy. *Nucleic Acids Res* 39:W475-W478.
- 589 38. Pritchard L, Glover RH, Humphris S, Elphinstone JG, Toth IK. 2016. Genomics and taxonomy
590 in diagnostics for food security: soft-rotting enterobacterial plant pathogens. *Analytical*
591 *Methods* 8:12-24.
- 592 39. Eren AM, Maignien L, Sul WJ, Murphy LG, Grim SL, Morrison HG, Sogin ML. 2013.
593 Oligotyping: differentiating between closely related microbial taxa using 16S rRNA gene data.
594 *Methods Ecol Evol* 4:1111-1119.
- 595 40. Madeira F, Park YM, Lee J, Buso N, Gur T, Madhusoodanan N, Basutkar P, Tivey ARN,
596 Potter SC, Finn RD, Lopez R. 2019. The EMBL-EBI search and sequence analysis tools APIs
597 in 2019. *Nucleic Acids Res* 47:W636-W641.
- 598 41. Kwong WK, Moran NA. 2013. Cultivation and characterization of the gut symbionts of honey
599 bees and bumble bees: description of *Snodgrassella alvi* gen. nov., sp. nov., a member of the
600 family *Neisseriaceae* of the *Betaproteobacteria*, and *Gilliamella apicola* gen. nov., sp. nov., a
601 member of *Orbaceae* fam. nov., *Orbales* ord. nov., a sister taxon to the order
602 '*Enterobacteriales*' of the *Gammaproteobacteria*. *Int J Syst Evol Microbiol* 63:2008-2018.
- 603 42. Bolyen E, Rideout JR, Dillon MR, Bokulich NA, Abnet CC, Al-Ghalith GA, Alexander H,
604 Alm EJ, Arumugam M, Asnicar F, Bai Y, Bisanz JE, Bittinger K, Brejnrod A, Brislawn CJ,
605 Brown CT, Callahan BJ, Caraballo-Rodríguez AM, Chase J, Cope EK, Da Silva R, Diener C,
606 Dorrestein PC, Douglas GM, Durall DM, Duvallet C, Edwardson CF, Ernst M, Estaki M,
607 Fouquier J, Gauglitz JM, Gibbons SM, Gibson DL, Gonzalez A, Gorlick K, Guo J, Hillmann B,
608 Holmes S, Holste H, Huttenhower C, Huttley GA, Janssen S, Jarmusch AK, Jiang L, Kaehler
609 BD, Kang KB, Keefe CR, Keim P, Kelley ST, Knights D, Koester I, Kosciulek T, Kreps J,
610 Langille MGI, Lee J, Ley R, Liu YX, Loftfield E, Lozupone C, Maher M, Marotz C, Martin
611 BD, McDonald D, McIver LJ, Melnik AV, Metcalf JL, Morgan SC, Morton JT, Naimey AT,
612 Navas-Molina JA, Nothias LF, Orchanian SB, Pearson T, Peoples SL, Petras D, Preuss ML,
613 Pruesse E, Rasmussen LB, Rivers A, Robeson MS 2nd, Rosenthal P, Segata N, Shaffer M,
614 Shiffer A, Sinha R, Song SJ, Spear JR, Swafford AD, Thompson LR, Torres PJ, Trinh P,
615 Tripathi A, Turnbaugh PJ, Ul-Hasan S, van der Hooft JJJ, Vargas F, Vázquez-Baeza Y,
616 Vogtmann E, von Hippel M, Walters W, Wan Y, Wang M, Warren J, Weber KC, Williamson
617 CHD, Willis AD, Xu ZZ, Zaneveld JR, Zhang Y, Zhu Q, Knight R, Caporaso JG. 2019.
618 Reproducible, interactive, scalable and extensible microbiome data science using QIIME 2. *Nat*
619 *Biotechnol* 37:852-857.
- 620 43. Callahan BJ, McMurdie PJ, Rosen MJ, Han AW, Johnson AJ, Holmes SP. 2016. DADA2:
621 High-resolution sample inference from Illumina amplicon data. *Nat Methods* 13:581-583.
- 622 44. Mirzajani A, Asharlous A, Kianpoor P, Jafarzadehpur E, Yekta A, Khabazkhoob M, Hashemi
623 H. 2019. Repeatability of curvature measurements in central and paracentral corneal areas of
624 keratoconus patients using Orbscan and Pentacam. *J Curr Ophthalmol* 31:382-386.
- 625 45. Soh EJY, Jaffé R, Ascher JS, Koch H, Medina LA, Moran NA, Kwong WK, Sing K-W. 2017.
626 Dynamic microbiome evolution. *Science Advances* 3:e1600513.
- 627 46. Chen S, Zhou Y, Chen Y, Gu J. 2018. fastp: an ultra-fast all-in-one FASTQ preprocessor.
628 *Bioinformatics* 34:i884-i890.

