

# The impact of DNA polymerase and number of rounds of amplification in PCR on 16S rRNA gene sequence data

**Running title:** Quantifying the effects of PCR conditions

Marc A Sze<sup>1</sup> and Patrick D Schloss<sup>1†</sup>

† To whom correspondence should be addressed: [pschloss@umich.edu](mailto:pschloss@umich.edu)

<sup>1</sup> Department of Microbiology and Immunology, University of Michigan, Ann Arbor, MI

## 1 **Abstract**

2 PCR amplification of 16S rRNA genes is a critical, yet under appreciated step in the generation  
3 of sequence data to describe the taxonomic composition of microbial communities. Numerous  
4 factors in the design of PCR can impact the sequencing error rate, the abundance of chimeric  
5 sequences, and the degree to which the fragments in the product represent their abundance in  
6 the original sample (i.e. bias). We compared the performance of high fidelity polymerases and  
7 varying number of rounds of amplification when amplifying a mock community and human stool  
8 samples. Although it was impossible to derive specific recommendations, we did observe general  
9 trends. Namely, using a polymerase with the highest possible fidelity and minimizing the number  
10 of rounds of PCR reduced the sequencing error rate, fraction of chimeric sequences, and bias.  
11 Evidence of bias at the sequence level was subtle and could not be ascribed to the fragments'  
12 fraction of bases that were guanines or cytosines. When analyzing mock community data, the  
13 amount that the community deviated from the expected composition increased with rounds of PCR.  
14 This bias was inconsistent for human stool samples. Overall the results underscore the difficulty  
15 of comparing sequence data that are generated by different PCR protocols. However, the results  
16 indicate that the variation in human stool samples is generally larger than that introduced by the  
17 choice of polymerase or number of rounds of PCR.

## 18 **Importance**

19 A steep decline in sequencing costs drove an explosion in studies characterizing microbial  
20 communities from diverse environments. Although a significant amount of effort has gone into  
21 understanding the error profiles of DNA sequencers, little has been done to understand the  
22 downstream effects of the PCR amplification protocol. We quantified the effects of the choice of  
23 polymerase and number of PCR cycles on the quality of downstream data. We found that these  
24 choices can have a profound impact on the way that a microbial community is represented in the  
25 sequence data. The effects are relatively small compared to the variation in human stool samples,  
26 however, care should be taken to use polymerases with the highest possible fidelity and to minimize

27 the number of rounds of PCR. These results also underscore that it is not possible to directly  
28 compare sequence data generated under different PCR conditions.

## 29 Introduction

30 16S rRNA gene sequencing is a powerful and widely used tool for surveying the structure of  
31 microbial communities (1–3). This approach has exploded in popularity with advances in sequencing  
32 throughput such that it is now possible to characterize numerous samples with thousands of  
33 sequences per sample. Many factors can impact how a natural community is represented by  
34 the sequencing data including the method of acquiring samples (4–8), storage conditions (4–6,  
35 9–12), extraction methods (13), amplification conditions (8, 14, 15), sequencing method (15–17),  
36 and analytical pipeline (15, 18–20). The increased sampling depth that is now available relative  
37 to previous Sanger sequencing-based methods is expected to compound the impacts of an  
38 investigator’s choices and the interpretation of their results.

39 One step in the generation of 16S rRNA gene sequence data that has been long known to have  
40 a significant impact on the description of microbial communities is the choice of conditions for  
41 PCR amplification (8, 14, 15). Factors such as the choice of primers have an obvious impact on  
42 which populations will be amplified (18, 21). However, a variety of PCR artifacts can also impact  
43 the perception of a community including the formation of chimeras (14, 22–24), misincorporation  
44 of nucleotides (23, 25, 26), preferential amplification of some populations over others leading to  
45 bias (24, 27–33), and accumulation of random amplification events leading to PCR drift (24, 27,  
46 32, 34). Many bioinformatic tools have been developed to identify chimeras; however, there are  
47 significant sensitivity and specificity tradeoffs (14, 35). Laboratory-based solutions to minimize  
48 chimera formation have also been proposed such as minimizing the amount of template DNA  
49 in the PCR, minimizing the number of rounds of PCR, minimizing the amount of shearing in the  
50 template DNA, using DNA polymerases that have a proof-reading ability, and emulsion PCR (14,  
51 23, 36). Others have attempted to account for PCR bias using modeling approaches (29, 37). In  
52 cases where such modeling approaches have been successful, it has been with relatively small  
53 communities with consistent composition (29). To minimize PCR drift, some investigators pool  
54 technical replicate PCRs hoping to average out the drift (34). Other factors that have been shown to  
55 impact the formation of PCR artifacts are outside the control of an investigator including the fraction  
56 of DNA bases that are guanines or cytosines, the variation in the length of the targeted region

57 across the community, the sequence of the DNA that flanks the template, and the genetic diversity  
58 of the community (28, 30–33). Early investigations of the factors that lead to the formation of PCR  
59 artifacts focused on analyzing binary mixtures of genomic DNA and 16S rRNA gene fragments to  
60 explore PCR biases and chimera formation. Although these studies were instrumental in forcing  
61 researchers to be cautious about the interpretation of their results, we have a poor understanding  
62 of how these factors affect the formation of PCR artifacts in more complex communities.

63 The influence that the choice of DNA polymerase has on the formation of PCR artifacts has not  
64 been well studied. There has been recent interest in how the choice of the hypervariable region  
65 and data analysis pipelines impact the sequencing error rate (15, 18–20); however, these studies  
66 use the same DNA polymerase in the PCR step and implicitly assume that the rate of nucleotide  
67 misincorporation from PCR are significantly smaller than those from the sequencing phase. There  
68 has been more limited interest in the impact that DNA polymerase choice has on the formation of  
69 chimeras (23, 38). A recent study found differences in the number of OTUs and chimeras between  
70 normal and high fidelity DNA polymerases (38). The authors of the study reduced the difference  
71 between two polymerases by optimizing the annealing and extension steps within the PCR protocol  
72 (38). Yet this optimization was specific for the community they were analyzing (i.e. captive and  
73 semi-captive red-shanked doucs) and assumed that if the two polymerases generate the same  
74 community structure that the community structure was correct. In fact, the community structure  
75 generated by both methods was not free of artifacts, but likely had the same artifacts. A challenge  
76 in these types of experiments is having *a priori* knowledge of the true community representation.  
77 A mock community with known composition allows researchers to quantify the sequencing error  
78 rate, fraction of chimeras, and bias (19); however, mock communities have a limited phylogenetic  
79 diversity relative to natural communities. Natural communities, in contrast, have an unknown  
80 community composition making absolute measurements impossible. They can be used to validate  
81 results from mock communities and to understand the relative impacts of artifacts on the ability to  
82 differentiate biological and methodological sources of variation. Given the large number of DNA  
83 polymerases available to researchers, it is unlikely that a specific recommendation is possible.  
84 Rather, the development of general best practices and understanding the impact of PCR artifacts  
85 on an analysis are needed.

86 This study investigated the impact of choice of high fidelity DNA polymerase and the number of  
87 rounds of amplification on the formation of PCR artifacts using a mock community and human  
88 stool samples. It was hypothesized that additional rounds of PCR would exacerbate the number  
89 of artifacts. We tested (i) the effect of the polymerase on the error rate of the bases represented  
90 in the final sequences, (ii) the fraction of sequences that appeared to be chimeras and the ability  
91 to detect those chimeras, (iii) the bias of preferentially amplifying one fragment over another in a  
92 mixed pool of templates, and (iv) inter-sample variation in community structure of samples amplified  
93 with the same polymerase across the amplification process. To characterize these factors we  
94 sequenced a mock community of 8 organisms with known sequences and community structure  
95 and human fecal samples with unknown sequences and community structures. We sequenced the  
96 V4 region of the 16S rRNA genes from a mock community by generating paired 250 nt reads on  
97 the Illumina MiSeq platform. This region and sequencing approach was used because it has been  
98 shown to result in a relatively low sequencing error rate and is a widely used protocol (18). To better  
99 understand the impact of DNA polymerase choice on PCR artifacts, we selected five high fidelity  
100 DNA polymerases and amplified the communities using 20, 25, 30, and 35 rounds of amplification.  
101 Collectively, our results suggest that the number of rounds and to a lesser extent the choice of DNA  
102 polymerase used in PCR impact the sequence data. The effects are consistent and are smaller  
103 than the biological differences between individuals.

## 104 **Results**

### 105 ***Sequencing errors vary by the number of cycles and the DNA polymerase used in PCR.***

106 The presence of sequence errors can confound the ability to accurately classify 16S rRNA gene  
107 sequences and group sequences into operational taxonomic units (OTUs). More importantly,  
108 sequencing errors themselves can alter the representation of the community. Therefore, it is  
109 important to minimize the number of sequencing errors. Using a widely-used approach that  
110 generates the lowest reported error rate, we quantified the error rate by sequencing the V4 region of  
111 the 16S rRNA genes from an 8 member mock community. We also removed any contigs that were  
112 at least three bases more similar to a chimera of two references than to a single reference sequence  
113 (18, 19, 39). Regardless of the polymerase, the error rate increased with the number of rounds of  
114 amplification (Figure 1). Using 30 rounds of PCR is a common approach across diverse types of  
115 samples. Among the data generated using 30 rounds of PCR the Accuprime polymerase had the  
116 highest error rate (i.e. 0.124%) followed by the Platinum (i.e. 0.094%), Phusion (i.e. 0.064%), KAPA  
117 (i.e. 0.062%), and Q5 (i.e. 0.060%) polymerases (Figure 1). When we applied a pre-clustering  
118 denoising step, which merged the counts of reads within 2 nt of a more abundant sequence (19),  
119 the error rates dropped considerably such that the Platinum polymerase had the highest error rate  
120 (i.e. 0.014%) followed by the Accuprime (i.e. 0.012%), Q5 (i.e. 0.0053%), Phusion (i.e. 0.0049%),  
121 and KAPA (i.e. 0.0049%) polymerases (Figure 1). Although specific recommendations are difficult  
122 to make because the phylogenetic diversity of the initial DNA template is likely to have an impact  
123 on the results, it is clear that using as few PCR cycles as necessary and a polymerase with the  
124 lowest possible error rate is a good guide to minimizing the impact of polymerase on the error rate.

