

# Contamination detection and microbiome exploration with GRIMER

Vitor C. Piro <sup>1,2,\*</sup> and Bernhard Y. Renard <sup>1</sup>

<sup>1</sup>Data Analytics and Computational Statistics, Hasso Plattner Institute, Digital Engineering Faculty, University of Potsdam, Potsdam 14482, Germany

<sup>2</sup>Department of Mathematics and Computer Science, Freie Universität Berlin, Berlin 14195, Germany

\*Correspondence address. Vitor C. Piro. Department of Mathematics and Computer Science, Freie Universität Berlin, Berlin 14195, Germany; E-mail: [vitor.piro@fu-berlin.de](mailto:vitor.piro@fu-berlin.de)

## Abstract

**Background:** Contamination detection is an important step that should be carefully considered in early stages when designing and performing microbiome studies to avoid biased outcomes. Detecting and removing true contaminants is challenging, especially in low-biomass samples or in studies lacking proper controls. Interactive visualizations and analysis platforms are crucial to better guide this step, to help to identify and detect noisy patterns that could potentially be contamination. Additionally, external evidence, like aggregation of several contamination detection methods and the use of common contaminants reported in the literature, could help to discover and mitigate contamination.

**Results:** We propose GRIMER, a tool that performs automated analyses and generates a portable and interactive dashboard integrating annotation, taxonomy, and metadata. It unifies several sources of evidence to help detect contamination. GRIMER is independent of quantification methods and directly analyzes contingency tables to create an interactive and offline report. Reports can be created in seconds and are accessible for nonspecialists, providing an intuitive set of charts to explore data distribution among observations and samples and its connections with external sources. Further, we compiled and used an extensive list of possible external contaminant taxa and common contaminants with 210 genera and 627 species reported in 22 published articles.

**Conclusion:** GRIMER enables visual data exploration and analysis, supporting contamination detection in microbiome studies. The tool and data presented are open source and available at <https://gitlab.com/dacs-hpi/grimer>.

**Keywords:** Contamination, Microbiome, Visualization, Taxonomy

## Introduction

Microbiome studies enable, via high-throughput sequencing, the investigation of the composition of complex microbial communities from diverse environments. Microbiome studies usually yield large amounts of raw sequences for several samples that can be analyzed with an increasing number of computational methods and databases. Standards, protocols, and best practices for designing and performing a microbiome study have been improving and changing over the years [1, 2], and the field is in constant evolution due to higher availability and reduced costs of sequencing runs as well as with the increase in number of publicly available reference sequences and computational methods.

In early stages of a standard *in silico* microbiome analysis, raw or quality-filtered sequences are classified or clustered into specific groups and quantified to generate a profile for a given environmental sample. Marker gene, whole-metagenome, and meta-transcriptome analyses have their own set of tools and standards that should be carefully chosen to generate reliable measurements for each sample in the study [3]. This step can be computationally intensive but reduces the large amount of data into a concise table of measurements. Alternatively, genome assembly can be performed for metagenomics samples, allowing genome-resolved analysis. Although still a complex task, gene prediction, taxonomic, and functional analyses are improved with metagenome-assembled genomes, resulting in overall better measurements [4].

After measurements are obtained, hypotheses are validated through data mining and statistical analysis. This step is mostly exploratory and specific to the hypotheses and research questions pursued, and the required analyses are difficult to be fully automatized. It is also very important to take in consideration the compositionality of data at this stage when working with the microbiome [5]. Several comprehensive and generalized analytical packages [6–8] and web platforms (Table 1) are available to perform a large number of microbiome analysis: basic data summaries, diversity and functional analysis, microbial interactions, and differential abundance, among others. Additionally, interactive tools for analytical and visual exploration are extremely helpful in this stage to better understand the data distribution and its properties and to guide further investigations to follow. In the last decade, several applications were developed with focus on visualization of microbiome data (Table 2). A comparison among many of those methods and their functionalities can be found in a recent review [9].

At this stage of a study, contamination detection should be considered. Contamination side effects have gained attention in recent years due to the controversial detection of a placental microbiome [31–33]. However, the issue is not new, and contamination has been known and reported for decades in the literature [34]. Contamination is characterized by exogenous DNA in a given sample introduced externally or internally. External contamination can come from diverse sources: DNA extraction kits, labo-

Received: October 31, 2022. Revised: February 6, 2023. Accepted: March 1, 2023

© The Author(s) 2023. Published by Oxford University Press GigaScience. This is an Open Access article distributed under the terms of the Creative Commons Attribution License (<http://creativecommons.org/licenses/by/4.0/>), which permits unrestricted reuse, distribution, and reproduction in any medium, provided the original work is properly cited.

**Table 1:** Web resources to process, analyze, and visualize microbiome data

Name	Website	Reference
MG-RAST	<a href="https://www.mg-rast.org/">https://www.mg-rast.org/</a>	[10]
MGNify	<a href="https://www.ebi.ac.uk/metagenomics/">https://www.ebi.ac.uk/metagenomics/</a>	[11]
MicrobiomeDB	<a href="https://microbiomedb.org/">https://microbiomedb.org/</a>	[12]
Nephele	<a href="https://nephele.niaid.nih.gov/">https://nephele.niaid.nih.gov/</a>	[13]
Qiita	<a href="https://qiita.ucsd.edu/">https://qiita.ucsd.edu/</a>	[14]

**Table 2:** Interactive analysis and visualization tools for microbiome data published in the last 10 years

Name	Focus	Platform	Website	Year	Reference
METAGENassist	Comparative metagenomics	Web	<a href="http://www.metagenassist.ca/METAGENassist/">http://www.metagenassist.ca/METAGENassist/</a>	2012	[15]
VAMPS	Microbial population structures	Web	<a href="https://vamps2.mbl.edu/">https://vamps2.mbl.edu/</a>	2014	[16]
Shiny-phyloseq	Microbiome analysis	Locally hosted (R)	<a href="https://joey711.github.io/shiny-phyloseq/">https://joey711.github.io/shiny-phyloseq/</a>	2015	[17]
MetaCoMET	Microbiome analysis	Web	<a href="https://probes.pw.usda.gov/MetaCoMET/">https://probes.pw.usda.gov/MetaCoMET/</a>	2016	[18]
BusyBee Web	Metagenomics binning and analysis	Web	<a href="https://ccb-microbe.cs.uni-saarland.de/busybee">https://ccb-microbe.cs.uni-saarland.de/busybee</a>	2017	[19]
MicrobiomeAnalyst	Microbiome analysis	Web	<a href="https://www.microbiomeanalyst.ca">https://www.microbiomeanalyst.ca</a>	2017	[20]
Burrito	Taxonomy and function analysis	Web	<a href="http://elbo-spice.cs.tau.ac.il/shiny/burrito/">http://elbo-spice.cs.tau.ac.il/shiny/burrito/</a>	2018	[21]
Pavian	Metagenomics analysis	Locally hosted (R)	<a href="https://github.com/fbreitwieser/pavian">https://github.com/fbreitwieser/pavian</a>	2019	[22]
GenePiper	Microbiome analysis	Locally hosted (R)	<a href="https://github.com/raytonghk/genepiper">https://github.com/raytonghk/genepiper</a>	2020	[23]
animalcules	Microbiome analysis	Locally hosted (R)	<a href="https://github.com/compbioed/animalcules">https://github.com/compbioed/animalcules</a>	2021	[24]
MicrobiomeExplorer	Microbiome analysis	Locally hosted (R)	<a href="https://github.com/zoecastillo/microbiomeExplorer">https://github.com/zoecastillo/microbiomeExplorer</a>	2021	[25]
microViz	Microbiome analysis	Locally hosted (R)	<a href="https://github.com/david-barnett/microViz/">https://github.com/david-barnett/microViz/</a>	2021	[26]
Namco	Microbiome analysis	Web	<a href="https://exbio.wzw.tum.de/namco/">https://exbio.wzw.tum.de/namco/</a>	2021	[27]
OpenContami	Contaminant detection	Web	<a href="https://openlooper.hgc.jp/opencontami/">https://openlooper.hgc.jp/opencontami/</a>	2021	[28]
wiSDOM	Microbiome analysis	Web or Locally hosted (R)	<a href="https://github.com/lunching/wiSDOM">https://github.com/lunching/wiSDOM</a>	2021	[29]
Mian	Microbiome analysis	Web	<a href="https://miandata.org/">https://miandata.org/</a>	2022	[30]
GRIMER	Contaminant detection	CLI + standalone file	<a href="https://github.com/pirov/grimer">https://github.com/pirov/grimer</a>	2022	this work

ratory reagents, surfaces and equipment, ultra-pure water, residuals from previous sequencing runs, and microbes from laboratory technicians [2, 35, 36]. Internal contamination can be defined as a undesired exchange of genetic material between samples, and it is usually referred as well-to-well contamination, cross-contamination, or sample "bleeding" as well as index switching in multiplexed sequencing libraries [37].

Contamination may affect most sequencing projects to some degree, especially low-biomass samples [38]. The composition of an environmental sample is mostly unknown before sequencing, increasing the complexity of detecting contamination when compared to a defined isolate genome and targeted sequencing project. Low-biomass samples (e.g., meconium, blood, human tissues) yield little to no DNA to be amplified and sequenced, an ideal scenario for exogenous contaminants to outcompete and dominate the biological signal.

It is important that contamination is acknowledged, accounted for, and discovered at the earliest stage of a study prior to statistical analysis, to not bias measurements and to

ensure that bias is not propagated into databases [39, 40]. Inclusion of negative and positive control samples is the recommended way to measure, detect, and mitigate contamination [2, 38, 41]. Negative controls should be included in the study design for every sample, extraction, or amplification batch. Once provided, controls should be carefully analyzed *in silico*, and results obtained should be applied to biological samples in terms of prevalence (e.g., observations in negative controls) but also based on the frequency in relation to DNA concentration [42, 43].

