

SHAMAN: a user-friendly website for metataxonomic analysis from raw reads to statistical analysis

Stevenn Volant

Institut Pasteur

Pierre Lechat

Institut Pasteur

Perrine Woringer

Institut Pasteur

Laurence Motreff

Institut Pasteur

Christophe Malabat

Institut Pasteur

Sean Kennedy

Institut Pasteur

Amine Ghozlane (✉ amine.ghozlane@pasteur.fr)

Institut Pasteur <https://orcid.org/0000-0001-7670-6235>

Software article

Keywords: Metagenomics, Differential analysis, Visualization, Web application

Posted Date: February 11th, 2020

DOI: <https://doi.org/10.21203/rs.2.23213/v1>

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

[Read Full License](#)

Version of Record: A version of this preprint was published at BMC Bioinformatics on August 10th, 2020.

See the published version at <https://doi.org/10.1186/s12859-020-03666-4>.

SOFTWARE ARTICLE

SHAMAN: a user-friendly website for metataxonomic analysis from raw reads to statistical analysis

Stevann Volant¹, Pierre Lechat¹, Perrine Woringer¹, Laurence Motreff², Christophe Malabat¹, Sean Kennedy² and Amine Ghozlane^{1,2*}

*Correspondence:

amine.ghozlane@pasteur.fr

¹Hub de Bioinformatique et

Biostatistique – Département

Biologie Computationnelle,

Institut Pasteur, USR 3756 CNRS,

28 rue du Docteur Roux, 75015

Paris, France

Full list of author information is available at the end of the article

Abstract

Background: Comparing the composition of microbial communities among groups of interest (e.g., patients vs healthy individuals) is a central aspect in microbiome research. It typically involves sequencing, data processing, statistical analysis and graphical representation of the detected signatures. Such an analysis is normally obtained by using a set of different applications that require specific expertise for installation, data processing and in some case, programming skills.

Results: Here, we present SHAMAN, an interactive web application we developed in order to facilitate the use of (i) a bioinformatic workflow for metataxonomic analysis, (ii) a reliable statistical modelling and (iii) to provide among the largest panels of interactive visualizations as compared to the other options that are currently available. SHAMAN is specifically designed for non-expert users who may benefit from using an integrated version of the different analytic steps underlying a proper metagenomic analysis. The application is freely accessible at <http://shaman.pasteur.fr/>, and may also work as a standalone application with a Docker container (aghozlane/shaman), conda and R. The source code is written in R and is available at <https://github.com/aghozlane/shaman>. Using two datasets (a mock community sequencing and published 16S rRNA metagenomic data), we illustrate the strengths of SHAMAN in quickly performing a complete metataxonomic analysis.

Conclusions: We aim with SHAMAN to provide the scientific community with a platform that simplifies reproducible quantitative analysis of metagenomic data.

Keywords: Metagenomics; Differential analysis; Visualization; Web application

1

2

3 Background

4 Quantitative metagenomic techniques have been broadly deployed to identify asso-
5 ciations between microbiome and environmental or individual factors (e.g., disease,
6 geographical origin, etc.). Analyzing changes in the composition and/or the abun-
7 dance of microbial communities yielded promising biomarkers, notably associated
8 with liver cirrhosis[1], diarrhea[2], colorectal cancer[3], or associated with various
9 pathogenic[4] or probiotic effects[5] on the host.

10 In metataxonomic studies, a choice is made prior to sequencing in order to specifi-
11 cally amplify one or several regions of the rRNA (usually the 16S or the 18S rRNA
12 genes for procaryotes/archaea and the ITS, the 23S or the 28S rRNA gene for eu-
13 karyotes) so that the composition of microbial communities may be characterized
14 with affordable techniques.

15 A typical workflow includes successive steps: (i) OTU (Operational Taxonomic
16 Unit) picking (dereplication, denoising, chimera filtering and clustering)[6], (ii) OTU
17 quantification in each sample and (iii) OTU annotating with respect to a reference
18 taxonomic database. This process may require substantial computational resources
19 depending on both the number of samples involved and the sequencing depth. Sev-
20 eral methods are currently available to complete these tasks, such as Mothur[7],
21 Usearch[8], DADA2[9] or Vsearch[10]. The popular Qiime[11] simplifies these tasks
22 (i to iii) and visualizations by providing a python-integrated environment. Schemat-
23 ically, once data processing is over, both a contingency table and a taxonomic table
24 are obtained. They contain the abundance of OTUs in the different samples and the
25 taxonomic annotations of OTUs, respectively. The data are normally represented
26 in the standard BIOM format[12].

27 Statistical analysis is then performed to screen significant variation in micro-
28 bial abundance. To this purpose, several R packages were developed, such as
29 Metastats[13] or Metagenomeseq[14]. It is worth noticing that other approaches
30 which were originally designed for RNA-seq, namely DESeq2[15] and EdgeR[16],
31 are also commonly used to carry out metataxonomic studies[17, 18]. They provide
32 an R integrated environment for statistical modelling in order to test the effects
33 of a particular factor on OTU abundance. Nevertheless using all of these different
34 methods requires a technical skills in Unix, R and experience in processing metage-
35 nomics data. To this end, we developed SHAMAN in order to provide a method
36 that simplifies the analysis of metataxonomic data, especially for users who are not
37 familiar with the technicalities of bioinformatic and statistical methods that are
38 commonly applied in this field.

39 SHAMAN is an all-inclusive approach to estimate the composition and abundance
40 of OTUs, based on raw sequencing data, and to perform the statistical analysis of
41 the processed files. First, the user can submit raw data in FASTQ format and de-
42 fine the parameters of the bioinformatic workflow. The output returns a BIOM file
43 for each database used as reference for annotation, a phylogenetic tree in Newick
44 format as well as FASTA-formatted sequences of all OTUs that were identified. The
45 second step consists in performing statistical analysis. The user has to provide a
46 "target" file that associates each sample with one or several explanatory variables.
47 These variables are automatically detected in the target file. An automatic filtering
48 of the contingency matrix of OTUs may be activated in order to remove features
49 with low frequency. Setting up the contrasts to be compared was also greatly sim-
50 plified. It consists in filling in a form that orients the choices of users when defining
51 the groups of interest. Several options to visualize data are available at three im-
52 portant steps of the process: quality control, bio-analysis and contrast comparison.
53 At each step, a number of common visual displays are implemented in SHAMAN
54 to explore data. In addition, SHAMAN also includes a variety of original displays
55 that is not available in other applications such as an abundance tree to visualize
56 count distribution according to the taxonomic tree and variables, or the logit plot
57 to compare feature p-values in two contrasts. Figures may be tuned to emphasize
58 particular statistical results (e.g., displaying significant features in a given contrast,
59 performing intersection between contrasts), to be more specific (e.g. feature abun-
60 dance in a given modality) or to improve the aesthetics of the graph (by changing

61 visual parameters). Figures fit publication standards and the corresponding file can
62 be easily downloaded.

63 Several web applications were developed to analyze data of metataxonomic studies,
64 notably, FROGS[19] as well as Qiita[20] for bioinformatic data processing, Shiny-
65 phyloseq[21] for statistical analysis, Metaviz[22] and VAMPS2[23] that make a par-
66 ticular focus on data visualization. While these interfaces propose related function-
67 alities, the main specificity of SHAMAN is to combine of all these steps in a single
68 user-friendly application. Last, SHAMAN may register a complete analysis which
69 may be of particular interest for matters of reproducibility.

70 **Implementation**

71 SHAMAN is implemented in R using the shiny-dashboard framework. The appli-
72 cation is divided into three main components (Fig. S1): a bioinformatic workflow
73 to process the raw FASTQ-formatted sequences, a statistical workflow to normalize
74 and further analyse data, as well as a visualization platform.