### 125 ***The fraction of sequences identified as being chimeric varies by the number of cycles 126 and the DNA polymerase used in PCR.***

127 Chimeric PCR products can significantly confound  
128 downstream analyses. Although numerous bioinformatic tools exist to identify and remove chimeric  
129 sequences with high specificity, their sensitivity is relatively low and can be reduced by the presence  
130 of sequencing errors (14, 35). Because the true sequences of the organisms in the mock community  
131 were known, we generated all possible chimeras between pairs of V4 16S rRNA gene fragments  
and used these possible chimeric sequences to screen the sequences generated under the different

132 PCR conditions to detect chimeras. The number of chimeras increased with rounds of amplification  
133 (Figure 2A). Interestingly, the fraction of chimeric sequences from the mock community varied by the  
134 type of polymerase used. After 30 rounds of PCR, the Platinum polymerase had the highest chimera  
135 rate (i.e. 18.2%) followed by the Q5 (i.e. 8.1%), Phusion (i.e. 7.5%), KAPA (i.e. 2.3%), and Accuprime  
136 (i.e. 0.9%) polymerases. To explore the characteristics of the chimeras further, we analyzed those  
137 chimeras formed after 35 cycles. Because of the uneven number of chimeras generated across  
138 the five polymerases, we subsampled the frequency of the chimeras to have the same number of  
139 chimeras per polymerase the Q5, Phusion, Accuprime, and Platinum polymerases; the chimeric  
140 sequence yield with the KAPA polymerase was significantly lower than the other polymerases and  
141 was omitted from our initial comparison. As has been shown previously (14), chimera formation was  
142 not random. Among the chimeras that were generated in mock community samples, 4.4% of the  
143 chimeras were found across all four polymerases. These chimeras represented between 67.6 and  
144 74.5% of the chimeras generated with each polymerase; they represented 40.4% of the chimeric  
145 sequences generated using the KAPA polymerase. These results indicate that the mechanisms  
146 leading to the formation of chimeras are largely independent of the properties of the polymerase,  
147 but are more likely due to the properties of the sequences.

148 Because our chimera screening procedure could only be applied to mock communities, we used  
149 the UCHIME algorithm to model the chimera screening approach that is used in most sequence  
150 curation pipelines. By comparing the output of UCHIME to our approach of screening for chimeras  
151 using all possible chimeras generated from the mock community sequences, we were able to  
152 calculate the UCHIME's sensitivity and specificity (Figure 2A). The specificity for all polymerases  
153 was above 95.4% and showed a weak association with the number of cycles used (Figure 2A).  
154 There was considerable inter-polymerase and inter-round of amplification variation in the sensitivity  
155 of UCHIME to detect the chimeras from the mock community. This suggested that the residual error  
156 rate after pre-clustering the sequence data did not compromise the sensitivity of UCHIME to detect  
157 chimeras. The sensitivity of UCHIME varied between 50 and 87.0% when at least 25 cycles were  
158 used. The generalizability of these results is limited because we used a single mock community  
159 with limited genetic diversity. Although we did not know the true chimera rate for our four human  
160 stool samples, we were able to calculate the fraction of sequences that UCHIME identified as being



161 chimeric (Figure 2B). These results followed those from the mock communities: additional rounds  
162 of amplification significantly increased the rate of chimeras and there was variation between the  
163 polymerases that we used. Although it was not possible to identify the features of a polymerase  
164 that resulted in higher rates of chimeras, it is clear that using the smallest number of PCR cycles  
165 possible will minimize the impact of chimeras.

166 **At the sequence level, PCR amplification bias is subtle.** Since researchers began using PCR  
167 to amplify 16S rRNA gene fragments there has been concern that amplifying fragments from  
168 a mixed template pool could lead to a biased representation in the pool of products and would  
169 confound downstream analyses (24, 27–33). The mock community was generated by mixing equal  
170 amounts of genomic DNA from 8 bacteria resulting in uneven representation of the *rrn* operons  
171 across the bacteria as each bacterium had a different genome size and varied in the number of  
172 operons in its genome. The vendor of the mock community subjects each lot of genomic DNA to  
173 shotgun sequencing to more accurately quantify the actual abundance of each organism in the  
174 community. It should be noted that this approach to quantifying abundance is also not without  
175 its own biases (40), but does provide an alternative approach to characterizing the structure of  
176 the mock community. We compared the vendor reported relative abundance of the 16S rRNA  
177 genes from each bacterium in the mock community to the data we generated across rounds of  
178 amplification and polymerase (Figure 3). Interestingly, for some bacteria, their representation  
179 became less biased with additional rounds of PCR (e.g. *L. fermentum*), while others became more  
180 biased (e.g. *E. faecalis*), and others had little change (e.g. *B. subtilis*). Contrary to prior reports  
181 (28), the percentage of bases in the V4 region that were guanines or cytosines was not predictive  
182 of the amount of bias. Across the strains there was no variation in the length of their V4 regions  
183 and they each had the same sequence in the region that the primers annealed. One of the bacteria  
184 represented in the mock community, *S. enterica*, had 6 identical copies of the V4 region and 1  
185 copy that differed from those by one nucleotide. The dominant copy had a thymidine and the rare  
186 copy had a guanine. We used the sequence data to calculate the ratio of the dominant to rare  
187 variants from *S. enterica* expecting a ratio near 6 (Figure S1). The Accuprime, Phusion, Platinum,  
188 and Q5 polymerases converged to a ratio of 5.4; however, the ratio for the KAPA polymerase was  
189 above 6 for all rounds of PCR (6.1-7.4) and the ratio for Q5 was below 6 for all rounds of PCR

190 (5.3-5.5). Given the subtle nature of the variation in the relative abundances of each 16S rRNA  
191 gene fragment, it was not possible to create generalizable rules that would explain the bias.

192 ***At the community level, the effects of PCR amplification bias grow with additional rounds***

193 ***of PCR.*** Because the variation in bias between polymerases and across rounds of PCR could be  
194 artificially inflated due to sequencing errors and chimeras, we analyzed the alpha and beta diversity  
195 of the mock community data at different phases of the sequence curation pipeline (Figure 4). First,  
196 we removed the chimeras from the mock community data as described above and mapped the  
197 individual reads to the OTUs that the 16S rRNA gene fragments would cluster into if there were no  
198 sequencing errors. This gave us a community distribution that reflected the distribution following  
199 PCR without any artifacts (Figure 4A; “No errors or chimeras”). Although the richness did not  
200 change, the Shannon diversity increased with the number of rounds of PCR for all polymerases  
201 except the KAPA polymerase, for which the diversity decreased. These data suggest that PCR  
202 had the effect of making the community distribution more even than it was originally, except for  
203 the data generated using the KAPA polymerase where the evenness decreased. Next, we used  
204 the observed sequence errors, but removed chimeras by comparing sequences to all possible  
205 chimeras between mock community sequences, and clustered the reads to OTUs (Figure 4A;  
206 “Residual errors, complete chimera removal”). The richness and diversity metrics trended higher  
207 with higher error rates and number of rounds of PCR. Finally, we used the observed sequence  
208 data and the UCHIME algorithm to identify chimeras (Figure 4A; “Residual errors, chimera removal  
209 with VSEARCH”). Again, the richness and diversity metrics trended higher with higher error rates  
210 and number of rounds of PCR. These comparisons demonstrated that although the bias at the  
211 sequence level was subtle, PCR introduces bias at the community level that is exacerbated by  
212 errors and chimeras when sequences are clustered into OTUs. When we measured the Bray-Curtis  
213 distance between the communities observed after 25 rounds of amplification and those at 30 and  
214 35, distances between 25 and 35 rounds were higher than between 25 and 30 rounds for each of  
215 the polymerases by an average of 0.022 units (Figure 4B). The Platinum polymerase varied the  
216 most across rounds of amplification (25 vs 30 rounds: 0.13; 25 vs 35 rounds: 0.15). For any number  
217 of cycles, the median Bray-Curtis distance between polymerases ranged between 0.074 and 0.11.  
218 Although the distances between samples were small, the ordination of these distances showed a

219 clear change in community structure with increasing rounds of PCR (Figure 4C). This observation  
220 was supported by our statistical analysis, which revealed that the effect of the number of rounds of  
221 PCR ( $R^2=0.21$ ,  $P<0.001$ ) was comparable to the choice of polymerase ( $R^2=0.20$ ,  $P<0.001$ ). These  
222 results demonstrate that subtle differences in relative abundances can have an impact on overall  
223 community structure. This variation underscores the importance of only comparing sequence data  
224 that have been generated using the same PCR conditions.

