# B- Training on Galaxy: Metabarcoding

June 2021 - Webinar

# FROGS Practice on 16S data

LUCAS AUER, MARIA BERNARD, LAURENT CAUQUIL, MAHENDRA MARIADASSOU, GÉRALDINE PASCAL & OLIVIER RUÉ

.

# Overview

# Objectives

Analyses of mibrobial communities

High-throughput sequencing of 16S/18S RNA amplicons

Illumina data, sequenced at great depth

Bioinformatics data processing

Abundance table

**WITH**
operational taxonomic units (OTUs) and
their taxonomic affiliation.

# Objectives: a count table

|  | Affiliation | Sample 1 | Sample 2 | Sample 3 | Sample 4 | Sample 5 | Sample 6 |
|---|---|---|---|---|---|---|---|
| OTU1 | Species A | 0 | 100 | 0 | 45 | 75 | 18645 |
| OTU2 | Species B | 741 | 0 | 456 | 4421 | 1255 | 23 |
| OTU3 | Species C | 12786 | 45 | 3 | 0 | 0 | 0 |
| OTU4 | Species D | 127 | 4534 | 80 | 456 | 756 | 108 |
| OTU5 | Species E | 8766 | 7578 | 56 | 0 | 0 | 200 |

# Why FROGS was developed ?

Most solutions are often designed for specialists making access difficult for the whole community (command lines).

We developed the pipeline **FROGS**: « ***Find Rapidly OTU with Galaxy Solution*** » usable with command lines or within interface.

# Who is in the current FROGS group?

**Maria BERNARD**  **Olivier RUÉ**

**Lucas AUER**  **Laurent CAUQUIL**

**Patrice Déhais**

Developers

Biology experts

Galaxy support

**Mahendra MARIADASSOU**

**Géraldine PASCAL**

Statistical expert

Coordinator
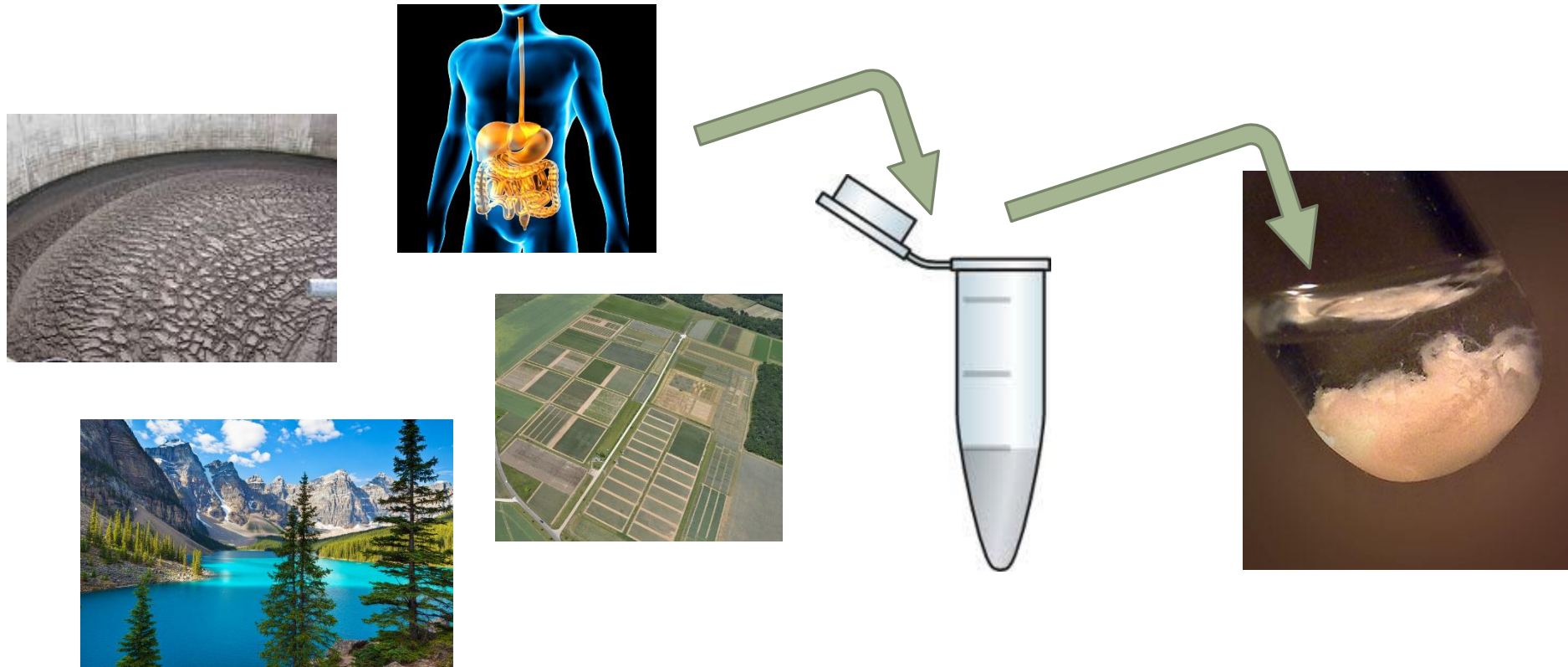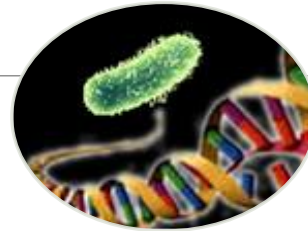
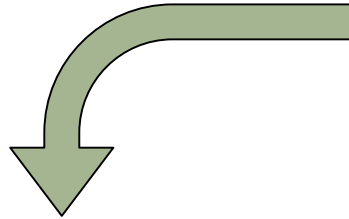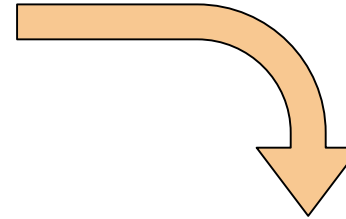# Material

# Sample collection and DNA extraction

# « Meta-omics » using next-generation sequencing (NGS)



DNA

RNA

| Metagenomics | Metatranscriptomics |
| --- | --- |

| Amplicon sequencing | Shotgun sequencing | RNA sequencing |
| --- | --- | --- |

Wolfe *et al.*, 2014

Almeida *et al.*, 2014

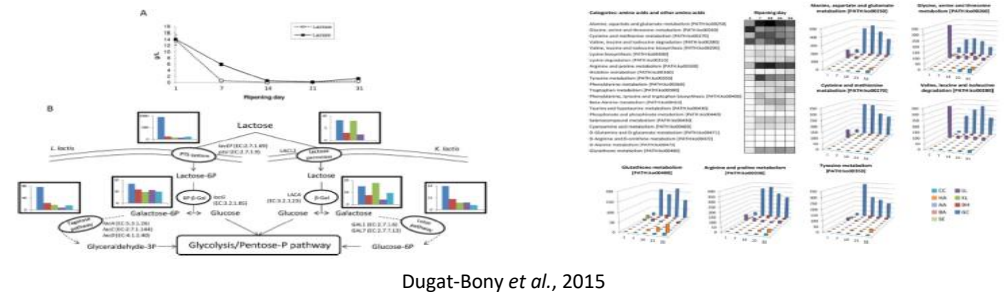Dugat-Bony *et al.*, 2015

Who is here?

