

1 **Improved microbial community characterization of 16S rRNA via metagenome hybridization**
2 **capture enrichment**

3 Megan S. Beaudry^{1,#,*}, Jincheng Wang^{1,2,#,‡}, Troy J. Kieran¹, Jesse Thomas^{1,3,§}, Natalia J. Bayona-
4 Vásquez^{1,4,δ}, Bei Gao¹¶, Alison Devault⁵, Brian Brunelle⁵, Kun Lu^{1,⋈}, Jia-Sheng Wang^{1,2}, Olin E.
5 Rhodes, Jr.³, Travis C. Glenn^{1,2,4},

6 ¹Department of Environmental Health Science, University of Georgia, Athens, GA 30602, USA

7 ²Interdisciplinary Toxicology Program, University of Georgia, Athens, GA 30602, USA

8 ³Savannah River Ecology Laboratory, University of Georgia, Aiken, SC 29808, USA

9 ⁴Institute of Bioinformatics, University of Georgia, Athens, GA 30602, USA

10 ⁵Daicel Arbor Biosciences, 5840 Interface Dr., Suite 101, Ann Arbor, MI 48103, USA

11 #equal contributions (co-first authors)

12 ‡current address: Department of Biochemistry and Microbiology, Rutgers University, New Brunswick,
13 NJ 08854, USA

14 §current address: Center for Disease Control, Atlanta, GA 30329, USA

15 δcurrent address: Department of Biology, Oxford College of Emory University, 801 Emory Street,
16 Oxford, GA, 30054, USA

17 ¶current address: School of Marine Sciences, Nanjing University of Information Science and
18 Technology, Nanjing, 210044, China

19 ⋈current address: Department of Environmental Sciences and Engineering, University of North
20 Carolina, Chapel Hill, NC 27599, USA

21 * **Correspondence:**

22 Megan Beaudry

23 Megan.Beaudry@uga.edu

24 **Keywords: amplicon, microbial diversity, microbiome, mock communities, next generation**
25 **sequencing, shotgun libraries, target enrichment**

26 **Abstract**

27 Environmental microbial diversity is often investigated from a molecular perspective using 16S
28 ribosomal RNA (rRNA) gene amplicons and shotgun metagenomics. While amplicon methods are
29 fast, low-cost, and have curated reference databases, they can suffer from amplification bias and are
30 limited in genomic scope. In contrast, shotgun metagenomic methods sample more genomic regions
31 with fewer sequence acquisition biases. However, shotgun metagenomic sequencing is much more
32 expensive (even with moderate sequencing depth) and computationally challenging. Here, we

33 develop a set of 16S rRNA sequence capture baits that offer a potential middle ground with the
34 advantages from both approaches for investigating microbial communities. These baits cover the
35 diversity of all 16S rRNA sequences available in the Greengenes (v. 13.5) database, with no
36 sequence having < 80% sequence similarity to at least one bait for all segments of 16S. The use of
37 our baits provide comparable results to 16S amplicon libraries and shotgun metagenomic libraries
38 when assigning taxonomic units from 16S sequences within the metagenomic reads. We demonstrate
39 that 16S rRNA capture baits can be used on a range of microbial samples (i.e., mock communities
40 and rodent fecal samples) to increase the proportion of 16S rRNA sequences (average >400-fold) and
41 decrease analysis time to obtain consistent community assessments. Furthermore, our study reveals
42 that bioinformatic methods used to analyze sequencing data may have a greater influence on
43 estimates of community composition than library preparation method used, likely in part to the extent
44 and curation of the reference databases considered.

45 **1 Introduction**

46 The study of microbes is critically important, as they have many essential roles in ecosystem
47 function, disease pathology, host physiology, and possibly assessing infectious disease outbreaks
48 (Dueker et al., 2018; Gallardo-Escárate et al., 2020). As microbial communities can often be highly
49 diverse and complex, it can be challenging to identify rare taxa in complex environmental samples
50 (e.g., soil, freshwater, etc.) with traditional and modern techniques (i.e., culturing, 16S amplicons, or
51 metagenomic shotgun libraries). Advances in sequencing technologies have transformed traditional
52 microbiology. Microbial communities that were previously considered indiscernible or unstudied,
53 can now be investigated at greater depths than ever before from many different environmental
54 systems (Gilmour et al., 2010; Kustin et al., 2019).

55 For decades, the 16S small subunit ribosomal RNA (rRNA) gene has been the gold standard marker
56 for microbial molecular taxonomic research (Woese and Fox, 1977; Meola et al., 2015), as this
57 highly conserved gene contains nine rapidly evolving hypervariable regions that aid in species
58 identification (Yuan et al., 2015). Amplicon sequencing, targeting the 16S rRNA, is a cost-effective
59 and high-throughput method used to study aquatic, terrestrial, food- and host-associated microbial
60 communities (Logares et al., 2014; Polka et al., 2015; Jiang et al., 2016; Jouselin et al., 2016;
61 Jouglin et al., 2019; Suenami et al., 2019; Ziegler et al., 2019). However, studies relying on 16S
62 rRNA amplicon sequencing have limitations and biases. Relevant biases in 16S rRNA amplicon
63 sequencing are associated with DNA extraction, amplification via PCR, sequencing, and sequence

64 analysis (Kennedy et al., 2014; Knight et al., 2018). Specifically, PCR biases include primer bias
65 (Klindworth et al., 2013; Kelly et al., 2019) and varying GC content (Aird et al., 2011). Additional
66 limitations associated with amplicon sequencing include challenges in the taxonomic characterization
67 of microbial communities, as well as accuracy and availability of reference databases (Kennedy et al.,
68 2014; Poretsky et al., 2014; Ritari et al., 2015; Knight et al., 2018). Furthermore, the selection of the
69 hypervariable region used for the amplicon analysis (i.e., V1-V3; V3-V4; V4; etc) can lead to
70 differences in bacterial identification (Vetrovsky and Baldrian, 2013; Martinez-Porchas et al., 2016).

71 In more recent years, metagenomic shotgun sequencing has aimed to characterize taxonomic profiles
72 of unique clade-specific marker genes to provide a balanced view of community composition and
73 function (Neelakanta and Sultana, 2013; Knight et al., 2018). However, metagenomic sequencing has
74 its own limitations; genomic DNA may contain non-target DNA (e.g., human DNA), which can
75 affect downstream analysis (e.g., mis-assemblies of sequence contigs, spurious reads) thus leading to
76 inaccurate conclusions (Schmieder and Edwards, 2011; Gasc and Peyret, 2018). Also, metagenomic
77 libraries are more expensive, take longer to prepare, and are much more complex than amplicon
78 libraries, requiring more computational effort (Sekse et al., 2017). In particular, it is difficult to
79 identify low abundance genetic traits and rare taxa using metagenomic libraries, and extensive deep
80 sequencing is often required to do so (Lasa et al., 2019). In summary, shotgun sequencing is less
81 biased and yields data on many genomic regions, but the main tradeoffs are high costs of library
82 preparation, sequencing, and analysis.

83 Mock communities can be used to help establish ground truth in microbial diversity studies, in
84 particular when comparing different library preparation methods (Costea et al., 2017; Rausch et al.,
85 2019). Rausch et al., 2019 provided a comparison of 16S rRNA amplicon sequencing and
86 metagenomic sequencing, and revealed similar community makeup (i.e., abundance and taxa
87 diversity) of their shallow mock community regardless of library type. Conversely, other studies have
88 found key differences in abundance and taxa of mock communities attributed to wet-laboratory
89 methods (Costea et al., 2017; Rausch et al., 2019). However, some of these differences may be
90 attributed to varying bioinformatic tactics.

91 In terms of bioinformatic analyses, advantages and limitation of methods, reference databases, and
92 software have been vastly described for both 16S rRNA and metagenomic strategies (Truong et al.,
93 2015; Callahan et al., 2016a; Costea et al., 2017; Escobar-Zepeda et al., 2018; Rausch et al., 2019).
94 The variation among these can lead to a lack of sensitivity and specificity that may contribute to
95 wrong classifications and/or no classification at a specific taxonomic level, and erroneous abundance

96 assignments (Escobar-Zepeda et al., 2018). In particular, it can be challenging to analyze
97 environmental samples, as most reference databases are based on human commensals (Dueholm et
98 al., 2020).

99 Both strategies (i.e., 16S rRNA amplicon and metagenomic shotgun libraries) present their own
100 challenges and variations in analyses (Knight et al., 2018), but metagenomic shotgun libraries tend to
101 perform at a higher sensitivity and specificity than 16S rRNA amplicon data (Escobar-Zepeda et al.,
102 2018). For metagenomic data, programs like MetaPhlan2 may be used to classify and estimate the
103 relative abundance of microbial cells by mapping reads against marker sequences to classify the
104 sequences at the sub-species to higher taxonomic levels (i.e., marker-gene approach) (Segata et al.,
105 2012; Truong et al., 2015). Whereas 16S rRNA amplicon data is commonly analyzed by inferring
106 representative sequences using a variety of methods, some of which are influenced by fragment size
107 and 16S region (Edgar, 2013; Callahan et al., 2016a; Callahan et al., 2016b). Furthermore, some
108 methods used to assign operational taxonomic units may result in limited resolution at lower
109 taxonomic levels (e.g., genus and species levels), as even organisms that share 98.75% sequences
110 may be different species (Mysara et al., 2017). Reference databases for 16S rRNA are much more
111 extensive than those for metagenomic analyses, which is key for superior analysis, particularly in
112 environmental samples (Escobar-Zepeda et al., 2018). However, variation in taxonomic classification
113 and abundance has also been associated with the use of different reference databases (Jovel et al.,
114 2016; Rausch et al., 2019).

115 Hybridization capture (also known as sequence capture, target capture, or targeted sequence capture)
116 is an enrichment technique that uses a set of biotinylated DNA or RNA baits that are complementary
117 to DNA sequences of interest to increase the proportion of DNA fragments of interest within DNA
118 libraries, subsequently characterizing the DNA by massively parallel sequencing (Lasa et al., 2019).
119 Hybridization capture assays have been designed previously for the 16S rRNA gene, using 15-1,402
120 baits (Gasc and Peyret, 2018; Barrett et al., 2020). Additional hybridization capture bait sets have
121 been designed for a variety of microbial projects, such as sets of defined pathogens or particular
122 genes, including virulence genes for *Vibrio* spp. that infect oysters (Lasa et al., 2019), bifidobacterial
123 in the gut of mammals (Lugli et al., 2019), and antibiotic resistance genes (Guiton et al., 2019).
124 Importantly, unlike other culture independent techniques, hybridization capture provides greater
125 phylogenetic resolution and increased sensitivity, while requiring fewer sequencing reads (Lasa et al.,
126 2019; Barrett et al., 2020). More specifically, 16S rRNA capture baits provide a cost-effective way to
127 identify bacteria in diverse environmental samples and identify rare taxa.

128 Here, we present a hybridization capture method (i.e., 16S-cap) to enrich metagenomic shotgun
129 libraries for DNA sequences of 16S rRNA genes. Our protocol improves on the existing methods by
130 including many more baits that better cover known sequence variation in 16S databases, taking
131 advantage of the extensive reference databases and ease of analyses of 16S rRNA sequences for
132 taxonomic classification and decreasing bias introduced from primer affinity, while reducing
133 sequencing costs per sample compared to unenriched metagenomic libraries. For microbes, targeted
134 sequence capture techniques for 16S rRNA have shown more accurate representation of microbial
135 communities compared to traditional methods (i.e., 16S rRNA amplicons, shotgun libraries) (Gasc
136 and Peyret, 2018). We provide a comparison of traditional methods for assessing composition of
137 microbial communities (i.e., 16S rRNA amplicons and metagenomic shotgun libraries) with our 16S-
138 cap method to characterize *in silico* mock, *in vitro* mock, and real microbial communities from
139 genomic data.

