1  Full title:

2  **New insights in host-associated microbial diversity with broad and**

3  **accurate taxonomic resolution.**

4

5  Running title:

6  Broad and precise microbiome structure resolution

7

8  Matthew T. Agler[1], Alfredo Mari[1], Nina Dombrowski[1], Stéphane Haquard[1] and

9  Eric M. Kemen[1*]

10

11  [1] Max Planck Institute for Plant Breeding Research, Cologne, Germany

12

13  **\* Corresponding author:**

14      Eric M. Kemen
15      Max Planck Research Group Fungal Biodiversity
16      Max Planck Institute for Plant Breeding Research
17      Carl-von-Linne Weg 10
18      50829 Cologne
19      Germany
20      E-mail: kemen@mpipz.mpg.de
21      Tel: +49-(0)221-5062-317
22      Fax: +49-(0)221-5062-353

23

24  **Total word count: 5265**

25  Introduction: 1137

26  Materials and Methods: 1504

27  Results: 1383

28  Discussion: 1161

29  Acknowledgement: 80

30

31  **Figures:**

32  Figure 1, Figure 2, Figure 3, Figure 4 and Figure 5, Figure 6 (all in color)

33

34  **Supporting Information:**

35      SI, File S1 (oligonucleotide sequences), File S2 (Data for Figure 6)

36

## Summary

- Deep microbiome profiling has sparked a revolution in biology, recontextualizing mechanisms such as macroorganismal adaptation and evolution. Amplicon sequencing has been critical for characterization of highly diverse microbiomes, but several challenges still hinder their investigation: (1) Poor coverage of the full diversity, (2) Read depth losses and (3) Erroneous diversity inflation/distortion.

- We developed a modular approach to quickly profile at least 8 interchangeable loci in a single sequencing run, including a simple and cost-effective way to block amplification of non-targets (like host DNA). We further correct observed distortion in amplified diversity by phylogenetically grouping erroneous OTUs, creating a phylogeny-based unit we call OPUs.

- Our approach achieves full, accurate characterization of a mixed-kingdom mock community of bacteria, fungi and oomycetes at high depth even in non-target contaminated systems. The OPU concept enables much more accurate estimations of alpha and beta diversity trends than OTUs and overcomes disagreements between studies caused by methodology. Leveraging the approach in the *Arabidopsis thaliana* phyllosphere, we generated to our knowledge the most complete microbiome survey to date.

- Microbiomes are extremely diverse, extending well beyond bacteria and fungi. Our method makes new questions in a variety of fields tractable with accurate, systems-based overviews of microbial community structures.

- 

**Keywords:**

microbiome, erroneous diversity, microbial diversity, holobiont, amplicon sequencing, bacteria, fungi, oomycete, protist

## Introduction

73  A revolution in biology is currently underway as our understanding of various
74  systems is brought into the context of newly characterized structures and
75  roles of symbiotic microbial consortia. This transformation is the result of
76  growing research on microbiota associated with various abiotic or biotic
77  systems (270 vs. 3494 publications had the words "Microbial community" or
78  "Microbiome" in the title in 2005 vs. 2015 according to a PubMed search on
79  Mar 22, 2016). Strong interest in this field is not surprising considering that
80  research is turning up important roles of the community context of
81  microorganisms in systems as diverse as biotechnological transformations
82  (Werner *et al.*, 2011) and plant and animal health and fitness (Hehemann *et
83  al.*, 2010, Mills *et al.*, 2013, Panke-Buisse *et al.*, 2015, Rolli *et al.*, 2015).

84  A typical approach employed by microbiome researchers is first to
85  characterize microbial community structures in a system of interest. To do so,
86  many studies rely on amplicon sequencing of phylogenetically informative
87  genomic loci to generate microbiota profiles. These profiles are then linked to
88  specific experimental parameters, host phenotypes or performance
89  measurements. Community profiling based on ribosomal gene phylogeny
90  dates to Pace and colleagues (Stahl *et al.*, 1985) who in 1985 reported on
91  isolating and sequencing the 5S rRNA gene in environmental samples to
92  identify abundant but uncultured bacteria. The technology has come a long
93  way: As with many recent important developments in biology, rapid and
94  inexpensive DNA sequencing technology has been an enabling force in
95  microbiome research. Its democratization, however, is due to development of
96  highly parallel library indexing which made high-throughput amplicon
97  sequencing extremely inexpensive on a per sample basis (Hamady *et al.*,
98  2008).

100  Today, with for example the MiSeq platform, amplicon libraries are routinely
101  and rapidly generated from hundreds of samples and sequenced together in a
102  single run. This process generates millions of sequences up to 600 bp in
103  length (Caporaso *et al.*, 2012), enabling extremely deep profiling of targeted
104  microbial groups. In the first experiment of the current study, we used the
105  Illumina MiSeq to characterize *A. thaliana* root compartments more deeply

3

106    than we could previously with 454 pyrosequencing (Schlaeppi *et al.*, 2014) in

107    hope of gaining new insights. We show that better diversity recovery with the

108    Illumina protocol, not read depth, enabled better differentiation of soil and

109    rhizosphere compartments. In addition to the need to maximize microbial

110    diversity coverage, we identified two other problems limiting characterization

111    of diverse microbial communities: (1) Losses to read depth because of non-

112    target amplification and (2) Artificial inflation/distortion of diversity due to

113    erroneous OTUs.

114    Limited profiling of diversity extends well beyond bacteria, since microbiomes

115    are often composed of species from all kingdoms of life. These cohabiting

116    members interact with the environment and influence one another via direct

117    associations (Fisher & Mehta, 2014) or indirectly via a host (Hajishengallis,

118    2015). To resolve these interactions and model microbial community

119    dynamics, robust systems approaches are needed (Lima-Mendez *et al.*,

120    2015). For example, analysis of modularity in microbial correlation networks

121    (i.e., co-occurring groups of microbes) has revealed rice root-associated

122    prokaryotes involved in methane cycling (Edwards *et al.*, 2015) as well as

123    modules of fungi and bacteria that together correlate with certain soil

124    parameters (de Menezes *et al.*, 2015). To improve the usefulness of such

125    approaches, some studies are profiling a larger diversity like bacteria and

126    fungi simultaneously (Marupakula *et al.*, 2016). Such approaches can reveal,

127    for example, keystone species that underlie microbial community structures

128    because they interact heavily, linking external abiotic and biotic sources of

129    variation to the community (Berry & Widder, 2014). Recent studies in

130    phyllosphere microbial communities (Agler *et al.*, 2016) and in ocean samples

131    (Chow *et al.*, 2014) have emphasized that keystone microbes participate

132    heavily in inter-kingdom interactions. Thus, broad coverage of diversity is

133    critical to pinpoint these important microbes in community surveys.

134    Parallel amplification and sequencing of multiple loci is one way to cover more

135    diversity and for this approach many well-characterized targets are available.

136    Among other target loci, structures of communities can be probed *via* the 16S

137    rRNA gene (Baker *et al.*, 2003), the internal transcribed spacer (ITS) region 1

138    or 2 (Blaalid *et al.*, 2013) or the 18S rRNA gene (Hugerth *et al.*, 2014) for

139    prokaryotic, fungal, and other eukaryotic microbes, respectively. For most of

140 these targets, many possible primer sets are available, each bringing their

141 own biases and specificities. Therefore, including multiple loci from a single

142 gene target can be advantageous and provide complementary information

143 (Wang *et al.*, 2016). Whatever the target choice, modular methods are

144 needed to quickly adapt methodology to specific questions because

145 represented diversity varies considerably between different microbial

146 communities. Previously, we developed a method to prepare, sequence and

147 analyze two loci from each of bacteria (16S), fungi (ITS) and oomycetes (ITS)

148 in parallel and used it to evaluate microbial structure and interactions in the

149 *Arabidopsis thaliana* phyllosphere (Agler *et al*, 2016). Here, we use a mixed

150 mock community of microorganisms to optimize it for high throughput,

151 resolution and accuracy.

152 We also address the two other major barriers to full characterization of

153 microbial communities. The first is amplification of host or non-target

154 organismal DNA such as mitochondrial 16S, chloroplast 16S or non-target

155 genomic ITS sequences that can lead to major loss of useful reads (Bulgarelli

156 *et al.*, 2012, Ihrmark *et al.*, 2012). We developed "blocking oligos" that

157 inexpensively and nearly completely eliminated non-target amplification in

158 mock communities with simulated "host" contamination. We also show that

159 their use does not bias results and that they can be adapted to block

160 amplification of undesirable microbial targets. Second, we addressed false

161 trends in recovered microbial community diversity arising because of

162 sequencing errors (Kunin *et al.*, 2010). Here, we introduce the "operational

163 phylogenetic unit" (OPU) – phylogenetic groupings of erroneous OTUs. We

164 show that this method resolves differences between our 454 and Illumina

165 methods caused by errors. We also demonstrate that OPUs can be used to

166 generate a phylogenetic beta diversity distance metric even for fungal ITS

167 reads and that they are a much more accurate direct measure of species

168 richness than OTUs.