However, due to the complexity and diverse possible sources of contamination, detection and mitigation are not a trivial tasks. Several approaches to identify and exclude background contamination in microbial studies have been proposed. These are based on exclusion of organisms detected in negative controls, use of replicates to find possible contaminants, removal of low abundant signals, negative correlation between organism abundance and bacterial load, clustering analysis, and others [44, 45]. Each approach has strengths and weaknesses based on the study design, data type, and control availability. Further, many studies do

not include or have a limited number of control samples due to the required increase in costs. Hornung et al. [41] reported that based on publications from the 2018 issues of *Microbiome* and *The ISME Journal*, only 30% cited the use of negative controls and only 10% positive controls. Moreover, Harrison et al. [46] reported that out of 50 selected publications from 2019 and 2020, only 15 used some type of negative control and 10 of positive control to account for reagent contamination. There was also no observed increase in positive or negative controls usage in the literature from 2015 to 2020, based on selected publications. Additionally, the detection of recurring contaminants in extraction kits and reagents (also called “kitome”) is known to be an issue [47] but remains under-explored, mainly for not being properly cataloged, centralized, or automated.

To overcome some of those challenges, we propose GRIMER, a tool to analyze, visualize, and explore microbiome studies with a focus on contamination detection. Based on a table of observations per sample, GRIMER generates an offline and interactive dashboard to automate data analysis, transformations, and plots and generates a set of charts integrating evidence for better decision-making and contamination detection. Additionally, we compiled an extensive list of common contaminants containing 210 genera and 627 species reported in 22 published articles. These data are integrated into the report. GRIMER is an effortless step once quantification is done, turning measurement tables into a interactive and dynamic report in seconds. GRIMER is open source, and the code is available at the GitHub repository [48]. Installation and usage instructions as well as an user manual are available in the repository. The tool is independent of analysis methods, does not rely on web or local servers, and generates standalone and shareable interactive dashboards.

## Methods

GRIMER analyzes and annotates multisample studies based on count tables and generates a report with several interactive plots to better explore the data and to facilitate contamination detection. GRIMER integrates several sources, references, analyses, and external tools and brings them together in one concise dashboard.

The output of GRIMER is a self-contained HTML file that can be visualized in any modern web browser. It works independently from any actively running server or web service. Once generated, it can be used and shared as an offline document. It has the advantages of a static report and a complex dashboard being portable and interactive. This feature makes it very convenient to distribute (e.g., as an email attachment), keep track of changes in analytical pipelines, and reproduce analyses in different environments.

GRIMER is independent of any quantification method and only requires a contingency table with raw counts of observations/components for each sample/composition in the study. Observations are usually, but not limited to, taxonomic entries (e.g., genus, species, strains), operational taxonomic units (OTUs), amplicon sequence variants, or sequence features. A count of unclassified or unassigned observations is also supported to generate normalized values. Additional files and data can be provided to expand GRIMER reports: study metadata, a taxonomy database, multiple control samples, the DNA concentration, custom contaminants, and reference groups of interest. The more information provided, the more complete and interactive the final report will be.

## Annotation

GRIMER annotates observations and samples linking data with external data sources.

Sample annotations are based on a user-provided study metadata, where each sample is described in 1 or more fields and variables. Those fields can contain either numeric or categorical values and are useful for grouping and clustering analyses as well as detection of batches and control/treatment effects.

Observation annotations are based on external lists of taxonomic entries, which can be used, for example, to link findings to common contaminants or connect analyses outcomes with known environments or biomes. Those entries can be easily provided by the user in a simple list of names or taxonomic identifiers in a formatted and annotated file (more information can be found in the GRIMER repository).

## Contamination references

We compiled an extensive list of possible contaminant taxa reported in several studies (Table 3). The studies selected were obtained from cross-references in review articles [38] and individual selected findings in the literature, usually focusing on contamination detection or mitigation. Articles were manually curated and more studies can potentially be added to the list, which is dynamically maintained. Contributions are welcome through the GRIMER repository [48]. The studies selected are very diverse in terms of sequencing technology, methodology used, and environment studied. Contamination in those studies can originate from diverse sequencing kits and reagents as well as the lab environment or other unknown sources. The idea behind compiling this list is to detect which taxa are the most recurrently identified as contaminant in diverse conditions, providing a guideline and consensus for further studies. Entries on this list are not strictly considered a contaminant and should not be used alone to define contamination in a study. However, it serves as an additional evidence supporting it, especially if entries are highly recurrent (Table 4) and corroborate with additional lines of evidence. Those contaminants were reported mainly at genus or species level in different formats, names, and taxonomies. We manually curated and converted them into the NCBI taxonomy [49] nomenclature for standardized usage.

Additionally, we compiled another list of common organisms found in probable external contamination sources: taxa commonly occurring in human skin, oral and nasal cavities, and face and other human limbs. Those were reported as possible sources of contamination [38]. Reference organisms names were obtained from BacDive [69], eHOMD [70], and further publications [71].

## MGNify

Additionally to the contamination references, a summary generated from the MGNify repository [11] is provided with counts of occurrences for each observation in thousands of microbiome studies, grouped by biome. MGNify is a resource to analyze microbiome data in an automated and standardized way. Thousands of analyzed studies are publicly available with related metadata. We mined this repository with the provided open API [72] and collected all taxonomic classifications available for every study. For each study, we collected the latest taxonomic classification based on the highest pipeline version available. If multiple classifications from different sources were present, we selected the largest one by file size. For each study output, the top 10 top most abundant organisms were linked to the study respective biome(s) definition, and a final count of top organisms by biome is generated.

**Table 3:** Summary of common contaminants taxa extracted from the literature. The complete list of taxa per study can be found in the GRIMER repository [48].

Organism group	Genus	Species	Reference
Bacteria	6	0	1998 Tanner et al. [50]
Bacteria	0	10	2002 Kulakov et al. [51]
Bacteria	4	0	2003 Grahn et al. [52]
Bacteria	16	0	2006 Barton et al. [53]
Bacteria	11	1	2014 Laurence et al. [54]
Bacteria	92	0	2014 Salter et al. [35]
Bacteria	7	0	2015 Jervis-Bardy et al. [42]
Bacteria	28	0	2015 Jousselein et al. [55]
Bacteria	77	127	2016 Glassing et al. [36]
Bacteria	23	0	2016 Lauder et al. [56]
Bacteria	6	0	2016 Lazarevic et al. [57]
Bacteria	62	0	2017 Salter et al. [58]
Bacteria	0	122	2018 Kirstahler et al. [59]
Bacteria	34	0	2018 Stinson et al. [60]
Bacteria	18	0	2019 Stinson et al. [61]
Bacteria	52	2	2019 Weyrich et al. [62]
Bacteria	8	26	2019 de Goffau et al. [63]
Bacteria	15	93	2020 Nejman et al. [64]
Viruses	0	1	2015 Kjartansdóttir et al. [65]
Viruses	0	1	2015 Mukherjee et al. [66]
Viruses	0	291	2019 Asplund et al. [67]
Eukaryota	0	3	2016 Czurda et al. [68]
Eukaryota	0	1	PRJNA168
Total (unique)	210	627	—

**Table 4:** Top 8 most reported taxa from Table 3 at genus and species level. If multiple child nodes of organisms are reported in the same study, they are counted here just once.

Genus	# reported	Species	# reported
<i>Pseudomonas</i>	13	<i>Cutibacterium acnes</i>	4
<i>Stenotrophomonas</i>	13	<i>Pseudomonas fluorescens</i>	4
<i>Ralstonia</i>	12	<i>Stenotrophomonas maltophilia</i>	4
<i>Bradyrhizobium</i>	11	<i>Acinetobacter baumannii</i>	3
<i>Methylobacterium</i>	11	<i>Bradyrhizobium elkanii</i>	3
<i>Acinetobacter</i>	10	<i>Corynebacterium tuberculostearicum</i>	3
<i>Corynebacterium</i>	10	<i>Rhodococcus fascians</i>	3
<i>Sphingomonas</i>	10	<i>Streptococcus mitis</i>	3

GRIMER uses this resource to annotate observations and links how many times each identified taxon was present in other biomes. This gives another level of evidence for the possible origin of certain taxa in a study, compared to thousands of other microbiome studies. For example, in the current version, the genus *Ralstonia*, a commonly reported contaminant, appeared in 30 environmental aquatic biome studies and 14 engineered bioreactor studies (out of a total of 79 studies) while the human-related bacterial genus *Prevotella* appears mostly in host-associated biomes (89% of occurrences). All 5 levels of biome classification are available for each taxonomic entry.

### Input data

GRIMER requires only a contingency table to generate the full report, either in a text/tabular format (observations and samples either in rows or columns with a header) or a BIOM file [73]. Further data can be provided to extend the report:

- Metadata: annotate samples and give further technical information. The metadata should be tabular and categorical, and numerical fields are supported.

- Taxonomy: GRIMER will automatically parse a given taxonomic annotation or generate one based on the provided observations. Data will be summarized in many taxonomic levels, and plots will be created accordingly. Taxonomy is fully automated for several commonly used taxonomies (NCBI, GTDB, SILVA, GreenGenes, OTT).
- Controls: 1 or more groups of control samples can be provided in a simple text file. Those samples will be further used to summarize data and annotate plots.
- References: custom sources of contamination or any references can be provided in addition to the precompiled ones described above.

