

# Feature matrix normalization, transformation and calculation of $\beta$ -diversity in metagenomics: Theoretical and applied perspectives on your decisions

1

2 **Casper Sahl Poulsen<sup>1\*</sup>, Frank Møller Aarestrup<sup>1</sup>, Christian Brinch<sup>1</sup> & Claus Thorn Ekstrøm<sup>2</sup>**

3 <sup>1</sup>Research Group for Genomic Epidemiology, National Food Institute, Technical University of  
4 Denmark, Kongens Lyngby, Denmark.

5 <sup>2</sup>Section of Biostatistics, Department of Public Health, University of Copenhagen, Copenhagen,  
6 Denmark.

7

8 **\* Correspondence:**

9 Casper Sahl Poulsen

10 cspoulsen@hotmail.com

11 **Keywords: Metagenomics, Data manipulation, Visualization, Dissimilarity index,  $\beta$ -diversity,**  
12 **normalization, transformation**

13

14 **Abstract**

15 Microbial metagenomics utilising next generation sequencing is a powerful experimental approach  
16 enabling detailed and potentially complete descriptions of the microbial world around and within us.  
17 Selecting how to perform feature data normalization, transformation and calculate  $\beta$ -diversity is a  
18 critical step in the analysis of metagenomic data, but also a step for which a multitude of methods are  
19 available. Researchers need to have a broad overview and understand the many methods that exist in  
20 the field and the consequences from applying them. In this perspectives article, some of the most widely  
21 used metagenomic feature data normalizations, transformations and  $\beta$ -diversity metrics are discussed  
22 in the context of multivariate visualizations. We provide a framework that other researchers can utilize  
23 to evaluate how robust their test data are when applying different normalizations, transformations and  
24  $\beta$ -diversity metrics, and visually compare the results of the methods. We constructed an *in silico* test  
25 dataset to evaluate the setup and clarify how the theoretical discussion is transferable to this data. We  
26 urge other researchers to implement their own test data, normalization, transformation,  $\beta$ -diversity  
27 metric and visualization methods, in the hope that it will advance better decision making both in study  
28 design and analysis strategy.

29

30 **1 The lack of consensus on how to perform data normalization, transformation and**  
31 **calculate  $\beta$ -diversity**

32 Next generation sequencing (NGS) is applied heavily in microbiome research, enabling both  
33 taxonomic and functional descriptions of microbiomes (1,2). Metagenomic data need to be processed

34 before analysis to make sample comparisons possible due to the differences in sequencing depth.  
35 Furthermore, there is an increasing awareness that such data are compositional and should be processed  
36 accordingly (3–7).

37 The choice of how to perform data normalization, transformation and computation of  $\beta$ -diversity can  
38 have a substantial impact on the results from the subsequent data analysis especially since metagenomic  
39 data typically are sparse, because some features are not present or their abundance are below the limit  
40 of detection. Classically, metagenomic feature data are either relativized to some sample characteristics  
41 such as the number of total reads, bacterial reads etc., or from a compositional, more transparent  
42 measure according to the number of assigned reads that is also known as total sum scaling (TSS). When  
43 relativizing, the precision of measurement is lost, considering that data are heteroscedastic direct  
44 comparison of samples is flawed if methods assume equal variance. (8). Therefore, rarefying can be  
45 performed, but it has been argued against due to the loss of power (8). Relativizing is highly influenced  
46 by the most abundant features, alternatively, the median, quantile normalization or cumulative sum  
47 scaling (CSS) can be used (9,10). Methods developed for normalizing data such as trimmed mean of  
48 M-values (TMM) and relative log expression (RLE), can be relevant if most features are not changing  
49 between samples (11,12). The compositional data analysis framework provides an additional approach  
50 to analyse metagenomic data with a multitude of possibilities for estimating zeroes and visualization  
51 (13–16). There are arguably advantages and disadvantages to applying all of the different methods  
52 described (7,9,17,18).

