

26 novel and existing analytics, visualization methods, and machine learning models. For example,
27 traditional microbiome analyses such as alpha/beta diversity and differential abundance analysis
28 are enhanced in the toolkit, while new methods such as biomarker identification are introduced.
29 Powerful interactive and dynamic figures generated by *animalcules* enable users to understand
30 their data and discover new insights. *animalcules* can be used as a standalone command-line R
31 package or users can explore their data with the accompanying interactive R Shiny interface.

32 **Conclusions:** We present *animalcules*, an R package for interactive microbiome analysis
33 through either an interactive interface facilitated by R Shiny or various command-line functions.
34 It is the first microbiome analysis toolkit that supports the analysis of all 16S rRNA, DNA-based
35 shotgun metagenomics, and RNA-sequencing based metatranscriptomics datasets. *animalcules*
36 can be freely downloaded from GitHub at <https://github.com/compbiomed/animalcules> or
37 installed through Bioconductor at
38 <https://www.bioconductor.org/packages/release/bioc/html/animalcules.html>.

39 **Keywords:** Microbiome analysis, Visualization, Interactive toolkit, Biomarker identification
40

41 **Background**

42 The complex role of the gut microbiota in shaping human health and disease has been
43 intensely investigated and explored in recent years, largely due to the availability of culture-
44 independent molecular-based high-throughput sequencing technologies. It is estimated that every
45 human host coexists with an average of 500-1000 different bacterial species [1–3] and research
46 has discovered that the microbiome is associated with host lifestyle and diet [4, 5] as well as
47 many diseases such as obesity, type 2 diabetes [6] and cancer [7]. Burgeoning sequencing
48 technology brings not only more data and capacity for microbiome research, but also new

49 challenges for data analytics and interpretation. Improved tools and methods for microbiome
50 data analytics will enhance our ability to understand the roles of microbes in diverse
51 environments, particularly understanding how they interact with each other as well as their
52 human hosts.

53 Current microbiome analysis typically consists of two important components: upstream
54 community profiling (e.g. what is the abundance of all microbes in each sample?) and
55 downstream high-level analysis (e.g. alpha/beta diversity analysis, differential abundance
56 analysis) [8]. In recent years, evolving data analytics, visualization, and machine learning
57 methods have been gradually applied to the development of many software tools and web servers
58 for microbiome data analysis covering these two components. [9–13]. However, new techniques
59 and sequencing technologies have steepened the learning curve for scientific researchers
60 applying new methods for microbiome data analysis and interpretation [14]. Furthermore,
61 existing tools are mostly dedicated to one aspect of analysis and/or are restricted to analyzing
62 one type of microbiome data. For example, while there are many tools and workflows for
63 analyzing 16S rRNA data, there are no existing tools and pipelines tailored for comprehensively
64 addressing the analytical needs of RNA-based metatranscriptomics.

65 **Table 1** gives a summary of the functions of these tools with respect to the analysis needs
66 of microbiome data. For marker gene-based data such as 16S rRNA, QIIME II [15] and Mothur
67 [16] provide a user interface and a plethora of analytic and visualization tools, but do not provide
68 support for metagenomic and metatranscriptomic data. Vegan [17] provides a wide variety of
69 functions for metagenomic data visualization, but lacks a user-interface, and tools for host and
70 microbial read alignment, differential expression, etc. BioBakery [18] provides a comprehensive
71 suite of tools for most metagenomic analysis needs for microbial communities, but relies on a

72 small set of markers to identify species, and does not address host or microbial expression.
73 Phyloseq [19] has a Shiny interface with tools for annotation, visualization, and diversity
74 analysis, but does not provide abundance analysis, and is no longer actively maintained by its
75 developers. None of these methods are comprehensive or specifically address the needs for
76 multiple types of 16S rRNA, metagenomic or metatranscriptomic data. Therefore, there are no
77 existing toolkits that contain a complete workflow for microbiome data analysis and
78 interpretation (with or without a graphical user interface).

79 Here we present *animalcules*, an interactive analysis and visualization toolkit for
80 microbiome data. *animalcules* supports the importing of microbiome profiles in multiple formats
81 such as a species count table, an organizational taxonomic unit (OTU) counts table, or Biological
82 Observation Matrix (BIOM) format [20]. These formats could be generated from common
83 microbiome data sources and analytical tools including 16S rRNA, metagenomics, and
84 metatranscriptomic data. Once data is uploaded, *animalcules* provides a useful data summary and
85 filtering function where users can view and filter their dataset using sample metadata, microbial
86 prevalence or relative abundance. Filtering the data in this way can significantly reduce the time
87 spent performing preprocessing and downstream analysis tasks. For data visualizations, such as
88 relative abundance bar charts and 3D dimension reduction plots (PCA/PCoA/tSNE/UMAP),
89 *animalcules* supports interactive operations where users can check the sample/microbe
90 information on each data point and adjust the figure format as needed, which is helpful for
91 recognizing elements or data patterns when the sample size or number of microbes is large.
92 Aside from common diversity analysis, differential abundance analysis, and dimension
93 reduction, *animalcules* supports biomarker identification by training a logistic regression or
94 random forest model with cross-validated biomarker performance evaluation. *animalcules*

95 provides a graphical user interface (GUI) through R/Shiny, which can be used even by users
 96 without prior programming knowledge, while experienced programmers can choose the
 97 command-line based R package or a combination of both.

98

99 **Table 1.** Comparison of *animalcules* and other popular microbiome analysis tools.

	Biobakery	Vegan	mothur	qiime2	phyloseq	animalcules
Filtering and Data Summary	✓		✓	✓	✓	✓
Interactive Visualization				✓		✓
Dimension Reduction			✓	✓	✓	✓
Differential Abundance Analysis	✓		✓	✓	✓	✓
Diversity Analysis		✓	✓	✓	✓	✓
Support for 16S rRNA Data	✓	✓	✓	✓	✓	✓
Support for total RNA-seq Data	✓	✓				✓
Biomarker Identification				✓		✓
Interface and Command-Line				✓		✓
Language/Platform	Web	R	R	Python	R	R

100

101

Implementation

102 Data Structures and Software Design

103 All data handling tasks and functions in *animalcules* are based upon and work with the
 104 MultiAssayExperiment (MAE) data structure [21]. The MAE class is a standard data structure
 105 for multi-omics experiments with efficient data retrieval and manipulation methods that support
 106 the linkage of samples across multiple assays. The MAE object has three key components:
 107 colData (contains subject or cell line level metadata), ExperimentList (stores data for one or

108 more assays), and sampleMap (relates experiments and samples). In *animalcules*, three tables
109 (sample metadata table, microbe count table, and taxonomy table) as well as the mapping
110 relationship between them are stored in the MAE class. It ensures correct alignment of assays
111 and subjects, and provides coordinated subsetting of samples and features. Additionally, it is
112 easy to convert to or from a MAE object from the SummarizedExperiment class, which has been
113 applied in many Bioconductor packages, enabling smooth interaction between other tools [21].
114 One important advantage of applying the MAE class in the microbiome research field is its
115 extensible design supporting many multi-omics layers of data. Multi-omics is becoming a trend
116 in the field, e.g., studying host-microbe interactions by combining host gene expression data and
117 microbial abundance data. *animalcules* is the first software tool for microbiome analysis to
118 integrate the MAE object and takes advantage of its unique properties by allowing the user to
119 store microbial data, host transcriptomics, metabolomics, as well as taxonomy information
120 within the same object (currently *animalcules* only supports microbiome analysis, but the MAE
121 structure enables future development that can address these data types. Additionally, the MAE
122 enables integration with other tools that do manage these data types, e.g. host transcriptomics).
123 The MAE object can also store processed versions of various assays (e.g. dimension-reduced
124 data) which allows for efficient manipulation and analysis downstream. This approach advances
125 standard microbiome analysis and data sharing by efficiently integrating the various multi-omics
126 datasets required.

127 Lastly, because all of the data is integrated within a single R object, users can serialize
128 the data to a single file which can be used for further analysis or share with other researchers. For
129 example, after processing and analyzing their data through the Shiny application, users can
130 export their datasets in the form of a serialized MAE object file, which can be later uploaded to

131 Shiny or imported in R for further exploration through the *animalcules* command line functions
132 or other methods. Integrating the MAE object brings efficiency, scalability, and reproducibility
133 to microbiome analysis through *animalcules*.

134

135 **Installation and Usage**

136 *animalcules* requires R $\geq 4.0.0$ and can be installed through Github or Bioconductor.

137 After loading the *animalcules* library in R, users can choose between launching the R Shiny GUI
138 (via the `run_animalcules()` function), or using the available command-line functions directly. In
139 the GUI, users can choose from the following tabs: Upload (select an example dataset, upload a
140 new dataset, or load a previously uploaded dataset), Summary and Filter (understand the data
141 distribution and filter the data by microbial features or sample phenotypes), Abundance (relative
142 abundance bar charts, heatmaps, and individual microbes boxplots), Diversity (statistical tests
143 and boxplots for alpha diversity and beta diversity), Dimension Reduction (PCA, PCoA, tSNE,
144 and UMAP), Differential Abundance (microbial differential abundance between sample groups),
145 and Biomarker (identify predictive microbial biomarkers). Common R functions in the package
146 are summarized in **Table 2**. A detailed tutorial on how to use the command-line version of
147 *animalcules* for microbiome data analysis can be found at

148 <https://compbiomed.github.io/animalcules-docs/articles/animalcules.html>.

149

150

151 **Table 2.** Table of exported functions and their descriptions available through the *animalcules* R package.

Data and Interface

`run_animalcules()`

Initiates a local instance of the *animalcules* Shiny application

Data Summary and Manipulation

<code>filter_summary_bar_density()</code>	Visualize sample/microbe data with a bar plot (categorical) or density plot (continuous)
<code>filter_summary_pie_box()</code>	Visualize sample/microbe data with a pie chart (categorical) or box plot (continuous)
<code>filter_categorize()</code>	Convert continuous variables into a various number of factors
<code>counts_to_logcpm()</code>	Covert counts table to a log counts per million table
<code>counts_to_relabu()</code>	Covert counts table to a relative abundances table
<code>upsample_counts()</code>	Up-sample counts table to a higher taxon level
<code>find_taxonomy()</code>	Find taxonomy for unlimited ids
<code>find_taxon_mat()</code>	Find taxonomy information matrix for unlimited ids
<code>mae_pick_samples()</code>	Isolate or discard samples from a multi-assay experiment object
<code>mae_pick_organisms()</code>	Isolate or discard microbes from a multi-assay experiment object

Sample Level Visualization

<code>relabu_barplot()</code>	Generate stacked bar plots of sample and group level microbe relative abundances
<code>relabu_boxplot()</code>	Generate box plots comparing organism prevalence across groups of samples
<code>relabu_heatmap()</code>	Generate a sample by microbe heatmap of counts
<code>dimred_pca()</code>	Return a 2D/3D scatter plot for dimensionality reduction through PCA
<code>dimred_pcoa()</code>	Return a 2D/3D scatter plot for dimensionality reduction through PCoA
<code>dimred_umap()</code>	Return a 2D/3D scatter plot for dimensionality reduction through UMAP
<code>dimred_tsne()</code>	Return a 2D/3D scatter plot for dimensionality reduction through t-SNE

Alpha and Beta Diversity

<code>diversities()</code>	Return alpha diversity
----------------------------	------------------------

do_alpha_div_test()	Compute various statistical tests for alpha diversity
alpha_div_boxplot()	Generate box plots comparing alpha diversity across groups of samples
diversity_beta_test()	Compute various statistical tests for beta diversity
diversity_beta_boxplot()	Generate box plots comparing beta diversity across groups of samples
diversity_beta_heatmap()	Generate a heatmap comparing beta diversity across groups of samples

Differential Abundance Analysis

differential_abundance()	Performs differential abundance analysis across groups of samples
--------------------------	---

Biomarker Discovery

find_biomarker()	Identifies microbes as potential biomarkers for groups of samples
------------------	---

152

153 Data Upload and Output

154 *animalcules* offers multiple options for importing data into the GUI or working with the
 155 MAE object for command line analysis. These include simple tab-delimited OTU or count
 156 matrices, typically generated by other tools such as QIIME II [15] or PathoScope [22], or using a
 157 MAE object available in the user’s session or in a file from a previous session of *animalcules*.
 158 Regardless of how the data is imported, the assay/OTU data will be available in the “Assay
 159 Viewer” section of the Upload tab.

