# Bioinformatics for microbiome research
# Day 2: microbial community analysis

Jyväskylä Summer School 2023
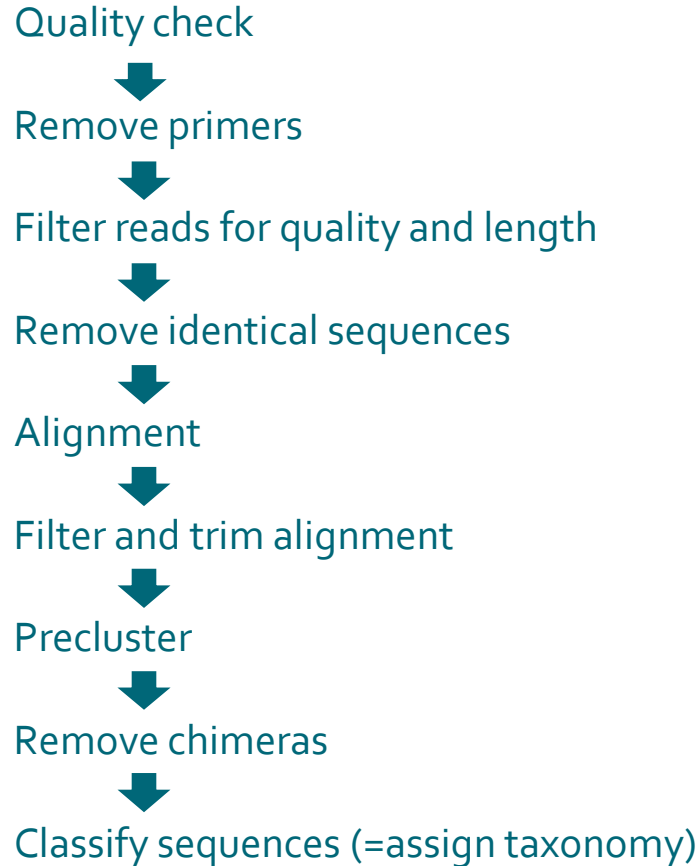
Heli Juottonen and Eija Korpelainen (slides modified from: Jesse Harrison)
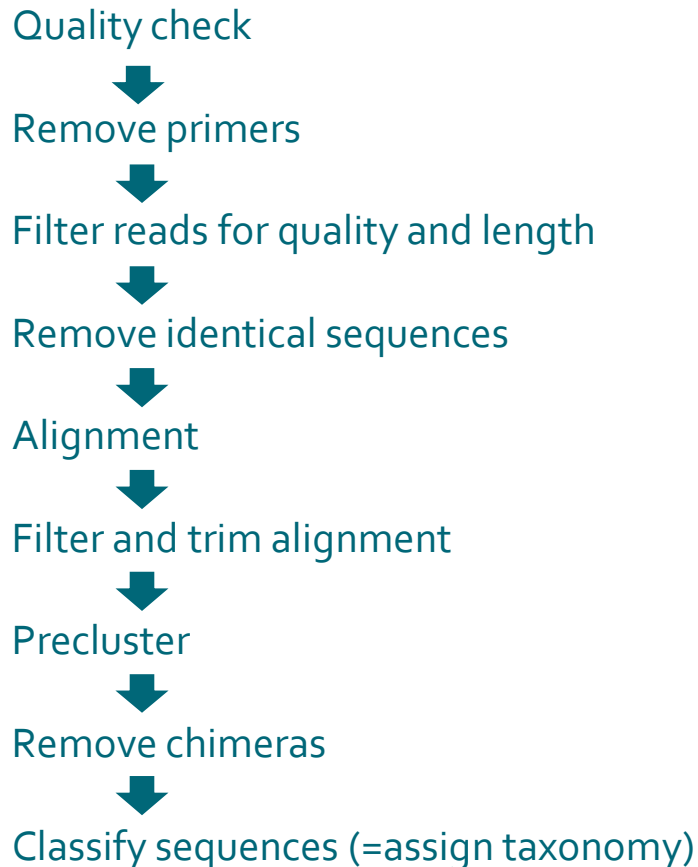CSC – IT Center for Science Ltd.
chipster@csc.fi

CSC
ICT Solutions for Brilliant Minds

# Outline of Day 1: Preprocessing reads

Quality check

⬇

Remove primers

⬇

Filter reads for quality and length

⬇

Remove identical sequences

⬇

Alignment

⬇

Filter and trim alignment

⬇

Precluster

⬇

Remove chimeras

⬇

Classify sequences (=assign taxonomy)

