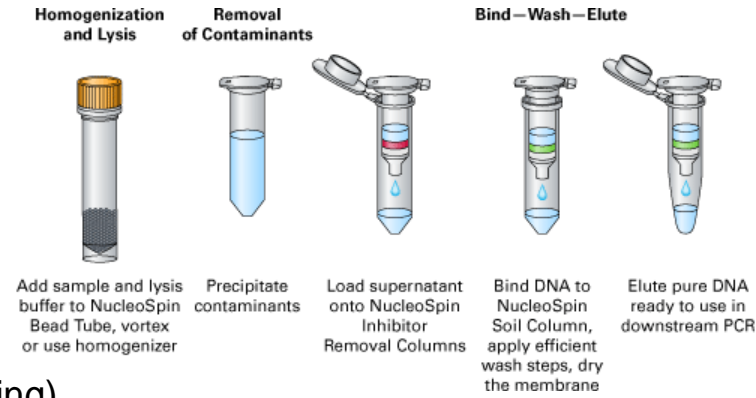# Bioinformatics and Microbiome Analysis MB140P94

# Amplicon data

**Tomáš Větrovský, Iñaki Odriozola and Petr Baldrian**
**Laboratory of Environmental Microbiology**
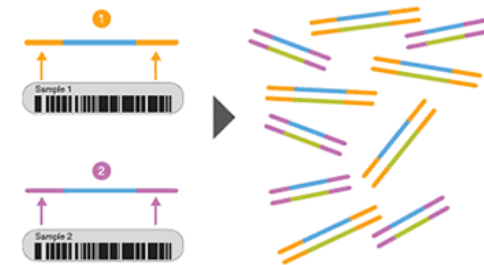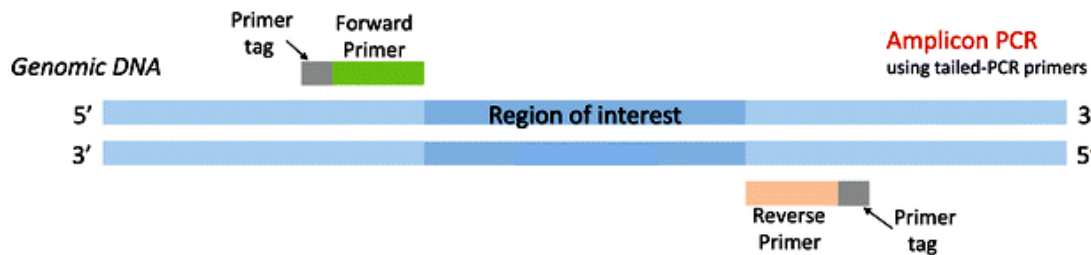**Institute of Microbiology of the CAS**

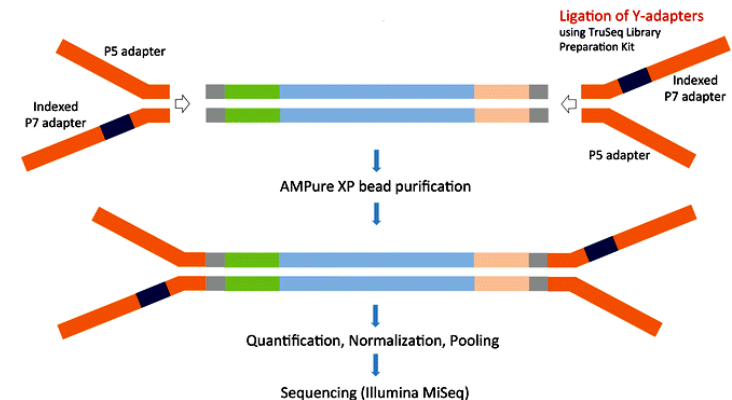# Illumina amplicon sequence library preparation

1. DNA isolation

Homogenization and Lysis — Removal of Contaminants — Bind—Wash—Elute

Add sample and lysis buffer to NucleoSpin Bead Tube, vortex or use homogenizer | Precipitate contaminants | Load supernatant onto NucleoSpin Inhibitor Removal Columns | Bind DNA to NucleoSpin Soil Column, apply efficient wash steps, dry the membrane | Elute pure DNA ready to use in downstream PCR

2. PCR with barcoded primers (multiplexing)

Primer tag — Forward Primer

Genomic DNA

5'  Region of interest  3'
3'  5'

Reverse Primer — Primer tag
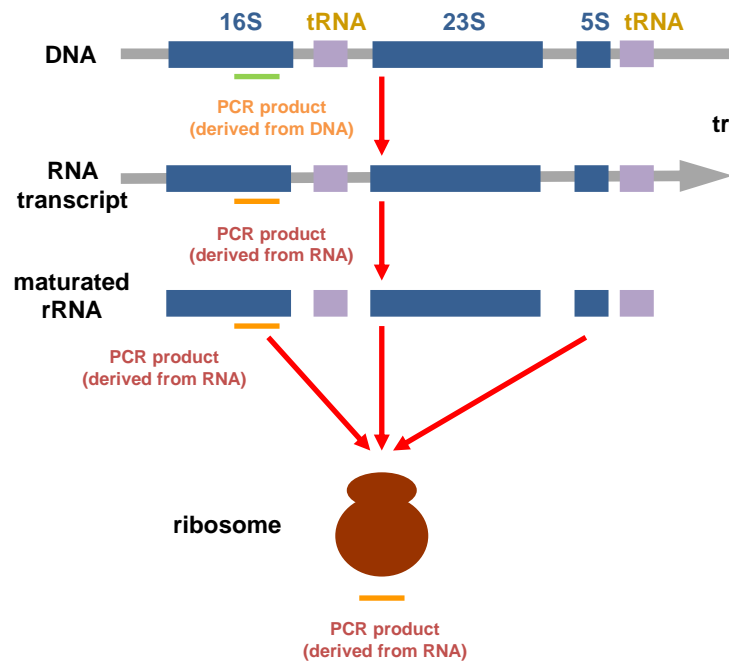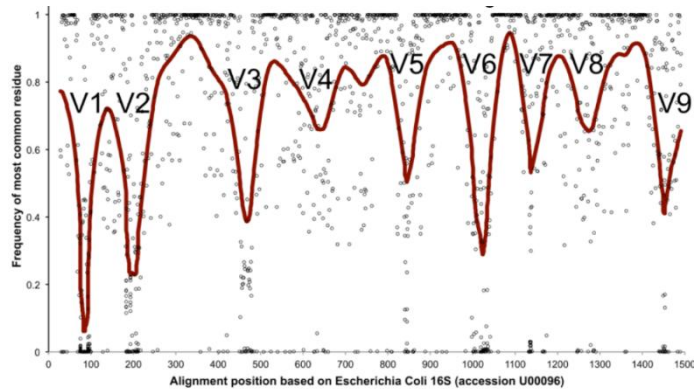
Amplicon PCR using tailed-PCR primers

Sample 1

Sample 2

3. Ligation of sequencing adapters – to attach short oligonucleotides (60bp) to your DNA fragments, these oligonucleotides are used to attach to the sequencing flow cell and they are also used as barcode of library

Ligation of Y-adapters using TruSeq Library Preparation Kit

P5 adapter

Indexed P7 adapter

Indexed P7 adapter

P5 adapter

AMPure XP bead purification

4. Quantification of the library by qPCR – to quantify of the exact amount of ligated fragments

Quantification, Normalization, Pooling

Sequencing (Illumina MiSeq)

# Most used marker genes

Bacterial 16S ribosomal RNA gene

Fungal internal transcribed spacer (ITS)



1-15 copies of rDNA per genome

X0–X00? copies of ITS per genome
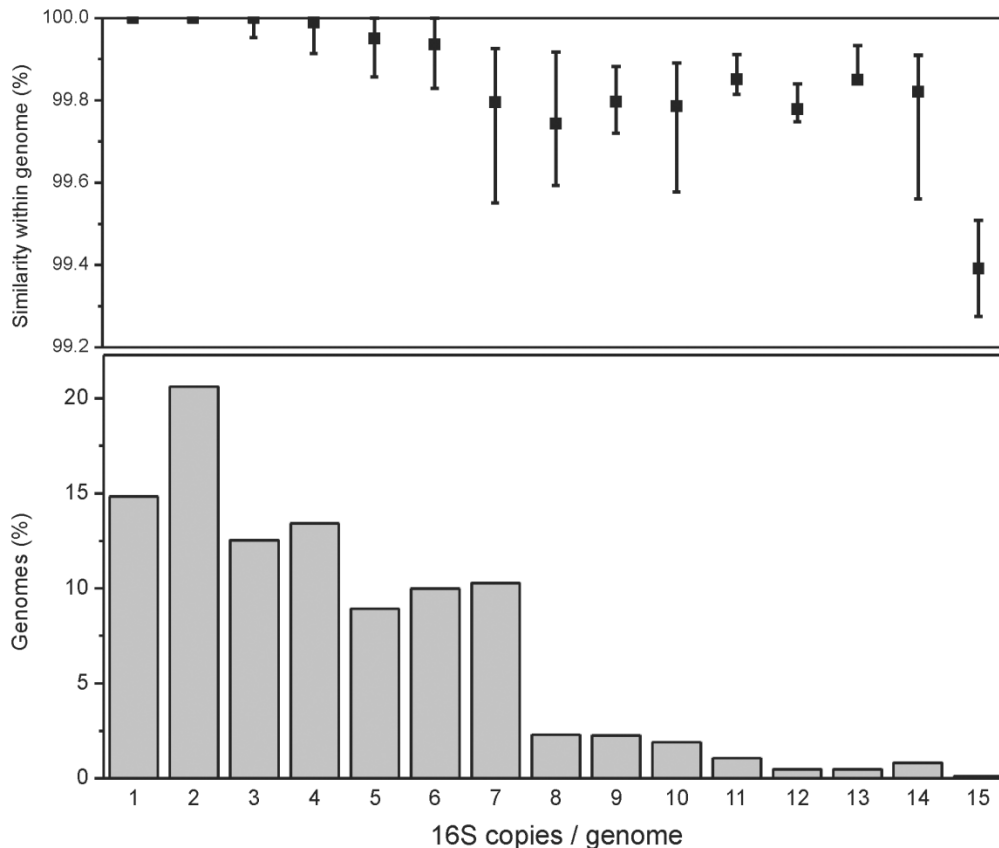
# 16S rDNA gene vs. alternative (low-copy) markers

**Pros: highly populated reference databases**

**Cons: muticopy nature of bacterial 16S rDNA gene**
- possibility of high intragenomic variability - diversity over estimation (number of OTUs)
- relative abundance estimation is skew -> normalisation by 16S copy number of closest taxon



16S rRNA within-genome similarity and copy numbers in bacterial genomes.

Upper panel: the similarity of genomes with various copy numbers: the values indicated represent the first, the second and the third quartile.

Lower panel: distribution of 16S rRNA copy numbers per genome in 1,690 sequenced bacterial genomes.

T. Větrovský & P. Baldrian - PloS one, 2013

# Sequencing platforms for amplicon sequencing (most used)



**454 Pyrosequencing**

Not supported anymore
(most studies 2009-2012)

Errors in homopolymeric regions

long reads
(up to 700 bp)

**Illumina**

The most used

Error rate less than 1%

pair-end data
(Illumina)

**IonTorrent**

Very chaep sequencing

Lot of errors

Medium read size
200-600 bp

**PacBio**

Still rare

Repeated sequncing of the same region

extra long reads
(up to 10.000 bp)

**Sequencing platforms for amplicon sequencing (based on GlobalFungi data sources)**

# Searching for true number of taxa (OTUs) in the data



Quince, Christopher, et al. "Accurate determination of microbial diversity from 454 pyrosequencing data." *Nature methods* 6.9 (2009): 639.

# Searching for true number of taxa (OTUs) in the data



Quince, Christopher, et al. "Accurate determination of microbial diversity from 454 pyrosequencing data." *Nature methods* 6.9 (2009): 639.

This is still a big concern to all microbial community analysis pipelines!

