# Metagenomics processing

E5444 Analysis of Sequencing Data

Vojtěch Bartoň, 2023

# Metagenomics

- **Microbial Community Genetics:** Metagenomics studies the genetics of entire microbial communities in environmental samples.

- **Genomic Snapshot:** It doesn't require isolating or culturing individual microorganisms, providing a holistic view of diverse microbes.

- **Health Insights:** It helps us study the human microbiome and its impact on health without traditional culturing methods

- **Applications in Ecology and Biotech:** Metagenomics informs ecosystem understanding and aids in biotech discoveries like enzymes and antibiotics

# Metagenomics: WMGS x 16S

**WMGS**

Captures the entire genetic content of all microorganisms.

High-resolution data for functional analysis and taxonomic identification.

Suitable for complex ecosystems and novel gene discovery.

Used in environmental and clinical metagenomics.

**16S**

Targets a specific 16S gene marker in bacteria and archaea.

Lower resolution, primarily for taxonomic identification.

Commonly used for microbial community profiling and diversity studies.
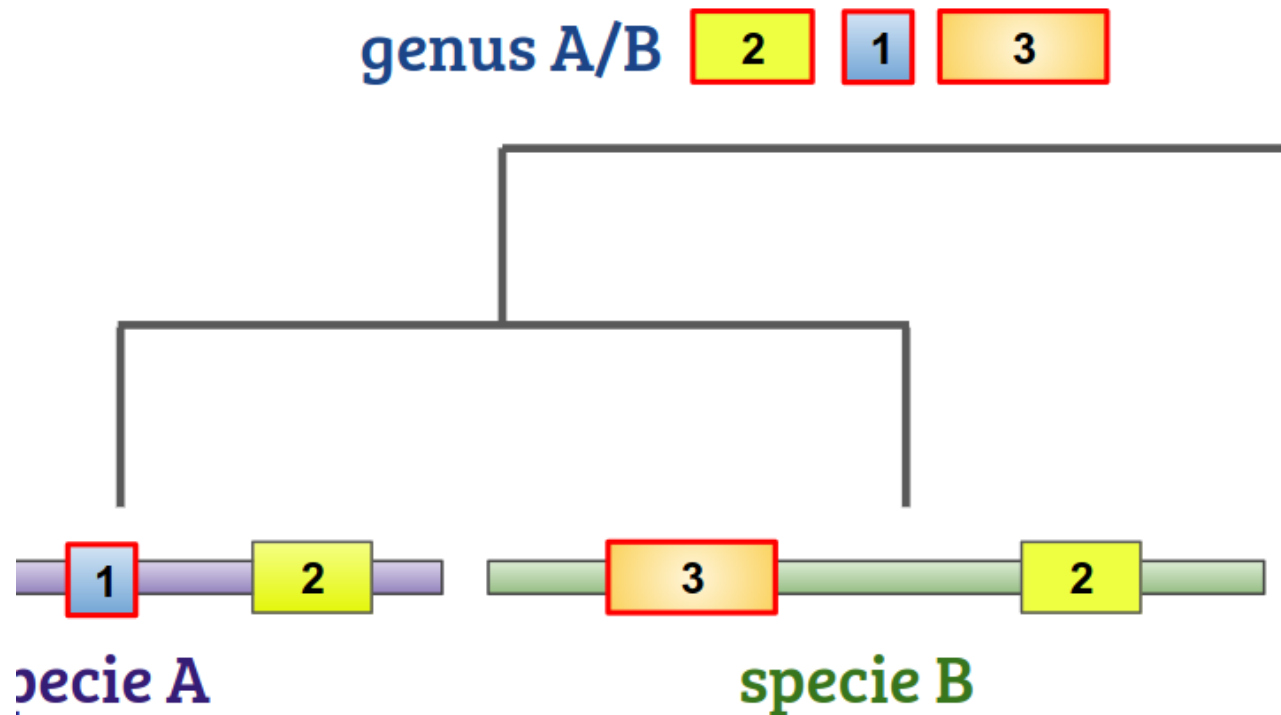
Cost-effective method for taxonomic analysis.

# MetaPhlAn

- **Meta**genomic **Ph**ylogenetic **An**alysis
- **Taxonomic Profiling:** MetaPhlAn identifies and quantifies microbes in samples.
- **Marker Gene Approach:** It uses unique genetic markers for speedy and accurate identification.
- **Efficiency:** Known for fast analysis of large datasets.
- **Applications:** Widely used in microbiome research and clinical metagenomics.