What can they do?

What are they doing?

# The gene encoding the small subunit of the ribosomal RNA

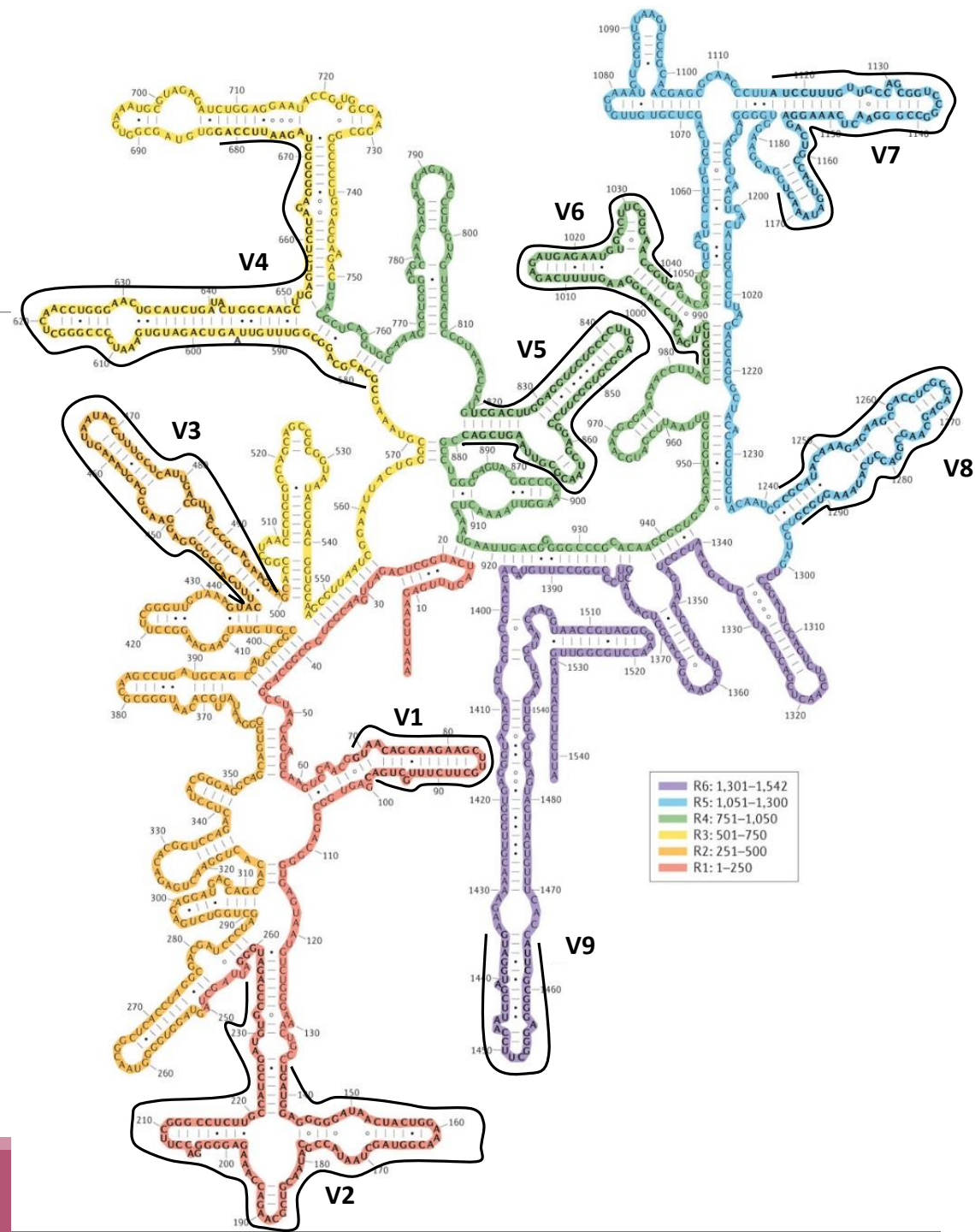The most widely used gene in **molecular phylogenetic** studies

Ubiquist gene : **16S rDNA** in prokaryotes ; **18S rDNA** in eukaryotes

**Gene encoding a ribosomal RNA :** non-coding RNA (not translated), part of the small subunit of the ribosome which is responsible for the translation of mRNA in proteins

Not submitted to lateral gene transfer

Availability of databases facilitating comparison
(Silva v138.1 - 2021: available SSU/LSU sequences to over **10,700,000**)

Secondary structure of the 16S rRNA of *Escherichia coli*

In red, fragment R1 including regions V1 and V2;
in orange, fragment R2 including region V3;
in yellow, fragment R3 including region V4;
in green, fragment R4 including regions V5 and V6;
in blue, fragment R5 including regions V7 and V8;
and in purple, fragment R6 including region V9.

*Uniting the classification of cultured and uncultured bacteria and archaea using 16S rRNA gene sequences*
*Pablo Yarza, et al.*
*Nature Reviews Microbiology 12, 635–645 (2014) doi:10.1038/nrmicro3330*

11

# The gene encoding the small subunit of the ribosomal RNA



0  100  200  300  400  500  600  700  800  900  1000  1100  1200  1300  1400  1500 bp

V1  V2  V3  V4  V5  V6  V7  V8  V9

**CONSERVED REGIONS:** unspecific applications

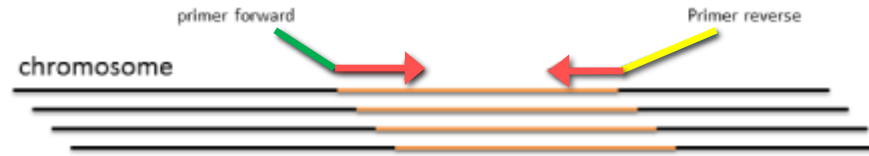**VARIABLE REGIONS:** group or species-specific applications

# Other targets

Bacterial lineages vary in their genomic contents, which suggests that different genes might be needed to resolve the diversity within certain taxonomic groups.

