

1 **Title: Genome wide association study reveals plant loci controlling heritability of the**  
2 **rhizosphere microbiome.**

3  
4 Authors: Siwen Deng<sup>1,2</sup>, Daniel Caddell<sup>2</sup>, Jinliang Yang<sup>3,4</sup>, Lindsay Dahlen<sup>1,5</sup>, Lorenzo Washington<sup>1</sup>,  
5 Devin Coleman-Derr<sup>1,2\*</sup>

6  
7 Affiliations:

8 <sup>1</sup>Department of Plant and Microbial Biology, University of California, Berkeley, CA, USA

9 <sup>2</sup>Plant Gene Expression Center, USDA-ARS, Albany, CA, USA

10 <sup>3</sup>Department of Agronomy and Horticulture, University of Nebraska-Lincoln, Lincoln, NE, USA

11 <sup>4</sup>Center for Plant Science Innovation, University of Nebraska-Lincoln, Lincoln, NE, USA

12 <sup>5</sup>Current affiliation: Department of Plant Sciences, University of California, Davis, CA, USA

13  
14 \*Author for correspondence:

15 Devin Coleman-Derr

16 Tel: 1-510-559-5911

17 Email: [colemanderr@berkeley.edu](mailto:colemanderr@berkeley.edu)

18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44

## 45 **Abstract**

46 Host genetics has recently been shown to be a driver of plant microbiome composition. However,  
47 identifying the underlying genetic loci controlling microbial selection remains challenging.  
48 Genome wide association studies (GWAS) represent a potentially powerful, unbiased method to  
49 identify microbes sensitive to host genotype, and to connect them with the genetic loci that  
50 influence their colonization. Here, we conducted a population-level microbiome analysis of the  
51 rhizospheres of 200 sorghum genotypes. Using 16S rRNA amplicon sequencing, we identify  
52 rhizosphere-associated bacteria exhibiting heritable associations with plant genotype, and identify  
53 significant overlap between these lineages and heritable taxa recently identified in maize.  
54 Furthermore, we demonstrate that GWAS can identify host loci that correlate with the abundance  
55 of specific subsets of the rhizosphere microbiome. Finally, we demonstrate that these results can  
56 be used to predict rhizosphere microbiome structure for an independent panel of sorghum  
57 genotypes based solely on knowledge of host genotypic information.

58  
59 **Keywords:** Rhizosphere, host genetics, microbiome, GWAS, heritability, amplicon sequencing,  
60 sorghum

## 61 **Introduction**

62 Recent work has shown that root-associated microbial communities are in part shaped by host  
63 genetics<sup>1-4</sup>. A study comparing the root microbiomes of a broad range of cereal crops has  
64 demonstrated a strong correlation between host genetic differences and microbiome composition<sup>4</sup>,  
65 suggesting that a subset of the plant microbiome may be influenced by host genotype across a  
66 range of plant hosts. In maize, these genotype-sensitive, or “heritable”, microbes are  
67 phylogenetically clustered within specific taxonomic groups<sup>5</sup>; however, it is unclear whether the  
68 increased genotype sensitivity in these lineages is unique to the maize microbiome or is common  
69 to other plant hosts as well.

70  
71  
72 Despite consistent evidence of the interaction between host genetics and plant microbiome  
73 composition, identifying specific genetic elements driving host-genotype dependent microbiome  
74 acquisition and assembly in plants remains a challenge. Recent efforts guided by *a priori*  
75 hypotheses of gene involvement have begun to dissect the impact of individual genes on  
76 microbiome composition<sup>6,7</sup>. However, these studies are limited to a small fraction of plant genes  
77 predicted to function in microbiome-related processes. Additionally, many plant traits expected to  
78 impact microbiome composition and activity, such as root exudation<sup>8</sup> and root system architecture<sup>9</sup>,  
79 are inherently complex and potentially governed by a very large number of genes. For these  
80 reasons, there is a need for alternative, large-scale and unbiased methods for identifying the genes  
81 that regulate host-mediated selection of the microbiome.

82  
83 Genome-wide association studies (GWAS) represent a powerful approach to map loci that are  
84 associated with complex traits in a genetically diverse population. Though pioneered for use in  
85 human genetics, to date the majority of GWAS have been conducted in plants<sup>10</sup>, and it has become  
86 an increasingly popular tool for studying the genetic basis of natural variation and traits of  
87 agricultural importance. When inbred lines are available, GWAS can be particularly useful; once  
88 genotyped, these lines can be phenotyped multiple times, making it possible to study many  
89 different traits in many different environments<sup>11</sup>. While GWAS is typically used in the context of  
90 a single quantitative phenotypic trait, analyses of multivariate molecular traits, such as

91 transcriptomic or metabolomic data, have also been conducted<sup>12,13</sup>. More recently, several attempts  
92 have been made to use host-associated microbiome census data as an input to GWAS, which in  
93 theory will allow for the identification of host genetic loci controlling microbiome composition<sup>14,15</sup>.

94  
95 In plants, a recent study in *Arabidopsis thaliana* used phyllosphere microbial community data as  
96 the phenotypic trait in a GWAS to demonstrate that plant loci responsible for defense and cell wall  
97 integrity affect microbial community variation<sup>16</sup>. Several other recent phyllosphere studies  
98 performed GWAS to identify genetic factors controlling microbiome associations with mixed  
99 degrees of success<sup>16-18</sup>. However, to our knowledge, use of GWAS in conjunction with the root  
100 associated microbiome has yet to be explored. In the context of the root microbiome, selection of  
101 sample type (rhizosphere or endosphere) and host system may be critical factors that determine  
102 the success of such effort. Previous work comparing the root microbiomes of diverse cereal crops  
103 have offered conflicting evidence as to whether host genotypic distance correlates most strongly  
104 with microbial communities distance within root endospheres or rhizospheres<sup>3,4</sup>. These data suggest  
105 that the sample type exhibiting the strongest correlation between genotype and microbiome  
106 composition may differ for each host, and that an initial evaluation of the degree of correlation  
107 between genotype and microbiome phenotype across sample types may be informative.

108  
109 In the context of the root microbiome, we propose *Sorghum bicolor* (L.) as an ideal plant system  
110 for GWAS-based dissection of host-genetic control of microbiome composition. Sorghum is a  
111 heavy producer of root exudates<sup>19</sup>, and the sorghum microbiome has been shown to house an  
112 unusually large number of host-specific microbes<sup>4</sup>. Additionally, there is a wide range of natural  
113 adaptation in traditional sorghum varieties from across Africa and Asia, and a collection of  
114 breeding lines generated from U.S. sorghum breeding programs, both of which provide a rich  
115 source of phenotypic and genotypic variation<sup>20</sup>. Several genome sequences of sorghum varieties  
116 have been completed, and variation in nucleotide diversity, linkage disequilibrium, and  
117 recombination rates across the genome have been quantified<sup>21</sup>, providing an understanding of the  
118 genomic patterns of diversification in sorghum. Finally, sorghum is an important cereal crop grown  
119 throughout the world as a food, feedstock, and biofuel, enabling direct integration of resulting  
120 discoveries into an agriculturally-relevant system.

121  
122 In this study, we dissect the host-genetic control of bacterial microbiome composition in the  
123 sorghum rhizosphere. Using 16S rRNA sequencing, we profiled the microbiome of a panel of 200  
124 diverse genotypes of field grown sorghum. We aim to demonstrate that a large fraction of the plant  
125 microbiome responds to host genotype, and that this subset shares considerable overlap with  
126 lineages shown to be susceptible to host genetic control in another plant host. Additionally, we  
127 tested whether GWAS can be used to identify specific genetic loci within the host genome that are  
128 correlated with the abundance of specific heritable lineages, and whether differences in  
129 microbiome composition can be predicted solely from genotypic information. Collectively, this  
130 work demonstrates the utility of GWAS for analysis of host-mediated control of rhizosphere  
131 microbiome phenotypes.

## 132 **Results**

133 **Diverse sorghum germplasm show rhizosphere is ideal for microbiome-based GWAS.** In this  
134 study, the relationship between host genotype and microbiome composition was explored through  
135 a field experiment involving 200 genotypes selected from the Sorghum Association Panel (SAP)  
136 germplasm collection<sup>20</sup> (Supplemental Table 1). As prior studies suggest that the strength of the  
137 correlation between host genotype and microbiome composition may vary by sample type in a  
138 host-dependent manner<sup>34</sup>, we first sought to determine whether leaf, root, or rhizosphere samples  
139 were most suitable for downstream GWAS in sorghum. Using a subset of 24 genotypes from our  
140 collection of 200 (Figure 1a, Supplemental Table 1), the microbiome composition of leaf, root,  
141 and rhizosphere sample types was analyzed using paired-end sequencing of the V3–V4 region of  
142 the ribosomal 16S rRNA on the Illumina MiSeq platform (Illumina Inc., San Diego, CA, USA).  
143 The resulting dataset demonstrated comparatively high levels of microbial diversity within both  
144 root and rhizosphere samples (Figure 1b) and strong clustering of above and below ground sample  
145 types (Figure 1c). Three independent Mantel's tests (9,999 permutations) were used to evaluate  
146 the degree of correlation between host genotypic distance and microbiome composition for leaf,  
147 root, and rhizosphere sample types (Figure 1d); of the three compartments, only rhizosphere  
148 exhibited a significant Mantel's correlation ( $R^2=0.13$ ,  $Df=1$ ,  $p=0.02$ ). Based on these results,  
149 subsequent investigation of the microbiomes of the full panel of 200 lines, including heritability  
150 and GWAS analyses, was performed using rhizosphere samples.

151  
152 To investigate host genotype dependent variation in the sorghum rhizosphere microbiome, the  
153 rhizospheres of 600 field grown plants (including three replicates of each of 200 genotypes) were  
154 profiled using V3-V4 16S rRNA amplicon sequencing. After removing rare OTUs with less than  
155 3 reads in at least 20% of the samples and normalizing to an even read depth of 18,000 reads per  
156 sample, the data set included 1,189 high-abundance OTUs representing 29 bacterial phyla.  
157 Compositional analysis of the resulting microbiome dataset exhibited profiles consistent with  
158 recent microbiome studies involving the sorghum rhizosphere<sup>12,22,23</sup> from a variety of field sites, with  
159 Proteobacteria, Actinobacteria and Acidobacteria comprising the top three dominant phyla  
160 (Supplemental Figure 1).

161  
162 **Sorghum and maize rhizospheres exhibit strong overlap in heritable taxa.** A recent study of  
163 two separate maize microbiome datasets suggests that specific bacterial lineages are more sensitive  
164 to the effect of host genotype than others<sup>5</sup>. To determine if a bacterial lineage's responsiveness to  
165 host genetics is a trait conserved across different plant hosts that diverged more than 11 million  
166 years ago<sup>24</sup>, the broad sense heritability ( $H^2$ ) of individual OTUs in our sorghum dataset was  
167 evaluated.  $H^2$ , which quantifies the proportion of variance that is explained by genetic rather than  
168 environmental effects, ranged from 0 to 66% for individual OTUs (Supplemental Table 2). By  
169 comparison,  $H^2$  for individual OTUs in the first of two experiments across 27 inbred maize lines  
170 had a maximum of 23% (performed in 2010), while the second exhibited a maximum of 54%  
171 (performed in 2015)<sup>5</sup>.

