

1 **Optimization of 16S amplicon analysis using mock communities: implications for**  
2 **estimating community diversity**

3

4 Andrew Krohn<sup>1,2</sup>, Bo Stevens<sup>1</sup>, Adam Robbins-Pianka<sup>4</sup>, Matthew Belus<sup>5</sup>, Gerard J. Allan<sup>1,2</sup>,  
5 Catherine Gehring<sup>1,3</sup>

6

7 <sup>1</sup>Department of Biological Sciences, Northern Arizona University, Flagstaff, AZ

8 <sup>2</sup>NAU Environmental Genetics and Genomics Laboratory, Flagstaff, AZ

9 <sup>3</sup>Merriam-Powell Center for Environmental Research, Flagstaff, AZ

10 <sup>4</sup>Department of Computer Science, University of Colorado Boulder, Boulder, CO

11 <sup>5</sup>Anschutz Medical Campus, University of Colorado Denver, Aurora, CO

12

13 Corresponding Author:

14 Andrew Krohn<sup>1,2</sup>

15

16 Email address: [alk224@nau.edu](mailto:alk224@nau.edu)

17

18

19

20

21

22

23

24

25

26

27

28

29

30

31

**32 Abstract:**

33 Diversity of complex microbial communities can be rapidly assessed by community  
34 amplicon sequencing of marker genes (*e.g.*, 16S), often yielding many thousands of DNA  
35 sequences per sample. However, analysis of community amplicon sequencing data requires  
36 multiple computational steps which affect the outcome of a final data set. Here we use mock  
37 communities to describe the effects of parameter adjustments for raw sequence quality filtering,  
38 picking operational taxonomic units (OTUs), taxonomic assignment, and OTU table filtering as  
39 implemented in QIIME 1.9.1. We demonstrate a workflow optimization based upon this  
40 exploration which we also apply to environmental samples. We found that quality filtering of  
41 raw data and filtering of OTU tables had large effects on observed OTU diversity. While all  
42 taxonomy assigners performed with similar accuracy, an appropriate choice of similarity  
43 threshold for defining OTUs depended on the method used for OTU picking. Our “default”  
44 analysis in QIIME overestimated mock community diversity by at least a factor of ten, compared  
45 to the optimized analysis which correctly characterized the taxonomic composition of the mock  
46 communities while still overestimating OTU diversity by about a factor of two. Though observed  
47 relative abundances of mock community member taxa were approximately correct, most were  
48 still represented by multiple OTUs. Low-frequency OTUs conspecific to constituent mock  
49 community taxa were characterized by multiple substitution and indel errors and the presence of  
50 a low quality base call resulting in sequence truncation during quality filtering. Low quality base  
51 calls were observed at “G” positions most of the time, and were also associated with a preceding  
52 “TTT” trinucleotide motif. Environmental diversity estimates were reduced by about 40% from  
53 2508 to 1533 OTUs when comparing output from the default and optimized workflows. We  
54 attribute this reduction in observed diversity to the removal of erroneous sequences from the data  
55 set. Our results indicate that both strict quality filtering of raw sequencing data and careful  
56 filtering of raw OTU tables are important steps for accurate estimation of microbial community  
57 diversity.

58

**59 Introduction:**

60 Over the past decade, community amplicon sequencing has become the preferred method  
61 for profiling diversity in microbial communities. Briefly, the technique uses the polymerase  
62 chain reaction (PCR) to amplify a pool of PCR products from an environmental sample to be

63 resolved by high throughput DNA sequencing. Similar sequences are binned together into  
64 operational taxonomic units (OTUs) which are compared against a database to obtain taxonomic  
65 classifications. Amplicon sequencing is flexible in that a community can be profiled for different  
66 genes which may represent markers better suited for certain microbial constituents (*e.g.*, 16S for  
67 prokaryotes, ITS for fungi), while profiling with functional genes can offer a better  
68 understanding of community traits (*e.g.*, Bentzon-Tilia *et al.*, 2015). While communities were  
69 originally profiled on 454 pyrosequencing instruments (Sogin *et al.*, 2006), amplicon sequencing  
70 has been adapted to newer instrumentation including sequencers from Illumina (Caporaso *et al.*,  
71 2012) and Pacific Biosciences (Fichot & Norman, 2013). Illumina sequencing is currently the  
72 most popular option due to several factors including cost, throughput, instrument availability,  
73 and the existence of multiple protocols for amplification and sequencing of marker gene pools on  
74 this platform (Caporaso *et al.*, 2012; Bokulich & Mills, 2013; Kozich *et al.*, 2013; Fadrosh *et al.*,  
75 2014).

76         Accurate determination of community diversity and taxonomic content are often primary  
77 aims of community amplicon sequencing projects. Systematic errors experienced during sample  
78 preparation such as PCR and sequencing errors can contribute to overestimates of diversity  
79 (Kunin, 2010). Additionally, signal cross-talk during index sequence cycles on Illumina  
80 sequencers can lead a researcher to falsely conclude that an organism is present in a sample  
81 (Kircher, Sawyer & Meyer, 2012; Nelson *et al.*, 2014). In the face of such potential  
82 complications, careful analysis is warranted to ensure that diversity estimates are not inflated and  
83 that data are properly filtered to avoid Type II errors. Several comprehensive tools exist for  
84 processing such data including mothur (Schloss *et al.*, 2009), QIIME (Caporaso *et al.*, 2010a),  
85 and UPARSE (Edgar, 2013). Many stand-alone tools are also available for performing specific  
86 bioinformatic tasks which may or may not be implemented in QIIME, mothur or UPARSE. It  
87 may be beneficial in some cases to perform separate bioinformatic steps with different software  
88 packages in order to obtain the most accurate community representation for a given ecosystem.  
89 However, it is up to the individual researcher to have a comprehensive understanding of the  
90 production and processing of amplicon sequencing data in order to make the best decisions  
91 during data processing.

92         Automated quality filtering is among the first steps performed in any sequencing project  
93 and is a necessity for managing modern DNA sequencing data sets. To achieve the status of

94 “finished,” genome sequencing projects require consensus base quality scores where the  
95 likelihood of an incorrect base call is less than 1 in 100,000 (q50), whereas assemblies using  
96 unfiltered data are considered “standard draft” and are expected to contain errors (Chain &  
97 Grafham, 2009). The default parameters in QIIME 1.9.1 require a minimum quality score of q4  
98 as recommended by Bokulich *et al.* (2013), and such data should be similarly treated as “draft”  
99 data. More reads are retained for downstream analysis, but a low quality score requirement also  
100 introduces an unknown degree of sequencing error as base quality scores may vary widely across  
101 a single sequencing run. Thus, data generated on runs with higher average error rates are more  
102 likely to overestimate alpha diversity if quality scores are not strictly controlled. While  
103 inconsistent qualities from sequencing runs can be effectively controlled via quality filtering,  
104 default quality filtering in QIIME retains reads that may be variably trimmed to a range of 75-  
105 100% of the original sequence length. Because the quality of different sequences may decrease  
106 nonuniformly across a sequencing run, variable read lengths may also contribute to an inflated  
107 estimate of OTU richness if reads are not dereplicated or sorted by size prior to clustering.

108       Quality-filtered amplicon sequencing data are clustered into OTU definitions, a  
109 computational process for which numerous programs are available. CD-HIT (Fu *et al.*, 2012),  
110 UCLUST (Edgar, 2010), BLAST (Altschul, 1990), and Swarm (Mahé *et al.*, 2014) are popular  
111 options that are all available in QIIME. Reference-based analysis techniques, such as BLAST,  
112 are known to incur biases according to the choice of reference database (Nelson *et al.*, 2014), but  
113 can easily be parallelized for more efficient computation. UCLUST can utilize a reference  
114 database, perform database-independent *de novo* clustering, or, as with the open-reference  
115 strategy currently implemented in QIIME, a combination of both methods (Navas-Molina *et al.*,  
116 2013). Pure *de novo* analysis is preferred by many as the approach least likely to impose a bias  
117 on the final outcome. One popular option for *de novo* OTU clustering is CD-HIT, but as this  
118 program cannot be parallelized it can be time-prohibitive when used with larger data sets.  
119 Swarm, another *de novo* OTU clustering program, allows for portions of the *de novo* clustering  
120 process to be parallelized, thus eliminating database-specific effects while also optimizing  
121 computational requirements. All OTU picking programs require the researcher to choose a  
122 similarity or distance threshold beyond which two sequences must be considered as separate  
123 OTUs. If present at this stage, PCR or sequencing errors may contribute to OTU inflation to an  
124 unknown degree. In addition to ensuring the data is properly filtered, one can also utilize a

125 conservative clustering threshold in order to avoid overestimation of community diversity (*e.g.*,  
126  $\leq 97\%$ ; Kunin *et al.*, 2010).

