1    **Title: Comparing faster evolving *rplB* and *rpsC* versus SSU rRNA for improved microbial**

2    **community resolution**

3

4    **Authors:** Jiarong Guo[a], James R. Cole[a], C. Titus Brown[b], James M. Tiedje[a]

5

6    **Author affiliation:**

7

8    [a]Center for Microbial Ecology, Michigan State University

9    [b]Department of Population Health and Reproduction, University of California, Davis

10

11    **Corresponding author:**

12    James M. Tiedje

13    tiedjej@msu.edu

14

17

18

19

20

21

22

23    **Abstract:**

24          Amplicon sequencing of the SSU rRNA gene is standard for microbial ecology but has

25    several drawbacks including limited resolving power for taxa below the level of genus and

26    variable multiplicity presenting difficulties in quantifying different organisms. Many conserved

27    protein-coding core genes are single copy and evolve faster than the SSU rRNA gene but their

28    use has been precluded by the lack of universal primers for amplicon sequencing. Recent

29    advances in gene targeted assembly methods for large shotgun metagenomes make their use

30    feasible. To evaluate this approach, we compared the variation of two single copy ribosomal

31    protein genes, *rplB* and *rpsC,* with the SSU rRNA gene for all completed bacterial genomes in

32    NCBI RefSeq. As expected, among pairwise comparisons of all species that belong to the same

33    genus, 94.9% and 91.0% of the pairs of *rplB* and *rpsC,* respectively, showed more variation than

34    did their SSU rRNA sequences. To circumvent primer bias and lack of universal primer issues of

35    amplicon methods, we used a gene targeted assembler, Xander, to assemble *rplB* and *rpsC* from

36    shotgun metagenomic data. When tested on rhizosphere samples of three crops -- corn, an

37    annual, and *Miscanthus* and switchgrass, both perennials -- both genes separated all three

38    communities while SSU rRNA gene could only separate the annual from the two perennial

39    communities in ordination analyses. Furthermore, the Xander assemblies of *rplB* and *rpsC*

40    yielded significantly higher numbers of OTUs (alpha diversity) than SSU rRNA gene recovered

41    from short reads and from amplicon data. These results confirm the better resolving ability of

42    these faster evolving marker genes for comparative microbiome studies.

43

44    **Importance:**

45      High resolution marker genes are central to determining diversity of communities and

46      differences between or among communities. Many ecologically determinative features occur at

47      genetic levels not resolved by the relatively conserved SSU rRNA gene; hence marker genes are

48      needed with finer community resolution. Further, if they were single copy, counting would be

49      more accurate than for the variable copy SSU rRNA genes. The rapid advancement of shotgun

50      sequencing and metagenome assembly has enabled us to avoid the need for and the inevitable

51      bias of primers, to recover single-copy protein-coding genes directly from shotgun metagenomes.

52      Targeting a few genes for assembly, like those coding for ribosomal proteins, samples more

53      organisms and speeds the analysis over using whole genome assemblies for this purpose.

54

**Introduction:**

55

56      Shaped by 3.5 billion years of evolution, microorganisms are estimated to comprise up to one

57      trillion species and the majority of genetic diversity in the biosphere (1). However, our

58      understanding of this diversity is limited because of this huge number, and that the majority are

59      yet to be cultured and their physiology or functions characterized. Since the pioneering work of

60      Carl Woese in the late 1970s, the SSU rRNA gene has been the dominant marker used in

61      microbial community structure analyses (2–5). While it has been extremely useful to advancing

62      understanding of the microbial world, it does have important limitations, namely that it is highly

63      conserved and that there are usually multiple copies and some with intra-genomic variations

64      making this gene problematic for taxonomic identification at species and ecotypes levels and

65      incapable of reflecting community distinctions at ecologically meaningful levels (6–8).

66      With the accelerated accumulation of microbial genomes in NCBI in recent years (9), whole

67      genome-based comparison is now feasible and a more accurate method for species and strain

68    identification (9–15). However, whole genome-based comparison is computationally more

69    expensive compared to marker gene comparison, and it is not yet possible to reliably obtain

70    genome sequences of many members of natural microbial communities. Hence, marker gene

71    analysis remains useful. Single copy protein coding housekeeping genes stand out as the best

72    candidates. First, their single copy status provides more accurate species and strain counting,

73    identification and OTU clustering than the SSU rRNA gene. Second, they are present in virtually

74    all members of the three domains of life. Third, protein coding genes evolve faster than rRNA

75    genes not only because rRNA genes are more conserved due to their critical role in ribosome

76    function (16), but also because of the redundancy in the genetic code, especially at the third

77    codon position (6).

78        Here, we evaluate two single copy protein coding genes, *rplB* (50S ribosomal large subunit

79    protein L2), and *rpsC* (30S ribosomal small subunit protein S3) as potential housekeeping genes

80    for phylogenetic markers for microbial community analyses. Earlier studies showed the potential

81    of protein coding genes over SSU rRNA genes as higher resolution phylogenetic markers for

82    microbial diversity analyses using both genomic data (111 genomes) and metagenomic data (< 6

83    Gbp by Sanger sequencing) (6, 8). We revisited this comparison with the now much larger data

84    set - all completed bacterial genomes (~4500 with one contig) and then tested the resolving

85    power of these two genes versus SSU rRNA gene among different crop rhizospheres using large

86    shotgun metagenomic data (~1TB). The novelty of our analyses is the application of gene

87    targeted assembly to recover single copy protein coding genes from shotgun metagenomic data

88    (17) and the use of *de novo* OTU-based diversity analyses, commonly used in microbial diversity

89    analyses, rather than just taxonomic identification as previous studies (6, 8).

90

4

**Methods:**

91

92    Bacterial genome assembly information from NCBI

93    (ftp://ftp.ncbi.nlm.nih.gov/genomes/refseq/bacteria/assembly_summ

94    ary.txt) was used to construct the link to download each genome based on the instructions

95    described in this link

96    (http://www.ncbi.nlm.nih.gov/genome/doc/ftpfaq/#allcomplete).

97    Command line "wget" was then used to retrieve the genome sequences with links obtained from

98    the above step.

99    For extracting genes from genomes, the SSU rRNA gene HMM (Hidden Markov Model)

100   from SSUsearch (18) was used to recover rRNA genes. Aligned *rplB* and *rpsC* nucleotide

101   sequences of the "training set" retrieved from the RDP FunGene database (19) were used to build

102   the HMM models using hmmbuild command in HMMER (version 3.1b2) (20). The nhmmer

103   command in HMMER was then used to identify SSU rRNA, *rplB* and *rpsC* sequences from

104   bacteria genomes obtained from NCBI using score cutoff (-T) of 60. Next, nhmmer hits of least

105   90% of the length of the HMM model were accepted as the target gene. For the purpose of

106   comparing SSU rRNA and *rplB* and *rpsC* gene distances, one copy of the SSU rRNA gene was

107   randomly picked from each genome. Pairwise comparison among gene sequences was done

108   using vsearch (version 1.1.3) with "--allpairs_global --acceptall" (21). Three species of

109   environmental interest, *Rhizobium leguminosarum*, *Pseudomonas putida* and *Escherichia coli,*

110   were chosen for closer comparison of *rplB* and *rpsC* pairwise distances and SSU rRNA gene

111   distances.

