# Microbiota-Analysis

Jacqueline Wyss, 21.06.2023

# Overview

1. Definitions

2. Sample processing

3. Data analysis

      - Alpha diversity

      - Beta diversity

      - Taxonomy

      - Differential abundance

      - Correlation

      - Models/classification

# Introduction:
# Microbiology vs Microbiota-Studies



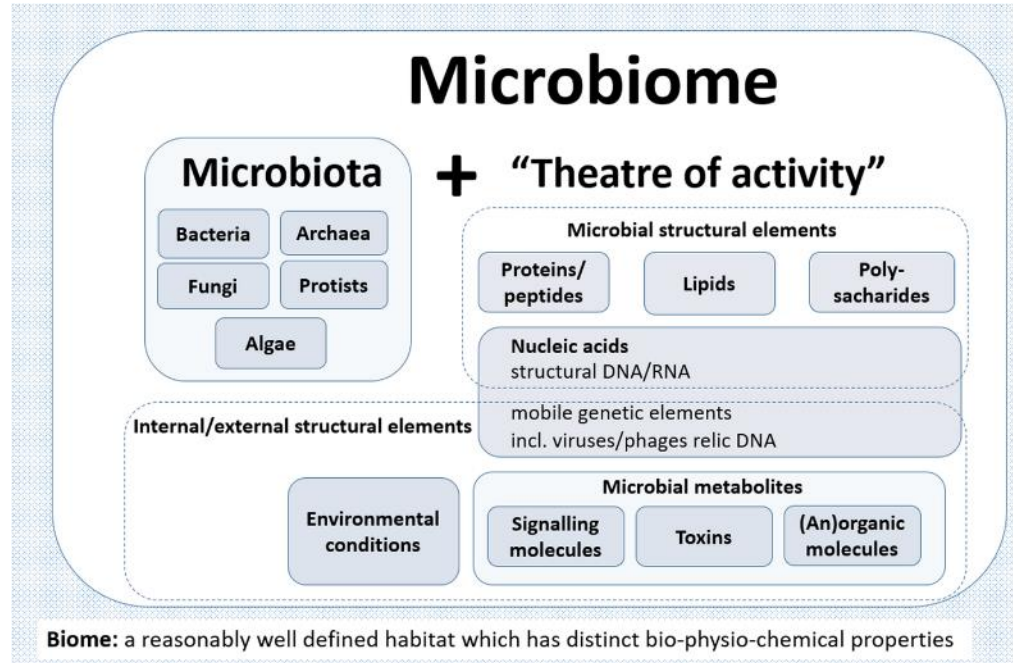Culturing based

Sequencing based

# Introduction: Definitions

Microbiota:

Microbiome:

Meta
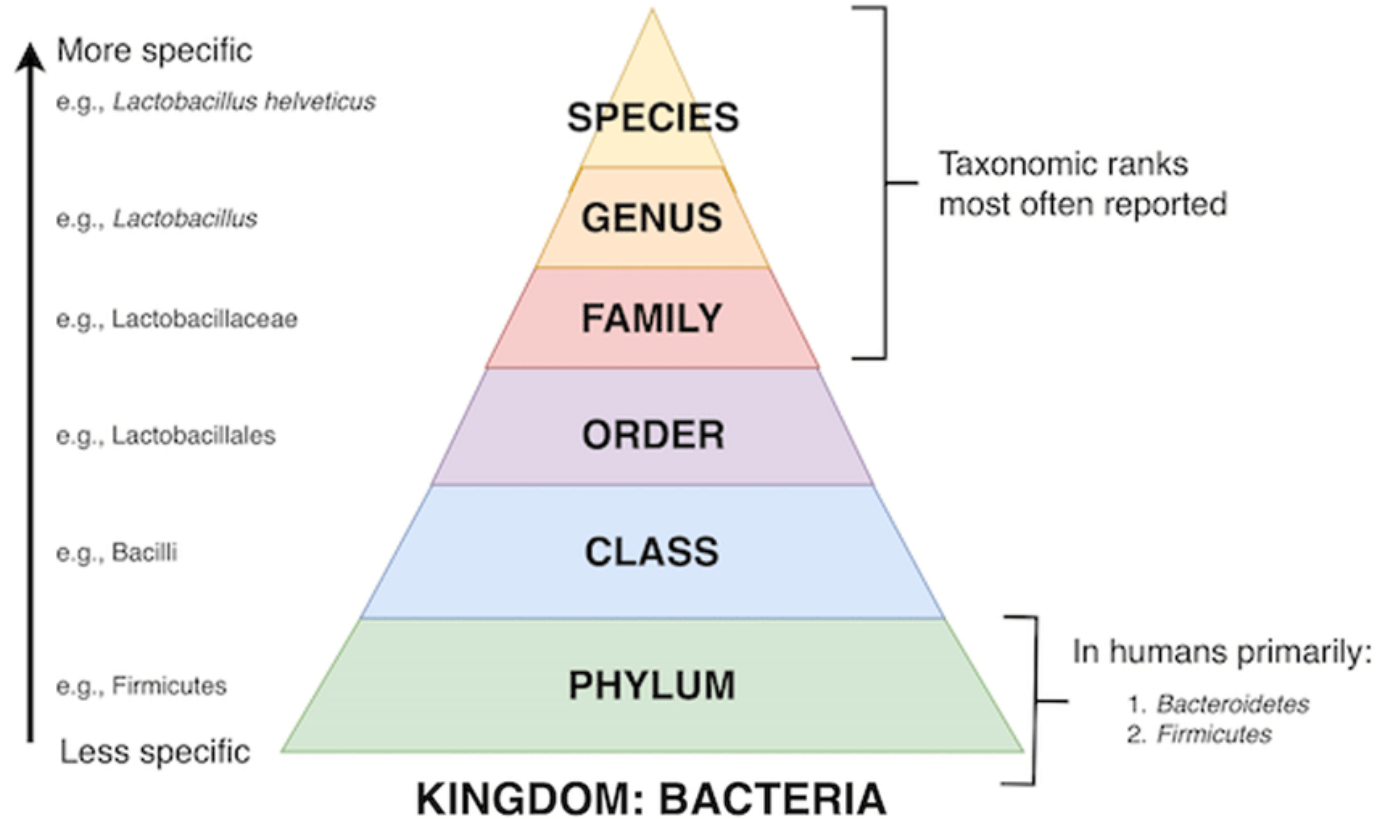
- Genomics

- Transcriptomics

- Metabolomics



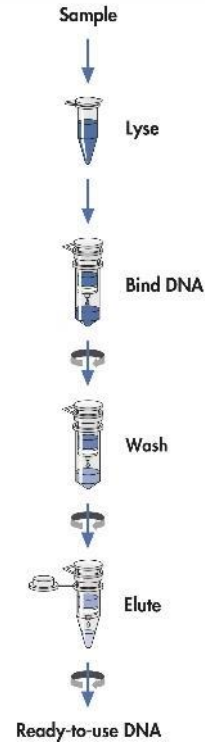Berg et al. Microbiome 2020

# Introduction: Definitions

Taxonomy:

- Species: Strains

- Substrains

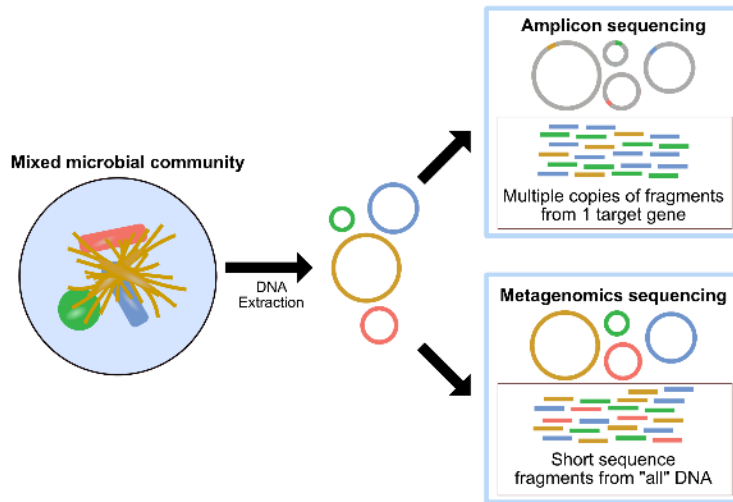Taxa: no definition of the

taxonomic level

More specific

e.g., *Lactobacillus helveticus*

e.g., *Lactobacillus*

e.g., Lactobacillaceae

e.g., Lactobacillales

e.g., Bacilli

e.g., Firmicutes

Less specific

**SPECIES**

**GENUS**

**FAMILY**

**ORDER**

**CLASS**

**PHYLUM**

**KINGDOM: BACTERIA**

Taxonomic ranks most often reported

In humans primarily:
1. *Bacteroidetes*
2. *Firmicutes*

Ameringen et al. Depression and Anxiety 2019

# Workflow:
# Sample preparation



Sample

Lyse

Bind DNA

Wash

Elute

Ready-to-use DNA

DNA is extracted from the samples

DNA can then be used for sequencing

# Workflow:
# Amplicon vs full metagenomic sequencing



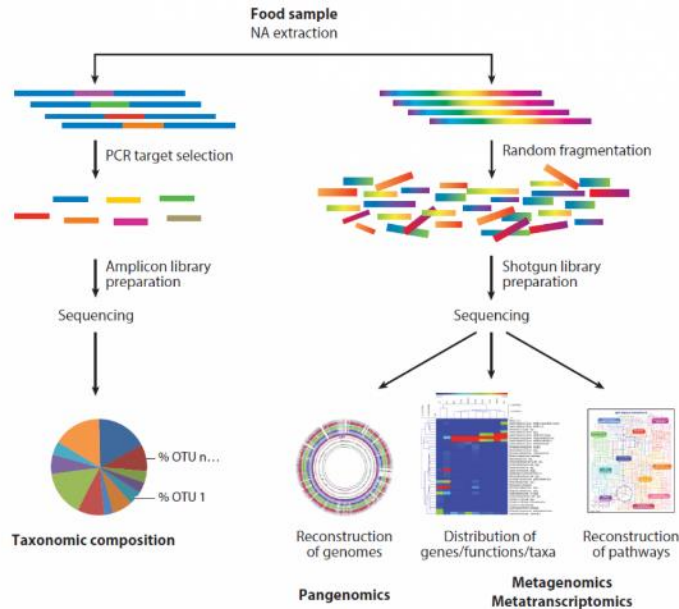https://astrobiomike.github.io

# Workflow:
# Amplicon vs full metagenomic sequencing



In microbiota studies most often a hypervariable region of the 16S subunit is used as a target gene for amplification