75 **Metataxonomic pipeline**

76 The bioinformatic workflow implemented in SHAMAN relies on the Galaxy
77 platform[24] that provides modular and scalable analyses. SHAMAN includes a
78 daemon-program (written in Python) using bioblend[25] to communicate with
79 Galaxy. It is worth noticing that previous studies, e.g. performed on mosquito
80 microbiota[26], showed that some non-annotated OTUs turned to be sequences of
81 the host organism. To overcome such issues, the user can optionally filter out reads
82 that align with the host genome and the PhiX174 genome (used as a control in Illu-
83 mina sequencers). The latter task is performed with Bowtie2 v2.2.6[27]. By default,
84 quality of reads is checked with AlienTrimmer[28] v0.4.0, a software for trimming
85 off contaminant sequences and clipping. Paired-end reads are then merged with
86 Pear[29] v0.9.10.1. OTU picking, taxonomic annotation and OTU quantification
87 are performed using Vsearch[10] v2.3.4.0, a software which is both accurate and
88 efficient [6]. The process also includes several steps of dereplication, singleton re-
89 moval and chimera detection. By default, clustering is performed with a threshold
90 of 97% in sequence identity. The input amplicons are aligned against the set of
91 detected OTUs to create a contingency table containing the number of amplicons
92 assigned to each OTU. The taxonomic annotation of OTUs is performed based on
93 various databases, i.e., with SILVA[30] rev. 132 SSU (for 16S, 18S rRNA genes) and
94 LSU (for 23S and 28S rRNA genes), Greengenes[31] (for 16S, 18S rRNA genes) and
95 Underhill rev. 1.6.1[32], Unite rev. 8.0[33] and Findley[34] for ITS rRNA sequences.
96 These databases are kept up-to-date every two month with biomaj.pasteur.fr. OTU
97 annotations are filtered according to their identity with the reference[35]. Phylum
98 annotations are kept when the identity between the OTU sequence and reference
99 sequence is $\geq 75\%$, $\geq 78.5\%$ for classes, $\geq 82\%$ for orders, $\geq 86.5\%$ for families,
100 $\geq 94.5\%$ for genera and $\geq 98\%$ for species. In addition, a taxonomic inference
101 made based on a naive Bayesian approach, RDP classifier[36] v2.12, is systemati-
102 cally provided. By default, RDP annotations are included whenever the annotation
103 probability is ≥ 0.5 . All the above-mentioned thresholds may be tuned by the user.
104 A phylogenetic analysis of OTUs is provided: multiple alignments are obtained with

105 Mafft[37] v7.273.1, filtering of regions that are insufficiently conserved is processed
106 using BMGE[38] v1.12 and finally, FastTree[39] v2.1.9 is used to infer the phylo-
107 genetic tree. Based on the latter tree, a Unifrac distance[40] may be computed in
108 SHAMAN to compare microbial communities. The outcomes of the overall workflow
109 are stored in several files: a BIOM file (per reference database), a phylogenetic tree
110 as well as a summary file describing the number of elements passing the different
111 steps of the workflow. The data are associated to a key that is unique to a project.
112 Such a key allows to automatically re-load all the results previously obtained in a
113 given project.

114 Statistical workflow

115 The statistical analysis in SHAMAN is based on DESeq2 which is a method to
116 model OTU counts with a negative binomial distribution. It is known as one of the
117 most accurate approach to detect differentially abundant bacteria in metagenomic
118 data[17, 18]. Relying on a robust estimation of variation in OTUs, the DESeq2
119 method shows suitable performances with datasets characterized by a relatively
120 low number of observations per group together with a high number of OTUs.
121 This method typically requires the following input files: a contingency table, a
122 taxonomic table and a target file describing the experimental design. These data are
123 processed to generate a meta-table that assign to each OTU a taxonomic annotation
124 and a raw count per sample.

125 Normalization

126 Normalization of the raw counts is one of the key issues when analyzing microbiome
127 experiments. The uniformity of the sequencing depth is affected by sample prepa-
128 ration and dye effects[41]. Normalizing data is therefore expected to increase the
129 accuracy of comparisons. It is done by adjusting the abundance of OTUs across sam-
130 ples. Four different normalization methods are currently implemented in SHAMAN.
131 For the sake of consistency, all of these methods are applied at the OTU level.
132 A first method is the relative log expression (RLE) normalization and is imple-
133 mented in the DESeq2 package. It consists in calculating a size factor of each sample,
134 i.e., a multiplication factor that increases or decreases the OTU counts in samples.
135 It is defined as the median ratio, between a given count and the geometric mean
136 of each OTU. Such a normalization was shown to be suited for metataxonomic
137 studies[17]. In practice, many OTUs are found in a few samples only, which trans-
138 late into sparse count matrices[14]. In this case, the RLE method may lead to a
139 defective normalization - as only a few OTU are taken into account - or might
140 be impossible if all OTUs show a null abundance in one sample at least. In the
141 Phyloseq[42] R package, the decision was made to replace the null abundance by a
142 count of 1. In SHAMAN, we decided to include two new normalization methods.
143 They are modified versions of the original RLE so that they better account for
144 matrix sparsity (number of zero-valued elements divided by the total number of el-
145 ements). In the *non-null normalization* (1) cells with null values are excluded from
146 the computation of the geometric mean. This method therefore takes all OTUs into
147 account when estimating the size factor. In the second method that we coined as
148 the *weighted non-null normalization* (2), weights are introduced so that OTUs with

149 a big number of occurrences have a higher influence when calculating the geometric
150 mean.

151 Assume that $C = (c_{ij})_{1 \leq i \leq k; 1 \leq j \leq n}$ is a contingency table where k and n are the
152 number of features (e.g. OTUs) and the number of samples, respectively. Here, c_{ij}
153 represents the abundance of the feature i in the sample j . The size factor of sample
154 j is denoted by s_j .

$$155 \quad s_j^{(1)} = \text{median}_i \frac{c_{ij}}{(\prod_{k \in S_i} c_{ik})^{1/n_i}}, \quad (1)$$

156

$$157 \quad s_j^{(2)} = w.\text{median}_i \frac{c_{ij}}{(\prod_{k \in S_i} c_{ik})^{1/n_i}}, \quad (2)$$

158 where S_i stands for the subset of samples with non null values for the feature j
159 and n_i is the size of this subset. The function *w.median* corresponds to a weighted
160 median.

161 An alternative normalization technique is the *total counts*[43] which is convenient
162 for highly unbalanced OTU distribution among samples.

163 Using a simulation-based approach, we addressed the question of the performance of
164 the *non-null* and the *weighted non-null normalization* techniques when the matrix
165 sparsity and the number of observations increase. We compared these new methods
166 to those normally performed with DESeq2 and Phyloseq. To do so, we normalized
167 500 matrices with varied sparsity levels (i.e., 0.28, 0.64 and 0.82) and a different
168 number of observations (i.e., 4, 10 and 30). We calculated the average coefficient of
169 variation (CVmean)[44] for each normalization method (Fig. S2). Considering that
170 these OTUs are assumed to have relatively constant abundance within the simu-
171 lations, the coefficient of variation is expected to be lower when the normalization
172 is more efficient. Overall, the *non-null* and the *weighted non-null* normalization
173 methods exhibited a lower coefficient of variation as compared to the other meth-
174 ods, when sparsity in the count matrix is high and the number of observations
175 is increased. These differences were clear especially when comparing DESeq2 and
176 Phyloseq to the *weighted non-null normalization* (sparsity ratio of 0.28, 0.64 and
177 0.82, with 30 samples; t-tests $p < 0.001$) (Fig. S2).

178 *Contingency table filtering*

179 In metataxonomic studies, contingency tables are often very sparse and after sta-
180 tistical analysis, some significant differences among groups may not be of great
181 relevance. This may arise when a feature, distributed in many samples with a low
182 abundance, is slightly more abundant in a group of comparison. These artifacts
183 are generally excluded by DESeq2 with an independent filtering. Furthermore, if
184 a feature is found in a few samples only, it may lead to non-reliable results when
185 its abundance is high (when it is not 0). Such distributions may also impact the
186 normalization process as well as the dispersion estimates. In order to avoid misin-
187 terpretation of results, we propose an optional extra-step of filtering: by excluding
188 features characterized by a low abundance and/or a low number of occurrence in
189 samples (e.g. features occurring in less than 20% of the samples). To set a by-default
190 abundance threshold, SHAMAN search for an inflection point at which the curve

191 between the number of observations and the abundance of feature changes from
 192 being linear to concave. This process is performed with linear regression in the
 193 following manner:

- 194 1 We define I the interval $\left[\min_j (\sum_i c_{ij}) ; \frac{\sum_{ij} c_{ij}}{k} \right]$.
- 195 2 For each $x \in I$, we compute $h(x)$ defined as the number of observations with
 196 a total abundance higher than x .
- 197 3 We compute the linear regression between $h(x)$ and x .
- 198 4 The intercept is set as the default threshold.