112   The shotgun data are from DNA from seven field replicates of rhizosphere samples of three

113   biofuel crops: corn (C) *Zea maize*, switchgrass (S) P*anicum virgatum*, and *Miscanthus* x gigantus

5

114    (M) that had been grown for 5 years. Shotgun sequence data for the 21 samples were downloaded

115    from the JGI web portal (http://genome.jgi.doe.gov/); JGI Project IDs are listed in Table S1. Raw

116    reads were quality trimmed using fastq-mcf in EA-Utils (verison 1.04.662)

117    (http://code.google.com/p/ea-utils) "-l 50 -q 30 -w 4 -k 0 -x 0 --max-ns 0 -X". Overlapping

118    paired-end reads were merged by FLASH (version 1.2.7) (22) with "-m 10 -M 120 -x 0.20 -r 140

119    -f 250 -s 25" described in (18).

120    SSU rRNA gene amplicon data (JGI project ID: 1025756) from the same DNA used for

121    shotgun sequence were trimmed the same way as shotgun data (described above). Paired ends

122    were joined by FLASH (-m 10 -M 150 -x 0.08 -p 33 -r 200 -f 300 -s 25) (22) and primer

123    sequences were removed by cutadapt (-f fasta --discard-untrimmed) (23). For community

124    analyses, the open reference OTU picking method in QIIME was used for clustering and Bray-

125    Curtis index was used for beta-diversity index (24).

126    For SSU rRNA gene analyses with shotgun data, SSU rRNA gene fragments and those

127    aligned to the V4 region (*E. coli* position: 577 - 727) of each sample were identified using the

128    SSUsearch pipeline (18) and clustered using RDP's McClust tool (25) at a distance of 0.05 and

129    minimal overlap of 25 bp, following the tutorial in SSUsearch (http://microbial-ecology-

130    protocols.readthedocs.io/en/latest/SSUsearch/overview.html).

131    Both *rplB* and *rpsC* sequences were assembled using Xander with

132    "MAX_JVM_HEAP=500G, FILTER_SIZE=40, K_SIZE=45, genes = *rplB* and *rpsC*,

133    MIN_LENGTH=150, THREADS=9" (17). Data for each crop were assembled separately. The

134    assembled *rplB* or *rpsC* sequences (nucleotide and protein) from the three crops were pooled and

135    clustered using RDP's McClust tool (25). For each gene, a table of OTU counts of each sample

136    was made based on mean k-mer coverage of the representative sequence of each OTU (provided

137    in "*_coverage.txt" output file from Xander). Further, diversity analyses were done with the

6

138    vegan package in R using functions "rda" for ordination and "diversity" for Shannon diversity

139    index, respectively, from the OTU (count) tables. An implementation of this pipeline is publicly

140    available at https://doi.org/10.5281/zenodo.1438073.

141        To assess how many potential target gene reads of *rplB* and *rpsC* were assembled by Xander,

142    we did a six-frame translation of the short reads (nucleotide sequences) into protein sequences by

143    transeq in EMBOSS tool (26). We then searched HMMs against the protein sequences and the

144    hits with bit score > 40 (e-value < $6.2 * 10^{-6}$) were treated as reads from the target gene.

145    Meanwhile, "*_match_reads.fa", a collection of reads that share a k-mer (k=45) with assembled

146    sequences, output from Xander, provided the reads assembled by Xander. Then we compared the

147    fold coverage of reads found by hmmsearch and reads used by Xander, by estimating fold

148    coverage of each read with median kmer coverage using khmer package (27, 28).

149    **Results:**

150        A total of 4,457 of complete bacteria genomes defined as one sequence were downloaded.

151    SSU rRNA gene copy number ranged from 1 to 16 with a mean of 4 and 99.9% of genomes have

152    single copies of *rplB* and *rpsC* (Table S2). Both of these genes were present in 4,440 of the

153    complete genomes. When evaluating intra-genomic variation among copies of SSU rRNA genes

154    in completed genomes of *R. leguminosarum*, *P. putida* and *E. coli*, *E. coli* had the largest

155    variation with a minimum of 95.4% identity (Fig. S1). For the pairwise comparison between

156    genomes, one copy of each gene was randomly picked as a representative for genomes with

157    multiple copies.

158        For the selected taxa, Rhizobiales, Pseudomonadales, *Rhizobium*, and *Pseudomonas*, *rplB*

159    and *rpsC* had similar variations and both had larger variation among the genomes than SSU

160    rRNA genes within their corresponding order (among genera), and genus (among species) (Fig. 1

7

161    and 2). When comparing all species of completed genomes that belong to the same genus, we

162    found SSU rRNA gene has an identity range of 63.2% to 100.0% and a median of 95.2%, *rplB*

163    has an identity range of 43.2% to 100.0% and a median of 87.2%, *rpsC* has an identity range of

164    46.0% to 100.0% and a median of 90.3%. Between *rplB* and SSU rRNA gene, 88,993 pairs

165    (94.9% of total) has larger variation in *rplB*, 3,573 pairs have larger variation in SSU rRNA gene,

166    and 1,167 pairs have the same variation (Fig. 3A); 77,885 pairs (91.0% of total) has larger

167    variation in *rpsC*, 6,074 pairs have larger variation in SSU rRNA gene, and 1,622 pairs have the

168    same variation (Fig. 3B); 54,755 pairs (63.7%) has larger variation in *rplB*, 28,393 pairs have

169    larger variation in *rpsC* gene, and 2,808 pairs have the same variation for *rplB* and *rpsC* (Fig.

170    3C).

171    　　We compared SSU rRNA genes with *rplB and rpsC* to test the ability of shotgun data to

172    resolve community differences among plant rhizospheres. We chose these two genes as they had

173    a suitable length for Xander assembly, were long enough for resolving power, and had HMMs

174    that were both specific and sensitive for fragment recovery due to their uniqueness in sequence as

175    parts of the ribosome, and both have been used as phylogenetic marker in other shotgun

176    metagenomic studies (29, 30). On average, 0.04% of total reads were identified as SSU rRNA

177    gene fragments and 0.004% of total reads aligned to the 150 bp of V4 region of the gene with

178    SSUsearch (18). Another 0.01% and 0.008% of total reads were identified as *rplB* and *rpsC,*

179    respectively, by Xander (Table S3). To test the sensitivity of Xander, we found that the number

180    of potential *rplB* and *rpsC* reads assembled were 49.5% and 47.9%, respectively, of those defined

181    by hmmsearch with bit score cutoff of 40 (Table S4) and have much higher fold coverage than

182    the rest of reads (excluding shared reads) in hmmsearch hits (Figure S2).