140 **2. Materials & Methods**

141 **2.1 Samples and DNA Extraction**

142 We used two commercial standard genomic DNA mock community collections to characterize
143 simple communities (HM-276D, BEI Resources, Manassas, VA; D6306, Zymo Research, Irvine,
144 CA). For complex communities, we used a subset of fecal samples from previous studies that
145 examined the impacts of environmental xenobiotic agents on the gut microbial communities of
146 rodent models (Gao et al., 2017; Wang et al., 2018). The first study examined carbamate insecticide
147 in male C57BL/6 mice (i.e., *Mus musculus*) (Gao et al., 2017), and the second examined green tea
148 polyphenols in female Sprague-Dawley rats (i.e., *Rattus norvegicus*) (Wang et al., 2018). DNA was
149 extracted using Qiagen Fast DNA Stool Mini Kit (QIAGEN, Valencia, CA, USA) or PowerSoil
150 DNA Isolation Kit (Mo Bio Laboratories, Carlsbad, CA, USA). Details on experimental design and
151 extractions are previously described (Gao et al., 2017; Wang et al., 2018).

152 **2.2 16S rRNA Amplicon Metabarcoding**

153 The primer pairs targeting the V3 and V4 16S regions (S-D- Bact-0341-b-S-17 and S-D-Bact-0785-
154 a-A-21) (Klindworth et al., 2013) were used for amplification of the 16S rRNA gene in rat fecal
155 samples and mock communities; and the primer pair targeting the V4 region (515-F and 806-R)
156 (Caporaso et al., 2012) was used on the mouse fecal samples. We created indexed fusion primers
157 with TruSeq compatible sequencing oligos as previously described using the *Adapterama I* and

158 *Adapterama II* systems (Glenn et al., 2019a; Glenn et al., 2019b) to generate amplicon libraries using
159 two rounds of PCR (Method 5 of Table 3 from Glenn et al. 2019b). For the first PCR, we prepared
160 individual 25 μ L PCR reactions for each sample using KAPA HiFi reagents (KAPA Biosystems,
161 Wilmington, MA, USA). Each PCR reaction mix included 5 μ L 5x KAPA HiFi buffer, 0.75 μ L 10
162 mM dNTPs, 0.5 μ L KAPA HiFi HotStart, 1.5 μ L 5 μ M forward indexed-fusion primer, 1.5 μ L 5 μ M
163 reverse indexed-fusion primer, and 1 μ L of 20 ng/ μ L DNA. PCR conditions were as follows: initial
164 denaturation at 95°C for 3 min; 15-18 cycles of 95°C for 20 sec, 60°C for 30 sec, and 72°C for 30
165 sec; final extension at 72°C for 5 min.

166 In preparation for the second PCR, we normalized individually indexed PCR products with a
167 SequalPrep Normalization Plate Kit (Invitrogen, Carlsbad, CA, USA) according to manufacturer's
168 protocols or by pooling them together based on agarose gel band brightness. These pools served as
169 the template for a second limited cycle PCR. Each 25 μ L PCR reaction mix included: 5 μ L 5x KAPA
170 HiFi buffer, 0.75 μ L 10 mM dNTPs, 0.5 μ L KAPA HiFi HotStart, 2.5 μ L of 5 μ M forward iTru5
171 primer, 2.5 μ L of 5 μ M reverse iTru7 primer, and 5 μ L of product from the first PCR. The following
172 were used as PCR conditions: initial denaturation at 95°C for 2 min; 10 cycles of 95°C for 20 sec,
173 60°C for 15 sec, and 72°C for 30 sec; final extension at 72°C for 5 min. These PCR products were
174 purified with Sera-Mag magnetic beads (Thermo-Scientific, Waltham, MA, USA). We quantified the
175 final products with a Qubit 2.0 Fluorometer (Thermo-Scientific, Waltham, MA, USA) and pooled
176 them in equal molar ratios for sequencing. Samples were sequenced using an Illumina MiSeq v2 600
177 cycle kit (Illumina, San Diego, CA, USA) at the Georgia Genomics and Bioinformatics Core
178 (Athens, GA, USA).

179 **2.3 Metagenomic Libraries**

180 Extracted DNA was sheared on a Bioruptor UCD-300 (Diagenode, Denville, NJ, USA) to an average
181 size of about 500 bp. We input ~100 ng of fragmented DNA into each reaction of a KAPA
182 HyperPrep Kit (KAPA Biosystems, Wilmington, MA, USA) following manufacturer's protocol at
183 half volume reaction size with 14 PCR cycles using iTru adaptors and indexed primers (Glenn et al.,
184 2019b). Samples were sequenced on an Illumina HiSeq 3000 with PE150 reads (Oklahoma Medical
185 Research Foundation, Oklahoma City, OK, USA).

186 **2.4 16S rRNA Bait Design**

187 We used Prokka v1.11 with default settings, to annotate and extract all 16S rRNA sequences in
188 GreenGenes v13.5 to ensure that only 16S rRNA regions were represented in the final bait set
189 (Seemann, 2014). Stretches of up to 25 Ns were replaced with T bases to facilitate probe design
190 across short unknown regions. We then used USEARCH v8.1 (Edgar, 2010) to sort by length (large
191 to short) and cluster (query coverage 90%, identity 90%) sequences, retaining one centroid from each
192 cluster. We then designed 120mer baits with flexible ~50% overlap. These baits were then clustered
193 using USEARCH (query coverage 75%, identity 78%), and one centroid per cluster retained.

194 **2.5 16S rRNA Hybridization Capture Enrichments**

195 Metagenomic libraries were combined into 500 ng pools of eight samples for rodents or two samples
196 for mock communities. Target enrichments of each pool were performed using myBaits kit (Arbor
197 Biosciences CAT # 308616, Ann Arbor, MI, USA) using the designed 16S rRNA Capture Baits
198 following manufacturer's protocol (v3.01) with a 24-hour 65°C hybridization. Following
199 hybridization, we used Dynabeads M-280 Streptavidin magnetic beads (Life Technologies, Carlsbad,
200 CA, US) for capturing and washing each biotinylated bait library. We then performed a post-
201 enrichment amplification using Illumina P5/P7 primers (Illumina, San Diego, CA, USA) and KAPA
202 HiFi HotStart reagents (KAPA Biosystems, Wilmington, MA, USA) using 98°C for 45 seconds,
203 followed by 16-22 cycles of 98°C for 20 seconds, 60°C for 30 seconds, and 72°C for 60 seconds,
204 ending with a final extension of 72°C for five minutes. PCR products were cleaned 1:1 with Sera-
205 Mag beads (Glenn et al., 2019a), quantified on Qubit and pooled in equimolar ratios for sequencing
206 paired-end 150 bp and 300 bp reads on Illumina HiSeq 3000 (Oklahoma Medical Research
207 Foundation, Oklahoma City, OK, USA) and MiSeq (Georgia Genomics Bioinformatics Core, Athens,
208 GA, USA), respectively.

209 **2.6 Simulating 16S rRNA Target Enrichment Data**

210 Three metagenomes (i.e., Lindgreen synthetic metagenome (Lindgreen et al., 2016); Zymo Mock
211 Community DS6306 genomes; and BEI Mock Community HM-276D) were used to simulate 16S
212 rRNA capture data. In summary, a fasta file containing our 120mer bait set was mapped to each
213 metagenome fasta file (Supplementary Data 1-3) using Burrows-Wheeler aligner (bwa) v.0.7.17 (Li
214 and Durbin, 2009). Samtools v1.9 (Li et al., 2009) was used to convert the obtained sam file into a
215 bam file. This mapping process is meant to simulate what would be an error- and bias-free
216 hybridization process. Following this, we obtained the mapping coordinates of the baits on the
217 reference metagenomes and extracted the sequences + 200 bp to the upstream and downstream of the

218 first position, if possible. Here, we sought to simulate a hybridization of the bait to the core of an
219 ~500 bp fragment while obtaining the flanking regions typically captured from use of biotinylated
220 baits.
221 The software ART 2016.06.05 (Huang et al., 2012) was then used to simulate > 200,000 paired-end
222 150 bp fastq reads from these extended reference sequences from each metagenome. These fastq files
223 were mapped to Greengenes 97% similarity database v.13.8 using BBmap v. 38.50 (Bushnell, 2014).
224 For each metagenome, we recorded the number of paired reads mapped to Greengenes, number of
225 forward reads, number of reverse reads and percentage average total mapped, and compared these
226 results with those from real samples also mapped to the Greengenes database (see below) (Altschul et
227 al., 1990).

228 **2.7 Data Processing and Analysis**

229 After obtaining demultiplexed Illumina pair-end raw sequences, we used library specific pipelines to
230 process the data (Figure 1). For 16S rRNA amplicon libraries, primers were removed using cutadapt
231 v1.15 (Martin, 2011). Following this, DADA2 (v1.8) was used for quality trimming and filtering, de-
232 replication and sequence-variant inference, merging paired-end reads, construction of feature tables,
233 low relative abundance filtering of 0.5%, removal of chimeras, and taxonomy assignment (Callahan
234 et al., 2016a). The taxonomy assignment was based on 97% clustered OTU based on Greengenes
235 v13.8 database in the DADA2 pipeline.

236 [Insert Figure 1]

237 For 16S-cap libraries, the resulting quality filtered reads were mapped to the 97% clustered OTU
238 based on Greengenes v13.8 database using BBmap v37.78 (Bushnell, 2014). The resulting mapping
239 information was filtered, and a hit was recorded if both ends of paired read hit the same reference, or
240 only one end of the paired read hit a reference. A low relative abundance filter of 0.5% was applied.

241 Also, we assessed the presence of non-target reads in the quality-filtered dataset by 1) running
242 MetaPhlan2 v2.7.8 (Segata et al., 2012; Truong et al., 2015), and 2) mapping to the rat and mice
243 genomes using Burrows-Wheeler aligner (bwa) v.0.7.17 (Li and Durbin, 2009).

244 For unenriched metagenomic libraries, Trimmomatic v0.36 (Bolger et al., 2014) was used for quality
245 trimming using a sliding window of three nucleotides with an average Q > 20, and minimum length
246 of 75 nucleotides. Reads that passed initial quality filtering (including both paired reads and orphan
247 reads) were fed to MetaPhlan2 v2.7.14 for taxonomy assignment (Segata et al., 2012; Truong et al.,
248 2015). A low relative abundance filter of 0.5% was applied. To further compare to the results from

249 16S-cap analysis, we performed the same 16S mapping steps to the GreenGenes database as
250 described for 16S-cap libraries for the unenriched libraries.

251 Data was analyzed using R statistical software (R Development Core Team, 2010). Duncan's
252 multiple range test was used to compare abundance estimates between library types. Additionally, for
253 all samples, abundance estimates were used to construct a Bray-Curtis dissimilarity matrix, which
254 was then used to generate a principle coordinate analysis (PCoA).

255 **3 Results**

256 **3.1 16S rRNA Capture Bait Design**

257 The 1,262,986 sequences comprising Greengenes v13.5 were annotated and 1,261,075 16S rRNA
258 sequences were retained. A total of 117 sequences containing consecutive runs of 25 or more
259 ambiguous bases (Ns) were removed. A total of 18,649 centroidal sequences were obtained from
260 USEARCH clustering. From these sequences, 413,480 120mer baits were designed. These baits were
261 then clustered using USEARCH, retaining one centroid per cluster, for a total of 37,745 baits.

