

# 1 How to normalize metatranscriptomic count 2 data for differential expression analysis

3 Heiner Klingenberg<sup>1</sup> and Peter Meinicke<sup>1</sup>

4 <sup>1</sup>Abteilung für Bioinformatik, Institut für Mikrobiologie und Genetik, Universität  
5 Göttingen, Goldschmidtstr. 1, 37077 Göttingen, Germany

6 Corresponding author:

7 Peter Meinicke<sup>1</sup>

8 Email address: peter@gobics.de

## 9 ABSTRACT

### 10 BACKGROUND

11 Differential expression analysis on the basis of RNA-Seq count data has become a standard tool in  
12 transcriptomics. Several studies have shown that prior normalization of the data is crucial for a reliable  
13 detection of transcriptional differences. Until now it is not clear whether and how the transcriptomic  
14 approach can be used for differential expression analysis in metatranscriptomics. The potential side  
15 effects that may result from direct application of transcriptomic tools to metatranscriptomic count data  
16 have not been studied so far.

### 17 METHODS

18 We propose a model for differential expression in metatranscriptomics that explicitly accounts for variations  
19 in the taxonomic composition of transcripts across different samples. As a main consequence the correct  
20 normalization of metatranscriptomic count data requires the taxonomic separation of the data into  
21 organism-specific bins. Then the taxon-specific scaling of organism profiles yields a valid normalization  
22 and allows to recombine the scaled profiles into a metatranscriptomic count matrix. This matrix can  
23 then be analyzed with statistical tools for transcriptomic count data. For taxon-specific scaling and  
24 recombination of scaled counts we provide a simple R script.

### 25 RESULTS

26 When applying transcriptomic tools for differential expression analysis directly to metatranscriptomic data  
27 the organism-independent (global) scaling of counts implies a high risk of falsely predicted functional  
28 differences. In simulation studies we show that incorrect normalization not only tends to loose significant  
29 differences but especially can produce a large number of false positives. In contrast, taxon-specific  
30 scaling can equalize the variation of relative library sizes from different organisms and therefore shows  
31 a reliable detection of significant differences in all simulations. On real metatranscriptomic data the  
32 results from taxon-specific and global scaling can largely differ. In our study, global scaling shows a high  
33 number of extra predictions which are not supported by single transcriptome analyses. Inspection of the  
34 scaling error suggests that these extra predictions may actually correspond to artifacts of an incorrect  
35 normalization.

### 36 CONCLUSIONS

37 As in transcriptomics, a proper normalization of count data is also essential for differential expression  
38 analysis in metatranscriptomics. Our model implies a taxon-specific scaling of counts for normalization  
39 of the data. The application of taxon-specific scaling consequently removes taxonomic composition  
40 variations from functional profiles and therefore effectively prevents the risk of false predictions due to  
41 incorrect normalization.

## 42 BACKGROUND

43 Metagenome analysis can provide a comprehensive view on the metabolic potential of a microbial  
44 community (Eisen, 2007; Simon and Daniel, 2009). In addition to the static functional profile of the  
45 metagenome, metatranscriptomic RNA sequencing (RNA-Seq) can highlight the multi-organism dynamics  
46 in terms of the corresponding expression profiles (Poretsky et al., 2005; Frias-Lopez et al., 2008; Gilbert  
47 et al., 2008; Urich et al., 2008). In particular, metatranscriptomics makes it possible to investigate the  
48 functional response of the community to environmental changes (Gilbert et al., 2008; Poretsky et al.,  
49 2009).

50 In single organism transcriptome studies, differential expression analysis based on RNA-Seq data  
51 has become an established tool (Marioni et al., 2008; Trapnell et al., 2012). For the analysis, first,  
52 quality-checked sequence reads are mapped to the organisms genome for transcript identification. Then  
53 the transcript counts are compared between different experimental conditions to identify statistically  
54 significant differences. Several studies have shown that read count normalization has a great impact on  
55 the detection of significant differences (Bullard et al., 2010; Dillies et al., 2013; Lin et al., 2016). The  
56 aim of the count normalization is to make the expression levels comparable across different samples  
57 and conditions. This is an essential prerequisite for distinguishing condition-dependent differences from  
58 spurious variation of expression levels.

59 In metatranscriptomics, already the transcript identification step can be challenging. In many cases,  
60 RNA-Seq faces a mixture of organisms for which no reference genome sequence is available. Several  
61 strategies have been suggested: de novo transcriptome assembly combined with successive homology-  
62 based annotation (Celaj et al., 2014), the direct functional annotation of reads by classification according  
63 to some protein database (Huson et al., 2011; Nacke et al., 2014; Hesse et al., 2015) or parallel sequencing  
64 of the corresponding metagenome with successive mapping of RNA-Seq reads to assembled and annotated  
65 contigs (Mason et al., 2012; Franzosa et al., 2014; Ye and Tang, 2016). For the subsequent comparison  
66 of counts between different conditions no standard protocol exists for differential expression analysis  
67 on metatranscriptomic data. Several studies and tools apply methods that have been developed for  
68 differential expression analysis in transcriptomics to metatranscriptomic count data (McNulty et al., 2013;  
69 Martinez et al., 2016; Macklaim et al., 2013). However, the question under which conditions established  
70 models from single organism transcriptomics also apply to organism communities has not been addressed  
71 sufficiently so far.

72 Here we present an extended statistical model for count data from metatranscriptomic RNA-Seq  
73 experiments. Theoretical considerations as well as studies on simulated and real count data show that  
74 correct normalization of the data is crucial and in general requires an organism-specific rescaling of  
75 expression profiles. This implies that a valid differential analysis should only include data that can be  
76 attributed to single organisms. The application of differential expression analysis to mixed-species data  
77 without prior separation can be found in several metatranscriptomic studies (Nacke et al., 2014; McNulty  
78 et al., 2013; De Filippis et al., 2016) as well as in dedicated pipelines for metatranscriptome analysis  
79 (Martinez et al., 2016; Westreich et al., 2016). Our results suggest that inadequate normalization of  
80 metatranscriptomic count data always bears the risk of serious errors in differential expression analysis  
81 and should be avoided consequently.

## 82 NORMALIZATION

83 To clarify our arguments for an alternative normalization of metatranscriptomic data we need to explain  
84 the statistical nature of the normalization problem. We first follow the approach of Anders and Huber  
85 (Anders and Huber, 2010) for single organism RNA-Seq count data and start with a basic model for the  
86 mean of the observed counts. The expected (mean) count  $\mathbb{E}[Y_{ij}]$  for gene (feature)  $i$  and sample  $j$  arises  
87 from a product of the per-gene quantity  $\lambda_{ic_j}$  under condition  $c_j$  and a size factor  $s_j$ :

$$\mathbb{E}[Y_{ij}] = \lambda_{ic_j} s_j \quad (1)$$

88 The factor  $\lambda_{ic_j}$  is proportional to the mean concentration of feature  $i$  under condition  $c_j$ . The size factor  $s_j$   
89 represents the sampling depth or library size. Usually, both factors are unknown. If we assume the  $i$ -th  
90 feature to be non-differentially expressed (NDE) we can represent the corresponding row of the count  
91 matrix by

$$\mathbb{E} \left[ Y_{i\bullet}^{(\text{NDE})} \right] = \lambda_i \mathbf{s} \quad (2)$$

92 where the relative feature abundance is equal for all samples and all size factors have been comprised in  
 93 the row vector  $\mathbf{s}$ . Thus, for NDE features the size factors are proportional to the expected counts. If we  
 94 knew which features are actually NDE, we would be able to estimate the required size (scaling) factors  
 95 for normalization from the corresponding counts.

96 With the common choice of a mean scaling factor of 1 the scaling factors can be estimated by

$$\hat{\mathbf{s}} = \frac{Y_{i\bullet}^{(\text{NDE})}}{\frac{1}{n} \sum_j Y_{ij}^{(\text{NDE})}} \quad (3)$$

97 where  $n$  is the number of samples. The denominator in the above equation corresponds to the arithmetic  
 98 mean of the counts for feature  $i$ . If the library sizes of the samples strongly diverge, possibly by orders of  
 99 magnitude, the geometric mean is more suitable (Anders and Huber, 2010).

100 To make the expected counts of different samples comparable, i.e. in order to compare the feature  
 101 concentrations, the columns of the data matrix are divided by the sample-specific scaling factors prior to  
 102 the testing for significant differences. As above, it is common to choose an average scaling factor 1. Thus,  
 103 if we would actually know an NDE feature beforehand, in principle, we could use it to estimate the scaling  
 104 factors. Usually this is not the case and we need to make some assumptions. A common assumption that  
 105 is used in current tools is that most of the features are NDE. Then it is possible to estimate the scaling  
 106 factors by some robust statistics. In DESeq for each sample the putative scaling factors from all features  
 107 are calculated and then the median of all these values is used as an estimator of the sample-specific scaling  
 108 factor (Anders and Huber, 2010). The median is highly robust, with a breakdown point of 50% and  
 109 therefore the estimator can be used if at least half of the data corresponds to NDE features. Without any  
 110 distinction between DE and NDE features the scaling factors have also been estimated from the count  
 111 sums of all samples. However, the potential shortcomings of this total count normalization have widely  
 112 been discussed (Anders and Huber, 2010; Robinson and Oshlack, 2010; Sonesson and Delorenzi, 2013).

113 In metatranscriptomics the situation is more complicated because for each organism we can have a  
 114 different scaling factor. So we have to extend the above sampling model to an  $N$ -organism mixture that  
 115 includes a matrix  $\mathbf{S}$  of organism-specific scaling factors  $s_{jk}$ :

$$\mathbb{E} [Y_{ij}] = \sum_{k=1}^N \lambda_{ijk} s_{jk} \quad (4)$$

116 where  $i, j, k$  are the feature, sample and organism indices, respectively. We omitted the condition dependency  
 117 ( $c_j$ ) for convenience.

118 In analogy to equation (2) for NDE features we have the following model for a feature row  $i$  of the  
 119 count matrix:

$$\mathbb{E} \left[ Y_{i\bullet}^{(\text{NDE})} \right] = \lambda_i^T \mathbf{S}^T \quad (5)$$

120 where the column vector  $\lambda_i$  contains all organism-specific rates for feature  $i$  and  $\lambda_i^T$  indicates transposition  
 121 of this vector.

122 Application of the above single-organism scheme for estimation of scaling factors is only valid if the  
 123 matrix of scaling factors has the following form:

$$\mathbf{S} = [\alpha_1 \mathbf{s}, \alpha_2 \mathbf{s}, \dots, \alpha_K \mathbf{s}] = \boldsymbol{\alpha}^T \quad (6)$$

124 where  $\boldsymbol{\alpha}$  is a column vector of organism-specific abundances and  $\mathbf{s}$  contains the sample-specific scaling  
 125 factors, now in a column vector, which is equal for all organisms. Then we can write

$$\mathbf{S}\lambda_i = \mathbf{s}\alpha^T \lambda_i = \tilde{\lambda}_i \mathbf{s} \quad (7)$$

126 where  $\tilde{\lambda}_i$  results from the dot product of the organism and the feature rates. This corresponds to equation  
127 (2) and allows to apply DESeq or other tools for single organism differential expression analysis to the  
128 metatranscriptomic count matrix. However, the underlying assumption that  $\mathbf{S}$  has column rank 1, i.e. all  
129 column vectors are collinear, would be hard to justify in practice. Implicitly we would assume that the  
130 relative contributions of all organisms are constant across all samples. In general, this assumption is not  
131 met, because for a real metatranscriptome, the organism composition of transcripts cannot be controlled  
132 and will be different for different samples. In our approach to normalization of metatranscriptomic  
133 counts we preprocess the data according to an organism-specific rescaling of separated counts so that the  
134 recombined count data actually meet the former assumption.