225 ***The choice of polymerase or the number of rounds of amplification have little impact on the***  
226 ***relative interpretation of community-wide metrics of diversity.*** We expected that the biases  
227 that we observed at the population and community levels using mock community data would  
228 be small relative to the expected differences between biological samples. To study this further,  
229 we calculated alpha and beta-diversity metrics using the human stool samples for each of the  
230 polymerases and rounds of amplification. We calculated the number of observed OTUs and  
231 Shannon diversity for each condition and stool sample (Figure 5A). Although there were clear  
232 differences between PCR conditions, the relative ordering of the stool samples did not meaningfully  
233 vary across conditions. When we characterized the variation between rounds of amplification  
234 using human stool samples, the distance between the 25 and 30 rounds and 25 and 35 rounds  
235 varied considerably between samples and polymerases (Figure 5B). In general the inter-round  
236 variation was lowest for the data generated using the KAPA and Accuprime polymerases. The data  
237 generated using the Platinum polymerase was consistent across rounds, but overall, it was more  
238 biased than the other polymerases. Considering the average distance across the four samples  
239 varied between 0.39 and 0.56, regardless of the polymerases and number of rounds of amplification,  
240 any bias due to amplification is unlikely to obscure community-wide differences between samples.  
241 In support of this was our principle coordinates analysis of the Bray-Curtis distances, which revealed  
242 distinct clusters by stool sample (Figure 5C). Within each cluster there were no obvious patterns  
243 related to the polymerase or number of rounds of PCR. Our statistical analysis revealed statistically  
244 significant differences in the community structures with the stool sample explaining the most  
245 variation ( $R^2=0.79$ ,  $P<0.001$ ), followed by the number of rounds of PCR ( $R^2=0.044$ ,  $P<0.001$ ) and  
246 the choice of polymerase ( $R^2=0.033$ ,  $P<0.001$ ). Together, these results indicate that for a coarse  
247 analysis of communities, the choice of number of rounds of amplification or polymerase are not

248 important, but that they must be consistent across samples. It is difficult to develop a specific  
249 recommendation based on the level of bias across rounds of PCR or polymerases; however, the  
250 general suggestion is to use as few rounds of amplification as possible.

251 ***There is little evidence of a relationship between polymerase or number of rounds of***  
252 ***amplification on PCR drift.*** There have been concerns that the same template DNA subjected  
253 to the same PCR conditions could result in different representations of communities because of  
254 random drift over the course of PCR. To test this, we determined the average Bray-Curtis distance  
255 between replicate reactions using the same polymerase and number of rounds of amplification  
256 (Figure 6). Using the mock community data there were no obvious trends. The average Bray-Curtis  
257 distance within a set of conditions varied by 0.062 to 0.11 units. Although we did not generate  
258 technical replicates of each of the stool samples, the inter-sample variation for each set of  
259 conditions was consistent and varied between 0.50 and 0.56 units. These data suggest that  
260 amplicon sequencing is robust to random variation in amplification and that any differences are  
261 likely to be smaller than what is considered biologically relevant.

## 262 Discussion

263 Our results suggest that the number of rounds of PCR and to a lesser degree the choice of DNA  
264 polymerase impact the analysis of 16S rRNA gene sequence data from bacterial communities.  
265 Although it was not possible to make direct connections between PCR conditions and specific  
266 sources of bias, we were able to identify general recommendations that reduce the amount of  
267 error, chimera formation, and bias. Researchers should strive to minimize the number of rounds  
268 of PCR and should use a high fidelity polymerase. Although specific PCR conditions impact the  
269 precise interpretation of the data, the effects were consistent and were smaller than the biological  
270 differences between the samples we tested. Based on these observations, amplicons must be  
271 generated by consistent protocols to yield meaningful comparisons. When comparing across  
272 studies, values like richness, diversity, and relative abundances must be made in relative and not  
273 absolute terms. Furthermore, care must be taken to not directly compare or pool samples from  
274 different studies. Instead, it is important to statistically model the study-based variation as has been  
275 done in recent meta-analyses that compared relative effect sizes or pooled data using a mixed  
276 effects statistical model (41, 42).

277 The observed sequencing error rates and alpha diversity metrics followed the manufacturers'  
278 measurements of their polymerases' fidelity (Figure 1). Accuprime and Platinum have fidelity that  
279 are approximately 10-times higher than that of Taq whereas the fidelity of Phusion, Q5, and KAPA  
280 are more than 100 times higher. Among these polymerases, the KAPA polymerase consistently  
281 resulted in a lower error rate, lower chimera rate, and lower bias across rounds of PCR for the mock  
282 community samples. Furthermore, among the human samples, the KAPA polymerase consistently  
283 had the lowest detected chimera rate and inter-cycle bias. These benefits were most accentuated  
284 at 35 cycles. However, in our experience and despite efforts to optimize the yield with the KAPA  
285 polymerase, the reactions typically had a high proportion of primer-dimer products and low yield of  
286 correctly-sized products. Although the error rate with the Accuprime polymerase was not as low as  
287 that with KAPA, we consider it to be an acceptable alternative. Considering polymerase development  
288 is an active area of commercial development with potential new polymerases becoming available,  
289 it is important for researchers to understand how changing the polymerase impacts downstream

290 analyses for their type of samples.

291 Over the past 20 years, a large literature has attempted to document various PCR biases and  
292 underscored the fact that data based on amplification of DNA from a mixed community are not a true  
293 representation of the actual community. In addition to obvious biases imposed by primer selection,  
294 other factors inherent in PCR can influence the representation of communities. Factors that can  
295 lead to preferential amplification of one fragment over another have included guanine and cytosine  
296 composition, length, flanking DNA composition, amount of DNA shearing, and number of rounds of  
297 PCR (24, 27–33). These factors may become exacerbated if PCR is performed on multiple samples  
298 that vary in their concentration (43). In addition, environmental and reagent contaminants can also  
299 have a significant impact on the analysis of low biomass samples (44). Less well understood is the  
300 effect of phylogenetic diversity on bias and chimera formation. Communities with low phylogenetic  
301 diversity may be more prone to chimera formation since chimeras are more likely to form among  
302 closely related sequences (14, 35). The interaction of these various influences on PCR artifacts are  
303 complex and difficult to tease apart. Minimizing the level of DNA shearing, controlling for template  
304 concentration across samples, and using the fewest number of rounds of PCR with a polymerase  
305 that has the highest possible fidelity are strategies that can be employed to minimize the formation  
306 of chimeras. Although care should always be taken when choosing a polymerase for 16S rRNA  
307 gene sequencing, our observations show that variation among polymerases is smaller than the  
308 actual biological variation in fecal communities between individuals.

309 Even with these strategies it is impossible to remove all PCR artifacts. Beyond the imperfections of  
310 the best polymerases, sometimes difficult to lyse organisms require stringent lysis steps and low  
311 biomass samples require additional rounds of PCR. A host of bioinformatics tools are available for  
312 removing residual sequencing errors (18, 45–47). Other tools are available for removing chimeras  
313 (14, 35) where there is a trade off between the sensitivity of detecting chimeras and the specificity  
314 of correctly calling a sequence a chimera. In recent years, parameters for these algorithms have  
315 been changed to increase their sensitivity with little evaluation of the effects on the specificity of  
316 the algorithms (45, 47). Others recommend removing any read that has an abundance below a  
317 specified threshold as a tool to remove PCR and sequencing artifacts (e.g. removing all sequences  
318 that only appear once) (20, 45–47). This method must be approached with caution as such

319 approaches are likely to introduce a different bias of the community representation and ignore the  
320 fact, as we showed, that artifacts may be quite abundant and reproducible. Ultimately, researchers  
321 must test their hypotheses with multiple methods to validate the claims they reach with any one  
322 method (48). All methods have biases and limitations and we must use complementary methods to  
323 obtain robust results.

## 324 **Materials & Methods**

325 **Mock community.** The ZymoBIOMICS™ Microbial Community DNA Standard (Zymo, CA, USA)  
326 was used for mock communities and the bacterial component was made up of *Pseudomonas*  
327 *aeruginosa*, *Escherichia coli*, *Salmonella enterica*, *Lactobacillus fermentum*, *Enterococcus*  
328 *faecalis*, *Staphylococcus aureus*, *Listeria monocytogenes*, and *Bacillus subtilis* at equal genomic  
329 DNA abundance ([https://web.archive.org/web/20171217151108/http://www.zymoresearch.com:](https://web.archive.org/web/20171217151108/http://www.zymoresearch.com:80/microbiomics/microbial-standards/zymbiomics-microbial-community-standards)  
330 [80/microbiomics/microbial-standards/zymbiomics-microbial-community-standards](https://web.archive.org/web/20171217151108/http://www.zymoresearch.com:80/microbiomics/microbial-standards/zymbiomics-microbial-community-standards)). The actual  
331 relative abundance for each bacterium was obtained from Zymo's certificate of analysis for the  
332 lot (Lot: ZRC187325), which they determined using shotgun metagenomic sequencing ([https:](https://github.com/SchlossLab/Sze_PCRSeqEffects_mSphere_2019/data/references/ZRC187325.pdf)  
333 [//github.com/SchlossLab/Sze\\_PCRSeqEffects\\_mSphere\\_2019/data/references/ZRC187325.pdf](https://github.com/SchlossLab/Sze_PCRSeqEffects_mSphere_2019/data/references/ZRC187325.pdf)).

334 **Human samples.** Fecal samples were obtained from 4 individuals who were part of an earlier  
335 study (49). These samples were collected using a protocol approved by the University of Michigan  
336 Institutional Review Board. For this study, the samples were de-identified. DNA was extracted from  
337 the fecal samples using the MOBIO™ PowerMag Microbiome RNA/DNA extraction kit (now Qiagen,  
338 MD, USA).