160

161 Six of the data importing options are described below:

- 162 1. *Count Table or OTU File (without taxonomy)*: This is the simplest option that enables the
 163 upload of an OTU or count table that has genomes/OTUs in the rows and samples in the
 164 columns. All functions and tools can be used for filtering, visualization, analysis of the

165 data, except the individual microbiomes or OTUs cannot be aggregated at different
166 levels.

167 2. *Count Table or OTU File (with taxonomy)*: This option provides an extension of the
168 previous but allows for associating the OTUs with taxonomy information and the
169 aggregation of microbes at different levels (e.g. species, genus, phylum, etc.). This
170 information can be provided as a separate table, with a row for each OTU in the table. In
171 addition, users can provide NCBI taxonomy IDs or NCBI accession numbers [23] and
172 *animalcules* will automatically generate the taxonomy table using the tools available in
173 the *taxize* R package [24]. The taxonomy table will be stored as a separate assay in the
174 MAE object, but will be linked to the rows of the OTU table through internal functions.
175 The taxonomy table will be available in the “Assay Viewer” section in the Upload tab.

176 3. *animalcules Object File*: Users can also directly upload a MAE object into the toolkit or
177 workflow. A MAE object could be generated from a previous *animalcules* session (stored
178 as an .rds file), converted from the output of any pre-processing pipeline, or generated
179 from some other source. This option allows for the efficient storage and re-upload of data
180 from a previous session, or enables the interaction between the command-line version
181 and the GUI version of *animalcules*. For example, users can conduct part of the analysis
182 in the GUI, save the results, and continue their analysis using command-line tools (inside
183 and outside of *animalcules*), and then re-upload the data to the GUI for further analysis or
184 visualization. This feature enables compatibility and interactivity that is not available in
185 other microbiome GUI or command-line toolkits.

186 4. *Pathoscope Output Files*: *animalcules* enables the direct upload of files generated from
187 the *PathoScope* pipeline [22]. These files are generally single tab-delimited tables for

188 each sample in the dataset, and contain NCBI taxonomy IDs for individual microbes.
189 *animalcules* combines and converts these files into a MAE object, and uses the *taxize*
190 package to generate the taxonomy table.

191 5. *BIOM Format File*: The standard BIological Observation Matrix (BIOM) format is a
192 commonly used format for representing samples by observation contingency tables [20].
193 The BIOM format is commonly used by QIIME II pipeline tools. We used the
194 *biomformat* R package [25] for uploading a BIOM file into *animalcules* as well as
195 outputting a BIOM file from *animalcules*. This enables interactivity between *animalcules*
196 and other microbiome analysis tools such as QIIME II.

197 6. *Example Data*: In *animalcules*, we have three pre-defined example datasets, including a
198 simulated dataset, a Tuberculosis 16S rRNA profiling dataset, and an Asthma
199 metatranscriptomic dataset. These example datasets allow users to try all the features and
200 functions in *animalcules* before users upload their own data, making it easy to learn how
201 to use *animalcules* and understand what analyses they can perform.

202

203 **Data Filtering and Summary**

204 The *animalcules* Shiny interface provides summary statistics to help users efficiently and
205 effectively assess data quality and filter low-quality microbes and samples. Users can visualize
206 the total number of reads for each organism through a scatter and density plot and filter
207 organisms based on average read number, relative abundance, or prevalence. Additionally, users
208 can visualize sample covariates through a pie and bar plot for categorical covariates or a scatter
209 and density plot for continuous covariates (**Figure 1**). Samples can be filtered based on one or
210 more covariates. Finally, users have the option to discard specific samples and/or organisms. As

211 samples and organisms are removed through any of the filtering methods, summary statistics and
212 plots are automatically refreshed to display any changes that may occur. If changes have been
213 made, users may download the modified data for later use. Visualizations of sample and microbe
214 data before and after filtering are generated with `animalcules::filter_summary_bar_density()` and
215 `animalcules::filter_summary_pie_box()` functions. For users who wish to inspect their data
216 before or after filtering, *animalcules* enables users to view and download five types of assays
217 generated including a count table, relative abundance table, logCPM table, taxonomy table and
218 annotation table. In addition, these tables can also be accessed directly from the MAE object
219 through standard R command line tools.

220

221 **Data Visualization**

222 A typical analysis involves visualization of microbe abundances across samples or groups
223 of samples. *animalcules* implements three common types of visualization plots including stacked
224 bar plots, heatmaps, and box plots. The stacked bar plots, generated with
225 `animalcules::relabu_barplot()` are used to visualize the relative abundance of microbes at a given
226 taxonomic level in each sample, represented as a single bar (**Figure 2**). Bars can be color-labeled
227 by one or more sample attributes and samples can also be aggregated by these attributes via
228 summing microbe abundances within groups. This is an efficient way for researchers to identify
229 sample- or group-level patterns at various taxonomic levels. Users also have the option to sort
230 the bars by sample attributes or by the abundance of one or more organisms. There is also a
231 convenient method for isolating or removing samples. With this tool, users can quickly scan
232 through different combinations of sample attributes and taxon levels for differential abundance in

233 one or more groups, outliers in terms of community profile, as well as sample clusters not
234 represented by known attributes.

235 Alternatively, users can investigate these questions through the heatmap visualization,
236 which represents a sample-by-organisms matrix that can be visualized at different taxonomic
237 levels. Many of the previously mentioned options are also compatible with the heatmap such as
238 color-labeling samples, sorting matrix rows by attributes or organisms and isolating or discarding
239 organisms and samples. After identifying potential differentially abundant microbes, users can
240 use the boxplot visualization to directly compare the abundance of one or more organisms
241 between categorical attributes. Organisms can be chosen from a given taxonomic level and
242 abundance can be represented as either counts, logCPM, or relative abundance. This plot can
243 also be generated in the command line using the `animalcules::relabu_heatmap()` function.

244

245 **Diversity Analysis**

246 Alpha diversity, which describes the richness and evenness of a microbial community, is
247 a vital indicator and measurement in microbiome analysis [26]. *animalcules* provides an
248 interactive box plot comparison of alpha diversity between selected groups of samples. Both
249 taxonomy levels and alpha diversity metrics (e.g. Shannon, Gini Simpson, Inverse Simpson) can
250 be changed and diversity can be calculated at multiple taxonomic levels [27, 28]. Alpha diversity
251 values for each sample can be output into the MAE object or as separate tables or files. Users can
252 also conduct alpha diversity statistical tests including Wilcoxon rank-sum test, T-test and
253 Kruskal-Wallis test [29, 30]. The alpha diversity boxplot as well as the statistical tests could be
254 generated in the command line using the `animalcules::alpha_div_boxplot()` function and
255 `animalcules::do_alpha_div_test()` function.

256 On the other hand, one can use distances between each microbial community sample, or
257 so-called beta diversity, as another key metric to consider for each analysis. Users can plot the
258 beta diversity heatmap by selecting different beta diversity dissimilarity metrics including Bray-
259 Curtis [31] or Jaccard index [32]. Users can also conduct beta diversity statistical testing between
260 groups including PERMANOVA [33], Wilcoxon rank-sum test, or Kruskal-Wallis test (**Figure**
261 **3**). The beta diversity comparison boxplot as well as the statistical tests can be generated in the
262 command line using the `animalcules::diversity_beta_boxplot()` function and
263 `animalcules::diversity_beta_test()` function.

264

265 **Dimension Reduction**

266 A crucial step in any data analysis workflow is to visualize and summarize highly
267 variable data in a lower-dimensional space (**Figure 4**). In *animalcules*, we implement four
268 commonly used dimensionality reduction techniques including Principal Components Analysis
269 (PCA), Principal Coordinates Analysis (PCoA), t-Distributed Stochastic Neighbor Embedding (t-
270 SNE) and Uniform Manifold Approximation and Projection (UMAP) [34–37] Both PCA and
271 PCoA project samples onto a new set of axes whereby a maximum amount of variation is
272 explained by the first, second, and third axes while t-SNE and UMAP are non-linear methods for
273 mapping data to a lower-dimensional embedding. Dimension reduction values for the dataset can
274 be output into the MAE object or as separate tables or files.

275 The original data used in each dimensionality reduction method can be either counts,
276 logCPM, or relative abundance, and can be visualized using a 2D or 3D (if two dimensions of
277 explained variance are inadequate) scatter plot. Data points can be colored by continuous sample
278 attributes and shaped by categorical attributes. With multiple dimensionality reduction

279 techniques and methods for data normalization, users can rapidly visualize the global and local
280 structure of their data, identify clustering patterns across one or more conditions, as well as
281 detect sample outliers. Dimensionality reduction can also be carried out in the command line
282 using the `animalcules::dimred_pca()` function for PCA, `animalcules::dimred_pcoa()` function for
283 PCoA, `animalcules::dimred_tsne()` function for t-SNE, and `animalcules::dimred_umap()` function
284 for UMAP.

285

286 **Differential Abundance Analysis**

287 There are many available tools for differential abundance estimation and inference. For
288 example, GLM (Generalized Linear Model) based methods including DESeq2 [38], edgeR [39],
289 and limma [40] model count based microbiome data or gene expression data by a negative
290 binomial distribution (DESeq2 and edgeR) or using log-counts (per million) and a Gaussian
291 distribution (limma) assumption. Core microbes that have different abundance in different
292 groups could be identified. Here in *animalcules*, we provide a DESeq2-based differential
293 abundance analysis (**Figure 5**). With the command-line function
294 `animalcules::differential_abundance()`, which by default uses the “DESeq2” method. Users can
295 choose the target variable, covariate variable, taxonomy level, minimum count cut-off, and an
296 adjusted p-value threshold. The analysis report will output not only the adjusted p-value and
297 log₂-fold-change of the microbes, but also the percentage, prevalence, and the group size-
298 adjusted fold change. Besides using DESeq2, in *animalcules* we also support differential
299 abundance analysis with limma, which requires users to specify in the command-line function as:
300 `animalcules::differential_abundance(method='limma')`.