199 (see Appendix 1 for more information). This extra-filtering aims at refining the
 200 first filtering processed with DESeq2 and may lead to a significant decrease of the
 201 computation time. The impact of filtering steps may be visually assessed with plots
 202 displaying the features that will be included in the analysis and those that will be
 203 discarded.

204 *Statistical modelling*

205 The statistical model relies on the variables that are available in the file of experi-
 206 mental design. By default, all variables are included in the model but the end-user
 207 can edit this selection and further add interactions between variables of interest.
 208 In addition, other variables such as batches or clinical data (e.g., age, sex, etc.)
 209 may be used as covariates. SHAMAN then automatically checks whether the model
 210 is statistically suitable (i.e., whether all the model parameters may be estimated).
 211 When it is not the case, an warning message appears and a "how to" box proposes
 212 a practical way to solve the issue. In SHAMAN, statistical models may be fitted
 213 at any taxonomic levels. Normalized counts are summed up within a given a taxo-
 214 nomic level.

215 To extract features that exhibit significant differential abundance (between two
 216 groups), the user must define a contrast vector. Both a guided mode and an
 217 expert mode are available in SHAMAN. In the guided mode, the user spec-
 218 ifies the groups to be compared using a dropdown menu. This mode is only
 219 available for DESeq2 v1.6.3 which is implemented in DESeq2shaman package
 220 (<https://github.com/ghozlane/DESeq2shaman>). In advanced comparisons, the
 221 user may define a contrast vector by specifying coefficients (e.g., -1, 0, 1) assigned
 222 to each variable.

223 **Visualization**

224 After running a statistical analysis, many displays are available:

225 (i) Diagnostic plots (such as barplots, boxplots, PCA, PCoA, NMDS and hierarchi-
 226 cal clustering) help the user examine both raw and normalized data. For instance,
 227 these plots may reveal clusters, sample mislabelling and/or batch effects. Scatter-
 228 plots of size factors and dispersion (i.e., estimates that are specific to DESeq2) are
 229 useful when assessing both the relevance and robustness of statistical models. PCA-
 230 and PCoA-plots associated with a PERMANOVA test may be used as preliminary
 231 results in the differential analysis as they may reveal global effects among groups
 232 of interest.

233 (ii) Significant features are gathered in a table including, the base mean (mean of
 234 the normalized counts), the fold change (i.e., the factor by which the average abun-
 235 dance changes from one group to the other), as well as the corresponding adjusted

236 p-values. The user may view tables for any contrasts and can export it into several
237 formats. Volcano plots and bar charts of p-values and log₂ fold change are also
238 available this section.

239 (iii) A global visualization section provides a choice of 9 interactive plots such as
240 barplots, heatmaps and boxplots to represent differences in abundance across groups
241 of interest. Diversity plots display the distribution of various diversity indices: alpha,
242 beta, gamma, Shannon, Simpson and inverse Simpson. Scatterplots and network
243 plots show association between feature abundance with other variables from the
244 target file. To explore variations of abundance across the taxonomic classification,
245 we included an interactive abundance tree and a Krona plot[45]. Rarefaction curves
246 are of great use to further consider the number of features in samples with respect
247 to the sequencing depth.

248 (iv) In the comparison section, plots displaying comparisons among contrasts may
249 be created. It includes several options such as, Venn diagram or upsetR graph[46]
250 (displaying subsets of common features across contrast), heatmap, a logit plot[47]
251 (showing the log₂ fold-change values in each feature), a density plot and a multiple
252 Venn diagram to summarize the number of features captured by each contrast. All
253 these graphs can be exported into four format (eps, png, pdf and svg).

254 Results

255 Comparison of SHAMAN with other available tools for meta-taxonomic analyses
256 A brief qualitative assessment of the strengths and limits of SHAMAN was done
257 in comparison with other similar web interfaces (Table 1). We first identified a list
258 of important considerations that have practical implications for the user such as
259 processing of raw sequencing data, statistical workflow, visualization, data storage
260 and accessibility. For each similar web interface, we then evaluated whether it met
261 those criteria. Besides that SHAMAN presents a number of advantages, we think
262 that such nested solution is essential for a careful interpretation of the results. Any
263 results in SHAMAN may be cross-checked with a quantification or an annotation
264 performed at an earlier stage. Furthermore several applications presented in Table
265 1, impose the burden to import/export R objects which requires skills in R pro-
266 gramming. It may also represent a source of issues for reproducibility, notably in
267 terms of compatibility of the packages over time.

268 User cases

269 To illustrate how SHAMAN works, we performed the analysis of two sequencing ex-
270 periments: a mock sequencing and a published dataset, afribiota dataset[48]. In both
271 analyses, we submitted the raw FASTQ files and provided a target file containing
272 sample information (needed for statistical analysis).

273 *Zymo Mock dataset*

274 The mock sequencing (EBI ENA code PRJEB33737) of the ZymoBIOMICS™
275 Microbial Community DNA was performed with an Illumina MiSeq. The Zymo
276 mock community is composed with 8 phylogenetically distant bacterial strains,
277 3 of which are gram-negative and 5 of which are gram-positive. DNA of two

278 yeast strains that are normally present in this community were not ampli-
279 fied. Genomic DNA from each bacterial strain was mixed in equimolar pro-
280 portions (<https://www.zymoresearch.com/zymbiomics-community-standard>). We
281 compared the impact of both the number of amplification cycles (25 and 30 cycles)
282 and the amount of DNA loaded in the flow cell (0.5ng and 1ng), on the microbial
283 abundance. Each sample was sequenced 3 times (experimental plan provided in
284 supplementary materials). Sequencing report provided by the sequencing facilities
285 indicated the presence of contaminants. To handle this issue, we filtered out the
286 genera occurring in less than 12 samples and outliers with a reduced log abundance
287 as compared to the other genera (Fig. S3). This process selected the 8 bacterial
288 stains of Zymo mock (Fig. 1). We then defined a statistical model that included
289 DNA amount and the number of amplification cycle as main effects and an in-
290 teraction between these variables. The statistical comparison showed a significant
291 impact of the number of amplification cycle compared to DNA amount. We found
292 no differential features between 0.5 ng and 1 ng DNA for each possible number of
293 cycle (25 and 30 cycles), while the comparison of number of amplification cycle for
294 each given amount of DNA showed significant impact on the abundance of mock
295 bacteria (Table S1, S2). These results are in agreement with previous studies that
296 presented the PCR-induced bias on equivalent mix[49, 50].

297 *Afribiota dataset*

298 The second dataset included samples of microbial communities in stunted chil-
299 dren aged 2-5y living in sub-Saharan Africa [48]. Three groups (nutritional sta-
300 tus) of individuals were considered: NN=non stunted, MCM=moderately stunted,
301 MCS=severely stunted. Samples originated from the small intestine fluids (gastric
302 and duodenal) and feces. The authors performed the bioinformatic treatment with
303 QIIME framework and the statistical analysis with several R packages including
304 Phyloseq for the normalization and DESeq2 for the differential analysis. 541 sam-
305 ples were available on EBI ENA (code PRJEB27868).

306 Using SHAMAN, raw reads were filtered against Human HG38 and PhiX174
307 genomes. A total of 2386 OTUs were calculated and 76% were annotated with
308 SILVA database at genus level. The sparsity rate of the contingency table was high
309 with 0.84. In consequence, we used the weighted non-null normalization which is
310 particularly efficient when the matrix highly sparse (Fig. S2).