- https://github.com/biobakery/MetaPhlAn

# Metaphlan: input data

- **Shotgun** Whole Metagenome Sequencing

- Sequences (fasta, fastq)
- Database of taxonomically known sequences
  - Unique regions

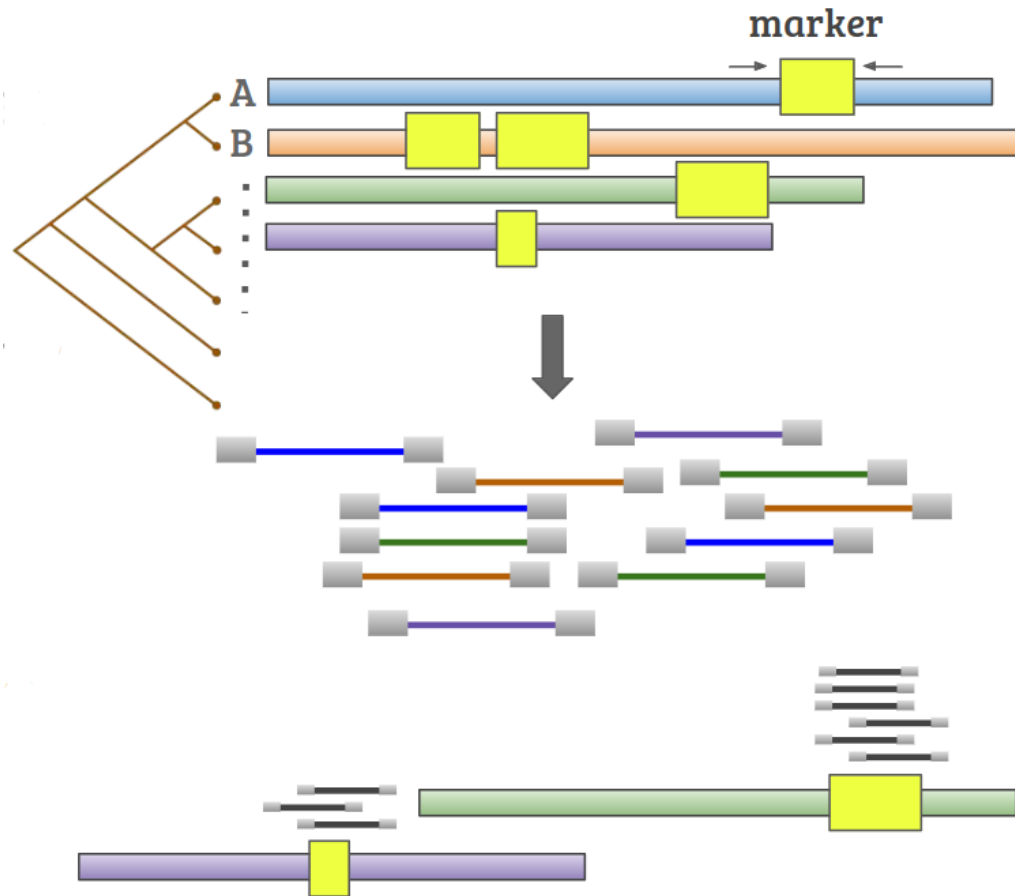# Metaphlan: marker genes



- **Clades**
  - Groups of genomes (organisms) believed to have evolved from a common ancestor
- **Clade-specific marker-genes**
  - Strongly conserved within the clade's genomes
  - Not similar to any sequence in other clades (of the same level)
  - Unique markers change as the clade level grows
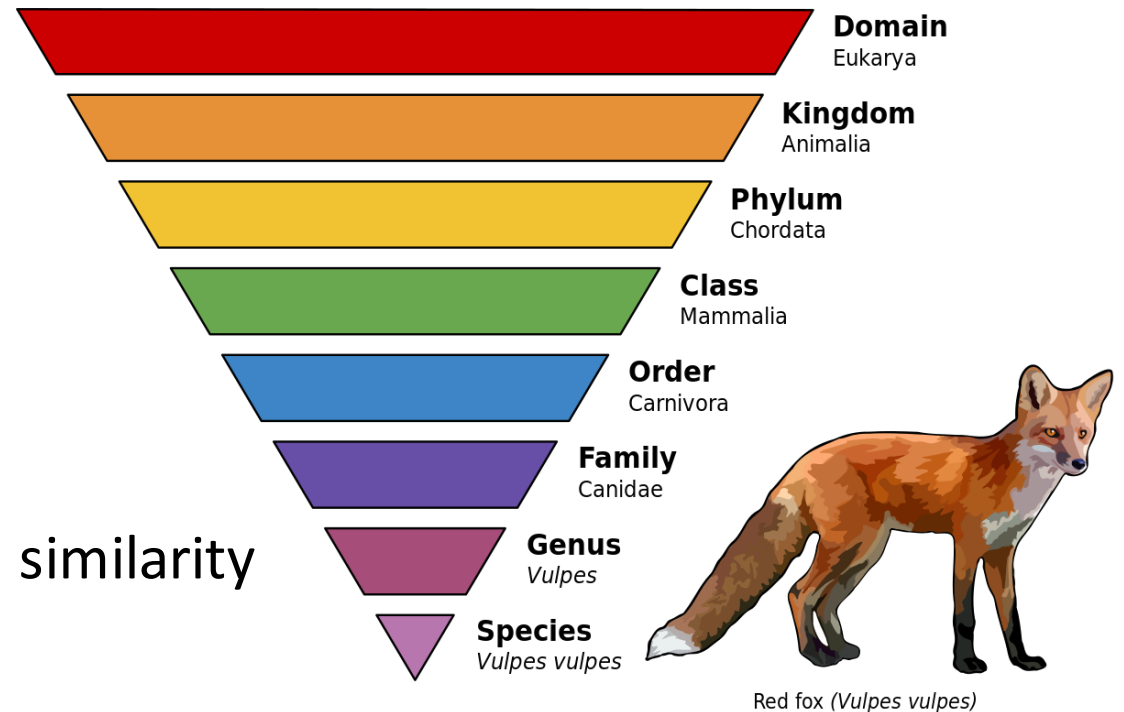  - They also accumulate in a way (direct vs indirect)...

