

Analysing microbiome intervention design studies: Comparison of alternative multivariate statistical methods

Maryia Khomich (✉ maryia.khomich@uib.no)

University of Bergen <https://orcid.org/0000-0002-6840-5739>

Ingrid Måge (✉ ingrid.mage@nofima.no)

Nofima <https://orcid.org/0000-0003-0364-0225>

Ida Rud

Nofima <https://orcid.org/0000-0002-1758-292X>

Ingunn Berget

Nofima <https://orcid.org/0000-0003-1027-1472>

Research Article

Keywords: gut microbiome, dietary intervention trials, differential abundance, multivariate, ANOVA, ASCA, FFMANOVA, ALDEx2, ANCOM, DESeq2, PERMANOVA, ANOSIM, SIMPER

Posted Date: September 16th, 2021

DOI: <https://doi.org/10.21203/rs.3.rs-910076/v1>

License:  This work is licensed under a Creative Commons Attribution 4.0 International License.

[Read Full License](#)

1 **Analysing microbiome intervention design studies: Comparison of** 2 **alternative multivariate statistical methods**

3 Maryia Khomich^{1,2*}, Ingrid Måge^{3*}, Ida Rud¹ and Ingunn Berget³

4

5 ¹Department of Food safety and quality, Division Food Science, Nofima – Norwegian Institute
6 of Food, Fisheries and Aquaculture Research, P.O. Box 210, 1431 Ås, Norway

7 ²Department of Clinical Science, University of Bergen, P.O. Box 7804, 5020 Bergen, Norway

8 ³Department of Raw materials and process optimisation, Division Food Science, Nofima –
9 Norwegian Institute of Food, Fisheries and Aquaculture Research, P.O. Box 210, 1431 Ås,
10 Norway

11 *Corresponding authors

12

13 Maryia Khomich: marykhomich@gmail.com, maryia.khomich@uib.no; ORCID: 0000-0002-
14 6840-5739

15 Ingrid Måge: ingrid.mage@nofima.no; ORCID: 0000-0003-0364-0225

16 Ida Rud: ida.rud@nofima.no; ORCID: 0000-0002-1758-292X

17 Ingunn Berget: ingunn.berget@nofima.no; ORCID: 0000-0003-1027-1472

18

19 **Abstract**

20 The diet plays a major role in shaping gut microbiome composition and function in both
21 humans and animals, and dietary intervention trials are often used to investigate and understand
22 these effects. A plethora of statistical methods for analysing the differential abundance of
23 microbial taxa exists, and new methods are constantly being developed, but there is a lack of
24 benchmarking studies and clear consensus on the best multivariate statistical practices. This
25 makes it hard for a biologist to decide which method to use. We compared the outcomes of
26 generic multivariate ANOVA (ASCA and FFMANOVA) against statistical methods commonly
27 used for community analyses (PERMANOVA and SIMPER) and methods designed for
28 analysis of count data from high-throughput sequencing experiments (ALDEx2, ANCOM and
29 DESeq2). The comparison is based on both simulated data and five published dietary
30 intervention trials representing different subjects and study designs. We found that the methods

31 testing differences at the community level were in agreement regarding both effect size and
32 statistical significance. However, the methods that provided ranking and identification of
33 differentially abundant operational taxonomic units (OTUs) gave incongruent results, implying
34 that the choice of method is likely to influence the biological interpretations. The generic
35 multivariate ANOVA tools have the flexibility needed for analysing multifactorial experiments
36 and provide outputs at both the community and OTU levels; good performance in the simulation
37 studies suggests that these statistical tools are also suitable for microbiome data sets.

38

39 **Keywords**

40 gut microbiome, dietary intervention trials, differential abundance, multivariate,
41 ANOVA, ASCA, FFMANOVA, ALDEx2, ANCOM, DESeq2, PERMANOVA, ANOSIM,
42 SIMPER

43

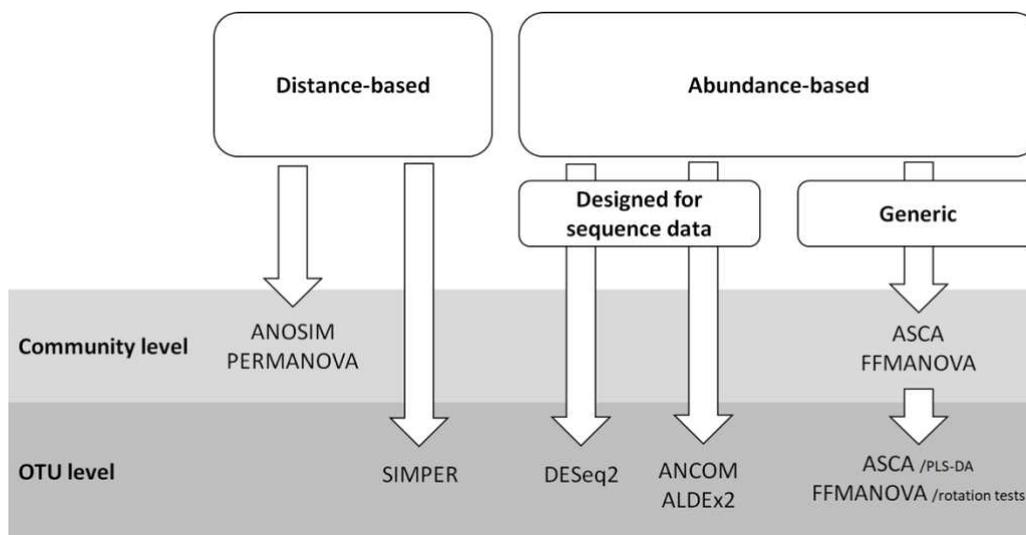
44 **Introduction**

45 The microbiome has emerged as an important link to health and disease [1]. Microbiome
46 analysis methods are rapidly advancing, in particular in areas such as compositional data
47 analysis, multi-omics and data integration [2, 3]. A clear understanding of the type of data being
48 analysed is crucial, given the growing number of studies uncovering the key role of
49 microbiome, its composition and functions following diet intervention or medical treatment [4].
50 At present, analysis of complex microbial data benefits from adapting the multivariate statistical
51 toolbox from ecology and environmental sciences, and a proper choice of statistical tools is
52 becoming increasingly important [5-7]. However, a lack of benchmarking studies and clear
53 consensus on the best multivariate statistical practices make comparisons across microbiome
54 data sets difficult [2, 8]. New methods are often tested by simulation studies, but there is always
55 a concern that simulations can be biased towards the tested statistical model and cannot mimic
56 the complexity of real microbiome data [6, 9]. Moreover, newly introduced tools are often
57 optimised, whereas the comparison of several statistical methods implies the use of standard or
58 default parameters [6, 9]. It is therefore of interest to compare existing methods on real data
59 sets of different complexity, in addition to simulation studies, to better understand how choice
60 of method affects the results.

61 Different statistical methods have different properties, and the choice of method should
62 depend on the scientific question, experimental design, data characteristics and expected

63 relationships among the variables. Furthermore, the choice of method is often biased by the
 64 research groups' tradition and familiarity with specific "toolboxes". Main differences between
 65 existing statistical approaches for analysing microbiome data are related to: (1) explorative
 66 versus confirmative; (2) univariate versus multivariate; (3) parametric versus nonparametric;
 67 (4) linear versus nonlinear; (5) compositional versus non-compositional; (6) distance-based
 68 versus count/abundance-based; and (7) incorporating phylogenetic information into the analysis
 69 or not [10-12]. Here, we explore statistical methods for analysing microbiome data from
 70 designed experiments with a focus on dietary intervention trials. In contrast to observational
 71 studies, these are usually small in sample size but performed in (semi)-controlled environments
 72 and tailored to a specific research hypothesis. The studies often include multiple experimental
 73 factors, possibly with more than two levels, and it is therefore natural to turn to analysis of
 74 variance (ANOVA)-like methods. Notably, most published analytical tools in microbiome
 75 research are essentially univariate [6], which led us to the conclusion that comparison of
 76 alternative multivariate statistical tools is sorely missing. From a biologist's point of view, it is
 77 also important that the methods are easy to interpret, both at the multivariate (microbial
 78 community) and univariate (microbial taxa or operational taxonomic units, OTUs) levels (see
 79 Fig 1).

80



81

82 **Fig 1. A diagram of statistical methods used in the study.**

83

84 **Distance-based methods**

85 The distance-based methods are multivariate since multiple variables (microbial OTUs)
 86 are used to calculate pairwise distances between samples. Among distance-based methods,

87 permutational multivariate analysis of variance (PERMANOVA) is the most widely used and
88 more powerful than the analysis of similarities (ANOSIM) to detect changes in community
89 structure [13-15]. Both methods may be implemented with any dissimilarity metric. Among
90 abundance-based beta diversity indices, Bray-Curtis is the most common choice for count data
91 [16, 17]. The most widely applied phylogenetic beta diversity indices are UniFrac-type metrics
92 [17-19]. However, UniFrac is unsuitable as a distance metric for studies with a small sample
93 size, which is usually the case for dietary intervention trials [20, 21]. Both PERMANOVA and
94 ANOSIM test differences at the community level but do not provide any information at the
95 OTU level. Similarity percentage analysis (SIMPER) works at the univariate level by
96 computing the relative contribution of each analysed microbial taxon (i.e. OTU) to the overall
97 average Bray-Curtis dissimilarities by pairwise comparison of two or more groups [15]. To the
98 best of our knowledge, no such method exists for the other distance metrics.

99 Distance-based methods have their strengths and weaknesses that are important to
100 account for beforehand. ANOSIM cannot deal with multifactorial designs, and both ANOSIM
101 and PERMANOVA may have problems detecting differences unless they are present in taxa
102 with high variability [22]. Newer methods aimed at assigning more interpretable effect sizes
103 are under constant development [6, 23]. For example, the more flexible PERMANOVA-S is an
104 extension to existing distance-based methods that can adjust for covariates and simultaneously
105 incorporate multiple distance metrics [24]. However, these methods do not consider
106 covariance/correlation between microbial species, and they encounter significant power loss if
107 *all* microbial species are used for distance calculations [25].

108

109 **Abundance-based methods**

110 The abundance-based methods can be either univariate (analysing each OTU
111 individually) or multivariate (focusing on covariance structure between OTUs). There are two
112 main approaches to deal with the special nature of abundance data: (1) application of methods
113 that consider the distribution of count data or (2) compositional data analysis (CoDa) based on
114 log-ratio transformed count data [26, 27]. Statistical methods designed for high-throughput
115 sequencing data are ANOVA-like differential expression analysis (ALDEx2) [26], analysis of
116 composition of microbiomes (ANCOM) [28], edgeR [29] and DESeq/DESeq2 [30]. edgeR and
117 DESeq2 model count data directly using generalized linear models with the negative binomial
118 distribution and the logistic link, respectively, whereas ALDEx2 and ANCOM use the log-ratio
119 transformation prior to univariate assessment of statistical significance for individual OTUs.

120 DESeq2 and edgeR are based on the same modelling approach but differ in normalisation,
121 outlier handling, and other adjustable parameters; these methods had similar performance in
122 simulation studies [30]. Thus, we decided to include only one of the methods – DESeq2 –
123 because differences between DESeq2 and edgeR are at a different conceptual level rather than
124 the other methods discussed. ALDEx2 uses a Dirichlet-multinomial probability distribution to
125 estimate abundances from count data and calculates the false discovery rate (FDR) based on
126 Monte Carlo simulations (see Fig. 3 in [26] for details). In ANCOM, the compositional nature
127 of the data is considered by testing the log-ratio for all pairs of OTUs, and then counting the
128 number of tests where the log-ratio is significantly different from zero. This number (W-stat)
129 can be used to obtain a ranking of OTUs most likely to differ between the groups. The newly
130 published ANCOM-BC corrects the bias induced by differences in sampling fractions and
131 provides p-values and confidence intervals for the differential abundance of each OTU [31].
132 The FDR and power were shown to be similar for both ANCOM and ANCOM-BC, and
133 therefore we limit the present study to ANCOM.