339 **PCR protocol.** Five high fidelity DNA polymerases were tested including AccuPrime™  
340 (ThermoFisher, MA, USA), KAPA HIFI (Roche, IN, USA), Phusion (New England Biolabs, MA,  
341 USA), Platinum (ThermoFisher, MA, USA), and Q5 (New England Biolabs, MA, USA). Manufacturer  
342 recommendations were followed except for the annealing and extension times, which were  
343 selected based on previously published protocols (18, 38). Primers targeting the V4 region of  
344 the 16S rRNA gene were used with modifications to generate MiSeq amplicon libraries (18)  
345 ([https://github.com/SchlossLab/MiSeq\\_WetLab\\_SOP/](https://github.com/SchlossLab/MiSeq_WetLab_SOP/)). The 16S rRNA gene targeting regions of  
346 the primers annealed to *E. coli* positions 515 to 533 (GTGCCAGCMGCCGCGGTAA) and 787 to  
347 806 (GGACTACHVGGGTWTCTAAT). The number of rounds of PCR used for each sample and  
348 polymerase started at 15 and increased by 5 rounds up to 35 cycles. Insufficient PCR product was  
349 generated using 15 rounds and has not been included in our analysis.

350 **Library generation and sequencing.** Each PCR condition (i.e. combination of polymerase and



351 number of rounds of PCR) were replicated four times for the mock community and one time for each  
352 fecal sample. Libraries were generated as previously described (18) ([https://github.com/SchlossLab/  
353 MiSeq\\_WetLab\\_SOP/](https://github.com/SchlossLab/MiSeq_WetLab_SOP/)). The libraries were sequenced using the Illumina MiSeq sequencing platform  
354 to generate paired 250-nt reads.

355 **Sequence processing.** The mothur software program (v 1.41) was used for all sequence  
356 processing steps (50). The protocol has been previously published (18) ([https://www.mothur.org/  
357 wiki/MiSeq\\_SOP](https://www.mothur.org/wiki/MiSeq_SOP)). Briefly, paired reads were assembled using mothur's make.contigs command to  
358 correct errors introduced by sequencing (18). Any assembled contigs that contained an ambiguous  
359 base call, mapped to the incorrect region of the 16S rRNA gene, or appeared to be a contaminant  
360 were removed from subsequent analyses. Sequences were further denoised using mothur's  
361 pre.cluster command to merge the counts of sequences that were within 2 nt of a more abundant  
362 sequence. The VSEARCH implementation of UCHIME was used to screen for chimeras (35, 51).  
363 At various stages in the sequence processing pipeline for the mock community data, the mothur  
364 seq.error command was used to quantify the sequencing error rate as well as the true chimera  
365 rate. This command uses the true sequences from the mock community to generate all possible  
366 chimeras and removes any contigs that were at least three bases more similar to a chimera than to  
367 a reference sequence. The command then counts the number of substitutions, insertions, and  
368 deletions in the contig relative to the reference sequence and reports the error rate without the  
369 inclusion of chimeric sequences (19). UCHIME's sensitivity was calculated as the percentage of  
370 true chimeras that were detected as chimeras when using UCHIME. Its specificity was calculated  
371 as the percentage of non-chimeric sequences that were detected as being non-chimeric by  
372 UCHIME. The reference sequences and *rrn* operon copy number for each bacterium were  
373 obtained from the ZymoBIOMICS™ Microbial Community DNA Standard protocol ([https://  
374 //web.archive.org/web/20181221151905/https://www.zymoresearch.com/media/amasty/amfile/  
375 attach/\\_D6305\\_D6306\\_ZymoBIOMICS\\_Microbial\\_Community\\_DNA\\_Standard\\_v1.1.3.pdf](https://web.archive.org/web/20181221151905/https://www.zymoresearch.com/media/amasty/amfile/attach/_D6305_D6306_ZymoBIOMICS_Microbial_Community_DNA_Standard_v1.1.3.pdf)).  
376 Sequences were assigned to operational taxonomic units (OTUs) at a threshold of 3% dissimilarity  
377 using the OptiClust algorithm (52). To adjust for unequal sequencing when measuring alpha and  
378 beta diversity, all samples were rarefied for downstream analysis. The Good's coverage for the  
379 samples was routinely greater than 95%.

380 **Statistical analysis.** All analysis was done with the R (v 3.5.1) software package (53). Data  
381 transformation and graphing were completed using the tidyverse package (v 1.2.1). The distance  
382 matrix data was analyzed using the adonis function within the vegan package (v 2.5.4).

383 **Reproducible methods.** The data analysis code for this study can be found at [https://github.com/](https://github.com/SchlossLab/Sze_PCRSeqEffects_mSphere_2019)  
384 SchlossLab/Sze\_PCRSeqEffects\_mSphere\_2019. The raw sequences are available at the SRA  
385 (Accession SRP132931).

### 386 **Acknowledgements**

387 We appreciate the willingness of the donors to provide stool samples. We also thank Judy Opp  
388 and April Cockburn for their assistance in sequencing the samples as part of the Microbiome Core  
389 Facility at the University of Michigan. Additional thanks to members of the Schloss lab and Dr. Marcy  
390 Balunas for reading earlier drafts of the manuscript and providing helpful critiques. Support for MAS  
391 came from the Canadian Institute of Health Research and NIH grant UL1TR002240 and support for  
392 PDS came from NIH grants P30DK034933, R01CA215574, and U19AI09087.

## 393 References

- 394 1. **Gilbert JA, Jansson JK, Knight R.** 2018. Earth microbiome project and global systems biology.  
395 *mSystems* **3**. doi:10.1128/msystems.00217-17.
- 396 2. **Human Microbiome Consortium.** 2012. Structure, function and diversity of the healthy human  
397 microbiome. *Nature* **486**:207–214. doi:10.1038/nature11234.
- 398 3. **Schloss PD, Girard RA, Martin T, Edwards J, Thrash JC.** 2016. Status of the archaeal and  
399 bacterial census: An update. *mBio* **7**. doi:10.1128/mbio.00201-16.
- 400 4. **Luo T, Srinivasan U, Ramadugu K, Shedden KA, Neiswanger K, Trumble E, Li JJ, McNeil**  
401 **DW, Crout RJ, Weyant RJ, Marazita ML, Foxman B.** 2016. Effects of specimen collection  
402 methodologies and storage conditions on the short-term stability of oral microbiome taxonomy.  
403 *Applied and Environmental Microbiology* **82**:5519–5529. doi:10.1128/aem.01132-16.
- 404 5. **Bassis CM, Nicholas M. Moore, Lolans K, Seekatz AM, Weinstein RA, Young VB, Hayden**  
405 **MK.** 2017. Comparison of stool versus rectal swab samples and storage conditions on bacterial  
406 community profiles. *BMC Microbiology* **17**. doi:10.1186/s12866-017-0983-9.
- 407 6. **Gorzelak MA, Gill SK, Tasnim N, Ahmadi-Vand Z, Jay M, Gibson DL.** 2015. Methods for  
408 improving human gut microbiome data by reducing variability through sample processing and  
409 storage of stool. *PLOS ONE* **10**:e0134802. doi:10.1371/journal.pone.0134802.
- 410 7. **Dominianni C, Wu J, Hayes RB, Ahn J.** 2014. Comparison of methods for fecal microbiome  
411 biospecimen collection. *BMC Microbiology* **14**:103. doi:10.1186/1471-2180-14-103.
- 412 8. **Santos QMB-d los, Schroeder JL, Blakemore O, Moses J, Haffey M, Sloan W, Pinto**  
413 **AJ.** 2016. The impact of sampling, PCR, and sequencing replication on discerning changes  
414 in drinking water bacterial community over diurnal time-scales. *Water Research* **90**:216–224.  
415 doi:10.1016/j.watres.2015.12.010.
- 416 9. **Sinha R, Chen J, Amir A, Vogtmann E, Shi J, Inman KS, Flores R, Sampson J, Knight R,**  
417 **Chia N.** 2015. Collecting fecal samples for microbiome analyses in epidemiology studies. *Cancer*

418 Epidemiology Biomarkers & Prevention **25**:407–416. doi:10.1158/1055-9965.epi-15-0951.

419 **10. Amir A, McDonald D, Navas-Molina JA, Debelius J, Morton JT, Hyde E, Robbins-Pianka**  
420 **A, Knight R.** 2017. Correcting for microbial blooms in fecal samples during room-temperature  
421 shipping. *mSystems* **2**:e00199–16. doi:10.1128/msystems.00199-16.

422 **11. Lauber CL, Zhou N, Gordon JI, Knight R, Fierer N.** 2010. Effect of storage conditions on  
423 the assessment of bacterial community structure in soil and human-associated samples. *FEMS*  
424 *Microbiology Letters* **307**:80–86. doi:10.1111/j.1574-6968.2010.01965.x.

425 **12. Song SJ, Amir A, Metcalf JL, Amato KR, Xu ZZ, Humphrey G, Knight R.** 2016. Preservation  
426 methods differ in fecal microbiome stability, affecting suitability for field studies. *mSystems*  
427 **1**:e00021–16. doi:10.1128/msystems.00021-16.