183    　　Beta diversity analyses of all three genes showed that the rhizosphere communities of the

184    annual crop, corn, were different from those of the two perennial grasses, *Miscanthus* and

8

185    switchgrass, but only *rplB* and *rpsC* distinguished the communities of the two perennial grasses

186    (Fig. 4). This was true whether the analysis was at the nucleotide or protein level. The alpha

187    diversity of the corn rhizosphere communities was significantly lower than those of *Miscanthus*

188    and switchgrass rhizospheres by all three measures except for Chao1 index with *rpsC* and SSU

189    rRNA gene (Fig. 5). When comparing among genes, the numbers of OTUs from *rplB* and *rpsC*

190    are also significantly higher than SSU rRNA gene (Fig. 5). Since SSUSearch returns shorter

191    fragments than Xander assembled genes, we also evaluated whether the longer fragments of SSU

192    rRNA from amplicon data ~ 250 to 300 bp, could distinguish the two perennial grass

193    communities, and they could not (Fig. 3E).

194

195    **Discussion:**

196       We confirmed the advantages of *rplB* and *rpsC* over the SSU rRNA gene as a more resolving

197    phylogenetic marker using updated large genomic data (~4500 complete genomes) (Fig. 1, 2 and

198    3). We also demonstrated that *rplB* and *rpsC* can be assembled from large shotgun metagenomes

199    and showed that they provided higher community resolution by separating *Miscanthus* and

200    switchgrass rhizosphere samples while the SSU rRNA gene did not (Fig. 4).  The two perennial

201    grasses would be expected to have more similar microbiome than the annual since the latter is re-

202    established each year while the fibrous perennial grass roots are more similar and not physically

203    disturbed annually and thus do not have full regrowth at a new random site each year.

204    In large genomic data analyses, *rplB* and *rpsC* show advantages in following three aspects:

205       First, SSU rRNA gene, a multiple copy gene, poses difficulties for interpreting species

206    abundance, while *rplB* and *rpsC* do not have the same issue as it is single copy genes in > 99.9%

207    of complete genomes (Table S2). Additionally, variations among multiple SSU copies can cause

208    multiple OTUs (sequence clusters) from the same species (Figure S1) and thus leads to

9

209    overestimation of species richness (31). Since a single copy of the *rplB* and *rpsC* genes is

210    contained in every cell in a community, the relative abundance of *rplB* and *rpsC* gene sequences

211    provides a reference for estimating the fraction of organisms possessing other genes.

212         Second, *rplB* and *rpsC* are better able to differentiate closely related species based on their

213    lower sequence similarities compared to the SSU rRNA gene in pairwise comparisons among

214    genomes (Fig. 1, 2, and 3). This is consistent with the crucial role SSU rRNA plays in translation

215    (ensuring translation accuracy) (16), also confirmed by another study showing SSU rRNA genes

216    (along with LSU rRNA genes, tRNA and ABC transporter genes) to be the most conserved genes

217    (32).

218         Third, SSU rRNA genes in genomes are also more prone to assembly errors (chimera) than

219    single copy genes due to their higher overall nucleotide identity and the presence of highly

220    conserved regions interspersed in SSU rRNA genes. Note that these erroneous sequences might

221    be further collected by databases and used as references for taxonomy, alignment, and chimera

222    detection, and thus have an impact on common microbial ecology diversity analyses. Switching

223    to a single copy gene that is less prone to assembly error can mitigate the above problem.

224         Finally, this method provides for higher resolution community diversity analyses in large

225    shotgun metagenomes, leveraging a scalable gene targeted assembler, Xander. Assembly is

226    desirable for short read data to correctly identify the gene and provide enough length for

227    resolving power, a major objective in ecology studies. Assembly misses the rarer species that do

228    not have enough sequencing depth in metagenomes, confirmed by the higher fold coverage of

229    reads used in assemblies compared to the other reads in hmmsearch hits (Fig. S2). We did find

230    that the number of reads used in assemblies are about half of the reads identified as the targeted

231    genes by hmmsearch (Table S4). The hmmsearch though could also have recovered some false

232    positives due to mistaken short-read identification and thus overestimated the total gene number.

10

233    However, *rplB* and *rpsC* yield significantly higher alpha diversity (Fig. 5) than SSU rRNA gene

234    despite missing rare members. Thus they reveal more diversity among abundant members than

235    SSU rRNA gene, which offsets and exceeds the diversity of the rare members that are not

236    assembled, further confirming their higher resolution.

237        We choose two protein coding genes to be sure our results were not gene specific, and both

238    gave very similar results at both the nucleotide and protein levels. At least from extensive

239    completed genomes, most of these two genes are single copy making quantitative (ratio)

240    comparisons with other genes more consistent. For future use, *rplB* might have slight advantage

241    over *rpsC* since it is longer, about 830 bp on average vs 660 bp of *rpsC*, providing a bit more

242    resolving power, which is consistent with results in genome comparisons showing *rplB* has lower

243    median sequence identity than *rpsC* (Fig 3).

244        It is of course possible to find in reference databases the best match to the assembled

245    sequence of these marker genes and potentially have finer taxonomic resolution than provided by

246    SSU rRNA. But, the reference database is only from sequenced genomes and hence is very

247    unbalanced and incomplete compared to 16S rRNA databases (17) so this use is not generally

248    beneficial at this time.

249        Although sequencing depth needed varies depending on community diversity, we estimate it

250    based on our rhizosphere soil samples as a practical guide. The reads from *rplB* are around 0.01%

251    of total (Table S3). Assuming a fold coverage of 3000 of *rplB* for each sample, to be comparable

252    to 3000 amplicons in planning amplicon-based studies, one needs about 25 Gbp (3000 * 830 /

253    0.01%) of shotgun metagenome (830 bp is the average gene length of *rplB*). The major

254    requirement for using this method beyond sufficient shotgun sequence depth is an access to a

255    high performance computer since large memory (> 250 Gb recommended for soil samples) is

256    needed to run Xander.

11

257

**Conclusion:**

259     We demonstrated that *rplB* and *rpsC*, single copy protein coding genes can provide finer

260     resolution of taxa and hence better distinguish among communities than the more commonly

261     used SSU rRNA gene and also provide finer scale *de novo* (OTU) diversity analysis. This method

262     does require shotgun sequence of sufficient depth, so is currently more costly than amplicon

263     based analyses, but as sequencing costs decline, capacity and access increase, read length grows,

264     and genome reference databases grow, single copy protein coding genes such *rplB* and *rpsC* have

265     the potential to complement or even replace the SSU rRNA gene as a phylogenetic marker and

266     better reflect ecology of communities.

