
PhyloToAST Documentation

Release 1.4.0rc1

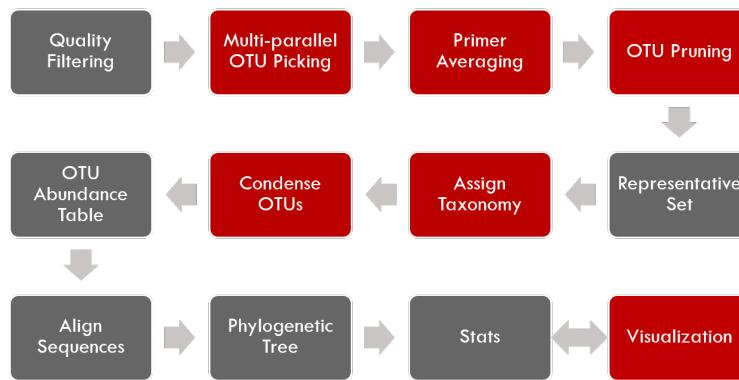
Shareef Dabdoub

Aug 09, 2018

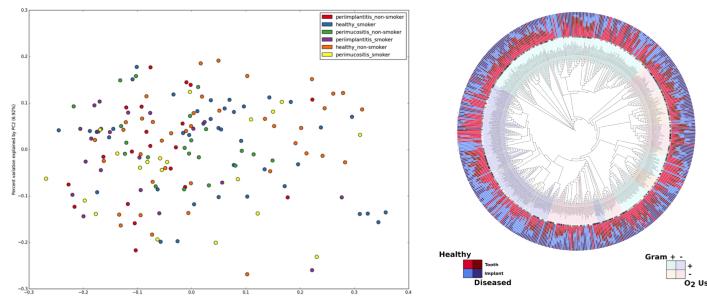
Contents

1 Installation	3
2 Using PhyloToAST	5
3 Citing PhyloToAST	37
4 Publications using PhyloToAST	39
5 References	41
6 Indices and tables	43

The PhyloToAST project is a collection of python scripts that modifies the original QIIME¹ pipeline:



by adding or modifying several steps (above in red) including support for PBS-based cluster-computing, multiple primer support², enhanced support for species-specific analysis, and additional visualization tools.



¹ **QIIME allows analysis of high-throughput community sequencing data.** J Gregory Caporaso, Justin Kuczynski, Jesse Stombaugh, Kyle Bittinger, Frederic D Bushman, Elizabeth K Costello, Noah Fierer, Antonio Gonzalez Pena, Julia K Goodrich, Jeffrey I Gordon, Gavin A Huttley, Scott T Kelley, Dan Knights, Jeremy E Koenig, Ruth E Ley, Catherine A Lozupone, Daniel McDonald, Brian D Muegge, Meg Pirrung, Jens Reeder, Joel R Sevinsky, Peter J Turnbaugh, William A Walters, Jeremy Widmann, Tanya Yatsunenko, Jesse Zaneveld and Rob Knight; Nature Methods, 2010; doi: [10.1038/nmeth.f.303](https://doi.org/10.1038/nmeth.f.303)

² **Target Region Selection Is a Critical Determinant of Community Fingerprints Generated by 16S Pyrosequencing.** Kumar PS, Brooker MR, Dowd SE, Camerlengo T (2011) Target Region Selection Is a Critical Determinant of Community Fingerprints Generated by 16S Pyrosequencing. PLoS ONE 6(6): e20956. doi: [10.1371/journal.pone.0020956](https://doi.org/10.1371/journal.pone.0020956)

CHAPTER 1

Installation

PhyloToAST is available on the Python Package Index (PyPI), and can be easily installed with pip:

```
$ pip install phylotoast
```


CHAPTER 2

Using PhyloToAST

2.1 API

API scripts of PhyloToAST.

2.1.1 biom_calc module

This module provides methods for calculating various metrics with regards to each OTU in an input OTU abundance table.

arcsine_sqrt_transform

Takes the proportion data from relative_abundance() and applies the variance stabilizing arcsine square root transformation:

$$X = \sin^{-1}(\sqrt(p))$$

```
usage: phylotoast.biom_calc.arcsine_sqrt_transform(rel_abd)
```

rel_abd:

Refers to a dictionary keyed on SampleIDs, and the values are dictionaries keyed on OTUID's and their values represent the relative abundance of that OTUID in that SampleID. rel_abd is the output of relative_abundance() function.

return:

Returns a dictionary keyed on SampleIDs, and the values are dictionaries keyed on OTUID's and their values represent the transformed relative abundance of that OTUID in that SampleID.

mean_otu_pct_abundance

Calculate the mean OTU abundance percentage.

```
usage: phylotoast.biom_calc.mean_otu_pct_abundance(rel_abd, otuIDs)
```

rel_abd:

Refers to a dictionary keyed on SampleIDs, and the values are dictionaries keyed on OTUID's and their values represent the relative abundance of that OTUID in that SampleID. rel_abd is the output of relative_abundance() function.

otuIDs:

A list of OTUID's for which the percentage abundance needs to be measured.

return:

A dictionary of OTUID and their percent relative abundance as key/value pair.

MRA

Calculate the mean relative abundance.

```
usage: phylotoast.biom_calc.MRA(biomf)
```

biomf:

A BIOM file.

return:

A dictionary keyed on OTUID's and their mean relative abundance for a given number of sampleIDs.

raw_abundance

Calculate the total number of sequences in each OTU or SampleID.

```
usage: phylotoast.biom_calc.raw_abundance(biomf, sampleIDs=None, sample_abd=True)
```

biomf:

A BIOM file.

sampleIDs:

A list of column id's from BIOM format OTU table. By default, the list has been set to None.

sample_abd:

A boolean operator to provide output for OTUID's or SampleID's. By default, the output will be provided for SampleID's.

return:

Returns a dictionary keyed on either OTUID's or SampleIDs and their respective abundance as values.

relative_abundance

Calculate the relative abundance of each OTUID in a Sample.

```
usage: phylotoast.biom_calc.relative_abundance(biomf)
```

biomf:

A BIOM format.

return:

Returns a dictionary keyed on SampleIDs, and the values are dictionaries keyed on OTUID's and their values represent the relative abundance of that OTUID in that SampleID.

transform_raw_abundance

Function to transform the total abundance calculation for each sample ID to another format based on user given transformation function.

```
usage: phylotoast.biom_calc.transform_raw_abundance(biomf, fn=math.log10,  
→sampleIDs=None, sample_abd=True)
```

biomf:

A BIOM file.

fn:

Mathematical function which is used to transform smax to another format. By default, the function has been given as base 10 logarithm.

sampleIDs:

A list of column id's from BIOM format OTU table. By default, the list has been set to None.

sample_abd:

A boolean operator to provide output for OTUID's or SampleID's. By default, the output will be provided for SampleID's.

return:

Returns a dictionary similar to output of raw_abundance function but with the abundance values modified by the mathematical operation. By default, the operation performed on the abundances is base 10 logarithm.

2.1.2 otu_calc module

otu_name

Determine a simple Genus-species identifier for an OTU, if possible. If OTU is not identified to the species level, name it as Unclassified (family/genus/etc...).

```
usage: phylotoast.otu_calc.otu_name(tax)
```

tax:

QIIME-style taxonomy identifiers, e.g. [‘k__Bacteria’, u’p__Firmicutes’, u’c__Bacilli’, ...]

return:

Returns genus-species identifier based on identified taxonomical level.

load_core_file

For core OTU data file, returns Genus-species identifier for each data entry.

```
usage: phylotoast.otu_calc.load_core_file(core_fp)
```

core_fp:

A file containing core OTU data. Output text file from QIIME's compute_core_microbiome.py script.

return:

Returns genus-species identifier based on identified taxonomical level.

assign_otu_membership

Determines the OTUIDs present in each sample.

```
usage: phylotoast.otu_calc.assign_otu_membership(biomfile)
```

biomfile:

BIOM table object from the biom-format library.

return:

Returns a dictionary keyed on Sample ID with sets containing the IDs of OTUIDs found in each sample.

2.1.3 util module

ensure_dir

Check to make sure the supplied directory path does not exist, if so, create it.

```
usage: phylotoast.util.ensure_dir(d)
```

d:

It is the full path to a directory.

return:

Does not return anything, but creates a directory path if it doesn't exist already.

file_handle

Takes either a file path or an open file handle, checks validity and returns an open file handle or raises an appropriate Exception.

```
usage: phylotoast.util.file_handle(fnh, mode='rU')
```

fnh:

It is the full path to a file, or open file handle.

mode:

The way in which this file will be used, for example to read or write or both. By default, file will be opened in rU mode.

return:
Returns an opened file for appropriate usage.

gather_categories

Find the user specified categories in the map and create a dictionary to contain the relevant data for each type within the categories. Multiple categories will have their types combined such that each possible combination will have its own entry in the dictionary.

```
usage: phylotoast.util.gather_categories(imap, header, categories=None)
```

imap:
The input mapping file data keyed by SampleID.

header:
The header line from the input mapping file. This will be searched for the user-specified categories.

categories:
The list of user-specified categories from the mapping file.

return:
A sorted dictionary keyed on the combinations of all the types found within the user-specified categories. Each entry will contain an empty DataCategory namedtuple. If no categories are specified, a single entry with the key ‘default’ will be returned.

parseFASTA

Parse the records in a FASTA-format file by first reading the entire file into memory.

```
usage: phylotoast.util.parseFASTA(fastaFNH)
```

fastaFNH:
The data source from which to parse the FASTA records. Expects the input to resolve to a collection that can be iterated through, such as an open file handle.

return:
FASTA records containing entries for id, description and data.

parse_map_file

Opens a QIIME mapping file and stores the contents in a dictionary keyed on SampleID (default) or a user-supplied one. The only required fields are SampleID, BarcodeSequence, LinkerPrimerSequence (in that order), and Description (which must be the final field).

```
usage: phylotoast.util.parse_map_file(mapFNH)
```

mapFNH:
Either the full path to the map file or an open file handle.

return:

A tuple of header line for mapping file and a map associating each line of the mapping file with the appropriate sample ID (each value of the map also contains the sample ID). An OrderedDict is used for mapping so the returned map is guaranteed to have the same order as the input file.

parse_taxonomy_table

Greengenes provides a file each OTU a full taxonomic designation. This method parses that file into a map with (key,val) = (OTU, taxonomy).

```
usage: phylotoast.util.parse_taxonomy_table(idtaxFNH)
```

idtaxFNH:

Either the full path to the map file or an open file handle.

return:

A map associating each OTU ID with the taxonomic specifier. An OrderedDict is used so the returned map is guaranteed to have the same order as the input file.

parse_unifrac

Parses the unifrac results file into a dictionary.

```
usage: phylotoast.util.parse_unifrac(unifracFN)
```

unifracFN:

The path to the unifrac results file.

return:

A dictionary with keys: ‘pcd’ (principle coordinates data) which is a dictionary of the data keyed by sample ID, ‘eigvals’ (eigenvalues), and ‘varexp’ (variation explained).

parse_unifrac_v1_8

Function to parse data from older version of unifrac file obtained from Qiime version 1.8 and earlier.

```
usage: phylotoast.util.parse_unifrac_v1_8(unifrac, file_data)
```

unifrac:

The path to the unifrac results file.

file_data

Unifrac data lines after stripping whitespace characters.

return:

A dictionary with keys: ‘pcd’ (principle coordinates data) which is a dictionary of the data keyed by sample ID, ‘eigvals’ (eigenvalues), and ‘varexp’ (variation explained).

parse_unifrac_v1_9

Function to parse data from newer version of unifrac file obtained from Qiime version 1.9 and later.

```
usage: phylotoast.util.parse_unifrac_v1_9(unifrac, file_data)
```

unifrac:

The path to the unifrac results file.

file_data

Unifrac data lines after stripping whitespace characters.

return:

A dictionary with keys: ‘pcd’ (principle coordinates data) which is a dictionary of the data keyed by sample ID, ‘eigvals’ (eigenvalues), and ‘varexp’ (variation explained).

split_phylogeny

Return either the full or truncated version of a QIIME-formatted taxonomy string.

```
usage: phylotoast.util.split_phylogeny(p, level='s')
```

p:

A QIIME-formatted taxonomy string: k_Foo; p_Bar; ...

level:

The different level of identification are kingdom (k), phylum (p), class (c),order (o), family (f), genus (g) and species (s). The default level of identification is species.

return:

A QIIME-formatted taxonomy string up to the classification given by param level.

storeFASTA

Parse the records in a FASTA-format file by first reading the entire file into memory.

```
usage: phylotoast.util.storeFASTA(fastaFNH)
```

fastaFNH:

The data source from which to parse the FASTA records. Expects the input to resolve to a collection that can be iterated through, such as an open file handle.

return:

FASTA records containing entries for id, description and data.

write_map_file

Given a list of mapping items (in the form described by the parse_mapping_file method) and a header line, write each row to the given input file with fields separated by tabs.

```
usage: phylotoast.util.write_map_file(mapFNH, items, header)
```

mapFNH:

Either the full path to the map file or an open file handle.

items:

The list of row entries to be written to the mapping file.

header:

The descriptive column names that are required as the first line of the mapping file.

return:

None.

2.1.4 graph_util

plot_kde

Plot a smoothed (by kernel density estimate) histogram.

```
usage: phylotoast.graph_util.plot_kde(data, ax, title=None, color='r', fill_bt=True)
```

data:

An array containing the data to be plotted.

ax:

The Axes object to draw to.

title:

The plot title.

color:

The color of the histogram line and fill. Note that the fill will be plotted with an alpha of 0.35.

fill_bt:

Specify whether to fill the area beneath the histogram line.

2.2 Data Handling

PhyloToAST scripts for data analysis and manipulation.

2.2.1 assign_taxonomy_by_blast_result.py

Assign taxonomy to a rep set of OTUs that were chosen by BLAST from an annotated database.

```
usage: assign_taxonomy_by_blast_result.py [-h] -i REP_SET_FP -t ID_TO_
                                          ↵TAXONOMY_FP [-o ASSIGNED_TAXONOMY_FP] [-v]
```

Required arguments

-i REP_SET_FP, --rep_set_fp REP_SET_FP

The set of representative sequences.

-t ID_TO_TAXONOMY_FP, --id_to_taxonomy_fp ID_TO_TAXONOMY_FP
 Path to tab-delimited file mapping sequences to assigned taxonomy.

Optional arguments

-o ASSIGNED_TAXONOMY_FP, --assigned_taxonomy_fp ASSIGNED_TAXONOMY_FP
 The path to the result file. By default outputs to assigned_taxonomy.txt

-h, --help
 Show the help message and exit

-v, --verbose
 Print detailed information about script operation.

2.2.2 barcode_filter.py

From an input FASTA file, filter all sequences with barcodes matching those in an input mapping file.

```
usage: barcode_filter.py [-h] -i INPUT_FASTA_FN -m MAPPING_FN [-q QUALITY_
FN] [-o OUTPUT_PREFIX] [-v]
```

Required arguments

-i INPUT_FASTA_FN, --input_fasta_fn INPUT_FASTA_FN
 The sequence data file to be filtered.

-m MAPPING_FN, --mapping_fn MAPPING_FN
 The mapping file containing the barcodes you want filtered sequenced to contain.

Optional arguments

-q QUALITY_FN, --quality_fn QUALITY_FN
 The quality data file. If you plan to use quality data with split_libraries.py, you have to filter the quality data as well.

-o OUTPUT_PREFIX, --output_prefix OUTPUT_PREFIX
 The prefix for the output filtered data

-h, --help
 Show the help message and exit

-v, --verbose
 Print detailed information about script operation.

2.2.3 biom_relative_abundance.py

Convert a BIOM file of OTU abundance data into a CSV of relative abundance data.

```
usage: biom_relative_abundance.py [-h] [-i INPUT_BIOM_FP] [-o OUTPUT_TSV_FP]
[-v]
```

Required arguments

-i INPUT_BIOM_FP, --input_biom_fp INPUT_BIOM_FP
The BIOM file path.

Optional arguments

-o OUTPUT_CSV_FP, --output_csv_fp OUTPUT_CSV_FP
A CSV table of relative OTU abundance data.

--stabilize_variance
Apply the variance-stabilizing arcsine square root transformation to the OTU proportion data.

-h, --help
Show the help message and exit

-v, --verbose
Print detailed information about script operation.

2.2.4 condense_workflow.py

This workflow script will run all three steps of the OTU condensing pipeline automatically with the default output file settings.

```
usage: condense_workflow.py [-h] -i ASSIGNED_TAXONOMY_FN -r REP_SET_FN -s
                            ↪SEQS_OTUS_FN [-L {k,p,c,o,f,g,s}] [-v]
```

Required arguments

-i ASSIGNED_TAXONOMY_FN, --assigned_taxonomy_fn ASSIGNED_TAXONOMY_FN
The taxonomy file output by the assign_taxonomy script.

-r REP_SET_FN, --rep_set_fn REP_SET_FN
The set of representative sequences.

-s SEQS_OTUS_FN, --seqs_otus_fn SEQS_OTUS_FN
The list of OTU IDs and their associated sequence IDs.

Optional arguments

-L {k,p,c,o,f,g,s}, --phylogenetic_level {k,p,c,o,f,g,s}
Set the phylogenetic level at which to define OTUs for condensing and downstream processing. Defaults to species level.

-h, --help
Show the help message and exit

-v, --verbose
Print detailed information about script operation.

2.2.5 extract_shared_or_unique_otuids.py

Parse a BIOM format file and obtain a list of unique OTUIDs found in each category in mapping file.

```
usage: extract_uniques.py [-h] [-p PREFIX] input_biom_fp output_dir mapping_
                           file category_column
```

Required arguments

`input_biom_fp`

BIOM format file path.

`mapping_file`

Mapping file with category information.

`category_column`

Column in mapping file specifying the category/condition of all samples.

Optional arguments

`-o OUTPUT_DIR, --output_dir OUTPUT_DIR`

Path to save category unique OTUIDs.

`-p PREFIX, --prefix PREFIX`

Provide specific text to prepend the output file names. By default, the ‘unique’ will be added in front of output filenames.

`-r REVERSE, --reverse REVERSE`

Get shared OTUIDs among all unique combinations of groups and write out the results to path provided to this option.

`-h, --help`

Show the help message and exit

2.2.6 filter_biom.py

Filter biom file on both ‘sample’ and ‘observation’ axes, given a list of sampleIDs to retain.

```
usage: filter_biom.py [-h] [-fo FILTER_OTUIDS_FNH] input_biom_fnh output_
                           biom_fnh mapping_fnh
```

Required arguments

`input_biom_fnh`

BIOM file path.

`output_biom_fnh`

Filtered biom output file.

`mapping_fnh`

Mapping file with sampleIDs to retain in it. The ‘#SampleID’ column will be used to get the list of all ids to retain.

Optional arguments

- fo** FILTER_OTUIDS_FNH, **--filter_otuids_fnh** FILTER_OTUIDS_FNH
Path to file to write out the list of OTUIDs not present in any SampleIDs in mapping file. This output is usually used to filter out unwanted otuids from “.tre” file. If not given, the discarded OTUIDs list will be saved in the current working directory.
- h, --help**
Show the help message and exit

2.2.7 filter_rep_set.py

Step 2 of the condensing process. Filter the representative sequence set to include only those sequences that map to unique OTUs.

```
usage: filter_rep_set.py [-h] -r REP_SET_FN -u UNIQUE_OTUS_FN [-o OUTPUT_  
- FILTERED REP_SET_FN] [-v]
```

Required arguments

- r** REP_SET_FN, **--rep_set_fn** REP_SET_FN
The set of representative sequences.
- u** UNIQUE_OTUS_FN, **--unique_otus_fn** UNIQUE_OTUS_FN
The condensed assigned taxonomy file.

Optional arguments

- o** OUTPUT_FILTERED REP_SET_FN, **--output_filtered_rep_set_fn** OUTPUT_FILTERED REP_SET_FN
The filtered representative set. By default outputs to condensed_rep_set.fna
- h, --help**
Show the help message and exit
- v, --verbose**
Print detailed information about script operation.

2.2.8 merge_otu_results.py

Distributing sequence data across the cluster for OTU picking results in a set of result files that need to be merged into a single pick otus result.

```
usage: merge_otu_results.py [-h] [-o OUTPUT_FN] [-v] pick_otus_results [pick_  
- otus_results ...]
```

Required arguments

- pick_otus_results**
The result files from multiple runs of a pick otus script that need to be merged.

Optional arguments

- o, --output_fn** OUTPUT_FN
The name of the file the merged results will be written to.
- h, --help**
Show the help message and exit.
- v, --verbose**
Print detailed information about script operation.

2.2.9 multi_parallel_pick_otus.py

Generate PBS scripts for submission to the OSC to run the QIIME parallel blast pick OTUs script on multiple input sequence data sets.

```
usage: osc_parallel_pick_otus.py [-h] -i INPUT_FNA [INPUT_FNA ...] [-t _  
-WALLTIME] [-n JOB_NAME] [-v]
```

Required arguments

- i** INPUT_FNA [INPUT_FNA ...], **--input_fna** INPUT_FNA [INPUT_FNA ...]
The names of the sequence files that will be have PBS scripts generated to process them. The expected input is from the split_sequence_data.py script (e.g. 0.fna, 1.fna, ..., n.fna).
- t** WALLTIME, **--walltime** WALLTIME
The maximum running time to specify to the OSC queuing system for each script.
- n** JOB_NAME, **--job_name** JOB_NAME
A descriptive name for the job script that will appear when checking the job status. Max length is 15 characters, but ‘_#’ will be appended to the name you provide to differentiate among all the jobs, so this parameter will be truncated if necessary to accommodate for the number of input files.
- h, --help**
Show the help message and exit
- v, --verbose**
This will cause the program to print the full path for each output file to the command line. This can be used for informational purposes or to pipe (l) to the PBS multi-submission script to automate job submission as soon as the scripts are created.