# Amplicons pipeline workflow

```
input data                          joining paired-ends
(paired-end FASTQ files)
         │                                   │
         │                                   ▼
         ▲                           quality filtering
         │                                   │
         │                                   ▼
   DNA sequencing             demultiplex samples,
         ▲                    orient sequences
         │                    and remove primers
         │                                   │
                                             ▼
   samples properties         remove too short and too long sequences
                                             │
                                             ▼
                             clustering to OTUs
                             and removal of chimeric sequences
```
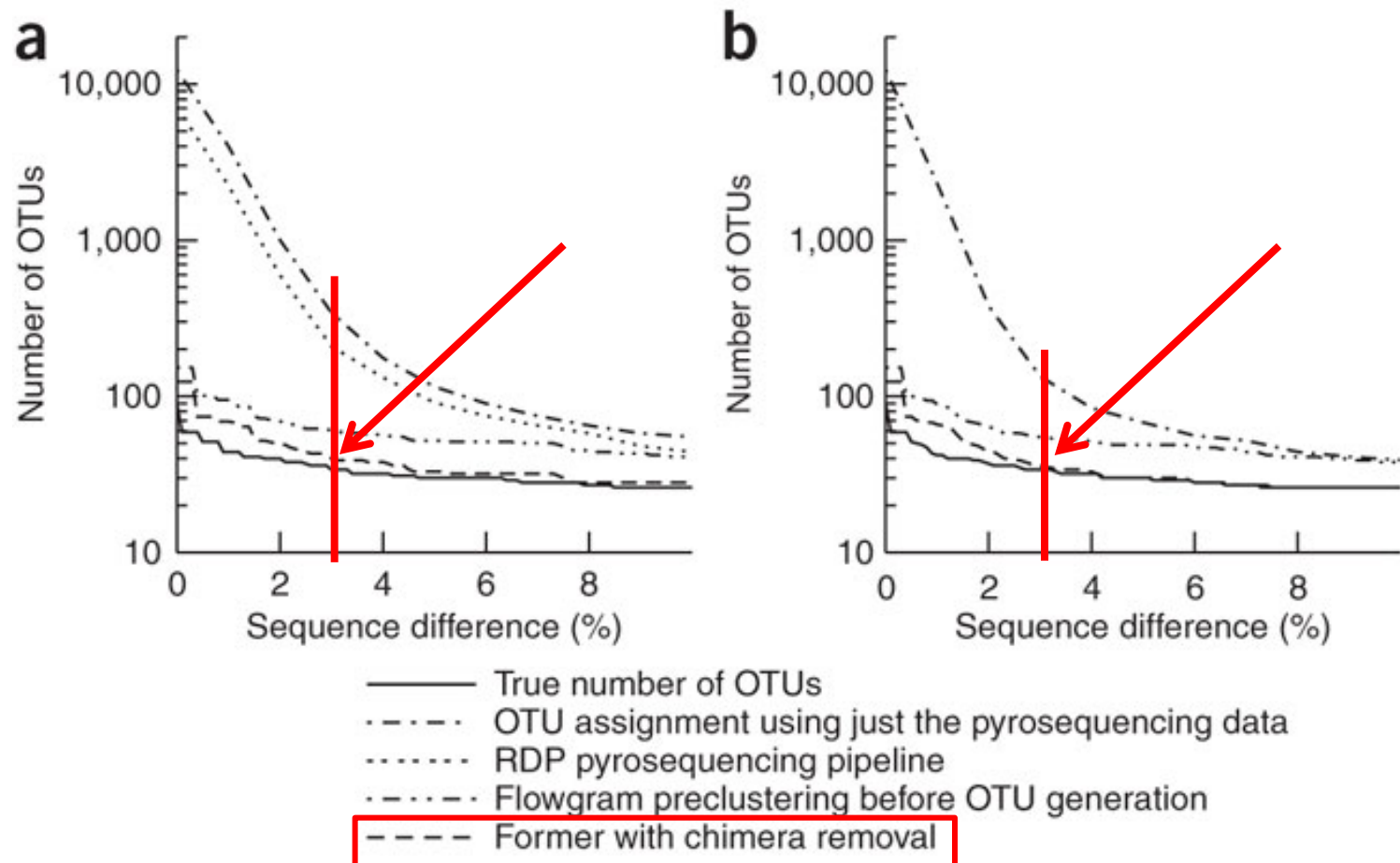
getting the representative sequences from the clusters

construction of OTU table

create phylogeny tree

Identification of OTUs

normalize table by sample size (16S copy number)

normalize (subsample) all samples to the same number of reads

metadata

beta diversity estimates (Fast Parallel UniFrac)

alpha diversity estimates

**analysis with R phyloseq library**

# GUI based alternative for Windows (http://www.biomed.cas.cz/mbu/lbwrf/seed/)

## SEED 2: a user-friendly platform for amplicon high-throughput sequencing data analyses



- editing of sequences and their titles
- sorting
- quality trimming
- pair-end joining
- grouping of sequences based on sequence motifs or sequence titles
- batch processing of sequence groups
- denoising
- chimera removal
- ITS extraction
- sequence alignments and clustering
- OTU table construction
- construction of consensus sequences
- creation of local databases for BLAST
- searching either local databases or the whole NCBI
- retrieval of taxonomical classification from the NCBI
- calculation of diversity parameters
- many more...

Větrovský, Baldrian & Morais (2018) SEED 2: a user-friendly platform for amplicon high-throughput sequencing data analyses. Bioinformatics, bty071, 2018

# SEED is alternative to

QIIME 2™ is a next-generation microbiome bioinformatics platform that is extensible, free, open source, and community developed.
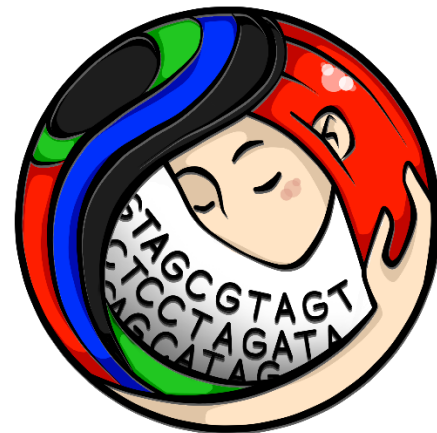
https://qiime2.org/



- **command line, Unix (Linux) based**

https://mothur.org/

mothur - one of the most widely used tools for analyzing 16S rRNA gene sequence data



- **command line, multiplatform**

## RAW DATA - R1 & R2 FASTQ

- Joining of pair-end data
- Quality filtering / Sample determination / Sequence trimming
- Chimera removal
- Preparing for clustering e.g.: fungal ITS extraction
- Clustering to OTUs
- Getting of the representative sequences from the clusters
- Identification of OTUs
- Construction of OTU table
- Estimation of diversity indices
- PROCESSING OF THE RESULTS

# RAW DATA

## BAC_R1.fastq - first sequence

@M03794:8:000000000-AJCUU:1:2114:9990:17907 1:N:0:7
AACAGCCGGACTACTGGGGGTTTCTAATCCTGTTTGCTCCCCACGCTTTCGTGCCTCAGTGTCAATGACCGTGTAGC
AAGCTGCCTTCGCAATTGGTGTTCTATGTCATATCTAAGCATTTCACCGCTACATGACATATTCCGCTTACCTCCAC
GATATTCAAGACTAATAGTATCAATGGCAGTTCCCAAGTTAAGCTCGGGGATTTCACCACGGACTTACTAGCCCACC
TACGCACCCTTTAAACCCAGT
+
BCCCCFCCCCCCGGGGGGGGFGGHHHCHHHHHHHHHHHHHGG2FGGGGGHGHGGHHHFHHHHHHHHHHHHHGGH
HHHHHHHHHHHHHHHGGGGHHHHGHGHHHHHHHHHFHHHHHHHHGHHHHHGGGGGHHHGHHHHHHHHGG
GGHHHHHHHGCGFDHHGGHHHHHHHHGG>GGGHHG/GHGHHHHHHFHFHGHHGHHG?DGGGGGGGGGGGGGAC
GGGGGGGGGGFGGFAAFF;@ADFFFFFFB/FFFFF;

## BAC_R2.fastq - first sequence

@M03794:8:000000000-AJCUU:1:2114:9990:17907 2:N:0:7
ACGAAGTGTGCCAGCAGCCGCGGTAATACGGAGGGTGCAAGCGTTATCCGGATTCACTGGGTTTAAAGGGGTGC
GTAGGTGGGCTAGTAAGTCCGTGGTGAAATCCCCGAGCTTAACTTGGGAACTGCCATTGATACTATTAGTCTTGA
ATATCGTGGAGGTAAGCGGAATATGTCATGTAGCGGTGAAATGCTTAGATATGACATAGAACACCAATTGCGAAG
GCAGCTTGCTACACGGTCATTGACACTG
+
BBBBBBBFFBFDGFFFGGGCGGCGGGHHHGDEDGGGFGGHHHGGGGGFDEE?EEGBGFHHHHGFFGGFFGGFE
FGDFDGFFEGGHHHGFEGFHFHEEFFFHHHHFHGGGFGFHHHFHHHHHHBGHHFBCHFHBGGGHHHHGHHHHH
HGFFHHHGHGD<GGAGHHHGGG@CFHHFCGHH:CCG?AAAGFEFEFGGGBFFFFFFGFFBBFGFFBDE?BBBB.@9
-9..A.B/:AFFFEF.@;AAF//99FFF/

# Joining of pair-end data (Illumina)

**Joining of pair-end data**

**Quality filtering Sample determination Sequence trimming**

**Chimera removal**

**Preparing for clustering e.g.: fungal ITS extraction**

**Clustering to OTUs**

**Getting of the representative sequences from the clusters**

**Identification of OTUs**

**Construction of OTU table**

**Estimation of diversity indices**

**PROCESSING OF THE RESULTS**

reads quality visualisation
usually: 0-40
darker = lower

**reads quality is dropping at the ends**



reconstructed sequence are based on bases with higher quality

PROGRAM: **FastqJoin**
set a minimal overlap length and overlap precision

# Sequence quality and quality filtering

| RAW DATA - R1 & R2 FASTQ |
|---|
| Joining of pair-end data |
| **Quality filtering Sample determination Sequence trimming** |
| Chimera removal |
| Preparing for clustering e.g.: fungal ITS extraction |
| Clustering to OTUs |
| Getting of the representative sequences from the clusters |
| Identification of OTUs |
| Construction of OTU table |
| Estimation of diversity indices |
| PROCESSING OF THE RESULTS |

FASTQ format

```
@FORJUSP02AJWD1
CCGTCAATTCATTTAAGTTTTAACCTTGCGGCCGTACTCCCCAGGCGGT
+
AAAAAAAAAAAA::99@::::??@@::FFAAAAACCAA::::BB@@?A?
```

Label

Sequence

Q scores (as ASCII chars)

Base=T, Q=':'=25

$$Q_{illumina} = -10 \times \log_{10}\left(\frac{P_e}{1 - P_e}\right),$$

where $P_e$ is the probability of identifying a base incorrectly. For Sanger and other platforms, the formula is as follows [8]:

$$Q_{PHRED} = -10 \times \log_{10}(P_e).$$

$$Q_{illumina} = 10 \times \log_{10}\left(10^{\left\{\frac{Q_{PHRED}}{10}\right\}} + 1\right)$$

Table 2. Phred quality scores are logarithmically linked to error probabilities (http://en.wikipedia.org/wiki/Phred_quality_score)

| Phred quality score | Probability of incorrect base call | Base call accuracy |
|---|---|---|
| 10 | 1 in 10 | 90% |
| 20 | 1 in 100 | 99% |
| 30 | 1 in 1000 | 99.90% |
| 40 | 1 in 10 000 | 99.99% |
| 50 | 1 in 100 000 | 99.999% |
| 60 | 1 in 1 000 000 | 99.9999% |

## Flowchart (left column)

RAW DATA - R1 & R2 FASTQ

↓

Joining of pair-end data

↓

**Quality filtering**
**Sample determination**
**Sequence trimming**

↓

Chimera removal

↓

Preparing for clustering e.g.: fungal ITS extraction

↓

Clustering to OTUs

↓

Getting of the representative sequences from the clusters

↓

Identification of OTUs

↓

Construction of OTU table

↓

Estimation of diversity indices

↓

PROCESSING OF THE RESULTS

## Multiple samples in one library (Multiplexing)

| name | FWDprimer | REVprimer |
|------|-----------|-----------|
| SAMPLE001 | 515F_T103 | 806R_T007 |
| SAMPLE002 | 515F_T002 | 806R_T052 |

