1    **Impact of sequence variant detection and bacterial DNA extraction methods on the**

2    **measurement of microbial community composition in human stool**

3

4    Riley Hughes[1], Zeynep Alkan[2], Nancy L. Keim[1,2] and Mary E. Kable[1,2]

5

6    **Affiliations:**

7    [1]University of California, Davis, Department of Nutrition, One Shields Ave., Davis, CA 95616,

8    USA

9    [2]USDA, Agricultural Research Service, Western Human Nutrition Research Center, Immunity

10   and Disease Prevention, 430 West Health Sciences Dr., Davis, CA, 95616, USA

11

12   **Corresponding Author:**

13   Mary E. Kable

14   mary.kable@ars.usda.gov

15 **ABSTRACT**

16 **Background:** The human gut microbiome has been widely studied in the context of human

17 health and metabolism, however the question of how to analyze this community remains

18 contentious. This study compares new and previously well established methods aimed at

19 reducing bias in bioinformatics analysis (QIIME 1 and DADA2) and bacterial DNA extraction of

20 human fecal samples in 16S rRNA marker gene surveys.

21 **Results:** Analysis of a mock DNA community using DADA2 identified more chimeras (QIIME

22 1: 0.70% of total reads vs DADA2: 1.96%), fewer sequence variants , (QIIME 1: 1297.4 $\pm$ 98.88

23 vs. DADA2: 136.27 $\pm$ 11.35, mean $\pm$ SD) and correct taxa at a higher resolution of classification

24 (i.e. genus-level) than open reference OTU picking in QIIME 1. Additionally, the extraction of

25 whole cell mock community bacterial DNA using four commercially available kits resulted in

26 varying DNA yield, quality and bacterial community composition. Of the four kits compared,

27 ZymoBIOMICS DNA Miniprep Kit provided the greatest yield, with a slight enrichment of

28 *Enterococcus.* However, QIAamp Fast DNA Stool Mini Kit resulted in the highest DNA quality.

29 Mo Bio PowerFecal DNA Kit had the most dramatic effect on the mock community

30 composition, resulting in an increased proportion of members of the family *Enterobacteriaceae*

31 and genus *Eshcerichia* as well as members of genera *Lactobacillus* and *Pseudomonas.* The

32 presence of a sterile fecal matrix had a slight, but inconsistent effect on the yield, quality and

33 taxa identified after extraction with all four DNA extraction kits. Extraction of bacterial DNA

34 from native stool samples revealed a distinct effect of the DNA stabilization reagent DNA/RNA

35 Shield on community composition, causing an increase in the detected abundance of members of

36 orders *Bifidobacteriales*, *Bacteroidales*, *Turicibacterales*, *Clostridiales* and *Enterobacteriales*.

37    **Conclusion:** These results confirm that the DADA2 algorithm is superior to sequence clustering

38    by similarity to determine microbial community structure. Additionally, commercially available

39    kits used for bacterial DNA extraction from fecal samples have some effect on the proportion of

40    high abundance members detected in a microbial community, but it is less significant than the

41    effect of using DNA stabilization reagent, DNA/RNA Shield.

42

43    **INTRODUCTION**

44    Marker gene surveys utilizing PCR amplification of a short region of the bacterial 16S

45    rRNA gene from bacterial DNA extracted from environmental samples are becoming

46    increasingly affordable, leading to their ubiquitous implementation in nearly every aspect of

47    biological sciences research [1-8]. However, this method can be heavily affected by technical

48    bias, which is induced at each step in the experimental protocol required to generate marker gene

49    data including; sample handling, bacterial DNA extraction, PCR amplification, sequencing and

50    bioinformatics anlysis [9, 10]. PCR conditions and primer choice can impact biases during the

51    amplification process, which has downstream effects on library preparation and formation of

52    chimeric sequences [11, 12]. However, the two most important sources of technical bias, which

53    can be relatively easily controlled, in marker gene surveys are DNA extraction and

54    bioinformatics analysis [13].

55    Clustering into operational taxonomic units (OTUs) has been one of the primary

56    bioinformatic methods used to group and identify bacterial taxa in samples in metagenomics and

57    marker gene based sequencing analyses. This method utilizes percent sequence similarity to

58    group sequences into operational taxonomic units (OTUs). The common similarity threshold

59    used to define these OTUs is 97%, which is based on a study showing that most strains have

60    97% 16S rRNA sequence similarity [14]. From each OTU cluster, a single sequence is selected

61    as the "representative sequence" and is classified based on a reference database. All sequences

62    within the OTU cluster are then given the same taxonomic classification. OTU clustering offers a

63    computational benefit, reducing millions of reads into only thousands of OTUs, allowing for

64    rapid analysis of datasets [15].

65         However, the OTU clustering method has long been understood to have a number of

66    drawbacks [15]. For instance, percent sequence similarity can overestimate the evolutionary

67    similarity between sequences, leading to inappropriate clustering of sequences. Additionally, the

68    standard 97% sequence similarity used to define species is an approximation and varies between

69    taxa [16]. Higher rate of false positives (i.e. identification of taxa not present in the sample) as

70    well as poor sequence and taxonomic resolution have also been cited as issues with OTU

71    clustering [17, 18]. With the development of a number of new algorithms for sequence variant

72    identification including Devisive Amplicon Denoising Algorithm (DADA2), unoise2, minimum

73    entropy decomposition and Deblur [19-23], additional criticisms have surfaced regarding the

74    OTU clustering method [24, 25] and the need to conduct and publish independent direct

75    comparisons of methods has arisen.

76         Before sequences can even be analyzed and results affected by OTU clustering vs.

77    sequence variant detection, DNA extraction methods can heavily influence the proportion of

78    bacterial taxa detected in an environmental sample. Previous studies investigating the impact of

79    various DNA extraction methods on 16S rRNA analyses of stool microbial communities each

80    lack the combined use of a mock community in the relevant stool matrix background.

81    Additionally, the number of optimizable steps in DNA extraction protocols results in a near

82    infinte number of possible ways to execute this type of experiment. Most notably, previous

83    efforts to compare DNA extraction methods have indicated that the bead beating protocol tends

84    to be the source of greatest variation between kits [26-28], yet few if any have held this variable

85    constant during comparison. Finally, as technology evolves, new DNA extraction kits and

86    bioinformatics methods are constantly being developed. Therefore, the need to compare and

87    analyze new methods remains.

88        In this study we perform two important comparisons. First, we examine DADA2's core

89    denoising algorithm relative to the open reference OTU clustering method used in QIIME 1 to

90    confirm which method results in a more accurate classification of the taxa present in a predefined

91    mock community of bacteria. Second, we use a whole cell mock community in a sterile feces

92    background to compare four relevant DNA extraction methods [10, 13, 29] with standardized

93    speed and duration of bead beating.

94

95    **METHODS**

96    **Preparation of stool samples**

97        Whole stool samples were collected at home by three human subjects, placed in a cooler

98    containing ice and brought to the Western Human Nutrition Research Center within 12 h of

99    generation. Upon arrival at the facility, each sample was stored briefly at 4℃ until it could be

100   homogenized in a stomacher for three minutes and flash frozen on dry ice. These samples were

101   thawed, combined in equal amounts, mixed, and then divided into 2 pools.  The first pool, which

102   will be referred to as "native stool", contained 1 g of stool from each subject, combined by

103   homogenization in a stomacher twice for 5 min.  A portion of this 3 g mixture was set aside in

104   100 mg aliquots for DNA extraction. The remaining 1 g of the mixture was combined with 9 mL

105   of nucleotide stabilization reagent (DNA/RNA Shield, Zymo Research, Irvine, CA) by vortexing

106    and incubated at room temperature overnight before 250 mg aliquots were weighed out for DNA

107    extraction.