311 Two analyses were performed, a global analysis that included duodenal, gastric as
312 well as feces samples and a more specific analysis including fecal samples only. Sta-
313 tistical models included the following variables, age, gender, country of origin and
314 nutritional status. Overall our results obtained when using SHAMAN were highly
315 consistent with those of Vonaesch *et al.* [48]. We detected a significant change in the
316 community composition between gastric and duodenal samples compared to feces
317 samples at Genus level (Fig. 2a) (PERMANOVA, $P=0.001$). The most abundant
318 genera were reported in Fig. S4. α -Diversity was not affected by stunting (Fig.2b).
319 We looked for a distinct signature of stunting in the feces. We report in the volcano
320 plot (Fig.2c) genera with differential abundance between stunt samples compared
321 to non-stunt (complete list available in Table S3). Twelve microbial taxa, corre-
322 sponding to members of the oropharyngeal core microbiota, were overrepresented

323 in feces samples of stunted children as compared with non-stunted children; more
324 particularly *Porphyromonas*, *Neisseria* and *Lactobacillus* (Fig.2d). These findings
325 were in agreement with the conclusions of the AfriBiota consortium while being
326 obtained within a few minutes of interaction with the SHAMAN interface.

327 Conclusions

328 SHAMAN enables user to run most of the classical metagenomics methods and
329 makes use of statistical analyses to provide support to each visualization. The pos-
330 sibility to deploy SHAMAN locally constitutes an important feature when the data
331 cannot be submitted on servers for privacy issues or insufficient internet access.
332 SHAMAN also simplifies the access to open computational facilities, making a care-
333 ful use of the dedicated server, galaxy.pasteur.fr.

334 During its development, we felt a strong interest of the metagenomics community.
335 We recorded 82 active users per month in 2019 (535 unique visitors in total) and 800
336 downloads of the docker application. We expect that SHAMAN will help researcher
337 performing a quantitative analysis of metagenomics data.

338 Data availability

339 Sequence reads of Zymo Mock have been deposited in the European Nucleotide Archive,
340 <https://www.ebi.ac.uk/ena/> accession no. PRJEB33737. The datasets generated, analysed during the current study
341 and the simulation script of the supplementary figure 2 are available in the microbiome repository:
342 10.6084/m9.figshare.11815860.

343 Software availability

344 Project name: SHAMAN
345 Project home page: <http://shaman.pasteur.fr>, <https://github.com/aghazlane/shaman>
346 Operating system: Platform independent
347 Programming language: R
348 Other requirements: Python 3
349 License: GNU GPL V3
350 Any restrictions to use by non-academics: No

351 Competing interests

352 The authors declare that they have no competing interests.

353 Author's contributions

354 SV, PL, PW, CM and AG developed SHAMAN; LM and SK performed the mock sequencing; SV, CM and AG
355 wrote the publication.

356 Acknowledgements

357 We thank Pascal Campagne for his comments, Hugo Varet for helpful discussions about DESeq2, Fabien Mareuil for
358 the help to deploy SHAMAN computation on Galaxy and Youssef Ghorbal for the maintenance of the databank, as
359 well as the IT System Department of Institut Pasteur, who manages installation and update of tools on TARS
360 cluster.

361 Author details

362 ¹Hub de Bioinformatique et Biostatistique – Département Biologie Computationnelle, Institut Pasteur, USR 3756
363 CNRS, 28 rue du Docteur Roux, 75015 Paris, France. ²Biomics – Département Génomes et Génétique, Institut
364 Pasteur, 28 rue du Docteur Roux, 75015 Paris, France.

365 References

- 366 1. Qin, N., Yang, F., Li, A., Pifti, E., Chen, Y., Shao, L., Guo, J., Le Chatelier, E., Yao, J., Wu, L., Zhou, J., Ni,
367 S., Liu, L., Pons, N., Batto, J.M., Kennedy, S.P., Leonard, P., Yuan, C., Ding, W., Chen, Y., Hu, X., Zheng,
368 B., Qian, G., Xu, W., Ehrlich, S.D., Zheng, S., Li, L.: Alterations of the human gut microbiome in liver
369 cirrhosis. *Nature* **513**, 59–64 (2014). doi:10.1038/nature13568
- 370 2. Pop, M., Walker, A.W., Paulson, J., Lindsay, B., Antonio, M., Hossain, M.A., Oundo, J., Tamboura, B., Mai,
371 V., Astrovskaya, I., Corrada Bravo, H., Rance, R., Stares, M., Levine, M.M., Panchalingam, S., Kotloff, K.,
372 Ikumapayi, U.N., Ebruke, C., Adeyemi, M., Ahmed, D., Ahmed, F., Alam, M.T., Amin, R., Siddiqui, S.,
373 Ochieng, J.B., Ouma, E., Juma, J., Mailu, E., Omere, R., Morris, J.G., Breiman, R.F., Saha, D., Parkhill, J.,
374 Nataro, J.P., Stine, O.C.: Diarrhea in young children from low-income countries leads to large-scale alterations
375 in intestinal microbiota composition. *Genome biology* **15**, 76 (2014). doi:10.1186/gb-2014-15-6-r76