53 Several R packages have implemented the techniques described above such as “vegan”, “edgeR”,  
54 “DESeq2”, “phyloseq” and “compositions” (10,13,19–21). From these packages, we have identified  
55 228 combinations of normalizing and transforming data and calculating  $\beta$ -diversity metrics. This is not  
56 an exhaustive list of possible methods to apply, and therefore processing metagenomic data is a task  
57 where tradition and ease of implementation are important factors governing researchers’ decisions.  
58 Understanding the more advanced methods, for instance, to perform compositional data analysis is  
59 most likely also a reason for these methods to not have become common as observed in other fields  
60 (22).

61 The aim of the present study is to provide theoretical as well as applied analytical perspectives on  
62 normalization and transformation of metagenomic data in the context of calculating  $\beta$ -diversity that is  
63 used for statistical inference and multivariate visualizations. We have constructed an *in silico* dataset  
64 to visualize how data processing affects metagenomic analysis. The dataset was used to investigate if  
65 methods are robust according to sequencing depth and the influence of changes in data structure.  
66 Furthermore, a visualization of a dissimilarity matrix containing the Procrustes test results for all  
67 selected methods compared pairwise provides a comparison of how the methods resemble each other.  
68 We provided the code used to generate the analysis in the hope that other researchers can use it as a  
69 tool for assessing the effect and sensitivity of using different transformation, normalization and  $\beta$ -  
70 diversity methods by incorporating their own test data, favourite methods or visualization techniques.

71

## 72 **2 Theoretical perspectives on data normalization, transformation and $\beta$ -diversity** 73 **calculation in metagenomics**

74 In this section, perspectives are provided on feature data transformation, normalization and  $\beta$ -diversity  
75 metrics where we, for the latter, have focused on Euclidean distance, Manhattan distance and Bray-  
76 Curtis dissimilarity due to their widespread use and acceptance in metagenomics. In terms of between

77 sample comparisons, normalization is primarily performed to take sequencing depth into account and,  
78 transformation is performed to weigh how the differences between features should be emphasized.

79 One of the most common methods to account for sequencing depth is TSS. When calculating  $\beta$ -  
80 diversity, this method is driven by the features with the largest differences between samples that are  
81 typically also the most abundant features because it scales linearly in absolute values. Multivariate  
82 visualisations and statistical tests therefore depend on the differences in the most abundant features.  
83 To deemphasize this effect, log transformations or square root transformations, such as the Hellinger  
84 transformation, are used. Sometimes these transformations are applied post TSS; however, if this is the  
85 case, it should be emphasized that sample values no longer add up to the same total sum anymore. If  
86 detection is the primary focus of analysis, making a presence absence (PA) transformation can be  
87 justified because this removes the effect of abundance. From a practical viewpoint, PA requires high  
88 specificity when mapping reads, commonly at the cost of sensitivity, control of contamination during  
89 sample processing and is not robust according to sequencing depth unless different detection limits are  
90 implemented prior to transformation.

91 Rarefying (or subsampling) provides data where precision of a measurement is the same across  
92 samples, typically performed by rarefying to the level of the sample with the fewest assigned reads, at  
93 the cost of sensitivity. The loss of precision is usually not a problem if sequencing depth is even, but a  
94 similar argument can be made in this case when relativizing. Data are still heteroscedastic and therefore  
95 should be modelled accordingly, i.e. when performing differential abundance analysis (8,9). Both  
96 relativizing and rarefying do not take the compositionality of data into account, but perform well if the  
97 most abundant features between samples are relatively constant, which is rarely known. If, on the other  
98 hand, most features are not changing between compared groups of samples, the median, RLE and TMM  
99 offer a better solution (9). This is also why RLE and TMM were implemented in DESeq2 and edgeR,  
100 respectively, for the analysis of expression data. In expression data, it is often a good assumption, for  
101 instance in a clinical study, that treatment only changes expression of a few genes (11,21). In  
102 metagenomics, this assumption could be met, but from our experience, working in the field and with  
103 spike-in organisms, this is rarely the case.

