

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

3

4 Andrew Krohn^{1,2}, Bo Stevens³, Adam Robbins-Pianka⁴, Matthew Belus⁵, Gerard J. Allan^{1,2},
5 Catherine Gehring^{1,6}

6

7 ¹Department of Biological Sciences, Northern Arizona University, Flagstaff, AZ

8 ²NAU Environmental Genetics and Genomics Laboratory, Flagstaff, AZ

9 ³School of Earth Sciences and Environmental Sustainability, Northern Arizona University,
10 Flagstaff, AZ

11 ⁴Department of Computer Science, University of Colorado Boulder, Boulder, CO

12 ⁵Anschutz Medical Campus, University of Colorado Denver, Aurora, CO

13 ⁶Merriam-Powell Center for Environmental Research, Flagstaff, AZ

14

15 Corresponding Author:

16 Andrew Krohn^{1,2}

17

18 Email address: alk224@nau.edu

19

20

21

22

23

24

25

26

27

28

29

30

31

32

33 Abstract:

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

60

61

62

63

64 Introduction:

65 Over the past decade, amplicon sequencing of marker gene fragments has become the
66 preferred method for profiling the diversity of microbial communities. Briefly, the technique
67 uses the polymerase chain reaction (PCR) to amplify a pool of PCR products from an
68 environmental sample to be resolved by high throughput DNA sequencing. Similar sequences
69 are binned together into operational taxonomic units (OTUs) and compared against a database to
70 obtain taxonomic classifications. Amplicon sequencing is flexible in that a community can be
71 profiled for different genes which may represent markers specifically suited for identification of
72 certain microbial constituents (e.g., 16S for bacteria and archaea, ITS for fungi). Similarly,
73 profiling with functional genes can offer a better understanding of community traits (e.g.,
74 Bentzon-Tilia et al., 2015). While communities were originally profiled with high-throughput
75 sequencing on 454 pyrosequencing instruments (Sogin et al., 2006), amplicon sequencing has
76 been adapted to newer instrumentation including sequencers from Illumina (Caporaso et al.,
77 2012) and Pacific Biosciences (Fichot & Norman, 2013). Illumina sequencing is currently the
78 most popular option due to several factors including cost, throughput, instrument availability,
79 and the existence of multiple protocols for amplification and sequencing of marker gene pools on
80 this platform (Caporaso et al., 2012; Bokulich & Mills, 2013; Kozich et al., 2013; Fadrosh et al.,
81 2014).

82 Accurate determination of community diversity and taxonomic content are often primary
83 aims of community amplicon sequencing projects. Systematic errors experienced during sample
84 preparation, such as PCR and sequencing errors, can contribute to overestimation of diversity
85 (Kunin, 2010). Additionally, signal cross-talk during index sequence cycles on Illumina
86 sequencers can lead to false identification of an organism in a sample (Kircher, Sawyer &
87 Meyer, 2012; Nelson et al., 2014). In the face of such potential complications, careful analysis is
88 necessary to ensure that diversity estimates are not inflated and that data are properly filtered to
89 avoid Type II errors. Several comprehensive tools exist for processing such data including
90 mothur (Schloss et al., 2009), QIIME (Caporaso et al., 2010a), and UPARSE (Edgar, 2013).
91 Many stand-alone tools are also available for performing specific bioinformatic tasks which may
92 or may not be implemented in QIIME, mothur or UPARSE. It may be beneficial in some cases to
93 perform separate bioinformatic steps with different software packages in order to obtain the most

94 accurate community representation for a given ecosystem. For instance, the use of various pre-
95 processing tools (e.g., error correction, chimera filtering) may improve the outcome for a given
96 data set. In this instance, the average researcher would require greater familiarity with the
97 production and processing of amplicon sequencing data in order to make the best decisions
98 during data processing.

99 Automated quality filtering is among the first steps performed in any sequencing project
100 and is a necessity for managing modern DNA sequencing data sets. To achieve the status of
101 “finished,” genome sequencing projects require consensus base quality scores where the
102 likelihood of an incorrect base call is less than 1 in 100,000 (q50), whereas assemblies using
103 unfiltered data are considered “standard draft” and are expected to contain errors (Chain &
104 Grafham, 2009). The default parameters in QIIME 1.9.1 require a minimum quality score of q4
105 as recommended by Bokulich et al. (2013), and should be similarly treated as “draft” data. More
106 reads are retained for downstream analysis, but a low quality score requirement also introduces
107 an unknown degree of sequencing error as base quality scores may vary widely across a single
108 sequencing run. Thus, data generated on runs with higher average error rates are more likely to
109 overestimate alpha diversity if quality scores are not strictly controlled (at the expense of
110 sequencing depth). Inconsistent qualities from sequencing runs can be effectively controlled via
111 quality filtering, and default quality filtering in QIIME retains reads that may be variably
112 trimmed to a range of 75-100% of the original sequence length. Because the quality of different
113 sequences may decrease non-uniformly across a sequencing run, variable read lengths may also
114 contribute to an inflated estimate of OTU richness if reads are not de-replicated or sorted by size
115 prior to clustering. Various error correction algorithms are available for processing Illumina data
116 (e.g., Kelley, Schatz & Salzberg, 2010; Medvedev et al., 2011; Nikolenko, Korobeynikov &
117 Alekseyev, 2013), the use of which may result in an increased number of reads retained
118 following quality filtering. Callahan et al. (2016) recently demonstrated a data processing
119 workflow that utilized error correction with good success, where the number of expected taxa
120 approximately equaled the number of observed OTUs, though we do not explore the use of error
121 correction techniques here. Chimera filtering, commonly performed following quality filtering, is
122 essential to remove PCR artifacts and further improves sequencing data quality.

123 Quality-filtered amplicon sequencing data are clustered into OTU definitions, a
124 computational process for which numerous programs are available. CD-HIT (Fu et al., 2012),
125 UCLUST (Edgar, 2010), BLAST (Altschul, 1990), and Swarm (Mahé et al., 2014) are popular
126 options that are all available in QIIME. Reference-based analysis techniques, such as BLAST,
127 are known to incur biases according to the choice of reference database (Nelson et al., 2014), but
128 can easily be parallelized for more efficient computation. UCLUST can utilize a reference
129 database, perform database-independent *de novo* clustering, or, as with the open-reference
130 strategy currently implemented in QIIME, a combination of both methods (Navas-Molina et al.,
131 2013). Pure *de novo* analysis is preferred by many as the approach least likely to impose a bias
132 on the final outcome. One popular option for *de novo* OTU clustering is CD-HIT, but as this
133 program cannot be parallelized it can be time-prohibitive when used with larger data sets.
134 Swarm, another *de novo* OTU clustering program, allows for portions of the *de novo* clustering
135 process to be parallelized, thus eliminating database-specific effects while also optimizing
136 computational requirements. All OTU picking programs require the researcher to choose a
137 similarity or distance threshold beyond which two sequences must be considered as separate
138 OTUs. If present at this stage, PCR or sequencing errors may contribute to OTU inflation to an
139 unknown degree. In addition to ensuring the data are properly filtered, one can also utilize a
140 conservative clustering threshold in order to avoid overestimation of community diversity (i.e.,
141 $\leq 97\%$; Kunin et al., 2010).

142 Taxonomic assignment, achieved through comparison of OTU definition sequences to a
143 reference database, can also be performed in a variety of ways. Popular methods include
144 BLAST, UCLUST, and RDP (Wang et al., 2007), and each are available in QIIME. In 2008, Liu
145 et al. reported that RDP provided the most accurate taxonomic assignments. Presently, other
146 techniques continue to be utilized by various amplicon sequencing analysis pipelines (e.g.,
147 Giongo et al., 2010; Gweon et al., 2015), revealing a lack of consensus among researchers.
148 Considering that improved taxonomic accuracies may be observed when sequences obtained for
149 study organisms are more similar to those populating the reference database, the relative success
150 of each algorithm may be context-dependent. For environmental data sets, accuracies of
151 taxonomic assignments are estimated by means of a confidence or quality value relevant to the
152 utilized technique (e.g., e-value for BLAST). Careful assessment of taxonomic accuracies can

153 only be done when the sequence content of a given sample can be anticipated. This can be
154 achieved with synthetic mock communities created *in silico* by extracting sequences from a
155 database (e.g., Bellemain et al., 2010) or using genomic mock communities that combine DNA
156 extracts from cultured organisms. Neither scenario is likely to provide an outcome that is directly
157 comparable to the natural complexities of environmental communities, yet both can offer a
158 measure of accuracy for taxonomic assignment methods.

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

171 Considering samples independently, Kircher, Sawyer & Meyer (2012) observed an
172 indexing inaccuracy rate of 0.3%, citing cluster mixing during sequencing as a mechanism by
173 which single-indexed Illumina sequences are likely attributed incorrectly to a particular sample.
174 For certain applications, their result argues that such data must be filtered at 0.3% by sample in
175 order to avoid Type II errors. Another common practice is to remove singleton OTUs (by sample
176 or by table) under the assumption that such OTUs represent errors generated during sequencing
177 (see Dickie, 2010). However, errors introduced during early PCR cycles may be faithfully
178 replicated many times so as to appear as valid OTUs, causing overestimation of OTU richness
179 even after singleton filtering (Nguyen et al., 2015). As an alternative, Nguyen et al. (2015)
180 suggest the removal of low-count or low-proportion OTUs by sample at a threshold informed by
181 mock community data. Mock communities used in this way may also identify certain sequence
182 motifs prone to error, which may help to identify whether novel OTUs observed in

183 environmental data should be considered suspect. Unfortunately, such controls are not available
184 for many data sets and artificial communities may not perform similarly to environmental
185 communities during sample preparation and analysis. Because samples are amplified
186 independently, PCR errors are likely to be present in the form of private OTUs observed only in
187 a single sample, so removal of unshared OTUs may be another effective precaution against
188 overestimation of diversity due to sequencing error.

189 As these examples illustrate, accurate filtering of an OTU table is not straightforward.
190 The sequence misattribution rate reported by Kircher, Sawyer & Meyer (2012) is vastly different
191 than the filtering threshold of 0.005% recommended by Bokulich et al. (2013), though their
192 recommendation was to filter across the entire OTU table. Since many amplicon sequencing
193 studies report relatively few taxa present above 0.3% per sample, filtering by sample at this
194 threshold (Kircher threshold) will exclude many valid taxa. The presence of misattributed
195 sequences may also diminish the efficacy of private OTU removal to eliminate PCR errors,
196 though dual-indexing of samples should reduce or eliminate sequence misattribution events
197 (Kircher, Sawyer & Meyer, 2012). Singleton filtering, however applied, is unlikely to be
198 thorough enough to remove errors that are either replicated during the PCR process or systematic
199 errors from the sequencing process. For single- or dual-indexed Illumina data, filtering at 0.005%
200 across the entire table (Bokulich threshold) may represent a viable compromise between
201 confident assignment of sequences to samples and the stringency that one imposes on filtering
202 the final table.