108        A "sterile" fecal sample was prepared as previously described [30] from the second pool,

109    which contained 3 g of stool from each subject. Briefly, the 9 g mixture was stirred together with

110    90 mL of boiling 30% hydrogen peroxide ($H_2O_2$) for 15 min. The boiled stool mixture was then

111    passed over a 0.22 µm vacuum filter (Sarstedt, Nümbrecht, Germany) to collect particulate

112    matter. Fecal particulate matter retained on the filter was then washed with sterile phosphate

113    buffered saline (DPBS, pH=7.0-7.3, ThermoFisher, Waltham, MA) in 100 mL batches until

114    $H_2O_2$ was no longer detected in the filtrate using detection strips (MQuant Peroxide Test,

115    MilliporeSigma, St. Louis, MO). This required 1 L of PBS. A total of 3.1 g of dry particulate

116    matter was collected from the filter surface and suspended in 4.5 mL sterile PBS to create a

117    sterile fecal matix. To create a mock stool sample with a known bacterial community, a 1.5g

118    aliquot of this sterile feces was homogenized in a stomacher for 2 min together with 1.125 mL of

119    commercially available whole cell mock community (ZymoBIOMICS Microbial Community

120    Standard, Zymo Research, Irvine, CA, lot number ZRC183430) and this mixture was set aside in

121    173 mg aliquots for DNA extraction (75 µL of mock community per 100 mg of stool). The

122    microbial strains included in the standard along with their theoretical relative abundances are

123    listed in Table S1. The remainder of the sterile fecal matrix was portioned into 100 mg aliquots

124    as blank controls for DNA extraction.

125

126    **Experimental design and bacterial DNA extraction**

127        A total of six sample types were prepared for DNA extraction; (1) kit blank with no

128    sample, (2) 75 µL mock community alone, (3) 100 mg sterile feces alone, (4) sterile feces with

6

129    mock community added totaling 173 mg as described above, (5) 100 mg native stool and (6) 25

130    mg native stool suspended in nucleotide stabilization reagent (DNA/RNA Shield, Zymo

131    Research, Irvine, CA) totaling 250 mg. Twelve aliquots of each sample type (three per kit) were

132    homogenized 5 times in bead tubes from three of the DNA extraction kits or bead tubes prepared

133    separately (described below) at 6.5 m/s for 1 min using a homogenizer (FastPrep-24 Classic

134    Instrument, MP Biomedicals).  Samples were rested on ice for three minutes between each

135    shaking interval. Bacterial DNA was then extracted using (1) QIAamp Fast DNA Stool Mini Kit,

136    (2) MO BIO PowerFecal DNA Kit, (3) ZR Fecal DNA Kit and (4) ZymoBIOMICS DNA

137    Miniprep Kit. For the QIAamp Fast DNA Stool Mini Kit, which contains no bead tubes, sterile 2

138    mL screw cap tubes containing 300 mg of 0.1 mm dimeter zirconia/silica beads (BioSpec

139    Products, Bartlesville, OK) were prepared separately and sterilized by autoclaving. After

140    homogenization by bead beating, the manufacturer's protocol was followed for each kit with the

141    following exceptions:

142        *All kits* - Wash and elution buffers were incubated on the column for 10 min prior to

143    centrifugation. Before the addition of elution buffer, columns were centrifuged for three minutes

144    with caps open in order to completely remove wash buffers.

145        *QIAamp Fast DNA Stool Mini Kit* – The protocol for "Isolation of DNA from Stool for

146    Pathogen Detection" in the QIAamp Fast DNA Stool Mini Handbook (03/2014) was used with

147    few modifications. Briefly, after bead-beating stool samples in 1 mL InhibitEx buffer, the

148    samples were heated at 95ºC for 5 minutes.  Centrifugation to remove particulate matter was

149    performed for 3 min on the whole sample and for an additional 3 min on the resulting

150    supernatant. A larger portion than recommended, 400 µL, of the clarified sample was transferred

151    to a new tube containing 30 µL proteinase K. Additional lysis was performed as described in the

152    manufacterer's protocol.  However, only 200 µL of lysate was added to the QIAamp spin

153    column. DNA was eluted with 30 µL buffer ATE.

154        *MO BIO (now QIAamp) PowerFecal DNA Kit –*DNA was eluted with 50 µL buffer C6.

155        *ZR Fecal DNA Kit (now Quick-DNA Fecal/Soil Microbe Miniprep Kit) –* DNA was

156    eluted in 50 µL DNAse free water.

157        *ZymoBIOMICS DNA Miniprep Kit –* DNA was eluted in 50 µL DNAse free water.

158

159    **Amplification and sequencing of 16S rRNA**

160        The 16S rRNA V4 region was amplified as previously described [31] using primers F515

161    and R806 [32]. A unique eight nucleotide Hamming code sequence was included on the 5' end

162    of F515 [33, 34] for amplification of each sample. Each 50 µL reaction mixture was composed

163    of 20 ng of template DNA, 1.5 U Ex Taq DNA polymerase (TaKaRa, Otsu, Japan), 100 nM of

164    forward primer, 100 nM reverse primer, 500 nM magnesium chloride, 200 nM dNTPs and 1X Ex

165    Taq buffer. Amplification was performed in triplicate for each sample with one cycle at 94°C for

166    3 min followed by 25 cycles of 94°C for 45 s, 50°C for 60 s, and 72°C for 60 s. A final extension

167    step was performed at 72°C for 10 min.  Equal volumes of each PCR reaction (40 µL) were

168    pooled and gel purified with the Wizard SV Gel and PCR cleanup system (Promega, Madison,

169    WI). Ligation of NEXTflex adapters (Bioo Scientific, Austin, TX) and 300-bp paired end

170    sequencing on an Illumina MiSeq instrument with MiSeq Reagent Kit v3 (Illumina) was

171    performed at the University of California, Davis (http://dnatech.genomecenter.ucdavis.edu/).

172        In order to eliminate the bias introduced by PCR amplification and sequencing from our

173    downstream analyses, a commercially available mock microbial community DNA standard

174    (ZymoBIOMICS™ Microbial Community DNA Standard, lot number ZRC187324), a sample

175    which we will refer to as Mock DNA was amplified and sequenced in the same manner as all

176    other experimental samples.  The DNA standard is a mixture of genomic DNA extracted and

177    quantified from pure cultures of eight bacterial and two fungal strains with the same theoretical

178    composition as the whole cell mock community described above. Metagenomic sequencing, was

179    performed by Zymo Research as part of their product quality assesment to determine the percent

180    relative abundance of the microbial strains in both the DNA and whole cell standards. Their

181    reported results are listed in Table S1.

182

183    **16S rRNA gene sequence analysis**

184        A summary of the methods used for analysis is described in Table 1. FASTQ files were

185    analyzed using QIIME version 1.9.1 [35], which will hereafter be referred to as QIIME 1, or

186    DADA2 version 1.4.0 [20]. R version 3.4.0 was used for all analyses. For the QIIME 1 analysis,

187    referred to throughout this manuscript, barcodes were extracted and the split_libraries_fastq.py

188    script was used for demultiplexing and quality filtering. Demultiplexing was performed only

189    with barcodes containing no sequencing errors, and quality filtering was performed at a Phred

190    quality threshold of 29. Chimeric sequences were identified with identify_chimeric_seqs.py

191    using usearch [36] and removed. The remaining DNA sequences were grouped into OTUs with

192    97% matched sequence identity by the use of pick_open_reference_otus.py. The default for open

193    reference OTU picking in QIIME is to use the first read as the representative sequence to form

194    the OTU clusters. In order to more closely imitate the DADA2 pipeline, this default behavior

195    was changed to use the most abundant sequence by passing a parameter file using the function

196    pick_rep_set.py (method most_abundant). Otherwise default parameters were used. Greengenes

197     13.8 was used as the reference database [37] for chimera checking, OTU picking, and taxonomy

198     assignment.

199

200     **Table 1.** Summary of bioinformatic methods

| Step | DADA2 1.4.0 | QIIME 1.9.1 |
|---|---|---|
| Input | demultiplexed fastq files + mapping file | fastq files + mapping file |
| Pre-Processing | - Filter and trim (trunclen=190, otherwise standard parameters)<br>- Dereplication | - Extract barcodes and remove primers<br>- Split libraries (demultiplex and quality filter – Q=29, otherwise default parameters) |
| Pick OTUs/variants | Sequence-variant inference (Sample inference/Denoising) | Open reference OTU picking (usearch, pick_rep_set method most_abundant, otherwise default parameters) |
| Remove chimeras | Remove bimeras | Remove chimeras* (usearch) |
| Assign taxonomy | Greengenes v13_8_99 | Greengenes v13_8_99 |

201     *Chimera/bimera removal comes before OTU picking in QIIME 1 but after Sample Inference in
202     DADA2
203

204          DADA2's denoising algorithm is based on pairwise comparison of sequences and uses

205     quality scores of the reads as well as the probability of various copy errors (transition

206     probabilities) that could be introduced during replication and sequencing. See Callahan et al. [20]

207     for full documentation of the core DADA2 algorithm. Methods used for DADA2 analyisis were

208     adapted from the DADA2 Pipline Tutorial (1.4) and DADA2 Frequently Asked Questions,

209     which are both currently available in the DADA2 GitHub documentation. Brielfly, prior to

210     analyses in DADA2, samples were demultiplexed using the QIIME 1.9.1 split_libraries_fastq.py

211     script with the following modifications from default parameters: -r (max bad run length) 999, -n

212     (max length of sequence) 999, -q (Phred quality threshold) 0, -p (min number of high quality

213     bases as fraction of read length) 0.0001 and --store_demultiplexed_fastq. This removed the

214     majority of quality filtering that is typically  implemented by the QIIME 1 pipeline using this

215    script and created individual fastq files for each sample. The demultiplexed files were used as the

216    input for DADA2.