172  
173 To explore whether microbes with high heritability in the sorghum dataset are phylogenetically  
174 clustered, we partitioned the 1,189 OTUs into heritable ( $n=347$ ) and non-heritable fractions  
175 ( $n=842$ ) using an  $H^2$  cutoff score of 0.15 (Figure 2a, Supplemental Table 3). Several bacterial  
176 orders, including Verrucomicrobiales, Flavobacteriales, Planctomycetales, and Burkholderiales,  
177 were observed to have significantly greater numbers of OTUs that are heritable, as compared to

178 the non-heritable OTU fraction (Fisher's exact test,  $p < 0.05$ , Figure 2a, Supplemental Table 3).  
179 Notably, all 6 Flavobacteriales OTUs were present in the heritable fraction (Figure 2b); by  
180 contrast, 40 other bacterial orders were only observed within the non-heritable fraction. Another  
181 bacterial order, Bacillales, contained a smaller number of OTUs in the heritable than non-heritable  
182 fraction, but the percentage of read counts attributable to its heritable OTUs was approximately  
183 eight-fold greater than those in the non-heritable fraction, suggesting that its heritable members  
184 are abundant organisms within the rhizosphere (Figure 2b). Collectively, these data demonstrate  
185 that a specific subset of bacterial lineages are enriched for members susceptible to host genotypic  
186 selection.

187  
188 We hypothesized that despite the considerable evolutionary distance between maize and sorghum,  
189 the bacterial lineages containing OTUs most responsive to host genotypic effects in maize would  
190 likely also contain OTUs exhibiting such susceptibility within sorghum. To test this, we compared  
191 the top 100 most heritable OTUs from both maize datasets (referred to as NAM 2010 and NAM  
192 2015) and the sorghum dataset described above, resulting in a combined dataset of 300 OTUs  
193 spanning 65 bacterial orders. After removing bacterial orders not observed in the sorghum dataset  
194 ( $n=18$ ), we noted that more than half were observed in at least two of the datasets, and  
195 approximately one third ( $n=15$ ) contained heritable OTUs in all three datasets (Figure 3a). To  
196 determine if this overlap was significantly greater than is expected by chance, we performed  
197 permutational resampling of 10,000 sets of randomly chosen sorghum OTUs for comparison.  
198 Notably, we found that the overlap between the heritable sorghum fraction with both the individual  
199 maize heritable fractions and the combined heritable maize OTUs to be significant, compared with  
200 the resampled sorghum OTUs (NAM 2010  $n=17$ ,  $p=0.0099$ , NAM 2015  $n=19$ ,  $p=0.0016$ ,  
201 combined  $n=15$ ,  $p=0.0344$ )(Figure 3a). Collectively, these results demonstrate that there is a  
202 conservation between the bacterial orders most sensitive to genotype across both maize and  
203 sorghum.

204  
205 In an effort to identify the bacterial lineages with the greatest propensity for high heritability, we  
206 calculated the number of heritable OTUs in each of the shared heritable bacterial orders identified  
207 above. We noted that among bacterial orders containing the greatest number of heritable OTUs  
208 across all three datasets were several that represent large lineages frequently observed within the  
209 root microbiome; (e.g. Actinomycetales) (Figure 3b). We hypothesized that this result is likely  
210 driven in part by the overall frequency of these lineages within the rhizosphere microbiome, with  
211 more common lineages resulting in a greater fraction of heritable microbes due to their ubiquity.  
212 To help account for this, we normalized the frequency of heritable sorghum OTUs ( $n=100$ ) by  
213 total sorghum OTU counts ( $n=1,189$ ) belonging to each order (Figure 3c, Supplemental Table 4).  
214 These results demonstrate that while the prevalence of Actinomycetales and Myxococcales among  
215 heritable microbes is consistent with their general prevalence in the overall dataset,  
216 Burkholderiales and two other lineages, including the Verrucomicrobia and Planctomycetes,  
217 exhibited a significant enrichment (Fisher's exact test,  $p < 0.001$ ) in the heritable fraction not  
218 expected to be influenced by abundance alone.



219 **Genome-wide association reveals genetic loci correlated with rhizosphere microbial**  
220 **abundance.** Recent work in the leaf microbiome has demonstrated the potential utility of GWAS  
221 for uncovering host loci correlated with microbiome composition<sup>18</sup>. Here, we sought to use GWAS  
222 with rhizosphere microbiome datasets using both global properties of the OTU dataset and the  
223 abundances of individual OTUs. For overall community composition, a subset of principal  
224 components (PCs) were selected from an analysis of the abundance patterns of the 1,189 OTUs.  
225 To prioritize individual PCs for inclusion in our GWAS analysis, we determined the heritability  
226 scores of each of the top ten PCs, which explained 75% of the total variance in our dataset  
227 (Supplemental Figure 2a). PCs with  $H^2$  equal to or greater than 0.25 (PC1, PC3, PC5, PC9, and  
228 PC10, Supplemental Figure 2a) were subjected to GWAS (Supplemental Figure 2b). The GWAS  
229 analysis performed for PC1, which explained 21% percent of total variance and had the second  
230 highest heritability ( $H^2=0.35$ ), revealed a significant correlation between community composition  
231 and a locus of approximately 1.15 Mb on chromosome 4 with a moderately stringent threshold of  
232  $-\log_{10}(p=10^{-4})$  (Figure 4a, Supplemental Figure 2b). Additionally, GWAS analyses that used PC5  
233 and PC10 as inputs, both revealed an identifiable peak on chromosome 6, though it was slightly  
234 below the threshold of significance (Supplemental Figure 2b).

235  
236 As principal components are derived from linear combinations of the abundance of individual  
237 OTUs within the dataset, it is unclear whether the correlations observed on chromosomes 4 and 6  
238 are driven by one common or two different sets of microbial lineages. To address this, we  
239 performed separate GWAS analyses using the abundances of each single OTU in our dataset as  
240 input (Figure 4b, Supplemental Figure 2c). From these analyses, we identified two distinct sets of  
241 39 and 10 OTUs with significant correlations with the loci on chromosomes 4 and 6, respectively,  
242 and only a single OTU belonging to the order Burkholderiales that was shared between the two  
243 loci (Supplemental Figure 2c), demonstrating that different sorghum loci influence the abundance  
244 patterns of different groups of microbes.

245  
246 To explore the relationship between the identified peak on chromosome 4 (Figure 4a) and the  
247 bacterial taxa with significant GWAS correlations at this locus (Figure 4b), we first sought to  
248 understand how relative abundance for these 40 OTUs varied across the sorghum panel. An  
249 analysis of the SNP data at this locus revealed two allele groups, the major allele containing 343  
250 sorghum genotypes and the minor allele containing 14 genotypes. Next, we observed that the  
251 majority of OTUs that were more prevalent in sorghum genotypes containing the major allele  
252 belonged to monoderm lineages, while the majority of OTUs more prevalent in the minor allele  
253 group belonged to diderm lineages (Figure 4b), suggesting that host genetic mechanisms at this  
254 locus are interacting with basal bacterial traits.

255  
256 To explore which genetic mechanisms might be driving the correlations observed on Chromosome  
257 4, we examined tissue specific expression patterns from publicly available RNA-Seq datasets  
258 obtained from phytozome v12.1<sup>25</sup> for all 27 genes in the 1.15 Mb interval (Figure 4c, Supplemental  
259 Table 5). Of these candidates, several were observed to exhibit strong root specific expression  
260 patterns, including three annotated candidates: gamma carbonic anhydrase-like 2, a putative Beta-  
261 1,4 endoxylanase, and disease resistance protein RGA2 (Figure 4c).

262  
263 **Sorghum genotypic data can predict microbiome composition.** To validate that allelic variation  
264 at the candidate locus on chromosome 4 contributes to differences in rhizosphere composition, we

265 conducted a follow up experiment with eighteen additional sorghum lines, including genotypes  
266 not present in the original study. To help disentangle phylogenetic-relatedness from locus-specific  
267 effects, we selected sorghum genotypes that spanned the diversity panel; additionally, for each  
268 minor allele genotype (n=9), we included a phylogenetically related major allele line (n=9) (Figure  
269 1a). Following two weeks of growth in a mixture of calcined clay and field soil in the growth  
270 chamber, we collected the rhizosphere microbiomes of each genotype and microbiome  
271 composition was analyzed using 16S rRNA amplicon sequencing as in the main study. A canonical  
272 analysis of principal coordinates (CAP) ordination constrained on genotypic group separated the  
273 rhizospheres of genotypes belonging to major and minor allele groups into distinct clusters (Figure  
274 5a, PERMANOVA  $F=2.66$ ,  $Df=1$ ,  $p=0.0061$ ), with genotype explaining approximately 7.5% (CAP1)  
275 of variance in the dataset.

276  
277 To identify which taxa drive the clustering observed in our CAPs analysis, and to compare this to  
278 taxa responsive to the chromosome 4 allele group in our main experiment, we performed an  
279 indicator species analysis on the validation dataset. A comparison of the significant indicator  
280 OTUs ( $p<0.05$ ) from each allele group in the validation dataset (n=65) demonstrated similar trends  
281 in abundance of indicator OTUs as observed in the main experiment (Figure 4b), with OTUs  
282 belonging to monoderm and diderm lineages enriched in the major and minor allele-containing  
283 lines, respectively. Interestingly, while most diderm lineages were more prevalent in the minor  
284 allele-containing lines, several diderm lineages including Gemmatimonadales, Acidobacteriales,  
285 and Sphingobacteriales contained OTUs that were more abundant within major allele lines.  
286 Notably, this pattern was observed in both the main experiment (Figure 4b) and validation  
287 experiment (Figure 5b). Collectively, this experiment supports the findings of our main  
288 experiment, in which allelic variation at a locus located on chromosome 4 was shown to correlate  
289 with the abundance of specific bacterial lineages.

