

SHAMAN USER'S GUIDE

SHiny Application for Metagenomic ANalysis

19/03/2018

The screenshot shows the 'Welcome to SHAMAN' page. On the left, there's a sidebar with links: Home, Tutorial, Download/Install, Raw data, and Upload your data. The main content area has a title 'Welcome to SHAMAN' and a paragraph about the application. Below that is a diagram of the workflow:

```
graph LR; A[Input files] --> B[Count, Metadata target file or BOW]; B --> C[Normalization at OTU level, modified PLF approach]; C --> D[Merging normalized counts at the user selected level]; D --> E[Filtering the features, optional step]; E --> F[Run DESeq2]; F --> G[Statistical Analysis]; G --> H[Define contrast vectors (comparisons), Get differential abundance, Differential Analysis]; H --> I[Diagnostic plots]; I --> J[Visualization plots].
```

On the right, there's a 'What's new in SHAMAN' section with three entries:

- March 30th 2017 - Krona, Phylogeny and bug fixes**: KRONA and phylogenetic tree plots are now available in visualization. Several new distance are available in PCOA. The import float count matrices is now ok. We have finally debugged the export of the relative abundance/normalized matrices.
- Dec 9th 2016 - Phylogenetic tree and stress plot**: You can now upload a phylogenetic tree to calculate the unfrac distance (only available at the OTU level). The stress plot has been added to evaluate the goodness of fit of the NMDS.
- Nov 22th 2016 - New visualization and bug fix**: We have implemented a new visualization called tree abundance. Some bugs have been fixed (thanks Carine Rey from ENS).
- Oct 12th 2016 - Filtering step and bugs fix**: We have added a filtering step before the differential analysis.

At the bottom left, there are logos for Institut Pasteur and CTRIA.

shaman.pasteur.fr

Auteurs : Stevenn Volant, Amine Ghozlane

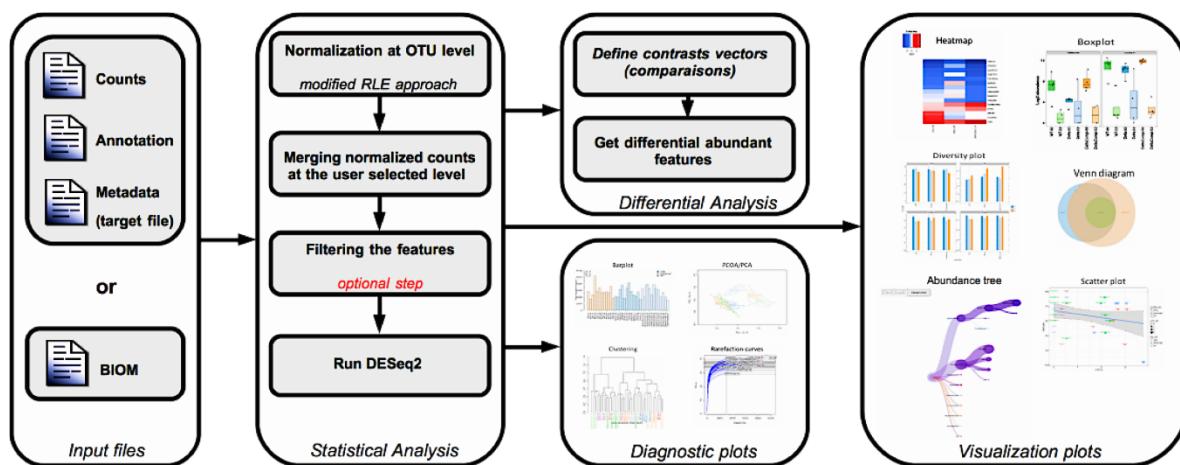
Introduction	3
Installation.....	4
Quick installation instructions.....	4
Tutorial	6
Homepage.....	6
Tutorial and Download	7
Raw data.....	8
<i>Read preparation.....</i>	8
<i>OTU processing</i>	9
Load count and annotation data.....	17
Statistical Analysis	21
Build the statistical model	21
<i>Loading the experimental design.....</i>	22
<i>Define the model.....</i>	23
Model options and normalisation.....	25
Data filtering (optionnel)	28
Define a contrast vector.....	29
Vérification du modèle statistique	30
Résultats de l'analyse différentielle	36
Visualizations	38
Visualisations des résultats	40
<i>Composition globale.....</i>	40
<i>Fold-change.....</i>	42
<i>Links between variables</i>	42
<i>Abondance et taxonomie.....</i>	43
Result comparisons	43
Bibliographie	46
Annexe A	47

INTRODUCTION

SHAMAN is a shiny application for differential analysis of metagenomic data (16S, 18S, 23S, 28S, ITS and WGS) including bioinformatics treatment of raw reads for targeted metagenomics, statistical analysis and result visualization with a large variety of plots (barplot, boxplot, heatmap, ...).

The statistical analysis performed by SHAMAN is based on DESeq2 R package [Anders and Huber 2010] which robustly identifies the differential abundant feature as suggested in [McMurdie and Holmes 2014, Jonsson2016].

SHAMAN is compatible with standard formats for metagenomic analysis (.csv, .tsv, .biom) and generated figures can be downloaded in several formats. Hereafter is the global workflow of the SHAMAN application:



INSTALLATION

Quick installation instructions

SHAMAN is available for local installation using Docker and R. This installation covers only the statistical analysis. The bioinformatics treatment is deported to the Institut Pasteur galaxy instance for performance reason.

- Docker install

Docker is the easiest way to use SHAMAN locally. It is a controlled virtual environment where every package required for SHAMAN are already installed and which have no impact on your local R installation.

First, download and install docker from <https://www.docker.com/>. Docker is available for Windows, Mac and Linux.

Run:

```
# Download shaman
docker pull aghozlane/shaman
# Execute shaman, port 80 and 5438 need to be available
docker run --rm -p 80:80 -p 5438:5438 aghozlane/shaman
```

Then, connect with your web browser to: <http://0.0.0.0/> or <http://localhost/>

If port 80 is already allocated, run:

```
docker run --rm -p 3838:80 -p 5438:5438 aghozlane/shaman
```

Then connect to <http://0.0.0.0:3838/> or <http://localhost:3838/> .

You can update your local version of SHAMAN with:

```
docker pull aghozlane/shaman
```

- Docker install

SHAMAN is available for R>3.1.2. The installation, download and execution can all be performed with a small R script:

```
# Load shiny packages
if(!require('shiny')){
    install.packages('shiny')
    library(shiny)
}
system("Rscript -e library(\"shiny\"); runGitHub(\"pierreLec/KronaRShy\", port=5438)",
wait=FALSE)
# Install dependencies, download last version from github,
# and run SHAMAN in one command:
runGitHub('aghozlane/shaman')
```

Of note, the R version has an impact on the contrast definition. DESeq2 contrast are harder to define on R>3.2 and some SHAMAN feature might be deactivated.

TUTORIAL

Homepage

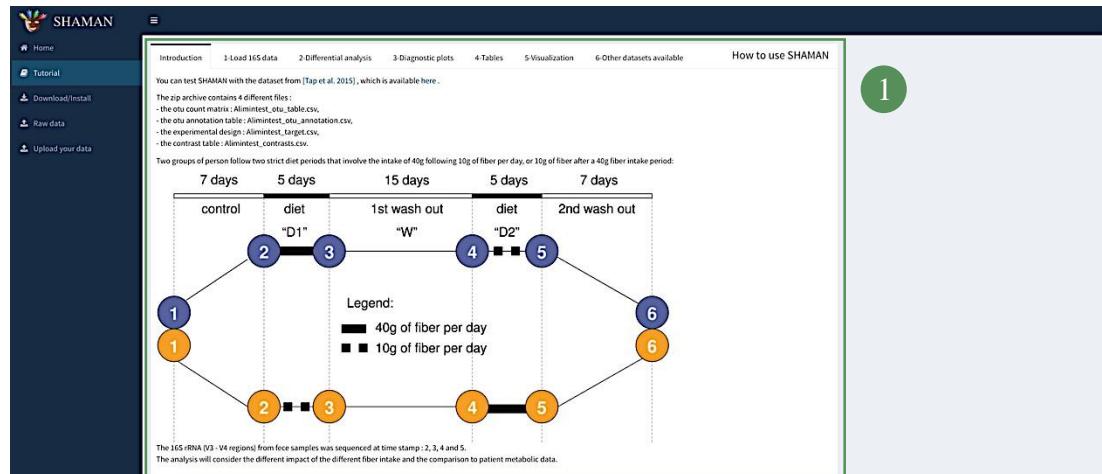
SHAMAN is available at <http://shaman.pasteur.fr/>. Hereafter is the homepage:

The screenshot shows the SHAMAN application homepage. On the left, there is a sidebar with a logo and five tabs: Home, Tutorial, Download/Install, Raw data, and Upload your data. A green circle labeled '1' is positioned next to the 'Upload your data' tab. The main content area has a dark header bar with the SHAMAN logo and a menu icon. Below the header, there are three sections: 'About', 'Authors', and 'Citing SHAMAN'. The 'About' section contains a detailed description of the application's purpose and methodology, mentioning DESeq2, RLE approach, PCOA, and ggplot2/D3.js packages. A green circle labeled '2' is positioned next to the 'About' section. The central part of the page displays a flowchart of the global workflow. It starts with 'Counts', 'Annotation', 'Metadata (target file)', or 'BIOM' input files. These lead to 'Normalization at OTU level (modified RLE approach)' and 'Merging normalized counts at the user selected level'. An optional step 'Filtering the features' leads to 'Run DESeq2'. This leads to 'Statistical Analysis' (which includes 'Define contrasts vectors (comparisons)', 'Get differential abundant features', and 'Differential Analysis'), 'Diagnostic plots' (including ChIP-seq, RNA-seq, and Metagenomic plots), and finally 'Visualization plots' (including Bar chart, Box plot, Violin plot, Abundance tree, and Stress plot). A red box highlights the 'Statistical Analysis' and 'Visualization plots' sections. A blue circle labeled '3' is positioned next to the 'Visualization plots' section. On the right side, there is a 'Welcome to SHAMAN' message and a 'What's new in SHAMAN' section. The 'What's new' section lists three updates: 'March 30th 2017 - Krona, Phylogeny and bug fixes', 'Dec 9th 2016 - Phylogenetic tree and stress plot', and 'Nov 22th 2016 - New visualization and bug fix'. Each update includes a brief description and a screenshot of the new feature.

- 1 Toolbar. 5 tab are available: Home, Tutorial, Download/Install, Raw data, Upload your data.
- 2 Description of the application
- 3 SHAMAN news.

Tutorial and Download

« Tutorial » and « Download/Install » panels are.



1

The screenshot shows the SHAMAN web interface. The left sidebar has links for Home, Tutorial, Download/Install (selected), Raw data, and Upload your data. The main content area is titled 'Download / Install SHAMAN'. It provides instructions for Docker install and R install (RC). The Docker install section includes steps for installing Docker on Ubuntu, Windows, and Mac, and running SHAMAN with commands like 'docker pull aghofzane/shaman && docker run -m -p 80:80 -p 5438:5438 aghofzane/shaman'. The R install (RC) section includes steps for installing the Docker toolbox and running 'Docker Quickstart Terminal'. It also provides instructions for connecting to the service via a web browser and updating the Docker container with 'docker pull aghofzane/shaman'.