217        Quality profiles of the reads were analyzed using the DADA2 function,

218    plotQualityProfile, to determine positions at which read quality greatly diminished. Reads were

219    then filtered and trimmed at the identified positions (truncLen=190) using the filterAndTrim

220    function with standard parameters (maxN=0, truncQ=2,and maxEE=2). Dereplication was then

221    used to identify all unique sequences present in the data set and determine the abundance of each

222    sequence. DADA2 also retains a summary of the quality information associated with each unique

223    sequence, using this to inform the error model of the subsequent denoising step, increasing its

224    accuracy [20]. DADA2's error model automatically filters out singletons, removing them before

225    the subsequent sample inference step. Quality of the error estimation was then visualized using

226    the plotErrors function to ensure good fit. Sample inference was performed using the inferred

227    error model and chimeric sequences were removed using the removeBimeraDenovo function. It

228    is relevant to note that DADA2 implements bimera removal after sample inference has been

229    performed, whereas removal of chimeras in QIIME 1 occurs before the OTU picking step. The

230    Greengenes 13.8 database was used to assign taxonomy using the assignTaxonomy function.

231

232    **Statistical Analysis**

233        OTU or sequence variant counts and rarefaction curves were determined on sequence

234    count files (referred to as sequence table and OTU table in DADA2 and QIIME 1 respectively)

235    generated by each analysis pipeline. These were determined using a count of the number of rows

236    in each output file that contained non-zero values, referred to as non-zero OTU/SV counts, for

237    each sample.

238   Analysis of the relative proportion of each bacterial taxa was made after the data were

239 rarefied at a sequencing depth of 50,000 sequences per sample for both QIIME 1 and DADA2.

240 The rarefied sequence variant counts were summed by taxonomic identification and differential

241 abundances between experimental groups were determined using LefSe [38]. This method

242 involves the Kruskal-Wallis (KW) sum-rank test between classes of data followed by (unpaired)

243 Wilcoxon rank-sum test to conduct pairwise tests among subclasses. LDA is then used to

244 estimate the effect size for each of the identified taxa. We used LEfSe (Galaxy Version 1.0) with

245 default paramters ($\alpha$ KW = 0.05; $\alpha$ Wilcoxon = 0.05; LDA score threshold = 2.0) as well as using

246 the all-against-all strategy for multi-class analysis. All other comparisons were made using either

247 Welch's t-test or Kruskal-Wallis (KW).

248

249 **RESULTS**

250 **The DADA2 denoising algorithm improves accuracy of bacterial community measurement.**

251   QIIME 1 and DADA2 were compared using 18,651,434 sequences generated by Illumina

252 MiSeq sequencing of 6 individual PCR amplifications of a microbial community DNA standard

253 (Mock DNA). After demultiplexing and quality filtering using QIIME 1, 790,502 total sequences

254 remained. Of these, 5,532 chimeras were identified using usearch, accounting for only 0.70% of

255 total sequences.  On the other hand, the trimming, denoising and dereplication steps of DADA2

256 resulted in 368 sequences (or inferred variants), which could be considered more equivalent to a

257 representative set of sequences picked by open reference OTU picking. Out of these sequences,

258 160 bimeras were identified, representing  43.48% of inferred variants, but only 1.96% of total

259 reads after dereplication, and filtering (1,354,268 reads), which is still nearly double the

260 percentage detected using usearch in QIIME 1.

261     QIIME 1 identified a much larger number of OTUs/SVs than DADA2 in Mock DNA

262     (QIIME 1: 1145.5 ± 68.73 vs. DADA2: 123.5 ± 8.12, mean ± SD) (**Figure 1A**). However,

263     DADA2 still greatly overestimated the number of non-zero variants relative to the expected

264     number of bacterial species present in the Mock DNA samples. Low abundance sequences

265     identified by DADA2 were investigated further. The Hamming distance of low abundant

266     sequences relative to more abundant sequence-variants they were split from fell in a range from

267     1 to 80, and quality scores at nucleotide positions used to determine a particular low abundance

268     sequence was unique relative to the more abundant sequence it was split away from were above

269     29. However, when BLAST was used to compare these low abundance sequences to those

270     available in the National Center for Biotechnology Information (NCBI) nucleotide database,

271     87% of uniqe sequences in the Mock DNA samples were exact matches (100% query cover,

272     100% identity)  to bacterial taxa that tend to be abundant in human stool samples, such as genera

273     *Bifidobacterium, Turicibacter, and Blautia*.

274          Rarefaction curves representing the discovery rate of unique sequences, potentially

275     attributed to new taxonomic units, as a function of sequencing effort (i.e. number of sequences)

276     [39], reflected the differences in non-zero OTU/SV counts between QIIME 1 and DADA2

277     (**Figure 1B and C**). As sequencing effort increases, QIIME 1 open reference OTU clustering

278     results in the detection of continually increasing numbers of unique sequences in Mock DNA

279     samples. However, the number of unique sequences detected by DADA2 does not increase with

280     sequencing effort in the same way as for QIIME 1, instead the number of unique sequences

281     detected levels out at approximately 50,000 sequences per sample.

282          While QIIME 1 identified many more OTUs/SVs than DADA2, rarefaction at 50,000

283     sequences per sample followed by removal of low abundance taxa (<1%) into a category termed

284    "Other", showed a similar taxonomic profile of the Mock DNA samples detected by both QIIME

285    1 and DADA2 (**Figure 2**). However, DADA2 identified correct taxa at a higher resolution of

286    classification (i.e. genus-level) with less redundancy (i.e. identification of the same taxa at

287    different levels of taxonomic classification, such as f_*Bacillaceae* and g_*Bacillus*) than QIIME 1.

288    More specifically, eight taxa were present at greater than 1% relative abundance as detected by

289    DADA2. Seven out of these eight were correctly identified to the genus level. The last variant

290    was correctly identified at the family level (e.g. f_*Enterobacteriaceae* includes *Salmonella*

291    *enterica*). QIIME 1 identified nine taxa present at greater than 1% relative abundance. Out of

292    these nine, four were redundant at different levels of phylogenetic resolution. These included

293    f__*Bacillaceae* and g__*Bacillus* as well as f__*Pseudomonadaceae* and g__*Pseudomonas*. All

294    nine taxa identified were present in the Mock DNA community (no spurious identification), but

295    two taxa remained classified only to the family level (f__*Enterobacteriaceae* and

296    f__*Listeriaceae*) (**Figure 2**). LefSe analysis showed significant differences in the majority of

297    taxa identified excluding only g__*Enterococcus* and g__*Staphylococcus*. Because of this

298    increased accuracy in taxonomic identification, the remainder of comparisons examining DNA

299    extraction kits were analyzed using DADA2.