**spacer** is not presented in native sequences, it is used to prevent overestimation of any taxa

**TAG** **SPACER** **ORIGINAL PRIMER**

515F_T103   AATATACGTGTGCCAGCMGCCGCGGTAA
515F_T002   ACGAAGTGTGCCAGCMGCCGCGGTAA

806R_T007   AGCCACCGGACTACHVGGGTWTCTAAT
806R_T052   ATCCTCCCGGACTACHVGGGTWTCTAAT

AATATAC          AGCCA

forward tags          reverse tags

ACGAA          ATCCTC

## Process flow (left column)

- **RAW DATA - R1 & R2 FASTQ**
- Joining of pair-end data
- **Quality filtering / Sample determination / Sequence trimming**
- Chimera removal
- Preparing for clustering e.g.: fungal ITS extraction
- Clustering to OTUs
- Getting of the representative sequences from the clusters
- Identification of OTUs
- Construction of OTU table
- Estimation of diversity indices
- **PROCESSING OF THE RESULTS**

# Sample determination (de-multiplexing)
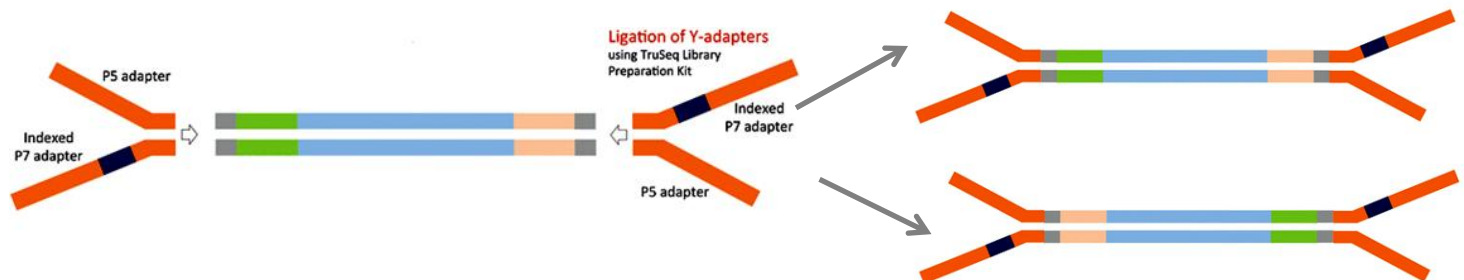
## Demultiplex samples

put sample names to sequence titles

>**SAMPLE034**|M03794:8:0000000
00-AJCUU:1:2114:9990:17907
CCTGTTTGCTCCCCACGCTTTC
GTGCCTCAGTGTCAATGACCGT
GTAGCAAGCTGCA…

## Orient sequences

P5 adapter

Indexed P7 adapter

**Ligation of Y-adapters** using TruSeq Library Preparation Kit

Indexed P7 adapter

P5 adapter

**cca 50 % of the reads are reverse complement oriented due to ligation of library adapters**

## Remove primers

since primer sequences are not native to the sample, they need to be removed before clustering to OTUs

# Sample determination (de-multiplexing) and removing barcodes and primers

| RAW DATA - R1 & R2 FASTQ |
|---|

| Joining of pair-end data |
|---|

| **Quality filtering Sample determination Sequence trimming** |
|---|

| Chimera removal |
|---|

| Preparing for clustering e.g.: fungal ITS extraction |
|---|

| Clustering to OTUs |
|---|

| Getting of the representative sequences from the clusters |
|---|

| Identification of OTUs |
|---|

| Construction of OTU table |
|---|

| Estimation of diversity indices |
|---|

| PROCESSING OF THE RESULTS |
|---|

GAGCGTGA      gITS7_T02
TGGCGTGA      gITS7_T03
GGTGCGTGA     gITS7_T06
AAAGCGTGA     gITS7_T08
TCTAGCGTGA    gITS7_T10

search for the sample barcodes at the beginning of reads

| Sequence | Query | RESULT |
|---|---|---|
| GAGCGTGA | gITS7_T02 | 14687 |
| TGGCGTGA | gITS7_T03 | 22568 |
| GGTGCGTGA | gITS7_T06 | 16835 |
| AAAGCGTGA | gITS7_T08 | 19258 |
| TCTAGCGTGA | gITS7_T10 | 20585 |

remove barcode after putting the name of sample to sequences header…

## Sequence trimming

RAW DATA - R1 & R2 FASTQ

↓

Joining of pair-end data

↓

Quality filtering
Sample determination
Sequence trimming

↓

Chimera removal

↓

Preparing for clustering
e.g.: fungal ITS
extraction

↓

Clustering to OTUs

↓

Getting of the
representative
sequences from
the clusters

↓

Identification of
OTUs

↓

Construction of
OTU table

↓

Estimation of
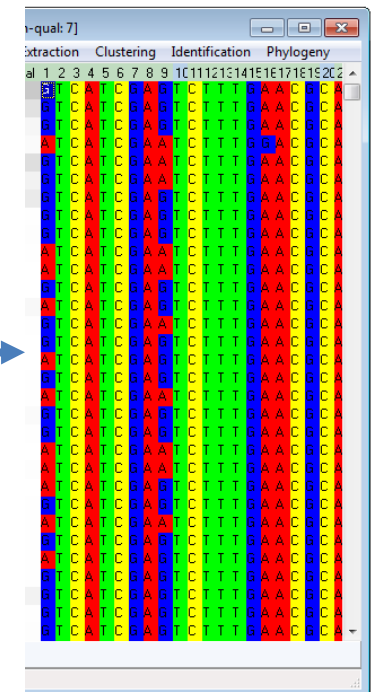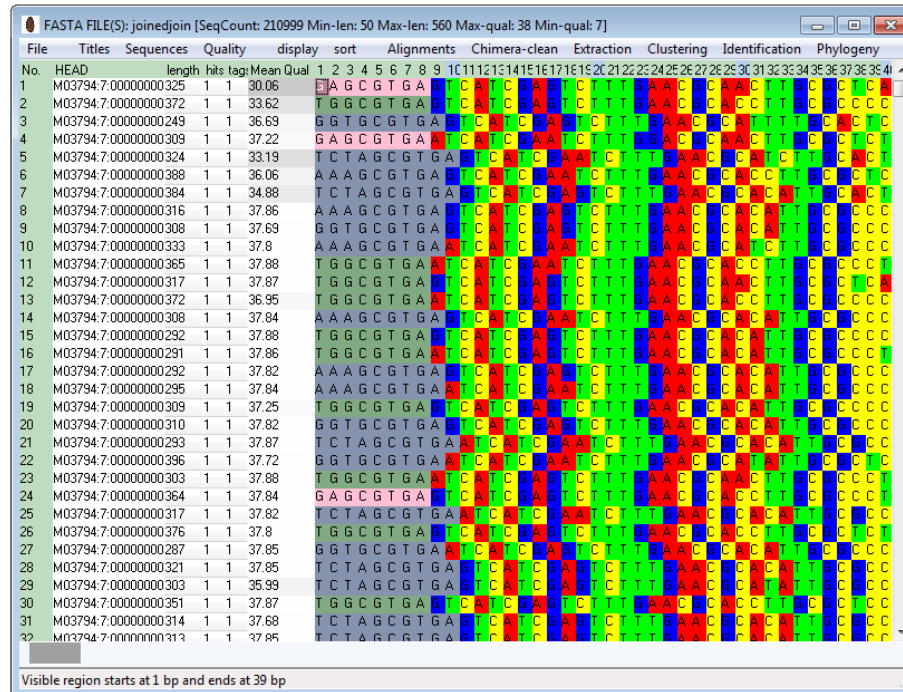diversity indices

↓

PROCESSING OF THE
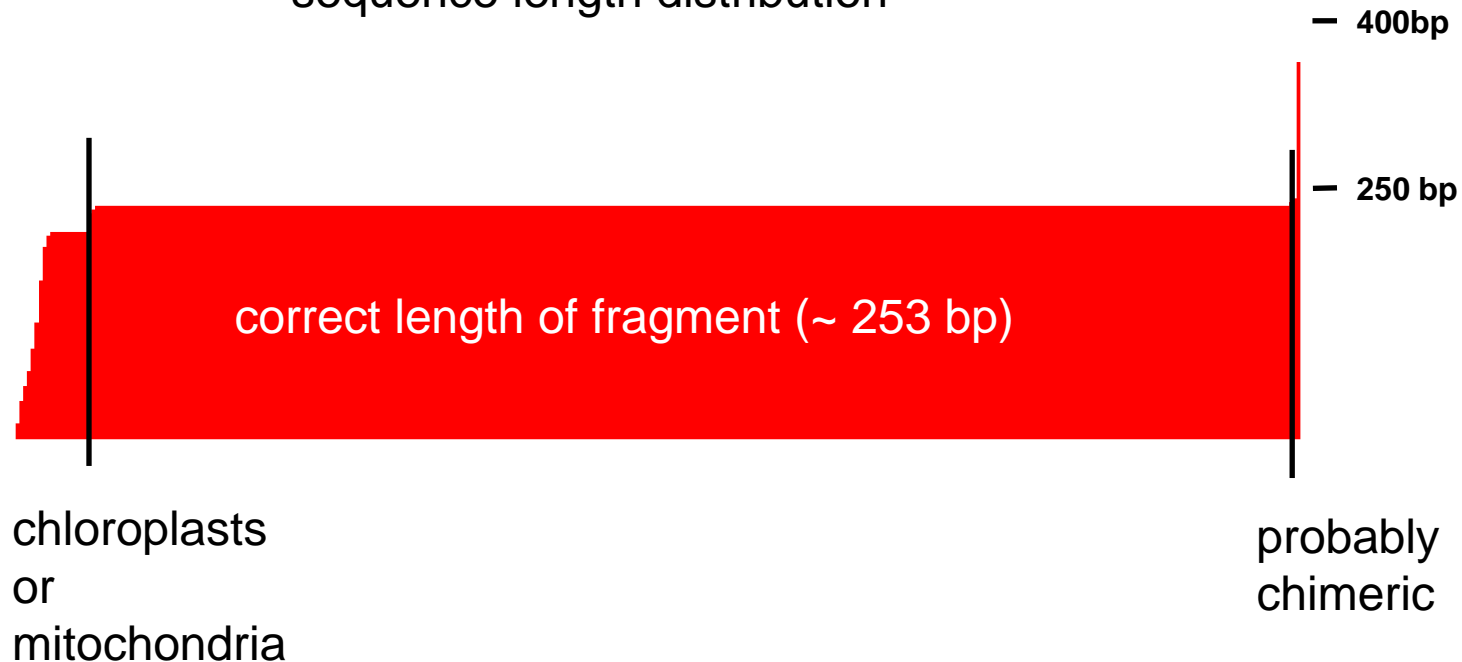RESULTS

removing sequences with aberrated length - depends on the marker gene

Too short – nonspecific PCR products/erroneus sequences
Too long - nonspecific PCR products/chimeric sequences

sequence length distribution

— 400bp

— 250 bp

correct length of fragment (~ 253 bp)

chloroplasts
or
mitochondria

probably
chimeric

Chimera removal

**Process flowchart (left column):**
- RAW DATA - R1 & R2 FASTQ
- Joining of pair-end data
- Quality filtering / Sample determination / Sequence trimming
- Chimera removal
- Preparing for clustering e.g.: fungal ITS extraction
- Clustering to OTUs
- Getting of the representative sequences from the clusters
- Identification of OTUs
- Construction of OTU table
- Estimation of diversity indices
- PROCESSING OF THE RESULTS

**Diagram labels:**
- Extension is aborted.
- Single stranded DNA template (strain A)
- Next round of PCR
- (2) Some mismatches occur here, but annealing happens
- (1) This incomplete product acts as a primer
- (3) Extension occurs using strain B as a template.
- Single stranded DNA template (strain B)
- Next round of PCR
- A chimera is formed.
- Strain A
- Strain B