https://astrobiomike.github.io

# Workflow:
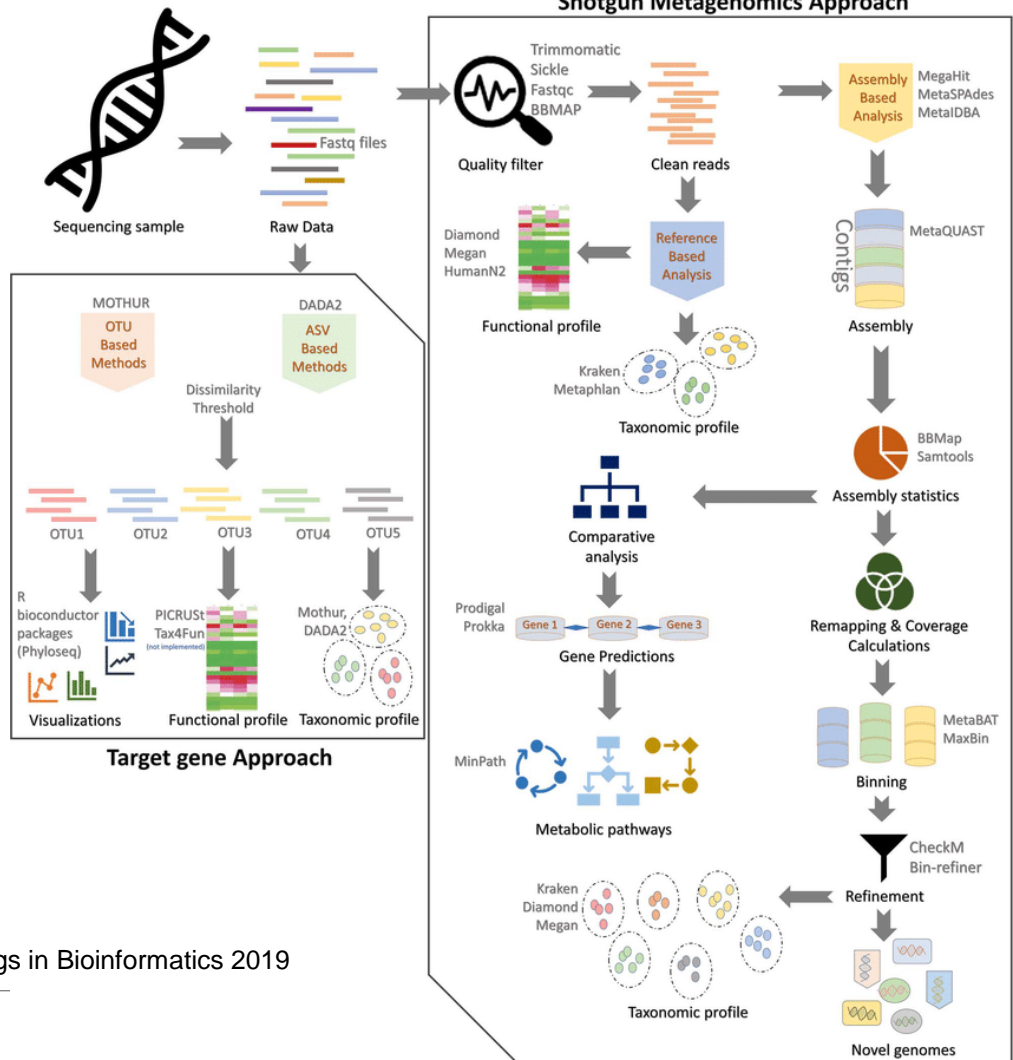# Amplicon vs full metagenomic sequencing



16S sequencing = amplicon sequencing

We mainly get information about the taxa composition

De Filippis F., et al, 2018

# Workflow: Sequencing analysis

Raw sequencing data is analysed by bioinformatic pipelines and packages



Bharti et al. Briefings in Bioinformatics 2019

# Raw sequencing data:
# FASTQ files

@DMJVU:00004:00006

TTCAATTGGCATTAGATACCCTGGTAGTCCACGCCGTAAACGATGTCGAC
TTGGAGGTTGTGCCCTTG

+

;75505057;66ACDCCCCC?DC6<;;;;666;606666,6666666666606066060666--
)-)-

@DMJVU:00004:00009
CTAGGAACCGCATTAGATACCCTGGTAGTCCACGCCGTAAACGATGAGT
GCTAGGTGTTAGACCCTTTCCGGGGTTTAGTGCCGCAGCTAACGCATTA
AGCACTCCGCCTGGGGAGTACGACCGCAAGGTTGAAACTCAAAGGAATT
GACGG

