

1 **Title:**

2 A Sensitivity Analysis of Methodological Variables Associated with Microbiome Measurements

3 **Authors:**

4 Samuel P. Forry<sup>1</sup>, Stephanie L. Servetas<sup>1</sup>, Jennifer N. Dootz<sup>1</sup>, Monique E. Hunter<sup>1</sup>, Jason G. Kralj<sup>1</sup>, James J.

5 Filliben<sup>2</sup>, Scott A. Jackson<sup>1</sup>

6 **Affiliation:**

7 1. Complex Microbial Systems Group, National Institute of Standards and Technology (NIST),

8 Gaithersburg, MD

9 2. Multimodal Information Group, National Institute of Standards and Technology (NIST),

10 Gaithersburg, MD

11 **Abstract**

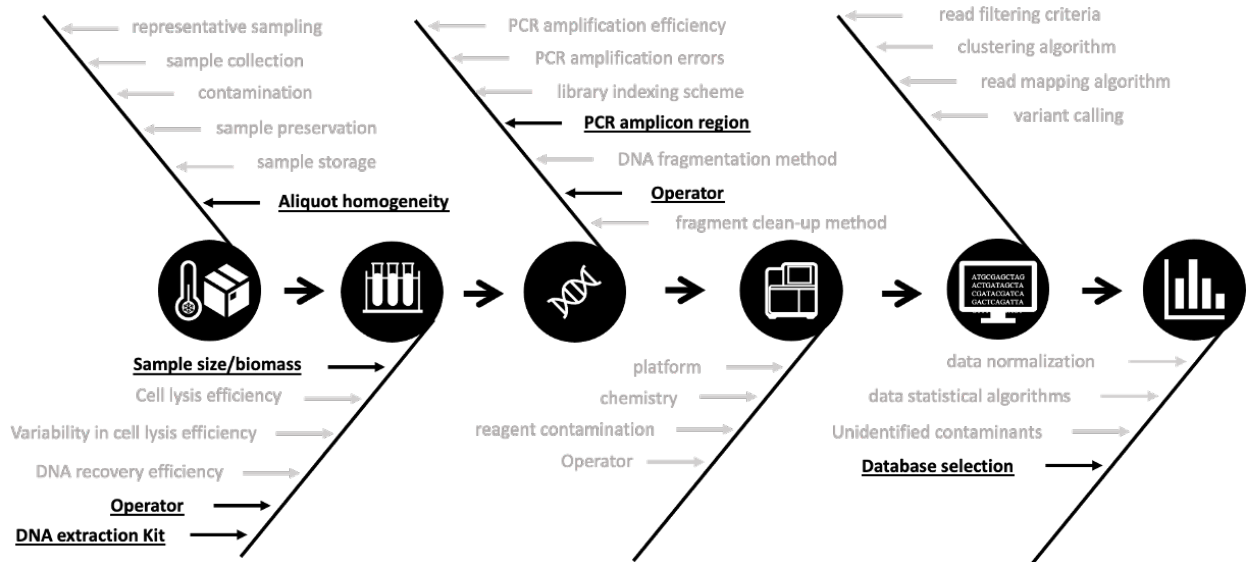
12 The experimental methods employed during metagenomic sequencing analyses of microbiome samples  
13 significantly impact the resulting data and typically vary substantially between laboratories. In this study,  
14 a full factorial experimental design was used to compare the effects of a select set of methodological  
15 choices (sample, operator, lot, extraction kit, variable region, reference database) on the analysis of  
16 biologically diverse stool samples. For each parameter investigated, a main effect was calculated that  
17 allowed direct comparison both between methodological choices (bias effects) and between samples  
18 (real biological differences). Overall, methodological bias was found to be similar in magnitude to real  
19 biological differences, while also exhibiting significant variations between individual taxa, even between  
20 closely related genera. The quantified method biases were then used to computationally improve the  
21 comparability of datasets collected under substantially different protocols. This investigation  
22 demonstrates a framework for quantitatively assessing methodological choices that could be routinely

23 performed by individual laboratories to better understand their metagenomic sequencing workflows  
24 and to improve the scope of the datasets they produce.

## 25 **Introduction**

26 The democratization of access to next-generation sequencing (NGS) has made this technology the  
27 workhorse of modern microbiology. NGS-based metagenomic sequencing (MGS) enables analysis of the  
28 genomic content of microbial communities, reporting both the taxonomic identification and count of  
29 sequencing reads and even enumeration of constituent microbes. In recent years, MGS has increasingly  
30 found its way into the clinic and other regulated spaces such as quality control testing where the  
31 measurement results can have important consequences.(1-3) Despite widespread adoption of MGS  
32 capabilities, the analytical performance remains uncertain, particularly when comparing between  
33 protocols.(4)

34 While MGS measurements may appear straightforward, the results actually depend on a complex series  
35 of steps and sample manipulations, typically including: sample storage and transport, DNA extraction  
36 and purification, DNA fragmentation and barcoding, library preparation and NGS, bioinformatic analysis  
37 to identify and enumerate unique microbial signatures, and taxonomic assignment against a database,  
38 as well as additional bioinformatic analyses dependent on particular experimental designs and project  
39 goals. Further, within each of these steps, many methodological choices must be specified for the  
40 particular protocol employed. Each of these steps, as well as the individual methodological choices, can  
41 introduce bias and variability that then accumulate through the measurement process. (5-7) The  
42 methodological contributions of bias and variability can be visually depicted using an Ishikawa diagram  
43 (Figure 1).



44

45 ***Figure 1 (link): Ishikawa diagram of bias associated with an amplicon metagenomic***

46 ***measurement workflow.*** The central image depicts a typical metagenomic analysis: sample

47 *handling and storage, DNA extraction, sequencing library prep, next-generation sequencing,*

48 *sequence data processing, and bioinformatic analysis.* At each step, some of the various

49 *methodological choices that can introduce variability and bias are listed. Bold and underlined*

50 *steps indicate methodological choices explicitly varied in the current study.*

51 In general, MGS datasets exhibit high precision (low variability) for locked-down workflows (8, 9), and

52 this makes MGS particularly useful within a research project where protocols can be easily controlled.

53 Given the myriad options for any single step (as well as a dynamic field where procedures continue to

54 evolve quickly), it is unlikely that two different labs would implement the exact same measurement

55 workflow from start to finish without substantial coordination. Hence, the lack of reproducibility

56 between laboratories is perhaps not surprising even when the results are highly reproducible within a

57 laboratory.(10) While there are limited opportunities through the MGS workflow for intermediate

58 characterization of the impact of each methodological decision, there are experimental designs aimed at

59 assessing the effect sizes of particular methodological decisions at the conclusion of the MGS

60 measurements.(11) One commonly used approach is a full factorial design, evaluating at least two  
61 specified options ('Level') for each tested methodological parameter.(12-14) This experimental design  
62 allows for a rigorous statistical assessment of the effect of individual parameters and provides  
63 quantitative data for direct comparisons.

64 We previously described the Mosaic Standards Challenge (MSC), an international interlaboratory study  
65 designed to assess the impact of methodological variability on MGS results.(15) The MSC employed five  
66 biologically distinct human fecal reference materials and a comprehensive standardized metadata  
67 reporting sheet that allowed participants to share exhaustive details about the protocols and methods  
68 used for analyzing the fecal materials. The MSC employed an "open protocol" design that allowed (and  
69 even encouraged) participants to follow divergent protocols as defined for their lab's routine MGS  
70 measurements. This strategy showcased the diversity of methodologies commonly employed across  
71 MGS measurement labs and highlighted the challenge of comparing MGS results between different  
72 protocols. In the current effort, we used the same five fecal samples from the MSC, but employed a fully  
73 specified experimental design where a small subset of methodological choices were systematically  
74 varied while the rest of the MGS measurement workflow was held constant.

75 The goal of this work was to demonstrate a rigorous experimental approach for evaluating bias in  
76 metagenomics workflows. Specifically, we used a full-factorial design (plus replication) for each fecal  
77 sample (n=5) to quantify the effects of different operators (n=2), reference material lot (n=2), DNA  
78 extraction kit (n=2), marker gene target (n=2), and reference database (n=3). The magnitude of effect for  
79 each selected methodological choice was quantified, and these main effects were compared between  
80 protocol choices and to the biological variability between fecal samples. Ratiometric analysis was  
81 employed to account for MGS measurement compositionality. In combination with spike-in species  
82 added to each fecal sample, this analysis allowed direct comparisons between samples. Among the

83 methodologies explored, the largest effects were observed for extraction kit selection, where the effect  
84 varied substantially between individual taxa, even within taxonomic clades.