203 In this study, we used simple genomic mock communities and an environmental data set
204 to describe the effects of parameter adjustments for methods implemented in QIIME 1.9.1
205 (Caporaso et al., 2010a) on sequence quality filtering, OTU picking, taxonomic assignment, and
206 OTU table filtering. We focused on QIIME because of its popularity and flexibility for
207 processing amplicon sequencing data sets. We hypothesized that observed OTU diversity will be
208 inflated due to the presence of PCR and/or sequencing artifacts, and that such effects will be
209 observable in simple genomic mock communities under the expectation that one OTU should be
210 observed per constituent taxon. Using five mock communities consisting of 4-8 taxa each, we
211 developed a modified protocol for the analysis of 16S community amplicon sequencing data, and
212 demonstrate the method on an environmental data set. By carefully controlling each of the steps

213 that we investigated, we were able to describe mock community compositions more correctly
214 than with a default workflow.

215

216 **Materials and Methods:**

217

218 *Mock communities*

219 DNA was extracted from axenic cultures of *Pseudomonas aeruginosa* (Proteobacteria),
220 *Proteus vulgaris* (Proteobacteria), *Klebsiella pneumoniae* (Proteobacteria), *Escherichia coli*
221 (Proteobacteria), *Bacillus megaterium* (Firmicutes), *Lactococcus lactis* (Firmicutes),
222 *Staphylococcus aureus* (Firmicutes), and *Micrococcus luteus* (Actinobacteria) using a PowerSoil
223 DNA Extraction Kit (MoBio Laboratories, Carlsbad, CA). DNA was quantified by PicoGreen
224 (Life Technologies, Carlsbad, CA) fluorescence, and normalized to approximately 0.75 ng/μL.
225 Five mock communities containing different ratios of bacterial taxa were constructed from the
226 extracted DNA. Community 0 contained equal volumes of DNA from each taxon; Community
227 1a contained 8% *M. luteus*, 42% *B. megaterium*, 42% *L. lactis*, and 8% *S. aureus*; Community 1b
228 contained 42% *M. luteus*, 8% *B. megaterium*, 8% *L. lactis*, and 42% *S. aureus*; Community 2a
229 contained 8% *E. coli*, 8% *K. pneumoniae*, 42% *P. vulgaris*, and 42% *P. aeruginosa*; Community
230 2b contained 42% *E. coli*, 42% *K. pneumoniae*, 8% *P. vulgaris*, and 8% *P. aeruginosa*. Final
231 concentrations for each mock community were determined to be ~ 0.75 ng/μL (Table S1).
232 Expected compositions of mock communities were corrected for genome size and copy number
233 against the CBS Genome Atlas Database (Hallin & Ussery, 2004).

234

235 *Environmental samples*

236 Environmental samples with an expected environmental contrast were collected from the
237 Northern Arizona University Pinyon Pine Common Garden near Sunset Crater National
238 Monument, AZ. During garden installation in October 2009, soil samples were collected from
239 holes dug to plant seedlings (“pre-tree” treatment). Soil core samples were taken from the same
240 seedlings in December 2010 (“post-tree” treatment). The top 2 centimeters (cm) of soil were
241 brushed aside prior to taking cores. A 2.5 cm diameter metal corer was placed 2 cm from the
242 seedling base and driven to a depth of 10 cm. Samples were kept on ice in the field and stored at
243 -20 °C until DNA extraction. DNA was extracted from homogenized soil cores using a

244 PowerSoil DNA Extraction Kit. Only samples which produced a clean ribosomal PCR product
245 were included in this study, resulting in unequal sample sizes between pre-tree (n = 13) and post-
246 tree (n = 28) groups. A random number generator was used to select a subset of post-tree samples
247 (n = 13) for comparisons of data with equal sample sizes. Samples were normalized to c. 1 ng/ μ L
248 prior to PCR amplification for library construction.

249 The environmental samples presented here are meant only to allow a demonstration of
250 the effects of a mock community-based workflow optimization on real environmental data.
251 Though we expect the presence of a seedling to create additional niche space which would
252 increase observed diversity, no background soil control samples were collected in order to
253 properly test this hypothesis. Nonetheless, the two sets of soil samples can be expected to vary
254 because of the presence or absence of a seedling and also due to differences in the time of
255 sampling, both year and season.

256

257 *Library construction and sequencing*

258 Amplicons were produced in a two-step protocol as suggested by Berry et al. (2011).
259 Briefly, samples were amplified in triplicate PCR reactions for the 16S V4 region using the
260 universal bacterial/archaeal primers 515F and 806R (Bates et al., 2011). First round reactions
261 were performed in triplicate in 384 well plates. The 8 μ L volumes contained the following: 1 μ M
262 each primer (Eurofins MWG Operon, LLC), 200 μ M each dNTP (Phenix Research, Candler,
263 NC), 0.01 U/ μ L Phusion Hot Start II DNA Polymerase (Life Technologies), 1X HF Phusion
264 Buffer (Life Technologies), 3 mM MgCl₂, 6% glycerol, and 1 μ L normalized template DNA.
265 Cycling conditions were: 2 minutes at 95°C followed by 20 cycles of 30 seconds at 95°C, 30
266 seconds at 55°C, 4 minutes at 60°C. Triplicate reactions for each sample were pooled by
267 combining 4 μ L from each, and 2 μ L was used to check for results on a 1% agarose gel. The
268 remainder was diluted 10-fold and used as template in a second PCR reaction in which 12 base
269 Golay indexed sequencing tails (Caporaso et al., 2012) were added. Second round reaction
270 conditions were identical to the first round except only one reaction was conducted per sample
271 and only 15 total cycles were performed. Indexed PCR products were purified using a 1:1 ratio
272 of 18% polyethylene glycol and carboxylated magnetic beads as described in Rohland & Reich
273 (2012), quantified by PicoGreen fluorescence, and an equal mass of each sample was combined
274 into a final sample pool. The pool was purified and concentrated, and subsequently quantified by

275 quantitative PCR against Illumina DNA Standards (Kapa Biosystems, Wilmington, MA).
276 Sequencing was carried out on a MiSeq Desktop Sequencer (Illumina Inc, San Diego, CA)
277 running in paired end 2x150 mode.

278

279 *Sanger sequencing of mock community members*

280 The 16S gene for each mock community member was sequenced by the Sanger method
281 to a minimum depth of 2 in order to provide an accurate sequence for assessing taxonomic
282 assignment methods. Briefly, PCR products were produced using primers 27F (Lane, 1991) and
283 806R or 515F and 1492R (Turner et al., 1999). Products were bead-purified with 18% PEG and
284 used as template in sequencing reactions containing 0.25 µL BigDye Terminator v3.1 (Life
285 Technologies), 1X BigDye Terminator Sequencing Buffer (Life Technologies), 3 µM primer and
286 1.5 mM additional MgCl₂. Cycling conditions were: 2 minutes at 95°C followed by 60 cycles of
287 5 seconds at 95°C, 5 seconds at 50°C, 2 minutes at 60°C. Sequencing products were bead-
288 purified with a 3:1 ratio of 25% PEG, resuspended in water, and sequenced on either a 3730xl or
289 a 3130 Genetic Analyzer (Life Technologies). Chromatograms were processed in Staden
290 Package v1.7 (Staden, Beal & Bonfield, 2000) and the resulting sequences used to augment the
291 Greengenes database so that an exact match for each expected OTU would be present during
292 taxonomy assignment. Taxonomic identity for each sequence was confirmed by comparing
293 against the non-redundant database at NCBI using the online BLAST tool (Altschul et al., 1990).
294 Sequences were deposited to GenBank with accession numbers KY007579-KY007586.

295

296 *Data processing and statistical analysis*

297 All bioinformatics were carried out on a Mac Pro (Apple, Inc.) running Ubuntu Linux
298 14.04 LTS (Canonical Ltd.) or the Monsoon high-performance computing cluster at Northern
299 Arizona University (<https://nau.edu/hpc/>) running CentOS 6.6 (The CentOS Project). Figures
300 were generated in Veusz v1.24 (<http://home.gna.org/veusz/>) or Geneious v8.1 (Biomatters Ltd.).
301 As contaminating PhiX Control sequence can complicate sequencing projects (Mukherjee et al.,
302 2015), we calculated the amount of PhiX Control among our demultiplexed data and removed it
303 prior to sample processing. This task was performed with the `akutils phix_filtering`
304 command in `akutils` v1.2 (Krohn, 2016; <https://github.com/alk224/akutils-v1.2>) which maps raw

305 data against the Enterobacteria phage phiX174 sensu lato complete genome sequence
306 (NC_001422.1) using Smalt 0.7.6 (<http://www.sanger.ac.uk/resources/software/smalt/>).

307 Overlapping paired end reads were aligned using the `akutils join_paired_reads`
308 command in `akutils` which employs the `fastq-join` command from `ea-utils` (Aronesty, 2011).
309 Demultiplexing and quality filtering of raw, joined data (mean length = 253 bp) was carried out
310 in QIIME with the `split_libraries_fastq.py` script using default parameters, or with
311 more strict requirements of a minimum quality threshold of q_{20} ($q = 19$), allowing 0-3 low-
312 quality base calls ($r = 1-3$), and requiring at least 95% of each read to be high quality ($p = 0.95$).
313 Chimeras were removed by the UCHIME method (Edgar et al., 2011) as implemented in `vsearch`
314 1.1.1 (Rognes et al., 2016) using either the `-uchime_denovo` or `-uchime_ref` option against
315 the Gold reference database (<http://drive5.com/uchime/gold.fa>). OTU picking and taxonomy
316 assignments were performed using the `akutils pick_otus` command in `akutils` which calls
317 standard functions in QIIME. After manual inspection of sequence divergence among congeneric
318 mock community members, sequences were dereplicated on the first 100 bases using the
319 `prefix_suffix` OTU picker in QIIME. OTU picking was performed with multiple similarity or
320 distance thresholds using common OTU picking algorithms (CD-HIT, UCLUST and BLAST at
321 97%, 95%, 92%, 90%, 85%, and Swarm at d_1 , d_2 , d_3 , d_4 , d_5). BLAST was used only for closed
322 reference analysis, UCLUST for open reference analysis, and CD-HIT and Swarm for *de novo*
323 analyses. Taxonomy was assigned using BLAST, RDP, and UCLUST options with default
324 settings available in QIIME 1.9.1 (UCLUST option in QIIME actually uses the USEARCH
325 algorithm for database matching steps). Reference-based OTU picking steps and taxonomic
326 assignments were conducted against the Greengenes 97% database (McDonald et al., 2012)
327 which had been formatted to include only the V4 region using the `akutils`
328 `format_database` command in `akutils`. Sequence alignments and phylogenetic trees were
329 produced using the `akutils align_and_tree` command in `akutils` which aligns sequences
330 using PyNAST (Caporaso et al., 2010b) and generates phylogenies with FastTree (Price, Dehal
331 & Arkin, 2009). Diversity analyses were conducted using the `akutils core_diversity`
332 command in `akutils`.

333 In order to facilitate assessment of optimal workflow steps, we first sought to establish a
334 method of filtering the final OTU tables by eliminating OTUs resulting from mixed clusters. To

335 this end, we processed the mock and environmental data sets through a default QIIME workflow
336 (see below) to assess taxonomic components, and compared methods for filtering OTU tables to
337 remove contaminating taxa from the mock data. An ideal filtering method should remove
338 erroneous OTUs that arise either from sequencing error or cluster mixing. Table filtering was
339 carried out using either the Kircher threshold (0.3% by sample; Kircher, Sawyer & Meyer,
340 2012), the Bokulich threshold (0.005% by table; Bokulich et al., 2013), singletons removed by
341 table (`mc2`), or singletons removed by sample (`n2`). Private OTUs were assumed to be errors and
342 were also removed in the `n2` tables. Filtered OTU tables were grouped according to filtering
343 method, and differences in the amount of OTUs classified as contaminating taxa was assessed by
344 one-way ANOVA. Tukey's HSD test was used to determine which groups were statistically
345 distinct.

