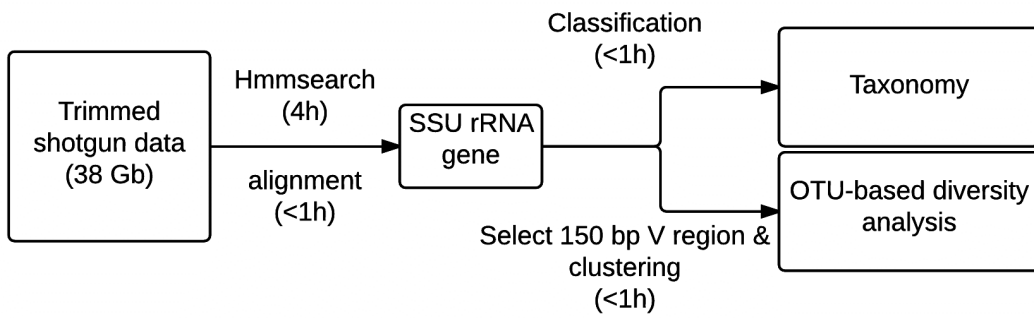1    **Supplemental Materials**

2    **TABLE S1** Higher fungi/bacteria ratio and percent of AMF fungi are in rhizosphere sample (M1)

3    than bulk sample (SB1).

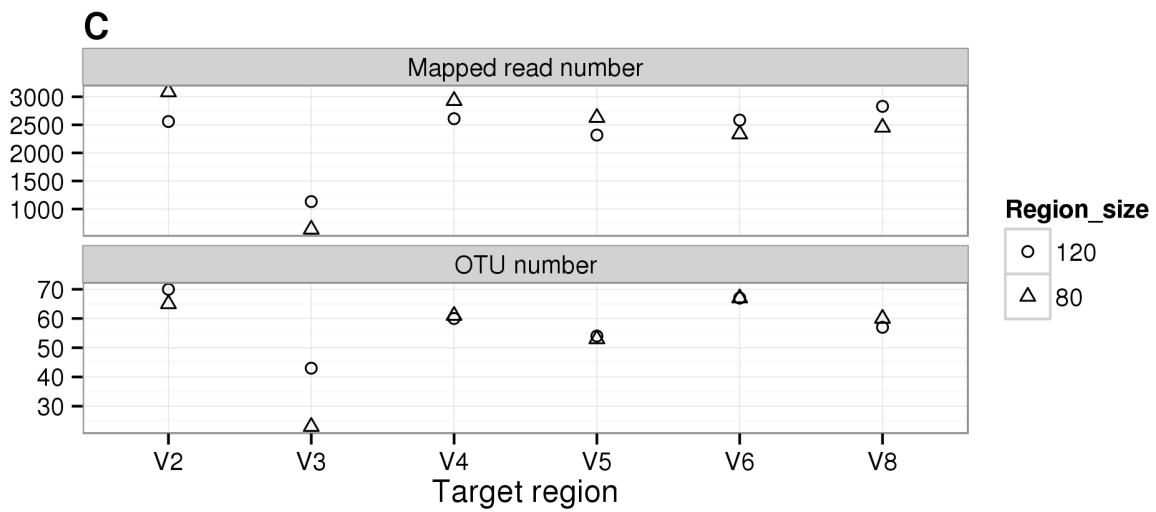|          | % Bacteria | % Fungi | % AMF in Fungi | Fungi/Bacteria |
|----------|------------|---------|----------------|----------------|
| SB1-SSU  | 97.00%     | 0.36%   | 0.00%          | 0.0037         |
| SB1-LSU  | 96.75%     | 0.42%   | 0.00%          | 0.0044         |
| M1-SSU   | 92.38%     | 2.60%   | 0.18%          | 0.0281         |
| M1-LSU   | 92.94%     | 2.48%   | 0.18%          | 0.0267         |