134 The drawback of univariate methods is that they treat all taxa as independent variables
135 without considering the covariance between the OTUs. Such methods may fail to detect
136 community-level differences [32]. A classical generalisation of ANOVA to multiple variables
137 (MANOVA) cannot be used when the number of variables exceeds the number of samples, as
138 it suffers from the problem of a singularity of covariance matrices and assumptions that are not
139 fulfilled [33, 34]. Novel statistical ANOVA-like methods include fifty-fifty multivariate
140 analysis of variance (FFMANOVA) [35] and ANOVA-simultaneous component analysis
141 (ASCA) [33]. Both methods are based on principal component analysis (PCA), and they can
142 handle multiple collinear responses. In FFMANOVA, the multivariate effects are estimated by
143 a modified variant of classical MANOVA, and OTU-level p-values are obtained by rotation
144 tests which adjust the p-values for multiple testing [36]. For ASCA, the multivariate effects are
145 calculated from combined sums-of-squares from all OTUs, and significance is assessed by
146 permutation testing. ASCA also provides scores and loadings related to each experimental
147 factor, which can be visualised in the same way as for PCA to better understand covariance
148 patterns within the data. The contribution of each OTU can be quantified by the loadings or by
149 partial least squares discriminant analysis (PLS-DA) for pairwise comparisons. ASCA has
150 recently gained popularity in metabolomics [37-39], and both ASCA and FFMANOVA have
151 successfully been applied to microbiome data [40-44].

152 Linear discriminant analysis effect size (LEfSe) is a stepwise approach that combines
153 univariate analysis with multivariate discriminant analysis [45]. LEfSe has found wide
154 application in microbiome research due to its easy to-use-and-interpret visualization [46, 47],
155 but it is not adapted to experimental designs with several multilevel factors and is therefore not
156 considered in this study.

157

158 **Method comparison**

159 An overview of the different methods compared in this study is given in Fig 1 and
160 Table 1. ANOSIM and PERMANOVA provide results only at the community level, while
161 SIMPER, DESeq2, ANCOM and ALDEx2 report results for single OTUs. ASCA and
162 FFMANOVA are generic methods and the only methods that provide results at both the
163 community and OTU level.

164 The aim of the method comparison was to investigate how different strategies for
165 statistical modelling affect biological inference. At the community level, methods were
166 compared with respect to effect sizes (expressed as percentage of explained variance) and
167 corresponding p-values. At the OTU level, comparison of methods is complex because some
168 methods provide results for an omnibus test of differences between factor levels (FFMANOVA
169 and ANCOM), whereas the other methods provide ranking for specific pairwise comparisons
170 (ASCA and PLS-DA, SIMPER) or contrasts/model coefficients (ALDEx2 and DESeq2). Even
171 so, a biologist will make inferences based on the output provided by the chosen method, and in
172 this context, it is relevant to compare the ranking statistics although the tests are not the same.
173 In our study, the ranking of OTUs was compared by Spearman's rank correlation and by
174 investigation of scatterplots between the different ranking metrics. For the simulated data,
175 where we know which OTUs are differentially abundant, True Positive Rate (TPR) and True
176 Negative Rate (TNR) were also evaluated.

177 We focused exclusively on designed experiments, which are usually smaller in sample
178 size and are more controlled in contrast to observational studies. We used five published data
179 sets as a basis for the comparisons (S1 Table). The following criteria for studies to be included
180 were considered: (1) at least two-factorial experimental design with a minimum of two-factor
181 levels; (2) either human or animal gut microbiome surveys; and (3) a taxonomic assignment at
182 the OTU level reported. Diet is the main factor of interest in all five studies, and we restricted
183 our comparisons to this factor.

184 The simulated data were based on data set 1 (S1 Table) using the same study design and
185 OTU counts as a starting point. Four different scenarios were simulated to investigate how the
186 methods perform in situations with varying effect sizes and different numbers of differentially
187 abundant OTUs. In all scenarios, one of the diet levels was manipulated to be significantly
188 different from the others, and there was no effect of the second experimental factor (dose). See
189 Methods section for further details.
190

191 **Table 1. An overview of statistical methods and their properties.**

Method	Method name	Number of experimental factors allowed	Parametric	Multivariate	Univariate	Provides output at community level	Statistics for ranking OTUs	Reference
ALDEx2	ANOVA-like differential expression tool for high-throughput sequencing data	any	yes	no	yes	no	p-values or effect sizes	[26]
ANCOM	Analysis of composition of microbiomes	main factor + covariates	yes	no	yes	no	W-stat for the main variable	[28]
ANOSIM	Analysis of similarities	one	no	yes	no	yes	no	[15]
ASCA	ANOVA-simultaneous component analysis	any	yes	yes	no	yes	loadings or PLS-DA regression coefficients	[33]
DESeq2	Differential gene expression analysis based on the negative binomial distribution	any	yes (GLM)	no	yes	no	p-values or effect sizes (coefficients)	[30]
FFMANOVA	Fifty-fifty multivariate ANOVA	any	yes	yes	yes (rotation tests)	yes	p-values	[35]
PERMANOVA	Permutational multivariate analysis of variance	any	no	yes	no	yes	no	[14]
SIMPER	Similarity percentage	two-group comparison	no	yes	no	no	permutation p-values	[15]

192

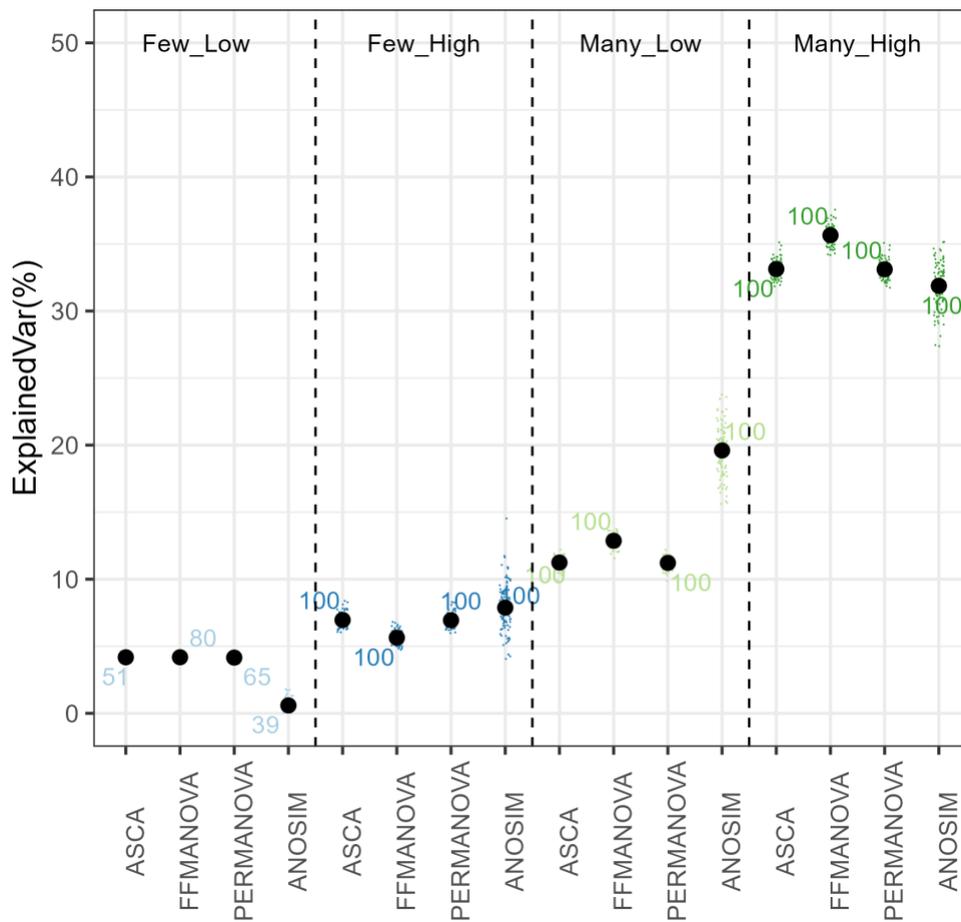
193 ANOVA – analysis of variance; GLM – generalized linear model; OTU – operational taxonomic unit; PLS-DA – partial least squares discriminant analysis.

194 **Results**

195 **Community level**

196 Four of the methods can be used to test the association between diet and the overall
197 microbiome composition, the distance-based ANOSIM and PERMANOVA, and abundance-
198 based ASCA and FFMANOVA. The results for each of the four simulated scenarios are shown
199 in Fig 2, and the results across real data sets are summarised in Table 2.

200



201

202 **Fig 2. Explained variance for simulated data and the relative number of simulations**
203 **where the simulated effect was detected.**

204 **Table 2. Community-level method comparison across five experimental data sets.**

Factor/ predictor	FFMANOVA ^{clr}		ASCA ^{clr}		PERMANOVA ^{clr}		ANOSIM ^{3, clr}		Factor/ predictor
	Effect size (explained variance), %	p-value ¹	Effect size (explained variance), %	p-value ²	Effect size (explained variance), %	p-value ²	Effect size (explained variance), %	p-value ²	
Moen et al., 2016 (data set 1) Model: OTU ~ fiber*dose							Moen et al., 2016 (data set 1) Model: OTU ~ fiber_dose		
fiber	34.31	< 0.001	37.39	< 0.001	37.26	0.001	82.19	0.001	fiber
dose	3.96	< 0.001	4.14	< 0.001	4.12	0.001			dose
fiber:dose	5.85	< 0.001	5.99	< 0.001	5.96	0.002			fiber:dose
residuals	55.75		52.69		52.67		17.81		residuals
Lai et al., 2016 (data set 2) Model: OTU ~ diet*exercise							Lai et al., 2016 (data set 2) Model: OTU ~ diet_exercise		
diet	27.31	< 0.001	30.50	< 0.001	30.85	0.001	99.11	0.001	diet
exercise	12.18	< 0.001	14.08	< 0.001	13.99	0.001			exercise
diet:exercise	7.93	< 0.001	7.51	< 0.001	7.47	0.002			diet:exercise
residuals	52.31		47.74		47.70		0.89		residuals
Le Sciellour et al., 2018 (data set 3) Model: OTU ~ diet*period + subject							Le Sciellour et al., 2018 (data set 3) Model: OTU ~ diet_period		
diet	1.78	< 0.001	1.99	< 0.001	2.14	0.001	13.33	0.001	diet
period	3.29	< 0.001	3.45	< 0.001	3.51	0.001			period
diet:period	1.32	< 0.001	1.41	0.006	1.36	0.009			diet:period
subject	26.53	< 0.001	26.42	0.073	27.28	0.046	86.67		subject
residuals	66.13		65.71		65.71				residuals
Wang et al., 2016 (data set 4) Model: OTU ~ diet + time + diet:time + subject							Wang et al., 2016 (data set 4) Model: OTU ~ diet_time		
diet	2.49	0.067	2.11	0.038	2.27	0.036	4.32	0.143	diet
time	2.00	0.007	1.96	0.082	1.77	0.246			time

diet:time	5.15	0.824	4.41	0.777	4.45	0.701			diet:time
subject	50.44	< 0.001	54.23	< 0.001	64.30	0.001	95.68		subject
residuals	31.23		27.25		27.21				residuals
Birkeland et al., 2020 (data set 5) Model: OTU ~ treatment:day + subject									Birkeland et al., 2020 (data set 5) Model: OTU ~ treatment_day
treatment:day	1.75	< 0.001	1.27	0.107	1.27	0.132	-0.03	0.998	treatment:day
subject	69.38	< 0.001	73.85	0	73.85	0.001			subject
residuals	28.88		24.87		24.87				residuals

205

206 Distance-based ANOSIM and PERMANOVA and abundance-based ASCA and FFMANOVA were compared with respect to effect sizes
 207 (expressed as percentage of explained variance) and corresponding p-values.

208 ¹based on the 50-50 F-test, 999 permutations

209 ²based on 999 permutations

210 ³based on a combined factor with no interaction in the model (limitation of ANOSIM)

211 ^{clr}centred log-ratio transformed data as input

212

213 **Simulated data.** As expected, the multivariate effect size (explained variance) is lowest
214 (around 5%) for the “*Few-Low*” scenario and highest (30-35%) for the “*Many-High*” scenario.
215 The multivariate effect was significant in 100% of the simulations in three scenarios with the
216 highest effect size. For the “*Few-Low*” scenario, FFMANOVA performed best by detecting
217 the effect in 80% of the data sets. For the “*Many*” simulations explained variance was slightly
218 higher with FFMANOVA than with ASCA and PERMANOVA, whereas for the “*Few-High*”
219 simulations the opposite trend was observed. ANOSIM was less consistent than the other
220 methods, with higher within-scenario variation and higher differences in effect size between
221 the two “*Low*” scenarios.