267

276

277

278

279

280

281

282

283

284

285

286

287  **References:**

288

289  1.  Locey KJ, Lennon JT. 2016. Scaling laws predict global microbial diversity. Proc Natl Acad
290      Sci USA 113:5970–5975.

291  2.  Woese CR, Fox GE. 1977. Phylogenetic structure of the prokaryotic domain: the primary
292      kingdoms. Proc Natl Acad Sci USA 74:5088–5090.

293  3.  Lane DJ, Pace B, Olsen GJ, Stahl DA, Sogin ML, Pace NR. 1985. Rapid determination of
294      16S ribosomal RNA sequences for phylogenetic analyses. Proc Natl Acad Sci USA 82:6955–
295      6959.

296  4.  Huse SM, Dethlefsen L, Huber JA, Mark Welch D, Welch DM, Relman DA, Sogin ML.
297      2008. Exploring microbial diversity and taxonomy using SSU rRNA hypervariable tag
298      sequencing. PLoS Genet 4:e1000255.

299  5.  Caporaso JG, Lauber CL, Walters WA, Berg-Lyons D, Huntley J, Fierer N, Owens SM,
300      Betley J, Fraser L, Bauer M, Gormley N, Gilbert JA, Smith G, Knight R. 2012. Ultra-high-

301     throughput microbial community analysis on the Illumina HiSeq and MiSeq platforms. ISME

302     J 6:1621–1624.

303  6.  Case RJ, Boucher Y, Dahllof I, Holmstrom C, Doolittle WF, Kjelleberg S. 2007. Use of 16S

304     rRNA and rpoB genes as molecular markers for microbial ecology studies. Appl Env

305     Microbiol 73:278–288.

306  7.  Wu M, Eisen JA. 2008. A simple, fast, and accurate method of phylogenomic inference.

307     Genome Biol 9:R151.

308  8.  Roux S, Enault F, Bronner G, Debroas D. 2011. Comparison of 16S rRNA and protein-

309     coding genes as molecular markers for assessing microbial diversity (Bacteria and Archaea)

310     in ecosystems. FEMS Microbiol Ecol 78:617–628.

311  9.  Land M, Hauser L, Jun SR, Nookaew I, Leuze MR, Ahn TH, Karpinets T, Lund O, Kora G,

312     Wassenaar T, Poudel S, Ussery DW. 2015. Insights from 20 years of bacterial genome

313     sequencing. Funct Integr Genomics 15:141–161.

314  10. Goris J, Konstantinidis KT, Klappenbach JA, Coenye T, Vandamme P, Tiedje JM. 2007.

315     DNA-DNA hybridization values and their relationship to whole-genome sequence

316     similarities. Int J Syst Evol Microbiol 57:81–91.

317  11. Luo C, Walk ST, Gordon DM, Feldgarden M, Tiedje JM, Konstantinidis KT. 2011. Genome

318     sequencing of environmental Escherichia coli expands understanding of the ecology and

319     speciation of the model bacterial species. Proc Natl Acad Sci USA 108:7200–7205.

320    12. Varghese NJ, Mukherjee S, Ivanova N, Konstantinidis KT, Mavrommatis K, Kyrpides NC,

321        Pati A. 2015. Microbial species delineation using whole genome sequences. Nucleic Acids

322        Res 43:6761–6771.

323    13. Rodriguez-R LM, Castro JC, Kyrpides NC, Cole JR, Tiedje JM, Konstantinidis KT. 2018.

324        How Much Do rRNA Gene Surveys Underestimate Extant Bacterial Diversity? Appl Environ

325        Microbiol 84:e00014-18.

326    14. Scortichini M, Marcelletti S, Ferrante P, Firrao G. 2013. A Genomic Redefinition of

327        Pseudomonas avellanae species. PLOS ONE 8:e75794.

328    15. Rodriguez-R LM, Gunturu S, Harvey WT, Rosselló-Mora R, Tiedje JM, Cole JR,

329        Konstantinidis KT. 2018. The Microbial Genomes Atlas (MiGA) webserver: taxonomic and

330        gene diversity analysis of Archaea and Bacteria at the whole genome level. Nucleic Acids

331        Res 46:W282–W288.

332    16. Carter AP, Clemons WM, Brodersen DE, Morgan-Warren RJ, Wimberly BT, Ramakrishnan

333        V. 2000. Functional insights from the structure of the 30S ribosomal subunit and its

334        interactions with antibiotics. Nature 407:340–348.

335    17. Wang Q, Fish JA, Gilman M, Sun Y, Brown CT, Tiedje JM, Cole JR. 2015. Xander:

336        employing a novel method for efficient gene-targeted metagenomic assembly. Microbiome

337        3:32.

338    18. Guo J, Cole JR, Zhang Q, Brown CT, Tiedje JM. 2015. Microbial community analysis with

339        ribosomal gene fragments from shotgun metagenomes. Appl Environ Microbiol AEM.02772-

340        15.

15

341    19. Fish JA, Chai B, Wang Q, Sun Y, Brown CT, Tiedje JM, Cole JR. 2013. FunGene: the

342        functional gene pipeline and repository. Front Microbiol 4.

343    20. Eddy SR. 2009. A new generation of homology search tools based on probabilistic inference.

344        Genome Inf 23:205–211.

345    21. Rognes T, Flouri T, Nichols B, Quince C, Mahé F. 2016. VSEARCH: a versatile open source

346        tool for metagenomics. PeerJ 4:e2584.

347    22. Magoc T, Salzberg SL. 2011. FLASH: fast length adjustment of short reads to improve

348        genome assemblies. Bioinformatics 27:2957–2963.

349    23. Martin M. 2011. Cutadapt removes adapter sequences from high-throughput sequencing

350        reads. EMBnet.journal 17:10–12.

351    24. Kuczynski J, Stombaugh J, Walters WA, Gonzalez A, Caporaso JG, Knight R. 2012. Using

352        QIIME to analyze 16S rRNA gene sequences from microbial communities. Curr Protoc

353        Microbiol Chapter 1:Unit 1E.5.

354    25. Cole JR, Wang Q, Fish JA, Chai B, McGarrell DM, Sun Y, Brown CT, Porras-Alfaro A,

355        Kuske CR, Tiedje JM. 2014. Ribosomal Database Project: data and tools for high throughput

356        rRNA analysis. Nucleic Acids Res 42:D633–642.