428 **13. Costea PI, Zeller G, Sunagawa S, Pelletier E, Alberti A, Levenez F, Tramontano M,**  
429 **Driessen M, Hercog R, Jung F-E, Kultima JR, Hayward MR, Coelho LP, Allen-Vercoe E,**  
430 **Bertrand L, Blaut M, Brown JRM, Carton T, Cools-Portier S, Daigneault M, Derrien M,**  
431 **Druesne A, Vos WM de, Finlay BB, Flint HJ, Guarner F, Hattori M, Heilig H, Luna RA,**  
432 **Hylckama Vlieg J van, Junick J, Klymiuk I, Langella P, Chatelier EL, Mai V, Manichanh C,**  
433 **Martin JC, Mery C, Morita H, O'Toole PW, Orvain C, Patil KR, Penders J, Persson S, Pons N,**  
434 **Popova M, Salonen A, Saulnier D, Scott KP, Singh B, Slezak K, Veiga P, Versalovic J, Zhao**  
435 **L, Zoetendal EG, Ehrlich SD, Dore J, Bork P.** 2017. Towards standards for human fecal sample  
436 processing in metagenomic studies. *Nature Biotechnology*. doi:10.1038/nbt.3960.

437 **14. Haas BJ, Gevers D, Earl AM, Feldgarden M, Ward DV, Giannoukos G, Ciulla D, Tabbaa D,**  
438 **Highlander SK, Sodergren E, Methe B, DeSantis TZ, Petrosino JF, Knight R, and BWB.** 2011.  
439 Chimeric 16S rRNA sequence formation and detection in sanger and 454-pyrosequenced PCR  
440 amplicons. *Genome Research* **21**:494–504. doi:10.1101/gr.112730.110.

441 **15. Sinha R, Abu-Ali G, Vogtmann E, Fodor AA, Ren B, Amir A, Schwager E, Crabtree J, Ma**  
442 **S, Abnet CC, Knight R, White O, Huttenhower C.** 2017. Assessment of variation in microbial  
443 community amplicon sequencing by the microbiome quality control (MBQC) project consortium.  
444 *Nature Biotechnology*. doi:10.1038/nbt.3981.

- 445 16. **Meisel JS, Hannigan GD, Tyldsley AS, SanMiguel AJ, Hodkinson BP, Zheng Q, Grice**  
446 **EA.** 2016. Skin microbiome surveys are strongly influenced by experimental design. *Journal of*  
447 *Investigative Dermatology* **136**:947–956. doi:10.1016/j.jid.2016.01.016.
- 448 17. **Caporaso JG, Lauber CL, Walters WA, Berg-Lyons D, Lozupone CA, Turnbaugh PJ,**  
449 **Fierer N, Knight R.** 2010. Global patterns of 16S rRNA diversity at a depth of millions of  
450 sequences per sample. *Proceedings of the National Academy of Sciences* **108**:4516–4522.  
451 doi:10.1073/pnas.1000080107.
- 452 18. **Kozich JJ, Westcott SL, Baxter NT, Highlander SK, Schloss PD.** 2013. Development of a  
453 dual-index sequencing strategy and curation pipeline for analyzing amplicon sequence data on  
454 the MiSeq illumina sequencing platform. *Applied and Environmental Microbiology* **79**:5112–5120.  
455 doi:10.1128/aem.01043-13.
- 456 19. **Schloss PD, Gevers D, Westcott SL.** 2011. Reducing the effects of PCR amplification and  
457 sequencing artifacts on 16S rRNA-based studies. *PLoS ONE* **6**:e27310. doi:10.1371/journal.pone.0027310.
- 458 20. **Bokulich NA, Subramanian S, Faith JJ, Gevers D, Gordon JI, Knight R, Mills DA,**  
459 **Caporaso JG.** 2012. Quality-filtering vastly improves diversity estimates from illumina amplicon  
460 sequencing. *Nature Methods* **10**:57–59. doi:10.1038/nmeth.2276.
- 461 21. **Parada AE, Needham DM, Fuhrman JA.** 2015. Every base matters: Assessing small subunit  
462 rRNA primers for marine microbiomes with mock communities, time series and global field samples.  
463 *Environmental Microbiology* **18**:1403–1414. doi:10.1111/1462-2920.13023.
- 464 22. **Wang GCY, Wang Y.** 1996. The frequency of chimeric molecules as a consequence of PCR  
465 co-amplification of 16S rRNA genes from different bacterial species. *Microbiology* **142**:1107–1114.  
466 doi:10.1099/13500872-142-5-1107.
- 467 23. **Potapov V, Ong JL.** 2017. Examining sources of error in PCR by single-molecule sequencing.  
468 *PLOS ONE* **12**:e0169774. doi:10.1371/journal.pone.0169774.
- 469 24. **Kebschull JM, Zador AM.** 2015. Sources of PCR-induced distortions in high-throughput  
470 sequencing data sets. *Nucleic Acids Research* gkv717. doi:10.1093/nar/gkv717.

- 471 25. **McInerney P, Adams P, Hadi MZ.** 2014. Error rate comparison during polymerase chain  
472 reaction by DNA polymerase. *Molecular Biology International* **2014**:1–8. doi:10.1155/2014/287430.
- 473 26. **Cline J.** 1996. PCR fidelity of pfu DNA polymerase and other thermostable DNA polymerases.  
474 *Nucleic Acids Research* **24**:3546–3551. doi:10.1093/nar/24.18.3546.
- 475 27. **Acinas SG, Sarma-Rupavtarm R, Klepac-Ceraj V, Polz MF.** 2005. PCR-induced  
476 sequence artifacts and bias: Insights from comparison of two 16S rRNA clone libraries  
477 constructed from the same sample. *Applied and Environmental Microbiology* **71**:8966–8969.  
478 doi:10.1128/aem.71.12.8966-8969.2005.
- 479 28. **Polz MF, Cavanaugh CM.** 1998. Bias in template-to-product ratios in multitemplate PCR.  
480 *Applied and Environmental Microbiology* **64**:3724–3730.
- 481 29. **Brooks JP, David J Edwards, Harwich MD, Rivera MC, Fettweis JM, Serrano MG, Reis**  
482 **RA, Sheth NU, Huang B, Girerd P, Strauss JF, Jefferson KK, Buck GA.** 2015. The truth about  
483 metagenomics: Quantifying and counteracting bias in 16S rRNA studies. *BMC Microbiology* **15**.  
484 doi:10.1186/s12866-015-0351-6.
- 485 30. **Suzuki MT, Giovannoni SJ.** 1996. Bias caused by template annealing in the amplification of  
486 mixtures of 16S rRNA genes by PCR. *Applied and environmental microbiology* **62**:625–630.
- 487 31. **Chandler D, Fredrickson J, Brockman F.** 1997. Effect of pcr template concentration on  
488 the composition and distribution of total community 16S rDNA clone libraries. *Molecular Ecology*  
489 **6**:475–482.
- 490 32. **Wagner A, Blackstone N, Cartwright P, Dick M, Misof B, Snow P, Wagner GP, Bartels J,**  
491 **Murtha M, Pendleton J.** 1994. Surveys of gene families using polymerase chain reaction: PCR  
492 selection and pcr drift. *Systematic Biology* **43**:250–261.
- 493 33. **Hansen MC, Tolker-Nielsen T, Givskov M, Molin S.** 1998. Biased 16S rDNA pcr amplification  
494 caused by interference from dna flanking the template region. *FEMS Microbiology Ecology*  
495 **26**:141–149.

- 496 34. **Kennedy K, Hall MW, Lynch MDJ, Moreno-Hagelsieb G, Neufeld JD.** 2014. Evaluating  
497 bias of illumina-based bacterial 16S rRNA gene profiles. *Applied and Environmental Microbiology*  
498 **80**:5717–5722. doi:10.1128/aem.01451-14.
- 499 35. **Edgar RC, Haas BJ, Clemente JC, Quince C, Knight R.** 2011. UCHIME improves sensitivity  
500 and speed of chimera detection. *Bioinformatics* **27**:2194–2200. doi:10.1093/bioinformatics/btr381.
- 501 36. **Williams R, Peisajovich SG, Miller OJ, Magdassi S, Tawfik DS, Griffiths AD.** 2006.  
502 Amplification of complex gene libraries by emulsion PCR. *Nature Methods* **3**:545–550.  
503 doi:10.1038/nmeth896.
- 504 37. **Edgar RC.** 2017. UNBIAS: An attempt to correct abundance bias in 16S sequencing, with  
505 limited success. doi:10.1101/124149.
- 506 38. **Gohl DM, Vangay P, Garbe J, MacLean A, Hauge A, Becker A, Gould TJ, Clayton JB,**  
507 **Johnson TJ, Hunter R, Knights D, Beckman KB.** 2016. Systematic improvement of amplicon  
508 marker gene methods for increased accuracy in microbiome studies. *Nature Biotechnology*  
509 **34**:942–949. doi:10.1038/nbt.3601.
- 510 39. **Quince C, Lanzen A, Davenport RJ, Turnbaugh PJ.** 2011. Removing noise from  
511 pyrosequenced amplicons. *BMC Bioinformatics* **12**:38. doi:10.1186/1471-2105-12-38.
- 512 40. **Nayfach S, Pollard KS.** 2016. Toward accurate and quantitative comparative metagenomics.  
513 *Cell* **166**:1103–1116. doi:10.1016/j.cell.2016.08.007.
- 514 41. **Sze MA, Schloss PD.** 2016. Looking for a signal in the noise: Revisiting obesity and the  
515 microbiome. *mBio* **7**. doi:10.1128/mbio.01018-16.
- 516 42. **Sze MA, Schloss PD.** 2018. Leveraging existing 16S rRNA gene surveys to identify  
517 reproducible biomarkers in individuals with colorectal tumors. *mBio* **9**. doi:10.1128/mbio.00630-18.
- 518 43. **Multinu F, Harrington SC, Chen J, Jeraldo PR, Johnson S, Chia N, Walther-Antonio MR.**  
519 **2018.** Systematic bias introduced by genomic DNA template dilution in 16S rRNA gene-targeted  
520 microbiota profiling in human stool homogenates. *mSphere* **3**. doi:10.1128/msphere.00560-17.