169 Finally, we leverage all of these benefits to profile microbes associated with

170 leaves of wild *A. thaliana* plants. To demonstrate full modularity and that any

171 desirable locus could be included, we expand to target microbial eukaryotes

172 with two loci of the 18S rRNA gene (8 loci total). The result is deep profiling of

173 all 8 loci and to our knowledge the most complete picture to date of diversity

174 in a microbial community. We provide all tools and information needed for
175 researchers to analyze up to 50 samples with the 8 loci described here or to
176 expand the system for their needs. Together, these simple solutions will
177 enable researchers to rapidly, accurately and nearly completely characterize
178 many microbial communities in a single sequencing run. We expect these
179 methods to broaden the applicability and impact of amplicon sequencing
180 experiments.
181

182 **Materials and Methods**

183 ***Comparison of methods for amplicon sequencing***

184 We first tested an amplification and sequencing protocol for the Illumina
185 MiSeq that includes amplicon generation in two PCR steps (**Fig. S1**). We
186 prepared bacterial 16S rRNA gene libraries from 3 bulk soil, 3 plant
187 rhizosphere and 2 plant root samples from the 'Eifel' natural site experiment
188 (Experiment 1, **Table S1**) and sequenced them as described in **SI methods**.
189 We combined Illumina data with data from the same samples previously
190 generated using 454 technology (Schlaeppi *et al*, 2014) and generated OTUs
191 as described in **SI methods**. We then summarized OTUs by taxonomy and
192 generated plots at the phylum level (all taxa) or the family level (the 20 most
193 abundant taxa). We compared the cumulative OTU discovery vs. depth
194 between the two technologies by rarifying tables at read depth intervals of 10
195 between 0 and 100 and 100 between 100 and 3000 and counting the number
196 of unique and shared OTUs generated by each technology. Finally, we
197 compared the ability of the technologies to discriminate between
198 compartments of *Arabidopsis thaliana* roots by plotting boxplots of Bray-Curtis
199 or weighted UniFrac distances between sample classes using phyloseq
200 (McMurdie & Holmes, 2013) and ggplots2.

201 ***Optimizing modular, multi-locus library preparation and sequencing***

202 We expanded the method used in the first experiment to target multiple loci in
203 a single sequencing run (Experiment 2, **Table S1** and **Table S2**). Accuracy
204 was tested by amplifying mixed kingdom mock communities (**Table S3**) in 6
205 separate PCR reactions targeting two loci from phylogenetically informative
206 regions of each of bacteria (16S rRNA V3-V4 and V5-V7), fungi (ITS1 and 2)

207  and oomycetes (ITS1 and 2). We tested effects of library preparation

208  methodology by performing PCR in one step (35 cycles) or two steps (10 then

209  25 cycles or 25 then 10 cycles) (**Fig. S1b** and **Table S1**). For two-step

210  preparations, the primers used in the first step consisted of unmodified

211  universal amplification primers (**Fig. 1a**). For single-step preparations and for

212  the second step in two-step preparations, primers were a concatenation of the

213  Illumina adapter P5 (forward) or P7 (reverse), an index sequence (reverse

214  only), a linker region, and the universal primer for the region being amplified

215  (**Fig. 1b**, **Fig. S1a** and **File S1**). Details of all PCR steps can be found in the

216  **SI methods**. Libraries were purified, quantified, and combined in equimolar

217  concentrations. Sequencing was on a single Illumina MiSeq lane (Illumina,

218  Inc.) by adapting the approach of Caporaso et al. (Caporaso *et al*, 2012) for

219  multiple loci (**Fig. S1c** and **File S1**). This recovers ~8% more high quality

220  bases than protocols relying on Illumina sequencing primers (calculated in **SI**

221  **Note**).

222  Details on generating OTU tables and taxonomy from raw multi-locus data

223  can be found in the **SI Methods**. We summarized bacterial, fungal and

224  oomycete OTU tables by taxonomic rank, converted abundances to relative

225  values and plotted the family-level taxonomic distribution directly from this

226  data with the package ggplots2 in R. To calculate distances of samples from

227  expected, we added the expected distributions (**Table S3**) to the OTU tables

228  and summarized taxa at the family level. After removing "host"-derived reads,

229  we calculated Bray-Curtis distances between samples using Vegan (Oksanen

230  *et al*., 2013). We plotted distances from expected distributions in boxplots

231  using ggplots2. Each box represents three "replicate" libraries generated with

232  the same mixed kingdom mock community template but with differing

233  amounts of "host" DNA added.

### *Avoiding non-target template amplification with "blocking oligos"*

235  To make the method applicable to host-associated studies, we addressed

236  non-target amplification in library preparation. In short, primers specific to the

237  known, undesirable template (hereafter "blocking oligos") are designed to bind

238  nested inside the universal primer binding sites (Experiment 3, **Table S1**).

239  Thereby, most amplicons made in the first PCR step from non-target template

240  are short and lack the universal primer sequences. These cannot be

241   elongated in the second PCR step and subsequently are not sequenced (**Fig.**
242   **1b**, **Fig. S1b** and **S1c**). Blocking oligos were designed for the *A. thaliana*
243   chloroplast (16S rRNA V3-V4 region) or mitochondria (16S rRNA V5-V7
244   region) and the *A. thaliana* ITS1 and ITS2 regions by adapting the approach
245   of Lundberg et al. (Lundberg *et al.*, 2013) originally for PNA clamps. See **SI**
246   **Methods** and **Fig. S2** for details of design and use in library preparation and
247   **File S1** their sequences. To analyze the percent reduction in host plant-
248   associated reads when blocking oligos were employed, we considered the
249   relative abundance of reads associated with the class "Chloroplast" or the
250   order "*Rickettsiales*" in the 16S OTU tables and reads in the kingdom
251   "*Viridiplantae*" in the ITS OTU tables in samples with *A. thaliana* DNA and with
252   and without blocking oligos.

253   ***Clustering OTUs by phylogeny into OPUs***

254   For the 454/Illumina comparison and multi-kingdom mock community data,
255   OTUs were clustered into phylogenetically closely related groups that we
256   called operational phylogenetic units (OPUs, **Fig. 1c,** Experiment 4 in **Table**
257   **S1**). In short, OTUs were divided at the rank of family, combined with
258   sequences from the taxonomy reference databases and a phylogenetic tree
259   was built for each by alignment with MUSCLE (Edgar, 2004). UPGMA trees
260   for each family were created with the R function hclust. The tree was
261   dynamically split into clusters using the hybrid method in cutreeDynamic in the
262   dynamicTreeCut package (Langfelder *et al.*, 2008) in R. This method was
263   designed to identify clusters in trees similar to hierarchical clustering but
264   without predetermined clustering depths. It dynamically identifies groups of
265   tips in a dendrogram that form clusters using both the tree and the distance
266   matrix that the tree is based on. A set of user-defined parameters define the
267   cluster detection sensitivity and we found that setting the minimum cluster
268   size to 15 and the deepSplit parameter to 3 was effective for OTU clustering.
269   We then generated a map of the OTUs in each OPU and generated an OPU
270   abundance table.

271   For the 454/Illumina data, overlap of OPU generation between technologies
272   and Bray-Curtis distance plots were generated exactly as described above for
273   OTUs. For mock communities, species richness estimates were based on
274   data from the evenly distributed mock community template with *A. thaliana*

275    "host" contamination, amplified in 2 steps (10 cycle / 25 cycle) with blocking

276    oligos. We used QIIME 1.8.0 (Caporaso *et al.*, 2010) to calculate the number

277    of observed species in 10 rarefactions at 30 evenly spaced depths based on

278    the OPU table, the OTU table, and the tables of OTUs grouped taxonomically

279    at levels species, genus, family and order. The maximum depth was based on

280    the OPU read depth since a few reads were discarded during OPU generation

281    (Bacteria V3/V4: 2530, V5/V6/V7: 34780, Fungi ITS1: 9930, ITS2: 48400,

282    Oomycete ITS1: 26820, ITS2: 5000). We plotted the average number of

283    observed species against sequencing depth for the bacterial 16S V3-V4

284    dataset and the ratios of observed:expected OTUs and OPUs for all datasets.

285    ***Characterizing* A. thaliana *phyllosphere microbiota***

286    We used the multi-locus approach to characterize the phyllosphere

287    microbiome of *A. thaliana* leaves infected with the oomycete pathogen *Albugo*

288    *laibachii* with near-complete taxonomic coverage (Experiment 5 in **Table S1**

289    and **Table S2**). Whole leaves (defined as a single whole rosette) or

290    endophytic fractions of leaves (defined as in (Agler *et al*, 2016)) were

291    collected in the wild (a total of 18 samples - 9 whole leaf, 9 endophyte) and

292    were immediately frozen on dry ice. DNA extraction was performed as

293    described previously (Agler *et al*, 2016). Library preparation, sequencing and

294    analysis was performed as described above. To more completely cover

295    eukaryotic microbial diversity, we expanded the 6 loci method to 8 with two

296    additional 18S rRNA gene loci (V4-V5 and V8-V9, see **File S1**).

297    To reduce *A. thaliana* or *A. laibachii* amplification in the 18S region we

298    designed additional blocking oligos for both of these organisms (**File S1**). We

299    tested them by preparing 18S amplicons from two mock communities

300    consisting of *A. thaliana* (97% or 87%), *A. laibachii* (0 or 10%)*, Sphingomonas*

301    sp. (1.5%), *Bacillus* sp. (1.5%) and 0.001% to 1% of target *Saccharomyces*

302    cerevisiae. (**Table S4**).

303    To provide a complete and concise picture of the diversity of microbiota

304    inhabiting *A. thaliana*, we combined the data from all samples. To visualize

305    data, we assigned taxonomy to OTUs and generated two phylogenetic trees

306    where branches represent unique genera. Trees were generated from the

307    taxonomic lineages (*not* OTU sequence similarity) with the ape package in R

308    and output as newick files (Paradis *et al.*, 2004). Therefore, OTUs from taxa

309   not represented in the databases are simply grouped as "Unclassified". These

310   were uploaded to iTOl v3.1 (Letunic & Bork, 2016) to color branches by

311   taxonomy or by targeted regions. The first tree, for Eukaryotes, includes data

312   from the 18S and ITS targeted regions. The second tree includes data from

313   the 16S targeted regions.