300

301    **DNA yield and quality vary among extraction kits.**

302          The efficiency of four commercial DNA extraction kits was assessed using commercially

303    available whole cell mock community (Mock Community) and the whole cell mock community

304    spiked into sterilized fecal matrix (Mock Community in Sterilized Feces). There was a

305    significant difference among the kits in DNA yield (KW Mock Community $P = 0.02488$, Mock

306    Community in Sterilized Feces $P = 0.01556$) and quality (KW Mock Community p = 0.03781,

14

307    Mock Community in Sterilized Feces $P = 0.04358$) from both sample types. ZR Fecal and

308    ZymoBIOMICS delivered the highest quantity of DNA for both the whole cell Mock

309    Community alone (ZR Fecal average = 59.3 ng/uL, ZymoBIOMICS average = 58.8 ng/uL) and

310    Mock Community in Sterile Feces (ZR Fecal average = 39.9ng/uL,  ZymoBIOMICS average =

311    31.8 ng/uL). However, QIAamp delivered the highest quality DNA from both sample types

312    (A260/A280 = 2.5 and 1.86 respectively) (**Figure. 3A and B**). The DNA yield and quality were

313    also affected by the presence of the sterile feces matrix. Both decreased in the presence of the

314    matrix for each kit, except for QIAamp. However, the difference in yield was only significant in

315    the ZR Fecal (Welch's t-test $P = 0.02699$) and ZymoBIOMICS ($P = 0.008911$) kits and the

316    difference in quality was only significant in the ZR Fecal kit ($P = 0.03097$ ).

317         The yield obtained from blank samples followed the same trend as the yield obtained

318    from mock community samples.  It was significantly higher for ZR Fecal and ZymoBIOMICS

319    than it was for the other two protocols, reaching levels greater than 10 ng/uL for each of the two

320    kits. However, the number of bacterial sequences detected after PCR and sequencing in the

321    blanks were not significantly different among kits (**Figure 3C**, KW Blank p-value = 0.09234).

322

323    **Measurement of bacterial community composition is affected by DNA extraction protocol.**

324         In addition to DNA yield and quality, the proportion of bacterial taxa measured after

325    extraction with each kit was determined.  The relative proportions of taxa expected to most

326    closely represent reality were determined using the Mock DNA standard described above.

327    Weighted UniFrac distances between extracted samples and the Mock DNA samples were

328    visualized by principal componants analysis (**Figure 4A**) and summarized in boxplots (**Figure**

329    **4B**). Samples extracted with the Mo Bio kit had the greatest combined distance from Mock DNA

330  (mean=0.0429, median=0.0409) compared to the other kits (Mo Bio mean = 0.0107, median =

331  0.0121; ZymoBIOMICS mean = 0.0039, median = 0.0034; ZR Fecal mean = 0.0097, median =

332  0.0078). Distances were significantly affected by the presence of a sterile fecal matrix in all kits

333  examined (Mo Bio $P$ =1.36e-05; QIAamp $P < 2.2e-16$; ZymoBIOMICS $P = 1.887e-6$; ZR Fecal

334  $P = 0.0131$). In the case of the Mo Bio and ZR Fecal kits, the presence of a stool matrix

335  decreased the distance from Mock DNA (Mo Bio mean with matrix = 0.0378, mean without

336  matrix = 0.0479; ZR Fecal mean with matrix = 0.0078, mean without matrix = 0.0117).

337  However, the opposite phenomenon occurred for the Qiagen kit protocol (mean with matrix =

338  0.0162, mean without matrix = 0.0053) and the ZymoBIOMICS kit (mean with matrix = 0.0064,

339  mean without matrix = 0.0039).

340       LEfSe analysis identified the greatest number of significantly different taxa in Mock

341  community samples extracted with the Mo Bio kit.  This included an increased proportion of

342  members of the family *Enterobacteriaceae* and genus *Eshcerichia* as well as members of genera

343  *Lactobacillus* and *Pseudomonas* (Figure. 5). Mo Bio also enriched the "Other" category,

344  indicating enrichment in several other low abundance taxa. Relative to the Mock DNA,

345  decreased abundance of members of the phylum *Firmicutes*, including  order *Bacillales* and class

346  *Bacilli* and genus *Listeria,* though not genus *Bacillus* were detected in all extracted samples.

347  Members of the gram positive genus *Staphylococcus* were also proportionally decreased in the

348  extracted samples relative to Mock DNA. Mock community samples extracted by

349  ZymoBIOMICS showed significant enrichment of  genus *Enterococcus*.

350

351  **Use of nucleotide stabilization reagent significantly affects measurement of microbial**

352  **community composition.**

16

353     After assessing the performance of different pipelines and extraction kits on the mock

354     community, we looked to confirm the relative efficiency of each kit and further investigate the

355     effect of the nucleotide stabilization reagent, DNA/RNA Shield, using a representative pool of

356     natural or native stool samples (Native Stool and Native Stool with DNA/RNA Shield). DNA

357     yield was significantly different among kits for extraction from pooled native stool samples

358     similar to observations for the mock community samples above (KW p-value = 0.0329 native

359     stool; 0.01556 native stool in shield). Additionally, the presence of stabilization reagent affected

360     the amount of DNA recovered by each kit. For both kits from Zymo Research (ZR Fecal and

361     ZymoBIOMICS), the amount of DNA recovered per gram of stool was significantly increased

362     (p-value = 0.0002916 and 0.01315) in the presence of stabilization reagent (**Figure 6A**). This

363     was not true for the other two protocols which showed a decrease. Although, the decrease was

364     only significant for the QIAamp kit (p-value = 0.003795). The quality of DNA recovered was

365     also significantly different among kits for extraction of the Native Stool and Native Stool with

366     DNA/RNA Shield (KW p-value = 0.02871; 0.01879), with QIAamp again providing the highest

367     quality DNA (**Figure 6B**). However, the quality of DNA was only significantly affected by the

368     presence of DNA/RNA Shield during extraction with the Mo Bio PowerFecal Kit (p-value =

369     6.435e-05).

370     Principal coordinate analysis of weighted UniFrac distances showed that samples

371     clustered by stabilization reagent first and by DNA extraction method second (**Figure 7A and

372     B**). The impact of stabilization reagent on community composition was again greatest for the Mo

373     Bio kit (**Figure 7C**). However, analysis of the relative abundance of bacterial taxa present after

374     extraction with each kit showed significant differences in relative proportion of taxa enriched

375     between samples with and without DNA/RNA Shield across all extraction kits (**Figure. 8**)  Order

17

376    *Clostridiales* including family *Ruminococcaceae* and genera *Clostridium, Oscillospira,*

377    *Ruminococcus,* and *SMB53* as well as order *Bifidobacteriales* including genus *Bifidobacterium*

378    were significantly increased in native stool with DNA/RNA Shield. Order *Bacteroidales*

379    including families *Rikenellaceae* and *Porphyromonadaceae* and genera *Bacteroides* and

380    *Parabacteroides*; order *Turicibacterales* including genus *Turicibacter*; and order

381    *Enterobacteriales* including genus *Escherichia* were also enriched in samples with DNA/RNA

382    Shield. While not significant at the order level, other members of *Firmicutes* and *Actinobacteria*,

383    including genera *Dorea*, *Faecalibacterium*, *Eggerthella*, *Roseburia*, *Collinsella*, *Coprococcus*,

384    and *Blautia* were decreased in the presence of stabilization reagent.

385

386    **DISCUSSION**

387         The determination of microbial community structure composition in environmental

388    samples can be heavily affected by technical bias. As new methods are developed to deal with

389    errors induced by DNA extraction, sequencing and other analysis methods, it remains necessary

390    to empirically compare and validate each method using microbial standards. Here we have

391    shown that DADA2 provides a more accurate assessment of the microbial community both in

392    terms of the number of sequence-variants detected as well as the identity and phylogenetic

393    resolution of taxa present. Additionally, if bead beating speed and duration are held constant, the

394    commercially available kit used for bacterial DNA extraction from fecal samples has minimal

395    effects on the proportion of high abundance members detected in a microbial community, except

396    in the case of chemical incompatibility, which may be present between the Mo Bio kit and the

397    DNA stabilization reagent, DNA/RNA Sheild.

398     The reduced number of unique sequences, identifiable to a higher taxonomic resolution

399     detected using DADA2 relative to the QIIME 1 OTU clustering method was likely due to the

400     method of error detection employed by DADA2, which statistically determines the most likely

401     sequencing errors in a particular data set and then adjusts for them rather than rounding out by an

402     allowable percent error (typically 97%). However, a number of low abundant taxa were also

403     identified using DADA2 that were not present in the reference sequences for the mock

404     community used for analysis.  It should be noted that these taxa were detected without a stringint

405     quality filter settting applied to the filterAndTrim function in DADA2.  Therefore, it is possible

406     that their number could be reduced further with a more stringint quality filter setting.