346 An optimal workflow was chosen by assessing diversity estimates and taxonomic
347 identities assigned to mock community data. The optimal OTU picking algorithm was
348 determined as the method that yielded the correct diversity result over the broadest range of
349 similarity or distance thresholds. Taxonomic accuracy was determined by seeding the
350 Greengenes database with the expected sequences from the mock community constituent taxa
351 prior to analysis, and inspecting the results. OTU tables from the optimal workflow across the
352 accurate range of similarity thresholds were filtered at each of the four thresholds described
353 above. Our "default QIIME workflow" was identical to the optimal workflow with the following
354 changes: the `split_libraries_fastq.py` command was performed with default settings;
355 OTU picking was performed with the `pick_open_reference_otus.py` command;
356 taxonomic assignment was performed with UCLUST; OTU tables were filtered with the
357 Bokulich threshold. Results from the optimal workflow were compared to the result obtained
358 from our default workflow. Environmental data was then processed using the best workflow
359 determined from this process and compared to the default result.

360 Diversity analyses for mock community data were calculated on OTU tables that had
361 been rarefied to 10,000 reads, or 5,000 reads for environmental data. Comparison of observed
362 mock community composition to the *a priori* expectation (Table S1) was conducted with
363 Spearman's rank correlation using species-level assignments. Comparison of observed OTU
364 diversity between environmental sample groupings was performed with nonparametric t-tests. A

365 random subset of post-tree samples from the environmental data ($n = 13$) was selected to
366 determine if unequal sample sizes were contributing to observed OTU diversity. Distance
367 matrices were calculated from environmental data for weighted UniFrac distance (Lozupone &
368 Knight, 2005). Tests of differences of total beta diversity were carried out on distance matrices
369 using PERMANOVA (Anderson, 2001), and differences in multivariate dispersion were detected
370 with PERMDISP (Anderson, Ellingsen & McArdle, 2006).

371 Representative sequences for the optimized mock community result were extracted from
372 the output data. When multiple OTU definition sequences represented the same taxonomic
373 identity, they were aligned with Mafft v7.123b (Kato & Standley, 2013) using the L-INS-i
374 setting. The lower abundance OTU for each multi-OTU taxon was assumed to be erroneous and
375 base differences compared to the major OTU were characterized. Trinucleotide motifs preceding
376 each base difference and terminal truncation position were tabulated. Because 2x150 sequencing
377 data does not fully overlap for 515F-806R amplicons (mean length = 253 bp), terminal base and
378 preceding trimers were considered in the context of the second read. Environmental data
379 processed through the optimal workflow was also investigated for terminal truncation positions
380 and preceding trinucleotide motifs. Because we have no reliable reference sequence for many
381 environmental OTUs, we investigated only OTUs that shared a taxonomic designation with at
382 least one other OTU, and had been truncated by more than 3 bases during quality filtering. For
383 mock and environmental data, motif and terminal base representations were tested against the
384 assumption of random occurrence with Chi-square tests.

385 We attempted to determine actual sequencing error rates for data used in either the
386 default QIIME workflow or our optimized workflow. Mock community reads were
387 demultiplexed in QIIME with `split_libraries_fastq.py` under default or strict quality
388 filtering, utilizing the `--store_demultiplexed_fastq` option. Demultiplexed fastq files
389 were imported into Geneious and aligned against the Sanger sequencing data for each mock
390 community member, requiring a 95% similarity in order to exclude contaminant sequences from
391 the alignments. The resulting alignments were exported in SAM format and SAM "NM" flags
392 were calculated in SAMtools v1.19 (Li et al., 2009). The `sam-stats` command in ea-utils was
393 used to calculate mismatch rates ("snp rate" field).

394

395 Results:

396 The sequencing run clustered at 1119 k/mm² (+/- 70) and resulted in 17.96 million total
397 reads passing filter, an overall error rate of 0.36%, and 91% of reads exceeded q30. PhiX
398 Sequencing Control v3 sequences (Illumina, Inc.) constituted 8.31% of the total run (percent
399 aligned). Once demultiplexed, mock community data contained 4.35% PhiX (103,070/2,371,510
400 reads) while the environmental data contained 4.10% PhiX (259,366/6,332,586 reads). Mock
401 community data demultiplexed under default parameters were determined to have an average
402 error rate of 0.3661% while stringent quality filtering yielded an improved error rate of 0.0990%.
403 Actual error rates varied for each taxon (Table S2). Denovo chimera detection found zero
404 chimeric reads, while reference-based detection consistently identified chimeras at a rate of
405 about 1% for each data set. As the more conservative option, we chose to utilize reference-based
406 chimera detection for the remainder of this study. Sample metadata is available in Table S11.
407 Raw sequencing data for samples used in this study are publicly available in the NCBI Sequence
408 Read Archive (study accession SRP091609; BioProject PRJNA348617).

409 Under default QIIME assessment, the mock community data showed substantial OTU
410 inflation; where there should have been just 8 OTUs, there were 127 (Table S3). When the
411 environmental data set was processed through the same workflow, 73 OTUs were classified at
412 the family level as Sphingomonadaceae. Together, these OTUs made up 5.3% of environmental
413 sequences, and Sphingomonadaceae was the most abundant classification observed at the family
414 level (Table S4). Three OTUs representing about 0.13% of the mock community data set were
415 also classified as Sphingomonadaceae, a designation which should be absent from the mock data.
416 This result led us to surmise that sequences from the environmental data set were contaminating
417 the mock communities during sequencing. Such sample cross-talk presumably arises from the
418 cluster mixing effect described by Kircher, Sawyer & Meyer (2012) where the index read from a
419 flowcell cluster is spuriously attributed to a neighboring cluster. The mock data also contained 3
420 OTUs classified as Planococcaceae (<0.03%) and 1 OTU classified as Methylobacteriaceae
421 (<0.01%), again corresponding with OTUs observed within the environmental data.
422 Sphingomonadaceae sequences were observed across all five mock communities, whereas
423 Planococcaceae was only associated with communities 0, 1a, and 1b, suggesting that cluster-
424 mixing events may occur non-randomly. Methylobacteriaceae was present as just a single read

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

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

439 Default quality filtering and OTU picking in QIIME resulted in overestimation of mock
440 community diversity regardless of how the final OTU table was filtered (Figure 1a-d; Figure S1:
441 Default mock community rarefactions). Diversity estimates were inflated up to 35 times when
442 singletons were removed by table, compared to nearly 3.5 times when filtering with the Kircher
443 threshold. Despite the reduction of OTU inflation by an order of magnitude, these results indicate
444 that revisions to initial processing steps may yield improved results. We therefore sought to
445 establish an optimized workflow that would produce the correct number of OTUs for an input of
446 known constituents. Using data that had been filtered according to strict standards during the
447 `split_libraries_fastq.py` step in QIIME ($q = 19, r = 0, p = 0.95$), a correct result was
448 achieved for each of the OTU picking algorithms tested. However, each algorithm differed in
449 which similarity threshold was required for the optimal result (Table 1). Closed reference OTU
450 picking with BLAST overestimated diversity above a similarity threshold of 92%. Open
451 reference OTU picking with UCLUST overestimated diversity at every threshold except 95%
452 similarity. *De novo* OTU picking using CD-HIT at thresholds below 92% and Swarm resolutions
453 below *d4* underestimated diversity. Swarm yielded the correct result over the broadest range of
454 tested distance thresholds (*d1-d4*), and offers other attractive features that made it stand out

455 among the tested OTU pickers (e.g., *de novo* picking, multi-threaded analysis). Thus, Swarm was
456 chosen as the optimal OTU picking method for the remainder of the study. We chose *d4* distance
457 as the optimal threshold as it was the most conservative setting to yield a correct result.

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

470 A perfect result for analysis of our mock communities requires stringent quality filtering
471 of the raw data. Default quality filtering in QIIME 1.9.1 was established according to Bokulich et
472 al. (2013). This imposes a minimum Phred quality score of 4 ($q = 3$), truncates sequences after
473 three bases are observed below this threshold ($r = 3$), and retains truncated reads that represent a
474 minimum of 75% of the original sequence length ($p = 0.75$). In contrast, we performed strict
475 quality filtering using $q = 19$, $r = 0$, and $p = 0.95$. This more stringent filtering protocol ensures
476 that data used for analysis are of much higher quality with approximately uniform read lengths.
477 An important consequence of such stringency is that much of the raw data is discarded. Of the
478 2,373,247 raw mock community sequences, default quality filtering retained 2,020,542 reads
479 (85.1%), whereas stringent parameters retained just 657,544 reads (27.7%). Holding constant $q =$
480 19 and $p = 0.95$, we found that increasing r during quality filtering had a profound effect on the
481 amount of data retained (Figure 2a). Allowing $r = 1$ resulted in an increase of data retention from
482 approximately 27% ($r = 0$) to over 56%. When $r = 2$ and $r = 3$, increases in data retention
483 showed diminishing returns, with 70% and 75% of the data retained, respectively. However, we
484 also found that allowing $r > 0$ will generally cause an inaccurate estimate of the number of

485 OTUs, depending on the criteria used for OTU picking (Figure 2b). With the $d1$ resolution,
486 increasing r will create a proportional inflation in the number of OTUs determined by Swarm. At
487 $d2$ resolution, allowing $r = 1$ still correctly described our simple mock community whereas
488 allowing $r = 2$ or $r = 3$ caused diversity to be overestimated. At resolutions $d3$ and $d4$, allowing r
489 > 0 caused underestimates of diversity. This suggests that the best result is obtained with the
490 most stringent quality filter, which we selected for our optimal workflow ($q = 19$). Similar results
491 using more data may be possible by allowing a small amount of errors (e.g., $r = 1$) and picking
492 OTUs with a more conservative similarity or distance threshold (e.g., Swarm at $d2$ resolution).

493 The Kircher threshold was effective at removing contaminating OTUs in our mock
494 community data thus yielding a near-perfect result (Figure 3). However, we anticipated that such
495 filtering could be too stringent for environmental analysis given the low per-sample OTU
496 frequencies commonly reported (e.g., Sogin et al., 2006). We compared the expected mock
497 community results to those observed with either default settings, or optimized settings for quality
498 filtering, OTU picking and taxonomy assignment, using each of the final OTU table filtering
499 methods we tested. For all comparisons, Spearman's rank correlation yielded significant p-values
500 (< 0.001), so we present only correlation values and 95% confidence intervals (CI) here. When
501 comparing the default analysis to the expected outcome, Spearman's r showed a negative
502 correlation ($r = -0.3494$; CI = $[-0.4280, -0.2655]$). Optimized results exhibited strong positive
503 correlations regardless of filtering threshold used. Lower values for Spearman's r occurred when
504 diversity was overestimated and when contaminants were present. Correlation with the expected
505 outcome improved as filtering stringency increased with every filtering method producing a
506 dramatic improvement over the default workflow (mc2: $r = 0.8075$, CI = $[0.7663, 0.8420]$; n2: r
507 = 0.8702 , CI = $[0.8344, 0.8987]$; Bokulich threshold: $r = 0.9135$, CI = $[0.8841, 0.9357]$; Kircher
508 threshold: $r = 0.9646$, CI = $[0.9495, 0.9752]$). The Bokulich threshold was chosen as our optimal
509 OTU table filtering method because it yielded the best correlation without being overly strict.