357    26. Rice P, Longden I, Bleasby A. 2000. EMBOSS: The European Molecular Biology Open

358        Software Suite. Trends Genet 16:276–277.

359    27. Crusoe MR, Alameldin HF, Awad S, Boucher E, Caldwell A, Cartwright R, Charbonneau A,

360        Constantinides B, Edvenson G, Fay S, Fenton J, Fenzl T, Fish J, Garcia-Gutierrez L, Garland

P, Gluck J, González I, Guermond S, Guo J, Gupta A, Herr JR, Howe A, Hyer A, Härpfer A, Irber L, Kidd R, Lin D, Lippi J, Mansour T, McA'Nulty P, McDonald E, Mizzi J, Murray KD, Nahum JR, Nanlohy K, Nederbragt AJ, Ortiz-Zuazaga H, Ory J, Pell J, Pepe-Ranney C, Russ ZN, Schwarz E, Scott C, Seaman J, Sievert S, Simpson J, Skennerton CT, Spencer J, Srinivasan R, Standage D, Stapleton JA, Steinman SR, Stein J, Taylor B, Trimble W, Wiencko HL, Wright M, Wyss B, Zhang Q, zyme en, Brown CT. 2015. The khmer software package: enabling efficient nucleotide sequence analysis. F1000Research.

28. Brown CT, Howe A, Zhang Q, Pyrkosz AB, Brom TH. 2012. A Reference-Free Algorithm for Computational Normalization of Shotgun Sequencing Data. ArXiv12034802 Q-Bio.

29. Sharon I, Kertesz M, Hug LA, Pushkarev D, Blauwkamp TA, Castelle CJ, Amirebrahimi M, Thomas BC, Burstein D, Tringe SG, Williams KH, Banfield JF. 2015. Accurate, multi-kb reads resolve complex populations and detect rare microorganisms. Genome Res 25:534–543.

30. Hug LA, Castelle CJ, Wrighton KC, Thomas BC, Sharon I, Frischkorn KR, Williams KH, Tringe SG, Banfield JF. 2013. Community genomic analyses constrain the distribution of metabolic traits across the Chloroflexi phylum and indicate roles in sediment carbon cycling. Microbiome 1:22.

31. Sun DL, Jiang X, Wu QL, Zhou NY. 2013. Intragenomic heterogeneity of 16S rRNA genes causes overestimation of prokaryotic diversity. Appl Env Microbiol 79:5962–5969.

380    32. Isenbarger TA, Carr CE, Johnson SS, Finney M, Church GM, Gilbert W, Zuber MT, Ruvkun

381       G. 2008. The most conserved genome segments for life detection on Earth and other planets.

382       Orig Life Evol Biosph 38:517–533.

383

384

385

386

387

388

389

390

391

392

393    **Figures**

394    Figure 1: Pairwise comparisons among all genomes of Rhizobiales (panels A, C, E) and of all

395    *Rhizobium* (panels B, D, F) using the SSU rRNA gene, *rplB* and *rpsC*. SSU rRNA gene identities

396    are higher than *rplB* and *rpsC*, and *rplB* and *rpsC* have similar sequence identities in most

397    genome pairs. The dashed line is y = x. Data below y=x line indicate the gene on X axis is more

398    conserved. The dot size indicates the number of pairwise comparisons with those values.

399

400

18

401     Figure 2: Pairwise comparison among all genomes of Pseudomonadales (panels A, C, E) and of

402     all *Pseudomonas* (panels B, D, F) using the SSU rRNA gene, *rplB* and *rpsC*. SSU rRNA gene

403     identities are higher than *rplB* and *rpsC* in most genome pairs, and *rplB* and *rpsC* have similar

404     sequence identities. The dashed line is y = x. Data below y=x line indicate gene on X axis is

405     more conserved. The dot size indicates the number of pairwise comparisons with those values.

406

407

408     Figure 3: Pairwise comparison among all completed genomes of species in the same genus using

409     the SSU rRNA gene, *rplB,* and *rpsC*. SSU rRNA gene identities are larger than *rplB* in most

410     genomes. The diagonal dashed line is y = x and data below the line indicates the gene on X axis

411     is more conserved. The dot intensity is the number of comparisons with those values. Subplot A,

412     B, and C are comparisons of *rplB* and SSU rRNA gene (93,733 pairwise comparisons), *rpsC* and

413     SSU rRNA gene (85,581 pairwise comparisons), *rplB* and *rpsC* (85,956 pairwise comparisons).

414

415

416     Figure 4: Comparison of the SSU rRNA gene, *rpsC* and *rplB* in beta diversity analyses

417     (ordination) using large soil metagenome sequences from seven field replicates. All genes show

418     that the microbial community of the corn (C) rhizosphere is significantly different from those of

419     *Miscanthus* (M) and switchgrass (S) while *rplB* and *rpsC* at both the nucleotide (n) and protein

420     (p) levels separate microbial communities of *Miscanthus* and switchgrass. The SSU rRNA gene

421     does not separate *Miscanthus* and switchgrass with either shotgun (SSU.sg) or amplicon

422     (SSU.am) data. "**" indicates p < 0.01 in PERMANOVA test.

423

424

425    Figure 5: Comparison of the SSU rRNA gene, *rplB*, and *rpsC* in alpha diversity analyses (Chao1,

426    Shannon, and OTU number by protein, "_p" and nucleotide, "_n") using large soil metagenome

427    sequence. All genes show that the microbial community of the corn (C) rhizosphere has

428    significantly less alpha diversity than those of *Miscanthus* (M) and switchgrass (S) except for

429    *rpsC* and SSU rRNA gene with Chao1 index ($p < 0.01$). Wilcoxon test was used to compare

430    SSU_sg against each of the other genes including rplB_n, rplB_p, rpsC_n, rpsC_p and SSU_am.

431    For Chao1 index, rplB_n and rpsC_n show significantly higher abundance than SSU_sg in

432    *Miscanthus* and switchgrass; For Shannon index and OTU number, rplB_n and rpsC_n show

433    significantly higher abundance than SSU_sg in all three crops ("***" is $p < 0.001$, "**" is $p <$