104 When calculating  $\beta$ -diversity, the length of the straight line between two points can be calculated, this  
105 is also known as the Euclidean distance. This method would be straightforward if the points were in  
106 Euclidean space, but metagenomic data are compositional and points are therefore confined to a  
107 simplex. When calculating Euclidean distances, the differences are squared, consequently, the greatest  
108 differences are further emphasized relative to using Manhattan distance or Bray-Curtis dissimilarity.  
109 This could be counterbalanced by performing a log or Hellinger transformation. Manhattan distance is  
110 the sum of absolute differences. Manhattan distance is also the numerator of Bray-Curtis dissimilarity  
111 that is then scaled to the sum of total features in the two samples. The Manhattan distance also does  
112 not account for the compositionality of data.

113

### 114 **3 Another approach to data normalization in metagenomics**

115 A solution to the challenges described above is to use a compositional analysis framework. Using  
116 centered log ratio (CLR) transformation, where the log of each feature is compared relative to the  
117 geometric mean, or the isometric log ratio (ILR) transformation, where orthogonal basis functions are  
118 used to span the simplex space somewhat analogous to the CLR transformation, in the context of  
119 calculating  $\beta$ -diversity (23,24). Performing both methods enables real-space calculations and

120 consequently Euclidean distances when calculating  $\beta$ -diversity. The methods are simple in principle,  
121 but zeroes have to be imputed and this represents a major challenge when dealing with metagenomic  
122 data that are typically sparse (4,25,26). One often-used solution is to detect features with a zero and  
123 then remove the features from all samples. This option is recommended when features are low abundant  
124 in the others samples, but in metagenomic studies, a feature might be relatively highly abundant in one  
125 sample and not present in another. Another approach is to add a pseudo-count, multiplicative simple  
126 replacement or a Bayesian approach (15,27,28). Nonetheless, in all imputations of zeroes, there is no  
127 way of knowing the difference between a “true” zero representing a feature that is not present and a  
128 zero that is below the detection limit. Imputation in this situation is therefore limited to assigning a  
129 value below the detection limit, even though the feature might not be present (27,28).

130 From a mathematical perspective, we expect the compositional methods to offer a desirable  
131 characteristic in that data are not constrained to the simplex, but considering sparsity, which is  
132 commonly an artefact of metagenomic data, zeroes have to be imputed.

133

#### 134 **4 Seeing is believing - *In silico* comparison of data normalization, transformation and $\beta$ -** 135 **diversity calculation in metagenomics**

136 To provide applied analytical perspectives, an *in silico* dataset was constructed to reflect typical  
137 challenges in metagenomic data including sparsity and differences in sequencing depth. A reference  
138 (Ref), equivalent to a sample, was created consisting of abundance profiles of 70 different organisms  
139 (i.e. number of sequence reads mapping to a given organism). The sample consisted of counts from the  
140 following abundance levels:

- 141 • **High** (1 random sampling between 1000-5000),
- 142 • **Medium high** (3 random samplings between 100-999 with replacement),
- 143 • **Medium** (9 random samplings between 5-99 with replacement),
- 144 • **Low** (27 random samplings between 0-4 with replacement), and
- 145 • **Not present** (30 zeroes).

146 The test data contained 70 different features (i.e. organisms), but this was a trade-off to make the  
147 analysis run on a desktop computer.

148 Eleven other samples were created, all variations of the reference:

- 149 • Multiplying counts with 2 (SF2) and 10 (SF10),
- 150 • Changing counts to zeroes in each of the different abundance levels (SwHato0, SwMHato0,  
151 SwMato0, SwLato0),
- 152 • Switching the highly abundant feature with one in each of the other abundance groups  
153 (SwHaMHa, SwHaMa, SwHaLa, SwHaNP), and
- 154 • Reversing the reference (RevRef) to create a very dissimilar sample only sharing a few low  
155 abundant features.

156 These artificial samples represent potential differences that are of interest to assess the effect of  
157 sequencing depth and structural differences in data. The full computer code documents the exact  
158 construction of the samples and their variations  
159 (<https://github.com/csapou/DataProcessinginMetagenomics>).