314   ***Data Availability and Figure Regeneration***

315   Raw sequencing data is being made publicly available *via* Qiita

316   (https://qiita.ucsd.edu/) study number 10408 and is currently available for

317   direct                                         download                                         at:

318   http://bioinfo.mpipz.mpg.de/download/MethodPaper_Share/.         All        modified

319   databases, OTU tables and metadata files, as well as scripts and instructions

320   to generate OPUs and recreate the main figures are available at Figshare

321   (https://figshare.com/s/07b3493d1f6442d34dfd).

322

323   **Results**

324   ***Pattern recovery depends on diversity coverage but is obscured by***

325   ***erroneous OTUs***

326   We first re-analyzed the 454-generated bacterial 16S data from (Schlaeppi *et*

327   *al*, 2014), confirming that rhizosphere and soil compartments from *A. thaliana*

328   roots were weakly distinct (**Fig. 2a**). We hypothesized that because of their

329   relatively high alpha diversity, higher read depth was needed to differentiate

330   microbiota between the compartments. Thus, we reanalyzed the same set of

331   samples at higher depth (86,406-211,907 reads/sample vs. 12,699-20,844

332   reads/sample previously) with our protocol for the Illumina platform (**Fig. 1a**

333   **and 1b and Table S1**). Bray-Curtis distances, which consider all OTUs

334   equally, suggested that the Illumina method indeed better distinguished

335   rhizosphere and soil compartments (**Fig. 2a**). Surprisingly, this was depth-

336   independent and was also true for differences between other compartments.

337   This was apparently driven by widely divergent OTU profiles with only 35% of

338   all OTUs observed in both datasets (**Fig. 2b**, 700 and 1424 OTUs were

339   unique to 454 and Illumina, respectively). Huge numbers of unique OTUs

340   suggested that either: (1) Differences in the methods of library generation and

341   sequencing resulted in little overlap or (2) OTUs were inflated by error. To

10

342  check this, we calculated between-sample weighted UniFrac distances, which

343  gives less importance to differences caused by closely related, likely

344  erroneous OTUs (**Fig. 2a)**. With this metric only soil and rhizosphere

345  compartments were better differentiated in the Illumina dataset, suggesting a

346  mixture of real differences and erroneous OTUs leading to false diversity

347  trends. True differences could be due to higher sensitivity with the Illumina

348  method to the phyla *Verrucomicrobia*, *TM7* and *Chloroflexi*, which were more

349  abundant in that dataset (**Fig. S3**). However, there were apparently too many

350  errors to locate soil/rhizosphere differential OTUs with certainty. In any case,

351  the role of improved taxonomic resolution in detecting fine differences

352  between datasets motivated us to expand to target loci beyond prokaryotes.

353

354  ***A fully modular, multi-locus approach to improve insight into microbial***

355  ***diversity***

356  We previously (Agler *et al*, 2016) adapted our 2-step Illumina amplicon library

357  generation protocol to simultaneously target 6 genomic loci, two from each of

358  bacteria, fungi and oomycetes. Here, we optimized the protocol by extensively

359  testing variations of it on a mock community consisting of microbes from the

360  three kingdoms (**Table S3**). We found that for all three kingdoms, taxa

361  distributions (shown at the order level in **Fig. 3a and Fig. S4a**) were similar to

362  expected. Mocks with staggered distributions of microorganisms were

363  generally closest to expected (**Fig. 3b and Fig. S4b)** because the effect of

364  underestimated taxa was sometimes stronger in even communities (e.g., the

365  order *Mucorales* was not efficiently recovered by fungal ITS1 primers **Fig. 3a**

366  **and Tables S5-S7**). Recovered community structures were reproducible,

367  since the distance of technical replicates from the expected distribution was

368  consistent (**Fig. 3b and Fig. S4b**). 2-step amplification recovered microbial

369  community structures that were closer to the expected than 1-step

370  amplification although the trend was not significant in all datasets. Further,

371  leaving the bulk of PCR cycling for the second step (10 cycles followed by 25

372  cycles) tended to give the most accurate results. These close-to-expected

373  taxonomy distributions were based on OTUs grouped by taxonomy. At the

374  OTU level we again observed inflated diversity due to erroneous OTUs. OTUs

375  overestimated species richness by on average 257.5%, 2575% and 387.5%

376    (only considering OTUs in the expected taxa) for bacteria, fungi and
377    oomycetes, respectively (**Tables S5-S7**).

378    We tested applicability of our method to host-associated microbiomes by
379    mixing 90% *A. thaliana* "host" DNA and 10% mock communities (**Fig. 3**). Non-
380    target host-derived DNA amplification accounted for up to 94% of reads
381    (chloroplast-derived in the 16S V3-V4 dataset) and much less but still
382    significant amounts in other target regions (**Fig. 3a and Fig. S4**). Therefore,
383    we developed and implemented "blocking oligos" to reduce amplification of
384    non-target DNA template. This method largely recovered read depth by
385    eliminating 60 - 90% of chloroplast contamination in bacterial 16S
386    communities and nearly all of the small amount of contamination in fungal ITS
387    communities (**Fig. 4a**). Importantly, employing blocking oligos did not change
388    the recovered distribution of taxa (each of the 2-step amplification boxplots in
389    **Fig. 3b** included a replicate with blocking oligos but all had the same distance
390    to expected). Thus, whereas extensive host contamination would obscure all
391    but the most abundant microbes, blocking oligos enable deeper amplicon
392    sequencing to uncover rare microbiota.

393

394    ### ***Recognizing true diversity trends with phylogenetic OTU clustering***

395    Prolific generation of erroneous OTUs strongly distorted true diversity
396    patterns. Since erroneous OTUs derive from the same true sequence (**Fig.
397    1c**), they should cluster closely in phylogenetic trees. Using this principle we
398    grouped OTUs generated in the 454/Illumina comparison into a unit that we
399    call operational phylogenetic units (OPUs). This approach reduced the total
400    from 3268 OTUs to 293 OPUs. OPUs properly grouped divergent erroneous
401    OTUs generated with 454 and Illumina since overlap between them increased
402    from 35% (OTUs, **Fig. 2b**) to 90% (OPUs, **Fig. 5a**). There were only 11 and
403    21 OPUs unique to 454- and Illumina, respectively. Bray-Curtis distances
404    based on OPUs (**Fig. 5b**) closely resembled UniFrac distances based on
405    OTUs (**Fig. 2a**) where Illumina only better distinguished soil and rhizosphere
406    compartments. We determined that 16 OTUs (maximum 150.8 reads/sample)
407    and 10 OPUs (maximum 213.7 reads/sample) significantly contributed to
408    observed soil/rhizosphere differentiation (**Table S8**). Significant OTUs and
409    OPUs were in agreement, since both were dominated by the phylum

12

410   *Chloroflexi* (4 of 6 OTUs with > 20 reads/sample and the most abundant OPU,

411   **Table S8**).

412   We used the mock community dataset to check if OPUs can provide accurate

413   diversity estimates in all target loci. Indeed, all OPU rarefaction curves quickly

414   reached an asymptote close to expected species richness (**Fig. 5c**).

415   Considering only expected families, species richness was estimated at 110%,

416   87.5% and 87.5% of expected for bacteria, fungi and oomycetes, respectively

417   (**Tables S5-S7**). For fungal and oomycete ITS2 data, richness was estimated

418   perfectly in all families, and the same was true for most families in other

419   datasets (**Tables S5-S7**). Considering all discovered OPUs and OTUs

420   (including non-targets and contaminants), phylogenetic grouping reduced the

421   average total number of observed units from 86 to 30.5 (OTUs to OPUs) for

422   bacteria (20 expected), 121.5 to 7 for fungi (4 expected) and 26 to 7.5 for

423   oomycetes (4 expected) (**Fig. 5d and Tables S5-S7**). Comparatively, OTUs

424   and even OTUs grouped by taxonomy extensively overestimated diversity and

425   their non-asymptotic rarefaction curve suggests continued inflation with

426   deeper sequencing (**Fig. 5c**). Overall, OPUs contribute to drastically improved

427   microbial diversity profiles from amplicon sequencing data.

428

429   ***Towards a complete survey of complex host-associated microbiomes***

430   Next, we leveraged the full modularity of our approach to provide a near-

431   complete survey of prokaryotic and eukaryotic diversity in *A. thaliana* leaves

432   collected in several locations using 8 loci (2 for each of bacteria, general

433   eukaryotes, fungi and oomycetes). Since the leaves of *A. thaliana* are often

434   infected by the obligate biotroph pathogen *A. laibachii*, templates from leaves

435   are dominated by *Arabidopsis* and *Albugo* genomic DNA. To overcome non-

436   target amplification of these two organisms by universal 18S primers, another

437   set of blocking oligos were designed. We tested the oligos by preparing 18S

438   amplicons from mock community templates (**Table S4**) containing bacterial,

439   *A. thaliana*, *A. laibachii* and target *S. cerevisiae* genomic DNA. Quantification

440   (qPCR) of target levels in the 18S libraries showed that blocking non-targets

441   increased target levels between ~57x (1% target template) and ~57,000x

442   (0.001% target template) (tested in the 18S V4-V5 region, **Fig. 4b**).