434    $0.01$, "*" is $p < 0.05$, "ns" is $p > 0.05$).
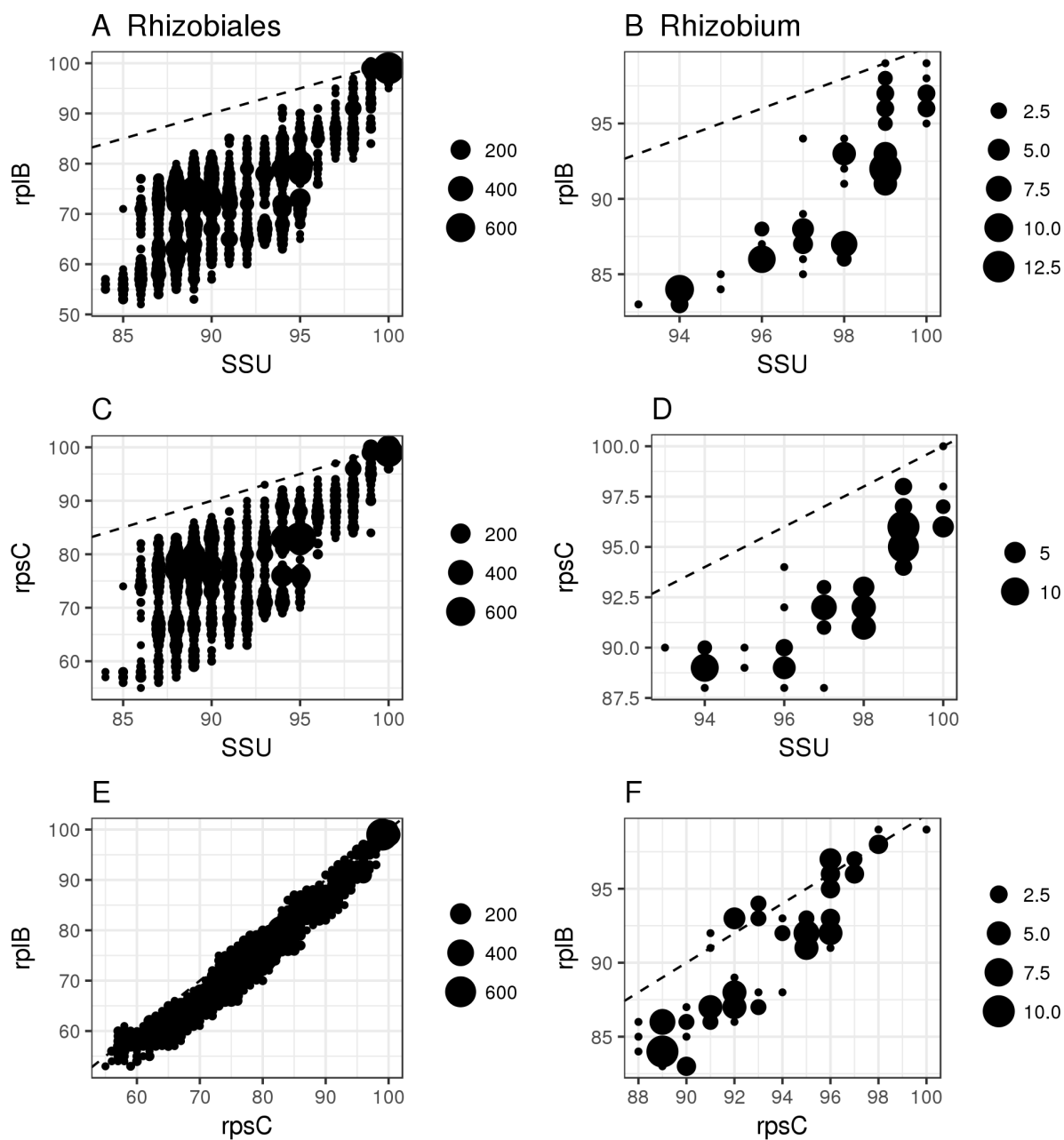
435

436

# Figures

Figure 1: Pairwise comparisons among all genomes of Rhizobiales (panels A, C, E) and of all *Rhizobium* (panels B, D, F) using the SSU rRNA gene, *rplB* and *rpsC*. SSU rRNA gene identities are higher than *rplB* and *rpsC*, and *rplB* and *rpsC* have similar sequence identities in most genome pairs. The dashed line is y = x. Data below y=x line indicate the gene on X axis is more conserved. The dot size indicates the number of pairwise comparisons with those values.
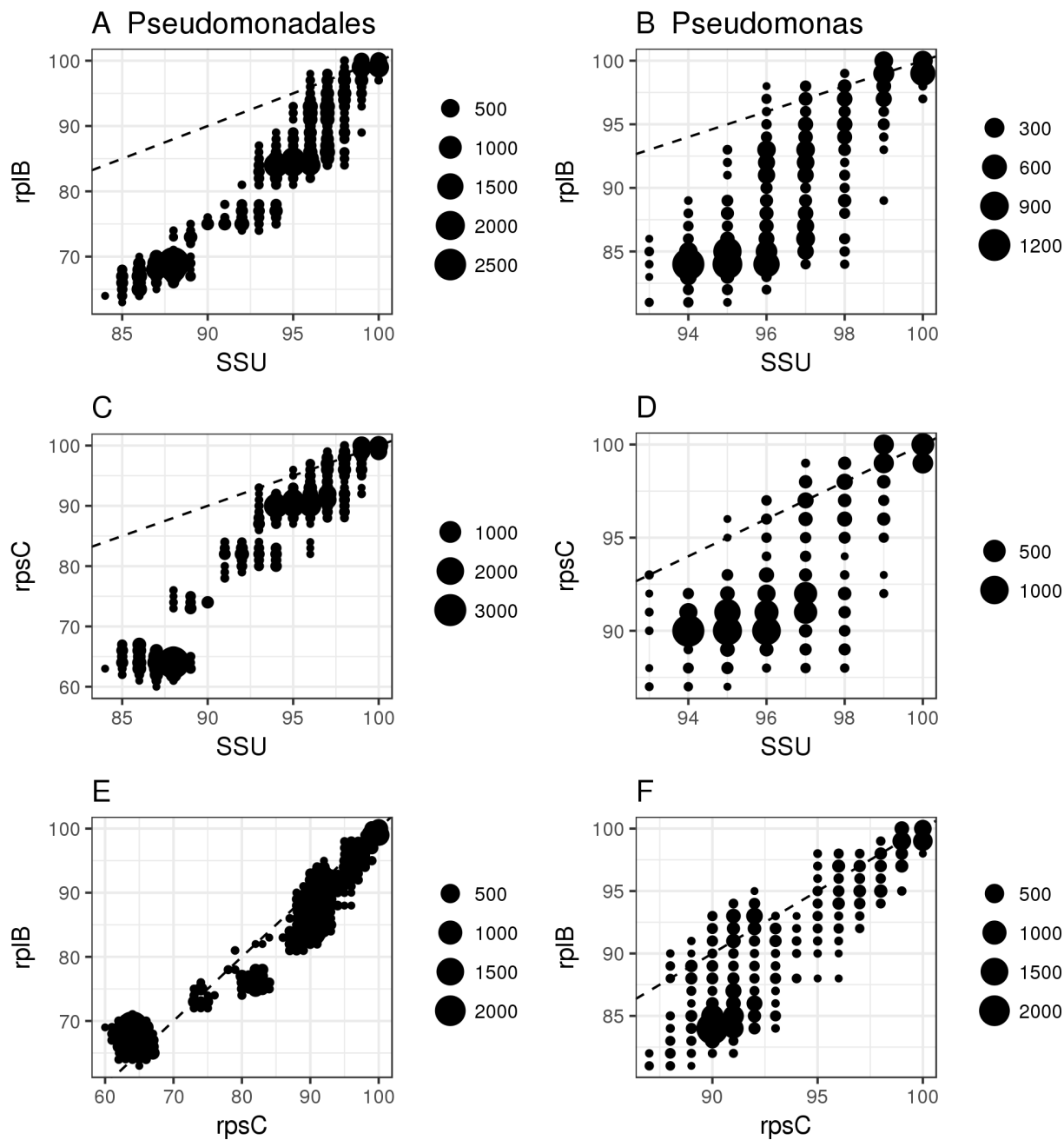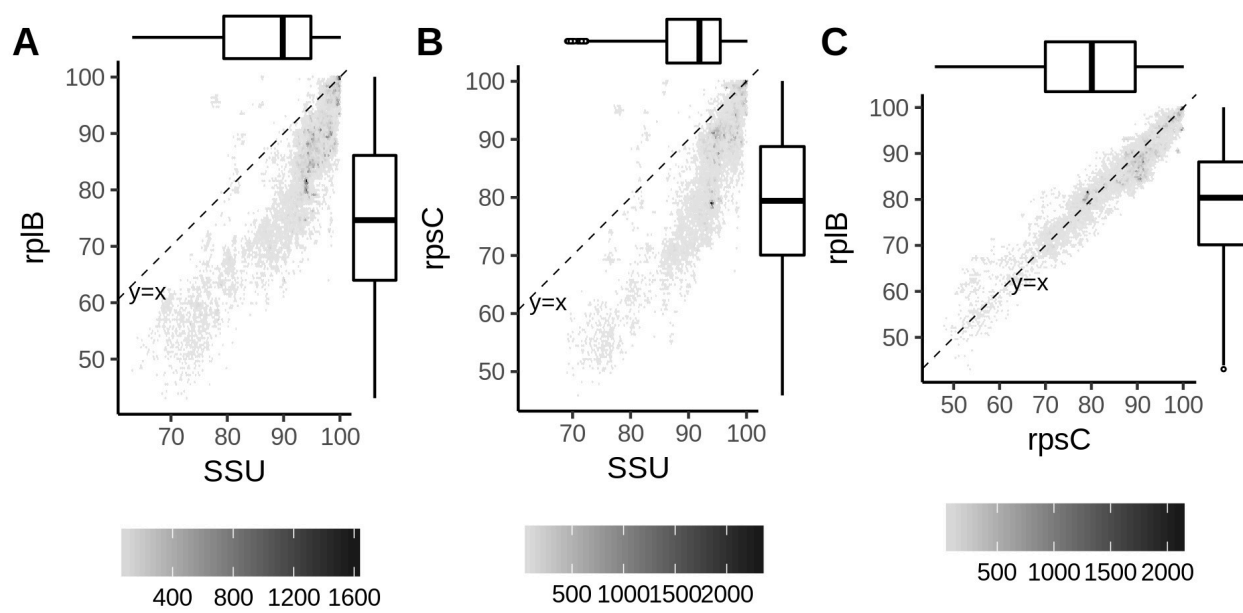
Figure 2: Pairwise comparison among all genomes of Pseudomonadales (panels A, C, E) and of all *Pseudomonas* (panels B, D, F) using the SSU rRNA gene, *rplB* and *rpsC*. SSU rRNA gene identities are higher than *rplB* and *rpsC* in most genome pairs, and *rplB* and *rpsC* have similar sequence identities. The dashed line is y = x. Data below y=x line indicate gene on X axis is more conserved. The dot size indicates the number of pairwise comparisons with those values.