+

;;5D0;7606=66;@CD???B:;;7;B@@B7<<<A>A;;;;;6@@;6>=;4==@0;;;05;49C
*44*4.4==*44*4444EAC?5:8//)--
//:3@39@<;;;7;C6;;@B7@CBACC::4:@@<@=B5:::/::::2*2.239.-----

Sample identifier

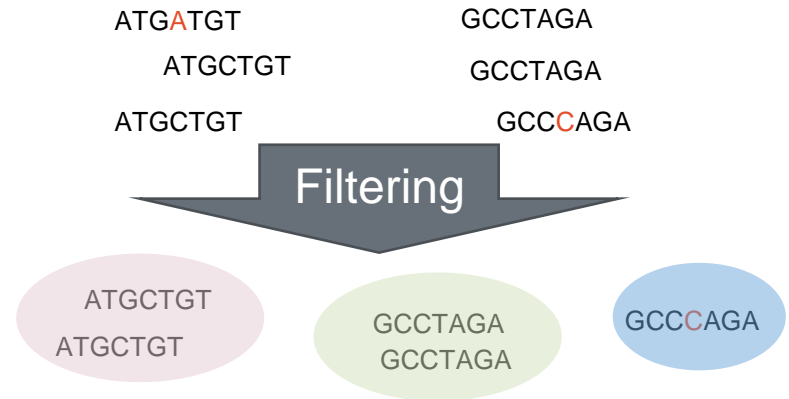Sequencing result

Quality score

# Workflow:
# Grouping reads

## OTU: operational taxonomic units

Grouping by consensus (e.g. 97%)

ATGATGT

ATGCTGT

ATGCTGT

GCCTAGA

GCCTAGA

GCCCAGA

## ASV: amplicon sequencing variants

Denoising and quality filtering before grouping exact matches

ATGATGT        GCCTAGA

ATGCTGT        GCCTAGA

ATGCTGT        GCCCAGA

Filtering

ATGCTGT

ATGCTGT

GCCTAGA

GCCTAGA

GCCCAGA

# Workflow:
# Scripts and Server



Sequences are filtered, demultiplexed grouped and assigned to known sequences in taxonomy databases

# Workflow:
# Feature table

| | Sample 1 | Sample 2 | Sample 3 |
|---|---|---|---|
| Taxa 1 | 0 | 0 | 2 |
| Taxa 2 | 1 | 0 | 0 |
| Taxa 3 | 10 | 2 | 15 |
| Taxa 4 | 0 | 1 | 0 |

Feature tables are sometimes filtered (rare taxa, low abundance, rarefaction)

INSELGRUPPE

# Microbiota data:
# analytic challenges



| Taxonomic table | | | | |
|---|---|---|---|---|
| | \multicolumn{4}{c}{Sample ID} |
| | S1 | S2 | S3 | S4 |
| OTU_1 | | | | |
| OTU_2 | | | | |
| OTU_3 | | | | |
| OTU_4 | | | | |
| OTU_5 | | | | |
| OTU_n | | | | |

| Functional table | | | | |
|---|---|---|---|---|
| | \multicolumn{4}{c}{Sample ID} |
| | S1 | S2 | S3 | S4 |
| KO_01 | | | | |
| KO_02 | | | | |
| KO_03 | | | | |
| KO_04 | | | | |
| KO_05 | | | | |
| KO_n | | | | |

Microbiota data is:

**multivariable** (= because every taxon is a variable)

**sparse** (= have many zeros, because taxa are often only present in few samples)

**compositional** (= are not absolute measurements but are proportions)

Liu et al. Protein & Cell 2021

# Microbiota as an ecosystem



beta-diversity

alpha-diversity

alpha-diversity

# Diversity measures
# alpha diversity

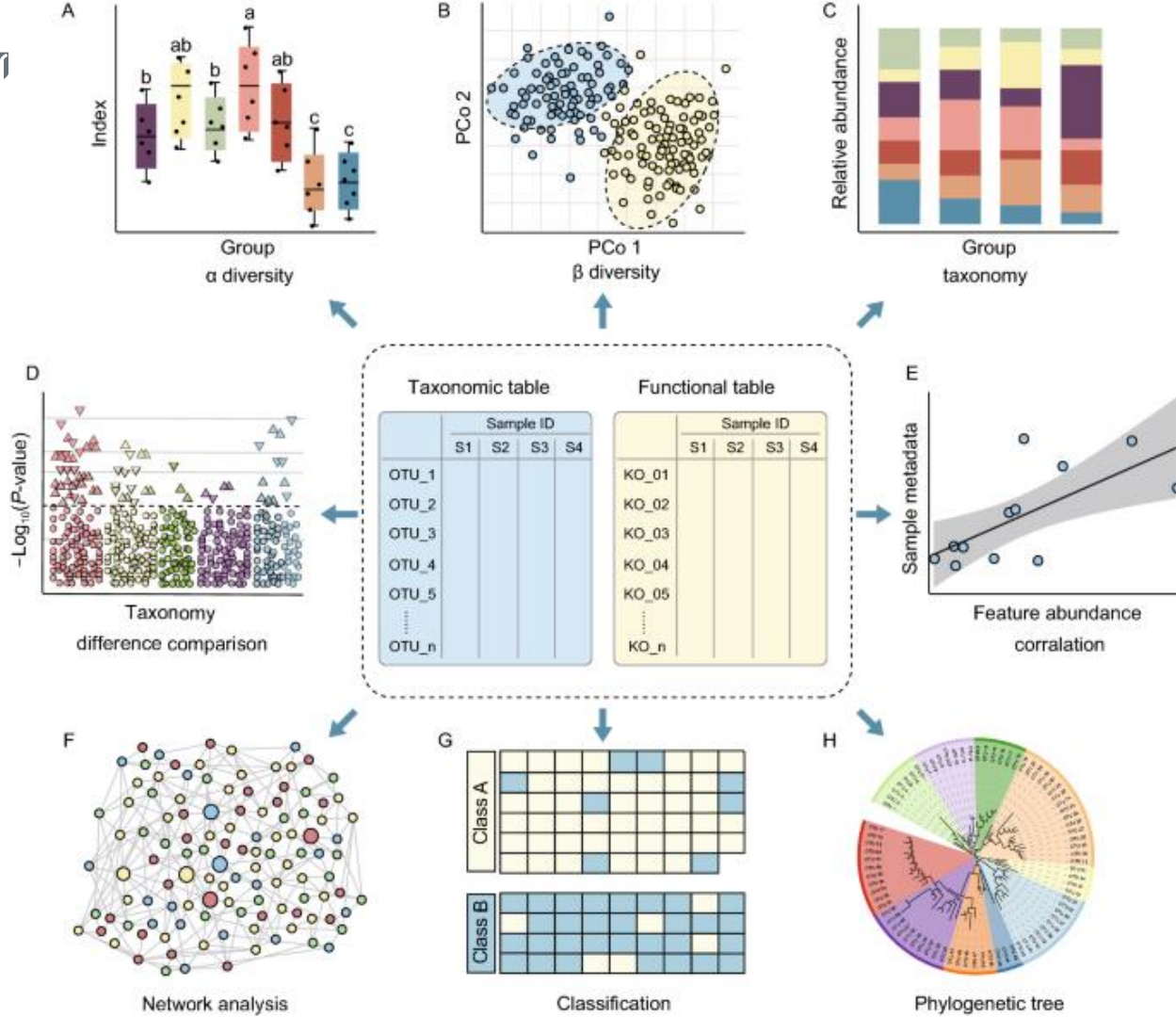Most simple: Species richness= total number of species at a site

Most popular alpha diversity measures:

- Simpson: between 0 and 1

$$D = 1 - \sum_{i=1}^{S} \frac{n_i(n_i - 1)}{n(n - 1)}$$

- Shannon: max= ln S

$$H' = -\sum_i p_i \cdot \ln p_i \quad \text{mit } p_i = \frac{n_i}{N}$$

- Others: Chao1, ACE, ..

Simpson diversity example calculation:
1*0/9*8 = 0
2*1/9*8 = 0.027
3*2/9*8 = 0.083
3*2/9*8 = 0.083
1-(2*0.083+0.027) = 0.8

# Alpha diversity Visualisation

Often visualised as boxplots per groups

Alpha diversity differences statistically evaluated like any other measurement:

- Wilcox-test, ANOVA

# Diversity measures
# beta diversity

Popular beta-diversity measures:

- Euclidean distance

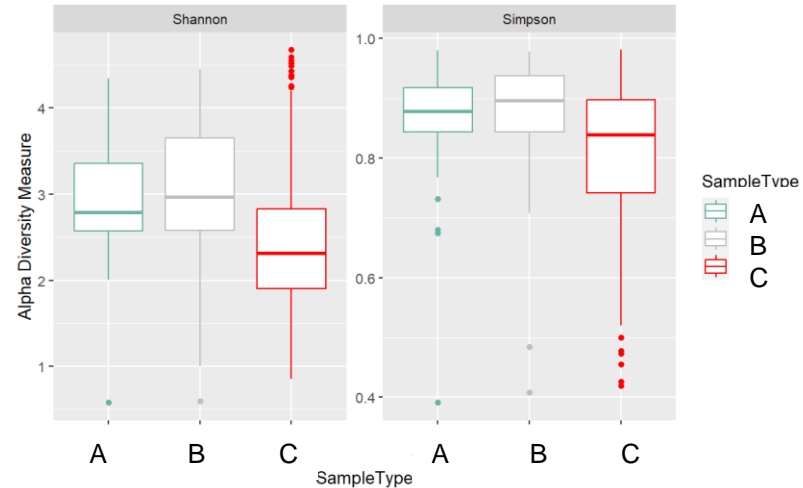- Aitchinson distance: central log transform + euclidean distance

- Bray-curtis dissimilarity

- UniFrac + weighted UniFrac: includes phylogenetic information

- Jaccard, Chao, Mahalobani, etc.