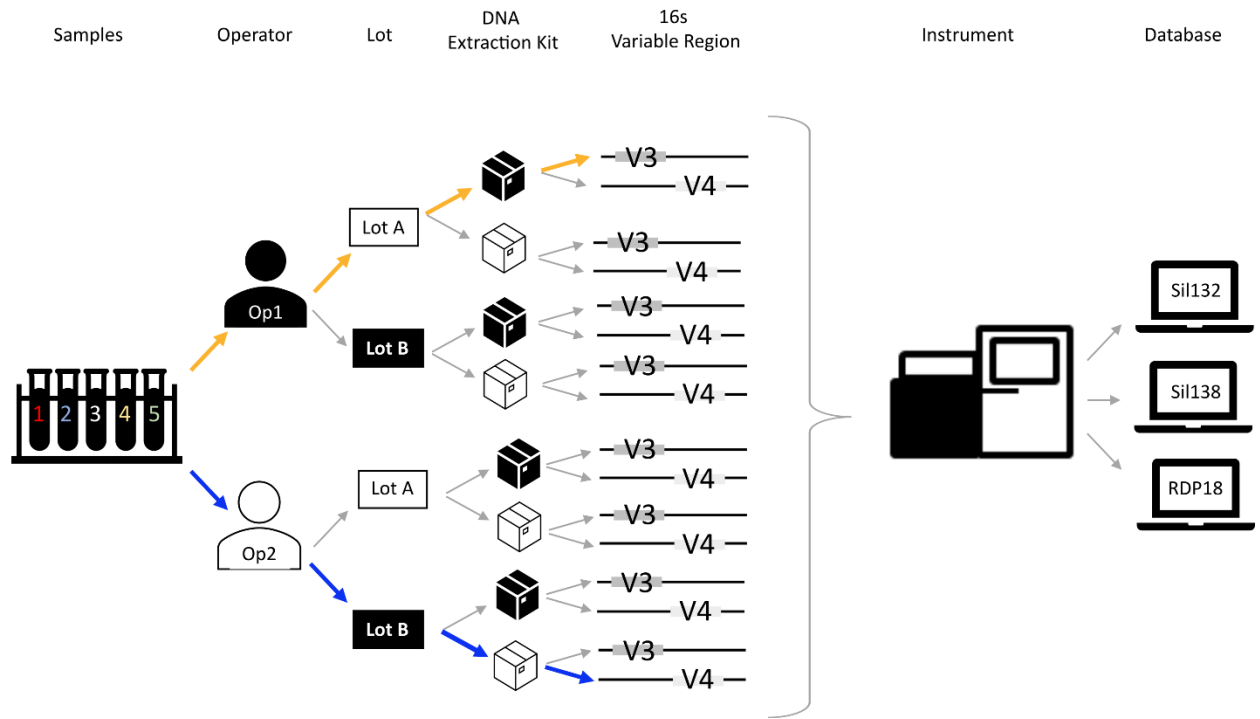
## 85 **Results and Discussion**

### 86 Metagenomic Sequencing Methodologies and Experimental Design

87 Using an Ishikawa diagram, a subset of parameters were selected (shown in bold and underlined) for  
88 this initial proof-of-concept study (Figure 1). The experimental design included parameters that were  
89 denoted ‘biological’ (i.e., actual biological differences between disparate stool samples) or  
90 ‘methodological’ (i.e., bias associated with the measurement protocol) to allow for direct comparison of  
91 the variability from biological differences to the variability due to methodological differences. For this  
92 study, 5 unique human fecal samples were selected representing the biological variability. For each  
93 sample, the methodological parameters of operator (n = 2), lot (n = 2), extraction kit (n = 2), 16S variable  
94 region (n = 2), database (n = 3), were chosen (with the indicated number of levels). Thus, 240 MGS  
95 datasets (270 with replicates) were generated from 48 distinct workflows (Figure 2). The orange and  
96 blue arrows in Figure 2 at the top and bottom arms of the study, respectively, represented orthogonal  
97 workflows where unique methods were chosen for each parameter. The two orthogonal workflows  
98 were replicated to assess precision, and the orange and blue lines were used here and in subsequent  
99 figures to showcase data originating from the two completely distinct workflows.

100 This work represents a proof-of-concept and was intended to demonstrate a full factorial analysis of a  
101 metagenomic sequencing workflow. Ideally, the implementation of a full-factorial experimental design  
102 would specify ‘high’ and ‘low’ values for each selected parameter to bound the range of normal  
103 experimental conditions. In the case of the metagenomics workflows evaluated here, the selected  
104 experimental parameters of interest exhibited discrete, non-numerical, non-continuous levels (e.g.,  
105 Operator, Lot, Extraction Kit), and the selected options (levels) may not bound the normal range of

106 experimental conditions. Herein, a rigorous statistical approach was demonstrated for evaluating bias in  
107 metagenomics workflows that others could utilize in their own laboratories with their particular  
108 protocols.



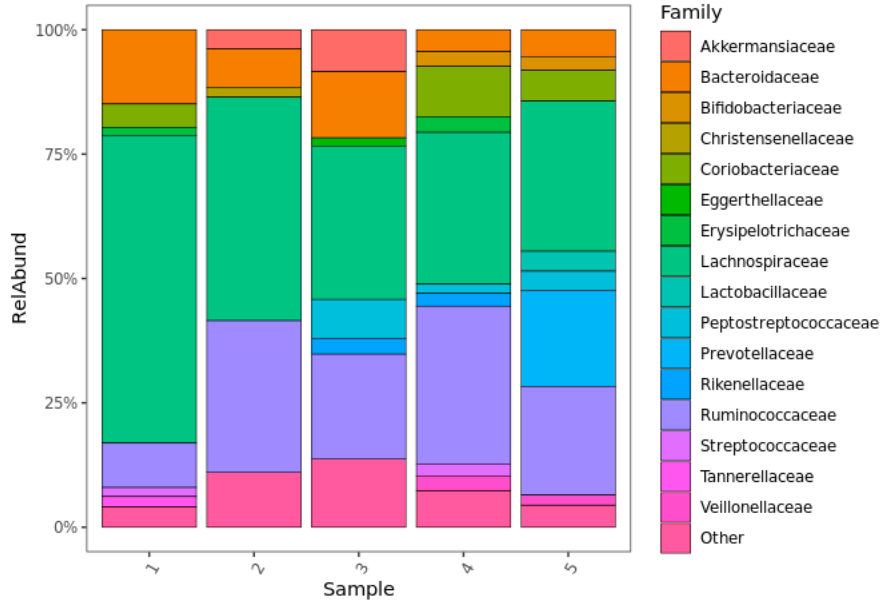
109

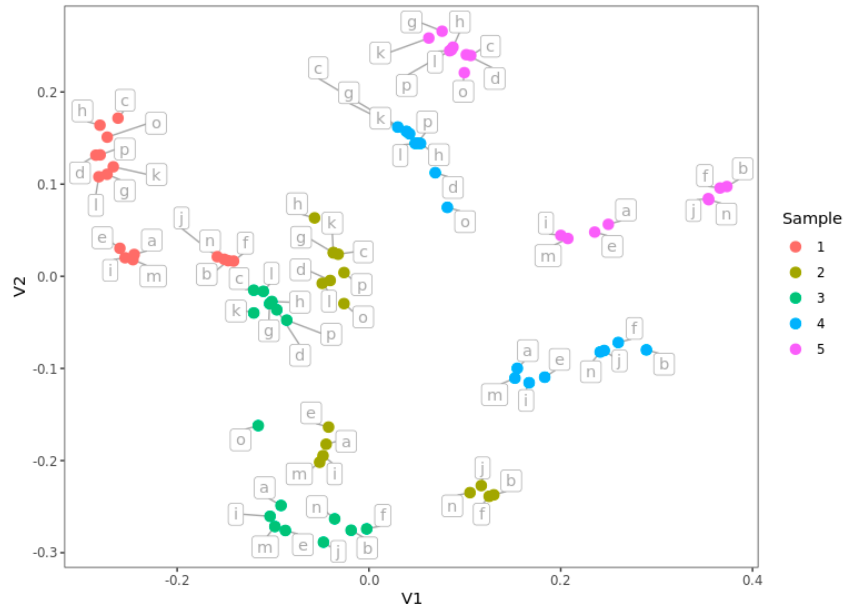
110 **Figure 2 (link): Systematic experimental design enables quantitative bias assessment.** The  
111 graphic above illustrates the levels compared for each step of the metagenomic workflow. In  
112 total, 5 different human stool samples were handled by 2 operators, preparing samples from 2  
113 lots, for MGS analysis using 2 DNA extraction kits and 2 variable regions, with taxonomic  
114 assignment coming from 3 databases. The top (orange) and bottom (blue) arms of the  
115 experimental design represent completely orthogonal protocol choices through all wet-bench  
116 operations, and these two protocols were replicated (2x) for each sample and database to  
117 explore reproducibility. This resulted in a total of 240 MGS data files (270 with replicates) that  
118 were analyzed.

119 MGS: Compositional Analysis

120 A typical metagenomic analysis may include bar charts to show relative abundance and community  
121 members. As such, the composition of each sample for a single workflow was visualized in a relative  
122 abundance bar chart at the Genus taxonomic level (Figure 3A). The differences between stool samples  
123 were evident, and good reproducibility was observed between replicate analyses. For example, while  
124 *Lachnospiraceae* and *Ruminococcaceae* make up the plurality of each sample, their relative proportions  
125 vary, and some taxa (e.g., *Rikenellaceae*, *Prevotellaceae*, *Akkermansiaceae*) appeared variably present,  
126 occurring above the 1.5 % cutoff in some samples but not others. A comparison of MGS bar charts from  
127 alternate workflows showed that differences in sample composition could be correlated with  
128 methodological choices (Figure SI-1); however, using bar charts to assess methodological parameters  
129 can be difficult to implement systematically. While differences between samples can be readily  
130 observed, it is difficult to determine the magnitude of these differences or to correlate them with  
131 particular methodological decisions. Also, comparing the compositional changes between many  
132 individual bar charts proves to be cumbersome and challenging to automate.

133 Another common approach is to evaluate shifts in composition by calculating similarities (or  
134 dissimilarities such as Bray-Curtis) between MGS datasets and then generating a principal coordinate  
135 analysis (PCoA) plot from the resulting distance metrics (Figure 3B). In this approach, a single point was  
136 generated for each unique combination of sample and analysis workflow, and clusters of points were  
137 then correlated with methodological parameters. For example, with the data plotted in Figure 3B, data  
138 from a common sample (same colors) generally clustered, with some overlap and sub-clustering  
139 associated with methodological parameters (denoted with letter flags). To compare between specific  
140 methodological choices, additional faceting generally allowed for their relative impact to be assessed  
141 (Figure SI-2a-e). This pattern recognition approach facilitated the evaluation of multiple datasets  
142 simultaneously; however, the results remained mostly qualitative in nature, as the apparent significance  
143 of each evaluated parameter depended on the ordination conditions of that dataset.





