# Effects of library size variance, sparsity, and compositionality on the analysis of microbiome data

Sophie J Weiss, Zhenjiang Xu, Amnon Amir, Shyamal Peddada, Kyle Bittinger, Antonio Gonzalez, Catherine Lozupone, Jesse R. Zaneveld, Yoshiki Vazquez-Baeza, Amanda Birmingham, Rob Knight

**Background:** Data from 16S amplicon sequencing present challenges to ecological and statistical interpretation. In particular, library sizes often vary over several ranges of magnitude, and the data contains many zeroes. Also, since researchers sample a small fraction of the ecosystem, the observed sequences are relative abundances and therefore the data is compositional. Here we evaluate methods developed in the literature to address these three challenges in the context of normalization and ordination analysis, which is commonly used to visualize overall differences in bacterial composition between sample groups, and differential abundance analysis, which tests for significant differences in the abundances of microbes between sample groups. **Results.** *Effects of normalization on ordination:* Most normalization methods successfully cluster samples according to biological origin when many microbes differ between the groups. For datasets in which clusters are subtle and/or sequence depth varies greatly between samples, or for metrics in which rare microbes play an important role, rarefying outperforms other techniques. For abundance-based metrics, rarefying as well as alternatives like DESeq and metagenomeSeq's cumulative sum scaling (CSS), seem to correctly cluster samples according to biological origin. With these normalization alternatives, clustering by sequence depth as a confounding variable must be checked for, especially for low library sizes. *Effects of differential abundance testing model choice:* We build on previous work to evaluate each statistical method using rarefied as well as unrarefied data. When the mean library sizes in the differential abundance groups differ by more than 2-3x, or the library sizes differ in distribution, our simulation studies reveal that each statistical method improved in its false positive rate when samples were rarefied. However, when the difference in library size mean is less than 2-3x, and the library sizes are similarly distributed, rarefying results in a loss of power for all methods. In this case, DESeq2 has the highest power to compare groups, especially for less than 20 samples per group. MetagenomeSeq's fitZIG is a faster alternative to DESeq2, although it does worse for smaller sample sizes (<50 samples per group) and tends to have a higher false positives rate. For larger sample sizes (>50 samples), rarefying paired with a non-parametric test, such as the Mann-Whitney test, can also yield equally high sensitivity. Based on these

results, we recommend a stepwise procedure in which sample groups are first tested for significant differences in library size. If there is a significant difference, we recommend rarefying with a non-parametric test. Otherwise, DESeq2 and/or fitZIG offer increased sensitivity, especially for rare OTUs and small sample numbers. **Conclusions.** These findings help guide which technique to use depending on the data characteristics of a given study.

2 **Effects of library size variance, sparsity, and compositionality on the analysis of**
3 **microbiome data**
4

5 Sophie J. Weiss[1], Zhenjiang Zech Xu[2], Amnon Amir[2], Shyamal Peddada[3], Kyle Bittinger[4],
6 Antonio Gonzalez[2], Catherine Lozupone[5], Jesse R. Zaneveld[6], Yoshiki Vázquez-Baeza[2],
7 Amanda Birmingham[7], Rob Knight[2,8a]
8
9

10
11 *[1]Department of Chemical and Biological Engineering, University of Colorado at Boulder,*
12 *Boulder, CO 80309*
13 *[2]Departments of Pediatrics, University of California San Diego, La Jolla, CA 92093*
14 *[3]Biostatistics and Computational Biology Branch, NIEHS, NIH*
15 *[4]Department of Microbiology, University of Pennsylvania, Philadelphia, PA 18014*
16 *[5]Department of Medicine, University of Colorado, Denver 80204*
17 *[6] Department of Microbiology, Oregon State University, 226 Nash Hall, Corvallis, OR 97331*
18 *[7]Center for Computational Biology and Bioinformatics, Dept. of Medicine, University of*
19 *California San Diego, La Jolla, CA 92093*
20 *[8]Department of Computer Science & Engineering, University of California San Diego, La Jolla,*
21 *CA 92093*
22
23
24
25 [a]To whom correspondence should be addressed

26 **Corresponding author**
27 Rob Knight, University of California San Diego, 9500 Gilman Drive, MC 0763 La Jolla, CA
28 92093
29 robknight@ucsd.edu   t:858-246-1184        f:858-246-1981
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46

47  ABSTRACT
48  **Background:** Data from 16S amplicon sequencing present challenges to ecological and
49  statistical interpretation. In particular, library sizes often vary over several ranges of magnitude,
50  and the data contains many zeroes. Also, since researchers sample a small fraction of the
51  ecosystem, the observed sequences are relative abundances and therefore the data is
52  compositional. Here we evaluate methods developed in the literature to address these three
53  challenges in the context of normalization and ordination analysis, which is commonly used to
54  visualize overall differences in bacterial composition between sample groups, and differential
55  abundance analysis, which tests for significant differences in the abundances of microbes
56  between sample groups.
57  **Results.** *Effects of normalization on ordination:* Most normalization methods successfully
58  cluster samples according to biological origin when many microbes differ between the groups.
59  For datasets in which clusters are subtle and/or sequence depth varies greatly between samples,
60  or for metrics in which rare microbes play an important role, rarefying outperforms other
61  techniques. For abundance-based metrics, rarefying as well as alternatives like DESeq and
62  metagenomeSeq's cumulative sum scaling (CSS), seem to correctly cluster samples according to
63  biological origin. With these normalization alternatives, clustering by sequence depth as a
64  confounding variable must be checked for, especially for low library sizes. *Effects of differential
65  abundance testing model choice:* We build on previous work to evaluate each statistical method
66  using rarefied as well as unrarefied data. When the mean library sizes in the differential
67  abundance groups differ by more than 2-3x, or the library sizes differ in distribution, our
68  simulation studies reveal that each statistical method improved in its false positive rate when
69  samples were rarefied.  However, when the difference in library size mean is less than 2-3x, and
70  the library sizes are similarly distributed, rarefying results in a loss of power for all methods.  In
71  this case, DESeq2 has the highest power to compare groups, especially for less than 20 samples
72  per group. MetagenomeSeq's fitZIG is a faster alternative to DESeq2, although it does worse for
73  smaller sample sizes (<50 samples per group) and tends to have a higher false positives rate. For
74  larger sample sizes (>50 samples), rarefying paired with a non-parametric test, such as the
75  Mann-Whitney test, can also yield equally high sensitivity. Based on these results, we
76  recommend a stepwise procedure in which sample groups are first tested for significant
77  differences in library size. If there is a significant difference, we recommend rarefying with a
78  non-parametric test.  Otherwise, DESeq2 and/or fitZIG offer increased sensitivity, especially for
79  rare OTUs and small sample numbers.
80  **Conclusions.** These findings help guide which technique to use, depending on the data
81  characteristics of a given study.
82
83  INTRODUCTION
84      Although data produced by high-throughput sequencing has proven extremely useful for
85  understanding microbial communities, the interpretation of these data is complicated by several
86  statistical challenges. To ease data interpretation, data are often normalized to account for the
87  sampling process and differences in sequencing efforts. Ordination analysis, such as principal
88  coordinates analysis (PCoA) (Gower 1966), is subsequently applied to these normalized data to
89  visualize broad trends of how similar or different bacteria are in certain sample types, such as
90  healthy *vs*. sick patients). Samples containing similar bacteria will group, or cluster, close
91  together, while differences in bacterial composition will cause separation in PCoA space. Next,

92  researchers may wish to determine, through statistical testing, which specific bacteria are
93  significantly differentially abundant between two sample type clusters.
94
95       For example, patients with *Clostridium difficile* infection cluster separately from healthy
96  patients in PCoA plots, and these overall differences in community composition are driven by
97  differences in microbial relative abundances (Kelly et al. 2014; Shankar et al. 2014; Weingarden
98  et al. 2015). Restoration of each intestinal bacteria type to healthy levels leads to patient
99  recovery, and causes samples from treated patients to overlap with healthy individuals in PCoA
100 plots. Significant changes in certain bacterial species abundances has also been linked to
101 inflammatory bowel diseases (Gevers et al. 2014), diarrhea (Pop et al. 2014), obesity (Ley et al.
102 2005; Ridaura et al. 2013; Turnbaugh et al. 2009), HIV (Lozupone et al. 2013a), diet (David et
103 al. 2014), culture, age, and antibiotic use (Lozupone et al. 2013b), among many other factors.
104 However, the veracity of these discoveries depends upon how well the chosen normalization and
105 differential abundance testing techniques address the statistical challenges posed by the
106 underlying community sequence data.
107
108      Following initial quality control steps to account for errors in the sequencing process,
109 microbial community sequencing data is typically organized into large matrices where the
110 columns represent samples, and rows contain observed counts of clustered sequences commonly
111 known as Operational Taxonomic Units, or OTUs, that represent bacteria types.  These tables are
112 often referred to as OTU tables. Several features of OTU tables can cause erroneous results in
113 downstream analyses if unaddressed. First, the microbial community in each biological sample
114 may be represented by very different numbers of sequences, reflecting differential efficiency of
115 the sequencing process rather than true biological variation. This problem is exacerbated by the
116 observation that the full range of species is rarely saturated, such that more bacterial species are
117 observed with more sequencing. (Similar trends by sequencing depth hold for discovery of genes
118 in shotgun metagenomic samples (Qin et al. 2010; Rodriguez & Konstantinidis 2014)). Thus,
119 samples with relatively few sequences can have inflated beta (β, or between sample) diversity,
120 because authentically shared OTUs are erroneously scored as unique to samples with more
121 sequences (Lozupone et al. 2011). Second, most OTU tables are sparse, meaning that they
122 contain a high proportion of zero counts (Paulson et al. 2013). This sparsity means that the
123 counts of rare OTUs are uncertain, since they are at the limit of sequencing detection ability
124 when there are many sequences per sample (i.e. large library size), and are undetectable when
125 there are few sequences per sample. Third, each sample is only a small percentage of its original
126 environment, constraining the total number of rRNA sequences to a constant sum; in such
127 "compositional" data, researchers do not know the absolute counts of each type of OTU but only
128 their relative abundances in relation to each other (Aitchison 1982; Friedman & Alm 2012;
129 Lovell D 2010). Uneven sampling depth, sparsity, and compositionality represent serious
130 challenges for interpreting these data.  No normalization method or differential abundance
131 testing method simultaneously addresses all of these challenges.  Thus, investigators must
132 choose methods based on relevant features of the dataset under consideration.
133
134 **Normalization**
135      Normalization is critical to address variability in sampling depths and number of zeros.
136 Microbial ecologists in the era of high-throughput sequencing have commonly normalized their
137 OTU matrices by rarefying, or drawing without replacement from each sample such that all