222 **Data set 1.** The effect of dietary fibres inulin (IN), cellulose (CE) or brewers spent grain
223 (BSG) on the overall caecal microbiota composition in mice from a study by Moen *et al.* [42]
224 accounted for 34-37% of the explained variance according to the FFMANOVA, ASCA and
225 PERMANOVA. In general, the three methods produced similar results, with slightly smaller p-
226 values by FFMANOVA and ASCA.

227 **Data set 2.** Lai *et al.* [48] investigated the effect of diet (the main variable), exercise
228 and their interaction on the overall faecal microbiota in sedentary and exercised mice fed high
229 fat or normal fat diet (four groups in total). Similarly, all tested methods, except ANOSIM,
230 produced congruent results (effect of diet 27-31%), with smaller p-values by FFMANOVA and
231 ASCA.

232 **Data set 3.** In a longitudinal study by Le Sciellour *et al.* [49] the authors tested the effect
233 of dietary fibre content on faecal microbiota in growing-finishing pigs fed alternately a low-
234 fibre and a high-fibre diet during four successive 3-week periods. In this survey, the effect of
235 diet was small (2%) compared to the effect of diet in data sets 1 and 2. Similarly, FFMANOVA
236 and ASCA reported slightly lower p-values than PERMANOVA, but all three methods agreed
237 on the effect of diet.

238 **Data set 4.** In a longitudinal study by Wang *et al.* [50] the objectives were to determine
239 the impact of beta glucan on the composition of faecal microbiota in mildly
240 hypercholesterolemic individuals. The individuals received for 5 weeks either a treatment
241 containing 3 g high molecular weight (HMW), 3 g low molecular weight (LMW), 5 g LMW
242 barley beta glucan or wheat and rice (control group) [50]. The effect of diet accounted for ~2%
243 of the explained variance reported by FFMANOVA, ASCA and PERMANOVA. Diet was
244 significant on a 5% level for PERMANOVA and ASCA ($p = 0.036$ and $p = 0.038$, respectively)
245 and on a 10% level for FFMANOVA ($p = 0.067$). However, different conclusions were drawn

246 with respect to time where significant result at the community level was obtained only for
247 FFMANOVA ($p = 0.007$).

248 **Data set 5.** Birkeland *et al.* [44] assessed the effect of prebiotic fibres or a control
249 supplement on faecal microbiota composition in human subjects with type two diabetes. The
250 interaction effect of treatment and day accounted for 1-2% of the explained variance according
251 to the FFMANOVA, ASCA and PERMANOVA. Significant results at the community level
252 were obtained only for FFMANOVA ($p < 0.001$).

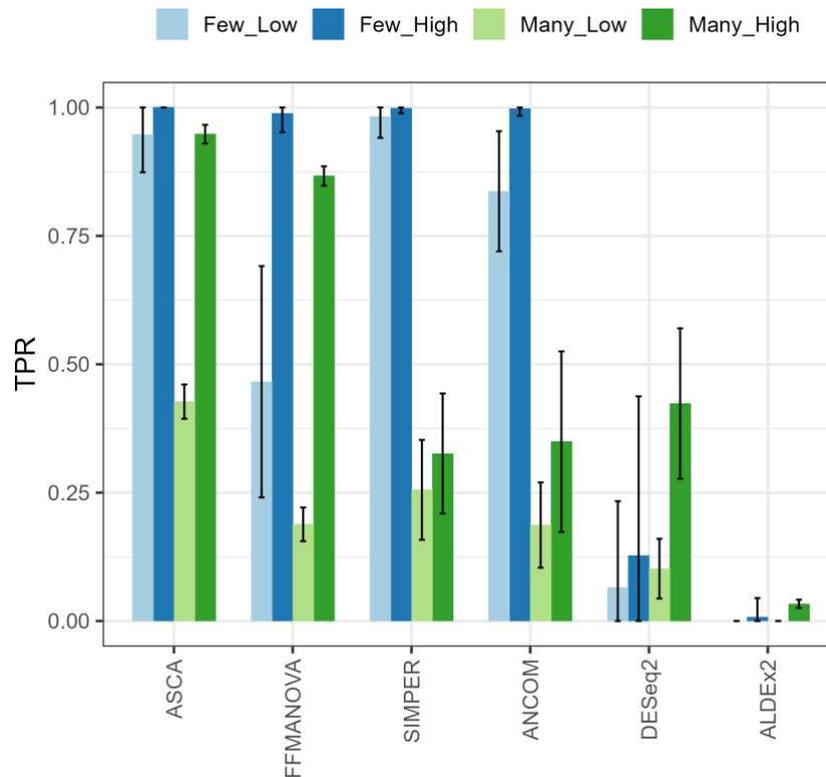
253 To summarise, PERMANOVA, FFMANOVA and ASCA gave almost identical results
254 regarding effect sizes and statistical significance across studies (Table 2). They all revealed that
255 there was a considerable difference in effect sizes of the main factor of interest (diet) between
256 animal (2-37%) and human (1-2%) dietary interventions, with effect sizes being very small in
257 human studies. In addition, three of the studies had crossover designs allowing for estimation
258 of interindividual variation. This variation was considerably higher for trials involving human
259 subjects (54-74%) compared to the animal study (26%). ANOSIM provided the most different
260 results and was not able to reveal the same biological insights since the multifactorial nature of
261 the studies cannot be taken into consideration by this approach.

262

263 **OTU level**

264 Six of the methods, namely SIMPER, ASCA, FFMANOVA, ANCOM, ALDEx2 and
265 DESeq2 can be used to make biological inferences for individual OTUs. The methods give
266 different outputs which can be used to identify differentially abundant OTUs and/or rank the
267 OTUs according to effect sizes (see Methods for details).

268 **Simulated data.** The True Positive Rate (TPR) for the four simulated scenarios is shown
269 in Fig 3. The True Negative Rate (TNR) was close to 100% for all methods and is therefore not
270 shown. ASCA provided the overall best results in terms of the TPR in all four scenarios.
271 FFMANOVA, SIMPER and ANCOM were all highly sensitive in the scenario with few
272 significant OTUs and a high effect size ("*Few-High*"). FFMANOVA was also highly sensitive
273 in the scenario with many significant OTUs and a high effect size ("*Many-High*"); SIMPER
274 and ANCOM performed best in the scenario with few significant OTUs with a lower effect size
275 ("*Few-Low*"). Both ALDEx2 and DESeq2 detected very few OTUs in any of the scenarios and
276 therefore had very low TPR.



277

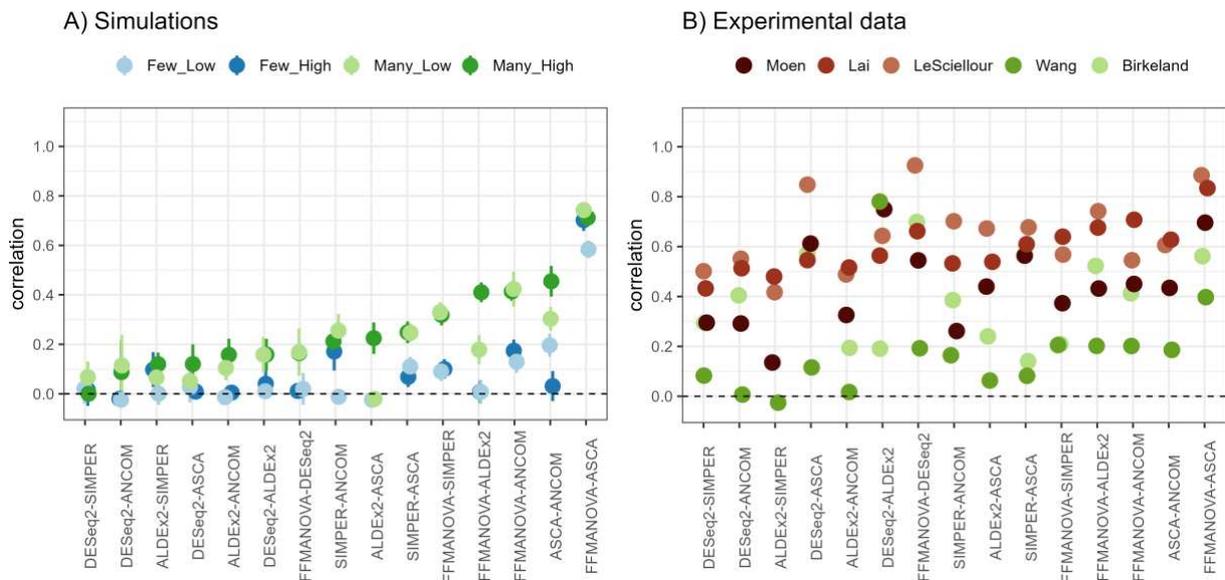
278 **Fig 3. Sensitivity (True Positive Rate) for the four scenarios in the simulation study.**

279

280 **Experimental data.** Summary tables and scatterplots comparing ranking for different
 281 methods on the experimental data sets are given as S1 File and S1 Fig, respectively. The number
 282 of significant OTUs detected by the methods is summarised in S2 Table. As expected, many
 283 significant OTUs were discovered in the studies with large multivariate effect sizes (data sets
 284 Moen and Lai), whereas few OTUs were found in the studies with low multivariate effect sizes.
 285 The highest number of significant OTUs was identified by FFMANOVA, with almost twice
 286 the numbers detected by ALDEx2. ANCOM differed considerably between the study designs.
 287 ASCA recovered fewer OTUs than the other methods for the Moen data set, but more OTUs
 288 than the other methods for the other data sets.

289 **Correlation between the methods.** Agreement between the methods was investigated
 290 by calculating Spearman's rank correlation between all pairs of output metrics (Fig 4). In the
 291 simulated data, FFMANOVA and ASCA had higher agreement than any other pair of methods,
 292 with correlations ranging from 0.6 to 0.75. In addition, the correlations were generally higher
 293 for scenarios with many differentially abundant OTUs for all methods. The results from the
 294 experimental data sets also showed that FFMANOVA and ASCA had highest agreement, with
 295 correlations ranging from 0.4 to 0.9. However, correlations varied considerably between the
 296 data sets. Highest agreements were observed for the animal studies, which are more controlled,

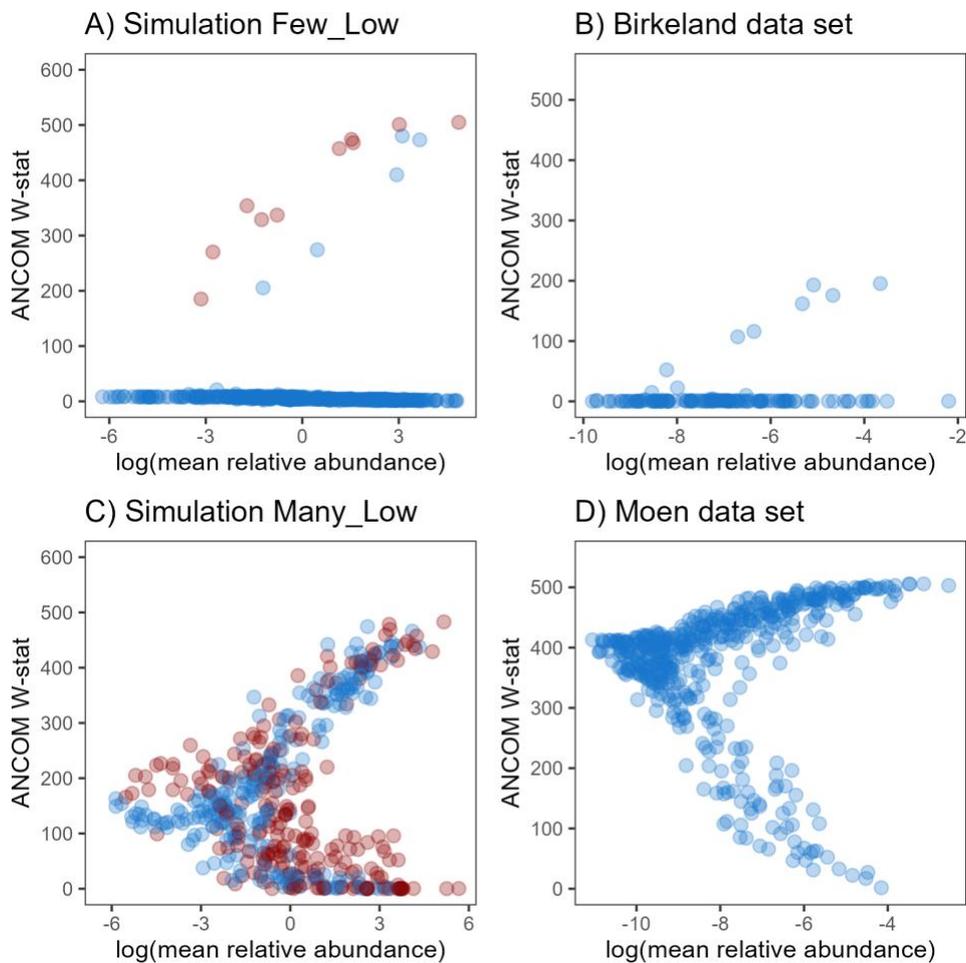
297 and where interindividual variation is smaller compared to studies involving human subjects.
 298 The Moen and Wang data sets had a lower correlation between the methods than the other data
 299 sets for many comparisons, which could be explained by the fact that there are more than two
 300 levels of the diet factor, and omnibus tests will therefore not completely agree with the pairwise
 301 comparisons.
 302



303
 304 **Fig 4. Spearman's correlation (Y-axis) calculated for pairwise comparison of statistical**
 305 **methods (X-axis) for (A) simulated data and (B) five experimental data sets.** Each point
 306 represents Spearman's rank correlation coefficient between OTU ranking metrics from the two
 307 methods compared.