145

Symbol	Operator	Lot	ExtractionKit	VRegion
a	J	A	Q	V3
b	M	A	Q	V3
c	J	A	Z	V3
d	M	A	Z	V3
e	J	B	Q	V3
f	M	B	Q	V3
g	J	B	Z	V3
h	M	B	Z	V3
i	J	A	Q	V4
j	M	A	Q	V4
k	J	A	Z	V4
l	M	A	Z	V4
m	J	B	Q	V4
n	M	B	Q	V4
o	J	B	Z	V4
p	M	B	Z	V4

146

147 **Figure 3: (a) Bar charts and (b) PcoA plots from 5 human stool samples.** (3A) Family level  
 148 relative abundance bar charts show the compositional diversity among 5 human stool samples  
 149 for a single workflow. For this plot, Families present at <1.5% relative abundance are grouped as  
 150 “Other”. Correlation between the observed compositions and specified protocol choices (see  
 151 Figure SI-1) can help identify particularly sensitive steps in an MGS workflow. (3B) Bray-Curtis

152 *dissimilarities can also be used to identify sensitive MGS protocol choices. Points are colored here*  
153 *to indicate fecal sample and letter flags show additional protocol choices for each dataset, as*  
154 *denoted in the table. Further clustering (see Figure SI-2) can be used to rank protocol choices*  
155 *with respect to their impact on the observed composition. Though these two plots are*  
156 *conventionally used, they are limited in their utility to identify and address biases.*

### 157 MGS: Ratiometric Analysis

158 An additional challenge when comparing MGS results between samples is the compositional nature of  
159 the datasets. That is, the measured relative abundance value for any taxa depended on both its actual  
160 abundance in the sample and on the sum of the observed abundances for all other taxa present in that  
161 sample. An alternative approach for comparing MGS results has focused on ratiometric comparisons,  
162 and these ratios of taxa have been shown to be *independent* of sample composition. (5, 8, 16-18)

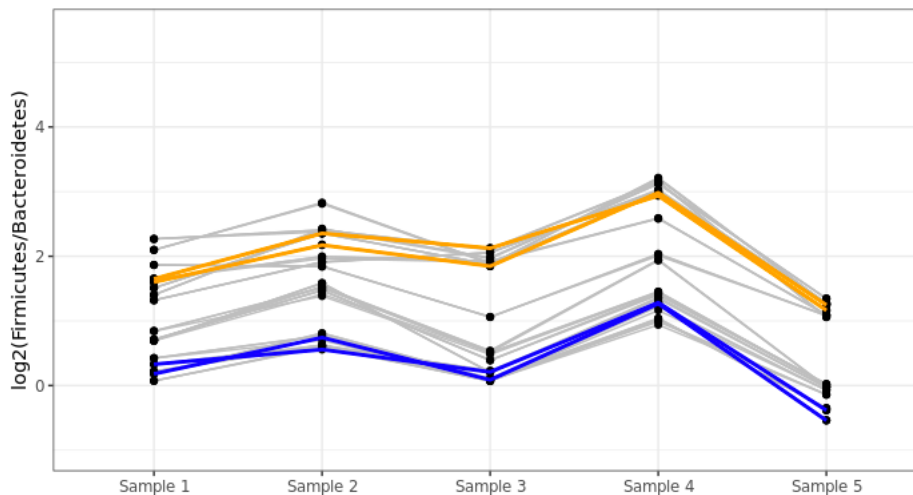
163 Whereas the initial relative abundance values could be misleading, ratios of relative abundance between  
164 two taxa within a sample have been shown to be more reliable, varying with just the actual taxa  
165 abundances and measurement biases.(8)

166 This concept has been implemented previously with native taxa, such as calculating the  
167 *Firmicutes:Bacteroidetes* (F:B) ratios (among others), with varying levels of success (10, 19-21). To  
168 demonstrate ratiometric analysis, we examined the F:B ratio within our experimental design (Figure 4).

169 The F:B ratio provided a single calculated value for each unique combination of stool sample and  
170 experimental workflows (black points), with the samples analyzed under identical workflows connected  
171 by light grey lines. For reference, the two completely orthogonal workflows are colored orange and blue  
172 (and match the top and bottom arms of the experimental design in Figure 2).

173 As expected, due to real variation in the actual F:B ratio between biologically distinct samples, similar  
174 trends ('M' shape) were observed across samples for each of the different protocols evaluated. A

175 Kruskal-Wallis Test confirmed that there was a statistically significant difference in F:B ratio between the  
176 stool samples ( $\chi^2(4) = 77.8, p = 1e-15$ ). However, within each sample, a greater than 4-fold difference in  
177 F:B ratio was observed between different workflows, reflecting the impact of methodological choices on  
178 the observed F:B ratio. Indeed, this workflow-dependent effect (bias) was at least as significant as the  
179 biological differences observed between distinct stool samples. In a rapidly evolving field like MGS, this  
180 kind of workflow-dependent measurement bias significantly complicates comparisons between studies.



181

182 ***Figure 4: Impact of methodological parameters on the Firmicutes:Bacteroidetes ratio.*** The  
183 *Firmicutes:Bacteroidetes* ratio was determined for each sample (black points) under every  
184 combination of the methodological parameters investigated in this study (48 combinations for  
185 each of the 5 samples). A grey line connects the dots between fecal samples processed using the  
186 same protocol; dots and lines overlapped for protocols where the varied parameters didn't  
187 change the observed *Firmicutes:Bacteroidetes* ratio. The blue and orange lines denote the results  
188 from the top and bottom legs of our experimental design flowchart (shown in Figure 2),  
189 respectively, and represent orthogonal choices for every methodological parameter (operator,  
190 lot, DNA extraction kit, variable region, database). There are two blue and two orange lines, as

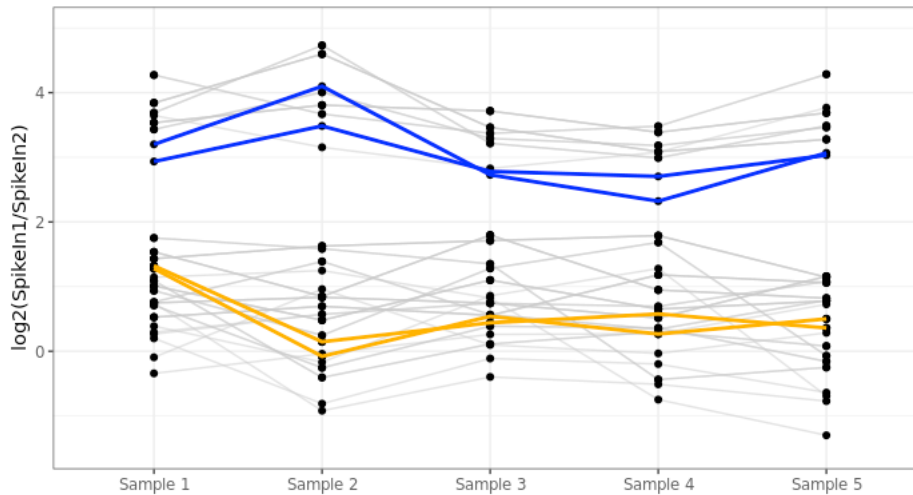
191 *each of these workflows were run in duplicate to demonstrate reproducibility when all*  
192 *parameters are held consistent.*

### 193 MGS: Internal Standards

194 Given that native taxa are expected to vary naturally between the stool samples, as observed here with  
195 Firmicutes and Bacteroidetes, an alternative approach for ratiometric analysis is to add exogenous  
196 strains (spike-ins) to serve as internal controls. In this study, two strains were uniformly added during  
197 the batch production of these stool samples. *Leifsonia xyli* (Gram positive) and *Aliivibrio fischeri*  
198 (formerly known as *Vibrio fischeri*, Gram negative), a plant pathogen and marine organism, respectively,  
199 were selected as they are unlikely to be found in human stool samples and could therefore be uniquely  
200 identified taxonomically. Unlike the native taxa, these spike-in organisms were added at consistent  
201 abundances across all stool samples.

202 Similar to the F:B ratio plotted in Figure 4, the *Leifsonia:Aliivibrio* (L:A) ratio was calculated for each  
203 unique combination of stool sample and experimental protocol (Figure 5). As before, a high degree of  
204 bias was observed between individual workflows (vertical distribution of points for each sample).  
205 However, unlike the previous analysis, the L:A ratio did not exhibit trends across samples ('M' shape  
206 from Figure 4). Supporting this lack of biological variability, a Kruskal-Wallis Test failed to find a  
207 statistically significant difference in L:A ratio between the stool samples ( $\chi^2(4) = 4.8$ ,  $p = 0.31$ ). Instead,  
208 observed variations in the L:A ratio between samples analyzed with a common workflow (i.e., following  
209 a single line trace) were attributed to measurement variability. Considering the differences in the  
210 relative abundances of *Firmicutes* and *Bacteroidetes* ( $\approx 88\%$ ) versus *Leifsonia* and *Aliivibrio* ( $\approx 0.6\%$ ), the  
211 increased measurement variability observed in Figure 4 was unsurprising. Since the spike-in organism  
212 abundances were detected consistently between samples, we concluded that either of these spike in  
213 organisms were suitable for use as internal standards for ratiometric analyses. For the purposes of

214 demonstrating ratiometric analysis we chose to use *Leifsonia*. As such, comparing the  
215 *Bacteroidetes:Leifsonia* ratio between different samples handled under a constant workflow allowed  
216 actual changes in Bacteroidetes to be quantified independent of any other native taxa.



217

218 ***Figure 5: Analysis of internal standards added uniformly to all stool samples.*** Two exogenous  
219 organisms, *Aliivibrio* and *Leifsonia*, were added as internal standards to the stool samples. The  
220 ratio of *Aliivibrio:Leifsonia* was calculated for each sample (black points) under every  
221 combination of the methodological parameters investigated in this study (48 combinations), as  
222 described in Figure 4. A grey line connects the dots between fecal samples processed using the  
223 same protocol; dots and lines overlapped for protocols where the varied parameters didn't  
224 change the observed *Aliivibrio:Leifsonia* ratio. The blue and orange lines denote the results from  
225 the top and bottom legs of our experimental design flowchart (shown in Figure 2), respectively,  
226 and represent orthogonal choices for every methodological parameter (operator, lot, DNA  
227 extraction kit, variable region, database). There are two blue and two orange lines, as each of  
228 these workflows were run in duplicate to demonstrate reproducibility when all parameters are  
229 held consistent.