138 samples have the same number of total counts. Samples with total counts below the defined
139 threshold are excluded, sometimes leading researchers to face difficult trade-offs between
140 sampling depth and the number of samples evaluated.  To ensure the proper total sum is chosen,
141 rarefaction curves can be constructed (Gotelli & Colwell 2001). These curves plot the number of
142 counts sampled (rarefaction depth) vs. the expected value of species diversity. Rarefaction
143 curves provide guidance that allows users to avoid negatively impacting the species diversity
144 found in samples by choosing too low a rarefaction depth.  The origins of rarefying sample
145 counts are mainly in sample species diversity measures, or alpha diversity (Brewer &
146 Williamson 1994; Gotelli & Colwell 2001).  However, more recently rarefying has been used in
147 the context of β-diversity (Horner-Devine et al. 2004; Jernvall & Wright 1998).  Rarefying
148 samples for normalization is now the standard in microbial ecology, and is present in all major
149 data analysis toolkits for this field (Caporaso et al. 2010; Jari Oksanen 2015; McMurdie &
150 Holmes 2013; Schloss et al. 2009). While rarefying is not an ideal normalization method, as it
151 reduces statistical power by removing some data, and was not designed to address
152 compositionality, alternatives to rarefying have not been sufficiently developed until recently.
153
154     Normalization alternatives to rarefying all involve some type of transformation, the most
155 common of which are scaling or log-ratio transformations. Effects of scaling methods depend on
156 the scaling factor chosen; often, a particular quantile of the data is used for normalization, but
157 choosing the correct quantile is difficult (Anders & Huber 2010; Bullard et al. 2010; Dillies et al.
158 2013; Paulson et al. 2013; Robinson & Oshlack 2010), and scaling can overestimate or
159 underestimate the prevalence of zero fractions, depending on whether zeroes are left in or thrown
160 out of the scaling (Agresti & Hitchcock ; Friedman & Alm 2012). This is because putting all
161 samples of varying sampling depth on the same scale ignores the differences in sequencing
162 depth, and therefore resolution of species, between the samples.  For example, a rare species
163 having zero counts in a small rRNA sample can have a small fractional abundance in a large
164 rRNA sample (unless further mathematical modeling beyond simple proportions is applied to
165 correct for this). Scaling can also distort OTU correlations across samples, again due to zeroes,
166 differences in sequencing depth, and sum constraints (Aitchison 1982; Buccianti et al. 2006;
167 Friedman & Alm 2012; Lovell D 2010; Pearson 1896).
168
169     While rarefying and some scaling techniques, such as total sum scaling (proportions),
170 treat OTU sequence counts as absolute environmental abundances, the counts are compositional
171 and only a fraction from the original environment, making only their relative ratios known
172 (Friedman & Alm 2012; Lovell D 2010). In contrast, log ratio transformations correct for
173 compositionality by exploiting this relative ratio information, and can also alleviate some noise
174 in the data (Aitchison 1982; Buccianti et al. 2006; Friedman & Alm 2012; Lovell D 2010).
175  However, because the log transformation cannot be applied to zeros (which are often well over
176 half of microbial data counts (Paulson et al. 2013)), sparsity is extremely problematic for
177 methods that rely on this transformation. One approach to this issue is to replace zeros with a
178 small value, known as a pseudocount. Despite active research on selection of pseudocount values
179 for scaling methods (Egozcue et al. 2003; Greenacre 2011), the choice of pseudocount values
180 can dramatically change the results (Costea et al. 2014; Paulson et al. 2014).
181
182 **Differential Abundance Testing**

183     For OTU differential abundance testing between conditions (e.g. case vs. control), a
184 common approach is to first rarify the count matrix to a fixed depth and then apply a non-
185 parametric test  (e.g. Mann-Whitney test for tests of two classes; Kruskal-Wallis test for tests of
186 multiple groups). Non-parametric tests are often preferred because most OTU counts are not
187 normally distributed (Wagner et al. 2011). However, this approach does not account for the fact
188 that the OTU counts are compositional. Also, nonparametric tests such as the Kruskal-Wallis test
189 do not fare well in terms of power when the data are sparse, but perform well when the data are
190 not sparse (Paulson et al. 2013). Recently, promising parametric models that make stronger
191 assumptions about the data have been developed in the subfields of transcriptomics ('RNA-Seq')
192 and metagenomic sequencing. These may additionally be useful for microbial marker gene data
193 (Anders & Huber 2010; Anders et al. 2013; Law et al. 2014; Love MI 2014; McMurdie &
194 Holmes 2014; Paulson et al. 2013; Robinson et al. 2010; Robinson & Smyth 2008). Such models
195 have greater detection power if their assumptions about the data are correct; however, studies of
196 these models on RNA-Seq data have shown that they can yield poor results (Rapaport et al.
197 2013) if relevant assumptions are not valid.
198
199     These parametric models are composed of a generalized linear model (GLM) that
200 assumes a distribution (Cameron & Trivedi), and there is debate about which distribution to use
201 (Auer & Doerge 2010; Cheung 2002; Connolly et al. 2009; Holmes et al. 2012; McMurdie &
202 Holmes 2014; Paulson et al. 2013; Rapaport et al. 2013; Soneson & Delorenzi 2013; White et al.
203 2009; Yu et al. 2013). In the genomics field, the negative binomial (NB) GLM has replaced the
204 Poisson GLM to allow for estimating overdispersion (Anders & Huber 2010; Anders et al. 2013;
205 Robinson et al. 2010). This model type was also one of the first in the RNA-Seq field, and
206 developed for use with a low number of replicates. NB models accommodate low replication by
207 assuming that OTUs of similar mean expression strength have similar variance in their sample
208 count distributions, estimating model parameters using this assumption, and then leveraging the
209 GLM to perform exact statistical tests. Also, while allowing for some overdispersion, the NB
210 often yields a poor fit in the case of a large number of zeroes, which is very typical in
211 microbiome data (Cheung 2002; Paulson et al. 2013). Zero-inflated GLMs, the most promising
212 of which is the zero-inflated Gaussian (ZIG), attempts to overcome this limitation (Paulson et al.
213 2013). The ZIG tries to address compositionality, sparsity and unequal sampling depth by
214 separately modeling 'structural' zero counts generated by e.g. under-sequencing and zeros
215 generated by the biological distribution of taxa. Log transformation of the non-zero counts yields
216 the Gaussian. However, this mixture model distribution is designed for continuous data rather
217 than discrete microbiome data. Hence, it is expected to do best in study designs that have large
218 sample sizes and high sequencing depths, and thus best approximate continuous distributions.
219
220     Here, we evaluate some of the most widely used or promising techniques for analyzing
221 sequencing data in the context of microbial ecology, with a focus on normalization and OTU
222 differential abundance testing. In addition to these widely used or promising methods, we also
223 test the naïve approaches of no normalization, and proportions (i.e. total sum scaling) for
224 comparison purposes.  Such comparisons are important, because while potential issues with
225 many methodologies are known, the balance of sensitivity and specificity for these methods in
226 situations commonly facing microbial ecologists is currently largely unknown.  Recent work in
227 this area (McMurdie & Holmes 2014), provides insights into the performance of parametric
228 normalization and differential abundance testing approaches for microbial ecology studies.

229  However, the work is primarily focused on estimating proportions from discrete data. Here we
230  update and expand these recent findings using both real and simulated datasets exemplifying the
231  additional combined challenges of uneven library sizes, sparsity, and compositionality.
232
233  MATERIALS AND METHODS
234  **Normalization**
235       The basic test of how well broad differences in microbial sample composition are
236  detected, as assessed by clustering analysis, was conducted as in 'Simulation A' from McMurdie
237  and Holmes (McMurdie & Holmes 2014).  Briefly, the 'Ocean' and 'Feces' microbiomes (the
238  microbial data from ocean and human feces samples, respectively) from the 'Global Patterns'
239  dataset (Caporaso et al. 2011b) were used as templates, modeled with a multinomial, and taken
240  to represent distinct classes of microbial community because they have few OTUs in common.
241   These two classes were mixed in many defined proportions (the 'effect size') in independent
242  simulations in order to generate simulated samples of varying clustering difficulty. Samples were
243  generated in sets of 40, as in McMurdie and Holmes (McMurdie & Holmes 2014). We also
244  tested smaller and larger sample sizes but saw little difference in downstream results. Additional
245  sets of 40 samples were simulated for varying library sizes (1000, 2000, 5000, and 10000
246  sequences per sample). These simulated samples were then used to assess normalization methods
247  by the proportion of samples correctly classified into the two clusters by the partitioning around
248  medioids (PAM) algorithm (Kaufman L. 1990; Reynolds A 2006).
249
250       McMurdie and Holmes (McMurdie & Holmes 2014) evaluated clustering accuracy with
251  five normalization methods (none, proportion, rarefying with replacement as in the multinomial
252  model (Colwell et al. 2012), DESeqVS (Anders & Huber 2010), and UQ-logFC (in the edgeR
253  package) (Robinson et al. 2010)) and six beta diversity metrics (Euclidean, Bray-Curtis (Bray &
254  Curtis 1957), PoissonDist (Witten 2011), top-MSD (Robinson et al. 2010), unweighed UniFrac
255  (Lozupone & Knight 2005), and weighted UniFrac (Lozupone et al. 2007)).  We modified the
256  normalization methods to those in Table S1 (none, proportion, rarefying without replacement as
257  in the hypergeometric model (Colwell et al. 2012), CSS (Paulson et al. 2013), logUQ (Bullard et
258  al. 2010), DESeqVS (Anders & Huber 2010), and edgeR-TMM (Robinson & Oshlack 2010))
259  and the beta diversity metrics to those in Fig2 and Fig. S1 (binary Jaccard, Bray-Curtis (Bray &
260  Curtis 1957), Euclidean, unweighed UniFrac (Lozupone & Knight 2005), and weighted UniFrac
261  (Lozupone et al. 2007)), thus including more recent normalization methods (Bullard et al. 2010;
262  Paulson et al. 2013), and only those beta diversity metrics that are most common in the literature.
263  We amended the rarefying method to the hypergeometric model (Colwell et al. 2012), which is
264  much more common in microbiome studies (Caporaso et al. 2010; Schloss et al. 2009).
265  Negatives in the DESeq normalized values (Anders & Huber 2010) were set to zero as in
266  McMurdie and Holmes (McMurdie & Holmes 2014), and a pseudocount of one was added to the
267  count tables (McMurdie & Holmes 2014). McMurdie and Holmes (McMurdie & Holmes 2014)
268  penalized the rarefying technique for dropping the lowest fifteenth percentile of sample library
269  sizes in their simulations by counting the dropped samples as 'incorrectly clustered'. Because the
270  15th percentile was used to set rarefaction depth, this capped clustering accuracy at 85%. We
271  instead quantified cluster accuracy among samples that were clustered following normalization
272  to exclude this rarefying penalty (Fig. S1). Conversely, it has since been confirmed that low-
273  depth samples contain a higher proportion of contaminants (rRNA not from the intended sample)
274  (Kennedy et al. 2014; Salter et al. 2014). Because the higher depth samples that rarefying keeps

275 may be higher quality and therefore give rarefying an unfair advantage, Fig. 2 compares
276 clustering accuracy for all the techniques based on the same set of samples remaining in the
277 rarefied dataset.
278
279          On the real datasets, non-parametric multivariate ANOVA (PERMANOVA) (Anderson
280 2001) was calculated by fitting a Type I sequential sums of squares model (y ~ Library_Size +
281 Biological_Effect).  Thus, we control for library size differences before assessing the effects on
282 the studied biological effect. All data was retrieved from QIITA (https://qiita.microbio.me).
283
284
285 **Differential Abundance Testing**
286          The simulation test for how well truly differentially abundant OTUs are recognized by
287 various parametric and non-parametric tests was conducted as in 'Simulation B' in McMurdie
288 and Holmes (McMurdie & Holmes 2014), with a few changes.  The basic data generation model
289 remained the same, but the creation of 'true positive' OTUs was either made symmetrical
290 through duplication or moved to a different step, to avoid introducing compositionality artifacts
291 (see below) depending on the simulation. The 'Global Patterns' (Caporaso et al. 2011b) dataset
292 was again used, because it was one of the first studies to apply high-thoughput sequencing to a
293 broad range of environments, which includes 9 environment types from 'Ocean', to 'Soil'; all
294 simulations were evaluated for all environments. Additionally, we verified the results on the
295 'Lean' and 'Obese' microbiomes from a different study (Piombino et al. 2014). As in McMurdie
296 and Holmes, significant changes were controlled for multiple comparisons using the Benjamini
297 & Hochberg (Benjamini & Hochberg 1995) False Discovery Rate (FDR) threshold of 0.05.
298
299          A simple overview of the two methods used for simulating differential abundance is
300 presented in Fig. S5a.  In McMurdie and Holmes' (McMurdie & Holmes 2014) 'Original'
301 simulation (second row), the distribution of counts from one environment (e.g. 'Ocean') was
302 modeled off of a multinomial template (first row) for two similar groups ('Ocean_1' and
303 'Ocean_2'), ensuring a baseline of all 'true negative' OTUs. Following the artificial inflation of
304 specific OTUs in the 'Ocean_1' samples to create 'true positives', fold-change estimates for
305 every other OTU are affected. Thus, 'true negatives' are possible 'true positives.' This is because
306 the counts in an OTU table are compositional, or relative abundances constrained to a sum. To
307 control for this we inflate OTUs by pairs of differentially abundant OTUs in both the 'Ocean_1'
308 and 'Ocean_2' samples (third row), creating a new 'Balanced' simulation.
309
310          We also tested the effect of differentially abundant organisms dominating one type of
311 community by drawing from a multinomial distribution where solely that organism's template
312 value is increased. This 'Compositional' approach is explained in Fig. S5b, and the results are
313 shown in Fig. S7.  In Fig S7, the environmental abundances of 25% of the OTUs in one group
314 are increased.
315
316          Besides the above procedural changes to the McMurdie and Holmes (McMurdie &
317 Holmes 2014) simulation, we also modified the rarefying technique from sampling with
318 replacement (multinomial) to sampling without replacement (hypergeometric - as in the previous
319 Normalization simulations) (Colwell et al. 2012). The testing technique was modified from a
320 two-sided Welch t-test to non-parametric Mann-Whitney test, which is widely used and more