160 To limit the number of combinations of normalization, transformation and  $\beta$ -diversity metrics in  
161 figures, we selected 36 methods. We included Euclidean distance, Manhattan distance and Bray-Curtis  
162 dissimilarity as  $\beta$ -diversity metrics, since these metrics are popular in metagenomics. The selected  
163 transformation and normalization steps were based on tradition in the field of microbial ecology (TSS,  
164 rarefying, PA and CSS). We also included Hellinger and log transformation both before and after TSS.  
165 Some methods are implemented to normalize RNA-expression data (TMM and DESeq (poscount  
166 argument in estimating SizeFactors)). For methods that adhere to the compositional data analysis  
167 framework, we included six methods that use Euclidean distances. Zeroes were estimated with  
168 multiplicative simple replacement or adding a pseudo-count of one prior to TSS, then performing both  
169 CLR and ILR. We also included TSS and then added a pseudo-count of the minimum divided by ten  
170 before performing CLR and ILR.

171 All statistical analysis and visualization of data were performed in R version 3.4.4, and data  
172 transformation, normalization and calculation of  $\beta$ -diversity were performed using the packages  
173 described above. To visualize the dissimilarities and distances between the different samples, we  
174 created a heatmap with accompanying dendrograms using complete linkage clustering of Euclidean  
175 distances based on the full-scale distance and dissimilarity matrices. Heatmaps were generated using  
176 the ‘pheatmap’ package by extracting the  $\beta$ -diversity to the reference sample.  $\beta$ -diversity values were  
177 made comparable in each of the methods by scaling to the max value. To compare all the distance and  
178 dissimilarity matrices pairwise, a Procrustes approach was used based on randomization tests (29,30).  
179 A dissimilarity matrix of the processing methods was created by subtracting the Procrustes correlations  
180 from one. Metric multidimensional scaling of the dissimilarity matrix was performed by running the  
181 capscale function unconstrained from ‘vegan’. The generation of the principal coordinates analysis  
182 (PCoA) plot of the first two dimensions, density plot of the correlations, stress plot containing a  
183 scatterplot of the distance observed in the PCoA as a function of the “true”  $\beta$ -diversity calculated and  
184 the scree plot showing the variance in the principal components were performed with ‘ggplot2’ (31).  
185 The code to generate test data and perform data processing is provided at  
186 <https://github.com/csapou/DataProcessinginMetagenomics> with additional principal component  
187 analysis (PCA) plots and PCoA plots of all individual methods and randomly generated samples.

188 From Figure 1 we find that samples scaled by a factor of 2 or 10 had a low  $\beta$ -diversity relative to the  
189 reference sample, indicating that the methods we selected were able to control the effect of sequencing  
190 depth, which is a bare minimum for applying them to this type of data. Some inconsistency was  
191 observed when performing log or log-ratio transformations. This effect can be reduced in this case by  
192 estimating zeroes at a lower level. The reverse sample representing a dissimilar community was also  
193 generally the one with the highest  $\beta$ -diversity relative to the reference.  $\beta$ -diversity metrics generally  
194 cluster containing either Euclidean distances or Bray-Curtis dissimilarity together with Manhattan  
195 distance. Bray-Curtis dissimilarity and Manhattan distance cluster when performing TSS, because the  
196 denominator evaluates to 2 when calculating the Bray-Curtis dissimilarity and is therefore just a factor  
197 of two scaling of the Manhattan distance. Transformation and normalization methods also cluster to  
198 some extent.

199 In Figure 2, where the processing methods are compared pairwise in full scale, the classical methods  
200 show a spectrum between abundance-driven processing exemplified by TSS and PA (Fig. 2A). In  
201 between these extremes, variations of Hellinger and log transformations are plotted as we expected  
202 from the theoretical discussion. The methods adhering to the compositional data analysis framework  
203 do not cluster, emphasizing the need for further investigations into the effect of estimating zeroes. The  
204 methods developed for normalization of gene expression data to perform differential abundance  
205 analysis are likely to perform badly with this *in silico* data because they assume that large proportions

206 of features are constant between samples. Comparing communities that are highly different, for  
207 example, the reverse sample in this dataset makes these methods inappropriate. The validation plots in  
208 the form of stress plot and scree plot show that the observed dissimilarity correlates with the ordination  
209 distance, and a large proportion of the variance is explained in the first two axes, respectively (Fig. 2C-  
210 D).