The genes that have been proposed for this task include those encoding :

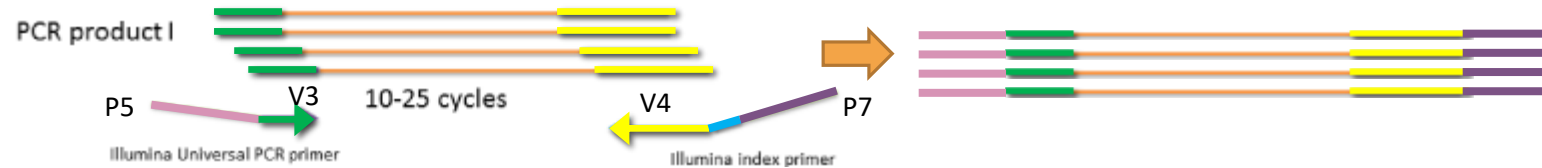- 23S rRNA,
- DNA gyrase subunit B (gyrB),
- RNA polymerase subunit B (rpoB),
- TU elongation factor (tuf),
- DNA recombinase protein (recA),
- protein synthesis elongation factor-G (fusA),
- dinitrogenase protein subunit D (nifD),
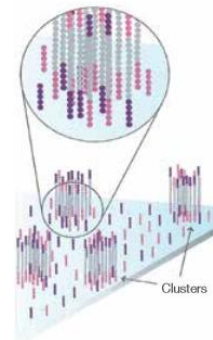- Internal Transcribed Spacer (ITS) for Fungi.

# Steps for Illumina sequencing

- 1st step : one PCR

- 2nd step: one PCR

- 3rd step: on flow cell, the cluster generations

- 4th step: sequencing

# Amplification and sequencing

« Universal » primer sets are used for **PCR amplification** of the phylogenetic biomarker

The primers contain adapters used for the sequencing step and barcodes (= tags = MIDs) to distinguish the samples (multiplexing = sequencing several samples on the same run)

| Adapter A | BarcodeSequence | LinkerPrimerSequence | Target Sequence | ReversePrimer | Adapter B |

example: V3                    Desired Sequence                    example: V4

# Illumina sequencing



Illumina video:
https://www.youtube.com/watch?v=fCd6B5HRaZ8

# Cluster generation
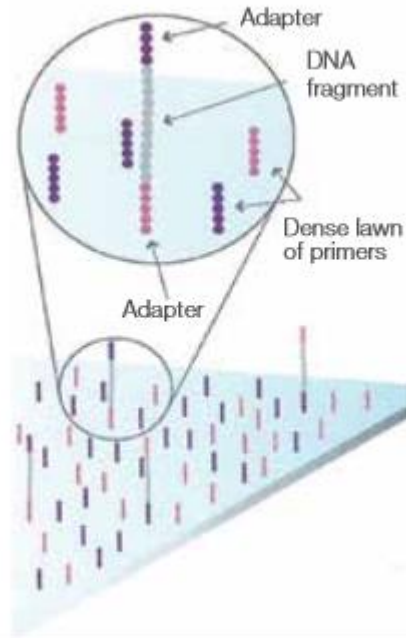
## Prepare Genomic DNA Sample



Randomly fragment genomic DNA and ligate adapters to both ends of the fragments.
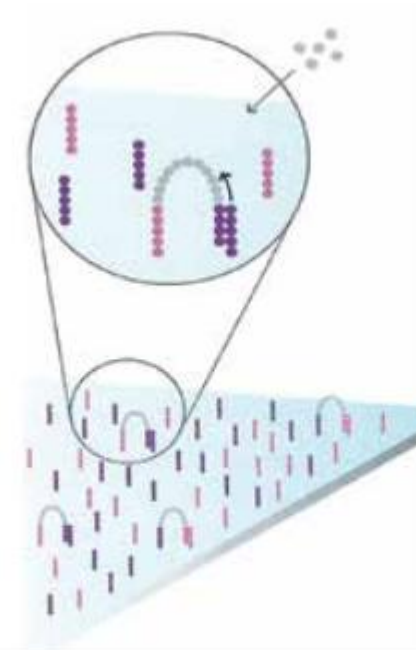
## Attach DNA to Surface



Bind single-stranded fragments randomly to the inside surface of the flow cell channels.

Attach DNA to surface

## Bridge Amplification



Add unlabeled nucleotides and enzyme to initiate solid-phase bridge amplification.

Bridge amplification

# Cluster generation

Fragments Become Double Stranded    Denature the Double-Stranded Molecules        Complete Amplification



The enzyme incorporates nucleotides to build double-stranded bridges on the solid-phase substrate.

Fragments become double stranded

Denaturation leaves single-stranded templates anchored to the substrate.

Denature the double-stranded molecule

Several million dense clusters of double-stranded DNA are generated in each channel of the flow cell.

Cycle of new strand synthesis and denaturation to make multiple copies of the same sequence (amplification)
Reverse strands are washed

# Sequencing by synthesis

### Determine First Base



The first sequencing cycle begins by adding four labeled reversible terminators, primers, and DNA polymerase.

Light signal is more strong in cluster

### Image First Base



After laser excitation, the emitted fluorescence from each cluster is captured and the first base is identified.

### Determine Second Base



The next cycle repeats the incorporation of four labeled reversible terminators, primers, and DNA polymerase.

# Sequencing by synthesis

### Image Second Chemistry Cycle



After laser excitation, the image is captured as before, and the identity of the second base is recorded.

### Sequencing Over Multiple Chemistry Cycles



GCTGA...

The sequencing cycles are repeated to determine the sequence of bases in a fragment, one base at a time.

Barcode is read, so cluster is identified.
After first sequencing (250 or 300 nt of Reverse strand), fragment form bridges again and Forward strand can be sequenced also.

DNA

Constant region

Divergent region

PCRs

Illumina index

Illumina adapter

Illumina adapter

Sequencing

Read 1

Read 2

Index 1

# Identification of bacterial populations may be not discriminating



0                        RNAr 16S gene total                        1500

F          V3                    V4          R

Amplicon

Constant regions

Divergent regions

# Amplification and sequencing

Sequencing is generally perform on Roche-454 (obsolete now) or Illumina MiSeq platforms or Oxford Nanopore Technology platform.