GRIMER will parse and process the data provided and run a set of analyses:

- General data summary by observation and samples, linking references, taxonomy, and metadata
- Filtering and transformation: observations and samples can be filtered to reduce noise or small counts. Transformations are applied (log, centered log-ratio, normalization) to account

for the compositionality of the data and improve some visualizations.

- Hierarchical clustering: 1 or more metrics and methods can be used to perform the clustering. The combination of all of them is executed and available in the report. For this analysis, zeros are replaced by small counts defined by the user.
- Correlation: symmetric proportionality coefficient (rho correlation) [74, 75] is calculated for top abundant observations in the study.
- DECONTAM [43]: R package with a simple method to detect contaminating taxa/observations based on 2 main assumptions: frequency of contaminant taxa inversely correlate with DNA concentrations, and contaminant taxa are more prevalent in control samples than in biological samples. DECONTAM uses linear models based on the assumptions and frequencies of the data and outputs a score for each observation to define contamination. If DNA concentration is not provided, total counts are used instead as an indirect concentration value replacement.
- MGnify: each taxon reported will be linked to the respective MGnify entry, reporting most common biome occurrences.

## GRIMER report

GRIMER will generate a report/dashboard with visualizations to better understand the distribution of observation counts among samples and the connection with external annotations, metadata, and taxonomy. Currently, GRIMER reports contain 4 main panels: Overview (Fig. 6), Samples (Fig. 1), Heatmap (Fig. 4), and Correlation (Fig. 5). Some of them were previously suggested to be adequate for contamination detection [45] and are commonly used in standard microbiome analysis. Every panel has 1 or more visualization and widgets to select, filter, group, and modify its contents. Panels can be reported independently.

An individual summary for each observation and its relation to annotations and distribution among samples can be found in the Overview panel (Fig. 6). Here, all evidence related to a specific observation is integrated for further examination. Each observation provided in the study is listed and summarized in a tabular format. Once selected, the distribution of counts for the specific observation for each sample can be observed in a bar plot. Information of annotations, MGnify biomes, and DECONTAM output are also available in the same interface. The DECONTAM output indicates if the observation is classified as a contaminant with a score and a plot showing the frequency of the selected observation against the DNA concentration for all samples containing that observation. Linear models showing the expected values for contamination and noncontamination values are also plotted. If provided, taxonomic lineages are integrated in the table, and plots and observation are decomposed and summarized into taxonomic levels. The Overview panel also roughly summarizes samples contents in the bar plot, with general classification metrics. Those bars can be transformed, annotated, grouped, and sorted to connect observation values to overall sample distribution.

In-depth evaluation of individual samples can be performed in the Samples panel (Fig. 1). Normalized distribution of top observations for each sample can be visualized in the bar plot to easily compare the overall distribution of observations among samples, with options for grouping and sorting by metadata. Automated selection of groups of samples is also possible by counts and metadata.

Several transformations can be applied to the data (normalization, log, center log ratio) to be further visualized in the Heatmap

panel (Fig. 4). Hierarchical clustering and grouping and sorting options can be independently selected for samples and observations to enable pattern detection (e.g., batch effects, treatment effects). Dendrograms are plotted when clustering options are selected. Annotation bars are plotted around the heatmap showing dynamic selection of sample annotations (metadata) and observation annotations (references, controls, and DECONTAM output). Metadata are automatically colored to reflect categories (distinct colors) and numeric (sequential colors) fields. Multiple metadata fields can be select interactively. Observation annotation values are normalized and plotted in the same color scale for easier interpretation. One heatmap is generated for each taxonomic level.

Correlation between observations is plotted as a matrix (Fig. 5). Positive or negative correlations among observations can point toward concurrent signals in the microbiome analysis. Observations present in multiple samples in similar ratios are positively correlated, and the opposite configures negative correlation. Once a signal is observed, the correlation matrix can indicate co-occurrence of observations and help to identify further candidates (e.g., cluster of co-occurring contaminants at similar ratios).

## Implementation

GRIMER is written in Python and Javascript and outputs a report file in HTML format. All visualizations and layouts are created with the the Bokeh library [76]. Bokeh plots, tables, and charts automatically provide a set of tools for interaction (e.g., zoom, selection) with an option to export the current selection to an image file. Many plots have interactive tool-tips, showing more information about the data under the mouse cursor. Help buttons are also included, explaining the plots and options.

Further libraries were used to analyze samples and generate the report: pandas [77] for general parsing and data structures, scipy [78] for hierarchical clustering, and scikit-bio [79] for transformations. Scripts to download and generate MGnify annotations and update reference sources are provided in the GRIMER repository [48].

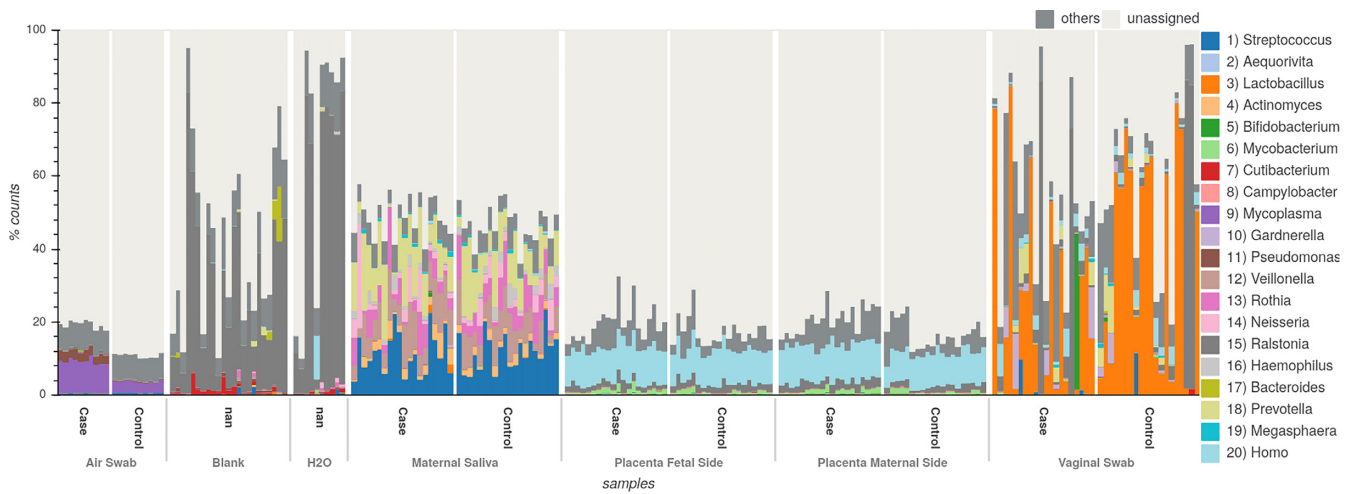
GRIMER automatically handles taxonomic entries using MultiTax [80]. GRIMER will automatically parse given taxonomies or download and convert any taxonomic ID or name internally and decompose results in taxonomic ranks. Currently supported taxonomies are NCBI, GTDB, Silva, GreenGenes, and OpenTree Taxonomy. Reference lists are currently only available based on the NCBI Taxonomy.

## Results

We reanalyzed publicly available studies to demonstrate the use of GRIMER reports in real case scenarios and what types of analyses are possible. In some examples, we try to reproduce analyses and in other cases point to new evidence that may have been overlooked. We encourage the readers to download [81] or open live examples of GRIMER reports [82] and interactively visualize the results being described to fully understand the capabilities of the report. All reports presented below were generated using GRIMER version 1.1.0.

## Detecting contamination

The attempt to detect and describe a possible human placental microbiome has motivated several studies and investigations [47, 63, 83, 84]. Leiby et al. [85] published a detailed and well-designed study contributing to the subject. Placental samples for term (con-



**Figure 1:** Bar plot with relative abundance of top 20 genera for the placenta study. Bars are grouped by sample type and case/control. There is a stark difference in composition between air, vaginal, and saliva samples to placental samples and controls (blank, H<sub>2</sub>O).

control) and preterm (case) newborns were collected for the maternal and fetal sides. Additionally, positive control samples were obtained from the mothers (saliva and cervicovaginal fluid) as well as negative control samples (air from the sample processing room, empty tubes, and PCR-grade water). The study was performed in both marker gene sequencing (amplicon) and metagenomics (MGS). The authors could not distinguish a unique placental microbiome that differs from the contamination background. We reanalyzed the samples in a standard pipeline with QIIME2 [6] for amplicon data and ganon for MGS data [86], generated a GRIMER report for both, and searched for the previously detected contamination.