127 Taxonomic assignment, achieved through comparison of OTU definition sequences to a  
128 reference database, can also be performed in a variety of ways. Popular methods include  
129 BLAST, UCLUST, and RDP (Wang *et al.*, 2007), and each are available in QIIME. In 2008, Liu  
130 *et al.* reported that RDP provided the most accurate taxonomic assignments. Presently, other  
131 techniques continue to be utilized by various amplicon sequencing analysis pipelines (*e.g.*,  
132 Giongo *et al.*, 2010; Gweon *et al.*, 2015), revealing a lack of consensus among researchers.  
133 Considering that improved taxonomic accuracies may be observed when sequences obtained for  
134 study organisms are more similar to those populating the reference database, it seems plausible  
135 that the relative success of each algorithm can be context-dependent. For environmental data  
136 sets, accuracies of taxonomic assignments are estimated by means of a confidence value relevant  
137 to the utilized technique (*e.g.*, e-value for BLAST). Careful assessment of taxonomic accuracies  
138 can only be done when the sequence content of a given sample can be anticipated. This can be  
139 achieved with synthetic mock communities created *in silico* by extracting sequences from a  
140 database (*e.g.*, Bellemain *et al.*, 2010) or using genomic mock communities that combine DNA  
141 extracts from cultured organisms. Neither scenario is likely to provide an outcome that is directly  
142 comparable to the natural complexities of environmental communities, yet both can offer a test  
143 of accuracy for taxonomic assignment methods.

144 Once quality filtered sequences have been clustered and taxonomically classified, they  
145 are compiled into an OTU table with count data for each observation. As OTUs defined from  
146 erroneous sequences may persist even to this point in the analysis, the resulting OTU table must  
147 be filtered prior to conducting diversity analyses, and the filtering approach can have a profound  
148 effect on the final result (Bokulich *et al.*, 2013). Although Bokulich *et al.* (2013) suggested the  
149 inclusion of mock communities on sequencing runs to assess the overall run quality and improve  
150 diversity assessments, they also provide a general recommendation to quality filter the final table  
151 by removing OTUs that represent less than 0.005% of the total read abundance. This has proven  
152 to be a useful guideline for numerous studies in which mock communities were not included.  
153 However, this practice ignores the independence of each sample and will treat samples  
154 differently according to sequencing depth such that low read count samples will be more  
155 severely filtered than samples with higher read counts.

156           Considering samples independently, Kircher, Sawyer & Meyer (2012) observed an  
157 indexing inaccuracy rate of 0.3%, citing cluster mixing during sequencing as a mechanism by  
158 which single-indexed Illumina sequences are likely attributed incorrectly to a particular sample.  
159 For certain applications, their result argues that such data must be filtered at this level in order to  
160 avoid Type II errors. Another common practice is to remove singleton OTUs (by sample or by  
161 table) under the assumption that such OTUs represent errors generated during sequencing (see  
162 Dickie, 2010). However, errors introduced during early PCR cycles may be faithfully replicated  
163 many times so as to appear as valid OTUs, causing overestimation of OTU richness even after  
164 singleton filtering (Nguyen *et al.*, 2015). As an alternative, Nguyen *et al.* (2015) suggest the  
165 removal of low-count or low-proportion OTUs by sample at a threshold informed by mock  
166 community data. Mock communities used in this way may also identify certain sequence motifs  
167 prone to error, which may help to identify whether novel OTUs observed in environmental data  
168 should be considered suspect. Unfortunately, such controls are not available for many data sets  
169 and artificial communities may not perform similarly to environmental communities during  
170 sample prep and analysis. Because samples are amplified independently, PCR errors are likely to  
171 be present in the form of private OTUs observed only in a single sample, so removal of unshared  
172 OTUs may be another effective precaution against overestimation of diversity due to sequencing  
173 error.

174           As these examples illustrate, proper filtering of an OTU table is not a straightforward  
175 task. The sequence misattribution rate reported by Kircher, Sawyer & Meyer (2012) is orders of  
176 magnitude above the filtering threshold of 0.005% recommended by Bokulich *et al.* (2013),  
177 though their recommendation was to filter across the entire OTU table. Since many amplicon  
178 sequencing studies report relatively few taxa present above 0.3% per sample, filtering by sample  
179 at this threshold (Kircher threshold) will exclude many valid taxa. The presence of misattributed  
180 sequences may also diminish the efficacy of private OTU removal to eliminate PCR errors,  
181 though dual-indexing of samples should reduce or eliminate sequence misattribution events  
182 (Kircher, Sawyer & Meyer, 2012). Singleton filtering, however applied, is unlikely to be  
183 thorough enough to remove errors that are either replicated during the PCR process, or represent  
184 systematic errors from the sequencing process. For single- or dual-indexed Illumina data,  
185 filtering at 0.005% across the entire table (Bokulich threshold) may represent a viable

186 compromise between confident assignment of sequences to samples and the stringency that one  
187 imposes on filtering the final table.

188 In this study we used simple genomic mock communities and an environmental data set  
189 to describe the effects of parameter adjustments for methods implemented in QIIME 1.9.1  
190 (Caporaso *et al.*, 2010a) on sequence quality filtering, OTU picking, taxonomic assignment, and  
191 OTU table filtering. We hypothesized that observed OTU diversity is dramatically inflated due to  
192 the presence of PCR and/or sequencing artifacts, and that such effects should be observed in  
193 simple genomic mock communities. Using five mock communities consisting of 4-8 taxa each,  
194 we developed a modified protocol for the analysis of 16S community amplicon sequencing data,  
195 and demonstrate the method on an environmental data set. By carefully controlling each of the  
196 steps that we investigated, we were able to describe mock community compositions more  
197 correctly than with a default workflow.

198

## 199 **Materials and Methods:**

200

### 201 *Mock communities*

202 DNA was extracted from axenic cultures of *Pseudomonas aeruginosa* (Proteobacteria),  
203 *Proteus vulgaris* (Proteobacteria), *Klebsiella pneumoniae* (Proteobacteria), *Escherichia coli*  
204 (Proteobacteria), *Bacillus megaterium* (Firmicutes), *Lactococcus lactis* (Firmicutes),  
205 *Staphylococcus aureus* (Firmicutes), and *Micrococcus luteus* (Actinobacteria) using a PowerSoil  
206 DNA Extraction Kit (MoBio Laboratories, Carlsbad, CA). DNA was quantified by PicoGreen  
207 (Life Technologies, Carlsbad, CA) fluorescence, and normalized to approximately 0.75 ng/μL.  
208 Five mock communities containing different ratios of bacterial taxa were constructed from the  
209 extracted DNA. Community 0 contained equal volumes of DNA from each taxon; Community  
210 1a contained 8% *M. luteus*, 42% *B. megaterium*, 42% *L. lactis*, and 8% *S. aureus*; Community 1b  
211 contained 42% *M. luteus*, 8% *B. megaterium*, 8% *L. lactis*, and 42% *S. aureus*; Community 2a  
212 contained 8% *E. coli*, 8% *K. pneumoniae*, 42% *P. vulgaris*, and 42% *P. aeruginosa*; Community  
213 2b contained 42% *E. coli*, 42% *K. pneumoniae*, 8% *P. vulgaris*, and 8% *P. aeruginosa*. Final  
214 concentrations for each mock community were determined to be ~ 0.75 ng/μL (Table S1: mock  
215 community construction). Expected compositions of mock communities were corrected for  
216 genome size and copy number against the CBS Genome Atlas Database (Hallin & Ussery, 2004).