211

## 212 **5 Multivariate visualization in metagenomics - a step forward**

213 We hope that the theoretical perspectives together with the visualizations provided demonstrate that  
214 data normalization, transformation and calculation of  $\beta$ -diversity have a substantial impact on the  
215 analysis and multivariate visualization of metagenomic data. We consider the public source code as a  
216 resource that other researchers can utilize to implement their own favourite methods for processing  
217 metagenomic data (<https://github.com/csapou/DataProcessinginMetagenomics>). Here, we also provide  
218 all of the 228 methods that we have identified with additional randomly generated samples. To perform  
219 a sensitivity analysis of the effect of using different data normalization and transformation strategies  
220 in the context of calculating  $\beta$ -diversity, a density plot is provided for all of the Procrustes test  
221 correlations. From the analysis on our test dataset we see that there are two peaks with approximately  
222 the same height and the lower one is centred around a correlation of 0.5, indicating that data processing  
223 is important for this test data (Fig. 2B). On the other hand, performing this analysis on another test  
224 dataset might reveal high correlations between all methods. This would indicate that the conclusions  
225 derived from the data are robust to withstand applying different normalizations, transformations and  
226  $\beta$ -diversity metrics.

227 Other relevant modifications include the removal of the reversed sample from the analysis to look at  
228 the subtle differences between similar samples. With the large number of combinations of  
229 normalizations, transformations and  $\beta$ -diversity metrics to select from, we discourage other researchers  
230 from implementing their real data to circumvent pipeline-hacking analogous to p-hacking (32). A better  
231 option for the users would be to implement their own relevant test dataset, and from this analysis, and  
232 together with theoretical considerations, select one or a few processing methods before analysing their  
233 real data. We hope that the code provided also eases the implementation of new methods. Generation  
234 of dendograms in the heatmap and the PCoA of Procrustes test results were run using default settings,  
235 and an investigation could also be initiated to assess how this might influence the results. Again, we  
236 urge others to implement their own favourite methods.

237 We would like to highlight other aspects of good scientific practice in metagenomics and refer readers  
238 to articles on study design (33–35), sample processing (36–39) and other aspects of metagenomic data  
239 analysis primarily focusing on differential abundance analysis of features (3,7–9,40–42).

240

## 241 **6 Acknowledgments**

242 The authors wish to thank Jeffrey Skiby for language editing.

243

## 244 **7 Author Contributions Statement**

245 CP and CE conceived the ideas and wrote the paper. CP analysed the data, made the figures and  
246 performed the literature review. FA and CB revised the manuscript. All authors read and approved the  
247 manuscript.

248

## 249 **8 Conflict of Interest Statement**

250 No conflict of interest

251

## 252 **9 Funding**

253 CP has received funding from the European Union's Horizon 2020 research and innovation programme  
254 under grant agreement No. 643476 (COMPARE).