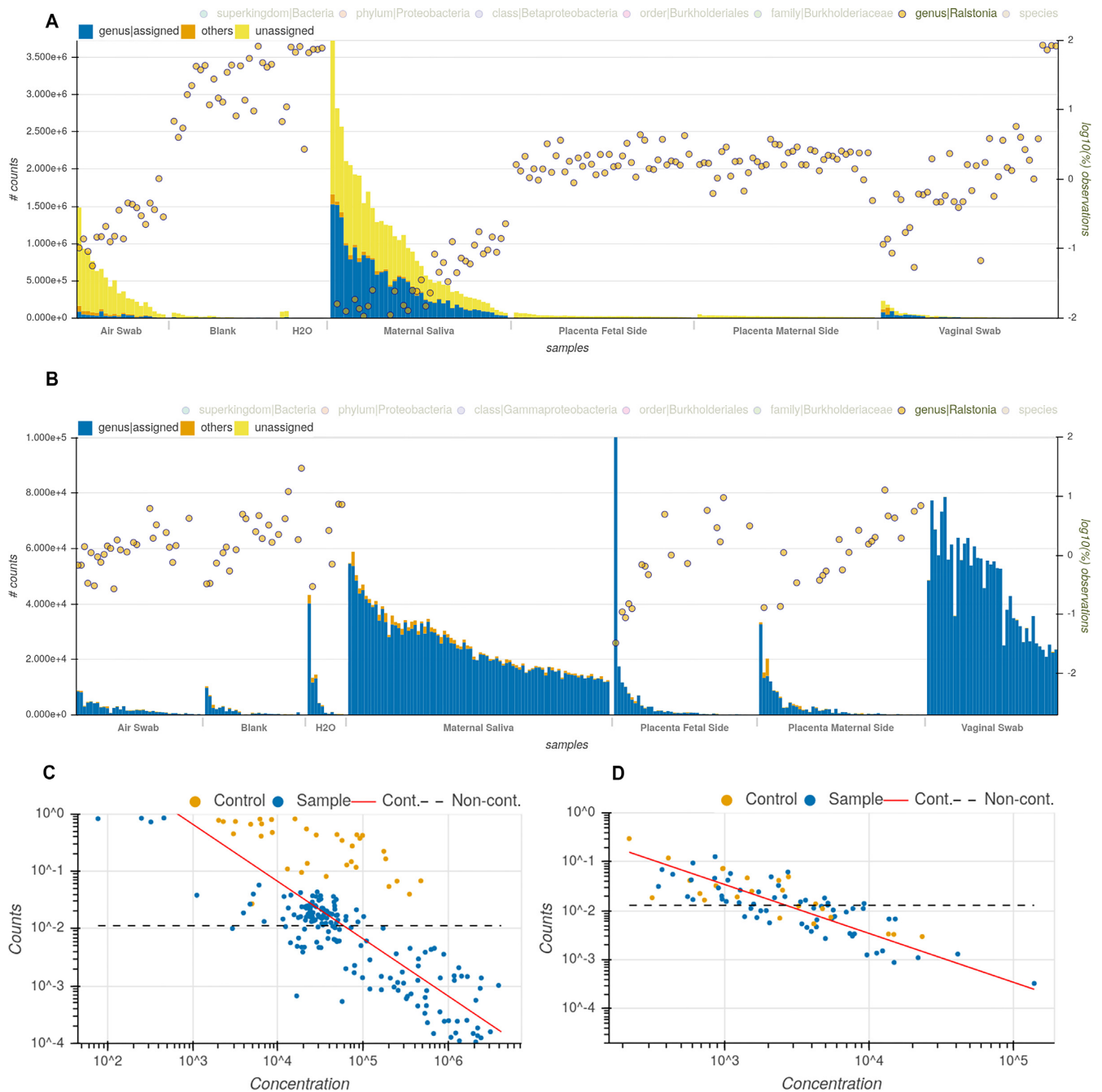
In the MGS report, the bar plot (Fig. 1) shows a stark difference in signal between sample types but a smaller difference in case and control groups. The *Ralstonia* genus is present in 96% of the all samples with an average abundance of 8.24%. Reads assigned for this genus were found in all negative control samples and H<sub>2</sub>O samples. *Ralstonia* was also reported in 12 studies as a common contaminant, based on our compiled contaminant list (Table 4), and it was classified as a contaminant by the DECONTAM method, based on the correlation of frequencies and the total number of reads per sample. Further, the abundance of this genus is higher in negative controls and placental samples as well as in samples with low number of reads, probably related to their low biomass, as depicted in the Fig. 2. Those results are in line with the ones reported in the original publication [85], even though the data were reanalyzed with a different set of tools, parameters, and reference databases. All evidence described pointing to *Ralstonia* as a contaminant was automatically generated by GRIMER and can be directly extracted from the Overview panel from the report. Besides human reads, *Ralstonia insidiosa* is the most prevalent species in this study. For the amplicon data, a similar pattern can be detected for the *Ralstonia* genus based on amplicon sequence variants (Fig. 2).

Further, all other taxa present can easily be verified for the same patterns. *Pseudomonas* show similar distribution and was also reported originally as probable contaminants in the placental samples. *Corynebacterium*, *Cutibacterium*, and *Mycobacterium*, although less prevalent, are further taxa with very similar patterns that could be potential contaminants and were not reported in the original publication.

### Multiple microbiome studies exploration

Definitive and robust conclusions from low-biomass studied environment are only possible with a set of controls and protocols to deal with contamination. KatharoSeq [87] is a well-designed protocol to better handle contamination in high-throughput low-biomass DNA microbial studies for amplicon sequencing or shotgun metagenomics. The protocol has guidelines for positive and negative controls implementation at the DNA extraction and library construction steps as well as computational approaches to define and exclude samples that did not achieve minimal amount of signal to be used. In their publication [87], the authors validate the protocol sequencing and analyzing with 3 low-biomass environments: the Jet Propulsion Laboratory spacecraft assembly facility (SAF), rooms of a neonatal intensive care unit (NICU), and an endangered abalone-rearing facility (abalone). A set of low biomass (LBM) negative controls to compare extraction kits is also included in the study.

We downloaded the OTU table and metadata from KatharoSeq evaluations for the 16S ribosomal RNA (rRNA) analyses available in Qiita [14] in the following configuration: reads trimmed at 150 bp and classified using closed-reference OTUs clustered at 97% similarity annotated with the greengenes taxonomy. A GRIMER report was generated for the raw table with all samples without any filtration. The heatmap generated for the annotated species level (Fig. 3) shows a distinct and clear pattern between environments and the LBM. As reported in the publication, abalone samples have a higher richness (here as species annotated OTUs) as well as the highest average number of reads per sample. It is possible to identify potential contaminants in the study by looking for observations prevalent across environments and the relation to its annotations. Using this analysis, we detected *Cutibacterium acnes*, which is reported as a common contaminant and human-related species, present among all 4 environments studied as well as highly frequent in negative and positive controls. Even though DECONTAM did not identify this taxon as a contaminant, related data still hold strong evidence for contaminant of *C. acnes* in this study. Furthermore, *Staphylococcus aureus* and *Staphylococcus epidermidis*, known as human-related bacteria, were detected in high abundances in both NICU and SAF environments—areas with low and high human exposure, respectively. However, both species were also relatively highly present in negative con-



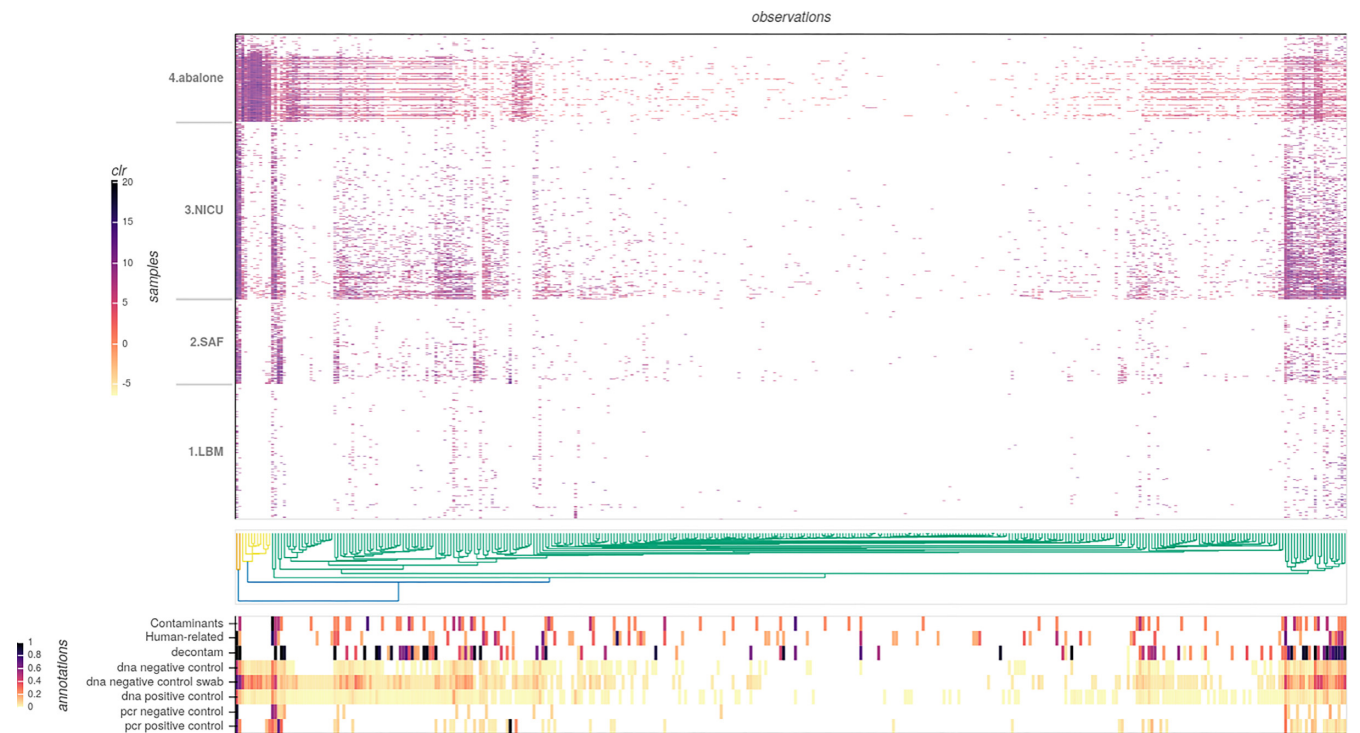
**Figure 2:** Evidence supporting *Ralstonia* as a contaminant taxon in the placenta study. (A, B) Right y-axis shows normalized abundance of genus *Ralstonia* in log scale for each sample in the MGS (A) and amplicon (B) data. Bars (left y-axis) summarized counts at the genus level for each sample. Samples are grouped by sample type and sorted by number of reads (x-axis). The yellow circles show abundance of *Ralstonia*, which is higher in the control samples (blank and H<sub>2</sub>O) as well as increased in real samples with low read count. (C, D) DECONTAM plots for *Ralstonia* genus for the MGS (C) and amplicon (D) data. DECONTAM plots show that taxa counts follow the expected distribution for contamination based on the number of reads per sample (red line).

controls, the abalone environment, and LBM samples. Additionally, both were positively classified as contamination by DECONTAM, indicating that besides human exposure, those organisms could be driven by an external source of contamination.