321 appropriate because the OTU distributions in microbiome data usually deviate from normality.
322 The techniques used (Table S2) differ only by the addition of another RNA-Seq method, Voom
323 (Law et al. 2014). Finally, we corrected the FPR definition (McMurdie & Holmes 2014) from
324 FP/(TP + FP) to FP/(TN + FP), where FP = number of false positive OTUs, TP = number of true
325 positive OTUs, and TN = number of true negative OTUs. This new simulation code can be
326 found in the supplemental R files (Differential_abundance.R, and
327 Differential_abundance_with_compositionality.R).
328
329 **Power Curve Calculations**
330     Similar to Table S1 in McMurdie and Holmes [27], we considered a very simplistic set-
331 up to evaluate the effect of rarefying on power when comparing two groups, labeled A and B. As
332 in McMurdie and Holmes [27], we considered the extreme case of a microbial population
333 consisting of only 2 species (or 2 OTUs), with OTU1 + OTU2 = library size.  For power
334 calculations, we assumed that the amount of OTU1 in group B is 85% of the amount of OTU1 in
335 group A.  Thus, it is enough to quantify the proportion of OTU1 in group A and library sizes of
336 groups A and B to specify the whole system.
337
338     We considered varied patterns of proportions of OTU1 in group A ranging from very rare
339 to common (0.5% to 50%).  The library size of group A was fixed at either 500, 1000 or 10,000
340 sequences per sample.  Meanwhile, the library size of group B was always taken to be at least as
341 large as that of group A and was either 10,000 or 100,000 sequences per sample. Various
342 rarefied percentages of the group B library size were considered.  The percent-rarefied
343 calculation for the first set of power curves is exemplified below using a library size of 500 for
344 library A, and an unrarefied library size of 10,000 for B:
345
346 Library size for A     Library size for B
347 ---------------------     ----------------------
348
349 500                          10,0000 (unrarefied case)
350 500                           5,000   (50% rarefied)
351 500                           1,000   (90% rarefied)
352 500                            500   (95% rarefied)
353
354     For each scenario of proportion of OTU1 and library sizes, power was computed using
355 Fisher's exact test.  Power calculations were done using the statistical software SAS.  Power
356 calculation results are provided in Fig. 5.
357
358 **Software Package Versions**
359     R version 3.1.0 (Team 2014) was used with Bioconductor (Gentleman et al. 2004)
360 packages phyloseq version 1.10.0, DESeq version 1.16.0, DESeq2 version 1.4.5, edgeR version
361 3.6.8, metagenomeSeq version 1.7.31, and Limma version 3.20.9.  Also, we used python-based
362 QIIME version 1.9.0, with Emperor (Vazquez-Baeza et al. 2013).
363
364
365 RESULTS AND DISCUSSION
366 **Normalization**

367    When there is a strong biological signal, and normalization is done properly, PCoA can
368    yield clear clustering and insight into microbial community differences (Fig. 1a). However, low-
369    depth samples can lead to poor cluster resolution (Fig. 1b), both by reducing information on
370    community structure, and by being more readily influenced by contamination (Kennedy et al.
371    2014; Salter et al. 2014). Furthermore, if no data normalization is applied, or the normalization
372    method fails to properly correct for differences in sequencing efficacy, the original library size of
373    the samples can confound biological differences (Fig. 1c). This is because samples of lower
374    sequencing depth fail to detect rare taxa. Highly sequenced samples will thus appear more
375    similar to each other than to shallow sequenced samples because they are scored as sharing the
376    same rare taxa.
377
378    To assess all the normalization methods (Table S1), we conducted simulations in the
379    context of results that are highly critical of the rarefying technique (McMurdie & Holmes 2014).
380    Briefly, only necessary modifications (Methods) were made to the code of McMurdie and
381    Holmes (McMurdie & Holmes 2014), making our approach easily comparable. If rarefying is
382    not penalized for the fifteenth percentile lowest depth samples that are thrown out, it can do
383    better than other techniques (Fig. S1). This practice of removing low depth samples from the
384    analysis is supported by the recent discovery that small biomass samples are of poorer quality
385    and may contain contaminating sequences (Kennedy et al. 2014; Salter et al. 2014). Furthermore,
386    alternatives to rarefying also recommend discarding low-depth samples, especially if they cluster
387    separately from the rest of the data (Love MI 2014; Paulson et al. 2013). If all other techniques
388    are run only on the same samples as rarefying, rarefying still does well (Fig. 2). These results
389    demonstrate that previous microbiome ordinations using rarefying as a normalization method
390    likely drew correct conclusions, even if some low depth samples were removed. However, these
391    results also suggest that CSS (Paulson et al. 2013) and DESeq's variance-stabilizing
392    transformation (Anders & Huber 2010) are promising alternatives for normalization prior to
393    PCoA analysis, especially for weighted distance metrics. For unweighted metrics that are based
394    on species presence and absence, like binary Jaccard and unweighted UniFrac, DESeq's
395    variance-stabilizing transformation performs poorly. This is because the negatives resulting from
396    DESeq's log-like transformation are set to zero (as in McMurdie and Holmes (McMurdie &
397    Holmes 2014)), which ignores rare species.
398
399    No good solution exists for the negatives output by the DESeq technique. DESeq was
400    developed mainly for use with Euclidean metrics (Lozupone & Knight 2005; Lozupone et al.
401    2007), for which negatives are not a problem; however, this issue yields misleading results for
402    ecologically useful non-Euclidean measures, like Bray-Curtis (Bray & Curtis 1957)
403    dissimilarity. Also, the negatives pose a problem to UniFrac's (Lozupone & Knight 2005;
404    Lozupone et al. 2007) branch length. The alternative to setting the negatives to zero, or adding
405    the absolute value of the lowest negative value back to the normalized matrix, will not work with
406    distance metrics that are not Euclidean because it amounts to multiplying the original matrix by a
407    constant due to DESeq's log-like transformation. Also, the addition of a constant (or
408    pseudocount; here, one) to the count matrix prior to CSS (Paulson et al. 2013), DESeq (Anders
409    & Huber 2010), and logUQ (Bullard et al. 2010) transformation as a way to avoid log(0) is not
410    ideal, and clustering results have been shown to be very sensitive to the choice of pseudocount,
411    due to the nonlinear nature of the log transform (Costea et al. 2014; Paulson et al. 2014). This
412    underscores the need for a better solution to the zero problem so that log-like approaches

413 inspired by Aitchison can be used (Aitchison 1982), and is especially critical since microbial
414 matrices are almost always much more than half sparse (Paulson et al. 2013).

415

416       While simulations are a useful initial check, real datasets are often much more complex.
417 Therefore, all normalization methods were also examined on real data to check for result and
418 methodological consistency. To perform an initial, detailed comparison of normalization
419 methods, we selected the data set from Caporaso *et al.* (Caporaso et al. 2012). The data included
420 a wide variety of samples, representing both environmental and host-associated sources. To
421 provide an extreme example of differences in sequencing depth, we artificially decreased the
422 library size by 90% for half the samples in the data set. The samples selected for library size
423 reduction were chosen randomly, and the same artificially altered data was used in all
424 normalization comparisons.

425

426       Using the data set from Caporaso *et al.* (Caporaso et al. 2012), we observed substantial
427 biases/confounding of results due to sequencing depth. In ordination of unweighted UniFrac
428 distance by PCoA, the soil samples were split into two groups along the first principal coordinate
429 when no normalization was used (Fig. 3a). Soil samples appearing in the group to the left had
430 more reads than those appearing in the group to the right. Similarly, the two stool samples in the
431 data set were arranged close to soil samples with similar library size. When the data was
432 rarefied prior to ordination, soil and stool samples were arranged along the first two coordinates
433 according to sample type rather than library size (Fig. 3b). Other methods of normalization
434 preserved the characteristic pattern seen in the non-normalized data, where soil and stool
435 samples were separated into groups according to library size (Fig. 3c-f).

436

437       Normalization did not affect conclusions drawn from non-parametric multivariate
438 ANOVA (PERMANOVA) (Anderson 2001), but we did observe differences in the effect size
439 estimated for sample type, and library size ($R^2$). Without normalization, the estimated effect size
440 of sample type for unweighted UniFrac distance was $R^2=0.40$. When the data was rarefied prior
441 to computing distances, the estimated effect size increased to $R^2=0.56$. Other methods of
442 normalization produced effect sizes similar to the non-normalized result. Although the true
443 effect size is not known for this data set, the environment of origin is known to be a dominant
444 effect in the determination of bacterial species observed (Lozupone & Knight 2007). Without
445 normalization, there is a large effect ($R^2=0.14$) corresponding to original library size, which is a
446 known artifact of the sequencing process. Rarefying helps to remove the effect of sequencing
447 depth ($R^2=0.045$), whereas other normalization techniques do not remove this signal artifact,
448 again resembling the non-normalized data.

449

450       As another example, we selected the inflammatory bowel disease (IBD) data set from
451 Gevers *et al.* (Gevers et al. 2014). In contrast to the previous data set, all samples here were
452 taken from a single environment type, namely human stool, and were extremely low depth,
453 having an average of 375 sequences per sample. In an ordination of unweighted UniFrac
454 distance with no normalization, there is again strong clustering by library size, with a group of
455 samples with low sequencing depth appearing slightly separate from the other samples (Fig.
456 S2a). Samples in the low-depth group are either dominated by a lack of species detected due to
457 few sequences, thus artificially inflating the β–diversity, or constitute different bacterial species
458 than the main group of stool samples, which should raise suspicion of potential problems from

459   contamination or poor quality PCR products.  Furthermore, the first principal coordinate in Fig
460   S2a is more strongly correlated with library size ($R^2$=0.055, Fig S2b) and poorly correlated with
461   disease state ($R^2$=0.022), with sampling depth explaining twice the variance of the studied
462   biological effect. Subsampling the data to uniform library size increased the correlation with
463   disease state ($R^2$=0.036), while other methods did not ($R^2$=0.022 for proportion, DESeq, and
464   CSS).  Because the average library size is so low for this study, the library size also affects
465   weighted UniFrac, where there is again low effect size for this gastrointestinal disorder.  Thus,
466   extremely low depth samples still need to be discarded from rarefying alternatives, especially if
467   they are suspected of yielding a poor representation of the true bacterial community due to
468   experimental factors.
469

470        PCoA plots using ecologically common metrics for all of the normalization techniques
471   on a few key real datasets representing a gradient (Lauber et al. 2009), distinct body sites
472   (Costello et al. 2009), and time series (Caporaso et al. 2011a) are shown in Supplemental Figures
473   S3-S4. Most measures do well in these cases where there is strong separation between the
474   categories. Clustering according to sequence depth is less of a problem in these datasets since
475   they have strong clustering patterns, however, some clustering according to depth persists.  For
476   example, in the 'Moving Pictures of the Human Microbiome' dataset (Caporaso et al. 2011a),
477   there is some clustering by sequence depth within each of the four main clusters when
478   normalization alternatives to rarefying are applied. It is noteworthy that CSS normalization
479   results appear robust to the distance metric used, including even Euclidean distance (results not
480   shown), which have been reported to perform poorly on highly sparse matrices (Legendre &
481   Gallagher 2001).
482