13

443   Here, we present the most complete picture of the *A. thaliana*-associated
444   microbiome (and to our knowledge any microbiome) ever assembled in a
445   single amplicon sequencing run (**Fig. 6**) based on combined diversity in 24
446   leaf samples collected from three wild locations. As expected, the 18S rRNA
447   gene primers recovered a wide diversity of fungal and non-fungal eukaryotic
448   microbiota, including various algae, cercozoa and amoebozoa (**Fig. 6a**). They
449   even suggested that insects and helminthes are or were present on the
450   leaves (**File S2**). The fungal and oomycete ITS datasets complemented the
451   broader 18S data with more specificity in those groups – together, these two
452   accounted for 44% of tree tips (observed genera, **Fig. 6a**). The prokaryote
453   trees further demonstrate complementarity for primer sets targeting the same
454   groups of microbes (**Fig. 6b**). Here, 42% of observed genera were discovered
455   by both primer sets, with complementary diversity discovery especially in the
456   phyla *Cyanobacteria* (V3-V4 dataset) and *Firmicutes* (V5-V7 dataset).
457

## Discussion

459   Amplicon sequencing of phylogenetically or functionally informative loci has
460   become an indispensable technique in a variety of biology-related fields
461   because its targeted approach (compared to untargeted approaches like
462   metagenomics) allows the most accurate annotation possible by using
463   specialized databases (DeSantis *et al.*, 2006). It has revealed that microbial
464   community structuring is more complex than previously thought and
465   suggested extensive interactions between (a)biotic factors and microbes (de
466   Menezes *et al*, 2015) and between microbes even across kingdoms (Agler *et*
467   *al*, 2016, Lima-Mendez *et al*, 2015). To understand these interactions,
468   microbiome researchers need to be able to more completely characterize
469   diversity in a single sequencing run. The current method enables this by
470   targeting at least 8 loci in parallel. This drastic increase in resolution critically
471   overcomes an inherent uncertainty in systems-scale investigations of factors
472   contributing to microbial community structures.
473   A key technique enabling the current advances is employment of a two-step
474   amplicon library preparation as opposed to a single step amplification. Many
475   commonly used protocols (e.g., the Earth Microbiome Project 16S protocol

476  based on (Caporaso *et al.*, 2011)) recommend using large, concatenated
477  primers and one-step amplification. The advantages in technical
478  reproducibility of a two-step approach that excludes concatenated primers in
479  the first step were already described by Berry et al. (Berry *et al.*, 2011) for T-
480  RFLP or 454 amplicon sequencing. For Illumina sequencing, concatenated
481  primer bias was addressed with a 3-step approach: A 2-step amplification plus
482  adapter ligation (Herbold *et al.*, 2015). That approach also allowed
483  characterization of multiple gene regions, but a "head" sequence was
484  concatenated to universal primers in the first amplification step. Our use of
485  only universal primers in the first step, therefore, probably explains why mock
486  community structure recovery was very accurate and replicable. Additionally,
487  our approach adds all required adapters for sequencing in the second
488  amplification step, eliminating problems associated with adapter ligation
489  (Sambrook *et al.*, 1989).

490  Another major problem in amplicon sequencing is associated with using
491  "universal" primers that in host-associated amplified non-target species,
492  sacrificing read depth and masking diversity (Hanshew *et al.*, 2013).
493  Previously, peptide nucleic acid "clamps" were used that were highly specific
494  to non-target templates and which physically block their amplification
495  (Lundberg *et al*, 2013). These clamps work efficiently in single-step
496  amplifications, but their production is expensive, limiting rapid development
497  and deployment of multiple clamps for new loci or for blocking several non-
498  targets. Other approaches, like using oligonucleotide clamps that physically
499  block the universal primer binding site (Vestheim & Jarman, 2008) are not
500  applicable here because target and non-target binding sites are too highly
501  conserved. Alternatively, blocking oligos are versatile and cheap and can be
502  extensively tested at very low costs and adapted to virtually any target. Some
503  universal primers targeting fungal ITS amplify fungal targets more efficiently
504  than host, which explains why we had only minor host contamination in mixed
505  mock community libraries. However, when fungal templates are less abundant
506  they can significantly amplify plant ITS (Ihrmark *et al*, 2012). Because
507  blocking oligos did not bias results, it is beneficial to always include them
508  when relative abundance of target and non-target DNA is unknown.

509    One of the most persistent problems identified so far in amplicon sequencing
510    is vast inflation of OTU diversity, mostly caused by sequencing errors (Kunin
511    *et al*, 2010). Despite useful approaches to remove erroneous reads (Bokulich
512    *et al.*, 2013, Reeder & Knight, 2010) or to reduce erroneous OTUs (Edgar,
513    2013) false diversity trends are commonly observed (Sinclair *et al.*, 2015).
514    Errors explain the popularity of phylogenetics-based tools for beta diversity
515    estimation based on the Kantorovich-Rubinstein metric, such as UniFrac
516    (Evans & Matsen, 2012, Lozupone & Knight, 2005) or alpha diversity metrics
517    like Faith's PD (Faith, 1992). These weight differences in samples caused by
518    distantly related OTUs more heavily since erroneous OTUs should be
519    phylogenetically closely related. They have been very successful in identifying
520    real differences between samples even when sequencing error is high.
521    However, they are generally not applicable to loci like ITS, where extreme
522    variability makes drawing phylogenetic relationships between all sequences
523    questionable (Schoch *et al.*, 2012). Additionally, the assumptions of
524    phylogenetic approaches do not hold when distantly related microbes occupy
525    similar niches. For example, the basidiomycete yeast-like *Pseudozyma* spp.
526    are phenotypically and ecologically much more similar to species like
527    *Dioszegia* sp. (Inácio *et al.*, 2005) than to plant pathogenic members of its
528    close relative *Ustilago* sp (Lefebvre *et al.*, 2013). Therefore, complementary
529    approaches are needed that are sensitive to shifts in abundance among
530    closely related taxa but which accurately delineate true and erroneous taxa.
531    The OPU approach addresses this problem because they are in principle like
532    phylogenetic diversity metrics – very closely related (likely erroneous) OTUs
533    are grouped into a unit which can be used to generate standard beta or alpha
534    diversity metrics. Therefore, the results are less abstract than UniFrac or
535    Faith's PD (which lacks a taxonomic unit) and should be more sensitive to
536    changes in abundance of closely related taxa. The term OPU was discussed
537    elsewhere (Pernthaler & Amann, 2005) in the context of using phylogenetic
538    grouping of organisms to move away from a specific percent identity as a
539    working taxonomic unit but not as a systematic way to group erroneous
540    OTUs. This concept was implemented in approaches to dynamically group
541    amplicon reads by phylogeny based on tree cutting. Here clusters of reads
542    were identified by training on a subset of data with known taxonomies (White

543    *et al.*, 2010) or by known differences in substitution errors between or within

544    species (Zhang *et al.*, 2013). General applicability of these approaches is

545    unclear because of the major computational resources needed to cluster raw

546    reads and because inferring phylogeny among all reads is questionable at

547    some highly divergent loci like ITS (Schoch *et al*, 2012). Our implementation

548    on the other hand uses pre-clustering of reads into OTUs and taxonomic

549    groups. In this way large datasets are not a barrier because parallelization

550    can be maximized according to available resources. Further, OTUs split by a

551    taxonomic rank, e.g., family, are closely related and phylogenetic relationships

552    can be determined even at highly divergent ITS loci. Therefore, diversity

553    metrics based on OPUs represent a much needed phylogenetic method for

554    loci that are not conserved enough to build alignments for example for

555    UniFrac distances.

556    The realization of the immense complexity of biological systems – and our

557    inability to adequately describe them - has led to many important, unresolved

558    issues. For example, there is ongoing debate about what it means to view

559    macroorganisms as holobionts, since symbiotic microbiota affect host health

560    and fitness (Brucker & Bordenstein, 2013, Sharma *et al.*, 2014). Unanswered

561    questions also linger, like what causes host genotype-independent taxonomic

562    conservation of plant root microbiomes over broad geographic distances

563    (Hacquard *et al.*, 2015). The tools described here will significantly increase

564    the ability of researchers to accurately resolve microbial communities,

565    addressing one of the primary limitations to progress. Although challenges

566    remain, we expect this approach to equip researchers to make better

567    hypotheses and to address seemingly intractable questions. These advances

568    will thereby assist in increasing discovery of the important roles of microbiota.

569

## Acknowledgements

579

## **Author contribution**

581    MA and EK conceptualized the project, designed and developed the

582    methodology and wrote the manuscript. For the technology comparison, ND

583    and SH performed experiments and MA, ND and SH analyzed the data. AM

584    adapted the method to the 18S rRNA region and assessed plant-associated

585    microbial communities. MA performed all other experiments and designed and

586    wrote the scripts to process and analyze the data.

587

## References


Earth Microbiome Project 16S rRNA Amplification Protocol. Accessed on: November 22, 2015. http://www.earthmicrobiome.org/emp-standard-protocols/16s/

Agler MT, Ruhe J, Kroll S, Morhenn C, Kim S-T, Weigel D *et al* (2016). Microbial Hub Taxa Link Host and Abiotic Factors to Plant Microbiome Variation. *PLoS Biol* 14: e1002352.

Baker GC, Smith JJ, Cowan DA (2003). Review and re-analysis of domain-specific 16S primers. *J Microbiol Meth* 55: 541-555.

Berry D, Ben Mahfoudh K, Wagner M, Loy A (2011). Barcoded Primers Used in Multiplex Amplicon Pyrosequencing Bias Amplification. *Applied and Environmental Microbiology* 77: 7846-7849.