# Metaphlan: Overview



- Reference genomes and their taxonomy
- Find clade-specific marker genes

- Sequence your sample

- Map to marker genes
- Count taxonomic units

# Metaphlan: Reference database

- **ChocoPhlAn**
- Acquire reference genomes
  - De novo assembly
  - Cultured species
  - **Uniprot core data**
- Acquire taxonomy
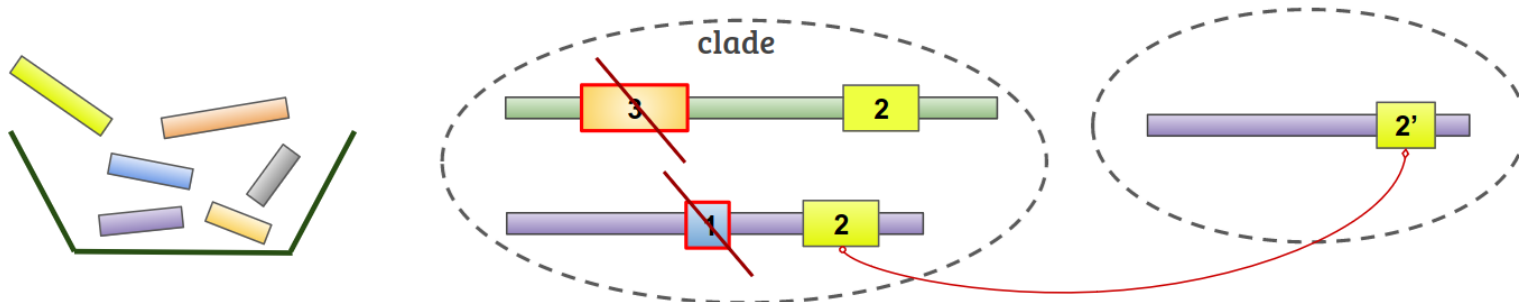  - Hierarchical clustering tree based on similarity
  - **Ncbi taxonomy**

**Domain**
Eukarya

**Kingdom**
Animalia

**Phylum**
Chordata

**Class**
Mammalia

**Order**
Carnivora

**Family**
Canidae

**Genus**
*Vulpes*

**Species**
*Vulpes vulpes*

Red fox *(Vulpes vulpes)*

# Metaphlan: Reference database

- The general process:
  - Each genome → bag-of-genes representation
  - Only conserved genes in the clade are saved
  - Inter-clade uniqueness index elimination
  - Single-copy genes were preferred of multi-copy genes

- Properties of the markers
  - Gene level
  - 5.1M filtered genes
  - 27K species-level genome bins
  - Not necessarily continuous (bag-of-genes)
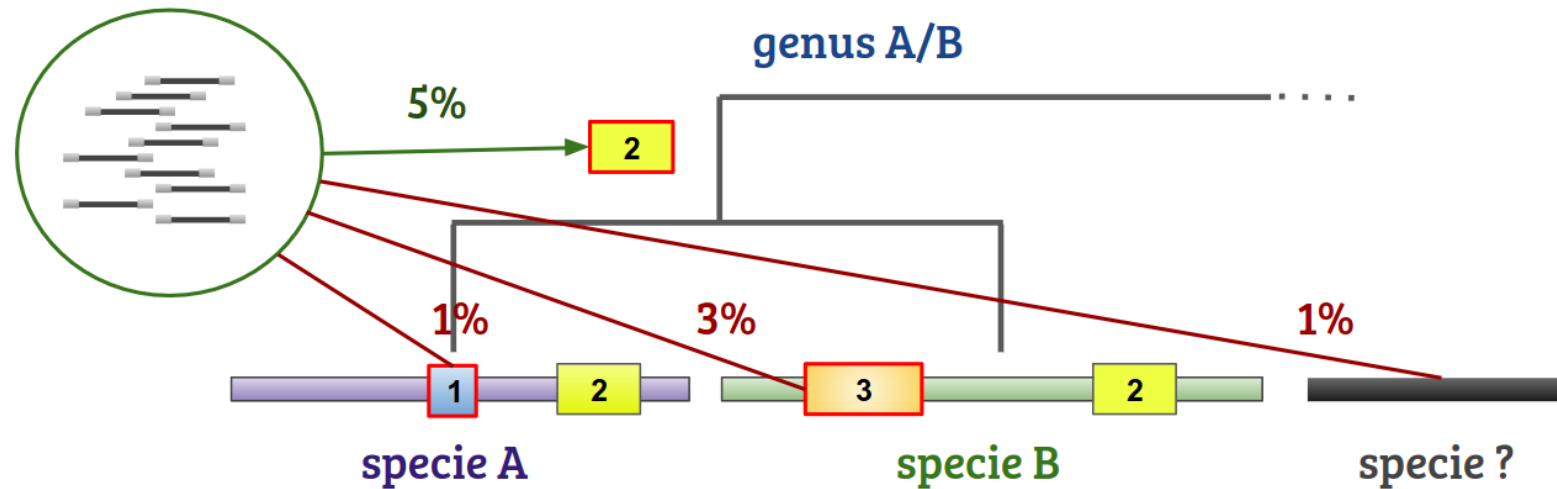  - ~4% of the total genome length
  - ~260 markers per specie