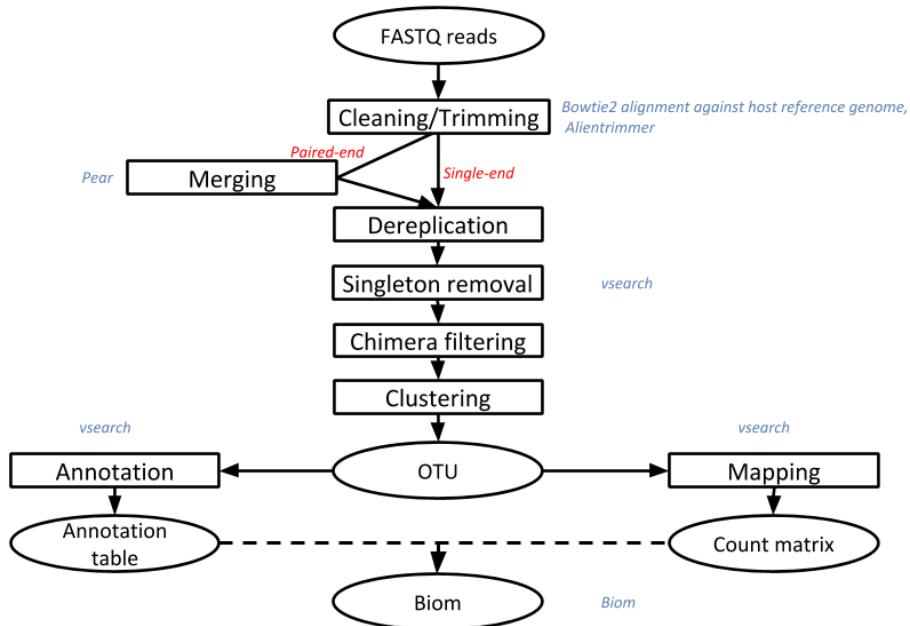
2

1 Online tutorial: it provides a description of SHAMAN usage for a set of 16S data. The dataset used is provided at download.

2 Installation guide. It is possible to install shaman with docker and R.

Raw data

SHAMAN provide access to a bioinformatics workflow for analyzing targeted metagenomics data. This workflow is based on *de novo* clustering. Operational Taxonomic Unit (OTU) are built from reads in a given experiment and annotated by alignment against reference databases. The aim is to identify among high quality amplicon, OTU sequences that will be representative of one species and considered for annotation and quantification. The workflow can be summaries as follow:



Read preparation

Cleaning

Cleaning step consists in the alignment of reads against the host and PhiX174 reference genome. Reads that did not align are considered for further analysis. This step is facultative but recommended. Very few amount of host DNA or Phi phage (used for calibration in every Illumina sequencing) can be identified in samples. A bowtie2 alignment with parameters (--sensitive) allows to eliminate these reads.

Trimming

READ	Sequence	GATTACA	...	TTA
	Quality	3031	...	161514
READ trimmed	Sequence	GATTACA	...	T
	Quality	3031	...	16

The trimming step consist to remove nucleotides sequences (adaptors, primers and non-confident nucleotides) in both 5' and 3' read ends. A list of every adaptors and primer used used in Illumina, Solid, Ion torrent, Truseq, and Nextera adaptors is

already included in SHAMAN. But user can specify his own adaptaters. This step is performed by Alientrimmer [Criscuolo 2013].

Merging

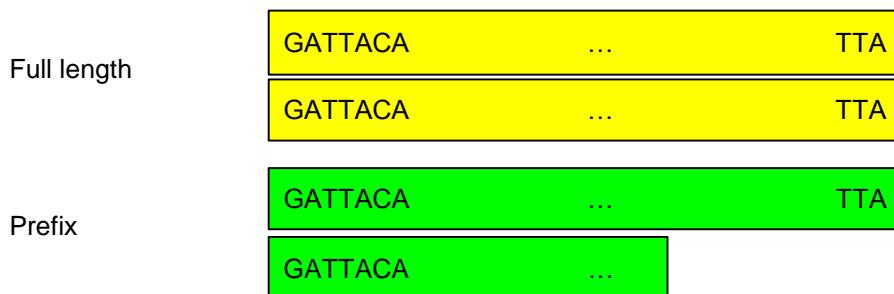


When reads are paired, a merging step is performed using Pear [Zhang 2013]. Reads are expected to overlap a given area of the ribosomal RNA. Sequence obtained are called amplicons.

OTU processing

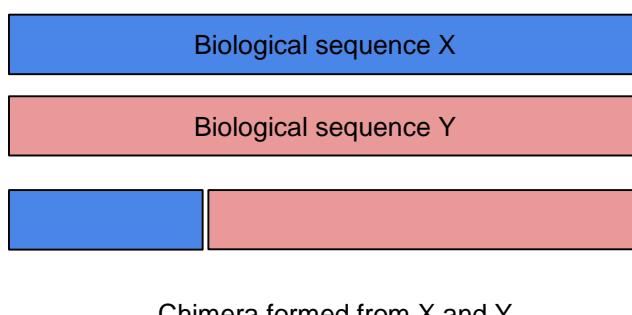
The OTU clustering step is performed with Vsearch [Rognes 2016].

Dereplication



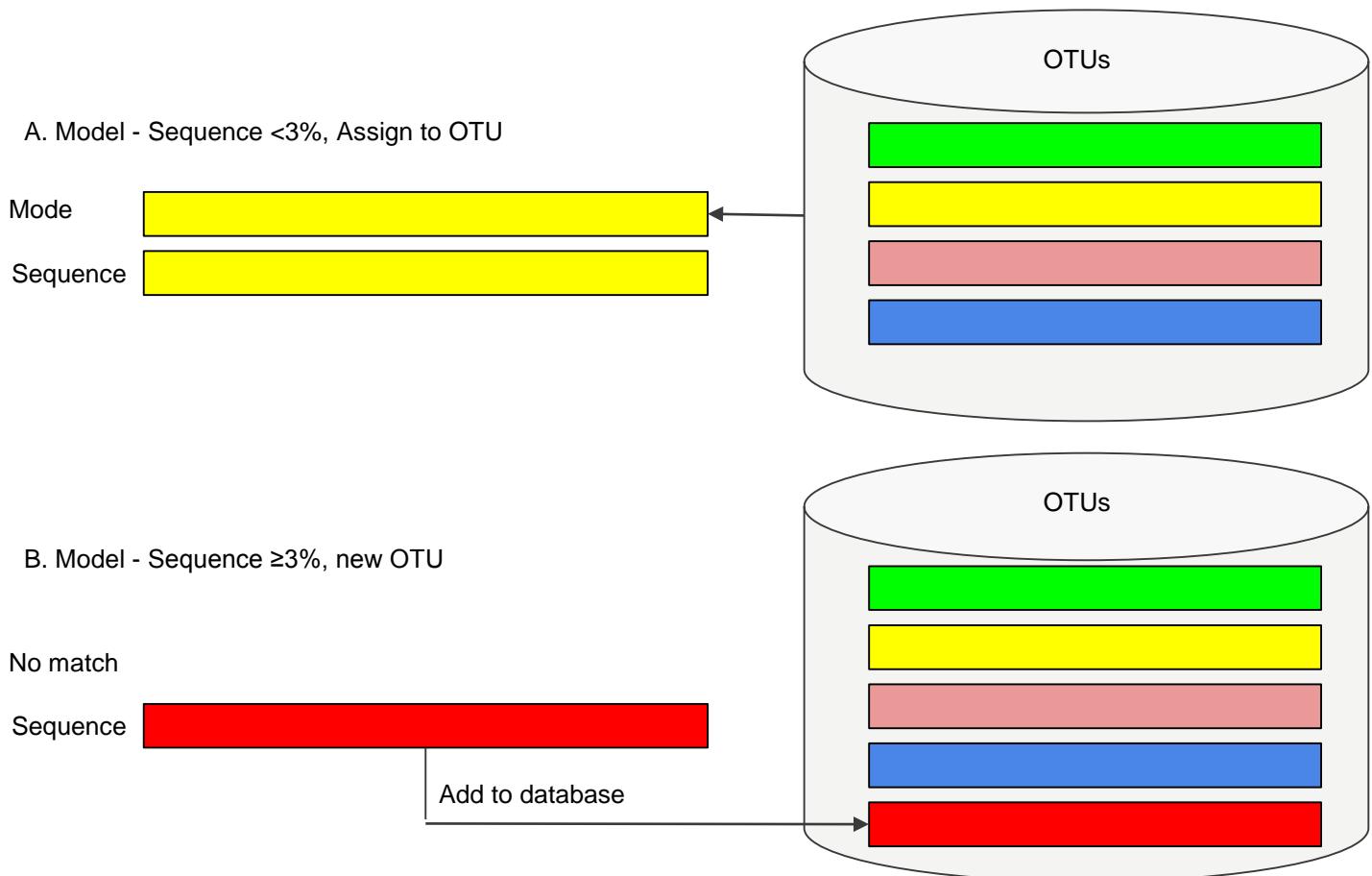
The dereplication consist in selecting one representative sequence among identical amplicon sequence. Two approach are available: the full length and the prefix dereplication. Full length approach will group sequence that are identical for their entire length. Prefix will also group sequence of different length. Only the longest sequence will be kept. The number of identical sequence present in a given sample is critical. Sequence are re-ordered by this “dereplicated abundance” which will drive the OTU clustering (see clustering section). Sequences with no identical respective are considered as singleton and are excluded of the OTU clustering process.

Chimera filtering



Chimera sequence are sequences that are composed by two other sequences. These sequences are filtered by a de novo approach consisting in a comparison with the other sequence in the dataset.

Clustering



The OTU clustering is performed with an Abundance-Greedy Clustering (AGC) [Westcott, Schloss, 2016 PeerJ; Rideout 2014; Schmidt et al. 2015]. The goal is to identify a set of correct biological sequences independently of sequencing error and ribosomal RNA single polymorphism. The algorithm can be described as follow:

Initial n groups ordered by their “dereplicated abundance”

Each step:

Pick a group and compare to the reference

If close to the reference:

Add in reference cluster

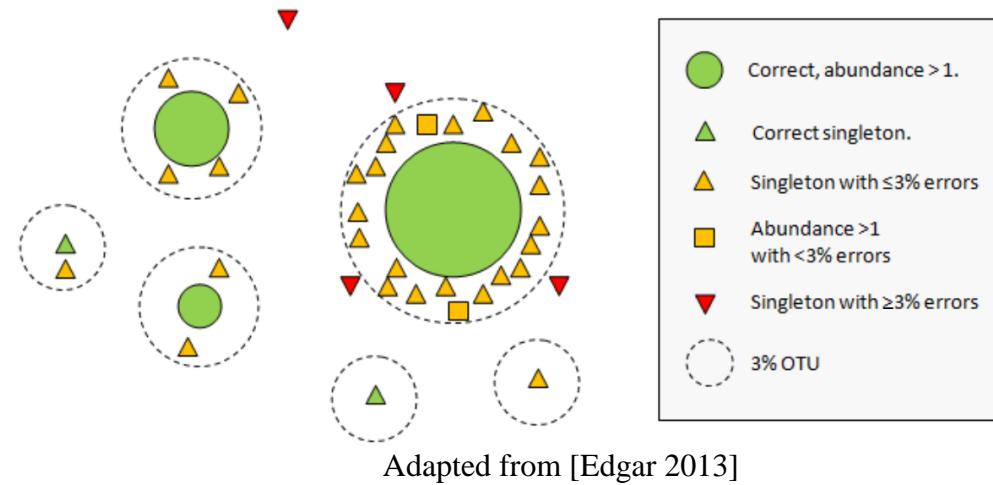
Otherwise:

Add it as a reference

AGC algorithm cost in the worst case $O(n^2)$ compared to $O(n^3)$ for hierarchical clustering approach implemented in Mothur for instance. The comparison is classically computed at threshold of 97% or 99%.

For 200 sequences, it represents 40000 operations for AGC compared to 8e+06 operations for hierarchical clustering. Sequences obtained after AGC are called OTU.

OTU abundance



The OTU abundance is calculated by aligning amplicon sequence against OTU sample per sample. The counts obtained are reproduced in a table called the count table.

OTU taxonomical annotation

The taxonomical annotation is performed by a fast Needleman-Wunsh alignment against ribosomal RNA databases. Annotation are filtered based on parameters identified in [Yarza 2014], as follow:

	Genus	Family	Order	Class	Phylum
Number of taxa	568	201	85	39	23
Median sequence identity	96.4% (96.2, 96.55)	92.25% (91.65, 92.9)	89.2% (88.25, 90.1)	86.35% (84.7, 87.95)	83.68% (81.6, 85.93)
Minimum sequence identity	94.8% (94.55, 95.05)	87.65% (86.8, 88.4)	83.55% (82.25, 84.8)	80.38% (78.55, 82.5)	77.43% (74.95, 79.9)
Threshold sequence identity	94.5%	86.5%	82.0%	78.5%	75.0%

Six databases are available in SHAMAN workflow:

- FINDLEY

Used for the taxonomical annotation of ITS sequences.

Findley, K., et al., Topographic diversity of fungal and bacterial communities in human skin. Nature, 2013, 498(7454), 367-370.

http://www.mothur.org/wiki/Findley_ITS_Database

- GREENGENES

Used for the taxonomical annotation of 16S, 18S sequences.