262 **3.2 Sequencing Summary Statistics**

263 A summary of average sequence statistics for each sample and library preparation type is given in
264 Table 1. For the 16S rRNA amplicon data, the number of total raw read pairs per sample ranged from
265 49,828 for the Zymo mock community to 136,184 for the BEI mock community, with rodent fecal
266 samples having intermediate depth. More reads (~77%) remained from the rodent fecal samples after
267 the denoising steps through the rigorous DADA2 pipelines versus the mock communities. Low
268 percentages of high quality reads remained following filtering for both the BEI and Zymo mock
269 communities (38.7% and 48.8% respectively). For the BEI mock community, initial index matching
270 in R2 reads caused ~30% loss of data (versus less than 5% typically observed in other samples) and
271 DADA2 quality trimming lost another ~30% of data. For the Zymo mock community, the loss of
272 data was mainly due to chimeric filtering (~30% of data loss).

273 For the unenriched libraries, the highest number of total raw read pairs ranged from 4,985,957 in the
274 Zymo mock community to 28,219,552 in the insecticide-treated mouse feces. The percentage of
275 reads retained after filtering was greater than 65% for all unenriched libraries. The average
276 percentage of reads mapped to GreenGenes ranged from 0.1% to 0.2% in the BEI and Zymo mock
277 communities.

278 For 16S-cap libraries, the PE150 reads had higher numbers of reads on average per sample type than
279 PE300 reads. The highest number of raw reads (i.e., 11,474,476) was obtained for the insecticide-
280 treated mouse feces with PE150 reads. The percentage of reads after filtering were greater than 70%
281 for all 16S-cap libraries. The average percentage of mapped reads was greater than 50% for all 16S-
282 cap libraries, with the highest percentage of mapping in the 16S-cap BEI mock community
283 sequenced with PE300 at 75.7%. On average among all sample types, the proportion of on target
284 reads was increased 435-fold when compared to unenriched libraries (range 283 – 499 fold increase,
285 Supplemental Table 2).

286 [Insert Table 1]

287 **3.3 16S rRNA Target Enrichment Simulated Reads**

288 Summary information for simulated reads is given in Table 2. We observed a higher percentage of
289 total mapped reads in our simulated mock communities than for the real data from those communities
290 (Table 2). For example, the real data from the Zymo mock community had an average total mapping
291 of 78.15% to GreenGenes, compared to 91.43% from the simulated data. Similarly, the BEI mock
292 community had an average total mapping of 78.62% for the real data, compared to 92.37% for the
293 simulated data.

294 [Insert Table 2]

295 **3.4 Validation on Mock Community Samples**

296 We initially prepared amplicon libraries, unenriched metagenomic libraries, and enriched our
297 metagenomic libraries using the target enrichment bait set (i.e., 16S-cap) we developed using two
298 mock communities (Table 1). At the phylum level both samples appear to provide accurate
299 identification of the microbes with good estimates of abundance, regardless of library type or data
300 analysis method used (Figure 2). The 16S-cap samples and metagenomic samples generate one
301 detection of a false positive phyla in the mock community samples (Figure 2). Additionally, in both
302 the unenriched and 16S-cap libraries analyzed with a 16S mapping approach, Cyanobacteria was
303 found in low abundance even though it was not expected to be present in the mock community.
304 However, when analyzing the unenriched library using marker gene approach, Cyanobacteria was not
305 found and instead Ascomycota was identified.

306 [Insert Figure 2]

307 At the genus level, 16S-cap and unenriched libraries reflect more accurate microbial community
308 composition and abundance for most taxa (Figure 3). The 16S-cap and unenriched libraries with 16S

309 mapping missed three genera: *Escherichia*, *Listeria*, and *Bacillus* for both mock community samples.
310 However, after identifying presumably false positive genera with above 1% abundance in the
311 samples analyzed with 16S mapping software, three families with no genus identification,
312 *Enterobacteriaceae*, *Listeriaceae*, *Bacillaceae*, were found, suggesting these are likely the missing
313 genera. In comparison, 16S rRNA amplicon-based analysis identified nearly all genera in mock
314 samples, however, its estimates of abundance for *Actinomyces*, *Propionibacterium*, *Pseudomonas*,
315 and *Rhodobacter* all greatly deviate from the nominal compositions. The unenriched metagenomic
316 libraries analyzed with a marker-gene approach were able to identify all 18 genera in the mock
317 communities, however its estimate of *Bacillus* abundance in both mock communities deviate from
318 the nominal composition (Figure 3).

319 [Insert Figure 3]

320 [Insert Figure 4]

321 In the BEI mock community libraries, relative abundance estimates in the 16S-cap libraries were
322 more accurate than the amplicon and unenriched libraries as measured by fold change being very
323 close to 1 (Figure 4). In the amplicon library, several genera (i.e., *Pseudomonas*, *Actinomyces*,
324 *Propionibacterium*, and *Rhodobacter*) are beyond the 2-fold change of their nominal compositions.
325 In particular one genus, *Rhodobacter*, proved to be challenging for all three library preparation
326 methods for accurate estimation of relative abundance. Duncan's multiple range test revealed that
327 there were significant differences ($p\text{-value} \leq 0.05$) between the BEI mock community amplicon and
328 16S-cap libraries, whereas the unenriched libraries were not found to be significantly different than
329 the amplicon or 16S-cap libraries. For the Zymo mock community libraries, relative abundance
330 estimates in the 16S-cap libraries are more accurate than relative abundance estimates for the
331 amplicon library. However, Duncan's multiple range test did not detect a significant difference
332 between the three library types (i.e., amplicon, unenriched, and enriched) (Figure 4).

333 **3.5 Validation on Fecal Samples**

334 Principle coordinate analysis was performed on mock community samples and additional samples
335 from laboratory mice and rats to further validate the 16S-cap method. When Bray-Curtis was used to
336 construct the dissimilarity matrix, which considers abundance estimates, we found that regardless of
337 analyses at the level of family (Figure 5A, left) or genus (Figure 5B, right) similar themes emerged.
338 We observed that the mock community samples were similar to each other regardless of library type.
339 Conversely, in the mouse and rat samples, we found that the unenriched libraries analyzed with a

340 marker-gene approach grouped together separately from amplicon, unenriched, and 16S-cap libraries,
341 all of which were analyzed with the 16S mapping approach.

342 [Insert Figure 5]

343 A comparison of Bray-Curtis distance was performed for rodent fecal samples at the level of family
344 and genus (Figure 6). This analysis revealed similar trends regardless of sample type or taxonomic
345 rank. The 16S-cap and unenriched libraries analyzed with 16S mapping approach showed to be the
346 most similar to each other, with a dissimilarity rate below 0.25. Bray-Curtis dissimilarity was higher
347 when comparing the amplicon libraries to both 16S-cap and unenriched libraries. When comparing
348 the unenriched libraries analyzed with two different analysis strategies (i.e., mapping reads to
349 GreenGenes vs gene-marker approach), we observed the highest degree of dissimilarity at both the
350 family and genus levels with dissimilarity rates at approximately 0.75. Post-hoc analysis revealed
351 that there were significant differences when comparing the unenriched and 16S-cap libraries to all
352 other library types, regardless of sample type or taxonomic rank (Figure 6).

353 [Insert Figure 6]

354 **4 Discussion**

355 Given the limitations of 16S rRNA amplicon and shotgun metagenomic libraries outlined in the
356 introduction, we sought to provide an alternative method to identify microbial community
357 composition by creating a 16S rRNA hybridization capture assay (i.e., 16S-cap). Our study revealed
358 two important things: 1) our 16S-cap method is an efficient way to obtain sequences from the
359 complete 16S rRNA gene to accurately reflect microbial community composition and abundance and
360 2) bioinformatic analysis methods greatly influence community composition in environmental
361 samples, regardless of library type. In our study we observed that sequences from 16S-cap were not
362 significantly different than sequences from unenriched shotgun libraries when analyzed using similar
363 bioinformatic methods and databases. However, we did find that the 16S-cap assay requires far fewer
364 reads, thus allowing enriched libraries to be characterized on benchtop sequencers, including
365 Illumina MiSeq instruments, at reasonable cost while overcoming the previously mentioned
366 limitations with direct 16S rRNA approaches and metagenomic approaches. These limitations
367 include selection and drift bias in PCR during amplicon library preparation and the potential for non-
368 target DNA (e.g., human DNA) in metagenomic libraries, which can lead to errors in downstream
369 analyses.

370 Enrichment for genes of interest is an important technique in characterizing complex environmental
371 samples. Previous studies have found other capture enrichment methods to increase the proportion of
372 on target reads from ~0.1% in unenriched shotgun libraries to ~60% in enriched libraries (Gasc and
373 Peyret, 2018). Similarly, we found 0.1 – 0.2% of unenriched libraries to map to the 16S rRNA,
374 whereas 58-76% of the enriched reads mapped to the 16S rRNA (Table 1). On average we achieved a
375 435-fold increase in reads mapped to the 16S rRNA in our 16S-cap libraries compared to the
376 unenriched libraries (Supplementary Table 2). *In silico* simulations of 16S-cap revealed that under
377 ideal conditions, 88-92% mapping to the 16S rRNA from mock communities could be achieved.
378 Therefore, our 16S-cap enrichment process helps to achieve a very high percentage of on-target
379 reads, but not quite as high as theoretically possible.

380 Our 16S-cap method identified several species that were not expected in the theoretical targets of the
381 mock communities, which may be attributed to several factors. First, the lack of genus identification
382 may be due to the mapping methods or clustering level used in data analysis rather than the library
383 preparation method. Both the 16S-cap and unenriched libraries analyzed with a 16S mapping method
384 failed to identify three genera *Escherichia*, *Listeria*, and some *Bacillus* in the mock communities.
385 However, there are three families, *Enterobacteriaceae*, *Listeriaceae*, and *Bacillaceae*, in the false
386 positive genera with >1% abundance that align with our missing genera. Thus, it appears that reads
387 for these three genera appear to be present, but are not being assigned appropriately at the genus
388 level. By assigning these unidentified genera as *Escherichia*, *Listeria*, and *Bacillus* respectively, the
389 16S-cap library is highly accurate in terms of taxonomic classification and abundance. Taxonomic
390 misassignment is a known problem with 16S mapping bioinformatic methods, and new software is
391 in development (Schloss and Westcott, 2011; Pollock et al., 2018; Zinger et al., 2019; Djemiel et al.,
392 2020). Additional work on the mapping and assignment processes used here, as well as comparisons
393 of newly developed and commonly used bioinformatic software is beyond the scope of this paper, but
394 warranted in future work. Other taxa identified that were not expected in the mock communities may
395 be due to reagent contamination or index hopping during sequencing. Several studies show that
396 contaminating DNA is common in laboratory reagents and DNA extraction kits (Salter et al., 2014;
397 Weiss et al., 2014; Kim et al., 2017; Eisenhofer et al., 2019; Zinter et al., 2019). Furthermore, studies
398 recommend sequencing negative controls consisting of ‘blank’ extractions and library preparations to
399 identify contamination by bacterial species (Salter et al., 2014; Knight et al., 2018). Conversely, false
400 positives of extremely low abundance (i.e., 0.1% or less) may be due to misassigned data that can

401 occur during Illumina sequencing. This phenomena is often referred to as index hopping (van der
402 Valk et al., 2019).

403 We compared theoretical target values of the BEI resources and Zymo mock communities to all three
404 library types (i.e., amplicon, unenriched, and 16S-cap) (Figures 3, 4). We find that the 16S-cap
405 libraries are representative of the target abundance values of the mock communities (Figure 3). Post-
406 hoc analysis revealed that the 16S rRNA amplicon library and 16S-cap library made from the BEI
407 mock community were significantly different from each other (p -value ≤ 0.05) based on relative
408 abundance. A PCoA revealed that in the mouse and rat samples the unenriched libraries analyzed
409 with a marker-gene approach grouped together separately from 16S rRNA amplicon libraries and
410 16S-cap and unenriched libraries analyzed with taxonomic binning approach (Figure 5). Thus,
411 enrichment and amplicon sequencing result in similar library composition, as do 16S-cap and
412 unenriched libraries analyzed with a 16S taxonomic binning approach. This indicates that our 16S-
413 cap method may be less biased than 16S amplification, but that analysis methods or the reference
414 database may greatly influence community composition results. Walsh et al. (2018) analyzed
415 different species classifiers using marker gene approaches and taxonomic binning, and found that the
416 results of the marker gene approach (i.e., MetaPhlAn2) were different from taxonomic binning
417 methods. Taxonomic binning methods are influenced by the size of the reference genome, whereas
418 marker gene approaches are not (Droge and McHardy, 2012; Balvociute and Huson, 2017; Walsh et
419 al., 2018). The use of hybridization capture baits may help alleviate some of these issues.