# Chimera removal

Joining of pair-end data

Quality filtering
Sample determination
Sequence trimming

Chimera removal

Preparing for clustering
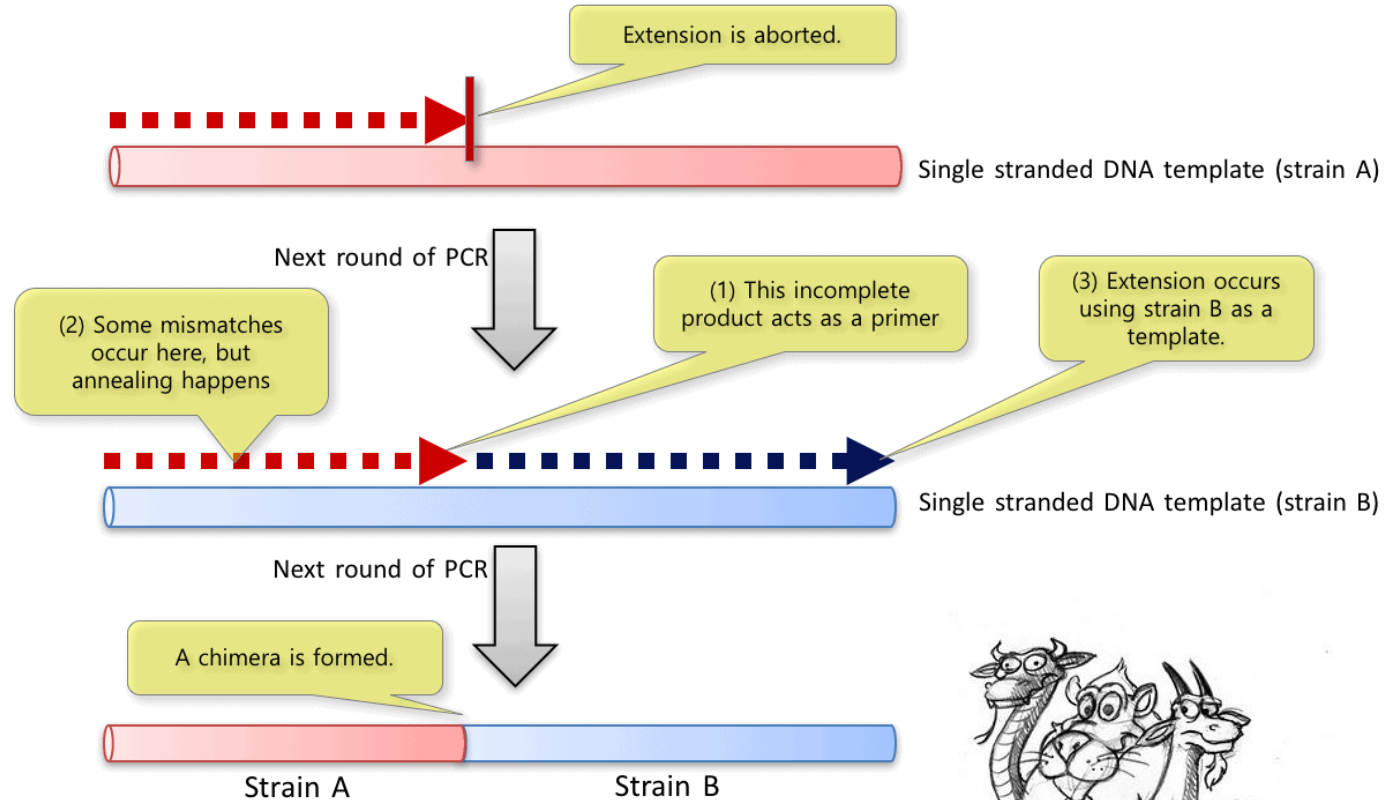e.g.: fungal ITS
extraction

Clustering to OTUs
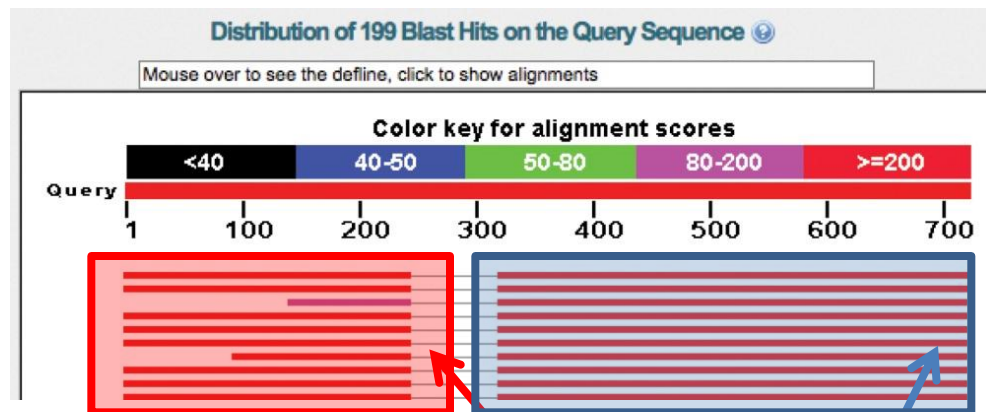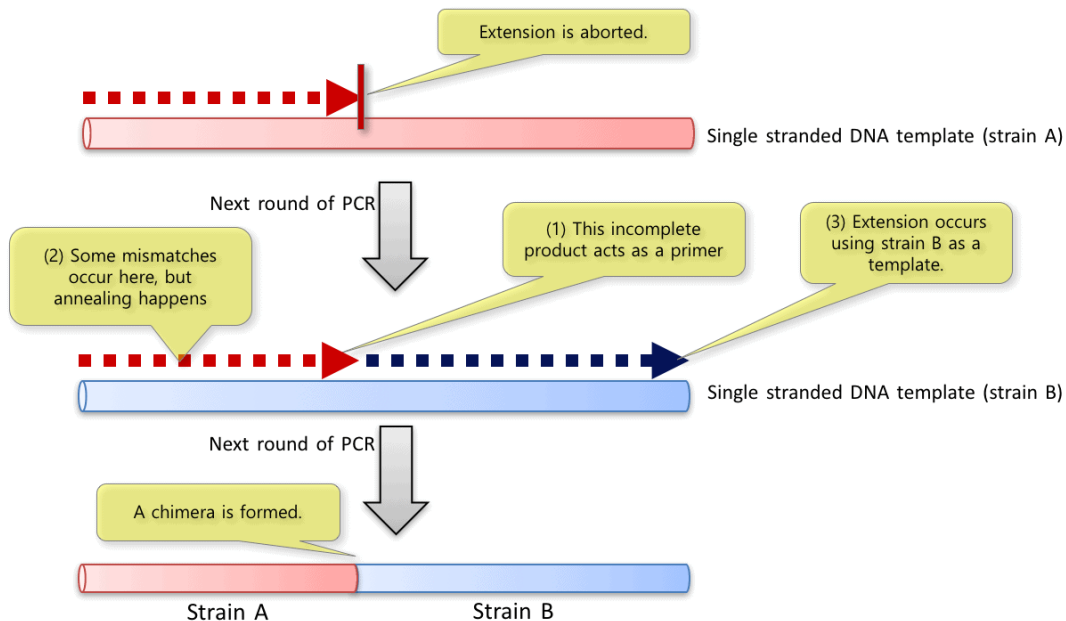
Getting of the
representative
sequences from
the clusters

Identification of
OTUs

Construction of
OTU table

Estimation of
diversity indices

PROCESSING OF THE
RESULTS

Extension is aborted.

Single stranded DNA template (strain A)

Next round of PCR

(2) Some mismatches occur here, but annealing happens

(1) This incomplete product acts as a primer

(3) Extension occurs using strain B as a template.

Single stranded DNA template (strain B)

Next round of PCR

A chimera is formed.

Strain A    Strain B

Distribution of 199 Blast Hits on the Query Sequence

Mouse over to see the defline, click to show alignments

Color key for alignment scores

| <40 | 40-50 | 50-80 | 80-200 | >=200 |

Query

1    100    200    300    400    500    600    700

Parental sequence 1    AAAAAAAAAAAAAAAAAAAA
Parental sequence 2    TTTTTTTTTTTTTTTTTTTT
Chimeric sequence      AAAAAAAATATTTTTTTTTT

# Chimera removal

**chimeric sequences are derived from parental sequences**
**-> program search for parents.**

## Flowchart (left column)

- RAW DATA - R1 & R2 FASTQ
- Joining of pair-end data
- Quality filtering / Sample determination / Sequence trimming
- **Chimera removal**
- Preparing for clustering e.g.: fungal ITS extraction
- Clustering to OTUs
- Getting of the representative sequences from the clusters
- Identification of OTUs
- Construction of OTU table
- Estimation of diversity indices
- PROCESSING OF THE RESULTS

## *de novo* approach:

- no reference database
- usually based on measuring of several sequence parts abundances – parents should be more abundant then its offspring (chimeras)

## reference based approach:

- pair-wise alignment with reference sequence database

candidate parents

query

A
Q
B
*Local*

A
Q
B
*Local-X*

A
Q
B
*Global-X*

Pair-wise alignments with all reference sequences

TGGGCGTAAAGCGCGCGTAGG   *Read*
TGtagcggtgaaacGCGTAGa   *Refseq 1*

TGGGCGTAAA   GCGCGCGTAGG   *Read*
TGGGCGTAAA   taGttacTgat   *Refseq 2*

+3                          +1

TGGGCGTAAA   GCGCGCGTAGG   *Read*
caagCatcAt   GCGCGCGTcGG   *Refseq 3*

Model with minimum score = nr. diffs + 3 x crossovers = 4.

TGGGCGTAAA   GCGCGCGTAGG   *Read*
TGGGCGTAAA   GCGCGCGTcGG   *Model m=2, d=1*

PROGRAMS: **UCHIME, UPARSE**

## Flowchart (left column)

- **RAW DATA - R1 & R2 FASTQ**
- Joining of pair-end data
- Quality filtering / Sample determination / Sequence trimming
- **Chimera removal**
- Preparing for clustering e.g.: fungal ITS extraction
- Clustering to OTUs
- Getting of the representative sequences from the clusters
- Identification of OTUs
- Construction of OTU table
- Estimation of diversity indices
- **PROCESSING OF THE RESULTS**

# Chimera removal

**chimeric sequences are derived from parental sequences -> program search for parents.**

## *de novo* approach:

- no reference database
- usually based on measuring of several sequence parts abundances – parents should be more abundant then its offspring (chimeras)

## reference based approach:

- pair-wise alignment with reference sequence database

## Problems

- algorithms are not optimal
- computation cost could be high
- problems with highly similar sequences

## Problems

- no appropriate reference database for environmental samples
- variable quality of reference databases

PROGRAMS: **UCHIME, UPARSE**

# Fungal ITS extraction

```
RAW DATA - R1 & R2
FASTQ
      ↓
Joining of pair-end data
      ↓
Quality filtering
Sample determination
Sequence trimming
      ↓
Chimera removal
      ↓
Preparing for clustering
e.g.: fungal ITS
extraction
      ↓
Clustering to OTUs
      ↓
Getting the
representative
sequences from
clusters
      ↓
Identification of
OTUs
      ↓
Construction of
OTU table
      ↓
Estimation of
diversity indices
      ↓
PROCESSING OF THE
RESULTS
```



**ITSx**
Improved software for detection and extraction of ITS1 and ITS2 from ribosomal ITS sequences of fungi and other eukaryotes for use in environmental sequencing

relies on **HMMER**
searching sequence databases for sequence homologs, and for making sequence alignments. It implements methods using probabilistic models called profile hidden Markov models (profile HMMs).

using the extracted variable ITS improves resolution when the OTUs are created

**http://microbiology.se/software/itsx/**

# Clustering to OTUs

**OTU (Operational taxonomic unit)** group of similar sequences grouped based on some similarity threshold
usually 97% similarity (16S, ITS) represents Species

## Heuristic

- comparison of each sequence with representative sequence ("seed")
- depends on sequence order

FAST

## Hierarchical

comparison of each sequence with each other (tree contruction)

SLOW

## Model based

- probabilistic, iterative
- uses more information than the sequence identity

VERY SLOW

---

Flowchart (left column):

- RAW DATA - R1 & R2 FASTQ
- Joining of pair-end data
- Quality filtering / Sample determination / Sequence trimming
- Chimera removal
- Preparing for clustering e.g.: fungal ITS extraction
- **Clustering to OTUs**
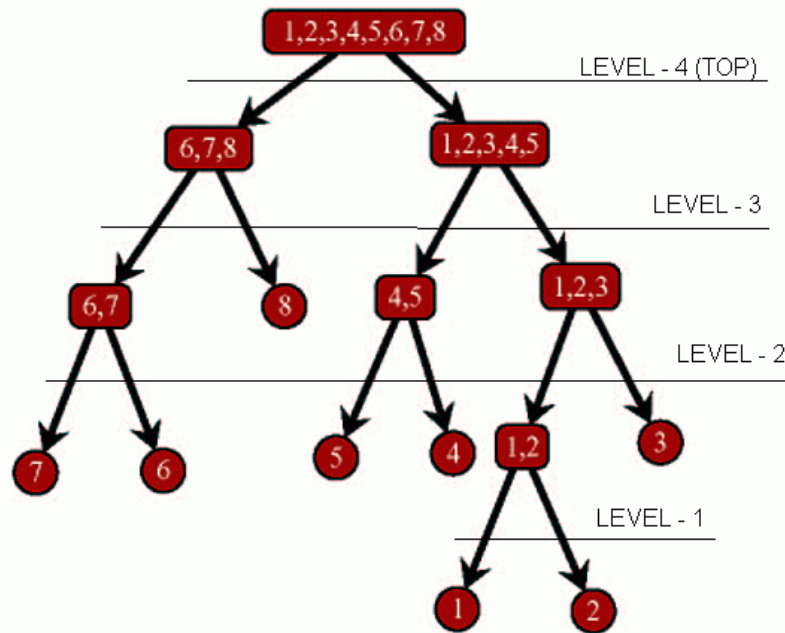- Getting the representative sequences from clusters
- Identification of OTUs
- Construction of OTU table
- Estimation of diversity indices
- PROCESSING OF THE RESULTS