- 629 47. Rognes T, Flouri T, Nichols B, Quince C, Mahé F. 2016. VSEARCH: a versatile open source
630 tool for metagenomics. *PeerJ* 4:e2584.
- 631 48. Zhang X, Li X, Su Q, Cao Q, Li C, Niu Q, Zheng H. 2019. A curated 16S rRNA reference
632 database for the classification of honeybee and bumblebee gut microbiota. *Biodiversity Science*
633 27:557-566.
- 634 49. Li H, Durbin R. 2010. Fast and accurate long-read alignment with Burrows-Wheeler transform.
635 *Bioinformatics* 26(5):589-595.
636



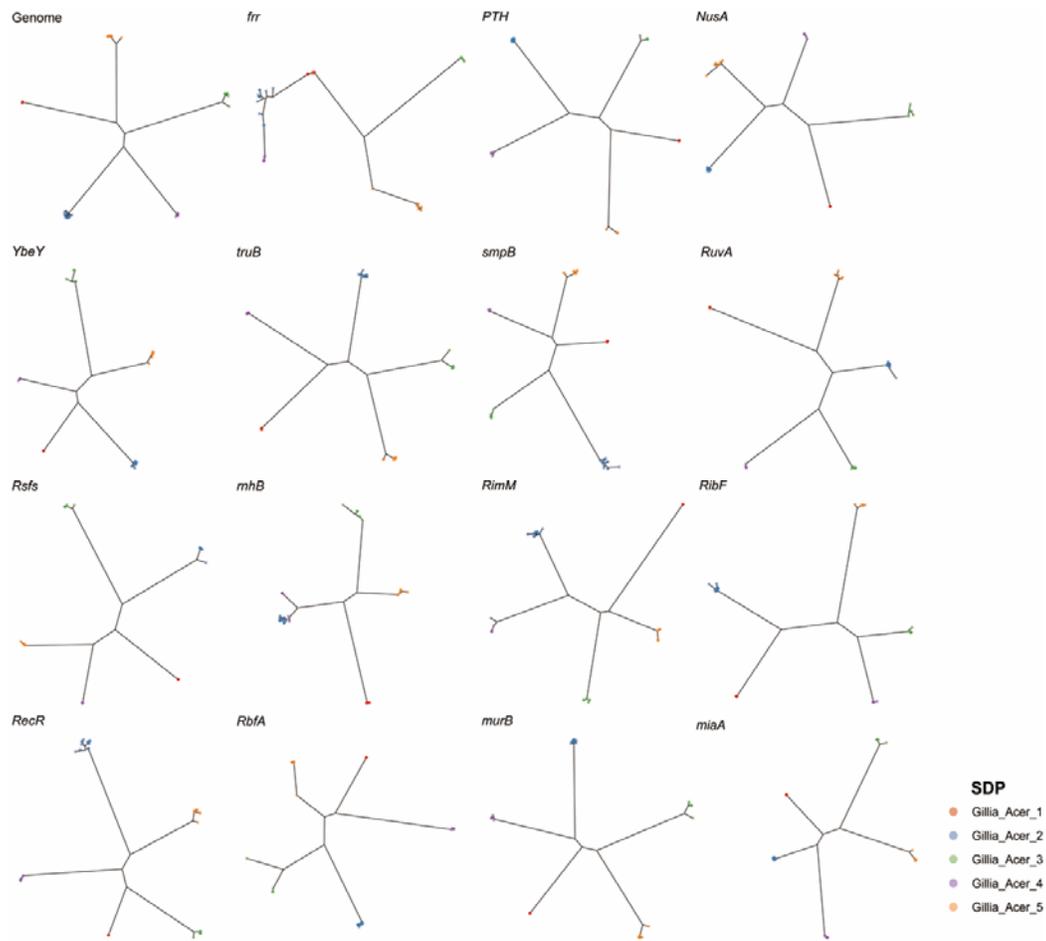
637

638

639 **FIG S1** The Shannon entropy across the remain marker genes of all *A. cerana Gilliamella*. The Shannon entropy
640 value is subsequently averaged by a 20-bp slide-window at a 5-bp step. Dash lines represent the mean Shannon
641 entropy values cross all sequences.