217

218 *Environmental samples*

219 Environmental samples with an expected environmental contrast were collected from the  
220 Northern Arizona University Pinyon Pine Common Garden near Sunset Crater National  
221 Monument, AZ. During garden installation in October 2009, soil samples were collected from  
222 holes dug to plant seedlings (“pre-tree” treatment). Soil core samples were taken from the same  
223 seedlings in December 2010 (“post-tree” treatment). The top 2 centimeters (cm) of soil were  
224 brushed aside prior to taking cores. A 2.5 cm diameter metal corer was placed 2 cm from the  
225 seedling base and driven to a depth of 10 cm. Samples were kept on ice in the field and stored at  
226 -20 °C until DNA extraction. DNA was extracted from homogenized soil cores using a  
227 PowerSoil DNA Extraction Kit. Only samples which produced a clean ribosomal PCR product  
228 were included in this study, resulting in unequal sample sizes between pre-tree (n = 13) and post-  
229 tree (n = 28) groups. A random number generator was used to select a subset of post-tree samples  
230 (n = 13) for comparisons of data with equal sample sizes. Samples were normalized to c. 1 ng/μL  
231 prior to PCR amplification for library construction.

232 The environmental samples presented here are meant only to allow a demonstration of  
233 the effects of a mock community-based workflow optimization on actual data. Though we expect  
234 the presence of a tree to create additional niche space which would increase observed diversity,  
235 no background soil control samples were collected. Observed differences, though likely to be  
236 real, could be influenced in part or in total by interannual environmental variations. Additionally,  
237 pre-tree and post-tree samples were collected during different months of the year, so seasonal  
238 differences could also contribute to the outcome.

239

240 *Library construction and sequencing*

241 Amplicons were produced in a two-step protocol as suggested by Berry *et al.* (2011).  
242 Briefly, samples were amplified in triplicate PCR reactions for the 16S v4 region using the  
243 universal prokaryotic primers 515F and 806R (Bates *et al.*, 2011). First round reactions were  
244 performed in triplicate in 384 well plates. The 8 μL volumes contained the following: 1 μM each  
245 primer (Eurofins MWG Operon, LLC), 200 μM each dNTP (Phenix Research, Candler, NC),  
246 0.01 U/μL Phusion Hot Start II DNA Polymerase (Life Technologies), 1X HF Phusion Buffer  
247 (Life Technologies), 3 mM MgCl<sub>2</sub>, 6% glycerol, and 1 μL normalized template DNA. Cycling



248 conditions were: 2 minutes at 95°C followed by 20 cycles of 30 seconds at 95°C, 30 seconds at  
249 55°C, 4 minutes at 60°C. Triplicate reactions for each sample were pooled by combining 4 µL  
250 from each, and 2 µL was used to check for results on a 1% agarose gel. The remainder was  
251 diluted 10-fold and used as template in a second PCR reaction in which 12 base Golay indexed  
252 tails (Caporaso *et al.*, 2012) were added. Second round reaction conditions were identical to the  
253 first round except only one reaction was conducted per sample and only 15 total cycles were  
254 performed. Indexed PCR products were purified using carboxylated magnetic beads as described  
255 in Rohland & Reich (2012), quantified by PicoGreen fluorescence, and an equal mass of each  
256 sample was combined into a final sample pool. The pool was purified and concentrated, and  
257 subsequently quantified by quantitative PCR against Illumina DNA Standards (Kapa  
258 Biosystems, Wilmington, MA). Sequencing was carried out on a MiSeq Desktop Sequencer  
259 (Illumina Inc, San Diego, CA) running in paired end 2x150 mode.

260

#### 261 *Data processing and statistical analysis*

262 All bioinformatics were carried out on a Mac Pro (Apple, Inc.) running Ubuntu Linux  
263 14.04 LTS (Canonical Ltd.) or the monsoon high-performance computing cluster at Northern  
264 Arizona University (<https://nau.edu/hpc/>) running CentOS 6.6 (The CentOS Project). Figures  
265 were generated in Veusz v1.24 (<http://home.gna.org/veusz/>) or Geneious v8.1 (Biomatters Ltd.).  
266 As contaminating PhiX Control sequence can complicate sequencing projects (Mukherjee *et al.*,  
267 2015), we calculated the amount of PhiX Control among our demultiplexed data and removed it  
268 prior to sample processing. This task was performed with the `akutils phix_filtering`  
269 command in `akutils` v1.2 (Krohn, 2016; <https://github.com/alk224/akutils-v1.2>) which maps raw  
270 data against the Enterobacteria phage phiX174 sensu lato complete genome sequence  
271 (NC\_001422.1) using Smalt 0.7.6 (<http://www.sanger.ac.uk/resources/software/smalt/>).

272 Paired end reads were joined using the `akutils join_paired_reads` command in  
273 `akutils` which employs `fastq-join` from `ea-utils` (Aronesty, 2011). Demultiplexing and quality  
274 filtering of raw, joined data (mean length = 253 bp) was carried out in QIIME with the  
275 `split_libraries_fastq.py` script using default parameters, or with more strict  
276 requirements of a minimum quality threshold of q20 ( $q = 19$ ), allowing 0-3 low-quality base  
277 calls ( $r = 1-3$ ), and requiring at least 95% of each read to be high quality ( $p = 0.95$ ). Chimeras  
278 were removed using `vsearch` 1.1.1 (Rognes *et al.*, 2015; <https://github.com/torognes/vsearch>)

279 against the Gold reference database (<http://drive5.com/uchime/gold.fa>). OTU picking and  
280 taxonomy assignments were performed using the `akutils pick_otus` command in `akutils`.  
281 After manual inspection of sequence divergence among congeneric mock community members,  
282 sequences were de-replicated on the first 100 bases using the `prefix_suffix` OTU picker in  
283 QIIME. OTU picking was performed with multiple similarity thresholds using common OTU  
284 picking algorithms (CD-HIT, UCLUST and BLAST at 97%, 95%, 92%, 90%, 85%, and Swarm  
285 at  $d1$ ,  $d2$ ,  $d3$ ,  $d4$ ,  $d5$ ). BLAST was used only for closed reference analysis, UCLUST for open  
286 reference analysis, and CD-HIT and Swarm for *de novo* analyses. Taxonomy was assigned using  
287 BLAST, RDP, and UCLUST using default settings available in QIIME 1.9.1. Reference-based  
288 OTU picking steps and taxonomy assignments were conducted against the Greengenes 97%  
289 database (McDonald *et al.*, 2012) which had been formatted to include only the v4 region using  
290 the `akutils format_database` command in `akutils`. Sequence alignments and phylogenetic  
291 trees were produced using the `akutils align_and_tree` command in `akutils` which aligns  
292 sequences using PyNAST (Caporaso *et al.*, 2010b) and generates phylogenies with FastTree  
293 (Price, Dehal & Arkin, 2009). Diversity analyses were conducted using the `akutils`  
294 `core_diversity` command in `akutils`.

295 In order to facilitate assessment of optimal workflow steps, we first sought to establish a  
296 method of filtering the final OTU tables by eliminating OTUs resulting from mixed clusters. To  
297 this end, we processed the mock and environmental data sets through a default QIIME workflow  
298 (see below) to assess taxonomic components, and compared methods for filtering OTU tables,  
299 which would remove contaminating taxa. An ideal filtering method should remove erroneous  
300 OTUs that arise either from sequencing error or cluster mixing. Table filtering was carried out  
301 using either the Kircher threshold (0.3% by sample; Kircher, Sawyer & Meyer, 2012), the  
302 Bokulich threshold (0.005% by table; Bokulich *et al.*, 2013), singletons removed by table (`mc2`),  
303 or singletons removed by sample (`n2`). Private OTUs were assumed to be errors and were also  
304 removed in the `n2` tables. Filtered OTU tables were grouped according to filtering method, and  
305 differences in the amount of OTUs classified as contaminating taxa was assessed by one-way  
306 ANOVA. Tukey's HSD test was used to determine which groups were statistically distinct.