## 230 Quantifying Parameter Effects

231 In Figure 5, the systematically varied workflows resulted in significant variation in the ratiometric MGS  
232 results within each sample. The full-factorial experimental design (Figure 2) was structured to balance  
233 each decision point in the workflow and facilitate independent characterization of each methodological  
234 choice. A Parameter Effect was calculated for each evaluated step in the workflow (i.e., parameter) and  
235 the method selected (i.e., level) by calculating the fold change of ratiometric results from workflows  
236 using the specified method versus the average results across all workflows (Equation 1)

$$237 \quad \text{Parameter Effect}_{(\text{Parameter}=\text{Level})} = \log_2 \left( \frac{GM(\text{Ratio}(\text{Taxa:Lf})_{\text{Parameter}=\text{Level}})}{GM(\text{Ratio}(\text{Taxa:Lf})_{\text{All Levels}})} \right) \quad (1)$$

238 where *GM* denotes the geometric mean. This calculated parameter effect allowed individual workflow  
239 steps and specific methodological choices to be directly compared.(11, 12)

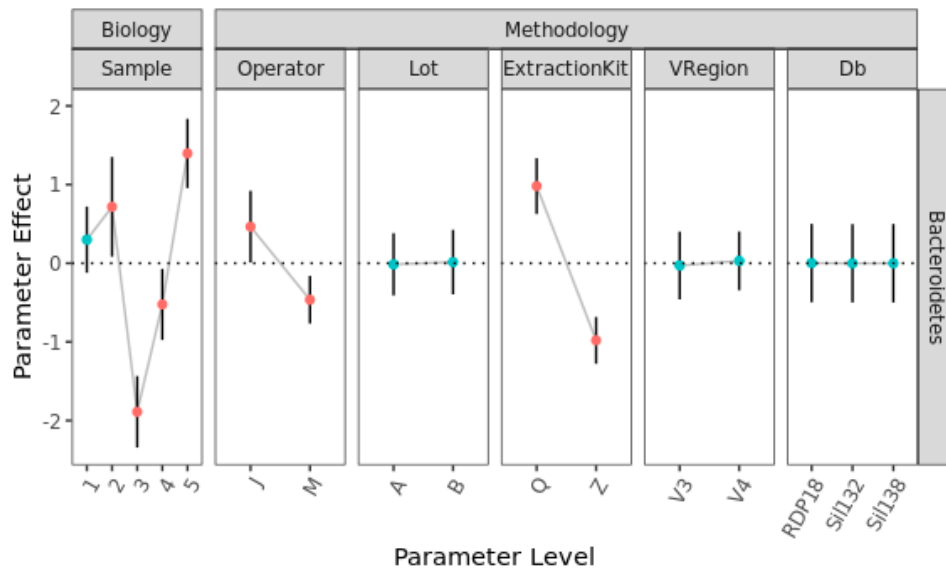
240 As expected, when the Parameter Effects for the ratio of the two internal standard taxa were calculated,  
241 (Figure SI-3) no significant effect was measured between samples due to the even amount of the  
242 internal standards added to each sample; however, several of the methodological parameters did  
243 exhibit statistically significant Parameter Effects for the ratio of the two internal standards. It was  
244 important to note that the Parameter Effect calculation (Equation 1) was based on ratios with the  
245 internal standard and could not attribute observed effects to either of the ratioed taxa individually.  
246 When we subsequently considered taxa native to the stool, a large Parameter Effect denoted a large  
247 bias differential between the taxa-of-interest and the internal standard, while a small or negligible  
248 Parameter Effect indicated a similar response between the taxa-of-interest and the internal standard.  
249 The Parameter Effect calculations for the phyla Bacteroidetes revealed significant effects for several  
250 parameters (Figure 6), additionally the differences between biological samples (e.g. aliquots) were also  
251 generally significant for Bacteroidetes. One benefit of the Parameter Effect calculation in the context of

252 a full-factorial design was its ability to quantify the effects of methodological parameters even in the  
253 context of biological variations. This was confirmed by calculating the Parameter Effects for  
254 *Bacteroidetes* within each stool sample individually (Figure SI-4). The within-sample Parameter Effects  
255 matched the aggregate Parameter Effects, albeit with a larger (less precise) confidence interval due to  
256 the smaller number of samples. In the aggregate Parameter Effects calculation, the extraction kit and  
257 operator parameters exhibited statistically significant effects, while the parameters of lot, 16S variable  
258 region, and database did not exhibit significant effects.

259 The Parameter Effects for Extraction Kit were similar in magnitude to the real biological differences  
260 between samples, potentially complicating comparisons between samples that were analyzed using  
261 different methods. For example, comparing stool samples 4 and 5 using Extraction Kit Z, while holding all  
262 other methods constant, revealed an actual difference in *Bacteroidetes:Leifsonia* ratio (B:L) of 4.1-fold.  
263 However, evaluating only sample 4 using Extraction Kits Q and Z revealed an apparent difference in B:L  
264 of 4.4-fold, just due to the methodological difference. When the measured B:L was compared between  
265 an analysis of stool sample 4 using extraction kit Q and an analysis of stool sample 5 using extraction kit  
266 Z, the results were essentially identical (1.05-fold difference) because the effects of the biological and  
267 methodological differences were of similar magnitude and in opposite directions. Alternately, analyzing  
268 sample 4 with extraction kit Z and sample 5 with extraction kit Q yielded an apparent change in B:L of  
269 10.3-fold because the biological and methodological differences compounded.

270 In addition to *Bacteroidetes*, Parameter Effects were calculated for multiple independent phyla across  
271 the 5 stool samples (Figure SI-5). As expected for diverse stool samples, the native phyla varied in  
272 different ways between the five fecal samples, reflecting their disparate biological compositions.  
273 However, the calculated Parameter Effects for methodological parameters also varied significantly  
274 between phyla. For instance, Proteobacteria and *Bacteroidetes* exhibited significant Parameter Effects  
275 for Extraction Kit, while Actinobacteria and Firmicutes did not exhibit significant Parameter Effects for

276 Extraction Kit (Figure SI-5). It was notable (while not unexpected) that these calculated Parameter  
277 Effects were taxa specific.



278

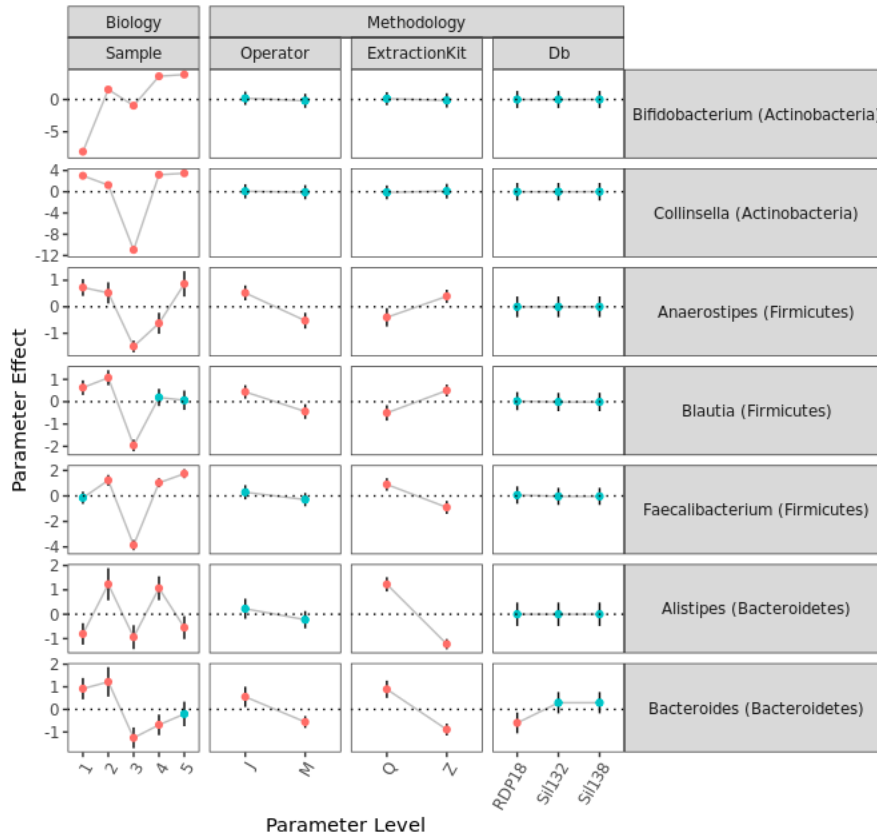
279 **Figure 6: Quantitative comparison between methodological parameters.** The  
280 *Bacteroidetes:Leifsonia* ratio was calculated for each dataset to enable comparison between  
281 protocols. The effect of each parameter was calculated by dividing the average ratio for all  
282 datasets at the denoted parameter level by the average ratio across all parameter levels, as  
283 shown in Equation 1. This parameter effect was plotted as a fold change on a log2 scale, such  
284 that the horizontal line at 0 denotes the null hypothesis of no effect. The magnitude of the effect  
285 of protocol choices (e.g., Extraction Kit) can be directly compared between parameters and  
286 parameter levels. Data error bars showed 99 % confidence intervals, and the points that are  
287 statistically significantly different from the mean ( $p < 0.01$ ) were colored red.

288 Quantifying Parameter Effects: Selected genera

289 While the analyses of native microbes presented thus far have focused on the highest taxonomic  
290 classification (i.e., phyla), much published research has linked beneficial or deleterious biological

291 outcomes with more specific taxonomic clades.(22) Thus, we calculated Parameter Effects for a panel of  
292 genera previously posited as having particular interest for gut health and representing multiple genera  
293 within specified phyla (Figure 7). Overall, the methodological factors of lot and variable region did not  
294 exhibit significant Parameter Effects for any of the genera evaluated here (and have been omitted from  
295 the plot), while the factors: Operator, Extraction Kit and Database were sometimes significant.