642

643



644

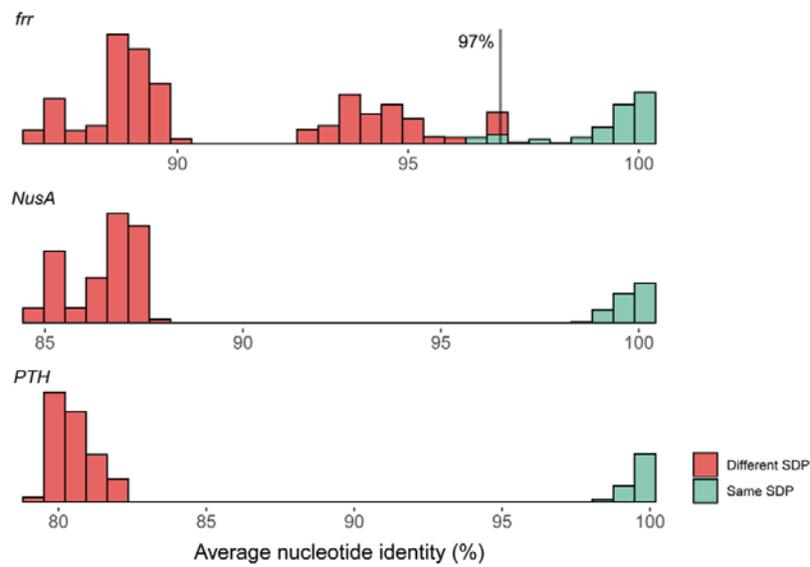
645

646

647

FIG S2 All but *mhB* of the 15 marker genes produce five SDPs for *A. cerana Gilliamella* phylotype in concert with the whole-genome result.

648



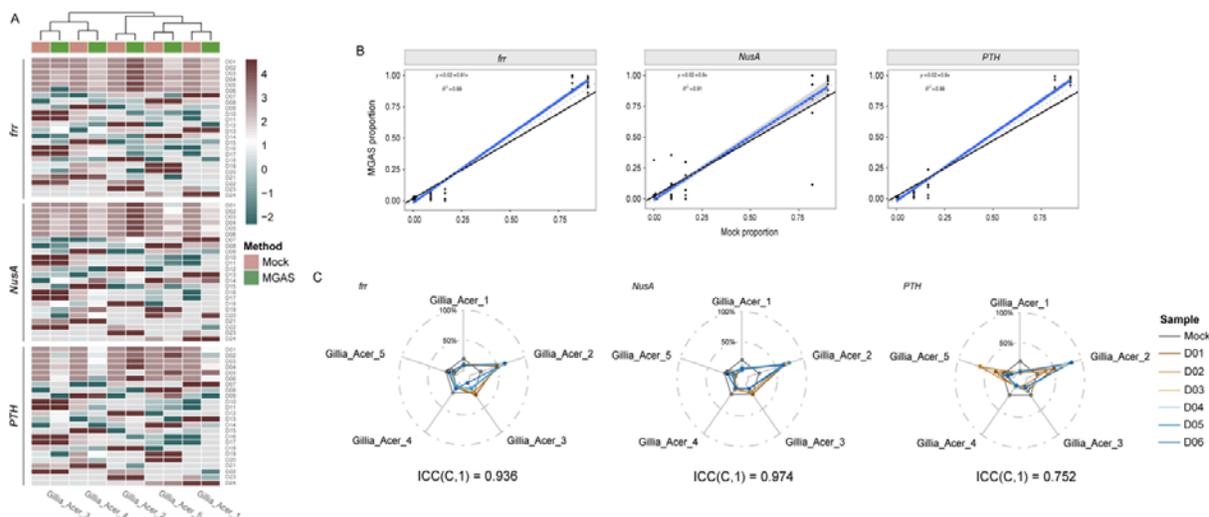
649

650

651 **FIG S3** Histograms of average nucleotide identity values of the 3 marker genes from comparisons between strains
652 belonging to the same SDPs (green) or different SDPs (red). Vertical black line indicates the threshold for
653 bacterial SDPs taxonomy for the present method.