290

## 291 Discussion

292 **Host selection of plant rhizosphere microbiomes.** Previous GWAS of plant-associated  
293 microbiome traits have often been conducted with leaf samples, and have not always been  
294 successful in identifying loci that correlate with microbiome phenotypes<sup>16-18</sup>. In this study, we  
295 compared the overall correlation between host genotype and bacterial microbiome distances across  
296 leaf, root, and rhizosphere of *Sorghum bicolor*, and demonstrate that of the three, the rhizosphere  
297 represents the most promising compartment for conducting experiments to untangle the heritability  
298 of the sorghum microbiome. Notably, the degree of correlation between sorghum phylogenetic  
299 distance and microbiome distance was highest in the rhizosphere and lowest in the leaves. This  
300 greater correlation observed in the rhizosphere could be in part due to the phyllosphere's relative  
301 compositional simplicity. Even *Arabidopsis* rosette leaves, which are in close proximity to soil,  
302 harbor a distinct and relatively simple bacterial community compared to the root<sup>26</sup>. By contrast, the  
303 rhizosphere represents a highly diverse and populated subset of the soil microbiome, and  
304 potentially offers a greater pool of microbes upon which the host may exert influence<sup>27</sup>.  
305 Alternatively, the rhizosphere's greater correlation with microbiome composition could be caused  
306 by the plant's relatively weaker ability to select epiphytes in its aboveground microbiome; while  
307 the arrival of phyllosphere colonists is largely thought to be driven by wind and rainfall dispersal<sup>28</sup>,  
308 root exudation is known to control chemotaxis and other colonization activities of select members  
309 of the surrounding soil environment. This provides an additional mechanism for host selection of  
310 its microbial inhabitants prior to direct interaction with the plant surface<sup>28,29,30</sup>. It is worth noting that

311 sorghum is known to be an atypically strong producer of root exudates<sup>19</sup>, and consequently it is  
312 possible that other plant hosts may demonstrate the greatest selective influence within tissues other  
313 than the rhizosphere. Future efforts to investigate host control of the microbiome through GWAS  
314 or related techniques would benefit from careful selection of sample type following pilot studies  
315 designed to explore heritability across different host tissues.

316 **Heritable rhizosphere microbes are phylogenetically clustered and similar across hosts.**

317 Within the rhizosphere, we demonstrate that microbiome constituents vary in broad sense  
318 heritability, and heritable taxa show a strong overlap with heritable lineages identified in maize,  
319 spanning fifteen different bacterial orders<sup>5</sup>. In particular, three of these orders, Verrucomicrobiales,  
320 Burkholderiales, and Planctomycetales were significantly enriched in the heritable fraction of our  
321 dataset. As members of Burkholderiales can form symbioses with both plant and animal hosts<sup>31,32</sup>,  
322 and some colonize specific members of a host genus or species<sup>33</sup>, it is feasible that such strong  
323 relationships necessitated additional genetic discrimination between hosts. Within *Burkholderia*  
324 spp., this could be facilitated by their relatively large pan-genome, with diversity driven by large  
325 multi-replicon genomes and abundant genomic islands<sup>34</sup>.

326  
327 These observations suggest that evaluating bacterial heritability may identify new lineages for  
328 which close or symbiotic but previously undetected associations with plant hosts exist. For  
329 example, we observed several lineages with high heritability that are common in soil, yet prior  
330 evidence of plant-microbe interactions in the literature is lacking, including Verrucomicrobiales  
331 and Planctomycetales. Interestingly, heritability in these lineages might be facilitated by the  
332 presence of a recently discovered shared bacterial microcompartment gene cluster present in both  
333 Planctomycetes and Verrucomicrobia, which confers the ability to degrade certain plant  
334 polysaccharides<sup>35</sup>. Indeed, microbiome composition is known to be driven in part by variations in  
335 polysaccharide containing sources including plant cell wall components and root exudates<sup>36</sup>.  
336 Additional experimentation with bacterial mutants lacking this genetic cluster could be useful for  
337 revealing its role in shaping plant microbe interactions.

338 **Sorghum loci are responsible for controlling the rhizobiome.** Our GWAS correlated host  
339 genetic loci and the abundance of specific bacteria within the host microbiome, as well as overall  
340 rhizosphere community structure. To our knowledge, this is the first example of such work in a  
341 crop rhizosphere. We identified two loci with strong associations with the microbiome structure.  
342 The most significant maps to a locus on chromosome 4 containing several candidate genes with  
343 root specific expression.

344  
345 One candidate gene located near the center of this locus encodes a beta 1,4 endo xylanase.  
346 Xylanases are responsible for the degradation of xylan into xylose, and are one of the primary  
347 catabolizers of hemicellulose, a major component of the plant cell wall<sup>37</sup>. As a result, beta 1,4 endo  
348 xylanases may play a role in shaping the degree of plasticity in the barrier between the root and  
349 surrounding rhizosphere environments, in turn influencing the release of cell wall or apoplast  
350 derived metabolites into the rhizosphere environment<sup>38</sup>. Alternatively, altered xylanase activity  
351 could lead to shifts in carbohydrate profiles within the cell wall, leading to heightened plant  
352 immune responses<sup>39,40</sup>; the catabolic byproducts of microbially-produced xylanase used in pathogen  
353 invasion are in part responsible for triggering innate immune responses in plants, and various



354 components of the plant immune signalling network have been shown to influence microbiome  
355 structure<sup>6,7</sup>.

356

357 Another candidate gene within the chromosome 4 locus, that also displays root-specific  
358 expression, is predicted to encode gamma carbonic anhydrase-like 2. In plants, carbonic  
359 anhydrases (CA) participate in aerobic respiration, and facilitate the reversible hydration of CO<sub>2</sub> to  
360 bicarbonate<sup>41,42</sup>. Previous studies have implicated CA activity in plant-microbe interactions<sup>43</sup>; an  
361 important role for CA was first observed in root nodules of legumes inoculated with *Rhizobium*<sup>44,45</sup>.  
362 CAs have since been implicated in disease resistance as well, having both antioxidant activity and  
363 salicylic acid binding capability<sup>46-48</sup>. Collectively, these studies suggest that a loss or alteration of  
364 function of CA could impact the composition of the rhizosphere microbiome. Future validation  
365 experiments using genetic mutants within this and other candidate genes can be used to help  
366 elucidate the underlying genetic element(s) responsible for modulation of the rhizosphere  
367 microbiome.

### 368 **Conclusion**

369 Although the underlying host genetic causes of shifts in the microbiome are not well understood,  
370 candidate driven approaches have implicated disease resistance<sup>6,7</sup>, nutrient status<sup>7,49,50</sup>, sugar  
371 signaling<sup>51</sup>, and plant age<sup>32,53</sup> as major factors. Non-candidate approaches to link host genetics and  
372 microbiome composition, such as GWAS, have the potential to discover novel mechanisms that  
373 can be added to this list. Here we show that GWAS can predict microbiome structure based on  
374 host genetic information, building on previous studies that have observed inter- and intra-species  
375 variation in microbiomes<sup>1,4,5,16,36,54-56</sup>. Collectively, our study adds to a growing list of evidence that  
376 genetic variation within plant host genomes modulates their associated microbiome. We anticipate  
377 that GWAS of plant microbiome association will promote a comprehensive understanding of the  
378 host molecular mechanisms underlying the assembly of microbiomes and facilitate breeding  
379 efforts to promote beneficial microbiomes and improve plant yield.

380

## 381 **Methods**

382 **Germplasm selection.** In order to ensure that microbiome profiling was performed on a  
383 representative subset of the broad genetic diversity present in the 378 member Sorghum  
384 Association Panel (SAP)<sup>30</sup>, subsets of 200 genotypes were randomly sampled from the SAP 10,000  
385 times and an aggregate nucleotide diversity score was calculated for each using the R package  
386 “PopGenome”<sup>37</sup>. From these data, the subset of 200 lines with the maximum diversity value was  
387 selected (Figure 1a, Supplemental Table 1). For the pilot experiment used to determine the  
388 appropriate sample type for GWAS, a subset of 24 lines was selected that included genotypes from  
389 a wide range of phylogenetic distances (Figure 1a, Supplemental Table 1). The phylogenetic tree  
390 of sorghum accessions was generated using the online tool: Interactive Tree Of Life (iTOL) v5<sup>38</sup>.

391 **Field experimental design and root microbiome sample collection.** The experimental field used  
392 in this study is an agricultural field site located in Albany, California (37.8864°N, 122.2982°W),  
393 characterized by a silty loam soil with pH 5.2<sup>4</sup>. Germplasm for the US SAP panel used in this  
394 study<sup>30</sup> were obtained from GRIN ([www.ars-grin.gov](http://www.ars-grin.gov)). To ensure a uniform starting soil  
395 microbiome for all sorghum seedlings and to control their planting density, seeds were first sown  
396 into a thoroughly homogenized field soil mix in a growth chamber with controlled environmental  
397 factors (25 °C, 16hr photoperiods) followed by transplantation to the field site. To prepare the soil  
398 for seed germination, 0.54 cubic meters of soil was collected at a depth of 0 to 20 cm from the  
399 field site subsequently used for planting, and homogenized by separately mixing 4 equally sized  
400 batches with irrigation water in a sterilized cement mixer followed by manual homogenization on  
401 a sterilized tarp surface. Soil was then transferred to sterilized 72-cell plant trays. To prepare seeds  
402 for planting, seeds were surface-sterilized through soaking 10 min in 10% bleach + 0.1% Tween-  
403 20, followed by 4 washes in sterile water. Following planting, sorghum seedlings were watered  
404 with approximately 5 ml of water using a mist nozzle every 24 hrs for the first three days, and  
405 bottom watered every three days until the 12th day, then transplanted to the field.

406  
407 The field consisted of three replicate blocks, with each block containing 200 plots for each of 200  
408 selected genotypes. Six healthy sorghum seedlings of each genotype were transplanted to their  
409 respective plots, separated by 15.2cm, and thinning to three seedlings per plot was performed at  
410 two weeks post transplanting. Plots were organized in an alternating pattern with respect to the  
411 irrigation line to maximize the distance between each plant (Supplemental Figure 3). Plants were  
412 watered for one hour, three times per week, using drip irrigation with 1.89 L/hour rate flow  
413 emitters. Manual weeding was performed three times per week throughout the growing season. To  
414 ensure that the genotypes were at a similar stage of development and that the host-associated  
415 microbiome had sufficient time to develop, collection of plant-associated samples was performed  
416 nine weeks post germination. Only the middle plant within each plot was harvested to help mitigate  
417 potential confounding plant-plant interaction effects resulting from contact with roots from  
418 neighboring plants of other genotypes. Rhizosphere, leaf, and root samples were collected as  
419 described previously<sup>39</sup>.

420  
421 **DNA extraction, PCR amplification, and Illumina sequencing.** DNA extractions, PCR  
422 amplification of the V3-V4 region of the 16S rRNA gene, and amplicon pooling were performed  
423 as described previously<sup>39</sup>. In brief, DNA extractions for all samples were performed using  
424 extraction kits (MoBio PowerSoil DNA Isolation Kit, MoBio Inc., Carlsbad, CA) following the  
425 manufacturer’s protocol. Amplification of the V3-V4 region of the 16S rRNA gene was performed

426 using dual-indexed 16s rRNA Illumina iTags primers 341F (5'-CCTACGGGNBGCASCAG-3')  
427 and 785R (5'-GACTACNVGGGTATCTAATCC-3'). An aliquot of the pooled amplicons was  
428 diluted to 10 nM in 30µL total volume before submitting to the QB3 Vincent J. Coates Genomics  
429 Sequencing Laboratory facility at the University of California, Berkeley for sequencing using  
430 Illumina Miseq 300bp pair-end with v3 chemistry. Sequences were returned demultiplexed, with  
431 adaptors removed.