2.2.10 multi_qsub.py

Submit multiple PBS job scripts to the queuing system (qsub) and store the output job IDs.

```
usage: multi_qsub.py [-h] [-t] job_scripts [job_scripts ...]
```

Required arguments

- job_scripts**
The job script files to submit to the queuing system.

Optional arguments

- h, --help**
Show the help message and exit
- t, --test**
Only print each of the qsub commands instead of actually running the commands.

2.2.11 network_plots_gephi.py

Create network plots based on correlation matrix.

```
usage: network_plots_gephi.py [-h] [-go GEXF_OUT] [-fp FIL_PCT] [-w STATS_OUT_FNH]_
                                biom_file mapping_file condition_column in_corr_mat cat_name
```

Required Arguments

biom_file

The biom-format file.

mapping_file

Mapping file for reading sampleIDs and their groups.

condition_column

Column name in mapping file denoting the categories.

in_corr_mat

Correlation matrix file. The format for the tab-separated file should be: Category -> Variable -> by Variable -> Correlation

cat_name

Category to be plotted.

Optional Arguments

- go GEXF_OUT, --gexf_out GEXF_OUT**
The directory to output the PCoA plots to.
- scaling_factor SCALING_FACTOR**
Graph information written to this Graph Exchange XML Format file. This file can be input to Gephi.
- fp FIL_PCT, --fil_pct FIL_PCT**
Specify the minimum value of correlation strength to display. By default, all correlations greater than or equal to 0.75 will be shown.
- w STATS_OUT_FNH, --stats_out_fnh STATS_OUT_FNH**
Write out graph statistics - degree and betweenness centrality calculations for each node.
- h, --help**
Show this help message and exit

2.2.12 otu_condense.py

Step 1 of the condensing process. Take a taxonomy table from the assign_taxonomy QIIME script and prune all redundant taxonomy strings

```
usage: otu_condense.py [-h] -i INPUT_ASSIGNED_TAXONOMY [-p PRUNED_OUTPUT_
->FILE] [-n NON_UNIQUE_OUTPUT_FILE] [-l {k,p,c,o,f,g,s}] [-v]
```

Required arguments

- i** INPUT_ASSIGNED_TAXONOMY, **--input_assigned_taxonomy** INPUT_ASSIGNED_TAXONOMY
The taxonomy file output by the assign_taxonomy script.

Optional arguments

- p** PRUNED_OUTPUT_FILE, **--pruned_output_file** PRUNED_OUTPUT_FILE
The output file for the pruned taxonomy list. Defaults to condensed_assigned_taxonomy.txt
- n** NON_UNIQUE_OUTPUT_FILE, **--non_unique_output_file** NON_UNIQUE_OUTPUT_FILE
The file will contain a list of pruned OTU IDs associated with the OTU IDs they replaced. Defaults to nonunique_otu_matrix.txt
- l** {k,p,c,o,f,g,s}, **--phylogenetic_level** {k,p,c,o,f,g,s}
Set the phylogenetic level at which to define OTUs for condensing and downstream processing. Defaults to species level.
- h, --help**
Show the help message and exit
- v, --verbose**
Print detailed information about script operation.

2.2.13 otu_to_tax_name.py

Convert a list of OTU IDs to a list of OTU IDs paired with Genus_species identifiers.

```
usage: otu_to_tax_name.py [-h] -i OTU_ID_FP -t TAXONOMY_FP [-o OUTPUT_FP]
```

Required arguments

- i** OTU_ID_FP, **--otu_id_fp** OTU_ID_FP
Either a text file containing a list (one per line) of OTU IDs, or a tab-separated (classic) BIOM-format file.
- t** TAXONOMY_FP, **--taxonony_fp** TAXONOMY_FP
A file associating OTU ID with a full taxonomic specifier.

Optional arguments

- o** OUTPUT_FP, **--output_fp** OUTPUT_FP
For a list input, a new file containing a list of OTU IDs and their corresponding short taxonomic identifiers separated by tabs. For a BIOM file input, a new mapping file with all the OTU IDs replaced by the short identifier.
- h, --help**
Show the help message and exit

2.2.14 pick_otus_condense.py

Step 3 of the condensing process. Condense the QIIME pick_otus.py script output by moving the sequences associated with non-unique OTUs to OTU IDs that were identified as unique.

```
usage: pick_otus_condense.py [-h] -s SEQS_OTUS -n NON_UNIQUE_OTU_MATRIX [-o CONDENSED_SEQS_OTUS_FILE] [-v]
```

Required arguments

-s SEQS_OTUS, **--seqs_otus** SEQS_OTUS

The list of OTU IDs and their associated sequence IDs.

-n NON_UNIQUE_OTU_MATRIX, **--non_unique_otu_matrix** NON_UNIQUE_OTU_MATRIX

The list of unique OTU IDs associated with the OTU IDs they replaced.

-o CONDENSED_SEQS_OTUS_FILE, **--condensed_seqs_otus_file** CONDENSED_SEQS_OTUS_FILE

The condensed set of OTU IDs and the matching sequences. By default outputs to condensed_seqs_otus.txt

Optional arguments

-h, --help

Show the help message and exit

-v, --verbose

Print detailed information about script operation.

2.2.15 primer_average.py

Combine multi-primer pick OTUs results files into a single results file while at the same time averaging sequence counts per sample for OTUs shared between the primer-set results. See reference: Kumar PS et al. (2011) doi:10.1371/journal.pone.0020956

```
usage: primer_average.py [-h] --p1 P1 --p2 P2 [-o OUTPUT_FP] [-v]
```

Required arguments

--p1 P1

Primer-set 1 seqs_otus results files.

--p2 P2

Primer-set 2 seqs_otus results files.

Optional arguments

-o OUTPUT_FP, **--output_fp** OUTPUT_FP

The combined seqs_otus file that has been averaged by shared OTU entries. Default: combined_seqs_otus.txt

-h, --help

Show the help message and exit

-v, --verbose

Print detailed information about script operation.

2.2.16 prune_otus.py

Parse the OTU-sequence data in two steps. First remove any OTUs that occur in less than a user-defined percent of samples (default 5%). Second, remove any OTUs that make up less than a user-defined percentage of the overall sequences (default 0.01%)

```
usage: prune_otus.py [-h] -i SEQS_OTUS_FN -t ID_TO_TAXONOMY_FN [-p PERCENT_OF_SAMPLES] [-s PERCENT_OF_SEQUENCES] [-l {k,p,c,o,f,g,s}] [-o OUTPUT_PRUNED_OTUS_FN] [-v] [-output_removed_otus_fn OUTPUT_REMOVED_OTUS_FN]
```

Required arguments

- i** SEQS_OTUS_FN, **--seqs_otus_fn** SEQS_OTUS_FN
The output from the pick OTUs step, e.g. seqs_otus.txt
- t** ID_TO_TAXONOMY_FN, **--id_to_taxonomy_fn** ID_TO_TAXONOMY_FN
Path to tab-delimited file mapping sequences to assigned taxonomy.

Optional arguments

- p** PERCENT_OF_SAMPLES, **--percent_of_samples** PERCENT_OF_SAMPLES
OTUs that occur in less than this percent of samples will be removed. Default is 5 percent.
- s** PERCENT_OF_SEQUENCES, **--percent_of_sequences** PERCENT_OF_SEQUENCES
OTUs that occur in less than this percent of total sequences will be removed. Default is 0.01 percent.
- l** {k,p,c,o,f,g,s}, **--phylogenetic_level** {k,p,c,o,f,g,s}
Set the phylogenetic level at which to join OTUs for consideration in pruning. Default is 'g'(group).
- o** OUTPUT_PRUNED_OTUS_FN, **--output_pruned_otus_fn** OUTPUT_PRUNED_OTUS_FN
The main output file that will contain the remaining OTUs and sequence IDs.
- output_removed_otus_fn** OUTPUT_REMOVED_OTUS_FN
The file to write out the set of OTUs that were removed by the filter.
- h, --help**
Show the help message and exit
- v, --verbose**
Print detailed information about script operation.

2.2.17 restrict_repset.py

Take a subset BIOM table (e.g. from a core calculation) and a representative set (repset) FASTA file and create a new repset restricted to the OTUs in the BIOM table.

```
usage: restrict_repset.py [-h] -i BIOM_FP -r REPSET_FP [-o REPSET_OUT_FP]
```

Required arguments

- i** BIOM_FP, **--biom_fp** BIOM_FP
Path to a biom-format file with OTU-Sample abundance data.
- r** REPSET_FP, **--repset_fp** REPSET_FP
Path to a FASTA-format file containing the representative set of OTUs.

Optional arguments

- o REPSET_OUT_FP, --repset_out_fp** REPSET_OUT_FP
Path to the new restricted repset file.
- h, --help**
Show the help message and exit

2.2.18 split_sequence_data.py

Split an input FASTA-formatted sequence file into a user-specified number of smaller files such that the sequence data is evenly distributed among them.

```
usage: split_sequence_data.py [-h] -i INPUT_FASTA_FN [-n NUM_OUTPUT_FILES] [-o OUTPUT_DIR] [-v]
```

Required arguments

- i INPUT_FASTA_FN, --input_fasta_fn** INPUT_FASTA_FN
The sequence data file to be split up into a series of smaller files.
- n NUM_OUTPUT_FILES, --num_output_files** NUM_OUTPUT_FILES
The number of files the input data should be split into.

Optional arguments

- o OUTPUT_DIR, --output_dir** OUTPUT_DIR
The location to write the split data files.
- h, --help**
Show the help message and exit
- v, --verbose**
Print detailed information about script operation.

2.2.19 transpose_biom.py

Transpose a BIOM-format file so that the matrix is sample by species.

```
usage: transpose_biom.py [-h] -i INPUT_BIOM_FP -m MAPPING [-c MAP_CATEGORY] [-o OUTPUT_BIOM_FP] [-v]
```

Required arguments

- i INPUT_BIOM_FP, --input_biom_fp** INPUT_BIOM_FP
The BIOM-format file.
- m MAPPING, --mapping** MAPPING
The mapping file specifying group information for each sample.
- o OUTPUT_BIOM_FP, --output_biom_fp** OUTPUT_BIOM_FP
The BIOM-format file to write.