Manhattan distance

Euclidean distance

$$BC_{ij} = 1 - \frac{2C_{ij}}{S_i + S_j}$$

C= common species
S= total species at a site

$$\left( \frac{sum\ of\ unshared\ branch\ lengths}{sum\ of\ all\ tree\ branch\ lengths} \right) = fraction\ of\ total\ unshared\ branch\ lengths$$

**Bet**
**Vis**

Choic



Distance metrics for dataset

# Beta diversity
# Pitfall: Batch effects



Bray-curtis-dissimilarity between samples

# Beta diversity
# Pitfall: Batch effects

# Beta diversity
# Ordination

Multivariate data can be plotted in a 2D way using dimensionality reduction techniques

- PCA:  for any multivariate data set, characterizing the main axes of variance in the data

- PCoA: main axes of variance for distance matrices

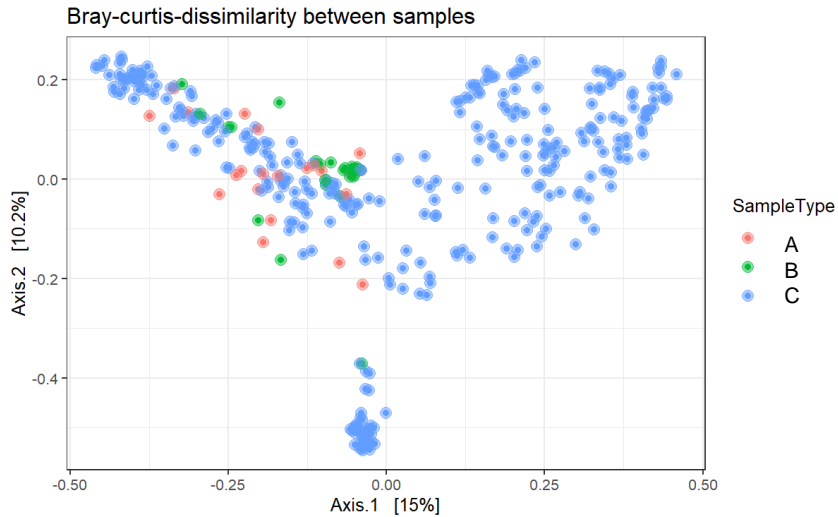# Principal components
# Idea

# Beta diversity
# Ordination

Multivariate data can be plotted in a 2D way using dimensionality reduction techniques

- PCA: for any multivariate data set, characterizing the main axes of variance in the data
- PCoA: main axes of variance for distance matrices


- NMDS (Non-Metric Multidimensional Scaling): iterative process, axes can not be interpreted


- Other: correspondance analysis (CA), redundancy analysis (RDA): uses metadata variables for scaling

# Beta diversity
# Visualisation

Choice of distance method is variable, ultimate goal is to see separation of groups

# Beta diversity
# Visualisation

Significant differences between groups need to tested by methods assessing multivariable data

- ANOSIM, ADONIS, PERMANOVA, MRPP

# Taxonomy

**Taxa composition**



Taxonomy results are not absolute, taxa composition are proportions

# Taxonomy

**Taxa absolute numbers**



This can lead to distortions if the overall bacterial count fluctuates between samples

# Taxonomy
# Phylogenetic trees



Khan academy

Phylogenetic trees can be built with reference databases or with information from the sequencing reads

# Taxonomy
# Phylogenetic trees



Phylogenetic trees can be built with reference databases or with information from the sequencing reads

Almeida et al. Nature, 2019

# Analysis:
# Correlation analysis



Wikipedia: Correlation

Data can be assessed for correlation structure

- Between microbiota

- Between microbiota and metadata

Because every taxa is tested seperately results have to be adjusted by multiple testing correction

# Correlation analysis: Visualisation



Koliada et al. BMC Microbiology 2017

Correlations with one predictor can be depticted as a regression

Multiple correlations can be depicted as a network analysis or in a correlation plot

# Correlation analysis: Visualisation



Correlations with one predictor can be depticted as a regression

Multiple correlations can be depicted as a network analysis or in a heatmap

Zhu et al. Frontiers in Cellular and Infection Microbiology, 2019

# Analysis
# Differential abundance

Differental abundance tells us
which taxa differ significantly
between groups

# Differential abundance Visualisation

Volcano plot



Arnoriaga-Rodríguez et al., Gut, 2021

# Analysis
# Differential abundance

Methods are not
standardised
and can output different
results

Nearing et al. Nature communications 2022

# Classification:
# Model fitting

If we have some hypothesis about the underlying structure of the data we can try to fit a model

- e.g. classification into two disease groups

- Model parameters can give information about which features (=e.g. taxa) are important

# Classification:
# Model fitting

If we have some hypothesis about the underlying structure of the data we can try to fit a model

- e.g. classification into two disease groups

- Model parameters can give information about which features (=e.g. taxa) are important

Examples:

- Multiple linear regression

- Random forest

- Partial Least-Squares Discriminant Analysis (PLS-DA)