Read quantity: ~10 000 reads per sample (454), ~30 000 reads per sample (MiSeq), up to several Tera of data (ONT).

Sequence lengths: >650 bp (Roche-454), 2 x 250 bp or 2 x 300 bp (MiSeq), Longest read > 2Mb (ONT)

# Methods

# FROGS Pipeline



**Pre-process**

FROGS Pre-process ✖
- TAR archive file
- dereplicated_file (fasta)
- count_file (tabular)
- summary_file (html)

**Clustering**

FROGS Clustering swarm ✖
- Sequences file
- Count file
- seed_file (fasta)
- abundance_biom (biom1)
- swarms_composition (tabular)

**Chimera**

FROGS Remove chimera ✖
- Sequences file
- Abundance file
- non_chimera_fasta (fasta)
- out_abundance_biom (biom1)
- out_abundance_count (tabular)
- summary_file (html)

**OTU Filters**

FROGS OTU Filters ✖
- Sequences file
- Abundance file
- output_fasta (fasta)
- output_biom (biom1)
- output_excluded (tabular)
- output_summary (html)

**Affiliation**

FROGS Affiliation OTU ✖
- OTU seed sequence
- Abundance file
- biom_affiliation (biom1)
- summary (html)

**Affiliation Filters**

FROGS Affiliation Filters ✖
- Sequences file
- Abundance file
- output_biom (biom1)
- output_fasta (fasta)
- output_impacted (tabular)
- output_multihit (tabular)
- output_summary (html)

25 tools in total

26

# FROGS Pipeline

**FROGS Pre-process ✗**
- TAR archive file
- dereplicated_file (fasta)
- count_file (tabular)
- summary_file (html)

**Pre-process**

**FROGS Clustering swarm ✗**
- Sequences file
- Count file
- seed_file (fasta)
- abundance_biom (biom1)
- swarms_composition (tabular)

**Clustering**

**FROGS Remove chimera ✗**
- Sequences file
- Abundance file
- non_chimera_fasta (fasta)
- out_abundance_biom (biom1)
- out_abundance_count (tabular)
- summary_file (html)

**Chimera**

**FROGS OTU Filters ✗**
- Sequences file
- Abundance file
- output_fasta (fasta)
- output_biom (biom1)
- output_excluded (tabular)
- output_summary (html)

**OTU Filters**

**FROGS Affiliation OTU ✗**
- OTU seed sequence
- Abundance file
- biom_affiliation (biom1)
- summary (html)

**Affiliation**

**FROGS Affiliation Filters ✗**
- Sequences file
- Abundance file
- output_biom (biom1)
- output_fasta (fasta)
- output_impacted (tabular)
- output_multihit (tabular)
- output_summary (html)

**Affiliation Filters**

**FROGS BIOM to TSV ✗**
- Abundance file
- Sequences file (optional)
- tsv_file (tabular)
- multi_affi_file (tabular)

**Convert to TSV**

# FROGS Pipeline

**Pre-process**

FROGS Pre-process
- TAR archive file
- dereplicated_file (fasta)
- count_file (tabular)
- summary_file (html)

**Clustering**

FROGS Clustering swarm
- Sequences file
- Count file
- seed_file (fasta)
- abundance_biom (biom1)
- swarms_composition (tabular)

**Chimera**

FROGS Remove chimera
- Sequences file
- Abundance file
- non_chimera_fasta (fasta)
- out_abundance_biom (biom1)
- out_abundance_count (tabular)
- summary_file (html)

**OTU Filters**

FROGS OTU Filters
- Sequences file
- Abundance file
- output_fasta (fasta)
- output_biom (biom1)
- output_excluded (tabular)
- output_summary (html)

**ITSX filter**

FROGS ITSx
- Sequences file
- Abundance file
- out_fasta (fasta)
- out_abundance_biom (biom1)
- out_excluded (fasta)
- summary_file (html)

**Affiliation**

FROGS Affiliation OTU
- OTU seed sequence
- Abundance file
- biom_affiliation (biom1)
- summary (html)

**Affiliation Filters**

FROGS Affiliation Filters
- Sequences file
- Abundance file
- output_biom (biom1)
- output_fasta (fasta)
- output_impacted (tabular)
- output_multihit (tabular)
- output_summary (html)

**Convert to TSV**

FROGS BIOM to TSV
- Abundance file
- Sequences file (optional)
- tsv_file (tabular)
- multi_affi_file (tabular)

# FROGS Tools for Bioinfomatics analyses

# FROGS Tools for Statistic analyses

# What kind of data ?
## 2 Histories

16S fastq sequences in an archive tar.gz

Food environment

chaillou_withprimers_64renamedsamples_V1V3_10000seq_R1R2.tar.gz

ITS data

METABARFOOD project

ITS.tar.gz

# Demultiplexing tool

skip

# Barcoding ?

# Demultiplexing

Sequence demultiplexing in function of barcode sequences :

- In forward
- In reverse
- In forward and reverse

Remove unbarcoded or ambiguous sequences

# Demultiplexing forward

Adapter A          Primer Fwd                                    Primer Rv

Barcode Fwd              Amplicon sequence targeted              Adapter B

Single-end sequencing

Paired-end sequencing

R1                          R1                                              R2

# Demultiplexing reverse

Adapter A          Primer Fwd                                                                                      Primer Rv          Adapter B

Amplicon sequence targeted                                                    Barcode Rv

Single end
sequencing

Paire end sequencing

R1                                                                                      R1                          R2

# Demultiplexing forward and reverse

The tool parameters depend on the input data type

FROGS Demultiplex reads (version 1.1.0)

**Barcode file:**

1: barcode.tabular ▾

This file describes barcodes and samples (one line by sample tabulated separated from barcode sequence(s)). See Help section

**Single or Paired-end reads:**

Single ▾

Select between paired and single end data

You have only R1 seq.

**Select fastq dataset:**

▾

Specify dataset of your single end reads

**barcode mismatches:**

0

Number of mismatches allowed in barcode

**barcode on which end ?:**

Forward ▾
Forward
Reverse
Both ends

Execute

Where is the barcode seq on the reads?