Optional arguments

-c MAP_CATEGORY, --map_category MAP_CATEGORY
A mapping category, such as TreatmentType, that will be used to split the data into separate BIOM files; one for each value found in the category.

-h, --help
Show the help message and exit

-v, --verbose
Print detailed information about script operation.

2.3 Visualization

PhyloToAST scripts used for visualizing data.

2.3.1 diversity.py

Calculate the alpha diversity of a set of samples using one or more metrics and output a kernal density estimator-smoothed histogram of the results.

```
usage: diversity.py [-h] [-d DIVERSITY [DIVERSITY ...]] [--plot_title PLOT_TITLE] [--image_type IMAGE_TYPE] [--save_calculations SAVE_CALCULATIONS] [--show_significance] [--show_available_metrics] -m MAP_FILE -i BIOM_FP -c CATEGORY --color_by COLOR_BY -o OUT_DIR
```

Required arguments

-m MAP_FILE, --map_file MAP_FILE
QIIME mapping file.

-i BIOM_FP, --biom_fp BIOM_FP
BIOM table file name

-c CATEGORY, --category CATEGORY
Specific category from the mapping file.

--color_by COLOR_BY
A column name in the mapping file containing hexadecimal (#FF0000) color values that will be used to color the groups. Each sample ID must have a color entry.

-o OUT_DIR, --out_dir OUT_DIR
The directory all plots will be saved to.

Optional arguments

-h, --help
show this help message and exit

-d DIVERSITY [DIVERSITY ...], --diversity DIVERSITY [DIVERSITY ...]
The alpha diversity metric. Default value is ‘shannon’, which will calculate the Shannon entropy. Multiple metrics can be specified (space separated). The full list of metrics is available at: <http://scikit-bio.org/docs/latest/generated/skbio.diversity.alpha.html>.

--plot_title PLOT_TITLE

The name of a PDF file the pathway map will be written to.

-p IMAGE_TYPE, **--image_type** IMAGE_TYPE

The type of image to save: PNG, SVG, etc.

--save_calculations SAVE_CALCULATIONS

Path and name of text file to store the calculated diversity metrics.

--show_significance

Display significance testing results. The results will be shown by default.

--show_available_metrics

Supply this parameter to see which alpha diversity metrics are available for usage. No calculations will be performed if this parameter is provided.

2.3.2 iTol.py

Create files appropriate for use in the iTOL visualization program by using the abundance data from a biom-format file and groups specified in a QIIME mapping file. The program also modifies a Newick-format phylogenetic tree file to use proper taxonomic names instead of OTU IDs for useful display in iTOL.

```
usage: iTol.py [-h] -i OTU_TABLE -m MAPPING [-t INPUT_TREE] [-e OUTPUT_TRE] [-o_
                ↪OUTPUT_ITOL_TABLE] [-c MAP_CATEGORIES] [-a {MRA,NMRA,raw}]
```

Required arguments

-i OTU_TABLE, **--otu_table** OTU_TABLE

The biom-format file with OTU-Sample abundance data.

-m MAPPING, **--mapping** MAPPING

The mapping file specifying group information for each sample.

Optional arguments

-t INPUT_TREE, **--input_tree** INPUT_TREE

A phylogenetic tree in Newick format to be modified by exchanging the OTU ID node names for taxonomic names.

-e OUTPUT_TRE, **--output_tre** OUTPUT_TRE

The output Newick-format tree (.tre) file

-o OUTPUT_ITOL_TABLE, **--output_itol_table** OUTPUT_ITOL_TABLE

Other than a phylogenetic tree, the main input to iTOL is a dataset file containing some representation of the abundance of every OTU across the specified data groups. This program provides multiple calculation methods. See the **-analysis_metric** option for details.

-c MAP_CATEGORIES, **--map_categories** MAP_CATEGORIES

Any mapping categories, such as treatment type, that will be used to group the data in the output iTOL table. For example, one category with three types will result in three data columns in the final output. Two categories with three types each will result in six data columns. Default is no categories and all the data will be treated as a single group.

-a {MRA, NMRA, raw}, **--analysis_metric** {MRA, NMRA, raw}

Specifies which metric is calculated on the abundance data in the OTU table. Available options: MRE - mean relative abundance (Abundance data is normalized by total sample abundance, then averaged across OTU),

NMRE - normalized mean relative abundance (MRE normalized by the total MRE across the groups as specified in `-map_categories`), raw (outputs the actual sequence abundance data for each OTU).

--stabilize_variance

Apply the variance-stabilizing arcsine square root transformation to the OTU proportion data. Recommended for usage with `-a NMRA` or `-a MRA`.

-h, --help

Show the help message and exit.

Workflow for generating useful phylogenetic trees using PhyloToAST

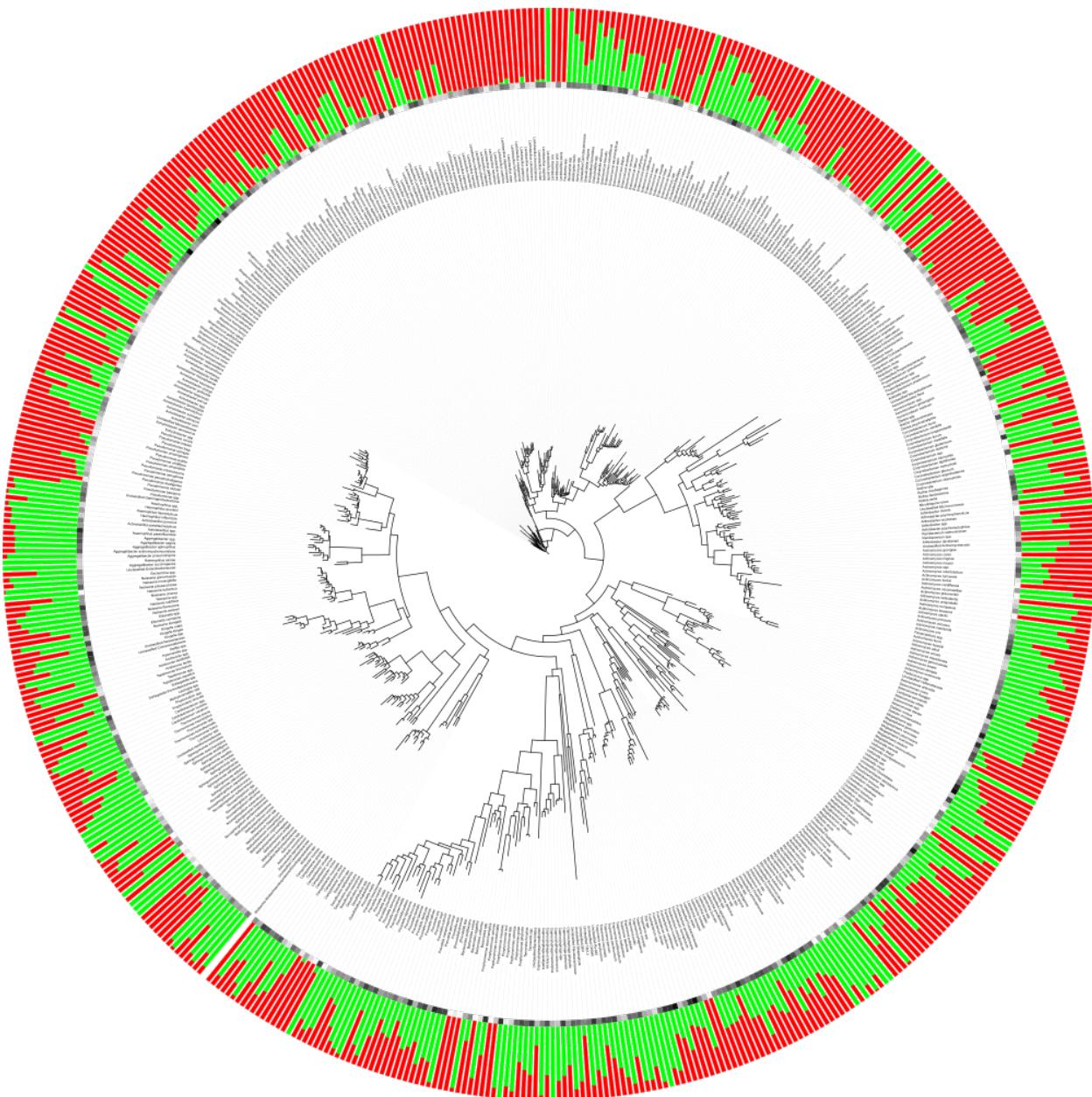
Step 1 : Obtain .tre file from QIIME's [make_phylogeny.py](#) script.

Step 2 : Run `iTol.py` script with `-a NMRA` analysis metric. This file will denote the multibar graph around the circular phylogenetic tree.

Step 3 : Run `iTol.py` script with `-a raw` analysis metric. This file will denote the gradient graph around the circular phylogenetic tree.

Step 4 : Upload modified .tre file from `iTol.py` script to [iTOL website](#). Add your dataset files and obtain the final phylogenetic tree figure.

Example output image



NOTE : Please refer to [iTOL help page](#) for changing dataset parameters.

2.3.3 LDA.py

This script calculates and returns LDA plots based on normalized relative abundances or distance matrices (for e.g. unifrac distance matrix).

```
usage: LDA.py [-h] -m MAP_FP -g GROUP_BY [-c COLORS] [-ot OTU_TABLE] [-dm DIST_MATRIX_FILE] [--save_lda_input SAVE_LDA_INPUT] [--plot_title PLOT_TITLE] [-o OUT_FP] [-d {2,3}] [--z_angles Z_ANGLES Z_ANGLES] [--figsize FIGSIZE FIGSIZE] [--font_size FONT_SIZE] [--label_padding LABEL_PADDING] [--annotate_points] [--ggplot2_style]
```

Required arguments

- m MAP_FP, --map_fp MAP_FP**
Metadata mapping file.
- g GROUP_BY [GROUP_BY ...], --group_by GROUP_BY [GROUP_BY ...]**
A column name in the mapping file containing categorical values that will be used to identify groups. Each sample ID must have a group entry. Default is no categories and all the data will be treated as a single group.

Optional arguments

- c COLORS, --colors COLORS**
A column name in the mapping file containing hexadecimal (#FF0000) color values that will be used to color the groups. Each sample ID must have a color entry.
- ot OTU_TABLE, --otu_table OTU_TABLE**
Input biom file format OTU table.
- dm DIST_MATRIX_FILE, --dist_matrix_file DIST_MATRIX_FILE**
Input distance matrix file.
- save_lda_input SAVE_LDA_INPUT**
Save a CSV-format file of the transposed LDA-input table to the file specified by this option.
- plot_title PLOT_TITLE**
Plot title. Default is no title.
- o OUT_FP, --out_fp OUT_FP**
The path and file name to save the plot under. If specified, the figure will be saved directly instead of opening a window in which the plot can be viewed before saving.
- d {2,3}, --dimensions {2,3}**
Choose whether to plot 2D or 3D.
- z_angles Z_ANGLES Z_ANGLES**
Specify the azimuth and elevation angles for a 3D plot.
- figsize FIGSIZE FIGSIZE**
Specify the ‘width height’ in inches for LDA plots. By default, figure size is 14x8 inches.
- font_size FONT_SIZE**
Sets the font size for text elements in the plot.
- label_padding LABEL_PADDING**
Sets the spacing in points between the each axis and its label.
- annotate**
If specified, each data point will be labeled with its sample ID. Currently, only works for 2D plots. Default is False.
- ggplot2_style**
Apply ggplot2 styling to the figure. Default is False.
- h, --help**
Show the help message and exit.