# Clustering to OTUs

## programs

| RAW DATA - R1 & R2 FASTQ |
|---|
| ↓ |
| Joining of pair-end data |
| ↓ |
| Quality filtering / Sample determination / Sequence trimming |
| ↓ |
| Chimera removal |
| ↓ |
| Preparing for clustering e.g.: fungal ITS extraction |
| ↓ |
| Clustering to OTUs |
| ↓ |
| Getting the representative sequences from clusters |
| ↓ |
| Identification of OTUs |
| ↓ |
| Construction of OTU table |
| ↓ |
| Estimation of diversity indices |
| ↓ |
| PROCESSING OF THE RESULTS |

| Tool | Distance calculation | Clustering algorithm | Reference |
|---|---|---|---|
| DOTUR (+ MUSCLE) | $MSA_{denovo}$ | Hierarchical | Schloss et al. 2005 |
| MOTHUR | $MSA_{profile}$ | Hierarchical | Schloss et al. 2009 |
| ESPRIT | PSA | Hierarchical | Sun et al. 2009 |
| SLP | PSA | Hierarchical | Huse et al. 2009 |
| ESPRIT-TREE | PSA | Hierarchical | Cai et al. 2011 |
| JMOTU | PSA | Hierarchical | Jones et al. 2011 |
| CD-HIT | PSA | Heuristic | Li et al. 2006 |
| USEARCH/UPARSE | PSA | Heuristic | Edgar et al. 2010/2013 |
| GRAMCLUSTER | PSA | Heuristic | Russell et al. 2010 |
| DNACLUST | PSA | Heuristic | Ghodsi et al. 2011 |
| CRUNCHCLUST | PSA | Heuristic | Hartmann et al. 2012 |
| DYSC | PSA | Heuristic | Zheng et al. 2012 |
| MS-CLUST | PSA | Heuristic | Chen et al. 2013 |
| TBC | PSA | Heuristic | Lee et al. 2012 |
| TSC | PSA | H&H combination | Jiang et al. 2012 |
| CROP | PSA | Model-based (BC) | Hao et al. 2011 |
| BEBAC | PSA | Model-based (BC) | Cheng et al. 2012 |
| DBC454 | composition | Model-based | Pagni 2013 |
| DBC | PSA | Model-based | Preheim et al. 2013 |
| M-PICK | graphical | Model-based | Wang et al. 2013 |

```
PSA - Pairwise Sequence Alignment
MSA - Multiple Sequence Alignment
```

## Clustering to OTUs (hierarchical)

### Flowchart

- RAW DATA - R1 & R2 FASTQ
- Joining of pair-end data
- Quality filtering / Sample determination / Sequence trimming
- Chimera removal
- Preparing for clustering e.g.: fungal ITS extraction
- **Clustering to OTUs**
- Getting the representative sequences from clusters
- Identification of OTUs
- Construction of OTU table
- Estimation of diversity indices
- PROCESSING OF THE RESULTS



all pairwise comparisons are performed and OTUs are delineated at fixed distance level

**linking method is an important driver of the outcome:**

- **single-linkage clustering** (SL) clusters may be merged together due to single sequences being close to each other, even though many of the sequences in each cluster may be very distant to each other
- **complete-linkage clustering** (CL) tends to find compact clusters of approximately equal diameters. With CL, all objects in a cluster are similar to each other
- **average-linkage clustering** (AL) can be seen as an intermediate between single and complete linkage clustering, resulting in more homogeneous clusters than those obtained by the single-linkage method

# Clustering to OTUs (heuristic - USEARCH)

**RAW DATA - R1 & R2 FASTQ**

Joining of pair-end data

Quality filtering
Sample determination
Sequence trimming

Chimera removal

Preparing for clustering
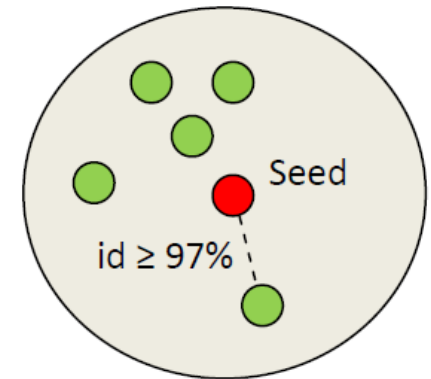e.g.: fungal ITS
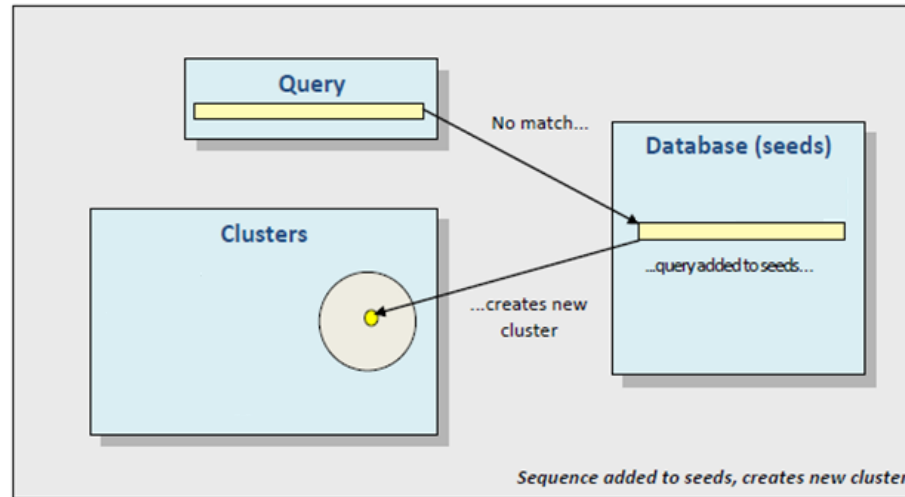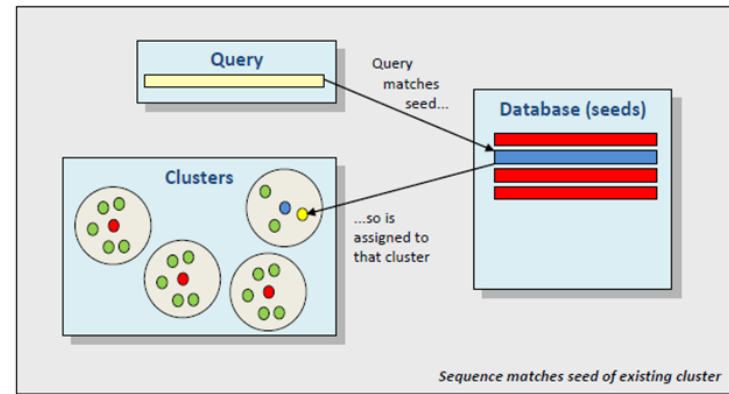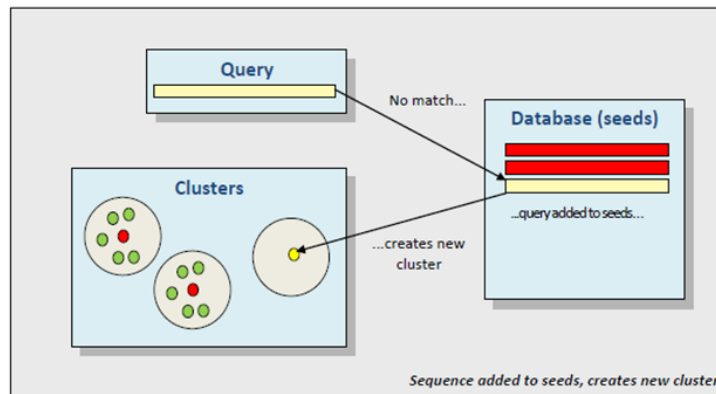extraction

Clustering to OTUs

Getting the
representative
sequences from
clusters

Identification of
OTUs
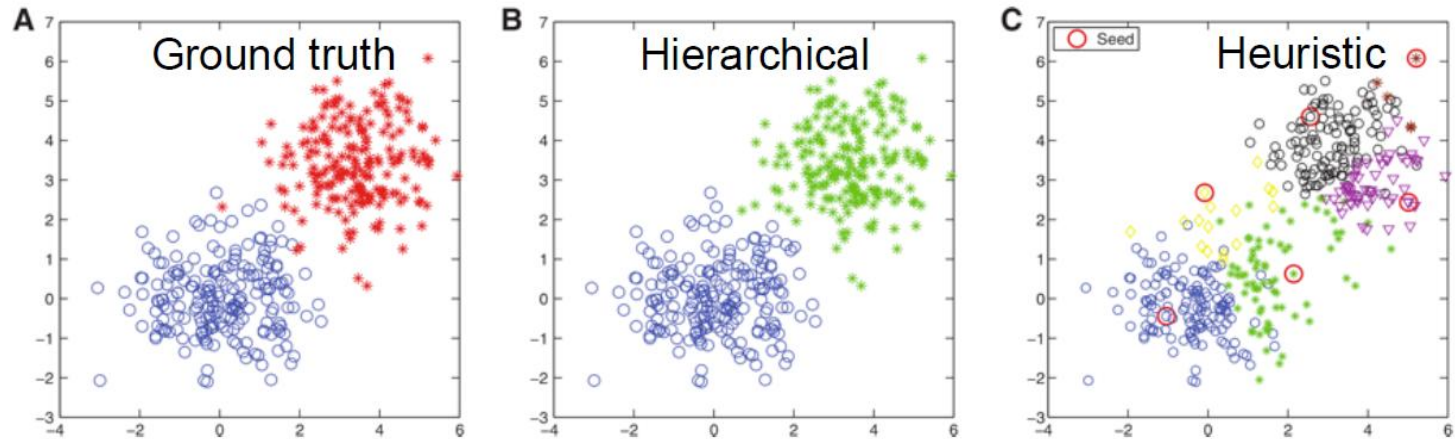
Construction of
OTU table
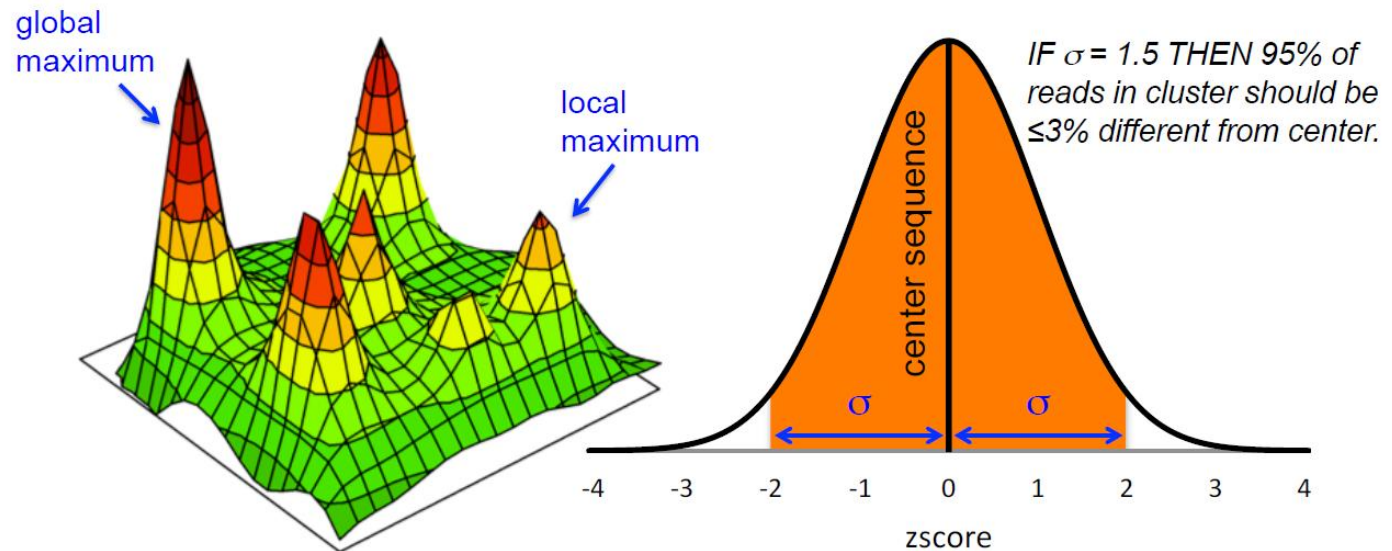
Estimation of
diversity indices

**PROCESSING OF THE
RESULTS**

Query

Clusters

Database (seeds)

No match...