308
 309 **Impact of OTU abundance.** Gut microbiome data are typically represented by a few
 310 dominant OTUs (relative abundance >1%) and a majority of low-abundant OTUs. One is
 311 therefore interested in ranking OTUs independent of their average abundance to detect
 312 biologically relevant changes in low-abundant OTUs. For all methods, except SIMPER,
 313 correlations in the range 0.1-0.4 were observed between the ranking statistics and the (log)
 314 mean relative abundance. However, it was not consistent between the data sets which method
 315 gave higher correlations. Correlations between the relative abundance and the ranking statistics
 316 were highest for the Moen data set and lowest for the Birkeland data set (S2 Table). ANCOM
 317 differed from the other methods as the ranking of OTUs was highly dependent on the abundance
 318 (Fig 5 and S1 Fig). In particular, the results indicated that highly abundant OTUs had either
 319 very high or very low W-stat, while low-abundant OTUs always had medium-to-low W-stat.

320 To the best of our knowledge, this finding has not been reported before and shows that ANCOM
321 is not able to identify changes in low-abundant OTUs. In the simulation study, clear differences
322 between the scenarios with “*Few*” or “*Many*” differentially abundant OTUs can be observed
323 (Fig 5A and 5C), and similar patterns are reflected for the experimental data (Fig 5B and 5D).
324 It has been shown that power of ANCOM drops when the number of differentially abundant
325 OTUs exceeds 25% [28]. In our “*Many*” simulations, 50% of the OTUs are differentially
326 abundant, which justifies why ANCOM performed poorly for these scenarios. The dependency
327 on the relative abundance can be more problematic also for the “*Many*” scenarios (Fig 5C) and
328 data sets with many differentially abundant OTUs as, for instance, the Moen data set (Fig 5D).
329



330

331

332 **Fig 5. Mean relative abundance (log-scale) plotted versus ANCOM W-stat.** (A) The “*Few-*
333 *Low*” simulation scenario and (B) Birkeland data set, (C) “*Many-Low*” simulation scenario and
334 (D) Moen data set.

335

336 Discussion

337 Diet is considered to be an important driver of microbiome variation [51]. However, in
338 observational population-based studies, diet consistently accounts for only a small proportion
339 of microbiome variation, and this is partly due to large interindividual differences in
340 microbiome composition, small sample sizes and limitations in study designs such as
341 potentially insufficient washout periods in crossover studies [51-55]. In general, higher
342 interindividual variation is observed in gut microbiome of human subjects compared to animal
343 species [56]. This was confirmed by our results, and it can partly explain the lower diet related
344 effect on the gut microbiomes in the two human studies. Variation in how much a diet can
345 influence the microbiome is also dependent on the nutritional differences of the compared diets.
346 Nevertheless, use of animal models to study a causal role of the gut microbiome in health and
347 disease is an established practice although animal models lack the specific interactions present
348 in the complex system of a human organism [57].

349 Univariate and multivariate analyses provide information at different levels. Biologists
350 will often find that the outputs of univariate analyses are easier to interpret compared to those
351 generated by multivariate analyses, though assumptions are similar for both method types [10].
352 In general, multivariate methods provide a more holistic overview of differences between
353 samples and account for correlations and interactions between the variables, whereas univariate
354 methods are well suited to point out the differences for specific microbial groups. Therefore,
355 the two levels of analysis provide complementary information, and it is generally of biological
356 interest to report differences at both levels.

357 FFMANOVA and ASCA consider the covariance between all OTUs. It could therefore
358 be expected that these methods would be better at detecting scenarios with many differentially
359 abundant OTUs due to the *consistency at large* phenomenon, i.e. that many OTUs carry the
360 same information and such methods collectively are able to detect small effects that are not
361 significant at the univariate level [58, 59]. At the community level, these methods performed
362 similarly to the distance-based methods. At the OTU level, they had considerably higher
363 sensitivity to identify true positive OTUs than the other methods in the “Many” OTU scenarios.

364 FFMANOVA and ASCA depend on the relative scaling of the OTUs, while for the
365 distance-based methods it depends on the chosen distance measure. It is a common practice in
366 many areas to scale all variables to equal variance thus giving them an equal weight in the
367 model, but other options are also possible, see, for instance, van den Berg *et al.* [60]. The clr-
368 transformation puts the variables at comparable levels, and the need for scaling is less obvious.

369 However, the highly abundant OTUs might still have slightly higher variance and scaling
370 should be considered depending on the data characteristics and the biological interpretation.
371 With scaling, all OTUs have the same contribution in the analysis, whereas without scaling the
372 more abundant OTUs will dominate the analyses and the inferences will be related to the more
373 abundant OTUs.

374 It is a known fact that microbial sequence data are zero-inflated, and rare OTUs should
375 be removed prior to downstream statistical analyses. We have observed that the threshold for
376 filtering out OTUs can significantly affect the results both at the community and OTU levels.
377 This can be exemplified in the human Birkeland data set, where stricter OTU-filtering
378 performed in the present study resulted in no significant treatment effect by ASCA, in contrast
379 to the original publication [44].

380 The tools tested in the present study vary in flexibility. SIMPER allows only pairwise
381 comparisons and ANOSIM provides multigroup analysis, but it is not suited for multifactorial
382 study designs. The other methods can employ more complex models with multiple factors with
383 varying number of levels, and corresponding interactions. In experiments with repeated
384 measurements, the subject effect is often included as a random factor. Neither of the methods
385 discussed here can do this, hence the subject effect was included as a fixed factor. These aspects
386 should be considered when selecting methods because different study designs might require
387 different types of statistical models and tests. Newer developments of ASCA, namely ASCA+
388 [61] and LiMM-PCA [62], increase flexibility when there are unbalanced designs or random
389 (i.e. subject) effects, respectively. Even so, the longitudinal modelling of large microbiome data
390 sets in combination with multiple covariates is only starting to emerge [6].

391 Although Spearman's rank correlation indicated good agreement for the animal studies
392 (Fig 4), little overlap between lists of "significant" results could be detected (results not shown).
393 There can be several reasons for this. One important aspect is that different criteria must be
394 used to define "significance" or generally "importance". FFMANOVA, SIMPER, DESeq2 and
395 ALDEx2 provide p-values, while more heuristic tools must be applied with ASCA and
396 ANCOM. For some methods, such as FFMANOVA and ANCOM, the ranking is related to *all*
397 levels of the experimental factors, whereas the other methods use only pairwise comparisons
398 (SIMPER) or comparison against a selected control/reference level (ALDEx2 and DESeq2).
399 ASCA provides a test related to all levels at the community level, while the PLS-DA step relates
400 to pairwise comparisons at the OTU level. This will introduce a bias for data sets with multilevel
401 factors where more than one level is different from each other. Nevertheless, we chose to

402 compare these rankings as this is the output that is available to the user. In the simulations only
403 one level of one factor was designed to have an effect, hence the “omnibus” and the “specific”
404 tests are more directly comparable. In experiments with multilevel factors, additional
405 information can be obtained by looking at the clr-difference between the groups of interest in
406 addition to the ranking statistics.

407 Moreover, differences in sample collection, sample preparation and sequencing
408 contribute to additional variability, which, in turn, affects the validity of the results [63] and
409 complicates comparisons across studies with similar interventions.

410 Past benchmarking studies [9, 64, 65] have reported varying results from different tools,
411 which was also confirmed in our study. Currently, there is no consensus for the best existing
412 tool for detecting differentially abundant microbial taxa, and there is no reason to believe that
413 one single method is best in all cases. Based on our simulations, the generic multivariate tools,
414 ASCA and FFMANOVA, performed best. We anticipate that these results will inform future
415 studies with more complex settings where more than one factor has an effect or interactions
416 between experimental factors are included in the model.

417 In addition to performance, ease of use is an important aspect when selecting the
418 appropriate tool. FFMANOVA and ASCA are based on standard statistical tools, namely PCA
419 and ANOVA. Some of the tools designed for microbiome high-throughput sequencing data, on
420 the other hand, can be difficult for non-statisticians, and it can be questioned if the users are
421 able to interpret all parameters correctly, even if the methods are supported by comprehensive
422 documentation. It is always good scientific practice to compare and report outputs from several
423 methods. There are no standards on how to report multiple modelling results, and there is a high
424 risk of “fishing for significance” when several methods are applied [66]. Before designing the
425 experiment, researchers should be aware of the different properties of the statistical methods
426 and consider whether it is most important not to miss out on any possible findings or to obtain
427 robust results. In the latter case, OTUs should be reported as differentially abundant only if they
428 were flagged as “significant” by several methods [66].

429

430 **Conclusion**

431 In the present study, we compared the performance of several multivariate ANOVA-
432 like statistical methods taking four simulated scenarios and five real dietary intervention
433 microbiome data sets as examples. At the community level, all the different methods came to

434 similar conclusions; at the OTU level, the agreement between the methods considerably varied.
435 ANCOM provided output metrics that were dependent on the average abundance, making it
436 impossible to detect differences in low-abundant OTUs. At the OTU level, the ranking of OTUs
437 obtained with different methods correlated better for animal studies than for human studies,
438 possibly due to lower interindividual variation in animal studies. Based on the simulation results
439 we advise applying FFMANOVA and ASCA for overall and pairwise comparisons of
440 microbiome data, respectively, also because these methods provide output at both the
441 community and OTU levels, can handle several design factors, as well as other data types
442 common in microbiome research.

443

444 **Methods**

445 **Experimental design and data characteristics.** The dietary intervention data sets were
446 selected based on the study design, with a minimum of two independent variables (for example,
447 diet and dose; see S1 Table for details). Prior to the statistical analyses, the data were filtered
448 to keep the OTUs that were present: (1) with relative abundance more than 0.005% in an
449 individual, and (2) in at least 50% of the individuals and in one of the groups. For each study
450 in total 507, 561, 560, 397 and 216 OTUs passed this filter and were subsequently used in
451 downstream analyses (S1 Table and S2 File). All statistical analyses were performed in R
452 version 4.1.0 [67] unless otherwise specified and were run in 999 permutations.

453 **Simulated data** were generated using the R package metaSPARSim [68] developed for
454 simulating 16S rDNA data. The simulation was a two-step process: (1) modelling the
455 abundance (expected counts for each experimental group) using a Gamma distribution; (2)
456 modelling the within-group variability using a Multivariate Hypergeometric (MHG)
457 distribution taking the output from step 1 as input parameters for the distribution. In addition,
458 metaSPARSim contains functions for parameter estimation from observed data. We used data
459 set 1 [42] (S1 Table) to estimate real starting parameters for the simulations. Raw counts were
460 generated for the same number of OTUs as in the experimental data set, and subsequently pre-
461 processed and analysed in the same manner as for the experimental data (described below).
462 Four different scenarios were simulated, with a varying number of the differentially abundant
463 OTUs (“*Few*” versus “*Many*”) and the effect sizes (“*Low*” versus “*High*”). 100 data sets were
464 generated and analysed for each of the four scenarios:

465 1) “*Few – Low*”: 10 randomly selected OTUs were assigned random log₂ fold changes from
466 a uniform distribution with boundaries [3,4].

