**Projection with Public Data (PPD)**

# Goal

- To compare users' 16S rRNA data with published datasets by processing and normalization them together, and projecting into 3D PCoA plot for visual comparative analysis

Starting from marker gene abundance data (OTU table, BIOM file, mothur output)

**Marker Data Profiling (MDP)**

**Shotgun Data Profiling (SDP)**

Starting from gene list or gene abundance data annotated by KO, EC or COG

**Click here to start**

Visually exploring your 16S rRNA data with a public data in a 3D PCoA plot

**Projection with Public Data (PPD)**

**Taxon Set Enrichment Analysis (TSEA)**

Starting with a list of taxa of interest (strains, species or higher level taxa)

# 1. Data Upload



**Note**: Please check "Format" section of MicrobiomeAnalyst web server for details about each format.

# 2. a) Data Integrity Check



| | |
|---|---|
| Data type: | OTU abundance table |
| File format: | text |
| OTU annotation: | greengene_id |
| OTU number: | 5138 |
| OTU with ≥ 2 counts: | 5138 |
| Sample number: | 73 |
| Number of experimental factors: | 1 |
| Total read counts: | 372013 |
| Average counts per sample: | 5096 |
| Maximum counts per sample: | 37579 |
| Minimum counts per sample: | 962 |

Text Summary    Graphic Summary

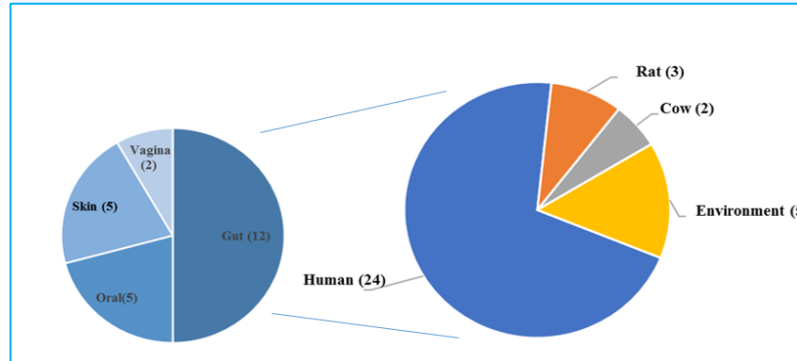**Click "Proceed" if all information is correct.**

- Provides the information summary of users uploaded data and also check whether all files necessary for analysis are present or not.

# 2. b) Graphic Summary



- Provides user the information about library size or total number of reads present in of each sample and help in identifying the potential outliers due to undersampling or sequencing errors.

# 3. Reference study selection



- In MicrobiomeAnalyst, human reference studies data has been downloaded from publicly accessible database (**QIITA**).

- These reference datasets has been primarily separated based on **the biom or site of sample collection**. Other than human, three others biom (including cow, mouse and environment) related datasets are available.

- Further, the datasets collected from human are partitioned based on **sampling body sites (4 body sites: gut ,oral  etc.)** , **sequencing platform** (3 platforms: Illumina HiSeq 2000 , 454 GS FLX etc.), target region ( 4 regions : V1, V2-V3 etc.) and **factors shaping microbiota**. (8 factors: **diet**, **host genetics**, **cultural traditions and geography**, **age**, **pregnancy** etc.)

- Currently, we have around **34** datasets derived from **~20** large-scale studies.

- Users should try to match experimental protocols used in their data generation with the reference data for meaningful comparison.

# 3. Reference study selection



Step 1 : Select the sampling site or biom of interest of reference study (studies from 4 biom)

| Studies | Target region | Sequence platform | No. of samples | Ref. |
|---|---|---|---|---|
| Healthy_whole_body | V2 | 454 GS FLX | 45 | Costello et al. 2009 |
| Dense_timeseries | V4 | Illumina HiSeq 2000 | 467 | Caporaso et al. 2011 |
| HMP_V35 | V3-5 | 454 GS FLX Titanium | 371 | HMP 2012 Consortium |
| HMP_V13 | V1-3 | 454 GS FLX Titanium | 204 | HMP 2012 Consortium |
| Global_gut | V4 | Illumina HiSeq 2000 | 528 | Yatsunenko et al. 2012 |
| Family_study | V2 | Illumina HiSeq 2000 | 169 | Song et al. 2013 |
| Diet_enterotype | V2 | 454 GS FLX Titanium | 85 | Wu et al. 2011 |
| Pregnant_women | V2 | 454 GS FLX and GS FLX Titanium | 667 | Koren et al. 2011 |
| Newborns_and_mothers | V2 | 454 GS FLX | 80 | Dominguez-Bello et al. 2010 |
| US_infant_timeseries | V2 | 454 GS FLX | 61 | Koenig et al. 2011 |
| Obese_twins | V2 | 454 GS FLX | 281 | Turnbaugh et al. 2009 |
| IBD_twins | V2 | 454 GS FLX | 114 | Willing et al. 2010 |