Berry D, Widder S (2014). Deciphering microbial interactions and detecting keystone species with co-occurrence networks. *Frontiers Microbiol* 5.

Blaalid R, Kumar S, Nilsson RH, Abarenkov K, Kirk PM, Kauserud H (2013). ITS1 versus ITS2 as DNA metabarcodes for fungi. *Mol Ecol Resour* 13: 218-224.

Bokulich NA, Subramanian S, Faith JJ, Gevers D, Gordon JI, Knight R *et al* (2013). Quality-filtering vastly improves diversity estimates from Illumina amplicon sequencing. *Nat Meth* 10: 57-59.

Brucker RM, Bordenstein SR (2013). The Hologenomic Basis of Speciation: Gut Bacteria Cause Hybrid Lethality in the Genus Nasonia. *Science* 341: 667-669.

Bulgarelli D, Rott M, Schlaeppi K, Ver Loren van Themaat E, Ahmadinejad N, Assenza F *et al* (2012). Revealing structure and assembly cues for Arabidopsis root-inhabiting bacterial microbiota. *Nature* 488: 91-95.

Caporaso JG, Kuczynski J, Stombaugh J, Bittinger K, Bushman FD, Costello EK *et al* (2010). QIIME allows analysis of high-throughput community sequencing data. *Nat Meth* 7: 335-336.

Caporaso JG, Lauber CL, Walters WA, Berg-Lyons D, Lozupone CA, Turnbaugh PJ *et al* (2011). Global patterns of 16S rRNA diversity at a depth of millions of sequences per sample. *Proc Nat Acad Sci* 108: 4516-4522.

Caporaso JG, Lauber CL, Walters WA, Berg-Lyons D, Huntley J, Fierer N *et al* (2012). Ultra-high-throughput microbial community analysis on the Illumina HiSeq and MiSeq platforms. *ISME J* 6: 1621-1624.

19

636
637 **Chow C-ET, Kim DY, Sachdeva R, Caron DA, Fuhrman JA (2014). Top-down**
638 **controls on bacterial community structure: microbial network analysis of**
639 **bacteria, T4-like viruses and protists. *ISME J* 8: 816-829.**
640
641 **de Menezes AB, Prendergast-Miller MT, Richardson AE, Toscas P, Farrell M,**
642 **Macdonald LM *et al* (2015). Network analysis reveals that bacteria and**
643 **fungi form modules that correlate independently with soil parameters.**
644 ***Environ Microbiol* 17: 2677-2689.**
645
646 **DeSantis TZ, Hugenholtz P, Larsen N, Rojas M, Brodie EL, Keller K *et al***
647 **(2006). Greengenes, a Chimera-Checked 16S rRNA Gene Database and**
648 **Workbench Compatible with ARB. *Appl Environ Microbiol* 72: 5069-5072.**
649
650 **Edgar RC (2004). MUSCLE: multiple sequence alignment with high accuracy**
651 **and high throughput. *Nucleic Acids Res* 32: 1792-1797.**
652
653 **Edgar RC (2013). UPARSE: highly accurate OTU sequences from microbial**
654 **amplicon reads. *Nat Meth* 10: 996-998.**
655
656 **Edwards J, Johnson C, Santos-Medellín C, Lurie E, Podishetty NK, Bhatnagar**
657 **S *et al* (2015). Structure, variation, and assembly of the root-associated**
658 **microbiomes of rice. *Proc Natl Acad Sci U S A*.**
659
660 **Evans SN, Matsen FA (2012). The phylogenetic Kantorovich–Rubinstein**
661 **metric for environmental sequence samples. *J R Stat Soc Series B Stat***
662 ***Methodol* 74: 569-592.**
663
664 **Faith DP (1992). Conservation evaluation and phylogenetic diversity. *Biol***
665 ***cons* 61: 1-10.**
666
667 **Fisher CK, Mehta P (2014). Identifying Keystone Species in the Human Gut**
668 **Microbiome from Metagenomic Timeseries Using Sparse Linear**
669 **Regression. *PLoS ONE* 9: e102451.**
670
671 **Hacquard S, Garrido-Oter R, González A, Spaepen S, Ackermann G, Lebeis S**
672 ***et al* (2015). Microbiota and Host Nutrition across Plant and Animal**
673 **Kingdoms. *Cell Host Microbe* 17: 603-616.**
674
675 **Hajishengallis G (2015). Periodontitis: from microbial immune subversion**
676 **to systemic inflammation. *Nat Rev Immunol* 15: 30-44.**
677
678 **Hamady M, Walker JJ, Harris JK, Gold NJ, Knight R (2008). Error-correcting**
679 **barcoded primers for pyrosequencing hundreds of samples in multiplex.**
680 ***Nat Meth* 5: 235-237.**
681
682 **Hanshew AS, Mason CJ, Raffa KF, Currie CR (2013). Minimization of**
683 **chloroplast contamination in 16S rRNA gene pyrosequencing of insect**
684 **herbivore bacterial communities. *J Microbiol Meth* 95: 149-155.**

685

686    Hehemann J-H, Correc G, Barbeyron T, Helbert W, Czjzek M, Michel G
687    (2010). Transfer of carbohydrate-active enzymes from marine bacteria to
688    Japanese gut microbiota. *Nature* 464: 908-912.
689

690    Herbold CW, Pelikan C, Kuzyk O, Hausmann B, Angel R, Berry D *et al* (2015).
691    A flexible and economical barcoding approach for highly multiplexed
692    amplicon sequencing of diverse target genes. *Frontiers Microbiol* 6: 731.
693

694    Hugerth LW, Muller EEL, Hu YOO, Lebrun LAM, Roume H, Lundin D *et al*
695    (2014). Systematic Design of 18S rRNA Gene Primers for Determining
696    Eukaryotic Diversity in Microbial Consortia. *PLoS ONE* 9: e95567.
697

698    Ihrmark K, Bödeker ITM, Cruz-Martinez K, Friberg H, Kubartova A, Schenck
699    J *et al* (2012). New primers to amplify the fungal ITS2 region – evaluation
700    by 454-sequencing of artificial and natural communities. *FEMS Microbiol*
701    *Ecol* 82: 666-677.
702

703    Inácio J, Portugal L, Spencer-Martins I, Fonseca Á (2005). Phylloplane
704    yeasts from Portugal: Seven novel anamorphic species in the Tremellales
705    lineage of the Hymenomycetes (Basidiomycota) producing orange-
706    coloured colonies. *FEMS Yeast Res* 5: 1167-1183.
707

708    Kunin V, Engelbrektson A, Ochman H, Hugenholtz P (2010). Wrinkles in the
709    rare biosphere: pyrosequencing errors can lead to artificial inflation of
710    diversity estimates. *Environ Microbiol* 12: 118-123.
711

712    Langfelder P, Zhang B, Horvath S (2008). Defining clusters from a
713    hierarchical cluster tree: the Dynamic Tree Cut package for R.
714    *Bioinformatics* 24: 719-720.
715

716    Lefebvre F, Joly DL, Labbé C, Teichmann B, Linning R, Belzile F *et al* (2013).
717    The Transition from a Phytopathogenic Smut Ancestor to an Anamorphic
718    Biocontrol Agent Deciphered by Comparative Whole-Genome Analysis.
719    *Plant Cell* 25: 1946-1959.
720

721    Letunic I, Bork P (2016). Interactive tree of life (iTOL) v3: an online tool for
722    the display and annotation of phylogenetic and other trees. *Nuc Acids Res*.
723

724    Lima-Mendez G, Faust K, Henry N, Decelle J, Colin S, Carcillo F *et al* (2015).
725    Determinants of community structure in the global plankton interactome.
726    *Science* 348.
727

728    Lozupone C, Knight R (2005). UniFrac: a New Phylogenetic Method for
729    Comparing Microbial Communities. *Appl Environ Microbiol* 71: 8228-8235.
730

731    Lundberg DS, Yourstone S, Mieczkowski P, Jones CD, Dangl JL (2013).
732    Practical innovations for high-throughput amplicon sequencing. *Nat Meth*
733    10: 999-1002.