296 Interestingly, while the two genera within the *Actinobacteria* phyla exhibited similar responses to the  
297 methodologies evaluated in this study, this was not the case within the *Firmicutes* or *Bacteroidetes*  
298 phyla. Within the *Firmicutes*, the *Faecalibacterium* genus was preferentially enriched by one extraction  
299 kit and exhibited no Parameter Effect for operator, while the *Blautia* and *Anaerostipes* genera were  
300 preferentially enriched by the other extraction kit and exhibited significant Parameter Effects for  
301 operator. Within *Bacteroidetes*, the *Alistipes* and *Bacteroides* genera exhibited similar Parameter Effects  
302 for extraction kit, but differed in the Parameter Effects for operator and database. These variations in  
303 the Parameter Effects within Phyla demonstrated significant phenotypic heterogeneity at the genus  
304 level and underscored the difficulty in identifying reference organisms as potential surrogates for other  
305 members within a taxonomic grouping. To further demonstrate the variation in Parameter Effect an  
306 expanded view of each taxonomic level (Phylum, Class, Order, Family) for the genera in Figure 7 is  
307 present in Figure S1-6. In some cases, trends stayed the same as you move down the taxonomic  
308 hierarchy (Figure S1-6e, operator) while in others significant differences were only observed at high  
309 taxonomic levels (Figure S1-6f, operator) or low taxonomic levels (Figure S1-6f, Extraction Kit). The scope  
310 of the current effort was limited to the specific parameters varied and the parameter levels employed;  
311 however, similar conclusions (i.e., taxa-specific Parameter Effects that can vary substantially within  
312 taxonomic groupings) should be expected for other MGS protocols.



313

314

315

316

317

318

319

320

321

**Figure 7: The effect of methodological parameter choices impacts the quantitation of some Genera associated in the literature with the health of the gut microbiome.** Like Figure 6, parameter effect was calculated for the ratio of various taxa of interest to an internal control (*Leifsonia*). The parameter effect was calculated by dividing the average results for the specified parameter level for the specified ratio by the average results across all factor levels for the specified ratio. Data error bars show a 99 % confidence interval, and statistically significant points ( $p < 0.01$ ) are colored red. The effects broken out by taxonomic levels are plotted in SI Figures 6a-6g.

322

Quantifying Parameter Effects: All Taxa

323

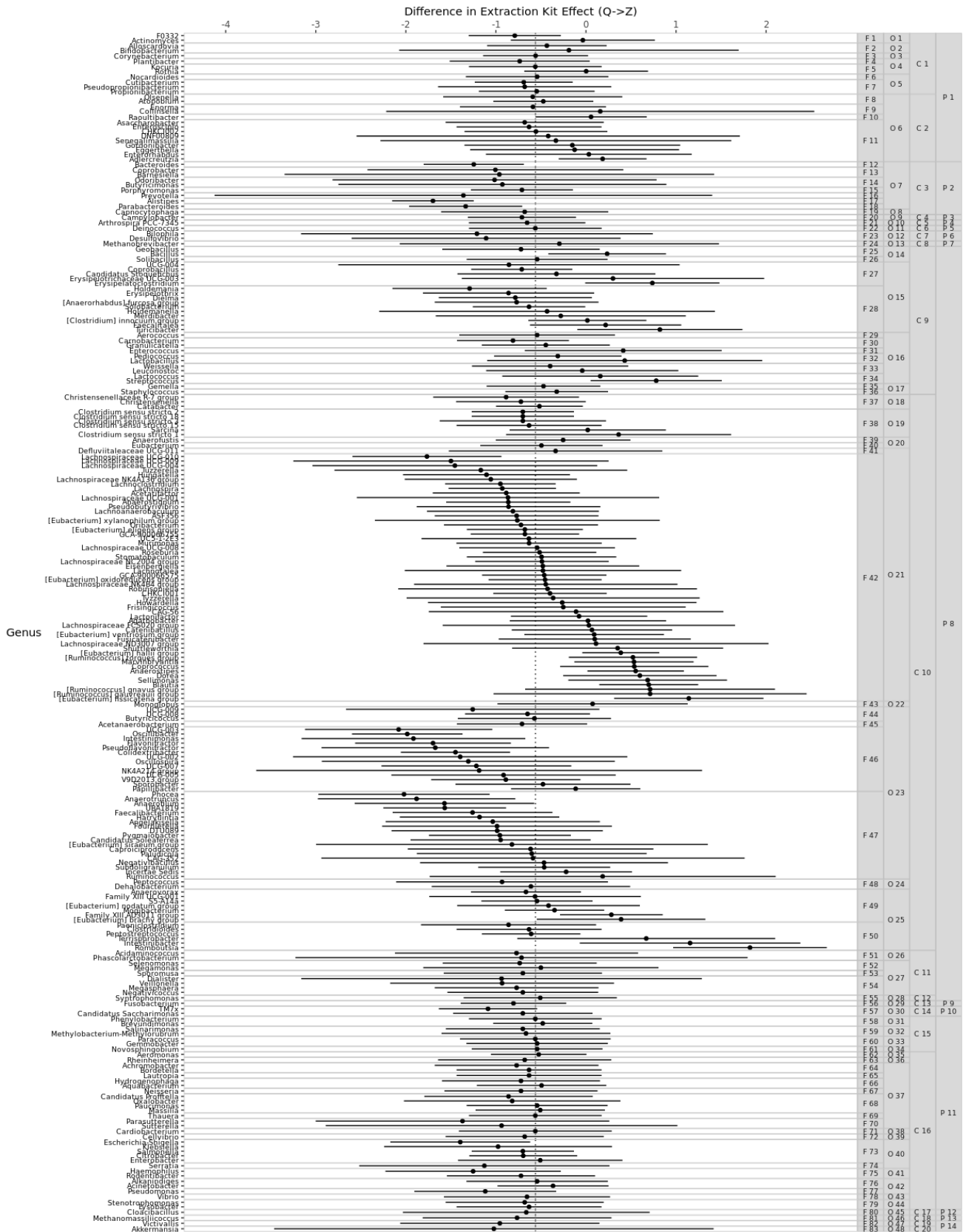
Using ratios with the internal standard, we calculated Parameter Effects for each individual taxa native

324

to the fecal samples. The methodological parameter which generally exhibited the most significant

325 effect was Extraction Kit, and we focused on that Parameter Effect to assess the degree of variation  
326 between taxa (Figure 8). Additionally, discrepancies in taxa names between databases or even database  
327 versions made systematic comparisons of individual taxa challenging, so this quantification of the  
328 Parameter Effect for extraction kit was performed within the Silva 138 database.

329 For some taxa, such as those in the family Clostridiaceae (Figure 8, F38), there appeared to be little or no  
330 effect of extraction kit, with all 99% confidence intervals overlapping with the mean and no statistical  
331 difference between the 6 genera. In comparison, the Peptostreptococcaceae family (Figure 8, F50)  
332 exhibits a wide range of extraction kit effects (and statistically significant differences) between the 6  
333 genera. In general, families with a greater number of genera present in the sample showed a greater  
334 range of extraction kit effects. While this observation was not surprising, it was notable that even closely  
335 related taxa could be impacted differently by methodological choices such as extraction kit. Figure 8  
336 would certainly look different in the details for other methodological choices or MGS workflows, but this  
337 demonstrated that in general the effects of methodological parameters (i.e, method bias) could vary  
338 substantially between taxa, even closely-related taxa, and cast doubt on the ability of well-characterized  
339 ‘reference species’ to serve as adequate surrogates.(23) This observed lack of consistency at the genera  
340 level has important implications for developing mock communities that will adequately stress a  
341 workflow (24). Instead, these data suggested that method bias for particular taxa-of-interest should be  
342 measured through experimental characterization of the protocols being considered.



<u>P- Phylum</u>	<u>C- Class</u>	<u>O- Order</u>	<u>F- Family</u>
1 Actinobacteria	1 Actinobacteria (Actinomycetota)	5 Propionibacteriales	7 Propionibacteriaceae
	2 Coriobacteriia	6 Coriobacteriales	11 Eggerthellaceae
8 Firmicutes	9 Bacilli	15 Erysipelotrichales	27 Erysipelatoclostridiaceae 28 Erysipelotrichaceae
	10 Clostridia	18 Christensenellales	37 Christensenellaceae
		19 Clostridiales	38 Clostridiaceae
		21 Lachnospirales	42 Lachnospiraceae
		23 Oscillospirales	44 Butyricoccaceae 46 Oscillospiraceae 47 Ruminococcaceae
		25 Peptostreptococcales- Tissierellales	49 Anaerovoracaceae 50 Peptostreptococcaceae
	11 Negativicutes	27 Veillonellales- Selenomonadales	54 Veillonellaceae
11 Proteobacteria	16 Gammaproteobacteria	37 Burkholderiales	68 Oxalobacteraceae
		40 Enterobacterales	73 Enterobacteriaceae