483        Thus, both simulations and real data suggest that rarefying remains a strong technique for
484   sample normalization prior to ordination and clustering, especially for presence/absence distance
485   metrics that have historically been very useful (such as binary Jaccard and unweighted UniFrac
486   (Lozupone & Knight 2005) distances), subtle effects, small library sizes, and large differences in
487   library size.  Of the other methods, and for weighted distance measures, we recommend
488   metagenomeSeq's CSS (Paulson et al. 2013) or DESeq's variance stabilizing transformation
489   (Anders & Huber 2010); however, the researcher must check for erroneous clustering according
490   to sequence depth.
491

**Differential Abundance Testing**

492
493        Differential abundance analysis is useful for testing whether certain microbes have higher
494   relative abundance in one condition vs. another (e.g. healthy vs. diseased patients). More
495   complex statistical methods specifically for RNA-Seq data have been developed and include
496   DESeq (Anders & Huber 2010), DESeq2 (Love MI 2014), edgeR (Robinson et al. 2010;
497   Robinson & Smyth 2008), and Voom (Law et al. 2014) (Table S2).  MetagenomeSeq (Paulson et
498   al. 2013) however, was developed specifically for microbial datasets, which usually contain
499   many more zeros than RNA-Seq data. These five methods incorporate more sensitive statistical
500   tests than the standard non-parametric tests such as the Wilcoxon rank-sum test, and they make
501   some distributional assumptions. Therefore, they hold great potential for better prediction of rare
502   OTU behavior.
503

504         Previous work in this area concluded that the newer differential abundance testing
505 models are worthwhile, and that the traditional practice of rarefying causes a high rate of false
506 positives (McMurdie & Holmes 2014). However, the latter conclusion was due to an artifact
507 within the simulation (see Methods, Fig. S5a-b). Instead, we found that rarefying does not cause
508 a high rate of false positives, but may lead to false negatives due to the decreased power that
509 results from throwing away some of the data (Fig. 4). The severity of the power decrease caused
510 by rarifying depends upon how much data has been thrown away. (This problem has been
511 known for a long time, leading to the general guideline to rarefy to the highest depth possible
512 without losing too many samples (Carcer et al. 2011).) In order to determine where the greatest
513 loss in power or information occurs when a dataset is rarefied, we constructed power curves
514 from a simple two-species simulation (Fig. 5). The greatest loss in power occurs for rare to
515 common OTUs (e.g. relative abundance ranging from 0.5% to 50%) depending on the library
516 size. This has also been observed in gene expression studies (Robles et al. 2012). Also,
517 consistent with other studies on subsampling (Carcer et al. 2011; Robles et al. 2012),
518 subsampling to library sizes close to the original does not have much effect on the results (50%
519 is treated as "close to the original" in this simplified example, but real microbiome studies are
520 much more complex and thus the real threshold is likely lower, and data-dependent). We also
521 observed that the performance of rarefying degrades faster for smaller library sizes.
522
523         Since simulations do not necessarily mirror reality, we again investigated the
524 performance of the techniques on real data. This was done for the techniques shown to be most
525 promising in the simulations: DESeq2 (Love MI 2014), edgeR (Robinson et al. 2010; Robinson
526 & Smyth 2008), metagenomeSeq (Paulson et al. 2013), and rarefying. Ranges of dataset sizes
527 were analyzed for environments that likely contain differentially abundant OTUs, as evidenced
528 by PCoA plots and significance tests (Fig. 6). Approximately 6 samples in each of the
529 categories of human skin vs. soil from Caporaso *et al.* (Caporaso et al. 2012), 28 samples in each
530 of the lean vs. obese categories from Piombino *et al.* (Piombino et al. 2014), and 500 samples in
531 the tongue vs. left palm categories from Caporaso *et al.* (Caporaso et al. 2011a) were tested.
532 Although we do not necessarily know which OTUs are true positives in these actual data, it is of
533 interest to investigate how the most promising techniques compare to each other. While
534 rarefying (at the 15th percentile as in McMurdie and Holmes (McMurdie & Holmes 2014)) finds
535 fewer OTUs as significant, the OTUs it does find to be differentially expressed are remarkably
536 stable. Agreeing with our modified simulation, it does not appear that rarefying causes a high
537 type I error. For example, in Fig. 6 there is high agreement between rarefying and the other
538 techniques. However, edgeR, which is known to be too lenient in its dispersion estimates (Love
539 MI 2014; Paulson et al. 2013), predicts a large number of significantly differentially abundant
540 OTUs relative to other methods, especially for studies with fewer samples (Fig. 6a), suggesting a
541 high false positive rate in agreement with RNA-Seq studies (Love MI 2014; Rapaport et al.
542 2013; Soneson & Delorenzi 2013).
543
544         We also used simulated data to investigate the situation in which the average library size
545 between the two groups was not approximately equal (Fig. 7). We found that of the newer
546 methods, metagenomeSeq's figZIG (Paulson et al. 2013) has a high sensitivity and a low false
547 positive rate (1-specificity) compared to the other techniques. However, the false positive rate is
548 still high. Rarefying achieves the lowest false positive rate, but at a cost to sensitivity. Thus, the
549 method employed by investigators may depend on the sensitivity of the analysis in question to

550 false negatives *vs.* false positives. We often place higher importance in reducing false positives,
551 but this will vary depending on experimental design. For example, study designs in which
552 community analysis is used as a pre-screening, and significant changes will be confirmed in
553 high-throughput follow-up experiments may allow greater tolerance of false positives. However,
554 while both fitZIG or rarefying followed by Wilcoxon rank sum tests in isolation may be
555 applicable for detecting differential abundance in particular situations, our results caution that
556 fitZIG should not be used on rarified data, as this combination of methods caused extremely high
557 false positive rates.
558
559     While the no-normalization or proportion approaches perform adequately where the
560 average library size is approximately the same between the two groups (Fig. 4), they do not
561 when one library is 10x larger than the other (Fig. 7). Therefore, we reiterate that neither the no-
562 normalization nor the naive proportion approach should be used for most statistical analyses. To
563 demonstrate this, we suggest the theoretical example of a data matrix with half the samples
564 derived from diseased patients and half from healthy patients.  If the samples from the healthy
565 patients have a 10x larger library size, OTUs of all mean abundance levels will be found to be
566 differentially abundant simply because they may have 10x the number of counts in the healthy
567 patient samples. (Such systematic bias can happen if, for example, healthy vs. diseased patients
568 are sequenced on separate sequencing runs or are being compared in a meta-analysis).  The same
569 warning applies for naive proportions, especially for rare OTUs that could be deemed
570 differentially abundant simply due to differences in sequencing depth.  This is seen even with
571 some filtering to remove very rare OTUs (Fig. 7).  We first observed a transition from the results
572 of Fig. 4 to Fig. 7 at around 2-3x difference in library sizes (Fig S6). Further, we investigated
573 uneven numbers of samples per class, with not much difference in results from Fig. 4.
574
575     While our previous simulations did not have compositionality, we next evaluated the
576 performance of the techniques with a compositional OTU table (see Methods, Fig. S5b). In
577 simulations where the abundances of 25% of the OTUs increased in one group, no method does
578 well in terms of false positive rate (Fig. S7). Proportion normalization again performs poorly in
579 the face of compositionality, which is present in all realistic datasets. For DESeq/DESeq2, poor
580 performance may be due to the model's assumption that differentially abundant OTUs are not a
581 large portion of the population (Dillies et al. 2013), or the model's overdispersion estimates
582 (Paulson et al. 2013).  Thus, compositionality is still a large unsolved problem in differential
583 abundance testing (Lovell et al. 2015), and we would urge caution in data sets where
584 compositionality may play a large role, e.g. when the alpha diversity of the samples is low
585 (Friedman & Alm 2012).
586
587 CONCLUSIONS
588     We built on the pioneering work of McMurdie and Holmes (McMurdie & Holmes 2014),
589 confirming that recently developed more complex techniques for normalization and differential
590 abundance testing hold potential. More testing of the approaches on experimental data is
591 necessary. Of methods for normalizing microbial data for ordination analysis, we found that
592 DESeq normalization (Anders & Huber 2010; Love MI 2014), which was developed for RNA-
593 Seq data and makes use of a log-like transformation, does not work well with ecologically useful
594 metrics, except weighted UniFrac (Lozupone et al. 2007).  In contrast, MetagenomeSeq's CSS
595 normalization (Paulson et al. 2013) was developed for microbial data and does not result in

596 troublesome negative output values. However, with techniques other than rarefying, library size
597 can be a confounding factor with very low library sizes (under approximately 1000 sequences
598 per sample), or if presence/absence metrics like unweighted UniFrac are used (Lozupone &
599 Knight 2005). Extremely low-depth samples should be removed regardless of normalization
600 technique, especially if it is suspected that they contain a higher proportion of contaminants
601 (Kennedy et al. 2014; Salter et al. 2014). Also, when using alternatives to rarefying, the
602 researcher must check that clustering by sequence depth does not obscure biologically
603 meaningful results. Therefore, rarefying is still an extremely useful normalization technique,
604 especially for presence/absence metrics. Rarefying can erase the artifact of sample library size
605 better than other normalization techniques, and results in a higher PERMANOVA effect size
606 ($R^2$) for the studied biological effect, especially for small (<1000 sequences per sample), and
607 uneven library sizes between groups. For both normalization and differential abundance testing,
608 we stress that no normalization and naive proportion approaches should not be used as they can
609 generate artifactual clusters based on sequencing depth, and may result in mistaken OTU
610 differential abundance significance or insignificance.

612       For differential abundance testing, we studied the methods using both simulations and
613 real data. The most promising of current techniques are based on GLMs with either the negative
614 binomial or zero-inflated Gaussian distributions. It appears that DESeq2 (Love MI 2014),
615 metagenomeSeq's fitZIG (Paulson et al. 2013), and rarefying are all acceptable techniques for
616 approximately even library sizes and numbers of samples per class. DESeq2 was designed for,
617 and is a good option for, increased sensitivity on smaller datasets; however computation time
618 becomes very slow for larger datasets, especially over 100 samples per category.
619 MetagenomeSeq's fitZIG is a faster option for larger library sizes, although it may have a higher
620 false positive rate. The fitZIG technique is designed for larger sample sizes, since more counts
621 per OTU enables more accurate approximation of a continuous distribution. Rarefying, paired
622 with traditional non-parametric tests to account for the non-normal distribution of microbial data,
623 is useful for all dataset sizes, with sensitivity approaching parametric models in larger
624 datasets. Rarefying yields fewer OTUs as significantly differentially abundant, but those OTUs
625 are robust, in the sense that they are almost always identified as significant by at least one other
626 differential abundance detection model. In the case of highly uneven library sizes per category
627 (greater than 2-3x library size difference), we recommend rarefying, which provides higher
628 specificity at a cost to sensitivity, or metagenomeSeq's fitZIG, giving higher sensitivity at a cost
629 to specificity, over the DESeq2 technique. In situations with highly compositional data, no
630 technique does well.

632       Prior to differential abundance analysis, we recommend checking for significant
633 differences in library size means and distribution between categories (e.g. healthy vs. sick); and
634 propose a Mann-Whitney test, although the subject could be investigated further. The Mann-
635 Whitney test works on the library sizes simulated for this study, as well as that of McMurdie and
636 Holmes (McMurdie & Holmes 2014). To check distributional differences, the library sizes of
637 one sample category can be multiplied by a factor (e.g. 2) to make the means comparable prior to
638 applying the Mann-Whitney test. If there is a significant difference in either mean or
639 distribution, we recommend rarefying paired with a non-parametric test; if not, alternatives to
640 rarefying may be used. For the parametric differential abundance techniques, it is recommended
641 that rare OTUs be filtered out of the matrix prior to differential abundance testing. However, we

advise OTU filtering after rarefying, and then applying non-parametric tests. Thanks to McMurdie and Holmes' previous work in this area (McMurdie & Holmes 2014), we recognize the potential of these newer techniques, and have incorporated DESeq2 (Love MI 2014) and metagenomeSeq (Paulson et al. 2013) normalization and differential abundance testing into QIIME version 1.9.0 (Caporaso et al. 2010), along with the traditional rarefying and non-parametric testing techniques.