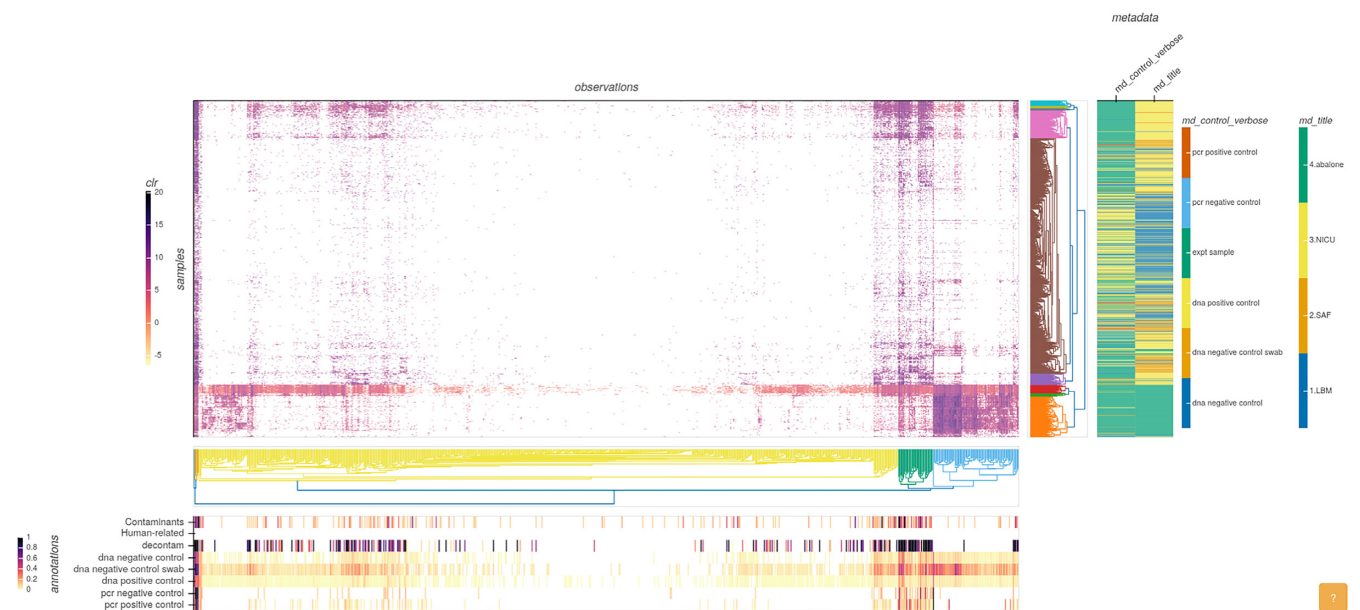
Species identification based on 16S rRNA is limited due to its low resolution: approximately 15% of the OTUs are annotated at species level and 69% at genus level in this study. The same analysis visualized at genus level gives an increased perception of the distribution of the data in this study. With a higher signal, it is possible to visualize how several clusters are formed and in many

cases agree in multiple levels of evidence supporting the possibility of contamination (Fig. 4).

Looking at the correlation between top observations reported (Fig. 5), a matrix of highly correlated genera can be detected. Such a pattern was previously reported to be an indication of contamination from reagent-derived sources since they are invariably present within samples in similar ratios [45]. Further inspection of those genera (*Glaciecola*, *Leucothrix*, *Mycoplasma*, *Oleibacter*, *Polaribacter*, *Pseudoalteromonas*, *Psychrilyobacter*, *Psychromonas*, *Shewanella*) shows that they are mainly from Aquatic/Marine biomes with

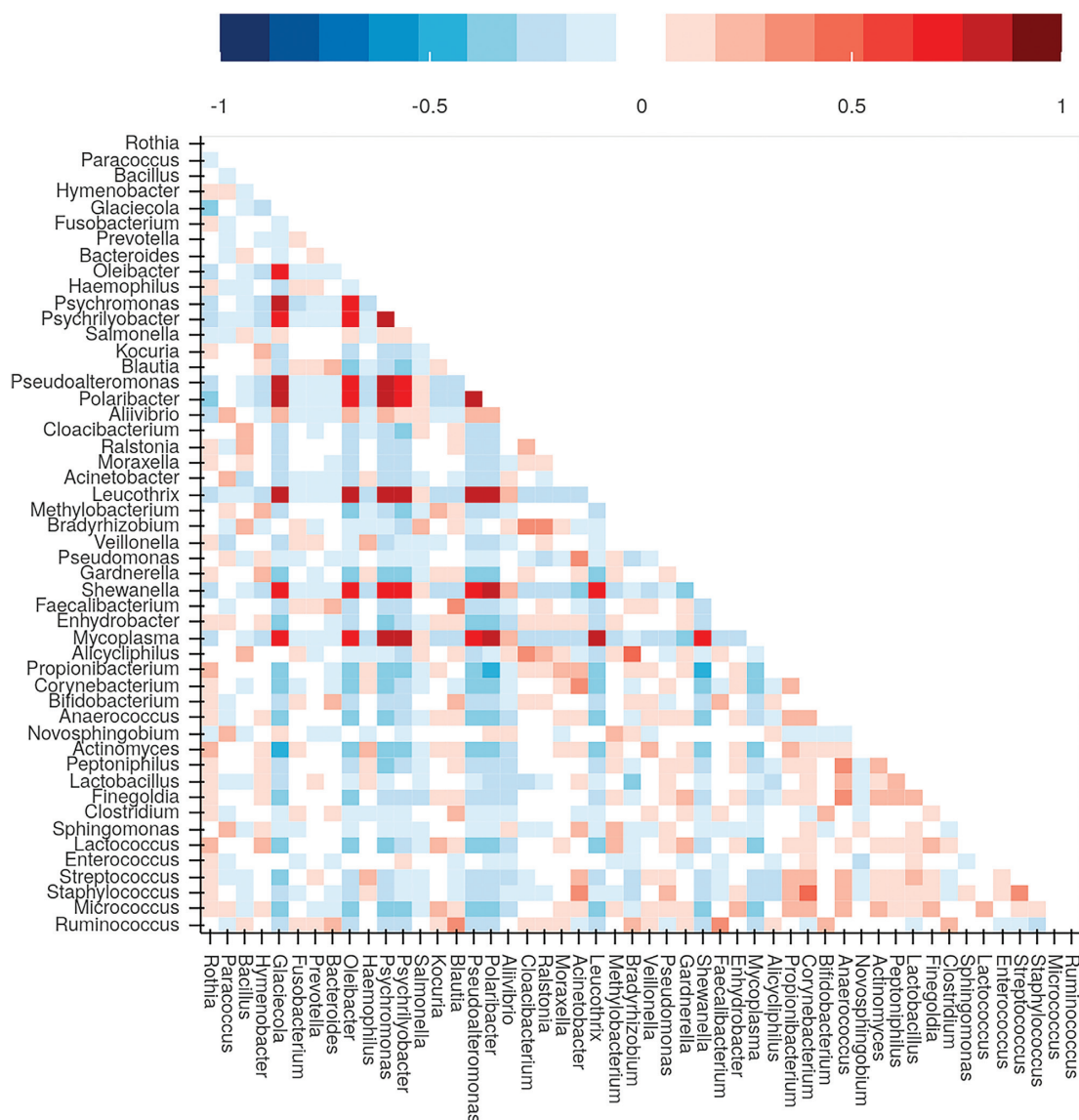


**Figure 3:** Heatmap visualization at species level for the KatharoSeq data. Samples are grouped by study type (y-axis) and clustered by observations (x-axis, euclidean distance metric, complete method). Data in the heatmap are center log ratio transformed. Bottom panel shows annotation related to the observations. “Contaminants” and “human-related” annotations are normalized counts against precompiled list of references described in this article. “Decontam” is the normalized DECONTAM *P* score. All “control” annotations show the proportion of the observation in the indicated group of control samples.



**Figure 4:** Heatmap visualization at genus level for the KatharoSeq data. Samples and observations axis are clustered and sorted based on the euclidean distance metric, complete method. Data in the heatmap are center log ratio transformed. Bottom panel shows annotation related to the observations. “Contaminants” and “human-related” annotations are normalized counts against precompiled list of references described in this article. “Decontam” is the normalized DECONTAM *P* score. All “control” annotations show the proportion of the observation in the indicated group of control samples. Metadata panel show color-coded sample information on study (*md\_title*) and type of sample (*md\_control\_verbose*). The annotation panel shows higher values on multiple sources of evidence for contamination relative to data clusters of the heatmap. Metadata panel shows how samples show independent patterns based on the environment (*md\_title*) and difference from controls (*md\_control\_verbose*).





**Figure 5:** Symmetric proportionality coefficient (rho correlation) between top 50 most abundant genera in the KatharoSeq data. Positive correlation values (between 0 and 1) are displayed in red. Negative correlation values (between  $-1$  and 0) are displayed in blue. Highly correlated matrix among 9 genera (dark red) points to reagent-derived contamination, when considered with other lines of evidence (Fig. 6).