432

433 **Amplicon sequence processing, OTU classification, and taxonomic assignment.** Sequencing  
434 data were analyzed using the iTagger pipeline to obtain OTUs<sup>60</sup>. In brief, after filtering 81,416,218  
435 16S rRNA raw reads for known contaminants (Illumina adapter sequence and PhiX), primer  
436 sequences were trimmed from the 5' ends of both forward and reverse reads. Low-quality bases  
437 were trimmed from the 3' ends prior to assembly of forward and reverse reads with FLASH<sup>61</sup>. The  
438 remaining 66,524,451 high-quality merged reads were clustered with simultaneous chimera  
439 removal using UPARSE<sup>62</sup>. After clustering, 37,867,921 read counts mapped to operational  
440 taxonomic units (OTUs) at 97% identity (Supplemental Table 6). Taxonomies were assigned to  
441 each OTU using the RDP Naïve Bayesian Classifier with custom reference databases<sup>63</sup>. For the 16S  
442 rRNA V3-V4 data, this database was compiled from the May 2013 version of the GreenGenes 16S  
443 database v13, trimmed to the V3-V4 region. After taxonomies were assigned to each OTU, OTUs  
444 were discarded if they were not assigned a Kingdom level RDP classification score of at least 0.5,  
445 or if they were not assigned to Kingdom Bacteria, which yielded 10,006 OTUs. In the downstream  
446 analyses, we removed low abundance OTUs because in many cases they are artifacts generated  
447 through the sequencing process. Samples with low read counts were also removed. To account for  
448 differences in sequencing read depth across samples, all samples were normalized to an even read  
449 depth of reads per sample random subsampling for specific analyses, or alternatively, by dividing  
450 the reads per OTU in a sample by the sum of usable reads in that sample, resulting in a table of  
451 relative abundance frequencies.

452

453 **Estimates of broad sense heritability of OTU abundance in rhizosphere.** To calculate the  
454 broad-sense heritability ( $H^2$ ) for individual OTU abundances, we fitted the following linear mixed  
455 model to OTU abundances of each individual OTU ( $n=1,189$ ) following a cumulative sum scaling<sup>64</sup>  
456 normalization procedure that adjusted for differences in sequencing depth and fit a normal  
457 distribution:

458

$$459 Y_{ijk} = u + G_i + R_j + B_{jk} + e$$

460

461 In this model for a given OTU,  $Y_{ijk}$  denotes the OTU abundance of the  $i^{\text{th}}$  genotype evaluated in the  
462  $k^{\text{th}}$  block of the  $j^{\text{th}}$  replicate;  $u$  denotes the overall mean;  $G_i$  is the random effect of the  $i^{\text{th}}$   
463 genotype;  $R_j$  is the random effect of the  $j^{\text{th}}$  replicate;  $B_{jk}$  is the random effect of the  $k^{\text{th}}$  block  
464 nested within the  $j^{\text{th}}$  replicate;  $e$  denotes the residual error. To account for the spatial effects in  
465 the field, additional spatial variables were fitted as random effects using 2-dimensional splines in  
466 the above model using an R add-on package “sommer”<sup>65</sup>.  $H^2$  was estimated as the amount of  
467 variance explained by the genotype term ( $V_G$ ) relative to the total variance ( $V_G + V_E/j$ ). Here  $j$  is the  
468 number of replications. To get the null distribution of  $H^2$ , each OTU was randomly shuffled 1,000  
469 times and then fitted to the same model as described above. Permutation p-value was calculated as  
470 the probability of the permuted  $H^2$  values bigger than the observed  $H^2$  value.

471

472 **Comparative analysis of heritable taxa between sorghum and maize datasets.** To identify the  
473 degree to which heritable taxa were shared between maize and sorghum, we compared the top 100  
474 most heritable OTUs from both maize datasets (referred to as NAM 2010 and NAM 2015) and the  
475 sorghum dataset generated in this study, resulting in a combined dataset of 300 OTUs spanning 65  
476 bacterial orders. As these three experiments were conducted at different field sites, a subset of the  
477 orders (n=18) containing heritable OTUs in the maize dataset were not detected in either the  
478 heritable or non-heritable fractions of the sorghum dataset and were excluded from subsequent  
479 comparative analyses. Of the remaining bacterial orders represented by these heritable OTUs, we  
480 determined the number (n=26) that contained heritable OTUs in at least two of the datasets, and  
481 the number (n=15) that contained heritable OTUs in all three datasets (Figure 3a). To evaluate  
482 whether the degree of overlap in heritable lineages is greater than what would be expected by  
483 chance, we performed a permutation test (n=10,000) in which we resampled 100 random OTUs  
484 from the 1,189 total sorghum OTUs and recomputed intersections with the two maize datasets. P-  
485 values are reported as the number of instances that these permutations returned a greater degree of  
486 overlap in these permutations divided by total number of permutations.

487  
488 **GWAS.** For each OTU, GWAS was conducted separately using the best linear unbiased predictors  
489 (BLUPs) obtained from the linear mixed model. Population structure was accounted for using  
490 statistical methods that allow us to detect both population structure (Q) and relative kinship (K) to  
491 control spurious association. The Q model ( $y = S\alpha + Qv + e$ ), the K model ( $y = S\alpha + Zu + e$ ), and  
492 the Q + K model ( $y = X\beta + S\alpha + Qv + Zu + e$ ) described previously<sup>66</sup>, were used in our study. In  
493 the model equations, y is a vector of phenotypic observation;  $\alpha$  is a vector of allelic effects; e is a  
494 vector of residual effects; v is a vector of population effects;  $\beta$  is a vector of fixed effects other  
495 than allelic or population group effects; u is a vector of polygenic background effects; Q is the  
496 matrix relating y to v; and X, S, and Z are incidence matrices of 1s and 0s relating y to  $\beta$ ,  $\alpha$ , and  
497 u, respectively. To account for the population structure and genetic relatedness, the first three  
498 principal components (PCs) and kinship matrix were calculated using the SNPs obtained from<sup>21</sup>  
499 and fitted into the MLM-based GWAS pipeline for each OTU using GEMMA<sup>67</sup>.

500 **GWAS validation experiment.** For the GWAS validation experiment, the 378 genotypes of the  
501 SAP were first subset into lines containing the major (n=343) and minor (n=14) allele for the two  
502 haplotypes found at the peak on chromosome 4 described in the text. Including the 178 genotypes  
503 not selected for the GWAS, a total of nine sorghum genotypes belonging to the minor allele were  
504 selected, with an effort to include genotypes spanning the phylogenetic tree. For each of these nine  
505 minor allele lines, another genotype containing the major allele with close overall genetic  
506 relatedness was selected, resulting in nine major and nine minor allele containing lines. Two  
507 replicates of each line were grown in growth chambers (33°C/28°C, 16h light/ 8h dark, 60%  
508 humidity) in a 10% vermiculite/ 90% calcined clay mixture rinsed with a soil wash prepared from  
509 a 2:1 ratio of field soil to water from the field site used in the GWAS. Plants were watered daily  
510 with approximately 5 ml of autoclaved Milli-Q water using a spray bottle for the first three days,  
511 followed by top watering with 15 ml of water every three days. An additional misting was  
512 performed to the soil surface every 24 hrs to prevent drying. Following two weeks of growth,  
513 plants were harvested and rhizosphere microbiomes extracted as described for the field  
514 experiment.

515

516 **Microbiome statistical analyses.** All statistical analyses of the amplicon datasets were performed  
517 in R using the normalized reduced dataset, unless stated otherwise. For alpha-diversity  
518 measurement, Shannon's Diversity was calculated as  $e^X$ , where X is Shannon's Entropy as  
519 determined with the diversity function in the R package `vegan`<sup>68</sup>. Principal coordinate analyses were  
520 performed with the function `pcoa` in the R package `ape`<sup>69</sup>, using the Bray-Curtis distance obtained  
521 from function `vegdist` in the R package `vegan`<sup>68</sup>. Mantel's tests were used to determine the  
522 correlation between host phylogenetic distances and microbiome distances using the `mantel`  
523 function in the R package `vegan`<sup>68</sup> with 9,999 permutations, and using Spearman's correlations to  
524 reduce the effect of outliers. Indicator species analyses were performed using the function `indval`  
525 in the R package `labdsv`<sup>70</sup>, with p-values based on permutation tests run with 10,000 permutations.  
526 To account for multiple testing performed for all 430 genera in our dataset, multiple testing  
527 correction was performed with an FDR of 0.05 using the `p.adjust` function in the base R package  
528 `stats`. Canonical Analysis of Principal Coordinates (CAP) was performed for the final validation  
529 experiment to test the amount of variance explained by genotypic group using the `capscale`  
530 function in the R package `vegan`<sup>68</sup>; an ANOVA like permutation test using the sum of all  
531 constrained eigenvalues was performed to determine the percent variance explained by each factor  
532 using the function `anova.cca` in the R package `vegan`<sup>68</sup>.

533  
534 **Analysis of sorghum RNA-seq datasets.** Publicly available sorghum RNA-Seq data for 27  
535 annotated genes in the 1.15 Mb interval of chromosome 4 (Sobic.004G153000 -  
536 Sobic.004G155900), were downloaded from phytozome v12.1<sup>25</sup> (Figure 4c, Supplemental Table  
537 5). Expression datasets were broadly grouped based on the tissue-type from which they were  
538 derived (root, leaf, or reproductive). To aid in the visualization of tissue specific expression of  
539 genes exhibiting large differences in absolute levels of gene expression, we normalized the  
540 Fragments Per Kilobase of transcript per Million mapped reads (FPKM) values for each gene in  
541 each tissue type by dividing by the average value of gene expression for that gene across all tissue  
542 types. We defined root-specific expression as genes that had a normalized FPKM less than 1 in no  
543 more than two root datasets, and a normalized FPKM greater than 1 in no more than two datasets  
544 of other tissue types (Figure 4c, Supplemental Table 5).

545  
546 **Data availability.** All datasets and scripts for analysis are available through github  
547 (<https://github.com/colemanderr-lab/Deng-2020>) and all short read data has been submitted to the  
548 NCBI SRA.

## 549 550 **Figure legends**

551 **Figure 1. Sample type and population selection.** **A** Phylogenetic tree representing the 378  
552 member sorghum association panel (SAP, inner ring), the subset of 200 lines selected for GWAS  
553 (2nd ring from the center, in blue), the 24 lines used for sample type selection (Pilot, 3rd ring from  
554 the center, in yellow), and the 18 genotypes used for GWAS validation containing either the  
555 Chromosome 4 minor allele (red) or major allele (brown) identified by GWAS (outer ring). **B**  
556 Shannon's Diversity values from 16S rRNA amplicon datasets for the leaf (green), root (yellow),  
557 and rhizosphere (red) sample types across all 24 genotypes used in the pilot experiment. **C**  
558 Principal coordinate analysis generated using Bray-Curtis distance for the 24 genotypes across leaf  
559 (green), root (yellow), and rhizosphere (red). **D** Mantel's R statistic plotted for each sample type  
560 indicating the degree of correlation between host genotypic distance and microbiome distance.