Workflow for generating LDA plots using PhyloToAST

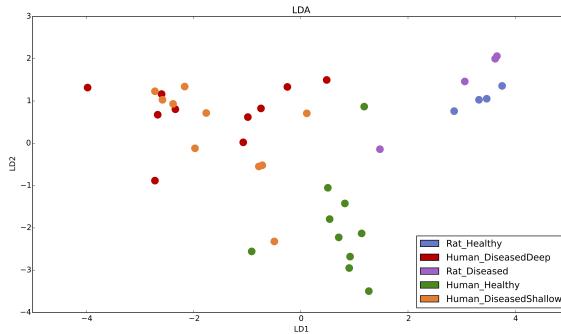
Step 1 : Create an all-pairs distance matrix for your sample data using the `beta_diversity.py` QIIME script. Different distance metrics can be calculated here: bray_curtis, morisita_horn, kulczynski, and many others.

Step 2 : Users can use either diversity distance matrix file or BIOM format file to run `LDA.py` script with all relevant parameters. If a distance matrix file is provided, LDA will use the distance matrix as an input, otherwise it will calculate OTU relative abundances from input BIOM format file and run LDA based on OTU relative abundances.

Example plots

2D LDA plot based on relative abundances with 5 metadata categories¹.

```
LDA.py -i table.biom -m mapping.txt -g Condition -c Colors
```



Citation:

2.3.4 LDA_bubble.py

This script returns LDA plots, with samples/dots sized by relative abundances of input OTU(s).

```
usage: LDA_bubble.py [-h] -i OTU_TABLE -m MAP_FP -g GROUP_BY -c COLOR_BY -ids OTU_IDS_FPS [-dm DIST_MATRIX_FILE] [--save_lda_input SAVE_LDA_INPUT] [-od OUTPUT_DIR] [--scale_by SCALE_BY] [-s SAVE_AS] [--ggplot2_style] [-v]
```

Required arguments

- i OTU_TABLE, --otu_table OTU_TABLE**
Input biom file format OTU table.
- m MAP_FP, --map_fp MAP_FP**
Metadata mapping file.
- g GROUP_BY, --group_by GROUP_BY**
A column name in the mapping file containing categorical values that will be used to identify groups. Each sample ID must have a group entry. Default is no categories and all the data will be treated as a single group.
- c COLOR_BY, --color_by COLOR_BY**
A column name in the mapping file containing hexadecimal (#FF0000) color values that will be used to color the groups. Each sample ID must have a color entry.
- ids OTU_IDS_FPS, --otu_ids_fps OTU_IDS_FPS**
Path to a file containing one OTU ID per line. One plot will be created for each OTU.

¹ Dabdoub, S. M. et al. **PhyloToAST: Bioinformatics tools for species-level analysis and visualization of complex microbial datasets**. Sci. Rep. 6, 29123; doi: 10.1038/srep29123 (2016).

Optional arguments

-dm DIST_MATRIX_FILE, --dist_matrix_file DIST_MATRIX_FILE
Input distance matrix file.

--save_lda_input SAVE_LDA_INPUT
Save a CSV-format file of the transposed LDA-input table to the file specified by this option.

-od OUTPUT_DIR, --output_dir OUTPUT_DIR
The directory to save the LDA bubble plots to. By default, plots will be saved in current working directory.

--scale_by SCALE_BY
Species relative abundance is multiplied by this factor in order to make appropriate visible bubbles in the output plots. Default scaling is 1000.

-s SAVE_AS, --save_as SAVE_AS
The type of image file for LDA plots. By default, plots will be saved in ‘svg’ format.

--ggplot2_style
Apply ggplot2 styling to the figure.

-v, --verbose
Displays species name as each is being plotted and stored to disk.

-h, --help
Show the help message and exit

2.3.5 PCoA.py

Create a series of 2D or 3D PCoA plots where the marker size varies by relative abundance of a particular OTU.

```
usage: PCoA.py [-h] -i COORD_FP -m MAP_FP -b COLORBY [-o OUT_FN] [-d {2,3}] [-t TITLE]
                [-s POINT_SIZE]
```

Required Arguments

-i COORD_FP, --coord_fp COORD_FP
Path to the principal coordinates result file (i.e., output from principal_coordinates.py).

-m MAP_FP, --map_fp MAP_FP
Path to the metadata mapping file.

-g GROUP_BY, --group_by GROUP_BY
Metadata category/categories (column headers) to group samples by in the plot. A single category can be specified as follows: -b Treatment. Multiple categories can also be specified: -b Treatment, Age. Finally, each specified category can be fixed to a single value: -b Treatment, Age, Gender=male. Note that in all cases, no spaces should be used). The program will create one group for each unique combination of values for the specified categories and put each sample in the appropriate group that matches its metadata. For example, if Treatment has two values (TreatmentA, TreatmentB) and Gender has two values (male, female), there are a total of 4 possible groups: TreatmentA and male, TreatmentA and female, TreatmentB and male, TreatmentB and female. In the output plot legend, multiple-category groups will have their values joined by an underscore: TreatmentA_male, TreatmentB_female.

Optional Arguments

-d {2,3}, --dimensions {2,3}
Choose whether to plot 2D or 3D. Default is a 2D plot.

-c COLORS, --colors COLORS
A column name in the mapping file containing hexadecimal (#FF0000) color values that will be used to color the groups. Each sample ID must have a color entry.

-s POINT_SIZE, --point_size POINT_SIZE
Specify the size of the circles representing each of the samples in the plot.

--pc_order PC_ORDER
Choose which Principle Coordinates are displayed and in which order, for example: 1,2 (Note the lack of any spaces around the comma).

--x_limits X_LIMITS X_LIMITS
Specify limits for the x-axis instead of automatic setting based on the data range. Should take the form: -x_limits -0.5 0.5

--y_limits Y_LIMITS Y_LIMITS
Specify limits for the y-axis instead of automatic setting based on the data range. Should take the form: -y_limits -0.5 0.5

--z_limits Z_LIMITS Z_LIMITS
Specify limits for the z-axis instead of automatic setting based on the data range. Should take the form: -z_limits -0.5 0.5

--z_angles Z_ANGLES Z_ANGLES
Specify the azimuth and elevation angles for a 3D plot.

-t TITLE, --title TITLE
Title of the plot.

--figsize FIGSIZE FIGSIZE
Specify the ‘width height’ in inches for PCoA plots. By default, figure size is 14x8 inches

--font_size FONT_SIZE
Sets the font size for text elements in the plot.

--label_padding LABEL_PADDING
Sets the spacing in points between the each axis and its label.

--annotate_points
If specified, each graphed point will be labeled with its sample ID.

--ggplot2_style
Apply ggplot2 styling to the figure.

-o OUT_FP, --out_fp OUT_FP
The path and file name to save the plot under. If specified, the figure will be saved directly instead of opening a window in which the plot can be viewed before saving.

-h, --help
Show the help message and exit.

Workflow for generating PCoA plots using PhyloToAST

Step 1 : Create an all-pairs distance matrix for your sample data using the `beta_diversity.py` QIIME script. Different distance metrics can be calculated here: bray_cuttis, morisita_horn, kulczynski, and many others.

Step 2 : Perform a principal coordinates analysis of the distance matrix from Step 1 using QIIME's `principal_coordinates.py` script.

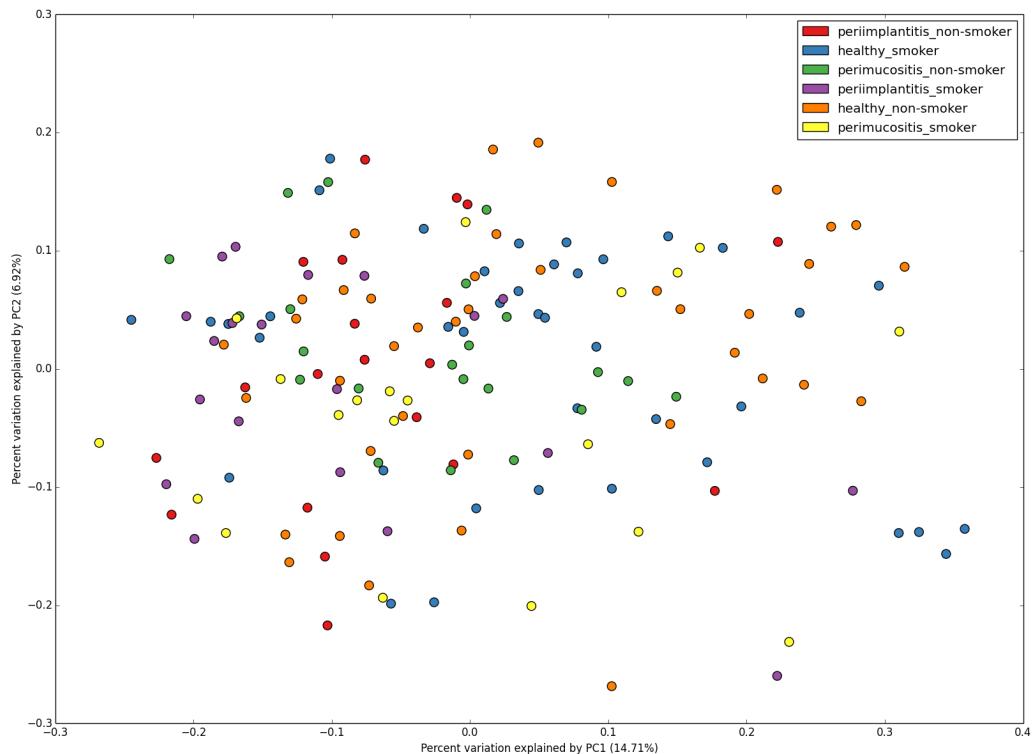
Alternate Step 1 and 2 combined for UniFrac PCoA: The `beta_diversity_through_plots.py` script produces the PCoA analysis of the UniFrac distances (weighted and unweighted) in one step.

Step 3 : Run PhyloToAST's `PCoA.py` with the input (`-i`) set to the output from Step 2.