- Neural networks

# Model fitting: pitfalls



Explain the Prediction

# Overfitting

Overfitting: If a model uses a lot of features it can be fitted very well to a training data set

- Extreme: one variable classifies one samples -> perfect fit

- Problem: **This model will not be informative for a new unrelated data**

Solution: Validation with new data

- Cross-validation

- Splitting of data before analysis

- Second evaluation cohort

Better

# Model evaluation: AUC-ROC-curves

Good AUC results

- Over 0.7: acceptable

- Over 0.8: excellent

- Over 0.9: outstanding

True positive rate (Precision) = TP/(TP+FN)

False positive rate= FP/(FP+TN)

Sensitivity (Recall): 1- TPR

Specificity: 1- FPR
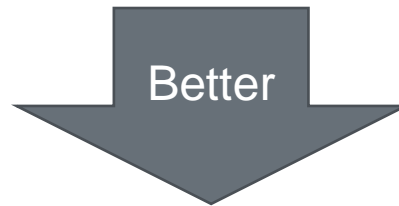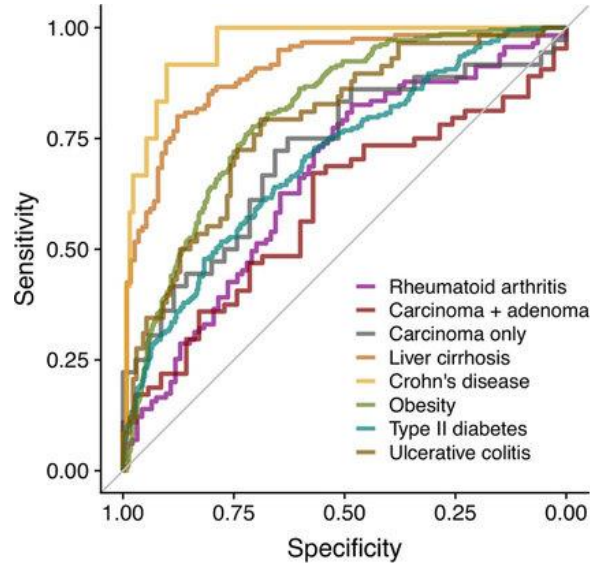


| Random Forest by Disease | | | | |
|---|---|---|---|---|
| Level | Disease | Color | OOB error | AUC |
| Module | Crohn's Disease | Yellow | 5.56% | 0.954 |
| Module | Liver cirrhosis | Orange | 17.09% | 0.902 |
| Module | Obesity | Green | 22.57% | 0.803 |
| Module | Ulcerative colitis | Brown | 25.26% | 0.783 |
| Module | Type II diabetes | Blue | 31.58% | 0.708 |
| Module | Rheumatoid arthritis | Purple | 35.58% | 0.664 |
| Module | Colorectal carcinoma | Red | 36.36% | 0.596 |
| Module | Carcinoma (without adenoma) | Grey | 35.21% | 0.715 |

Legend: Rheumatoid arthritis, Carcinoma + adenoma, Carcinoma only, Liver cirrhosis, Crohn's disease, Obesity, Type II diabetes, Ulcerative colitis

Armour et. al, mSystems, 2019

# Metagenomic analysis: Functional potential

Metagenomic sequencing analyses the whole genome of  bacteria
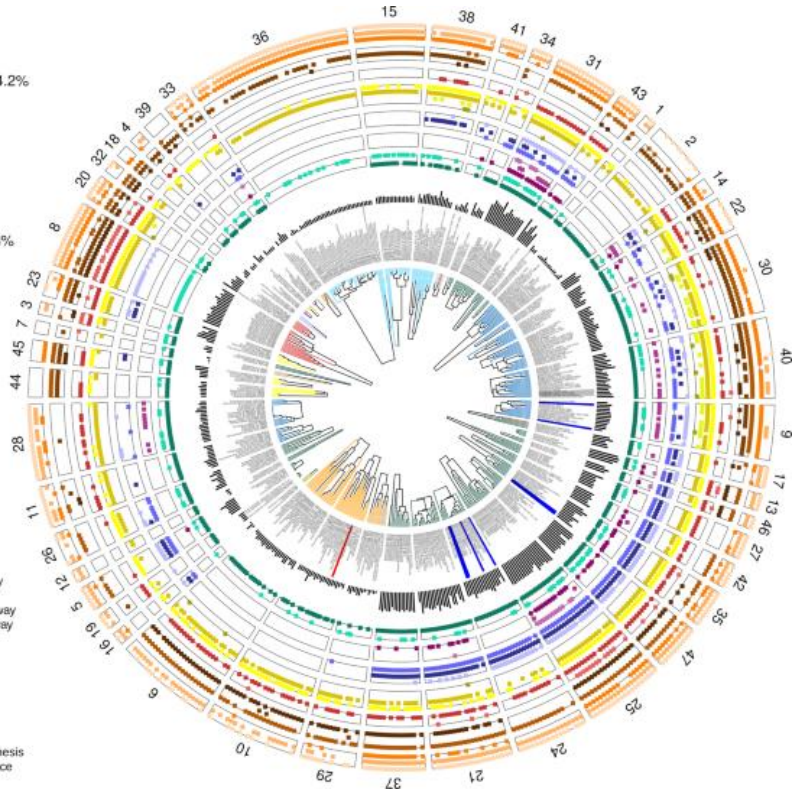
This includes bacterial enzymes:

- Analysis of functional potential



**Taxonomy:**
- Firmicutes – 4%
- Bdellovibrionota – 0.6%
- Gammaproteobacteria – 34.2%
- Actinobacteriota – 4.7%
- Deinococcota – 0.6%
- Thermotogota – 1.7%
- Spirochaetota – 0.6%
- Others – 3%
- Desulfobacterota – 1.1%
- Cyanobacteria – 15%
- Aquificota – 0.6%
- Alphaproteobacteria – 21.8%
- Bacteroidota – 11.4%
- Campylobacterota – 0.6%

**Interaction-traits:**
- Vitamin B1 biosynthesis
- Vitamin B7 biosynthesis
- Vitamin B12 biosynthesis
- Vitamin B1 transport
- Vitamin B7 transport
- Vitamin B12 transport
- Fe–Siderophore
- Fe–Siderophore transporter
- Auxin: Indole–3–pyruvate pathway
- Auxin: Tryptamine pathway
- Auxin: Indole–3–acetonitrile pathway
- Auxin: Indole–3–acetamide pathway
- Quorum sensing
- Chemotactic behavior
- Motility and adhesion apparatus
- Type III secretion system
- Type IV secretion system
- Type VI secretion system
- Antimicrobial compounds biosynthesis
- Antimicrobial compounds resistance

Zoccarato et. al, communications biology, 2022

# Metagenomic analysis: Functional potential

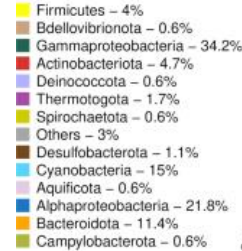Metagenomic sequencing analyses the whole genome of bacteria

This includes bacterial enzymes:
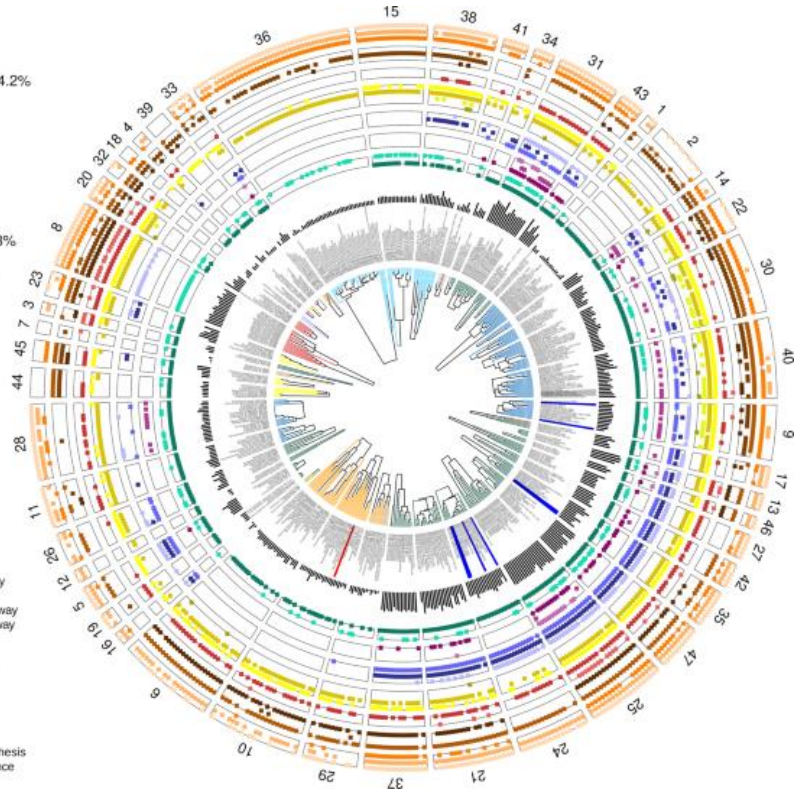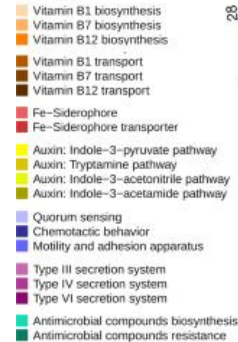
- Analysis of functional potential

Because of higher resolution also:

- Species level taxonomic profiling



Zoccarato et. al, communications biology, 2022

# Type of analysis depends on the research question

# What questions can we answer?

- We can identify if specific taxa are present

- We can compare if a sample has reduced alpha diversity (e.g. after antibiotic therapy)

- We can compare samples with others (e.g. is a sample more similar to other from heathly people or from a disease group)

- We can establish taxa signatures that are diagnostic or predicitve for clinical variables

- Metagenomic sequencing can give answers about metabolic functions of taxa: pathways, metabolites, bacterial signalling molecules, ..

# INSELGRUPPE

Questions?