654

655

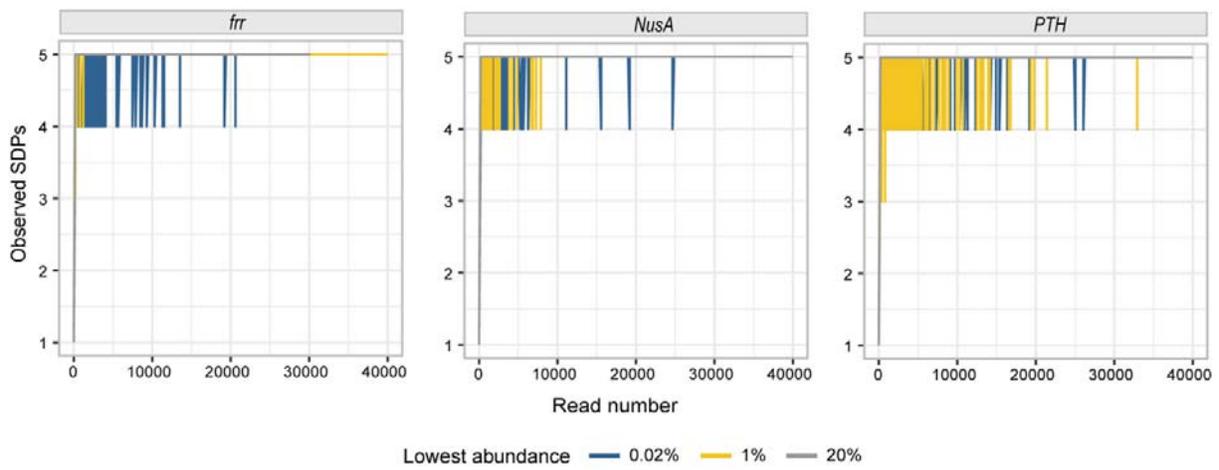


656

657

658 **FIG S4** MGAS accurately identifies the *A. cerana* *Gilliamella* SDPs in DNA mock samples. (A) Relative SDP
659 abundances in mock samples revealed by MGAS. The results shown in the heatmap are the logarithms of the
660 relative abundances percentage of the five representative strains of the five SDPs of *A. cerana* *Gilliamella*. Grey
661 box indicates a relative abundance at zero. (B) Spearman correlation of SDP abundances in *A. cerana* *Gilliamella*
662 communities revealed by sequencing against mock samples, $p < 2.2e-16$. The black line presents the linear
663 regression of the MGAS results against SDP abundances in mock samples. The blue solid and gray dashed lines
664 represent a 1: 1 line and the fitted exponential regression (with 95 % confidence interval shown in gray shade),
665 respectively. (C) Repeatability of relative abundance between replicates of DNA mock samples. $n = 6$, ICC(C,1)
666 is 0.936, 0.974 and 0.752 for *frr*, *NusA* and *PTH* genes, respectively.

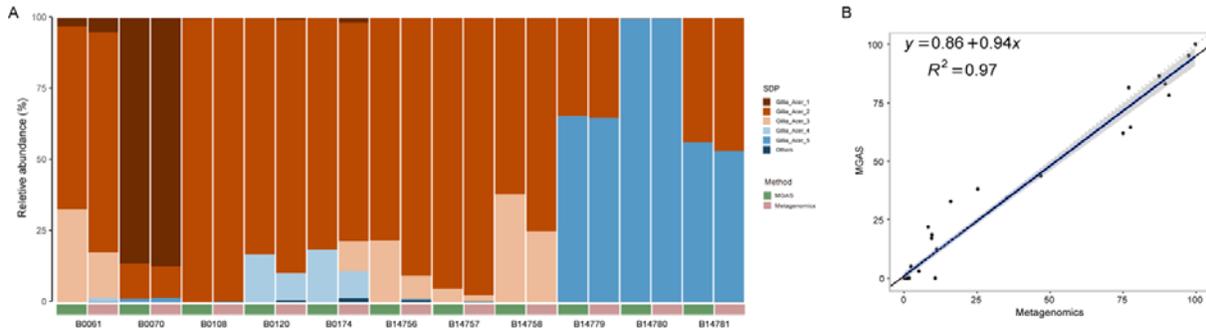
667



668
669
670
671
672

FIG S5 Rarefaction curves of detected bacterial SDPs in bacterial mock samples reach the saturation stage with increasing read numbers.

673



674

675

676

677

678

679

680

681

682

FIG S6 Amplicon sequencing with the *PTH* gene showed high congruence to metagenomic sequencing at SDP-level analyses. (A) Relative abundances of *Gilliamella* SDPs revealed by MGAS (*PTH* gene) and metagenomics sequencing of *A. cerana* gut communities. (B) Spearman correlation coefficient between MGAS and metagenomics results, with $R^2 = 0.97$, $p < 2.2e-16$. The black line presents the linear regression of the MGAS results in SDP abundances against those of metagenomics. The blue solid and gray dashed lines represent a 1: 1 line and the fitted exponential regression (with 95 % confidence interval shown in gray shade), respectively.