## 135 MATERIALS & METHODS

### 136 Taxon-specific scaling and global scaling

137 We propose a method to prepare metatranscriptomic data for differential expression analysis. The  
138 method is referred to as taxon-specific scaling. As an essential prerequisite our approach requires that  
139 the data is first partitioned according to the contributing organisms. Then the count data matrix from  
140 each partition is normalized separately. Here, established tools from transcriptomics can be used to  
141 estimate the corresponding scaling factors. Finally, the normalized count data matrices are summed  
142 up to provide normalized metatranscriptomic count data which can be analyzed in terms of differential  
143 expression (Fig. 1). Here all statistical models and tools for count-based differential expression analysis  
144 in transcriptomics can in principle be used to identify differentially expressed features.

145 If we denote the original count matrix for organism  $k$  as  $\mathbf{Y}_k$  and the associated vector of estimated  
146 scaling factors as  $\hat{\mathbf{s}}_k$  the normalized metatranscriptomic count matrix is computed by

$$\tilde{\mathbf{Y}} = \sum_k \mathbf{Y}_k \text{diag}^{-1}(\hat{\mathbf{s}}_k) \quad (8)$$

147 Here, the  $\text{diag}^{-1}$  operator transforms the scaling vector to a diagonal matrix with inverse scaling factors  
148 on the diagonal and zeros everywhere else. We provide an R script where we use DESeq2 for scaling  
149 factor estimation and identification of significant differences (see Additional File 1).

150 In principle, our method is computationally simple and the hard work has to be done beforehand in  
151 order to provide the partitioned data in terms of the organism-specific count matrices. This is the realm of  
152 binning methods and, in addition, may require sequence assembly tools to achieve a sufficient sequence  
153 length for reliable separation.

154 At this point, the question may arise why to get back to metatranscriptomic data when differential  
155 expression analysis could be performed for separate organisms or specific taxa. There are several reasons  
156 why the analysis of the recombined metatranscriptome data can be useful: first of all, the statistical power  
157 of organism-specific tests may be low due to decreased counts. If several organisms show the same slight  
158 difference, this difference may only become statistically significant when accumulating their normalized  
159 counts. Or a feature may show differences for single organisms but these differences may cancel out  
160 when correctly summarized. In this case the corresponding feature is not indicative for the experimental  
161 condition with regard to the whole community. Therefore the analysis of separate organism transcriptomes  
162 and the analysis of the rectified metatranscriptome data should be combined to provide a complete picture  
163 of the community response.

164 In our study and in the supplied R script we use DESeq2 to compute scaling factors and to identify  
165 significant differences on the basis of the normalized count matrices. We decided for DESeq2 for several  
166 reasons. It is an established tool in transcriptomics which has shown a good performance in compar-  
167 ative studies (Soneson and Delorenzi, 2013; Dillies et al., 2013) and which has already been used for  
168 metatranscriptome analysis (McNulty et al., 2013; Martinez et al., 2016; De Filippis et al., 2016). In  
169 particular, the estimation of scaling factors is robust and can be performed as a separate prior step apart  
170 from the computation of significant differences. The latter aspect is important for taxon-specific scaling  
171 which requires to apply the normalization independently. However, we would like to emphasize that our

172 arguments for the taxon-specific scaling approach do not depend on a particular statistical tool and in fact  
173 the main findings of our study can be reproduced with other tools, such as edgeR (Robinson et al., 2010),  
174 SAMseq(Li and Tibshirani, 2013) or limma(Ritchie et al., 2015). In some experiments we also used  
175 edgeR and total count (TC) normalization to study the impact of different transcriptomic scaling methods.

176 In contrast to taxon-specific scaling, global scaling performs the normalization of metatranscriptomic  
177 data without prior separation, i.e. sample-specific scaling factors are estimated from the original metatran-  
178 scriptomic counts. In general, taxon-specific and global scaling will result in distinct normalized count  
179 matrices which in turn can lead to largely differing results in differential expression analysis. We tried to  
180 show this on simulated and real count data as described in the following.

### 181 **Synthetic data generation and analysis tools**

182 The tool `compcodeR` (Soneson, 2014) was used to generate all simulated data. The tool generates count  
183 data based on a negative binomial distribution model with parameters estimated from real transcriptome  
184 data (Pickrell et al., 2010; Cheung et al., 2010). If not explicitly specified, the `compcodeR` parameters in  
185 the R function “`generate.org.mat`” are used (see Additional File 1). All analyses were performed with R  
186 version 3.3.0 and DESeq2 version 1.8.2, edgeR version 3.10.5.

### 187 **Simulated metatranscriptome**

188 A metatranscriptome arises from a mixture of various organisms, each with individual features. As a  
189 result, a metatranscriptome can include features covered by all taxa as well as features occurring only  
190 in few or a single organism. Generally, the count contributions from different organisms are not equal  
191 and vary across samples. We refer to this as the variation of the library size. Therefore, `compcodeR`  
192 was used to generate multiple data sets with different total count numbers to simulate the variation of  
193 organism-specific library sizes. Thereby, each generated data set mimics the contribution of a single  
194 organism. The data sets were then combined to simulate a metatranscriptomic count matrix.

195 As with all simulations, the data can only provide a coarse approximation of real metatranscriptomic  
196 counts which depends on particular parameters. Therefore, settings for the number of features and the  
197 number of total counts influence the results. Each organism is simulated with 100 differentially expressed  
198 features (DEF), 50 of them upregulated, and with 900 features that were non-differentially expressed  
199 (NDE).

200 Each data set consists of two conditions, A and B, with six samples (replicates) per condition and  
201 five organisms (Org1 to Org5) per sample. In the first three simulations, the different organism profiles  
202 are stacked, to exclude any interference between features from different organisms in the combined data.  
203 Accordingly, the final count matrix has 12 columns and 5000 rows that correspond to samples and features,  
204 respectively. The data generation process provides the necessary information to calculate the number  
205 of true positives (TP) and false positives (FP). The label  $L_i$  is DE or NDE according to feature  $i$  being  
206 differentially expressed or non-differentially expressed. The statistical test used to detect DEF, provides a  
207 p-value for each feature. The predicted label  $\hat{L}_i$  is DE if the adjusted p-value (Benjamini and Hochberg,  
208 1995) is below a threshold of 0.05 for feature  $i$ . The TP and FP counts are calculated for each organism  $k$   
209 individually.

$$210 \quad TP_k = |\{i : \hat{L}_i = DE \wedge L_i = DE\}| \quad (9)$$

$$211 \quad FP_k = |\{i : \hat{L}_i = DE \wedge L_i = NDE\}| \quad (10)$$

### 212 **Simulation I: “Without library size variation”**

213 In the first experiment we simulate the case where the library size (LS) for each of the five organisms does  
214 not vary across different samples. Although this is an unrealistic case we performed this simulation  
215 to verify that both normalization approaches work equally well under ideal conditions. In addition, we  
216 wanted to investigate how different organism abundances affect the identification of DEF. Each organism  
217 was assigned a fixed total count number across all samples, without variation in library size. We simulated  
218 organism Org1 with a base count of  $1e^7$  followed by organism Org2 with  $5e^6$ ,  $1e^6$  for organism Org3  
219 and organism Org4, Org5 with  $5e^5$ ,  $1e^5$  respectively. Because data is generated without variation for the  
number of counts per sample, no normalization is required, i.e. the correct scaling factors for all samples  
and all organisms are the same (= 1).

220 **Simulation II: “With library size variation”**

221 In the second experiment we simulated a more realistic situation, with varying LS for all included  
222 organisms.

223 Organism base counts are identical to the first simulation but the LS is randomly increased or reduced  
224 according to a random factor between 0.5 and 2. Due to the different library sizes of the samples, a prior  
225 normalization is required.

226 **Simulation III: “Condition dependent variation”**

227 In the third simulation, we investigated to what extent a condition dependent variation of LS can affect  
228 the normalization results. Under condition A we increase LS of Org1 by a random factor between 1.5 and  
229 2 while under condition B we decrease the LS by a random factor between 0.5 and 0.667. For Org2 the  
230 direction of change is reversed, with a random decrease under condition A and an increase for condition  
231 B. For Org1 and Org2 the same base count as for Simulation I and II is used and for Org3-5 all parameters  
232 from Simulation II are used.

233 **Simulation IV: “Mixed feature effects”**

234 In the previous simulations I-III, the generated count matrices from different organisms are stacked in  
235 the combined count matrix. By stacking the organism profiles, the provided feature labels from the data  
236 generation process can be used to validate the predictions. Now we want to analyze which effects can  
237 be observed if each feature can accumulate counts from different organisms. This case is common for  
238 real metatranscriptomic counts which arise from a mixture of organisms. Of particular interest are two  
239 effects that we refer to as “cancellation” and “boosting”. To simplify the analysis of these effects we  
240 restricted our simulation to a mixture of two organisms, which had the same base counts as Org1 and  
241 Org2 in Simulation I. The cancellation effect is observed when DEF that are significant in one or both  
242 organism datasets lose significance in the mixture. In contrast, the boosting effect is observed for DEF  
243 which are only significant in the combined dataset. We generated data for three samples per condition, to  
244 limit the variability of the superimposed count data. The first 100 features of the organism profiles were  
245 DE while the remaining 900 were NDE. The simulation is divided into part A, where we superimposed  
246 features with the same DE direction and part B, where the corresponding DE features of Org1 and Org2  
247 had opposite directions.

248 Note that the aim of Simulation IV was not to compare the different normalization approaches but  
249 instead to demonstrate the possible effects that may result from mixed organism count data. However,  
250 the simulation cannot be used to draw conclusions about the frequencies of the effects for real data. In  
251 particular, we expect the boosting effect to be much stronger for real data where organisms with a similar  
252 response may provide correlated features that can emphasize trends or differences between conditions  
253 when superimposing their counts.

254 **Part A** For this part of the simulation, we superimposed equally directed features of the two organisms.  
255 With 100 features selected as DE, the first 50 are “upregulated” followed by 50 “downregulated” features  
256 and 900 NDE features. This simulation was expected to show the boosting effect as well as the cancellation  
257 effect.

258 **Part B** In part B we tried to further increase the frequency of the cancellation effect. An important  
259 aspect of identifying DEF is the difference between the mean count values of the two conditions. To  
260 bring the mean count values of the mixture for the two conditions closer together, we added the sorted  
261 “upregulated” features of Org1 and the sorted “downregulated” features of Org2 and vice versa. The true  
262 mean values available from the data generation process are the basis for the sorting. For the generated  
263 DEF the sorting ensures that high count values in condition A from one organism are balanced with high  
264 count values in condition B from the other organism.

265 **Simulation V: “False positive control”**

266 In the final simulation, the aim was to investigate the effect of an LS shift within mixed count data  
267 without any DE features. Here we particularly wanted to measure the impact of the normalization on the  
268 false positive rate. This kind of analysis has also been proposed for transcriptomics to validate the false  
269 discovery control in differential expression analysis tools (Soneson and Delorenzi, 2013).

270 We used the parameters from Simulation III, but instead of six samples per condition, compcodeR  
271 generated 12 samples per condition where we only used the samples for the first condition. For Org3-5

272 we use the LS variation as used in Simulation III. For Org1 and Org2 the LS shift between different  
273 conditions in Simulation III is now applied between samples 1-6 (condition A) and 7-12 (condition B).  
274 The LS ranges were identical to those in Simulation III. The results of the differential expression analysis  
275 for global and taxon-specific scaling on the superimposed data were compared.

#### 276 **Error calculation of the scaling factor**

The scaling error  $E_k$  was estimated from the difference between the sample-specific scaling factors  $\hat{s}_{jk}$  and the actual "true" scaling factors  $s_{jk}$  for each organism  $k$  as provided by the simulation parameters. In both cases the factors are scaled to provide a unit mean across samples. To obtain a scaling error between 0 and 1, we compute the error by:

$$E_k = \frac{\sum_j |\hat{s}_{jk} - s_{jk}|}{2n} \quad (11)$$

277 where  $n$  is the number of samples. In addition, we used the logarithmic measure  $\log_2 \hat{s}_{jk}/s_{jk}$  to represent  
278 the directed error.