...query added to seeds...

...creates new
cluster

*Sequence added to seeds, creates new cluster*

Seed

id ≥ 97%

cluster definition

Query

Clusters

Database (seeds)

No match...

...query added to seeds...

...creates new
cluster

*Sequence added to seeds, creates new cluster*

Query

Clusters

Database (seeds)

Query
matches
seed...

...so is
assigned to
that cluster

*Sequence matches seed of existing cluster*

# Clustering to OTUs (hierarchical vs. heuristic)

| RAW DATA - R1 & R2 FASTQ |
|---|

↓

| Joining of pair-end data |
|---|

↓

| Quality filtering<br>Sample determination<br>Sequence trimming |
|---|

↓

| Chimera removal |
|---|

↓

| Preparing for clustering<br>e.g.: fungal ITS<br>extraction |
|---|

↓

| **Clustering to OTUs** |
|---|

↓

| Getting the<br>representative<br>sequences from<br>clusters |
|---|

↓

| Identification of<br>OTUs |
|---|

↓

| Construction of<br>OTU table |
|---|

↓

| Estimation of<br>diversity indices |
|---|

↓

| PROCESSING OF THE<br>RESULTS |
|---|



**Hierarchical clustering**
- is able to identify the real clusters (ideally)
- computationally expensive

x

**Heuristic clustering**
- computationally cheap
- often generates artificial clusters (overestimated diversity)

Joining of pair-end data

Quality filtering
Sample determination
Sequence trimming

Chimera removal

Preparing for clustering
e.g.: fungal ITS extraction

Clustering to OTUs

Getting the representative sequences from clusters

Identification of OTUs
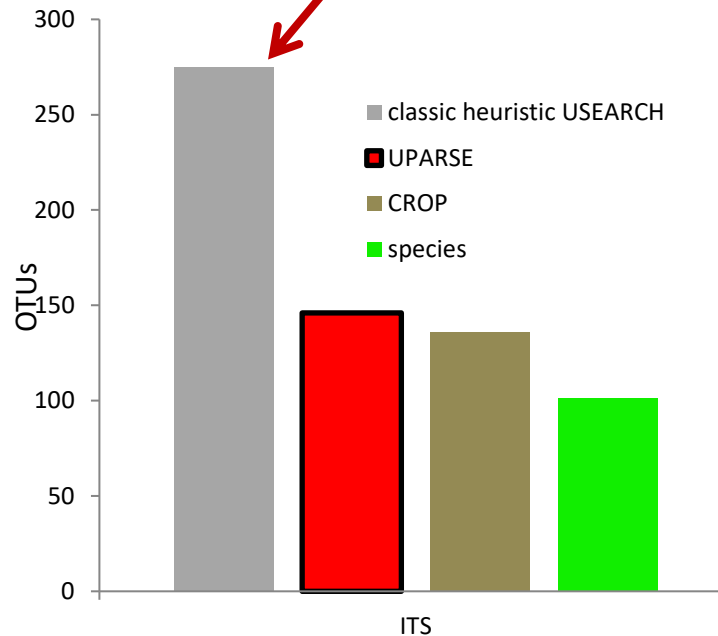
Construction of OTU table

Estimation of diversity indices

PROCESSING OF THE RESULTS

# Clustering to OTUs (model based)

## CROP (FILTER – PSA – BAYESÏAN)

Hao et al. 2011 (Bioinformatics): "If we consider the sequences as data points in a high-dimensional space [...], then the probability that a sequence belongs to a cluster becomes a function of the distance between the sequence and the center."

**CROP uses a mixture model to find subpopulations among all sequences under the assumption that they are independently drawn from a mixture of Gaussian distributions.**

global maximum

local maximum

IF $\sigma$ = 1.5 THEN 95% of reads in cluster should be ≤3% different from center.

center sequence

$\sigma$          $\sigma$

-4   -3   -2   -1   0   1   2   3   4

zscore

# Clustering to OTUs (comparison)

## RAW DATA - R1 & R2 FASTQ

- Joining of pair-end data
- Quality filtering / Sample determination / Sequence trimming
- Chimera removal
- Preparing for clustering e.g.: fungal ITS extraction
- **Clustering to OTUs**
- Getting the representative sequences from clusters
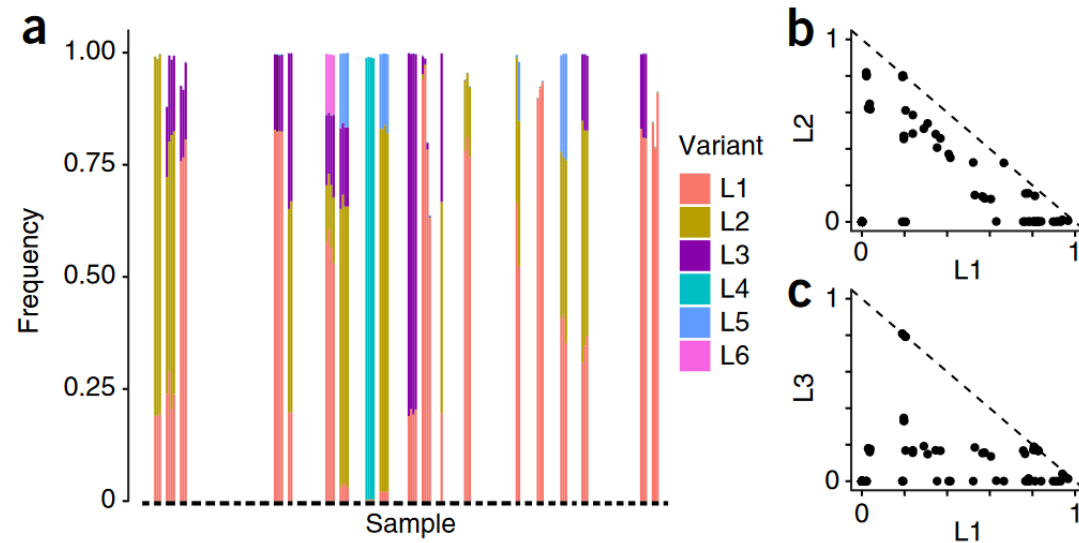- Identification of OTUs
- Construction of OTU table
- Estimation of diversity indices
- **PROCESSING OF THE RESULTS**

Problems with overestimation of obtained OTUs…



mock fungal community about 100 fungal species

solution
**UPARSE (USEARCH)** – improved heuristic algorithm which is able to recognize chimeric sequences

**http://drive5.com/uparse/**

Edgar, R.C. (2013) UPARSE: Highly accurate OTU sequences from microbial amplicon reads, *Nature Methods* [Pubmed:23955772,  dx.doi.org/10.1038/nmeth.2604].

## Clustering-independent methods

```
RAW DATA - R1 & R2
FASTQ
        ↓
Joining of pair-end data
        ↓
Quality filtering
Sample determination
Sequence trimming
        ↓
Chimera removal
        ↓
Preparing for clustering
e.g.: fungal ITS
extraction
        ↓
Clustering to OTUs
        ↓
Getting the
representative
sequences from
clusters
        ↓
Identification of
OTUs
        ↓
Construction of
OTU table
        ↓
Estimation of
diversity indices
        ↓
PROCESSING OF THE
RESULTS
```
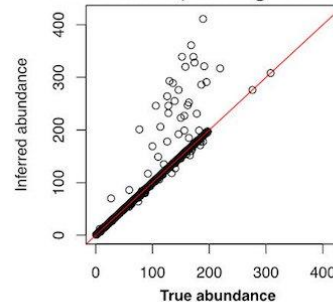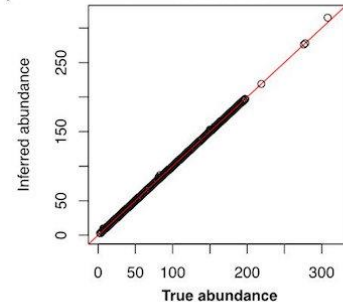
Callahan, Benjamin J., et al. "**DADA2**: high-resolution sample inference from Illumina amplicon data." *Nature methods* (2016).



L. crispatus sequence variants in the human vaginal community during pregnancy. DADA2 identified six L. crispatus 16S rRNA sequence variants present in multiple samples and a significant fraction of all reads.

Amir, Amnon, et al. "**Deblur** Rapidly Resolves Single-Nucleotide Community Sequence Patterns." *mSystems* 2.2 (2017).

**RAW DATA - R1 & R2 FASTQ**

**Joining of pair-end data**

**Quality filtering Sample determination Sequence trimming**

**Chimera removal**

**Preparing for clustering e.g.: fungal ITS extraction**

**Clustering to OTUs**

**Getting the representative sequences from clusters**

**Identification of OTUs**

**Construction of OTU table**

**Estimation of diversity indices**

**PROCESSING OF THE RESULTS**



**Accuracy: Simulated data**

3% OTUs (average linkage) | DADA2

TP: 978
FP: 272
FN: 77
cor: 0.935

TP: 1042
FP: 0
FN: 13
cor: 0.999

**Data:** Kopylova, et al. mSystems, 2016.

## Advantages

Resolution: DADA2 infers exact amplicon sequence variants (ASVs) from amplicon data, resolving biological differences of even 1 or 2 nucleotides.

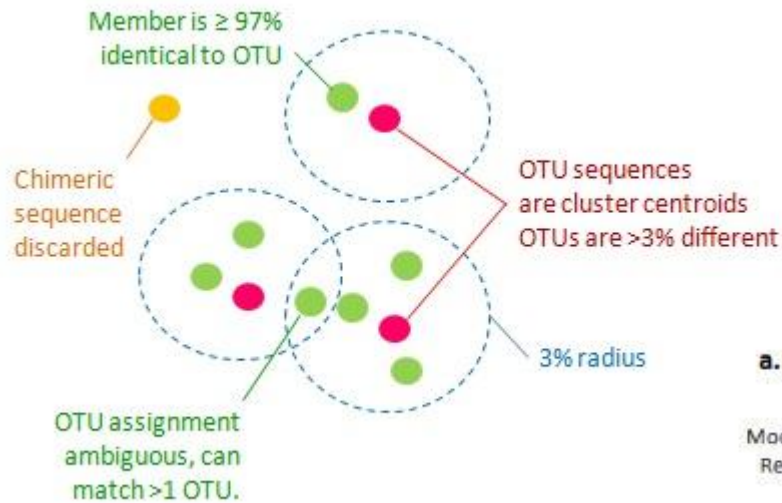Accuracy: DADA2 reports fewer false positive sequence variants than other methods report false OTUs.

Comparability: The ASVs output by DADA2 can be directly compared between studies, without the need to reprocess the pooled data.