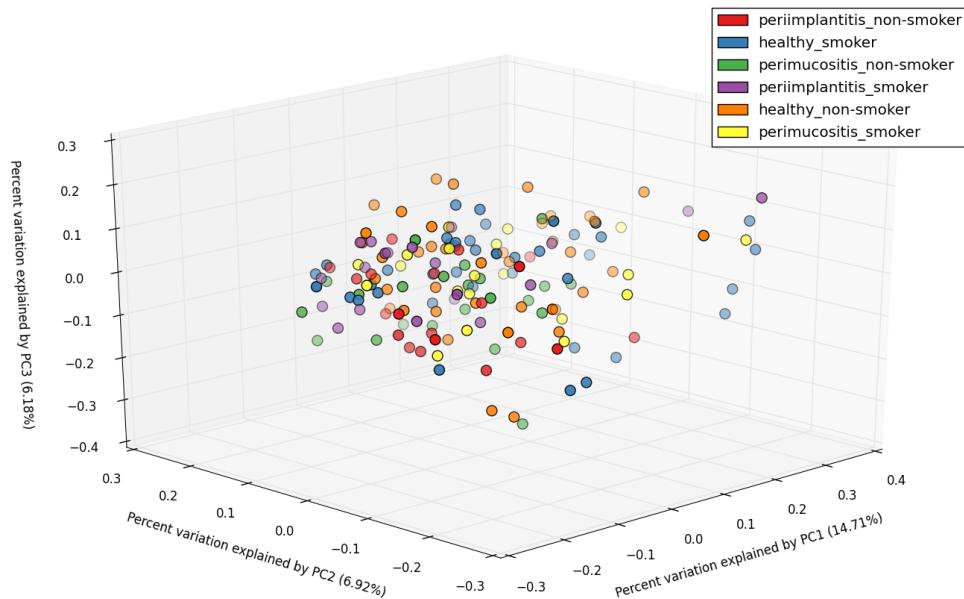
For minimum functionality, also set the mapping file (`-m`), and the grouping category column within the mapping file (`-b`). If you want to specify your own colors for the groups, also specify `-c` option. To get a 3D plot that is rotatable/zoomable specify `-d 3`.

Example plots

2D PCoA plot with 2 metadata categories - DiseaseState and SmokingStatus.



3D PCoA plot with 2 metadata categories - DiseaseState and SmokingStatus.



2.3.6 PCoA_bubble.py

Create a series of Principal Coordinate plots for each OTU in an input list where the plot points are varied in size by the relative abundance of the OTU (relative to either Sample or the total contribution of the OTU to the data set).

```
usage: PCoA_bubble.py [-h] -i OTU_TABLE -m MAPPING -pc PCOA_FP -b GROUP_BY [-c
                      ↪COLORS] -ids OTU_IDS_FP [-o OUTPUT_DIR] [-s SAVE_AS] [--scale_by SCALE_BY] [--
                      ↪ggplot2_style] [-v]
```

Required Arguments

- i OTU_TABLE, --otu_table OTU_TABLE**
The biom-format file with OTU-Sample abundance data.
- m MAPPING, --mapping MAPPING**
The mapping file specifying group information for each sample.
- pc PCOA_FP, --pcoa_fp PCOA_FP**
Principal Coordinates Analysis file. Eg. unweighted_unifrac_pc.txt, or any other output from principal_coordinates.py.
- b GROUP_BY, --group_by GROUP_BY**
Column name in mapping file specifying group information.
- c COLORS, --colors COLORS**
A column name in the mapping file containing hexadecimal (#FF0000) color values that will be used to color the groups. Each sample ID must have a color entry.
- ids OTU_IDS_FP, --otu_ids_fp OTU_IDS_FP**
Path to a file containing one OTU ID per line. One plot will be created for each OTU.

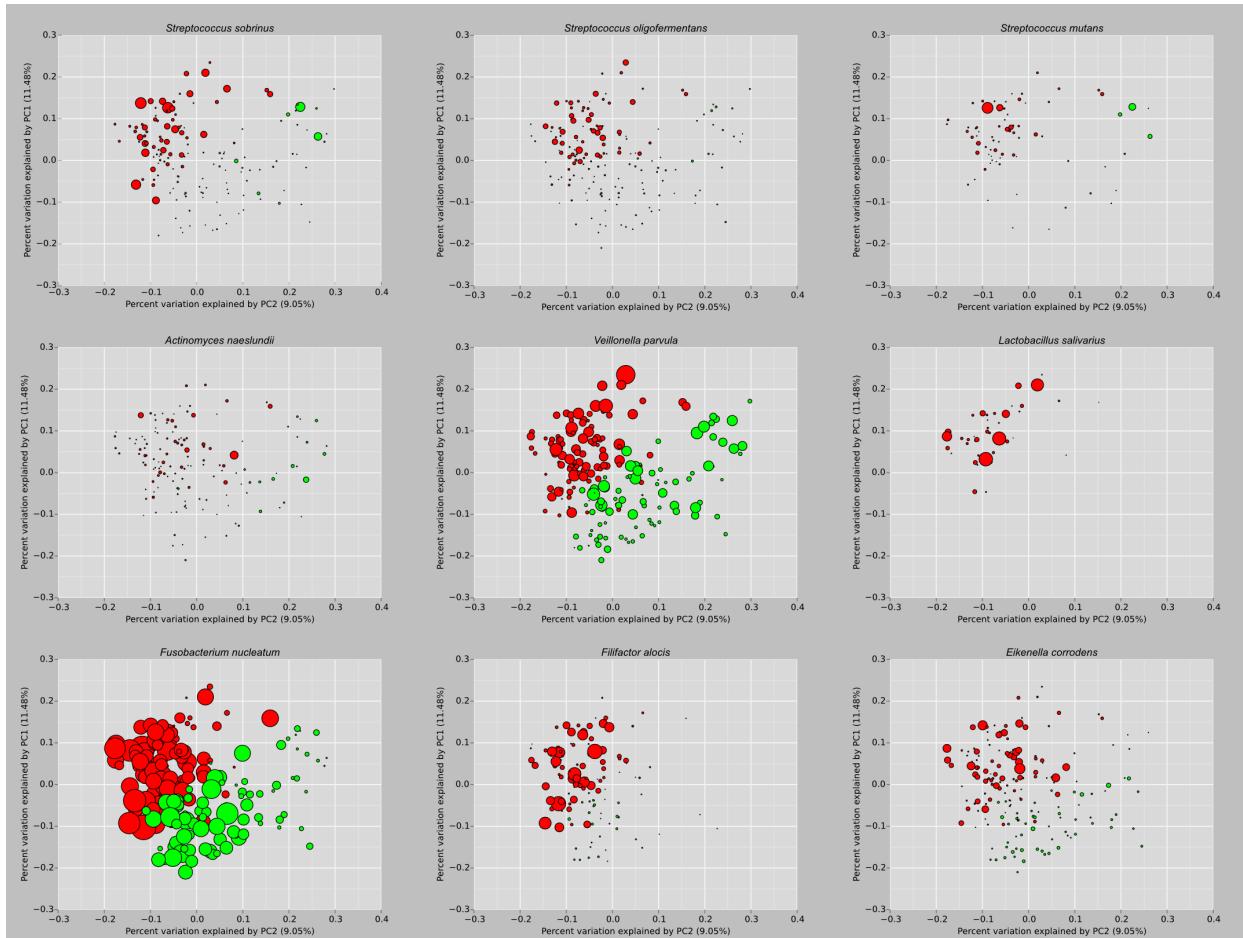
Optional Arguments

- o OUTPUT_DIR, --output_dir OUTPUT_DIR**
The directory to output the PCoA plots to.
- s SAVE_AS, --save_as SAVE_AS**
The type of image file for PCoA plots. By default, files will be saved in SVG format.
- scale_by SCALE_BY**
Species relative abundance is multiplied by this factor in order to make appropriate visible bubbles in the output plots. Default is 1000.
- ggplot2_style**
Apply ggplot2 styling to the figure.
- v, --verbose**
Displays species name as each is being plotted and saved.
- h, --help**
Show this help message and exit

Example plot

PCoA bubble plots of subgingival microbiome pathogens of smokers¹.

¹ **The Subgingival Microbiome of Clinically Healthy Current and Never Smokers.** Matthew R Mason, Philip M Preshaw, Haikady N Nagaraja, Shareef M Dabdoub, Anis Rahman and Purnima S Kumar; doi: [10.1038/ismej.2014.114](https://doi.org/10.1038/ismej.2014.114)



Citation:

2.4 Complete Script List

All available PhyloToAST scripts.

2.4.1 filter_keep_otus_by_sample

This filter allows for the removal of sequences not contained within a user-specified list of Sample IDs. This script examines each OTU and removes any sequences not originating from the specified set of allowed Sample IDs. Any empty OTUs that result are removed.

```
usage: filter_keep_otus_by_sample.py [-h] -i OTU_MAP -k SAMPLES_TO_KEEP_FP -o OUTPUT_
                                     ↵OTU_MAP_FP [-v]
```

Required Arguments

-i OTU_MAP, **--otu_map** OTU_MAP

Path to the input OTU map (i.e., the output from pick_otus.py)

- k** SAMPLES_TO_KEEP_FP, **--samples_to_keep_fp** SAMPLES_TO_KEEP_FP
Path to the file containing Sample IDs to keep in the new OTU map. One Sample ID per line.
- o** OUTPUT_OTU_MAP_FP, **--output_otu_map_fp** OUTPUT_OTU_MAP_FP
Path to the output filtered OTU map

Optional Arguments

- h, --help**
Show this help message and exit
- v, --verbose**
Specify for verbose description of the script output.