#### 279 **Metatranscriptome data**

280 For a real data study, we chose a metatranscriptome dataset from mice gut (McNulty et al., 2013). The  
281 experiment includes 12 different species (see Additional File 2 : Tab. 1) representing an artificial human  
282 gut microbial community which was inserted into germ free mice. In the original study the diet for  
283 the mice was changed at different time points. Metatranscriptomic data is available for 6 time points  
284 which provide the conditions for our analysis. The available processed count data was obtained from the  
285 European Bioinformatics Institute [<http://www.ebi.ac.uk/>, ArrayExpress, E-GEOD-48993] and contains  
286 gene names and the associated numbers.

287 Because the gene to Pfam (Finn et al., 2014) mapping is available for most organisms, we selected  
288 Pfam protein domains as features for the differential expression analysis. Each Pfam domain family is a  
289 feature in the resulting vector, including only Pfams observed at least once. We transformed the available  
290 RPKM values for the genes back to raw counts. For genes with multiple Pfam annotations, we add the  
291 raw count values of the gene to all associated Pfam features. From the available data, we constructed  
292 a count matrix for each condition and organism (Additional File 3). Here, each column constitutes a  
293 different sample and each row represents a particular feature. Because the count data from *Bacteroides*  
294 *cellulosilyticus* WH2 did not map to gene names, all related counts are excluded from the analysis.

295 A differential expression analysis for all pairwise combinations of distinguished conditions was  
296 performed to compare the results of global and taxon-specific scaling. We calculated the number of DEF  
297 predicted a) with both methods with the same fold change direction, b) with both methods but with an  
298 opposite fold change direction and c) with only one scaling method. In addition, we investigated the  
299 overlap between the single organism transcriptome analyses and the differential expression analysis for  
300 the mixture. We applied a significance threshold on the adjusted p-value of 0.05 for the prediction of DEF.

## 301 **RESULTS & DISCUSSION**

302 In the first part of our evaluation we examined the performance of taxon-specific and global scaling  
303 methods on simulated data. Because simulations I to III had been designed to provide a clear ground truth  
304 we were able to distinguish true positive predictions of DEF from falsely classified features. In the second  
305 part we show results on real metatranscriptomic count data. Here the ground truth is not known and  
306 therefore we restrict the analysis on the comparison of the results from the two normalization approaches.  
307 Because it is impossible to verify the correctness of predictions we focus on analyzing the agreement or  
308 disagreement on DEF detection in this case.

### 309 **Simulation I**

310 In this experiment, we measured the ability to detect DEF in a metatranscriptome without variation of  
311 organism-specific library sizes across different samples. This situation, in principle, does not require any  
312 normalization and therefore we expected taxon-specific ("tax") and global ("glo") scaling to yield similar  
313 results. This is confirmed by the resulting true positive (TP) predictions of DEF for the included organisms  
314 (Fig. 2). For both approaches the number of true positives is higher for more abundant organisms due to

315 an increased statistical power of the corresponding tests. The final profile includes 100 DE and 900 NDE  
316 features for each organism, resulting in 5000 features in total.

317 We repeated the analysis with edgeR and TC normalizations to estimate the scaling factors (see  
318 Additional File 2 : Fig. 1). For edgeR, the number of correctly identified DEF is lower for all organisms  
319 (see Additional File 2 : Fig. 1). For this particular data set, the library size (LS) was correctly adjusted by  
320 DESeq2 with scaling factors close to 1 for all samples. Again both normalization approaches performed  
321 equally well. A similar picture can be expected for a varying total LS of the metatranscriptome samples  
322 as long as the relative LS of the organisms does not vary across different samples.

### 323 **Simulation II**

324 When introducing organism-specific LS variation across samples the picture changes. For the global  
325 scaling approach the results show a decrease in the average TP rate for all organisms (see Fig. 3). This  
326 trend is also visible when edgeR and TC normalization are used for differential expression analysis (see  
327 Additional File 2 : Fig 2). On the other hand, with taxon-specific scaling the results are very similar to  
328 results from Simulation I (see Fig. 2). With this method more DEF are correctly identified than with  
329 global scaling. The difference in the number of correctly identified DEF for global scaling is dependent  
330 on the parameter settings for the LS variation. For a lower amplitude of the LS variation, the TP rate for  
331 global scaling increases (Additional File 2 : Fig. 3). The range of TP for the most abundant organisms  
332 Org1 and Org2 is broader with global scaling (see Fig. 3) which also shows a higher scaling error (SE) for  
333 organisms with a lower sequencing depth (see Additional File 2 : Fig. 4).

334 The receiver operating characteristics (ROC) shows a higher area under curve (AUC) value for taxon-  
335 specific scaling (0.8776) than for global scaling (0.8282). The curve for global scaling also shows a higher  
336 degree of variation across different simulation runs (Fig. 4 dotted lines).

### 337 **Simulation III**

338 With the inclusion of a condition dependent variation of the LS this simulation experiment can be viewed  
339 as a worst case study. For global scaling, the observed number of true positives is higher for all data sets  
340 compared to Simulation II (see Fig. 5). However, the number of false positive predictions explodes and  
341 even exceeds the total number of DEF (500) resulting in average TP and FP numbers of 228 ( $\pm 11$ ) and  
342 1523 ( $\pm 78$ ), i.e.  $\sim 35\%$  of all features are predicted to be DEF.

343 In particular, the biggest portion of FP accumulates in features from Org1 and Org2 (see Fig. 5).  
344 Inspecting the log<sub>2</sub> fold changes (see Fig. 6) a shift from the correct center of 0 upwards and downwards  
345 can be observed for Org1 and Org2, respectively. As a result, many DEF are identified with a wrong  
346 (opposite) direction and most of the false positive detections just reflect the direction of this shift. This  
347 situation implies a total loss of control over the false discovery rate. The results with edgeR and TC  
348 normalization show a similarly high FP rate (see Additional File 2 : Fig. 5).

349 As a consequence, the ROC curve collapses for global scaling (see Fig. 7) with an AUC of 0.6396.  
350 In contrast, taxon-specific scaling does not suffer from condition-dependent LS variation and the results  
351 compare well with those of Simulation I & II showing a similar shape of the ROC curve (AUC: 0.8785).  
352 With taxon-specific scaling, the average TP across all species is 237 which corresponds to a sensitivity of  
353  $\sim 47\%$ . For global scaling the total number of predicted DEF (TP + FP) is dependent on the amplitude of  
354 the condition dependent LS shift and increases for bigger shifts.

### 355 **Simulation IV**

356 With this simulation we wanted to analyze the effects that result from the superposition of counts from  
357 different species. In this case we do not compare the two normalization approaches. DEF that can not be  
358 detected for each organism separately may be identified as DE in the mixture (boosting effect) and DEF  
359 that can be detected for single organisms may disappear in the mixture (cancellation effect). Here we  
360 analyzed the frequencies of the different effects for the features that were labeled DE according to the  
361 data generation process. The presented numbers result from averaging over 100 iterations and for the  
362 observed effects these numbers sum up to the total number of DE-labeled features (100).

363 **Part A** First, we just added the data matrices of the two simulated organisms, i.e. the upregulated and  
364 downregulated feature counts are summed up. In the results, boosting as well as cancellation effects can  
365 be observed. The chance to observe the boosting effect is low because the features of the two organisms  
366 and conditions are not correlated. Furthermore, when counts with different orders of magnitude are added,



367 only slight changes in the mean counts can be observed. As a result, the boosting and cancellation effects  
368 only show median frequencies of 5 and 9, respectively (see Fig. 8). The main portion of features identified  
369 as DE in the mixture is also identified as DE for at least one of the organisms, followed by the fraction of  
370 features that are insignificant in all cases (see Fig. 8). Increasing the number of samples per condition  
371 from 3 to 6 further reduces the median number of boosted features (3) while also increasing the overall  
372 ability to correctly identify DE features (see Additional File 2 : Fig. 6).

373 **Part B** In the second part of the simulation, we changed the order of the features to intensify the  
374 cancellation effect. This was achieved by adding sorted upregulated features of one organism to sorted  
375 downregulated features of the other organism. The median number of features identified as DE for both  
376 cases, mixture and single organism analysis, is reduced to 7 and the boosting effect almost completely  
377 disappears (see Fig. 8). Due to our simulation setup, features identified as DE in one of the organisms but  
378 not identified in the mixture were the second most frequent (see Fig. 8) while features predicted as NDE  
379 for both organisms and for the mixture were most frequent.

380 Increasing the number of samples per condition from three to six resulted in a stronger cancellation  
381 effect and a higher number of features identified as DE for one of the organisms and in the mixture (see  
382 Additional File 2 : Fig. 6). Again, the increased number of samples improves the overall detection of  
383 DEF.

### 384 **Simulation V**

385 In the final simulation, we wanted to investigate the effect of the scaling on the false discovery rate  
386 for mixed organism count data. Therefore the data matrices (12 samples and 1000 features) for this  
387 experiment were generated without any DEF.

388 For taxon-specific scaling, the number of significant features is negligible with an average number of  
389 1.6, which is well within the range of the estimated false discovery rate (FDR). In contrast, for global  
390 scaling the average number of predicted DEF is  $628 \pm 22$ . While with taxon-specific scaling there is  
391 little if any significant difference, global scaling predicted  $\sim 63\%$  of all features to be DE. Only a small  
392 number of predicted DEF results from analysis of the organism transcriptomes with a mean of 1.2, again  
393 well within the FDR range for an adjusted p-value threshold of 0.05. The results from Simulation IV  
394 show that the prevalent effect of the superposition of uncorrelated data is the cancellation effect and the  
395 boosting effect could only be observed in a few cases. However, in Simulation IV we focused on DEF  
396 that were marked as differentially expressed by the data generation process. Therefore we can exclude  
397 that the high number of DEF predicted by global scaling, results from the boosting effect.

### 398 **Real metatranscriptome data**

399 While the objective of the simulation studies was to evaluate and compare the two normalization ap-  
400 proaches in terms of correctly identified DEF, we do not have a ground truth for the analysis of the real  
401 data. Therefore we focused on an analysis of the (dis-)agreement in results between both approaches  
402 without assessing the actual detection performance. The analyzed data comprises Pfam counts from 11  
403 organisms, 6 conditions according to different time points and 4 replicates per condition. An overview on  
404 predicted DEF in all pairwise condition comparisons is shown in Fig. 9.

405 For both approaches, the number of DEF peaked at “day 13” vs. “day 27” with 512 and 756 significant  
406 features for taxon-specific and global scaling. The number of DEF was low when conditions close together  
407 on the time line were compared (“day 15” vs. “day 16” or “day 29” vs. “day 30”).

408 For global scaling, the number of predicted DEF was generally higher than for taxon-specific scaling.  
409 For some of the comparisons, the number of extra predictions under global scaling was even higher  
410 than the number of shared predictions (see Fig. 9). The high number of extra predictions observed with  
411 global scaling is especially prevalent for the comparisons “day 15” vs “day 16” with 16 times more extra  
412 predictions than predictions shared with taxon-specific scaling and “day 13” vs “day 16” with 3 times  
413 more extra predictions than shared predictions.

414 We also compared the results of the mixture analysis for global scaling and taxon-specific scaling  
415 to the transcriptome analyses of the individual organisms. For 10 of the 15 comparisons, the majority  
416 of features predicted as DE with global scaling were not predicted as DE in any transcriptome (see  
417 Additional File 2 : Fig. 7). In contrast, with taxon-specific scaling only the comparisons “day 29” vs “day  
418 30” and “day 15” vs “day 16” show a higher fraction of significant features not predicted as DE in the  
419 transcriptomes. These two comparisons are also the ones with the smallest total number of predicted DEF.

420 When taking into account the direction of the differential expression, the number of DEF predicted by  
421 both methods but with contrary regulation direction is low. Here comparison “day 16” vs “day 30” shows  
422 the highest number of significant features with an opposite direction (5).

#### 423 **Analysis of “day 13” vs “day 27”**

424 As described in the original study, “day 13” and “day 27” each correspond to the final day of a particular  
425 diet. Because the number of predicted DEF for both scaling methods (global and taxon-specific) was the  
426 highest here, we analyze the results for this comparison in more detail.