685 (Continued Table S1)

<i>Sordgrassella</i>	<i>Sordgrassella</i>	B3809	Stead_Accr_Advs_1	<i>Apis cerana</i>	China Caohu Town, Mengzi County, Henghe Prefecture, Yunnan	2,23	21	JANFH00000000
<i>Sordgrassella</i>	<i>Sordgrassella</i>	B3837	Stead_Accr_Advs_1	<i>Apis cerana</i>	China Changqing National Nature Reserve, Yang County, Hanchong, Shaanxi	2,04	31	JANFH00000000
<i>Sordgrassella</i>	<i>Sordgrassella</i>	wk8237A	Stead_Accr_Advs_1	<i>Apis andrewmyerensis</i>	Singapore Hort Park	2,32	21	MEJH00000000
<i>Sordgrassella</i>	<i>Sordgrassella</i>	wk8237B	Stead_Accr_Advs_1	<i>Apis cerana</i>	Singapore Clementi Park	2,32	31	MEJH00000000
<i>Sordgrassella</i>	<i>Sordgrassella</i>	wk8298B	Stead_Accr_Advs_1	<i>Apis cerana</i>	Singapore Hort Park	2,34	42	MEJH00000000
<i>Sordgrassella</i>	<i>Sordgrassella</i>	Q02	Stead_Accr_1	<i>Apis mellifera</i>	USA: West Haven, CT	1,60	299	JAEH00000000
<i>Sordgrassella</i>	<i>Sordgrassella</i>	A-10-12	Stead_Accr_2	<i>Apis mellifera</i>	USA: Arizona	2,50	63	NAGY01000000
<i>Sordgrassella</i>	<i>Sordgrassella</i>	A-11	Stead_Accr_2	<i>Apis mellifera</i>	USA: Arizona	2,43	122	NAGY00000000
<i>Sordgrassella</i>	<i>Sordgrassella</i>	A-11-12	Stead_Accr_2	<i>Apis mellifera</i>	USA: Arizona	2,50	90	NAGY01000000
<i>Sordgrassella</i>	<i>Sordgrassella</i>	A-11-12	Stead_Accr_2	<i>Apis mellifera</i>	USA: Arizona	2,50	58	NAGY00000000
<i>Sordgrassella</i>	<i>Sordgrassella</i>	A12	Stead_Accr_2	<i>Apis mellifera</i>	USA: Arizona	2,40	214	NAGY00000000
<i>Sordgrassella</i>	<i>Sordgrassella</i>	A2	Stead_Accr_2	<i>Apis mellifera</i>	USA: Arizona	2,43	84	NAGY01000000
<i>Sordgrassella</i>	<i>Sordgrassella</i>	A-2-12	Stead_Accr_2	<i>Apis mellifera</i>	USA: Arizona	2,50	75	NAGY01000000
<i>Sordgrassella</i>	<i>Sordgrassella</i>	A3	Stead_Accr_2	<i>Apis mellifera</i>	USA: Arizona	2,43	109	NAGY01000000
<i>Sordgrassella</i>	<i>Sordgrassella</i>	A5	Stead_Accr_2	<i>Apis mellifera</i>	USA: Arizona	2,43	120	NAGY01000000
<i>Sordgrassella</i>	<i>Sordgrassella</i>	A5-24	Stead_Accr_2	<i>Apis mellifera</i>	USA: Arizona	2,49	172	NAGY01000000
<i>Sordgrassella</i>	<i>Sordgrassella</i>	A-9-24	Stead_Accr_2	<i>Apis mellifera</i>	USA: Arizona	2,50	62	NAGY01000000
<i>Sordgrassella</i>	<i>Sordgrassella</i>	Aw-18	Stead_Accr_2	<i>Apis mellifera</i>	USA: Arizona	2,50	88	NAGY01000000
<i>Sordgrassella</i>	<i>Sordgrassella</i>	Aw-20	Stead_Accr_2	<i>Apis mellifera</i>	USA: Arizona	2,50	65	NAGY01000000