307 An optimal workflow was chosen by assessing diversity estimates and taxonomic  
308 identities assigned to mock community data. The optimal OTU picking algorithm was  
309 determined as the method that yielded the correct diversity result over the broadest range of

310 similarity thresholds. Taxonomic accuracy was determined by seeding the Greengenes database  
311 with the expected sequences from the mock community constituent taxa prior to analysis, and  
312 inspecting the results. OTU tables from the optimal workflow across the accurate range of  
313 similarity thresholds were filtered at each of the four thresholds described above. Our “default  
314 QIIME workflow” was identical to the optimal workflow with the following changes: the  
315 `split_libraries_fastq.py` command was performed with default settings; OTU picking  
316 was performed with the `pick_open_reference_otus.py` command; taxonomic assignment  
317 was performed with UCLUST; OTU tables were filtered with the Bokulich threshold. Results  
318 from the optimal workflow were compared to the result obtained from our default workflow.  
319 Environmental data was then processed using the best workflow determined from this process  
320 and compared to the default result.

321         Diversity analyses for mock community data were calculated on OTU tables that had  
322 been rarefied to 10,000 reads, or 5,000 reads for environmental data. Comparison of observed  
323 mock community composition to the *a priori* expectation (Table S1) was conducted with  
324 Spearman’s rank correlation using species-level assignments. Comparison of observed OTU  
325 diversity between environmental sample groupings was performed with nonparametric t-tests. A  
326 random subset of post-tree samples from the environmental data ( $n = 13$ ) was selected to  
327 determine if unequal sample sizes were contributing to observed OTU diversity. Distance  
328 matrices were calculated from environmental data for weighted UniFrac distance (Lozupone &  
329 Knight, 2005). Tests of differences of total beta diversity were carried out on distance matrices  
330 using PERMANOVA (Anderson, 2001), and differences in multivariate dispersion were detected  
331 with PERMDISP (Anderson, Ellingsen & McArdle, 2006).

332         Representative sequences for the optimized mock community result were extracted from  
333 the output data. When multiple OTU definition sequences represented the same taxonomic  
334 identity, they were aligned with Mafft v7.123b (Kato & Standley, 2013) using the L-INS-i  
335 setting. The lower abundance OTU for each multi-OTU taxon was assumed to be erroneous and  
336 base differences compared to the major OTU were characterized. Trinucleotide motifs preceding  
337 each base difference and terminal truncation position were tabulated. Because 2x150 sequencing  
338 data does not fully overlap for 515F-806R amplicons (mean length = 253 bp), terminal base and  
339 preceding trimers were considered in the context of the second read. Environmental data  
340 processed through the optimal workflow was also investigated for terminal truncation positions

341 and preceding trinucleotide motifs. Because we have no reliable reference sequence for many  
342 environmental OTUs, we investigated only OTUs that shared a taxonomic designation with at  
343 least one other OTU, and had been truncated by more than 3 bases during quality filtering. For  
344 mock and environmental data, motif and terminal base representations were tested against the  
345 assumption of random occurrence with Chi-square tests.

346

#### 347 **Results:**

348         The sequencing run clustered at 1119 k/mm<sup>2</sup> (+/- 70) and resulted in 17.96 million total  
349 reads passing filter, an overall error rate of 0.36%, and 91% of reads exceeded q30. PhiX  
350 Sequencing Control v3 sequences (Illumina, Inc.) constituted 8.31% of the total run (percent  
351 aligned). Once demultiplexed, mock community data contained 4.35% PhiX (103,070/2,371,510  
352 reads) while the environmental data contained 4.10% PhiX (259,366/6,332,586 reads). Raw  
353 sequencing data for the samples used in this study and a QIIME-formatted mapping file are  
354 publicly available in the QIITA database (<https://qiita.ucsd.edu/>) under study ID number 10479.

355         Under default QIIME assessment, the mock community data showed substantial OTU  
356 inflation; where there should have been just 8 OTUs, there were 127 (Table S2, default mock  
357 analysis). When the environmental data set was processed through the same workflow, 73 OTUs  
358 were classified at the family level as Sphingomonadaceae. Together, these OTUs made up 5.3%  
359 of environmental sequences, and Sphingomonadaceae was the most abundant classification  
360 observed at the family level (Table S3, default environmental analysis). Three OTUs  
361 representing about 0.13% of the mock community data set were also classified as  
362 Sphingomonadaceae, a designation which should be absent from the mock data. This result led  
363 us to surmise that sequences from the environmental data set were contaminating the mock  
364 communities during sequencing. Such sample cross-talk presumably arises from the cluster  
365 mixing effect described by Kircher, Sawyer & Meyer (2012) where the index read from a  
366 flowcell cluster is spuriously attributed to a neighboring cluster. The mock data also contained 3  
367 OTUs classified as Planococcaceae (<0.03%) and 1 OTU classified as Methylobacteriaceae  
368 (<0.01%), again corresponding with OTUs observed within the environmental data.  
369 Sphingomonadaceae sequences were observed across all five mock communities, whereas  
370 Planococcaceae was only associated with communities 0, 1a, and 1b, suggesting that cluster-  
371 mixing events may occur non-randomly. Methylobacteriaceae was present as just a single read

372 among community 1b. Three mock community OTUs were observed at low levels in  
373 communities from which they should be absent, indicating additional cluster-mixing within the  
374 mock community data.

375 As the most prevalent non-target taxon observed among the mock community data, we  
376 sought to establish a method for filtering OTU tables that would eliminate the presence of  
377 contaminating Sphingomonadaceae reads. OTU tables generated for the mock communities by  
378 each of the OTU picking, taxonomy assignment and table filtering methods were compared for  
379 the presence of Sphingomonadaceae contaminants. Considering filtering method (mc2, n2,  
380 Kircher threshold, or Bokulich threshold) as the predictive variable, we found strong differences  
381 among them in removing non-target OTUs ( $F_{3,239} = 89.301$ ,  $p < 0.0001$ ). The least severe  
382 filtering method (mc2) retained the most Sphingomonadaceae OTUs (2.50 +/- 1.21) followed by  
383 n2 (2.45 +/- 1.21), and Bokulich threshold (1.85 +/- 0.86). Only the Kircher threshold completely  
384 removed Sphingomonadaceae contamination from the mock community OTU tables effectively  
385 (Tukey's HSD,  $p < 0.05$ ).

386 Default quality filtering and OTU picking in QIIME resulted in overestimation of mock  
387 community diversity regardless of how the final OTU table was filtered (Figure 1a-d; Figure S1:  
388 Default mock community rarefactions). Diversity estimates were inflated up to 35 times when  
389 singletons were removed by table, compared to nearly 3.5 times when filtering with the Kircher  
390 threshold. Despite the reduction of OTU inflation by an order of magnitude, these results indicate  
391 that revisions to initial processing steps may yield improved results. We therefore sought to  
392 establish an optimized workflow that would produce the correct number of OTUs for an input of  
393 known constituents. Using data that had been filtered according to strict standards during the  
394 `split_libraries_fastq.py` step in QIIME ( $q = 19$ ,  $r = 0$ ,  $p = 0.95$ ), a correct result was  
395 achieved for each of the OTU picking algorithms tested. However, each algorithm differed in  
396 which similarity threshold was required for the optimal result (Table 1). Closed reference picking  
397 with BLAST overestimated diversity above a similarity threshold of 92%. Open reference  
398 picking with UCLUST overestimated diversity at every threshold except 95% similarity. *De*  
399 *novo* picking using CD-HIT at thresholds below 92% and Swarm resolutions below  $d4$   
400 underestimated diversity. Swarm yielded the correct result over the broadest range of tested  
401 similarity thresholds ( $d1$ - $d4$ ), and offers other attractive features that made it stand out among the  
402 tested OTU pickers (*e.g.*, *de novo* picking, multi-threaded analysis). Thus, Swarm was chosen as

403 the optimal OTU picking method for the remainder of the study. We chose *d4* similarity as the  
404 optimal threshold as it was the most conservative setting to yield a correct result.