**Author Contributions**

     S.J.W, Z.Z.X., A.A., S.D.P., K.B., A.G., J.R.Z., and R.K. designed and conceived analyses. S.J.W. and J.N.P. wrote the QIIME scripts, and A.G.P. and Y.V.B. helped integrate the scripts into QIIME. S.J.W. wrote the initial manuscript, and all authors provided invaluable feedback and insights into analyses and the manuscript. All authors approved the final version of the manuscript.

Agresti A, and Hitchcock DB. *Bayesian inference for categorical data analysis : a survey*. Gainesville, Fla.?: University of Florida.

Aitchison J. 1982. The Statistical-Analysis of Compositional Data. *Journal of the Royal Statistical Society Series B-Methodological* 44:139-177.

Anders S, and Huber W. 2010. Differential expression analysis for sequence count data. *Genome Biol* 11.

Anders S, McCarthy DJ, Chen YS, Okoniewski M, Smyth GK, Huber W, and Robinson MD. 2013. Count-based differential expression analysis of RNA sequencing data using R and Bioconductor. *Nature Protocols* 8:1765-1786.

Anderson MJ. 2001. A new method for non-parametric multivariate analysis of variance. *Austral Ecology* 26:32-46.

Auer PL, and Doerge RW. 2010. Statistical design and analysis of RNA sequencing data. *Genetics* 185:405-416.

Benjamini Y, and Hochberg Y. 1995. Controlling the False Discovery Rate - a Practical and Powerful Approach to Multiple Testing. *Journal of the Royal Statistical Society Series B-Methodological* 57:289-300.

Bray JR, and Curtis JT. 1957. An Ordination of the Upland Forest Communities of Southern Wisconsin. *Ecological Monographs* 27:326-349.

Brewer A, and Williamson M. 1994. A New Relationship for Rarefaction. *Biodiversity and Conservation* 3:373-379.

Buccianti A, Mateu-Figueras G, and Pawlowsky-Glahn V. 2006. *Compositional data analysis in the geosciences : from theory to practice*. London: The Geological Society.

687    Bullard JH, Purdom E, Hansen KD, and Dudoit S. 2010. Evaluation of statistical methods for
688        normalization and differential expression in mRNA-Seq experiments. *BMC*
689        *Bioinformatics* 11:94.
690    Cameron AC, and Trivedi PK. *Regression analysis of count data*.
691    Caporaso JG, Kuczynski J, Stombaugh J, Bittinger K, Bushman FD, Costello EK, Fierer N,
692        Pena AG, Goodrich JK, Gordon JI, Huttley GA, Kelley ST, Knights D, Koenig JE, Ley RE,
693        Lozupone CA, McDonald D, Muegge BD, Pirrung M, Reeder J, Sevinsky JR, Turnbaugh
694        PJ, Walters WA, Widmann J, Yatsunenko T, Zaneveld J, and Knight R. 2010. QIIME
695        allows analysis of high-throughput community sequencing data. *Nat Methods* 7:335-
696        336.
697    Caporaso JG, Lauber CL, Costello EK, Berg-Lyons D, Gonzalez A, Stombaugh J, Knights D,
698        Gajer P, Ravel J, Fierer N, Gordon JI, and Knight R. 2011a. Moving pictures of the
699        human microbiome. *Genome Biol* 12:R50.
700    Caporaso JG, Lauber CL, Walters WA, Berg-Lyons D, Huntley J, Fierer N, Owens SM, Betley J,
701        Fraser L, Bauer M, Gormley N, Gilbert JA, Smith G, and Knight R. 2012. Ultra-high-
702        throughput microbial community analysis on the Illumina HiSeq and MiSeq
703        platforms. *ISME J* 6:1621-1624.
704    Caporaso JG, Lauber CL, Walters WA, Berg-Lyons D, Lozupone CA, Turnbaugh PJ, Fierer N,
705        and Knight R. 2011b. Global patterns of 16S rRNA diversity at a depth of millions of
706        sequences per sample. *Proc Natl Acad Sci U S A* 108 Suppl 1:4516-4522.
707    Carcer DA, Denman SE, McSweeney C, and Morrison M. 2011. Evaluation of subsampling-
708        based normalization strategies for tagged high-throughput sequencing data sets
709        from gut microbiomes. *Appl Environ Microbiol* 77:8795-8798.
710    Cheung YB. 2002. Zero-inflated models for regression analysis of count data: a study of
711        growth and development. *Stat Med* 21:1461-1469.
712    Colwell RK, Chao A, Gotelli NJ, Lin SY, Mao CX, Chazdon RL, and Longino JT. 2012. Models
713        and estimators linking individual-based and sample-based rarefaction,
714        extrapolation and comparison of assemblages. *Journal of Plant Ecology* 5:3-21.
715    Connolly SR, Dornelas M, Bellwood DR, and Hughes TP. 2009. Testing species abundance
716        models: a new bootstrap approach applied to Indo-Pacific coral reefs. *Ecology*
717        90:3138-3149.
718    Costea PI, Zeller G, Sunagawa S, and Bork P. 2014. A fair comparison. *Nat Methods* 11:359.
719    Costello EK, Lauber CL, Hamady M, Fierer N, Gordon JI, and Knight R. 2009. Bacterial
720        community variation in human body habitats across space and time. *Science*
721        326:1694-1697.
722    David LA, Maurice CF, Carmody RN, Gootenberg DB, Button JE, Wolfe BE, Ling AV, Devlin
723        AS, Varma Y, Fischbach MA, Biddinger SB, Dutton RJ, and Turnbaugh PJ. 2014. Diet
724        rapidly and reproducibly alters the human gut microbiome. *Nature* 505:559-563.
725    Dillies MA, Rau A, Aubert J, Hennequet-Antier C, Jeanmougin M, Servant N, Keime C, Marot
726        G, Castel D, Estelle J, Guernec G, Jagla B, Jouneau L, Laloe D, Le Gall C, Schaeffer B, Le
727        Crom S, Guedj M, and Jaffrezic F. 2013. A comprehensive evaluation of
728        normalization methods for Illumina high-throughput RNA sequencing data analysis.
729        *Brief Bioinform* 14:671-683.

730 Egozcue JJ, Pawlowsky-Glahn V, Mateu-Figueras G, and Barcelo-Vidal C. 2003. Isometric
731       logratio transformations for compositional data analysis. *Mathematical Geology*
732       35:279-300.
733 Friedman J, and Alm EJ. 2012. Inferring correlation networks from genomic survey data.
734       *PLoS Comput Biol* 8:e1002687.
735 Gentleman RC, Carey VJ, Bates DM, Bolstad B, Dettling M, Dudoit S, Ellis B, Gautier L, Ge YC,
736       Gentry J, Hornik K, Hothorn T, Huber W, Iacus S, Irizarry R, Leisch F, Li C, Maechler
737       M, Rossini AJ, Sawitzki G, Smith C, Smyth G, Tierney L, Yang JYH, and Zhang JH. 2004.
738       Bioconductor: open software development for computational biology and
739       bioinformatics. *Genome Biol* 5.
740 Gevers D, Kugathasan S, Denson LA, Vazquez-Baeza Y, Van Treuren W, Ren B, Schwager E,
741       Knights D, Song SJ, Yassour M, Morgan XC, Kostic AD, Luo C, Gonzalez A, McDonald
742       D, Haberman Y, Walters T, Baker S, Rosh J, Stephens M, Heyman M, Markowitz J,
743       Baldassano R, Griffiths A, Sylvester F, Mack D, Kim S, Crandall W, Hyams J,
744       Huttenhower C, Knight R, and Xavier RJ. 2014. The treatment-naive microbiome in
745       new-onset Crohn's disease. *Cell Host Microbe* 15:382-392.
746 Gotelli NJ, and Colwell RK. 2001. Quantifying biodiversity: procedures and pitfalls in the
747       measurement and comparison of species richness. *Ecology Letters* 4:379-391.
748 Gower JC. 1966. Some Distance Properties of Latent Root and Vector Methods Used in
749       Multivariate Analysis. *Biometrika* 53:325-&.
750 Greenacre M. 2011. Measuring Subcompositional Incoherence. *Mathematical Geosciences*
751       43:681-693.
752 Holmes I, Harris K, and Quince C. 2012. Dirichlet multinomial mixtures: generative models
753       for microbial metagenomics. *PLoS One* 7:e30126.
754 Horner-Devine MC, Lage M, Hughes JB, and Bohannan BJ. 2004. A taxa-area relationship for
755       bacteria. *Nature* 432:750-753.
756 Jari Oksanen FGB, Roeland Kindt, Pierre Legendre, Peter R. Minchin, R. B. O'Hara, Gavin L.
757       Simpson, Peter Solymos, M. Henry H. Stevens and Helene Wagner. 2015. vegan:
758       Community Ecology Package. *R package version 22-1*.
759 Jernvall J, and Wright PC. 1998. Diversity components of impending primate extinctions.
760       *Proc Natl Acad Sci U S A* 95:11279-11283.
761 Kaufman L. RP. 1990. *Finding Groups in Data: An introduction to Cluster Analysis*: JohnWiley
762       & Sons.
763 Kelly CR, Ihunnah C, Fischer M, Khoruts A, Surawicz C, Afzali A, Aroniadis O, Barto A,
764       Borody T, Giovanelli A, Gordon S, Gluck M, Hohmann EL, Kao D, Kao JY, McQuillen
765       DP, Mellow M, Rank KM, Rao K, Ray A, Schwartz MA, Singh N, Stollman N, Suskind
766       DL, Vindigni SM, Youngster I, and Brandt L. 2014. Fecal microbiota transplant for
767       treatment of Clostridium difficile infection in immunocompromised patients. *Am J*
768       *Gastroenterol* 109:1065-1071.
769 Kennedy K, Hall MW, Lynch MD, Moreno-Hagelsieb G, and Neufeld JD. 2014. Evaluating bias
770       of illumina-based bacterial 16S rRNA gene profiles. *Appl Environ Microbiol* 80:5717-
771       5722.
772 Lauber CL, Hamady M, Knight R, and Fierer N. 2009. Pyrosequencing-based assessment of
773       soil pH as a predictor of soil bacterial community structure at the continental scale.
774       *Appl Environ Microbiol* 75:5111-5120.