- 467 2) *“Few – High”*: 10 randomly selected OTUs were assigned random log₂ fold changes from
468 a uniform distribution with boundaries [8,9].
- 469 3) *“Many – Low”*: 254 randomly selected OTUs were assigned random log₂ fold changes
470 from a uniform distribution with boundaries [3,4].
- 471 4) *“Many – High”*: 254 randomly selected OTUs were assigned random log₂ fold changes
472 from a uniform distribution with boundaries [8,9].

473 **Zero-value replacements** were done prior to clr-transformation [27, 69] by applying
474 function *cMultRepl* with a setting *method* = ‘CZM’ in the R package *zCompositions* [70]. Zero
475 replacement is an ongoing and yet unsolved statistical problem in microbiome research, and
476 newer methods are constantly being developed and applied to both simulated and experimental
477 data sets [71-74].

478 **ANOSIM** and **PERMANOVA** were run using the functions *anosim* and *adonis* in the R
479 package *vegan* [75], with a setting *method* = *Euclidean* and the clr-transformed data as input.

480 **SIMPER** was run on filtered relative abundance data using function *simper* from the R
481 package *vegan* [75]. The p-values for each OTU between selected pairwise comparisons of diet
482 levels were used as a ranking metric; the p-value represents the probability of getting a larger
483 contribution to the Bray-Curtis dissimilarity in a random permutation of the group factor.

484 **FFMANOVA** was performed by using the function *ffmanova* implemented in the R
485 package *ffmanova* [35] on the clr-transformed and standardised data. Raw p-values were used
486 as a ranking metric.

487 **ASCA** was run on the clr-transformed and centred data using an in-house implementation
488 in MATLAB (R2018b, The MathWorks Inc.); p-values at the community level were calculated
489 by permutation tests (*n* = 999). PLS-DA models [76] with the ASCA diet effect matrix,
490 residuals as a predictor and factor levels as a response were used for pairwise comparison of
491 diet levels. Variable Importance in Prediction (VIP) values [76] were used to identify significant
492 OTUs, and the VIP threshold was set using the Uninformative Variable Elimination (UVE)
493 method [77]. The UVE procedure was repeated 100 times, and OTUs that were above the
494 threshold in >95% of the repetitions were defined as significant. For study designs with two
495 diet levels, the ASCA loadings were used to rank OTUs, while PLS-DA regression coefficients
496 were used for pairwise comparisons of multilevel diet factors.

497 **ANCOM** was run by using the ANCOM 2.0 source code implemented in R [28] using
498 filtered relative abundance data as input. The W-stat was used to rank OTUs indicating the

499 number of significantly different pairwise log-ratios while adjusting for FDR by applying a
500 Benjamini-Hochberg correction at a 0.05 level of significance.

501 **ALDEx2** was run using the functions *aldex.clr* and *aldex.glm* from the R package ALDEx2
502 v.1.18.0 [26]. We used raw counts as input and p-values for selected factor level contrasts to
503 rank OTUs.

504 **DESeq2** was run using the functions *DESeqDataSetFromMatrix* (to generate object),
505 *DESeq* (for analysis) and *results* (to extract results) from the Bioconductor package version
506 1.32.0 [30]. We used raw counts as input and raw p-values for selected level contrasts to rank
507 OTUs (default settings).

508 **Factor level comparisons.** For study designs with more than two diet groups, only one
509 pairwise comparison was analysed for the methods based on the pairwise group comparison,
510 namely ALDEx2, ASCA and SIMPER. The following pairs with the most contrasting outcomes
511 were compared: (1) BSG group vs. IN group [42]; (2) control group vs. 3 g HMW [50] and (3)
512 Placebo 6 weeks group vs. Fibre 6 weeks group [44].

513 **Differentially abundant OTUs** were identified by setting thresholds on the ranking metrics
514 for each method. For the methods providing p-values, the threshold was set to 0.01. For
515 ANCOM, the 60th percentile of the empirical distribution of the W-stat was used as a threshold.
516 For ASCA, OTUs selected in 95 out of 100 UVE-runs were identified as differentially
517 abundant.

518 **True Positive Rate (TPR) and True Negative Rate (TNR)** were calculated for the
519 simulated data. TPR (also called *Power* or *Sensitivity*) is calculated as TP/P , where TP is the
520 number of true differentially abundant OTUs identified by a statistical method and P is the
521 number of differentially abundant OTUs defined in the simulation setup. The TNR (also called
522 *Specificity*) is calculated as TN/N , where TN is the number of true non-differentially abundant
523 OTUs identified by a statistical method and N is the corresponding number defined in the
524 simulation setup.

525

526 **List of abbreviations**

527 ANOVA: Analysis of variance; ALDEx2: ANOVA-like differential expression tool for
528 high-throughput sequencing data; ANCOM: Analysis of composition of microbiomes;
529 ANOSIM: Analysis of similarities; ASCA: ANOVA-simultaneous component analysis;
530 DESeq2: Differential gene expression analysis based on the negative binomial distribution;

531 FDR: False discovery rate; FFMANOVA: Fifty-fifty multivariate ANOVA; OTU: Operational
532 taxonomic unit; PERMANOVA: Permutational multivariate analysis of variance; PLS-DA:
533 Partial least squares discriminant analysis; SIMPER: Similarity percentage.

534

535 **Acknowledgements**

536 We thank all the authors who were involved in data generation used in the present study.
537 All data sets are properly cited and referred to in S1 Table.

538

539 **Funding**

540 This work was funded by Nofima – Norwegian Institute of Food, Fisheries and
541 Aquaculture Research and Foundation for Research Levy on Agricultural Products (Projects
542 RCN 262306 and 262308). The funding sponsors had no role in the design of the study, the
543 collection, analyses and interpretation of data, in the writing of the manuscript or in the decision
544 to distribute the results.

545

546 **Author contributions statement**

547 IM, IB and IR conceived the project. MK, IM, IB and IR designed the project. MK, IM
548 and IB analysed the data. MK, IM, IB and IR interpreted the results. MK and IB generated the
549 figures and tables. MK was a major contributor in writing the manuscript, with help from IM
550 and IB. All authors reviewed and edited the manuscript and approved the final version.

551

552 **Competing interests' statement**

553 The authors declare no competing interests.

554

555 **Ethics approval and consent to participate**

556 Not applicable. The study is based on already published data sets listed in S1 Table.

557

558 **Availability of data and materials**

559 All data analysed in the present study are included in Supplementary Information.
560 Custom R scripts to run the experimental and simulated data analyses are available as S3-S5
561 Files.

562

563 **References**

564 1. Li H. Microbiome, metagenomics, and high-dimensional compositional data analysis.
565 *Annu Rev Stat Appl.* 2015;2:73-94. doi: [https://doi.org/10.1146/annurev-statistics-010814-](https://doi.org/10.1146/annurev-statistics-010814-020351)
566 [020351](https://doi.org/10.1146/annurev-statistics-010814-020351).

567 2. Blanco-Míguez A, Fdez-Riverola F, Sánchez B, Lourenço A. Resources and tools for
568 the high-throughput, multi-omic study of intestinal microbiota. *Brief Bioinform.*
569 2019;20(3):1032-56. doi: <https://doi.org/10.1093/bib/bbx156>.

570 3. Knight R, Vrbanac A, Taylor BC, Aksenov A, Callewaert C, Debelius J, et al. Best
571 practices for analysing microbiomes. *Nat Rev Microbiol.* 2018;16(7):410-22. doi:
572 <https://doi.org/10.1038/s41579-018-0029-9>.

573 4. Gloor GB, Reid G. Compositional analysis: a valid approach to analyze microbiome
574 high-throughput sequencing data. *Can J Microbiol.* 2016;62(8):692-703. doi:
575 <https://doi.org/10.1139/cjm-2015-0821>.

576 5. Lê Cao K-A, Costello M-E, Lakis VA, Bartolo F, Chua X-Y, Brazeilles R, et al.
577 MixMC: a multivariate statistical framework to gain insight into microbial communities. *PLoS*
578 *One.* 2016;11(8):e0160169. doi: <https://doi.org/10.1371/journal.pone.0160169>.

579 6. Mallick H, Ma S, Franzosa EA, Vatanen T, Morgan XC, Huttenhower C. Experimental
580 design and quantitative analysis of microbial community multiomics. *Genome Biol.*
581 2017;18(1):228. doi: <https://doi.org/10.1186/s13059-017-1359-z>.

582 7. Waldron L. Data and statistical methods to analyze the human microbiome. *mSystems.*
583 2018;3(2):e00194-17. doi: <https://doi.org/10.1128/mSystems.00194-17>.

584 8. Schloss PD. Identifying and overcoming threats to reproducibility, replicability,
585 robustness, and generalizability in microbiome research. *mBio.* 2018;9(3):e00525-18. doi:
586 <https://doi.org/10.1128/mBio.00525-18>.

587 9. Weiss S, Xu ZZ, Peddada S, Amir A, Bittinger K, Gonzalez A, et al. Normalization and
588 microbial differential abundance strategies depend upon data characteristics. *Microbiome.*
589 2017;5(1):27. doi: <https://doi.org/10.1186/s40168-017-0237-y>.

- 590 10. Paliy O, Shankar V. Application of multivariate statistical techniques in microbial
591 ecology. *Mol Ecol*. 2016;25(5):1032-57. doi: <https://doi.org/10.1111/mec.13536>.
- 592 11. Morton JT, Sanders J, Quinn RA, McDonald D, Gonzalez A, Vázquez-Baeza Y, et al.
593 Balance trees reveal microbial niche differentiation. *mSystems*. 2017;2(1):e00162-16. doi:
594 <https://doi.org/10.1128/mSystems.00162-16>.
- 595 12. Wang S, Cai TT, Li H. Hypothesis testing for phylogenetic composition: a minimum-
596 cost flow perspective. *Biometrika*. 2021;108(1):17-36. doi:
597 <https://doi.org/10.1093/biomet/asaa061>.
- 598 13. Anderson MJ. A new method for non-parametric multivariate analysis of variance.
599 *Austral Ecol*. 2001;26(1):32-46. doi: <https://doi.org/10.1111/j.1442-9993.2001.01070.pp.x>.
- 600 14. Anderson MJ, Walsh DC. PERMANOVA, ANOSIM, and the Mantel test in the face of
601 heterogeneous dispersions: what null hypothesis are you testing? *Ecol Monogr*.
602 2013;83(4):557-74. doi: <https://doi.org/10.1890/12-2010.1>.
- 603 15. Clarke KR. Non-parametric multivariate analyses of changes in community structure.
604 *Aust J Ecol*. 1993;18(1):117-43. doi: <https://doi.org/10.1111/j.1442-9993.1993.tb00438.x>.
- 605 16. Schroeder PJ, Jenkins DG. How robust are popular beta diversity indices to sampling
606 error? *Ecosphere*. 2018;9(2):e02100. doi: <https://doi.org/10.1002/ecs2.2100>.
- 607 17. Wong RG, Wu JR, Gloor GB. Expanding the UniFrac toolbox. *PLoS One*.
608 2016;11(9):e0161196. doi: <https://doi.org/10.1371/journal.pone.0161196>.
- 609 18. Lozupone C, Knight R. UniFrac: a new phylogenetic method for comparing microbial
610 communities. *Appl Environ Microbiol*. 2005;71(12):8228-35. doi:
611 <https://doi.org/10.1128/AEM.71.12.8228-8235.2005>.
- 612 19. Lozupone C, Lladser ME, Knights D, Stombaugh J, Knight R. UniFrac: an effective
613 distance metric for microbial community comparison. *ISME J*. 2011;5(2):169-72. doi:
614 <https://doi.org/10.1038/ismej.2010.133>.
- 615 20. Roldán Ahumada JA, Avendaño Garrido ML. A commentary on diversity measures
616 UniFrac in very small sample size. *Evol Bioinform*. 2019;15:1176934319843515. doi:
617 <https://doi.org/10.1177%2F1176934319843515>.
- 618 21. Schloss PD. Evaluating different approaches that test whether microbial communities
619 have the same structure. *ISME J*. 2008;2(3):265-75. doi: <https://doi.org/10.1038/ismej.2008.5>.