510 Output for the environmental data using either the default or optimized workflow was
511 examined for basic diversity statistics. Default analysis identified 2508 OTUs classified into 388
512 taxonomic assignments (OTUs per taxon: mean = 6.46, median = 2; Figure S3: Default
513 environmental rarefactions). The optimized analysis identified 1533 OTUs classified into 328
514 taxonomic assignments (OTUs per taxon: mean = 4.67, median = 2; Figure S4: Optimized

515 environmental rarefactions). By treatment, OTU diversity was reduced about twofold when
516 assessed via the optimized workflow and compared to the default results (Figure 4a-4b). In the
517 default analysis, pre-tree soils hosted 978.30 +/- 128.42 OTUs while post-tree soils had 1138.35
518 +/- 86.34 OTUs (nonparametric T-test = 4.578, $p < 0.001$). In the optimized analysis, pre-tree
519 soils contained 543.28 +/- 79.32 compared to 674.95 +/- 50.21 OTUs in post-tree soils
520 (nonparametric T-test = 6.277, $p < 0.001$). Differences in beta diversity were observed between
521 treatments for each workflow using weighted UniFrac distance matrices (Figure 4c-4d; default
522 PERMANOVA = 8.181, $p < 0.001$; optimized PERMANOVA = 9.355, $p < 0.001$). We also
523 noticed an increase in multivariate dispersion in the optimized workflow, though the differences
524 were not found to be significant in either case (default PERMDISP = 1.086, $p = 0.294$; optimized
525 PERMDISP = 2.160, $p = 0.158$). When data was processed with equivalent sample sizes, the
526 same patterns were observed for both alpha diversity (pre-tree = 545.76 +/- 78.84, post-tree =
527 679.00 +/- 47.86; nonparametric T-test = 5.004, $p < 0.001$) and beta diversity (PERMANOVA =
528 6.585, $p < 0.001$), though statistical power was slightly reduced, and multivariate dispersion
529 increased (PERMDISP = 3.248, $p = 0.071$), consistent with a reduction in sample size.

530 Of the 17 OTUs observed in the optimized mock result, the nine extra OTUs therein were
531 composed of three contaminants and six spurious OTUs representing sequence variants of the
532 target taxa. All extra OTUs were present at low levels ranging from 0.003% to 0.17% per sample
533 (Table S5). That sequence counts of contaminant OTUs were observed in all samples, but only
534 for select taxa, strongly suggests that cluster mixing events occur non-randomly during Illumina
535 sequencing. Species-level mock community observations from the optimized workflow describe
536 the eight constituent taxa at approximately the correct proportions. However, six of the eight taxa
537 were represented by two OTUs each. The main OTU for each taxon was present as 6.30% to
538 19.09% of the total community while the rates of lower frequency OTUs ranged from 0.01% to
539 0.05%. Manual inspection of conspecific OTU sequence alignments revealed multiple
540 substitution and indel positions within the first 100 bases which prevented these sequences from
541 dereplicating into the correct sequence during our workflow (Table S6). Additionally, these
542 sequence variants were shorter than the main constituent sequence by at least seven bases,
543 indicating that they derive from inherently lower quality reads. Inspection of trinucleotide motifs
544 preceding each substitution or indel position did not reveal any pattern relating to the observed

545 errors (Table S7). Consistent with the results of Schirmer et al. (2015), we observed a higher rate
546 of errors among A or C bases than G or T (error ratio = 1.67). Since A and C or G and T bases
547 share fluorescence excitation wavelengths during Illumina 4-channel sequencing-by-synthesis
548 (SBS), this result suggests that some of the errors we observed were indeed the result of
549 systematic errors during sequencing, although this study was not designed to distinguish between
550 such errors and those generated during PCR. Examining the terminal trinucleotide motif
551 immediately preceding truncation positions (Table S8) we observed “TTT” 83% of the time (χ^2_{63}
552 = 271.333, $p < 0.0001$). Additionally, the correct base at the truncation position was “G” 83% of
553 the time ($\chi^2_3 = 11.33$, $p = 0.0101$). An example alignment for the two OTUs representing *B.*
554 *megaterium* is presented in Figure 5a, illustrating the “TTT” motif preceding a “G” truncation
555 position (reverse complemented).

556 Truncation positions and preceding trimers were also characterized for environmental
557 data, resulting in 34 “suspect” OTUs (Table S9). Of these, 27 OTUs had been truncated at a “G”
558 position (79.41%; $\chi^2_3 = 54.235$, $p < 0.0001$), and just 10 possible trimers were represented
559 preceding the truncation position. The motifs “TTT” and “TTC” were substantially
560 overrepresented, being observed 14 (41.18%) and 7 (20.59%) times, respectively ($\chi^2_{63} =$
561 474.235 , $p < 0.0001$). An example alignment for 5 OTUs classified to the family level as
562 Sphingomonadaceae is presented in Figure 5b, and includes one such suspect OTU with a “TTT”
563 motif preceding a “G” truncation position (reverse complemented).

564

565 Discussion:

566 Our results show that amplicon sequencing data from Illumina MiSeq instruments
567 requires stringent quality filtering in order to provide the most accurate estimates of diversity.
568 Kunin et al. (2010) found that diversity was grossly overestimated for their mock community
569 data until a quality threshold of q27 was implemented. Similarly, Nelson et al. (2014) observed
570 high overestimation of mock community diversity (25-125 times expected) unless the data was
571 carefully controlled. Our optimal workflow still overestimated the OTU diversity of our simple
572 mock communities by a factor of about two. While this is still an overestimation, it is an
573 improvement over results obtained by default processing. Our optimized protocol yielded a
574 reasonable characterization of taxonomic content for mock communities (Table S5) and

575 environmental data (Table S10) alike, though it is important to recognize that mock community
576 results may not always generalize well to environmental samples.

577 Some authors have suggested that excessive OTU diversity may be at least partially
578 explained by the presence of unfiltered chimeric reads (Edgar, 2013), ribosomal paralogs (Pei et
579 al., 2010), or laboratory contaminants (Nelson et al., 2014). It seems worth noting that the level
580 of chimeric reads in our data was very low compared to rates observed by others (e.g., Schloss,
581 Gevers & Westcott, 2011; Edgar, 2013). We speculate this is due the use of a high-fidelity
582 polymerase and low cycling conditions during library construction, consistent with the results of
583 Gohl et al. (2016). As chimeras are thought to form primarily when incomplete products from
584 the previous cycle act as primers during the extension step (Haas et al., 2011), we made use of an
585 extra-long, low temperature extension of 4 minutes at 60 °C in an attempt to minimize this effect.
586 We tested the cycling conditions by amplifying serial dilutions of 16S products by qPCR (data
587 not shown) and found it yielded an efficiency of about 1, lending further support to the
588 possibility that our data is virtually chimera-free. Intragenomic ribosomal diversity is also an
589 unlikely explanation for OTU inflation in our mock community results. While structural changes
590 are often associated with diversity of the ribosomal operon (Lim, Furuta & Kobayashi, 2012),
591 these should have little impact on the sequence diversity of the 16S V4 region. In fact, Sun et al.
592 (2013) found that the V4-V5 region suffers from lower rates of intragenomic diversity compared
593 to other variable regions of the 16S rRNA gene. Using a quality cut-off of q20 across a 253 nt
594 sequence, paralogous sequences may remain, though we did not observe any such sequences at a
595 rate high enough to be considered as potential paralogs. Further, the observed proportion of each
596 constituent was quite close to expected proportions after accounting for genome size and 16S
597 rRNA gene copy numbers (Figure 3). All contaminants that we observed in the mock community
598 data could be directly attributed to taxa present in the environmental data set.

599 Schirmer et al. (2015) observed that error rates reported by Illumina MiSeq sequencers,
600 according to the PhiX Control v3, do not accurately reflect those of amplicon sequences. Their
601 conclusion that actual error rates were higher than those indicated by q-scores reported by the
602 MiSeq has important implications for the use of Illumina sequencing in estimating microbial
603 diversity. It is possible that newer imaging strategies (e.g., 2-channel SBS chemistry used by
604 Illumina NextSeq and MiniSeq instruments) will provide improved parity between the estimated

605 and actual error rates, but this will require careful testing. Interestingly, when we attempted to
606 determine actual error rates through alignment of mock community sequences (demultiplexed
607 under default settings) to their expected result, we observed a very close correlation compared to
608 the error reported by PhiX Control (0.36% vs. 0.37%). This result was not consistent across
609 different mock community constituent taxa, suggesting that error rates can be taxon-specific
610 (Table S2). We further note that data filtered under our strict filtering conditions, which
611 stipulated a minimum per-base quality of q20, yielded an average mismatch rate of just 0.099%
612 (q30), indicating that most of our data is of exceptional quality following quality filtering.

613 Of the non-target OTUs present in our optimized mock community result, one third were
614 contaminants arising from cluster mixing events during sequencing and two thirds were sequence
615 variants of the constituent OTUs which may have arisen during PCR, sequencing, or a
616 combination of the two. Cluster mixing can be controlled by dual-indexing of samples (Kircher,
617 Sawyer & Meyer, 2012), but errors arising during PCR or sequencing represent systematic errors
618 inherent to the procedure of amplicon sequencing which are difficult, if not impossible, to
619 completely eliminate irrespective of indexing strategy. Even though dual-indexing offers a clear
620 advantage over single indexing with regard to sample attribution, single-indexed protocols (e.g.,
621 Caporaso et al., 2012) remain popular and widely used. Single-indexed data still yields valuable
622 information and should not be discounted, as long as researchers are aware of the limitations.
623 Dual-indexed designs should be encouraged for new research projects (e.g., Kozich et al., 2013;
624 Fadrosh et al., 2014).

625 We echo the recommendation by others (e.g., Bokulich et al., 2013; Schirmer et al.,
626 2015) to include control mock community samples to guide data analysis. PhiX Control v3 is
627 still needed to improve sequence diversity for the purpose of cluster map generation
628 (<https://goo.gl/NpauDN>), but an alternative reference sequence could be used with onboard mock
629 communities to more directly estimate error profiles for community amplicon sequencing data.
630 PhiX sequence itself likely contributes little, if at all, to inflation of diversity estimates, and is
631 easily quantified and removed. Though such an effect is direct evidence of cluster mixing, the
632 rate of PhiX infiltration is likely much higher than the rate of sample mixing because PhiX
633 Control is unindexed, producing no fluorescent signal during indexing cycles. Spurious OTUs

634 defined from contaminating PhiX sequence may be more prevalent amid sequence data which
635 was accompanied by higher concentrations of PhiX Control v3 during sequencing.