301

302

303 **Biomarker Identification**

304 One unique feature of *animalcules* is the biomarker identification module. Users can
305 choose either logistic regression [41] or random forest [42] classification model to identify a
306 microbe biomarker. The feature importance score for each microbe will be provided (**Figure 6**),
307 in addition to AUC values and average cross-validation ROC curves for evaluating biomarker
308 prediction performance. The biomarker identification can also be conducted by the command-
309 line function `animalcules::find_biomarker()`.

310

311 **Results**

312 To illustrate the utility of *animalcules*, we include two example analyses using the pre-
313 loaded datasets packaged within *animalcules*; the first being an asthma metatranscriptomic
314 dataset, and the second a TB 16S rRNA dataset. For brevity, we do not explore all *animalcules*
315 functions in each analysis, but focus on and expand the relevant analyses for the scientific
316 questions for each example. Both analyses could be reproduced within the *animalcules* Rshiny
317 app by using the corresponding example datasets.

318

319 **Example 1: Asthma nasal swabs metatranscriptomic dataset**

320 The asthma metagenomic shotgun RNA sequencing dataset was generated from
321 participants of the AsthMap (Asthma Severity Modifying Polymorphisms) project and originally
322 reported in a research article characterizing asthma-associated microbial communities [43]. It
323 contains 14 total samples of nasal epithelial cells collected from 8 children and adolescents with
324 asthma and 6 healthy controls. The goal of this study was to further understand the relationship
325 between the microbiome and host inflammatory processes in asthmatic children.

326 To characterize the relationship between microbial communities and asthma, species-
327 level abundances were visualized by plotting the group-wise relative abundance of microbes
328 across asthma and control subjects. This plot can be generated with the
329 `animalcules::relabu_barplot()` function as well as under the *Abundance* tab of the Shiny
330 application. It is clear that *Moraxella catarrhalis* is overrepresented in asthmatics versus
331 controls, which was a major discovery in the original publication. This microbe - which is known
332 to cause infections in the respiratory system - could serve as a biomarker for early disease
333 detection, severity of disease, or potential for exacerbation. In addition, other dysbiosis to the
334 airway microbiome included differences in other genera such as *Corynebacterium aurimucosum*,
335 which is underrepresented in asthmatics versus controls (**Figure 7**).

336 To further investigate the overrepresentation and underrepresentation of *M. catarrhalis*
337 and *C. aurimucosum* respectively in asthmatics versus controls, we use boxplots, generated with
338 the `animalcules::boxplot()` function, to visualize the relative abundance in each group and to get
339 a better sense of the mean and variance of the distribution across samples. These plots confirm
340 the previous results by showing a drastic difference in abundance (**Figure 8**). Furthermore, we
341 employed DESeq2 to conduct a differential abundance analysis of microbe species for asthmatics
342 versus controls. This analysis shows that *M. catarrhalis* is significantly ($q = 1.78e-3$)
343 overrepresented ($\text{Log}_2\text{FC} = 5.9$) in asthmatics. It also shows that *C. aurimucosum* is
344 overrepresented ($\text{Log}_2\text{FC} = 2.66$) in controls, however not at a statistically significant level ($q =$
345 0.236). This table was generated with the `animalcules::differential_abundance()` feature.

346 Through the *animalcules* interface, we were able to rapidly visualize sample- and group-
347 level microbial communities between asthmatic and control samples and test for over- and

348 underrepresented organisms in asthmatics, identifying *M. catarrhalis* and *C. aurimucosum*
349 respectively.

350

351 **Example 2: Tuberculosis 16S rRNA profiling dataset**

352 This 16S rRNA TB dataset comes from a pilot TB study containing 12 subjects, 30
353 respiratory tract samples and 417 species of microbe [44] Among the 12 subjects, there are 6
354 patients with pulmonary tuberculosis and 6 healthy control individuals. Sample tissue type
355 includes sputum, oropharynx, and nasal respiratory tract. The goal of this study is to learn the
356 microbial community differences in the respiratory tract between healthy and TB patients, and to
357 evaluate the sample/tissue types that were most effective for exploring differences between the
358 microbiome of TB samples vs. controls.

359 We first conducted an overall assessment of the data, focusing on how the microbial
360 taxonomy affects sample variables such as disease status. We used the barplot function in
361 *animalcules* to visualize the taxonomic profile for each sample, colored by any annotation
362 variable (here we used disease information, where dark blue represents control and yellow
363 represents TB samples). In **Figure 9** we display the genus and phylum level abundances.

364 From the taxonomy barplot, we find different patterns that exist in TB vs. control
365 samples. At the genus level (**Figure 9A**), *Streptococcus* appears to have a higher relative
366 abundance in TB samples compared to the control samples. In the phylum level (**Figure 9B**), we
367 found *Firmicutes* to be more abundant in TB samples. Both figures were generated using
368 command-line function `animalcules::rebalu_barplot()`.

369 To obtain a quantitative understanding of the ecological diversity difference between TB
370 and control samples, we compared the alpha and beta diversity of our samples. For alpha

371 diversity, we compared the Shannon index in TB vs. control samples (see **Figure 10.A**).
372 *animalcules* automatically conducted a non-parametric Wilcoxon rank-sum test and a parametric
373 Welch two-sample T-test on these diversity measures. Here the Wilcoxon rank-sum test gives a
374 p-value of 0.0060, while Welch two-sample T-test gives a p-value 0.0077, thus showing a
375 significant difference in diversity between TB and control groups. From the boxplot, we observe
376 that the alpha diversity is higher in the control group. The alpha diversity boxplot was generated
377 by `animalcules::alpha_div_boxplot()`, and the statistical test was generated by
378 `animalcules::do_alpha_div_test()`.

379 As for comparing beta diversity, we plotted the Bray-Curtis distance to compare: within
380 the TB group, within the control group, and between the TB/control (**Figure 10.B**). The average
381 distance between the two groups is higher than two separate within-group distances, meaning
382 both TB samples and control samples are more similar to themselves. Furthermore, we
383 conducted a PERMANOVA test between the two groups, and it shows a significant difference
384 with a p-value of 0.003. The beta diversity comparison boxplot was generated by
385 `animalcules::diversity_beta_boxplot()`, and the PERMANOVA test was generated by
386 `animalcules::diversity_beta_test()`.

387 After exploring this TB dataset in terms of relative abundance and diversity analysis, we
388 were certain that there is a significant difference between TB and control groups in the
389 microbiome. Here, with the biomarker function in *animalcules*, we were able to build a microbial
390 biomarker that could help us predict TB status. Using a logistic regression model, 3-fold cross-
391 validation (CV), the number of CV repeats as 3, and top biomarker proportion as 0.05, we
392 identified an 8-genus biomarker for TB classification. Then we tested the biomarker performance
393 by using only the 8-genus biomarker for cross-validation, and the prediction performance ROC is

394 displayed in (**Figure 11**). We used `animalcules::find_biomarker()` to identify the biomarker, plot
395 the feature importance score barplot and the ROC curve. Here we have a very high AUC =
396 0.913, thus providing evidence that the microbiome could serve as a biomarker for TB
397 prediction, and our biomarker has a differentiating power between TB vs. healthy controls. This
398 result suggests that further evaluation of microbial biomarkers for TB is warranted. Previously,
399 people have been using transcriptomic biomarkers for TB diagnosis [45], our new finding of
400 using microbes as diagnosis biomarker can lead to a potential total RNA-seq directed TB disease
401 biomarker that involves both host transcriptome gene expression as well as microbial abundance,
402 which has the potential of higher accuracy for TB diagnosis because it considers both host and
403 microbial side, or even the host-microbe interaction in TB.

404 To summarize, with the help of *animalcules*, we explored and compared the microbial
405 community difference between TB and control samples. Our analysis shows that the microbial
406 community structure in the control group is more diverse and evenly distributed compared to the
407 one in the TB group. Also, the TB group as well as control group each has a specific microbial
408 composition that is shared within the group. Finally, we identified a subset of microbes that
409 indicate its differentiating power between TB vs. control samples, which can be used as a new
410 TB disease biomarker.

411

412

413

414

Discussion

415

416

A fundamental characteristic of *animalcules* is its seamless interaction with the user through dynamic visualization tools. This design logic is rooted in the fact that researchers in

417 microbiology must analyze their data at multiple levels (taxonomy) and multiple scales
418 (normalization), thus data visualization and analysis become complicated without an organized
419 analysis framework and workflow. *animalcules* solves this problem by providing a platform for
420 interactively exploring large datasets, making it easier for users to identify patterns inherent in
421 the dataset through appropriate analysis methods. Key analysis methods allow users to
422 investigate differences in grouped relative abundance patterns between multiple sample groups
423 in the phylum level, check the top abundant species in one specific sample group, or to check the
424 individual sample-wise microbiome composition at different taxon levels. Patterns identified can
425 be further tested through alpha/beta diversity statistical tests, differential abundance analysis, as
426 well as biomarker identification. Furthermore, *animalcules* utilizes the MAE object, an efficient
427 data structure for multi-omic sequencing data, which could be extended in the future to
428 incorporate host sequencing assays, and enable compact methods for analyzing host-microbe
429 interactions.

430

431 **Conclusion**

432 In this report, we present *animalcules*, an open-source R package and Shiny application
433 dedicated to microbiome analysis for both 16S rRNA and shotgun sequencing (metagenomics
434 and metatranscriptomics) data. We incorporate leading and novel methods in an efficient
435 framework for researchers to characterize and understand the microbial community structure in
436 their data, leading to valuable insights into the connection between the microbial community and
437 phenotypes of interest.

438

439 **Availability and requirements**

440 Project name: *animalcules*
441 Project home page: <https://compbiomed.github.io/animalcules-docs/>
442
443 Operating system(s): Linux, OS X, Windows
444 Programming language: R
445 License: GNU GPLv3

446

447 **Declarations**

448 **Acknowledgments**

449 This research was supported by grants from the NIH R01 GM127430 and U01CA220413. The
450 authors thank Lucas Schiffer (Boston University) as well as other members of the Johnson lab
451 for helpful comments.

452 **Authors' contributions**

453 YZ, AF, SM contributed to the software development. YZ, AF, TF, SM and WEJ contributed to
454 the writing and review of the manuscript before submission for publication. All authors read and
455 approved the final manuscript.

456 **Availability of data and materials**

457 animalcules is freely available on GitHub at <https://github.com/compbiomed/animalcules> or
458 Bioconductor at <https://bioconductor.org/packages/release/bioc/html/animalcules.html> and is
459 accompanied by comprehensive documentation and tutorials at

460 <https://compbiomed.github.io/animalcules-docs/>

461 **Competing interests**

462 None declared.

463 **Consent for publication**

464 Not applicable.

465

466 **Ethics approval and consent to participate**

467 Not applicable.

468

469

470

References

471

472 1. Gilbert J, Blaser M, Caporaso J, medicine JJ-N, 2018 undefined. Current understanding of the human

473 microbiome. *nature.com*. <https://www.nature.com/articles/nm.4517>. Accessed 6 May 2020.