407     Optimization parameters aside, many of these taxa were abundant in DNA that was extracted

408     from native stool samples at the same time as the mock community samples in this study.  This

409     indicates that some contamination of the Mock DNA sample occurred leading to a slightly

410     greater number of detected sequence variants than we expected. However, because the same

411     samples were analyzed using both pipelines, our conculsions regarding the improved accuracy of

412     DADA2 remain valid.

413     Subsequent to the selection of a bioinformatics pipeline for our analyses, we found that

414     DNA yield and quality varied among mock community samples and blanks from four

415     commercially available DNA extraction kits. Given that the number of bacterial sequences

416     detected in Zymo Research blanks were not significantly higher than in the other kits, it is

417     unclear why the DNA yield was high in blank samples extracted using these two kits. We

418     suppose that either, the chemistry involved in the Zymo Research kits results in absorbance at

419     A260, or that there is viral or fungal DNA contaminant in the kit, which was undetected by our

420     PCR protocol.

19

421    Within all four, both yield and quality were slightly impacted by the presence of sterile

422    fecal matrix. The trend for a reduction in yield in the presence of matrix in three out of the four

423    kits suggests that, as expected, the presence of physical impediments to bead-beating that tend to

424    be present in the stool matrix, primarily undigested food particles, likely inhibit the effectiveness

425    of the beads in disrupting bacterial cells [41, 42]. One exception was the QIAamp kit. However,

426    the composition of the microbial community in the context of sterile fecal matrix was more

427    dissimilar than mock community alone from the proportions predicted by our control. This did

428    not result in a significantly detectable change in the relative proportions of abundant taxa, but

429    insignificant increases in gram negative organisms and decreases in proportions of some gram

430    positive organisms were observed. This would be expected if decreased efficiency of bacterial

431    cell wall disruption by bead-beating occurred in the presence of the sterile fecal matrix. The Mo

432    Bio kit, on the other hand, displayed decreased yield in the presence of sterile fecal matrix, but

433    the microbial community composition tended to become more similar to the control than mock

434    community along. The garnet beads included in the Mo Bio kit were pulverized at the speed and

435    duration of shaking used in our protocol (see materials and methods).  It is therefore possible that

436    the very small broken particles of these beads disrupted bacterial cells so effectively that exposed

437    DNA was also pulverized and the presence of a fecal matrix helped prevent some of this

438    disruption.

439    Given that speed and duration of bead beating were held constant, the trends described

440    for yield and quality across the four extraction kits in the presence of sterile fecal matrix suggest

441    that the size, shape and composition of beads play role in the ability to sufficiently disrupt the

442    stool matrix and facilitate the detection of "realistic" proportions of bacterial taxa. A second

443    explanation for varying results across the kits, predominated by a slight decrease in nucleotide

444    quality, in the presence of the sterile fecal matrix might be PCR inhibitors, such as

445    carbohydrates, coming from the stool matrix, which are eliminated to differing degrees of

446    completeness by each kit and could also be affected by use of stabilization or preservation

447    reagent [43].

448        All four DNA extraction protocols showed a decreased relative abundance of the phylum

449    *Firmicutes* and genus *Staphylococcus* extracted from whole cell mock community relative to the

450    Mock DNA control sample. This may indicate that even the robust bead beating protocol used in

451    this study (see Methods) was not sufficient to fully lyse all gram positive organisms contained in

452    the stool samples.  However, as shown in Table S1, the relative abundances of taxa in the whole

453    cell mock community, as estimated by Zymo Research using metagenomics sequencing, differed

454    from that in the Mock DNA. It is possible that the difference, reported by Zymo, was caused by

455    biases introduced by the DNA extraction kit that they used to determine the abundance of their

456    own community. Therefore, we are unable to use the difference between the "measured" values

457    as the expected difference between the Mock DNA and mock microbial community samples as

458    this would simply be a comparison of their extraction and sequencing methods and our own.

459    Given that both sample types were prepared with the same theoretical proportions, our analysis

460    presumes that the Mock DNA is a close representation of the proportions in the whole cell mock

461    community. Under this assumption, the ZymoBIOMICS DNA Miniprep Kit was determined to

462    provide the closest representation of the "true" microbial community in a stool sample.  On the

463    other hand, the Mo Bio kit had the most distinct deviation from the expected microbial

464    community composition. In the mock microbial community, characterized by significant

465    increases in *Lactobacillus* and several gram negative organisms relative to the Mock DNA

466    control.

467     A native stool sample was used to determine the effect of DNA stabilization reagent on

468     the overall microbial community composition. An additional element contributing to the

469     differential community composition observed using the Mo Bio kit may be explained by

470     analyses which showed that native stool samples were most dramatically affected by the

471     presence of nucleotide stabilization reagent when extracted with Mo Bio. This indicates a

472     potential incompatability of the Mo Bio kit with the DNA/RNA Shield stabilization reagent,

473     which was also used to stabilize the commercially available Mock Microbial Community. This

474     putative chemical incompatibility may have affected the microbial community composition

475     observed in all DNA/RNA Shield-suspended samples extracted using the Mo Bio kit. This

476     includes the whole cell mock community samples, which are available only suspended in

477     DNA/RNA Shield. On the other hand, the increase in DNA yield per gram of stool in the

478     presence of stabilization reagent used together with the Zymo Research kits is perhaps

479     unsurprising because all components were manufactured by Zymo Research and were likely

480     optimized to be used together.  However, we have shown that the stabilization reagent can also

481     be used successfully with the QIAamp kit.  Although there is a decrease in yield, the reagent

482     does not cause a decrease in the DNA quality. It should be noted, however, that our analyses

483     show  the use of DNA/RNA Shield, alters the observed abundance of numerous taxa compared

484     to native stool and this should be taken into consideration when planning studies and  comparing

485     results from studies which differ in their use of stabilization reagent.

486     Development of best practices and standardized methods for  microbiota analysis  is

487     critical for the advancement of research in many fields including personalized nutrition, ecology,

488     and food science/safety. It will be necessary to perform similar experiments as new technologies

489     are developed in order to make informed choices when determining which methods will provide

490     the most accurate data.

491

492     **FIGURE LEGENDS**

493     **Figure 1. Variant counts resulting from QIIME 1 and DADA2 analyses.** A) Boxplot of

494     comparison between DADA2 and QIIME 1 OTU/sequence variant counts for Mock DNA. B)

495     Rarefaction curves showing differences in taxonomic discovery rate between DADA2 and

496     QIIME 1. Six individually amplified and sequenced Mock DNA samples were analyzed with

497     each pipeline.

498

499     **Figure 2. Relative taxonomic abundance of Mock Community DNA samples analyzed by**

500     **DADA2 and QIIME 1**. OTU and sequence variant counts were rarefied at 50,000 sequences per

501     sample for both groups. All taxa present at <1% abundance were grouped into the "Other"

502     category. Each bar represents six PCR amplifications of Mock DNA. +Signficantly enriched in

503     DADA2 analyzed samples. *Significantly enriched in QIIME 1 analyzed samples.

504

505     **Figure 3. Yield and quality of mock community DNA extracted by four commercial kits.**

506     Boxplots showing A) DNA yield (ng/uL), B) DNA quality (A260/A280) and C) raw sequence

507     counts obtained from whole cell mock community (Mock Community) or whole cell mock

508     community spiked into a sterile feces matrix (Mock Comm in Sterile Feces), sterile feces alone

509     (Sterile Feces) or no sample (Kit Blank) using each of four commercial DNA extraction kits

510     (MoBio, Qiagen, ZRFecal and ZymoBIOMICS). Three of each sample type are represented.

511

23

512    **Figure 4. Weighted UniFrac of distance between extracted mock community samples and**

513    **Mock DNA.** A) Principal coordinate analysis of weighted UniFrac distances among Mock

514    Community or Mock Community in Sterile Feces and Mock DNA (control) samples. B)

515    Boxplots summarizing the weighted UniFrac distance between Mock DNA and each extracted

516    sample type grouped by extraction kit (Mo Bio, QIAamp, ZR Fecal and ZymoBIOMICS).

517

518    **Figure 5. Relative taxonomic abundance of mock community taxa after extraction by four**

519    **commercial kits.** The proportions of taxa present in the Mock DNA sample are shown for

520    comparison (None). Each bar represents a summary of technical replicates (six Mock DNA

521    samples and three of each of the other sample types). MC is used to designate whole cell mock

522    community only and SF is used to designate mock community spiked into sterile feces matrix.