FROGS Demultiplex reads (version 1.1.0)

**Barcode file:**

1: barcode.tabular ▾

This file describes barcodes and samples (one line by sample tabulated separated from barcode sequence(s)). See Help section

**Single or Paired-end reads:**

Paired ▾

Select between paired and single end data

You have R1 and R2 seq.

**Select first set of reads:**

▾

Specify dataset of your forward reads

**Select second set of reads:**

▾

Specify dataset of your reverse reads

**barcode mismatches:**

0

Number of mismatches allowed in barcode

**barcode on which end ?:**

Forward ▾
Forward
Reverse
Both ends

at the begining of the forward end or of the reverse end or both?

Execute

FROGS Demultiplex reads ✖

Barcode file

Select fastq dataset

demultiplexed_archive (data)

undemultiplexed_archive (data)

summary (tabular)

**Demultiplexing**

FROGS Demultiplex reads Attribute reads to samples in function of inner barcode. (Galaxy Version 2.0.0)  ▾ Options

**Barcode file**

📄 🗗 📁   24: barcode_forward.tabular

This file describes barcodes and samples (one line by sample tabulated separated from barcode sequence(s)). See Help section

**Single or Paired-end reads**

Single  ▾

Select between paired and single-end data

**Select fastq dataset**

📄 🗗 📁   6: multiplex.fastq  ▾

Specify dataset of your single end reads

**Barcode mismatches**

0

Number of mismatches allowed in barcode

**Barcode on which end ?**

Forward  ▾

The barcode is placed either at the beginning of the forward end or of the reverse end or both?

✔ Execute

---

**Input example**

| MgArd0001 | ACAGCGT |
|-----------|---------|
| MgArd0009 | ACAGTAG |
| MgArd0017 | ACGTCAG |
| MgArd0029 | ACTCAGT |
| MgArd0038 | ACTCGTC |
| MgArd0046 | AGCAGTC |
| MgArd0054 | AGCTATG |
| MgArd0062 | AGCTCGC |
| MgArd0073 | AGTATCT |
| MgArd0081 | AGTCTGC |

if index is in only at forward: tabular file with 2 columns sample names + barcodes

# Advices

<div style="border:1px solid; background-color:#f7f0d0; padding:10px; text-align:center;">For your own data</div>

- Do not forget to indicate barcode sequence as they are in the fastq sequence file, especially if you have data multiplexed via the reverse strand.

- For the mismatch threshold, we advised you to let the threshold to 0, and if you are not satisfied by the result, try with 1. The number of mismatch depends on the length of the barcode, but often those sequences are very short so 1 mismatch is already more than the sequencing error rate.

- If you have different barcode lengths, you must demultiplex your data in different times beginning by the longest barcode set and used the "unmatched" or "ambiguous" sequence with smaller barcode and so on.

- If you have Roche 454 sequences in sff format, you must convert them with some program like [sff2fastq](#)

# Outputs

**9: FROGS Demultiplex reads: report**

**8: FROGS Demultiplex reads: undemultiplexed.tar.gz**

**7: FROGS Demultiplex reads: demultiplexed.tar.gz**

| 1 | 2 |
|---|---|
| #sample | count |
| ambiguous | 0 |
| MgArd0009 | 91 |
| MgArd0017 | 166 |
| MgArd0038 | 1208 |
| MgArd0029 | 193 |
| unmatched | 245 |
| MgArd0001 | 119 |
| MgArd0081 | 246 |
| MgArd0046 | 401 |
| MgArd0054 | 243 |
| MgArd0073 | 474 |
| MgArd0062 | 1127 |

With barcode mismatches >1 sequence can corresponding to several samples. Sequences that match at only one sample are affected to this sample but the others (ambiguous) are not re-affected to a sample.

Sequences without known barcode. So these sequences are non-affected to a sample.

A tar archive is created by grouping one (or a pair of) fastq file per sample with the names indicated in the first column of the barcode tabular file.

# Format: Barcode

BARCODE FILE is expected to be tabulated:
- first column corresponds to the sample name (unique, without space)
- second to the forward sequence barcode used (None if only reverse barcode)
- optional third is the reverse sequence barcode (optional)

Take care to indicate sequence barcode in the strand of the read, so you may need to reverse complement the reverse barcode sequence. Barcode sequence must have the same length.

Example of barcode file.
The last column is optional, like this, it describes sample multiplexed by both fragment ends.

| MgArd00001 | ACAGCGT | ACGTACA |
|---|---|---|

# Format : FastQ

FASTQ : Text file describing biological sequence in 4 lines format:

- first line start by "@" correspond to the sequence identifier and optionally the sequence description. "@Sequence_1 description1"
- second line is the sequence itself. "ACAGC"
- third line is a "+" following by the sequence identifier or not depending on the version
- fourth line is the quality sequence, one code per base. The code depends on the version and the sequencer

```
@HNHOSKD01ALD0H
ACAGCGTCAGAGGGGTACCAGTCAGCCATGACGTAGCACGTACA
+
CCCFFFFFFHHHHHHJJIJJJHHFF@DEDDDDDDD@CDDDDACDD
```

# How it works ?

For each sequence or sequence pair the sequence fragment at the beginning (forward multiplexing) of the (first) read or at the end (reverse multiplexing) of the (second) read will be compare to all barcode sequence.

If this fragment is equal (with less or equal mismatch than the threshold) to one (and only one) barcode, the fragment is trimmed and the sequence will be attributed to the corresponding sample.

Finally fastq files (or pair of fastq files) for each sample are included in an archive, and a report describes how many sequence are attributed for each sample.

# Pre-process tool

From demultiplex tool

454 or illumina

ITS

MiSeq Fastq R2

MiSeq Fastq R1

Already merged

FROGS Pre-process Illumina ✖

Archive file

dereplicated_file (fasta)

count_file (tabular)

summary_file (html)

**Pre-process**

# Pre-process

- Merging of R1 and R2 reads

- Delete sequences without good primers

- Finds and removes adapter sequences

- Delete sequence with not expected lengths

- Delete sequences with ambiguous bases (N)

- Dereplication


- + removing homopolymers (size = 8 ) for 454 data

- + quality filter for 454 data

## Example for:

- Illumina MiSeq data

- 1 sample

- Non joined

**FROGS Pre-process** merging, denoising and dereplication. (Galaxy Version r3.0-3.0) ▾ Options

**Sequencer**

Illumina ▾

Select the sequencing technology used to produce the sequences.

**Input type**

Files by samples ▾

Samples files can be provided in single archive or with two files (R1 and R2) by sample.