405 Taxonomic accuracy for Swarm-picked OTUs (*d4*) was assessed for the different  
406 taxonomy assigners using default parameters in QIIME 1.9.1. To control for reference database  
407 bias, we added representative sequences from each of the correct OTUs to our Greengenes  
408 reference with a unique identifier. We observed that BLAST returned the representative  
409 sequence 100% of the time, while RDP and UCLUST never found the exact match (Table 2).  
410 Even though RDP and UCLUST did not find optimal sequences, assignments were correct,  
411 though less specific in taxonomic depth. BLAST yielded similar results when the representative  
412 sequences were not present in the database (Table 2). While BLAST offers the advantage of  
413 obtaining the best sequence match when available in the database, RDP and UCLUST both offer  
414 an advantage in substantially reducing computational time while providing reasonable accuracy  
415 for most applications. For the analysis presented here, we chose BLAST as the optimal  
416 taxonomy assigner for its superior accuracy.

417 A perfect result for analysis of our mock communities requires stringent quality filtering  
418 of the raw data. Default quality filtering in QIIME 1.9.1 was established according to Bokulich *et*  
419 *al.* (2013). This imposes a minimum Phred quality score of 4 ( $q = 3$ ), truncates sequences after  
420 three bases are observed below this threshold ( $r = 3$ ), and retains truncated reads that represent a  
421 minimum of 75% of the original sequence length ( $p = 0.75$ ). In contrast, we performed strict  
422 quality filtering using  $q = 19$ ,  $r = 0$ , and  $p = 0.95$ . This more stringent filtering protocol ensures  
423 that data used for analysis are of much higher quality with approximately uniform read lengths.  
424 An important consequence of such stringency is that much of the raw data is discarded. Of the  
425 2,373,247 raw mock community sequences, default quality filtering retained 2,020,542 reads  
426 (85.1%), whereas stringent parameters retained just 657,544 reads (27.7%). Holding constant  $q =$   
427 19 and  $p = 0.95$ , we found that increasing  $r$  during quality filtering had a profound effect on the  
428 amount of data retained (Figure 2a). Allowing  $r = 1$  resulted in an increase of data retention from  
429 approximately 27% ( $r = 0$ ) to over 56%. When  $r = 2$  and  $r = 3$ , increases in data retention  
430 showed diminishing returns, with 70% and 75% of the data retained, respectively. However, we  
431 also found that allowing  $r > 0$  will generally cause an inaccurate estimate of the number of  
432 OTUs, depending on the criteria used for OTU picking (Figure 2b). With the *d1* resolution,  
433 increasing  $r$  will create a proportional inflation in the number of OTUs determined by Swarm. At



434  $d2$  resolution, allowing  $r = 1$  still correctly described our simple mock community whereas  
435 allowing  $r = 2$  or  $r = 3$  caused diversity to be overestimated. At resolutions  $d3$  and  $d4$ , allowing  $r$   
436  $> 0$  caused underestimates of diversity. This suggests that the best result is obtained with the  
437 most stringent quality filter, which we selected for our optimal workflow ( $q = 19$ ). Similar results  
438 using more data may be possible by allowing a small amount of errors (*e.g.*,  $r = 1$ ) and picking  
439 OTUs with a more conservative similarity threshold (*e.g.*, Swarm at  $d2$  resolution).

440 The Kircher threshold was effective at removing contaminating OTUs in our mock  
441 community data thus yielding a near-perfect result (Figure 3). However, we anticipated that such  
442 filtering could be too stringent for environmental analysis given the low per-sample OTU  
443 frequencies commonly reported (*e.g.*, Sogin *et al.*, 2006). We compared the expected mock  
444 community results to those observed with either default settings, or optimized settings for quality  
445 filtering, OTU picking and taxonomy assignment, using each of the final OTU table filtering  
446 methods we tested. For all comparisons, Spearman's rank correlation yielded significant p-values  
447 ( $< 0.001$ ), so we present only correlation values and 95% confidence intervals here. When  
448 comparing the default analysis to the expected outcome, Spearman's  $r$  showed a negative  
449 correlation ( $r = -0.3494$ ; CI =  $[-0.4280, -0.2655]$ ). Optimized results exhibited strong positive  
450 correlations regardless of filtering threshold used. Lower values for Spearman's  $r$  occurred when  
451 diversity was overestimated and when contaminants were present. Correlation with the expected  
452 outcome improved as filtering stringency increased with every filtering method producing a  
453 dramatic improvement over the default workflow (mc2:  $r = 0.8075$ , CI =  $[0.7663, 0.8420]$ ; n2:  $r$   
454  $= 0.8702$ , CI =  $[0.8344, 0.8987]$ ; Bokulich threshold:  $r = 0.9135$ , CI =  $[0.8841, 0.9357]$ ; Kircher  
455 threshold:  $r = 0.9646$ , CI =  $[0.9495, 0.9752]$ ). The Bokulich threshold was chosen as our optimal  
456 OTU table filtering method because it yielded the best correlation without being overly strict.

457 Output for the environmental data using either the default or optimized workflow were  
458 examined for basic diversity statistics. Default analysis identified 2508 OTUs classified into 388  
459 taxonomic assignments (OTUs per taxon: mean = 6.46, median = 2; Figure S3: Default  
460 environmental rarefactions). The optimized analysis identified 1533 OTUs classified into 328  
461 taxonomic assignments (OTUs per taxon: mean = 4.67, median = 2; Figure S4: Optimized  
462 environmental rarefactions). By treatment, OTU diversity was reduced about twofold when  
463 assessed via the optimized workflow and compared to the default results (Figure 4a-4b). In the  
464 default analysis, pre-tree soils hosted 978.30 +/- 128.42 OTUs while post-tree soils had 1138.35



465 +/- 86.34 OTUs (nonparametric T-test = 4.578,  $p < 0.001$ ). In the optimized analysis, pre-tree  
466 soils contained 543.28 +/- 79.32 compared to 674.95 +/- 50.21 OTUs in post-tree soils  
467 (nonparametric T-test = 6.277,  $p < 0.001$ ). Differences in beta diversity were observed between  
468 treatments for each workflow using weighted UniFrac distance matrices (Figure 4c-4d; default  
469 PERMANOVA = 8.181,  $p < 0.001$ ; optimized PERMANOVA = 9.355,  $p < 0.001$ ). We also  
470 noticed an increase in multivariate dispersion in the optimized workflow, though the differences  
471 were not found to be significant in either case (default PERMDISP = 1.086,  $p = 0.294$ ; optimized  
472 PERMDISP = 2.160,  $p = 0.158$ ). When data was processed with equivalent sample sizes, the  
473 same patterns were observed for both alpha diversity (pre-tree = 545.76 +/- 78.84, post-tree =  
474 679.00 +/- 47.86; nonparametric T-test = 5.004,  $p < 0.001$ ) and beta diversity (PERMANOVA =  
475 6.585,  $p < 0.001$ ), though statistical power was slightly reduced, and multivariate dispersion  
476 increased (PERMDISP = 3.248,  $p = 0.071$ ), consistent with a reduction in sample size.

477         Of the 17 OTUs observed in the optimized mock result, the nine extra OTUs therein were  
478 composed of three contaminants and six spurious OTUs representing sequence variants of the  
479 target taxa. All extra OTUs were present at low levels ranging from 0.003% to 0.17% per sample  
480 (Table S4). That sequence counts of contaminant OTUs were observed in all samples, but only  
481 for select taxa, strongly suggests that cluster mixing events occur non-randomly during Illumina  
482 sequencing. Species-level mock community observations from the optimized workflow describe  
483 the eight constituent taxa at approximately the correct proportions. However, six of the eight taxa  
484 were represented by two OTUs each. The main OTU for each taxon was present as 6.30% to  
485 19.09% of the total community while the rates of lower frequency OTUs ranged from 0.01% to  
486 0.05%. Manual inspection of conspecific OTU sequence alignments revealed multiple  
487 substitution and indel positions within the first 100 bases which prevented these sequences from  
488 dereplicating into the correct sequence during our workflow (Table S5: OTU sequence  
489 alignments by taxonomy). Additionally, these sequence variants were shorter than the main  
490 constituent sequence by at least seven bases, indicating that they derive from inherently lower  
491 quality reads. Inspection of trinucleotide motifs preceding each substitution or indel position did  
492 not reveal any pattern relating to the observed errors (Table S6: Error-associated sequence  
493 motifs). Consistent with the results of Schirmer *et al.* (2015), we observed a higher rate of errors  
494 among A or C bases than G or T errors (error ratio = 1.67). Since A and C or G and T bases  
495 share fluorescence excitation wavelengths during Illumina 4-channel sequencing-by-synthesis

496 (SBS), this result suggests that some of the errors we observed were indeed the result of  
497 systematic errors during sequencing, although this study was not designed to distinguish between  
498 such errors and those generated during PCR. Examining the terminal trinucleotide motif  
499 immediately preceding truncation positions (Table S7: Terminal-associated sequence motifs) we  
500 observed “TTT” 83% of the time ( $X^2_{63} = 271.333$ ,  $p < 0.0001$ ). Additionally, the correct base at  
501 the truncation position was “G” 83% of the time ( $X^2_3 = 11.33$ ,  $p = 0.0101$ ). An example  
502 alignment for the two OTUs representing *B. megaterium* is presented in Figure 5a, illustrating  
503 the “TTT” motif preceding a “G” truncation position (reverse complimented).