**Output so far:**
1. FASTA file of processed reads
2. count file (which read in which sample)
3. taxonomy file (taxonomy of each read)

# Outline of Day 1: Preprocessing reads

Quality check

⬇

Remove primers

⬇

Filter reads for quality and length

⬇

Remove identical sequences

⬇

Alignment

⬇

Filter and trim alignment

⬇

Precluster

⬇

Remove chimeras

⬇

Classify sequences (=assign taxonomy)

**Output so far:**
1. FASTA file of processed reads
2. count file (which read in which sample)
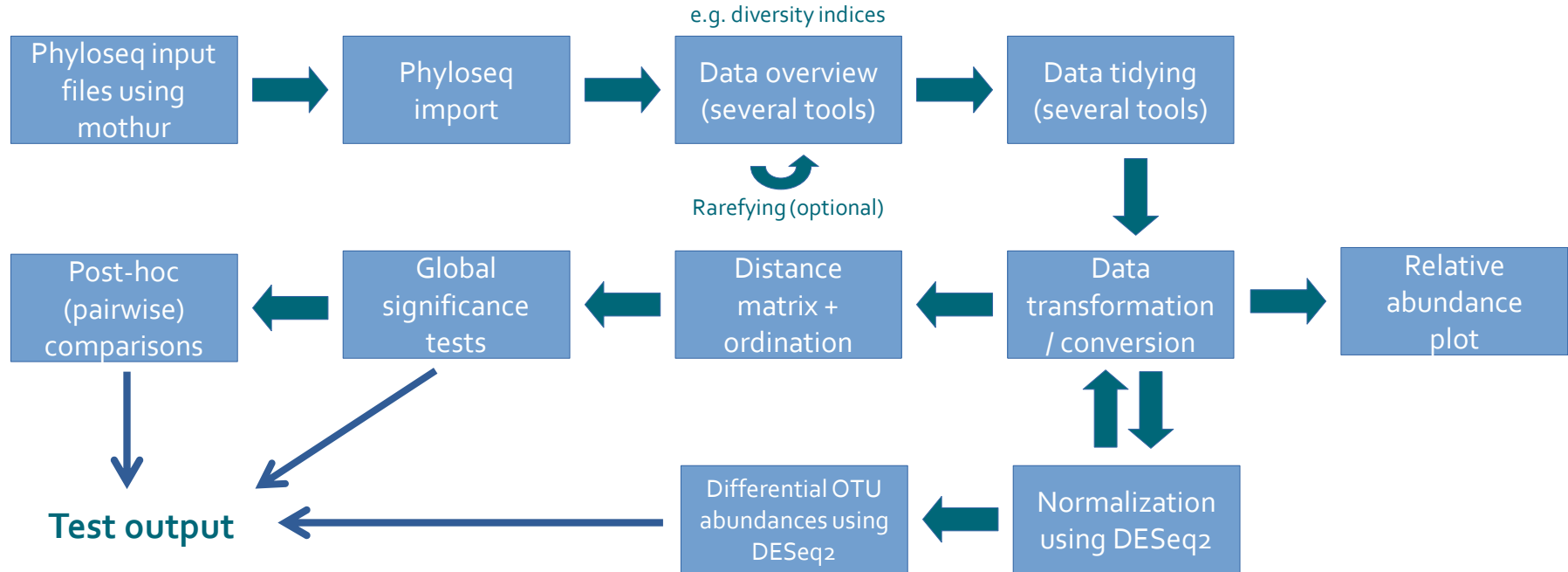3. taxonomy file (taxonomy of each read)

⬇

**Day 2: Community analysis**
- Clustering into OTUs
- Import into **Phyloseq**
- Data tidying & transformations
- Taxonomy plots
- Alpha diversity
- Beta diversity: ordinations & statistics

**Part 1: Tool overview and data importing**

# Workflow for microbial community analysis in Chipster

# Generating input files for phyloseq

Phyloseq is a multi-use R package for microbial community data processing and analysis

https://joey711.github.io/phyloseq/

# Generating input files for phyloseq

**Specifications for creating phyloseq input files:**

- type of data (16S/18S or ITS)

- % cut off for <u>OTU clustering</u>

- files produced by mothur
  - o final FASTA
  - o count file
  - o taxonomy file = taxonomy assignment of each <u>read</u>



Generate input files for phyloseq

**Parameters**                                    ↺ Reset All

Type of data                                      | 16S or 18S ▾ |
Indicate if you have ITS data as it is treated differently.