- 376 3. Zeller, G., Tap, J., Voigt, A.Y., Sunagawa, S., Kultima, J.R., Costea, P.I., Amiot, A., Böhm, J., Brunetti, F.,
377 Habermann, N., Herczeg, R., Koch, M., Luciani, A., Mende, D.R., Schneider, M.A., Schrotz-King, P.,
378 Tournigand, C., Tran Van Nhieu, J., Yamada, T., Zimmermann, J., Benes, V., Kloor, M., Ulrich, C.M., von
379 Knebel Doeberitz, M., Sobhani, I., Bork, P.: Potential of fecal microbiota for early-stage detection of colorectal
380 cancer. *Molecular systems biology* **10**, 766 (2014). doi:10.15252/msb.20145645
- 381 4. Quereda, J.J., Dussurget, O., Nahori, M.-A., Ghazlane, A., Volant, S., Dillies, M.-A., Regnault, B., Kennedy,
382 S., Mondot, S., Villoing, B., Cossart, P., Pizarro-Cerda, J.: Bacteriocin from epidemic listeria strains alters the
383 host intestinal microbiota to favor infection. *Proceedings of the National Academy of Sciences of the United*
384 *States of America* **113**, 5706–5711 (2016). doi:10.1073/pnas.1523899113
- 385 5. Veiga, P., Gallini, C.A., Beal, C., Michaud, M., Delaney, M.L., DuBois, A., Khlebnikov, A., van
386 Hylckama Vlieg, J.E., Punit, S., Glickman, J.N., *et al.*: Bifidobacterium animalis subsp. lactis fermented milk
387 product reduces inflammation by altering a niche for colitogenic microbes. *Proceedings of the National*
388 *Academy of Sciences* **107**(42), 18132–18137 (2010). doi:10.1073/pnas.1011737107
- 389 6. Westcott, S.L., Schloss, P.D.: De novo clustering methods outperform reference-based methods for assigning
390 16s rRNA gene sequences to operational taxonomic units. *PeerJ* **3**, 1487 (2015). doi:10.7717/peerj.1487
- 391 7. Schloss, P.D., Westcott, S.L., Ryabin, T., Hall, J.R., Hartmann, M., Hollister, E.B., Lesniewski, R.A., Oakley,
392 B.B., Parks, D.H., Robinson, C.J., Sahl, J.W., Stres, B., Thallinger, G.G., Van Horn, D.J., Weber, C.F.:
393 Introducing mothur: open-source, platform-independent, community-supported software for describing and
394 comparing microbial communities. *Applied and environmental microbiology* **75**, 7537–7541 (2009).
395 doi:10.1128/AEM.01541-09
- 396 8. Edgar, R.C.: Uparse: highly accurate OTU sequences from microbial amplicon reads. *Nature methods* **10**,
397 996–998 (2013). doi:10.1038/nmeth.2604
- 398 9. Callahan, B.J., McMurdie, P.J., Rosen, M.J., Han, A.W., Johnson, A.J.A., Holmes, S.P.: Dada2:
399 High-resolution sample inference from illumina amplicon data. *Nature methods* **13**, 581–583 (2016).
400 doi:10.1038/nmeth.3869
- 401 10. Rognes, T., Flouri, T., Nichols, B., Quince, C., Mahé, F.: Vsearch: a versatile open source tool for
402 metagenomics. *PeerJ* **4**, 2584 (2016). doi:10.7717/peerj.2584
- 403 11. Caporaso, J.G., Kuczynski, J., Stombaugh, J., Bittinger, K., Bushman, F.D., Costello, E.K., Fierer, N., Pena,
404 A.G., Goodrich, J.K., Gordon, J.I., *et al.*: Qiime allows analysis of high-throughput community sequencing data.
405 *Nature methods* **7**(5), 335–336 (2010). doi:10.1038/nmeth.f.303
- 406 12. McDonald, D., Clemente, J.C., Kuczynski, J., Rideout, J.R., Stombaugh, J., Wendel, D., Wilke, A., Huse, S.,
407 Hufnagle, J., Meyer, F., *et al.*: The biological observation matrix (biom) format or: how i learned to stop
408 worrying and love the ome-ome. *GigaScience* **1**(1), 7 (2012). doi:10.1186/2047-217X-1-7
- 409 13. Paulson, J.N., Pop, M., Bravo, H.C.: Metastats: an improved statistical method for analysis of metagenomic
410 data. *Genome Biology* **12**(1), 17 (2011)
- 411 14. Paulson, J.N., Stine, O.C., Bravo, H.C., Pop, M.: Differential abundance analysis for microbial marker-gene
412 surveys. *Nature methods* **10**, 1200–1202 (2013). doi:10.1038/nmeth.2658
- 413 15. Love, M.I., Huber, W., Anders, S.: Moderated estimation of fold change and dispersion for rna-seq data with
414 deseq2. *Genome Biology* **15**, 550 (2014). doi:10.1186/s13059-014-0550-8
- 415 16. Robinson, M.D., McCarthy, D.J., Smyth, G.K.: edgeR: a bioconductor package for differential expression analysis
416 of digital gene expression data. *Bioinformatics* **26**(1), 139–140 (2010). doi:10.1093/bioinformatics/btp616
- 417 17. McMurdie, P.J., Holmes, S.: Waste not, want not: why rarefying microbiome data is inadmissible. *PLoS*
418 *computational biology* **10**(4), 1003531 (2014). doi:10.1371/journal.pcbi.1003531
- 419 18. Viktor Jonsson, O.N. Tobias Österlund, Kristiansson, E.: Statistical evaluation of methods for identification of
420 differentially abundant genes in comparative metagenomics. *BMC Genomics* **17**(1), 1–14 (2016).
421 doi:10.1186/s12864-016-2386-y
- 422 19. Escudé, F., Auer, L., Bernard, M., Mariadassou, M., Cauquil, L., Vidal, K., Maman, S., Hernandez-Raquet, G.,
423 Combes, S., Pascal, G.: Frogs: find, rapidly, OTUs with galaxy solution. *Bioinformatics* **34**(8), 1287–1294 (2017).
424 doi:10.1093/bioinformatics/btx791
- 425 20. Gonzalez, A., Navas-Molina, J.A., Kosciolk, T., McDonald, D., Vázquez-Baeza, Y., Ackermann, G., DeReus,
426 J., Janssen, S., Swafford, A.D., Orchanian, S.B., Sanders, J.G., Shorenstein, J., Holste, H., Petrus, S.,
427 Robbins-Pianka, A., Brislawn, C.J., Wang, M., Rideout, J.R., Bolyen, E., Dillon, M., Caporaso, J.G.,
428 Dorrestein, P.C., Knight, R.: Qiita: rapid, web-enabled microbiome meta-analysis. *Nature methods* **15**, 796–798
429 (2018). doi:10.1038/s41592-018-0141-9
- 430 21. McMurdie, P.J., Holmes, S.: Shiny-phyloseq: Web application for interactive microbiome analysis with
431 provenance tracking. *Bioinformatics (Oxford, England)* **31**, 282–283 (2015). doi:10.1093/bioinformatics/btu616
- 432 22. Wagner, J., Chelaru, F., Kancherla, J., Paulson, J.N., Zhang, A., Felix, V., Mahurkar, A., Elmqvist, N.,
433 Corrada Bravo, H.: Metaviz: interactive statistical and visual analysis of metagenomic data. *Nucleic acids*
434 *research* **46**(6), 2777–2787 (2018). doi:10.1093/nar/gky136
- 435 23. Huse, S.M., Mark Welch, D.B., Voorhis, A., Shipunova, A., Morrison, H.G., Eren, A.M., Sogin, M.L.: Vamps: a
436 website for visualization and analysis of microbial population structures. *BMC bioinformatics* **15**, 41 (2014).
437 doi:10.1186/1471-2105-15-41
- 438 24. Afgan, E., Baker, D., Batut, B., van den Beek, M., Bouvier, D., Cech, M., Chilton, J., Clements, D., Coraor, N.,
439 Grünig, B.A., Guerler, A., Hillman-Jackson, J., Hiltmann, S., Jalili, V., Rasche, H., Soranzo, N., Goecks, J.,
440 Taylor, J., Nekrutenko, A., Blankenberg, D.: The galaxy platform for accessible, reproducible and collaborative
441 biomedical analyses: 2018 update. *Nucleic acids research* **46**, 537–544 (2018). doi:10.1093/nar/gky379
- 442 25. Sloggett, C., Goonasekera, N., Afgan, E.: Bioblend: automating pipeline analyses within galaxy and cloudman.
443 *Bioinformatics* **29**(13), 1685–1686 (2013). doi:10.1093/bioinformatics/btt199
- 444 26. Dickson, L.B., Jiolle, D., Minard, G., Moltini-Conclois, I., Volant, S., Ghazlane, A., Bouchier, C., Ayala, D.,
445 Paupy, C., Moro, C.V., *et al.*: Carryover effects of larval exposure to different environmental bacteria drive adult
446 trait variation in a mosquito vector. *Science advances* **3**(8), 1700585 (2017). doi:10.1126/sciadv.1700585
- 447 27. Langmead, B., Trapnell, C., Pop, M., Salzberg, S.L.: Ultrafast and memory-efficient alignment of short dna