504 Truncation positions and preceding trimers were also characterized for environmental  
505 data, resulting in 34 “suspect” OTUs (Table S8: Environmental terminal errors). Of these, 27  
506 OTUs had been truncated at a “G” position (79.41%;  $X^2_3 = 54.235$ ,  $p < 0.0001$ ), and just 10  
507 possible trimers were represented preceding the truncation position. The motifs “TTT” and  
508 “TTC” were substantially overrepresented, being observed 14 (41.18%) and 7 (20.59%) times,  
509 respectively ( $X^2_{63} = 474.235$ ,  $p < 0.0001$ ). An example alignment for 5 OTUs classified to the  
510 family level as Sphingomonadaceae is presented in Figure 5b, and includes one such suspect  
511 OTU with a “TTT” motif preceding a “G” truncation position (reverse complimented).

512

### 513 **Discussion:**

514 We have shown that amplicon sequencing data from Illumina MiSeq instruments should  
515 be stringently filtered in order to provide the most accurate estimates of diversity. Kunin *et al.*  
516 (2010) found that diversity was grossly overestimated for their mock community data until a  
517 quality threshold of q27 was implemented. Similarly, Nelson *et al.* (2014) observed high  
518 overestimation of mock community diversity (25-125 times expected) unless the data was  
519 carefully controlled. Our optimal workflow still overestimated the OTU diversity of our simple  
520 mock communities by a factor of about 2. This slight overestimation is a dramatic improvement  
521 over that obtained by default processing, and our optimized protocol yielded a reasonable  
522 characterization of taxonomic content for mock communities (Table S4) and environmental data  
523 (Table S9) alike.

524 Schirmer *et al.* (2015) observed that error rates as reported by the sequencer according to  
525 the PhiX Control v3 do not accurately reflect those of amplicon sequences. Their conclusion that  
526 actual error rates were higher than those indicated by q-scores reported by the MiSeq has

527 important implications for the use of Illumina sequencing in estimating microbial diversity. It is  
528 possible that newer imaging strategies (*e.g.*, 2-channel SBS chemistry used by Illumina NextSeq  
529 and MiniSeq instruments) will provide improved parity between the estimated and actual error  
530 rates, but this will require careful testing to determine empirically. Of the non-target OTUs  
531 present in our optimized mock community result, one third were contaminants arising from  
532 cluster mixing events during sequencing and two thirds were sequence variants of the constituent  
533 OTUs which may have arisen during PCR, sequencing, or a combination of the two. Cluster  
534 mixing can be controlled by dual-indexing of samples (Kircher, Sawyer & Meyer, 2012), but  
535 errors arising during PCR or sequencing represent systematic errors inherent to the procedure of  
536 amplicon sequencing which are difficult, if not impossible, to completely eliminate irrespective  
537 of indexing strategy. Even though dual-indexing offers a clear advantage over single indexing  
538 with regard to sample attribution, single-indexed protocols (*e.g.*, Caporaso *et al.*, 2012) remain  
539 popular and are widely used. Such data still yields valuable information and should not be  
540 discounted, as long as researchers are aware of the limitations. Dual-indexed designs should be  
541 encouraged for new research projects (*e.g.*, Kozich *et al.*, 2013; Fadrosch *et al.*, 2014).

542 We echo the recommendation by others (*e.g.*, Bokulich *et al.*, 2013; Schirmer *et al.*,  
543 2015) to include control mock community samples to guide data analysis. PhiX Control v3 is  
544 still needed to improve sequence diversity for the purpose of cluster map generation  
545 ([https://support.illumina.com/content/dam/illumina-](https://support.illumina.com/content/dam/illumina-marketing/documents/products/technotes/hiseq-phix-control-v3-technical-note.pdf)  
546 [marketing/documents/products/technotes/hiseq-phix-control-v3-technical-note.pdf](https://support.illumina.com/content/dam/illumina-marketing/documents/products/technotes/hiseq-phix-control-v3-technical-note.pdf)), but an  
547 alternative reference sequence could be used with onboard mock communities to more directly  
548 estimate error profiles for community amplicon sequencing data. PhiX sequence itself likely  
549 contributes little (if at all) to inflation of diversity estimates, and is easily quantified and  
550 removed. Though such an effect is direct evidence of cluster mixing, the rate of PhiX infiltration  
551 is likely much higher than the rate of sample mixing because PhiX Control is unindexed,  
552 producing no fluorescent signal during indexing cycles. Spurious OTUs defined from  
553 contaminating PhiX sequence may be more prevalent amid sequence data which was  
554 accompanied by higher concentrations of PhiX Control v3 during sequencing.

555 Though this study was not designed for careful investigation of errors generated during  
556 amplicon sequencing projects, we were able to observe that certain bases and motifs were more  
557 frequently associated with low-quality base calls than should be expected by chance. The

558 presence of a “TTT” or “TTC” motif immediately preceding a “G” position near the end of a  
559 sequence (near the start of the second read) was most frequently associated with an erroneous or  
560 suspect OTU (Table S7). Indeed, mock community diversity was inflated on account of this  
561 effect, but determining the source of such error requires more careful investigation than is  
562 possible here. In addition to the terminal truncation observations, we note that all other observed  
563 errors in the mock community sequences occurred within the first 100 bp of sequence, specific to  
564 the non-overlapping region of the first sequencing read (Figure 5a). It is likely that the errors we  
565 observed here would have occurred less frequently had we used fully-overlapping reads for this  
566 study. Importantly, the motif-specific patterns we observed were consistent between the mock  
567 and environmental data sets (Figure 5; Table S7; Table S8).

568 Estimates of alpha diversity are more sensitive than beta diversity calculations to the  
569 effects of cluster mixing and systematic errors. Input sequence quality had the most profound  
570 effect on alpha diversity estimates. Increasing the number of allowed low-quality reads (*r*  
571 parameter in `split_libraries_fastq.py`) increases the amount of data available for  
572 processing, but also changes observed diversity. For this reason, we suggest that diversity  
573 estimates should be performed only with data that has been stringently filtered for quality.  
574 Because errors in amplicon sequencing data may follow sequence-specific patterns (Schirmer *et*  
575 *al.*, 2015; this study), spurious OTUs may provide artificial support to the statistical separation of  
576 experimental treatments when derived from OTUs driving such differences. Alternatively,  
577 spurious OTUs arising from taxa which are not differentially represented among treatments  
578 could provide artificial noise, making it more difficult to detect real differences. In either  
579 scenario, careful quality filtering can diminish such effects.