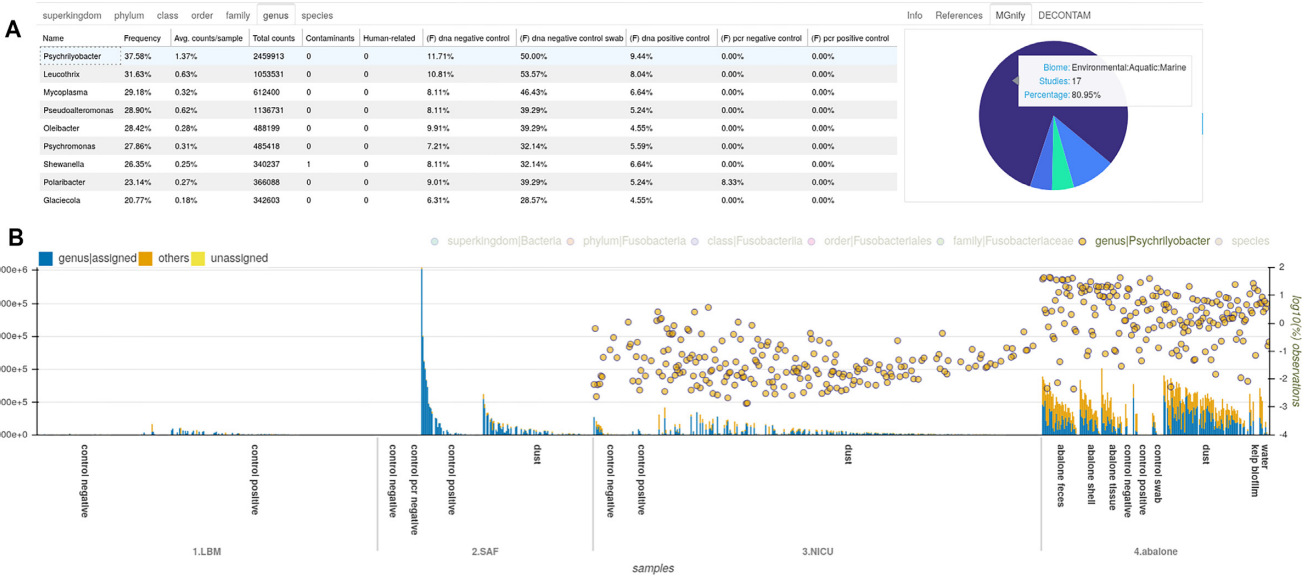
help of the matching results with the MGnify database (Fig. 6B). Further, they are more prevalent in negative controls (Fig. 6A), an evidence of DNA extraction kit or sample processing contaminants. Those organisms are highly frequent in the abalone study, which is a Marine environment, and some of them were also described in the original publication. Although in very low amounts, those groups were also reported present in NICU samples (Fig. 6C), pointing to possible well-to-well contamination.

## Discussion

GRIMER is an easy-to-use and accessible tool for specialists and nonspecialists that generates a concise interactive offline dashboard with a set of analyses, visualizations, and data connections from a simple table of counts. It automatically summarizes several levels of evidence to better understand the relation between observations, samples, metadata, and taxonomy. GRIMER reports are a valuable resource for investigating contamination, a problem that affects every microbiome study to some degree.

All the conclusion and visualizations presented in this work in the Results section were solely based on GRIMER reports, showing that microbiome analysis, contamination investigation, and detection are possible with the methodology proposed. The use of multiple sources of evidence to annotate observations improves the ability to better detect clear contaminants in microbiome studies as well as to point to probable groups of candidate contaminants.

In addition to the GRIMER software, we compiled and provided in this work a list of common taxa contaminants based on 22 publications (Table 3). Many of the reported contaminants are recurrent in diverse studies, pointing to a consensus for some taxa (Table 4) as a probable contaminant. Taxa in this list cannot be strictly considered a contaminant by itself. However, they can corroborate suspicious contamination discovered via several other lines of evidence without the extra effort of researching the literature. The presented list is not comprehensive but a first step to centralize and standardize recurring contaminants described in the literature. We expect this list to incrementally grow over time



**Figure 6:** GRIMER overview panel plots (A) listing of 9 highly correlated genera detected in Fig. 5. Samples have high incidence in DNA negative controls. MGNify plot showing proportion of biomes related to *Psychrilyobacter* in the whole MGNify database. (B) bar plot listing samples (x-axis), grouped by study and sample type and sorted by total number of reads. Bars represent the total number of counts for each sample and are annotated with the proportion assigned to genus level (left y-axis). Log-transformed abundance of *Psychrilyobacter* is displayed in yellow circles (right y-axis). This taxon is abundant in the abalone samples but has some signal in the NICU samples that are inversely correlated to the total amount of reads, pointing to potential contamination. The other 8 taxa show similar patterns in the report.

as more evidence of kit and laboratory contamination becomes available. The information of common contaminants is a valuable resource to aid contamination detecting, and we are willing to keep and extend it. Improvements to the list and suggestions of further candidate taxa can be provided via the GRIMER repository at [48]. As a future work, the list can be associated with study details as biome, extraction kit, and methodology to be further queried and integrated in more details.

Additionally to the aforementioned common contaminants, GRIMER can also use general lists of custom organisms to annotate samples. In this work and by default, human-related organisms commonly occurring in human skin and oral and nasal cavities, as well as face and other human limbs, are used since they can be external sources of contamination. Those lists can be easily provided as taxonomic identifiers or names to GRIMER. If the target study conflicts with any of those environments (e.g., study of human skin), one could simply remove the related entries from the configuration files. More details and examples on how to perform this can be found in the online documentation.

GRIMER works out-of-the-box with as little data as possible but can incrementally expand the reports when more data are provided and can be adapted for user necessities. GRIMER is fast and generates reports in a matter of seconds on a standard notebook. The outcome dashboard is lightweight and can handle hundreds to thousands of samples and observations. Report sizes usually vary from 1 to 10 MB and are highly compressible, since they are text-based HTML files. GRIMER reports with a higher number of samples (thousands) can grow significantly in size (10–100 MB) but still run normally. If report size is a limitation, many options can be adjusted to generated more compact files: reducing number of taxonomy ranks displayed, less combinations of analyses, and filtering very low-abundant observations, among others.

One of the core strengths of GRIMER is the taxonomy automation. It accepts taxonomic identifiers from several different taxonomies but also parses names and converts them to their respec-

tive identifiers. If only 1 taxonomic level is provided (e.g., species level), GRIMER can decompose and summarize the data in higher ranks. That means that users do not have to handle taxonomy and everything will work automatically. GRIMER was developed in a way that new visualizations can be included with little effort.

We listed and summarized a list of similar currently available methods published in the last 10 years (Table 2) as well as web platforms for complete analyses of microbiome data (Table 1). A list of functionalities between similar available tools is provided in [9], but a detailed comparison with GRIMER is out of the scope of this work. Most methods share some basic functions (e.g., taxonomic abundance analysis) but are diverse in many other aspects and were sometimes developed with specific goals (e.g., function analysis, biomarker identification). However, there is no comprehensive method that can provide a complete solution for the many possible analyses in a microbiome study. We believe that many of those tools, besides their overlapping functions, are complementary and can be used concurrently. GRIMER mainly shares features with pavian [22] in terms of general microbiome exploration and support to metagenomics data and with OpenContami [28] regarding contaminant detection. GRIMER, however, is unique in its output format. The vast majority of the currently available tools are web based, are hosted in a remote server, or rely on a local hosted web server to properly work (Table 2). This may be impractical for many nonspecialists and for long-term storage and reproducibility. GRIMER reports are portable and fully functional offline. This allows analysis to be accessible by many researchers with different backgrounds working together in the same study, increasing direct interaction with data. The portability also enables better documentation of results, reproducibility, and shareability. Further, web-based tools may disappear after some years of inactivity or lack of funding, and analysis may be lost, as it is the case for for some methods (Table 5). GRIMER reports are completely offline and will work as long as the report file is safely stored.

**Table 5:** Tools and web resources no longer available, supported, or inaccessible (as of 28 February 2022)

Name	Reason	Year	Reference
Community-analyzer	Website offline	2013	[88]
calypso	Website offline	2017	[89]
Metaviz	Web tool not responsive	2018	[90]
iMAP	No longer supported due to funding	2019	[91]
biomminer	Page not found	2020	[92]

Overall, we believe that GRIMER is a valuable contribution to the microbiome field and can facilitate data exploration, analysis, and contamination detection.

## Data Availability

In addition to the GitHub repository [48], an archival copy of the code and supporting data are available via the GigaScience repository, GigaDB [93].

GRIMER reports presented in this article are available for download via Zenodo [81] and can be used interactively online [82].

The datasets and metadata for the placenta study were obtained from ENA database with the Bioproject ID PRJNA451186.

The datasets and metadata for the KatharoSeq study were obtained from the qiita website [94] (log-in required).

## Availability and Requirements

Project name: GRIMER

Project homepage: <https://github.com/pirovc/grimer>

Operating system(s): Platform independent

Programming language: Python 3.5 or higher

Other requirements: bokeh 2.2.3 or higher

License: MIT License

Any restrictions to use by nonacademics: Use based on MIT license

RRID: SCR\_023265

biotools: grimer

## Abbreviations

bp: base pairs; MGS: metagenomics; NCBI: The National Center for Biotechnology Information; NICU: neonatal intensive care unit; OTU: operational taxonomic unit; rRNA: ribosomal RNA.

## Competing Interests

The authors declare no competing interests.