DeSantis, T. Z., et al. Greengenes, a chimera-checked 16S rRNA gene database and workbench compatible with ARB. Applied and environmental microbiology, 2006, 72(7), 5069-5072.

http://greengenes.secondgenome.com/downloads/database/13_5

- SILVA LSU, SSU

Used for the taxonomical annotation of 16S, 18S, 23S, 28S sequences.

Pruesse, E., et al., SILVA: a comprehensive online resource for quality checked and aligned ribosomal RNA sequence data compatible with ARB. Nucleic acids research, 2007, 35(21), 7188-7196.

<https://www.arb-silva.de/>

- UNDERHILL

Used for the taxonomical annotation of ITS sequences.

Tang J, Iliev I, Brown J, Underhill D and Funari V. Mycobiome: Approaches to Analysis of Intestinal Fungi. Journal of Immunological Methods, 2015, 421:112-21.

<https://riscweb.csmc.edu/microbiome/thf/>

- UNITE

Used for the taxonomical annotation of ITS sequences.

Abarenkov, K., et al., The UNITE database for molecular identification of fungi—recent updates and future perspectives. New Phytologist, 2010, 186(2), 281-285.

<https://unite.ut.ee/repository.php>

The databases are updated every 2 months.

This workflow is available in Raw data section, is installed on galaxy.pasteur.fr and makes profit on the Institut Pasteur cluster to run calculation. The usage of this workflow was simplified in this section.

The screenshot shows the SHAMAN workflow interface on a computer screen. The interface is divided into several sections:

- Top Left:** "SHAMAN" logo and navigation links: Home, Tutorial, Download/install, Raw data, Upload your data.
- Top Center:** "GET A KEY" button (disabled), "Require a valid email address".
- Top Right:** "FASTQ FILES" button (disabled), "Load the fastq files".
- Top Far Right:** "MATCHING PAIRS" button (disabled), "Only for paired-end sequencing".
- Left Column:** "About" section with "Type of data" radio buttons (16S, 18S, 23S/28S, ITS selected), "Email address" input field (yy@xx), "Get key" button, "Select the host" dropdown (None), "Specify the primer" checkbox, and "More workflow options" link.
- Middle Left:** "Directory containing the FastQ files" section with "Select your fastq files" button (disabled), "Browse..." button, and "No file selected" message.
- Middle Right:** "Enter the key" section with a large input field and "Check project" button.
- Bottom Left:** "Match the paired files (only for paired-end sequencing)" section with "Suffix R1 (Forward, do not include file extension)" dropdown (R1 selected), "Suffix R2 (Reverse, do not include file extension)" dropdown (R2 selected), "Match" button, and "Remove file(s)" link.
- Bottom Right:** "Summary of your analysis" section with "Type of data: ITS", "Paired-end sequencing: n", "Number of samples: 0", "Removed samples: 0", and "Save" button.

Four circular numbered callouts point to specific areas:

- Step 1: Top right of the "About" section.
- Step 2: "Select your fastq files" button in the middle left.
- Step 3: "Match the paired files" section at the bottom left.
- Step 4: "Summary of your analysis" section at the bottom right.

1

Main panel allows to describe data characteristics:

- type of sequencing (16S, 18S, 23S, 28S, ITS)
- Paired/single end sequencing
- Host of the microbiome
- Specify primers
- And workflow parameters

The user need to specify a mail and click on “get a key”. This key is the project identifier. It allows to follow calculation progress, view and download results.

2

Loading area for .fq, .fq.gz, .fastq and fastq.gz files

3

For paired end sequencing, files corresponding to R1 and R2 of a given sample need be paired.

A pattern is required to identify R1 and R2 files. The exclusion of this pattern must allow to get the same sample name.

Matched sample will be located next to its pair

4

Summary of the analysis.

For the example, we will load the mock sequencing () performed with an Illumina MiSeq sequencer performed at Institut Pasteur using microbial community DNA from Zymobiomics

(https://www.zymoresearch.com/media/amasty/amfile/attach/_D4300T_D4300_D4304_ZymoBIOMICS_DNA_Miniprep_Kit_1_1_2_LKN-SW_.pdf).

The screenshot shows the ZymoBIOMICS software interface. At the top, there are three green status bars: the first says "KEY CREATED! Your key is 30a47fca2dc3", the second says "FASTQ FILES 42 files are loaded", and the third says "PAIRS ARE MATCHED 21 samples are detected". A large green circle labeled "1" is positioned above the main interface.

About

Type of data: 16S 18S 23S/28S ITS
Email address *: amine.ghozlane@pasteur.fr
Specify the primer:
Paired-end sequencing ? Yes No
Select the host: None
More workflow options:

Directory containing the FastQ files

Select your fastq files: Browse... 42 files
Upload complete

List of the fastq files in the selected directory

```
Ing-25cycles_1_S38_L001_R1_001.fastq.gz
Ing-25cycles_1_S38_L001_R2_001.fastq.gz
Ing-25cycles_2_S39_L001_R1_001.fastq.gz
Ing-25cycles_2_S39_L001_R2_001.fastq.gz
Ing-25cycles_3_S40_L001_R1_001.fastq.gz
Ing-25cycles_3_S40_L001_R2_001.fastq.gz
Ing-25cycles_3_S50_L001_R1_001.fastq.gz
```

Match the paired files (only for paired-end sequencing)

Suffix R1 (Forward, do not include file extension): _R1
Suffix R2 (Reverse, do not include file extension): _R2

Ing-25cycles_1_S38_L001_R1_001.fastq.gz	Ing-25cycles_1_S38_L001_R2_001.fastq.gz
Ing-25cycles_2_S39_L001_R1_001.fastq.gz	Ing-25cycles_2_S39_L001_R2_001.fastq.gz
Ing-25cycles_3_S40_L001_R1_001.fastq.gz	Ing-25cycles_3_S40_L001_R2_001.fastq.gz
Ing-30cycles_1_S50_L001_R1_001.fastq.gz	Ing-30cycles_1_S50_L001_R2_001.fastq.gz

Match Remove file(s) Check and submit

A large red circle labeled "2" is positioned above the "Match the paired files" section.

1

When fastq file were successfully loaded, all panel turn green. The key 30a47fca2dc3 is now specific to my project.

2

Matched paired-end reads appears in the same order.

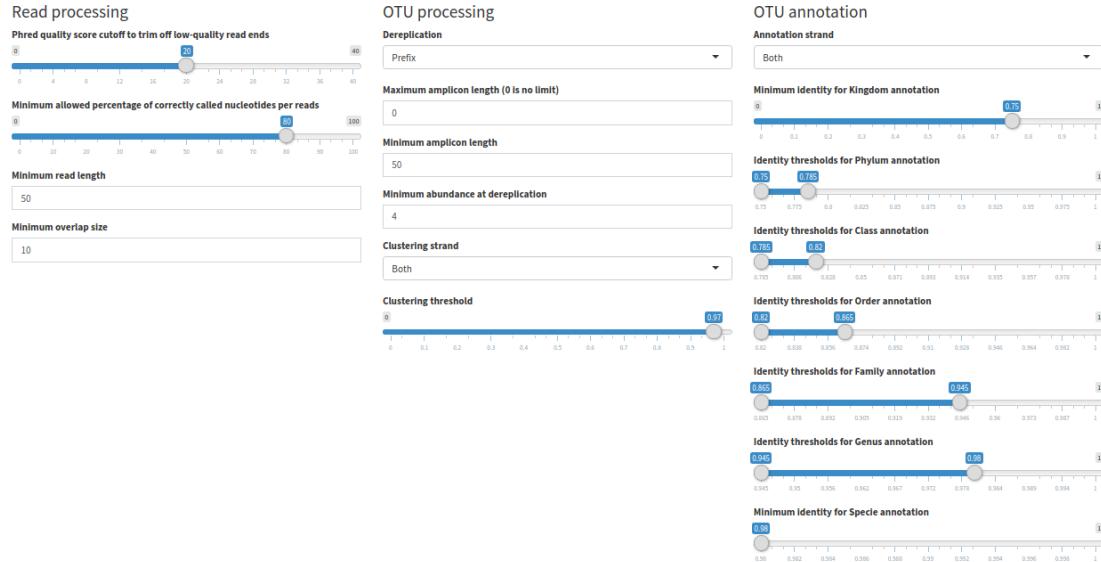
More options are available to specify primer used for sequencing and change workflow parameters when checking correspond boxes.

Specify the primer More workflow options

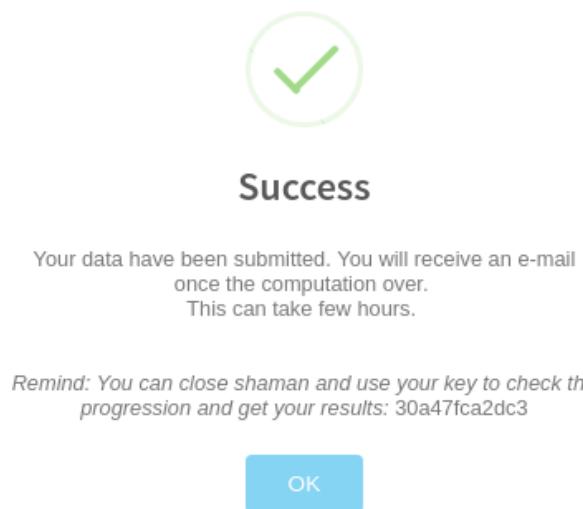
Primer sequences indicated in the following panels will be considered for read trimming.

Forward primer: TCGTCGGCAGCGCTCAGATGTGATAAGAGACAGCCCTACGGNGGCWGAG
Reverse primer: GTCTCGTGGGCTCGGAGATGTATAAGAGACAGGACTACHVGGGTATCTAATCC

Users can also change workflow parameters for read processing, OTU processing and OTU annotation with the following panel:



When parameters are identified, user can run calculation with the button “check and submit”. Successful submission will be notified by the following panel:



Calculation progress will be notified in the running interface with an interactive progress percentage:

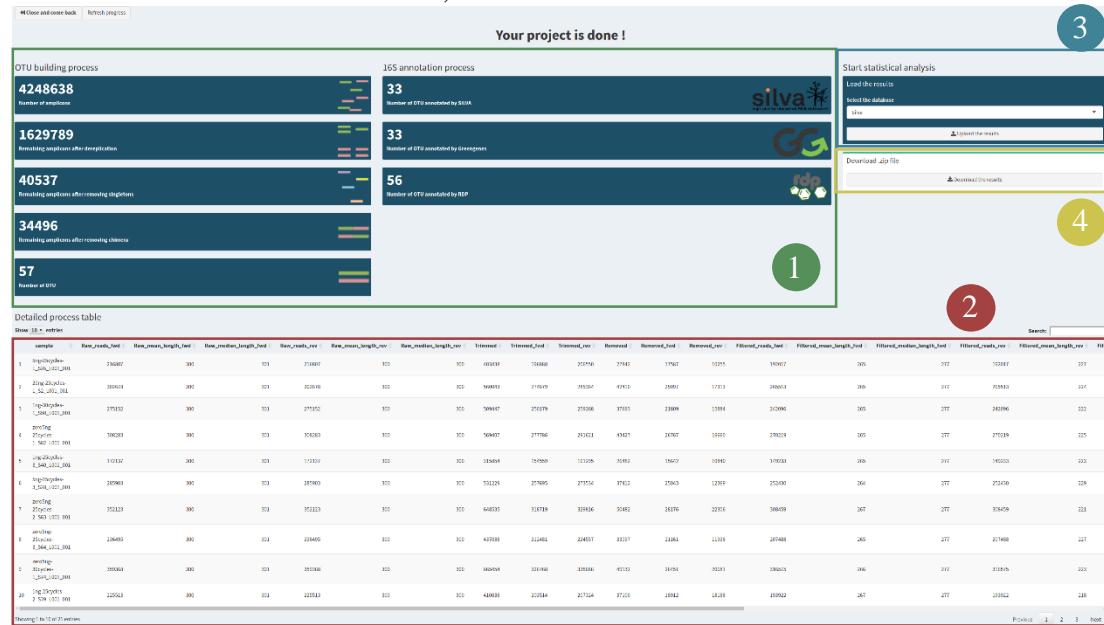


by entering the project key (remotely):



And by email.