- 620 22. Warton DI, Wright ST, Wang Y. Distance-based multivariate analyses confound
621 location and dispersion effects. *Methods Ecol Evol.* 2012;3(1):89-101. doi:
622 <https://doi.org/10.1111/j.2041-210X.2011.00127.x>.
- 623 23. Randolph TW, Zhao S, Copeland W, Hullar M, Shojaie A. Kernel-penalized regression
624 for analysis of microbiome data. *Ann Appl Stat.* 2018;12(1):540-66. doi:
625 <https://dx.doi.org/10.1214/17-AOAS1102>.
- 626 24. Tang Z-Z, Chen G, Alekseyenko AV. PERMANOVA-S: association test for microbial
627 community composition that accommodates confounders and multiple distances.
628 *Bioinformatics.* 2016;32(17):2618-25. doi: <https://doi.org/10.1093/bioinformatics/btw311>.
- 629 25. Wu C, Chen J, Kim J, Pan W. An adaptive association test for microbiome data. *Genome*
630 *Med.* 2016;8(1):56. doi: <https://doi.org/10.1186/s13073-016-0302-3>.
- 631 26. Fernandes AD, Reid JN, Macklaim JM, McMurrough TA, Edgell DR, Gloor GB.
632 Unifying the analysis of high-throughput sequencing datasets: characterizing RNA-seq, 16S
633 rRNA gene sequencing and selective growth experiments by compositional data analysis.
634 *Microbiome.* 2014;2(1):15. doi: <https://doi.org/10.1186/2049-2618-2-15>.
- 635 27. Gloor GB, Macklaim JM, Pawlowsky-Glahn V, Egozcue JJ. Microbiome datasets are
636 compositional: and this is not optional. *Front Microbiol.* 2017;8:2224. doi:
637 <https://doi.org/10.3389/fmicb.2017.02224>.
- 638 28. Mandal S, Van Treuren W, White RA, Eggesbø M, Knight R, Peddada SD. Analysis of
639 composition of microbiomes: a novel method for studying microbial composition. *Microb Ecol*
640 *Health Dis.* 2015;26(1):27663. doi: <https://doi.org/10.3402/mehd.v26.27663>.
- 641 29. Robinson MD, McCarthy DJ, Smyth GK. edgeR: a Bioconductor package for
642 differential expression analysis of digital gene expression data. *Bioinformatics.*
643 2010;26(1):139-40. doi: <https://doi.org/10.1093/bioinformatics/btp616>.
- 644 30. Love MI, Huber W, Anders S. Moderated estimation of fold change and dispersion for
645 RNA-seq data with DESeq2. *Genome Biol.* 2014;15(12):1-21. doi:
646 <https://doi.org/10.1186/s13059-014-0550-8>.
- 647 31. Lin H, Peddada SD. Analysis of compositions of microbiomes with bias correction. *Nat*
648 *Commun.* 2020;11(1):1-11. doi: <https://doi.org/10.1038/s41467-020-17041-7>.
- 649 32. Saccetti E, Hoefsloot HC, Smilde AK, Westerhuis JA, Hendriks MM. Reflections on
650 univariate and multivariate analysis of metabolomics data. *Metabolomics.* 2014;10(3):361-74.
651 doi: <https://doi.org/10.1007/s11306-013-0598-6>.

- 652 33. Smilde AK, Jansen JJ, Hoefsloot HC, Lamers R-JA, Van Der Greef J, Timmerman ME.
653 ANOVA-simultaneous component analysis (ASCA): a new tool for analyzing designed
654 metabolomics data. *Bioinformatics*. 2005;21(13):3043-8. doi:
655 <https://doi.org/10.1093/bioinformatics/bti476>.
- 656 34. Stahle L, Wold S. Multivariate analysis of variance (MANOVA). *Chemometr Intell Lab*
657 *1990;9:127-41*. doi: [https://doi.org/10.1016/0169-7439\(90\)80094-M](https://doi.org/10.1016/0169-7439(90)80094-M).
- 658 35. Langsrud Ø. 50–50 multivariate analysis of variance for collinear responses. *J R Stat*
659 *Soc - Ser D Stat*. 2002;51(3):305-17. doi: <https://doi.org/10.1111/1467-9884.00320>.
- 660 36. Langsrud Ø. Rotation tests. *Stat Comput*. 2005;15(1):53-60. doi:
661 <https://link.springer.com/content/pdf/10.1007/s11222-005-4789-5.pdf>.
- 662 37. Combrink M, Du Preez I, Ronacher K, Walzl G, Loots DT. Time-dependent changes in
663 urinary metabolome before and after intensive phase tuberculosis therapy: a
664 pharmacometabolomics study. *OMICS*. 2019;23(11):560-72. doi:
665 <https://doi.org/10.1089/omi.2019.0140>.
- 666 38. Gómez-Canela C, Prats E, Lacorte S, Raldúa D, Piña B, Tauler R. Metabolomic changes
667 induced by nicotine in adult zebrafish skeletal muscle. *Ecotox Environ Safe*. 2018;164:388-97.
668 doi: <https://doi.org/10.1016/j.ecoenv.2018.08.042>.
- 669 39. Trimigno A, Khakimov B, Savorani F, Tenori L, Hendrixson V, Čivilis A, et al.
670 Investigation of Variations in the Human Urine Metabolome amongst European Populations:
671 An Exploratory Search for Biomarkers of People at Risk-of-Poverty. *Mol Nutr Food Res*.
672 2019;63(1):1800216. doi: <https://doi.org/10.1002/mnfr.201800216>.
- 673 40. Bjerke GA, Rudi K, Avershina E, Moen B, Blom H, Axelsson L. Exploring the brine
674 microbiota of a traditional Norwegian fermented fish product (Rakfisk) from six different
675 producers during two consecutive seasonal productions. *Foods*. 2019;8(2):72. doi:
676 <https://doi.org/10.3390/foods8020072>.
- 677 41. Måge I, Steppeler C, Berget I, Paulsen JE, Rud I. Multi-way methods for understanding
678 longitudinal intervention effects on bacterial communities. *bioRxiv*: 363630v1 [Preprint].
679 2018:[cited 1 Sept 2021]. Available from: <https://www.biorxiv.org/content/10.1101/363630v1>.
- 680 42. Moen B, Henjum K, Måge I, Knutsen SH, Rud I, Hetland RB, et al. Effect of dietary
681 fibers on cecal microbiota and intestinal tumorigenesis in azoxymethane treated A/J Min/+
682 mice. *PLoS One*. 2016;11(5):e0155402. doi: <https://doi.org/10.1371/journal.pone.0155402>.

- 683 43. Moen B, Røssvoll E, Måge I, Møretrø T, Langsrud S. Microbiota formed on attached
684 stainless steel coupons correlates with the natural biofilm of the sink surface in domestic
685 kitchens. *Can J Microbiol.* 2016;62(2):148-60. doi: <https://doi.org/10.1139/cjm-2015-0562>.
- 686 44. Birkeland E, Gharagozlian S, Birkeland KI, Valeur J, Måge I, Rud I, et al. Prebiotic
687 effect of inulin-type fructans on faecal microbiota and short-chain fatty acids in type 2 diabetes:
688 a randomised controlled trial. *Eur J Nutr.* 2020;59:3325-38. doi:
689 <https://doi.org/10.1007/s00394-020-02282-5>.
- 690 45. Segata N, Izard J, Waldron L, Gevers D, Miropolsky L, Garrett WS, et al. Metagenomic
691 biomarker discovery and explanation. *Genome Biol.* 2011;12(6):R60. doi:
692 <https://doi.org/10.1186/gb-2011-12-6-r60>.
- 693 46. Chi C, Xue Y, Lv N, Hao Y, Liu R, Wang Y, et al. Longitudinal gut bacterial
694 colonization and its influencing factors of low birth weight infants during the first 3 months of
695 life. *Front Microbiol.* 2019;10:1105. doi: <https://doi.org/10.3389/fmicb.2019.01105>.
- 696 47. Girard C, Tromas N, Amyot M, Shapiro BJ. Gut microbiome of the canadian arctic
697 Inuit. *mSphere.* 2017;2(1):e00297-16. doi: <https://doi.org/10.1128/mSphere.00297-16>.
- 698 48. Lai Z-L, Tseng C-H, Ho HJ, Cheung CK, Lin J-Y, Chen Y-J, et al. Fecal microbiota
699 transplantation confers beneficial metabolic effects of diet and exercise on diet-induced obese
700 mice. *Sci Rep.* 2018;8(1):15625. doi: <https://doi.org/10.1038/s41598-018-33893-y>.
- 701 49. Le Sciellour M, Labussière E, Zemb O, Renaudeau D. Effect of dietary fiber content on
702 nutrient digestibility and fecal microbiota composition in growing-finishing pigs. *PLoS One.*
703 2018;13(10):e0206159. doi: <https://doi.org/10.1371/journal.pone.0206159>.
- 704 50. Wang Y, Ames NP, Tun HM, Tosh SM, Jones PJ, Khafipour E. High molecular weight
705 barley β -glucan alters gut microbiota toward reduced cardiovascular disease risk. *Front*
706 *Microbiol.* 2016;7:129. doi: <https://doi.org/10.3389/fmicb.2016.00129>.
- 707 51. Vangay P, Johnson AJ, Ward TL, Al-Ghalith GA, Shields-Cutler RR, Hillmann BM, et
708 al. US immigration westernizes the human gut microbiome. *Cell.* 2018;175(4):962-72.e10. doi:
709 <https://doi.org/10.1016/j.cell.2018.10.029>.
- 710 52. Hughes RL, Kable ME, Marco M, Keim NL. The role of the gut microbiome in
711 predicting response to diet and the development of precision nutrition models. Part II: results.
712 *Adv Nutr.* 2019;10(6):979-98. doi: <https://doi.org/10.1093/advances/nmz049>.

- 713 53. Johnson AJ, Vangay P, Al-Ghalith GA, Hillmann BM, Ward TL, Shields-Cutler RR, et
714 al. Daily sampling reveals personalized diet-microbiome associations in humans. *Cell Host*
715 *Microbe*. 2019;25(6):789-802. e5. doi: <https://doi.org/10.1016/j.chom.2019.05.005>.
- 716 54. Rothschild D, Weissbrod O, Barkan E, Kurilshikov A, Korem T, Zeevi D, et al.
717 Environment dominates over host genetics in shaping human gut microbiota. *Nature*.
718 2018;555(7695):210-5. doi: <https://doi.org/10.1038/nature25973>.
- 719 55. So D, Whelan K, Rossi M, Morrison M, Holtmann G, Kelly JT, et al. Dietary fiber
720 intervention on gut microbiota composition in healthy adults: a systematic review and meta-
721 analysis. *Am J Clin Nutr*. 2018;107(6):965-83. doi: <https://doi.org/10.1093/ajcn/nqy041>.
- 722 56. Nagpal R, Wang S, Solberg Woods LC, Seshie O, Chung ST, Shively CA, et al.
723 Comparative microbiome signatures and short-chain fatty acids in mouse, rat, non-human
724 primate, and human feces. *Front Microbiol*. 2018;9:2897. doi:
725 <https://doi.org/10.3389/fmicb.2018.02897>.
- 726 57. Baker DH. Animal models in nutrition research. *Nutr J*. 2008;138(2):391-6. doi:
727 <https://doi.org/10.1093/jn/138.2.391>.
- 728 58. Hui BS, Wold HOA. Consistency and consistency at large of partial least squares
729 estimates. In: Jöreskog KG, Wold HOA, editors. *Systems under indirect observation, part II*.
730 Amsterdam: North-Holland Publ. Co.; 1982. p. 119-30.
- 731 59. Schneeweiss H. Consistency at large in models with latent variables. In: Haagen K,
732 Bartholomew DJ, Deistler M, editors. *Statistical modelling and latent variables*. Amsterdam:
733 Elsevier; 1993. p. 299-320.
- 734 60. van den Berg RA, Hoefsloot HC, Westerhuis JA, Smilde AK, van der Werf MJ.
735 Centering, scaling, and transformations: improving the biological information content of
736 metabolomics data. *BMC Genomics*. 2006;7(1):142. doi: [https://doi.org/10.1186/1471-2164-7-](https://doi.org/10.1186/1471-2164-7-142)
737 [142](https://doi.org/10.1186/1471-2164-7-142).
- 738 61. Thiel M, Féraud B, Govaerts B. ASCA+ and APCA+: Extensions of ASCA and APCA
739 in the analysis of unbalanced multifactorial designs. *J Chemom*. 2017;31(6):e2895. doi:
740 <https://doi.org/10.1002/cem.2895>.
- 741 62. Martin M, Govaerts B. LiMM-PCA: Combining ASCA+ and linear mixed models to
742 analyse high-dimensional designed data. *J Chemom*. 2020;34(6):e3232. doi:
743 <https://doi.org/10.1002/cem.3232>.