Tabs: Human Gut | Human Skin | Human Oral | Human Vagina | Mouse | Cow | Environmental

Step 2 : Select the reference study you want to compare with your data. (Reference data: Greengenes id annotated)

Submit

Step 3 : Click "Submit" to proceed

# 4. Data processing

**This is the most critical step in which:**

- User uploaded data will be merge user with selected reference data on the basis of common features (taxonomic labels). Even though, reference data are annotated with **Greengenes** database , user can also upload **SILVA** annotated data. (internal mapping will be performed on the basis of common **GenBank** id)

- There has to be significant features overlap (20%) between user and reference data for meaningful comparison. (otherwise, you can't proceed)

- Merged data is processed and normalized together.

- User can choose from multiple distance matrices as well as number of features to keep for performing **PCoA** comparative analysis.

# 5. Visual exploring the result
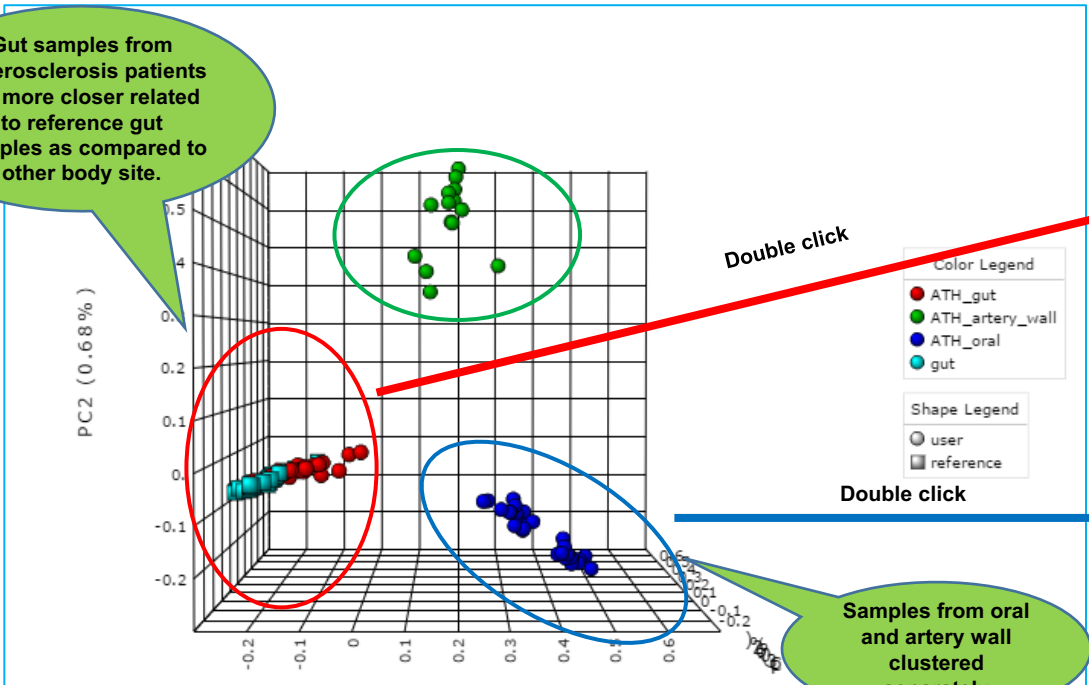


**Data Projection View**

Experimental factor — Group

Distance measure — Bray-Curtis Index

No. of features to use — Top 20% most abundant

Submit

**Distance measure to compute dissimilarity between samples**

**No. of features to use for computing dissimilarity (either all or top 20 % most abundant)**

Taxonomic level: Phylum — Update

Merge small taxa (as Others): percentage < 0.0 (0.0 - 1)