4

5

6

7

8

9

10

11

12

13

```
                                      Classification
                                         (<1h)
┌──────────────┐                                    ┌──────────────────┐
│              │  Hmmsearch     ┌──────────┐   ┌────▶│    Taxonomy      │
│   Trimmed    │    (4h)        │ SSU rRNA │   │     └──────────────────┘
│ shotgun data │───────────────▶│   gene   │───┤     ┌──────────────────┐
│   (38 Gb)    │  alignment     └──────────┘   └────▶│OTU-based diversity│
│              │    (<1h)                            │     analysis     │
└──────────────┘                  Select 150 bp V region &  └──────────────────┘
                                        clustering
                                         (<1h)
```
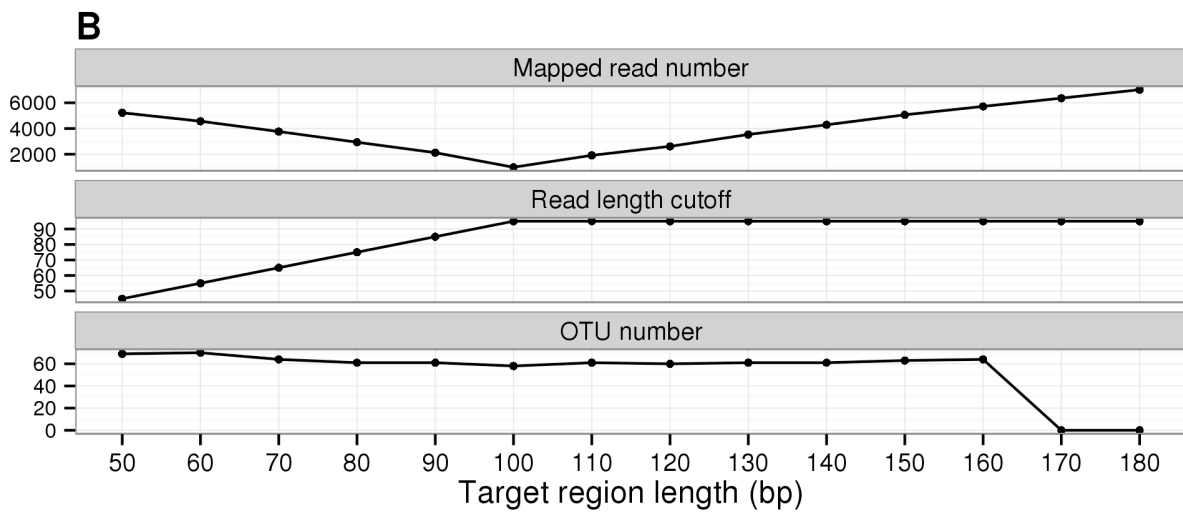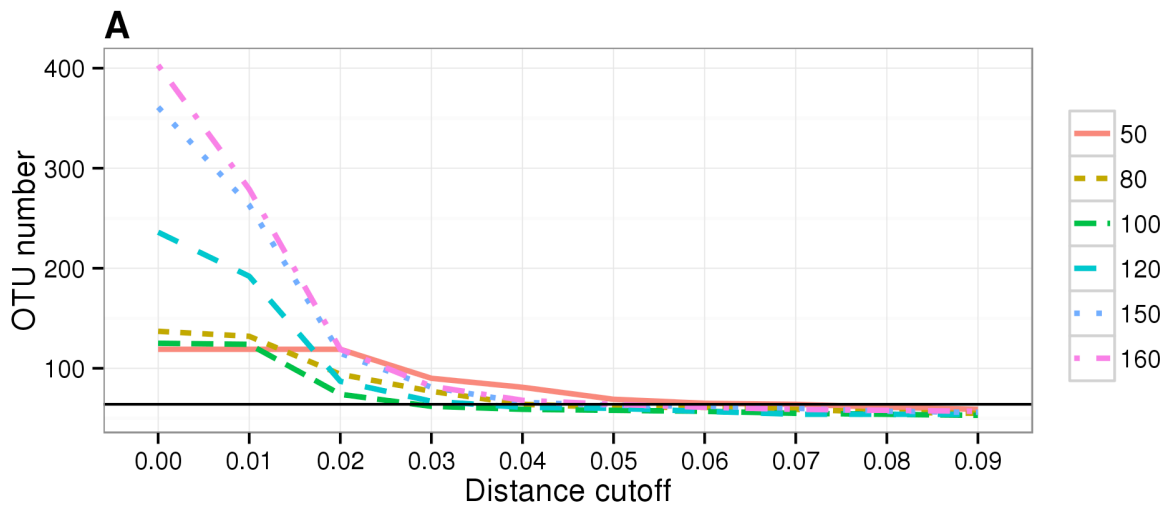
14

15   **FIG S1** Flowchart of SSUsearch pipeline. SSU rRNA gene fragments were retrieved by an

16   hmmsearch and alignment step, which could be further used for reference-based (supervised)

17   diversity analysis (taxonomy). Those fragments aligned to 150 bp of a variable region could be

18   used for OTU-based (unsupervised) diversity analysis. SSU rRNA gene identification

19   hmmsearch is the most time consuming steps. For 1 lane of trimmed HiSeq data (38 Gb) from

20   Miscanthus rhizosphere sample (M1), SSU rRNA gene identification took about 4 hours with

21   peak memory usage of about 4.5 Gb. In the analysis pipeline, where there was not a clear

22   performance difference between tools, we mostly used Mothur including databases. For (de

23   novo) OTU based analysis, SSU rRNA gene reference sequences with taxonomy information are

24   required for classification. Another smaller set of aligned references is required to align the gene

25   fragments from shotgun data. The SILVA database (the official one, not the one included with