## Funding

This work was financially supported by the by the Deutsche Forschungsgemeinschaft (DFG, German Research Foundation) project number 458163427 and by the German Ministry for Education and Research (Bundesministerium für Bildung und Forschung - BMBF) grant number 01KI1905D. We acknowledge support by the OpenAccess Publication Fund of Freie Universität Berlin.

## Authors' Contributions

V.C.P. and B.Y.R. conceptualized and developed the idea. V.C.P. wrote the software and main manuscript text, including all anal-

yses. B.Y.R. reviewed and contributed to the manuscript text. All authors reviewed the manuscript.

## Acknowledgments

We thank all partners in the ZooSeq project for helpful discussions and support in this project.

## References

- Pollock, J, Glendinning, L, Wisedchanwet, T, et al. The madness of microbiome: attempting to find consensus “best practice” for 16S microbiome studies. *Appl Environ Microbiol* 2018;**84**(7):e02627–17.
- Kim, D, Hofstaedter, CE, Zhao, C, et al. Optimizing methods and dodging pitfalls in microbiome research. *Microbiome* 2017;**5**(1):52.
- Knight, R, Vrbanac, A, Taylor, BC, et al. Best practices for analysing microbiomes. *Nat Rev Microbiol* 2018;**16**(7):410–22.
- Kayani, MuR, Huang, W, Feng, R, et al. Genome-resolved metagenomics using environmental and clinical samples. *Brief Bioinform* 2021;**22**(5):bbab030.
- Gloor, GB, Macklaim, JM, Pawlowsky-Glahn, V, et al. Microbiome datasets are compositional: and this is not optional. *Front Microbiol* 2017;**8**: 1–6.
- Bolyen, E, Rideout, JR, Dillon, MR, et al. Reproducible, interactive, scalable and extensible microbiome data science using QIIME 2. *Nat Biotechnol* 2019;**37**(8):852–7.
- Schloss, PD, Westcott, SL, Ryabin, T, et al. Introducing mothur: open-source, platform-independent, community-supported software for describing and comparing microbial communities. *Appl Environ Microbiol* 2009;**75**(23):7537–41.
- McMurdie, PJ, Holmes, S. phyloseq: An R package for reproducible interactive analysis and graphics of microbiome census data. *PLoS ONE* 2013;**8**(4):e61217.
- Peeters, J, Thas, O, Shkedy, Z, et al. Exploring the microbiome analysis and visualization landscape. *Front Bioinform* 2021;**1**:774631.
- Meyer, F, Bagchi, S, Chaterji, S, et al. MG-RAST version 4—lessons learned from a decade of low-budget ultra-high-throughput metagenome analysis. *Brief Bioinform* 2019;**20**(4):1151–9.
- Mitchell, AL, Almeida, A, Beracochea, M, et al. MGnify: the microbiome analysis resource in 2020. *Nucleic Acids Res* 2020;**48**(D1):D570–8.
- Oliveira, FS, Brestelli, J, Cade, S, et al. MicrobiomeDB: a systems biology platform for integrating, mining and analyzing microbiome experiments. *Nucleic Acids Res* 2018;**46**(D1):D684–91.
- Weber, N, Liou, D, Dommer, J, et al. Nephele: a cloud platform for simplified, standardized and reproducible microbiome data analysis. *Bioinformatics* 2018;**34**(8):1411–3.
- Gonzalez, A, Navas-Molina, JA, Kosciulek, T, et al. Qiita: rapid, web-enabled microbiome meta-analysis. *Nat Methods* 2018;**15**(10):796–8.

15. Arndt, D, Xia, J, Liu, Y, et al. METAGENassist: a comprehensive web server for comparative metagenomics. *Nucleic Acids Res* 2012;**40**(W1):W88–5.
16. Huse, SM, Mark Welch, DB, Voorhis, A, et al. VAMPS: a website for visualization and analysis of microbial population structures. *BMC Bioinform* 2014;**15**(1):41.
17. McMurdie, PJ, Holmes, S. Shiny-phyloseq: web application for interactive microbiome analysis with provenance tracking. *Bioinformatics* 2015;**31**(2):282–3.
18. Wang, Y, Xu, L, Gu, YQ, et al. MetaCoMET: a web platform for discovery and visualization of the core microbiome. *Bioinformatics* 2016;**32**(22):3469–70.
19. Laczny, CC, Kiefer, C, Galata, V, et al. BusyBee Web: metagenomic data analysis by bootstrapped supervised binning and annotation. *Nucleic Acids Res* 2017;**45**(W1):W171–W179.
20. Dhariwal, A, Chong, J, Habib, S, et al. MicrobiomeAnalyst: a web-based tool for comprehensive statistical, visual and meta-analysis of microbiome data. *Nucleic Acids Res* 2017;**45**(W1):W180–8.
21. McNally, CP, Eng, A, Noecker, C, et al. BURRITO: an interactive multi-omic tool for visualizing taxa–function relationships in microbiome data. *Front Microbiol* 2018;**9**:365.
22. Breitwieser, FP, Salzberg, SL. Pavian: interactive analysis of metagenomics data for microbiome studies and pathogen identification. *Bioinformatics* 2020;**36**(4):1303–1304.
23. Tong, WM, Chan, Y. GenePiper, a graphical user interface tool for microbiome sequence data mining. *Microbiol Resour Announc* 2020;**9**(1):e01195–19.
24. Zhao, Y, Federico, A, Faits, T, et al. animalcules: interactive microbiome analytics and visualization in R. *Microbiome* 2021;**9**(1):76.
25. Reeder, J, Huang, M, Kaminker, JS, et al. MicrobiomeExplorer: an R package for the analysis and visualization of microbial communities. *Bioinformatics* 2021;**37**(9):1317–8.
26. Barnett, DJm, Arts, ICw, Penders, J. microViz: an R package for microbiome data visualization and statistics. *J Open Source Softw* 2021;**6**(63):3201.
27. Dietrich, A, Matchado, MS, Zwiebel, M, et al. Namco: a microbiome explorer. *bioRxiv* 2021. <https://www.biorxiv.org/content/10.1101/2021.12.15.471754v1>.
28. Park, SJ, Nakai, K. OpenContami: a web-based application for detecting microbial contaminants in next-generation sequencing data. *Bioinformatics* 2021;**37**(18):3021–2.
29. Su, SC, Galvin, JE, Yang, SF, et al. wiSDOM: a visual and statistical analytics for interrogating microbiome. *Bioinformatics* 2021;**37**(17):2795–7.
30. Jin, BT, Xu, F, Ng, RT, et al. Mian: interactive web-based microbiome data table visualization and machine learning platform. *Bioinformatics* 2022;**38**(4):1176–8.
31. Fricke, WF, Ravel, J. Microbiome or no microbiome: are we looking at the prenatal environment through the right lens? *Microbiome* 2021;**9**(1):9.
32. Blaser, MJ, Devkota, S, McCoy, KD, et al. Lessons learned from the prenatal microbiome controversy. *Microbiome* 2021;**9**(1):8.
33. Walter, J, Hornef, MW. A philosophical perspective on the prenatal in utero microbiome debate. *Microbiome* 2021;**9**(1):5.
34. Rand, KH, Houck, H. Taq polymerase contains bacterial DNA of unknown origin. *Mol Cell Probes* 1990;**4**(6):445–50.
35. Salter, SJ, Cox, MJ, Turek, EM, et al. Reagent and laboratory contamination can critically impact sequence-based microbiome analyses. *BMC Biol* 2014;**12**(1):87.
36. Glassing, A, Dowd, SE, Galandiuk, S, et al. Inherent bacterial DNA contamination of extraction and sequencing reagents may affect interpretation of microbiota in low bacterial biomass samples. *Gut Pathogens* 2016;**8**(1):24.
37. Minich, JJ, Sanders, JG, Amir, A, et al. Quantifying and understanding well-to-well contamination in microbiome research. *mSystems* 2019;**4**(4):e00186–19.
38. Eisenhofer, R, Minich, JJ, Marotz, C, et al. Contamination in low microbial biomass microbiome studies: issues and recommendations. *Trends Microbiol* 2019;**27**(2):105–17.
39. Gruber, K. Here, there, and everywhere. *EMBO Rep* 2015;**16**(8):898–901.
40. Breitwieser, FP, Perteu, M, Zimin, A, et al. Human contamination in bacterial genomes has created thousands of spurious proteins. *Genome Res* 2019;**29**(6):954–960.
41. Hornung, BVH, Zwitter, RD, Kuijper, EJ. Issues and current standards of controls in microbiome research. *FEMS Microbiol Ecol* 2019;**95**(5).
42. Jervis-Bardy, J, Leong, LEX, Marri, S, et al. Deriving accurate microbiota profiles from human samples with low bacterial content through post-sequencing processing of Illumina MiSeq data. *Microbiome* 2015;**3**(1):19.
43. Davis, NM, Proctor, DM, Holmes, SP, et al. Simple statistical identification and removal of contaminant sequences in marker-gene and metagenomics data. *Microbiome* 2018;**6**(1):226.
44. Marsh, RL, Nelson, MT, Pope, CE, et al. How low can we go? The implications of low bacterial load in respiratory microbiota studies. *Pneumonia* 2018;**10**(1):7.
45. Goffau, MCD, Lager, S, Salter, SJ, et al. Recognizing the reagent microbiome. *Nat Microbiol* 2018;**3**(8):851–3.
46. Harrison, JG, Randolph, GD, Buerkle, CA. Characterizing microbiomes via sequencing of marker loci: techniques to improve throughput, account for cross-contamination, and reduce cost. *mSystems* 2021;**0**:e00294–21.
47. Olomu, IN, Pena-Cortes, LC, Long, RA, et al. Elimination of “kitome” and “splashome” contamination results in lack of detection of a unique placental microbiome. *BMC Microbiol* 2020;**20**(1):157.
48. Piro, VC. *pirovc/grimer*. 2023. <https://github.com/pirovc/grimer>.
49. Schoch, CL, Ciufu, S, Domrachev, M, et al. NCBI Taxonomy: a comprehensive update on curation, resources and tools. *Database* 2020;**2020**:baaa062.
50. Tanner, MA, Goebel, BM, Dojka, MA, et al. Specific ribosomal DNA sequences from diverse environmental settings correlate with experimental contaminants. *Appl Environ Microbiol* 1998;**64**(8):3110–3.
51. Kulakov, LA, McAlister, MB, Ogden, KL, et al. Analysis of bacteria contaminating ultrapure water in industrial systems. *Appl Environ Microbiol* 2002;**68**(4):1548–55.
52. Grahm, N, Olofsson, M, Ellnebo-Svedlund, K, et al. Identification of mixed bacterial DNA contamination in broad-range PCR amplification of 16S rDNA V1 and V3 variable regions by pyrosequencing of cloned amplicons. *FEMS Microbiol Lett* 2003;**219**(1):87–91.
53. Barton, HA, Taylor, NM, Lubbers, BR, et al. DNA extraction from low-biomass carbonate rock: an improved method with reduced contamination and the low-biomass contaminant database. *J Microbiol Methods* 2006;**66**(1):21–31.
54. Laurence, M, Hatzis, C, Brash, DE. Common contaminants in next-generation sequencing that hinder discovery of low-abundance microbes. *PLoS ONE* 2014;**9**(5):e97876.
55. Jousset, E, Clamens, AL, Galan, M, et al. Assessment of a 16S rRNA amplicon Illumina sequencing procedure for studying