# Metaphlan: Taxonomic profiling

- Map reads against reference database of marker genes
- Calculate relative abundance
  - Sum the total reads mapped to clade markers
  - Divide by marker's total length
  - Abundances in every clade-level sum up to 100%
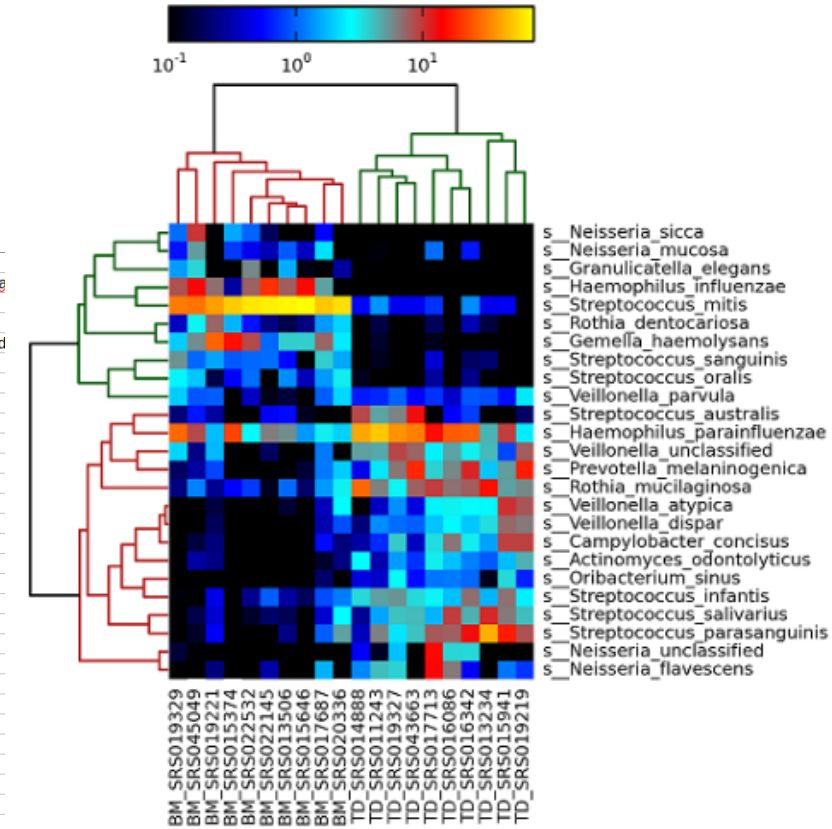
# Metaphlan: Taxonomic profiling