- 744 63. Hughes RL, Marco ML, Hughes JP, Keim NL, Kable ME. The role of the gut
745 microbiome in predicting response to diet and the development of precision nutrition models—
746 Part I: overview of current methods. *Adv Nutr.* 2019;10(6):953-78. doi:
747 <https://doi.org/10.1093/advances/nmz022>.
- 748 64. Hawinkel S, Mattiello F, Bijmans L, Thas O. A broken promise: microbiome differential
749 abundance methods do not control the false discovery rate. *Brief Bioinform.* 2019;20(1):210-
750 21. doi: <https://doi.org/10.1093/bib/bbx104>.
- 751 65. Thorsen J, Brejnrod A, Mortensen M, Rasmussen MA, Stokholm J, Al-Soud WA, et al.
752 Large-scale benchmarking reveals false discoveries and count transformation sensitivity in 16S
753 rRNA gene amplicon data analysis methods used in microbiome studies. *Microbiome.*
754 2016;4(1):62. doi: <https://doi.org/10.1186/s40168-016-0208-8>.
- 755 66. Boulesteix AL, Binder H, Abrahamowicz M, Sauerbrei W. On the necessity and design
756 of studies comparing statistical methods. *Biometrical J.* 2017;60(1):216-8. doi:
757 <https://doi.org/10.1002/bimj.201700129>.
- 758 67. R Development Core Team. R: A language and environment for statistical computing.
759 R Foundation for Statistical Computing, Vienna, Austria. <https://www.R-project.org/>. 2021.
- 760 68. Patuzzi I, Baruzzo G, Losasso C, Ricci A, Di Camillo B. metaSPARSim: a 16S rRNA
761 gene sequencing count data simulator. *BMC Bioinform.* 2019;20(9):1-13. doi:
762 <https://doi.org/10.1186/s12859-019-2882-6>.
- 763 69. Aitchison J. The analysis of compositional data. London: Chapman and Hall; 1986. 416
764 p.
- 765 70. Palarea-Albaladejo J, Martín-Fernández JA. zCompositions—R package for
766 multivariate imputation of left-censored data under a compositional approach. *Chemometr*
767 *Intell Lab.* 2015;143:85-96. doi: <https://doi.org/10.1016/j.chemolab.2015.02.019>.
- 768 71. Brill B, Amir A, Heller R. Testing for differential abundance in compositional counts
769 data, with application to microbiome studies. arXiv:190408937 [Preprint]. 2019:[cited 1 Sept
770 2021]. Available from: <https://arxiv.org/pdf/1904.08937.pdf>.
- 771 72. Costea PI, Zeller G, Sunagawa S, Bork P. A fair comparison. *Nat Methods.*
772 2014;11(4):359. doi: <https://doi.org/10.1038/nmeth.2897>.
- 773 73. Wang S. Robust differential abundance test in compositional data. arXiv: 210108765
774 [Preprint]. 2021:[cited 1 Sept 2021]. Available from: <https://arxiv.org/pdf/101.08765.pdf>.

- 775 74. Lubbe S, Filzmoser P, Templ M. Comparison of zero replacement strategies for
776 compositional data with large numbers of zeros. *Chemometr Intell Lab.* 2021;210:104248. doi:
777 <https://doi.org/10.1016/j.chemolab.2021.104248>.
- 778 75. Oksanen J, Blanchet FG, Friendly M, Kindt R, Legendre P, McGlinn D, et al. *vegan*:
779 Community ecology package (R package version 2.5-5) <https://cran.r-project.org/> and
780 <https://github.com/vegandevs/vegan/>. 2019.
- 781 76. Eriksson L, Byrne T, Johansson E, Trygg J, Vikström C. *Multi-and megavariable data*
782 *analysis basic principles and applications*. Malmö, Sweden: Umetrics Academy; 2013.
- 783 77. Centner V, Massart D-L, de Noord OE, de Jong S, Vandeginste BM, Sterna C.
784 Elimination of uninformative variables for multivariate calibration. *Anal Chem.*
785 1996;68(21):3851-8. doi: <https://doi.org/10.1021/ac960321m>.
- 786
- 787

Figures

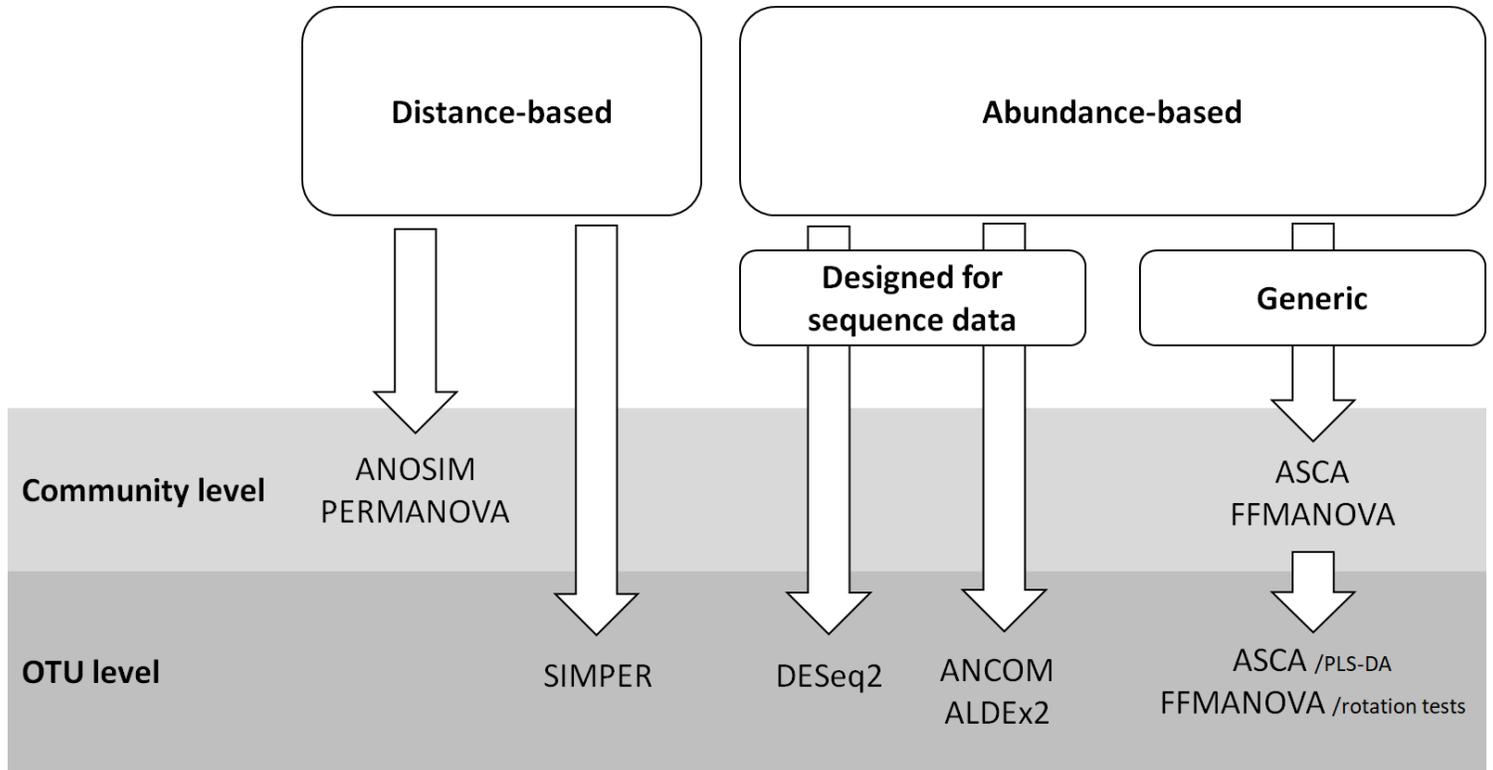


Figure 1

A diagram of statistical methods used in the study.

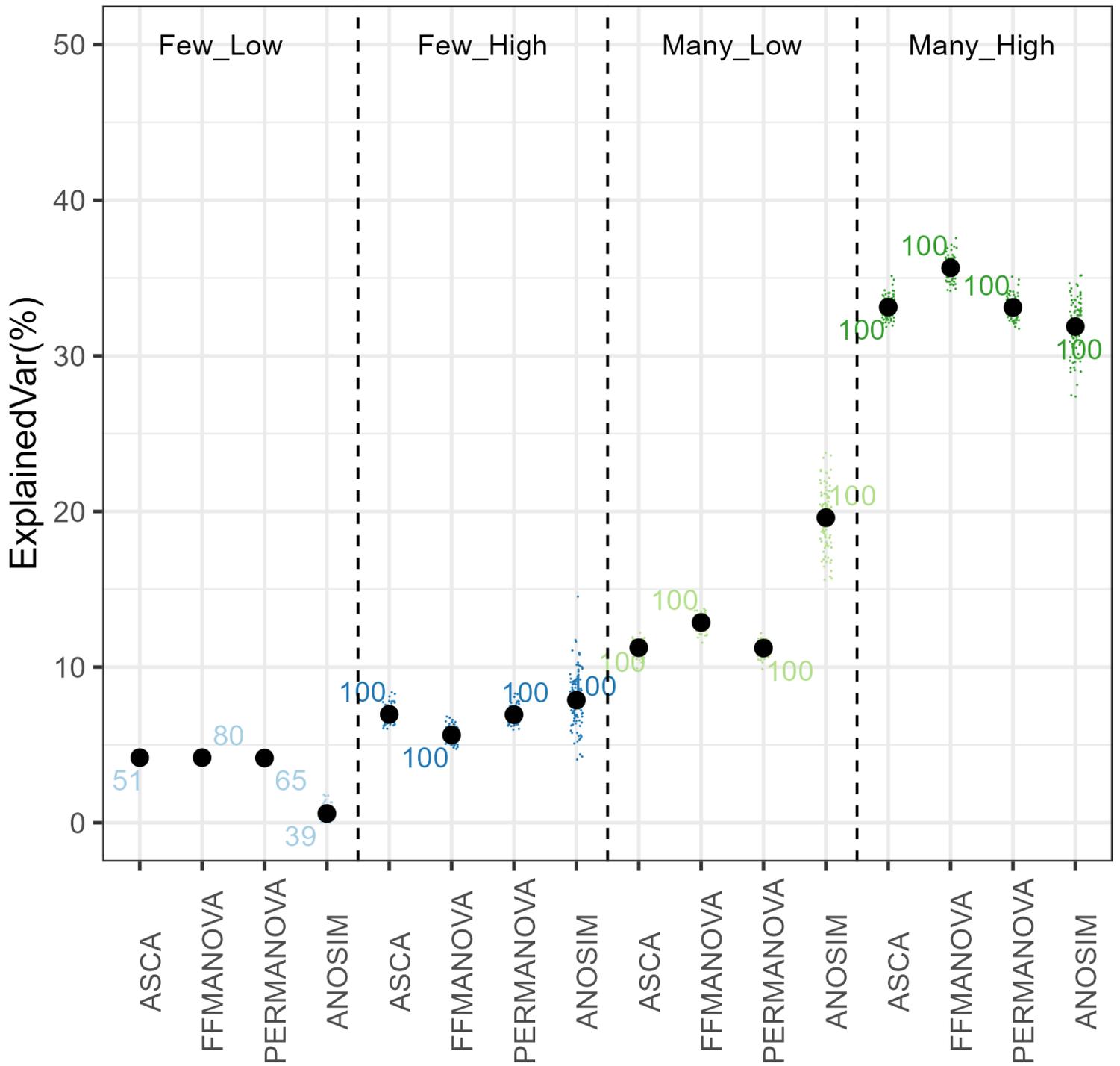


Figure 2

Explained variance for simulated data and the relative number of simulations where the simulated effect was detected.