- 521 44. **Salter SJ, Cox MJ, Turek EM, Calus ST, Cookson WO, Moffatt MF, Turner P, Parkhill**  
522 **J, Loman NJ, Walker AW.** 2014. Reagent and laboratory contamination can critically impact  
523 sequence-based microbiome analyses. *BMC Biology* **12**. doi:10.1186/s12915-014-0087-z.
- 524 45. **Callahan BJ, McMurdie PJ, Rosen MJ, Han AW, Johnson AJA, Holmes SP.** 2016. DADA2:  
525 High-resolution sample inference from illumina amplicon data. *Nature Methods* **13**:581–583.  
526 doi:10.1038/nmeth.3869.
- 527 46. **Edgar RC.** 2016. UNOISE2: Improved error-correction for illumina 16S and ITS amplicon  
528 sequencing. doi:10.1101/081257.
- 529 47. **Amir A, McDonald D, Navas-Molina JA, Kopylova E, Morton JT, Xu ZZ, Kightley EP,**  
530 **Thompson LR, Hyde ER, Gonzalez A, Knight R.** 2017. Deblur rapidly resolves single-nucleotide  
531 community sequence patterns. *mSystems* **2**. doi:10.1128/msystems.00191-16.
- 532 48. **Schloss PD.** 2018. Identifying and overcoming threats to reproducibility, replicability,  
533 robustness, and generalizability in microbiome research. *mBio* **9**. doi:10.1128/mbio.00525-18.
- 534 49. **Seekatz AM, Rao K, Santhosh K, Young VB.** 2016. Dynamics of the fecal microbiome  
535 in patients with recurrent and nonrecurrent clostridium difficile infection. *Genome Medicine* **8**.  
536 doi:10.1186/s13073-016-0298-8.
- 537 50. **Schloss PD, Westcott SL, Ryabin T, Hall JR, Hartmann M, Hollister EB, Lesniewski RA,**  
538 **Oakley BB, Parks DH, Robinson CJ, Sahl JW, Stres B, Thallinger GG, Horn DJV, Weber CF.**  
539 2009. Introducing mothur: Open-source, platform-independent, community-supported software  
540 for describing and comparing microbial communities. *Applied and Environmental Microbiology*  
541 **75**:7537–7541. doi:10.1128/aem.01541-09.
- 542 51. **Rognes T, Flouri T, Nichols B, Quince C, Mahé F.** 2016. VSEARCH: A versatile open source  
543 tool for metagenomics. *PeerJ* **4**:e2584. doi:10.7717/peerj.2584.
- 544 52. **Westcott SL, Schloss PD.** 2017. OptiClust, an improved method for assigning  
545 amplicon-based sequence data to operational taxonomic units. *mSphere* **2**:e00073–17.  
546 doi:10.1128/mspheredirect.00073-17.



547 53. **R Core Team**. 2018. R: A language and environment for statistical computing. R Foundation  
548 for Statistical Computing, Vienna, Austria.

549 **Figure 1. The error rate of assembled mock community sequence reads increases with the**  
550 **number of rounds of PCR; however, much of this error was eliminated by denoising and**  
551 **followed the relative error rates provided by the manufacturers.** Each line represents the  
552 mean of four replicates.

553 **Figure 2. The fraction of all denoised sequences that were identified as being chimeric**  
554 **increases with the number of rounds of PCR used and varied between polymerases.** (A)  
555 Sequencing of a mock community allowed us to identify the total fraction of sequences that were  
556 chimeric as well as the specificity and sensitivity of UCHIME to detect those chimeras. Each line  
557 represents the mean of four replicates. (B) Sequencing of four human stool samples after using  
558 one of five different polymerases again demonstrated increased rate of chimera formation with  
559 increasing number of rounds of PCR and variation across polymerases.

560 **Figure 3. The relative abundances of mock community sequence reads mapped to**  
561 **reference sequences differed subtly from the expected relative abundances as determined**  
562 **by shotgun metagenomic sequencing.** Bias did not increase with number of rounds of PCR or  
563 vary by polymerase or the guanine and cytosine content of the fragment. The expected relative  
564 abundance of each organism is indicated by the horizontal gray line. The percentage of bases  
565 that were guanines or cytosines within the V4 region of the 16S rRNA genes in each organism is  
566 indicated by the number in the lower left corner of each panel. Each line represents the mean of  
567 four replicates.

568 **Figure 4. Despite evidence of subtle PCR bias at the genome level, there was significant**  
569 **evidence of bias using community-wide metrics that grew with the number of rounds of**  
570 **PCR when using a mock community.** (A) With the exception of the KAPA polymerase data, the  
571 richness and Shannon diversity values increased with number of rounds of PCR and the inclusion of  
572 residual sequencing errors and chimeras. The horizontal black line indicates the expected richness  
573 and diversity. (B) Relative to the mock community sampled after 25 rounds of PCR, the distance  
574 to the communities sampled after 30 and 35 rounds of PCR increased for all polymerases. (C)  
575 The variation between samples demonstrated a significant change in the community driven by the  
576 number of rounds of PCR and the polymerase used. The ellipses represent bivariate normally

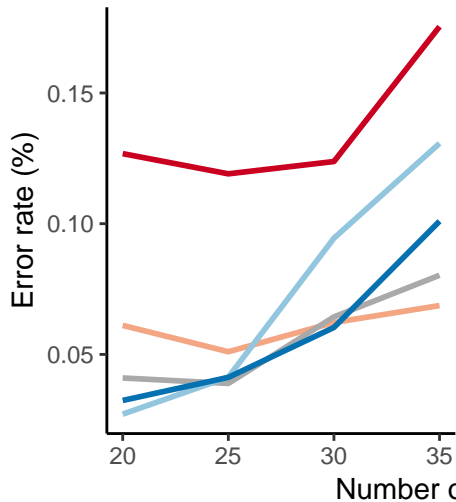
577 distributed 95% confidence intervals. The data in A and B represents the mean of four replicates.

578 **Figure 5. Sequencing of human stool samples indicated clear increase in bias with number**  
579 **of rounds of PCR, however, the bias appeared to be consistent within each sample.** (A) With  
580 the exception of data collected using the KAPA polymerase, the richness and Shannon diversity  
581 values increased with number of rounds of PCR. (B) Relative to the stool communities sampled  
582 after 25 rounds of PCR, the distance to the stool communities sampled after 30 and 35 rounds of  
583 PCR was inconsistent and there was little difference in variation for data collected using the KAPA  
584 polymerase. (C) The variation between stool samples was larger than the amount of variation  
585 introduced by varying the number of rounds of PCR or polymerase. The ellipses represent bivariate  
586 normally distributed 95% confidence intervals. Results for some samples at 20 cycles are not  
587 presented because it was not possible to obtain a sufficient number of reads for those polymerases.

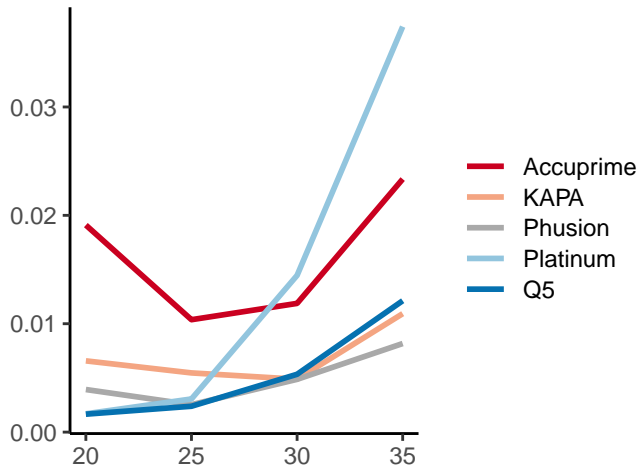
588 **Figure 6. The average distance between replicates of sequencing the same mock**  
589 **community or between the human stool samples (i.e. drift) did not vary by number of**  
590 **rounds of PCR or by polymerase. .**