When calculations are terminated, an interface is available to check



1 Information about the number of sequences obtained at each clustering step and annotated against ribosomal RNA databases.

2 Information about the read and OTU processing sample by sample

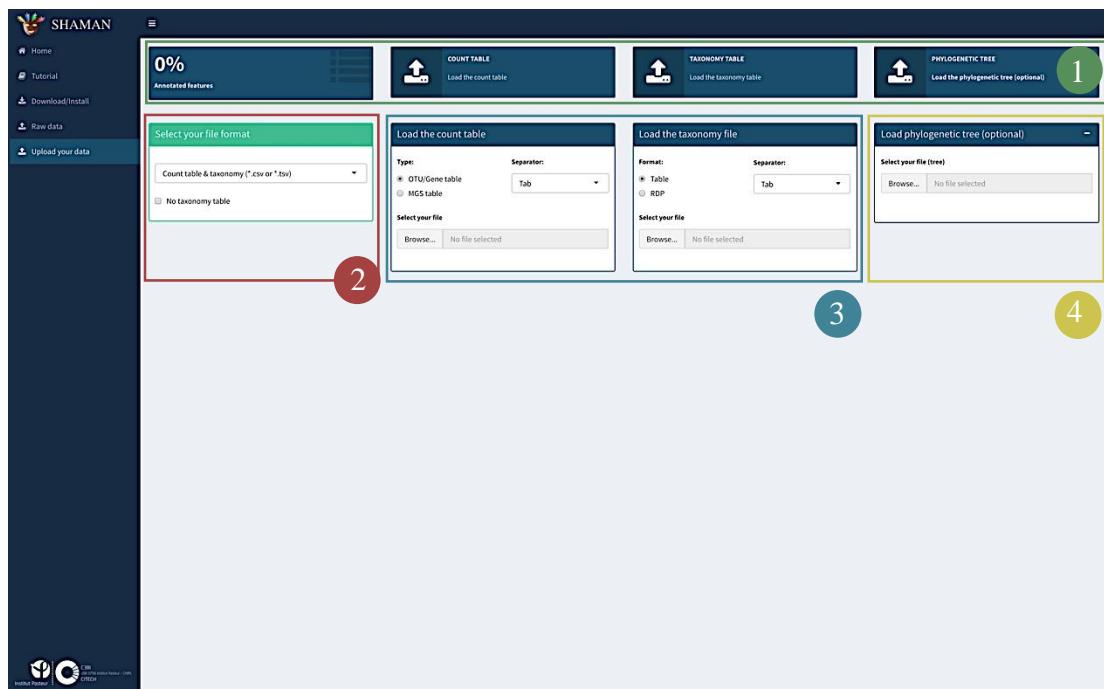
3 Panel allowing to load resulting data for corresponding databases.

4 Download all results.

A command line version of the workflow is available at <https://github.com/aghozlane/masque>.

Load count and annotation data

Processed data can be loaded in Shaman. User need a raw count table providing the number of read aligned against each OTU/gene for every sample and a matrix providing an annotation for each element. These information can also be provided as a biom file, as follow:



- 1 Information about the count table, the annotation table and the phylogenetic tree (optional)
- 2 Select the entry file format. User can choose to upload a count table and an annotation table or one file in biom format.

The user can also perform a study without an annotation file, in the case where the analysis is performed at the same level as the count file.
- 3 Load input files. The user can choose the type of data (OTU / Gene or MGS), the column separator and, for the RDP annotation, the probability threshold.
- 4 Load the phylogenetic tree .tree (optional)

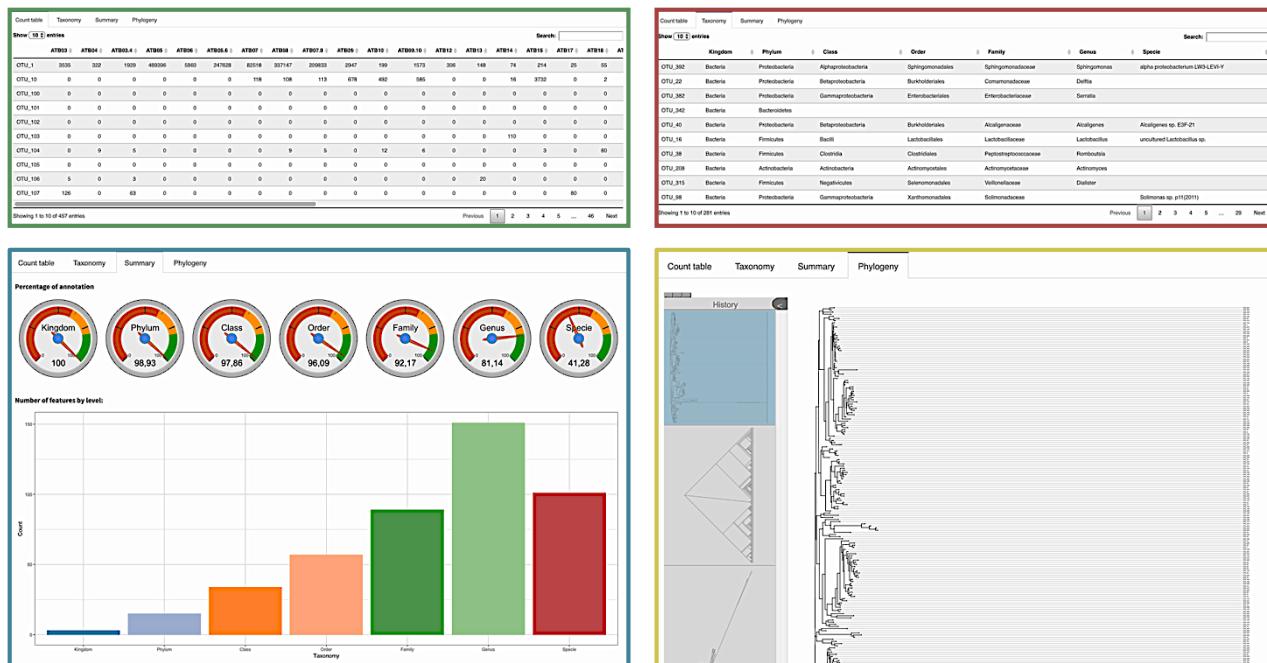
Once the data are loaded, SHAMAN displays the data of the user into two tables (count and annotation). Some graphical representations allowing to check the quality of the annotation are also available.

The screenshot shows the SHAMAN web application interface. At the top, there are three status boxes: 'Annotated features' (61.49%, green), 'COUNT TABLE' (Format seems to be OK, green), and 'TAXONOMY TABLE' (Format seems to be OK, green). Below these are four upload sections: 'Select your file format' (Count table & taxonomy (*.csv or *.tsv)), 'Load the count table' (OTU/Gene table selected, separator Tab, file chu_counts_table.tsv uploaded complete), 'Load the taxonomy file' (Table selected, separator Tab, file silva_annotation.tsv uploaded incomplete), and 'Load phylogenetic tree (optional)' (Browse phylo_tree.tree, Upload complete). A large central table displays the count table data, showing entries for OTU_1 through OTU_107 across various samples (ATB03 to ATB18). The table includes a search bar and navigation links for previous, next, and specific page numbers (1 to 46).

- 1 Successful file loading will end up to turn the boxes turn green and indicate the percentage of annotation.
- 2 Area presenting an overview of uploaded files and basic statistics. The details of the 4 tabs are given on page 19.

The count table must be in the following format: samples in column and individuals (OTU, ...) in row.

The 4 tabs of the previous box are the following



- 1 Count table with sample in column and OTU / genes in line.
- 2 Annotation table with different annotation levels in column and the OTU / genes on line.
- 3 Figure presenting the percentage of OTU / genes annotated at each annotation level and the number of sample at each level.
- 4 Figure of the phylogenetic tree (only if a phylogenetic .tree is loaded)

Once the data has been loaded, two new tabs appear "Statistical analysis" and "Visualization".

Statistical analysis - This tab is divided into 3 sub-tabs:

- Run differential analysis: allows to define the statistical model as well as some options, the taxonomic level and to create the vectors of contrasts for defining comparisons
- Diagnostic plots: many different visualizations (barplots, boxplots, clustering, PCA, PCoA, NMDS, ...) are available to control the design of the experiment, check the normalization and identify possible sequencing problems.
- Tables: get the results of the differential analysis for each defined contrast vector.

Visualization - This tab is divided into 2 sub-tabs:

- Global views: provides many interactive representations (barplots, heatmap, boxplots, krona, tree of abundance, curves of rarefaction, diversities) to visualize sample composition as well as the results of the analysis differential according to the experimental design.
- Plots comparison: a venn diagram and a heatmap of the log2-foldchange are available to compare the results of the analysis for each contrast (at least 2 contrast vectors have to be defined)

STATISTICAL ANALYSIS

SHAMAN statistics are based on DESeq2 package. This package requires the definition of a statistical model and to set up several parameters (not presented here). The different steps can be summarized as follow:

1. Data normalization (calculation of size factors)
2. Estimation of the dispersion based on a modelization of the average normalized counts and the empirical dispersion.
3. Adjustment of the generalized linear model based on the negative binomial distribution and a log link function.
4. Statistical test (Wald test) performed on model parameters
5. Outlier filtering
6. Correction of the multiple tests

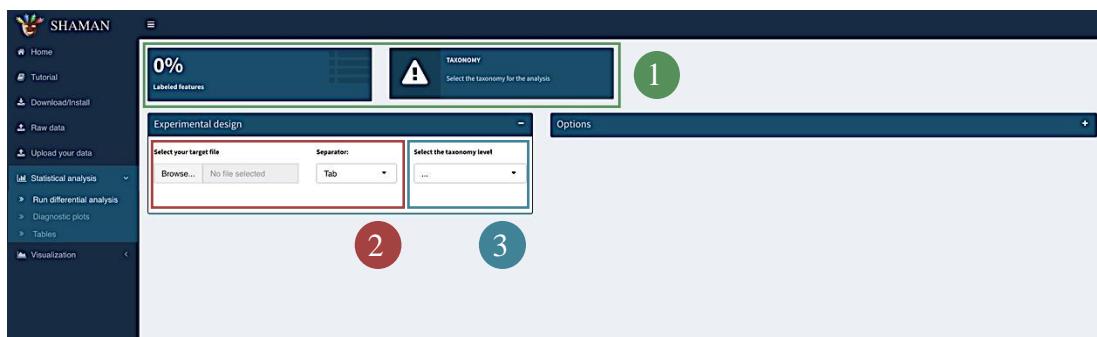
DESeq2 have been shown as one the best method to identify differentially abundant features.

Build the statistical model

In order to build the statistical model, the user must first provide the experimental design (target file) and select the taxonomic level. It is then necessary to select the variables of interest describing the biological phenomenon studied as well as the possible interactions between these variables. Before starting the analysis, the user can set the various options related to DESeq2 such as the independent filtering, the shape of the dispersion modeling, the method of correction of the multiple tests ... (for more information see [Love et al., 2014]).

All these steps must be carried out in the "Run differential analysis" sub-tab.

Loading the experimental design

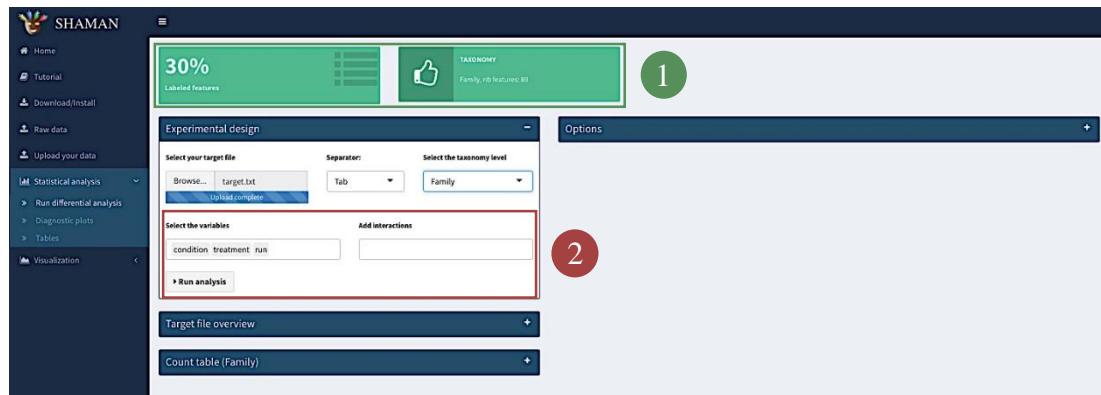


- 1 Information on experimental design and annotation level.
- 2 Area to load the experimental design. The user can also define the column separator used in the input file. This file must respect a certain format (see Appendix A for more information).
- 3 Choice of the annotation level. These levels correspond to those provided in the annotation file.