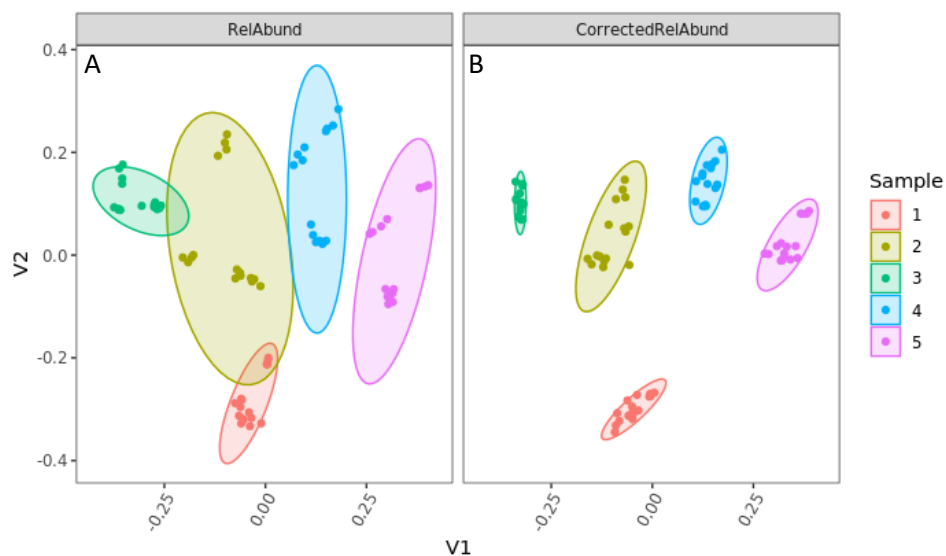
344 *Figure 8: The Effect of Extraction Kit on Different Genera. For each genera in the fecal samples,*  
345 *the difference in the Parameter Effects of the two Extraction Kits ( $Effect_{Kit=Z} - Effect_{Kit=Q}$ )*  
346 *was calculated. (Leifsonia was used as the internal standard.) Error bars show a 99 % confidence*  
347 *interval. Genus name is given on the left hand of the figure, with corresponding Phylum, class,*  
348 *order, and family groupings on the right. The table below the figure provides the taxonomic key*  
349 *for all Families with  $\geq 3$  Genera, and the full taxonomic key (249 genera) is included in the*  
350 *Supporting Information (SI Table 1).*

### 351 Correcting for Methodological Bias

352 Mathematically, quantification of Parameter Effect provided a mechanism for achieving comparability  
353 between differing protocols.(8) The observed differences between protocols in the observed relative  
354 abundances of individual taxa (i.e., method bias) could be computationally brought into agreement by  
355 correcting for the taxa-specific Parameter Effect. To validate this experimentally, Genus-level Parameter  
356 Effects were calculated for the methodological parameters of extraction kit and operator. These two  
357 parameters were selected because they had previously exhibited significant Parameter Effects for  
358 multiple tested Genera. Additionally, the Parameter Effects were only calculated using data from a  
359 single stool sample to simulate the ability of this approach to improve MGS characterization of  
360 previously untested samples. Stool sample 3 was used to calculate the Parameter Effects because it  
361 exhibited the most unique Genera. Using the calculated Parameter Effects, the observed abundances of  
362 Genera in all five samples were then computationally corrected to account for the protocol employed  
363 (i.e., which extraction kit and operator were specified). The initial- and corrected- relative abundances  
364 were used to calculate Bray-Curtis Dissimilarities between samples across all protocols as seen in the  
365 principal coordinate analysis plot colored by Sample, with the variability between protocols  
366 approximated with a 95 % data ellipse (Figure 9). The effect of the computational correction was evident

367 in the reduced size of the data ellipses for each sample, demonstrating a significant decrease in the  
368 observed dissimilarities between data collected under differing protocols.

369 It was notable that the improved comparability (reduced method bias) between protocols was observed  
370 for all five stool samples, even samples 1, 2, 4, and 5 that were not included in the Parameter Effect  
371 calculations. The ellipse centroids did not shift with computational correction because the Parameter  
372 Effect calculation was centered around a fold-change of 0 (no effect), and the arms of the experimental  
373 design were evenly balanced. It should also be noted that this computational correction did not  
374 necessarily improve the resulting data with respect to the actual biological composition of the samples.  
375 Rather, the data have been made more comparable between protocols by accounting for their taxa-  
376 specific biases. In this way, the dissimilarities within the datasets for each sample were reduced. This  
377 result strongly supported the idea that MGS data collected under differing protocols could be  
378 harmonized through the systematic analysis of common samples (e.g., reference materials).



379

380 ***Figure 9: Correcting for Extraction Kit and Operator Biases.*** Genus level Parameter Effects were  
381 *calculated for Extraction Kit and Operator using data stool sample 3, and the Relative*  
382 *Abundances of individual Genera in all five samples were computationally corrected to account*

383 *for the protocol used (i.e., which extraction kit and operator were specified). Bray-Curtis*  
384 *dissimilarities were calculated between all datasets and displayed on principal coordinate plots*  
385 *of (A) Relative Abundance or (B) corrected Relative Abundance, colored by fecal sample. 95 %*  
386 *data ellipses are shown to highlight the amount of dispersion observed between differing*  
387 *protocols for both conditions.*

## 388 **Conclusions**

389 The impact of methodological variables on MGS results is generally well recognized. In the current  
390 manuscript we combined a full factorial experimental design strategy with ratiometric analyses to  
391 quantify the impact of specific methodological choices. When the F:B ratio was calculated for individual  
392 stool samples, a >4-fold divergence in results was observed between different protocols (Figure 4).  
393 Replicate analyses confirmed that MGS results were reproducible run-to-run, so the observed variations  
394 were attributed to bias associated with specific methodological differences in the analysis protocols.  
395 Furthermore, this methodological bias was similar in magnitude to the observed differences between  
396 the biologically distinct samples.

397 Using the full factorial design, the Parameter Effect could be calculated to directly compare between  
398 specific biological and methodological parameters (Figure 6). Among the 5 methodological parameters  
399 evaluated in the current investigation, extraction kit and operator were statistically significant, while  
400 stool lot, 16S variable region, and database were not statistically significant. The observation that  
401 extraction kit introduced significant measurement bias was not novel; however, the benefit of the full-  
402 factorial experimental design was that the Parameter Effect for each variable could be quantified and  
403 directly compared to the Parameter Effect associated with biological variations. For example, for  
404 extraction kits specifically, the observed Bacteroidetes measurement bias was sufficiently large to mask  
405 real differences between distinct biological samples. This observation raises questions about attempts to

406 compare MGS data collected under different protocols. Somewhat surprisingly, operator also yielded a  
407 Parameter Effect that was weak but statistically significant and showed that human variability can still  
408 impact results even with locked-down protocols and identical lab equipment.

409 By combining the full factorial design with ratiometric analysis against internal control organisms,  
410 Parameter Effects could be calculated for individual taxa-of-interest native to the stool samples. For a  
411 panel of 7 genera posited in the literature to be associated with specific gut-health outcomes, we  
412 observed a large diversity of Parameter Effects (Figure 7), with some taxa exhibiting minimal  
413 measurement bias for the methods evaluated here, while several others exhibiting significant bias  
414 across multiple methodological decisions (Figure 7). Notably, Parameter Effects did not appear to group  
415 by taxonomic clade (e.g., phylum) and this somewhat unexpected finding warrants further investigation.  
416 Focusing specifically on the extraction kit, some related taxa exhibited similar Parameter Effects while  
417 other closely related taxa (e.g., genera within a family) exhibited quite divergent Parameter Effects. This  
418 diversity of methodological bias could have important implication for reference material development  
419 (i.e., mock communities) where microbial mixtures often contain only 10s of strains. These findings  
420 suggest that the generalizability of method performance based on a small subset of strains could be very  
421 limited and justifies efforts to develop more complex mock communities and natural surrogate  
422 reference materials that contain many more taxa and are compositionally similar to samples of interest.  
423 Indeed, using genera-specific Parameter Effects calculated for one of the samples investigated here,  
424 method bias between divergent protocols was significantly reduced for all samples.

425 The current report details some of the measurement challenges within the metagenomic space and  
426 presents a method to systematically assess sources of bias. However, there are a few limitations worth  
427 noting. As designed, this study only examined a small subset of methodological variables, while holding  
428 all others constant. This enabled the full-factorial experimental design and allowed Parameter Effects to

429 be quantified for those variables. This study was not meant to serve as a comprehensive review of  
430 available kits, library preparations, NGS technology, databases, etc. largely due to a second limitation of  
431 the full factorial design: it is resource intensive and expensive. A significant amount of starting material  
432 was needed to ensure a homogenous sample supply, a notable challenge for clinical and archival  
433 samples that are routinely limited. As reagent costs and personnel time requirements quickly scale in  
434 increasing numbers of methodological variables, a partial-factorial experimental design may be  
435 sufficient in many cases. Nevertheless, as technology advances, the quantitative determination of the  
436 bias associated with specific changes to the workflow can make future results backward compatible.

437 In summary, we presented an experimental design and analysis strategy that enables direct quantitative  
438 comparison of different biological and methodological parameters. We encourage MGS research groups  
439 to consider characterizing the impact of their in-house protocol choices on their MGS results. As more  
440 labs seek to improve the comparability of their MGS results between divergent protocols, this approach  
441 provides a quantitative strategy for assessing methodological bias and shows how data collected under  
442 differing protocols can be harmonized by widespread analysis of a common sample (e.g., reference  
443 material).

444 **Methods** ([Mosaic Materials-and-methods.docx](#))

#### 445 **Samples and Sample Handling**

446 Samples used in this study were previously described.<sup>(15)</sup> Briefly, five samples were generated from five  
447 different donors. Each sample was prepared by pooling and homogenizing 4 consecutive bowel  
448 movements from the donor, combining homogenized material with Omnigene Gut Solution, spiking with  
449  $10^8$  CFU/mL of *Aliivibrio- fischeri* (Gram negative, formerly known as *Vibrio fischeri*) and *Leifsonia xyli*  
450 (Gram positive), and then aliquoting 1 mL aliquots. The final concentration of the material was  
451 100mg/mL. Samples were received as a set of five, one sample from each donor (labeled 1-5). For this

452 study we ordered two sets (designated Lot A and Lot B) of the samples. Samples were stored at -80 °C  
453 until the experiments were conducted. On the day of the experiment, samples were vortexed vigorously  
454 and then 100 µL from each sample (1-5) was aliquoted into five 1.5-mL microcentrifuge tubes. This  
455 process was repeated for both lots. Operator 1 used three aliquots of each sample (1-5) from Lot A, and  
456 two aliquots from each sample (1-5) from Lot B. Operator 2 used two aliquots of each sample (1-5) from  
457 Lot A, and three aliquots from each sample (1-5) from Lot B.