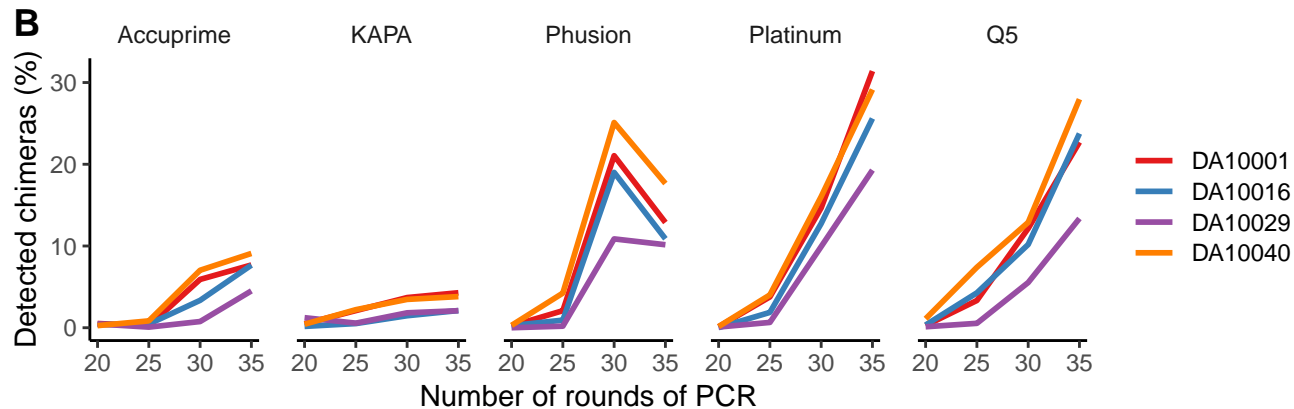
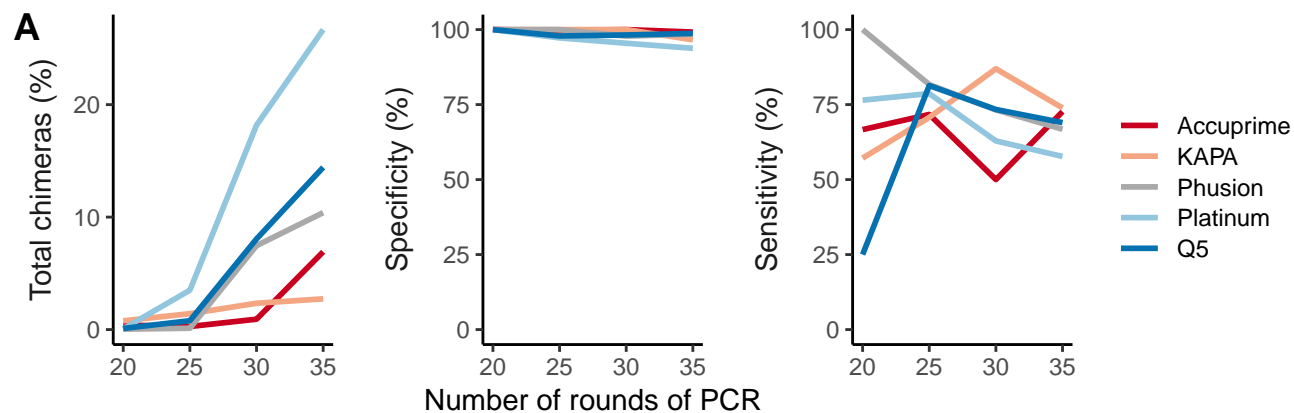
591 **Figure S1: With the exception of the sequence data generated using the KAPA polymerase,**  
592 **the ratio of the two *Salmonella enterica* V4 sequences from the mock community was lower**  
593 **than the expected ratio of 6:1.** The dominant and rare *S. enterica* V4 sequences differed by a  
594 single base. The horizontal gray line indicates the expected 6:1 ratio. Each line represents the  
595 mean of four replicates.

Raw contigs

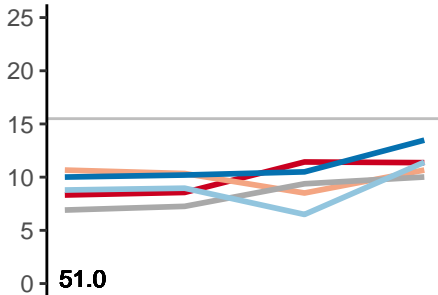


Denoised contigs

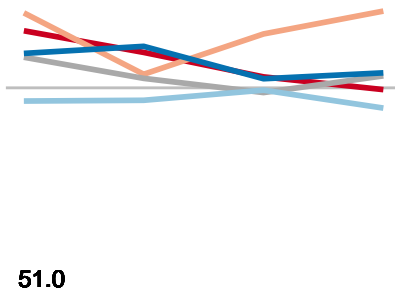




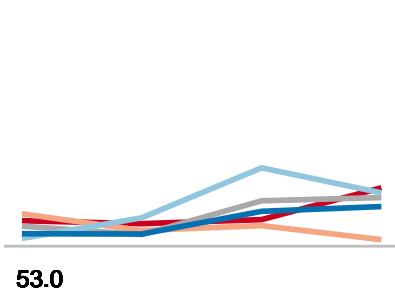
*Lactobacillus fermentum*



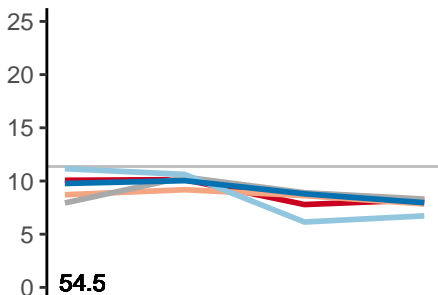
*Staphylococcus aureus*



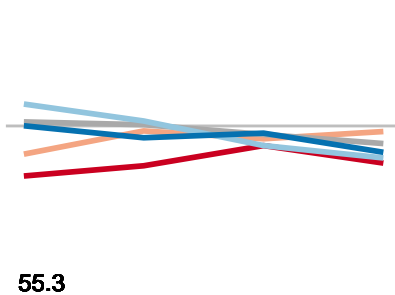
*Pseudomonas aeruginosa*



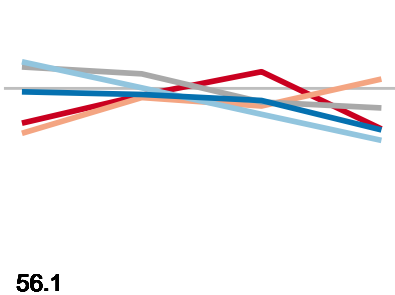
*Enterococcus faecalis*



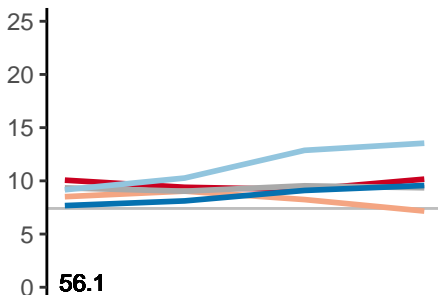
*Listeria monocytogenes*



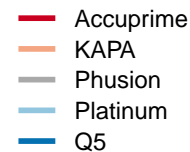
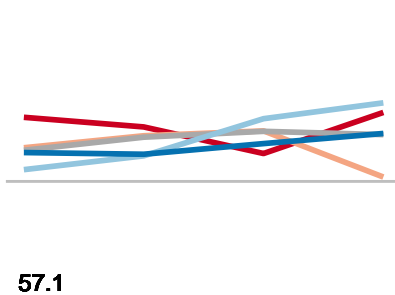
*Bacillus subtilis*



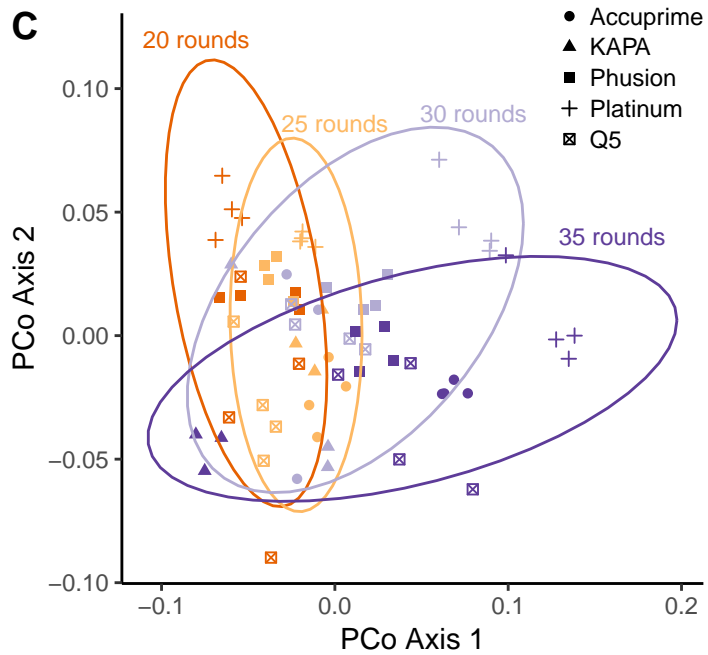
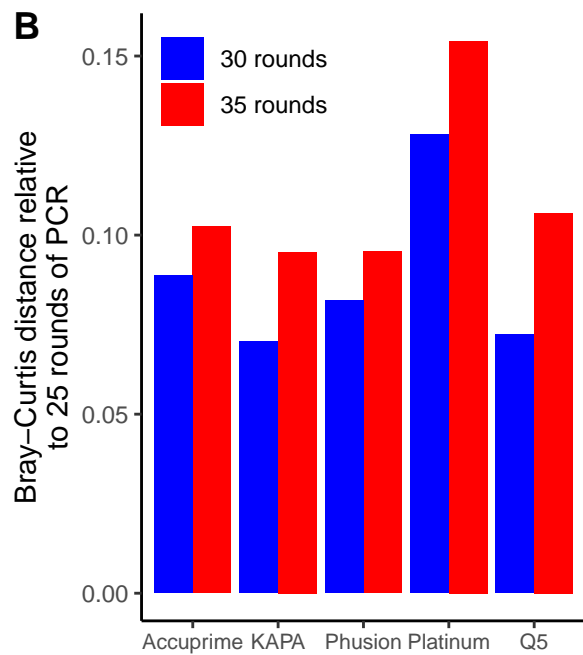
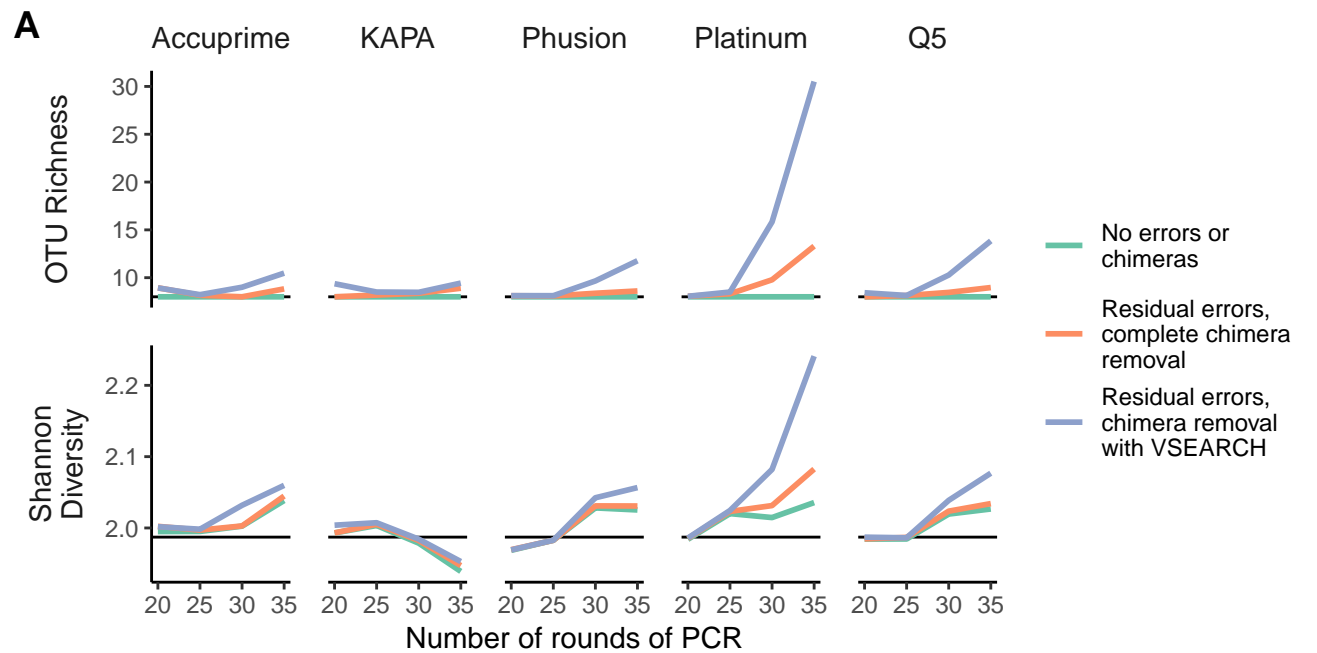
*Escherichia coli*