Cutoff                                            | 0.03 ▾ |
Dissimilarity threshold for OTU clustering, e.g. a cut-off value of 0.03
corresponds to 97% similarity

**Input files**

FASTA file                                        | chimeras.removed.fasta.gz ▾ |

Mothur count file                                 | chimeras.removed.count_table ▾ |

Sequences taxonomy assignment file                | sequences-taxonomy-assignmer ▾ |

# Generating input files for phyloseq

**Generated input files:**

- .shared file (mothur file format)
  - samples in rows, OTUs in columns
  - how many reads of each OTU in each sample (OTU table)

- consensus taxonomy file
  - taxonomy assignments of OTUs

- phenodata file

# Phenodata file: fill in sample information



The phenodata file is an editable table with

1) unique IDs for each sample

2) information on sample groupings

# Converting input files into a phyloseq object

# Converting input files into a phyloseq object

```
### Imported phyloseq object ###


phyloseq-class experiment-level object
otu_table()   OTU Table:          [ 1437 taxa and 16 samples ]
sample_data() Sample Data:        [ 16 samples by 6 sample variables ]
tax_table()   Taxonomy Table:     [ 1437 taxa by 6 taxonomic ranks ]


### Sample names ###


 [1] "HPc1_cut"    "HPc2_cut"    "HPc5_cut"    "HPc6_cut"    "HPps1_cut"
 [6] "HPps2_cut"   "HPps5_cut"   "HPps6_cut"   "KEKc3_cut"   "KEKc4_cut"
[11] "KEKc5_cut"   "KEKc6_cut"   "KEKps3_cut"  "KEKps4_cut"  "KEKps5_cut"
[16] "KEKps6_cut"


### Sample variables ###


[1] "sample"        "original_name" "site"          "individual"
[5] "bagging"       "honeybees"
```

Produces **a phyloseq object (.Rda)** and a text summary

The Rda file is used as **the input for downstream analyses**

OTU table, taxonomy table and sample data can be exported for use outside Chipster (Microbial amplicon data analyses / **Extract information from the Phyloseq object**)

# Part 2: Data tidying and alpha diversity

# Taxon-level clean-up tools

Under Microbial amplicon data analyses:

- Filter by taxonomic group
  - Remove non-specific sequences (keep e.g. Bacteria or Archaea only)

- Remove selected taxa
  - Remove chloroplast and/or mitochondrial sequences
  - (Manually remove specific taxa)