**Reads already contiged ?**

No ▾

The inputs contain 1 file by sample : R1 and R2 are already merged by pair.

**Samples**

1: Samples

**Name**

sampleA

The sample name.

**Reads 1**

📄 🗐 📁 1: http://genoweb.toulouse.inra.fr/~formation/15_FROGS/FROGS_ini/D TA/sampleA_R1.fastq ▾

R1 FASTQ file of paired-end reads.

**reads 2**

📄 🗐 📁 2: http://genoweb.toulouse.inra.fr/~formation/15_FROGS/FROGS_ini/D TA/sampleA_R2.fastq ▾

R2 FASTQ file of paired-end reads.

➕ Insert Samples

**Reads 1 size**

250

The maximum read1 size.

**Reads 2 size**

250

The maximum read2 size.

**mismatch rate.**

0.1

The maximum rate of mismatches in the overlap r

**Parameters for the merging**

**Merge software**

Vsearch ▾

Select the software to merge paired-end reads.

**Would you like to keep unmerged reads?**

Yes | No

No : Unmerged reads will be excluded; Yes : unmerged reads will be artificially combined with 100 N. (default No)

**Minimum amplicon size**

`340`

The minimum size for the amplicons.

**Maximum amplicon size**

`450`

The maximum size for the amplicons.

V4-16S variability
Mean size = 390 ncl.

**Sequencing protocol**

`Illumina standard` ▼

The protocol used for sequencing step: standard or custom with PCR primers as sequencing primers.

**5' primer**

`GTGCCAGCMGCCGCGGTAA`

The 5' primer sequence (wildcards are accepted). The orienta... ...ameters'.

Primer sequences

**3' primer**

`ATTAGAWACCCBDGTAGTCC`

The 3' primer sequence (wildcards are accepted). The orientation is detailed below in 'Primers parameters'.

✔ Execute

degenerate primer
are accepted
(IUPAC code)

Pre-process example 1

## Example for:

- Roche 454 data

- 1 sample

- Only one read (454 process)

**FROGS Pre-process** Step 1 in metagenomics analysis: denoising and dereplication. (Galaxy Version 1.5.0)          ▾ Options

**Sequencer**

454                                                                                                                    ▾

Select the sequencer family used to produce the sequences.

**Input type**

One file by sample                                                                                                     ▾

Samples files can be provided in single archive or with one file by sample.

**Samples**

1: Samples

**Name**

my_sample

The sample name.

**Sequence file**

📄 🗐 🗀   1: /work/formation/FROGS/454.fastq.gz                                                                         ▾

FASTQ file of sample.

➕ Insert Samples

**Minimum amplicon size**

380

The minimum size for the amplicons (with primers).

**Maximum amplicon size**

500

The maximum size for the amplicons (with primers).

**[V3 – V4] 16S variability**

**5' primer**

ACGGGAGGCAGCAG

The 5' primer sequence (wildcards are accepted). The orient... ...rameters'.

**3' primer**

AGGATTAGATACCCTGGTA

The 3' primer sequence (wildcards are accepted). The orientation is detailed below in 'Primers parameters'.

**Primer sequences**

✔ Execute

## Example for:

- Illumina MiSeq data
- 9 samples in 1 archive
- Joined
- Without sequenced PCR primers (Kozich protocol)



**FROGS Pre-process** Step 1 in metagenomics analysis: denoising and dereplication. (Galaxy Version 1.5.0)　▾ Options

**Sequencer**

| Illumina |

**Sequencing technology**

Select the sequencer family used to produce the sequences.

**Input type**

| Archive |

**One file per sample and all files are contained in a archive**

Samples files can be provided in single archive or with two files (R1 and R2) by sample.

**Archive file**

📄 🗐 📁 | 1: /work/project/frogs/Formation/100spec_90000seq_9samples_Hantagulumic.tar.gz |

The tar file containing the sequences file(s) for each sample.

**Reads already contiged ?**

| Yes |

**Paire-end sequencing all ready joined**

The archive contains 1 file by sample : Reads 1 and Reads 2 are already contiged by pair.

**Minimum amplicon size**

| 380 |

The minimum size for the amplicons.

**[V3 – V4] 16S variability**

**Maximum amplicon size**

| 500 |

The maximum size for the amplicons.

**Sequencing protocol**

| Custom protocol (Kozich et al. 2013) |

**No more primers**

The protocol used for sequencing step: standard or custom with PCR primers as sequencing primers.

✔ Execute

**Pre-process example 3**

# Which primers for 16S ?



| NGS platforms | 16S region | PCR primers | Estimated insert size to read (E. coli) | Sequencing |
|---|---|---|---|---|
| Illumina MiSeq PE (Pair End) | V3V4 | 341F & 805R | 427 bp | 250 bp x 2 or 300 bp x 2 |
| Illumina HiSeq/iSeq100 (Earth Microbiome Project) | V4 | 515FB & 806RB | 250 bp | 150 x 2 |

| Name of primer F=forward, R=reverse | Sequence |
|---|---|
| 8F | AGAGTTTGATCCTGGCTCAG |
| 27F | AGAGTTTGATCMTGGCTCAG |
| 336R | ACTGCTGCSYCCCGTAGGAGTCT |
| 337F | GACTCCTACGGGAGGCWGCAG |
| 337F | GACTCCTACGGGAGGCWGCAG |
| 341F | CCTACGGGNGGCWGCAG |
| 515FB | GTGYCAGCMGCCGCGGTAA |
| 518R | GTATTACCGCGGCTGCTGG |
| 533F | GTGCCAGCMGCCGCGGTAA |
| 785F | GGATTAGATACCCTGGTA |
| 805R | GACTACHVGGGTATCTAATCC |
| 806RB | GGACTACNVGGGTWTCTAAT |
| 907R | CCGTCAATTCCTTTRAGTTT |
| 928F | TAAAACTYAAAKGAATTGACGGG |
| 1100F | YAACGAGCGCAACCC |
| 1100R | GGGTTGCGCTCGTTG |
| 1492R | CGGTTACCTTGTTACGACTT |

# What does the Pre-process tool do?

- Merging of R1 and R2 reads with vsearch, flash or pear (only in command line)

- Delete sequences without good primers

- Finds and removes adapter sequences with cutadapt

- Delete sequence with not expected lengths

- Delete sequences with ambiguous bases (N)

- Dereplication

- + removing homopolymers (size = 8 ) for 454 data

- + quality filter for 454 data

**VSEARCH: a versatile open source tool for metagenomics.**
Rognes T, Flouri T, Nichols B, Quince C, Mahé F.
PeerJ. 2016 Oct 18;4:e2584. eCollection 2016.