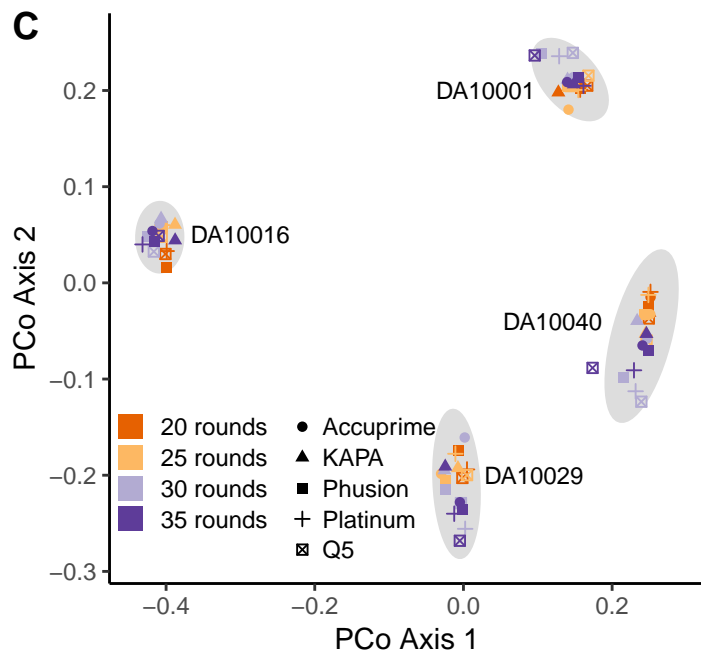
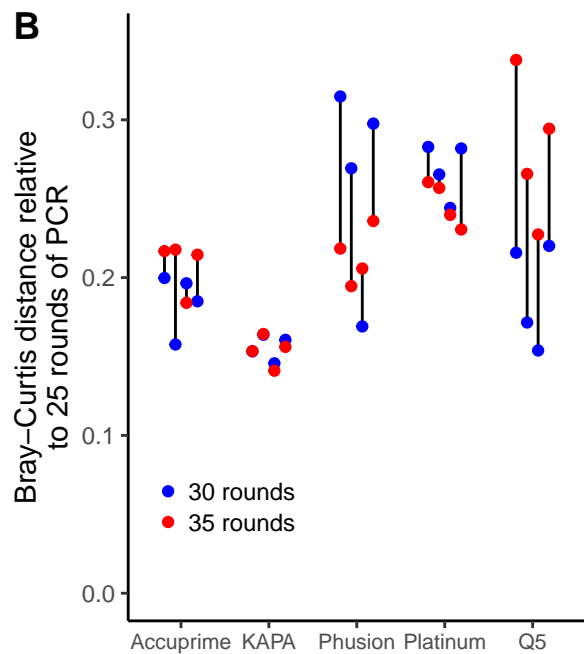
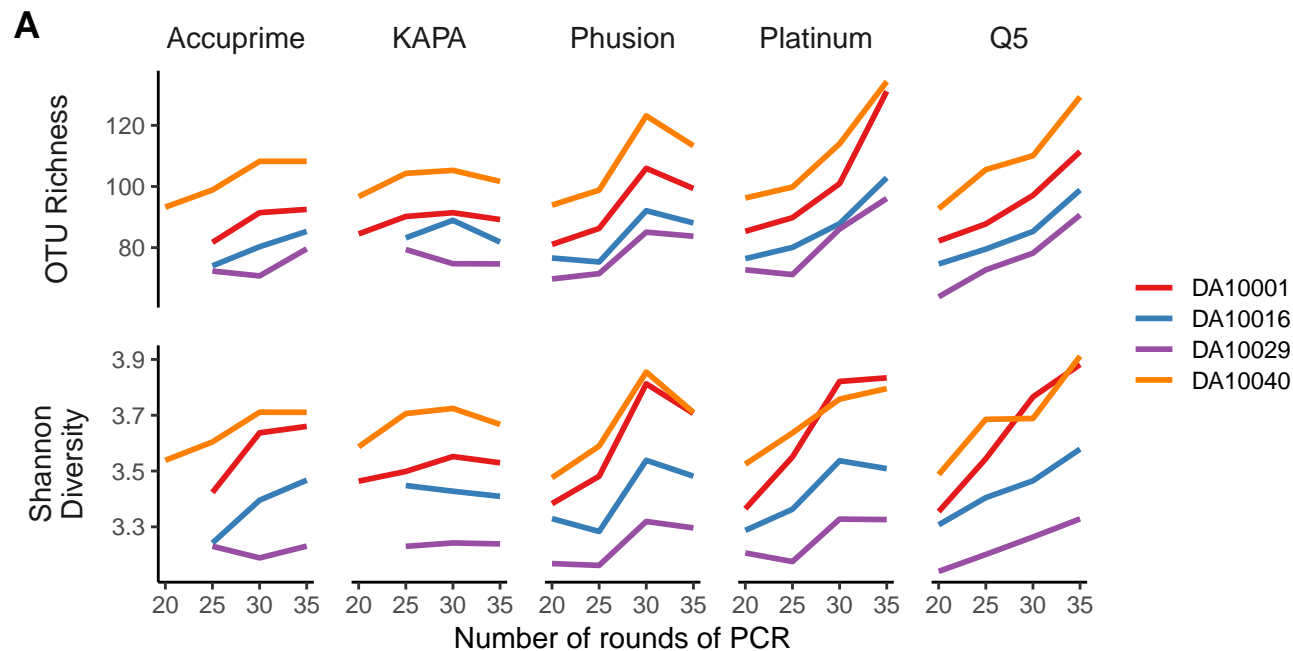


*Salmonella enterica*



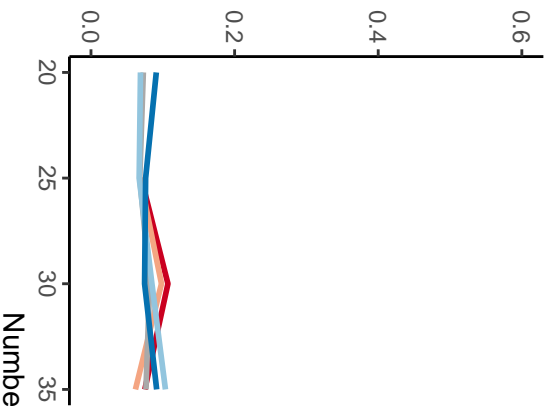
Number of rounds of PCR



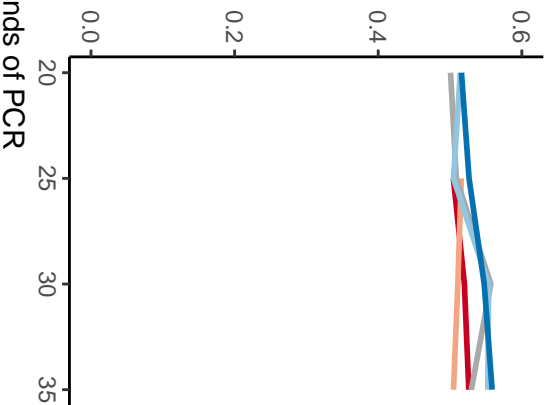




Mean intra-replicate Bray-Curtis distances for the mock community



Mean inter-sample Bray-Curtis distances for the stool samples



- Accuprime
- KAPA
- Phusion
- Platinum
- Q5