- Overview of taxon composition
  - user-specified level

# Data inspection: Sequence numbers, rarefaction curve and alpha diversity estimates



```
### Per-sample sequence no.s ###


  HPc1_cut   HPc2_cut   HPc5_cut   HPc6_cut  HPps1_cut  HPps2_cut  HPps5_cut
      5084       7467       5859       5218       2664       1585       4495
 HPps6_cut  KEKc3_cut  KEKc4_cut  KEKc5_cut  KEKc6_cut KEKps3_cut KEKps4_cut
      6161      10198       3000       4277       8703       6851       3789
KEKps5_cut KEKps6_cut
      7943      18524




### Alpha diversity estimates (observed OTUs, Chao1, Shannon's index, Pielou's evenness) ###



          Observed   Chao1 se.chao1  Shannon    pielou    sample
HPc1_cut       280 438.8095 40.27640 3.424902 0.6078137  HPc1_cut
HPc2_cut       295 475.0244 44.97236 3.318920 0.5836002  HPc2_cut
HPc5_cut       310 488.5789 45.77474 3.715068 0.6476111  HPc5_cut
HPc6_cut       332 536.1395 49.12739 4.019842 0.6924632  HPc6_cut
HPps1_cut      239 320.2222 23.13010 4.044690 0.7385587 HPps1_cut
HPps2_cut       95 118.0769 12.23199 3.474902 0.7630644 HPps2_cut
HPps5_cut      300 473.3182 42.63040 3.798141 0.6658987 HPps5_cut
HPps6_cut      405 669.0000 62.88645 4.397870 0.7325038 HPps6_cut
KEKc3_cut      267 462.5385 48.92997 2.625653 0.4699367 KEKc3_cut
```
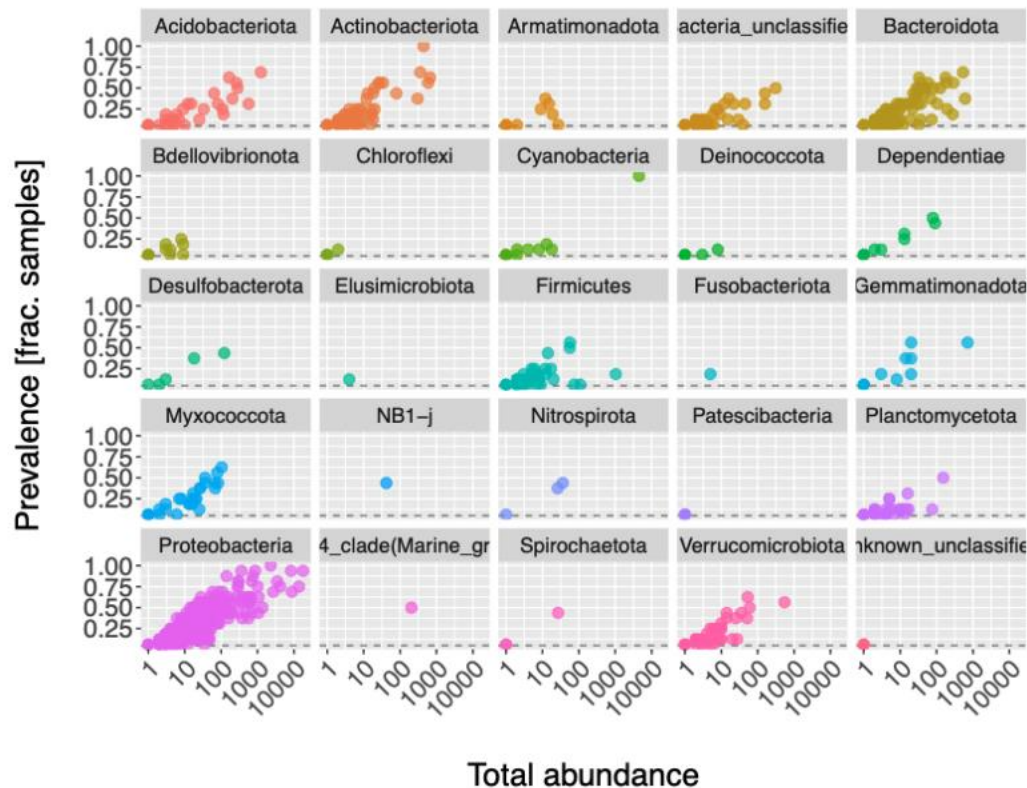
# Filtering low-abundance OTUs

- Filter out OTUs that occur in less than x % of samples
  - Proportional prevalence filtering

- Remove singletons and doubletons
  - Remove OTUs with 0-2 occurrences

# Visualizing low-abundance OTUs

Additional prevalence summaries

- Visualization of OTU prevalence at phylum level
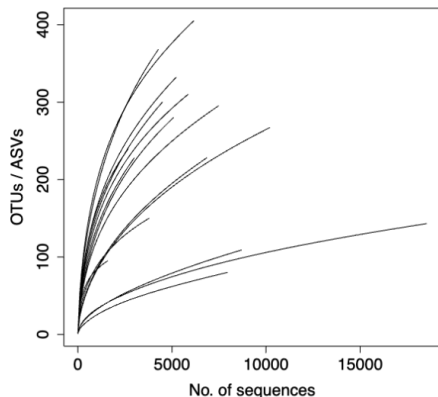
-  Text summary on low prevalence OTUs

# Alpha diversity

**Diversity within a habitat:** how many species are present, at which relative amounts?

- Species richness  - how many species?
  - observed number of OTUs
  - richness estimators (Chao1)

- Evenness – abundance distribution of species?
  - all OTUs equally abundant vs. a few dominant and a lot of rare OTUs
  - Pielou's evenness

- Many diversity indices **combine richness and evenness** (e.g. Shannon index)

# Rarefaction?



**Uneven sequence numbers** among samples can bias comparisons, especially with alpha diversity. Solutions:

- rarefying: equal number of reads picked from all samples (Rarefy OTU data to even depth)

- data transformation  (Transform OTU counts)

# Part 3: Transformations and ordinations

# Transformation of OTU data

Four options (August 2023)

# Relative abundance (%) bar plots

# Beta diversity

**Change in community composition among habitats**

    o Does the microbial community composition in treatment A differ from treatment B?