561



562 **Figure 2. Taxonomic classification of heritable rhizosphere microbes.** **A** The relative  
563 percentage of total OTUs belonging to each of the top 17 bacterial orders for all OTUs (left bar),  
564 non-heritable OTUs (middle bar), or heritable OTUs (right bar). Orders with significantly different  
565 numbers of OTUs in the heritable ( $H > 0.15$ ) as compared to the non-heritable fraction ( $H < 0.15$ ),  
566 as determined by Fisher's exact test ( $p < 0.05$ ), are indicated with asterisks. **B** Order-level  
567 scatterplot of the  $\log_2$  ratio between heritable and non-heritable OTU counts (x-axis) and read count  
568 abundance (y-axis). Circle sizes represent the total abundance represented by each bacterial order.  
569 Points within the dashed lines indicate merged bacterial orders that were present only in the  
570 heritable (upper right) or non-heritable (lower left) fractions.

571  
572 **Figure 3. Heritability of rhizosphere microbes across maize and sorghum.** **A** Proportional  
573 Venn diagram of bacterial orders containing heritable OTUs identified in this study (Sorghum  
574 SAP), compared with those found in a large-scale field study of maize nested association mapping  
575 (NAM) parental lines grown over two separate years, published in Walters et al., 2018<sup>5</sup>. The top  
576 100 heritable OTUs (based on  $H^2$ ) from each dataset were classified at the taxonomic rank of order  
577 to generate the Venn diagram. NAM heritable orders only present in the SAP non-heritable fraction  
578 are represented by the blue sections. Superscript letters indicate the frequency that a random  
579 subsampling of 100 sorghum OTUs (10,000 permutations) produced greater overlap with maize  
580 OTUs from either single year (a/b) or both (c). **B** Stacked barplot displaying cumulative counts (y-  
581 axis) of OTUs identified as heritable in any of the three datasets for all bacterial orders (x-axis)  
582 which have a total of at least three heritable OTUs. **C** The fraction of heritable sorghum OTUs  
583 relative to all sorghum OTUs within each order are displayed as a heatmap. Asterisks indicate  
584 orders enriched in heritable OTUs (Fisher's exact test,  $p < 0.001$ ).

585  
586 **Figure 4. A sorghum genetic locus is correlated with rhizosphere microbial abundance.** **A**  
587 Manhattan plot of PC1 community analysis GWAS. **B** Individual OTU GWAS of all OTUs with  
588 at least 5 SNPs above a threshold of  $-\log_{10}(p=10^{-2.5})$  in the 1.15 Mb window identified on the same  
589 chromosome 4 locus identified by PC1 GWAS (lower heatmap). Ratio of OTUs that associate with  
590 the sorghum major (red) or minor (blue) allele groups within this locus (upper heat map). OTUs  
591 were grouped based on the predicted presence of one or two membranes (monoderm or diderm)  
592 within each bacterial order and colored as in figure 2. **C** Tissue-specific gene expression data for  
593 sorghum genes within the chromosome 4 locus. Darker blue indicates higher expression  
594 (normalized FPKM). Asterisks indicate genes whose expression are predicted  
595 to be root-specific.

596  
597 **Figure 5. Sorghum genetic information can be used to predict rhizosphere microbiome**  
598 **composition.** **A** Canonical Analysis of Principal Coordinates of the rhizosphere microbiome for  
599 nine major allele genotypes (red) and nine minor allele genotypes (blue). **B** Ratio of indicator  
600 OTUs that associate with the sorghum major (red) or minor (blue) allele groups. OTUs were  
601 grouped based on the predicted presence of one or two membranes (monoderm or diderm), within  
602 each bacterial order, and colored as in figures 2 and 4.

## 603 604 **Acknowledgments**

605 We thank Dr. Sam Leiboff, Dr. Ling Xu, Edi Wipf, and Tuesday Simmons for their helpful  
606 discussions and critical readings of the manuscript. This research was funded by a grant from the  
607 US Department of Agriculture (2030-12210-002-00D).

608

609 **Author contributions.** S.D. conceived and designed the experiments, performed the experiments,  
610 analyzed the data, and prepared figures and/or tables; D.C. conceived and designed the  
611 experiments, analyzed the data, and prepared figures and/or tables; J.Y. conceived and designed  
612 the experiments, and analyzed the data; L.D. performed the experiments; L.W. performed the  
613 experiments and analyzed the data; D.C-D. conceived and designed the experiments, analyzed the  
614 data, and prepared figures and/or tables; All authors authored or reviewed drafts of the paper and  
615 approved the final draft.

616

## 617 **References cited**

- 618 1. Peiffer, J. A. *et al.* Diversity and heritability of the maize rhizosphere microbiome under  
619 field conditions. *Proc. Natl. Acad. Sci. U. S. A.* **110**, 6548–6553 (2013).
- 620 2. Schlaeppli, K., Dombrowski, N., Oter, R. G., Ver Loren van Themaat, E. & Schulze-Lefert,  
621 P. Quantitative divergence of the bacterial root microbiota in *Arabidopsis thaliana* relatives.  
622 *Proc. Natl. Acad. Sci. U. S. A.* **111**, 585–592 (2014).
- 623 3. Edwards, J. *et al.* Structure, variation, and assembly of the root-associated microbiomes of  
624 rice. *Proc. Natl. Acad. Sci. U. S. A.* **112**, E911–20 (2015).
- 625 4. Naylor, D., DeGraaf, S., Purdom, E. & Coleman-Derr, D. Drought and host selection  
626 influence bacterial community dynamics in the grass root microbiome. *ISME J.* **11**, 2691–  
627 2704 (2017).
- 628 5. Walters, W. A. *et al.* Large-scale replicated field study of maize rhizosphere identifies  
629 heritable microbes. *Proc. Natl. Acad. Sci. U. S. A.* **115**, 7368–7373 (2018).
- 630 6. Lebeis, S. L. *et al.* PLANT MICROBIOME. Salicylic acid modulates colonization of the  
631 root microbiome by specific bacterial taxa. *Science* **349**, 860–864 (2015).
- 632 7. Castrillo, G. *et al.* Root microbiota drive direct integration of phosphate stress and  
633 immunity. *Nature* **543**, 513–518 (2017).
- 634 8. Zhalnina, K. *et al.* Dynamic root exudate chemistry and microbial substrate preferences  
635 drive patterns in rhizosphere microbial community assembly. *Nat Microbiol* **3**, 470–480  
636 (2018).
- 637 9. Saleem, M., Law, A. D., Sahib, M. R., Pervaiz, Z. H. & Zhang, Q. Impact of root system  
638 architecture on rhizosphere and root microbiome. *Rhizosphere* **6**, 47–51 (2018).
- 639 10. Brachi, B., Morris, G. P. & Borevitz, J. O. Genome-wide association studies in plants: the  
640 missing heritability is in the field. *Genome Biol.* **12**, 232 (2011).
- 641 11. Atwell, S. *et al.* Genome-wide association study of 107 phenotypes in *Arabidopsis thaliana*  
642 inbred lines. *Nature* **465**, 627–631 (2010).
- 643 12. Wu, S. *et al.* Mapping the *Arabidopsis* Metabolic Landscape by Untargeted Metabolomics at  
644 Different Environmental Conditions. *Mol. Plant* **11**, 118–134 (2018).
- 645 13. Schaefer, R. J. *et al.* Integrating Coexpression Networks with GWAS to Prioritize Causal  
646 Genes in Maize. *Plant Cell* **30**, 2922–2942 (2018).
- 647 14. Davenport, E. R. *et al.* Genome-Wide Association Studies of the Human Gut Microbiota.  
648 *PLoS One* **10**, e0140301 (2015).

- 649 15. Wang, J. *et al.* Genome-wide association analysis identifies variation in vitamin D receptor  
650 and other host factors influencing the gut microbiota. *Nat. Genet.* **48**, 1396–1406 (2016).
- 651 16. Horton, M. W. *et al.* Genome-wide association study of *Arabidopsis thaliana* leaf microbial  
652 community. *Nat. Commun.* **5**, 5320 (2014).
- 653 17. Wallace, J. G., Kremling, K. A., Kovar, L. L. & Buckler, E. S. Quantitative Genetics of the  
654 Maize Leaf Microbiome. *Phytobiomes Journal* **2**, 208–224 (2018).
- 655 18. Roman-Reyna, V. *et al.* The rice leaf microbiome has a conserved community structure  
656 controlled by complex host-microbe interactions. *bioRxiv* 615278 (2019)  
657 doi:10.1101/615278.
- 658 19. Baerson, S. R. *et al.* A functional genomics investigation of allelochemical biosynthesis in  
659 *Sorghum bicolor* root hairs. *J. Biol. Chem.* **283**, 3231–3247 (2008).
- 660 20. Casa, A. M. *et al.* Community Resources and Strategies for Association Mapping in  
661 *Sorghum*. *Crop Sci.* **48**, 30–40 (2008).
- 662 21. Morris, G. P. *et al.* Population genomic and genome-wide association studies of  
663 agroclimatic traits in sorghum. *Proc. Natl. Acad. Sci. U. S. A.* **110**, 453–458 (2013).
- 664 22. Xu, L. *et al.* Drought delays development of the sorghum root microbiome and enriches for  
665 monoderm bacteria. *Proc. Natl. Acad. Sci. U. S. A.* **115**, E4284–E4293 (2018).
- 666 23. Oberholster, T., Vikram, S., Cowan, D. & Valverde, A. Key microbial taxa in the  
667 rhizosphere of sorghum and sunflower grown in crop rotation. *Sci. Total Environ.* **624**, 530–  
668 539 (2018).
- 669 24. Swigonova, Z. *et al.* On the tetraploid origin of the maize genome. *Comp. Funct. Genomics*  
670 **5**, 281–284 (2004).
- 671 25. Goodstein, D. M. *et al.* Phytozome: a comparative platform for green plant genomics.  
672 *Nucleic Acids Res.* **40**, D1178–86 (2012).
- 673 26. Bergelson, J., Mittelstrass, J. & Horton, M. W. Characterizing both bacteria and fungi  
674 improves understanding of the *Arabidopsis* root microbiome. *Sci. Rep.* **9**, 24 (2019).
- 675 27. Bodenhausen, N., Horton, M. W. & Bergelson, J. Bacterial communities associated with the  
676 leaves and the roots of *Arabidopsis thaliana*. *PLoS One* **8**, e56329 (2013).
- 677 28. Copeland, J. K., Yuan, L., Layeghifard, M., Wang, P. W. & Guttman, D. S. Seasonal  
678 community succession of the phyllosphere microbiome. *Mol. Plant. Microbe. Interact.* **28**,  
679 274–285 (2015).
- 680 29. Badri, D. V., Chaparro, J. M., Zhang, R., Shen, Q. & Vivanco, J. M. Application of natural  
681 blends of phytochemicals derived from the root exudates of *Arabidopsis* to the soil reveal  
682 that phenolic-related compounds predominantly modulate the soil microbiome. *J. Biol.*  
683 *Chem.* **288**, 4502–4512 (2013).
- 684 30. Zhang, N. *et al.* Effects of different plant root exudates and their organic acid components  
685 on chemotaxis, biofilm formation and colonization by beneficial rhizosphere-associated  
686 bacterial strains. *Plant Soil* **374**, 689–700 (2014).
- 687 31. Angus, A. A. *et al.* Plant-associated symbiotic Burkholderia species lack hallmark strategies  
688 required in mammalian pathogenesis. *PLoS One* **9**, e83779 (2014).