427 We mapped the Pfam-annotated features which were predicted as DE for global scaling and taxon-  
428 specific scaling to Gene Ontology (GO) terms and compared the results. With taxon-specific scaling 250  
429 GO terms with at least one DEF mapping were identified, while global scaling resulted in 311 GO terms.  
430 GO terms associated to biological processes with a high agreement between the two methods were for  
431 example cellular amino acid metabolic process where both methods identified 7 of the 9 associated Pfams  
432 as DE and carbohydrate metabolic process with 11 DEF shared between taxon-specific scaling and global  
433 scaling. In this category, taxon-specific scaling and global scaling predicted 5 additional DEF uniquely.

434 GO terms with predicted DEF from taxon-specific scaling alone included magnesium ion binding  
435 and fucose metabolic process with 3 predicted DEF each. For global scaling alone, we found DNA  
436 modification, molybdopterin cofactor biosynthetic process, and RNA modification with 3, 2 and 2  
437 predicted DEF respectively (see Additional File 4 for a complete list).

438 **Extra predictions** In the condition comparison “day 13” vs “day 27”, both normalization approaches  
439 shared 376 features predicted as DE. With taxon-specific scaling, 136 extra predictions were observable  
440 while global scaling resulted in 380 extra predictions. When comparing the results of both scaling methods  
441 to the single organism transcriptome analyses, global scaling and taxon-specific scaling resulted in 252  
442 and 69 predictions that were insignificant in all single analyses (see Fig. 9). Both methods shared an  
443 overlap of 53 DEF that were not detected in any of the transcriptome analyses. Thus, global scaling would  
444 suggest a boosting effect for  $\sim 13\%$  of the features that were insignificant in all transcriptome analyses. In  
445 contrast, with taxon-specific scaling, a putative boosting effect for only  $\sim 4\%$  of the transcriptomic NDE  
446 features can be observed. In the single organism transcriptomes, a total of 1302 features were predicted to  
447 be DEF at least once. Both methods lead to a similar cancellation rate of  $\sim 66\%$  for taxon-specific scaling  
448 and  $\sim 61\%$  for global scaling.

449 The fraction of shared DEF predictions between the two scaling methods is lower if the DEF are  
450 supported by a smaller number of transcriptome analyses. For features supported by one transcriptome the  
451 agreement was  $\sim 48\%$ , increasing to  $\sim 56\%$  for two transcriptomes. In the range of three to six supporting  
452 transcriptomes, the agreement increases to  $\sim 59\%$ ,  $\sim 80\%$ ,  $100\%$  and  $100\%$  respectively. While the relative  
453 agreement between taxon-specific and global scaling increases, the total number of features supported by  
454 multiple transcriptomes decreases (Additional File 2 : Fig. 7).

455 **Scaling error** In the simulations, the differences in the estimated scaling factors for the single organism  
456 profiles in comparison to the actual scaling factors were low with a scaling error of  $\sim 0.01$ . In Simulation  
457 II we found the scaling error to be high in general with global scaling (Additional File 2 : Fig. 3) and  
458 in Simulation III we showed the drastic increase of features falsely identified as DE with global scaling  
459 when two organisms with condition-dependent abundance shifts were combined.

460 To further investigate the increased number of DEF predicted by global scaling, we compared the  
461 scaling factors estimated for the single organism profiles with the estimated scaling factors for the global  
462 normalization in the comparison “day 13” vs “day 27”. For several organisms, a pattern emerged which  
463 showed a condition specific scaling error (Fig. 10). While the scaling factors for one condition are  
464 too small, the scaling factors in the other condition are too high. As a result, global scaling leads to  
465 condition-dependent errors which may cause extra predictions of DEF because an artificial shift between  
466 the two conditions is introduced.

467 Incorrect scaling is especially problematic for features, which mainly comprise counts from one  
468 organism or when counts from mixed organisms with the same scaling shift are analyzed. For a quantifi-  
469 cation we determined the number of features, for which counts from a single organism (or the mixture of  
470 organisms with the same scaling error direction) exceed  $80\%$  of the normalized counts for that feature.  
471 For extra predictions obtained only with global scaling without evidence from the transcriptome analyses,  
472 the counts from a single organism are the main contribution for 82 of 199 features. Additional 43 features  
473 are predicted from the summed counts from organisms with the same scaling shift.

474 **Single feature analysis** For the comparison “day 13” vs “day 27” we now show several examples for  
475 features that reflect particular mixed organism effects or the different behavior of the scaling methods.

476 With regard to the scaling error discussed above, several features actually collect counts from mainly  
477 one organism (see Fig. 11, 12 and Additional File 2 : Fig. 8). In addition, we found several features  
478 that show an observable boosting effect with global as well as taxon-specific scaling (e.g. Fig. 13 or  
479 Additional File 2 : Fig. 9). For these features, significant differences were only observed in the combined  
480 metatranscriptome. Some of the extra predictions obtained with taxon-specific scaling are results of  
481 a putative boosting effect (see Additional File 2 : Fig. 10 and Fig. 11). In contrast, some of the extra  
482 predictions resulting in a putative boosting effect were observable only with global scaling (see Fig. 11  
483 and Additional File 2 : Fig 8). In both cases, the incorrect scaling factors resulted in the detection of DE  
484 for mainly one organism. For other features, the combination of multiple incorrect scaling factors also  
485 predicted DE when global scaling was used (see Additional File 2 : Fig. 12).

486 The prevalent effect for both normalization methods was cancellation. Often, the DEF from multiple  
487 transcriptomes cancel each other out (see Fig. 14 and Additional File 2 : Fig. 13) and in some cases an  
488 organism switch could be observed (see Additional File 2 : Fig. 14 and 15).

489 An example for a contradicting expression direction is shown in Fig 15. For this feature, the incorrect  
490 scaling factors obtained by global scaling for Org2 and Org4 suggest a higher expression of this feature in  
491 “day 13”, while taxon-specific scaling predicted the expression to be higher in the “day 27” condition.

## 492 **CONCLUSIONS**

493 Differential expression analysis in metatranscriptomics is challenging. Metatranscriptomic count data  
494 from RNA-Seq experiments show two main modes of biological variation. The functional composition  
495 of transcripts reflects the activity of organisms and systematic changes might indicate a metabolic  
496 response to experimental conditions. The taxonomic composition of transcripts can change as well and a  
497 change may not necessarily be explainable in terms of controlled experimental conditions. In contrast to  
498 metagenomics, in metatranscriptomics the questions “who is there?” and “what are they doing?” are not  
499 necessarily connected and need to be answered separately. If the two questions are not separated, there is  
500 a considerable risk, to interpret variations in the taxonomic composition as functional changes. This may  
501 even happen if the functional profiles of all organisms stay the same under different conditions.

502 Normalization of metatranscriptomic data must have the goal to eliminate the influence of taxonomic  
503 variations from functional analysis. We argue that for a correct normalization the metatranscriptome needs  
504 to be decomposed to normalize the organism profiles independently. Then the metatranscriptomic count  
505 data may be recombined from the normalized profiles to look for any global trends in the superimposed  
506 count data. If differential expression tools are directly applied to the metatranscriptomic count matrix a  
507 high risk of erroneous results is encountered. Our simulations indicate that the main risk is not to miss  
508 some of the true differences but the real danger is to detect a large number of false functional differences  
509 which arise from taxonomic abundance variations across samples. In particular, if these variations are  
510 condition dependent the false positive rate can explode, circumventing all statistical control mechanisms  
511 for bounding the false discovery rate.

512 We would like to point out that our findings do not affect metatranscriptome studies that just aim  
513 to analyze the functional repertoire from RNA-Seq data. The question which functions or genes are  
514 expressed is much easier to answer than the question what is the functional response to a change of  
515 experimental conditions. However, it is important to note that our results do not only apply to the classic  
516 two conditions setup that we used throughout our study. Also for multiple conditions and time series a  
517 correct normalization is essential to separate functional from taxonomic trends in the metatranscriptomic  
518 composition variations.

## 519 **COMPETING INTERESTS**

520 The authors declare that they have no competing interests.

## 521 **AUTHOR’S CONTRIBUTIONS**

522 HK implemented the method and performed all experiments. PM designed the model and wrote the “Nor-  
523 malization” section of the manuscript. HK & PM designed the simulations, performed data interpretation  
524 and wrote the manuscript. All authors read and approved the final manuscript.

## 525 FUNDING

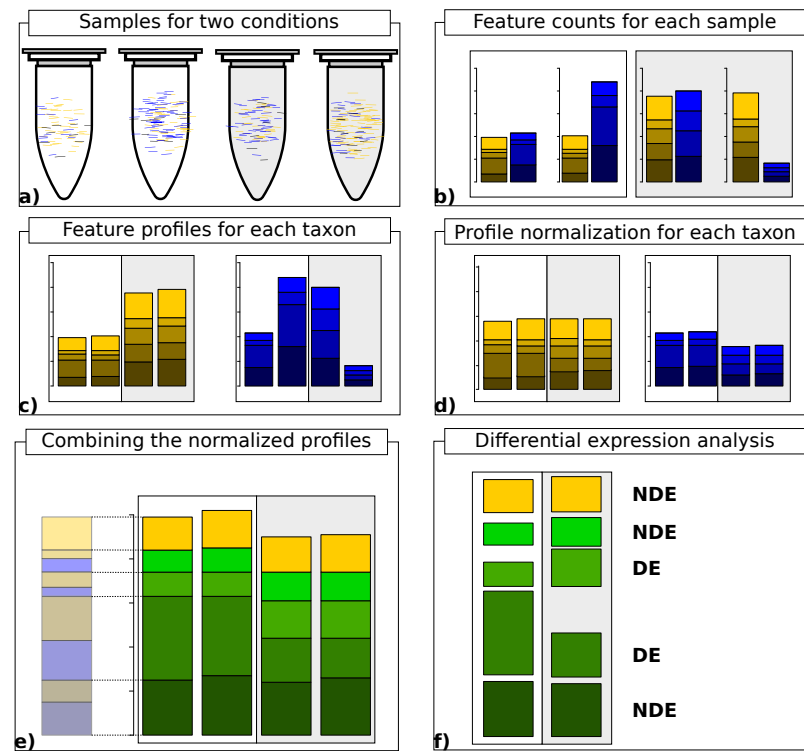
526 This work was partly funded by DFG (Project: “Computational models for metatranscriptome analysis”,  
527 Me3138).