636 Although this study was not designed for careful investigation of errors generated during
637 amplicon sequencing projects, we were able to observe that certain bases and motifs were more
638 frequently associated with low-quality base calls than should be expected by chance. The
639 presence of a “TTT” or “TTC” motif immediately preceding a “G” position near the end of a
640 sequence (near the start of the second read) was most frequently associated with an erroneous or
641 suspect OTU (Table S8). Indeed, mock community diversity was inflated on account of this
642 effect, but determining the source of such error requires more careful investigation than is
643 possible here, given that this study derives from a single MiSeq run with limited taxonomic
644 diversity. In addition to the terminal truncation observations, we note that all other observed
645 errors in the mock community sequences occurred within the first 100 bp of sequence, specific to
646 the non-overlapping region of the first sequencing read (Figure 5a). It is likely that the errors we
647 observed here would have occurred less frequently had we used fully-overlapping reads for this
648 study. Importantly, the motif-specific patterns we observed were consistent between the mock
649 and environmental data sets (Figure 5; Table S8; Table S9).

650 Estimates of alpha diversity are more sensitive than beta diversity calculations to the
651 effects of cluster mixing and systematic errors. Increasing the number of allowed low-quality
652 reads (*r* parameter in `split_libraries_fastq.py`) increases the amount of data available
653 for processing, but also changes observed diversity. For this reason, we suggest that alpha
654 diversity estimates should be performed only with data that has been stringently filtered for
655 quality. Because errors in amplicon sequencing data may follow sequence-specific patterns
656 (Schirmer et al., 2015; this study), spurious OTUs may provide artificial support to the statistical
657 separation of experimental treatments. Alternatively, spurious OTUs arising from taxa which are
658 not differentially represented among treatments could provide artificial noise, making it more
659 difficult to detect real differences. In either scenario, careful quality filtering can diminish such
660 effects.

661 Our results suggest that alpha diversity can be overestimated if sequencing error rates are
662 not carefully controlled. Here we observed this effect with a QIIME-based workflow, although
663 QIIME is just one of a variety of tools used in data analysis for such work. Because errors may

664 arise systematically during PCR or sequencing implies that a similar effect is likely to be
665 observed regardless of which analysis pipeline is used to assess the data. We made use of a high-
666 fidelity polymerase (Phusion Hot Start II) in contrast to many studies which continue to utilize
667 Taq polymerase, with which PCR-derived errors will be more prevalent. Lower fidelity will
668 promote more PCR-derived errors, and those generated during early cycles will be highly
669 perpetuated, an effect which would be more problematic under high-cycling conditions. This
670 effect was recently demonstrated by Gohl et al. (2016), who also showed that PCR-chimeras are
671 virtually absent from protocols utilizing low-cycling conditions. Because errors may follow
672 sequence-specific patterns, some diversity estimates may be particularly inflated for certain taxa,
673 which can further affect studies using taxonomic content to predict community function (e.g.
674 Langille et al., 2013). The use of phylogenetic metrics (e.g., phylogenetic diversity for alpha
675 diversity, UniFrac for beta diversity) during data analysis will likely diminish the effects of
676 complications associated with systematically-inflated OTU diversity. Though the quality-
677 filtering recommendations outlined by Bokulich et al., (2013) have subsequently provided
678 valuable guidance to numerous researchers, newer quality-filtering methods promise
679 improvements in accuracy and read retention (e.g., Puente-Sánchez, Aguirre & Parro, 2016).
680 Careful consideration of the results presented here and elsewhere (Kunin et al., 2010; Schirmer et
681 al., 2015) will improve upon our collective interpretation of microbial diversity across
682 environments.

683

684 **Conclusions:**

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

693

694 1. Remove PhiX Control v3 contamination with Smalt

- 695 2. Align read pairs with fastq-join
- 696 3. **Strict quality filter in QIIME ($q = 19, r = 0, p = 0.95$)**
- 697 4. Chimera filtering with vsearch
- 698 5. Sequence dereplication with prefix/suffix OTU picker
- 699 6. **Pick OTUs with Swarm ($d4$ resolution, adjust as necessary)**
- 700 7. **Assign taxonomy with BLAST (default settings)**
- 701 8. **Filter output table at the Bokulich threshold**
- 702

703 Our results were consistent with the hypothesis that mock community diversity would be
704 inflated due to the presence of PCR or sequencing errors in the data. By imposing more rigorous
705 quality filtering of raw sequencing data, much of this error is removed. The effects of remaining
706 errors can be minimized by utilizing a conservative similarity or distance threshold during OTU
707 picking. By characterizing mock communities at multiple thresholds, one can identify a
708 sufficiently conservative similarity or distance value ($d4$ in our case) which should offer
709 improved confidence when measuring environmental diversity. If mock communities are
710 unavailable, we advocate the use of a workflow based upon the above optimization. For studies
711 utilizing an alternative locus, we suggest adjusting the clustering threshold based on the length of
712 the amplicon (e.g., more conservative clustering for longer amplicons) until mock communities
713 can be employed to determine a more informed threshold.

714

715 **Acknowledgements:**

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

719

720 **References:**

721

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

724

725 Anderson, M. J. (2001). A new method for non-parametric multivariate analysis of variance.

726 *Austral Ecology*, 26, 32-46.

727

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

730

731 Aronesty, E. (2011). *ea-utils*: Command-line tools for processing biological sequencing data;
732 <https://expressionanalysis.github.io/ea-utils/>

733

734 Bates, S. T., Berg-Lyons, D., Caporaso, J. G., Walters, W. A., Knight, R., & Fierer, N. (2011).
735 Examining the global distribution of dominant archaeal populations in soil. *The ISME Journal*, 5,
736 908–17.

737

738 Bellemain, E., Carlsen, T., Brochmann, C., Coissac, E., Taberlet, P., & Kausrud, H. (2010). ITS
739 as an environmental DNA barcode for fungi: an in silico approach reveals potential PCR biases.
740 *BMC Microbiology*, 10, 189.

741

742 Bentzon-Tilia, M., Traving, S. J., Mantikci, M., Knudsen-Leerbeck, H., Hansen, J. L., Markager,
743 S., & Riemann, L. (2015). Significant N₂ fixation by heterotrophs, photoheterotrophs and
744 heterocystous cyanobacteria in two temperate estuaries. *The ISME Journal*, 9, 273–285.

745

746 Berry, D., Mahfoudh, K. B., Wagner, M., & Loy, A. (2011). Barcoded Primers Used in
747 Multiplex Amplicon Pyrosequencing Bias Amplification. *Applied and Environmental*
748 *Microbiology*, 77, 612–612.

749

750 Bokulich, N. A., & Mills, D. A. (2013). Improved selection of internal transcribed spacer-
751 specific primers enables quantitative, ultra-high-throughput profiling of fungal communities.
752 *Applied and Environmental Microbiology*, 79, 2519–2526.

753

754 Bokulich, N. A., Subramanian, S., Faith, J. J., Gevers, D., Gordon, J. I., Knight, R., Mills, D. A.,
755 & Caporaso, J. G. (2013). Quality-filtering vastly improves diversity estimates from Illumina
756 amplicon sequencing. *Nature Methods*, 10, 57–9.

757
758 Callahan, B. J., Mcmurdie, P. J., Rosen, M. J., Han, A. W., Johnson, A. J., & Holmes, S. P.
759 (2016). DADA2: High resolution sample inference from amplicon data. *Nature Methods*, *13*,
760 581-583.

761
762 Caporaso, J. G., Kuczynski, J., Stombaugh, J., Bittinger, K., Bushman, F. D., Costello, E. K.,
763 Fierer, N., Peña, A. G., Goodrich, J. K., Gordon, J. I., Huttley, G. A., Kelley, S. T., Knights, D.,
764 Koenig, J. E., Ley, R. E., Lozupone, C. A., McDonald, D., Muegge, B. D., Pirrung, P., Reeder,
765 J., Sevinsky, J. R., Turnbaugh, P. J., Walters, W. A., Widmann, J., Yatsunenko, T., Zaneveld, J.,
766 & Knight, R. (2010a). QIIME allows analysis of high-throughput community sequencing data.
767 *Nature Methods*, *7*, 335–336.

768
769 Caporaso, J. G., Bittinger, K., Bushman, F. D., Desantis, T. Z., Andersen, G. L., & Knight, R.
770 (2010b). PyNAST: A flexible tool for aligning sequences to a template alignment.
771 *Bioinformatics*, *26*, 266–267.

772
773 Caporaso, J. G., Lauber, C. L., Walters, W. A., Berg-Lyons, D., Huntley, J., Fierer, N., Owens,
774 S. M., Betley, J., Fraser, L., Bauer, M., Gormley, N., Gilbert, J. A., Smith, G., & Knight, R.
775 (2012). Ultra-high-throughput microbial community analysis on the Illumina HiSeq and MiSeq
776 platforms. *The ISME Journal*, *6*, 1621–1624.

777
778 Chain, P., & Grafham, D. (2009). Genome project standards in a new era of sequencing. *Science*,
779 *326*, 1–5.

780
781 Dickie, I. A. (2010). Insidious effects of sequencing errors on perceived diversity in molecular
782 surveys. *New Phytologist*, *188*, 916-918.

783
784 Edgar, R. C. (2010). Search and clustering orders of magnitude faster than BLAST.
785 *Bioinformatics*, *26*, 2460–1.

786

- 787 Edgar, R. C., Haas, B. J., Clemente, J. C., Quince, C., & Knight, R. (2011). UCHIME improves
788 sensitivity and speed of chimera detection, *Bioinformatics*, *27*, 2194–2200.
789
- 790 Edgar, R. (2013). UPARSE: highly accurate OTU sequences from microbial amplicon reads.
791 *Nature Methods*, *10*, 996–998.
792
- 793 Fadrosch, D. W., Ma, B., Gajer, P., Sengamalay, N., Ott, S., Brotman, R. M., & Ravel, J. (2014).
794 An improved dual-indexing approach for multiplexed 16S rRNA gene sequencing on the
795 Illumina MiSeq platform. *Microbiome*, *2*, 6.
796
- 797 Fichot, E. B., & Norman, R. S. (2013). Microbial phylogenetic profiling with the Pacific
798 Biosciences sequencing platform. *Microbiome*, *1*, 10.
799
- 800 Fu, L., Niu, B., Zhu, Z., Wu, S., & Li, W. (2012). CD-HIT: Accelerated for clustering the next-
801 generation sequencing data. *Bioinformatics*, *28*, 3150–3152.
802
- 803 Giongo, A., Crabb, D. B., Davis-Richardson, A. G., Chauliac, D., Mobberley, J. M., Gano, K.
804 A., Mukherjee, N., Casella, G., Roesch, L. F. W., Walts, B., Riva, A., King, G., & Triplett, E. W.
805 (2010). PANGEA: pipeline for analysis of next generation amplicons. *The ISME Journal*, *4*,
806 852–61.
807
- 808 Gohl, D. M., Vangay, P., Garbe, J., MacLean, A., Hauge, A., Becker, A., Gould, T. J., Clayton,
809 J. B., Johnson, T. J., Hunter, R., Knights, D., & Beckman, K. B. (2016). Systematic improvement
810 of amplicon marker gene methods for increased accuracy in microbiome studies. *Nature*
811 *Biotechnology*, *34*, 942-949.
812
- 813 Gweon, H. S., Oliver, A., Taylor, J., Booth, T., Gibbs, M., Read, D. S., Griffiths, R. I., &
814 Schonrogge, K. (2015). PIPITS: An automated pipeline for analyses of fungal ITS sequences
815 from the Illumina sequencing platform. *Methods in Ecology and Evolution*, *6*, 973-980.
816