474 2. Turnbaugh PJ, Ley RE, Hamady M, Fraser-Liggett CM, Knight R, Gordon JI. The Human Microbiome Project.

475 *Nature*. 2007;449:804–10. doi:10.1038/nature06244.

476 3. Lloyd-Price J, Abu-Ali G, Huttenhower C. The healthy human microbiome. *Genome Medicine*. 2016;8:1–11.

477 4. David LA, Materna AC, Friedman J, Campos-Baptista MI, Blackburn MC, Perrotta A, et al. Host lifestyle affects

478 human microbiota on daily timescales. *Genome Biol*. 2014;15:R89. doi:10.1186/gb-2014-15-7-r89.

479 5. David LA, Maurice CF, Carmody RN, Gootenberg DB, Button JE, Wolfe BE, et al. Diet rapidly and reproducibly

480 alters the human gut microbiome. *Nature*. 2014;505:559–63. doi:10.1038/nature12820.

481 6. Hartstra A V., Bouter KEC, Bäckhed F, Nieuwdorp M. Insights Into the Role of the Microbiome in Obesity and

482 Type 2 Diabetes. *Diabetes Care*. 2015;38:159–65. doi:10.2337/DC14-0769.

483 7. Schwabe RF, Jobin C. The microbiome and cancer. *Nat Rev Cancer*. 2013;13:800–12. doi:10.1038/nrc3610.

484 8. Knight R, Vrbanac A, Taylor BC, Aksenov A, Callewaert C, Debelius J, et al. Best practices for analysing

485 microbiomes. *Nature Reviews Microbiology*. 2018;16:410–22.

486 9. Saulnier DM, Riehle K, Mistretta TA, Diaz MA, Mandal D, Raza S, et al. Gastrointestinal microbiome signatures

487 of pediatric patients with irritable bowel syndrome. *Gastroenterology*. 2011;141:1782–91.

488 10. Qu K, Guo F, Liu X, Lin Y, microbiology QZ-F in, 2019 undefined. Application of machine learning in

489 microbiology. *ncbi.nlm.nih.gov*. <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC6482238/>. Accessed 6 May 2020.

490 11. Zhou YH, Gallins P. A review and tutorial of machine learning methods for microbiome host trait prediction.

491 *Front Genet*. 2019;10 JUN.

- 492 12. Arango-Argoty G, Garner E, Pruden A, Heath LS, Vikesland P, Zhang L. DeepARG: A deep learning approach
493 for predicting antibiotic resistance genes from metagenomic data. *Microbiome*. 2018;6.
- 494 13. Reiman D, Metwally A, Dai Y. Using convolutional neural networks to explore the microbiome. In: Proceedings
495 of the Annual International Conference of the IEEE Engineering in Medicine and Biology Society, EMBS. Institute
496 of Electrical and Electronics Engineers Inc.; 2017. p. 4269–72.
- 497 14. Allaband C, McDonald D, Vázquez-Baeza Y, Minich JJ, Tripathi A, Brenner DA, et al. *Microbiome 101:*
498 *Studying, Analyzing, and Interpreting Gut Microbiome Data for Clinicians*. *Clinical Gastroenterology and*
499 *Hepatology*. 2019;17:218–30.
- 500 15. Bolyen E, Rideout JR, Dillon MR, Bokulich NA, Abnet CC, Al-Ghalith GA, et al. Reproducible, interactive,
501 scalable and extensible microbiome data science using QIIME 2. *Nature Biotechnology*. 2019;37:852–7.
- 502 16. Schloss PD, Westcott SL, Ryabin T, Hall JR, Hartmann M, Hollister EB, et al. Introducing mothur: Open-
503 Source, Platform-Independent, Community-Supported Software for Describing and Comparing Microbial
504 Communities. *Appl Environ Microbiol*. 2009;75:7537–41. doi:10.1128/AEM.01541-09.
- 505 17. Dixon P. VEGAN, a package of R functions for community ecology. *Journal of Vegetation Science*.
506 2003;14:927–30.
- 507 18. Mciver LJ, Abu-Ali G, Franzosa EA, Schwager R, Morgan XC, Waldron L, et al. bioBakery: a meta’omic
508 analysis environment. doi:10.1093/bioinformatics/btx754.
- 509 19. Mcmurdie PJ, Holmes S. phyloseq: An R Package for Reproducible Interactive Analysis and Graphics of
510 Microbiome Census Data. *journals.plos.org*. 2013;8. doi:10.1371/journal.pone.0061217.
- 511 20. McDonald D, Clemente J, ... JK-, 2012 undefined. The Biological Observation Matrix (BIOM) format or: how
512 I learned to stop worrying and love the ome-ome. *academic.oup.com*. [https://academic.oup.com/gigascience/article-](https://academic.oup.com/gigascience/article-abstract/1/1/2047-217X-1-7/2656152)
513 [abstract/1/1/2047-217X-1-7/2656152](https://academic.oup.com/gigascience/article-abstract/1/1/2047-217X-1-7/2656152). Accessed 6 May 2020.
- 514 21. Ramos M, Schiffer L, Re A, Azhar R, Basunia A, Rodriguez C, et al. Software for the Integration of Multiomics
515 Experiments in Bioconductor. *Cancer Res*. 2017;77:e39–42. doi:10.1158/0008-5472.CAN-17-0344.
- 516 22. Hong C, Manimaran S, Shen Y, Perez-Rogers JF, Byrd AL, Castro-Nallar E, et al. PathoScope 2.0: a complete
517 computational framework for strain identification in environmental or clinical sequencing samples. *Microbiome*.
518 2014;2:33. doi:10.1186/2049-2618-2-33.
- 519 23. Federhen S. The NCBI Taxonomy database. doi:10.1093/nar/gkr1178.

- 520 24. Chamberlain SA, Szöcs E. Taxize: Taxonomic search and retrieval in R. *F1000Research*. 2013;2.
- 521 25. McMurdie P, 1.0 JP-P version, 2016 undefined. biomformat: an interface package for the BIOM file format.
- 522 26. Whittaker RH. EVOLUTION AND MEASUREMENT OF SPECIES DIVERSITY. *Taxon*. 1972;21:213–51.
- 523 27. Spellerberg IF, Fedor PJ. A tribute to Claude-Shannon (1916-2001) and a plea for more rigorous use of species
524 richness, species diversity and the “Shannon-Wiener” Index. *Glob Ecol Biogeogr*. 2003;12:177–9.
- 525 28. Jost L. *Entropy and diversity*. *Oikos*. 2006;113:363–75. doi:10.1111/j.2006.0030-1299.14714.x.
- 526 29. Kruskal WH, Wallis WA. Use of Ranks in One-Criterion Variance Analysis. *J Am Stat Assoc*. 1952;47:583–
527 621.
- 528 30. Mann H, statistics DW-T *annals of mathematical*, 1947 undefined. On a test of whether one of two random
529 variables is stochastically larger than the other. *JSTOR*.
- 530 [https://www.jstor.org/stable/2236101?casa_token=47wpm5LL1p8AAAAA:8eTiq-60-Km-](https://www.jstor.org/stable/2236101?casa_token=47wpm5LL1p8AAAAA:8eTiq-60-Km-02twhkibGHIq68tNENNLK06hpehy3dGEApYMZ6sIWvb8qn3M8TgHgZ_sZF-KPJ17wluojPglbzIXnoOiy17J17_3V7w1C2Imi8HQdnc)
531 [02twhkibGHIq68tNENNLK06hpehy3dGEApYMZ6sIWvb8qn3M8TgHgZ_sZF-](https://www.jstor.org/stable/2236101?casa_token=47wpm5LL1p8AAAAA:8eTiq-60-Km-02twhkibGHIq68tNENNLK06hpehy3dGEApYMZ6sIWvb8qn3M8TgHgZ_sZF-KPJ17wluojPglbzIXnoOiy17J17_3V7w1C2Imi8HQdnc)
532 [KPJ17wluojPglbzIXnoOiy17J17_3V7w1C2Imi8HQdnc](https://www.jstor.org/stable/2236101?casa_token=47wpm5LL1p8AAAAA:8eTiq-60-Km-02twhkibGHIq68tNENNLK06hpehy3dGEApYMZ6sIWvb8qn3M8TgHgZ_sZF-KPJ17wluojPglbzIXnoOiy17J17_3V7w1C2Imi8HQdnc). Accessed 6 May 2020.
- 533 31. Bray JR, Curtis JT. An Ordination of the Upland Forest Communities of Southern Wisconsin. *Ecol Monogr*.
534 1957;27:325–49.
- 535 32. Jaccard P. THE DISTRIBUTION OF THE FLORA IN THE ALPINE ZONE. *New Phytol*. 1912;11:37–50.
- 536 33. Anderson MJ. Permutational Multivariate Analysis of Variance (PERMANOVA). In: *Wiley StatsRef: Statistics*
537 *Reference Online*. John Wiley & Sons, Ltd; 2017. p. 1–15.
- 538 34. Pearson K. LIII. On lines and planes of closest fit to systems of points in space . London, Edinburgh, Dublin
539 *Philos Mag J Sci*. 1901;2:559–72.
- 540 35. Borg I, Groenen P. *Modern multidimensional scaling: Theory and applications*. 2005.
- 541 <https://books.google.com/books?hl=en&lr=&id=duTODldZzRcC&oi=fnd&pg=PR7&dq=Modern+multidimensiona>
542 [l+scaling:+Theory+and+applications&ots=SE4u8pOIuU&sig=X-YTeJ17yQgULfnXLG0y5oVRvF8](https://books.google.com/books?hl=en&lr=&id=duTODldZzRcC&oi=fnd&pg=PR7&dq=Modern+multidimensiona). Accessed 6
543 May 2020.
- 544 36. Van Der Maaten L, Hinton G. *Visualizing Data using t-SNE*. 2008.
- 545 <http://www.jmlr.org/papers/volume9/vandermaaten08a/vandermaaten08a.pdf>. Accessed 2 Feb 2019.
- 546 37. McInnes L, Healy J, Melville J. *UMAP: Uniform Manifold Approximation and Projection for Dimension*
547 *Reduction*. 2018. <http://arxiv.org/abs/1802.03426>. Accessed 6 May 2020.

548 38. Love MI, Huber W, Anders S. Moderated estimation of fold change and dispersion for RNA-seq data with
549 DESeq2. *Genome Biol.* 2014;15:550. doi:10.1186/s13059-014-0550-8.

550 39. Robinson M, McCarthy D, Bioinformatics GS-, 2010 undefined. edgeR: a Bioconductor package for differential
551 expression analysis of digital gene expression data. *academic.oup.com*.
552 <https://academic.oup.com/bioinformatics/article-abstract/26/1/139/182458>. Accessed 6 May 2020.

553 40. Smyth GK. limma: Linear Models for Microarray Data. In: *Bioinformatics and Computational Biology Solutions*
554 *Using R and Bioconductor*. Springer-Verlag; 2005. p. 397–420.