255

## 256 **10 References**

- 257 1. Kim M, Lee K-H, Yoon S-W, Kim B-S, Chun J, Yi H. Analytical Tools and Databases for  
258 Metagenomics in the Next-Generation Sequencing Era. *Genomics Inform* [Internet].  
259 2013;11(3):102. Available from:  
260 <https://synapse.koreamed.org/DOIx.php?id=10.5808/GI.2013.11.3.102>
- 261 2. Land M, Hauser L, Jun S-R, Nookaew I, Leuze MR, Ahn T-H, et al. Insights from 20 years of  
262 bacterial genome sequencing. *Funct Integr Genomics*. 2015;
- 263 3. Nayfach S, Pollard KS. Toward Accurate and Quantitative Comparative Metagenomics. *Cell*  
264 [Internet]. 2016;166(5):1103–16. Available from: <http://dx.doi.org/10.1016/j.cell.2016.08.007>
- 265 4. Gloor GB, Reid G. Compositional analysis: a valid approach to analyze microbiome high-  
266 throughput sequencing data. *Can J Microbiol*. 2016;62(8):692–703.
- 267 5. Gloor GB, Macklaim JM, Pawlowsky-Glahn V, Egozcue JJ. Microbiome datasets are  
268 compositional: And this is not optional. *Front Microbiol*. 2017;8:1–6.
- 269 6. Odintsova V, Tyakht A, Alexeev D. Guidelines to Statistical Analysis of Microbial  
270 Composition Data Inferred from Metagenomic Sequencing. *Curr Issues Mol Biol* [Internet].  
271 2017;24:17–36. Available from: <http://www.caister.com/cimb/abstracts/v24/17.html>
- 272 7. Weiss S, Xu ZZ, Peddada S, Amir A, Bittinger K, Gonzalez A, et al. Normalization and  
273 microbial differential abundance strategies depend upon data characteristics. *Microbiome*.  
274 2017;5(1):1–18.
- 275 8. McMurdie PJ, Holmes S. Waste Not, Want Not: Why Rarefying Microbiome Data Is  
276 Inadmissible. *PLoS Comput Biol*. 2014;10(4).
- 277 9. Pereira MB, Wallroth M, Jonsson V, Kristiansson E. Comparison of normalization methods  
278 for the analysis of metagenomic gene abundance data. *BMC Genomics*. 2018;19(1):1–17.

- 279 10. Oshlack A, Robinson MD. A scaling normalization method for differential expression analysis  
280 of RNA-seq data. *Genome Biol.* 2010;11.
- 281 11. Robinson MD, Oshlack A. A scaling normalization method for differential expression analysis  
282 of RNA-seq data. *Genome Biol.* 2010;11(R25).
- 283 12. Love MI, Huber W, Anders S. Moderated estimation of fold change and dispersion for RNA-  
284 seq data with DESeq2. *Genome Biol.* 2014;15(550).
- 285 13. Boogaart AKG Van Den, Tolosana-delgado R, Bren M, Boogaart MKG Van Den. Package ‘  
286 compositions.’ 2018;
- 287 14. Cao KAL, Costello ME, Lakis VA, Bartolo F, Chua XY, Brazeilles R, et al. MixMC: A  
288 multivariate statistical framework to gain insight into microbial communities. *PLoS One.*  
289 2016;
- 290 15. Palarea-Albaladejo J, Martín-Fernández JA. ZCompositions - R package for multivariate  
291 imputation of left-censored data under a compositional approach. *Chemom Intell Lab Syst*  
292 [Internet]. 2015;143:85–96. Available from: <http://dx.doi.org/10.1016/j.chemolab.2015.02.019>
- 293 16. Silverman JD, Washburne AD, Mukherjee S, David LA. A phylogenetic transform enhances  
294 analysis of compositional microbiota data. *Elife.* 2017;6:1–20.
- 295 17. Dillies MA, Rau A, Aubert J, Hennequet-Antier C, Jeanmougin M, Servant N, et al. A  
296 comprehensive evaluation of normalization methods for Illumina high-throughput RNA  
297 sequencing data analysis. *Brief Bioinform.* 2013;14(6):671–83.
- 298 18. Parks DH, Beiko RG. Measures of phylogenetic differentiation provide robust and  
299 complementary insights into microbial communities. *ISME J [Internet].* 2013;7(1):173–83.  
300 Available from: <http://dx.doi.org/10.1038/ismej.2012.88>
- 301 19. Oksanen J et al. *vegan: Community Ecology Package [Internet].* 2019. Available from:  
302 <https://cran.r-project.org/package=vegan>
- 303 20. McMurdie PJ, Holmes S. Phyloseq: An R Package for Reproducible Interactive Analysis and  
304 Graphics of Microbiome Census Data. *PLoS One.* 2013;
- 305 21. Anders S, Huber W. Differential expression analysis for sequence count data. *Genome Biol.*  
306 2010;11(R106).
- 307 22. Filzmoser P, Hron K, Reimann C. The bivariate statistical analysis of environmental  
308 (compositional) data. *Sci Total Environ [Internet].* 2010;408(19):4230–8. Available from:  
309 <http://dx.doi.org/10.1016/j.scitotenv.2010.05.011>
- 310 23. Aitchison J. *The Statistical Analysis of Compositional Data.* London: Chapman and Hall;  
311 1986.
- 312 24. Egozcue J, Pawłowsky Glahn V, Mateu-Figueras G, Barceló Vidal C. Isometric logratio for  
313 compositional data analysis. *Math Geol.* 2003;35(3):279–300.
- 314 25. Costea PI, Zeller G, Sunagawa S, Bork P. A fair comparison. *Nat Methods.* 2014;11(4):359.
- 315 26. Tsilimigras MCB, Fodor AA. *Annals of Epidemiology* Compositional data analysis of the