817 Hallin, P. F., & Ussery, D. W. (2004). CBS Genome Atlas database: A dynamic storage for
818 bioinformatic results and sequence data. *Bioinformatics*, *20*, 3682–3686.
819

820 Haas, B. J., Gevers, D., Earl, A. M., Feldgarden, M., Ward, D. V., Giannoukos, G., Ciulla, D.,
821 Tabbaa, D., Highlander, S. K., Sodergren, E., Methé, B., DeSantis, T. Z., Petrosino, J. F., Knight,
822 R. & Birren, B. W. (2011). Chimeric 16S rRNA sequence formation and detection in Sanger and
823 454-pyrosequenced PCR amplicons. *Genome Research*, *21*, 494–504.
824

825 Katoh, K., & Standley, D. M. (2013). MAFFT multiple sequence alignment software version 7:
826 Improvements in performance and usability. *Molecular Biology and Evolution*, *30*, 772-
827 780.reference
828

829 Kelley, D. R., Schatz, M. C., & Salzberg, S. L. (2010). Quake: quality-aware detection and
830 correction of sequencing errors. *Genome Biology*, *11*, R116.
831

832 Kircher, M., Sawyer, S., & Meyer, M. (2012). Double indexing overcomes inaccuracies in
833 multiplex sequencing on the Illumina platform. *Nucleic Acids Research*, *40*, e3.
834

835 Kozich, J. J., Westcott, S. L., Baxter, N. T., Highlander, S. K., & Schloss, P. D. (2013).
836 Development of a dual-index sequencing strategy and curation pipeline for analyzing amplicon
837 sequence data on the MiSeq Illumina sequencing platform. *Applied and Environmental*
838 *Microbiology*, *79*, 5112-5120.
839

840 Krohn, A. (in review). akutils-v1.2: Facilitating analyses of microbial communities through
841 QIIME. *The Journal of Open Source Software*, (in review).
842

843 Kunin, V., Engelbrekton, A., Ochman, H., & Hugenholtz, P. (2010). Wrinkles in the rare
844 biosphere: Pyrosequencing errors can lead to artificial inflation of diversity estimates.
845 *Environmental Microbiology*, *12*, 118–123.
846

847 Lane, D. J. (1991). 16S/23S rRNA sequencing. Pp. 115–176 in E. Stackebrandt and M.
848 Goodfellow, eds. *Nucleic acid techniques in bacterial systematics*. New York, NY: John Wiley.
849
850 Langille, M. G. I., Zaneveld, J., Caporaso, J.G., McDonald, D., Knights, D., Reyes, J. A.,
851 Clemente, J. C., Burkepille, D. E., Vega Thurber, R. L., Knight, R., Beiko, R. G., & Huttenhower,
852 C. (2013). Predictive functional profiling of microbial communities using 16S rRNA marker
853 gene sequences. *Nature Biotechnology*, *31*, 814-821.
854
855 Li H., Handsaker B., Wysoker A., Fennell T., Ruan J., Homer N., Marth G., Abecasis G., &
856 Durbin R. (2009). The Sequence Alignment/Map format and SAMtools. *Bioinformatics*, *25*,
857 2078–2079.
858
859 Lim, K., Furuta, Y., & Kobayashi, I. (2012). Large variations in bacterial ribosomal RNA genes.
860 *Molecular Biology and Evolution*, *29*, 2937–2948.
861
862 Liu, Z., DeSantis, T. Z., Andersen, G. L., & Knight, R. (2008). Accurate taxonomy assignments
863 from 16S rRNA sequences produced by highly parallel pyrosequencers. *Nucleic Acids Research*,
864 *36*, e120.
865
866 Lozupone, C., & Knight, R. (2005). UniFrac: a New Phylogenetic Method for Comparing
867 Microbial Communities. *Applied and Environmental Microbiology*, *71*, 8228–8235.
868
869 Mahé, F., Rognes, T., Quince, C., de Vargas, C., & Dunthorn M. (2014). Swarm: robust and fast
870 clustering method for amplicon-based studies. *PeerJ*, *2*, e593.
871
872 McDonald, D., Price, M. N., Goodrich, J., Nawrocki, E. P., DeSantis, T. Z., Probst, A.,
873 Andersen, G. L., Knight, R., & Hugenholtz, P. (2012). An improved Greengenes taxonomy with
874 explicit ranks for ecological and evolutionary analyses of bacteria and archaea. *The ISME*
875 *Journal*, *6*, 610–8.
876

- 877 Medvedev, P., Scott, E., Kakaradov, B., & Pevzner, P. (2011). Error correction of high-
878 throughput sequencing datasets with non-uniform coverage. *Bioinformatics*, *27*, 137–141.
879
- 880 Mukherjee, S., Huntemann, M., Ivanova, N., Kyrpides, N. C., & Pati, A. (2015). Large-scale
881 contamination of microbial isolate genomes by Illumina PhiX control. *Standards in Genomic
882 Sciences*, *10*, 1–4.
883
- 884 Navas-Molina, J. A., Peralta-Sánchez, J. M., González, A., McMurdie, P. J., Vázquez-Baeza, Y.,
885 Xu, Z., Ursell, L. K., Lauber, C., Zhou, H., Song, S. J., Huntley, J., Ackermann, G. L., Berg-
886 Lyons, D., Holmes, S., Caporaso, J. G., Knight, R. (2013). Advancing our understanding of the
887 human microbiome using QIIME. *Methods in Enzymology*, *531*, 371-444.
888
- 889 Nelson, M. C., Morrison, H. G., Benjamino, J., Grim, S. L., & Graf, J. (2014). Analysis,
890 optimization and verification of Illumina-generated 16S rRNA gene amplicon surveys. *PLoS
891 One*, *9*, e94249.
892
- 893 Nguyen, N. H., Smith, D., Peay, K., & Kennedy, P. (2015). Parsing ecological signal from noise
894 in next generation amplicon sequencing. *New Phytologist*, *205*, 1389-1393.
895
- 896 Nikolenko, S. I., Korobeynikov, A. I., & Alekseyev, M. A. (2013). BayesHammer: Bayesian
897 clustering for error correction in single-cell sequencing. *BMC Genomics*, *14*, 1–11.
898
- 899 Pei, A. Y., Oberdorf, W. E., Nossa, C. W., Agarwal, A., Chokshi, P., Gerz, E. A., Jin, Z., Lee, P.,
900 Yang, L., Poles, M., Brown, S. M., Sotero, S., DeSantis, T., Brodie, E., Nelson, K., Pei, Z.
901 (2010). Diversity of 16S rRNA genes within individual prokaryotic genomes. *Applied and
902 Environmental Microbiology*, *76*, 3886–3897.
903
- 904 Price, M. N., Dehal, P. S., & Arkin, A. P. (2009). Fasttree: Computing large minimum evolution
905 trees with profiles instead of a distance matrix. *Molecular Biology and Evolution*, *26*, 1641–
906 1650.
907

- 908 Puente-Sánchez, F., Aguirre, J., & Parro, V. (2016). A novel conceptual approach to read-
909 filtering in high-throughput amplicon sequencing studies. *Nucleic Acids Research*, *44*, e40.
910
- 911 Rognes, T., Flouri, T., Nichols, B., Quince, C., & Mahé, F. (2016). VSEARCH: a versatile open
912 source tool for metagenomics. *PeerJ*, *4*, e2584.
913
- 914 Rohland, N., & Reich, D. (2012). Cost-effective, high-throughput DNA sequencing libraries for
915 multiplexed target capture. *Genome Research*, *22*, 939–46.
916
- 917 Schirmer, M., Ijaz, U. Z., D’Amore, R., Hall, N., Sloan, W. T., & Quince, C. (2015). Insight into
918 biases and sequencing errors for amplicon sequencing with the Illumina MiSeq platform. *Nucleic
919 Acids Research*, *43*, 1–16.
920
- 921 Schloss, P. D., Westcott, S. L., Ryabin, T., Hall, J. R., Hartmann, M., Hollister, E. B.,
922 Lesniewski, R. A., Oakley, B. B., Parks, D. H., Robinson, C. J., Sahl, J. W., Stres, B., Thallinger,
923 G. G., Van Horn, D. J. & Weber, C. F. (2009). Introducing mothur: Open-source, platform-
924 independent, community-supported software for describing and comparing microbial
925 communities. *Applied and Environmental Microbiology*, *75*, 7537–7541.
926
- 927 Schloss, P. D., Gevers, D., Westcott, S. L. (2011). Reducing the effects of PCR amplification and
928 sequencing Artifacts on 16s rRNA-based studies. *PLoS ONE*, *6*, e27310.
929
- 930 Sogin, M. L., Morrison, H. G., Huber, J. A., Welch, D. M., Huse, S. M., Neal, P. R., Arrieta, J.
931 M. & Herndl, G. J. (2006). Microbial diversity in the deep sea and the underexplored “rare
932 biosphere”. *Proceedings of the National Academy of Sciences of the United States of America*,
933 *103*, 12115–12120.
934
- 935 Staden, R., Beal, K. F, Bonfield, J. K. (2000). The Staden package, 1998. *Methods in Molecular
936 Biology*, *132*, 115-130.
937

938 Sun, D. L., Jiang, X., Wu, Q. L., & Zhou, N. Y. (2013). Intragenomic heterogeneity of 16S
939 rRNA genes causes overestimation of prokaryotic diversity. *Applied and Environmental*
940 *Microbiology*, 79, 5962–5969.

941

942 Turner, S., Pryer, K. M., Miao, V. P., Palmer, J. D. (1999). Investigating deep phylogenetic
943 relationships among cyanobacteria and plastids by small subunit rRNA sequence analysis. *The*
944 *Journal of Eukaryotic Microbiology*, 46, 327–338.

945

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

949

950

951

952

953

954

955

956

957

958

959

960

961

962

963

964

965

966

967

968

969 Table 1: Mock community OTU picking results. Comparison of OTU picking methods using
970 mock communities with eight total taxa. Cell values represent the number of OTUs returned for a
971 given algorithm and similarity or distance threshold. Correct results are emphasized with bold
972 text.