- Unclassified case
  - Move up in the taxonomy tree

# Metaphlan: outputs and visualisations

- bugs_list.tsv

- Utility scripts



| #mpa_vJan21_CHOCOPhlAnSGB_202103 | | | |
|---|---|---|---|
| #/mnt/storage-brno3-cerit/share/310000-SCI/999990-rcxag/Projects/2_Infrastructure/1_Bioinformatics/Tools/conda_env/biobakery/bin/metaphlan /mnt/storage-brno3-cerit/share/310000-SCI/999990-rcxag/Projects/2_Infrastructure/1_Bioinforma | | | |
| #2616413 reads processed | | | |
| #SampleID | Metaphlan_Analysis | | |
| #clade_name | NCBI_tax_id | relative_abundance | add |
| k__Bacteria | 2 | 99.00119 | |
| k__Eukaryota | 2759 | 0.99881 | |
| k__Bacteria\|p__Actinobacteria | 2\|201174 | 90.92536 | |
| k__Bacteria\|p__Proteobacteria | 2\|1224 | 8.07582 | |
| k__Eukaryota\|p__Basidiomycota | 2759\|5204 | 0.99881 | |
| k__Bacteria\|p__Actinobacteria\|c__Actinobacteria | 2\|201174\|1760 | 90.92536 | |
| k__Bacteria\|p__Proteobacteria\|c__Gammaproteobacteria | 2\|1224\|1236 | 7.27117 | |
| k__Eukaryota\|p__Basidiomycota\|c__Malasseziomycetes | 2759\|5204\|1538075 | 0.99881 | |
| k__Bacteria\|p__Proteobacteria\|c__Betaproteobacteria | 2\|1224\|28216 | 0.80465 | |
| k__Bacteria\|p__Actinobacteria\|c__Actinobacteria\|o__Propionibacteriales | 2\|201174\|1760\|85009 | 90.42538 | |
| k__Bacteria\|p__Proteobacteria\|c__Gammaproteobacteria\|o__Alteromonadales | 2\|1224\|1236\|135622 | 7.27117 | |
| k__Eukaryota\|p__Basidiomycota\|c__Malasseziomycetes\|o__Malasseziales | 2759\|5204\|1538075\|162474 | 0.99881 | |
| k__Bacteria\|p__Proteobacteria\|c__Betaproteobacteria\|o__Burkholderiales | 2\|1224\|28216\|80840 | 0.80465 | |
| k__Bacteria\|p__Actinobacteria\|c__Actinobacteria\|o__Corynebacteriales | 2\|201174\|1760\|85007 | 0.49999 | |
| k__Bacteria\|p__Actinobacteria\|c__Actinobacteria\|o__Propionibacteriales\|f__Propionibacteriaceae | 2\|201174\|1760\|85009\|31957 | 90.42538 | |
| k__Bacteria\|p__Proteobacteria\|c__Gammaproteobacteria\|o__Alteromonadales\|f__Alteromonadaceae | 2\|1224\|1236\|135622\|72275 | 7.27117 | |
| k__Eukaryota\|p__Basidiomycota\|c__Malasseziomycetes\|o__Malasseziales\|f__Malasseziaceae | 2759\|5204\|1538075\|162474\|742845 | 0.99881 | |
| k__Bacteria\|p__Proteobacteria\|c__Betaproteobacteria\|o__Burkholderiales\|f__Comamonadaceae | 2\|1224\|28216\|80840\|80864 | 0.80465 | |
| k__Bacteria\|p__Actinobacteria\|c__Actinobacteria\|o__Corynebacteriales\|f__Lawsonellaceae | 2\|201174\|1760\|85007\|2805586 | 0.49999 | |
| k__Bacteria\|p__Actinobacteria\|c__Actinobacteria\|o__Propionibacteriales\|f__Propionibacteriaceae\|g__Cutibacterium | 2\|201174\|1760\|85009\|31957\|1912216 | 90.42538 | |
| k__Bacteria\|p__Proteobacteria\|c__Gammaproteobacteria\|o__Alteromonadales\|f__Alteromonadaceae\|g__Alishewanella | 2\|1224\|1236\|135622\|72275\|111142 | 7.27117 | |
| k__Eukaryota\|p__Basidiomycota\|c__Malasseziomycetes\|o__Malasseziales\|f__Malasseziaceae\|g__Malassezia | 2759\|5204\|1538075\|162474\|742845\|55193 | 0.99881 | |
| k__Bacteria\|p__Proteobacteria\|c__Betaproteobacteria\|o__Burkholderiales\|f__Comamonadaceae\|g__Delftia | 2\|1224\|28216\|80840\|80864\|80865 | 0.80465 | |
| k__Bacteria\|p__Actinobacteria\|c__Actinobacteria\|o__Corynebacteriales\|f__Lawsonellaceae\|g__GGB2722 | 2\|201174\|1760\|85007\|2805586\| | 0.49999 | |
| k__Bacteria\|p__Actinobacteria\|c__Actinobacteria\|o__Propionibacteriales\|f__Propionibacteriaceae\|g__Cutibacterium\|s__Cutibacterium_acnes | 2\|201174\|1760\|85009\|31957\|1912216\|1747 | 89.48138 | |
| k__Bacteria\|p__Proteobacteria\|c__Gammaproteobacteria\|o__Alteromonadales\|f__Alteromonadaceae\|g__Alishewanella\|s__Alishewanella_agri | 2\|1224\|1236\|135622\|72275\|111142\|553384 | 7.27117 | |
| k__Eukaryota\|p__Basidiomycota\|c__Malasseziomycetes\|o__Malasseziales\|f__Malasseziaceae\|g__Malassezia\|s__Malassezia_restricta | 2759\|5204\|1538075\|162474\|742845\|55193\|76775 | 0.99881 | |
| k__Bacteria\|p__Actinobacteria\|c__Actinobacteria\|o__Propionibacteriales\|f__Propionibacteriaceae\|g__Cutibacterium\|s__Cutibacterium_granulosum | 2\|201174\|1760\|85009\|31957\|1912216\|33011 | 0.944 | |
| k__Bacteria\|p__Proteobacteria\|c__Betaproteobacteria\|o__Burkholderiales\|f__Comamonadaceae\|g__Delftia\|s__Delftia_acidovorans | 2\|1224\|28216\|80840\|80864\|80865\|80866 | 0.80465 | |
| k__Bacteria\|p__Actinobacteria\|c__Actinobacteria\|o__Corynebacteriales\|f__Lawsonellaceae\|g__GGB2722\|s__GGB2722_SGB3663 | 2\|201174\|1760\|85007\|2805586\|\| | 0.49999 | |
| k__Bacteria\|p__Actinobacteria\|c__Actinobacteria\|o__Propionibacteriales\|f__Propionibacteriaceae\|g__Cutibacterium\|s__Cutibacterium_acnes\|t__SGB16955 | 2\|201174\|1760\|85009\|31957\|1912216\|1747\| | 89.48138 | |
| k__Bacteria\|p__Proteobacteria\|c__Gammaproteobacteria\|o__Alteromonadales\|f__Alteromonadaceae\|g__Alishewanella\|s__Alishewanella_agri\|t__SGB9784 | 2\|1224\|1236\|135622\|72275\|111142\|553384\| | 7.27117 k__Bacteria\|p__Proteobacteria\|c__Gammaproteobacteria\|o__Alteromonadales | |
| k__Eukaryota\|p__Basidiomycota\|c__Malasseziomycetes\|o__Malasseziales\|f__Malasseziaceae\|g__Malassezia\|s__Malassezia_restricta\|t__EUK76775 | 2759\|5204\|1538075\|162474\|742845\|55193\|76775\| | 0.99881 | |
| k__Bacteria\|p__Actinobacteria\|c__Actinobacteria\|o__Propionibacteriales\|f__Propionibacteriaceae\|g__Cutibacterium\|s__Cutibacterium_granulosum\|t__SGB16958 | 2\|201174\|1760\|85009\|31957\|1912216\|33011\| | 0.944 | |
| k__Bacteria\|p__Proteobacteria\|c__Betaproteobacteria\|o__Burkholderiales\|f__Comamonadaceae\|g__Delftia\|s__Delftia_acidovorans\|t__SGB12680 | 2\|1224\|28216\|80840\|80864\|80865\|80866\| | 0.80465 k__Bacteria\|p__Proteobacteria\|c__Betaproteobacteria\|o__Burkholderiales\|f__ | |
| k__Bacteria\|p__Actinobacteria\|c__Actinobacteria\|o__Corynebacteriales\|f__Lawsonellaceae\|g__GGB2722\|s__GGB2722_SGB3663\|t__SGB3663 | 2\|201174\|1760\|85007\|2805586\|\|\| | 0.49999 | |