If no annotation file has been provided, the only option is "OTU/gene", which corresponds to the level at which the counting table was calculated.

Define the model

Once the target file is loaded and the taxonomic level is selected, SHAMAN checks that the sample names match the names of the counting table and merges counts of OTUs / genes with the same annotation.



- 1 Successful loading of the target file will end up to turn the boxes in green and to indicate the percentage of sample present in the counting table for which SHAMAN has information coming from the design loaded by the user.
- 2 Selection of the variables of interest (coming from the design). The variables can be qualitative and quantitative. It is also possible to add interactions between the variables.

In the example, the model has 3 variables:

- condition: 2 modalities "WT" and "KO"
- treatment: 2 treatments "A" and "B"
- run: "batch" effect which allows to take into account a 2-run sequencing

In this case, it would be wise to add an interaction between the condition and treatment variables. Indeed, it can be assumed that the effect of the treatment depends on the condition and therefore wants to test the effects by subgroups. Here, to simplify the notation this effect has been neglected.

Caution: It is recommended to avoid using numbers for qualitative variables. For instance, for the "run" variable, the user should prefer the notation "r1" and "r2" to "1" and "2".

The screenshot displays two tables side-by-side. The left table, titled 'Target file overview', lists samples with columns for sampleID, condition, treatment, and run. The right table, titled 'Count table (Family)', lists bacterial families with columns for SN03.4, SN07.8, SN14, ATB03.4, ATB07.8, and ATB12. Both tables include search bars and navigation buttons.

sampleID	condition	treatment	run
SN03.4	WT	A	r1
SN07.8	WT	A	r1
SN14	WT	A	r2
ATB03.4	ATB03.4	B	r1
ATB07.8	ATB07.8	B	r1
ATB12	ATB12	B	r2
ATB13	ATB13	KO	r2
ATB15	ATB15	KO	r2
ATB05.6	ATB05.6	KO	r1
SN13	SN13	KO	r2

	SN03.4	SN07.8	SN14	ATB03.4	ATB07.8	ATB12
Acetobacteraceae	9580	41940	494	12296	21174	1
Aeromonadaceae	2868	108150	74	46	63	1
Alcaligenaceae	0	0	1	115	0	1
Burkholderiaceae	0	6	1	87	3179	1
Cannabacteriaceae	0	17	194	0	0	1
Comamonadaceae	4	1	74	462	1	1
Corynebacteriaceae	0	2	65	22	2	1
Enterobacteriaceae	175580	678417	128624	3164	219289	52
Family XI	4	18	5	256	24	
Flavobacteriaceae	14333	16782	20885	18519	9546	1983

1

Aperçu du fichier target chargé par l'utilisateur. Ce fichier contient les informations permettant de décrire chaque échantillon par une ou plusieurs variables (qualitatives et/ou quantitatives).

L'utilisateur a la possibilité de supprimer certains échantillons de mauvaise qualité et d'exporter ce nouveau fichier target.

2

Aperçu de la table des comptages normalisés « agrégée » au niveau taxonomique sélectionné par l'utilisateur.

L'utilisateur a la possibilité d'exporter la table des comptages normalisés et/ou les abondances relatives

Model options and normalisation

The screenshot shows the 'Options' tab of the DESeq2 software interface. It has three main tabs: 'Statistical model' (selected), 'Normalization', and 'Filtering'. Under 'Statistical model', there are four sections: 'Type of transformation' (radio buttons for 'VST' and 'rlog', 'VST' is selected), 'Independent filtering' (radio buttons for 'True' and 'False', 'True' is selected), 'p-value adjustement' (radio buttons for 'BH' and 'BY', 'BH' is selected), and 'Level of significance' (text input field containing '0.05'). Below these are sections for 'Cooks cut-off' (radio buttons for 'Auto', 'No cut-off', and 'Value', 'Auto' is selected), 'Local function' (radio buttons for 'Median' and 'Shorth', 'Median' is selected), and 'Relationship' (radio buttons for 'Parametric' and 'Local', 'Parametric' is selected).

- **Type of transformation**

Two transformations are available in DESeq2, VST (Variance Stabilizing Transformation) and rlog (regularized log transformation). The data being heteroscedastic, this transformation makes it possible to eliminate the dependence of the variance on the mean. It is only used to visualize and/or classify the data (modeling is done on the counts).

When the number of samples is large, it is recommended to use VST transformation which is faster than the rlog.

- **Independent filtering**

This filter based on the average counts over all the samples. It aims at filtering the individuals (species, genera, ...) which are very unlikely to be differentially abundant. The threshold used for the filter is determined such that the number of significant individuals reach its maximum for a given FDR.

- **p-value adjustment**

Two usual methods for multiple correction by FDR are proposed: Benjamini-Hochberg et Benjamini-Yekuteli.

- **Level of significance**

Level of significance for the statistical test. By default, this value is 5%.

- **Cooks cut-off**

The log-fold change value is strongly influenced by outliers. To detect outliers, the Cooks distance which measure the impact of removing a sample on the estimation of the parameters can be used. The 99th percentile of the Fisher

distribution ($F(p, m-p)$, with p the number of parameters and m the number of samples) is used as a threshold for the Cooks distance.

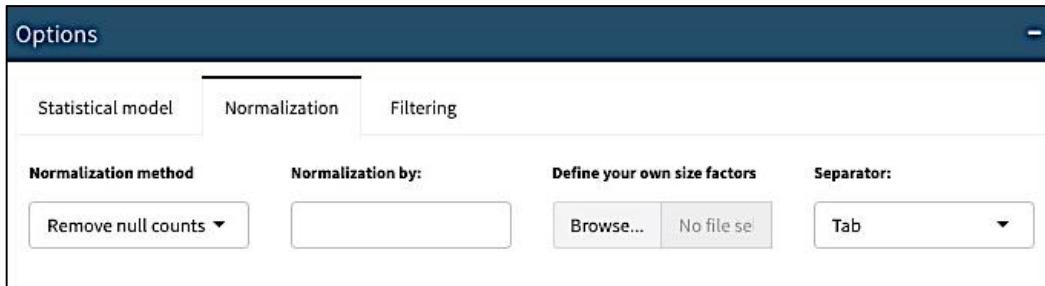
- **Local function**

To calculate the size factors (used to normalize the data), there are two options "median" or "shorth". The first one, "median" is the default method which aims at calculating the median of the ratio between counts and the geometric mean (see normalization section for more details). The second one, "shorth", calculates the average of the smallest interval that covers half of the values. This option is especially recommended for low counts.

- **Relationship**

For each feature, the dispersion is estimated by modeling the relationship between empirical dispersion and mean counts. To model this relation, three methods are proposed:

- ✓ a parametric regression of the form $\alpha_{tr}(\bar{\mu}) = \alpha_0 + \frac{a_1}{\bar{\mu}}$
- ✓ a local regression which makes it possible to obtain a better modeling when the form in "1/x" is not well adapted.
- ✓ the mean. This approach can be used when the number of individuals is small.



- **Normalization method**

SHAMAN proposes 4 normalization methods:

- ✓ Usual : default method of DESeq2 which consists in calculating the median of the ratio between the counts and the geometric mean

$$s_j = \text{median}_i \frac{K_{ij}}{\left(\prod_{k=1}^n K_{ik}\right)^{1/n}}$$

- ✓ Remove null counts : The calculation is done only with the non null counts

- $$s_j = \text{median}_i \frac{K_{ij}}{\left(\prod_{k \in S_i} K_{ik}\right)^{1/n_i}}$$
- ✓ Weighted : weighted version of the previous method (the weights are proportional to the number of samples with a non null value).
 - ✓ Total count : this method aims at calculating a size factor by dividing the total counts of each sample by the average of the totals over all the samples. This approach must be used when the composition of one condition to another is very different, i.e. when several species are present in one condition and absent in the other (and vice versa).

For metagenomic analysis, the matrices are very sparse(many 0). Therefore, it is recommended to avoid the use of the "usual" method.

- **Normalization by**

This option makes it possible to perform group normalization according to a variable of the experimental design.

This option can be useful when the user wants to visualize the results of several studies without modifying the normalization.

Caution: avoid selecting a variable of interest to normalize the data because it can create some biases.

- **Define your own size factors**

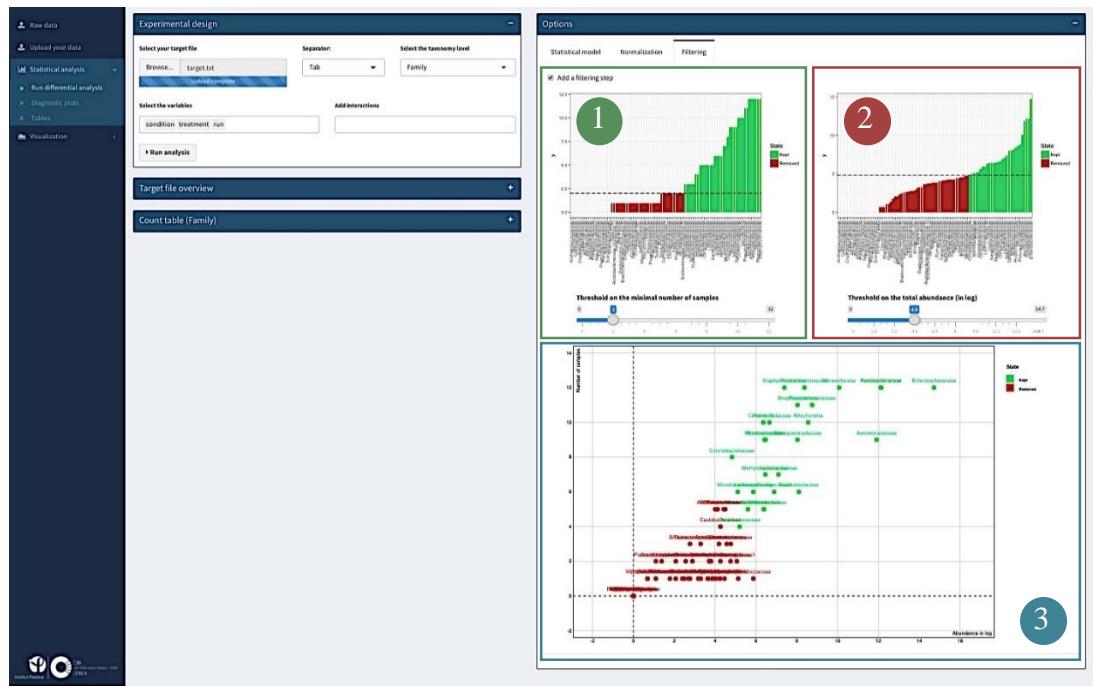
The user can load his own vector of size factors (obtained by another method for instance).

- **Separator**

File separator for the size factor file.

Data filtering (optionnel)

In some cases, when many elements have low counts and/or are only detected for a small number of samples, it may be useful to filter them before starting the analysis. This will not have a strong influence on the results of the differential analysis. Indeed, the independent filtering already filter these elements. For studies with many features, it allows to reduce the computation time.



- 1 Threshold on the occurrence in all samples. By default, the threshold set by SHAMAN corresponds to 20% of the samples.
- 2 Threshold on average abundance. To calculate a threshold automatically, SHAMAN performs a linear regression between the number of sample with an average abundance of at least x and the average abundance.
- 3 Figure representing the number of occurrence in all samples compared to the average abundance. The red dots correspond to the elements that will be filtered once the two filters are applied.

Define a contrast vector

$\beta_i^c = c^t \beta_i$ Once the statistical analysis is done, SHAMAN provides the list of parameters corresponding to the variables included in the model. From the estimation of these parameters, the user can test different effects. Denote by β_i the vector of parameters. Let c be a vector of contrasts, define as a linear combination of parameters. We then use a Wald test whose test statistic is distributed according to a standard normal distribution:

The matrix Σ_i corresponds to the variance-covariance matrix of the model $\frac{\beta_i^c}{\sqrt{c^t \Sigma_i c}} \sim N(0, 1)$ parameters.