Computational Scaling: The compute time of DADA2 scales linearly sample number, and memory requirements are essentially flat.
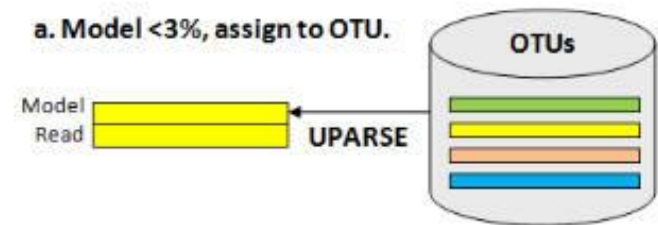
## Disadvantage

sequence variants are not representing the real sequences (they are estimated based on the errors modeling)
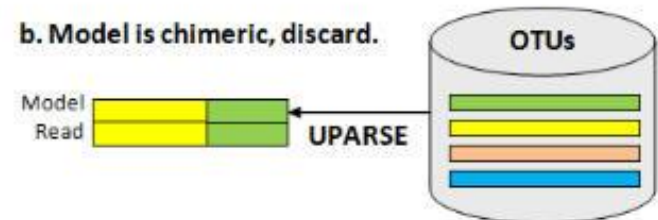
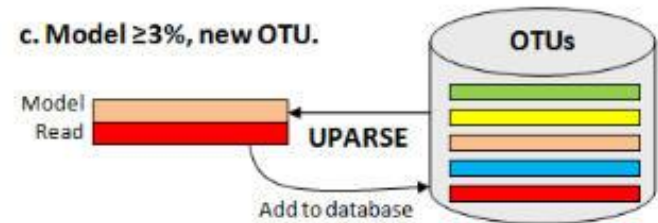# UPARSE: Clustering and chimera removal in the same time

**RAW DATA - R1 & R2 FASTQ**

**Joining of pair-end data**

**Quality filtering Sample determination Sequence trimming**

**Chimera removal**

**Preparing for clustering e.g.: fungal ITS extraction**

**Clustering to OTUs**

**Getting the representative sequences from clusters**

**Identification of OTUs**

**Construction of OTU table**

**Estimation of diversity indices**

**PROCESSING OF THE RESULTS**

Member is ≥ 97% identical to OTU

Chimeric sequence discarded

OTU sequences are cluster centroids OTUs are >3% different

3% radius

OTU assignment ambiguous, can match >1 OTU.

a. Model <3%, assign to OTU. — Model Read — UPARSE — OTUs

b. Model is chimeric, discard. — Model Read — UPARSE — OTUs

c. Model ≥3%, new OTU. — Model Read — UPARSE — OTUs — Add to database

Edgar, R.C. (2013) UPARSE: Highly accurate OTU sequences from microbial amplicon reads, *Nature Methods*

## Flowchart (left column)

- **RAW DATA - R1 & R2 FASTQ**
- Joining of pair-end data
- Quality filtering / Sample determination / Sequence trimming
- Chimera removal
- Preparing for clustering e.g.: fungal ITS extraction
- Clustering to OTUs
- Getting the representative sequences from clusters
- Identification of OTUs
- **Construction of OTU table**
- Estimation of diversity indices
- **PROCESSING OF THE RESULTS**

# Construction of OTU table

**OTU table** - matrix that gives the number of reads per sample per OTU

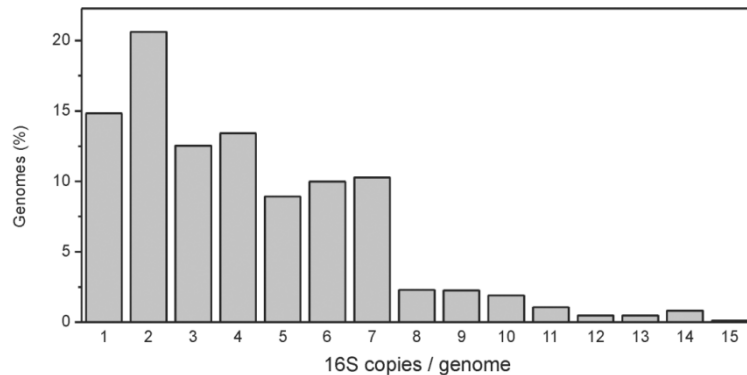| OTU_ID | SAMPLE_1 | SAMPLE_2 | SAMPLE_3 | SAMPLE_4 | SAMPLE_5 | SAMPLE_6 | SAMPLE_7 | SAMPLE_8 |
|---|---|---|---|---|---|---|---|---|
| CL00001 | 249 | 189 | 220 | 311 | 1 | 16 | 68 | 2 |
| CL00002 | 201 | 19 | 169 | 438 | 1 | 8 | 12 | 0 |
| CL00003 | 190 | 39 | 176 | 210 | 0 | 21 | 20 | 1 |
| CL00004 | 183 | 36 | 195 | 177 | 1 | 16 | 16 | 0 |
| CL00005 | 0 | 26 | 2 | 35 | 20 | 164 | 4 | 116 |
| CL00006 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 |
| CL00007 | 133 | 71 | 125 | 89 | 0 | 3 | 26 | 0 |
| CL00008 | 106 | 42 | 96 | 158 | 0 | 10 | 14 | 0 |
| CL00009 | 95 | 46 | 108 | 134 | 2 | 7 | 24 | 0 |
| CL00010 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 3 |
| CL00011 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 |

**OTU frequency does not correlate with species frequency**

This means, for example, that the most abundant OTU does not have to be the most abundant species – especially because of multi-copy nature of target genes as 16S and ITS
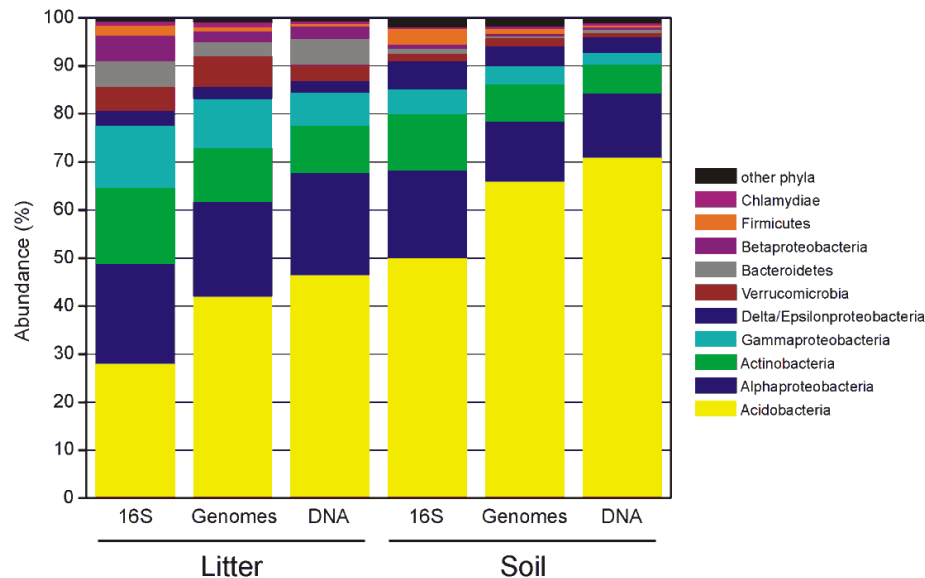
**Singleton counts are especially suspect**
- many OTU table entries are often singletons (have value 1) for smaller OTUs because the total count is distributed over several samples
- Small counts are more likely to be spurious, especially singletons, either because the OTU itself is spurious (e.g., an undetected chimera), or because of cross-talk
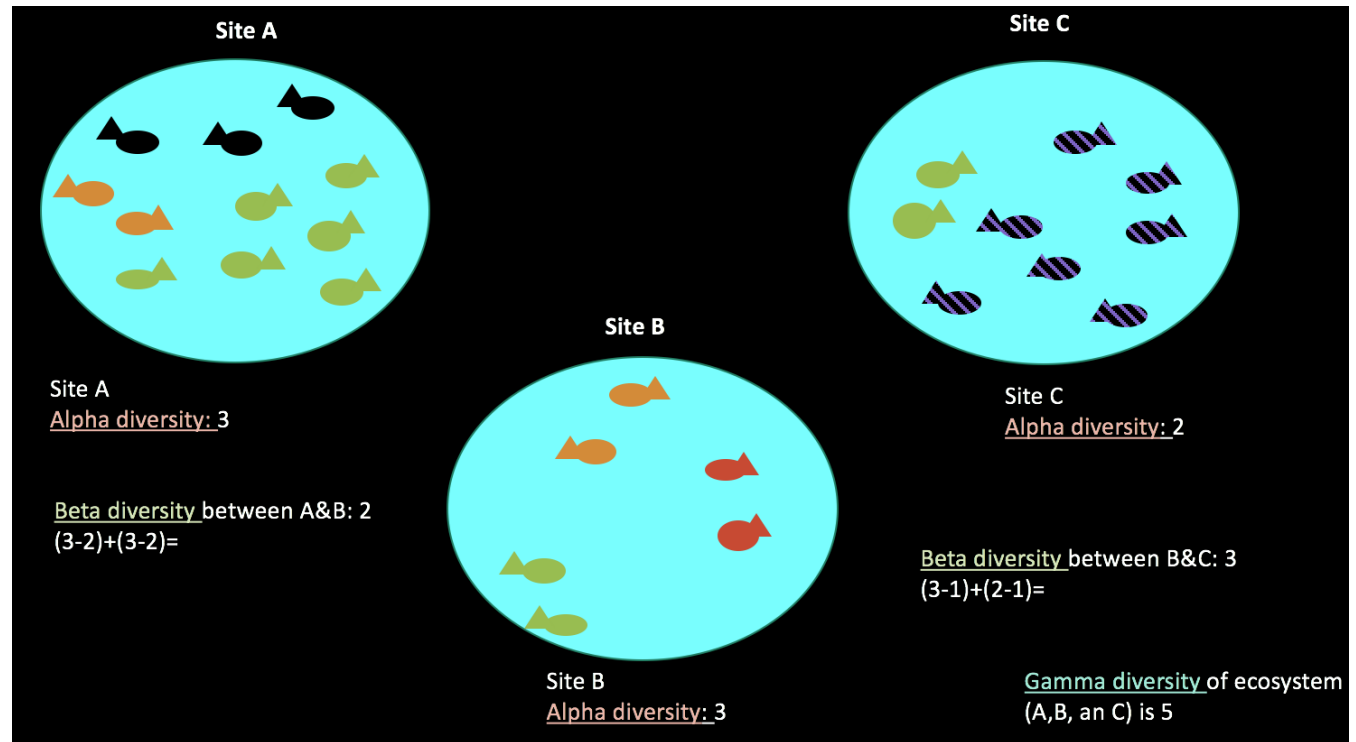
# Normalize OTU table by 16S copy number



**rrnDB**
A searchable database documenting variation in ribosomal RNA operons (rrn) in Bacteria and Archaea. Find information such as the 16S gene copy number of an organism by looking up its name under the NCBI or RDP taxonomy or by full-text search of rrnDB's records.



**Abundance of bacterial 16S rRNA sequences, genomes and DNA in forest litter and soil.**

Relative abundance of bacterial 16S rRNA sequences in the amplicon pool from Picea abies litter and soil (Baldrian et al., 2012), and estimates of the relative abundance of bacterial genomes and DNA. The estimates were calculated using the values of 16S rRNA copy numbers and genome sizes of the closest hits to each bacterial OTU.

T. Větrovský & P. Baldrian - PloS one, 2013 & Stoddard et al. (2015) https://rrndb.umms.med.umich.edu/

**RAW DATA - R1 & R2 FASTQ**

↓

**Joining of pair-end data**

↓

**Quality filtering Sample determination Sequence trimming**

↓

**Chimera removal**

↓

**Preparing for clustering e.g.: fungal ITS extraction**

↓

**Clustering to OTUs**

↓

**Getting of the representative sequences from the clusters**

↓

**Identification of OTUs**

↓

**Construction of OTU table**

↓

**Estimation of diversity indices**

↓

**PROCESSING OF THE RESULTS**

# Estimation of diversity indices



Site A
Alpha diversity: 3

Beta diversity between A&B: 2
(3-2)+(3-2)=

Site B
Alpha diversity: 3

Site C
Alpha diversity: 2

Beta diversity between B&C: 3
(3-1)+(2-1)=

Gamma diversity of ecosystem
(A,B, an C) is 5

- **Alpha-diversity**: diversity of organisms in one sample / environment
  - **Shannon index**
  - **Chao1**
  - **Observed OTUs (Richness)**

- **Beta-diversity**: differences in diversities across samples or environments
  - **UniFrac** (Lozupone et al, AEM, 2005) (phylogenetic)
  - **Bray-Curtis** dissimilarity measure (OTU abundance)
  - **Jaccard** similarity coefficient (OTU presence/absence)

# Alpha diversity

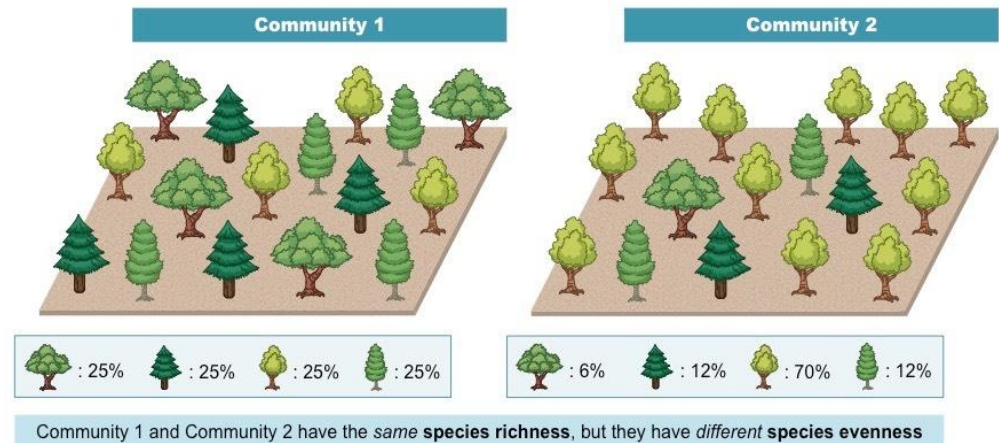$$H' = -\sum_{i=1}^{S} p_i \ln p_i$$

## Shannon index (Shannon entropy)
Then the Shannon entropy quantifies the uncertainty in predicting the species identity of an individual that is taken at random from the dataset.