# Metaphlan: application

- Shotgun sequencing

- Microbiome Profiling

- Metagenomics

- Metatranscriptomics


- As input for HUMAnN
  - profiling the abundance of microbial metabolic pathways and other molecular functions

# Metaphlan: pros & cons

**Pros**

- Rapid profiling

- Accuracy

- Versatile

- Quantitative output

**Cons**

- Limited functional information

- Reference-dependent

- Computational resources

- Interpreting unknowns

# QIIME2

| | A | B | C |
|---|---|---|---|
| 1 | sampleID | forwardReads | reverseReads |
| 2 | SRR10070130 | s3://ngi-igenomes/test-data/ampliseq/SRR10070130_1.fastq.gz | s3://ngi-igenomes/test-data/ampliseq/SRR10070130_2.fastq.gz |
| 3 | SRR10070131 | s3://ngi-igenomes/test-data/ampliseq/SRR10070131_1.fastq.gz | s3://ngi-igenomes/test-data/ampliseq/SRR10070131_2.fastq.gz |
| 4 | SRR10070132 | s3://ngi-igenomes/test-data/ampliseq/SRR10070132_1.fastq.gz | s3://ngi-igenomes/test-data/ampliseq/SRR10070132_2.fastq.gz |
| 5 | SRR10070133 | s3://ngi-igenomes/test-data/ampliseq/SRR10070133_1.fastq.gz | s3://ngi-igenomes/test-data/ampliseq/SRR10070133_2.fastq.gz |
| 6 | SRR10070134 | s3://ngi-igenomes/test-data/ampliseq/SRR10070134_1.fastq.gz | s3://ngi-igenomes/test-data/ampliseq/SRR10070134_2.fastq.gz |
| 7 | SRR10070141 | s3://ngi-igenomes/test-data/ampliseq/SRR10070141_1.fastq.gz | s3://ngi-igenomes/test-data/ampliseq/SRR10070141_2.fastq.gz |
| 8 | SRR10070149 | s3://ngi-igenomes/test-data/ampliseq/SRR10070149_1.fastq.gz | s3://ngi-igenomes/test-data/ampliseq/SRR10070149_2.fastq.gz |
| 9 | SRR10070150 | s3://ngi-igenomes/test-data/ampliseq/SRR10070150_1.fastq.gz | s3://ngi-igenomes/test-data/ampliseq/SRR10070150_2.fastq.gz |
| 10 | SRR10070151 | s3://ngi-igenomes/test-data/ampliseq/SRR10070151_1.fastq.gz | s3://ngi-igenomes/test-data/ampliseq/SRR10070151_2.fastq.gz |
| 11 | SRR10102392 | s3://ngi-igenomes/test-data/ampliseq/SRR10102392_1.fastq.gz | s3://ngi-igenomes/test-data/ampliseq/SRR10102392_2.fastq.gz |
| 12 | SRR10102393 | s3://ngi-igenomes/test-data/ampliseq/SRR10102393_1.fastq.gz | s3://ngi-igenomes/test-data/ampliseq/SRR10102393_2.fastq.gz |
| 13 | SRR10102394 | s3://ngi-igenomes/test-data/ampliseq/SRR10102394_1.fastq.gz | s3://ngi-igenomes/test-data/ampliseq/SRR10102394_2.fastq.gz |