## 528 REFERENCES

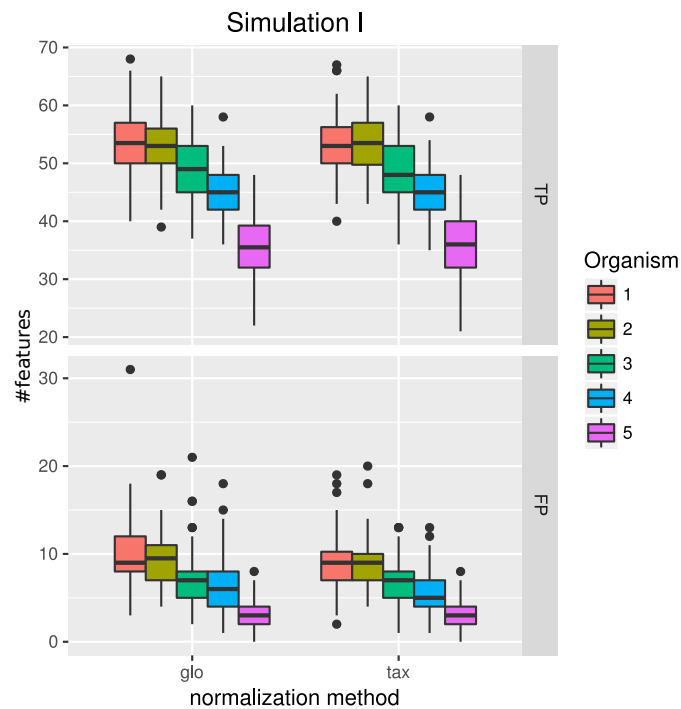
- 529 Anders, S. and Huber, W. (2010). Differential expression analysis for sequence count data. *Genome*  
530 *Biology*, 11(10):R106.
- 531 Benjamini, Y. and Hochberg, Y. (1995). Controlling the False Discovery Rate: A Practical and Powerful  
532 Approach to Multiple Testing. *Journal of the Royal Statistical Society. Series B (Methodological)*,  
533 57(1):289–300.
- 534 Bullard, J. H., Purdom, E., Hansen, K. D., and Dudoit, S. (2010). Evaluation of statistical methods for  
535 normalization and differential expression in mRNA-Seq experiments. *BMC Bioinformatics*, 11:94.
- 536 Celaj, A., Markle, J., Danska, J., and Parkinson, J. (2014). Comparison of assembly algorithms for  
537 improving rate of metatranscriptomic functional annotation. *Microbiome*, 2:39.
- 538 Cheung, V. G., Nayak, R. R., Wang, I. X., Elwyn, S., Cousins, S. M., Morley, M., and Spielman, R. S.  
539 (2010). Polymorphic *Cis*- and *Trans*-Regulation of Human Gene Expression. *PLoS Biology*, 8(9):1–14.
- 540 De Filippis, F., Genovese, A., Ferranti, P., Gilbert, J. A., and Ercolini, D. (2016). Metatranscriptomics re-  
541 veals temperature-driven functional changes in microbiome impacting cheese maturation rate. *Scientific*  
542 *Reports*, 6:21871.
- 543 Dillies, M. A., Rau, A., Aubert, J., Hennequet-Antier, C., Jeanmougin, M., Servant, N., Keime, C.,  
544 Marot, G., Castel, D., Estelle, J., Guernec, G., Jagla, B., Jouneau, L., Laloe, D., Le Gall, C., Schaeffer,  
545 B., Le Crom, S., Guedj, M., and Jaffrezic, F. (2013). A comprehensive evaluation of normalization  
546 methods for Illumina high-throughput RNA sequencing data analysis. *Briefings in Bioinformatics*,  
547 14(6):671–683.
- 548 Eisen, J. A. (2007). Environmental shotgun sequencing: its potential and challenges for studying the  
549 hidden world of microbes. *PLoS Biology*, 5(3):e82.
- 550 Finn, R. D., Bateman, A., Clements, J., Coghill, P., Eberhardt, R. Y., Eddy, S. R., Heger, A., Hetherington,  
551 K., Holm, L., Mistry, J., Sonnhammer, E. L., Tate, J., and Punta, M. (2014). Pfam: the protein families  
552 database. *Nucleic Acids Res.*, 42(Database issue):D222–230.
- 553 Franzosa, E. A., Morgan, X. C., Segata, N., Waldron, L., Reyes, J., Earl, A. M., Giannoukos, G., Boylan,  
554 M. R., Ciulla, D., Gevers, D., Izard, J., Garrett, W. S., Chan, A. T., and Huttenhower, C. (2014).  
555 Relating the metatranscriptome and metagenome of the human gut. *Proceedings of the National*  
556 *Academy of Sciences of the United States of America*, 111(22):E2329–2338.
- 557 Frias-Lopez, J., Shi, Y., Tyson, G. W., Coleman, M. L., Schuster, S. C., Chisholm, S. W., and Delong, E. F.  
558 (2008). Microbial community gene expression in ocean surface waters. *Proceedings of the National*  
559 *Academy of Sciences of the United States of America*, 105(10):3805–3810.
- 560 Gilbert, J. A., Field, D., Huang, Y., Edwards, R., Li, W., Gilna, P., and Joint, I. (2008). Detection of Large  
561 Numbers of Novel Sequences in the Metatranscriptomes of Complex Marine Microbial Communities.  
562 *PLoS ONE*, 3(8):e3042.
- 563 Hesse, C. N., Mueller, R. C., Vuyisich, M., Gallegos-Graves, L. V., Gleasner, C. D., Zak, D. R., and Kuske,  
564 C. R. (2015). Forest floor community metatranscriptomes identify fungal and bacterial responses to N  
565 deposition in two maple forests. *Front Microbiol*, 6:337.
- 566 Huson, D. H., Mitra, S., Ruscheweyh, H. J., Weber, N., and Schuster, S. C. (2011). Integrative analysis of  
567 environmental sequences using MEGAN4. *Genome Res.*, 21(9):1552–1560.
- 568 Li, J. and Tibshirani, R. (2013). Finding consistent patterns: a nonparametric approach for identifying  
569 differential expression in RNA-Seq data. *Stat Methods Med Res*, 22(5):519–536.
- 570 Lin, Y., Golovnina, K., Chen, Z. X., Lee, H. N., Negron, Y. L., Sultana, H., Oliver, B., and Harbison, S. T.  
571 (2016). Comparison of normalization and differential expression analyses using RNA-Seq data from  
572 726 individual *Drosophila melanogaster*. *BMC Genomics*, 17:28.
- 573 Macklaim, J. M., Fernandes, A. D., Di Bella, J. M., Hammond, J. A., Reid, G., and Gloor, G. B. (2013).  
574 Comparative meta-RNA-seq of the vaginal microbiota and differential expression by *Lactobacillus*  
575 *inners* in health and dysbiosis. *Microbiome*, 1(1):12.
- 576 Marionni, J. C., Mason, C. E., Mane, S. M., Stephens, M., and Gilad, Y. (2008). RNA-seq: an assessment of

- 577 technical reproducibility and comparison with gene expression arrays. *Genome Research*, 18(9):1509–  
578 1517.
- 579 Martinez, X., Pozuelo, M., Pascal, V., Campos, D., Gut, I., Gut, M., Azpiroz, F., Guarner, F., and  
580 Manichanh, C. (2016). MetaTrans: an open-source pipeline for metatranscriptomics. *Scientific Reports*,  
581 6:26447.
- 582 Mason, O. U., Hazen, T. C., Borglin, S., Chain, P. S., Dubinsky, E. A., Fortney, J. L., Han, J., Holman,  
583 H. Y., Hultman, J., Lamendella, R., Mackelprang, R., Malfatti, S., Tom, L. M., Tringe, S. G., Woyke,  
584 T., Zhou, J., Rubin, E. M., and Jansson, J. K. (2012). Metagenome, metatranscriptome and single-cell  
585 sequencing reveal microbial response to Deepwater Horizon oil spill. *ISME J*, 6(9):1715–1727.
- 586 McNulty, N. P., Wu, M., Erickson, A. R., Pan, C., Erickson, B. K., Martens, E. C., Pudlo, N. A., Muegge,  
587 B. D., Henrissat, B., Hettich, R. L., and Gordon, J. I. (2013). Effects of Diet on Resource Utilization by  
588 a Model Human Gut Microbiota Containing *Bacteroides cellulosilyticus* WH2, a Symbiont with an  
589 Extensive Glycobiome. *PLoS Biology*, 11(8):1–20.
- 590 Nacke, H., Fischer, C., Thürmer, A., Meinicke, P., and Daniel, R. (2014). Land use type significantly  
591 affects microbial gene transcription in soil. *Microbial Ecology*, 67(4):919–930.
- 592 Pickrell, J. K., Marioni, J. C., Pai, A. A., Degner, J. F., Engelhardt, B. E., Nkadori, E., Veyrieras, J.-B.,  
593 Stephens, M., Gilad, Y., and Pritchard, J. K. (2010). Understanding mechanisms underlying human  
594 gene expression variation with RNA sequencing. *Nature*.
- 595 Poretsky, R. S., Bano, N., Buchan, A., LeClerc, G., Kleikemper, J., Pickering, M., Pate, W. M., Moran,  
596 M. A., and Hollibaugh, J. T. (2005). Analysis of microbial gene transcripts in environmental samples.  
597 *Applied and Environmental Microbiology*, 71(7):4121–4126.
- 598 Poretsky, R. S., Hewson, I., Sun, S., Allen, A. E., Zehr, J. P., and Moran, M. A. (2009). Comparative  
599 day/night metatranscriptomic analysis of microbial communities in the North Pacific subtropical gyre.  
600 *Environmental Microbiology*, 11(6):1358–1375.
- 601 Ritchie, M. E., Phipson, B., Wu, D., Hu, Y., Law, C. W., Shi, W., and Smyth, G. K. (2015). limma powers  
602 differential expression analyses for RNA-sequencing and microarray studies. *Nucleic Acids Research*,  
603 43(7):e47.
- 604 Robinson, M. D., McCarthy, D. J., and Smyth, G. K. (2010). edgeR: a Bioconductor package for  
605 differential expression analysis of digital gene expression data. *Bioinformatics*, 26(1):139–140.
- 606 Robinson, M. D. and Oshlack, A. (2010). A scaling normalization method for differential expression  
607 analysis of RNA-seq data. *Genome Biology*, 11(3):R25.
- 608 Simon, C. and Daniel, R. (2009). Achievements and new knowledge unraveled by metagenomic ap-  
609 proaches. *Applied Microbiology and Biotechnology*, 85(2):265–276.
- 610 Sonesson, C. (2014). compcodeR – an R package for benchmarking differential expression methods for  
611 RNA-seq data. *Bioinformatics*, 30(17):2517–2518.
- 612 Sonesson, C. and Delorenzi, M. (2013). A comparison of methods for differential expression analysis of  
613 RNA-seq data. *BMC Bioinformatics*, 14:91.
- 614 Trapnell, C., Roberts, A., Goff, L., Pertea, G., Kim, D., Kelley, D. R., Pimentel, H., Salzberg, S. L.,  
615 Rinn, J. L., and Pachter, L. (2012). Differential gene and transcript expression analysis of RNA-seq  
616 experiments with TopHat and Cufflinks. *Nature Protocols*, 7(3):562–578.
- 617 Urich, T., Lanzen, A., Qi, J., Huson, D. H., Schleper, C., and Schuster, S. C. (2008). Simultaneous assess-  
618 ment of soil microbial community structure and function through analysis of the meta-transcriptome.  
619 *PLoS ONE*, 3(6):e2527.
- 620 Westreich, S. T., Korf, I., Mills, D. A., and Lemay, D. G. (2016). SAMSA: a comprehensive metatran-  
621 scriptome analysis pipeline. *BMC Bioinformatics*, 17(1):399.
- 622 Ye, Y. and Tang, H. (2016). Utilizing de Bruijn graph of metagenome assembly for metatranscriptome  
623 analysis. *Bioinformatics*, 32(7):1001–1008.