523    +Signficantly enriched in Mock DNA samples. *Significantly enriched in extracted samples.

524

525    **Figure 6. Yield and quality of pooled native stool bacterial community DNA extracted by**

526    **four commercial kits.** Boxplots showing A) DNA yield in ng/g of stool and B) quality

527    (A260/A280) obtained from pooled (three stool samples) native stool community (Native Stool)

528    and pooled native stool community suspended in nucleotide stabilization reagent (Native Stool in

529    DNA Shield) using each of four commercial DNA extraction kits (Mo Bio, QIAamp, ZR Fecal

530    and ZymoBIOMICS). Three of each sample type are represented.

531

532    **Figure 7. Relative composition of microbial communities in pooled native stool samples**

533    **with and without stabilization reagent.** Weighted UniFrac of distance between native stool

534    samples colored by A) Sample Type with (dark blue) or without (light blue) DNA/RNA Shield

24

535   B) Extraction Kit and C) Boxplot of the weighted UniFrac distance between Native Stool w/

536   DNA Shield samples and Native Stool w/o DNA Shield samples, separated by DNA Extraction

537   Kit.

538

539   **Figure 8. Relative abundance of taxa in pooled native stool.** Relative proportions of taxa are

540   shown at the order level for pooled native stool samples with (NatwSh) or without (Nat) pre-

541   incubation in nucleotide stabilization reagent for each extraction kit. *Significantly enriched in

542   samples without DNA shield. +Significantly enriched in samples with DNA shield.

543

544   **DECLARATIONS**

545   **Ethics approval and consent to participate**

546   The institutional review board of the University of California, Davis approved this study and all

547   participants provided written informed consent (clinicaltrails.gov registration number

548   NCT02298725).

549

550   **Consent for publication**

551   Not applicapble

552

553   **Availability of data and material**

554   All 16S rRNA sequences used in this analysis were deposited in the Qiita database

555   (https://qiita.ucsd.edu) under study ID 11427 and in the European Nucleotide Archive (ENA)

556   under accession number ERP104979.

557

558    **Competing interests**

559    The authors declare that they have no competing interests.

560

561    **Funding**

565

566    **Authors' contributions**

567    RH performed bioinformatics analyses and together with MEK wrote the main text of the

568    manuscript. ZA performed DNA extraction, PCR, compiled data and contributed to writing the

569    methods section. NK designed and managed the human study which provided stools for these

570    analyses and provided editorial input for the manuscript.

571

572    **Acknowledgements**

579

580    **REFERENCES**

581    1.    Manor, O. and E. Borenstein, *Systematic Characterization and Analysis of the Taxonomic*

582          *Drivers of Functional Shifts in the Human Microbiome.* Cell Host Microbe, 2017. **21**(2):

583          p. 254-267.

584    2.    Lloyd-Price, J., et al., *Strains, functions and dynamics in the expanded Human*

585          *Microbiome Project.* Nature, 2017. **550**(7674): p. 61-66.

586    3.    Howard, M.M., T.H. Bell, and J. Kao-Kniffin, *Soil microbiome transfer method affects*

587          *microbiome composition, including dominant microorganisms, in a novel environment.*

588          FEMS Microbiol Lett, 2017. **364**(11).

589    4.    Oliverio, A.M., M.A. Bradford, and N. Fierer, *Identifying the microbial taxa that*

590          *consistently respond to soil warming across time and space.* Glob Chang Biol, 2017.

591          **23**(5): p. 2117-2129.

592    5.    Ward, C.S., et al., *Annual community patterns are driven by seasonal switching between*

593          *closely related marine bacteria.* ISME J, 2017. **11**(6): p. 1412-1422.

594    6.    Neave, M.J., et al., *Endozoicomonas genomes reveal functional adaptation and plasticity*

595          *in bacterial strains symbiotically associated with diverse marine hosts.* Sci Rep, 2017. **7**:

596          p. 40579.

597    7.    Lemas, D.J., et al., *Alterations in human milk leptin and insulin are associated with early*

598          *changes in the infant intestinal microbiome.* Am J Clin Nutr, 2016. **103**(5): p. 1291-300.

599    8.    Noble, E.E., et al., *Early-Life Sugar Consumption Affects the Rat Microbiome*

600          *Independently of Obesity.* J Nutr, 2017. **147**(1): p. 20-28.

601    9.    Zielińska, S., et al., *The choice of the DNA extraction method may influence the outcome*

602          *of the soil microbial community structure analysis.* MicrobiologyOpen, 2017.

603    10.    Lozupone, C.A., et al., *Meta-analyses of studies of the human microbiota.* Genome Res,

604            2013. **23**(10): p. 1704-14.

605    11.    Aird, D., et al., *Analyzing and minimizing PCR amplification bias in Illumina sequencing*

606            *libraries.* Genome biology, 2011. **12**(2): p. R18.

607    12.    Schirmer, M., et al., *Insight into biases and sequencing errors for amplicon sequencing*

608            *with the Illumina MiSeq platform.* Nucleic acids research, 2015. **43**(6): p. e37-e37.

609    13.    Sinha, R., et al., *Assessment of variation in microbial community amplicon sequencing by*

610            *the Microbiome Quality Control (MBQC) project consortium.* Nat Biotechnol, 2017.

611    14.    Konstantinidis, K.T. and J.M. Tiedje, *Genomic insights that advance the species*

612            *definition for prokaryotes.* Proceedings of the National Academy of Sciences of the

613            United States of America, 2005. **102**(7): p. 2567-2572.

614    15.    Nguyen, N.-P., et al., *A perspective on 16S rRNA operational taxonomic unit clustering*

615            *using sequence similarity.* npj Biofilms and Microbiomes, 2016. **2**: p. 16004.

616    16.    Schloss, P.D. and S.L. Westcott, *Assessing and improving methods used in operational*

617            *taxonomic unit-based approaches for 16S rRNA gene sequence analysis.* Appl Environ

618            Microbiol, 2011. **77**(10): p. 3219-26.

619    17.    Kopylova, E., et al., *Open-Source Sequence Clustering Methods Improve the State Of the*

620            *Art.* mSystems, 2016. **1**(1).

621    18.    Tikhonov, M., R.W. Leach, and N.S. Wingreen, *Interpreting 16S metagenomic data*

622            *without clustering to achieve sub-OTU resolution.* The ISME journal, 2015. **9**(1): p. 68.

623    19.    Edgar, R.C., *UNOISE2: improved error-correction for Illumina 16S and ITS amplicon*

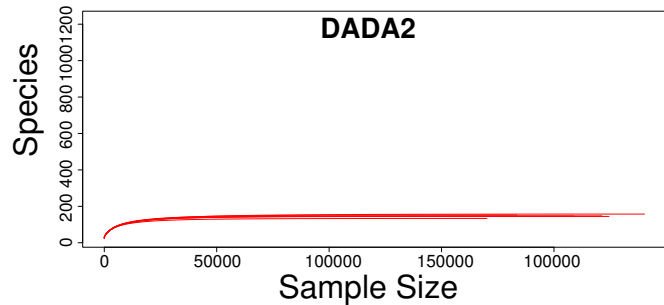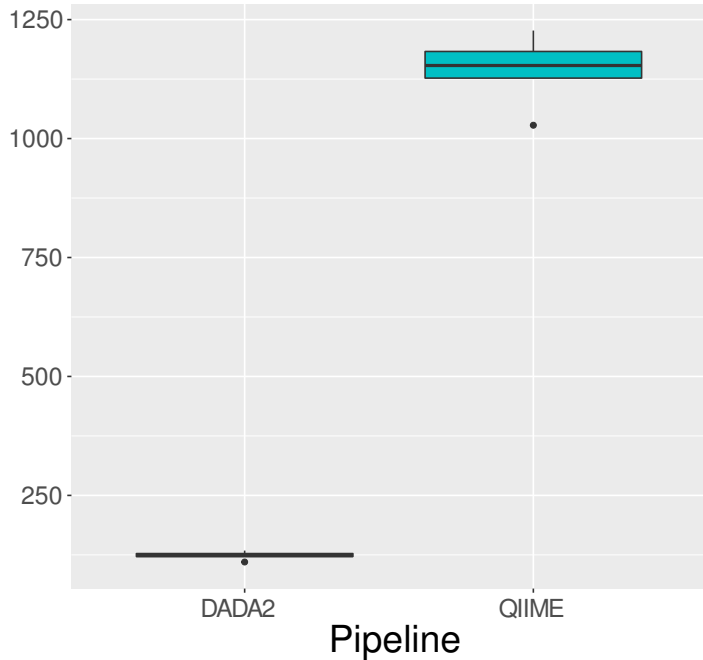624            *sequencing.* bioRxiv, 2016: p. 081257.