26   QIIME or Mothur) was used to build the HMM since the reference set from SILVA was more up

27   to date. Two scripts in QIIME were used to cluster (UCLUST, the default) and pick

28   representative sequences. Building the HMM is not part of the pipeline but using the built HMM

29   is. We found that complete-linkage clustering is faster and requires less memory with McClust

30   than with Mothur (dist.seqs and cluster). Additionally, we use two scripts in QIIME

31   (pick_otus.py and pick_rep_set.py) to select representative sequences for building the HMM due

32   to ease of use and the GreenGenes database is included for use with the Copyrighter copy number
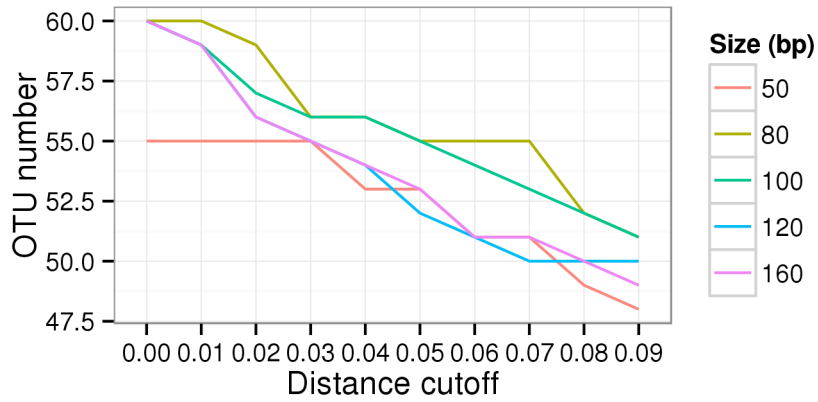
33   correction tool.

34
35

36

37

38

39

4

40    **FIG S2** Testing the effect of target region size and variable region on clustering on a synthetic

41    community with 64 species with read length at about 100 bp. Subfigure A shows a distance

42    cutoff of 4% or 5% is proper for all regions sizes from 50 bp to 160 bp in V4 (OTU number

43    approached the species number 64 as indicated by the black line).  Subfigure B shows more

44    details in the method used for subfigure A. Panel "Read Length cutoff" in B shows minimum

45    read length was set to the target region size minus 5 bp if the region size was less than 100 bp,

46    and 95 bp when the region size was longer than 100 bp. As a result, the number of reads aligned

47    decreased as the target region size increased until 100 bp, and then the number of reads aligned

48    increased with target region as shown in Panel "Mapped read number" in B. Panel "OTU

49    number" in C shows OTU number at distance cutoff of 0.05 and our method works well from a

50    50 bp region to 160 bp region in V4. Subfigure C tests our unsupervised method on multiple

51    hyper-variable regions (V2, V3, V4, V5, V6, V8) with region size of 120 bp (circle) and 80 bp

52    (triangle). Panel "Mapped read number" in C shows the number of reads mapped to each chosen

53    region. Panel "OTU number" in C shows the number of OTUs in each region at distance cutoff of

54    0.05. All regions have consistent mapped read number and OTUs except V3.