Example 1

If the user wants to test whether there is a treatment effect (A vs B), it corresponds to test whether the difference A-B is zero. To do so, a contrast vector composed of a 1 and a -1 for the parameters associated with the treatments A and B has to be created. The hypotheses of the test will be as follows:

$$\begin{cases} H_0: \beta_A - \beta_B = 0 \\ H_1: \beta_A - \beta_B \neq 0 \end{cases}$$

In a complex experimental design, many comparisons can be done by defining a set of contrast vectors.

Example 2 :

Assume that we have an experimental design with measurements for 3 treatments A, B and C. The user wants to know if the effect of treatment C corresponds to the average of treatments A and B. In this case , 3 parameters will be estimated (one for each treatment). To get the desired comparison, the user must defined the following contrast vector:

$$\mathbf{c}^t \boldsymbol{\beta} = \left[\begin{array}{ccc} \frac{1}{2} & \frac{1}{2} & -1 \end{array} \right] \begin{bmatrix} \beta_1 \\ \beta_2 \\ \beta_3 \end{bmatrix} = \frac{1}{2}\beta_1 + \frac{1}{2}\beta_2 - \beta_3 = 0$$

The screenshot displays four panels of a software interface for defining contrasts:

- Panel 1:** Contrasts (New). Shows a dropdown menu "Compare" set to "A" and another dropdown "To" set to "B". A button "+ Add" is visible.
- Panel 2:** Contrasts (advanced user). Shows a "Select a file of contrasts" input field with "Browse..." and "No file selected" buttons, and a "Separator:" dropdown set to "Space".
- Panel 3:** Define contrasts by yourself. Contains a "Contrast name" input field and a "+ Add contrast" button. Below are several numerical input fields for different variables: Intercept (0), conditionKO (0), conditionWT (0), treatmentA (0), treatmentB (0), runr1 (0), and runr2 (0).
- Panel 4:** Defined contrasts. Shows a list of contrasts: "KO_vs_WT" and "A_vs_B". It includes a "Remove" button and an "Export" button.

- 1 Area for automatic definition of contrast vectors based on model variables. This panel allows to define most of the usual contrast vectors.
- 2 Loading area of the contrast file created from SHAMAN. The variables must exactly match those that were used to create the contrast file.
- 3 Area for manual definition of contrast vectors. This area is reserved for advanced users knowing how to define a vector of contrasts from the parameters of the statistical model.
- 4 List of contrasts created (by one of the 3 possibilities). The user can export the contrasts vectors as a .txt file (for later use) and / or remove some contrasts.

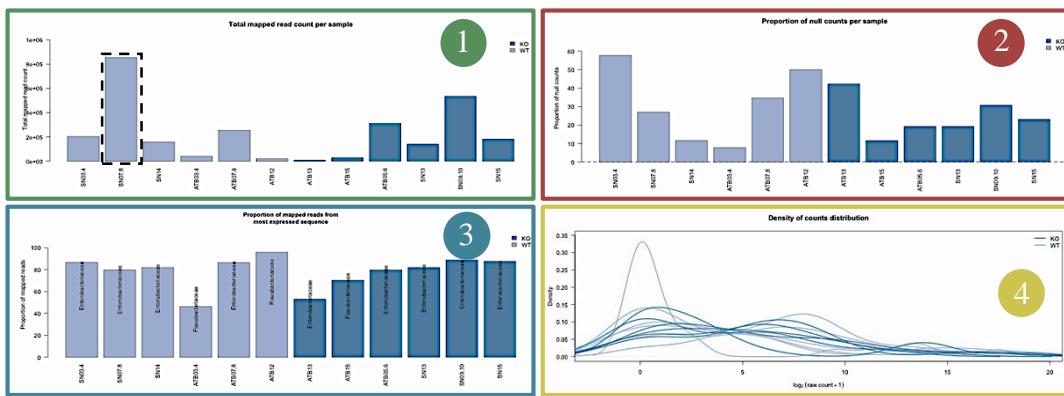
Assessing the statistical model

Once the model defined and analysis carried out, the user must check that the normalization, the estimation of the dispersion and the size factors are correct with the Total barplot, boxplots, the dispersion and the size factor vs total plots.

Different ordination methods (PCA, PCoA, NMDS) are also available to check that there has been no inversion of samples and that the biological effect studied is well observed in the data.

All of these representations are in the "Diagnostic Plots" sub-tab.





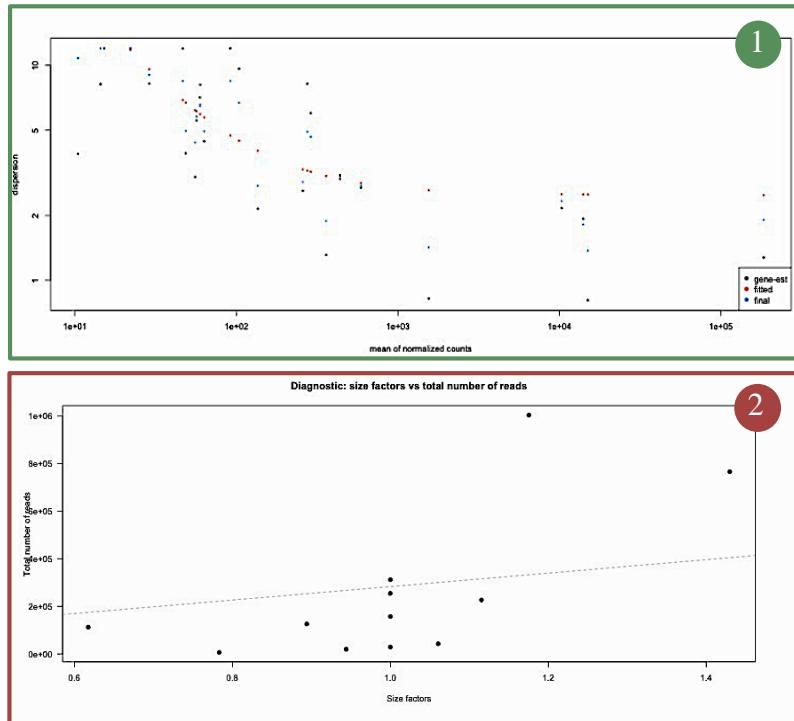
- 1 Diagramme en barres des comptages par échantillon
- 2 Diagramme en barres des comptages nuls par échantillon
- 3 Diagramme en barres des individus les plus abondants
- 4 Représentations de la densité des comptages en log2.

Interpretation:

These diagrams allow to identify possible issues linked to the sequencing depth. In the example presented, one of the "SN07-8" samples (in dashed lines) seems to have a large number of reads aligned with the others. This sample is also apparent when looking at the representation of densities.

In this context, the user must give attention to this sample and check if the whole bioinformatics process went well. It might be necessary to delete this sample from the analysis.

Note: The rarefaction curves presented in the "Visualization" section can also be used to determine if one sample must be deleted from the study.

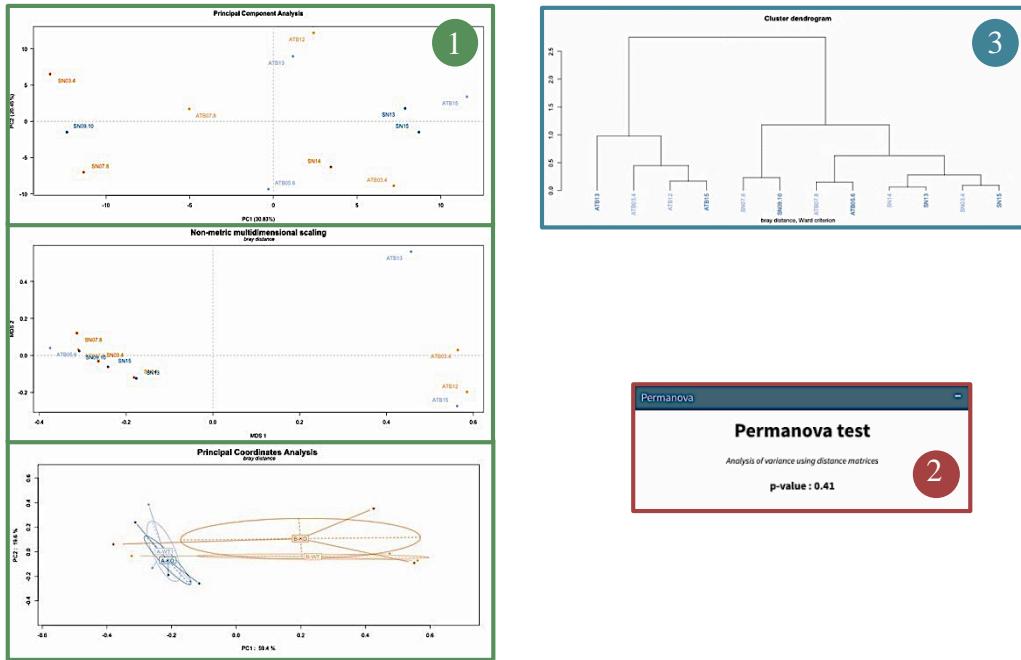


- 1 Représentation de la moyenne des comptages normalisés en fonction de la dispersion.
- 2 Représentation du nombre total de lecture alignés en fonction des size factors.

Interpretation :

These two representations allow to verify the normalization and the estimation of the dispersion. Concerning the estimation of the dispersion, the user must check that the model chosen corresponds to the shape of the point cloud. In the example presented here, the default function ("parametric") did not provide a satisfactory modeling, it was then decided to use the "local" option to obtain a modeling more consistent with the data.

The representation of the total number of reads aligned according to the size factors allows to check the normalization of the data and to identify possible outliers. In this graph, the size factors must be within a reasonable range (between 0.5 and 2.5 approximately) and the points must be close to the dashed line. Indeed, the size factors are supposed to correct the biases due to the variation of the sequencing depth, it is thus important to verify that there is indeed a link between size factors and sequencing depth.



- 1 Représentation des 2 premiers axes des 3 méthodes d'ordination (PCA, NMDS et PCoA) en fonction des variables d'intérêts sélectionnées par l'utilisateur.
- 2 Résultat du test de permanova. Ce test repose sur un calcul de distance entre les points. Non disponible pour la PCA.
- 3 Classification hiérarchique des échantillons.

Interpretation :

The ordination methods (PCA, NMDS and PCoA) have a double role, they allow to detect aberrant points or inversions of samples and to check if the studied biological effect brings a strong variability between the samples. In this example, we can see that the first axis of the PCoA (the one with the highest percentage of variance) allows to separate the 2 treatments A and B.

When samples are not sequenced in the same time, it is possible to see a "run" effect (samples would be clusterised according by sequencing batch). In this case, it is necessary to take this effect as an adjustment variable in the model (just add a "run" variable in the target file and incorporate it into the model). The experimental design must be considered before the analysis to avoid confusion problems between the variables (for example: to sequence all the samples of the treatment A together and, in a second time, those of the treatment B).

Regarding the p-value associated with the permanova test, it does not allow to determine that at least one of the groups is significantly different from the others (at a risk of 5%). When the test is significant, it means that at least one group is different from the others. When the biological effect studied is strong, the hierarchical classification makes it possible to identify possible inversion (s) of sample (s).

Differential analysis

Once the analysis is done, the results are available in tabular form in the "Table" tab.

ID	baseMean	FoldChange	log2FoldChange	pvalue_adjusted
Aeromonadaceae	10317.6	8.44996e+01	6.401	0.000131581648423139
Propionibacteriaceae	359.06	6.061851e-02	-4.044	0.00257263629427126
Enterobacteriaceae	184018.57	1.094042e+01	3.452	0.0098581723662584
Moraxellaceae	1547.76	8.210246e+00	3.037	0.0098581723662584
Burkholderiaceae	274.13	2.982792e-02	-5.067	0.0116530785397626
Nitochondria	438.25	7.214072e-02	-3.793	0.0201932

- 1 4 tables de résultats sont disponibles :
 - ✓ Significant: list of elements detected as differentially abundant.
 - ✓ Complete: list with all elements.
 - ✓ Up, down: lists of the elements detected as differentially abundant according to the sign of the fold change.
- 2 Result table:
 - ✓ baseMean: Average of normalized counts on all the samples.
 - ✓ FoldChange: Estimated fold change for selected contrast
 - ✓ Log2FoldChange : Fold change transformed in log2
 - ✓ Pvalue_adjusted: adjusted p-value (by default p-values are adjusted by BH procedure, other methods are available)
- 3 Selection of the desired comparison based on the user defined contrast vectors.
- 4 Export results. The user can choose the table to export as well as the column separator.