- 448 sequences to the human genome. *Genome biology* **10**(3), 25 (2009). doi:10.1186/gb-2009-10-3-r25
- 449 28. Criscuolo, A., Brisse, S.: Alientrimmer: a tool to quickly and accurately trim off multiple short contaminant
450 sequences from high-throughput sequencing reads. *Genomics* **102**(5-6), 500–506 (2013).
451 doi:10.1016/j.ygeno.2013.07.011
- 452 29. Zhang, J., Kobert, K., Flouri, T., Stamatakis, A.: Pear: a fast and accurate illumina paired-end read merger.
453 *Bioinformatics* **30**(5), 614–620 (2013). doi:10.1093/bioinformatics/btt593
- 454 30. Pruesse, E., Quast, C., Knittel, K., Fuchs, B.M., Ludwig, W., Peplies, J., Glöckner, F.O.: Silva: a
455 comprehensive online resource for quality checked and aligned ribosomal rna sequence data compatible with
456 arb. *Nucleic acids research* **35**(21), 7188–7196 (2007). doi:10.1093/nar/gkm864
- 457 31. DeSantis, T.Z., Hugenholtz, P., Larsen, N., Rojas, M., Brodie, E.L., Keller, K., Huber, T., Dalevi, D., Hu, P.,
458 Andersen, G.L.: Greengenes, a chimera-checked 16s rna gene database and workbench compatible with arb.
459 *Applied and environmental microbiology* **72**, 5069–5072 (2006). doi:10.1128/AEM.03006-05
- 460 32. Tang, J., Iliev, I.D., Brown, J., Underhill, D.M., Funari, V.A.: Mycobiome: approaches to analysis of intestinal
461 fungi. *Journal of immunological methods* **421**, 112–121 (2015). doi:10.1016/j.jim.2015.04.004
- 462 33. Abarenkov, K., Henrik Nilsson, R., Larsson, K.-H., Alexander, I.J., Eberhardt, U., Erland, S., Høiland, K.,
463 Kjølner, R., Larsson, E., Pennanen, T., Sen, R., Taylor, A.F.S., Tedersoo, L., Ursing, B.M., Vrålstad, T.,
464 Liimatainen, K., Peintner, U., Kõljalg, U.: The unite database for molecular identification of fungi—recent
465 updates and future perspectives. *The New phytologist* **186**, 281–285 (2010).
466 doi:10.1111/j.1469-8137.2009.03160.x
- 467 34. Findley, K., Oh, J., Yang, J., Conlan, S., Deming, C., Meyer, J.A., Schoenfeld, D., Nomicos, E., Park, M.,
468 Program, N.I.S.C.C.S., Kong, H.H., Segre, J.A.: Topographic diversity of fungal and bacterial communities in
469 human skin. *Nature* **498**, 367–370 (2013). doi:10.1038/nature12171
- 470 35. Yarza, P., Yilmaz, P., Pruesse, E., Glöckner, F.O., Ludwig, W., Schleifer, K.-H., Whitman, W.B., Euzéby, J.,
471 Amann, R., Rosselló-Móra, R.: Uniting the classification of cultured and uncultured bacteria and archaea using
472 16s rna gene sequences. *Nature reviews. Microbiology* **12**, 635–645 (2014). doi:10.1038/nrmicro3330
- 473 36. Wang, Q., Garrity, G.M., Tiedje, J.M., Cole, J.R.: Naive bayesian classifier for rapid assignment of rna
474 sequences into the new bacterial taxonomy. *Applied and environmental microbiology* **73**, 5261–5267 (2007).
475 doi:10.1128/AEM.00062-07
- 476 37. Katoh, K., Standley, D.M.: Mafft multiple sequence alignment software version 7: improvements in performance
477 and usability. *Molecular biology and evolution* **30**(4), 772–780 (2013). doi:10.1093/molbev/mst010
- 478 38. Criscuolo, A., Gribaldo, S.: Bmge (block mapping and gathering with entropy): a new software for selection of
479 phylogenetic informative regions from multiple sequence alignments. *BMC evolutionary biology* **10**(1), 210
480 (2010). doi:10.1186/1471-2148-10-210
- 481 39. Price, M.N., Dehal, P.S., Arkin, A.P.: Fasttree: computing large minimum evolution trees with profiles instead
482 of a distance matrix. *Molecular biology and evolution* **26**(7), 1641–1650 (2009). doi:10.1093/molbev/msp077
- 483 40. Lozupone, C., Knight, R.: Unifrac: a new phylogenetic method for comparing microbial communities. *Applied
484 and environmental microbiology* **71**, 8228–8235 (2005). doi:10.1128/AEM.71.12.8228-8235.2005
- 485 41. Sims, D., Sudbery, I., Iltott, N.E., Heger, A., Ponting, C.P.: Sequencing depth and coverage: key considerations
486 in genomic analyses. *Nature reviews. Genetics* **15**, 121–132 (2014). doi:10.1038/nrg3642
- 487 42. McMurdie, P.J., Holmes, S.: phyloseq: an r package for reproducible interactive analysis and graphics of
488 microbiome census data. *PLoS one* **8**(4), 61217 (2013). doi:10.1371/journal.pone.0061217
- 489 43. Evans, C., Hardin, J., Stoebel, D.M.: Selecting between-sample RNA-seq normalization methods from the
490 perspective of their assumptions. *Briefings in Bioinformatics* **19**(5), 776–792 (2017). doi:10.1093/bib/bbx008
- 491 44. Dillies, M.-A., Rau, A., Aubert, J., Hennequet-Antier, C., Jeanmougin, M., Servant, N., Keime, C., Marot, G.,
492 Castel, D., Estelle, J., *et al.*: A comprehensive evaluation of normalization methods for illumina high-throughput
493 rna sequencing data analysis. *Briefings in bioinformatics* **14**(6), 671–683 (2013). doi:10.1093/bib/bbs046
- 494 45. Ondov, B.D., Bergman, N.H., Phillippy, A.M.: Interactive metagenomic visualization in a web browser. *BMC
495 Bioinformatics* **12**(1) (2011). doi:10.1186/1471-2105-12-385
- 496 46. Conway, J.R., Lex, A., Gehlenborg, N.: Upsetr: an r package for the visualization of intersecting sets and their
497 properties. *Bioinformatics* **33**(18), 2938–2940 (2017). doi:10.1093/bioinformatics/btx364
- 498 47. Hourdel, V., Volant, S., O'Brien, D.P., Chenal, A., Chamot-Rooke, J., Dillies, M.-A., Brier, S.: Memhdx: an
499 interactive tool to expedite the statistical validation and visualization of large hdx-ms datasets. *Bioinformatics*
500 **32**(22), 3413–3419 (2016). doi:10.1093/bioinformatics/btw420
- 501 48. Vonaesch, P., Morien, E., Andrianonimiadana, L., Sanke, H., Mbecko, J.-R., Huus, K.E., Naharimananirina,
502 T., Gondje, B.P., Nigatoloum, S.N., Vondo, S.S., *et al.*: Stunted childhood growth is associated with
503 decompartmentalization of the gastrointestinal tract and overgrowth of oropharyngeal taxa. *Proceedings of the
504 National Academy of Sciences* **115**(36), 8489–8498 (2018). doi:10.1073/pnas.1806573115
- 505 49. Acinas, S.G., Sarma-Rupavtarm, R., Klepac-Ceraj, V., Polz, M.F.: Pcr-induced sequence artifacts and bias:
506 insights from comparison of two 16s rna clone libraries constructed from the same sample. *Appl. Environ.
507 Microbiol.* **71**(12), 8966–8969 (2005)
- 508 50. Sipos, R., Székely, A.J., Palatinszky, M., Révész, S., Márialigeti, K., Nikolausz, M.: Effect of primer mismatch,
509 annealing temperature and pcr cycle number on 16s rna gene-targetting bacterial community analysis. *FEMS
510 microbiology ecology* **60**, 341–350 (2007). doi:10.1111/j.1574-6941.2007.00283.x

Table 1 Comparison of SHAMAN with other web interface for metataxonomic analysis.

Category	SHAMAN	FROGS	Qiita	Shiny-phyloseq	Metaviz	Vamps
OTU processing	Yes	Yes	Yes	No	No	No
Normalization	Yes	Yes	No	No	No	No
Modelisation	Yes	Manova	No	D	M	No
Diversity analysis	Yes	Yes	Yes	Alpha	Alpha	Alpha
Phylogenetic analysis	Yes	Yes	Yes	Yes	No	Tree
Feature abundance plots	Yes	Yes	Yes	Yes	Yes	Yes
Ordination plots	Yes	Yes	Yes	Yes	Yes	Yes
Network plots	Yes	No	No	Yes	No	Yes
Geographic distribution plots	No	No	No	No	No	Yes
Statistics plots	Yes	No	NR	Yes	NR	NR
Interactive visualization	31	2;P	3	8	9	17
Raw data storage	No	No	Yes	No	No	Yes
Result storage	Yes	No	Yes	No	No	Yes
Online web Interface	Yes	No	Yes	No	Yes	Yes
R packaging	No	NR	NR	Yes	Yes	NR
Docker	Yes	No	No	No	Yes	No
Conda	Yes	Yes	Yes	No	Yes	No

D: Export from DESeq2, M: Export from Metagenomeseq, NR: Non relevant feature, P: Import/Export to Phyloseq, Number of unique interactive visualization are reported for each application in section 'Interactive visualization'

512 Figures

Figure 1 Barplot of taxa abundance of ZYMO MOCK samples. We summed the abundance of the OTU annotated at genera level with SILVA database and plotted the average abundance per condition.

Figure 2 Afribiota study of small intestine fluids and feces from stunt children compared to non stunt. (a) PCoA plot the Bray-Curtis dissimilarity index of the samples. Duodenal samples are colored in blue, light blue for Gastric and orange for Feces. PERMANOVA test based on the sample type yielded a P value of 0.001. (b) Alpha diversity analysis of non-stunt (NN), moderately stunted (MCM) and severely stunted (MCS). Overlapping confidence interval indicates that the diversity are not different between NN, MCM and MCS in duodenal, gastric and feces samples. (c) Volcano plot of differentially abundant genera in the feces of stunt children compared to non-stunt. We plot the log₂ fold change against the -log₁₀ adjusted p-value. Microbial taxa in red correspond to an increase of abundance and in blue to a decrease abundance. Labeled dots correspond to taxa from oropharyngeal core microbiota. (d) Log₂ abundance of differential abundant taxa from oropharyngeal core microbiota in stunt and non-stunt children feces.