| | A | B | C | D | E |
|---|---|---|---|---|---|
| 1 | ID | name | habitat | Riv_vs_Gro | Sed_vs_Soil |
| 2 | SRR10070130 | SRR10070130-Riverwater | Riverwater | Riverwater | |
| 3 | SRR10070131 | SRR10070131-Riverwater | Riverwater | Riverwater | |
| 4 | SRR10070132 | SRR10070132-Groundwater | Groundwater | Groundwater | |
| 5 | SRR10070133 | SRR10070133-Groundwater | Groundwater | Groundwater | |
| 6 | SRR10070134 | SRR10070134-Riverwater | Riverwater | Riverwater | |
| 7 | SRR10070141 | SRR10070141-Groundwater | Groundwater | Groundwater | |
| 8 | SRR10070149 | SRR10070149-Sediment | Sediment | | Sediment |
| 9 | SRR10070150 | SRR10070150-Sediment | Sediment | | Sediment |
| 10 | SRR10070151 | SRR10070151-Sediment | Sediment | | Sediment |
| 11 | SRR10102392 | SRR10102392-Soil | Soil | | Soil |
| 12 | SRR10102393 | SRR10102393-Soil | Soil | | Soil |
| 13 | SRR10102394 | SRR10102394-Soil | Soil | | Soil |

# QIIME2

# QIIME2: Output artifacts and visualisations

- *.qza - zip folder, containing data and metadata
- *.qzv - zip folder, containind data, metadata and visualisations
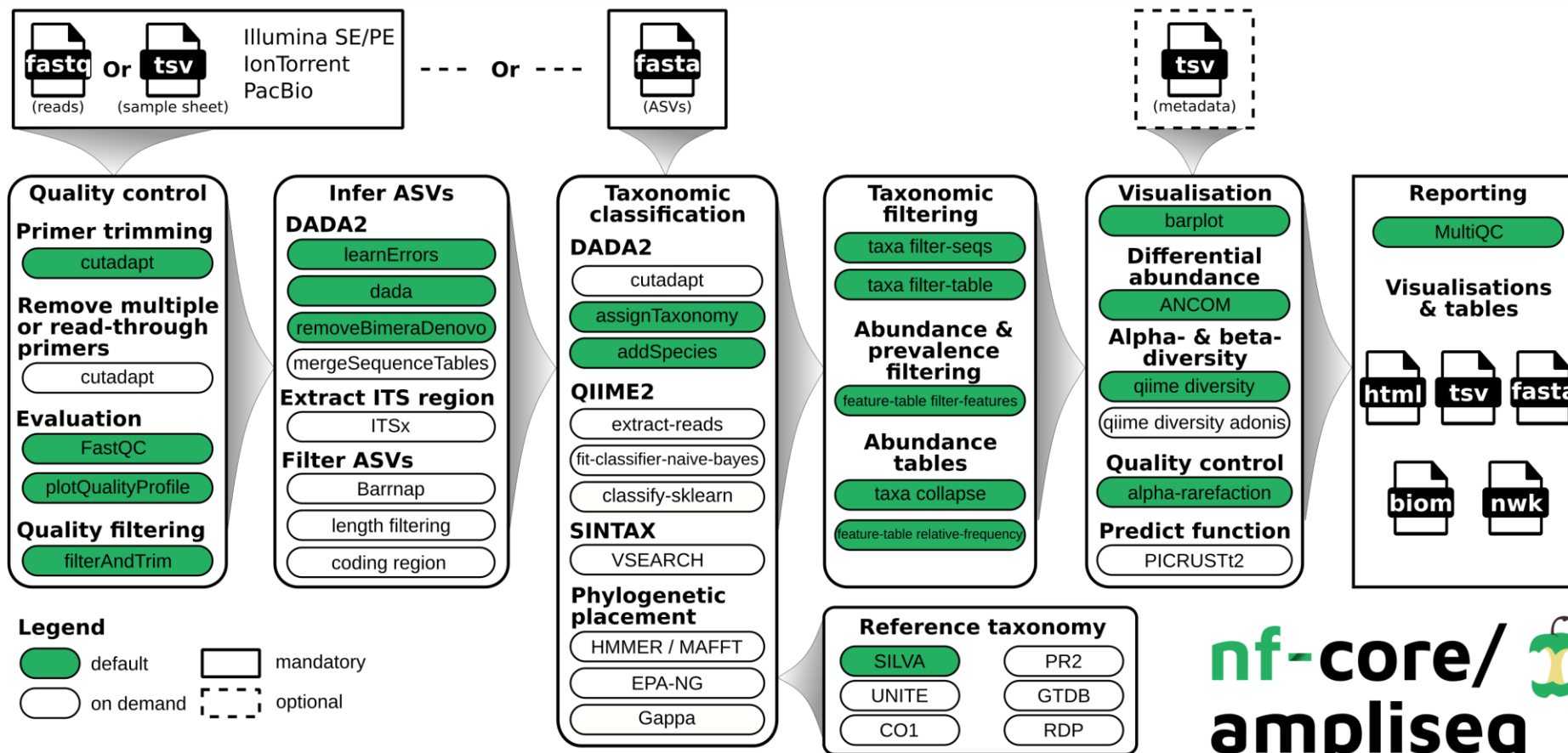
# QIIME2: pros & cons