420 Other groups have designed a more limited bait set to hybridize all known 16S rRNA gene sequences
421 by focusing on highly conserved regions and incorporating ambiguities(Gasc and Peyret, 2018).
422 When validating their bait set on a mock community, they found that they detected 24 of 26 genera
423 tested, and that two less abundant species (i.e., *Methanobrevibacter smithii* and *Methanococcus*
424 *aelocius* at 0.00006%) were missed. In addition, Cariou et al., (2018) tested hybridization capture
425 probes designed by Gasc & Peyret (2018) on a previously characterized pea aphid and found their
426 enriched libraries to be representative of the bacterial population. There are some key differences
427 between the design of our baits set and Gasc & Peyret (2018). Foremost, is the number of baits
428 included in the bait set. Our bait set included 37,745 120mer baits, whereas Gasc & Peyret bait set
429 include 15 baits that are 28- to 50-mer. Using more baits with more sequence variation among the
430 baits helps to capture a greater range of diverse targets and thus generates more accurate abundance
431 estimates of the full range of community members. Having a more extended bait set, as ours, may
432 allow to overcome some of the previous challenges, demonstrated by the ability to detect all genera

433 in our mock communities. These aspects are critical when studying environmental samples and
434 searching for rare taxa. In addition, the use of longer hybridization times or “double capture” (i.e.,
435 when captured product is captured again) can improve the percentage of on target reads and help
436 capture rare sequences. Future work to identify the optimal bait set(s) for various microbial
437 communities and research objectives should include a direct comparison of the Gasc & Peyret (2018)
438 bait set verses our bait set.

439 Preparing 16S-cap libraries can most readily be accomplished by using an existing enrichment kit,
440 which ranges in cost from \$1,500 - \$5,200 depending on the number of reactions purchased. To
441 reduce reagent costs and hands-on time, we have successfully pooled multiple samples (see section
442 2.5), which is commonly done (Glenn and Faircloth, 2016). For example, pooling samples in groups
443 of eight reduces capture costs from \$93.75 per sample to \$11.72 per sample (Supplementary Table
444 3). Larger numbers of samples can be pooled to further reduce costs, but there are tradeoffs (see
445 Glenn & Faircloth, 2016). Our baitset is commercially available from Arbor Biosciences in ready-to-
446 use kit format, and the bait sequences are freely available to the scientific community
447 (Supplementary Data 4). Thus, our baits can be modified and/or synthesized by any strategy any
448 researcher desires.

449 Sequencing 16S-cap libraries require less extensive sequencing than unenriched shotgun
450 metagenomic libraries, which reduces costs (Supplementary Tables 4, 5). For example, a 100-fold
451 16S-cap enrichment sequenced on an Illumina MiSeq Nano PE150 provides a cost-savings of
452 approximately \$315 compared to an unenriched metagenomic shotgun library requiring 1 million
453 reads (Supplementary Table 4). Indeed, 16S-cap makes it economically and logistically reasonable to
454 routinely screen for 16S segments from enriched shotgun metagenomic libraries on Illumina MiSeqs.
455 16S-cap decreases costs when using a production scale Illumina sequencer (e.g., Illumina NovaSeq)
456 to less than \$0.10 per sample when achieving a 100-fold enrichment (Supplementary Table 5).
457 However, because production scale sequencers produce 400 – 2,500 million read pairs, to achieve
458 low cost for samples needing relatively few reads, each run requires huge numbers of samples or a
459 mixture of some samples needing large numbers of reads (i.e., a mixture of projects; see Glenn et al.
460 2019a). Due to the limited savings possible on production sequencing costs (Supplementary Table 4),
461 the savings in data transfer, storage, and compute time may be more significant than savings in
462 sequencing costs.

463 In summary, our data demonstrates that the 16S-cap assay and unenriched shotgun metagenomic
464 libraries produce very similar community profiles. Importantly, our 16S-cap library is produced from

465 a metagenomic library, which eliminates primer (though not all PCR) biases. Additionally, our 16S-
466 cap assay provides a deeper community profile (i.e., more 16S reads that can be queried to a
467 database) with far fewer reads than the unenriched shotgun metagenomic libraries. In environmental
468 samples, we routinely achieved > 400-fold enrichment. Thus, expensive deep sequencing is
469 unnecessary for 16S-cap libraries because a few thousand reads provide the same number of 16S
470 rRNA sequences as millions of shotgun reads. By trading modest additional library preparation costs
471 for reduced sequencing costs (Supplementary Tables 3-5), 16S cap is economical and opens up the
472 possibility of adding deep taxonomic sampling to studies that are capturing other genes of interests
473 (e.g., antibiotic resistance genes (Guiton et al., 2019; Oladeinde et al., 2019; Thomas et al., 2020). In
474 comparison to amplicon libraries, the 16S-cap assay will be more expensive, however, it provides
475 superior microbial community resolution, increased accuracy of relative abundance, and greater
476 flexibility in terms of sequencer and kit choice. We believe that our bait set is a valuable tool to
477 efficiently and accurately identify microbial community composition and would be well-suited to be
478 used in combination with other bait sets targeting different genes of interest (e.g., antimicrobial
479 resistance baits).

480 *Figures*

481 **Figure 1.** Overview of data analysis methods on the three library types (i.e., 16s amplicon, 16s
482 hybridization bait capture, and metagenomic libraries).

483 **Figure 2.** Relative abundance of bacterial phyla in mock community controls sequenced and
484 analyzed using different methods. Phyla that are not among the nominal composition of the
485 respective mock communities are plotted as black dots next to z_Others. The black dot in the
486 enriched and unenriched library analyzed with 16S mapping software the assigned phyla was
487 Cyanobacteria. In unenriched libraries analyzed with a marker gene approach, the assigned phyla was
488 Ascomycota. Colored vertical bar in each panel represents the nominal abundance of respective
489 phylum. X-axis is plotted in log-scale to show the low abundance phylum. Row panel strips labels
490 identify the mock communities; column panel strips labels identify library type (i.e., amplicon,
491 enriched 16S-cap, unriched metagenomic library) and analyzing strategy (i.e., denoising,
492 16Smapping, and marker gene).

493 **Figure 3.** Relative abundance of bacterial genera in mock community controls sequenced and
494 analyzed using different methods. Panel A is the BEI mock community. Panel B is the Zymo mock
495 community. Genera that are not among the nominal composition of the respective mock communities

496 were plotted as black dots under z_Others. Colored vertical bar in each panel represents the nominal
497 abundance of respective genus. X-axis plotted in log-scale to show the low abundance genus. Row
498 panel strips labels identify the mock communities; column panel strips labels identify library type
499 (i.e., amplicon, enriched 16S-cap, unriched metagenomic library) and analyzing strategy (i.e.,
500 denoising, 16Smapping, and marker gene).

501 **Figure 4.** Fold change (i.e., upper or under) comparing the relative abundances of respective genera
502 in each library to its nominal abundance. Duncan's multiple range test was performed to compare
503 each library type for each mock community. Letters indicate whether significant differences were
504 detected.

505 **Figure 5.** PCoA plots were constructed using Bray-Curtis dissimilarity matrix at a family level
506 (panel A) and genus level (panel B). Each project is represented by a colored dot (i.e., orange = BEI
507 mock community, green = mouse samples, blue = rat samples, and purple = Zymo mock
508 community). Each library type, sequencing read length and data analysis method is represented by a
509 different shape (i.e., circle = amplicon library, square = 16S-cap enriched PE150 reads, diamond =
510 unenriched PE150 analyzed with 16S mapping and triangle = unenriched PE150 analyzed with
511 metagenome mapping). Numbers represent sample number.

512 **Figure 6.** A comparison of the Bray-Curtis distance metric was performed for each library type at a
513 genus level using box plots. Bray-Curtis distance is indicated on the y-axis. Library type is indicated
514 on the x-axis. Duncan's multiple range test was performed to compare each library type for each
515 mock community. Letters indicate whether significant differences were detected.

516 ***Permission to reuse and Copyright***

517 *Figures, tables, and images will be published under a Creative Commons CC-BY licence and*
518 *permission must be obtained for use of copyrighted material from other sources (including re-*
519 *published/adapted/modified/partial figures and images from the internet). It is the responsibility of*
520 *the authors to acquire the licenses, to follow any citation instructions requested by third-party rights*
521 *holders, and cover any supplementary charges.*

523 **Table 1.** A brief overview of the average summary statistics (i.e., number of samples, total raw read-pairs, average filtered/bar, average
524 mapped/filtered) for each sample type of each library type (i.e., 16S amplicon libraries, 16S-cap enriched, and unenriched).

Library Type	Read Length	Sample Type	N Samples	Total Raw Read-Pairs	Total Filtered Reads	Average Filtered/Raw (Mean±SD)	Average Mapped/Filtered (Mean±SD)
Amplicon-16S/V3V4	PE300	Rat feces	5	318,561	247,781	(77.3±6.2)%	NA
Amplicon-16S/V3V4	PE300	BEI Mock	1	136,184	52,734	38.7%	NA
Amplicon-16S/V3V4	PE300	Zymo Mock	1	49,828	24,301	48.8%	NA
Amplicon-16S/V4	PE250	Mice feces	8	526,754	389,000	(77.6±7.1)%	NA
Enriched	PE150	Mice feces	8	8,321,081	11,474,476	(70.1±5.4)%	(59.1±0.8)%
Enriched	PE150	Rat feces	5	6,450,541	9,470,428	(72.9±2.1)%	(57.8±4.1)%
Enriched	PE150	BEI Mock	1	5,345,638	8,203,396	76.7%	70.4%
Enriched	PE150	Zymo Mock	1	3,359,376	5,140,030	76.5%	70.1%
Enriched	PE300	Mice feces	8	1,050,608	1,573,122	(75.1±3.2)%	(59.9±2.1)%
Enriched	PE300	BEI Mock	1	737,309	1,108,481	75.2%	75.7%
Enriched	PE300	Zymo Mock	1	467,250	721,740	77.2%	73.8%
Unenriched	PE150	Mice feces	8	28,219,552	37,894,050	(68.6±6.4)%	0.1%
Unenriched	PE150	Rat feces	5	16,266,683	28,448,468	(87.4±0.9)%	0.1%
Unenriched	PE150	BEI Mock	1	6,263,379	8,889,636	71%	0.2%
Unenriched	PE150	Zymo Mock	1	4,985,957	7,001,503	70.2%	0.2%

526

527 **Table 2.** Summary statistics for simulated data and real data from mock communities, libraries were enriched for 16S using the 16S-cap

528 enrichment and sequenced on an Illumina MiSeq PE150 reads.