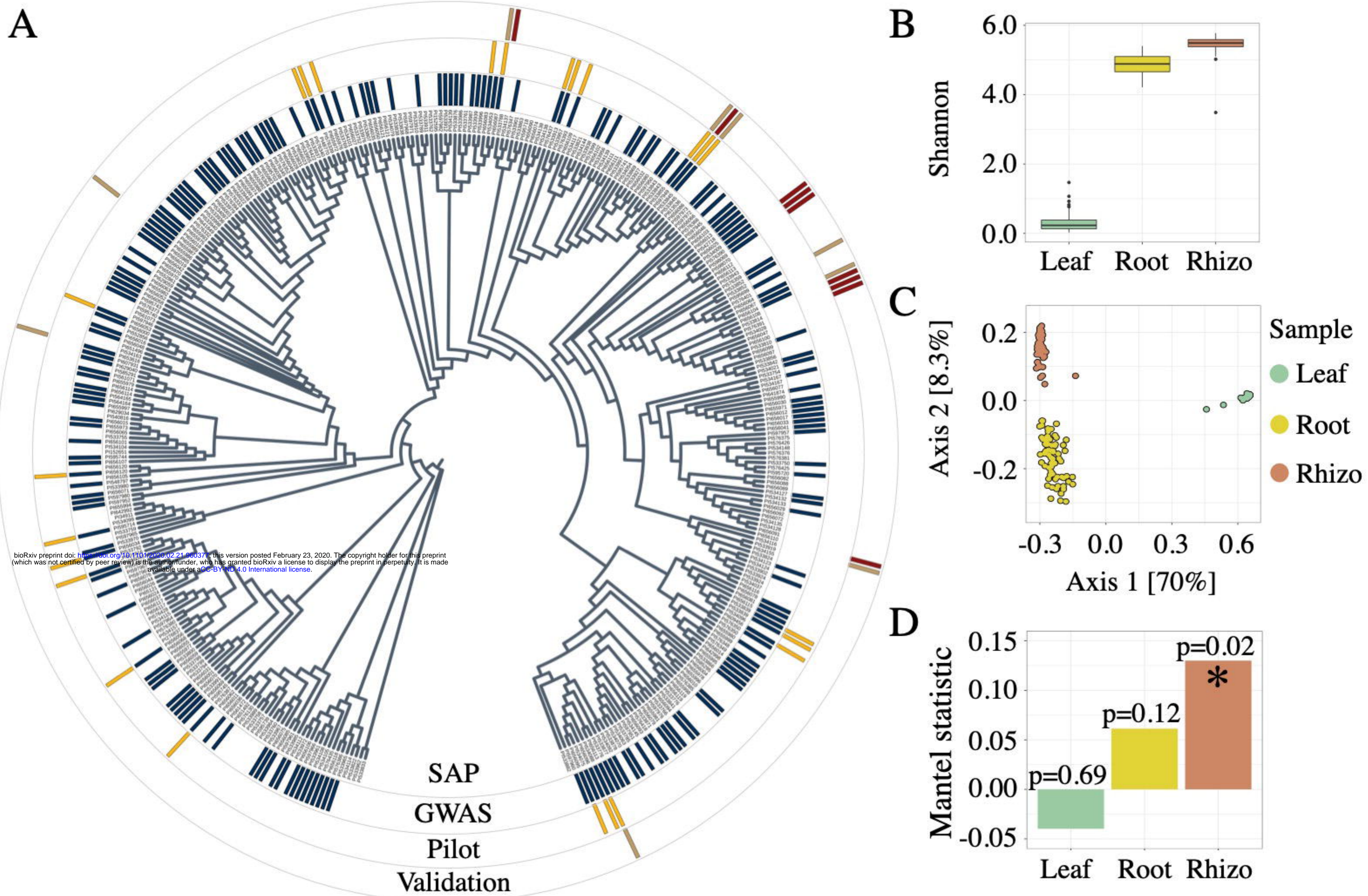
- 689 32. Kim, J. K. & Lee, B. L. Symbiotic factors in Burkholderia essential for establishing an  
690 association with the bean bug, Riptortus pedestris. *Arch. Insect Biochem. Physiol.* **88**, 4–17  
691 (2015).
- 692 33. Shu, L. *et al.* Symbiont location, host fitness, and possible coadaptation in a symbiosis  
693 between social amoebae and bacteria. *Elife* **7**, (2018).
- 694 34. Mannaa, M., Park, I. & Seo, Y.-S. Genomic Features and Insights into the Taxonomy,  
695 Virulence, and Benevolence of Plant-Associated Burkholderia Species. *Int. J. Mol. Sci.* **20**,  
696 (2018).
- 697 35. Erbilgin, O., McDonald, K. L. & Kerfeld, C. A. Characterization of a planctomycetal  
698 organelle: a novel bacterial microcompartment for the aerobic degradation of plant  
699 saccharides. *Appl. Environ. Microbiol.* **80**, 2193–2205 (2014).
- 700 36. Bulgarelli, D. *et al.* Revealing structure and assembly cues for Arabidopsis root-inhabiting  
701 bacterial microbiota. *Nature* **488**, 91–95 (2012).
- 702 37. Meents, M. J., Watanabe, Y. & Samuels, A. L. The cell biology of secondary cell wall  
703 biosynthesis. *Ann. Bot.* **121**, 1107–1125 (2018).
- 704 38. Sasse, J., Martinoia, E. & Northen, T. Feed Your Friends: Do Plant Exudates Shape the Root  
705 Microbiome? *Trends Plant Sci.* **23**, 25–41 (2018).
- 706 39. Claverie, J. *et al.* The Cell Wall-Derived Xyloglucan Is a New DAMP Triggering Plant  
707 Immunity in *Vitis vinifera* and *Arabidopsis thaliana*. *Front. Plant Sci.* **9**, 1725 (2018).
- 708 40. Hou, S., Liu, Z., Shen, H. & Wu, D. Damage-Associated Molecular Pattern-Triggered  
709 Immunity in Plants. *Front. Plant Sci.* **10**, 646 (2019).
- 710 41. Parisi, G. *et al.* Gamma carbonic anhydrases in plant mitochondria. *Plant Mol. Biol.* **55**,  
711 193–207 (2004).
- 712 42. DiMario, R. J., Clayton, H., Mukherjee, A., Ludwig, M. & Moroney, J. V. Plant Carbonic  
713 Anhydrases: Structures, Locations, Evolution, and Physiological Roles. *Mol. Plant* **10**, 30–  
714 46 (2017).
- 715 43. Floryszak-Wieczorek, J. & Arasimowicz-Jelonek, M. The multifunctional face of plant  
716 carbonic anhydrase. *Plant Physiol. Biochem.* **112**, 362–368 (2017).
- 717 44. Atkins, C. A. Occurrence and some properties of carbonic anhydrases from legume root  
718 nodules. *Phytochemistry* **13**, 93–98 (1974).
- 719 45. De La Peña, T. C., Frugier, F. & McKhann, H. I. A carbonic anhydrase gene is induced in  
720 the nodule primordium and its cell-specific expression is controlled by the presence of  
721 *Rhizobium* during development. *The Plant* (1997).
- 722 46. Slaymaker, D. H. *et al.* The tobacco salicylic acid-binding protein 3 (SABP3) is the  
723 chloroplast carbonic anhydrase, which exhibits antioxidant activity and plays a role in the  
724 hypersensitive defense response. *Proc. Natl. Acad. Sci. U. S. A.* **99**, 11640–11645 (2002).
- 725 47. Restrepo, S. *et al.* Gene profiling of a compatible interaction between *Phytophthora*  
726 *infestans* and *Solanum tuberosum* suggests a role for carbonic anhydrase. *Mol. Plant.*  
727 *Microbe. Interact.* **18**, 913–922 (2005).
- 728 48. Wang, Y.-Q. *et al.* S-nitrosylation of AtSABP3 antagonizes the expression of plant

- 729 immunity. *J. Biol. Chem.* **284**, 2131–2137 (2009).
- 730 49. Khan, G. A., Vogiatzaki, E., Glauser, G. & Poirier, Y. Phosphate Deficiency Induces the  
731 Jasmonate Pathway and Enhances Resistance to Insect Herbivory. *Plant Physiol.* **171**, 632–  
732 644 (2016).
- 733 50. Hiruma, K. *et al.* Root Endophyte *Colletotrichum tofieldiae* Confers Plant Fitness Benefits  
734 that Are Phosphate Status Dependent. *Cell* **165**, 464–474 (2016).
- 735 51. Yamada, K., Saijo, Y., Nakagami, H. & Takano, Y. Regulation of sugar transporter activity  
736 for antibacterial defense in Arabidopsis. *Science* **354**, 1427–1430 (2016).
- 737 52. Wagner, M. R. *et al.* Host genotype and age shape the leaf and root microbiomes of a wild  
738 perennial plant. *Nat. Commun.* **7**, 12151 (2016).
- 739 53. Edwards, J. A. *et al.* Compositional shifts in root-associated bacterial and archaeal  
740 microbiota track the plant life cycle in field-grown rice. *PLoS Biol.* **16**, e2003862 (2018).
- 741 54. Lundberg, D. S. *et al.* Defining the core Arabidopsis thaliana root microbiome. *Nature* **488**,  
742 86–90 (2012).
- 743 55. Haney, C. H., Samuel, B. S., Bush, J. & Ausubel, F. M. Associations with rhizosphere  
744 bacteria can confer an adaptive advantage to plants. *Nat Plants* **1**, (2015).
- 745 56. Fitzpatrick, C. R. *et al.* Assembly and ecological function of the root microbiome across  
746 angiosperm plant species. *Proc. Natl. Acad. Sci. U. S. A.* **115**, E1157–E1165 (2018).
- 747 57. Pfeifer, B., Wittelsbürger, U., Ramos-Onsins, S. E. & Lercher, M. J. PopGenome: an  
748 efficient Swiss army knife for population genomic analyses in R. *Mol. Biol. Evol.* **31**, 1929–  
749 1936 (2014).
- 750 58. Letunic, I. & Bork, P. Interactive Tree Of Life (iTOL) v4: recent updates and new  
751 developments. *Nucleic Acids Res.* **47**, W256–W259 (2019).
- 752 59. Simmons, T., Caddell, D. F., Deng, S. & Coleman-Derr, D. Exploring the Root Microbiome:  
753 Extracting Bacterial Community Data from the Soil, Rhizosphere, and Root Endosphere. *J.*  
754 *Vis. Exp.* (2018) doi:10.3791/57561.
- 755 60. Bolyen, E. *et al.* Reproducible, interactive, scalable and extensible microbiome data science  
756 using QIIME 2. *Nat. Biotechnol.* **37**, 852–857 (2019).
- 757 61. Magoč, T. & Salzberg, S. L. FLASH: fast length adjustment of short reads to improve  
758 genome assemblies. *Bioinformatics* **27**, 2957–2963 (2011).
- 759 62. Edgar, R. C. UPARSE: highly accurate OTU sequences from microbial amplicon reads. *Nat.*  
760 *Methods* **10**, 996–998 (2013).
- 761 63. Wang, Q., Garrity, G. M., Tiedje, J. M. & Cole, J. R. Naive Bayesian classifier for rapid  
762 assignment of rRNA sequences into the new bacterial taxonomy. *Appl. Environ. Microbiol.*  
763 **73**, 5261–5267 (2007).
- 764 64. Paulson, J. N., Colin Stine, O., Bravo, H. C. & Pop, M. Differential abundance analysis for  
765 microbial marker-gene surveys. *Nature Methods* vol. 10 1200–1202 (2013).
- 766 65. Covarrubias-Pazarán, G. Genome-Assisted Prediction of Quantitative Traits Using the R  
767 Package sommer. *PLoS One* **11**, e0156744 (2016).
- 768 66. Yu, J., Holland, J. B., McMullen, M. D. & Buckler, E. S. Genetic design and statistical