<i>Sordgrassella</i>	<i>Sordgrassella</i>	11	Stead_Accr_2	<i>Apis mellifera</i>	USA: Arizona	2,39	381	NAGY00000000
<i>Sordgrassella</i>	<i>Sordgrassella</i>	ES2406	Stead_Accr_2	<i>Apis mellifera</i>	Switzerland: Lucerne	2,45	15	QGL000000000
<i>Sordgrassella</i>	<i>Sordgrassella</i>	J21	Stead_Accr_2	<i>Apis mellifera</i>	USA: West Haven, CT	2,33	456	AVZ000000000
<i>Sordgrassella</i>	<i>Sordgrassella</i>	MS1-3	Stead_Accr_2	<i>Apis mellifera</i>	USA: Austin, TX	2,50	93	MEJH00000000
<i>Sordgrassella</i>	<i>Sordgrassella</i>	N-23	Stead_Accr_2	<i>Apis mellifera</i>	Norway: Ass	2,42	128	NAGY00000000
<i>Sordgrassella</i>	<i>Sordgrassella</i>	N9	Stead_Accr_2	<i>Apis mellifera</i>	Norway: Ass	2,40	129	NAGY01000000
<i>Sordgrassella</i>	<i>Sordgrassella</i>	N-S1	Stead_Accr_2	<i>Apis mellifera</i>	Norway: Ass	2,42	98	NAGY00000000
<i>Sordgrassella</i>	<i>Sordgrassella</i>	N-S2	Stead_Accr_2	<i>Apis mellifera</i>	Norway: Ass	2,42	73	NAGY01000000
<i>Sordgrassella</i>	<i>Sordgrassella</i>	N-S4	Stead_Accr_2	<i>Apis mellifera</i>	Norway: Ass	2,42	38	NAGY00000000
<i>Sordgrassella</i>	<i>Sordgrassella</i>	N-S5	Stead_Accr_2	<i>Apis mellifera</i>	Norway: Ass	2,42	77	NAGY00000000
<i>Sordgrassella</i>	<i>Sordgrassella</i>	N-W4	Stead_Accr_2	<i>Apis mellifera</i>	Norway: Ass	2,42	75	NAGY01000000
<i>Sordgrassella</i>	<i>Sordgrassella</i>	N-W7	Stead_Accr_2	<i>Apis mellifera</i>	Norway: Ass	2,42	62	NAGY01000000
<i>Sordgrassella</i>	<i>Sordgrassella</i>	OH1	Stead_Accr_2	<i>Apis mellifera</i>	USA: West Haven, CT	1,37	401	JAEH00000000
<i>Sordgrassella</i>	<i>Sordgrassella</i>	P14	Stead_Accr_2	<i>Apis mellifera</i>	USA: West Haven, CT	1,31	385	JAC000000000
<i>Sordgrassella</i>	<i>Sordgrassella</i>	PE00171	Stead_Accr_2	<i>Apis mellifera</i>	USA: West Haven, CT	2,52	77	MEJH00000000
<i>Sordgrassella</i>	<i>Sordgrassella</i>	PE00178	Stead_Accr_2	<i>Apis mellifera</i>	USA: West Haven, CT	2,52	135	MEJH00000000
<i>Sordgrassella</i>	<i>Sordgrassella</i>	wk82	Stead_Accr_2	<i>Apis mellifera</i>	USA: Austin, TX	2,53	Combing	CP07446
<i>Sordgrassella</i>	<i>Sordgrassella</i>	N-53	Stead_Accr_2	<i>Apis mellifera</i>	Norway: Ass	2,46	79	NAGY00000000
<i>Sordgrassella</i>	<i>Sordgrassella</i>	wk8332	Stead_Accr_2	<i>Apis mellifera</i>	Malaysia: Genting Highlands	2,49	80	MEJH00000000
<i>Sordgrassella</i>	<i>Sordgrassella</i>	wk8339	Stead_Accr_2	<i>Apis mellifera</i>	Malaysia: Genting Highlands	2,50	27	MEJH00000000
<i>Sordgrassella</i>	<i>Sordgrassella</i>	wk89	Stead_Accr_2	<i>Apis mellifera</i>	USA: West Haven, CT	2,55	15	MEJH00000000