Sample ID	Library Type	Avg. No. of (Simulated) Reads	No. of Simulated Reads	Matched Pairs	Matched Forward	Matched Reverse	Total Mapped	Percent of Avg. Total Mapped
Simulated Data								
Zymo Mock	Enriched-PE150	412,520	206,260	171708	190,964	186,216	377,180	91.43%
BEI Mock	Enriched-PE150	415,472	207736	176547	193998	189,777	383,775	92.37%
Lindgreen et al., 2016	Enriched-PE150	490,238	245119	188620	218911	213,918	432,829	88.29%
Real Data								
Zymo Mock	Enriched-PE150	3,904,480	1,952,240	1,314,654	1,548,323	1,503,225	3,051,548	78.15%
BEI Mock	Enriched-PE150	6,260,110	3,130,055	2,127,656	2,486,274	2,435,425	4,921,699	78.62%

529

530 ***Conflict of Interest***

531 The EHS DNA lab provides oligonucleotide aliquots and library preparation services at cost,
532 including some oligonucleotides and services used in this manuscript (baddna.uga.edu). Brian
533 Brunelle and Alison Devault are employed by, and thereby have financial interest in, Daicel Arbor
534 Biosciences, who provided the in-solution capture reagents used in this work.

535 ***Author Contributions***

536 TG conceived of the project. JW, JT, TK, BG, KL, and TG designed experiments. JW, TK, and BG
537 performed the experiments. AD and BB designed the baits. JW and NB analyzed the data. AD, BB,
538 KL, OR, JSW, and TG provided funding and resources. MB wrote the manuscript. JW, TK, NB
539 wrote sections of the manuscript. MB and JW produced figures and tables. All authors critically
540 reviewed, edited, and approve of this work.

541 ***Funding***

542 Funding for this grant was provided by the Center for Disease Control contract 200-2018-02889
543 (75D30118C02889), US Department of Energy Cooperative Agreement number DE-FC09-
544 07SR22506, National Institute of Health (R01ES024950, P30ES010126, and P42ES031007) and
545 United States Agency for International Development via Peanut and Mycotoxin Innovation
546 Laboratory (ECG-A00-13-00001-00). Daicel Arbor Biosciences provided the customized in-solution
547 capture reagents used in this work.

548 ***Acknowledgments***

549 We acknowledge the contributions of Marissa Howard, Amanda Sullivan, Allison Perry, and Laura
550 Rose.

551 ***Disclaimer***

552 This report was prepared as an account of work sponsored by agencies of the United States
553 Government. Neither the United States Government, nor any agency thereof, nor any of their
554 employees makes any warranty, express or implied, or assumes any legal liability or responsibility
555 for the accuracy, completeness, or usefulness of any information, apparatus, product, or process
556 disclosed or represents that its use would not infringe privately owned rights. Reference herein to any

557 specific commercial product, process, or service by trade name, trademark, manufacturer, or
558 otherwise does not necessarily constitute or imply its endorsement, recommendation, or favoring by
559 the United States Government or any agency thereof. The views and opinions of authors expressed
560 herein do not necessarily state or reflect those of the United States Government or any agency
561 thereof.

562

563

Bibliography

564

- 565 Aird, D., Ross, M.G., Wei-Sheng, C., Danielsson, M., Fennell, T., Russ, C., et al. (2011). Analyzing
566 and minimizing PCR amplification bias in Illumina sequencing libraries. *Genome Biology* 12,
567 1-14.
- 568 Altschul, S.F., Gish, W., Miller, W., Myers, E., and Lipman, D.J. (1990). Basic Local Alignment
569 Search Tool. *Journal of Molecular Biology* 215, 403-410.
- 570 Balvociute, M., and Huson, D.H. (2017). SILVA, RDP, Greengenes, NCBI and OTT - how do these
571 taxonomies compare? *BMC Genomics* 18(Suppl 2), 114. doi: 10.1186/s12864-017-3501-4.
- 572 Barrett, S.R., Hoffman, N.G., Rosenthal, C., Bryan, A., Marshall, D.A., Lieberman, J., et al. (2020).
573 Sensitive identification of bacterial DNA in clinical specimens by broad range 16S rRNA
574 enrichment. *J Clin Microbiol*, 1-30. doi: 10.1128/JCM.01605-20.
- 575 Bolger, A.M., Lohse, M., and Usadel, B. (2014). Trimmomatic: a flexible trimmer for Illumina
576 sequence data. *Bioinformatics* 30(15), 2114-2120. doi: 10.1093/bioinformatics/btu170.
- 577 Bushnell, B. (2014). BBMAP: A fast, accurate, splice-aware aligner. *Lawrence Berkeley National*
578 *Laboratory*.
- 579 Callahan, B.J., McMurdie, P.J., Rosen, M.J., Han, A.W., Johnson, A.J., and Holmes, S.P. (2016a).
580 DADA2: High-resolution sample inference from Illumina amplicon data. *Nat Methods* 13(7),
581 581-583. doi: 10.1038/nmeth.3869.
- 582 Callahan, B.J., McMurdie, P.J., Rosen, M.J., Han, A.W., Johnson, A.J.A., and Holmes, S.P. (2016b).
583 DADA2: High-resolution sample inference from Illumina amplicon data. *Nature Methods*
584 13(7), 581-583. doi: 10.1038/nmeth.3869.
- 585 Caporaso, J.G., Lauber, C.L., Walters, W.A., Berg-Lyons, D., Huntley, J., Fierer, N., et al. (2012).
586 Ultra-high-throughput microbial community analysis on the Illumina HiSeq and MiSeq
587 platforms. *ISME J* 6(8), 1621-1624. doi: 10.1038/ismej.2012.8.
- 588 Costea, P.I., Zeller, G., Sunagawa, S., Pelletier, E., Alberti, A., Levenez, F., et al. (2017). Towards
589 standards for human fecal sample processing in metagenomic studies. *Nat Biotechnol* 35(11),
590 1069-1076. doi: 10.1038/nbt.3960.
- 591 Djemiel, C., Dequiedt, S., Karimi, B., Cottin, A., Girier, T., El Djoudi, Y., et al. (2020). BIOCUM-
592 PIPE: a new user-friendly metabarcoding pipeline for the characterization of microbial
593 diversity from 16S, 18S and 23S rRNA gene amplicons. *BMC Bioinformatics* 21(1), 492. doi:
594 10.1186/s12859-020-03829-3.
- 595 Droge, J., and McHardy, A.C. (2012). Taxonomic binning of metagenome samples generated by
596 next-generation sequencing technologies. *Brief Bioinform* 13(6), 646-655. doi:
597 10.1093/bib/bbs031.
- 598 Dueholm, M.S., Andersen, K.S., McIlroy, S.J., Kristensen, J.M., Yashiro, E., Karst, S.M., et al.
599 (2020). Generation of Comprehensive Ecosystem-Specific Reference Databases with Species-
600 Level Resolution by High-Throughput Full-Length 16S rRNA Gene Sequencing and
601 Automated Taxonomy Assignment (AutoTax). *mBio* 11(5). doi: 10.1128/mBio.01557-20.

- 602 Dueker, M.E., French, S., and O'Mullan, G.D. (2018). Comparison of Bacterial Diversity in Air and
603 Water of a Major Urban Center. *Front Microbiol* 9, 2868. doi: 10.3389/fmicb.2018.02868.
- 604 Edgar, R.C. (2010). Search and clustering orders of magnitude faster than BLAST. *Bioinformatics*
605 26(19), 2460-2461. doi: 10.1093/bioinformatics/btq461.
- 606 Edgar, R.C. (2013). UPARSE: highly accurate OTU sequences from microbial amplicon reads. *Nat*
607 *Methods* 10(10), 996-998. doi: 10.1038/nmeth.2604.
- 608 Eisenhofer, R., Minich, J.J., Marotz, C., Cooper, A., Knight, R., and Weyrich, L.S. (2019).
609 Contamination in Low Microbial Biomass Microbiome Studies: Issues and
610 Recommendations. *Trends Microbiol* 27(2), 105-117. doi: 10.1016/j.tim.2018.11.003.
- 611 Escobar-Zepeda, A., Godoy-Lozano, E.E., Raggi, L., Segovia, L., Merino, E., Gutierrez-Rios, R.M.,
612 et al. (2018). Analysis of sequencing strategies and tools for taxonomic annotation: Defining
613 standards for progressive metagenomics. *Sci Rep* 8(1), 12034. doi: 10.1038/s41598-018-
614 30515-5.
- 615 Gallardo-Escárate, C., Valenzuela-Muñoz, V., Núñez-Acuña, G., Valenzuela-Miranda, D., Castellón,
616 F., Benavente-Cartes, B., et al. (2020). The wastewater microbiome: a novel insight for
617 COVID-19 surveillance. *Research Square*, 1-20. doi: 10.21203/rs.3.rs-62651/v1.
- 618 Gao, B., Bian, X., Mahbub, R., and Lu, K. (2017). Sex-Specific Effects of Organophosphate
619 Diazinon on the Gut Microbiome and Its Metabolic Functions. *Environ Health Perspect*
620 125(2), 198-206. doi: 10.1289/EHP202.
- 621 Gasc, C., and Peyret, P. (2018). Hybridization capture reveals microbial diversity missed using
622 current profiling methods. *Microbiome* 6(1), 61. doi: 10.1186/s40168-018-0442-3.
- 623 Gilmour, M.W., Graham, M., Van Domselaar, G., Tyler, S., Kent, H., Trout-Yakel, K.M., et al.
624 (2010). High-throughput genome sequencing of two *Listeria monocytogenes* clinical isolates
625 during a large foodborne outbreak. *BMC Genomics* 11, 120. doi: 10.1186/1471-2164-11-120.
- 626 Glenn, T.C., and Faircloth, B.C. (2016). Capturing Darwin's dream. *Mol Ecol Resour* 16(5), 1051-
627 1058. doi: 10.1111/1755-0998.12574.
- 628 Glenn, T.C., Nilsen, R.A., Kieran, T.J., Sanders, J.G., Bayona-Vasquez, N.J., Finger, J.W., et al.
629 (2019a). Adapterama I: universal stubs and primers for 384 unique dual-indexed or 147,456
630 combinatorially-indexed Illumina libraries (iTru & iNext). *PeerJ* 7, e7755. doi:
631 10.7717/peerj.7755.
- 632 Glenn, T.C., Pierson, T.W., Bayona-Vasquez, N.J., Kieran, T.J., Hoffberg, S.L., Thomas Iv, J.C., et
633 al. (2019b). Adapterama II: universal amplicon sequencing on Illumina platforms
634 (TaggiMatrix). *PeerJ* 7, e7786. doi: 10.7717/peerj.7786.
- 635 Guitor, A.K., Raphenya, A.R., Klunk, J., Kuch, M., Alcock, B., Surette, M.G., et al. (2019).
636 Capturing the Resistome: a Targeted Capture Method To Reveal Antibiotic Resistance
637 Determinants in Metagenomes. *Antimicrob Agents Chemother* 64(1), 1-37. doi:
638 10.1128/AAC.01324-19.
- 639 Huang, W., Li, L., Myers, J.R., and Marth, G.T. (2012). ART: a next-generation sequencing read
640 simulator. *Bioinformatics* 28(4), 593-594. doi: 10.1093/bioinformatics/btr708.
- 641 Jiang, Y., Xiong, X., Danska, J., and Parkinson, J. (2016). Metatranscriptomic analysis of diverse
642 microbial communities reveals core metabolic pathways and microbiome-specific
643 functionality. *Microbiome* 4, 2. doi: 10.1186/s40168-015-0146-x.
- 644 Jouglin, M., Blanc, B., de la Cotte, N., Bastian, S., Ortiz, K., and Malandrino, L. (2019). First
645 detection and molecular identification of the zoonotic *Anaplasma capra* in deer in France.
646 *PLoS One* 14(7), e0219184. doi: 10.1371/journal.pone.0219184.
- 647 Jousset, E., Clamens, A.L., Galan, M., Bernard, M., Maman, S., Gschloessl, B., et al. (2016).
648 Assessment of a 16S rRNA amplicon Illumina sequencing procedure for studying the
649 microbiome of a symbiont-rich aphid genus. *Mol Ecol Resour* 16(3), 628-640. doi:
650 10.1111/1755-0998.12478.