55

56

57

58

59

60

61

62

63



64

**FIG S3** Testing the effect of target region size on V4 of full-length SSU rRNA genes from a

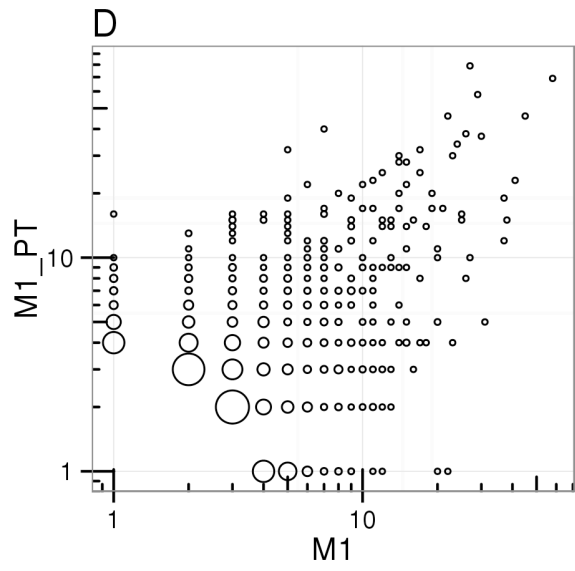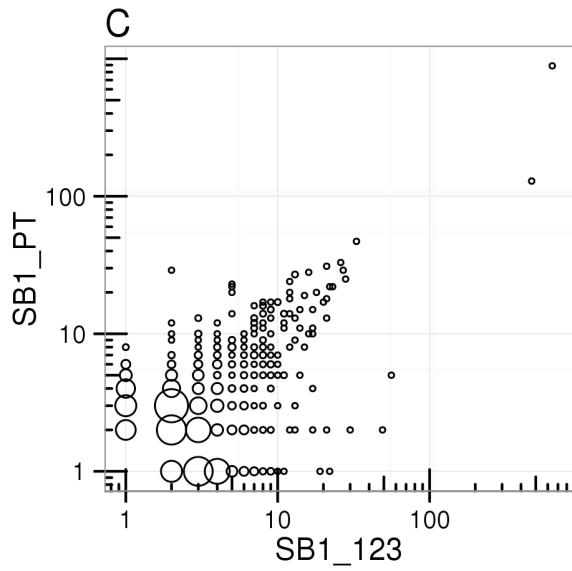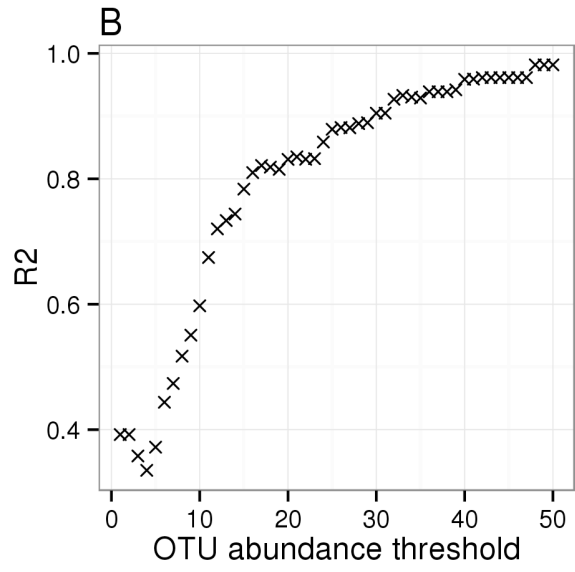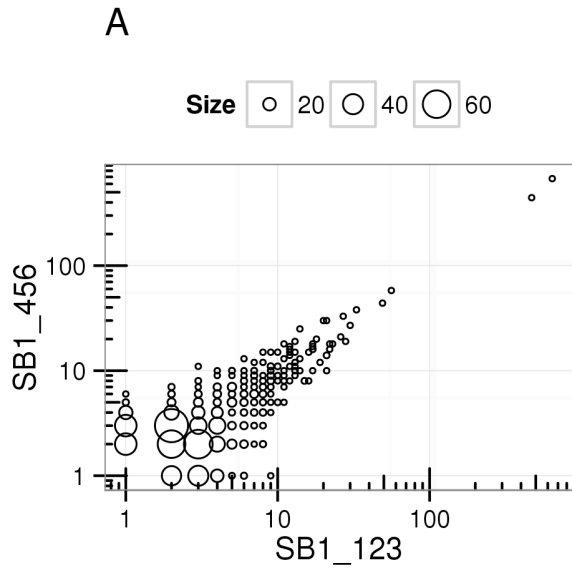synthetic community with 64 species.

67

68

69

70

71

72

73
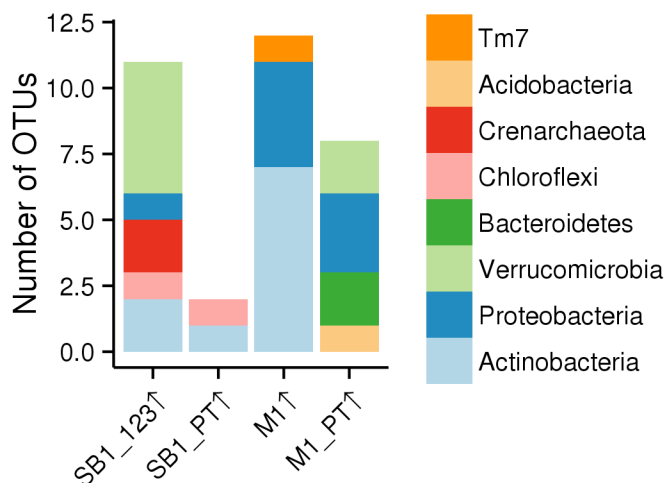
74

75

76

77

78

79

80

81



82

7

83    **FIG S4** Technical reproducibility test of our unsupervised clustering and comparison of OTU

84    abundances between paired shotgun and amplicon data. Subfigure A shows consistent OTU

85    abundance profiles in two technical replicates (Pearson's correlation coefficient is 0.997). X axis

86    shows number of reads in each OTU in replicate SB1_123, and y axis shows number of reads in

87    each OTU in replicate SB1_456. The size of circle is proportional to number of OTUs at the

88    same location in the plot (with the same counts in SB1_123 and also in SB1_456).  The

89    consistency of counts of each OTU in two replicates becomes better when the abundance of

90    OTUs are higher. Subfigure B shows progressive dropout analysis of two technical replicates of

91    shotgun data. There is significant correlation of counts of each OTU between technical replicates;