Figure 3: Pairwise comparison among all completed genomes of species in the same genus using the SSU rRNA gene, *rplB,* and *rpsC*. SSU rRNA gene identities are larger than *rplB* in most genomes. The diagonal dashed line is y = x and data below the line indicates the gene on X axis is more conserved. The dot intensity is the number of comparisons with those values. Subplot A, B, and C are comparisons of *rplB* and SSU rRNA gene (93,733 pairwise comparisons), *rpsC* and SSU rRNA gene (85,581 pairwise comparisons), *rplB* and *rpsC* (85,956 pairwise comparisons).
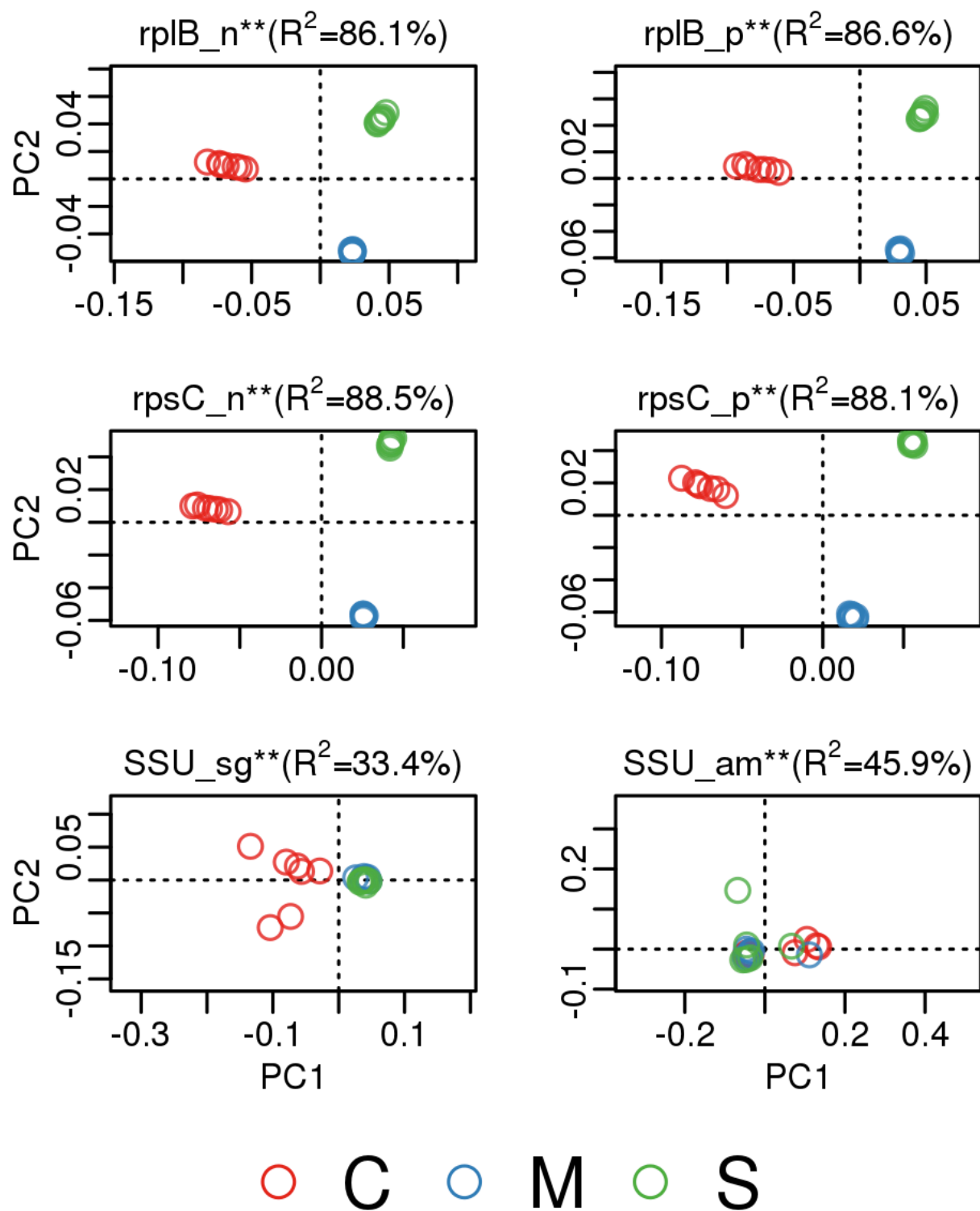
5