555 41. Jr DH, Lemeshow S, Sturdivant R. *Applied logistic regression*. 2013.
556 <https://books.google.com/books?hl=en&lr=&id=64JYAwwAAQBAJ&oi=fnd&pg=PR13&dq=logistic+regression&ots=DsjS909nlN&sig=iYEn6STGF9Q4T8dvq35h4032IDI>. Accessed 6 May 2020.

557
558 42. Breiman L. Random forests. *Mach Learn.* 2001;45:5–32.

559 43. Castro-Nallar E, Bendall ML, Pérez-Losada M, Sabuncyan S, Severance EG, Dickerson FB, et al. Composition,
560 taxonomy and functional diversity of the oropharynx microbiome in individuals with schizophrenia and controls.
561 *PeerJ.* 2015;3:e1140. doi:10.7717/peerj.1140.

562 44. Botero LE, Delgado-Serrano L, Cepeda ML, Bustos JR, Anzola JM, Del Portillo P, et al. Respiratory tract
563 clinical sample selection for microbiota analysis in patients with pulmonary tuberculosis. *Microbiome.* 2014;2.

564 45. Leong S, Zhao Y, Joseph NM, Hochberg NS, Sarkar S, Pleskunas J, et al. Existing blood transcriptional
565 classifiers accurately discriminate active tuberculosis from latent infection in individuals from south India.
566 *Tuberculosis.* 2018;109:41–51.

567
568
569
570
571
572
573
574
575

576

577

Figure and table titles and legend

578

579 **Table 1.** Comparison of *animalcules* and other popular microbiome analysis tools.

580 **Table 2.** Table of exported functions and their descriptions available through the *animalcules* R package.

581

582 **Figure 1. *animalcules* Data Filtering and Summary tab.** In the right panel, a table of data summary metrics, a
583 scatter/boxplot, and a density plot are displayed for continuous variables. For categorical variables,
584 *animalcules* will automatically identify and show the pie and bar plots instead.

585 **Figure 2. *animalcules* Abundance Tab.** In the subtab panel, users can select between a bar plot, heatmap, or box
586 plot. In the bar plot setting, in the left panel, users can select the color by variable, taxonomy level, and sort by
587 option. In the right panel, *animalcules* will show an interactive plot where users can mouse-hover to check the
588 identity of any color bar shown in the plot.

589 **Figure 3. *animalcules* Diversity tab.** In the subtab panel, the user could select between alpha diversity analysis and
590 beta diversity analysis. Here in the beta diversity analysis, the right panel controls what statistical test to use,
591 which condition to test on, and show statistical test results in a table as well as a boxplot.

592 **Figure 4. *animalcules* Dimension Reduction tab.** In the subtab panel, the user could select between PCA, PCoA, t-
593 SNE, and UMAP. Here in the PCA subtab, the user could choose the taxonomy level, color by variable, and in
594 advanced options, the user could also specify up to three PCs for visualization, shape by variable, and which
595 data type to use.

596 **Figure 5. *animalcules* Differential Abundance tab.** In the subtab panel, users select between DESeq2 and limma. In
597 the left panel, users specify taxonomy level, target condition, covariate variables, count cut-off, and adjusted
598 p-value threshold. In the right panel, a detailed differential abundance result table is shown.

599 **Figure 6. *animalcules* Biomarker tab.** In the left panel, users select taxonomy level, and target condition. In the
600 advanced options: number of cross-validations folds, number of cross-validations repeats, biomarker
601 proportion, and classification model. In the right panel, *animalcules* will show the biomarker list, importance
602 plot, and ROC plot.

603 **Figure 7. Relative abundance of microbial species bar plot.** A stacked bar plot representing the group-wise

604 relative abundance of microbial species in asthmatics (purple) and healthy controls (yellow).

605 **Figure 8. Relative abundance boxplot for differentially abundant species.** *Left.* A boxplot of relative abundance
606 of *M. catarrhalis* in asthmatics (green) and healthy controls (blue). *Right.* A boxplot of relative abundance of
607 *C. aurimucosum* in asthmatics (green) and healthy controls (blue).

608 **Figure 9. Sample-wise relative abundance bar plot.** A stacked bar plot representing the sample-wise relative
609 abundance of microbial species in TB (yellow) and healthy controls (blue). Figure A is the genus level and
610 figure B is the phylum level.

611 **Figure 10. TB example dataset diversity analysis.** **A.** Alpha diversity boxplot between control (red) group and TB
612 (blue) group. **B.** Beta diversity boxplot within the TB (blue), within the control(orange), and between
613 TB/control group (green).

614 **Figure 11. Biomarker ROC curve.** ROC shows AUC and cross-validation prediction performance of the identified
615 biomarker.

616

617

618

619

animalcules manuscript

figures

Filter Categorize Assay Dashboard

Filter By
 Metadata

Select a Condition
 AGE

Include
 1 99

Filter

Reset

Download Animalcules File

Download Biom File

Advanced Options

Summary Statistics

Number of Samples	50
Number of Covariates	4
Number of Organisms	100
Sample Mean Counts	24041
Sample Median Counts	4121
Organism Mean Counts	12021
Organism Median Counts	203

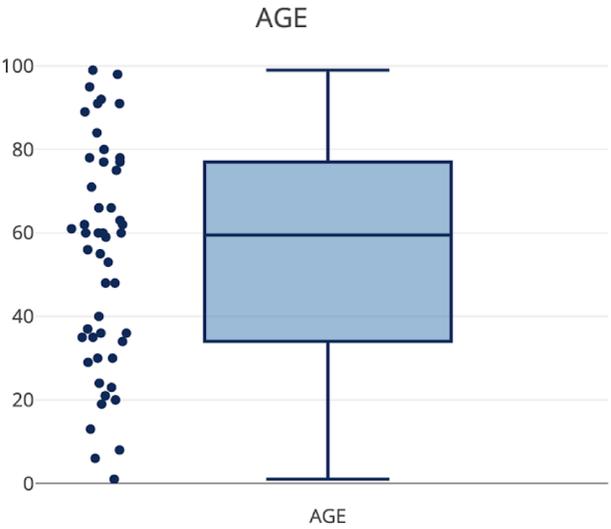
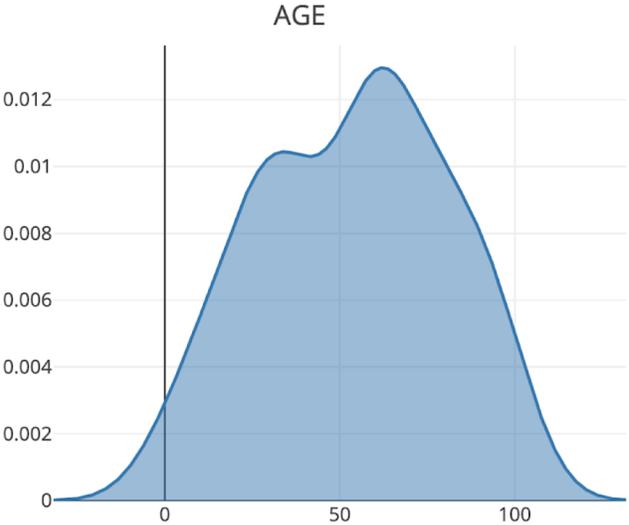


Figure 1. animalcules Data Filtering and Summary tab. In the right panel, a table of data summary metrics, a scatter/boxplot, and a density plot are displayed for continuous variables. For categorical variables, *animalcules* will automatically identify and show the pie and bar plots instead.

Color Samples by Condition

Group Samples by Condition

Tax Level

Sort By

No Sorting

Conditions

Organisms

Advanced Options

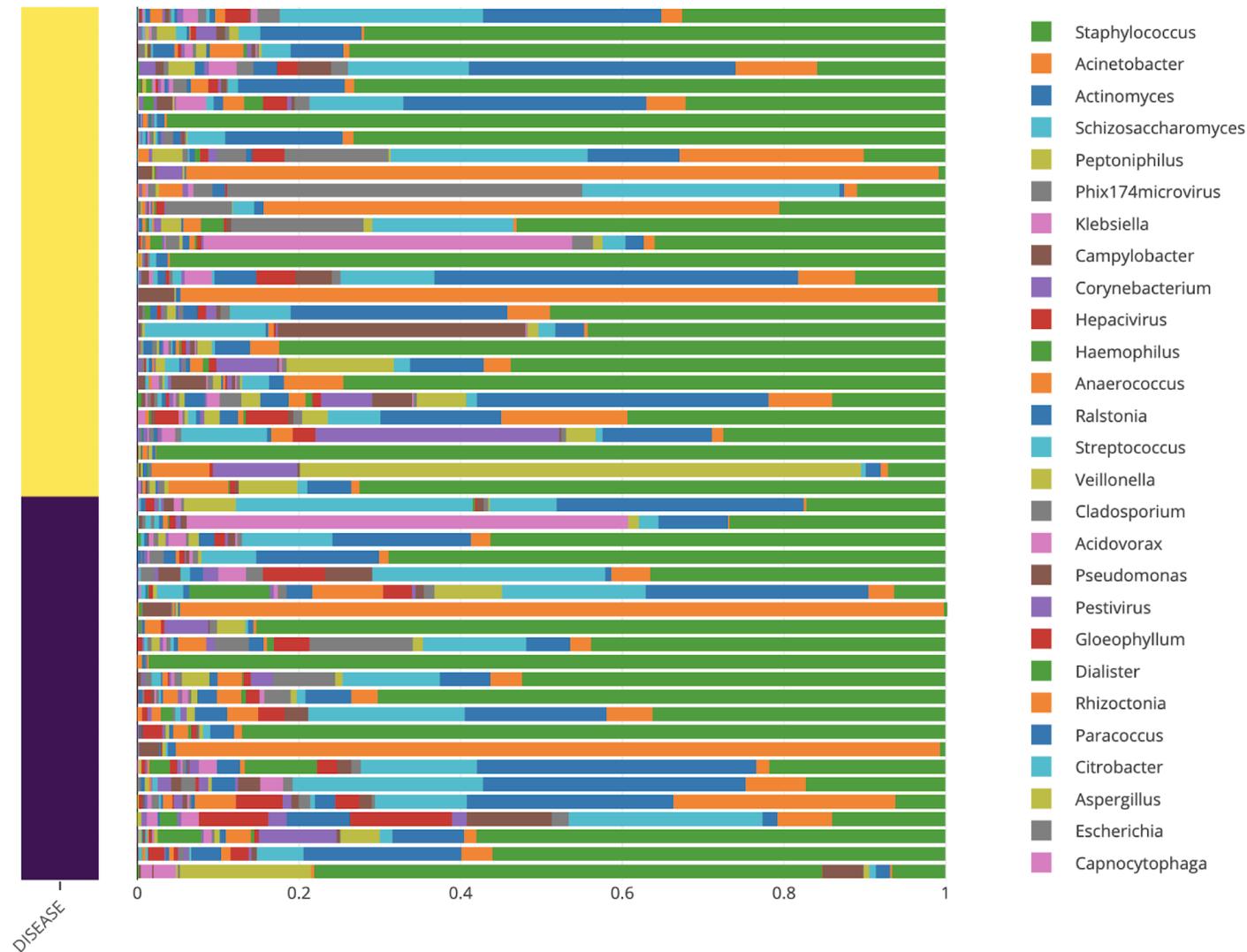


Figure 2. *animalcules* Abundance Tab. In the subtab panel, users can select between a bar plot, heatmap, or box plot. In the bar plot setting, in the left panel, users can select the color by variable, taxonomy level, and sort by option. In the right panel, *animalcules* will show an interactive plot where users can mouse-hover to check the identity of any color bar shown in the plot.