775  Law CW, Chen YS, Shi W, and Smyth GK. 2014. voom: precision weights unlock linear model
776      analysis tools for RNA-seq read counts. *Genome Biol* 15.
777  Legendre P, and Gallagher ED. 2001. Ecologically meaningful transformations for
778      ordination of species data. *Oecologia* 129:271-280.
779  Ley RE, Backhed F, Turnbaugh P, Lozupone CA, Knight RD, and Gordon JI. 2005. Obesity
780      alters gut microbial ecology. *Proc Natl Acad Sci U S A* 102:11070-11075.
781  Love MI HWaAS. 2014. Moderated estimation of fold change and dispersion for RNA-seq
782      data with DESeq2. *Genome Biol* 15.
783  Lovell D ea. 2010. Caution! compositions! can constraints on omics data lead analyses
784      astray? *CSIRO*:1-44.
785  Lovell D, Pawlowsky-Glahn V, Egozcue JJ, Marguerat S, and Bahler J. 2015. Proportionality:
786      a valid alternative to correlation for relative data. *PLoS Comput Biol* 11:e1004075.
787  Lozupone C, and Knight R. 2005. UniFrac: a new phylogenetic method for comparing
788      microbial communities. *Appl Environ Microbiol* 71:8228-8235.
789  Lozupone C, Lladser ME, Knights D, Stombaugh J, and Knight R. 2011. UniFrac: an effective
790      distance metric for microbial community comparison. *ISME J* 5:169-172.
791  Lozupone CA, Hamady M, Kelley ST, and Knight R. 2007. Quantitative and qualitative beta
792      diversity measures lead to different insights into factors that structure microbial
793      communities. *Appl Environ Microbiol* 73:1576-1585.
794  Lozupone CA, and Knight R. 2007. Global patterns in bacterial diversity. *Proc Natl Acad Sci*
795      *U S A* 104:11436-11440.
796  Lozupone CA, Li M, Campbell TB, Flores SC, Linderman D, Gebert MJ, Knight R, Fontenot AP,
797      and Palmer BE. 2013a. Alterations in the Gut Microbiota Associated with HIV-1
798      Infection. *Cell Host Microbe* 14:329-339.
799  Lozupone CA, Stombaugh J, Gonzalez A, Ackermann G, Wendel D, Vazquez-Baeza Y, Jansson
800      JK, Gordon JI, and Knight R. 2013b. Meta-analyses of studies of the human
801      microbiota. *Genome Research* 23:1704-1714.
802  McMurdie PJ, and Holmes S. 2013. phyloseq: An R Package for Reproducible Interactive
803      Analysis and Graphics of Microbiome Census Data. *PLoS One* 8.
804  McMurdie PJ, and Holmes S. 2014. Waste Not, Want Not: Why Rarefying Microbiome Data
805      Is Inadmissible. *PLoS Comput Biol* 10.
806  Paulson JN, Bravo HC, and Pop M. 2014. Reply to: "a fair comparison". *Nat Methods* 11:359-
807      360.
808  Paulson JN, Stine OC, Bravo HC, and Pop M. 2013. Differential abundance analysis for
809      microbial marker-gene surveys. *Nat Methods* 10:1200-1202.
810  Pearson K. 1896. Mathematical contributions to the theory of evolution: On a form of
811      spurious correlation which may arise when indices are used in the measurements of
812      organs. *Proc Roy Soc* 60:489-498.
813  Piombino P, Genovese A, Esposito S, Moio L, Cutolo PP, Chambery A, Severino V, Moneta E,
814      Smith DP, Owens SM, Gilbert JA, and Ercolini D. 2014. Saliva from obese individuals
815      suppresses the release of aroma compounds from wine. *PLoS One* 9:e85611.
816  Pop M, Walker AW, Paulson J, Lindsay B, Antonio M, Hossain MA, Oundo J, Tamboura B, Mai
817      V, Astrovskaya I, Corrada Bravo H, Rance R, Stares M, Levine MM, Panchalingam S,
818      Kotloff K, Ikumapayi UN, Ebruke C, Adeyemi M, Ahmed D, Ahmed F, Alam MT, Amin
819      R, Siddiqui S, Ochieng JB, Ouma E, Juma J, Mailu E, Omore R, Morris JG, Breiman RF,

820          Saha D, Parkhill J, Nataro JP, and Stine OC. 2014. Diarrhea in young children from
821          low-income countries leads to large-scale alterations in intestinal microbiota
822          composition. *Genome Biol* 15:R76.
823 Qin J, Li R, Raes J, Arumugam M, Burgdorf KS, Manichanh C, Nielsen T, Pons N, Levenez F,
824          Yamada T, Mende DR, Li J, Xu J, Li S, Li D, Cao J, Wang B, Liang H, Zheng H, Xie Y, Tap
825          J, Lepage P, Bertalan M, Batto JM, Hansen T, Le Paslier D, Linneberg A, Nielsen HB,
826          Pelletier E, Renault P, Sicheritz-Ponten T, Turner K, Zhu H, Yu C, Li S, Jian M, Zhou Y,
827          Li Y, Zhang X, Li S, Qin N, Yang H, Wang J, Brunak S, Dore J, Guarner F, Kristiansen K,
828          Pedersen O, Parkhill J, Weissenbach J, Bork P, Ehrlich SD, and Wang J. 2010. A
829          human gut microbial gene catalogue established by metagenomic sequencing.
830          *Nature* 464:59-65.
831 Rapaport F, Khanin R, Liang Y, Pirun M, Krek A, Zumbo P, Mason CE, Socci ND, and Betel D.
832          2013. Comprehensive evaluation of differential gene expression analysis methods
833          for RNA-seq data. *Genome Biol* 14:R95.
834 Reynolds A RG, Iglesia B, Rayward-Smith V. 2006. Clustering rules: A comparison of
835          partitioning and hierarchical clustering algorithms. *Journal of Mathematical*
836          *Modelling and Algorithms* 5:475-504.
837 Ridaura VK, Faith JJ, Rey FE, Cheng J, Duncan AE, Kau AL, Griffin NW, Lombard V, Henrissat
838          B, Bain JR, Muehlbauer MJ, Ilkayeva O, Semenkovich CF, Funai K, Hayashi DK, Lyle
839          BJ, Martini MC, Ursell LK, Clemente JC, Van Treuren W, Walters WA, Knight R,
840          Newgard CB, Heath AC, and Gordon JI. 2013. Gut microbiota from twins discordant
841          for obesity modulate metabolism in mice. *Science* 341:1241214.
842 Robinson MD, McCarthy DJ, and Smyth GK. 2010. edgeR: a Bioconductor package for
843          differential expression analysis of digital gene expression data. *Bioinformatics*
844          26:139-140.
845 Robinson MD, and Oshlack A. 2010. A scaling normalization method for differential
846          expression analysis of RNA-seq data. *Genome Biol* 11:R25.
847 Robinson MD, and Smyth GK. 2008. Small-sample estimation of negative binomial
848          dispersion, with applications to SAGE data. *Biostatistics* 9:321-332.
849 Robles JA, Qureshi SE, Stephen SJ, Wilson SR, Burden CJ, and Taylor JM. 2012. Efficient
850          experimental design and analysis strategies for the detection of differential
851          expression using RNA-Sequencing. *BMC Genomics* 13:484.
852 Rodriguez RL, and Konstantinidis KT. 2014. Estimating coverage in metagenomic data sets
853          and why it matters. *ISME J* 8:2349-2351.
854 Salter SJ, Cox MJ, Turek EM, Calus ST, Cookson WO, Moffatt MF, Turner P, Parkhill J, Loman
855          NJ, and Walker AW. 2014. Reagent and laboratory contamination can critically
856          impact sequence-based microbiome analyses. *BMC Biol* 12:87.
857 Schloss PD, Westcott SL, Ryabin T, Hall JR, Hartmann M, Hollister EB, Lesniewski RA,
858          Oakley BB, Parks DH, Robinson CJ, Sahl JW, Stres B, Thallinger GG, Van Horn DJ, and
859          Weber CF. 2009. Introducing mothur: open-source, platform-independent,
860          community-supported software for describing and comparing microbial
861          communities. *Appl Environ Microbiol* 75:7537-7541.
862 Shankar V, Hamilton MJ, Khoruts A, Kilburn A, Unno T, Paliy O, and Sadowsky MJ. 2014.
863          Species and genus level resolution analysis of gut microbiota in Clostridium difficile
864          patients following fecal microbiota transplantation. *Microbiome* 2:13.

865 Soneson C, and Delorenzi M. 2013. A comparison of methods for differential expression
866      analysis of RNA-seq data. *BMC Bioinformatics* 14:91.
867 Team RC. 2014. R: A language and environment for statistical computing. *R Foundation for*
868      *Statistical Computing, Vienna, Austria*.
869 Turnbaugh PJ, Hamady M, Yatsunenko T, Cantarel BL, Duncan A, Ley RE, Sogin ML, Jones
870      WJ, Roe BA, Affourtit JP, Egholm M, Henrissat B, Heath AC, Knight R, and Gordon JI.
871      2009. A core gut microbiome in obese and lean twins. *Nature* 457:480-484.
872 Vazquez-Baeza Y, Pirrung M, Gonzalez A, and Knight R. 2013. EMPeror: a tool for
873      visualizing high-throughput microbial community data. *Gigascience* 2:16.
874 Wagner BD, Robertson CE, and Harris JK. 2011. Application of two-part statistics for
875      comparison of sequence variant counts. *PLoS One* 6:e20296.
876 Weingarden A, Gonzalez A, Vazquez-Baeza Y, Weiss S, Humphry G, Berg-Lyons D, Knights D,
877      Unno T, Bobr A, Kang J, Khoruts A, Knight R, and Sadowsky MJ. 2015. Dynamic
878      changes in short- and long-term bacterial composition following fecal microbiota
879      transplantation for recurrent Clostridium difficile infection. *Microbiome* 3:10.
880 White JR, Nagarajan N, and Pop M. 2009. Statistical methods for detecting differentially
881      abundant features in clinical metagenomic samples. *PLoS Comput Biol* 5:e1000352.
882 Witten DM. 2011. Classification and Clustering of Sequencing Data Using a Poisson Model.
883      *Annals of Applied Statistics* 5:2493-2518.
884 Yu D, Huber W, and Vitek O. 2013. Shrinkage estimation of dispersion in Negative Binomial
885      models for RNA-seq experiments with small sample size. *Bioinformatics* 29:1275-
886      1282.
887
888
889
890
891
892 FIGURE CAPTIONS:
893
894 **Figure 1: Effect of sampling depth on ordination methods.** (a) Data rarefied at 500 sequences
895 per sample**.** (b, c) Data not normalized, with a random half of the samples subsampled to 500
896 sequences per sample and the other half to 50 sequences per sample. (b) is colored by
897 subject_ID, (c) is colored by sequences per sample. Non-parametric ANOVA (PERMANOVA)
898 effect sizes ($R^2$) roughly represent the percent variance that can be explained by the given
899 variable. Asterisk (*) indicates significance at $p < 0.01$. The distance metric of unweighted
900 UniFrac was used for all panels.

901 **Figure 2:  Comparison of common distance metrics and normalization methods across**
902 **library sizes when low-coverage samples are excluded.**
903 Clustering accuracy is shown for all combinations of five common distance metrics (panels
904 arranged from left to right) across four library depths (panels arranged from top to bottom; $N_L$,
905 median library size), six sample normalization methods (series within each panel), and several
906 effect sizes (x-axis within panels).  In all cases, samples below the 15th percentile of library size
907 were dropped from the analysis in order to isolate the effects of rarifying from the effects of
908 dropping low-coverage samples. The x-axis ('effect size') within each panel represents the
909 multinomial mixing proportions of the two sample classes *'Ocean'* and '*Feces*'.  A higher effect

910 size represents an easier clustering task. The y-axis ('accuracy') shows the accuracy of each
911 classifier, as assessed by the fraction of simulated samples correctly clustered.
912
913 **Figure 3: Rarefying clusters more according to biological origin, and diminishes the effect**
914 **of library size.** Rarefying exhibits a higher effect size ($R^2$) for biological origin, and a lower
915 effect size ($R^2$) of original library size. Unweighted UniFrac was used for clustering, and a
916 random half of samples were subsampled to 10 times fewer sequences per sample.  The 45-
917 degree line splits low from high depth samples in all but the rarefying technique. For each letter
918 (a-f), the left PCoA plot is colored according to the 'Canine Feces', etc. legend, and the right
919 PCoA plot is colored according to the 'High/Low Library Size' legend.
920
921 **Figure 4: Differential abundance detection performance.**
922 The AUC ('Area Under the Curve') version of the ROC ('Receiver Operator Characteristic')
923 curve is the ratio of sensitivity to (1-specificity), or true positive rate vs. false positive rate. A
924 higher AUC indicates better differential abundance detection performance. The 'effect size'
925 represents the fold-change of the 'true positive' OTUs from one condition (e.g. case) to another
926 (e.g. control). The right axis represents the median library size ($N_L$), while the shading on the
927 graph lines represents the number of samples per class.  'Model/None' represents data analyzed
928 with a parametric statistical model (e.g. DESeq), or no normalization.  Blue lines in, e.g. the
929 DESeq column represents the data was rarefied, then DESeq was applied.  Since the fitZIG
930 model depends upon original library size information, the model does poorly on rarefied data.
931
932 **Figure 5: The effect of rarefying on power for different OTU relative abundances and**
933 **library sizes.**
934 The detection power for differentially abundant OTUs of varying levels of relative abundance
935 (very rare to common). This is for two samples A and B.  For power calculations, we assumed
936 that OTU1 fraction of group B is 85% of the OTU1 fraction of group A. Library type A was
937 fixed, while library size B was subsampled at different percentages, creating the power curves
938 calculated with Fisher's exact test.
939
940 **Figure 6: Comparison of the most promising differential abundance detection techniques**
941 **on real datasets.**
942 Each table's diagonal represents the number of OTUs found significant (Benjamini & Hochberg
943 FDR < 0.05) by that technique.  The off-diagonal entries represent the number of shared
944 differentially abundant OTUs between two techniques. The bar charts represents the percentage
945 of differentially abundant OTUs shared by at least one other technique.
946
947 **Figure 7: Differential abundance detection performance where one sample group average**
948 **library size is 10 times the size of the other.**  Labels are the same as in Fig. 4.  A significant
949 difference from the results of Fig. 4 was first observed at 2-3-fold difference in library sizes (see
950 Fig. S6).
951
952 **Figure S1: Simulated clustering accuracy if rarefying is not penalized for removing the**
953 **lowest 15th percentile samples.**
954 The right axis represents the median library size ($N_L$), while the x-axis 'effect size' is the
955 multinomial mixing proportions of the two classes of samples, *'Ocean'* and *'Feces'*. See caption

956     for Fig. 2 for further details.