734

735 Marupakula S, Mahmood S, Finlay RD (2016). Analysis of single root tip
736 microbiomes suggests that distinctive bacterial communities are selected
737 by Pinus sylvestris roots colonized by different ectomycorrhizal fungi.
738 *Environ Microbiol*.
739
740 McMurdie PJ, Holmes S (2013). phyloseq: An R Package for Reproducible
741 Interactive Analysis and Graphics of Microbiome Census Data. *PLoS ONE* 8:
742 e61217.
743
744 Mills E, Shechtman K, Loya Y, Rosenberg E (2013). Bacteria appear to play
745 important roles in both causing and preventing the bleaching of the coral
746 Oculina patagonica. *Mar Ecol Prog Ser* 489: 155-162.
747
748 Oksanen J, Blanchet G, Kindt R, Legendre P, Minchin P, O'Hara R *et al*
749 (2013). vegan: Community Ecology Package version 2.0-10. Accessed on:
750 http://CRAN.R-project.org/package=vegan
751
752 Panke-Buisse K, Poole AC, Goodrich JK, Ley RE, Kao-Kniffin J (2015).
753 Selection on soil microbiomes reveals reproducible impacts on plant
754 function. *ISME J* 9: 980-989.
755
756 Paradis E, Claude J, Strimmer K (2004). APE: Analyses of Phylogenetics and
757 Evolution in R language. *Bioinformatics* 20: 289-290.
758
759 Pernthaler J, Amann R (2005). Fate of Heterotrophic Microbes in Pelagic
760 Habitats: Focus on Populations. *Microbiol Mol Biol Rev* 69: 440-461.
761
762 Reeder J, Knight R (2010). Rapid denoising of pyrosequencing amplicon
763 data: exploiting the rank-abundance distribution. *Nature methods* 7: 668-
764 669.
765
766 Rolli E, Marasco R, Vigani G, Ettoumi B, Mapelli F, Deangelis ML *et al* (2015).
767 Improved plant resistance to drought is promoted by the root-associated
768 microbiome as a water stress-dependent trait. *Environ Microbiol* 17: 316-
769 331.
770
771 Sambrook J, Fritsch EF, Maniatis T (1989). *Molecular cloning: A laboratory
772 manual*, vol. 2. Cold spring harbor laboratory press New York.
773
774 Schlaeppi K, Dombrowski N, Oter RG, Ver Loren van Themaat E, Schulze-
775 Lefert P (2014). Quantitative divergence of the bacterial root microbiota in
776 Arabidopsis thaliana relatives. *Proc Nat Acad Sci* 111: 585-592.
777
778 Schoch CL, Seifert KA, Huhndorf S, Robert V, Spouge JL, Levesque CA *et al*
779 (2012). Nuclear ribosomal internal transcribed spacer (ITS) region as a
780 universal DNA barcode marker for Fungi. *Proc Nat Acad Sci* 109: 6241-
781 6246.
782

783  Sharma R, Mishra B, Runge F, Thines M (2014). Gene loss rather than gene
784  gain is associated with a host jump from monocots to dicots in the smut
785  fungus Melanopsichium pennsylvanicum. *Genome Biol Evol*.
786
787  Sinclair L, Osman OA, Bertilsson S, Eiler A (2015). Microbial Community
788  Composition and Diversity via 16S rRNA Gene Amplicons: Evaluating the
789  Illumina Platform. *PLoS ONE* 10: e0116955.
790
791  Stahl DA, Lane DJ, Olsen GJ, Pace NR (1985). Characterization of a
792  Yellowstone hot spring microbial community by 5S rRNA sequences. *Appl*
793  *Environ Microbiol* 49: 1379-1384.
794
795  Vestheim H, Jarman S (2008). Blocking primers to enhance PCR
796  amplification of rare sequences in mixed samples - a case study on prey
797  DNA in Antarctic krill stomachs. *Frontiers Zool* 5: 12.
798
799  Wang S, Sun B, Tu J, Lu Z (2016). Improving the microbial community
800  reconstruction at the genus level by multiple 16S rRNA regions. *Journal of*
801  *Theoretical Biology* 398: 1-8.
802
803  Werner JJ, Knights D, Garcia ML, Scalfone NB, Smith S, Yarasheski K *et al*
804  (2011). Bacterial community structures are unique and resilient in full-
805  scale bioenergy systems. *Proc Nat Acad Sci* 108: 4158-4163.
806
807  White JR, Navlakha S, Nagarajan N, Ghodsi M-R, Kingsford C, Pop M (2010).
808  Alignment and clustering of phylogenetic markers-implications for
809  microbial diversity studies. *BMC bioinformatics* 11: 152.
810
811  Zhang J, Kapli P, Pavlidis P, Stamatakis A (2013). A general species
812  delimitation method with applications to phylogenetic placements.
813  *Bioinformatics* 29: 2869-2876.
814
815
816
817

818 **Figures**

819

820 **Figure 1. Strategy to increase taxonomic coverage and accuracy of**
821 **amplicon sequencing.** A. In the first step, 8 individual PCR reactions are
822 performed per sample each targeting a specific gene region. B. Blocking
823 oligos are employed in the first PCR step which are specific for non-target
824 templates so that these cannot be elongated to the final libraries with
825 concatenated primers in the second step. P5 and P7 are standard Illumina
826 adapter sequences. L indicates the linker sequence. C. Inflation of the number
827 of observed OTUs caused largely by sequencing error is addressed by
828 dividing them up by taxonomic lineages, building a phylogenetic tree and then
829 clustering closely related members into operational phylogenetic units
830 (OPUs).

831
832 **Figure 2. A comparison of 454- vs. Illumina-based amplicon sequencing**
833 **protocols shows little overlap of OTUs, suggesting high erroneous OTU**
834 **generation.** A. Bray-Curtis distances based on OTU relative abundances
835 suggest that data recovered with the Illumina protocol significantly better
836 distinguishes all pairs of compartments. Weighted UniFrac, however,
837 suggests that the Illumina protocol only better distinguishes rhizosphere and
838 soil compartments, implying that differences between others were due to
839 closely related and probably erroneous OTUs. B. Between 454 and Illumina
840 datasets, rarefaction curves of the number of OTUs discovered with
841 increasing read depth suggest that erroneous OTUs are very common since
842 most OTUs are unique to datasets produced by one or the other technology,
843 with only about one-third of all OTUs found in both datasets. The rarefaction
844 curves are separated into OTUs that are unique to Illumina, unique to 454 or
845 shared by both technologies.

846
847 **Figure 3**. **Reproducible and accurate characterization of mock**
848 **communities of bacteria, fungi, and oomycetes by amplicon sequencing.**
849 A. Observed taxa at the order level in sequenced mock communities closely
850 matched expected communities. The taxa "Other" is primarily non-target
851 amplification from *A. thaliana* "host" DNA that was added to test blocking

24

852    oligomers which prevent "host" DNA amplification. "NA" indicates a sample

853    where sequencing depth was too low after subsampling to be included. B.

854    Distance (Bray-Curtis distance based on relative abundance of order-level

855    taxa) of sequenced communities from the expected distribution where 0 is

856    identical and 1 is unrelated. Even or staggered revers to the distribution of the

857    organisms in the mock communities (see expected distributions in A). PCR

858    Steps refers to a 1-step (35 cycles with concatenated primers) or 2-step (10 or

859    25 cycles with standard primers followed by 25 or 10 cycles with extension

860    primers containing Illumina adapters) amplification protocol. Letters indicate p

861    < 0.1 (FDR-corrected) based on pairwise t-tests between groups.

862

863    **Figure 4. Employing blocking oligomers greatly reduces non-target and**

864    **increases target yield in amplicon libraries.** A. Near-complete reduction of

865    amplification of *A. thaliana* "host" non-target plastid 16S or ITS by employing

866    blocking oligos in preparation of mock community libraries. B. Relative

867    increase of target (*Saccharomyces* sp.) 18S V4-V5 region amplicons (qPCR

868    $2^{-\Delta Cq}$ values relative to measurement without blocking oligomers) in mock

869    community libraries prepared with blocking oligomers to reduce *A. thaliana*

870    and *A. laibachii* non-target amplification.

871

872    **Figure 5. Clustering of operational taxonomic units (OTUs) into**

873    **operational phylogenetic units (OPUs) by their phylogenetic relatedness**

874    **corrects erroneous diversity discovery.** A. Between 454 and Illumina

875    datasets, the number of shared and unique OPUs and the fraction of shared

876    OPUs demonstrates that OTU clustering greatly reduces erroneous dataset

877    disagreements compared to Fig. 2B. B. Bray-Curtis distances based on OPUs

878    displays similar trends as weighted UniFrac distances with the only significant

879    differences between rhizosphere and soil compartments. C. Rarefaction

880    curves of observed units at the OTU level, various taxonomic ranks, and for

881    OPUs for bacterial 16S V3/V4 amplicon data show that unclustered OTUs and

882    most taxonomic ranks greatly overestimate the expected diversity and the

883    curves do not reach an asymptote, while OPUs quickly reach an asymptote

884    close to the expected diversity. D. Numbers of observed OTUs or OPUs vs.

885    expected units for all target regions demonstrates near-expected numbers of

886 taxa in all target regions. Data in C. and D. is generated from the evenly
887 distributed mock community with *A. thaliana* "host" DNA and using host-
888 blocking oligomers and only considers OTUs and OPUs in the expected
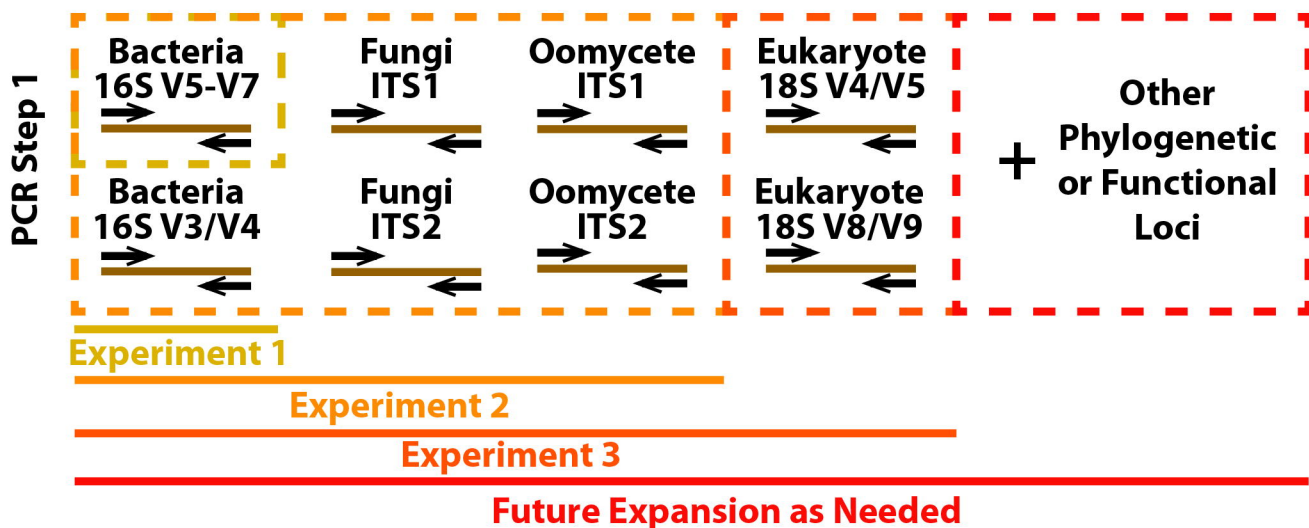889 taxonomic families.