Taxonomy Level

Color Samples by Condition

Sort By
 No Sorting
 Conditions
 Advanced Options

Select Test

Only variables with 2 levels are supported

Select condition

Number of permutations

	Df	SumOfSqs	R2	F	Pr(>F)
condition	1	0.228622865940151	0.0158332366277792	0.77222213390879	0.634
Residual	48	14.2108042275093	0.984166763372221		
Total	49	14.4394270934494	1		

Showing 1 to 3 of 3 entries

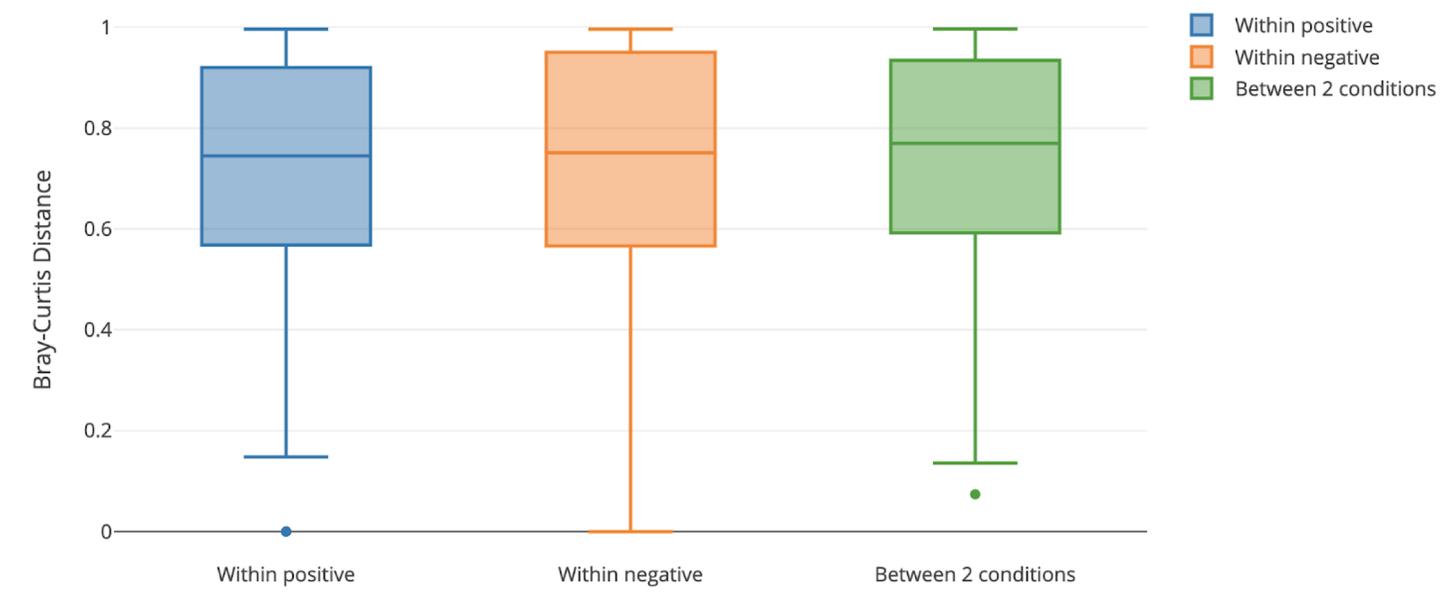


Figure 3. *animalcules* Diversity tab. In the subtab panel, the user could select between alpha diversity analysis and beta diversity analysis. Here in the beta diversity analysis, the right panel controls what statistical test to use, which condition to test on, and show statistical test results in a table as well as a boxplot.

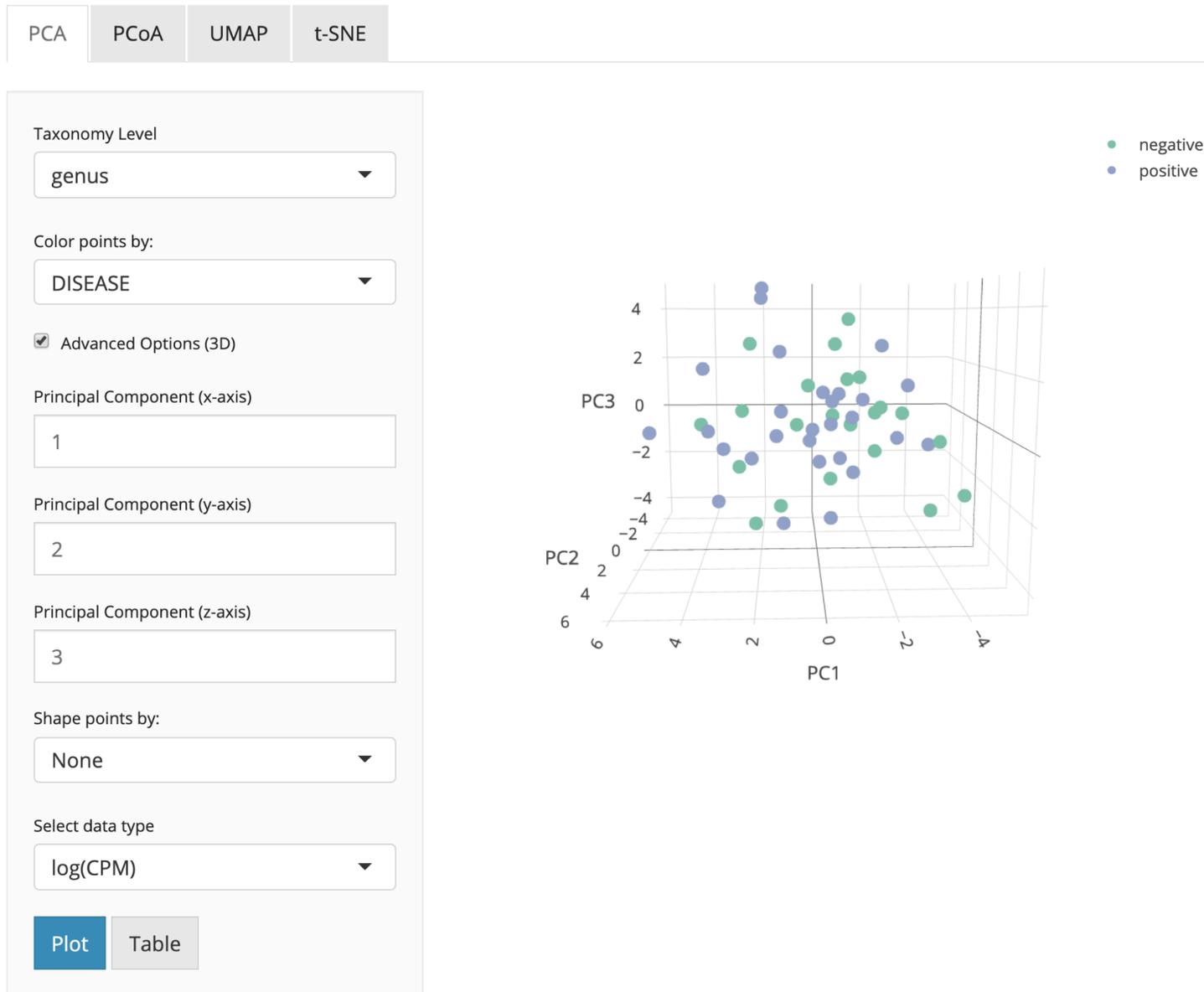


Figure 4. *animalcules* Dimension Reduction tab. In the subtab panel, the user could select between PCA, PCoA, t-SNE, and UMAP. Here in the PCA subtab, the user could choose the taxonomy level, color by variable, and in advanced options, the user could also specify up to three PCs for visualization, shape by variable, and which data type to use.

DESeq2
limma

Taxonomy Level
genus

Select condition
DISEASE

Advanced Options

Select (multiple) covariates

Minimum count cut-off
500

Choose padj cut-off
0.5

Run

Note: For multi-level target variable, all significant results will be printed if existed

Show 10 entries Search:

	microbe	padj	pValue	log2FoldChange	positive	negative	prevalence	Group Size adjusted fold change
1	Paracoccus	0.00529	0.000212	2.93	19/28	14/22	66.00%	1.07
2	Acinetobacter	0.073	0.00584	1.16	28/28	22/22	100.00%	1
3	Schizosaccharomyces	0.176	0.0211	-0.772	28/28	22/22	100.00%	1
4	Corynebacterium	0.323	0.0544	2.12	20/28	11/22	62.00%	1.43
5	Dialister	0.323	0.0646	-4	5/28	5/22	20.00%	1.27
6	Pestivirus	0.348	0.0961	-0.651	27/28	22/22	98.00%	1.04
7	Rhizoctonia	0.348	0.0974	-2.9	6/28	6/22	24.00%	1.27

Showing 1 to 7 of 7 entries Previous 1 Next

Figure 5. *animalcules* Differential Abundance tab. In the subtab panel, users select between DESeq2 and limma. In the left panel, users specify taxonomy level, target condition, covariate variables, count cut-off, and adjusted p-value threshold. In the right panel, a detailed differential abundance result table is shown.

Taxonomy Level
genus ▼

Select Target Condition:
DISEASE ▼

Advanced Options

If the dataset is too small or unbalanced, cross-validation can't be applied. You will see error messages like: NA/NaN/Inf in foreign function call.