$p_i$ – proportion of the population made up of species i
$S$ – number of species in sample

## Species evenness
Species evenness refers to how close in numbers each species in an environment is. Mathematically it is defined as a diversity index, a measure of biodiversity which quantifies how equal the community is numerically.



**Community 1** — : 25% : 25% : 25% : 25%
**Community 2** — : 6% : 12% : 70% : 12%

Community 1 and Community 2 have the *same* **species richness**, but they have *different* **species evenness**

$$J' = \frac{H'}{H'_{\max}} \qquad H'_{\max} = -\sum_{i=1}^{S} \frac{1}{S} \ln \frac{1}{S} = \ln S$$

## Chao1 index
Estimate diversity from abundance data (importance of rare OTUs)

$$S_{est} = S_{obs} + \left( \frac{f_1^{\,2}}{2 f_2} \right)$$

where $S_{obs}$ is the number of species in the sample, $f_1$ is the number of singletons and $f_2$ is the number of doubletons.

Shannon, C. E. (1948) A mathematical theory of communication.
The Bell System Technical Journal, 27, 379–423 and 623–656.
*Chao*, A.; Shen, T-J. (2003)

https://palaeo-electronica.org/2011_1/238/estimate.htm

# Getting of the representative sequences from the clusters

RAW DATA - R1 & R2 FASTQ

Joining of pair-end data

Quality filtering
Sample determination
Sequence trimming

Chimera removal

Preparing for clustering
e.g.: fungal ITS
extraction

Clustering to OTUs

Getting the
representative
sequences from
clusters

Identification of
OTUs

Construction of
OTU table

Estimation of
diversity indices

PROCESSING OF THE
RESULTS

**centroid**



T = identity threshold
centroid sequence
member sequence

**consensus**



**most abundant**

## Workflow (left column)

- **RAW DATA - R1 & R2 FASTQ**
- Joining of pair-end data
- Quality filtering / Sample determination / Sequence trimming
- Chimera removal
- Preparing for clustering e.g.: fungal ITS extraction
- Clustering to OTUs
- Getting the representative sequences from clusters
- **Identification of OTUs**
- Construction of OTU table
- Estimation of diversity indices
- **PROCESSING OF THE RESULTS**

# Taxonomic classification of OTUs

## Similarity-based

Find homology or minimum alignment distance

Tools:
- local alignments (e.g. BLAST, MEGAN, METAXA2, RTAX)
- global alignments (e.g. GAST)
- overlap alignments (e.g. SINA)

Pro/Con:
- good accuracy for similar sequences
- performs less well on distant lineages
- can be slow on large reference databases

## Composition-based

Detect specific features

Tools:
- kmer searches (e.g. NBC/RDP, UTAX, SINTAX)
- hidden Markov models (e.g. PHYMMBL, C16S)

Pro/Con:
- computationally efficient and fast
- performs well on distant lineages
- training required
- limited resolution for shorter sequences

## Phylogeny-based

Evolutionary model to determine best placement

Tool:
- ML, NJ, Bayesian (e.g. PPLACER, EPA)

Pro/Con:
- great accuracy for similar sequences
- classification in its evolutionary context
- computationally complex
- requires accurate reference tree
- difficult for non-coding regions

# Identification of OTUs



**RAW DATA - R1 & R2 FASTQ**

Joining of pair-end data

Quality filtering
Sample determination
Sequence trimming

Chimera removal

Preparing for clustering
e.g.: fungal ITS
extraction

Clustering to OTUs

Getting the
representative
sequences from
clusters

**Identification of OTUs**

Construction of OTU table

Estimation of diversity indices

**PROCESSING OF THE RESULTS**

All genes
**GenBank - *genetic sequence database, an annotated collection of all publicly available DNA sequence***
• largest ☺
• many errors ☹

**https://www.ncbi.nlm.nih.gov/genbank/**

Gen – GenBank
WGS – whole genome sequences

# Identification of OTUs

**RAW DATA - R1 & R2 FASTQ**

Joining of pair-end data

Quality filtering
Sample determination
Sequence trimming

Chimera removal

Preparing for clustering
e.g.: fungal ITS
extraction

Clustering to OTUs

Getting the
representative
sequences from
clusters

**Identification of OTUs**

Construction of
OTU table

Estimation of
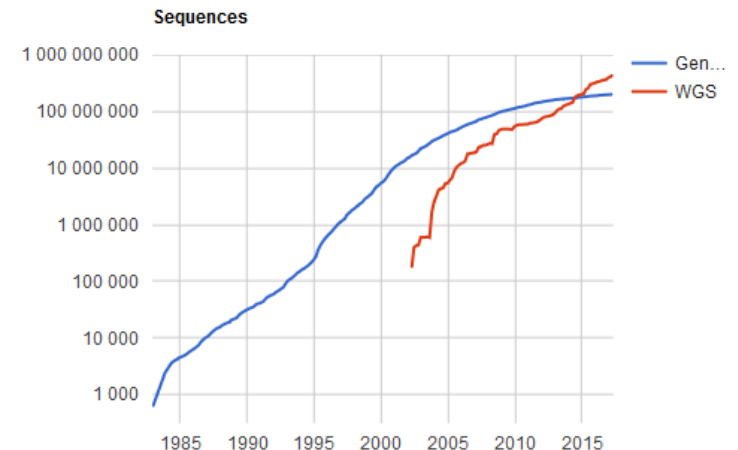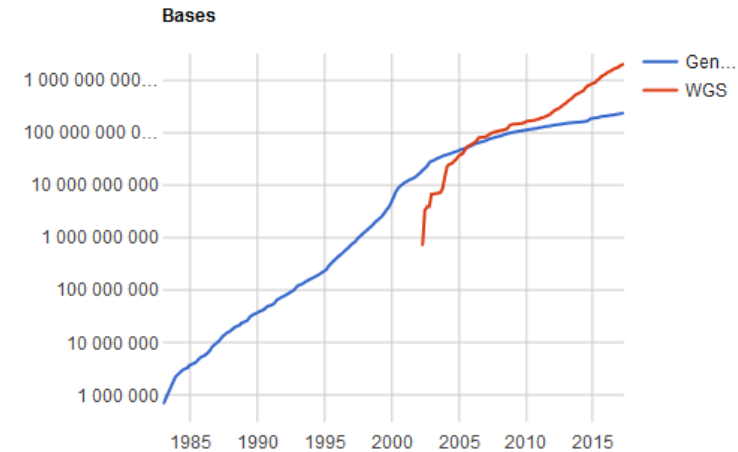diversity indices

**PROCESSING OF THE RESULTS**

Identification of bacteria
**RDP – *Ribosomal Database Project***
provides quality-controlled, aligned and annotated Bacterial and Archaeal 16S rRNA sequences, and Fungal 28S rRNA sequences, and a suite of analysis tools to the scientific community



**Hb** Hierarchy Browser  **Cl** Classifier  **Pm** Probe Match  **Fg** FunGene  **Al** Aligner  **Tb** Tree Builder  **Os** RDP Open Source  **Tu** Tutorials

**Mg** MlxS GoogleSheets  **Lc** Library Compare  **Sm** Sequence Match  **Rp** RDPipeline

**https://rdp.cme.msu.edu/**

**SILVA** - provides comprehensive, quality checked and regularly updated datasets of aligned small (16S/18S, SSU) and large subunit (23S/28S, LSU) ribosomal RNA (rRNA) sequences for all three domains of life (Bacteria, Archaea and Eukarya).

**silva**
high quality ribosomal RNA databases

**https://www.arb-silva.de/**

**GREENGENES**
The 16S rRNA Gene Database and Tools

**http://greengenes.secondgenome.com/**

# Identification of OTUs

**RAW DATA - R1 & R2 FASTQ**

↓

**Joining of pair-end data**

↓

**Quality filtering Sample determination Sequence trimming**

↓

**Chimera removal**

↓

**Preparing for clustering e.g.: fungal ITS extraction**

↓

**Clustering to OTUs**

↓

**Getting the representative sequences from clusters**

↓

**Identification of OTUs**

↓

**Construction of OTU table**

**Estimation of diversity indices**

↓

**PROCESSING OF THE RESULTS**

Identification of fungi
**UNITE - *Unified system for the DNA based fungal species linked to the classification***



Run Analysis   Search Pages   Resources   Statistics   Notes and News   Workbench

**unite** community

*Unified system for the DNA based fungal species linked to the classification Ver. 7.1*

Current version: **7.2**; Last updated: 2017-06-08 (**read more**)
Number of ITS sequences (UNITE+INSD): **741 222**; Number of UNITE fungal Species Hypotheses with DOIs at 1.5% threshold: **73 929** (**more statistics**)

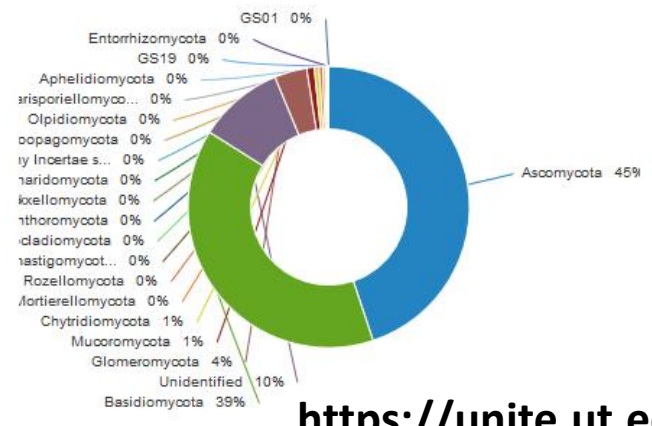| Threshold | 1.5 % ▾ | Include | All SH-s ▾ | Start typing taxon name here ... | Go | Reset | ❶ |

▸ Ascomycota (33,051)
▸ Basidiomycota (28,743)
▸ Unidentified (7,240)
▸ Glomeromycota (2,743)
▸ Mucoromycota (636)
▸ Chytridiomycota (388)
▸ Mortierellomycota (322)
▸ Rozellomycota (198)
▸ Neocallimastigomycota (90)
▸ Blastocladiomycota (42)
▸ Entomophthoromycota (41)
▸ Kickxellomycota (39)
▸ Monoblepharidomycota (22)
▸ Fungi phy Incertae sedis (19)
▸ Zoopagomycota (15)
▸ Olpidiomycota (13)

SH graph: Fungi

GS01 0%
Entorrhizomycota 0%
GS19 0%
Aphelidiomycota 0%
...arisporiellomycota... 0%
Olpidiomycota 0%
...oopagomycota 0%
...y Incertae s... 0%
...haridomycota 0%
...kxellomycota 0%
...nthoromycota 0%
...cladiomycota 0%
...astigomycot... 0%
Rozellomycota 0%
Mortierellomycota 0%
Chytridiomycota 1%
Mucoromycota 1%
Glomeromycota 4%
Unidentified 10%
Basidiomycota 39%
Ascomycota 45%

**https://unite.ut.ee/**

RAW DATA - R1 & R2 FASTQ

Joining of pair-end data

Quality filtering
Sample determination
Sequence trimming

Chimera removal

Preparing for clustering
e.g.: fungal ITS extraction

Clustering to OTUs

Getting the representative sequences from clusters

Identification of OTUs

Construction of OTU table

Estimation of diversity indices

PROCESSING OF THE RESULTS

# PROCESSING OF THE RESULTS

## Introduction to multivariate data analysis (Iñaki Class)