- 769 power of nested association mapping in maize. *Genetics* **178**, 539–551 (2008).
- 770 67. Zhou, X. & Stephens, M. Genome-wide efficient mixed-model analysis for association  
771 studies. *Nat. Genet.* **44**, 821–824 (2012).
- 772 68. Oksanen, J. *et al.* Vegan: community ecology package. software. (2016).
- 773 69. Paradis, E., Claude, J. & Strimmer, K. APE: Analyses of Phylogenetics and Evolution in R  
774 language. *Bioinformatics* **20**, 289–290 (2004).
- 775 70. Roberts, D. W. & Roberts, M. D. W. Package ‘labdsv’. *Ordination and Multivariate* (2016).

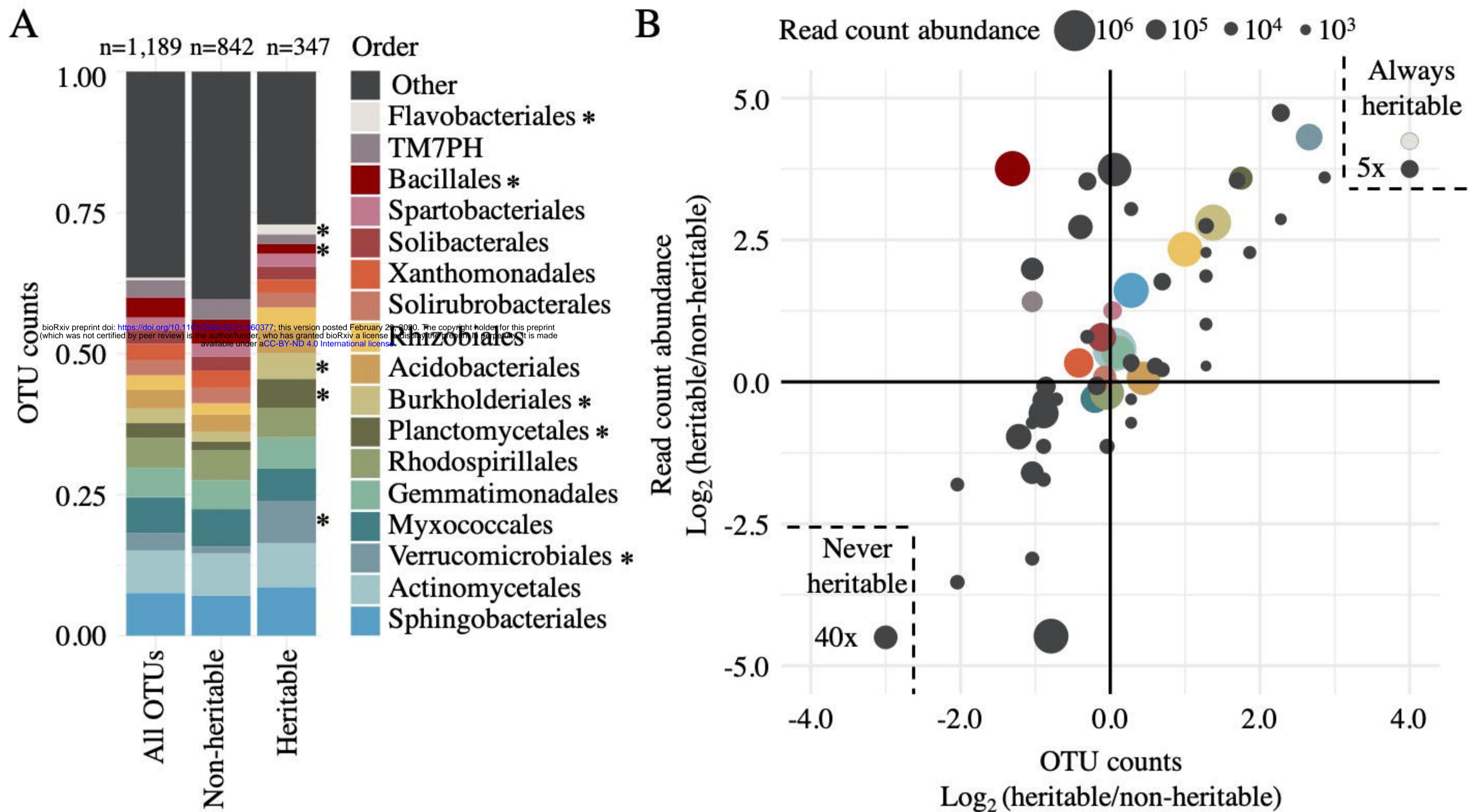




bioRxiv preprint doi: <https://doi.org/10.1101/2020.02.21.400371>; this version posted February 23, 2020. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted bioRxiv a license to display the preprint in perpetuity. It is made available under aCC-BY-ND 4.0 International license.

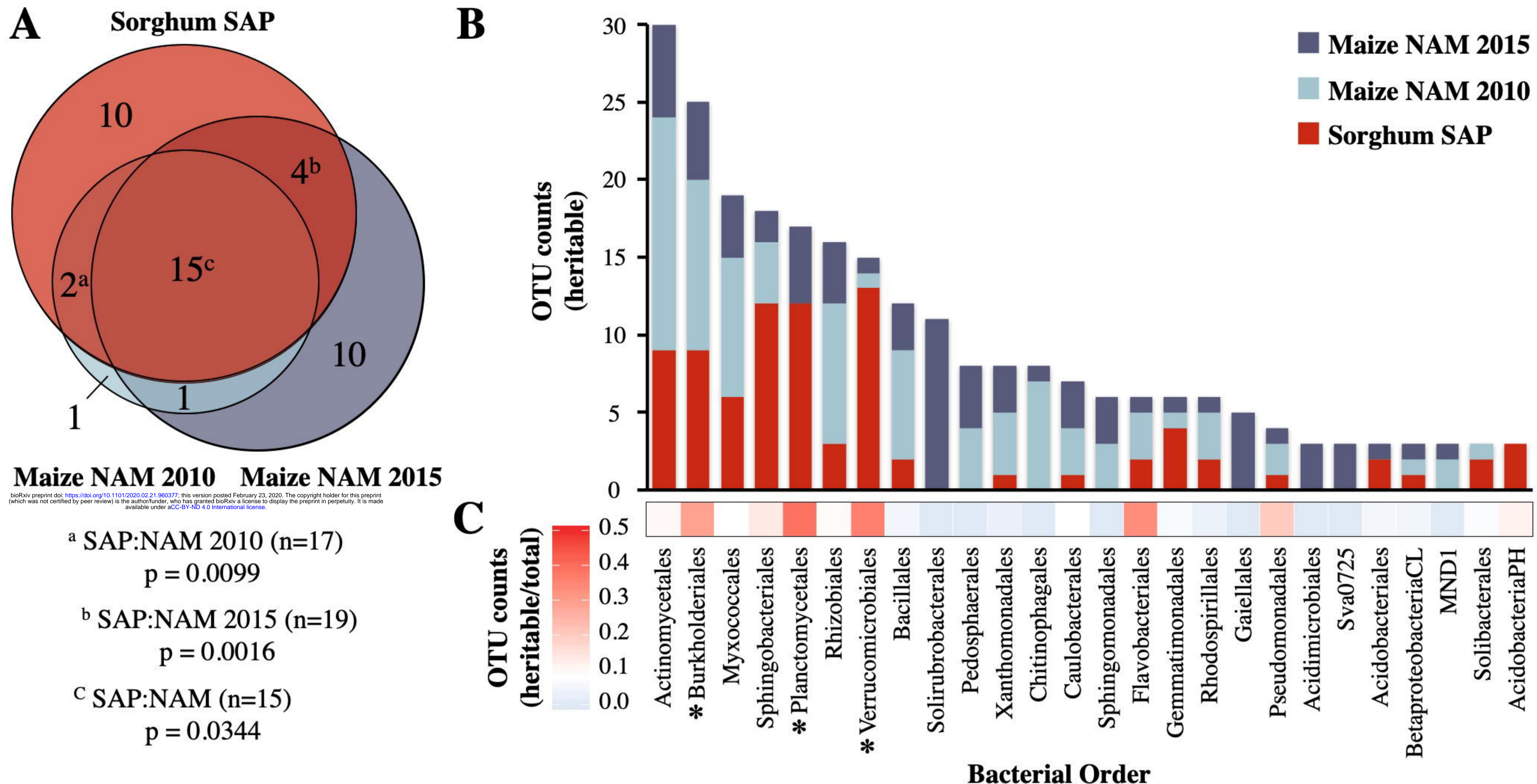
**Figure 1. Sample type and population selection.** **A** Phylogenetic tree representing the 378 member sorghum association panel (SAP, inner ring), the subset of 200 lines selected for GWAS (2nd ring from the center, in blue), the 24 lines used for sample type selection (Pilot, 3rd ring from the center, in yellow), and the 18 genotypes used for GWAS validation containing either the Chromosome 4 minor allele (red) or major allele (brown) identified by GWAS (outer ring). **B** Shannon's Diversity values from 16S rRNA amplicon datasets for the leaf (green), root (yellow), and rhizosphere (red) sample types across all 24 genotypes used in the pilot experiment. **C** Principal coordinate analysis generated using Bray-Curtis distance for the 24 genotypes across leaf (green), root (yellow), and rhizosphere (red). **D** Mantel's R statistic plotted for each sample type indicating the degree of correlation between host genotypic distance and microbiome distance.