92    X axis is the threshold of OTU abundance and y axis is the $R^2$ of linear regression of log

93    transformed OTU abundances in two replicates. OTUs with lower abundance than the thresholds

94    (x axis) were discarded before regression analysis. Subfigure C and D shows comparison of OTU

95    abundance profile between paired shotgun and amplicon data in bulk soil sample (SB1) and

96    rhizosphere sample (M1), respectively. There is inconsistency between shotgun data and

97    amplicon data in both samples. X axis shows number of SSU rRNA gene fragments in shotgun

98    data per OTU in log scale, and y axis shows number of amplicon sequences in each OTU in log

99    scale. The OTU abundance in both amplicon and shotgun data were increased by 1 to avoid 0

100   counts that can be displaced in log scale. The size of circle is proportional to number of OTUs

101   with the same abundance in both types of data. There are OTUs with significantly different

102   abundances in the two types of data (circles deviate from diagonal line). Pearson's correlation

103   between two types of data is 0.873 in SB1 and 0.581 in M1.

104

105

106



107

**FIG S5** Phyla of OTUs significantly different between shotgun data and amplicon data. SB1_123

are shotgun data and SB1_PT are amplicon data both from the same DNA from bulk soil sample.

M1 are shotgun data and M1_PT are amplicon data both from the same DNA from Miscanthus

rhizosphere sample. OTUs significantly different were defined as those with total abundance > 10

and fold change between two types of data > 5 or < 0.2. Verrucomicrobia was biased against in

bulk soil sample amplified by V6-V8 primer (SB1_PT) but biased for in rhizosphere sample

amplified with V4 primer (M1_PT). Actinobacteria was biased against in rhizosphere sample

(M1).

108
109
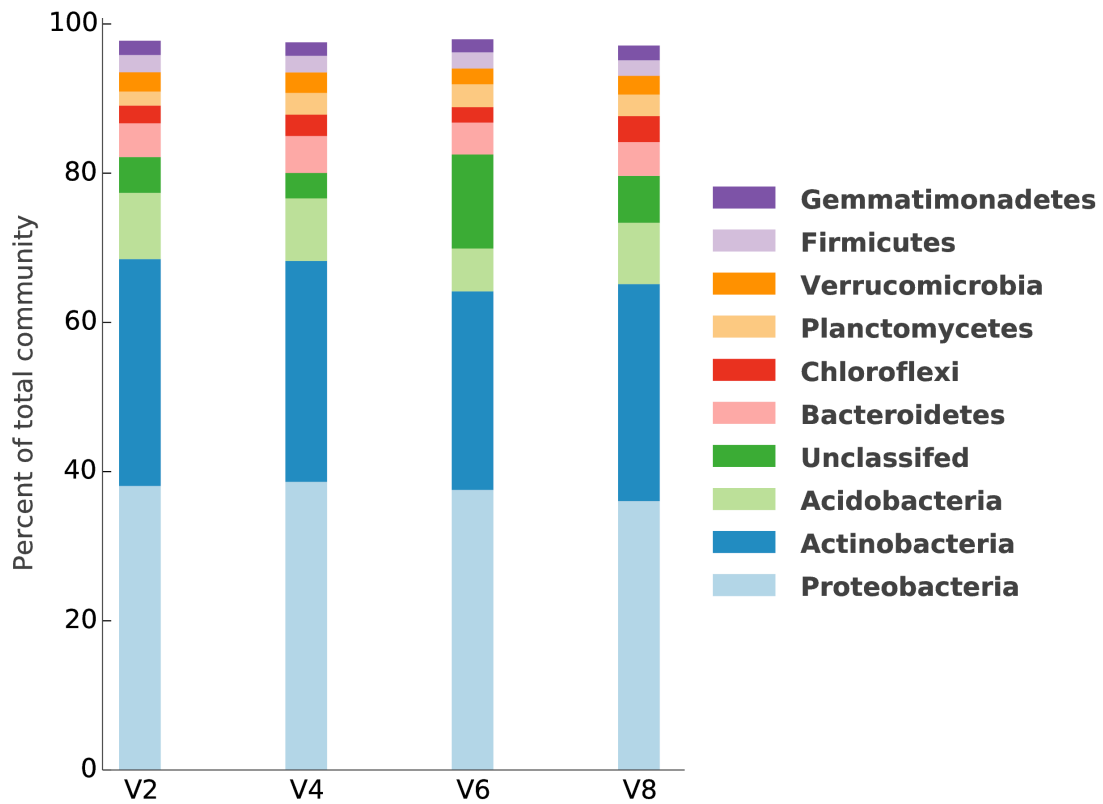110
111
112
113
114
115
116
117
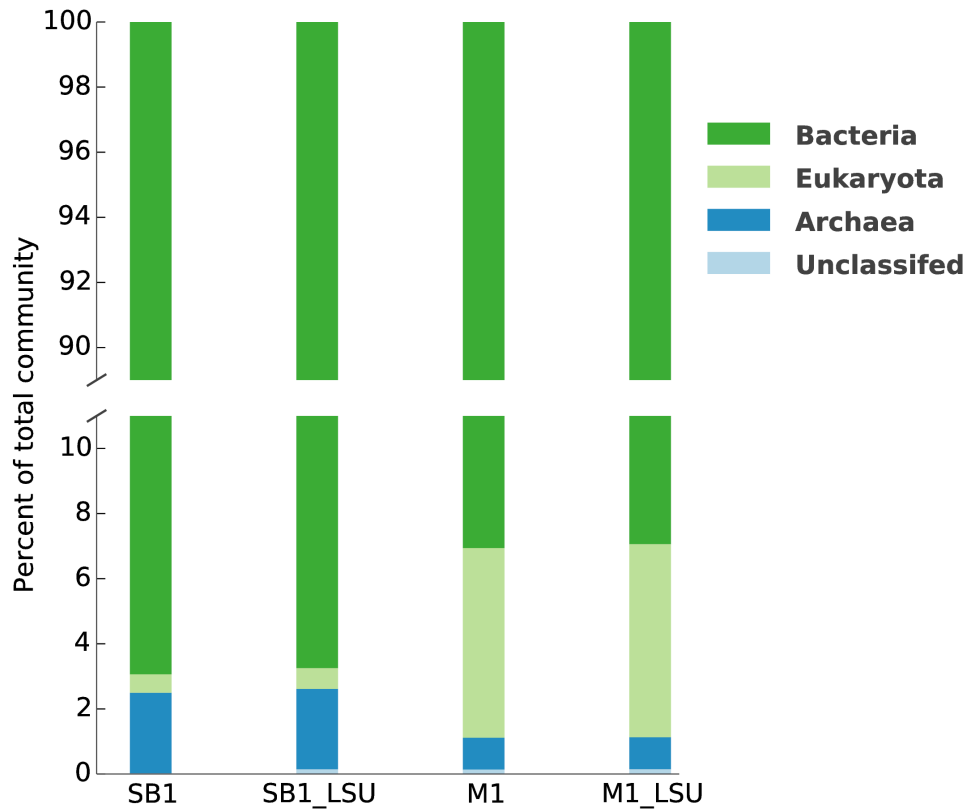118
119
120
121
122