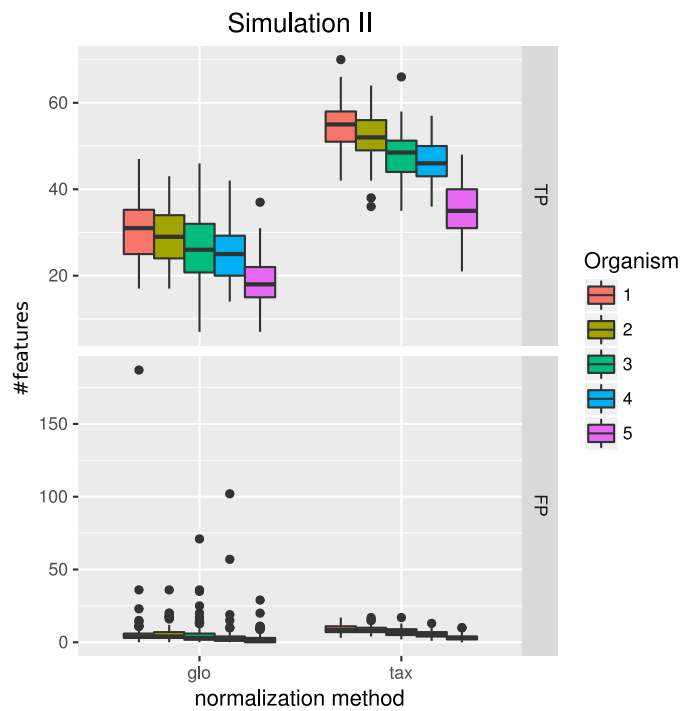
## 624 FIGURES



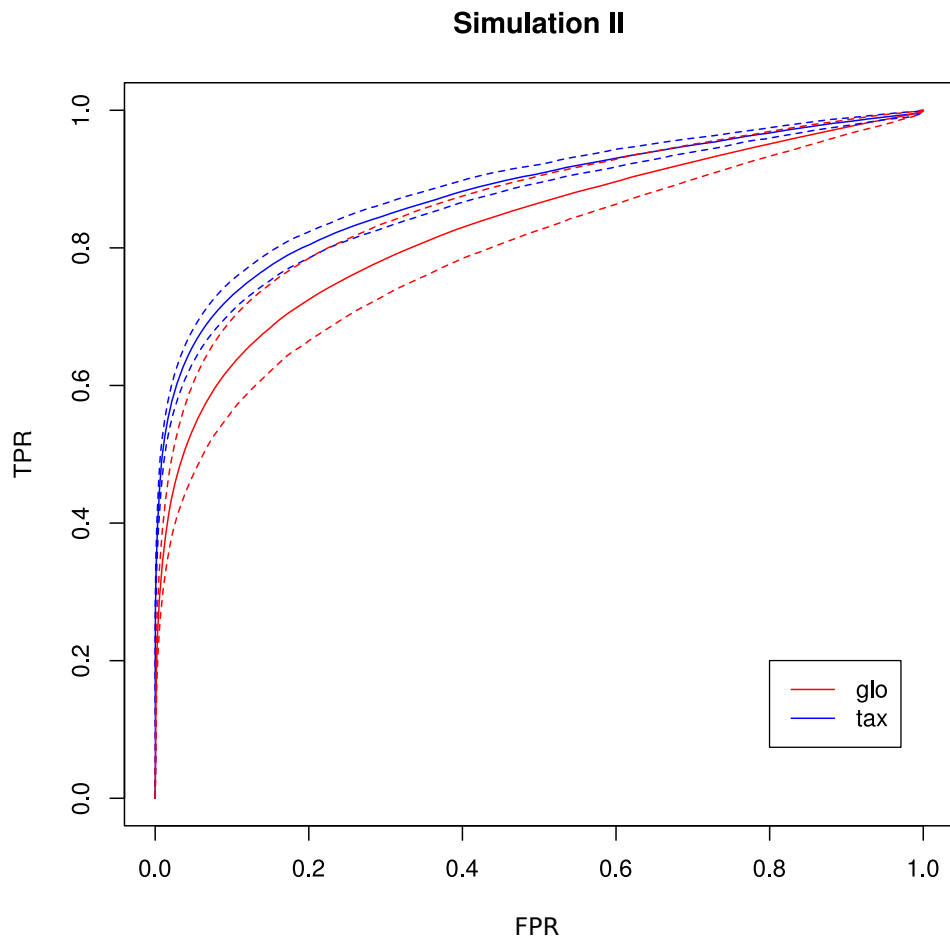
**Figure 1.** Workflow for taxon-specific normalization. a) Sequence samples from conditions A (white) and B (light gray). Assign each sequence read to taxonomic and feature categories. b) Compute feature profiles from the assignment counts. c) Obtain count matrix from taxon-specific feature profiles. d) Normalize feature profiles of each taxon-specific count matrix separately. e) Recombine normalized feature profiles of all taxa into a metatranscriptomic profile. f) Perform differential expression analysis on metatranscriptomic count matrix.



**Figure 2.** Simulation I. Number of true positive (TP) and false positive (FP) features identified with DESeq2 for global (glo) and taxon-specific (tax) scaling: Boxplots represent variation over 100 runs of the simulation.

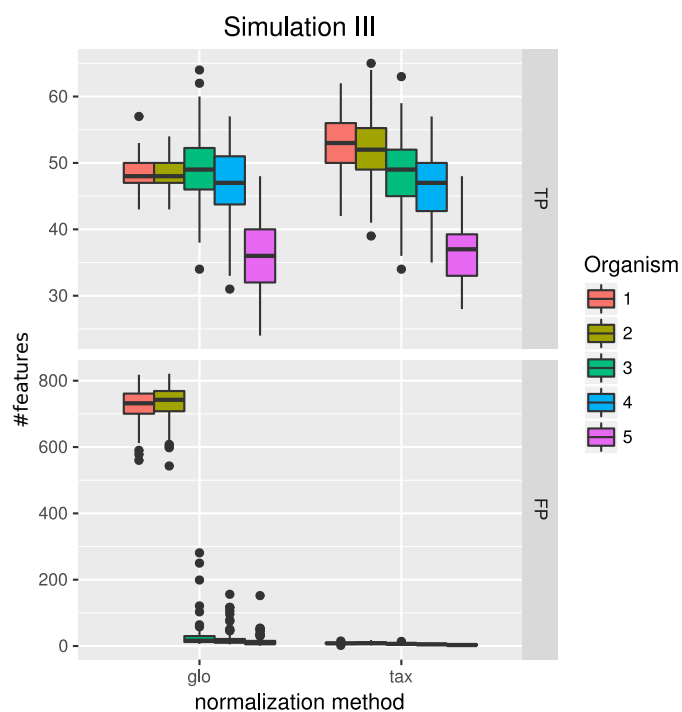


**Figure 3.** Simulation II. Number of true positive (TP) and false positive (FP) features identified with DESeq2 for global (glo) and taxon-specific (tax) scaling. FP boxplots appear compressed due to outliers. Organism order for FP is the same as for TP boxplots.



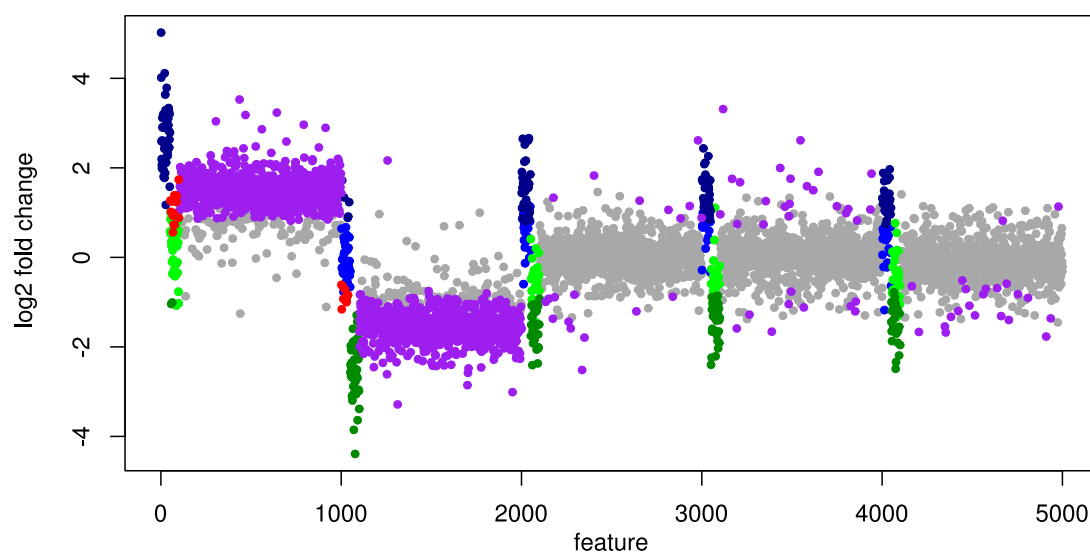
**Figure 4.** ROC curves for Simulation II. Average curve for taxon-specific scaling (blue) vs. average curve for global scaling (red) with false positive rate (FPR) on x-axis and true positive rate (TPR) on y-axis. Dotted lines above and below indicate the standard deviation for each method. The average area under curve (AUC) is 0.8776 for taxon-specific scaling and 0.8282 for global scaling.



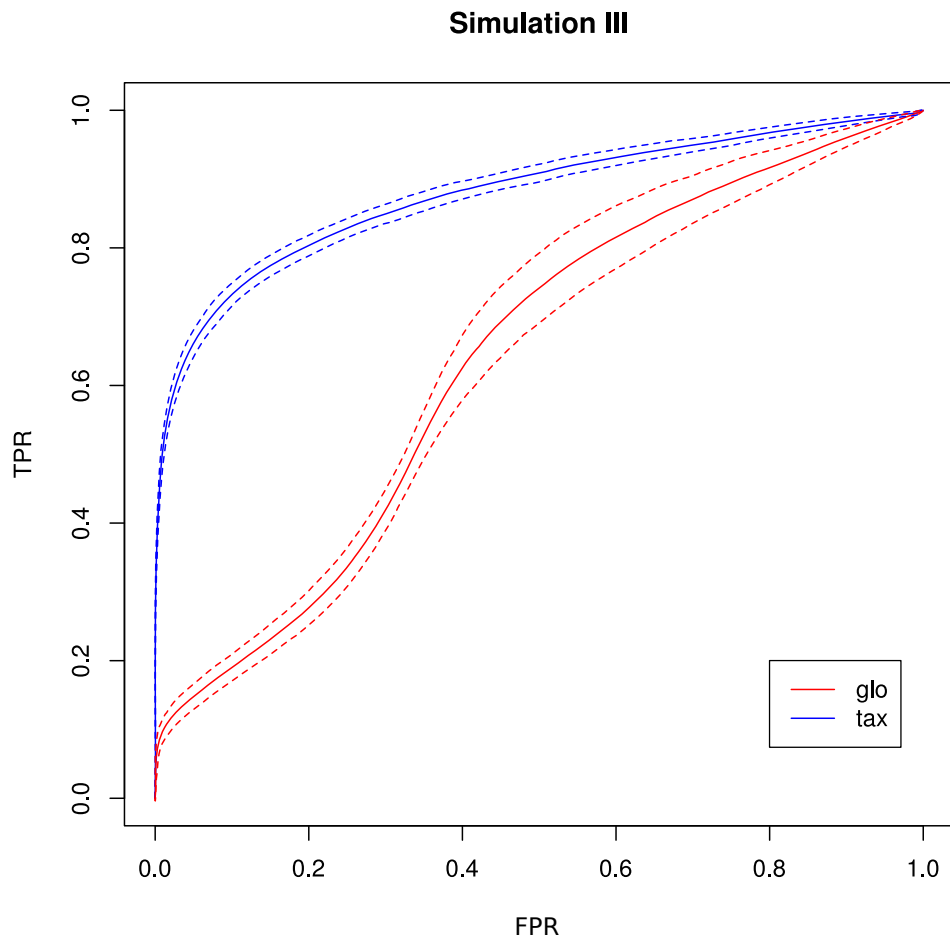


**Figure 5.** Simulation III. Number of true positive (TP) and false positive (FP) features identified with DESeq2 for global (glo) and taxon-specific (tax) scaling. Boxplots represent variation over 100 runs of the simulation.