513 **Additional Files**

514 Appendix 1: Mathematical definition of the contingency table filtering

515 Let us denote by \mathcal{F} the entire set of OTU. We propose to only consider a subset \mathcal{R} of OTUs, for the analysis:516 $\mathcal{R} = \mathcal{L}_1 \cap \mathcal{L}_2$, where \mathcal{L}_1 is defined by

$$517 \quad \mathcal{L}_1 = \left\{ f \in \mathcal{F} \mid \sum_j \mathbb{1}_{\{c_{fj} > 0\}} \geq l_1 \right\} \quad \text{with} \quad l_1 = \lfloor 0.8 \times k_{max} \rfloor,$$

518 where c_{fj} is the abundance of the feature f in the sample j while k_{max} is the maximum number of samples in519 which the feature is found. The subset \mathcal{L}_2 corresponds to the features with a not too small abundance and is

520 defined by

$$521 \quad \mathcal{L}_2 = \left\{ f \in \mathcal{F} \mid \sum_j c_{fj} \geq l_2 \right\},$$

522 where l_2 is the intercept of the linear regression between the variables $y_k = \sum_i \mathbb{1}_{\{\sum_j c_{ij} > x_k\}}$ and

$$523 \quad x_k \in \left[\min_{f \in \mathcal{F}} \left(\sum_j c_{ij} \right); \lambda \right]. \quad \lambda \text{ is a tuning parameter whose default value is } \left[\frac{\sum_{ij} c_{ij}}{n} \times 0.05 \right].$$

524 Supplementary Figure 1: SHAMAN workflow.

Figure S1 Shaman workflow. SHAMAN can start from raw reads or from processed data. In this last case, it needs at least three tables, the count matrix, the annotation table and the metadata which can either be provided into three different files or by using the BIOM format. The user can then select the variables of interest and add some batch effects. Once run, contrast vectors can easily be defined and interactive visualizations are available.

525 Supplementary Figure 2: Boxplots of the average coefficient of variation for DESeq2, non-null, weighted non-null
526 and Phyloseq normalization.

Figure S2 Boxplots of the average coefficient of variation for DESeq2, non-null, weighted non-null and Phyloseq normalization. The sparsity of the count matrices is chosen within $\{0.28, 0.64, 0.82\}$ and the number of observations vary within $\{4, 10, 30\}$. 500 normalizations were performed at each level of sparsity and for each number of observations. The results were analyzed by using a t-test, $***p < 0.001$. DESeq2 normalization did not converged when the matrix sparsity and the number of observations was high (e.g., with a sparsity of 0.64 and 30 observations).

527 Supplementary Figure 3: Filtering of Zymo mock contaminant genera according to occurrence in samples and
528 abundance.

Figure S3 Filtering of Zymo mock contaminant genera according to occurrence in samples and abundance. Genera occurring in at least 12 samples and with log abundance ≥ 5 were kept in the study (green dots). Cutibacterium, Pelomonas and Sphingomonas genera counts (red dots) were removed from the contingency table.

529 Supplementary Figure 4: Barplot of the 12 most abundant genera in duodenal, gastric and feces from children of
530 Central African Republic (CAR) and Madagascar.

Figure S4 Barplot of the 12 most abundant genera in duodenal, gastric and feces from children of Central African Republic (CAR) and Madagascar. Haemophilus and Neisseria genera were also observed in cultured samples.

531 Supplementary Table 1: Summary of taxa differentially abundant when compared 25 and 30 amplification cycles for
532 0.5ng DNA load.

ID	Base mean	Fold Change	Log2 fold change	P-value adjusted
Bacillus	23314	1.613	0.69	0.0.0058
Listeria	22410	1.633	0.708	0.0.0058
Salmonella	19308	1.767	0.822	0.0.0058
Pseudomonas	8696	1.609	0.686	0.0045
Staphylococcus	26546	1.496	0.58	0.0045
Escherichia-Shigella	13921	1.432	0.518	0.0051
Enterococcus	11024	1.417	0.504	0.0084

Table S1 Summary of taxa differentially abundant when compared 25 and 30 amplification cycles for 0.5ng DNA load Positive fold change indicates an increase of abundance of the taxa for the 25 amplification cycles.

533 Supplementary Table 2: Summary of taxa differentially abundant when compared 25 and 30 amplification cycles for
534 1ng DNA load.

Id	Base mean	Fold change	Log2 fold change	P-value adjusted
Enterococcus	11024	1.385	0.47	0.0303
Listeria	22410	1.427	0.514	0.0303
Salmonella	19308	1.473	0.559	0.0303
Staphylococcus	26546	1.381	0.465	0.0303
Escherichia-Shigella	13921	1.322	0.402	0.0393
Pseudomonas	8696	1.409	0.495	0.0414

Table S2 Summary of taxa differentially abundant when compared 25 and 30 amplification cycles for 1ng DNA load Positive fold change indicates an increase of abundance of the taxa for the 25 amplification cycles.

535 Supplementary Table 3: Summary of taxa differentially abundant when compared samples of stunted to non stunted
536 children.

Id	Base mean	Fold change	Log2 fold change	P-value adjusted
Porphyromonas	4.73	3.772	1.915	5.755e-11
Neisseria	7.16	2.789	1.48	1.605e-09
Lactobacillus	51.69	5.152	2.365	4.962e-08
Prevotella	5.04	2.354	1.235	5.836e-05
Ruminococcaceae UCG-009	3.46	0.419	-1.255	0.0003
Weissella	11.33	3.927	1.974	0.0003
Actinobacillus	3.98	2.672	1.418	0.0006
Rikenellaceae RC9 gut group	86.21	0.33	-1.598	0.0028
Ruminococcaceae UCG-011	502.47	1.85	0.887	0.0028
Streptococcus	237.2	1.723	0.785	0.0036
Christensenellaceae R-7 group	49.22	0.534	-0.906	0.0042
Aggregatibacter	4.31	2.222	1.152	0.0061
Granulicatella	0.81	1.991	0.993	0.0061
Ureaplasma	5.52	2.94	1.556	0.0061
Ruminococcaceae UCG-002	144.01	0.601	-0.735	0.0065
Streptobacillus	0.63	3.351	1.745	0.0065
Terrisporobacter	4.29	0.514	-0.961	0.0065
[Eubacterium] xylanophilum group	2.05	0.405	-1.305	0.0073
Fusobacterium	28.55	1.964	0.974	0.0124
Abiotrophia	0.45	2.733	1.45	0.0126
Campylobacter	89.02	1.739	0.798	0.0181
[Eubacterium] coprostanoligenes group	119.89	0.657	-0.606	0.0181
Capnocytophaga	0.19	3.005	1.587	0.0181
Ruminococcaceae UCG-005	211.21	0.683	-0.55	0.0192
Ruminococcaceae UCG-010	31.58	0.57	-0.81	0.0192
Veillonella	475.74	1.515	0.599	0.0218
Ruminococcaceae NK4A214 group	36.02	0.651	-0.62	0.0257
Morganella	0.42	4.099	2.035	0.0271
Haemophilus	260.89	1.494	0.58	0.0301
Family XIII AD3011 group	9.67	0.653	-0.614	0.0344
Lactococcus	35.73	2.161	1.112	0.0373
Methanobrevibacter	5.99	0.321	-1.638	0.0429
Turicibacter	10.03	0.534	-0.9	0.0429
Escherichia-Shigella	25.03	1.574	0.654	0.0444
Kingella	0.47	2.822	1.497	0.0488

Table S3 Summary of taxa differentially abundant when compared samples of stunted to non stunted children. Positive and negative fold changes indicate respectively an increase of abundance of the taxa in stunted and non-stunted children.