Bioinformatics (2011) 27 (21):2957-2963. doi:10.1093/bioinformatics/btr507
**FLASH: fast length adjustment of short reads to improve genome assemblies**
TanjaMagoc, Steven L. Salzberg

Bioinformatics (2014) 30 (5):614–620 doi.org/10.1093/bioinformatics/btt593
**PEAR: a fast and accurate Illumina Paired-End reAd mergeR**
J. Zhang, K. Kobert, T. Flouri, A. Stamatakis,

EMBnet Journal, Vol17 no1. doi : 10.14806/ej.17.1.200
**Cutadapt removes adapter sequences from high-throughput sequencing reads**
Marcel Martin

# How work reads merging ?

WITH VSEARCH

# The aim of Vsearch is to merge R1 with R2

Case of a sequencing of overlapping sequences: case of 16S V3-V4 amplicon MiSeq sequencing:

Imagine a real amplicon sequence of 400bp

400bp

Imagine a Miseq paired sequencing of 2x250bp

R1 : 250bp

R2 : 250bp

Reconstructing amplicon sequence is possible thanks to the overlap region

Merged sequence length : 400bp, with 100bp overlap

# The aim of Vsearch is to merge R1 with R2

Case of a sequencing of over-overlapping sequences:

Imagine a real amplicon sequence of 200bp

200bp

Imagine a Miseq paired sequencing of 2x250bp

R1 : 250bp

R2 : 250bp

FROGS takes in charge this case in trimming over bases

200bp

Merged sequence length : 200bp, with 100% overlap

# Practice:

# Exercise

Go to « 16S » history

Launch the pre-process tool on that data set

→ objective: understand Vsearch software

# 16S dataset presentation:

A real analysis provided by Stéphane Chaillou *et al.*

Comparison of meat and seafood bacterial communities.

8 environment types (EnvType) :
- Meat → Ground Beef, Ground veal, Poultry sausage, Diced bacon
- Seafood → Cooked schrimps, Smoked salmon, Salmon filet, Cod filet



Chaillou, S. et al (2015). Origin and ecological selection of core and food-specific bacterial communities associated with meat and seafood spoilage. ISME J, 9(5):1105-1118.

# 16S dataset presentation:

BH    VH    MV    DL    CD    SF    FS    FC

From Chaillou paper, we produced simulated data:

- 64 samples of 16S amplicons

- R1 and R2 overlapping reads of 300 bases.

- 8 replicates per condition

- with errors among the linear curve 2.54e-1 2.79e-1

- with 10% chimeras

- Primers for V1-V3:
  - 5' AGAGTTTGATCCTGGCTCAG 3'
  - 5' CCAGCAGCCGCGGTAAT 3'

Chaillou, S. et al (2015). Origin and ecological selection of core and food-specific bacterial communities associated with meat and seafood spoilage. ISME J, 9(5):1105-1118.

**FROGS Pre-process** merging, denoising and dereplication. (Galaxy Version 3.2.1)    ▾ Options

**Sequencer**

| Illumina | ▾ |

Select the sequencing technology used to produce the sequences.

**Input type**

| TAR Archive | ▾ |

Samples files can be provided in a single TAR archive or sample by sample (with one or two files each).

**TAR archive file**

| 📄  🗐  📁 | 1: http://genoweb.toulouse.inra.fr/~formation/15_FROGS/Webinar_data/chaillou_withprimers_64renamedsam... | ▾ |

The TAR file containing the sequences file(s) for each sample.

**Are reads already merged ?**

| No | ▾ |

The archive contains 1 file by sample : R1 and R2 pair are already merged in one sequence.

**Reads 1 size**

| 300 |

The maximum read1 size.

**Reads 2 size**

| 300 |

The maximum read2 size.

**Mismatch rate.**

| 0.1 |

The maximum rate of mismatch in the overlap region

**Merge software**

| Vsearch |                    Vsearch is recommended (in command line, prefer pear)

Select the software to merge paired-end reads.

**Would you like to keep unmerged reads?**

| Yes | No |

No : Unmerged reads will be excluded; Yes : unmerged reads will be artificially combined with 100 N. (default No)

**Minimum amplicon size**

400

The minimum size for the amplicons (with primers).

**Maximum amplicon size**

580

The maximum size for the amplicons (with primers).

**Sequencing protocol**

Illumina standard ▾

The protocol used for sequencing step: standard or custom with PCR primers as sequencing primers.

**5' primer**

AGAGTTTGATCCTGGCTCAG

The 5' primer sequence (wildcards are accepted). The orientation is detailed below in 'Primers parameters' help section.

**3' primer**

CCAGCAGCCGCGGTAAT

The 3' primer sequence (wildcards are accepted). The orientation is detailed below in 'Primers parameters' help section.

✔ Execute

Amplicon length distribution before trimming and filtering

**Minimum amplicon size**

400

The minimum size for the amplicons (with primers).

**Maximum amplicon size**

580

The maximum size for the amplicons (with primers).

**Sequencing protocol**

Illumina standard ▼

The protocol used for sequencing step: standard or custom with PCR primers as sequencing primers.

**5' primer**

AGAGTTTGATCCTGGCTCAG

The 5' primer sequence (wildcards are accepted). The orientation is detailed below in 'Primers parameters' help section.

**3' primer**

CCAGCAGCCGCGGTAAT

The 3' primer sequence (wildcards are accepted). The orientation is detailed below in 'Primers parameters' help section.