## 458 **DNA Extraction**

### 459 ***Genomic DNA Extraction Kit Q:***

460 Samples were extracted using two different methods. Method 1 was done with the QIAamp Fast DNA  
461 Stool Mini DNA prep protocol (Qiagen, Cat# 51504). Briefly, 1 mL of the Inhibit EX buffer was added to  
462 each of the samples and vortexed (Mo Bio Vortex – Genie 2 at max speed) for 1 minute. The samples  
463 were then centrifuged (Eppendorf Centrifuge 5417R) at 16,000×g for 1 minute to pellet the samples. A  
464 new 2 mL centrifuge tube was prepared containing 25 µL of proteinase K and 600 µL of the supernatant  
465 from the original tube was added to the tube containing proteinase K (supplied with kit). This was  
466 followed by the addition of 600 µL of Buffer AL; samples were vortexed for 15 seconds to mix. The  
467 samples were incubated for 10 minutes and 70 °C. After incubation, 600 µL of molecular biology grade  
468 ethanol (100 %, Sigma-Aldrich, E7023-1L) was added and vortexed. In 600 µL increments, the lysate was  
469 loaded onto the QIAamp spin column, and samples were then centrifuged 16,000×g for 1 minute. After  
470 centrifugation, columns were placed in new 2 mL collection tube and the flow-through was discarded.  
471 This process was repeated until all the lysate had been loaded onto the column. Next, the columns were  
472 washed first with 500 µL of Buffer AW1 (centrifuged 16,000×g for 1 minute), then 500 µL of Buffer AW2  
473 (centrifuged at 16,000×g for 3 minutes). Flow-through was discarded and tubes replaced between each  
474 wash. To remove any remaining wash solution, columns were placed in a new collection tube and

475 centrifuged at 16,000×g for 3 minutes. Columns were then placed in 1.5 mL labeled microcentrifuge  
476 tubes and 100 µL of Buffer ATE was added to the center of the membrane. The samples incubated for 5  
477 minutes and then centrifuged at 16,000×g for 1 minute to elute the DNA.

#### 478 ***Genomic DNA Extraction Kit Z:***

479 Method 2 was done with the ZR Fecal DNA miniprep (cat#D6010) from Zymo Research following the  
480 manufacturer protocol with minor modifications. Briefly, each sample was combined with lysis buffer  
481 and loaded onto the MoBio vortex Genie 2 with 2 mL tube adapter for 20 minutes at full speed. The lysis  
482 tube was centrifuged at 10,000×g for 1 minute. Then 400 µL of the lysate was added to the IV spin filter  
483 and centrifuged at 7,000×g for 1 minute. The process was repeated until all lysate had been passed  
484 through the spin filter. Filtrate was combined with 1,200 µL of the fecal DNA binding buffer and mixed  
485 well by pipetting. The mixture was transferred to the IIC column in 800 µL increments and centrifuged at  
486 10,000×g for 1 minute. The supernatant was discarded. The bound DNA was washed and eluted in 150  
487 µL of the elution buffer after incubation at room temperature for five minutes. The DNA was then  
488 transferred to the prepared IV-HRC spin filter and centrifuged at 8,000×g for three minutes.

489 Extracted DNA from both methods was quantified using DeNovix dsDNA High Sensitivity (Catalog  
490 number: KIT-DSDNA-HIGH-2) and analyzed on the DeNovix DS-11FX+ Spectrophotometer/ Fluorometer.

#### 491 **DNA Library Preparation**

492 All PCR reactions were carried out using Kapa HiFi HotStart ready mix 2x Master Mix (KapaBiosystems,  
493 cat# 07958935001) in 0.2 mL thin wall PCR plate (Fisher Scientific, AB0800150). The V34 and V4 variable  
494 regions were chosen for 16S amplicon sequencing on the fecal samples. V34 (340F (5'→3')  
495 TCGTCGGCAGCGTCAGATGTGTATAAGAGACAGCTACGGGNGGCWGCAG/ 806R (5'→3')  
496 GTCTCGTGGGCTCGGAGATGTGTATAAGAGACAGGGACTACNVGGGTWTCTAAT) and V4 (515F (5'→3')  
497 (TCGTCGGCAGCGTCAGATGTGTATAAGAGACAGGTGYCAGCMGCCGCGGTAA/ 806R (5'→3'))

498 GTCTCGTGGGCTCGGAGATGTGTATAAGAGACAGGGACTACNVGGGTWTCTAAT) were ordered as pre-  
499 mixed primer pairs and include the Illumina handles (underlined) on the 5' end (RxnReady Primer Pools,  
500 IDT). The 16S PCR amplification was carried out using the following reaction mixture: 12.5 µL of Kapa  
501 HiFi HotStart 2x Master Mix, 1 µL of 10 µM V34 primer pool, 12 ng of extracted DNA, and nuclease free  
502 water to bring the final volume to 25 µL. Amplification was carried out using the following protocol:  
503 Initial denature (95.0°C – 3:00), 25 cycles of denature (98.0 °C – 0:20), annealing (55.0 °C – 0:15),  
504 elongation (72.0 °C – 0:20), and a final extension (72.0 °C – 1:00). The 16S reactions were purified using  
505 SPRIselect beads (Beckman Coulter, cat# B23318) to select for amplicons. The beads were added at a  
506 0.8:1 ratio (20 µL of beads for 25 µL reaction), mixed by pipetting, and incubated on the bench for 1  
507 minute then place on magnet for 1 minute. Without disturbing beads, liquid was discarded and while  
508 keeping the plate on magnet, 180 µL of 80 % molecular grade ethanol (Sigma-Aldrich, cat# E7023-1L)  
509 was added. After 30 seconds, ethanol, was removed carefully to not disturb the beads (first with 200 µL  
510 pipette set to 180 µL, then with 20 µL pipette set to 20 µL to remove any excess ethanol. The plate was  
511 then removed from magnet and 40 µL of RNase/DNase free molecular biology grade water (Promega,  
512 cat# P1195) from was added and mixed by pipetting. Bead were incubated with the water for one  
513 minute, the place on magnet and allowed to sit for one minute. Careful not to disturb the beads, water  
514 was carefully removed and placed in a new 0.2 mL PCR plate.

515 Cleaned 16S amplicons were then barcoded using dual indexes from Illumina (Nextera XT Index kit V2,  
516 Illumina cat# 15052163) PCR reactions were set up similarly to the 16S reaction described above with  
517 the following exceptions, 2.5 µL of cleaned 16S amplicons were added as the DNA template and 5 µL of  
518 each indexing primer was used. The same thermocycler conditions were used as for the 16S PCR with  
519 the following modification: only 8 cycles were run. Indexed, 16S amplicons were then cleaned following  
520 the same SPRIselect protocol detailed above, with the following modifications: beads are added at a 1:1

521 ratio (25  $\mu$ L of beads for 25  $\mu$ L reaction. Following clean-up, samples were quantified by fluorescence  
522 (DeNovix dsDNA high sensitivity assay cat# Kit-DSDNA-HIGH-2).

### 523 **Sequencing**

524 10 ng of each sample were pooled for the sequencing. The samples were first pooled by operator. Pools  
525 were quantified by Fluorescence (DeNovix dsDNA high sensitivity assay cat# Kit-DSDNA-HIGH-2) and nM  
526 of each pool was calculated. DNA pools were stored at 4 °C until sequencing. On the day of sequencing,  
527 each operator diluted the pools to 4 nM. Pools were combined in equal volume and prepared for  
528 sequencing following the MiSeq System Denature and Dilute Libraries Guide (Document # 15039740  
529 v10, Protocol A). Denatured libraries were diluted to a final concentration of 12 pM and combined with  
530 5 % PhiX control (V3 cat# 15017666 from Illumina). Paired-end sequencing was performed on an  
531 Illumina MiSeq with 2 $\times$ 300 bp reads (MiSeq Reagent Kit v3 600-cycle, cat #: MS-102-3003).

### 532 **Data analysis**

533 Adapter trimming was done as part of the Illumina MiSeq Generate FASTQ workflow. Fastq files were  
534 imported to RStudio (R version 4.1.0) for processing. Cutadapt (2.8)(25) was used for primer trimming  
535 followed by DADA2 (1.20.0)(26) with taxonomic assignment. Briefly, Cutadapt identified and trimmed  
536 primers allowing for no more than 2 mismatches with the primer sequencing. Primer trimmed  
537 sequences were then fed into the DADA2 pipeline. Reads were filtered and trimmed using the following  
538 parameters: all reads (forwards and reverse) were trimmed to 250 bp, reads could contain no Ns, maxEE  
539 (expected error) of five for both forward and reverse reads, and the default truncQ (2). Default settings  
540 for error learning, dereplicating, merging, and chimera identification were used. The R code for the  
541 initial sequence processing and the statistical analyses shown herein are publicly available at  
542 <https://data.nist.gov/od/id/mds2-3092>.

### 543 **NIST Disclaimer**

544 Certain commercial equipment, instruments, or materials are identified in this paper to foster  
545 understanding. Such identification does not imply recommendation or endorsement by the National  
546 Institute of Standards and Technology, nor does it imply that the materials or equipment identified are  
547 necessarily the best available for the purpose. The reference materials used in this study were not  
548 certified by NIST and are not official [NIST Reference Materials](#).

#### 549 **Ethics approval and Consent to Participate**

550 All work was reviewed and approved by the U. S. National Institute of Standards and Technology (NIST)  
551 Research Protections Office. This study (protocol #: MML-2019-0135) was determined to be “not human  
552 subjects research” as defined in the Common Rule (45 CFR 46, Subpart A).

#### 553 **Consent for Publication**

554 Not applicable.

#### 555 **Availability of Data and Materials**

556 All metagenomic sequencing results and the code used for analyses in this manuscript are available  
557 online (<https://data.nist.gov/od/id/mds2-3092>). Remaining units of the fecal materials used for this  
558 project are available for purchase from The BioCollective.