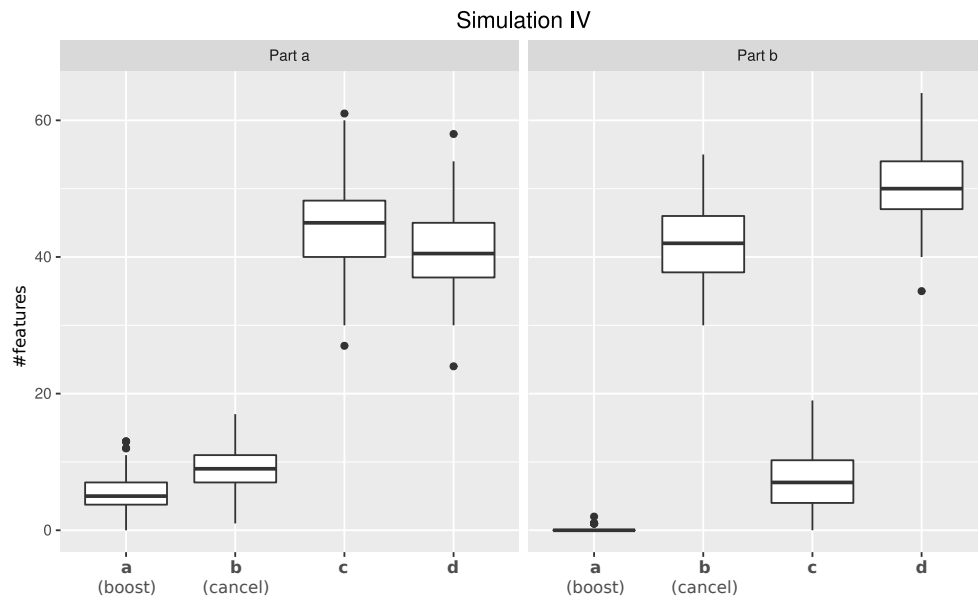
### DEF identification with global scaling method



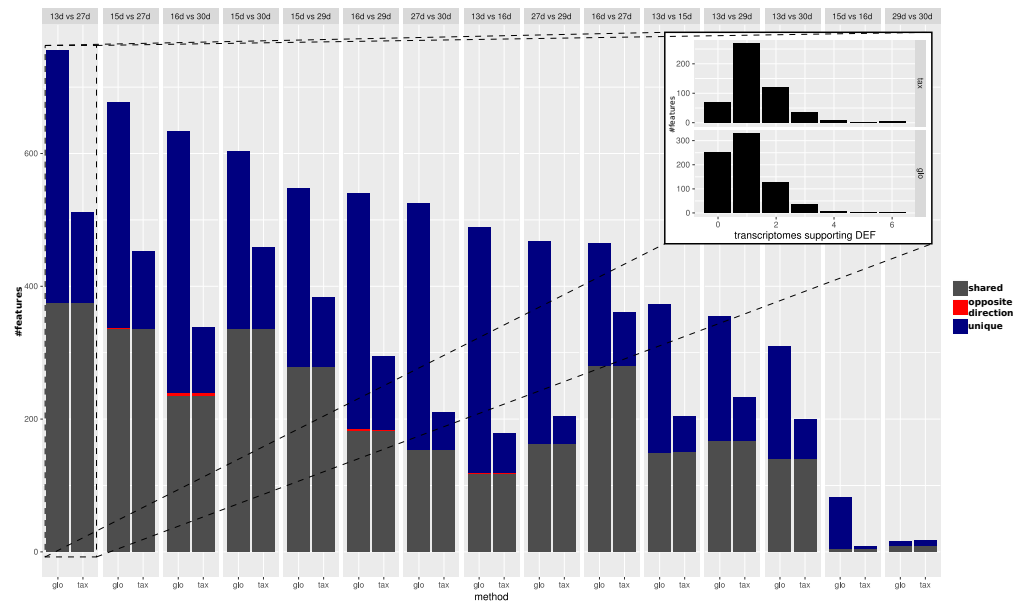
**Figure 6.** Simulation III: Log<sub>2</sub> fold changes. Log<sub>2</sub> fold changes of features for global normalization on one example data set. Along x-axis, features (dots) are ordered according to five stacked organism profiles, each with 1000 features of which the first 50 features are “upregulated”, and the next 50 features are “downregulated”. Gray dots correspond to correctly detected NDE features, light green dots to downregulated features which are missed and dark green dots to correctly identified downregulated DEF. Light blue dots correspond to missed upregulated DEF and dark blue dots to correctly identified upregulated DEF. Red dots mark DEF where global scaling leads to an incorrect direction. Purple dots correspond to NDE features which are incorrectly identified as significant features.



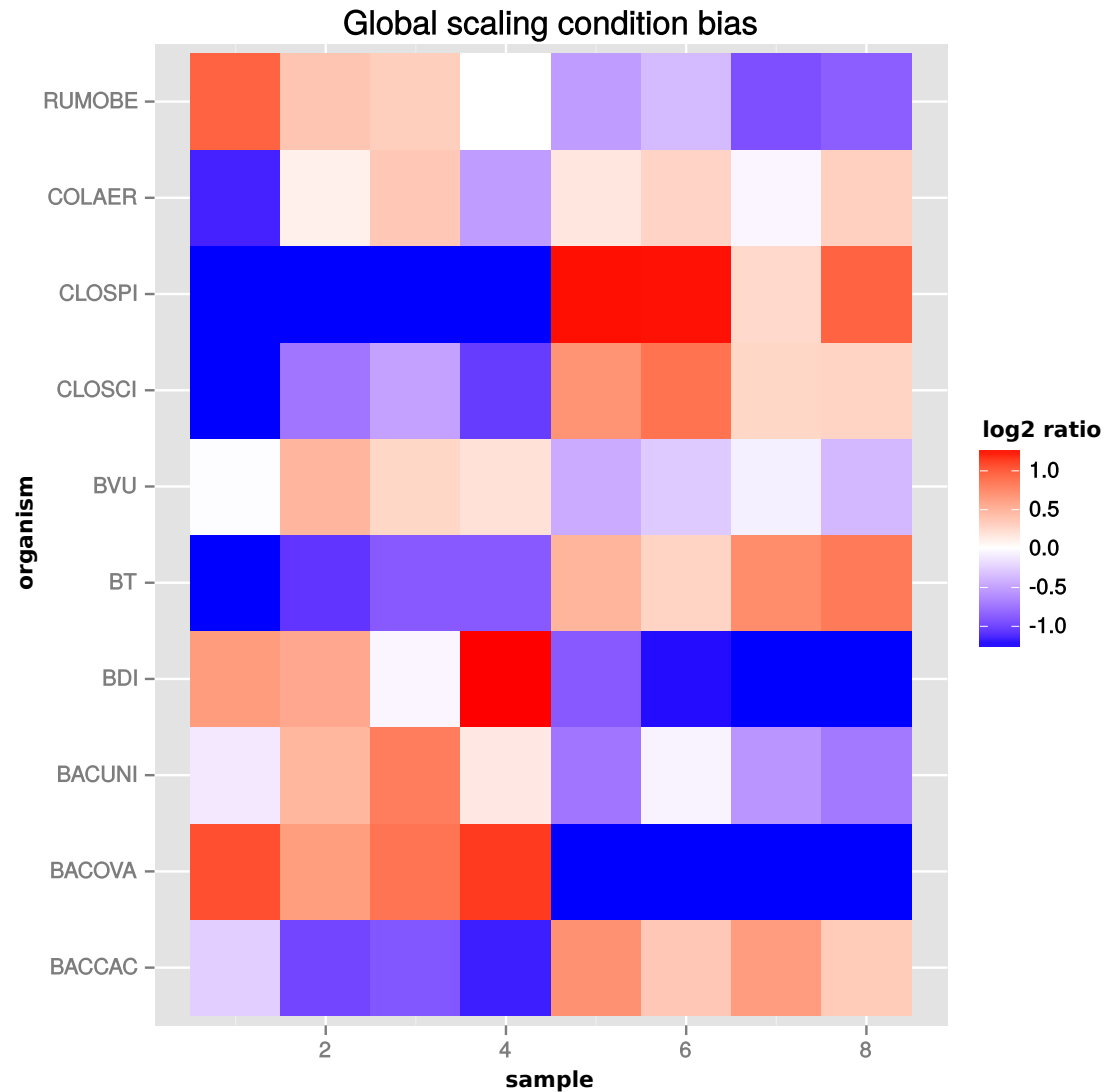
**Figure 7.** ROC curves for Simulation III. Average curve for taxon-specific scaling (blue) vs. average curve for global scaling (red) with false positive rate (FPR) on x-axis and true positive rate (TPR) on y-axis. Dotted lines above and below indicate the standard deviation for each method. The average area under curve (AUC) is 0.8785 for taxon-specific scaling and 0.6369 for global scaling.



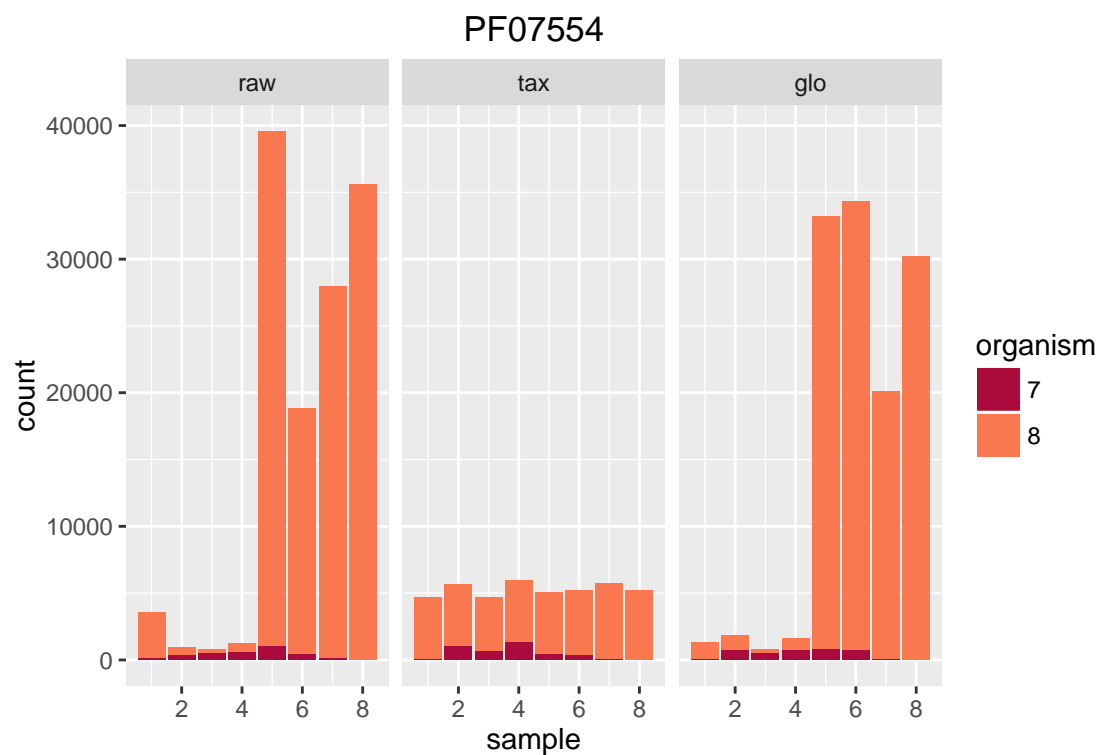
**Figure 8.** Simulation IV. Number of DEF that show a particular mixed organism effect. Effects are distinguished according to a) features not detected as DE in a single transcriptome but predicted as DE in the mixture (boosting), b) features detected as DE in at least one transcriptome but predicted as NDE in the mixture (cancellation), c) features identified as DE in at least one transcriptome and also predicted as DE in the mixture and d) features not detected as DE for both, transcriptomes and mixture. Boxplots represent variation over 100 runs of the simulation.



**Figure 9.** Predicted DEF for real data. Number of significant features from taxon-specific scaling (“tax”, right bar) and global scaling (“glo”, left bar) for different condition comparisons. Colors indicate shared significant features with same direction of difference (grey), shared significant features with opposite direction (red) and mutually exclusive features (purple) that are only found to be significant for one scaling method. Smaller figure: histogram for predicted DEF according to the number of single organism analyses that show a significant difference (x-axis). Upper part shows results for taxon-specific scaling and lower part for global scaling. For example, a high bar at “0” means that many features are found to be significant for the metatranscriptome which are not significant for any of the single transcriptome analyses