Run

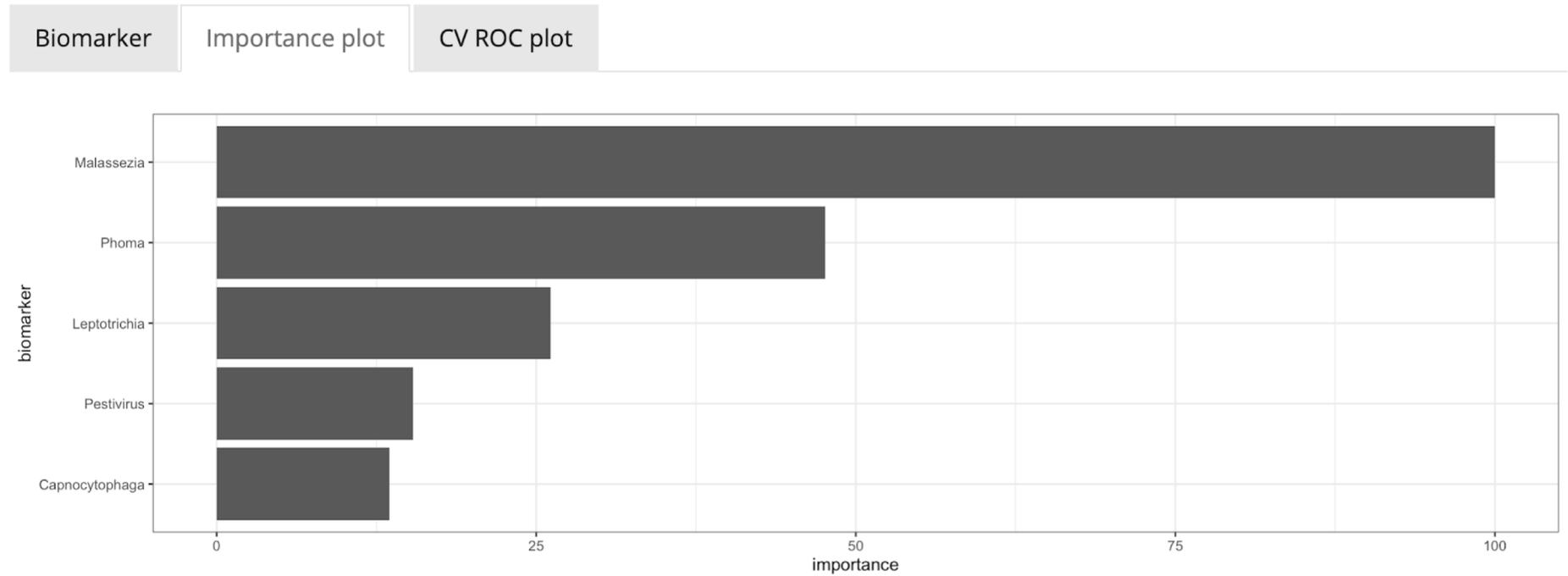


Figure 6. *animalcules* Biomarker tab. In the left panel, users select taxonomy level, and target condition. In the advanced options: number of cross-validations folds, number of cross-validations repeats, biomarker proportion, and classification model. In the right panel, *animalcules* will show the biomarker list, importance plot, and ROC plot.

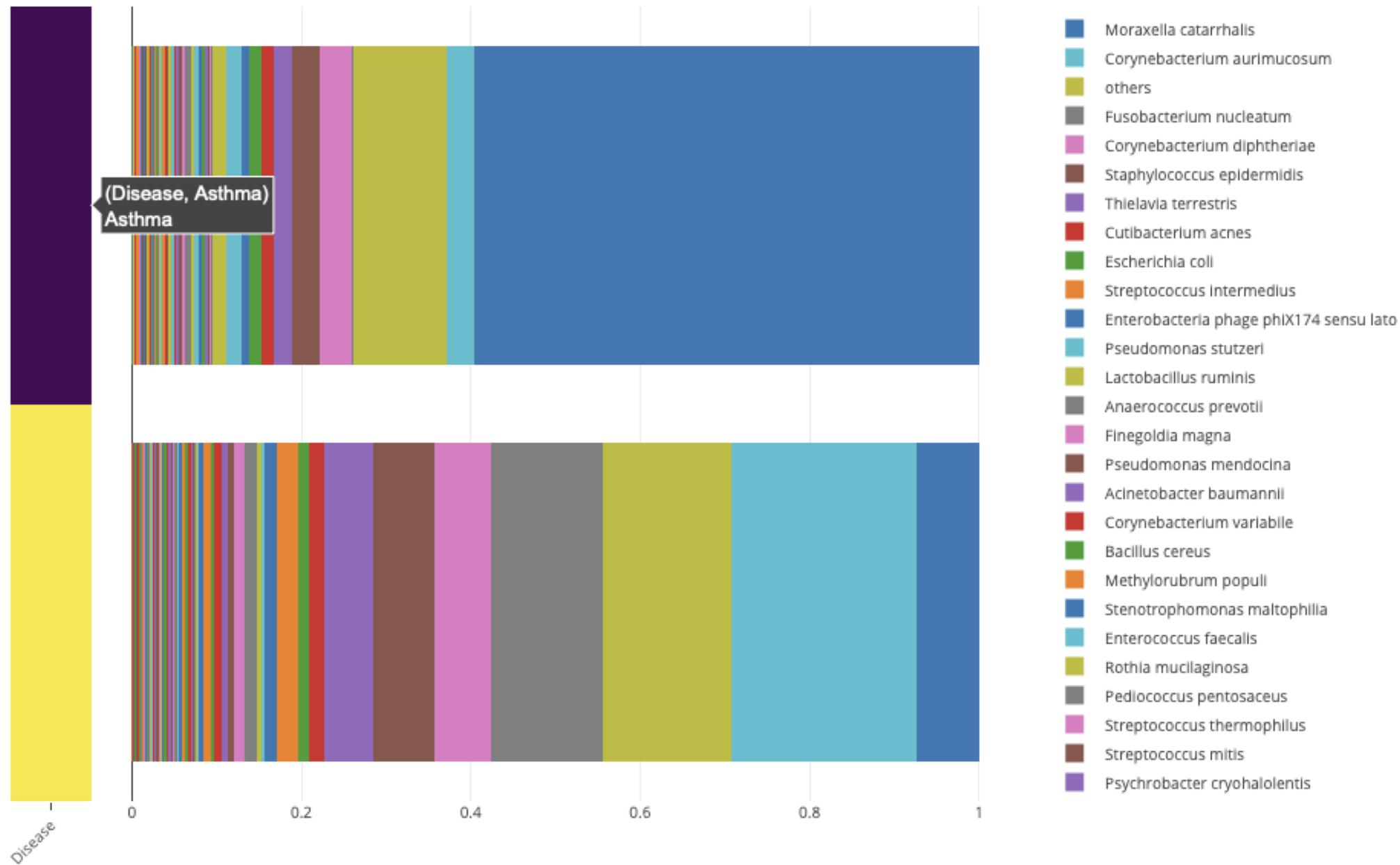


Figure 7. Relative abundance of microbial species bar plot. A stacked bar plot representing the group-wise relative abundance of microbial species in asthmatics (purple) and healthy controls (yellow).

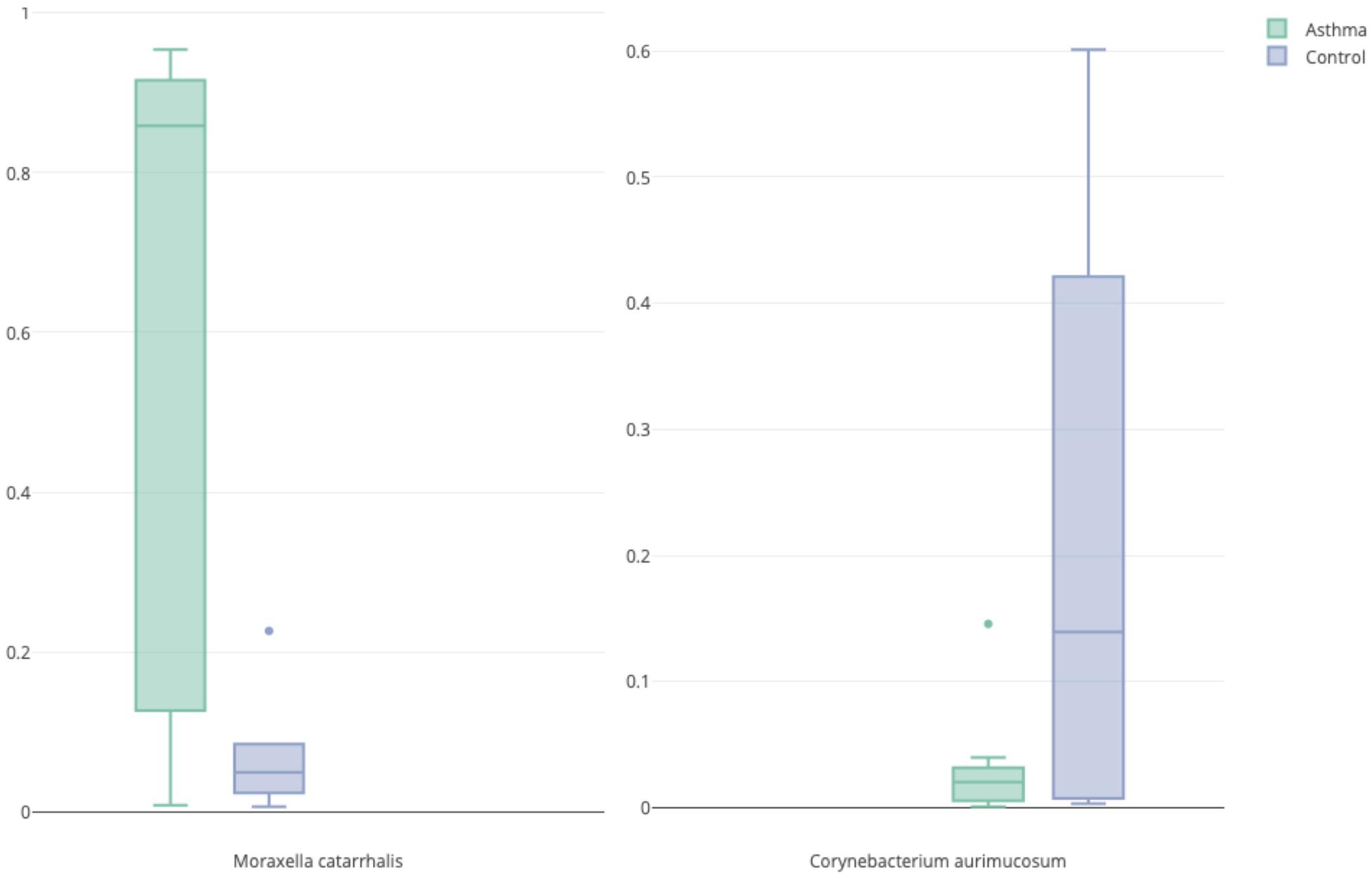


Figure 8. Relative abundance boxplot for differentially abundant species. *Left.* A boxplot of relative abundance of *M. catarrhalis* in asthmatics (green) and healthy controls (blue). *Right.* A boxplot of relative abundance of *C. aurimucosum* in asthmatics (green) and healthy controls (blue).