- Unlike in alpha diversity, OTU identity matters

- Both identity and relative abundance of OTUs usually included

- Quantified with **distance or dissimilarity measures**

# Distance matrices and ordinations

Distance measures available: **Bray-Curtis** or Euclidean

- o centered log ratio (CLR) transformation + Euclidean distance = **Aitchinson distance**

**Ordinations: visualizing beta diversity**

- o nMDS (non-metric multidimensional scaling)
  - o overall variation among samples displayed
- o db-RDA (distance-based redundancy analysis)
  - o focus on the variation explained by phenodata variable(s)



**Recommended for more information:**
Guide to Statistical Analysis in Microbial Ecology:
https://sites.google.com/site/mb3gustame/

# Non-metric multidimensional scaling (nMDS)

# Distance-based redundancy analysis (db-RDA)

Requires specifying one or more phenodata variables

# Distance-based redundancy analysis (db-RDA)



db−RDA

**Constrained ordination** = focus on the community variation explained by the phenodata variable(s)

**Recommended for more information:**
Guide to Statistical Analysis in Microbial Ecology:
https://sites.google.com/site/mb3gustame/

# Part 4: Statistics

# PERMANOVA (permutational multivariate analysis of variance)

Input: distance matrix (ps_dist.Rda)

- Global test: '**Does community structure differ between sample groups**?'
  - o Pairwise test: 'Which groups differ from one another?'

- Currently: several phenodata variables -> added sequentially -> order matters!

- Influenced by both **location** and **dispersion**
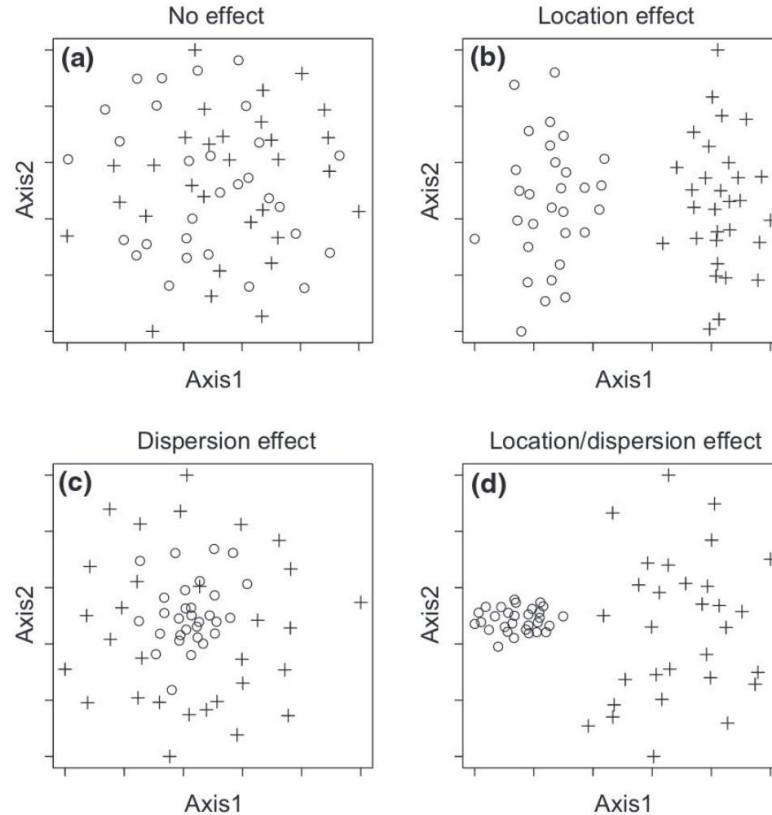
**Location vs. dispersion**

A significant PERMANOVA result can be due to:

- Location effect

- Dispersion effect

- Combination of both



No effect

Location effect

Dispersion effect

Location/dispersion effect

# PERMANOVA output

```
### Global PERMANOVA summary ###




$aov.tab
Permutation: free
Number of permutations: 999

Terms added sequentially (first to last)

             Df SumsOfSqs MeanSqs F.Model      R2 Pr(>F)
get(pheno1)   1    1.5086 1.50863  6.0983 0.30342  0.001 ***
Residuals    14    3.4634 0.24738         0.69658
Total        15    4.9720                 1.00000
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

$call
adonis(formula = ps_dist ~ get(pheno1), data = ps_df)
```