123

124



125

126 **FIG S6** Bacterial phylum profile comparison using different variable regions. Different variable

127 regions have similar taxonomy profiles, except that V6 has more unclassified sequences. The

128 minimum Pearson's correlation between the regions is 0.96. Classifications were done using SSU

129 rRNA gene fragments from Miscanthus rhizosphere soil sample (M1) and SILVA database as

130 reference.

131

132

**FIG S7** Taxonomy profile comparison at domain level using SSU and LSU rRNA genes. For

both the bulk soil sample (SB1) and rhizosphere sample (M1), SSU and LSU show consistent

domain level taxonomy distribution (Pearson's correlation coefficient = 1). "_LSU" indicates

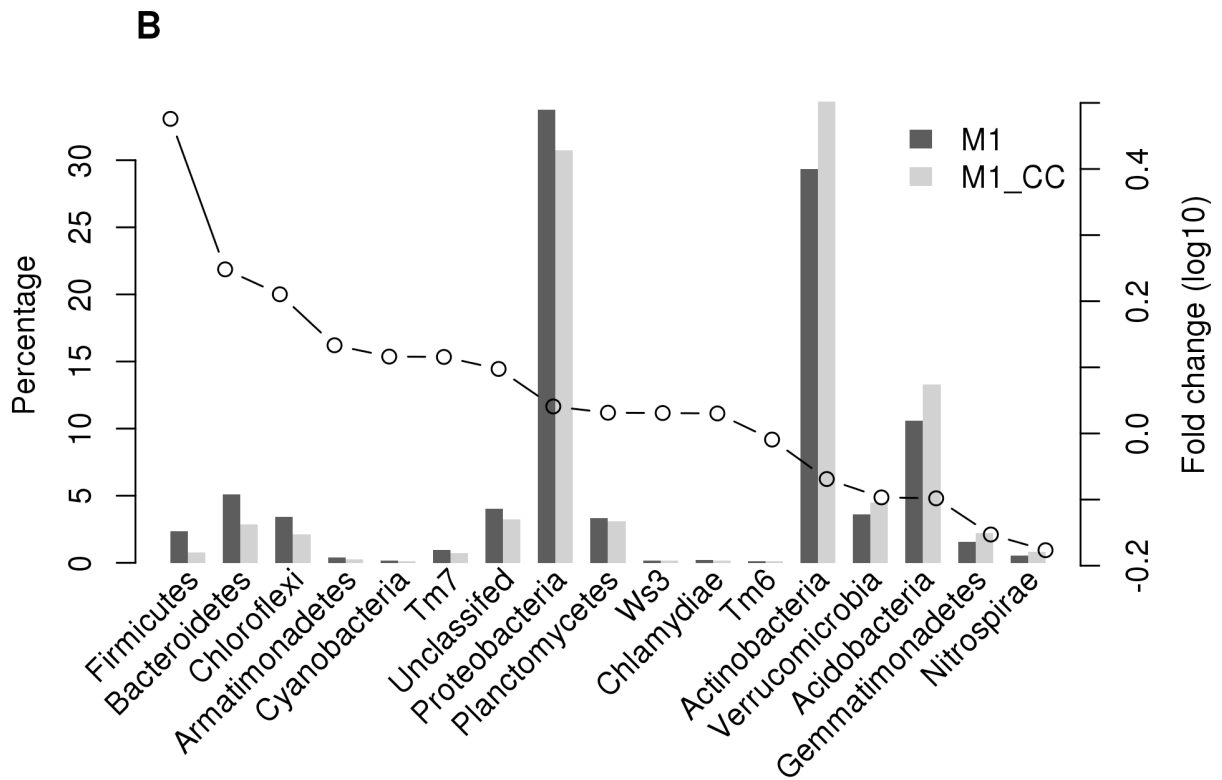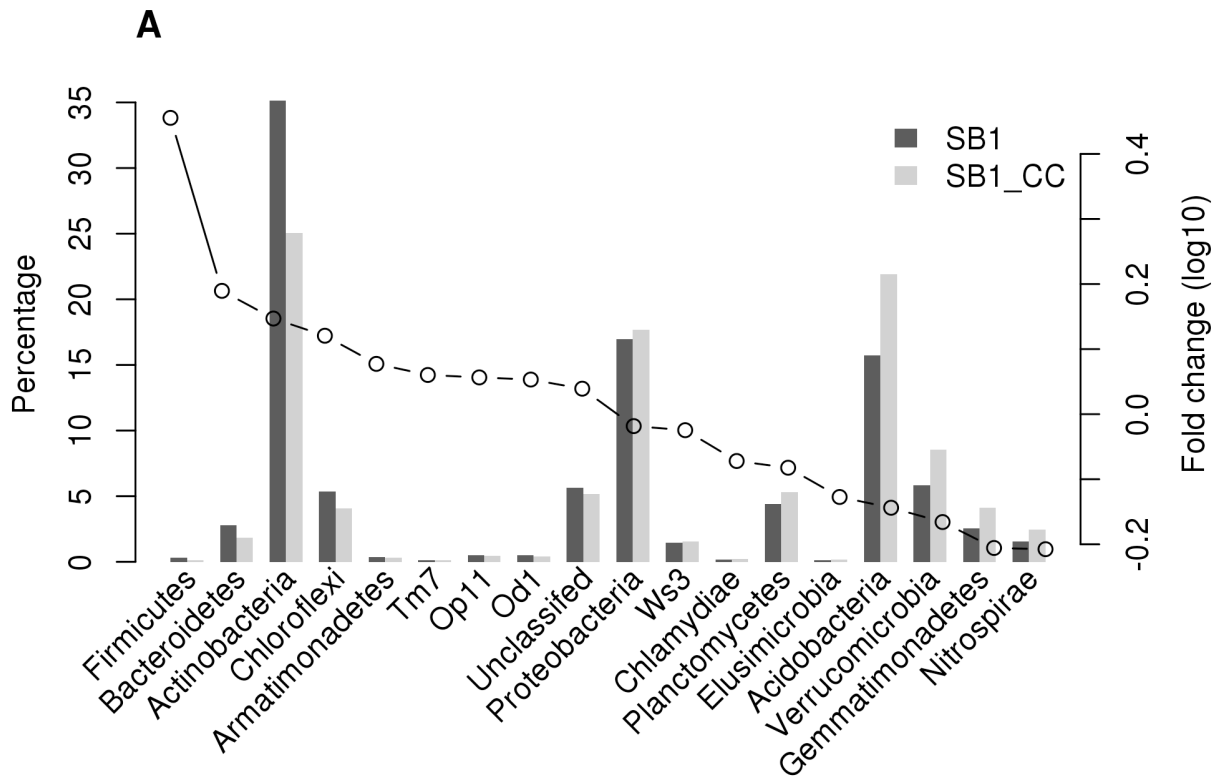taxonomy from LSU rRNA SILVA database and the rest are classified by SSU rRNA database.