Interpretation:

The tables presented above allows to quickly identify the elements that have a significant change in term of abundance between the conditions studied. The presented p-values correspond to the p-values of the Wald test (after correction for multiple comparisons) at the threshold of significance chosen (e.g. 5%).

BaseMean informs about the mean abundance element (species, gender, ...).

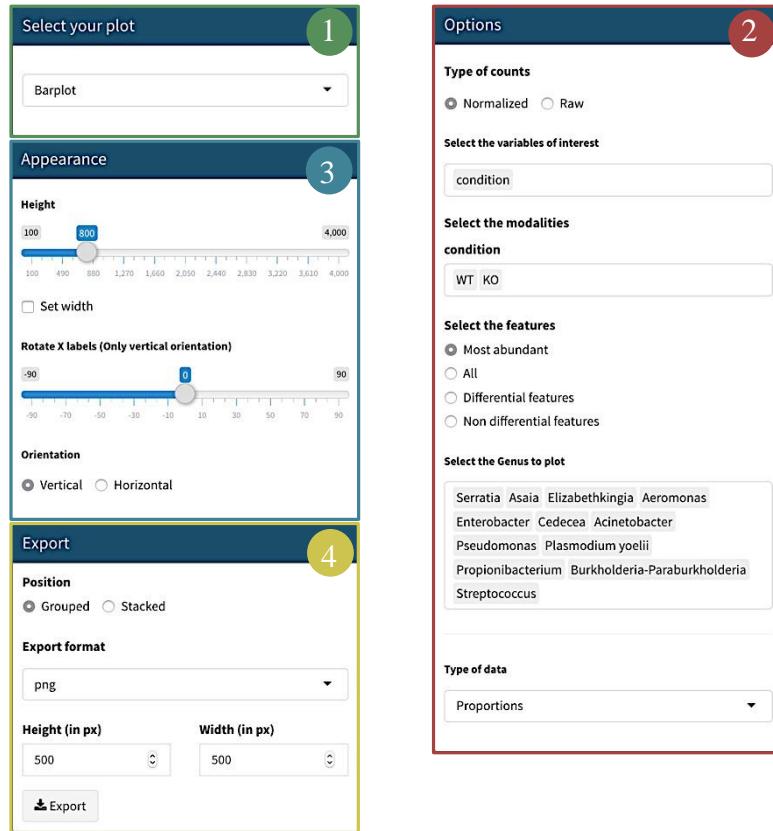
Note: the log2FoldChange are obtained from the parameter estimates and then the values are shrunk toward 0 for features with a high variance (strongly impacts the low counts).

VISUALIZATIONS

SHAMAN offers both a robust statistical analysis of data and a wide range of representations (barplots, heatmap, boxplots, abundance tree, scatterplot, diversity, rarefaction curves, krona and Venn digrams) allowing a complete analysis of the results. They are divided into two categories:

- ✓ Global views : représentations permettant de visualiser l'abondance des individus en fonction des variables définies dans le target.
- ✓ Comparison plots : représentations spécifiques permettant de comparer les résultats obtenus pour 2 contrastes ou plus.

Regardless of the type of visualization used, SHAMAN offers many options to the user (which will remain active from one representation to another). These options are in a column to the right of the visualizations and consists of 4 elements.



- 1 Menu déroulant permettant de sélectionner la visualisation souhaitée. 8 visualisations différentes sont proposées : barplot, heatmap, boxplot, Tree, scatterplot, diversity, rarefaction, Krona.
- 2 Liste des options disponibles pour la visualisation sélectionnée. Pour la majorité des visualisations disponibles, l'utilisateur a la possibilité de sélectionner les données qu'il souhaite représenter (normalisées or not), de définir les variables ainsi que de sélectionner la totalité ou seulement une partie des modalités et de sélectionner les éléments à représenter (les plus abondants, tous, seulement les différentiellement abondants ou non pour une des comparaisons définies).
- 3 Liste des options disponibles pour modifier l'apparence du graphique (dépend de la visualisation sélectionnée)
- 4 Export des résultats. L'utilisateur peut choisir le format (png, svg, eps, pdf) ainsi que la taille de la figure à exporter.

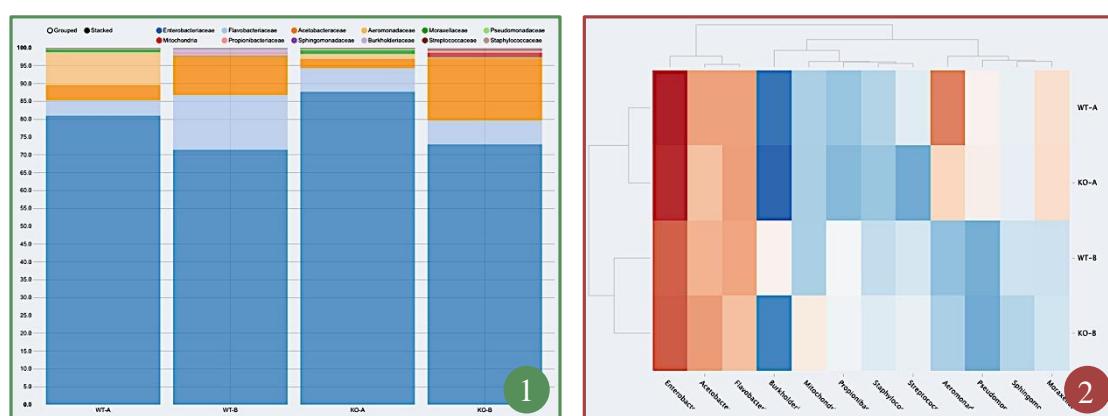
Visualisations des résultats

In the "Global views" tab, eight interactive visualizations are available to view the results.

Note: The interpretation of targeted metagenomics data must be done carefully. rRNA are present in several copy number in bacterial genomes [Vetrovsky and Baldrian 2013, Klappenbach 2001] from 1 to 15. In consequence, it is possible to analyze the abundance of in given genera between conditions but not to compare the abundance of two genera.

Overall composition

Barplot or heatmap are ideal visualizations for the analysis of the whole study. These representations make it possible to quickly visualize the differences between the conditions studied.



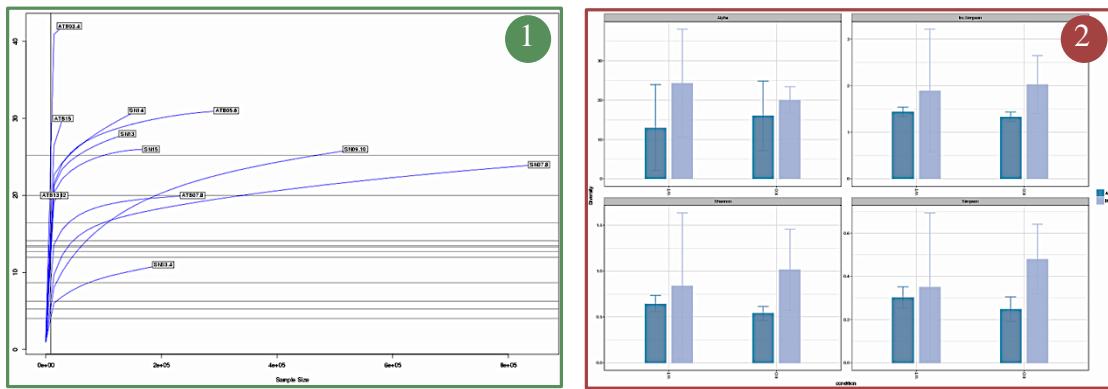
1 Barplot of the 12 most abundant elements. Each color represents the proportion of each element.

2 Heatmap for the 12 most abundant items per variable. The color depends on the average level of abundance. Dark red = very abundant; Dark blue = low abundance.

Interpretation:

The two previous figures show that the overall composition varies between the treatments (A and B) but not for the conditions (KO and WT) on the 12 most abundant. This observation is reinforced by the hierarchical clustering in the heatmap that groups the treatments between them.

The **rarefaction curves** and the **diversity measures** also provides an overview of the global composition according to the metadata.



- 1 Rarefaction curves for each sample. It corresponds to the number of elements detected according to the depth of sequencing.
- 2 Measure of diversity according to the variables of interest. SHAMAN offers 6 different measures: alpha, beta, gamma, shannon, simpson and inverse simpson.

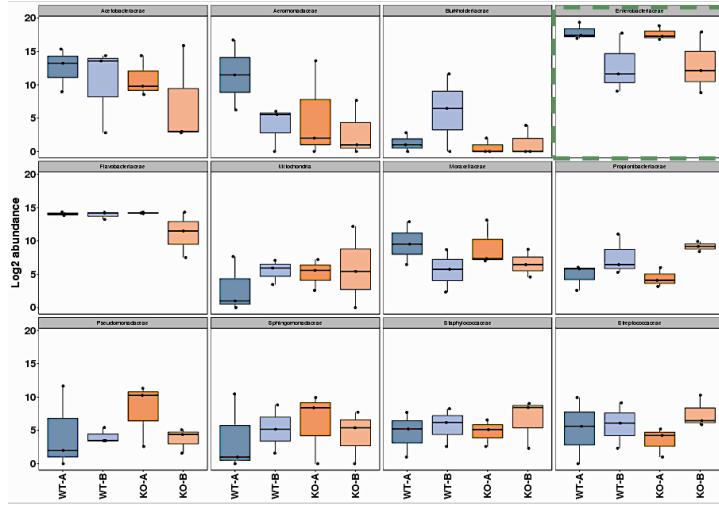
Interpretation:

The rarefaction curves highlight possible sequencing problems. The best case is to obtain rarefaction curves that all converge to a plateau, meaning that the sequencing depth was enough to discover, for instance, every genus that are present in the samples. In the example presented, the 3 samples closest to the ordinate axis show that the sequencing is not sufficient. In this case, it must be verified that normalization is not responsible for this result (by displaying the rarefaction curves for non-standardized data) and, if necessary, change the normalization method. If standardization is not responsible, it must be verified that there has been no problem upstream of the statistical analysis.

The diversities are useful to view the impact of the biological condition studied regarding to the composition of the samples. Confidence intervals make it possible to conclude whether the difference is significant or not (overlapping or not intervals). The user can also export the values used for the representations.

Fold-change

Boxplots are the usual representation to visualize the results of the differential analysis and log2 fold-change estimated.

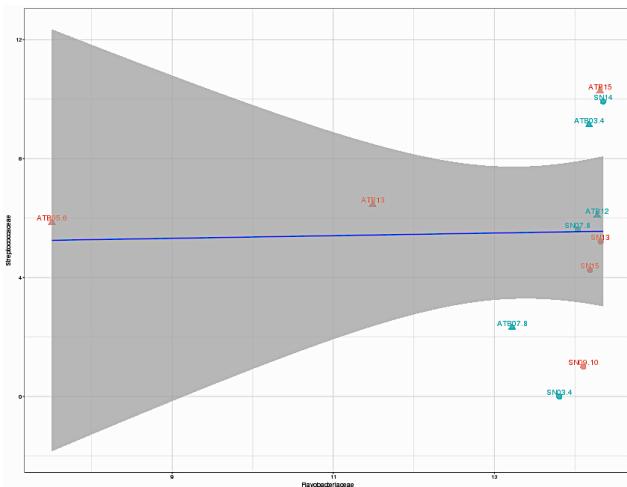


For each selected element, SHAMAN represents a boxplot for all the modalities of the variables of interest (with a different color). This representation highlights the differences identified during the differential analysis. For example, for the family Enterobacteriaceae (top right), it is clear that treatment B reduces the abundance of this family (whatever the condition). This effect was detected as significant in the result tables (3rd row of the table).

Links between variables

SHAMAN enables to cross abundance levels between two elements and to measure their correlation. It is also possible to model the abundance of a species with respect to diversity or to an auxiliary variable like age, bmi, ... which can be informative notably for the analysis of clinical data.