**Gut samples from Atherosclerosis patients are more closer related to reference gut samples as compared to other body site.**
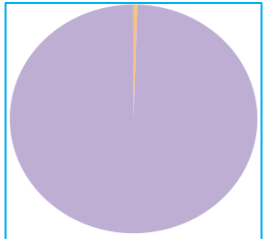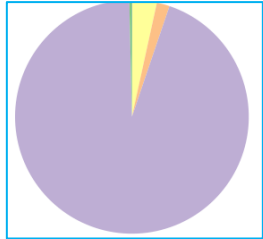
PC2 (0.68 %)

Double click

Color Legend
- ATH_gut
- ATH_artery_wall
- ATH_oral
- gut

Shape Legend
- user
- reference

Double click

**Samples from oral and artery wall clustered separately**

Gut

Oral

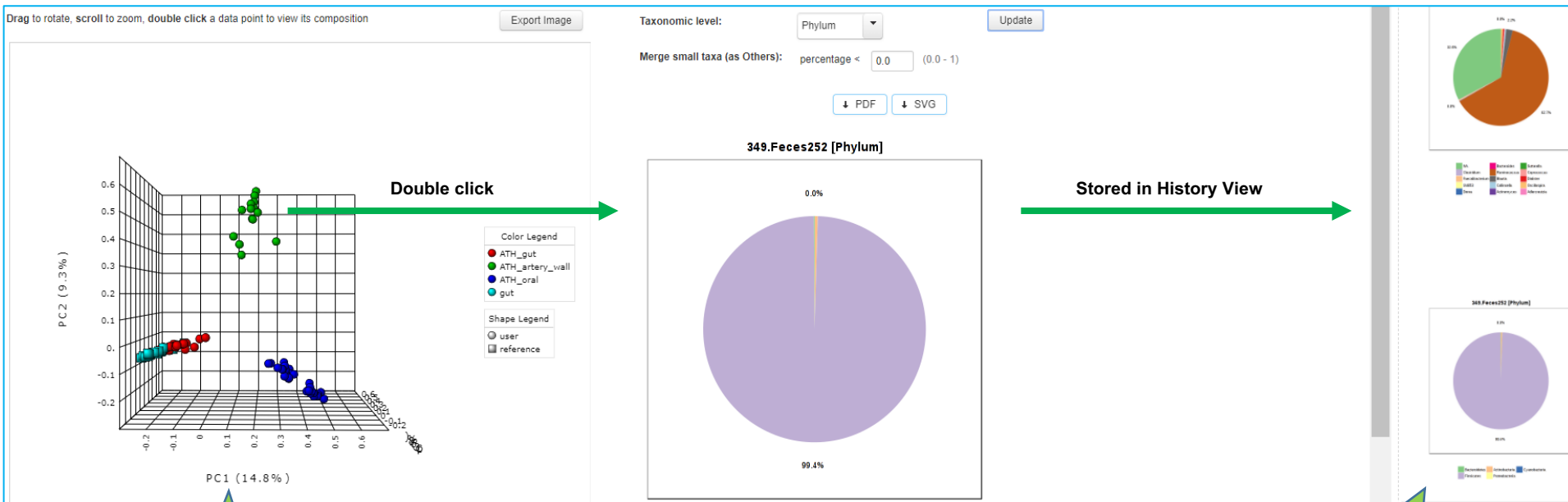**Phylum-level composition for samples**

- Result is visualized as **interactive 3D** plot in which user can found some interesting patterns in their data in comparison with reference (e.g. **healthy** vs **disease** samples clustering separately or **body sites** clustered separately)
- Users can double click individual data points (sample) to see the corresponding distributions of the taxa that underlying these difference at all possible taxonomic levels.

# 5. Visual exploring the result



**Drag** to rotate, **scroll** to zoom, **double click** a data point to view its composition

Export Image

Double click

PC2 (9.3%)

Color Legend
- ATH_gut
- ATH_artery_wall
- ATH_oral
- gut

Shape Legend
- user
- reference

PC1 (14.8%)

Taxonomic level: Phylum

Merge small taxa (as Others): percentage < 0.0 (0.0 - 1)

Update

↓ PDF   ↓ SVG

349.Feces252 [Phylum]

0.0%

99.4%

Stored in History View

349.Feces252 [Phylum]

**3 D PCoA plot**

**current selected or double clicked sample composition**

**History View**
**all the images of the samples that have been inspected.**

==END==