625  20.  Callahan, B.J., et al., *DADA2: High-resolution sample inference from Illumina amplicon*

626       *data.* Nat Methods, 2016. **13**(7): p. 581-3.

627  21.  Amir, A., et al., *Deblur Rapidly Resolves Single-Nucleotide Community Sequence*

628       *Patterns.* mSystems, 2017. **2**(2).

629  22.  Eren, A.M., et al., *Minimum entropy decomposition: unsupervised oligotyping for*

630       *sensitive partitioning of high-throughput marker gene sequences.* ISME J, 2015. **9**(4): p.

631       968-79.

632  23.  Rosen, M.J., et al., *Denoising PCR-amplified metagenome data.* BMC bioinformatics,

633       2012. **13**(1): p. 283.

634  24.  Edgar, R.C., *Accuracy of microbial community diversity estimated by closed- and open-*

635       *reference OTUs.* PeerJ, 2017. **5**: p. e3889.

636  25.  Callahan, B.J., P.J. McMurdie, and S.P. Holmes, *Exact sequence variants should replace*

637       *operational taxonomic units in marker gene data analysis.* bioRxiv, 2017: p. 113597.

638  26.  Vishnivetskaya, T.A., et al., *Commercial DNA extraction kits impact observed microbial*

639       *community composition in permafrost samples.* FEMS microbiology ecology, 2014.

640       **87**(1): p. 217-230.

641  27.  Bahl, M.I., A. Bergström, and T.R. Licht, *Freezing fecal samples prior to DNA extraction*

642       *affects the Firmicutes to Bacteroidetes ratio determined by downstream quantitative PCR*

643       *analysis.* FEMS microbiology letters, 2012. **329**(2): p. 193-197.

644  28.  Guo, F. and T. Zhang, *Biases during DNA extraction of activated sludge samples*

645       *revealed by high throughput sequencing.* Applied microbiology and biotechnology, 2013.
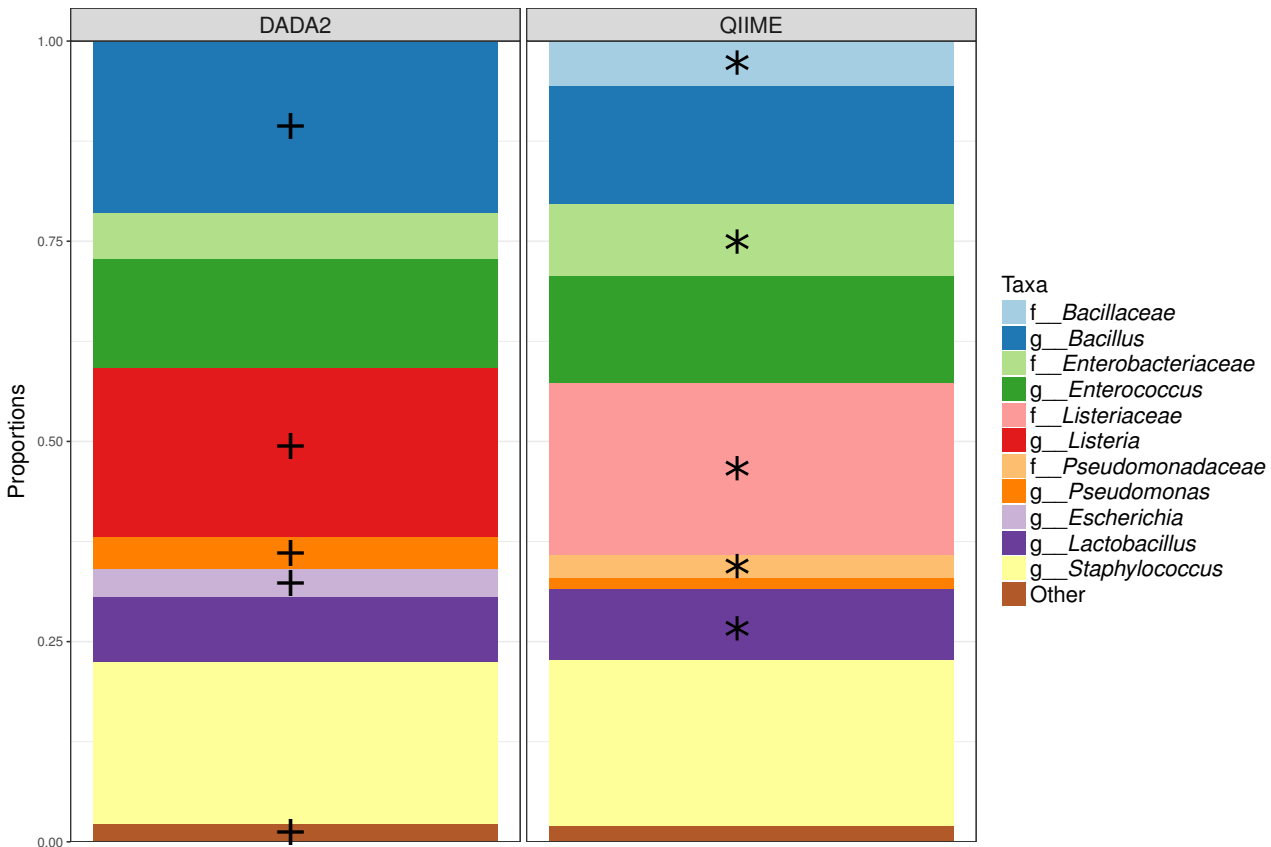
646       **97**(10): p. 4607-4616.

647    29.    Garcia-Mazcorro, J.F., et al., *Influence of whole-wheat consumption on fecal microbial*

648           *community structure of obese diabetic mice.* PeerJ, 2016. **4**: p. e1702.

649    30.    Plante, C.J. and S. Stinson, *Recolonization and cues for bacterial migration into mock¹*

650           *deposit-feeder fecal casts.* Aquatic microbial ecology, 2003. **33**(2): p. 107-115.

651    31.    Kable, M.E., et al., *The Core and Seasonal Microbiota of Raw Bovine Milk in Tanker*

652           *Trucks and the Impact of Transfer to a Milk Processing Facility.* MBio, 2016. **7**(4).

653    32.    Caporaso, J.G., et al., *Global patterns of 16S rRNA diversity at a depth of millions of*

654           *sequences per sample.* Proceedings of the National Academy of Sciences, 2011.

655           **108**(Supplement 1): p. 4516-4522.

656    33.    Hamady, M., et al., *Error-correcting barcoded primers for pyrosequencing hundreds of*

657           *samples in multiplex.* Nature methods, 2008. **5**(3): p. 235-237.

658    34.    Bokulich, N.A., et al., *Next-generation sequencing reveals significant bacterial diversity*

659           *of botrytized wine.* PloS one, 2012. **7**(5): p. e36357.

660    35.    Caporaso, J.G., et al., *QIIME allows analysis of high-throughput community sequencing*

661           *data.* Nature methods, 2010. **7**(5): p. 335-336.

662    36.    Edgar, R.C., *Search and clustering orders of magnitude faster than BLAST.*

663           Bioinformatics, 2010. **26**(19): p. 2460-2461.

664    37.    DeSantis, T.Z., et al., *Greengenes, a chimera-checked 16S rRNA gene database and*

665           *workbench compatible with ARB.* Applied and environmental microbiology, 2006. **72**(7):

666           p. 5069-5072.