580 Given the vast number of studies that have already utilized Illumina sequencing for  
581 community amplicon profiling, it seems likely that estimates of alpha diversity for a wide variety  
582 of environments could be inflated due to uncontrolled error rates. Here we observed this effect  
583 with a QIIME-based workflow, though QIIME is just one of a variety of tools used in data  
584 analysis for such work. That errors may arise systematically during PCR or sequencing implies  
585 that a similar effect is likely to be observed regardless of which analysis pipeline is used to  
586 assess the data. We also made use of a high-fidelity polymerase (Phusion Hot Start II) in contrast  
587 to many studies which continue to utilize Taq polymerase, with which PCR-derived errors will  
588 be more prevalent. Lower fidelity will promote more PCR-derived errors, and those generated

589 during early cycles will be highly perpetuated, an effect which would be more problematic under  
590 high-cycling conditions. Because errors may follow sequence-specific patterns, some diversity  
591 estimates may be particularly inflated for certain taxa, which can further affect studies using  
592 taxonomic content to predict community function (*e.g.* Langille *et al.*, 2013). The use of  
593 phylogenetic metrics (*e.g.*, phylogenetic diversity for alpha diversity, UniFrac for beta diversity)  
594 during data analysis will likely diminish the effects of complications associated with  
595 systematically-inflated OTU diversity. Though the quality-filtering recommendations outlined by  
596 Bokulich *et al.*, (2013) have subsequently provided valuable guidance to numerous researchers,  
597 careful consideration of the results presented here and elsewhere (Kunin *et al.*, 2010; Schirmer *et*  
598 *al.*, 2015) will improve upon our collective interpretation of microbial diversity across  
599 environments.

600

### 601 **Conclusions:**

602 In this study, we observed that each of the various workflow components tested (quality  
603 filtering, OTU picking, taxonomic assignment, and OTU table filtering) affect the outcome of an  
604 amplicon sequencing project. Though high quality output can be achieved through a variety of  
605 means, in this study the optimal result was achieved with a specific set of steps. We outline them  
606 here as a general recommendation for processing community amplicon data generated on MiSeq  
607 instruments through QIIME 1.9.1 (Caporaso *et al.*, 2010a). Analysis parameters can and should  
608 be adjusted as necessary for individual data sets. The optimal workflow as performed in this  
609 study was as follows (optimized steps in bold):

610

- 611 1. Remove PhiX Control v3 contamination with Smalt
- 612 2. Align read pairs with fastq-join
- 613 3. **Strict quality filter in QIIME ( $q = 19, r = 0, p = 0.95$ )**
- 614 4. Chimera filtering with vsearch
- 615 5. Sequence dereplication with prefix/suffix OTU picker
- 616 6. **Pick OTUs with Swarm ( $d4$  resolution, adjust as necessary)**
- 617 7. **Assign taxonomy with BLAST (default settings)**
- 618 8. **Filter output table at the Bokulich threshold**

619

620 Our results were consistent with the hypothesis that mock community diversity would be  
621 inflated due to the presence of PCR or sequencing errors in the data. By imposing more rigorous  
622 quality filtering of raw sequencing data, much of this error is removed. The effects of remaining  
623 errors can be minimized by utilizing a conservative similarity or distance threshold during OTU  
624 picking. By characterizing mock communities at multiple thresholds, one can identify a  
625 sufficiently conservative similarity or distance value ( $d_4$  in our case) which should offer  
626 improved confidence when measuring environmental diversity. If mock communities are  
627 unavailable, we advocate the use of a workflow based upon the above optimization. For studies  
628 utilizing an alternative locus, we suggest adjusting the clustering threshold based on the length of  
629 the amplicon (*e.g.*, more conservative clustering for longer amplicons) until mock communities  
630 can be employed to determine a more informed threshold.

631

### 632 **Acknowledgements:**

633 The authors would like to thank Linda Fitchett-Hewitt for providing axenic bacterial  
634 cultures, Dreux Patch for propagating them and extracting DNA, and the NAU Environmental  
635 Genetics and Genomics Laboratory.

636

### 637 **References:**

638

639 Altschul, S. F., Gish, W., Miller, W., Myers, E. W., & Lipman, D. J. (1990). Basic local  
640 alignment search tool. *Journal of Molecular Biology*, 215, 403–10.

641

642 Anderson, M. J. (2001). A new method for non-parametric multivariate analysis of variance.  
643 *Austral Ecology*, 26, 32-46.

644

645 Anderson, M.J., Ellingsen, K.E., McArdle, B.H. (2006). Multivariate dispersion as a measure of  
646 beta diversity. *Ecology Letters*, 9, 683–693.

647

648 Aronesty, E. (2011). *ea-utils*: Command-line tools for processing biological sequencing data;

649 <http://code.google.com/p/ea-utils>



650  
651 Bates, S. T., Berg-Lyons, D., Caporaso, J. G., Walters, W. A., Knight, R., & Fierer, N. (2011).  
652 Examining the global distribution of dominant archaeal populations in soil. *The ISME Journal*, 5,  
653 908–17.  
654  
655 Bellemain, E., Carlsen, T., Brochmann, C., Coissac, E., Taberlet, P., & Kauserud, H. (2010). ITS  
656 as an environmental DNA barcode for fungi: an in silico approach reveals potential PCR biases.  
657 *BMC Microbiology*, 10, 189.  
658  
659 Bentzon-Tilia, M., Traving, S. J., Mantikci, M., Knudsen-Leerbeck, H., Hansen, J. L., Markager,  
660 S., & Riemann, L. (2015). Significant N<sub>2</sub> fixation by heterotrophs, photoheterotrophs and  
661 heterocystous cyanobacteria in two temperate estuaries. *The ISME Journal*, 9, 273–285.  
662  
663 Berry, D., Mahfoudh, K. B., Wagner, M., & Loy, A. (2011). Barcoded Primers Used in  
664 Multiplex Amplicon Pyrosequencing Bias Amplification. *Applied and Environmental*  
665 *Microbiology*, 77, 612–612.  
666  
667 Bokulich, N. A., & Mills, D. A. (2013). Improved selection of internal transcribed spacer-  
668 specific primers enables quantitative, ultra-high-throughput profiling of fungal communities.  
669 *Applied and Environmental Microbiology*, 79(8), 2519–26.  
670  
671 Bokulich, N. A., Subramanian, S., Faith, J. J., Gevers, D., Gordon, J. I., Knight, R., Mills, D. A.,  
672 & Caporaso, J. G. (2013). Quality-filtering vastly improves diversity estimates from Illumina  
673 amplicon sequencing. *Nature Methods*, 10, 57–9.  
674  
675 Caporaso, J. G., Kuczynski, J., Stombaugh, J., Bittinger, K., Bushman, F. D., Costello, E. K.,  
676 Fierer, N., Peña, A. G., Goodrich, J. K., Gordon, J. I., Huttley, G. A., Kelley, S. T., Knights, D.,  
677 Koenig, J. E., Ley, R. E., Lozupone, C. A., McDonald, D., Muegge, B. D., Pirrung, P., Reeder,  
678 J., Sevinsky, J. R., Turnbaugh, P. J., Walters, W. A., Widmann, J., Yatsunenko, T., Zaneveld, J.,  
679 & Knight, R. (2010a). QIIME allows analysis of high-throughput community sequencing data.  
680 *Nature Methods*, 7, 335–336.



681  
682 Caporaso, J. G., Bittinger, K., Bushman, F. D., Desantis, T. Z., Andersen, G. L., & Knight, R.  
683 (2010b). PyNAST: A flexible tool for aligning sequences to a template alignment.  
684 *Bioinformatics*, 26, 266–267.

685  
686 Caporaso, J. G., Lauber, C. L., Walters, W. A., Berg-Lyons, D., Huntley, J., Fierer, N., Owens,  
687 S. M., Betley, J., Fraser, L., Bauer, M., Gormley, N., Gilbert, J. A., Smith, G., & Knight, R.  
688 (2012). Ultra-high-throughput microbial community analysis on the Illumina HiSeq and MiSeq  
689 platforms. *The ISME Journal*, 6, 1621–1624.

690  
691 Chain, P., & Grafham, D. (2009). Genome project standards in a new era of sequencing. *Science*,  
692 326, 1–5.

693  
694 Dickie, I. A. (2010). Insidious effects of sequencing errors on perceived diversity in molecular  
695 surveys. *New Phytologist*, 188, 916-918.

696  
697 Edgar, R. C. (2010). Search and clustering orders of magnitude faster than BLAST.  
698 *Bioinformatics*, 26, 2460–1.

699  
700 Edgar, R. (2013). UPARSE: highly accurate OTU sequences from microbial amplicon reads.  
701 *Nature Methods*, 10(10), 996–998.