890

891 **Figure 6. A comprehensive overview of high diversity microbiomes**
892 **inhabiting *A. thaliana* leaves revealed by parallel amplicon sequencing**
893 **of 8 loci targeting eukaryotes and prokaryotes microbes.** A. 6 loci
894 targeting eukaryotes: Two regions of the 18S rRNA gene (V4-V5 and V8-V9),
895 two regions of the fungal ITS (ITS 1 and 2) and two regions of the oomycete
896 ITS (ITS 1 and 2) revealed a diverse eukaryotic microbiota. The 18S loci
897 revealed the broadest diversity but was complemented by fungi and
898 oomycete-specific primer sets which had more detailed resolution within these
899 groups. "Target loci specificity" refers to the taxa identified with each target
900 group (Eukaryotes) or locus (Prokaryotes). B. 2 loci targeting prokaryotes:
901 Two regions of the 16S rRNA gene (V3-V4 and V5-V7) that amplify mostly
902 bacteria revealed a largely overlapping diversity profile complemented by
903 unique discovery of taxa from each of the two target regions.
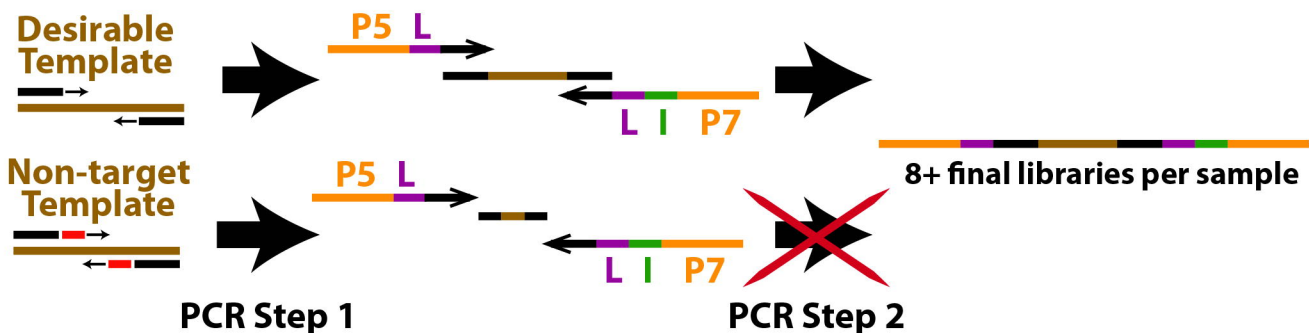
904

# A. Modular, multi-locus amplicon targeting approach

**8+ individual PCR reactions per sample**



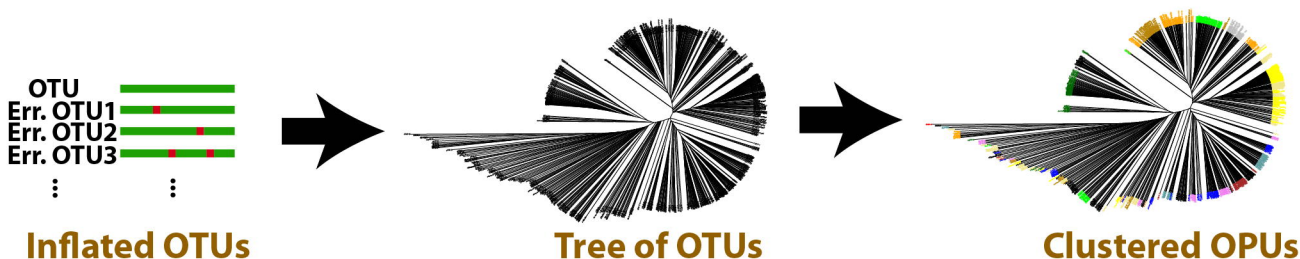# B. Blocking non-target host or microbial amplification

Specific non-target blocking oligos in PCR step 1 (for each of 8+ target regions)
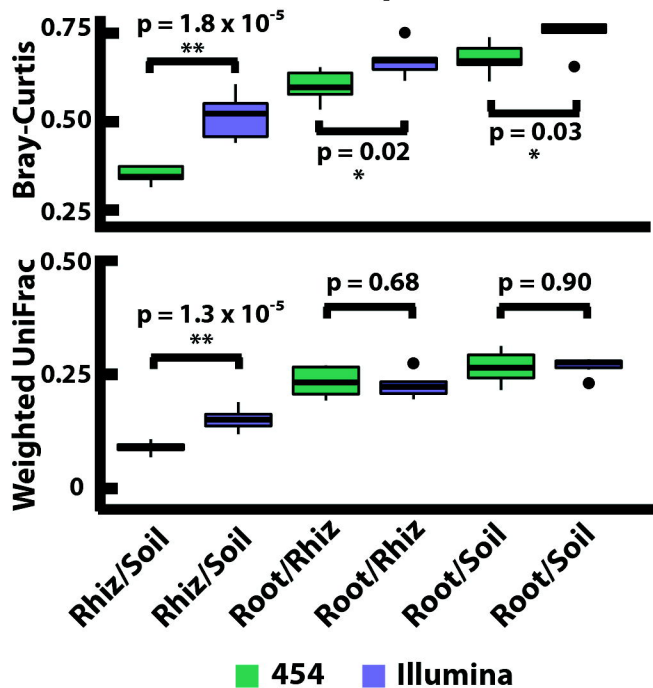Concatenated primers complementary to universal primer in PCR step 2



# C. Eliminating diversity inflation
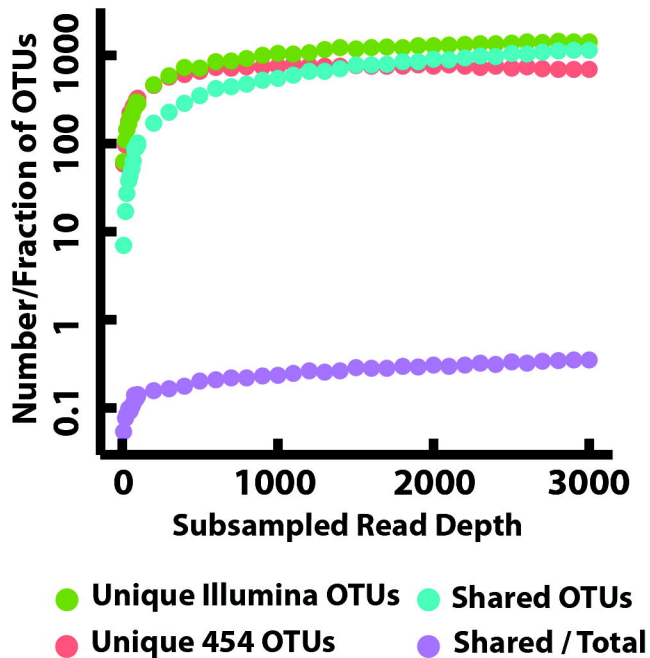
Cluster OTUs into "OPUs" by phylogenetic similarity
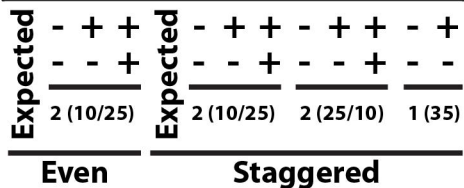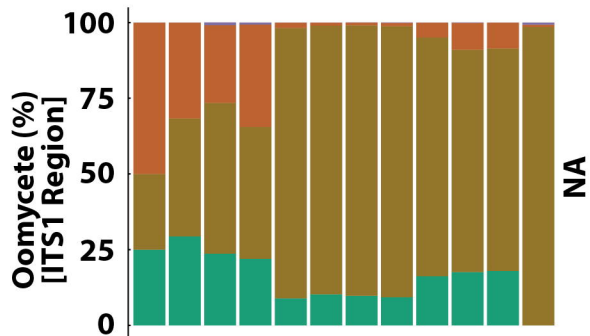
**A**

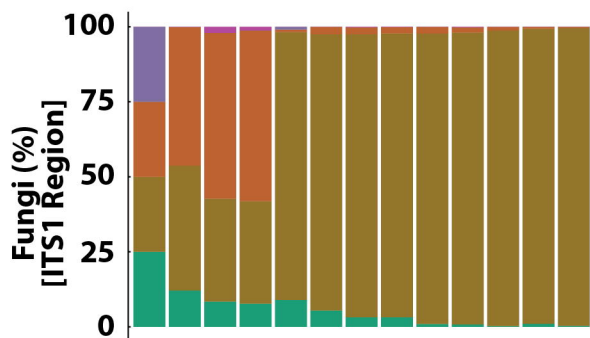**OTU-based Distances Between Compartments**