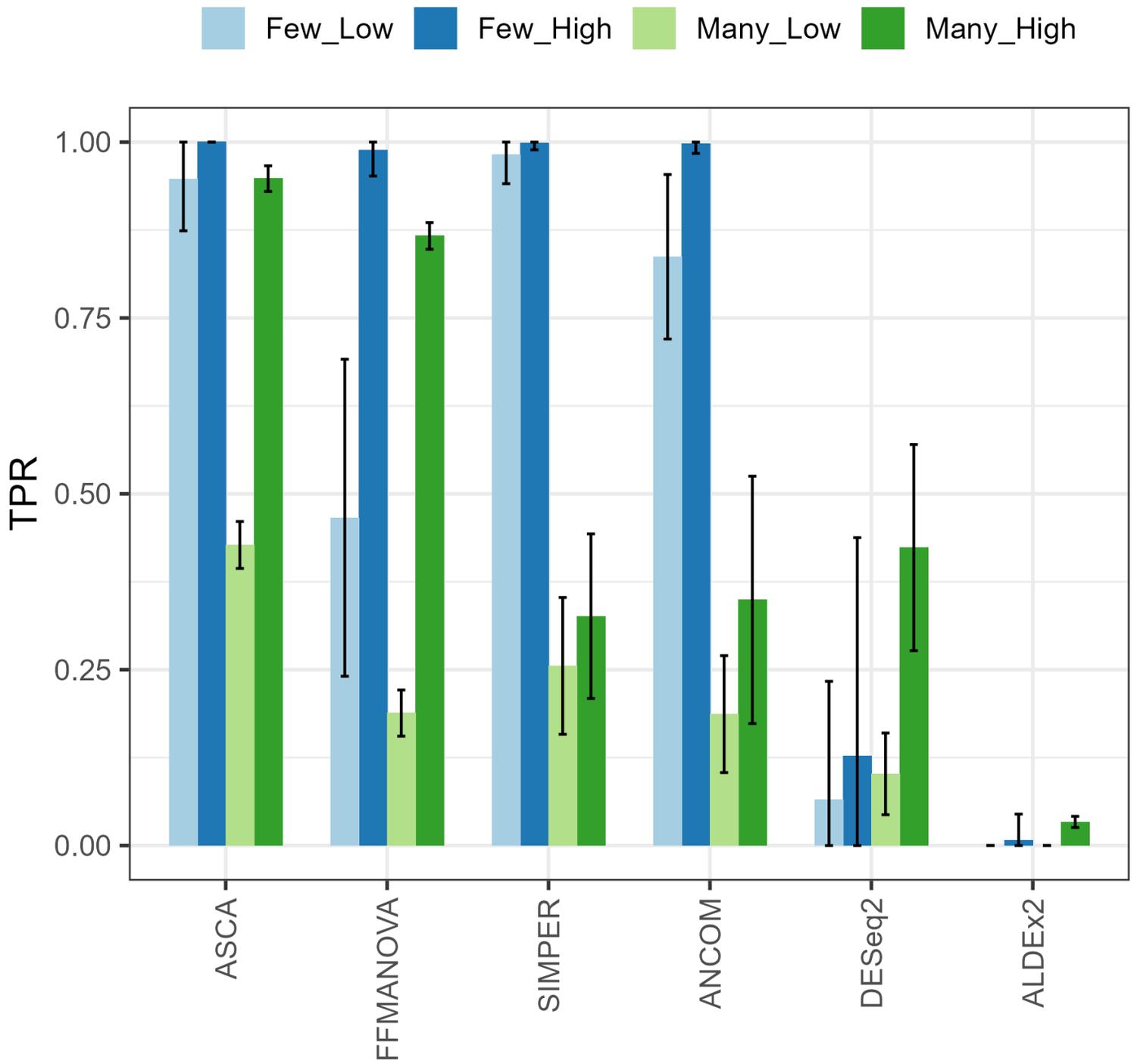


Figure 3

Sensitivity (True Positive Rate) for the four scenarios in the simulation study.

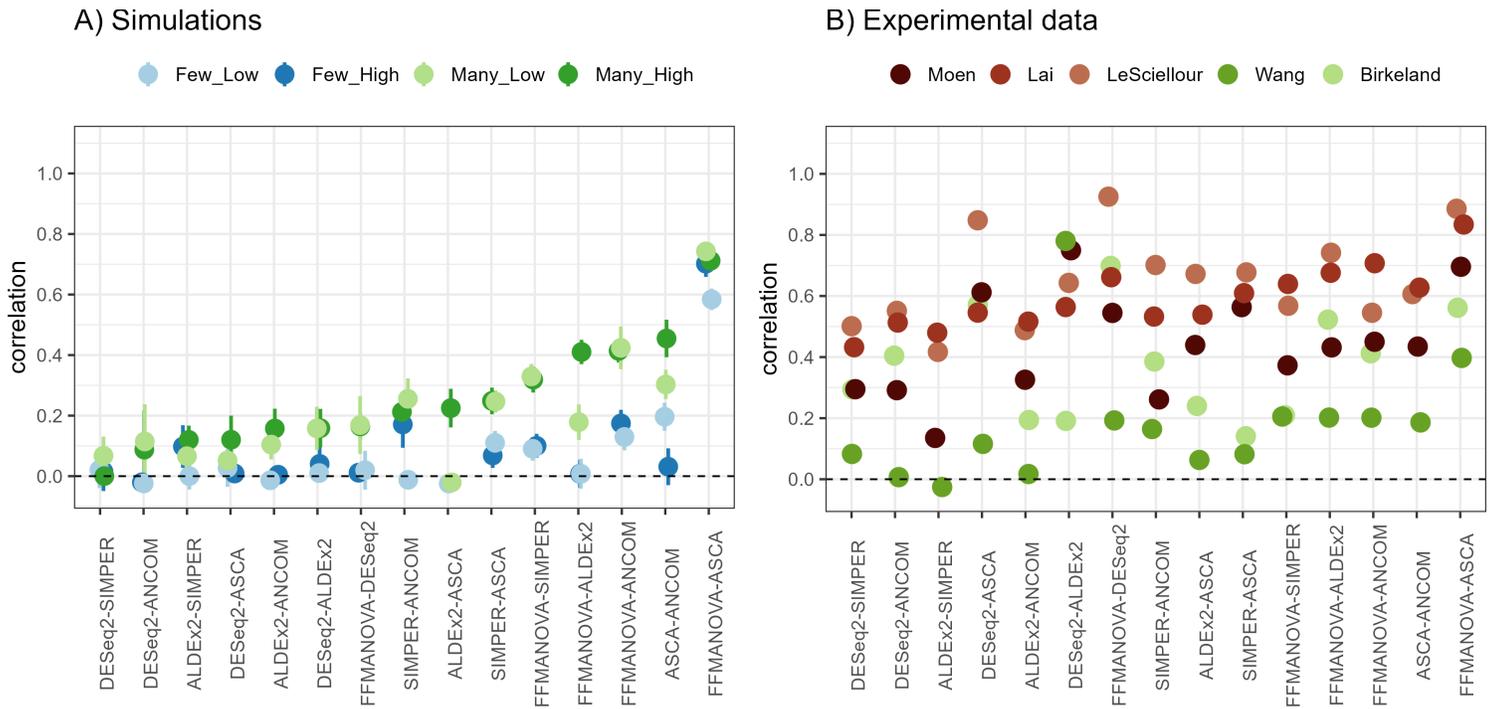


Figure 4

Spearman's correlation (Y-axis) calculated for pairwise comparison of statistical methods (X-axis) for (A) simulated data and (B) five experimental data sets. Each point represents Spearman's rank correlation coefficient between OTU ranking metrics from the two methods compared.

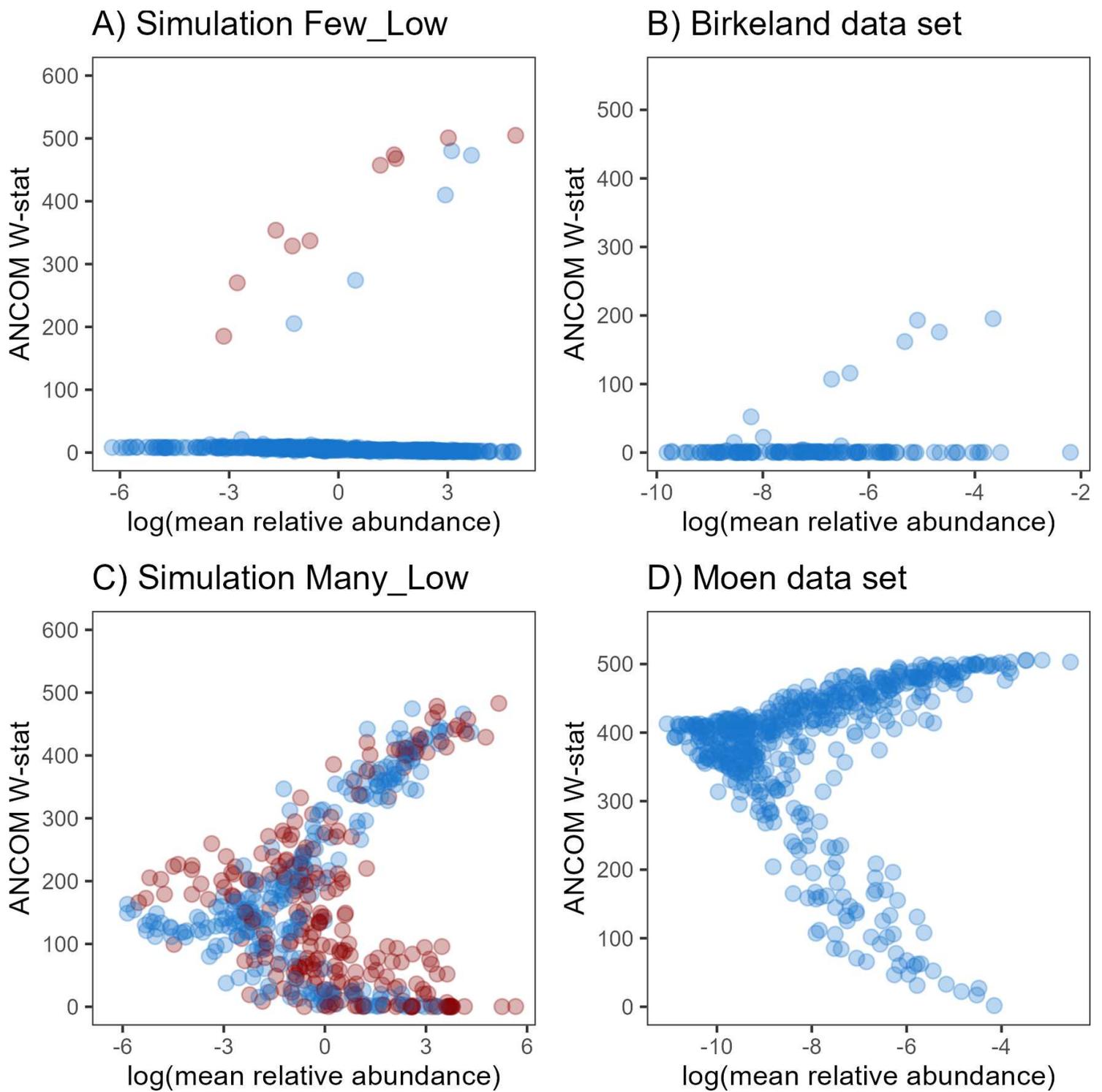


Figure 5

Mean relative abundance (log-scale) plotted versus ANCOM W-stat. (A) The “Few-Low” simulation scenario and (B) Birkeland data set, (C) “Many-Low” simulation scenario and (D) Moen data set.

Supplementary Files

This is a list of supplementary files associated with this preprint. [Click to download.](#)

- [KHOMICHetal2021preprintSupportingInformation.pdf](#)
- [S1Fig.pdf](#)
- [S1File.xlsx](#)
- [S1Table.xlsx](#)
- [S2Fig.pdf](#)
- [S2File.xlsx](#)
- [S2Table.xlsx](#)
- [S3File.rdata](#)
- [S4File.r](#)
- [S5File.r](#)