#### 559 **Competing Interests**

560 The authors declare no competing interests.

#### 561 **Authors Contributions**

562 S.P.F., S.L.S., J.G.K., J.J.F., and S.A.J. conceptualized the project and designed experiments. J.N.D. and  
563 M.E.H. analyzed samples. S.P.F. and S.L.S. analyzed data. S.P.F., S.L.S., J.G.K., and S.A.J. wrote the

564 manuscript and prepared figures. All authors approved the final manuscript except for J.J.F. who passed  
565 away before the completion of the project.

## 566 **Acknowledgements**

567 The authors wish to acknowledge Sheng Lin-Gibson, Kirsten H. Parratt, Alshae R. Logan, and Lisa M.  
568 Stabryla for critical feedback during manuscript review.

## 569 **References**

- 570 [1. Simner PJ, Miller S, Carroll KC. 2018. Understanding the Promises and Hurdles of Metagenomic](#)  
571 [Next-Generation Sequencing as a Diagnostic Tool for Infectious Diseases. \*Clinical Infectious\*](#)  
572 [Diseases 66:778-788.](#)
- 573 [2. Goldberg B, Sichtig H, Geyer C, Ledebner N, Weinstock GM. 2015. Making the Leap from](#)  
574 [Research Laboratory to Clinic: Challenges and Opportunities for Next-Generation Sequencing in](#)  
575 [Infectious Disease Diagnostics. \*Mbio\* 6.](#)
- 576 [3. Chiu CY, Miller SA. 2019. Clinical metagenomics. \*Nature Reviews Genetics\* 20:341-355.](#)
- 577 [4. Sczyrba A, Hofmann P, Belmann P, Koslicki D, Janssen S, Droge J, Gregor I, Majda S, Fiedler J,](#)  
578 [Dahms E, Bremges A, Fritz A, Garrido-Oter R, Jorgensen TS, Shapiro N, Blood PD, Gurevich A, Bai](#)  
579 [Y, Turaev D, DeMaere MZ, Chikhi R, Nagarajan N, Quince C, Meyer F, Balvociute M, Hansen LH,](#)  
580 [Sorensen SJ, Chia BKH, Denis B, Froula JL, Wang Z, Egan R, Don Kang D, Cook JJ, Deltel C,](#)  
581 [Beckstette M, Lemaitre C, Peterlongo P, Rizk G, Lavenier D, Wu YW, Singer SW, Jain C, Strous M,](#)  
582 [Klingenberg H, Meinicke P, Barton MD, Lingner T, Lin HH, Liao YC, et al. 2017. Critical](#)  
583 [Assessment of Metagenome Interpretation-a benchmark of metagenomics software. \*Nat\*](#)  
584 [Methods 14:1063-1071.](#)
- 585 [5. Lim MY, Song EJ, Kim SH, Lee J, Nam YD. 2018. Comparison of DNA extraction methods for](#)  
586 [human gut microbial community profiling. \*Syst Appl Microbiol\* 41:151-157.](#)

- 587 [6. Sinha R, Abu-Ali G, Vogtmann E, Fodor AA, Ren B, Amir A, Schwager E, Crabtree J, Ma S, Abnet](#)  
588 [CC, Knight R, White O, Huttenhower C. 2017. Assessment of variation in microbial community](#)  
589 [amplicon sequencing by the Microbiome Quality Control \(MBQC\) project consortium. Nature](#)  
590 [Biotechnology 35:1077-1086.](#)
- 591 [7. Zaiko A, Greenfield P, Abbott C, von Ammon U, Bilewitch J, Bunce M, Cristescu ME, Chariton A,](#)  
592 [Dowle E, Geller J, Gutierrez AA, Hajibabaei M, Haggard E, Inglis GJ, Lavery SD, Samuiloviene A,](#)  
593 [Simpson T, Stat M, Stephenson S, Sutherland J, Thakur V, Westfall K, Wood SA, Wright M, Zhang](#)  
594 [G, Pochon X. 2022. Towards reproducible metabarcoding data: Lessons from an international](#)  
595 [cross-laboratory experiment. Molecular Ecology Resources 22:519-538.](#)
- 596 [8. McLaren M, Willis A, Callahan B. 2019. Consistent and correctable bias in metagenomic](#)  
597 [sequencing experiments. Elife 8.](#)
- 598 [9. Shkoporov AN, Ryan FJ, Draper LA, Forde A, Stockdale SR, Daly KM, McDonnell SA, Nolan JA,](#)  
599 [Sutton TDS, Dalmaso M, McCann A, Ross RP, Hill C. 2018. Reproducible protocols for](#)  
600 [metagenomic analysis of human faecal phageomes. Microbiome 6:68.](#)
- 601 [10. Magne F, Gotteland M, Gauthier L, Zazueta A, Pessoa S, Navarrete P, Balamurugan R. 2020. The](#)  
602 [Firmicutes/Bacteroidetes Ratio: A Relevant Marker of Gut Dysbiosis in Obese Patients? Nutrients](#)  
603 [12:1474.](#)
- 604 [11. Trutna LS, P; del Castillo, E; Moore, T; Hartley, S; Hurwitz, A. 2012. In Tobias PT, L \(ed\),](#)  
605 [NIST/SEMATECH e-Handbook of Statistical Methods.](#)  
606 <https://www.itl.nist.gov/div898/handbook/>. doi:doi.org/10.18434/M32189.
- 607 [12. Galyean AA, Filliben JJ, Holbrook RD, Vreeland WN, Weinberg HS. 2016. Asymmetric flow field](#)  
608 [flow fractionation with light scattering detection - an orthogonal sensitivity analysis. J](#)  
609 [Chromatogr A 1473:122-132.](#)

- 610 [13. Scott JHJ. 2007. Accuracy issues in chemical and dimensional metrology in the SEM and TEM.](#)  
611 [Measurement Science and Technology 18:2755-2761.](#)
- 612 [14. Lee Y, Filliben JJ, Micheals RJ, Jonathon Phillips P. 2013. Sensitivity analysis for biometric](#)  
613 [systems: A methodology based on orthogonal experiment designs. Computer Vision and Image](#)  
614 [Understanding 117:532-550.](#)
- 615 [15. Forry SP, Servetas SL, Kralj JG, Soh K, Hadjithomas M, Cano R, Carlin M, de Amorim MG, Auch B,](#)  
616 [Bakker MG, Bartelli TF, Bustamante JP, Cassol I, Chalita M, Dias-Neto E, Duca AD, Gohl DM,](#)  
617 [Kazantseva J, Haruna MT, Menzel P, Moda BS, Neuberger-Castillo L, Nunes DN, Patel IR, Peralta](#)  
618 [RD, Saliou A, Schwarzer R, Sevilla S, Takenaka IKTM, Wang JR, Knight R, Gevers D, Jackson SA.](#)  
619 [2023. doi:10.1101/2023.04.28.538741.](#)
- 620 [16. Morton JT, Marotz C, Washburne A, Silverman J, Zaramela LS, Edlund A, Zengler K, Knight R.](#)  
621 [2019. Establishing microbial composition measurement standards with reference frames. Nat](#)  
622 [Commun 10:2719.](#)
- 623 [17. Fernandes AD, Reid JN, Macklaim JM, McMurrrough TA, Edgell DR, Gloor GB. 2014. Unifying the](#)  
624 [analysis of high-throughput sequencing datasets: characterizing RNA-seq, 16S rRNA gene](#)  
625 [sequencing and selective growth experiments by compositional data analysis. Microbiome 2:15.](#)
- 626 [18. Gloor GB, Wu JR, Pawlowsky-Glahn V, Egozcue JJ. 2016. It's all relative: analyzing microbiome](#)  
627 [data as compositions. Ann Epidemiol 26:322-9.](#)
- 628 [19. Houtman TA, Eckermann HA, Smidt H, de Weerth C. 2022. Gut microbiota and BMI throughout](#)  
629 [childhood: the role of firmicutes, bacteroidetes, and short-chain fatty acid producers. Scientific](#)  
630 [Reports 12.](#)
- 631 [20. Mariat D, Firmesse O, Levenez F, Guimaraes VD, Sokol H, Dore J, Corthier G, Furet JP. 2009. The](#)  
632 [Firmicutes/Bacteroidetes ratio of the human microbiota changes with age. BMC Microbiology 9.](#)

- 633 [21. Takezawa K, Fujita K, Matsushita M, Motooka D, Hatano K, Banno E, Shimizu N, Takao T, Takada](#)  
634 [S, Okada K, Fukuhara S, Kiuchi H, Uemura H, Nakamura S, Kojima Y, Nonomura N. 2021. The](#)  
635 [Firmicutes/Bacteroidetes ratio of the human gut microbiota is associated with prostate](#)  
636 [enlargement. Prostate 81:1287-1293.](#)
- 637 [22. Janssens Y, Nielandt J, Bronselaer A, Debunne N, Verbeke F, Wynendaele E, Van Immerseel F,](#)  
638 [Vandewynckel YP, De Tre G, De Spiegeleer B. 2018. Disbiome database: linking the microbiome](#)  
639 [to disease. BMC Microbiol 18:50.](#)
- 640 [23. McLaren MR, Nearing JT, Willis AD, Lloyd KG, Callahan BJ. 2022.](#)  
641 [doi:10.1101/2022.08.19.504330.](#)
- 642 [24. Sergaki C, Anwar S, Fritzsche M, Mate R, Francis RJ, MacLellan-Gibson K, Logan A, Amos GCA.](#)  
643 [2022. Developing whole cell standards for the microbiome field. Microbiome 10:123.](#)
- 644 [25. Martin M. 2011. Cutadapt removes adapter sequences from high-throughput sequencing reads.](#)  
645 [EMBnetjournal 17.](#)
- 646 [26. Callahan BJ, McMurdie PJ, Rosen MJ, Han AW, Johnson AJ, Holmes SP. 2016. DADA2: High-](#)  
647 [resolution sample inference from Illumina amplicon data. Nat Methods 13:581-3.](#)
- 648