Figures

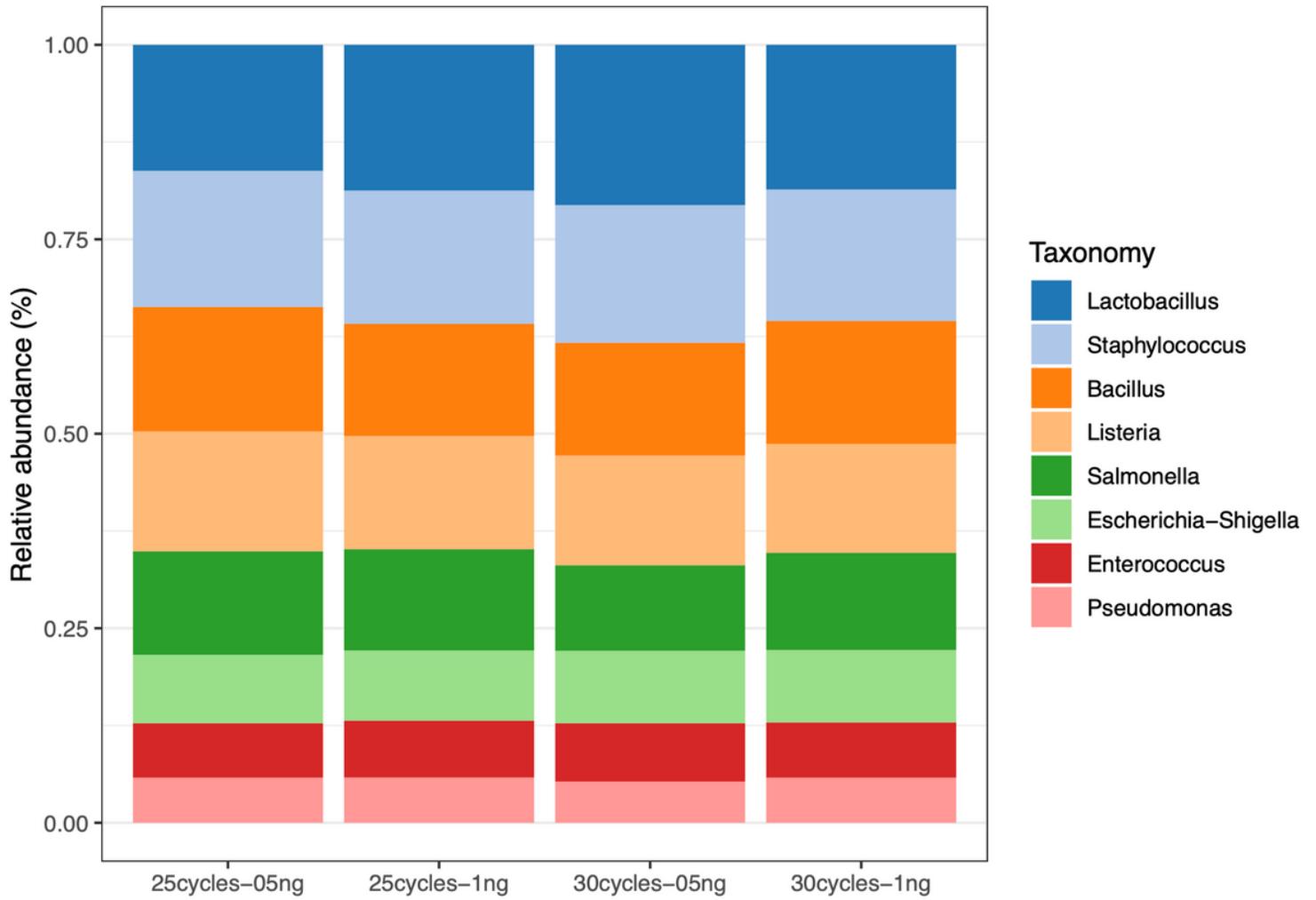


Figure 1

Barplot of taxa abundance of ZYMO MOCK samples. We summed the abundance of the OTU annotated at genera level with SILVA database and plotted the average abundance per condition.

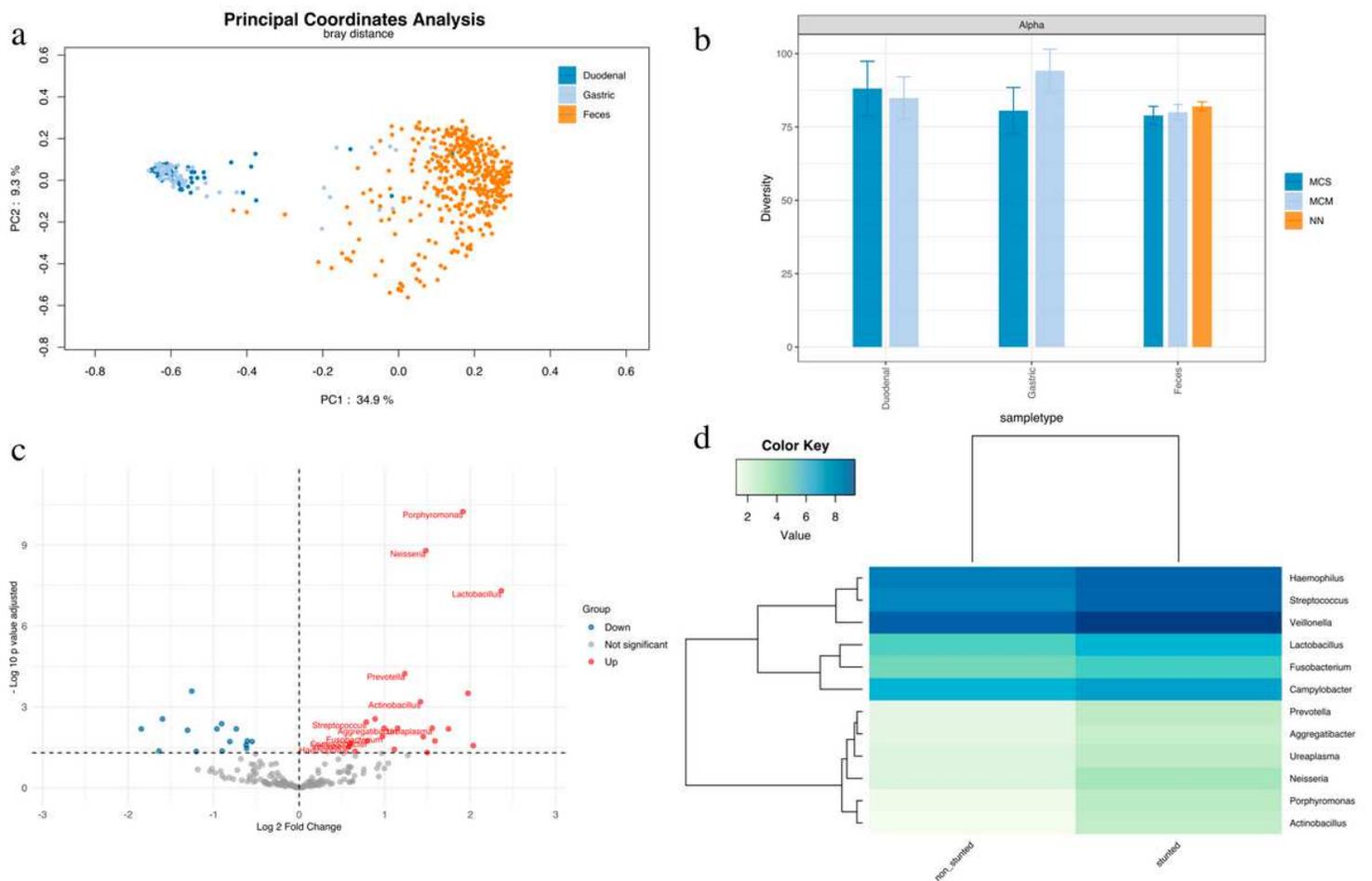


Figure 2

Afribiota study of small intestine fluids and feces from stunt children compared to non stunt. (a) PCoA plot the Bray-Curtis dissimilarity index of the samples. Duodenal samples are colored in blue, light blue for Gastric and orange for Feces. PERMANOVA test based on the sample type yielded a P value of 0.001. (b) Alpha diversity analysis of non-stunt (NN), moderately stunted (MCM) and severely stunted (MCS). Overlapping confidence interval indicates that the diversity are not different between NN, MCM and MCS in duodenal, gastric and feces samples. (c) Volcano plot of differentially abundant genera in the feces of stunt children compared to non-stunt. We plot the log₂ fold change against the -log₁₀ adjusted p-value. Microbial taxa in red correspond to an increase of abundance and in blue to a decrease abundance. Labeled dots correspond to taxa from oropharyngeal core microbiota. (d) Log₂ abundance of differential abundant taxa from oropharyngeal core microbiota in stunt and non-stunt children feces.

Supplementary Files

This is a list of supplementary files associated with this preprint. Click to download.

- [bmcarticle.bib](#)

- figS1.eps
- figS2.png
- figS3.png
- figS4.eps