CHAPTER 3

Citing PhyloToAST

Dabdoub, S. M. et al. PhyloToAST: Bioinformatics tools for species-level analysis and visualization of complex microbial datasets. Sci. Rep. 6, 29123, 2016; doi: [10.1038/srep29123](https://doi.org/10.1038/srep29123)

CHAPTER 4

Publications using PhyloToAST

Paropkari, A. D. et al. *Smoking, Pregnancy and the Subgingival Microbiome*. Sci. Rep. 6, 30388, 2016; doi: [10.1038/srep30388](https://doi.org/10.1038/srep30388)

Tsigarida and Dabdoub et al., *The Influence of Smoking on the Peri-Implant Microbiome*. Journal of Dental Research, 2015; doi: [10.1177/0022034515590581](https://doi.org/10.1177/0022034515590581)

Mason et al., *The subgingival microbiome of clinically healthy current and never smokers*. The ISME Journal, 2014; doi: [10.1038/ismej.2014.114](https://doi.org/10.1038/ismej.2014.114)

Dabdoub et al., *Patient-specific Analysis of Periodontal and Peri-implant Microbiomes*. Journal of Dental Research, 2013; doi: [10.1177/0022034513504950](https://doi.org/10.1177/0022034513504950)

CHAPTER 5

References

CHAPTER 6

Indices and tables

- genindex
- modindex
- search

Symbols

- annotate
 - command line option, 27
 - annotate_points
 - command line option, 30
 - color_by COLOR_BY
 - command line option, 23
 - figsize FIGSIZE FIGSIZE
 - command line option, 27, 30
 - font_size FONT_SIZE
 - command line option, 27, 30
 - ggplot2_style
 - command line option, 27, 29, 30, 33
 - label_padding LABEL_PADDING
 - command line option, 27, 30
 - output_removed_otus_fn
 - PUT_REMOVED_OTUS_FN
 - command line option, 21
 - p1 P1
 - command line option, 20
 - p2 P2
 - command line option, 20
 - pc_order PC_ORDER
 - command line option, 30
 - plot_title PLOT_TITLE
 - command line option, 23, 27
 - save_calculations SAVE_CALCULATIONS
 - command line option, 24
 - save_lda_input SAVE_LDA_INPUT
 - command line option, 27, 29
 - scale_by SCALE_BY
 - command line option, 29, 33
 - scaling_factor SCALING_FACTOR
 - command line option, 18
 - show_available_metrics
 - command line option, 24
 - show_significance
 - command line option, 24
 - stabilize_variance
 - command line option, 24
 - x_limits X_LIMITS X_LIMITS
 - command line option, 30
 - y_limits Y_LIMITS Y_LIMITS
 - command line option, 30
 - z_angles Z_ANGLES Z_ANGLES
 - command line option, 27, 30
 - z_limits Z_LIMITS Z_LIMITS
 - command line option, 30
 - L {k,p,c,o,f,g,s}, -phylogenetic_level {k,p,c,o,f,g,s}
 - command line option, 14
 - a {MRA,NMRA,raw}, {MRA,NMRA,raw}
 - command line option, 24
 - b GROUP_BY, -group_by GROUP_BY
 - command line option, 32
 - c CATEGORY, -category CATEGORY
 - command line option, 23
 - c COLORS, -colors COLORS
 - command line option, 27, 30, 32
 - c COLOR_BY, -color_by COLOR_BY
 - command line option, 28
 - c MAP_CATEGORIES, MAP_CATEGORIES
 - command line option, 24
 - c MAP_CATEGORY, MAP_CATEGORY
 - command line option, 23
 - d DIVERSITY [DIVERSITY ...], -diversity DIVERSITY [DIVERSITY ...]
 - command line option, 23
 - d {2,3}, -dimensions {2,3}
 - command line option, 27, 30
 - dm DIST_MATRIX_FILE, DIST_MATRIX_FILE
 - command line option, 27, 29
 - e OUTPUT_TRE, -output_tre OUTPUT_TRE
 - command line option, 24
 - fo FILTER_OTUIDS_FNH, -filter_otuids_fnh FILTER_OTUIDS_FNH
 - command line option, 24
- OUT-

	command line option, 16	
-fp FIL_PCT, -fil_pct FIL_PCT	command line option, 18	
-g GROUP_BY [GROUP_BY ...], -group_by GROUP_BY [GROUP_BY ...]	command line option, 27	
-g GROUP_BY, -group_by GROUP_BY	command line option, 28, 29	
-go GEXF_OUT, -gexf_out GEXF_OUT	command line option, 18	
-h, -help	command line option, 13–23, 25, 27, 29, 30, 33, 35	
-i ASSIGNED_TAXONOMY_FN, assigned_taxonomy_fn SIGNED_TAXONOMY_FN	AS-	
	command line option, 14	
-i BIOM_FP, -biom_fp BIOM_FP	command line option, 21, 23	
-i COORD_FP, -coord_fp COORD_FP	command line option, 29	
-i INPUT_ASSIGNED_TAXONOMY, input_assigned_taxonomy PUT_ASSIGNED_TAXONOMY	IN-	
	command line option, 19	
-i INPUT BIOM FP, -input_biom_fp PUT_BIOM_FP	IN-	
	command line option, 14, 22	
-i INPUT_FASTA_FN, -input_fasta_fn PUT_FASTA_FN	IN-	
	command line option, 13, 22	
-i INPUT_FNA [INPUT_FNA ...], -input_fna PUT_FNA [INPUT_FNA ...]	IN-	
	command line option, 17	
-i OTU_ID_FP, -otu_id_fp OTU_ID_FP		
	command line option, 19	
-i OTU_MAP, -otu_map OTU_MAP		
	command line option, 34	
-i OTU_TABLE, -otu_table OTU_TABLE		
	command line option, 24, 28, 32	
-i REP_SET_FP, -rep_set_fp REP_SET_FP		
	command line option, 12	
-i SEQS_OTUS_FN, -seqs_otus_fn SEQS_OTUS_FN		
	command line option, 21	
-ids OTU_IDS_FP, -otu_ids_fp OTU_IDS_FP		
	command line option, 28, 32	
-k SAMPLES_TO_KEEP_FP, -samples_to_keep_fp SAMPLES_TO_KEEP_FP		
	command line option, 34	
-l {k,p,c,o,f,g,s}, -phylogenetic_level {k,p,c,o,f,g,s}		
	command line option, 19, 21	
-m MAPPING, -mapping MAPPING		
	command line option, 22, 24, 32	
-m MAPPING_FN, -mapping_fn MAPPING_FN		
	command line option, 13	
		-m MAP_FILE, -map_file MAP_FILE
		command line option, 23
		-m MAP_FP, -map_fp MAP_FP
		command line option, 27–29
		-n JOB_NAME, -job_name JOB_NAME
		command line option, 17
		-n NON_UNIQUE_OTU_MATRIX,
		–non_unique_otu_matrix
		NON_UNIQUE_OTU_MATRIX
		command line option, 20
		-n NON_UNIQUE_OUTPUT_FILE,
		–non_unique_output_file
		NON_UNIQUE_OUTPUT_FILE
		command line option, 19
		-n NUM_OUTPUT_FILES, -num_output_files
		NUM_OUTPUT_FILES
		command line option, 22
		-o ASSIGNED_TAXONOMY_FP,
		assigned_taxonomy_fp AS-
		SIGNED_TAXONOMY_FP
		command line option, 13
		-o CONDENSED_SEQS_OTUS_FILE,
		condensed_seqs_otus_file CON-
		DENSED_SEQS_OTUS_FILE
		command line option, 20
		-o OUTPUT_BIOM_FP, -output_biom_fp OUT-
		PUT_BIOM_FP
		command line option, 22
		-o OUTPUT_CSV_FP, -output_csv_fp OUT-
		PUT_CSV_FP
		command line option, 14
		-o OUTPUT_DIR, -output_dir OUTPUT_DIR
		command line option, 15, 22, 33
		-o OUTPUT_FILTERED REP_SET_FN, -
		output_filtered_rep_set_fn OUT-
		PUT_FILTERED REP_SET_FN
		command line option, 16
		-o OUTPUT_FN, -output_fn OUTPUT_FN
		command line option, 17
		-o OUTPUT_FP, -output_fp OUTPUT_FP
		command line option, 19, 20
		-o OUTPUT_ITOL_TABLE, -output_itol_table OUT-
		PUT_ITOL_TABLE
		command line option, 24
		-o OUTPUT_OTU_MAP_FP, -output_otu_map_fp OUT-
		PUT_OTU_MAP_FP
		command line option, 35
		-o OUTPUT_PREFIX, -output_prefix OUT-
		PUT_PREFIX
		command line option, 13
		-o OUTPUT_PRUNED_OTUS_FN, -
		output_pruned_otus_fn OUT-
		PUT_PRUNED_OTUS_FN
		command line option, 21

-o OUT_DIR, --out_dir OUT_DIR
 command line option, 23

-o OUT_FP, --out_fp OUT_FP
 command line option, 27, 30

-o REPSET_OUT_FP, --repset_out_fp
 REPSET_OUT_FP
 command line option, 22

-od OUTPUT_DIR, --output_dir OUTPUT_DIR
 command line option, 29

-ot OTU_TABLE, --otu_table OTU_TABLE
 command line option, 27

-p IMAGE_TYPE, --image_type IMAGE_TYPE
 command line option, 24

-p PERCENT_OF_SAMPLES, --percent_of_samples
 PERCENT_OF_SAMPLES
 command line option, 21

-p PREFIX, --prefix PREFIX
 command line option, 15

-p PRUNED_OUTPUT_FILE, --pruned_output_file
 PRUNED_OUTPUT_FILE
 command line option, 19

-pc PCOA_FP, --pcoa_fp PCOA_FP
 command line option, 32

-q QUALITY_FN, --quality_fn QUALITY_FN
 command line option, 13

-r REPSET_FP, --repset_fp REPSET_FP
 command line option, 21

-r REP_SET_FN, --rep_set_fn REP_SET_FN
 command line option, 14, 16

-r REVERSE, --reverse REVERSE
 command line option, 15

-s PERCENT_OF_SEQUENCES, --percent_of_sequences
 CENT_OF_SEQUENCES
 command line option, 21

-s POINT_SIZE, --point_size POINT_SIZE
 command line option, 30

-s SAVE_AS, --save_as SAVE_AS
 command line option, 29, 33

-s SEQS_OTUS, --seqs_otus SEQS_OTUS
 command line option, 20

-s SEQS_OTUS_FN, --seqs_otus_fn SEQS_OTUS_FN
 command line option, 14

-t ID_TO_TAXONOMY_FN, --id_to_taxonomy_fn
 ID_TO_TAXONOMY_FN
 command line option, 21

-t ID_TO_TAXONOMY_FP, --id_to_taxonomy_fp
 ID_TO_TAXONOMY_FP
 command line option, 12

-t INPUT_TREE, --input_tree INPUT_TREE
 command line option, 24

-t TAXONOMY_FP, --taxonony_fp TAXONOMY_FP
 command line option, 19

-t TITLE, --title TITLE

command line option, 30

-t WALLTIME, --walltime WALLTIME
 command line option, 17

-t, --test
 command line option, 18

-u UNIQUE_OTUS_FN, --unique_otus_fn
 UNIQUE_OTUS_FN
 command line option, 16

-v, --verbose
 command line option, 13, 14, 16, 17, 19–23, 29, 33, 35

-w STATS_OUT_FNH, --stats_out_fnh
 STATS_OUT_FNH
 command line option, 18

A

ax:
 command line option, 12

B

biom_file
 command line option, 18

biomf:
 command line option, 6, 7

biomfile:
 command line option, 8

C

cat_name
 command line option, 18

categories:
 command line option, 9

category_column
 command line option, 15

color:
 command line option, 12

command line option
 --annotate, 27
 --annotate_points, 30
 --color_by COLOR_BY, 23
 --figsize FIGSIZE FIGSIZE, 27, 30
 --font_size FONT_SIZE, 27, 30
 --ggplot2_style, 27, 29, 30, 33
 --label_padding LABEL_PADDING, 27, 30
 --output_removed_otus_fn
 PUT_REMOVED_OTUS_FN, 21