To realize these different crossings, SHAMAN proposes a **scatterplot** allowing to represent all the couples of elements, to change the points (form and color) according to the variables of interest and to add a third variable on which depends the size of the points.

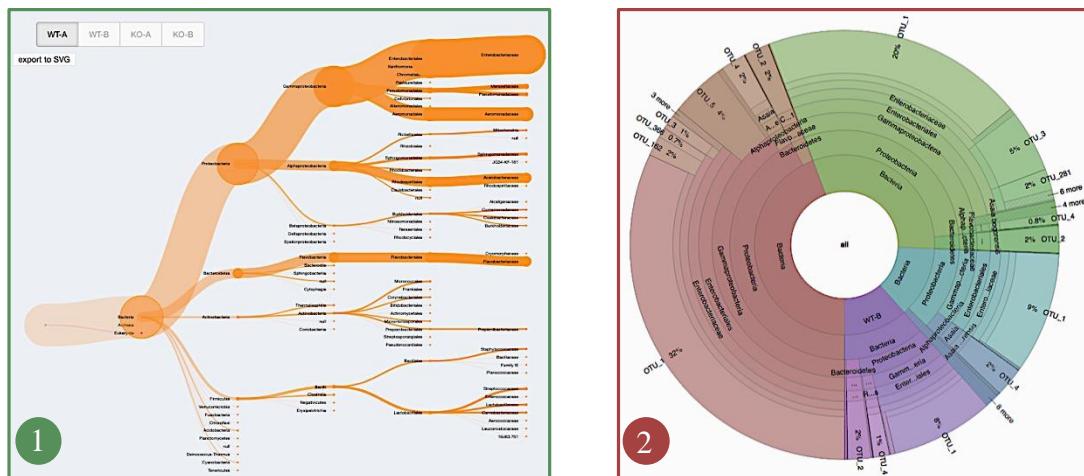


The user can also add a linear regression (blue line) with confidence area (gray area). The equation of the line, the coefficient of determination R^2 and the tests for slope and intercept are then available. The user can also access to the correlation table (Pearson or Spearman).

Abundance and taxonomy

All previous representations allow to visualize the results of the analysis at the level defined by the user during the definition of the statistical model.

Result comparisons

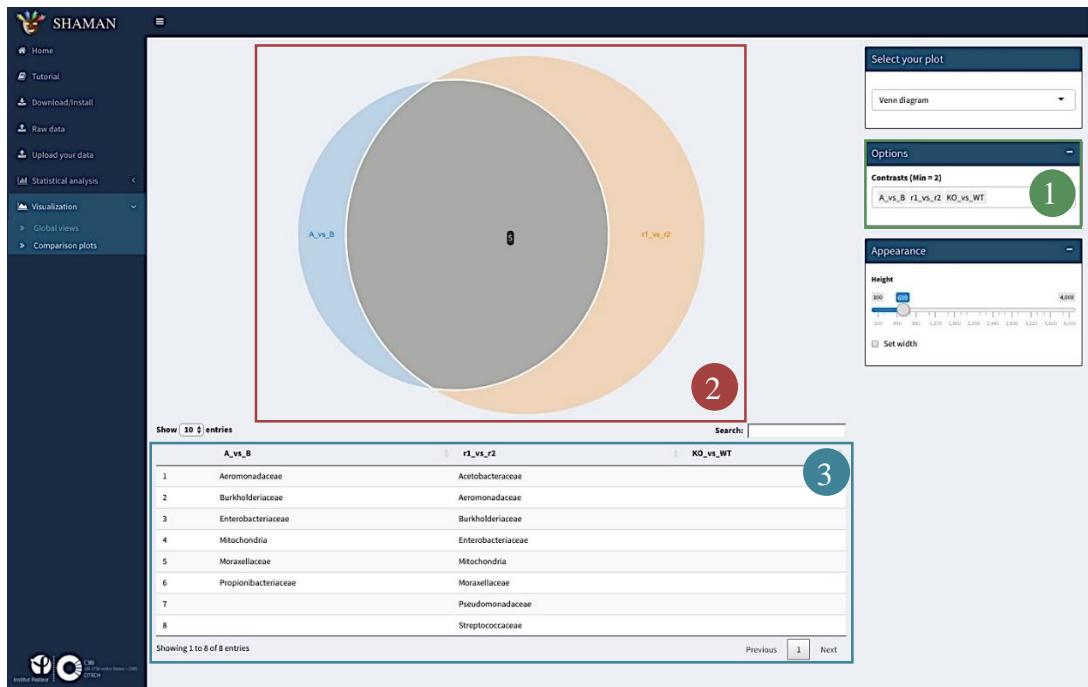


- 1 Abundance tree according to the taxonomic level. The user can choose the samples they wish to represent.
- 2 Representation in the form of a krona plot¹. This interactive visualization allows to navigate through different taxonomic levels.

The main objective of SHAMAN is to perform a quantitative study of the differences between experimental conditions. A more qualitative approach is proposed in SHAMAN to compare the results of different studies. This is available through 2 representations, the Venn

¹ <https://github.com/marbl/Krona/wiki>

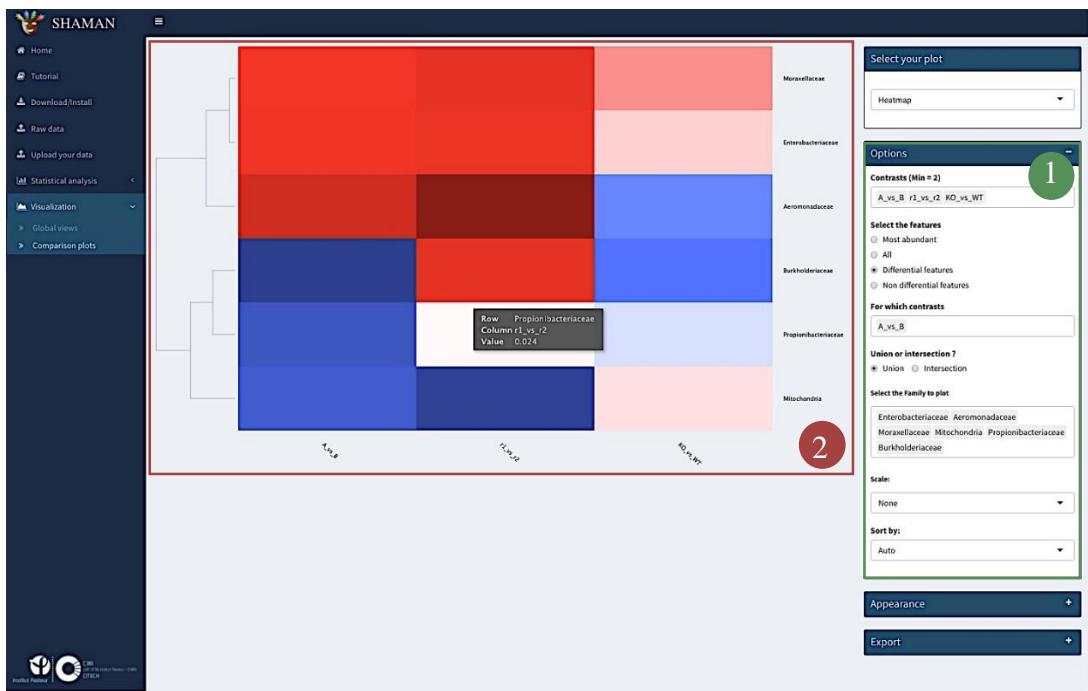
diagram that compare the elements detected as differentials for 2 or more comparisons and the heatmap of the log2 fold change that illustrate the strength of the difference. This type of representation allows to compare the results of several studies from a qualitative point of view. In many cases, a quantitative analysis can also be performed (with 2 comparisons) by defining a suitable contrast ("advanced user" mode).



- 1 Selection of contrasts that will be compared. At least two contrasts must be selected.
- 2 Venn diagram between differentially abundant elements for selected contrasts.
- 3 Lists elements for each contrast.

Interpretation:

The Venn diagram is interpreted by identifying the number of common elements between contrasts. Each circle is associated to the comparison of one contrast.



- 1 List of the different options proposed to draw the heatmap.
- 2 Representation of log2 fold changes as a heatmap. The color depends on the value of the log2 fold change. Dark red: high positive value; Dark blue: high negative value; White: null value.

Interpretation:

This representation allows to identify common elements between 2 or more contrasts. The shade of color gives the direction and the strength of the difference.

BIBLIOGRAPHIE

Anders and Huber, Differential expression analysis for sequence count data, *Genome Biology*, 2010

Jonsson V, Österlund T, Nerman O et al. Statistical evaluation of methods for identification of differentially abundant genes in comparative metagenomics. *BMC Genomics* 2016; 17: 78

Love, Huber and Anders, Moderated estimation of fold change and dispersion for RNA-Seq data with DESeq2, *Genome Biology*, 2014

McMurdie, P. J. & Holmes, S. Waste not, want not: why rarefying microbiome data is inadmissible. *PLOS Comput. Biol.* 10, e1003531 (2014)

How to create a target file for Shaman

format, contents and constraints

The target file is required for each analysis done with Shaman. It must contain all the available information on the samples (corresponding to the metadata) that will be used to build the statistical model and/or to visualize the data. To be loaded in Shaman the target file must respect some properties

- 1. The first column is dedicated to the sample name which must correspond exactly to the column names of the count matrix. At least 2 samples are required.** *Once the count matrix is loaded, check carefully the sample names in the “count table” tab. Sometimes some characters are modified with the loading.*
 - 2. At least one variable must be provided.** *In Example 1, two variables are provided (condition and treatment).*
 - 3. NA or missing values are not allowed.**
 - 4. A variable with the same value for each sample is not allowed.** *This kind of variable should be removed from the target file before loading.*
 - 5. The selected variables for the statistical model must not be collinear.** *It means that if one variable can be determined by another variable or a combination of variables the analysis cannot be done with all the variables. However, the user can use this variable for visualization. (See example 3).*
 - 6. Be careful, numeric variables will be considered as quantitative variable.** *For instance, do not use 1 and 2 to describe two different conditions but C1 and C2 or A and B. (see example 3)*
 - 7. Avoid using special characters such as / \ ? * : < > / + , [] - + () % @ " &**
-

Example 1: Target file with 2 variables (condition and treatment)

sampleID	condition	treatment
S1	WT	A
S2	WT	A
S3	KO	A
S4	KO	A
S5	WT	B
S6	WT	B
S7	KO	B
S8	KO	B

Error

The model matrix is not full rank. One or more variables or interaction terms are linear combinations of the others and must be removed.

Reminder: Your target file must contain at least 2 columns and 2 rows. NA's values are not allowed and the variables must not be collinear.

This is a usual example in which we have 2 variables to describe the samples (condition and treatment). For instance, the user will be able to define the following model: condition + treatment + condition:treatment and then get differentially abundant features between treatments A and B for each condition.

Example 2: Target file with collinearity problem

sampleID	condition	treatme	group
D	nt		
S1	WT	A	g1
S2	WT	A	g1
S3	KO	A	g2
S4	KO	A	g2
S5	WT	B	g3
S6	WT	B	g3
S7	KO	B	g4
S8	KO	B	g4

In this example, group = condition + treatment, so the variables are collinear.

Note that this file can be loaded in Shaman without error but the error will appear if the user tries to define a model with the three variables condition, treatment and group.

Example 3: Quantitative versus qualitative variable.

Target file		Model parameters	
sampleID		condition	
S1	1	condition	
S2	1	<input type="text" value="0"/>	
S3	2	<input type="text" value="0"/>	
S4	2	<input type="text" value="0"/>	
S5	3	<input type="text" value="0"/>	
S6	3	<input type="text" value="0"/>	
S7	4	<input type="text" value="0"/>	
S8	4	<input type="text" value="0"/>	
sampleID		condition	
S1	C1	conditionC1	
S2	C1	<input type="text" value="0"/>	
S3	C2	<input type="text" value="0"/>	
S4	C2	<input type="text" value="0"/>	
S5	C3	<input type="text" value="0"/>	
S6	C3	<input type="text" value="0"/>	
S7	C4	<input type="text" value="0"/>	
S8	C4	<input type="text" value="0"/>	

In case 1, condition is considered as a numeric variable which leads to only one parameter in the statistical model. It assumes that the difference between 1 and 3 is two times the difference between 1 and 2 and so on. In case 2, there is no order between the conditions.