- 316 microbiome : fundamentals , tools , and challenges. *Ann Epidemiol* [Internet].  
317 2016;26(5):330–5. Available from: <http://dx.doi.org/10.1016/j.annepidem.2016.03.002>
- 318 27. Martín-Fernández JA, Barceló-Vidal C, Pawlowsky-Glahn V. Dealing with Zeros and Missing  
319 Values in Compositional Data Sets Using Nonparametric Imputation. *Math Geol*.  
320 2003;35(3):253–78.
- 321 28. Martín-Fernández JA, Hron K, Templ M, Filzmoser P, Palarea-Albaladejo J. Bayesian-  
322 multiplicative treatment of count zeros in compositional data sets. *Stat Modelling*.  
323 2015;15(2):134–58.
- 324 29. Jackson D. PROTEST: A PROcrustean Randomization TEST of community environment  
325 concordance. *Ecoscience*. 1995;2(3):297–303.
- 326 30. Peres-Neto PR, Jackson DA. How well do multivariate data sets match? The advantages of a  
327 procrustean superimposition approach over the Mantel test. *Oecologia*. 2001;129(2):169–78.
- 328 31. Wickham H. *ggplot2: Elegant Graphics for Data Analysis*. New York: Springer-Verlag; 2016.
- 329 32. Simonsohn U, Nelson LD, Simmons JP. P-curve: A key to the file-drawer. *J Exp Psychol Gen*.  
330 2014;143(2):534–47.
- 331 33. Prosser JI. Replicate or lie. *Environ Microbiol*. 2010;12(7):1806–10.
- 332 34. Lennon JT. Replication, lies and lesser-known truths regarding experimental design in  
333 environmental microbiology. *Environ Microbiol*. 2011;13(6):1383–6.
- 334 35. ten Hoopen P, Finn RD, Bongo LA, Corre E, Fosso B, Meyer F, et al. The metagenomic data  
335 life-cycle: Standards and best practices. *Gigascience*. 2017;6(8):1–11.
- 336 36. Costea PI, Zeller G, Sunagawa S, Pelletier E, Alberti A, Levenez F, et al. Towards standards  
337 for human fecal sample processing in metagenomic studies. *Nat Biotechnol* [Internet].  
338 2017;35(11):1069–76. Available from: <http://dx.doi.org/10.1038/nbt.3960>
- 339 37. Wiehlmann L, Pienkowska K, Hedtfeld S, Dorda M, Tümmler B. Impact of sample processing  
340 on human airways microbial metagenomes. *J Biotechnol* [Internet]. 2017 May 20 [cited 2018  
341 Jan 2];250:51–5. Available from:  
342 <http://linkinghub.elsevier.com/retrieve/pii/S0168165617300123>
- 343 38. Jones MB, Highlander SK, Anderson EL, Li W, Dayrit M, Klitgord N, et al. Library  
344 preparation methodology can influence genomic and functional predictions in human  
345 microbiome research. *Proc Natl Acad Sci U S A* [Internet]. 2015 Nov 10 [cited 2018 Jan  
346 2];112(45):14024–9. Available from:  
347 <http://www.pnas.org/lookup/doi/10.1073/pnas.1519288112>
- 348 39. Knudsen BE, Bergmark L, Munk P, Lukjancenko O, Priemé A, Aarestrup FM, et al. Impact of  
349 Sample Type and DNA Isolation Procedure on Genomic Inference of Microbiome  
350 Composition. Jansson JK, editor. *mSystems* [Internet]. 2016 Oct 25 [cited 2018 Jan  
351 2];1(5):e00095-16. Available from:  
352 <http://msystems.asm.org/lookup/doi/10.1128/mSystems.00095-16>
- 353 40. Jonsson V, Österlund T, Nerman O, Kristiansson E. Statistical evaluation of methods for  
354 identification of differentially abundant genes in comparative metagenomics. *BMC Genomics*