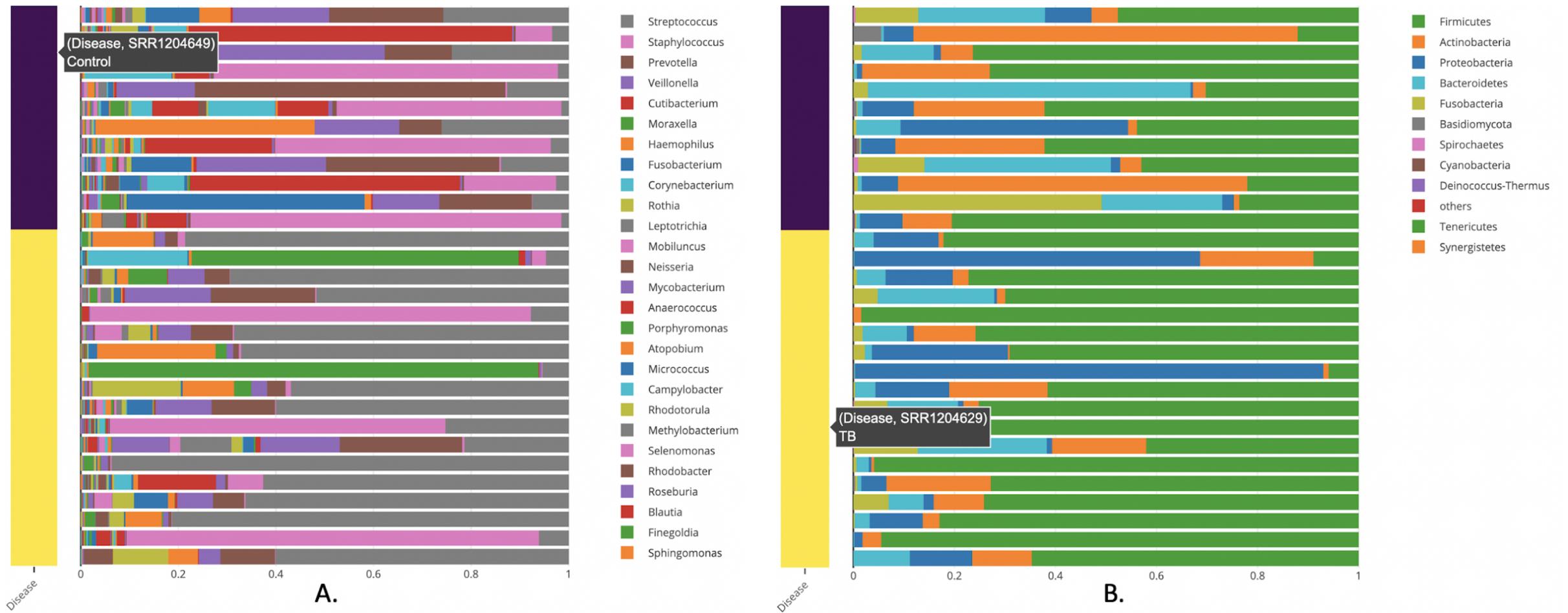
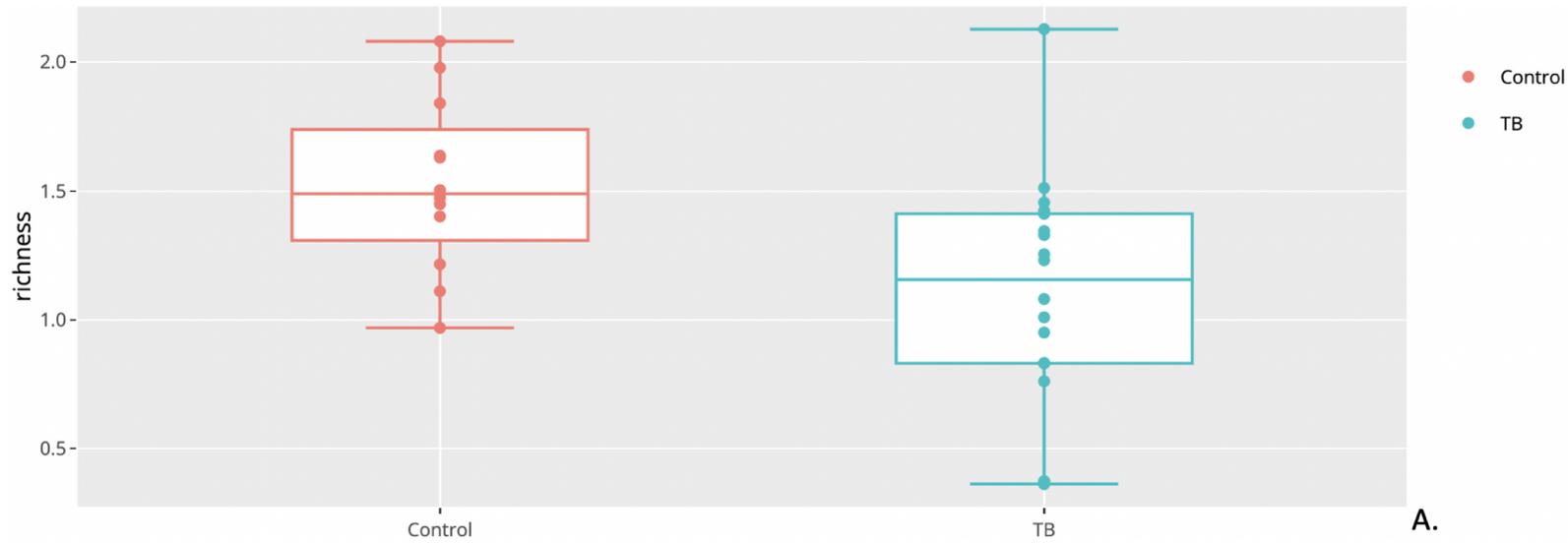
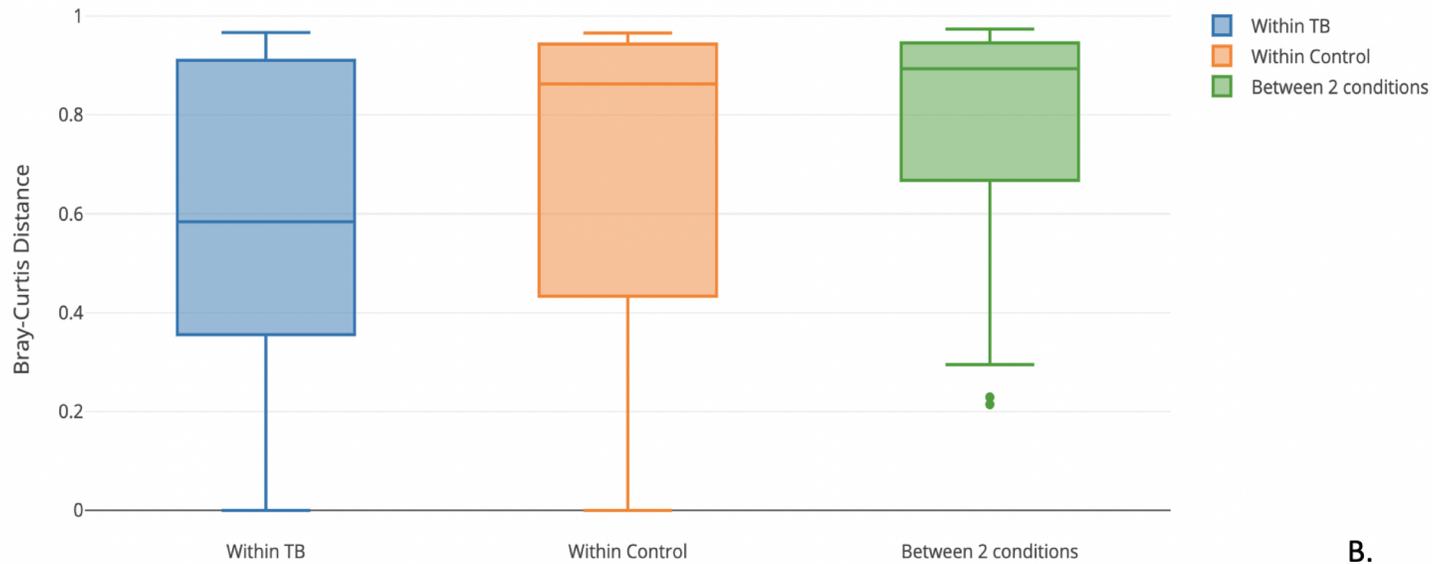


Figure 9. Sample-wise relative abundance bar plot. A stacked bar plot representing the sample-wise relative abundance of microbial species in TB (yellow) and healthy controls (blue). Figure A is the genus level and figure B is the phylum level.



A.



B.

Figure 10. TB example dataset diversity analysis. A. Alpha diversity boxplot between control (red) group and TB (blue) group. B. Beta diversity boxplot within the TB (blue), within the control (orange), and between TB/control group (green).

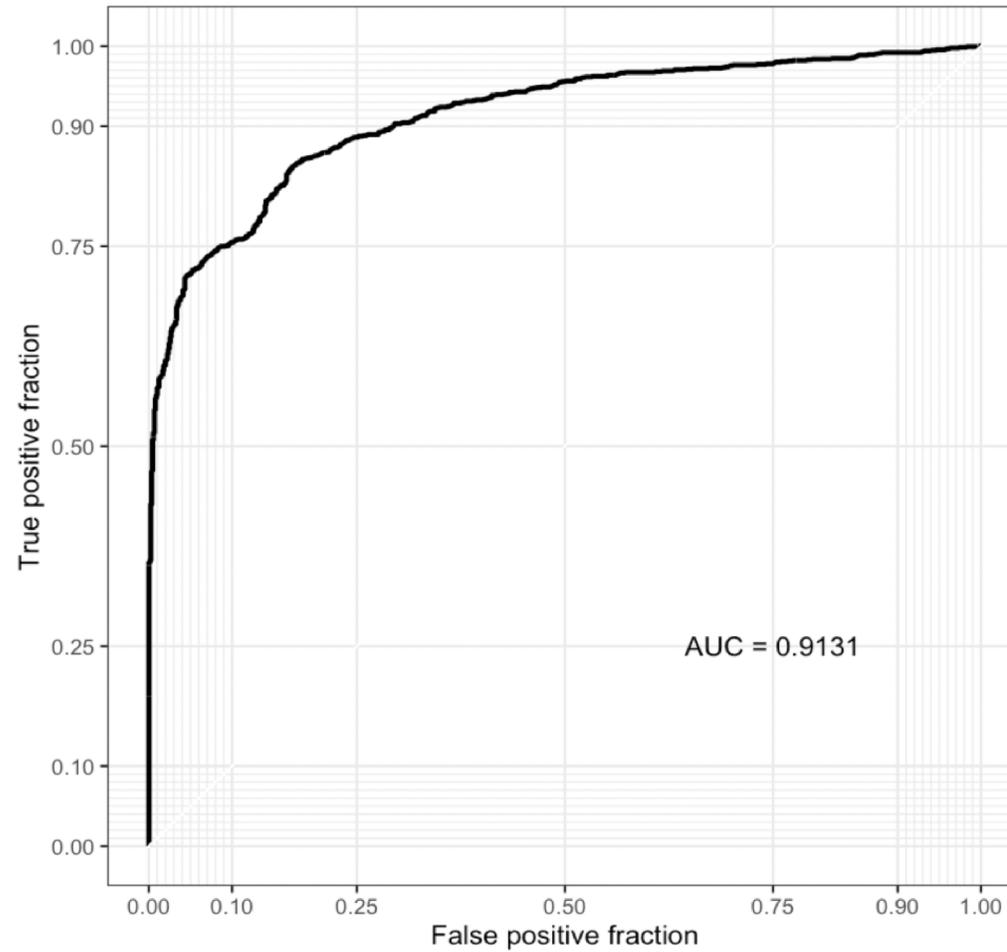


Figure 11. Biomarker ROC curve. ROC shows AUC and cross-validation prediction performance of the identified biomarker.