	97%/d1	95%/d2	92%/d3	90%/d4	85%/d5
BLAST	9	9	8	8	8
CD-HIT	8	8	8	5	5
UCLUST	11	8	9	9	11
Swarm	8	8	8	8	5

973
974
975
976
977
978
979
980
981
982
983
984
985
986
987
988
989
990
991

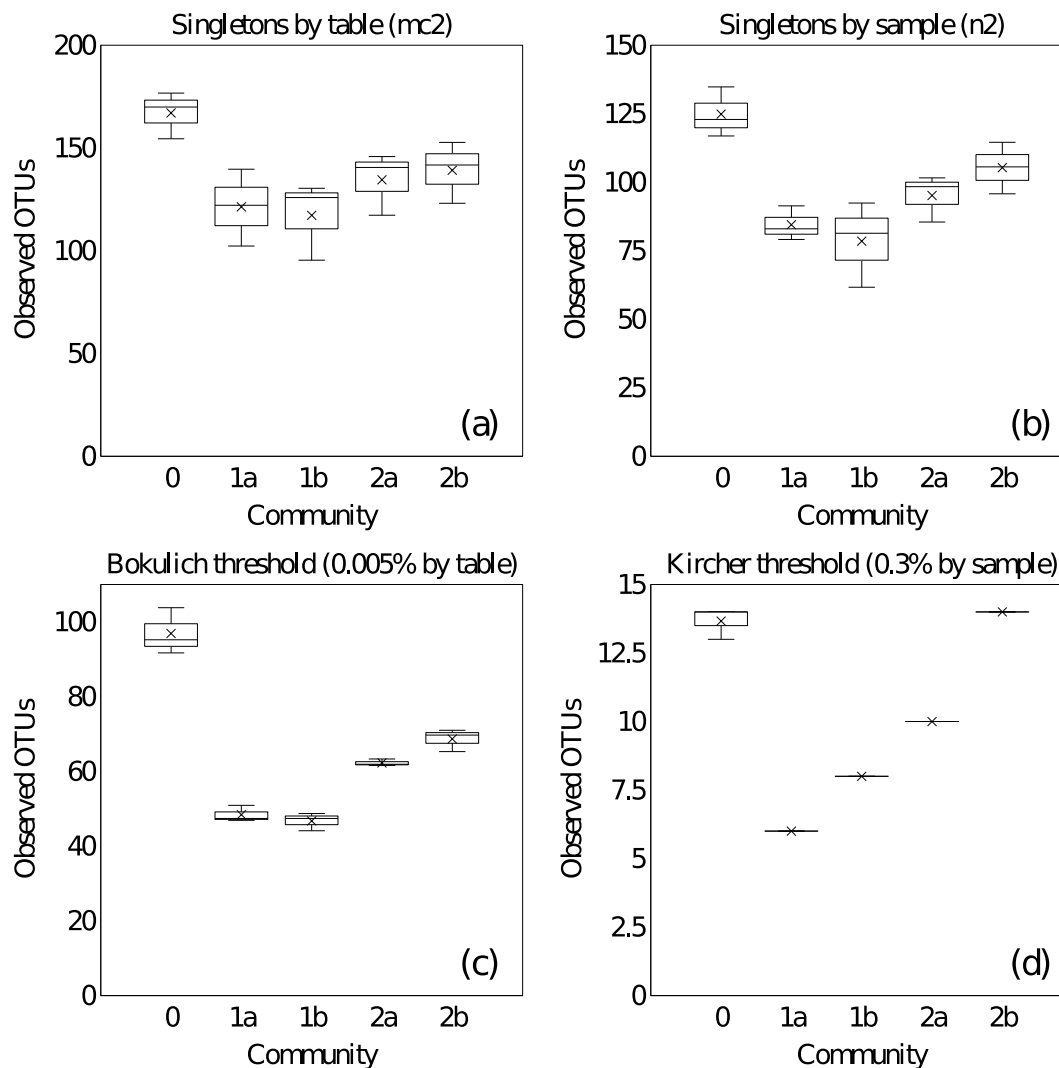
992 Table 2: Mock community taxonomic assignment results. Comparison of taxonomic assignment methods using mock communities
 993 with eight total taxa. Representative sequences were added to Greengenes to test each of the assigners for the ability to return the
 994 optimal sequence from a database search. Resulting taxonomic assignments for each method are listed with method-specific
 995 confidence values indicated in parentheses (BLAST: e-value, RDP: bootstrap confidence, and UCLUST: p-value). An asterisk (*)
 996 indicates an exact match against the correct representative sequence. The final column shows results with BLAST when the
 997 representative sequences are excluded from the database.

OTU ID	Expected result	BLAST	RDP	UCLUST	BLAST (Greengenes only)
denovo 0	<i>Escherichia coli</i>	* <i>Escherichia coli</i> (1e-139)	<i>Escherichia coli</i> (0.860)	Enterobacteriaceae (1.00)	Enterobacteriaceae (1e-139)
denovo 1	<i>Staphylococcus aureus</i>	* <i>Staphylococcus aureus</i> (6e-135)	<i>Staphylococcus</i> (0.980)	<i>Staphylococcus epidermidis</i> (0.67)	<i>Staphylococcus</i> (6e-135)
denovo 2	<i>Bacillus megaterium</i>	* <i>Bacillus megaterium</i> (2e-137)	<i>Bacillus cereus</i> (0.720)	<i>Bacillus</i> (0.67)	<i>Bacillus cereus</i> (2e-137)
denovo 3	<i>Klebsiella pneumoniae</i>	* <i>Klebsiella pneumoniae</i> (1e-139)	Enterobacteriaceae (1.000)	Enterobacteriaceae (0.67)	Enterobacteriaceae (1e-139)
denovo 4	<i>Proteus vulgaris</i>	* <i>Proteus vulgaris</i> (1e-139)	<i>Proteus</i> (0.960)	<i>Proteus</i> (0.67)	<i>Proteus</i> (2e-137)
denovo 5	<i>Lactococcus lactis</i>	* <i>Lactococcus lactis</i> (1e-136)	<i>Lactococcus</i> (0.640)	<i>Lactococcus</i> (1.00)	<i>Lactococcus</i> (1e-136)
denovo 6	<i>Micrococcus</i>	* <i>Micrococcus luteus</i>	<i>Micrococcus</i> (0.900)	<i>Micrococcus</i> (0.67)	<i>Micrococcus</i> (2e-

	<i>luteus</i>	(2e-137)			137)
denovo 7	<i>Pseudomonas aeruginosa</i>	* <i>Pseudomonas aeruginosa</i> (1e-139)	<i>Pseudomonas</i> (0.520)	Pseudomonadaceae (0.67)	<i>Pseudomonas</i> (1e-139)
Matches		8/8 (100%)	0/8 (0%)	0/8 (0%)	0/8 (0%)

998

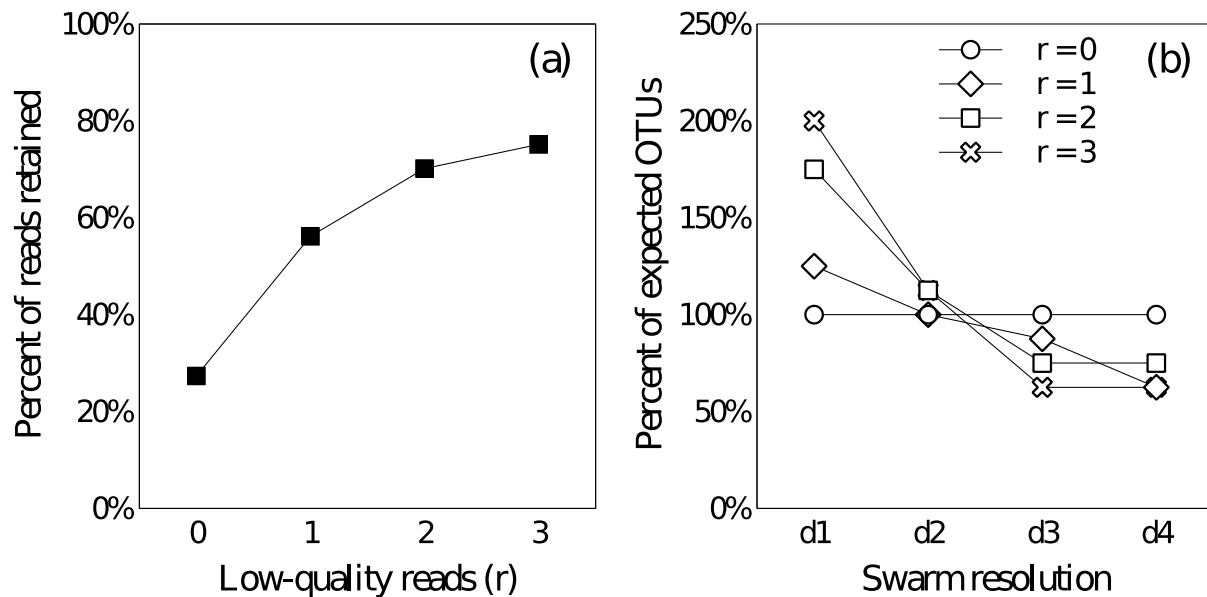
999 Figure 1: Alpha diversities of mock communities by default analysis. Alpha diversity
 1000 comparisons for mock communities using default analysis settings and rarefied to 10,000
 1001 sequences per sample. Mean values are represented by an “x” while median values are
 1002 represented by a straight line. Each plot depicts a different OTU table filtering method: (a)
 1003 singletons removed across the entire table (mc2), (b) private alleles and singletons removed per
 1004 sample (n2), (c) OTUs removed that do not exceed 0.005% of the total data (Bokulich), and (d)
 1005 OTUs removed that do not exceed 0.3% per sample (Kircher). In every case, diversity estimates
 1006 are inflated (expected values are Community 0: 8 OTUs; Communities 1a, 1b, 2a, 2b: 4 OTUs
 1007 each).