667    38.    Segata, N., et al., *Metagenomic biomarker discovery and explanation.* Genome biology,
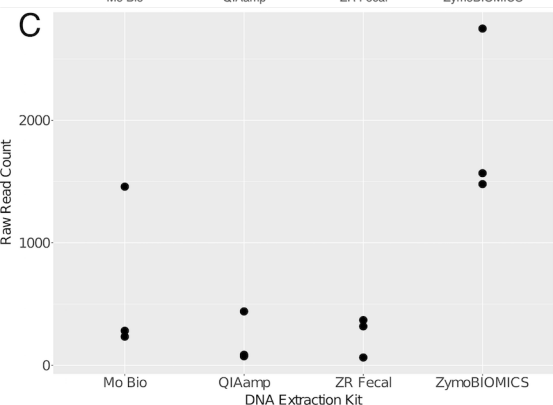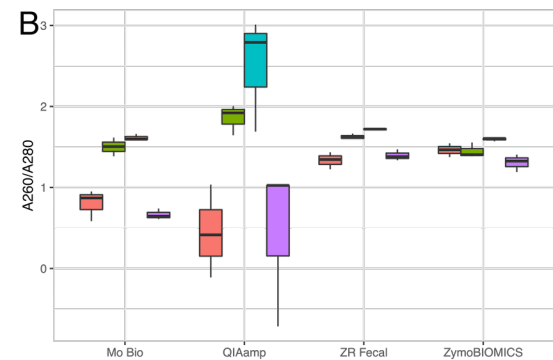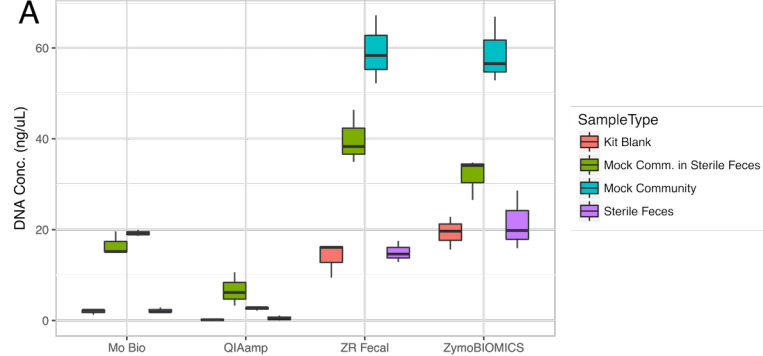
668           2011. **12**(6): p. R60.

669    39.    Li, K., et al., *Analyses of the microbial diversity across the human microbiome.* PloS one,

670            2012. **7**(6): p. e32118.

671    40.    Eren, A.M., et al., *Oligotyping: differentiating between closely related microbial taxa*

672            *using 16S rRNA gene data.* Methods in Ecology and Evolution, 2013. **4**(12): p. 1111-

673            1119.

674    41.    Leite, F.L., et al., *Comparison of fecal DNA extraction kits for the detection of*

675            *Mycobacterium avium subsp. paratuberculosis by polymerase chain reaction.* Journal of

676            veterinary diagnostic investigation, 2013. **25**(1): p. 27-34.

677    42.    Ruggieri, J., et al., *Techniques for Nucleic Acid Purification from Plant, Animal, and*

678            *Microbial Samples.* Sample Preparation Techniques for Soil, Plant, and Animal Samples,

679            2016: p. 41-52.

680    43.    Nechvatal, J.M., et al., *Fecal collection, ambient preservation, and DNA extraction for*

681            *PCR amplification of bacterial and human markers from human feces.* Journal of

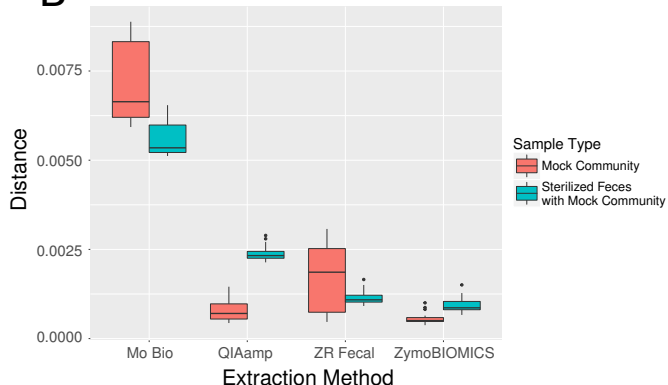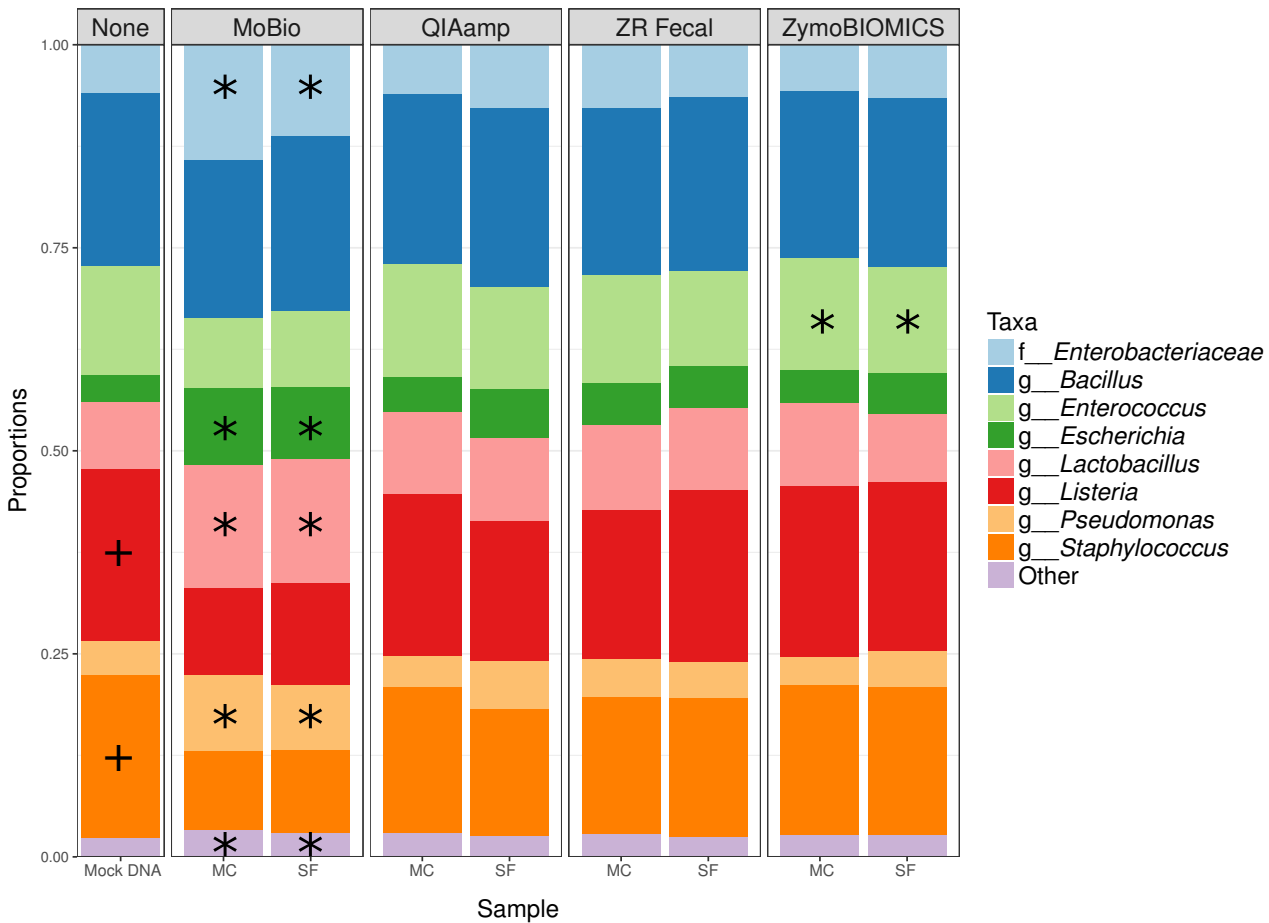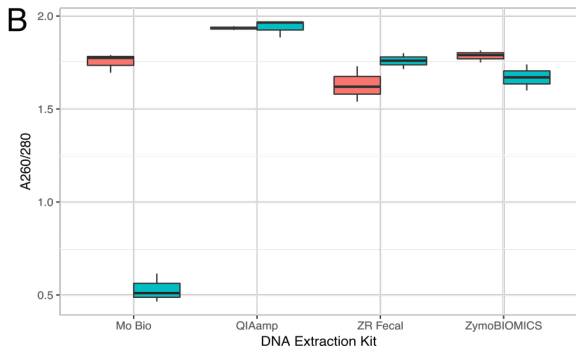682            microbiological methods, 2008. **72**(2): p. 124-132.


683