OUT-
 --p1 P1, 20
 --p2 P2, 20
 --pc_order PC_ORDER, 30
 --plot_title PLOT_TITLE, 23, 27

SAVE_CALCULATIONS, 24

SAVE_LDA_INPUT, 27, 29

SCALE_BY, 29, 33

SCALING_FACTOR, 18

```

--show_available_metrics, 24
--show_significance, 24
--stabilize_variance, 14, 25
--x_limits X_LIMITS X_LIMITS, 30
--y_limits Y_LIMITS Y_LIMITS, 30
--z_angles Z_ANGLES Z_ANGLES, 27, 30
--z_limits Z_LIMITS Z_LIMITS, 30
-L {k,p,c,o,f,g,s}, --phylogenetic_level
{k,p,c,o,f,g,s}, 14
-a {MRA,NMRA,raw}, --analysis_metric
{MRA,NMRA,raw}, 24
-b GROUP_BY, --group_by GROUP_BY, 32
-c CATEGORY, --category CATEGORY, 23
-c COLORS, --colors COLORS, 27, 30, 32
-c COLOR_BY, --color_by COLOR_BY, 28
-c MAP_CATEGORIES, --map_categories
MAP_CATEGORIES, 24
-c MAP_CATEGORY, --map_category
MAP_CATEGORY, 23
-d DIVERSITY [DIVERSITY ...], --diversity DI-
VERSITY [DIVERSITY ...], 23
-d {2,3}, --dimensions {2,3}, 27, 30
-dm DIST_MATRIX_FILE, --dist_matrix_file
DIST_MATRIX_FILE, 27, 29
-e OUTPUT_TRE, --output_tre OUTPUT_TRE, 24
-fo FILTER_OTUIDS_FNH, --filter_otuids_fnh FIL-
TER_OTUIDS_FNH, 16
-fp FIL_PCT, --fil_pct FIL_PCT, 18
-g GROUP_BY [GROUP_BY ...], --group_by
GROUP_BY [GROUP_BY ...], 27
-g GROUP_BY, --group_by GROUP_BY, 28, 29
-go GEXF_OUT, --gexf_out GEXF_OUT, 18
-h, --help, 13–23, 25, 27, 29, 30, 33, 35
-i ASSIGNED_TAXONOMY_FN,
--assigned_taxonomy_fn AS-
SIGNED_TAXONOMY_FN, 14
-i BIOM_FP, --biom_fp BIOM_FP, 21, 23
-i COORD_FP, --coord_fp COORD_FP, 29
-i INPUT_ASSIGNED_TAXONOMY,
--input_assigned_taxonomy IN-
PUT_ASSIGNED_TAXONOMY, 19
-i INPUT_BIOM_FP, --input_biom_fp IN-
PUT_BIOM_FP, 14, 22
-i INPUT_FASTA_FN, --input_fasta_fn IN-
PUT_FASTA_FN, 13, 22
-i INPUT_FNA [INPUT_FNA ...], --input_fna IN-
PUT_FNA [INPUT_FNA ...], 17
-i OTU_ID_FP, --otu_id_fp OTU_ID_FP, 19
-i OTU_MAP, --otu_map OTU_MAP, 34
-i OTU_TABLE, --otu_table OTU_TABLE, 24, 28,
32
-i REP_SET_FP, --rep_set_fp REP_SET_FP, 12
-i SEQS_OTUS_FN, --seqs_ots_fn
SEQS_OTUS_FN, 21
-ids OTU_IDS_FP, --otu_ids_fp OTU_IDS_FP, 28,
32
-k SAMPLES_TO_KEEP_FP, --samples_to_keep_fp SAM-
PLES_TO_KEEP_FP, 34
-l {k,p,c,o,f,g,s}, --phylogenetic_level
{k,p,c,o,f,g,s}, 19, 21
-m MAPPING, --mapping MAPPING, 22, 24, 32
-m MAPPING_FN, --mapping_fn MAPPING_FN,
13
-m MAP_FILE, --map_file MAP_FILE, 23
-m MAP_FP, --map_fp MAP_FP, 27–29
-n JOB_NAME, --job_name JOB_NAME, 17
-n NON_UNIQUE_OTU_MATRIX,
--non_unique_otu_matrix NON_UNIQUE_OTU_MATRIX, 20
-n NON_UNIQUE_OUTPUT_FILE,
--non_unique_output_file NON_UNIQUE_OUTPUT_FILE, 19
-n NUM_OUTPUT_FILES, --num_output_files
NUM_OUTPUT_FILES, 22
-o ASSIGNED_TAXONOMY_FP,
--assigned_taxonomy_fp AS-
SIGNED_TAXONOMY_FP, 13
-o CONDENSED_SEQS_OTUS_FILE,
--condensed_seqs_ots_file CON-
DENSED_SEQS_OTUS_FILE, 20
-o OUTPUT_BIOM_FP, --output_biom_fp OUT-
PUT_BIOM_FP, 22
-o OUTPUT_CSV_FP, --output_csv_fp OUT-
PUT_CSV_FP, 14
-o OUTPUT_DIR, --output_dir OUTPUT_DIR, 15,
22, 33
-o OUTPUT_FILTERED REP_SET_FN,
--output_filtered_rep_set_fn OUT-
PUT_FILTERED REP_SET_FN, 16
-o OUTPUT_FN, --output_fn OUTPUT_FN, 17
-o OUTPUT_FP, --output_fp OUTPUT_FP, 19, 20
-o OUTPUT_ITOL_TABLE, --output_itol_table
OUTPUT_ITOL_TABLE, 24
-o OUTPUT_OTU_MAP_FP, --output_otu_map_fp
OUTPUT_OTU_MAP_FP, 35
-o OUTPUT_PREFIX, --output_prefix OUT-
PUT_PREFIX, 13
-o OUTPUT_PRUNED_OTUS_FN,
--output_pruned_ots_fn OUT-
PUT_PRUNED_OTUS_FN, 21
-o OUT_DIR, --out_dir OUT_DIR, 23
-o OUT_FP, --out_fp OUT_FP, 27, 30
-o REPSET_OUT_FP, --repset_out_fp
REPSET_OUT_FP, 22
-od OUTPUT_DIR, --output_dir OUTPUT_DIR, 29
-ot OTU_TABLE, --otu_table OTU_TABLE, 27
-p IMAGE_TYPE, --image_type IMAGE_TYPE, 24

```

-p PERCENT_OF_SAMPLES, percent_of_samples
 CENT_OF_SAMPLES, 21
 -p PREFIX, --prefix PREFIX, 15
 -p PRUNED_OUTPUT_FILE, --pruned_output_file
 PRUNED_OUTPUT_FILE, 19
 -pc PCOA_FP, --pcoa_fp PCOA_FP, 32
 -q QUALITY_FN, --quality_fn QUALITY_FN, 13
 -r REPSET_FP, --repset_fp REPSET_FP, 21
 -r REP_SET_FN, --rep_set_fn REP_SET_FN, 14, 16
 -r REVERSE, --reverse REVERSE, 15
 -s PERCENT_OF_SEQUENCES, percent_of_sequences
 CENT_OF_SEQUENCES, 21
 -s POINT_SIZE, --point_size POINT_SIZE, 30
 -s SAVE_AS, --save_as SAVE_AS, 29, 33
 -s SEQS_OTUS, --seqs_otus SEQS_OTUS, 20
 -s SEQS_OTUS_FN, --seqs_otus_fn
 SEQS_OTUS_FN, 14
 -t ID_TO_TAXONOMY_FN, --id_to_taxonomy_fn
 ID_TO_TAXONOMY_FN, 21
 -t ID_TO_TAXONOMY_FP, --id_to_taxonomy_fp
 ID_TO_TAXONOMY_FP, 12
 -t INPUT_TREE, --input_tree INPUT_TREE, 24
 -t TAXONOMY_FP, --taxonomy_fp TAXONOMY_FP, 19
 -t TITLE, --title TITLE, 30
 -t WALLTIME, --walltime WALLTIME, 17
 -t, --test, 18
 -u UNIQUE_OTUS_FN, --unique_otus_fn
 UNIQUE_OTUS_FN, 16
 -v, --verbose, 13, 14, 16, 17, 19–23, 29, 33, 35
 -w STATS_OUT_FNH, --stats_out_fnh
 STATS_OUT_FNH, 18
 ax:, 12
 biom_file, 18
 biomf:, 6, 7
 biomfile:, 8
 cat_name, 18
 categories:, 9
 category_column, 15
 color:, 12
 condition_column, 18
 core_fp:, 8
 d:, 8
 data:, 12
 fastaFNH:, 9, 11
 file_data, 10, 11
 fill_bt:, 12
 fn:, 7
 fnh:, 8
 header:, 9, 12
 idtaxFNH:, 10
 imap:, 9

- in_corr_mat, 18
 input_biom_fnh, 15
 input_biom_fp, 15
 items:, 12
 job_scripts, 17
 level:, 11
 mapFNH:, 9, 12
 mapping_file, 15, 18
 mapping_fnh, 15
 mode:, 8
 otuIDs:, 6
 output_biom_fnh, 15
 p:, 11
 pick_otus_results, 16
 rel_abd:, 5, 6
 return:, 5–12
 Returns an opened file for appropriate usage., 8
 sample_abd:, 6, 7
 sampleIDs:, 6, 7
 tax:, 7
 title:, 12
 unifrac:, 10, 11
 unifracFNH:, 10
 condition_column
 command line option, 18
 core_fp:
 command line option, 8

D

d:
 command line option, 8

data:
 command line option, 12

F

fastaFNH:
 command line option, 9, 11

file_data
 command line option, 10, 11

fill_bt:
 command line option, 12

fn:
 command line option, 7

fnh:
 command line option, 8

H

header:
 command line option, 9, 12

I

idtaxFNH:
 command line option, 10

imap:
 command line option, 9

in_corr_mat
 command line option, 18

input_biom_fnh
 command line option, 15

input_biom_fp
 command line option, 15

items:
 command line option, 12

J

job_scripts
 command line option, 17

L

level:
 command line option, 11

M

mapFNH:
 command line option, 9, 12

mapping_file
 command line option, 15, 18

mapping_fnh
 command line option, 15

mode:
 command line option, 8

O

otuIDs:
 command line option, 6

output_biom_fnh
 command line option, 15

P

p:
 command line option, 11

pick_otus_results
 command line option, 16

R

rel_abd:
 command line option, 5, 6

return:
 command line option, 5–12

Returns an opened file for appropriate usage.
 command line option, 8

S

sample_abd:
 command line option, 6, 7

sampleIDs:

command line option, 6, 7

T

tax:
 command line option, 7

title:
 command line option, 12

U

unifrac:
 command line option, 10, 11

unifracFN:
 command line option, 10