702  
703 Fadrosch, D. W., Ma, B., Gajer, P., Sengamalay, N., Ott, S., Brotman, R. M., & Ravel, J. (2014).  
704 An improved dual-indexing approach for multiplexed 16S rRNA gene sequencing on the  
705 Illumina MiSeq platform. *Microbiome*, 2, 6.

706  
707 Fichot, E. B., & Norman, R. S. (2013). Microbial phylogenetic profiling with the Pacific  
708 Biosciences sequencing platform. *Microbiome*, 1(1), 10.

709  
710 Fu, L., Niu, B., Zhu, Z., Wu, S., & Li, W. (2012). CD-HIT: Accelerated for clustering the next-  
711 generation sequencing data. *Bioinformatics*, 28, 3150–3152.

712  
713 Giongo, A., Crabb, D. B., Davis-Richardson, A. G., Chauliac, D., Mobberley, J. M., Gano, K.  
714 A., Mukherjee, N., Casella, G., Roesch, L. F. W., Walts, B., Riva, A., King, G., & Triplett, E. W.  
715 (2010). PANGEA: pipeline for analysis of next generation amplicons. *The ISME Journal*, *4*,  
716 852–61.  
717  
718 Gweon, H. S., Oliver, A., Taylor, J., Booth, T., Gibbs, M., Read, D. S., Griffiths, R. I., &  
719 Schonrogge, K. (2015). PIPITS: An automated pipeline for analyses of fungal ITS sequences  
720 from the Illumina sequencing platform. *Methods in Ecology and Evolution*, *6*, 973-980.  
721  
722 Hallin, P. F., & Ussery, D. W. (2004). CBS Genome Atlas database: A dynamic storage for  
723 bioinformatic results and sequence data. *Bioinformatics*, *20*, 3682–3686.  
724  
725 Katoh, K., & Standley, D. M. (2013). MAFFT multiple sequence alignment software version 7:  
726 Improvements in performance and usability. *Molecular Biology and Evolution*, *30*, 772-  
727 780.reference  
728  
729 Kircher, M., Sawyer, S., & Meyer, M. (2012). Double indexing overcomes inaccuracies in  
730 multiplex sequencing on the Illumina platform. *Nucleic Acids Research*, *40*, e3.  
731  
732 Kozich, J. J., Westcott, S. L., Baxter, N. T., Highlander, S. K., & Schloss, P. D. (2013).  
733 Development of a dual-index sequencing strategy and curation pipeline for analyzing amplicon  
734 sequence data on the MiSeq Illumina sequencing platform. *Applied and Environmental*  
735 *Microbiology*, *79*, 5112-5120.  
736  
737 Krohn, A. (2016). akutils-v1.2: Facilitating analyses of microbial communities through QIIME.  
738 Zenodo. [10.5281/zenodo.56764](https://doi.org/10.5281/zenodo.56764)  
739  
740 Kunin, V., Engelbrektson, A., Ochman, H., & Hugenholtz, P. (2010). Wrinkles in the rare  
741 biosphere: Pyrosequencing errors can lead to artificial inflation of diversity estimates.  
742 *Environmental Microbiology*, *12*, 118–123.

743  
744 Langille, M. G. I., Zaneveld, J., Caporaso, J.G., McDonald, D., Knights, D., Reyes, J. A.,  
745 Clemente, J. C., Burkepile, D. E., Vega Thurber, R. L., Knight, R., Beiko, R. G., & Huttenhower,  
746 C. (2013). Predictive functional profiling of microbial communities using 16S rRNA marker  
747 gene sequences. *Nature Biotechnology*, *31*, 814-821.  
748  
749 Liu, Z., DeSantis, T. Z., Andersen, G. L., & Knight, R. (2008). Accurate taxonomy assignments  
750 from 16S rRNA sequences produced by highly parallel pyrosequencers. *Nucleic Acids Research*,  
751 *36*, e120.  
752  
753 Lozupone, C., & Knight, R. (2005). UniFrac: a New Phylogenetic Method for Comparing  
754 Microbial Communities. *Applied and Environmental Microbiology*, *71*, 8228–8235.  
755  
756 Mahé, F., Rognes, T., Quince, C., de Vargas, C., & Dunthorn M. (2014). Swarm: robust and fast  
757 clustering method for amplicon-based studies. *PeerJ*, *2*:e593.  
758  
759 McDonald, D., Price, M. N., Goodrich, J., Nawrocki, E. P., DeSantis, T. Z., Probst, A.,  
760 Andersen, G. L., Knight, R., & Hugenholtz, P. (2012). An improved Greengenes taxonomy with  
761 explicit ranks for ecological and evolutionary analyses of bacteria and archaea. *The ISME*  
762 *Journal*, *6*, 610–8.  
763  
764 Mukherjee, S., Huntemann, M., Ivanova, N., Kyrpides, N. C., & Pati, A. (2015). Large-scale  
765 contamination of microbial isolate genomes by Illumina PhiX control. *Standards in Genomic*  
766 *Sciences*, *10*, 1–4.  
767  
768 Navas-Molina, J. A., Peralta-Sánchez, J. M., González, A., McMurdie, P. J., Vázquez-Baeza, Y.,  
769 Xu, Z., Ursell, L. K., Lauber, C., Zhou, H., Song, S. J., Huntley, J., Ackermann, G. L., Berg-  
770 Lyons, D., Holmes, S., Caporaso, J. G., Knight, R. (2013). Advancing our understanding of the  
771 human microbiome using QIIME. *Methods in Enzymology*, *531*, 371-444.  
772

- 773 Nelson, M. C., Morrison, H. G., Benjamino, J., Grim, S. L., & Graf, J. (2014). Analysis,  
774 optimization and verification of Illumina-generated 16S rRNA gene amplicon surveys. *PloS*  
775 *One*, 9, e94249.  
776
- 777 Nguyen, N. H., Smith, D., Peay, K., & Kennedy, P. (2015). Parsing ecological signal from noise  
778 in next generation amplicon sequencing. *New Phytologist*, 205, 1389-1393.  
779
- 780 Price, M. N., Dehal, P. S., & Arkin, A. P. (2009). Fasttree: Computing large minimum evolution  
781 trees with profiles instead of a distance matrix. *Molecular Biology and Evolution*, 26, 1641–  
782 1650.  
783
- 784 Rognes, T., Mahé, F., Flouri, T., Quince, C., & Nichols, B. (2015). vsearch: VSEARCH version  
785 1.0.16. Zenodo. [10.5281/zenodo.15524](https://doi.org/10.5281/zenodo.15524)  
786
- 787 Rohland, N., & Reich, D. (2012). Cost-effective, high-throughput DNA sequencing libraries for  
788 multiplexed target capture. *Genome Research*, 22, 939–46.  
789
- 790 Schirmer, M., Ijaz, U. Z., D'Amore, R., Hall, N., Sloan, W. T., & Quince, C. (2015). Insight into  
791 biases and sequencing errors for amplicon sequencing with the Illumina MiSeq platform. *Nucleic*  
792 *Acids Research*, 43, 1–16.  
793
- 794 Schloss, P. D., Westcott, S. L., Ryabin, T., Hall, J. R., Hartmann, M., Hollister, E. B., ... Weber,  
795 C. F. (2009). Introducing mothur: Open-source, platform-independent, community-supported  
796 software for describing and comparing microbial communities. *Applied and Environmental*  
797 *Microbiology*, 75(23), 7537–7541.  
798
- 799 Sogin, M. L., Morrison, H. G., Huber, J. a, Mark Welch, D., Huse, S. M., Neal, P. R., ... Herndl,  
800 G. J. (2006). Microbial diversity in the deep sea and the underexplored “rare biosphere”.  
801 *Proceedings of the National Academy of Sciences of the United States of America*, 103(32),  
802 12115–20.  
803

804 Wang, Q., Garrity, G. M., Tiedje, J. M., & Cole, J. R. (2007). Naive Bayesian classifier for rapid  
805 assignment of rRNA sequences into the new bacterial taxonomy. *Applied and Environmental*  
806 *Microbiology*, 73, 5261–7.