**Pros**

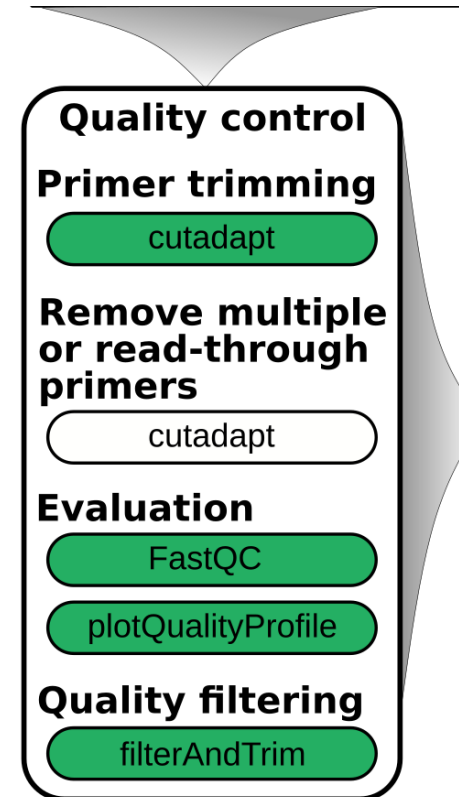- Comprehensive pipeline

- Plugins

- GUI

- Modularity

**Cons**

- Own data types

- Learning curve
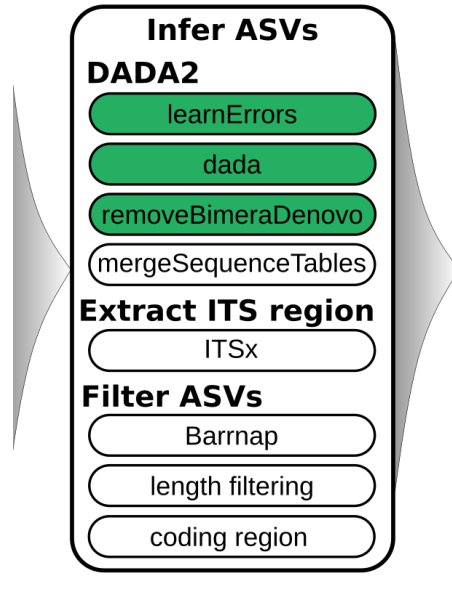
- Reference dependent

# nf-core/Ampliseq pipeline



CC-BY 4.0. Design originally by Zandra Fagernäs

# Ampliseq: quality control

- Data preprocessing

- Check reads quality

- Perform filt&trim

# Ampliseq: ASVs calculation



- DADA2
- Error estimation
- Chimera removal
- Contamination removal
- Filtering

# Ampliseq: Taxonomic classification



- Database dependent
- Infer species
- Confidence intervals
- Multiple assignment

# Ampliseq: Taxonomic filtering

**Taxonomic filtering**

- taxa filter-seqs
- taxa filter-table

**Abundance & prevalence filtering**

- feature-table filter-features

**Abundance tables**

- taxa collapse
- feature-table relative-frequency

- Filter specific taxa
- Abundance filtering

# Ampliseq: Post processing

- Visualisation

- Diversity computation

- Functional analysis

# Ampliseq: Visualisation

# Ampliseq: functional profiling

- Picrust2
- Phylogenetic Investigation of Communities by Reconstruction of Unobserved States

- KEGG and COG database
- Based on phylogeny
- Genes present in microbial genomes are similar amongst relatives
- When sufficient genome sequences are available, it is possible to predict which gene families are present in a given microbial OTU from phylogeny alone.

# Ampliseq: Report

# Ampliseq: pros & cons

**Pros**

- Standardized

- Easy to run

- Comprehensive analysis

- Community-driven

**Cons**

- Learning curve

- Resource intensive

- Needs setup

- Software versions dependent