686

TABLE S2 Information of the marker genes.

PhyEco marker	Gene	Length/bp
B000079	<i>frr</i>	558
B000041	<i>NusA</i>	1,476
B000103	<i>PTH</i>	642
B000063	<i>RbfA</i>	378
B000080	<i>RecR</i>	606
B000039	<i>rnhB</i>	627
B000096	<i>RibF</i>	939
B000086	<i>RimM</i>	531
B000062	<i>RsfS</i>	315
B000071	<i>RuvA</i>	609
B000065	<i>smpB</i>	483
B000032	<i>truB</i>	921
B000082	<i>miaA</i>	912
B000114	<i>murB</i>	1,011
B000081	<i>YebY</i>	468

TABLE S3 List of barcode sequences.

Barcode NO.	Forward seq (5'to 3')	Reverse seq (5'to 3')
B01	ATCACG	ACTGAT
B02	CGATGT	ATGAGC
B03	TTAGGC	ATTCTT
B04	TGACCA	CAAAAG
B05	ACAGTG	CAACTA
B06	GCCAAT	CACCGG
B07	CAGATC	CACGAT
B08	ACTTGA	CACTCA
B09	GATCAG	CAGGCG
B10	TAGCTT	CATGGC
B11	GGCTAC	CATTTT
B12	CTTGTA	CCAACA
B13	AGTCAA	CGGAAT
B14	AGTTCC	CTAGCT
B15	ATGTCA	CTATAC
B16	CCGTCC	CTCAGA
B17	GTAGAG	GACGAC
B18	GTCCGC	TAATCG
B19	GTGAAA	TACAGC
B20	GTGGCC	TATAAT
B21	GTTTCG	TCATTC
B22	CGTACG	TCCCGA
B23	GAGTGG	TCGAAG
B24	GGTAGC	TCGGCA

TABLE S4 Mixing ratio of mock samples.

SampleID	Mixing ratio					Barcode NO.
	B2776	B2889	B3801	B3172	B3788	
S01	20.00	20.00	20.00	20.00	20.00	B01
S02	25.00	25.00	0.00	25.00	25.00	B02
S03	0.00	25.00	25.00	25.00	25.00	B03
S04	25.00	0.00	25.00	25.00	25.00	B04
S05	25.00	25.00	25.00	0.00	25.00	B05
S06	25.00	25.00	25.00	25.00	0.00	B06
S07	33.33	33.33	0.00	33.33	0.00	B07
S08	0.00	33.33	0.00	33.33	33.33	B08
S09	0.00	0.00	33.33	33.33	33.33	B09
S10	33.33	0.00	33.33	0.00	33.33	B10
S11	33.33	33.33	33.33	0.00	0.00	B11
S12	24.39	24.39	2.44	24.39	24.39	B12
S13	24.94	24.94	0.25	24.94	24.94	B13
S14	24.99	24.99	0.02	24.99	24.99	B14
S15	20.00	20.00	20.00	20.00	20.00	B15
S16	1.00	9.00	50.00	10.00	30.00	B16
S17	50.00	1.00	30.00	9.00	10.00	B17
S18	30.00	50.00	10.00	1.00	9.00	B18
S19	10.00	30.00	9.00	50.00	1.00	B19
S20	9.00	10.00	1.00	30.00	50.00	B20
S21	0.00	0.00	10.00	90.00	0.00	B21
S22	0.00	0.00	1.00	99.00	0.00	B22
S23	0.00	0.00	0.10	99.90	0.00	B23
S24	0.00	0.00	50.00	50.00	0.00	B24
D01	20.00	20.00	20.00	20.00	20.00	B01
D02	20.00	20.00	20.00	20.00	20.00	B02
D03	20.00	20.00	20.00	20.00	20.00	B03
D04	20.00	20.00	20.00	20.00	20.00	B04
D05	20.00	20.00	20.00	20.00	20.00	B05
D06	20.00	20.00	20.00	20.00	20.00	B06
D07	90.00	9.00	0.90	0.09	0.01	B07
D08	9.00	0.90	0.09	0.01	90.00	B08
D09	0.90	0.09	0.01	90.00	9.00	B09
D10	0.09	0.01	90.00	9.00	0.90	B10
D11	0.09	0.01	90.00	9.00	0.90	B11
D12	0.01	90.00	9.00	0.90	0.09	B12
D13	90.00	9.00	0.90	0.09	0.01	B13
D14	9.00	0.90	0.09	0.01	90.00	B14
D15	0.90	0.09	0.01	90.00	9.00	B15
D16	0.09	0.01	90.00	9.00	0.90	B16
D17	0.09	0.01	90.00	9.00	0.90	B17

(Continued Table S4)

D18	0.01	90.00	9.00	0.90	0.09	B18
D19	0.00	0.00	1.64	16.39	81.97	B19
D20	0.00	0.00	1.64	16.39	81.97	B20
D21	0.00	1.64	16.39	81.97	0.00	B21
D22	1.64	16.39	81.97	0.00	0.00	B22
D23	16.39	81.97	0.00	0.00	1.64	B23
D24	81.97	0.00	0.00	1.64	16.39	B24

Note: B2776, B2889, B3801, B3172 and B3788 are the representative strain of Acer_Giliia_1 to Acer_Giliia_5, respectively.