- the microbiome of a symbiont-rich aphid genus. *Mol Ecol Res* 2015;**16**(3):628–40.
56. Lauder, AP, Roche, AM, Sherrill-Mix, S, et al. Comparison of placenta samples with contamination controls does not provide evidence for a distinct placenta microbiota. *Microbiome* 2016;**4**(1):29.
  57. Lazarevic, V, Gaia, N, Girard, M, Schrenzel, J. Decontamination of 16S rRNA gene amplicon sequence datasets based on bacterial load assessment by qPCR. *BMC Microbiol* 2016;**16**(1):73.
  58. Salter, SJ, Turner, C, Wathnanaworawit, W, et al. A longitudinal study of the infant nasopharyngeal microbiota: the effects of age, illness and antibiotic use in a cohort of South East Asian children. *PLoS Neglected Trop Dis* 2017;**11**(10):e0005975.
  59. Kirstahler, P, Bjerrum, SS, Friis-Møller, A, et al. Genomics-based identification of microorganisms in human ocular body fluid. *Sci Rep* 2018;**8**(1):4126.
  60. Stinson, LF, Keelan, JA, Payne, MS. Comparison of meconium DNA extraction methods for use in microbiome studies. *Front Microbiol* 2018;**9**: 270.
  61. Stinson, LF, Keelan, JA, Payne, MS. Identification and removal of contaminating microbial DNA from PCR reagents: impact on low-biomass microbiome analyses. *Lett Appl Microbiol* 2019;**68**(1):2–8.
  62. Weyrich, LS, Farrer, AG, Eisenhofer, R, et al. Laboratory contamination over time during low-biomass sample analysis. *Mol Ecol Res* 2019;**19**(4):982–6.
  63. de Goffau, MC, Lager, S, Sovio, U, et al. Human placenta has no microbiome but can contain potential pathogens. *Nature* 2019;**572**(7769):329–4.
  64. Nejman, D, Livyatan, I, Fuks, G, et al. The human tumor microbiome is composed of tumor type-specific intracellular bacteria. *Science* 2020;**368**(6494):973–80.
  65. Kjartansdóttir, KR, Friis-Nielsen, J, Asplund, M, et al. Traces of ATCV-1 associated with laboratory component contamination. *Proc Natl Acad Sci* 2015;**112**(9):E925–6.
  66. Mukherjee, S, Huntemann, M, Ivanova, N, et al. Large-scale contamination of microbial isolate genomes by Illumina PhiX control. *Standards Genomic Sci* 2015;**10**(1):18.
  67. Asplund, M, Kjartansdóttir, KR, Møllerup, S, et al. Contaminating viral sequences in high-throughput sequencing viromics: a linkage study of 700 sequencing libraries. *Clin Microbiol Infect* 2019;**25**(10):1277–85.
  68. Czurda, S, Smelik, S, Preuner-Stix, S, et al. Occurrence of fungal DNA contamination in PCR reagents: approaches to control and decontamination. *J Clin Microbiol* 2016;**54**(1):148–52.
  69. Reimer, LC, Vetcinina, A, Carbasse, JS, et al. BacDive in 2019: bacterial phenotypic data for high-throughput biodiversity analysis. *Nucleic Acids Res* 2019;**47**(D1):D631–6.
  70. Escapa, IF, Chen, T, Huang, Y, et al. New insights into human nostril microbiome from the expanded human oral microbiome database (eHOMD): a resource for the microbiome of the human aerodigestive tract. *mSystems* 2018;**3**(6): e00187–18.
  71. Byrd, AL, Belkaid, Y, Segre, JA. The human skin microbiome. *Nat Rev Microbiol* 2018;**16**(3):143–55.
  72. MGnify API. Archiving, analysis and integration of metagenomics data. EMBL-EBI. 2023. <https://www.ebi.ac.uk/metagenomics/api/v1/> [Accessed 21 Feb 2023].
  73. McDonald, D, Clemente, JC, Kuczynski, J, et al. The biological observation matrix (BIOM) format or: how I learned to stop worrying and love the ome-ome. *GigaScience* 2012;**1**(1): 7.
  74. Lovell, D, Pawlowsky-Glahn, V, Egozcue, JJ, et al. Proportionality: a valid alternative to correlation for relative data. *PLoS Comput Biol* 2015;**11**(3):e1004075.
  75. Erb, I, Notredame, C. How should we measure proportionality on relative gene expression data? *Theory Biosci* 2016;**135**(1–2):21–36.
  76. Team, BD. Bokeh: Interactive Data Visualization in the browser, from Python. 2023. <https://bokeh.org/> [Accessed 21 Feb 2023].
  77. McKinney, W. *Data structures for statistical computing in python*. Proceedings of the 9th Python in Science Conference, Austin, Texas; 2010. p. 56–61.
  78. Virtanen, P, Gommers, R, Oliphant, TE, et al. SciPy 1.0: fundamental algorithms for scientific computing in Python. *Nat Methods* 2020;**17**(3):261–72.
  79. scikit-bio Development Team, . scikit-bio: a bioinformatics library for data scientists, students, and developers. 2023. <http://scikit-bio.org> [Accessed 21 Feb 2023].
  80. Piro, VC. MultiTax GitHub. 2022. <https://github.com/pirovc/multitax>.
  81. Piro, VC, Renard, BY. Contamination detection and microbiome exploration with GRIMER. *Zenodo*. 2023. <https://dx.doi.org/10.5281/zenodo.7103846>.
  82. Piro, VC. *grimer-reports*. 2023.
  83. Silverstein, RB, Mysorekar, IU. Group therapy on in utero colonization: seeking common truths and a way forward. *Microbiome* 2021;**9**(1):7.
  84. Sterpu, I, Fransson, E, Hugerth, LW, et al. No evidence for a placental microbiome in human pregnancies at term. *Am J Obstet Gynecol* 2021;**224**(3):296.
  85. Leiby, JS, McCormick, K, Sherrill-Mix, S, et al. Lack of detection of a human placenta microbiome in samples from preterm and term deliveries. *Microbiome* 2018;**6**(1):196.
  86. Piro, VC, Dadi, TH, Seiler, E, et al. ganon: precise metagenomics classification against large and up-to-date sets of reference sequences. *Bioinformatics* 2020;**36**(Suppl 1):i12–20.
  87. Minich, JJ, Zhu, Q, Janssen, S, et al. KatharoSeq enables high-throughput microbiome analysis from low-biomass samples. *mSystems* 2018;**3**(3):e00218–17.
  88. Kuntal, BK, Ghosh, TS, Mande, SS. Community-analyzer: a platform for visualizing and comparing microbial community structure across microbiomes. *Genomics* 2013;**102**(4):409–18.
  89. Zakrzewski, M, Proietti, C, Ellis, JJ, et al. Calypso: a user-friendly web-server for mining and visualizing microbiome-environment interactions. *Bioinformatics* 2017;**33**(5):782–3.
  90. Wagner, J, Chelaru, F, Kancherla, J, et al. Metaviz: interactive statistical and visual analysis of metagenomic data. *Nucleic Acids Res* 2018;**46**(6):2777–87.
  91. Buza, TM, Tonui, T, Stomeo, F, et al. iMAP: an integrated bioinformatics and visualization pipeline for microbiome data analysis. *BMC Bioinformatics* 2019;**20**(1):374.
  92. Shamsaddini, A, Dadkhah, K, Gillevet, PM. BiomMiner: an advanced exploratory microbiome analysis and visualization pipeline. *PLoS One* 2020;**15**(6):e0234860.
  93. Piro, VC, Renard, BY. Supporting data for “Contamination Detection and Microbiome Exploration with GRIMER.” *GigaScience Database*. 2023. <http://dx.doi.org/10.5524/102359>.
  94. *KatharoSeq enables high-throughput microbiome analysis from low-biomass samples*. 2023. <https://qiita.ucsd.edu/study/description/10934>.