137

138

139

140

141

**A**



**B**

142

143    **FIG S8** Bacterial phylum level taxonomy summary before and after SSU rRNA gene copy

144    correction. Left vertical axis with bar plot shows percentage in total community, while right

145    vertical axis with line plot shows fold change after copy number correction. Taxa with relative

146    abundances of more than 0.1% before copy correction were chosen and were ordered based on

147    fold change. Subfigure A is for bulk soil sample (SB1) and B is for Miscanthus rhizosphere

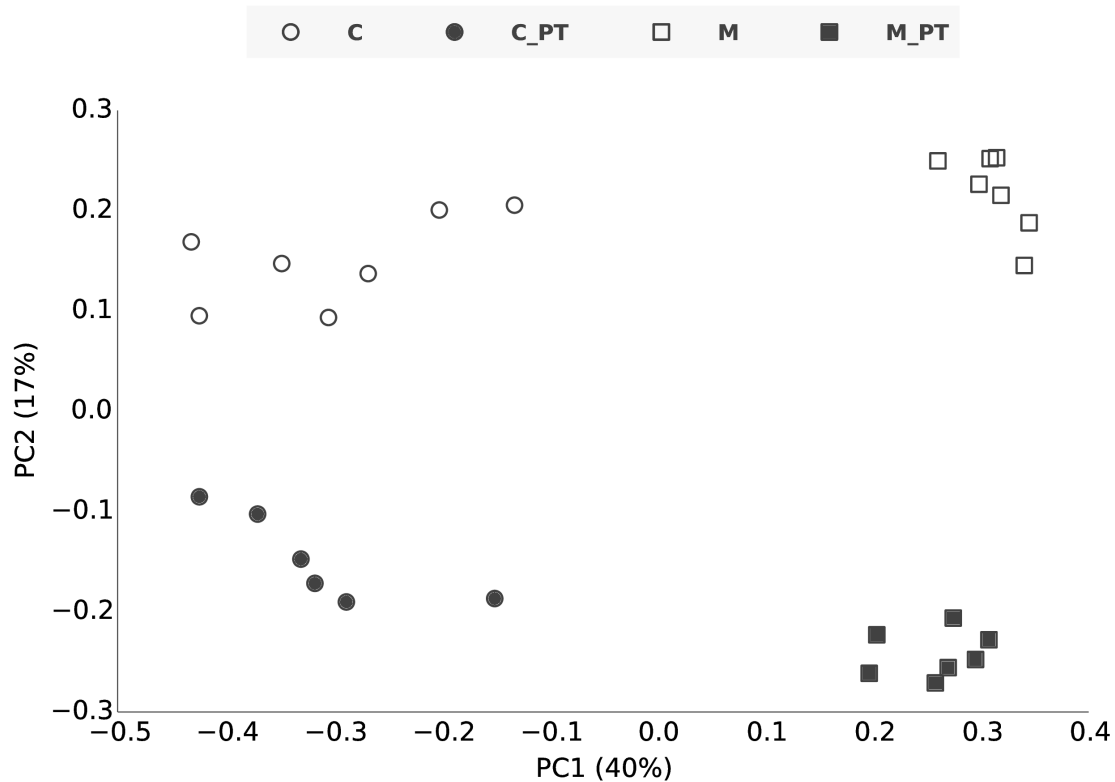148    sample (M1).

149

150

151

152

153

154

155

156

157

158

159

160

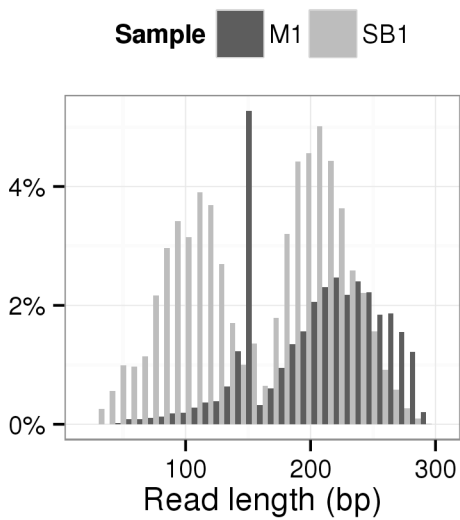161

162

163



164

165      **FIG S9** Ordination plot of amplicon and shotgun derived data after copy correction. The plot is

166      similar to the one without copy correction (Fig. 3). There were significant differences in

167      amplicon and shotgun derived data (y-axis) and of corn and Miscanthus rhizosphere samples (x-

168      axis), (AMOVA p-value < 0.001), after copy number correction. PCoA was applied to OTU table

169      resulting from *de novo* clustering with shotgun data and amplicon data using 150bp of V4 region.

170      The filled markers ("_PT") are amplicon data and the unfilled markers are shotgun data.

171

172

173

174



175

176  **FIG S10** Length distribution of trimmed reads after quality trimming and paired-end merging.

177  SB1 is the bulk soil data and M1 is the rhizosphere data. The reads with >150 bp result from the

178  merged paired ends, which benefits classification and clustering in downstream analyses. Reads

179  less than 150 bp are also used in the analysis and come from unmerged paired reads.

180

181

182

183

184

185

186