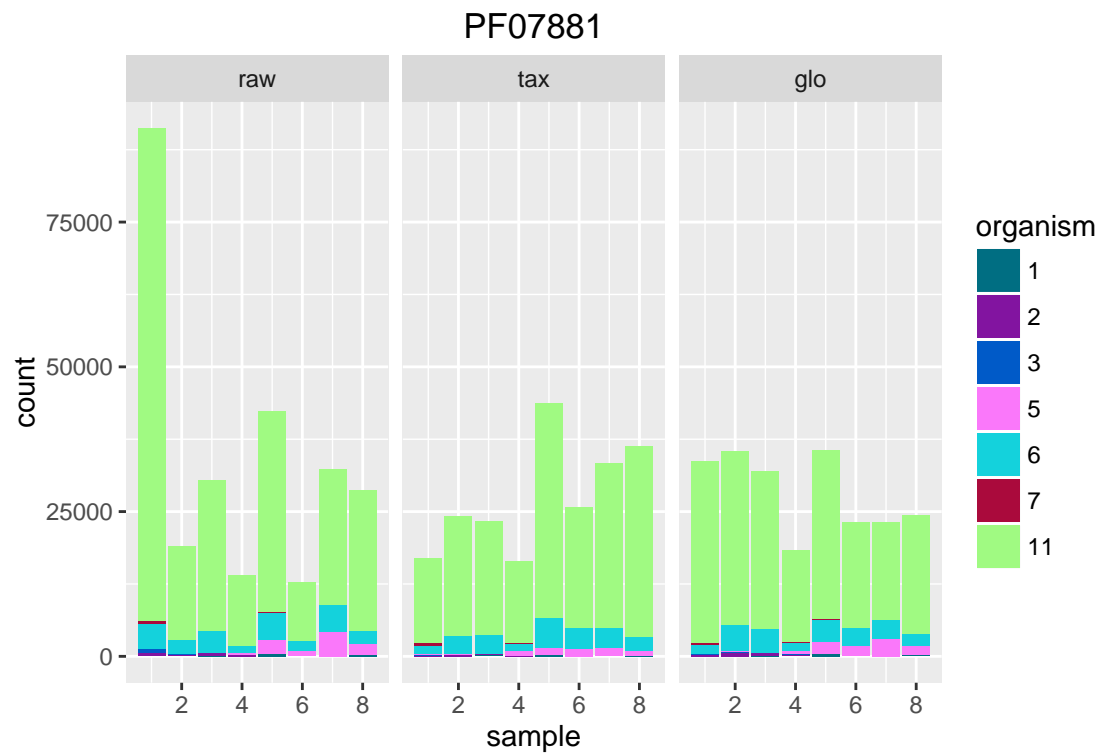


**Figure 10.** Global scaling condition bias. Direction of global scaling “error” in terms of the log<sub>2</sub>-ratio of scaling factors from transcriptomic and global scaling. Results for different organisms in the comparison of “day 13” vs. “day 27”. For symmetry of the color range the negative log<sub>2</sub>-ratio was capped at -1.25, with error scores below that threshold showing the same color (blue). Samples 1-4 are from condition A and samples 5-8 from condition B. For the species name abbreviations see Additional File 2 : Tab. 1. *D. longicatena DSM 13814* was not observed in that particular condition comparison.



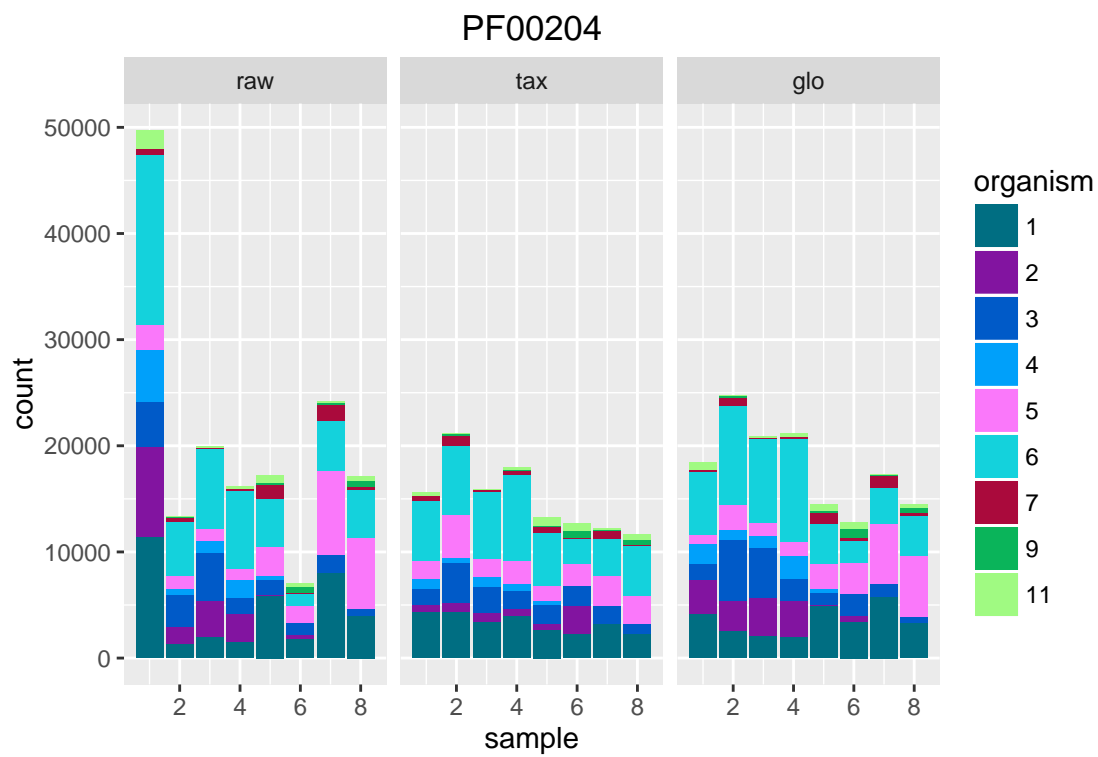
**Figure 11.** Single feature analysis (PF07554). Stacked bars in three parts (x-axis) show organism-specific counts before scaling (left), after taxon-specific scaling (middle) and after Taxon-specific and global scaling result in adjusted p-values 0.98 and  $7.23^{-56}$ , respectively.

global scaling (right). Extra prediction of DEF by (mis)scaling of one organism with global scaling.

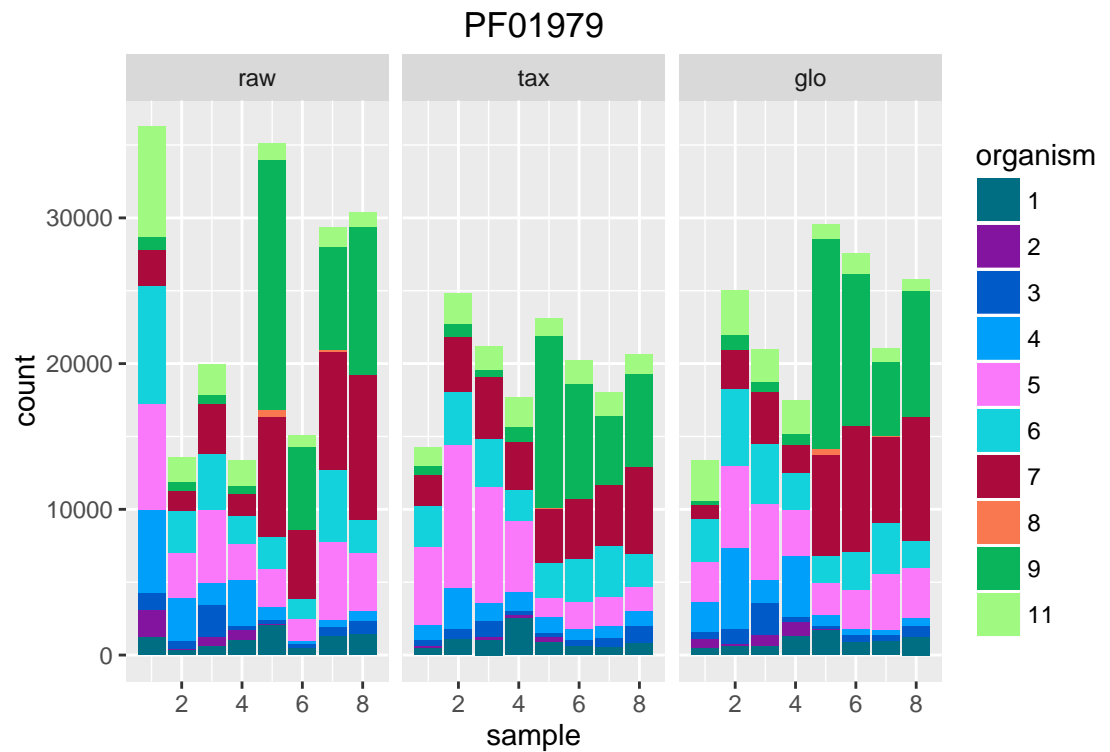


**Figure 12.** Single feature analysis (PF07881). Stacked bars in three parts (x-axis) show organism-specific counts before scaling (left), after taxon-specific scaling (middle) and after global scaling (right). Loss of significance due to (mis)scaling of profiles from mainly one organism. Feature is significant for organisms Org2, Org5 and Org11 in transcriptome analysis. Taxon-specific and global scaling result in adjusted p-values  $1.79e^{-3}$  and 0.66, respectively.

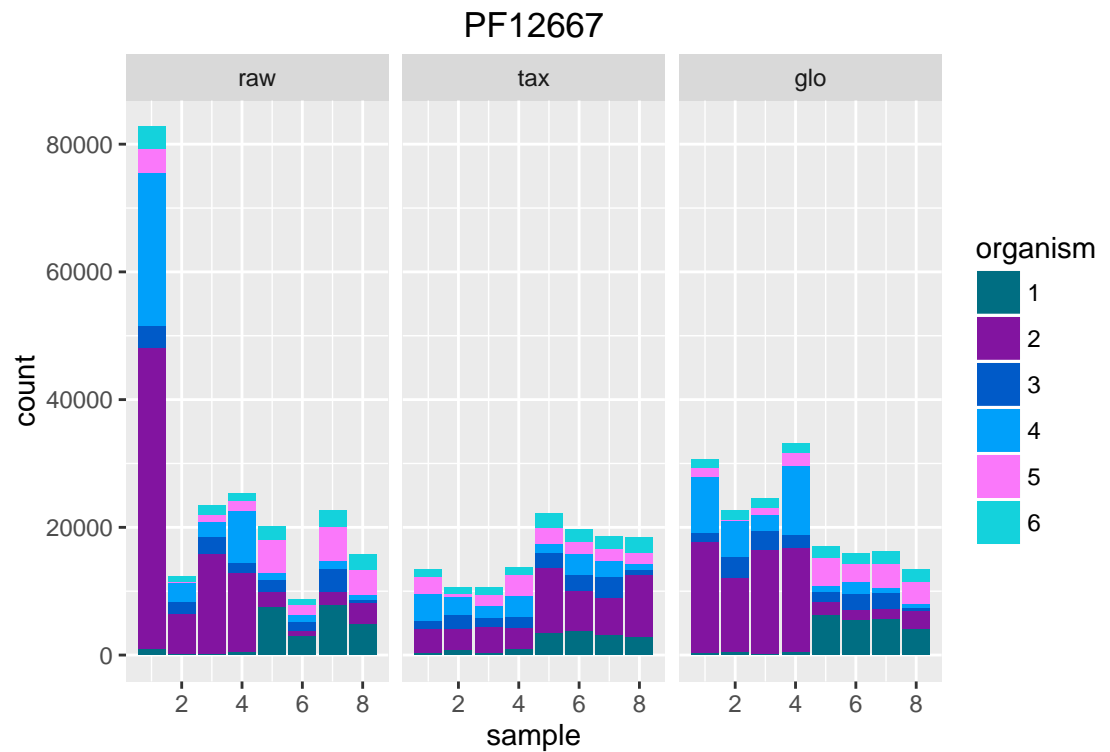




**Figure 13.** Single feature analysis (PF00204). Stacked bars in three parts (x-axis) show organism-specific counts before scaling (left), after taxon-specific scaling (middle) and after global scaling (right). With boosting effect observable for both methods, i.e. feature is insignificant in transcriptome analyses and significant for metatranscriptomic counts. For Org1 and Org6 the adjusted p-values are close to 0.05 for single analysis, with taxon-specific scaling an adjusted p-value  $p = 0.005$  is achieved.



**Figure 14.** Single feature analysis (PF01979). Stacked bars in three parts (x-axis) show organism-specific counts before scaling (left), after taxon-specific scaling (middle) and after global scaling (right). Feature is significant for transcriptome analysis but becomes insignificant for metatranscriptomic counts (cancellation effect), i.e. the adjusted p-value with both scaling methods is above 0.05.



**Figure 15.** Single feature analysis (PF12667). Stacked bars in three parts (x-axis) show organism-specific counts before scaling (left), after taxon-specific scaling (middle) and after global scaling (right). Significant feature with opposite direction for the two scaling methods. Taxon-specific and global scaling result in adjusted p-values  $6.91e^{-5}$  and  $1.33e^{-5}$ , respectively. The log<sub>2</sub> fold change for condition A in comparison to condition B is 0.70 for taxon-specific scaling and -0.82 for global scaling.