Bray-Curtis

p = 1.8 x 10⁻⁵
**

p = 0.02
*

p = 0.03
*

Weighted UniFrac

p = 1.3 x 10⁻⁵
**

p = 0.68

p = 0.90

Rhiz/Soil  Rhiz/Soil  Root/Rhiz  Root/Rhiz  Root/Soil  Root/Soil

■ 454   ■ Illumina

**B**

**454 vs. Illumina OTU Observation Curves**

Number/Fraction of OTUs

Subsampled Read Depth

● Unique Illumina OTUs   ● Shared OTUs
● Unique 454 OTUs   ● Shared / Total
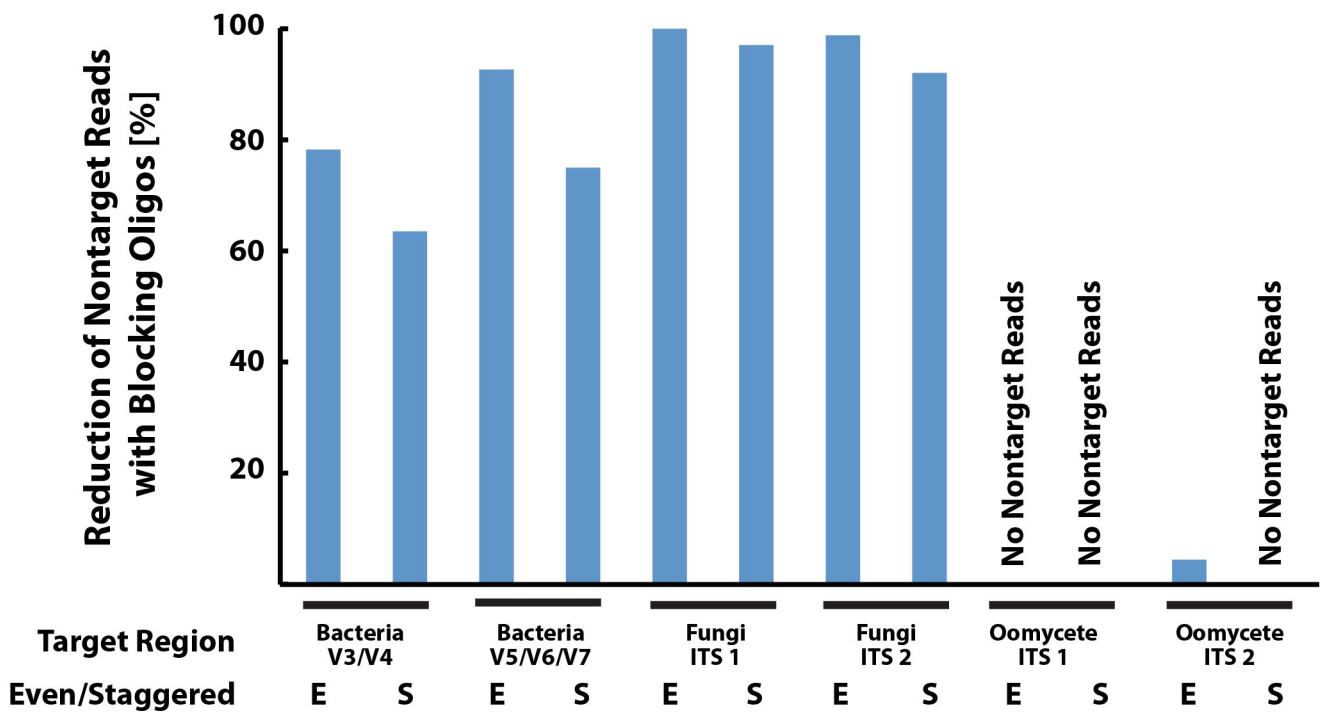
**A.** Microbial Relative Abundance Including *A. thaliana* amplicons

**B.** Distance to Expected Excluding *A. thaliana*

Expected Taxa:
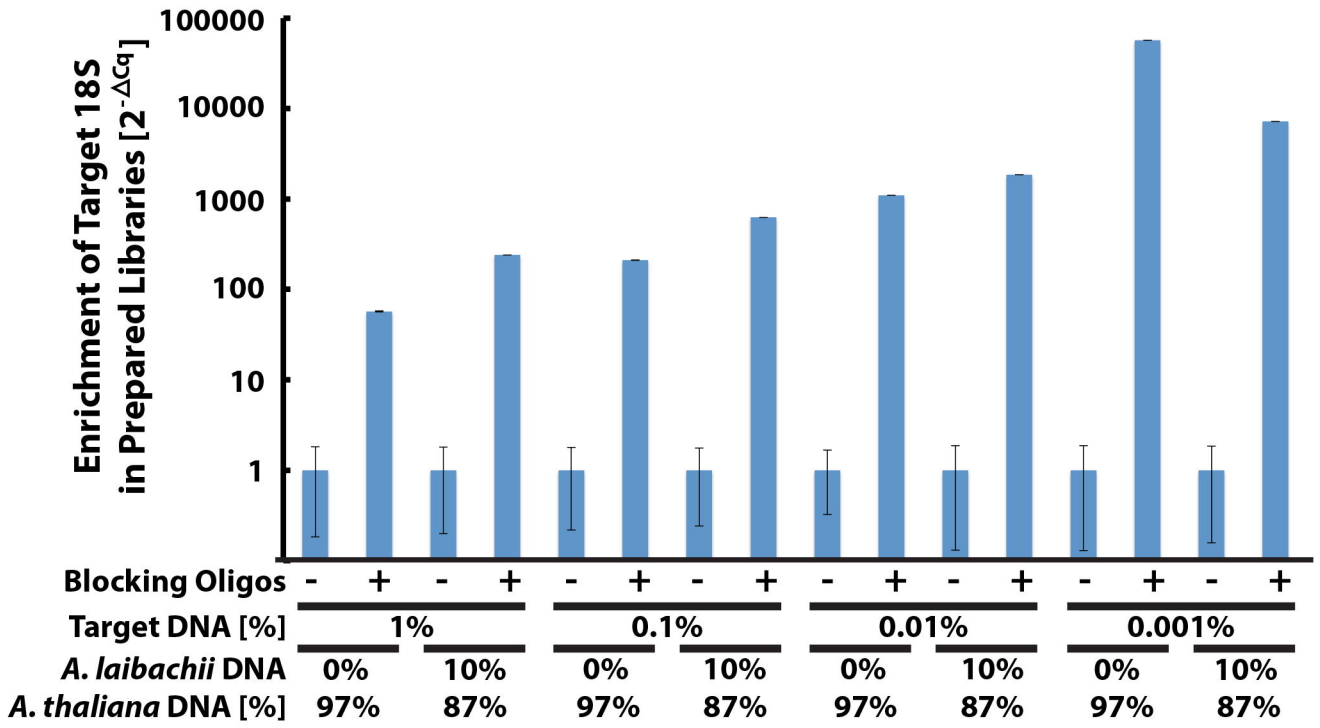- *Actinomycetales*
- *Bacteroidales*
- *Bacillales*
- *Lactobacillales*
- *Clostridiales*
- *Rhodobacteriales*
- *Neisseriales*
- *Campylobacterales*
- *Enterobacteriales*
- *Pseudomonadales*
- *Deinococcales*
- Other (Nontarget)

Bacteria (%) [V3/V4 rRNA Gene]

- *Saccharomycetales*
- *Sordariomycetes*
- *Usilaginales*
- *Mucorales*
- Other (Nontarget)

Fungi (%) [ITS1 Region]

- *Peronosporales*
- *Pythiales*
- *Saprolegiales*
- Other (Nontarget)

Oomycete (%) [ITS1 Region]

Nontarget DNA
Nontarget Blocking
PCR Steps
Community Ratios

Expected — 2 (10/25) — 2 (10/25) — 2 (25/10) — 1 (35)

Even — Staggered

**A.**



**B.**

**A.** 454 vs. Illumina OPU Observation Curves

OPUs
- Shared
- Unique Illumina
- Unique 454
- Shared/Total

X-axis: Subsampled Read Depth (Thousands)
Y-axis: Fraction/Number

**B.** OPU-based distances between compartments

Technology
- 454
- Illumina

Y-axis: Bray-Curtis

Rhiz/Soil: $p=1.12 \times 10^{-5}$
Root/Rhiz: $p=0.51$
Root/Soil: $p=0.51$

**C.** Mock Community Bacterial Alpha Diversity

Operational Unit
- Order
- Family
- Genus
- Species
- OTU (97% ID)
- OPU
- Actual

X-axis: Subsampled Read Depth (Thousands)
Y-axis: Observed Taxonomic Units

**D.** Mock Community Observed vs. Expected Units

Operational Unit
- OTU
- OPU

X-axis categories: Bacteria V3-V4, Bacteria V5-V7, Fungi ITS1, Fungi ITS2, Oomycete ITS1, Oomycete ITS2
Y-axis: Observed:Expected Taxonomic Units Ratio

**Recovered Diversity**

**Target loci specificity**

**A.**

**Eukaryotes**

Fungi — Plants
Green Algae — Metazoa
Cercozoa — Oomycetes
Red Algae — Amoebozoa

Oomycete ITS
Fungal ITS
Eukaryote 18S

**% of Tips**
1% — 0%
43% — 6%
50% — .3%

**B.**

**Prokaryotes**

Proteobacteria — Gemmatinomonadetes
Actinobacteria — Verrucomicrobia
Firmicutes — Acidobacteria
Bacteroidetes — Termi
Cyanobacteria — Unclassified

Bacteria 16S V3/V4
Bacteria 16S V5-V7

**% of Tips**
19%
39%
42%