957

958     **Figure S2: Low library size samples can diminish result quality, regardless of**
959     **normalization technique.**  We show the inflammatory bowel disease (IBD) dataset of Gevers et
960     al.(Gevers et al. 2014), which has an average library size 375 sequences per sample. **(a)**
961     Extremely low depth samples cluster in lower right hand corner of PCoA plots with no
962     normalization, or rarefying alternatives, unweighted UniFrac. (b) The original library size of
963     samples is a dominant effect, even influencing weighted UniFrac, with low library sizes and
964     subtle biological clustering for rarefying alternatives.  This diminishes if low library size
965     samples are removed from analysis.

966

967     **Figure S3: All normalization techniques on key microbiome datasets, Bray Curtis distance.**
968     Rows of panels show (from top to bottom) data from 88soils (Lauber et al. 2009), Body Sites
969     (Costello et al. 2009), Moving Pictures (Caporaso et al. 2011a).  88 soils is colored according to
970     a color gradient from low to high pH.  The Costello et al. body sites dataset is colored according
971     to body site: feces (blue), oral cavity (purple), the rest of the colors are external auditory canal,
972     hair, nostril, skin, and urine. Moving Pictures dataset: Left and Right palm (red/blue), tongue
973     (green), feces (orange).  It is important to note that all the samples in these datasets are
974     approximately the same depth, and there are very strong driving gradients.

975

976     **Figure S4: All normalization techniques on key microbiome datasets, unweighed UniFrac**
977     **distance.**  See Figure S3 caption for details.

978

979     **Figure S5: Simple example of the reasoning behind differential abundance simulations.** (a)
980     In actual OTU tables generated from sequencing data, the counts (left column) are already
981     compositional and therefore only relative (left column). Application of the 'effect size' to the
982     original 'Multinomial' template to create fold-change differences disturbs the distinction
983     between true positive (TP) and true negative (TN) OTUs in the 'Original' simulation, but not the
984     'Balanced' simulation. (c) Creation of a 'Compositional' OTU table from the 'Multinomial'
985     template, where the counts/relative abundances are intentionally blurred for the TN OTUs.

986

987     **Figure S6: Differential abundance detection performance where one sample group average**
988     **library size is 3 times the size of the other.**  Labels are the same as in Fig. 4.

989

990     **Figure S7: Differential abundance detection performance when the dataset is**
991     **compositional**. 25% of OTUs are differentially abundant.  Labels are the same as in Fig. 4.
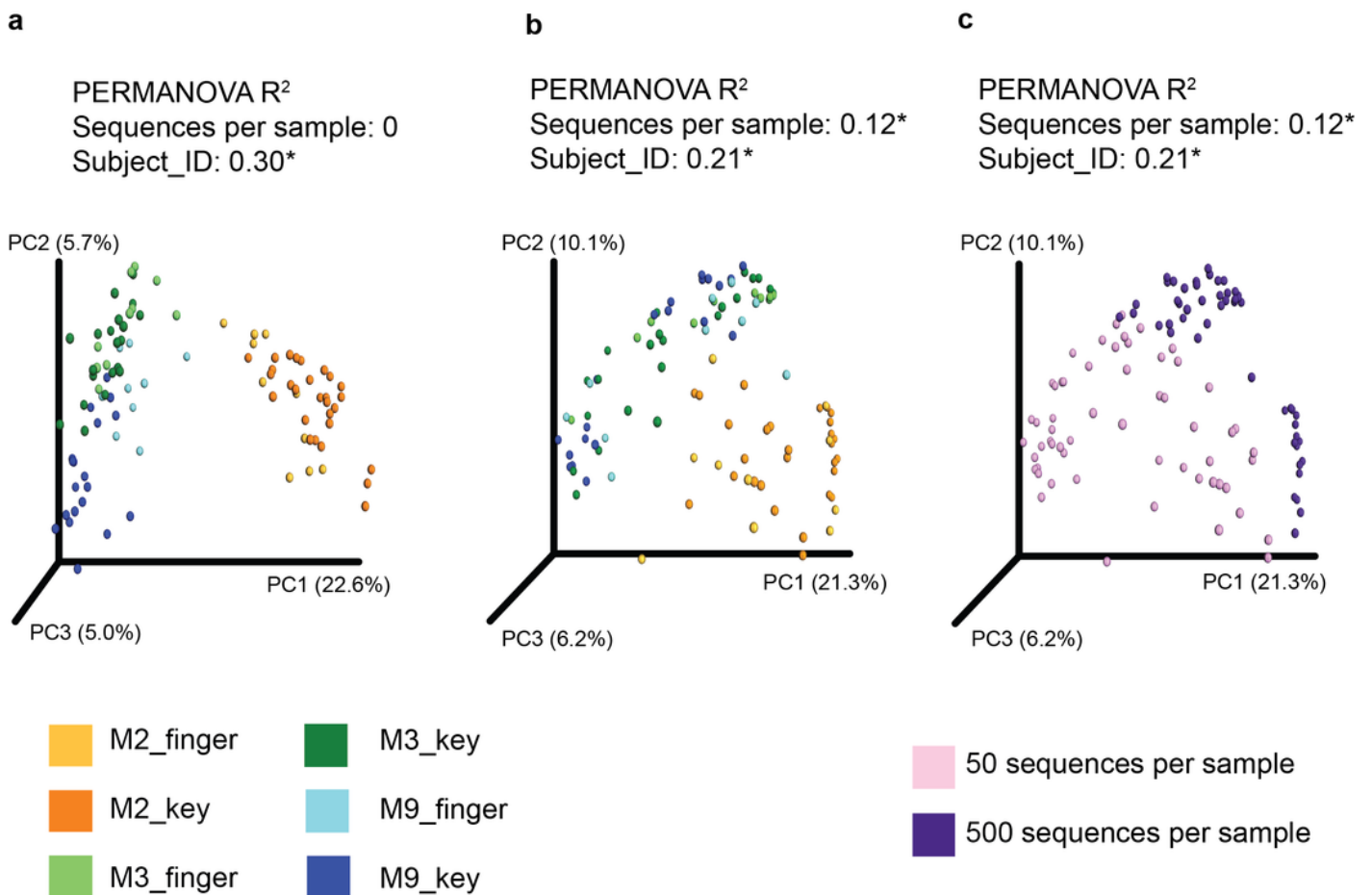
992

993

994

995

996

# 1
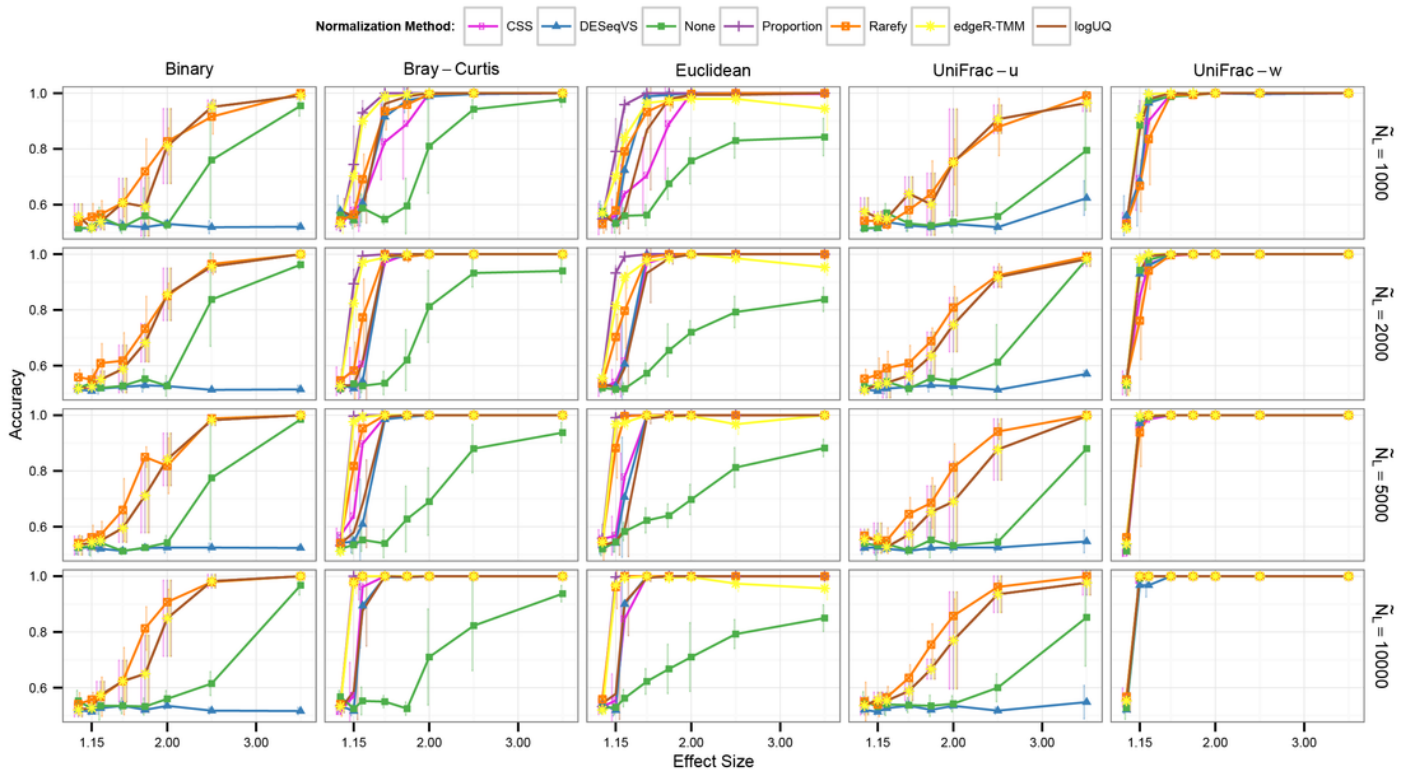
Effect of sampling depth on ordination methods

(a) Data rarefied at 500 sequences per sample. (b, c) Data not normalized, with a random half of the samples subsampled to 500 sequences per sample and the other half to 50 sequences per sample. (b) is colored by subject_ID, (c) is colored by sequences per sample. Non-parametric ANOVA (PERMANOVA) effect sizes ($R^2$) roughly represent the percent variance that can be explained by the given variable. Asterisk (*) indicates significance at p < 0.01. The distance metric of unweighted UniFrac was used for all panels.