Figure 4: Comparison of the SSU rRNA gene, *rpsC* and *rplB* in beta diversity analyses (ordination) using large soil metagenome sequences from seven field replicates. All genes show that the microbial community of the corn (C) rhizosphere is significantly different from those of *Miscanthus* (M) and switchgrass (S) while *rplB* and *rpsC* at both the nucleotide (n) and protein (p) levels separate microbial communities of *Miscanthus* and switchgrass. The SSU rRNA gene does not separate *Miscanthus* and switchgrass with either shotgun (SSU.sg) or amplicon (SSU.am) data. "**" indicates p < 0.01 in PERMANOVA test.
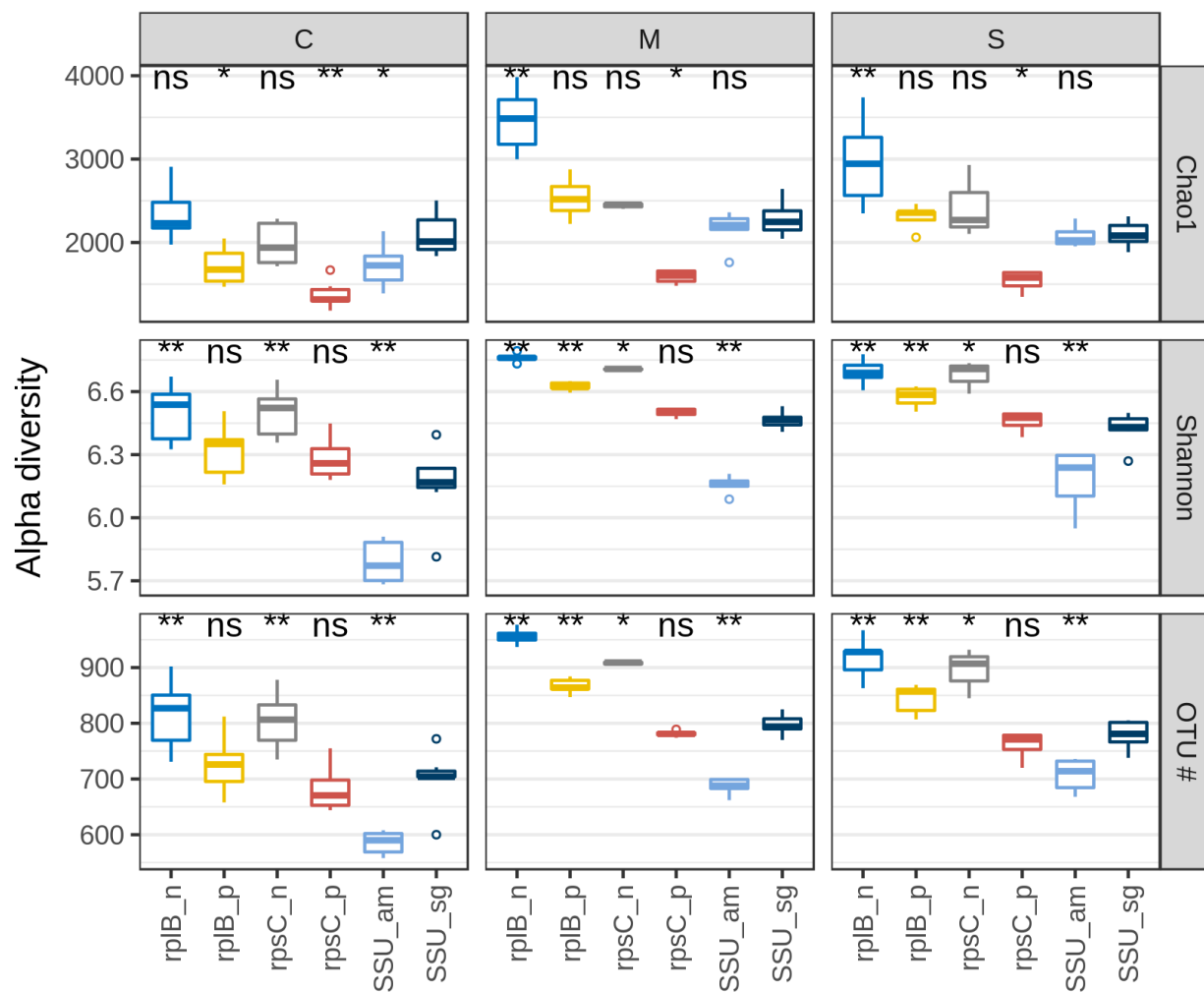
7

Figure 5: Comparison of the SSU rRNA gene, *rplB*, and *rpsC* in alpha diversity analyses (Chao1, Shannon, and OTU number by protein, "_p" and nucleotide, "_n") using large soil metagenome sequence. All genes show that the microbial community of the corn (C) rhizosphere has significantly less alpha diversity than those of *Miscanthus* (M) and switchgrass (S) except for *rpsC* and SSU rRNA gene with Chao1 index ($p < 0.01$). Wilcoxon test was used to compare SSU_sg against each of the other genes including rplB_n, rplB_p, rpsC_n, rpsC_p and SSU_am. For Chao1 index, rplB_n and rpsC_n show significantly higher abundance than SSU_sg in *Miscanthus* and switchgrass; For Shannon index and OTU number, rplB_n and rpsC_n show significantly higher abundance than SSU_sg in all three crops ("***" is $p < 0.001$, '**' is $p < 0.01$, '*' is $p < 0.05$, "ns" is $p > 0.05$).

9