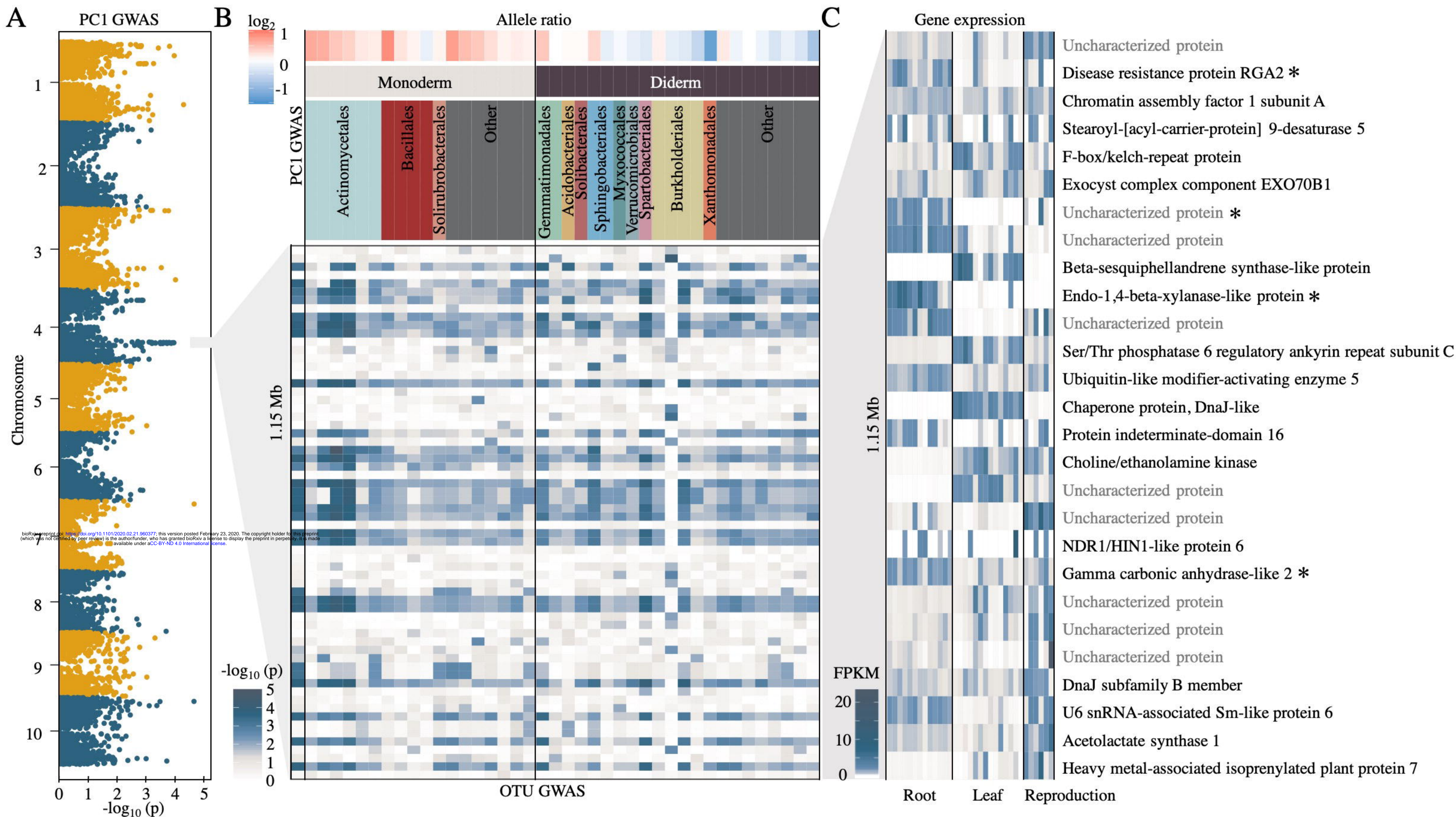
**Figure 2. Taxonomic classification of heritable rhizosphere microbes.** **A** The relative percentage of total OTUs belonging to each of the top 17 bacterial orders for all OTUs (left bar), non-heritable OTUs (middle bar), or heritable OTUs (right bar). Orders with significantly different numbers of OTUs in the heritable ( $H^2 > 0.15$ ) as compared to the non-heritable fraction ( $H^2 < 0.15$ ), as determined by Fisher's exact test ( $p < 0.05$ ), are indicated with asterisks. **B** Order-level scatterplot of the  $\log_2$  ratio between heritable and non-heritable OTU counts (x-axis) and read count abundance (y-axis). Circle sizes represent the total abundance represented by each bacterial order. Points within the dashed lines indicate merged bacterial orders that were present only in the heritable (upper right) or non-heritable (lower left) fractions.





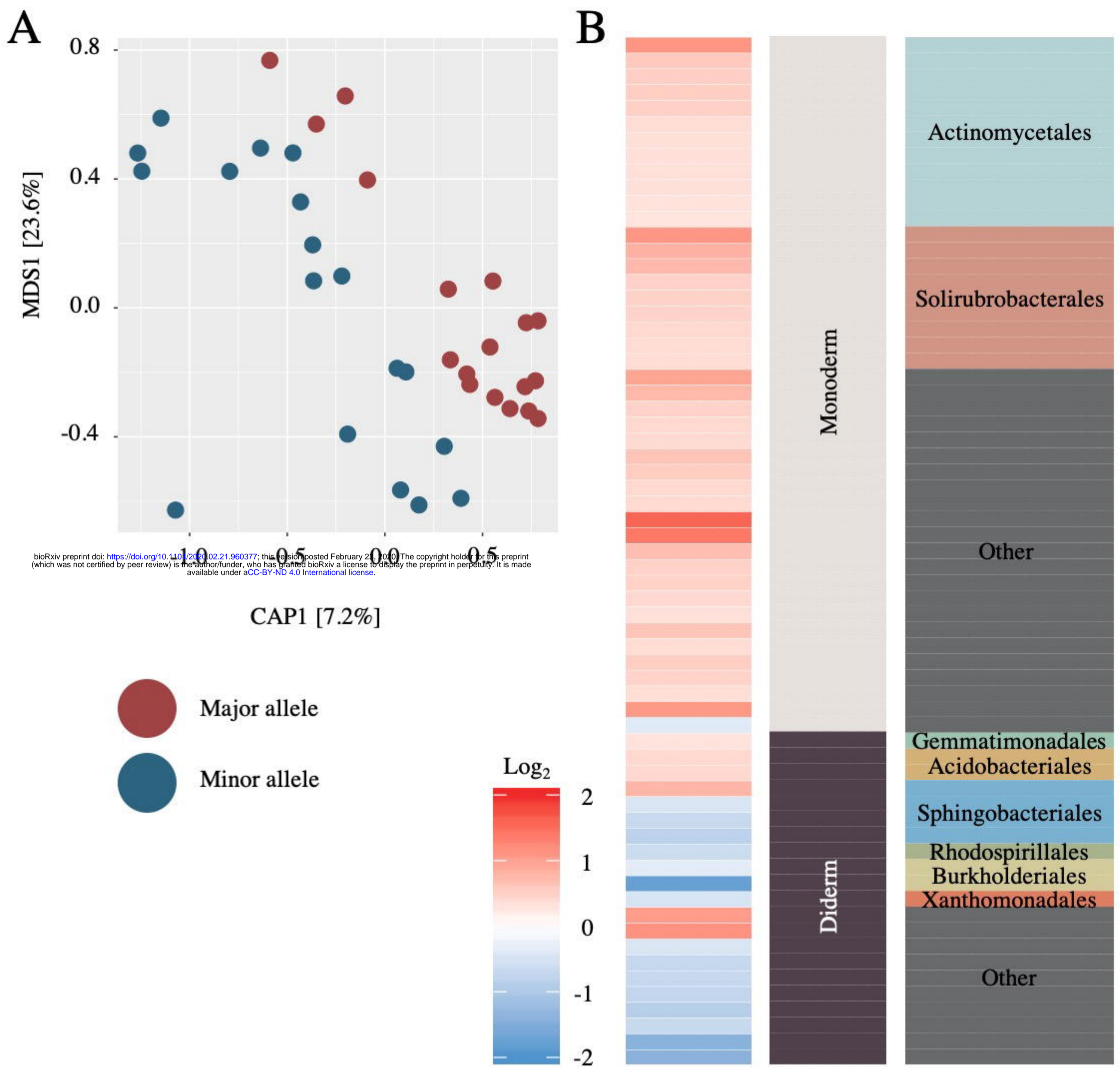
**Figure 3. Heritability of rhizosphere microbes across maize and sorghum.** **A** Proportional Venn diagram of bacterial orders containing heritable OTUs identified in this study (Sorghum SAP), compared with those found in a large-scale field study of maize nested association mapping (NAM) parental lines grown over two separate years, published in Walters et al., 2018<sup>5</sup>. The top 100 heritable OTUs (based on  $H^2$ ) from each dataset were classified at the taxonomic rank of order to generate the Venn diagram. NAM heritable orders only present in the SAP non-heritable fraction are represented by the blue sections. Superscript letters indicate the frequency that a random subsampling of 100 sorghum OTUs (10,000 permutations) produced greater overlap with maize OTUs from either single year (a/b) or both (c). **B** Stacked barplot displaying cumulative counts (y-axis) of OTUs identified as heritable in any of the three datasets for all bacterial orders (x-axis) which have a total of at least three heritable OTUs. **C** The fraction of heritable sorghum OTUs relative to all sorghum OTUs within each order are displayed as a heatmap. Asterisks indicate orders enriched in heritable OTUs (Fisher's exact test,  $p < 0.001$ ).





**Figure 4. A sorghum genetic locus is correlated with rhizosphere microbial abundance.** **A** Manhattan plot of PC1 community analysis GWAS. **B** Individual OTU GWAS of all OTUs with at least 5 SNPs above a threshold of  $-\log_{10}(p)=10^{-2.5}$  in the 1.15 Mb window identified on the same chromosome 4 locus identified by PC1 GWAS (lower heatmap). Ratio of OTUs that associate with the sorghum major (red) or minor (blue) allele groups within this locus (upper heatmap). OTUs were grouped based on the predicted presence of one or two membranes (monoderm or diderm) within each bacterial order and colored as in figure 2. **C** Tissue-specific gene expression data for sorghum genes within the chromosome 4 locus. Darker blue indicates higher expression (normalized FPKM). Asterisks indicate genes whose expression are predicted to be root-specific.





**Figure 5. Sorghum genetic information can be used to predict rhizosphere microbiome composition.** **A** Canonical Analysis of Principal Coordinates of the rhizosphere microbiome for nine major allele genotypes (red) and nine minor allele genotypes (blue). **B** Ratio of indicator OTUs that associate with the sorghum major (red) or minor (blue) allele groups. OTUs were grouped based on the predicted presence of one or two membranes (monoderm or diderm), within each bacterial order, and colored as in figures 2 and 4.