**Df:** Degrees of freedom

**F.Model:** Test statistic (pseudo-$F$)

**R2**: Variation explained by the model

**Pr(>F):** Statistical significance ($p$ value)

# PERMDISP: test for the homogeneity of multivariate dispersion

Input: distance matrix

- Test if a significant PERMANOVA result is due to dispersion, not (only) location

- Significant result -> be careful with PERMANOVA interpretation

# Post-hoc comparisons

Sample group comparisons **following significant global PERMANOVA:**

- Pairwise PERMANOVA (similar as global test but for sample pairs)

Dispersion comparisons **following significant PERMDISP:**

- Tukey's Honestly Significant Difference (HSD) test

- Both methods use a correction for multiple testing (Benjamin-Hochberg correction)

# DESeq2

- Originates from the RNAseq field

- Addresses the question: '**Which taxa are differentially abundant between sample groups?'**

- Enables inferences such as: 'Illness x is associated with a reduction in the abundance of beneficial gut microbes y and z'

- Input: untransformed data -> convert to DESeq2 format with Transform OTU counts (corrects for differences in sequencing depth)

- Results given as **log fold changes**

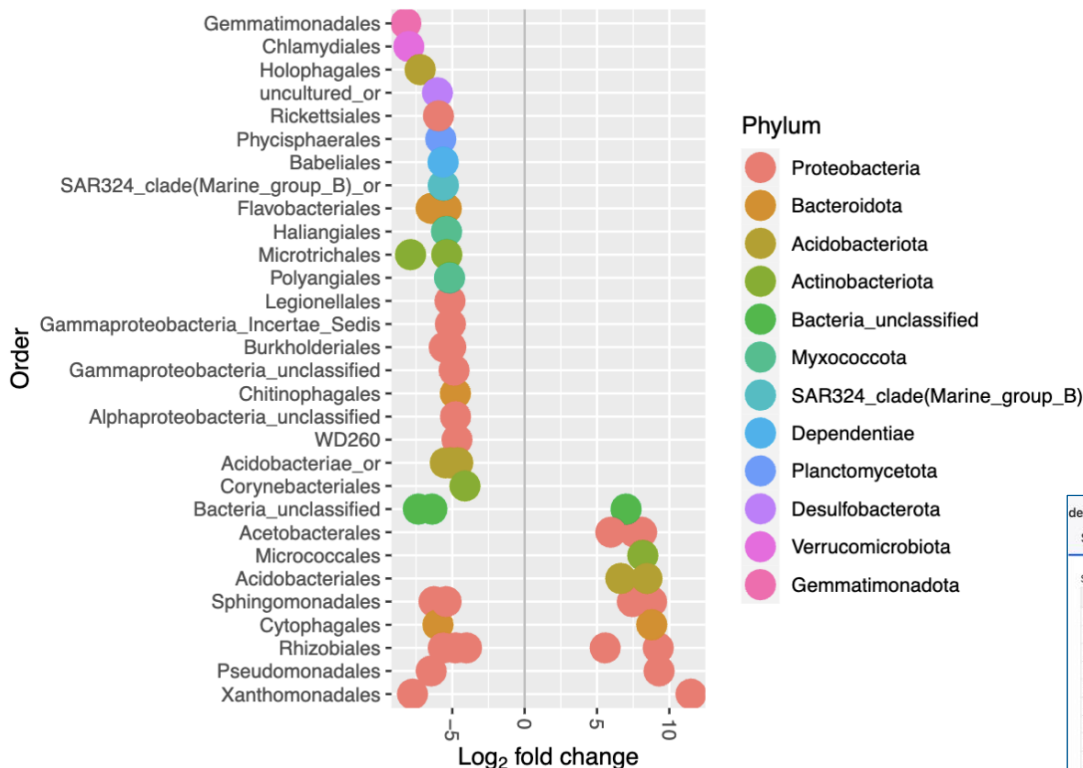More information: https://joey711.github.io/phyloseq-extensions/DESeq2.html

# DESeq2

Current tool configuration (August 2023):

- **Focused on comparison of two groups at a time**
  - o If selected phenodata column has >2 groups, specify a pair (Group 1 and Group 2)

- **Reference level selection**:
  - o Phenodata column with two groups -> first in alphabet is the reference level (e.g. 'b vs. a' or 'sick vs. healthy'=
  - o Phenodata column with >2 groups -> 'Group 2' is the reference level

# DESeq2



8-fold increase compared to reference level = log2 fold change 3 (because $2^3=8$)

each dot = OTU with adjusted p value < 0.01