- 651 Jovel, J., Patterson, J., Wang, W., Hotte, N., O'Keefe, S., Mitchel, T., et al. (2016). Characterization
652 of the Gut Microbiome Using 16S or Shotgun Metagenomics. *Front Microbiol* 7, 459. doi:
653 10.3389/fmicb.2016.00459.
- 654 Kelly, R.P., Shelton, A.O., and Gallego, R. (2019). Understanding PCR Processes to Draw
655 Meaningful Conclusions from Environmental DNA Studies. *Sci Rep* 9(1), 12133. doi:
656 10.1038/s41598-019-48546-x.
- 657 Kennedy, K., Hall, M.W., Lynch, M.D., Moreno-Hagelsieb, G., and Neufeld, J.D. (2014). Evaluating
658 bias of illumina-based bacterial 16S rRNA gene profiles. *Appl Environ Microbiol* 80(18),
659 5717-5722. doi: 10.1128/AEM.01451-14.
- 660 Kim, D., Hofstaedter, C.E., Zhao, C., Mattei, L., Tanes, C., Clarke, E., et al. (2017). Optimizing
661 methods and dodging pitfalls in microbiome research. *Microbiome* 5(1), 52. doi:
662 10.1186/s40168-017-0267-5.
- 663 Klindworth, A., Pruesse, E., Schweer, T., Peplies, J., Quast, C., Horn, M., et al. (2013). Evaluation of
664 general 16S ribosomal RNA gene PCR primers for classical and next-generation sequencing-
665 based diversity studies. *Nucleic Acids Res* 41(1), e1. doi: 10.1093/nar/gks808.
- 666 Knight, R., Vrbanac, A., Taylor, B.C., Aksenov, A., Callewaert, C., Debelius, J., et al. (2018). Best
667 practices for analysing microbiomes. *Nat Rev Microbiol* 16(7), 410-422. doi:
668 10.1038/s41579-018-0029-9.
- 669 Kustin, T., Ling, G., Sharabi, S., Ram, D., Friedman, N., Zuckerman, N., et al. (2019). A method to
670 identify respiratory virus infections in clinical samples using next-generation sequencing. *Sci*
671 *Rep* 9(1), 2606. doi: 10.1038/s41598-018-37483-w.
- 672 Lasa, A., di Cesare, A., Tassistro, G., Borello, A., Gualdi, S., Furones, D., et al. (2019). Dynamics of
673 the Pacific oyster pathobiota during mortality episodes in Europe assessed by 16S rRNA gene
674 profiling and a new target enrichment next-generation sequencing strategy. *Environ*
675 *Microbiol* 21(12), 4548-4562. doi: 10.1111/1462-2920.14750.
- 676 Li, H., and Durbin, R. (2009). Fast and accurate short read alignment with Burrows-Wheeler
677 transform. *Bioinformatics* 25(14), 1754-1760. doi: 10.1093/bioinformatics/btp324.
- 678 Li, R., Yu, C., Li, Y., Lam, T.W., Yiu, S.M., Kristiansen, K., et al. (2009). SOAP2: an improved
679 ultrafast tool for short read alignment. *Bioinformatics* 25(15), 1966-1967. doi:
680 10.1093/bioinformatics/btp336.
- 681 Lindgreen, S., Adair, K.L., and Gardner, P.P. (2016). An evaluation of the accuracy and speed of
682 metagenome analysis tools. *Sci Rep* 6, 19233. doi: 10.1038/srep19233.
- 683 Logares, R., Sunagawa, S., Salazar, G., Cornejo-Castillo, F.M., Ferrera, I., Sarmiento, H., et al.
684 (2014). Metagenomic 16S rDNA Illumina tags are a powerful alternative to amplicon
685 sequencing to explore diversity and structure of microbial communities. *Environ Microbiol*
686 16(9), 2659-2671. doi: 10.1111/1462-2920.12250.
- 687 Lugli, G.A., Duranti, S., Milani, C., Mancabelli, L., Turrone, F., Sinderen, D.V., et al. (2019).
688 Uncovering Bifidobacteria via Targeted Sequencing of the Mammalian Gut Microbiota.
689 *Microorganisms* 7(11), 1-11. doi: 10.3390/microorganisms7110535.
- 690 Martin, M. (2011). Cutadapt removes adapter sequences from high-throughput sequencing reads.
691 *EMBnet* 17(1), 10-12.
- 692 Martinez-Porchas, M., Villalpando-Canchola, E., and Vargas-Albores, F. (2016). Significant loss of
693 sensitivity and specificity in the taxonomic classification occurs when short 16S rRNA gene
694 sequences are used. *Heliyon* 2(9), e00170. doi: 10.1016/j.heliyon.2016.e00170.
- 695 Meola, M., Lazzaro, A., and Zeyer, J. (2015). Bacterial Composition and Survival on Sahara Dust
696 Particles Transported to the European Alps. *Front Microbiol* 6, 1454. doi:
697 10.3389/fmicb.2015.01454.

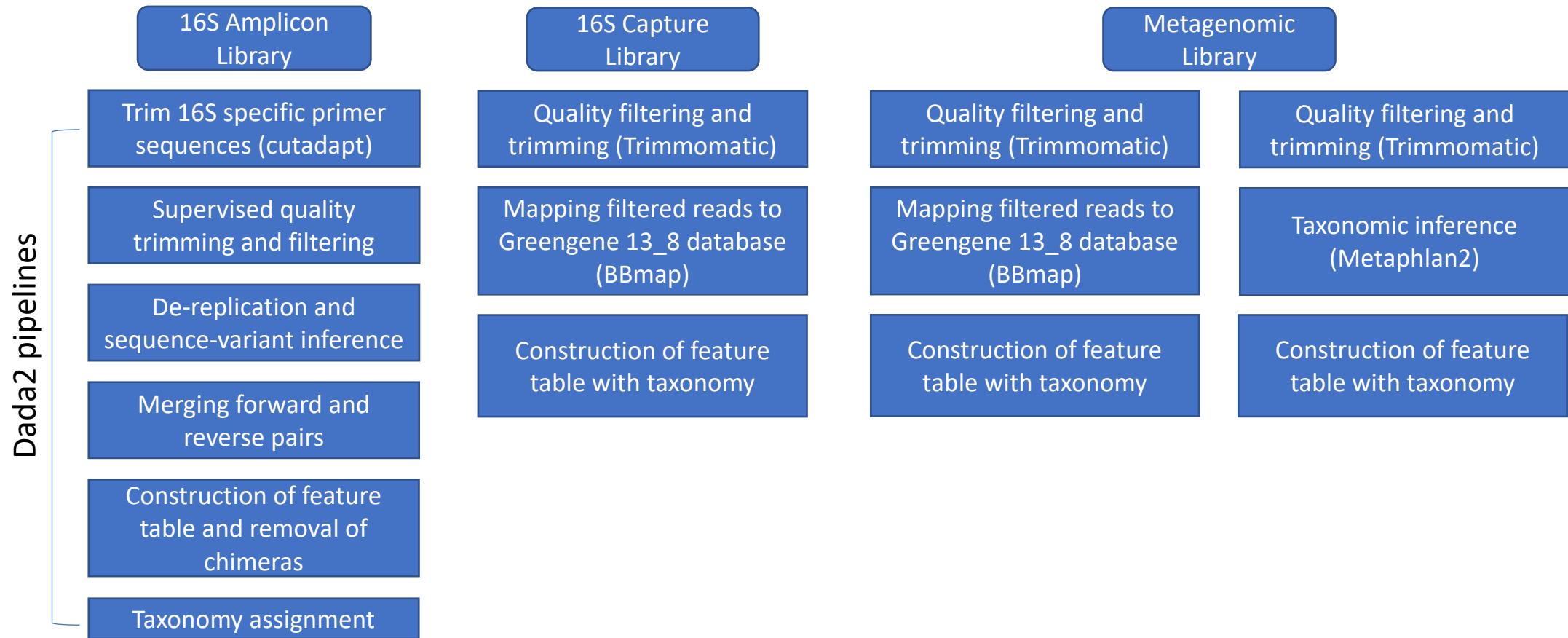
- 698 Mysara, M., Vandamme, P., Props, R., Kerckhof, F.M., Leys, N., Boon, N., et al. (2017).
699 Reconciliation between operational taxonomic units and species boundaries. *FEMS Microbiol*
700 *Ecol* 93(4), 1-12. doi: 10.1093/femsec/fix029.
- 701 Neelakanta, G., and Sultana, H. (2013). The use of metagenomic approaches to analyze changes in
702 microbial communities. *Microbiol Insights* 6, 37-48. doi: 10.4137/MBI.S10819.
- 703 Oladeinde, A., Cook, K., Lakin, S., Woyda, R., Abdo, Z., Looft, T., et al. (2019). Horizontal Gene
704 Transfer and Acquired Antibiotic Resistance in *Salmonella enterica* Serovar Heidelberg
705 following In Vitro Incubation in Broiler Ceca. *Applied and Environmental Microbiology*
706 85(22), e01903-01919.
- 707 Polka, J., Rebecchi, A., Pisacane, V., Morelli, L., and Puglisi, E. (2015). Bacterial diversity in typical
708 Italian salami at different ripening stages as revealed by high-throughput sequencing of 16S
709 rRNA amplicons. *Food Microbiol* 46, 342-356. doi: 10.1016/j.fm.2014.08.023.
- 710 Pollock, J., Glendinning, L., Wisedchanwet, T., and Watson, M. (2018). The Madness of
711 Microbiome: Attempting To Find Consensus "Best Practice" for 16S Microbiome Studies.
712 *Appl Environ Microbiol* 84(7). doi: 10.1128/AEM.02627-17.
- 713 Poretsky, R., Rodriguez, R.L., Luo, C., Tsementzi, D., and Konstantinidis, K.T. (2014). Strengths
714 and limitations of 16S rRNA gene amplicon sequencing in revealing temporal microbial
715 community dynamics. *PLoS One* 9(4), e93827. doi: 10.1371/journal.pone.0093827.
- 716 R Development Core Team (2010). (Vienna, Austria: R Foundation for Statistical Computing).
- 717 Rausch, P., Ruhlemann, M., Hermes, B.M., Doms, S., Dagan, T., Dierking, K., et al. (2019).
718 Comparative analysis of amplicon and metagenomic sequencing methods reveals key features
719 in the evolution of animal metaorganisms. *Microbiome* 7(1), 133. doi: 10.1186/s40168-019-
720 0743-1.
- 721 Ritari, J., Salojarvi, J., Lahti, L., and de Vos, W.M. (2015). Improved taxonomic assignment of
722 human intestinal 16S rRNA sequences by a dedicated reference database. *BMC Genomics* 16,
723 1056. doi: 10.1186/s12864-015-2265-y.
- 724 Salter, S.J., Cox, M.J., Turek, E.M., Calus, S.T., Cookson, W.O., Moffatt, M.F., et al. (2014).
725 Reagent and laboratory contamination can critically impact sequence-based microbiome
726 analyses. *BM Biology* 12(87), 1-12.
- 727 Schloss, P.D., and Westcott, S.L. (2011). Assessing and improving methods used in operational
728 taxonomic unit-based approaches for 16S rRNA gene sequence analysis. *Appl Environ*
729 *Microbiol* 77(10), 3219-3226. doi: 10.1128/AEM.02810-10.
- 730 Schmieder, R., and Edwards, R. (2011). Fast identification and removal of sequence contamination
731 from genomic and metagenomic datasets. *PLoS One* 6(3), e17288. doi:
732 10.1371/journal.pone.0017288.
- 733 Seemann, T. (2014). Prokka: rapid prokaryotic genome annotation. *Bioinformatics* 30(14), 2068-
734 2069. doi: 10.1093/bioinformatics/btu153.
- 735 Segata, N., Waldron, L., Ballarini, A., Narasimhan, V., Jousson, O., and Huttenhower, C. (2012).
736 Metagenomic microbial community profiling using unique clade-specific marker genes. *Nat*
737 *Methods* 9(8), 811-814. doi: 10.1038/nmeth.2066.
- 738 Sekse, C., Holst-Jensen, A., Dobrindt, U., Johannessen, G.S., Li, W., Spilberg, B., et al. (2017).
739 High Throughput Sequencing for Detection of Foodborne Pathogens. *Front Microbiol* 8,
740 2029. doi: 10.3389/fmicb.2017.02029.
- 741 Suenami, S., Konishi Nobu, M., and Miyazaki, R. (2019). Community analysis of gut microbiota in
742 hornets, the largest eusocial wasps, *Vespa mandarinia* and *V. simillima*. *Sci Rep* 9(1), 9830.
743 doi: 10.1038/s41598-019-46388-1.
- 744 Thomas, J.C.t., Oladeinde, A., Kieran, T.J., Finger, J.W., Jr., Bayona-Vasquez, N.J., Cartee, J.C., et
745 al. (2020). Co-occurrence of antibiotic, biocide, and heavy metal resistance genes in bacteria