1008

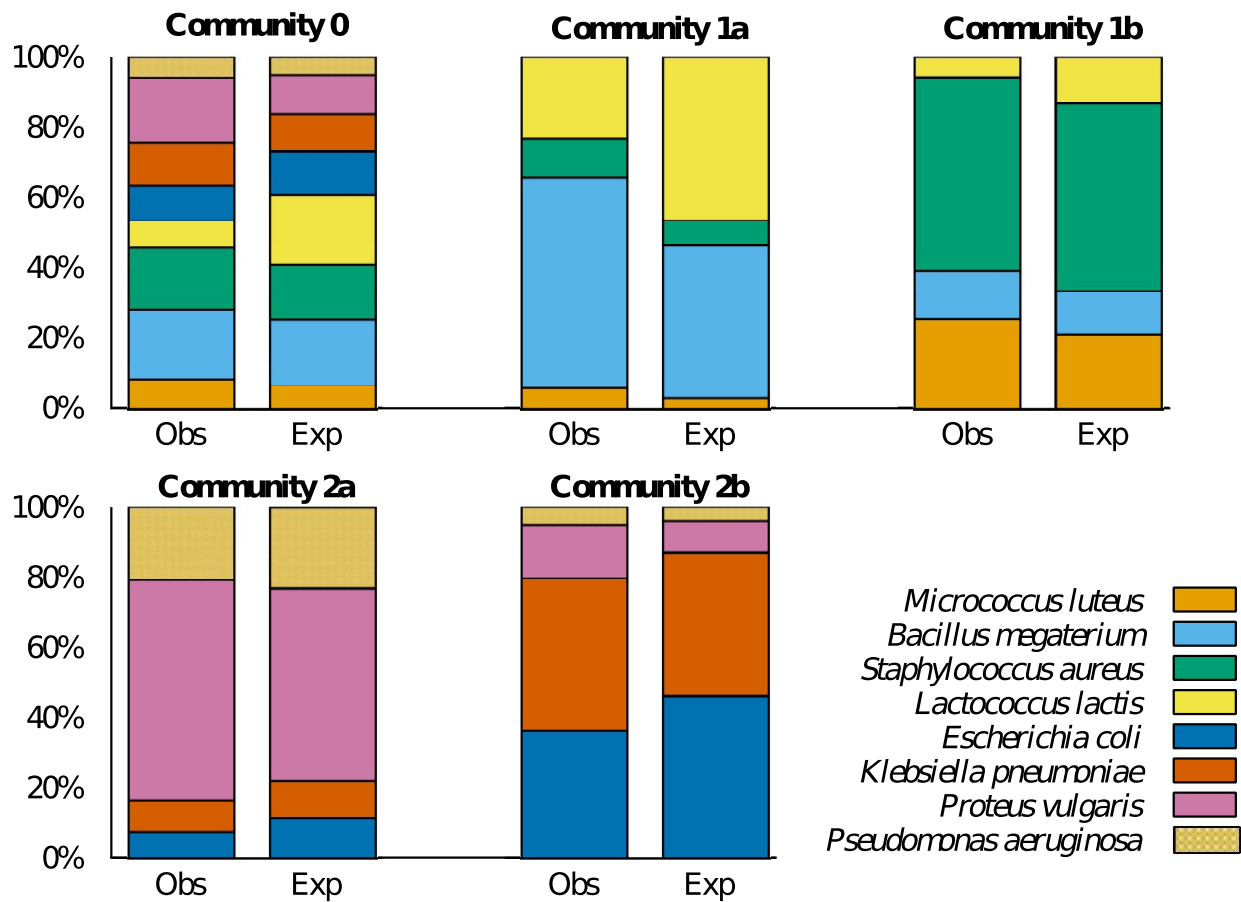
1009

1010 Figure 2: Effect of allowed low-quality reads on data retention and observed OTU diversity. The
1011 effect of adjusting (a) the r parameter (allowed low quality reads), in the `split_libraries_fastq.py`
1012 command in QIIME, on the percent of raw reads retained after data processing (Swarm d4 only),
1013 and (b) the observed number of OTUs in the mock community data compared to expected values
1014 (expressed as percent of total) across Swarm resolutions d1-d4.



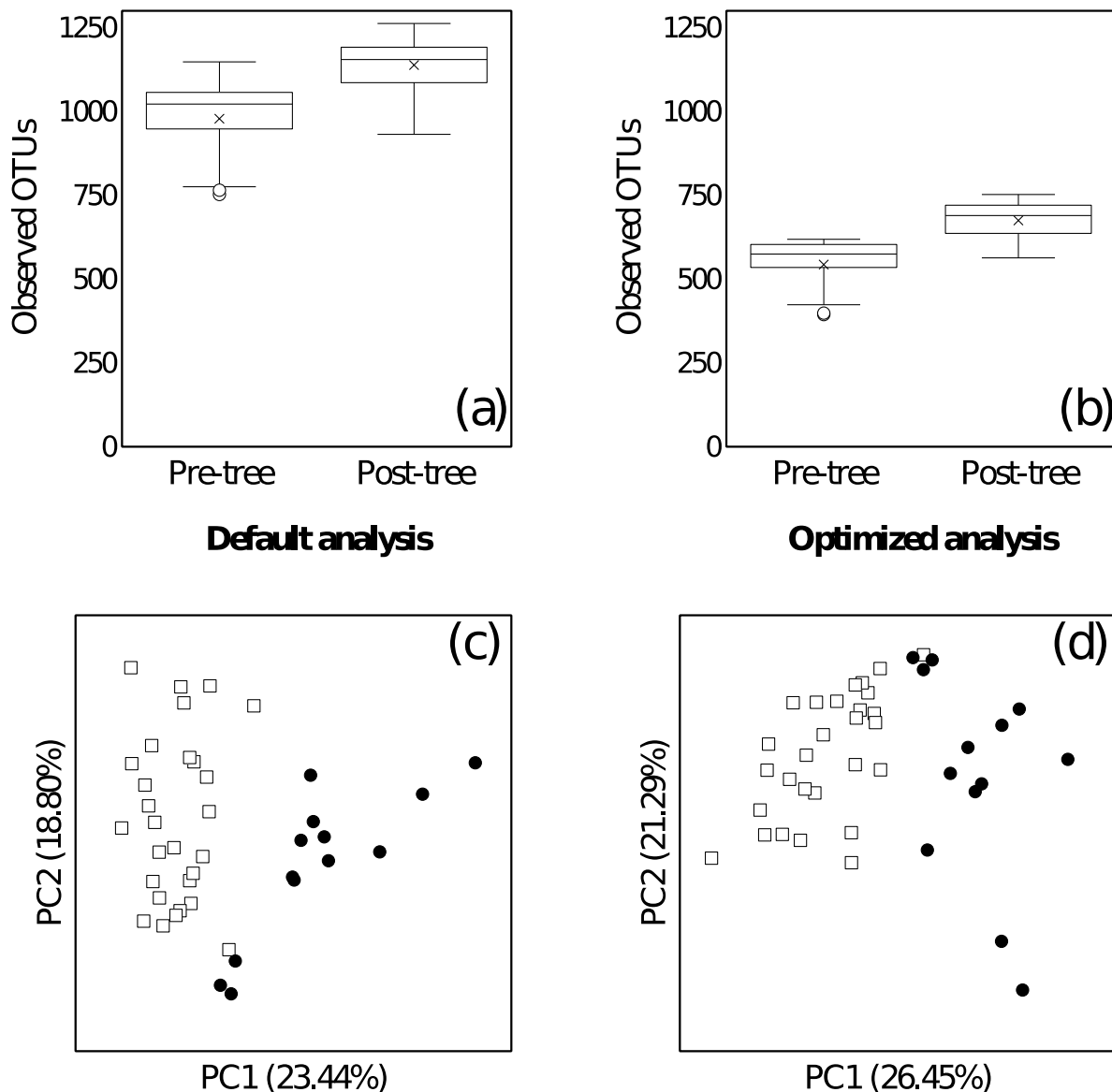
1015
1016
1017
1018
1019
1020
1021
1022
1023
1024
1025
1026
1027
1028
1029

1030 Figure 3: Observed versus expected mock community compositions. Observed mock community
 1031 compositions (Obs) plotted against expected values (Exp) for each separate community. Data
 1032 presented here are for the optimized analysis using the Kircher threshold for OTU table filtering,
 1033 and yielded a strong positive correlation to the expected result (Spearman's $r = 0.9646$, CI =
 1034 $[0.9495, 0.9752]$).



1035
 1036
 1037
 1038
 1039
 1040
 1041
 1042
 1043

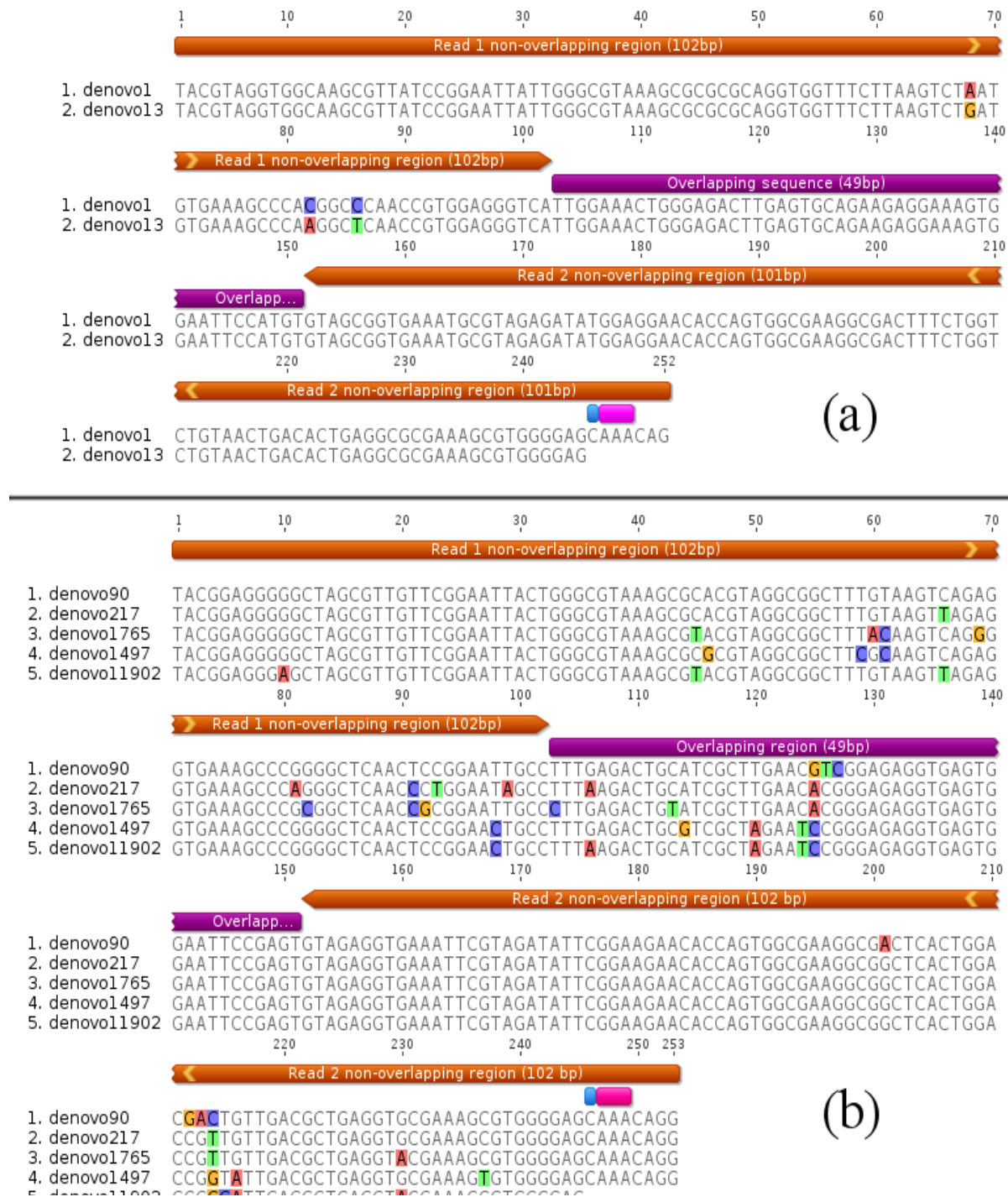
1044 Figure 4: Diversity analyses for environmental data rarefied to 5000 reads per sample. Alpha
1045 diversity as assessed by default (a) or optimized (b) workflow. Beta diversity assessed by default
1046 (c) or optimized (d) workflow (white squares: pre-tree; black circles: post-tree). All data show
1047 the same trends with strong statistical support (see text). Optimization has a strong effect on the
1048 interpretation of alpha diversity results while results are similar between workflows for beta
1049 diversity analyses



1050

1051

1052 Figure 5: Mafft alignments illustrating typical truncation of low-frequency OTUs for mock and
1053 environmental sequences. Mafft alignments for (a) two OTUs constituting mock community
1054 member *Bacillus megaterium*, and (b) five environmental OTUs classified to the family level as
1055 Sphingomonadaceae. Sequences are labelled with de novo OTU designations from the optimized
1056 workflow, and base positions are indicated above the sequence alignments. Highlighted bases
1057 indicate differences from the alignment consensus (not shown) and are colored according to
1058 identity. Overlapping and non-overlapping regions of the first and second reads are indicated
1059 above the alignments. Pink and blue positions indicate a “TTT” trimer preceding a “G”
1060 truncation position respectively. In each case, the bottom sequence represents a low-frequency
1061 OTU which was truncated during quality filtering.



1062