- 355 [Internet]. 2016;17(1):1–14. Available from: <http://dx.doi.org/10.1186/s12864-016-2386-y>
- 356 41. Thorsen J, Brejnrod A, Mortensen M, Rasmussen MA, Stokholm J, Al-Soud WA, et al. Large-  
357 scale benchmarking reveals false discoveries and count transformation sensitivity in 16S  
358 rRNA gene amplicon data analysis methods used in microbiome studies. *Microbiome*  
359 [Internet]. 2016 Dec 25 [cited 2019 Nov 28];4(1):62. Available from:  
360 <http://microbiomejournal.biomedcentral.com/articles/10.1186/s40168-016-0208-8>
- 361 42. Russel J, Thorsen J, Brejnrod AD, Bisgaard H, Sørensen SJ, Burmølle M. DAtest: a  
362 framework for choosing differential abundance or expression method. *bioRxiv* [Internet]. 2018  
363 Jan 2 [cited 2019 Nov 28];241802. Available from:  
364 <https://www.biorxiv.org/content/10.1101/241802v1>

365

## 366 11 Figure legends

367 **Figure 1: Heatmap visualizing the  $\beta$ -diversity to the reference sample using different strategies**  
368 **to normalize and transform data.** The  $\beta$ -diversity relative to the reference for each method was  
369 normalized according to the max value. Dendograms were created using complete linkage clustering  
370 of Euclidean distances. The rows in the heatmap represent different modifications to the reference,  
371 where RevRef represents reversing the reference, Sw represents switching, Ha represents high  
372 abundance, MHa represents medium high abundance, Ma represents medium abundance, La represents  
373 low abundance, NP represents not present, and SF represents scaling factor. The column labels in the  
374 heatmap contain extended explanations of zero estimation, where TSS represents total sum scaling,  
375 Rar represents rarefying, pa represents presence absence, CSS represents cumulative sum scaling, off  
376 represents an offset of zeroes, est represents a zero estimate using multiplicative simple replacement,  
377 ilr represents isometric log ratio transformation, clr represents centred log ratio transformation, and  
378 TMM represents trimmed mean of M-values.

379 **Figure 2: A: Principal coordinates analysis (PCoA) of the dissimilarity matrix containing**  
380 **pairwise comparisons of 1 - Procrustes correlations between methods, B: Density plot of**  
381 **correlations, C: Stress plot comparing observed dissimilarity to ordination distance in the**  
382 **PCoA and D: Scree plot of the percent of variation explained by the axes.** A dissimilarity matrix  
383 was created from all of the pairwise comparisons of metagenomics data analysis pipelines  
384 represented by one minus the Procrustes correlation. Redundancy analyses were performed  
385 unconstrained using the capscale function in vegan creating PCoA-, density, stress- and scree plot  
386 with ggplot2. Ellipses were added manually highlighting presence absence (pa), total sum scaling  
387 (TSS), and the compositional methods centred log ratio transformation (clr) and isometric log ratio  
388 transformation (ilr). In the processing (transformations and normalizations) legend, CSS represents  
389 cumulative sum scaling, and TMM represents trimmed mean of M-values. A small amount of jitter  
390 was added to distinguish clr and ilr.

391

## 392 1 Data Availability Statement

393 The datasets generated and analyzed for this study can be found at:  
394 <https://github.com/csapou/DataProcessinginMetagenomics>