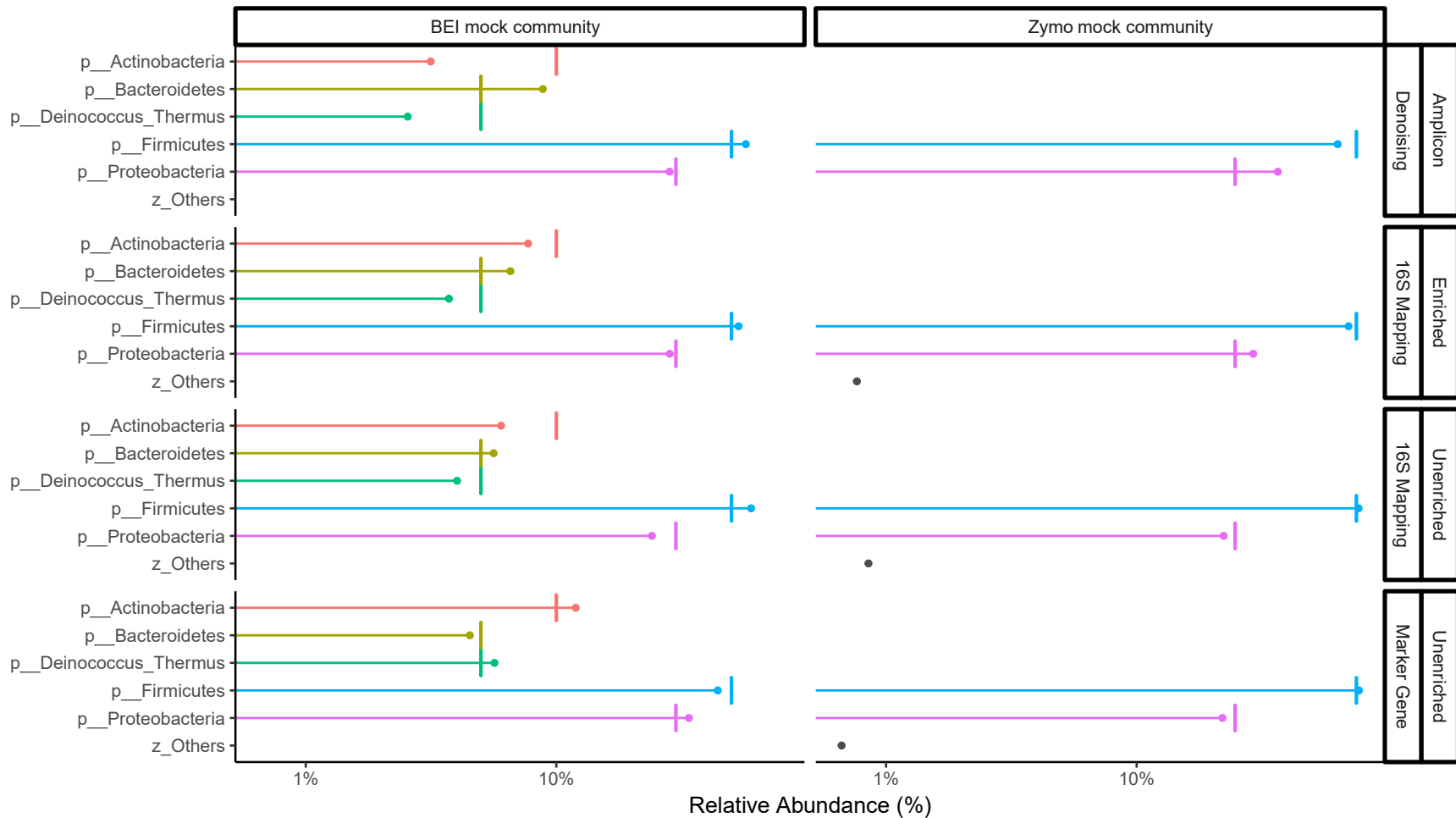
- 746 from metal and radionuclide contaminated soils at the Savannah River Site. *Microb*
747 *Biotechnol* 13(4), 1179-1200. doi: 10.1111/1751-7915.13578.
- 748 Truong, D.T., Franzosa, E.A., Tickle, T.L., Scholz, M., Weingart, G., Pasolli, E., et al. (2015).
749 MetaPhlan2 for enhanced metagenomic taxonomic profiling. *Nat Methods* 12(10), 902-903.
750 doi: 10.1038/nmeth.3589.
- 751 van der Valk, T., Vezzi, F., Ormestad, M., Dalen, L., and Guschanski, K. (2019). Index hopping on
752 the Illumina HiSeqX platform and its consequences for ancient DNA studies. *Mol Ecol*
753 *Resour*, 1171-1181. doi: 10.1111/1755-0998.13009.
- 754 Vetrovsky, T., and Baldrian, P. (2013). The variability of the 16S rRNA gene in bacterial genomes
755 and its consequences for bacterial community analyses. *PLoS One* 8(2), e57923. doi:
756 10.1371/journal.pone.0057923.
- 757 Walsh, A.M., Crispie, F., O'Sullivan, O., Finnegan, L., Claesson, M.J., and Cotter, P.D. (2018).
758 Species classifier choice is a key consideration when analysing low-complexity food
759 microbiome data. *Microbiome* 6(1), 50. doi: 10.1186/s40168-018-0437-0.
- 760 Wang, J., Tang, L., Zhou, H., Zhou, J., Glenn, T.C., Shen, C.L., et al. (2018). Long-term treatment
761 with green tea polyphenols modifies the gut microbiome of female sprague-dawley rats. *J*
762 *Nutr Biochem* 56, 55-64. doi: 10.1016/j.jnutbio.2018.01.005.
- 763 Weiss, S., Amir, A., Hyde, E.R., Metcalf, J.L., Song, S.J., and Knight, R. (2014). Tracking down the
764 sources of experimental contamination in microbiome studies. *Genome Biology* 15, 1-3.
- 765 Woese, C.R., and Fox, G.E. (1977). Phylogenetic structure of the prokaryotic domain: the primary
766 kingdoms. *Proc Natl Acad Sci U S A* 74(11), 5088-5090. doi: 10.1073/pnas.74.11.5088.
- 767 Yuan, C., Lei, J., Cole, J., and Sun, Y. (2015). Reconstructing 16S rRNA genes in metagenomic data.
768 *Bioinformatics* 31(12), i35-i43. doi: 10.1093/bioinformatics/btv231.
- 769 Ziegler, M., Grupstra, C.G.B., Barreto, M.M., Eaton, M., BaOmar, J., Zubier, K., et al. (2019). Coral
770 bacterial community structure responds to environmental change in a host-specific manner.
771 *Nat Commun* 10(1), 3092. doi: 10.1038/s41467-019-10969-5.
- 772 Zinger, L., Bonin, A., Alsos, I.G., Balint, M., Bik, H., Boyer, F., et al. (2019). DNA metabarcoding-
773 Need for robust experimental designs to draw sound ecological conclusions. *Mol Ecol* 28(8),
774 1857-1862. doi: 10.1111/mec.15060.
- 775 Zinter, M.S., Mayday, M.Y., Ryckman, K.K., Jelliffe-Pawlowski, L.L., and DeRisi, J.L. (2019).
776 Towards precision quantification of contamination in metagenomic sequencing experiments.
777 *Microbiome* 7(1), 62. doi: 10.1186/s40168-019-0678-6.

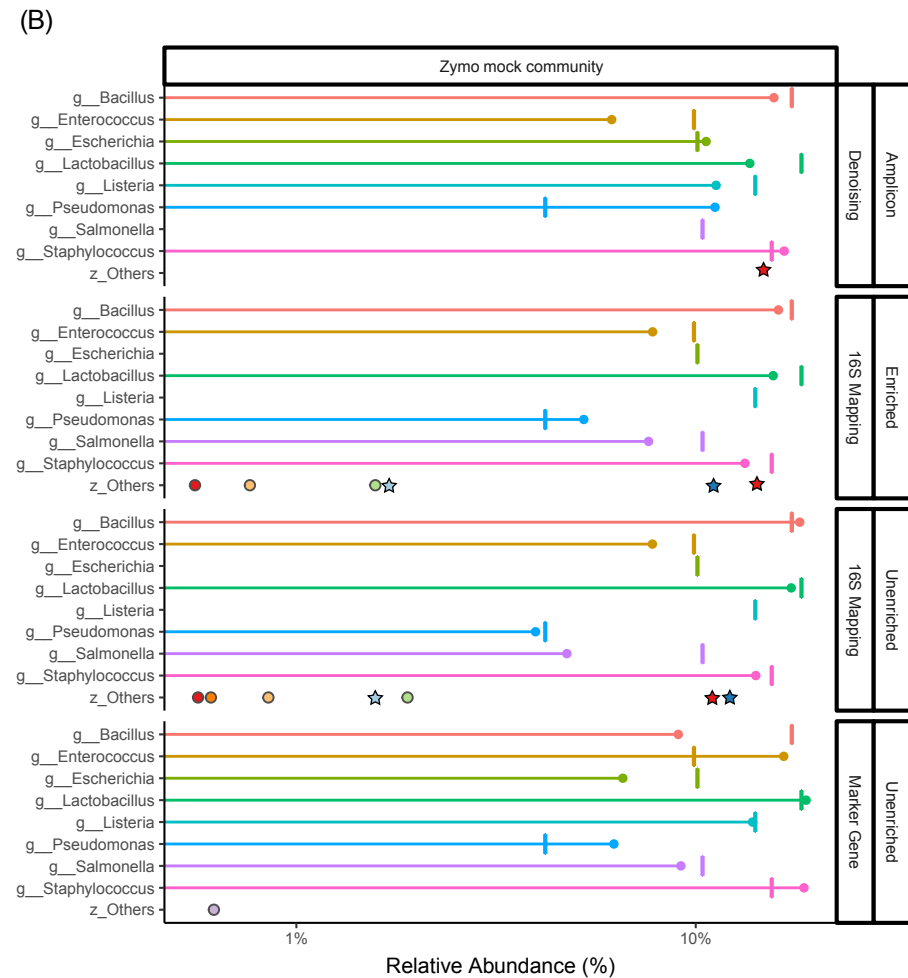
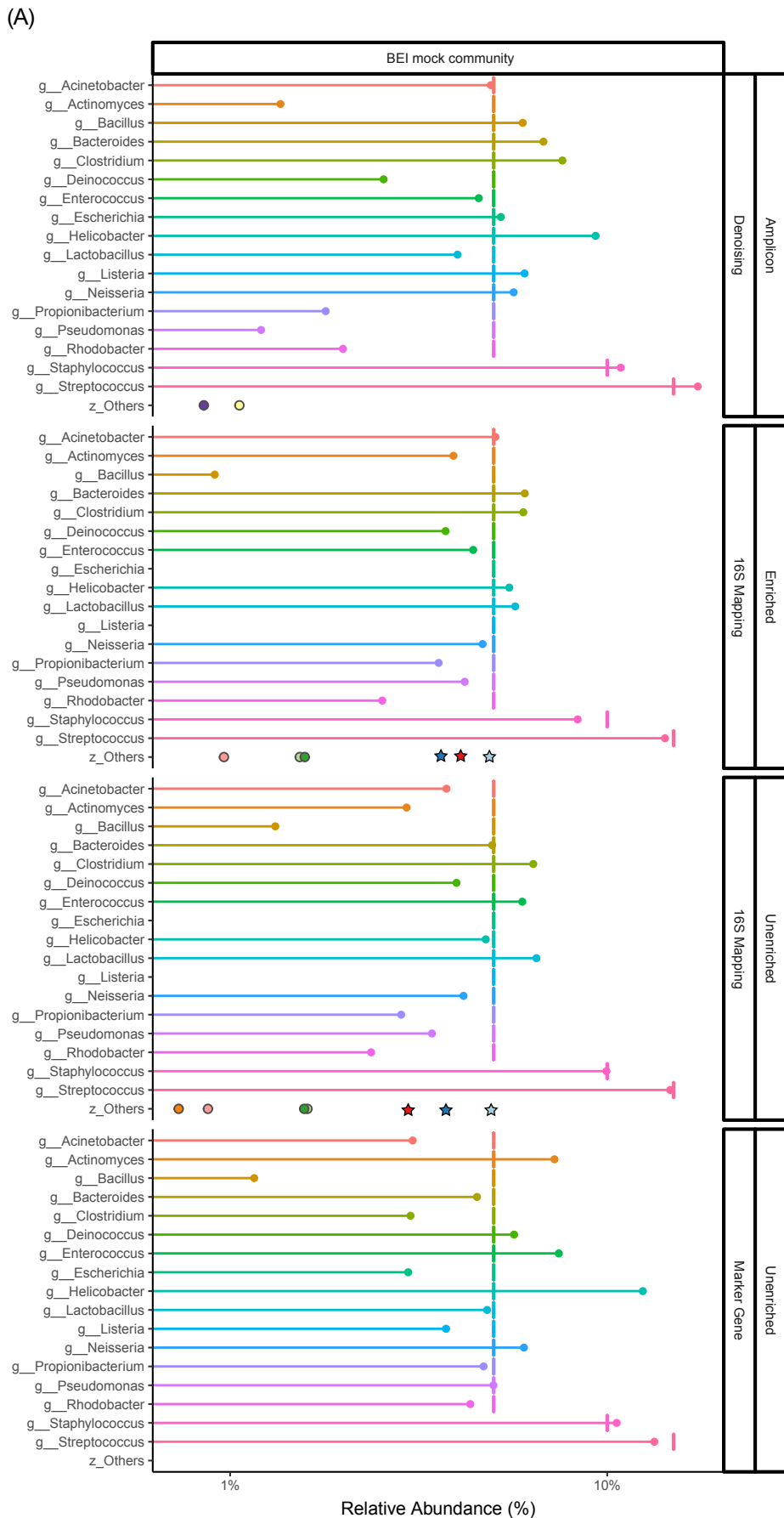
778 **Data Availability**