**a**

PERMANOVA $R^2$
Sequences per sample: 0
Subject_ID: 0.30*

PC2 (5.7%)
PC1 (22.6%)
PC3 (5.0%)

**b**

PERMANOVA $R^2$
Sequences per sample: 0.12*
Subject_ID: 0.21*

PC2 (10.1%)
PC1 (21.3%)
PC3 (6.2%)

**c**

PERMANOVA $R^2$
Sequences per sample: 0.12*
Subject_ID: 0.21*

PC2 (10.1%)
PC1 (21.3%)
PC3 (6.2%)

M2_finger    M3_key
M2_key       M9_finger
M3_finger    M9_key

50 sequences per sample
500 sequences per sample

# 2

Comparison of common distance metrics and normalization methods across library sizes when low-coverage samples are excluded.

Clustering accuracy is shown for all combinations of five common distance metrics (panels arranged from left to right) across four library depths (panels arranged from top to bottom; $N_L$, median library size), six sample normalization methods (series within each panel), and several effect sizes (x-axis within panels). In all cases, samples below the 15th percentile of library size were dropped from the analysis in order to isolate the effects of rarifying from the effects of dropping low-coverage samples. The x-axis ('effect size') within each panel represents the multinomial mixing proportions of the two sample classes '*Ocean*' and '*Feces*'. A higher effect size represents an easier clustering task. The y-axis ('accuracy') shows the accuracy of each classifier, as assessed by the fraction of simulated samples correctly clustered.

3

Rarefying clusters more according to biological origin, and diminishes the effect of library size.

Rarefying exhibits a higher effect size ($R^2$) for biological origin, and a lower effect size ($R^2$) of original library size. Unweighted UniFrac was used for clustering, and a random half of samples were subsampled to 10 times fewer sequences per sample. The 45-degree line splits low from high depth samples in all but the rarefying technique. For each letter (a-f), the left PCoA plot is colored according to the 'Canine Feces', etc. legend, and the right PCoA plot is colored according to the 'High/Low Library Size' legend.

Legend:
- Canine Feces (red)
- Human Feces (yellow)
- Canine Mouth (blue)
- Human Mouth (teal)
- Canine Skin (orange)
- Human Skin (purple)
- Soil (light blue)
- High Library Size (green)
- Low Library Size (magenta)

*** p<0.001    ** p<0.01    *p<0.05

**a  None: PERMANOVA R²**
Library_Size: 0.14***
Host_Bodysite: 0.40***

**b  rarefy: PERMANOVA R²**
Library_Size: 0.045*
Host_Bodysite: 0.56***

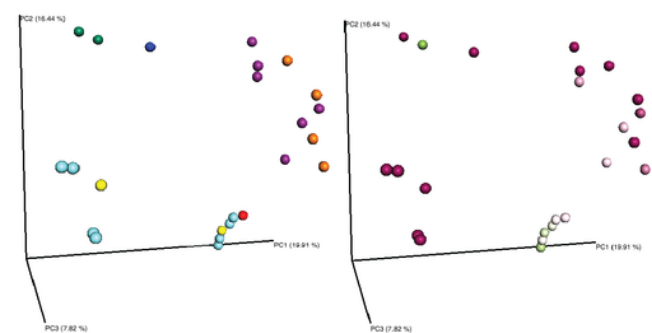**c  Proportion: PERMANOVA R²**
Library_Size: 0.14***
Host_Bodysite: 0.40***

**d  CSS: PERMANOVA R²**
Library_Size: 0.13***
Host_Bodysite: 0.42***

**e  DESeq: PERMANOVA R²**
Library_Size: 0.12***
Host_Bodysite: 0.46***

**f  edgeR-TMM: PERMANOVA R²**
Library_Size: 0.10***
Host_Bodysite: 0.46***

# 4

Differential abundance detection performance.

The AUC ('Area Under the Curve') version of the ROC ('Receiver Operator Characteristic') curve is the ratio of sensitivity to (1-specificity), or true positive rate vs. false positive rate. A higher AUC indicates better differential abundance detection performance. The 'effect size' represents the fold-change of the 'true positive' OTUs from one condition (e.g. case) to another (e.g. control). The right axis represents the median library size ($N_L$), while the shading on the graph lines represents the number of samples per class. 'Model/None' represents data analyzed with a parametric statistical model (e.g. DESeq), or no normalization. Blue lines in, e.g. the DESeq column represents the data was rarefied, then DESeq was applied. Since the fitZIG model depends upon original library size information, the model does poorly on rarefied data.
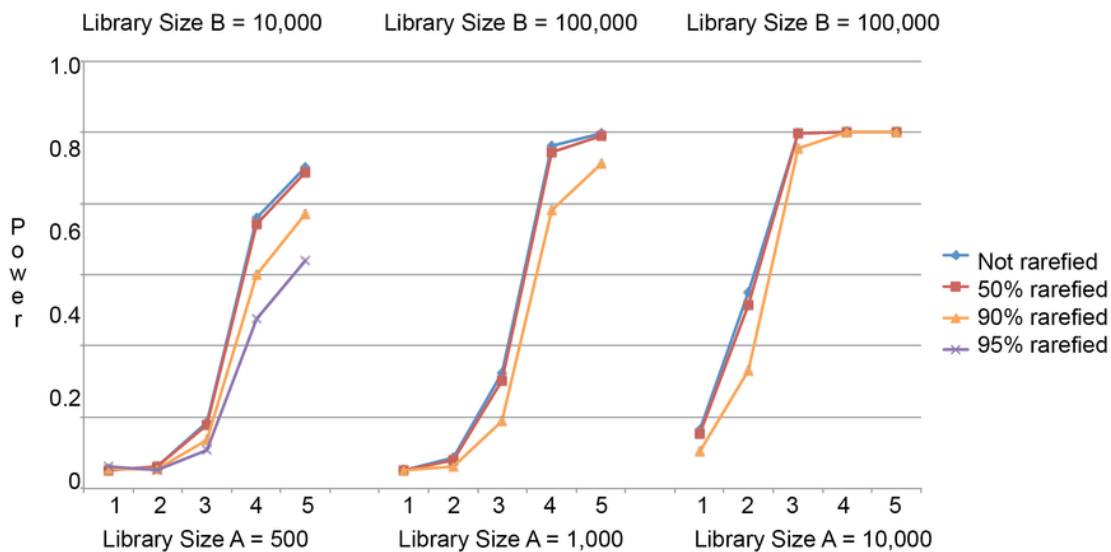
# 5

The effect of rarefying on power for different OTU relative abundances and library sizes.

The detection power for differentially abundant OTUs of varying levels of relative abundance (very rare to common). This is for two samples A and B. For power calculations, we assumed that OTU1 fraction of group B is 85% of the OTU1 fraction of group A. Library type A was fixed, while library size B was subsampled at different percentages, creating the power curves calculated with Fisher's exact test.
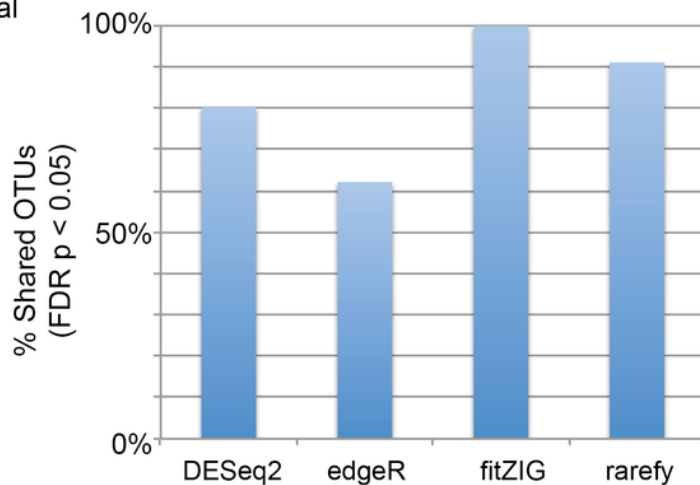
# 6

Comparison of the most promising differential abundance detection techniques on real datasets.

Each table's diagonal represents the number of OTUs found significant (Benjamini & Hochberg FDR < 0.05) by that technique. The off-diagonal entries represent the number of shared differentially abundant OTUs between two techniques. The bar charts represents the percentage of differentially abundant OTUs shared by at least one other technique.

**a** Caporaso *et al.* Ultra-high-throughput microbial community analysis on the Illumina HiSeq and MiSeq platforms *ISME* (2012).
~ 6 skin samples, 8 soil samples
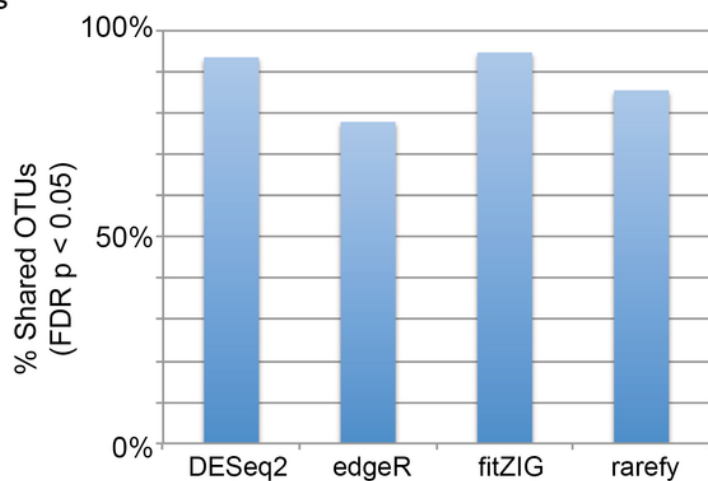mean sequences per sample: 1.3 million

|  | DESeq2 | edgeR | fitZIG | rarefy |
|---|---|---|---|---|
| DESeq2 | **1182** | 1135 | 934 | 376 |
| edgeR |  | **2357** | 1344 | 706 |
| fitZIG |  |  | **1445** | 680 |
| rarefy |  |  |  | **763** |



**b** Piombino *et al.* Saliva from Obese Individuals Suppresses the Release of Aroma Compounds from Wine. *PLoS One* (2014).
~28 samples per category (lean vs. obese)
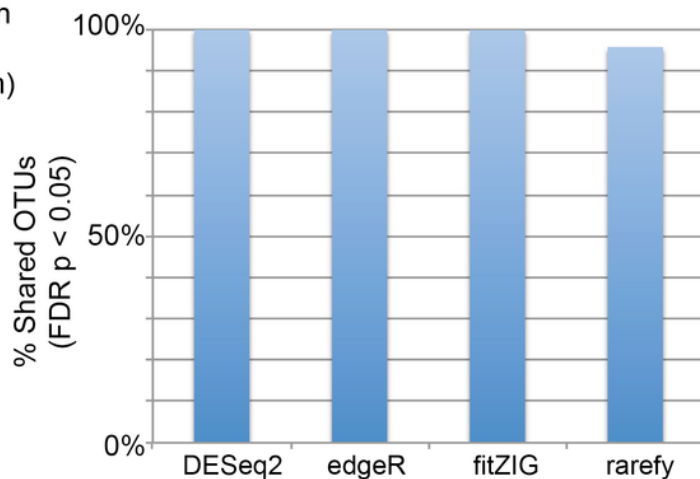mean sequences per sample: 75,580

|  | DESeq2 | edgeR | fitZIG | rarefy |
|---|---|---|---|---|
| DESeq2 | **232** | 206 | 164 | 122 |
| edgeR |  | **267** | 159 | 113 |
| fitZIG |  |  | **189** | 128 |
| rarefy |  |  |  | **145** |



**c** Caporaso *et al.* Moving Pictures of the Human Microbiome. *Genome Biol.* (2011).
~500 samples per category (tongue vs. left palm)
mean sequences per sample: 25,600

|  | DESeq2 | edgeR | fitZIG | rarefy |
|---|---|---|---|---|
| DESeq2 | **1070** | 1051 | 955 | 1038 |
| edgeR |  | **1122** | 995 | 1097 |
| fitZIG |  |  | **1038** | 1004 |
| rarefy |  |  |  | **1185** |

# 7

Differential abundance detection performance where one sample group average library size is 10 times the size of the other.

Labels are the same as in Fig. 4. A significant difference from the results of Fig. 4 was first observed at 2-3-fold difference in library sizes (see Fig. S6).