TABLE S5 Statistics of data outputs.

LibraryID	Raw reads	Clean reads	Raw base/G	Clean base/G	Effective rate/%	Q20/%	Q30/%	GC content/%
f1S01-f1S24	2,910,358	2,904,839	0.87	0.87	99.81	98.96	96.8	42.92
f2S01-f2S24	4,370,025	4,362,026	1.31	1.31	99.82	98.99	96.88	42.85
f3S01-f3S24	3,971,727	3,966,181	1.19	1.19	99.86	98.34	94.65	42.84
N1S01-N1S24	3,101,708	3,097,334	0.93	0.93	99.86	98.83	96.33	38.67
N2S01-N2S24	3,455,304	3,451,312	1.04	1.04	99.88	97.96	93.59	38.66
N3S01-N3S24	2,893,355	2,889,594	0.87	0.87	99.87	97.8	93.25	38.64
P1S01-P1S24	5,446,708	5,439,697	1.63	1.63	99.87	99.1	96.48	36.74
P2S01-P2S24	2,698,030	2,694,490	0.81	0.81	99.87	98.97	96.29	36.88
P3S01-P3S24	3,377,356	3,371,139	1.01	1.01	99.82	97.95	93.16	37.08
fD01-fD24	3,599,515	3,595,846	1.08	1.08	99.9	98.41	95.28	42.44
ND01-ND24	4,399,592	4,393,529	1.32	1.32	99.86	98.42	95.02	38.73
PD01-PD24	3,387,737	3,380,942	1.02	1.01	99.8	98.53	95.13	36.64
fB0061-fB14781	5,441,182	5,434,806	1.63	1.63	99.88	99.24	97.36	42.25
PB0061-PB14781	3,591,409	3,587,252	1.08	1.08	99.88	98.8	95.75	36.55

Note: 1. f1 - f3, N1 - N3 and P1 - P3 represent the three replicates for *frr*, *NusA* and *PTH* gene sequencing, and S01 - S24 are mock samples with different ratio of mixing bacterial cells shown in Table S4.

2. f, N and P represent the *frr*, *NusA* and *PTH* gene sequencing, and D01 - D24 are mock sample with different ratio of mixing bacterial DNA shown in Table S4.

3. f and P represent the *frr* and *PTH* gene sequencing, and BXXXX present *Apis cerana* gut sample

TABLE S6 Summary of read processing and data obtained from marker gene, 16S V4 amplicon and metagenomic sequencing of honey bee guts.

Gut ID	Raw PE reads				Joined and filtered reads				<i>Gilliamella</i> reads	
	16S	<i>frr</i>	<i>PTH</i>	Meta	16S	<i>frr</i>	<i>PTH</i>	Meta	16S	Meta
B0061	84,584	358,851	65,424	33,687,518	82,760	311,842	60,964	5,159,739	25,700	3,185
B0070	85,169	253,390	94,393	37,685,288	83,491	223,179	91,436	5,909,974	47,892	7,661
B0108	83,908	349,570	123,368	45,115,102	82,132	344,248	121,018	4,242,595	39,424	2,189
B0120	85,368	389,432	29,023	34,543,934	83,267	262,273	28,277	4,523,366	22,741	815
B0154	83,691	289,878	-	34,969,488	81,143	226,310	-	4,622,933	26,346	1,086
B0174	84,281	361,728	75,882	38,471,836	81,979	287,357	68,222	6,058,633	18,258	243
B14756	-	224,748	118,334	21,491,068	-	194,850	114,236	4,622,933	-	17,472
B14757	-	354,956	158,879	22,959,909	-	328,188	156,087	4,622,933	-	5,546
B14758	-	277,823	165,638	24,408,709	-	224,100	160,741	9,658,926	-	4,182
B14779	-	342,928	151,481	23,802,654	-	285,197	146,421	7,922,065	-	8,133
B14780	-	301,064	48,088	23,495,452	-	272,069	47,247	7,123,654	-	3,064
B14781	-	291,415	71,187	22,381,871	-	237,037	68,190	9,626,197	-	3,686

Note: “16S” indicates 16S V4; “Meta” indicates metagenomic; “-” indicates no test.