779 The datasets generated in this study can be found in the authors dropbox
780 <https://www.dropbox.com/sh/exg0kow6pyghlmx/AAAIn7R93EawGUDO7TQ6NDIY?dl=0> .

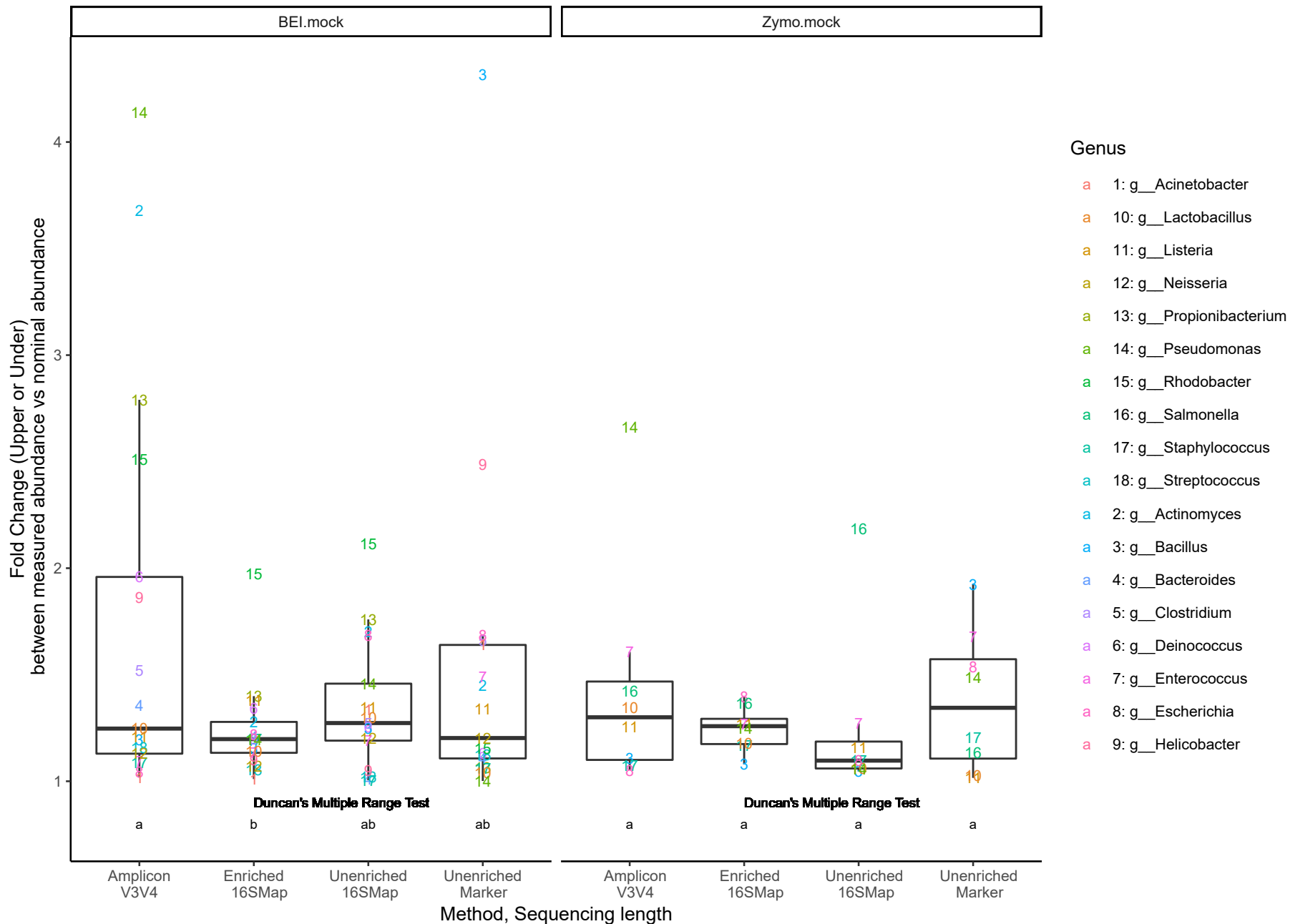
781



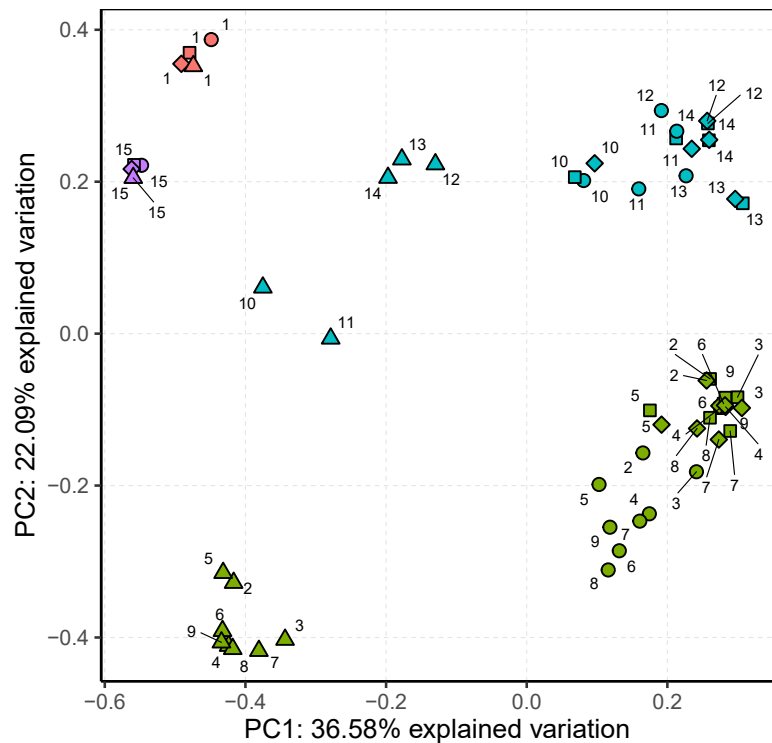




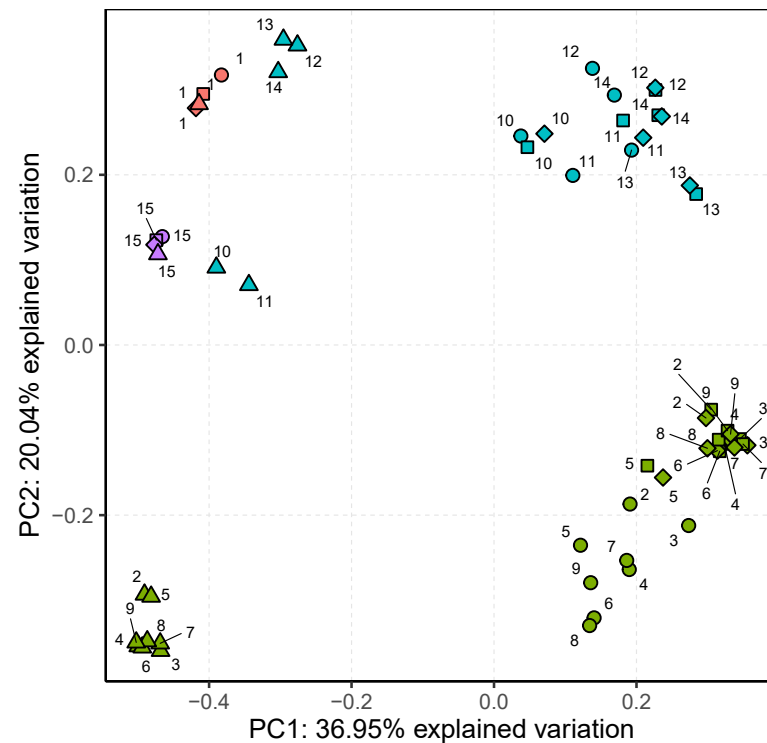
- Family
- ☆ f__Bacillaceae
 - ★ f__Enterobacteriaceae
 - ★ f__Listeriaceae
 - f__Carnobacteriaceae
 - f__Clostridiaceae
 - f__Acaryochloridaceae
 - f__Neisseriaceae
 - f__Planococcaceae
 - f__Prevotellaceae
 - f__S24-7
 - f__Saccharomycetaceae



(A) Family Level PCoA



(B) Genus Level PCoA



Projects

- BEI.mock
- Mouse
- Rat
- Zymo.mock

Library

- Amplicon_Denosing
- Enriched_16SMapping
- ◇ Unenriched_16SMapping
- △ Unenriched_Marker

Comparison of bray-curtis distance

Based on Genus level composition