✔ Execute

N.B.
Primers in 5' → 3' sens

Ex: read R1
@63_0 reference=otu_00517 position=1..300
AGAGTTTGATCCTGGCTCAGgatgaacgctagcggggaggcttaacacatgcaagccgaggggg
tagaattagcttgctaatttgagaccggcgcacgggtgcgtaacgcgtatgcaacttgccctactgaaaa
ggatagcccagagaaatttggattaatactttataatagactgaatggcatcatttagtttttgaaagattt
atcgcagtaggataggcatgcgtaagattagatagttggtgaggtaacggctcaccaagtcgacgatct
ttaggggggcctgagagggtgaacccccca

Ex: read R2
@63_0 reference=otu_00517 position=1..300 errors=5%G
ATTAGCGCGGCTGCTGGcacggagttagccggtgcttattcttctggtaccttcagctacttacac
gtaagtaggtttatccccagataaaagtagtttacaacccataaggccgtcatcctacacgcgggatggc
tggatcaggcttccacccattgtccaatattcctcactgctgcctcccgtaggagtctggtccgtgtctcag
taccagtgtgggggttcaccctctcaggcccccctaaagatcgtcgacttggtgagccgttacctcaccaa
ctatctaatcttacgcatgcct

R2 primer must be reverse transcribed

# Exercise

1. Do you understand how enter your primers ?

2. What is the « FROGS Pre-process:  dereplicated.fasta » file  ?

3. What is the «  FROGS Pre-process: count.tsv » file ?

4. Explore the file «  FROGS Pre-process: report.html  »

5. *Who loose a lot of sequences ?*

# Exercise

6.  How many sequences are there in the input file ?

7.  How many sequences did not have the 5' primer?

8.  How many sequences still are after pre-processing the data?

9.  How much time did it take to pre-process the data ?

10.  What is the length of your merged reads before preprocessing ?

11.  What can you tell about the samples, based on amplicon size distributions ?

Q1: Do you understand how enter your primers ?

**Minimum amplicon size**

400

The minimum size for the amplicons (with primers).

**Maximum amplicon size**

580

The maximum size for the amplicons (with primers).

**Sequencing protocol**

Illumina standard

The protocol used for sequencing step: standard or custom with PCR primers as sequencing primers.

**5' primer**

AGAGTTTGATCCTGGCTCAG

The 5' primer sequence (wildcards are accepted). The orientation is detailed below in 'Primers parameters' help section.

**3' primer**

CCAGCAGCCGCGGTAAT

The 3' primer sequence (wildcards are accepted). The orientation is detailed below in 'Primers parameters' help section.

✔ Execute

N.B.
Primers in 5' → 3' sens

R2 primer must be reverse transcribed
Use https://www.bioinformatics.nl/cgi-bin/emboss/revseq

Q2: What is the « FROGS Pre-process:  dereplicated.fasta » file  ?

Q3: What is the «  FROGS Pre-process: count.tsv » file ?



```
>06_5949;size=4 reference=otu_00680 position=1..300 errors=20%T
AGACCGGCGCACGGGTGCGTAACGCGTATGCAATCTGCCTTTCACAGAGGGATAGCCCAGAGAAATTTGGATTAATACCTCATA
>56_3551;size=1 reference=otu_00680 position=1..300 errors=21%A
AAGACCGGCGCACGGGTGCGTAACGCGTATGCAATCTGCCTTTCACAGAGGGATAGCCCAGAGAAATTTGGATTAATACCTCAT
>53_322;size=1 reference=otu_01408,otu_00680 amplicon=1..300,1..300 position=1..300
ATTGAACGGTGGCGGCATGCCTACGCGTATGCAATCTGCCTTTCACAGAGGGATAGCCCAGAGAAATTTGGATTAATACCTCAT
>56_2589;size=1 reference=otu_00680 position=1..300 errors=21%C
CAGACCGGCGCACGGGTGCGTAACGCGTATGCAATCTGCCTTTCACAGAGGGATAGCCCAGAGAAATTTGGATTAATACCTCAT
>56_7560;size=1 reference=otu_00680 position=1..300 errors=21%C
CAGACCGGCGCACGGGTGCGTAACGCGTATGCAATCTGCCTTTCACAGAGGGATAGCCCAGAGAAATTTGGATTAATACCTCAT
>36_626;size=1 reference=otu_00680 position=1..300 errors=21%C
CAGACCGGCGCACGGGTGCGTAACGCGTATGCAATCTGCCTTTCACAGAGGGATAGCCCAGAGAAATTTGGATTAATACCTCAT
>53_6128;size=1 reference=otu_00231,otu_00941,otu_00680 amplicon=1..300,1..300,1..30
CTGGCTCAGGATGAACGCCGTAACGCGTATGCAATCTGCCTTTCACAGAGGGATAGCCCAGAGAAATTTGGATTAATACCTCAT
>51_6860;size=1 reference=otu_00799,otu_00680 amplicon=1..300,1..300 position=1..300
GACGAAGGCGCACGGGTGCGTAACGCGTATGCAATCTGCCTTTCACAGAGGGATAGCCCAGAGAAATTTGGATTAATACCTCAT
```

| #id | BHT0.LOT01 | BHT0.LOT03 | BHT0.LOT04 | BHT0.LOT05 | BHT0.LOT06 | BHT0.LOT07 |
|---|---|---|---|---|---|---|
| 06_5949 | 0 | 0 | 0 | 0 | 0 | 0 |
| 56_3551 | 0 | 0 | 0 | 0 | 0 | 0 |
| 53_322 | 0 | 0 | 0 | 0 | 0 | 0 |
| 56_2589 | 0 | 0 | 0 | 0 | 0 | 0 |
| 56_7560 | 0 | 0 | 0 | 0 | 0 | 0 |
| 36_626 | 0 | 0 | 0 | 0 | 0 | 0 |
| 53_6128 | 0 | 0 | 0 | 0 | 0 | 0 |
| 51_6860 | 0 | 0 | 0 | 0 | 0 | 0 |
| 56_6896 | 0 | 0 | 0 | 0 | 0 | 0 |

| #id | BHT0.LOT01 | BHT0.LOT03 | BHT0.LOT04 | BHT0.LOT05 | BHT0.LOT06 | BHT0.LOT07 |
|---|---|---|---|---|---|---|
| 56_3997 | 0 | 0 | 0 | 0 | 0 | 0 |
| 59_6 | 0 | 0 | 0 | 0 | 191 | 111 |
| 59_5144 | 0 | 0 | 0 | 0 | 1 | 0 |
| 59_5852 | 0 | 0 | 0 | 0 | 1 | 0 |
| 60_1696 | 0 | 0 | 0 | 0 | 0 | 1 |
| 59_6656 | 0 | 0 | 0 | 0 | 1 | 0 |
| 59_1192 | 0 | 0 | 0 | 0 | 1 | 0 |

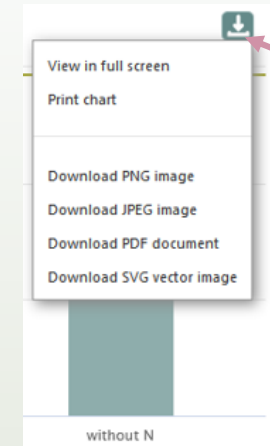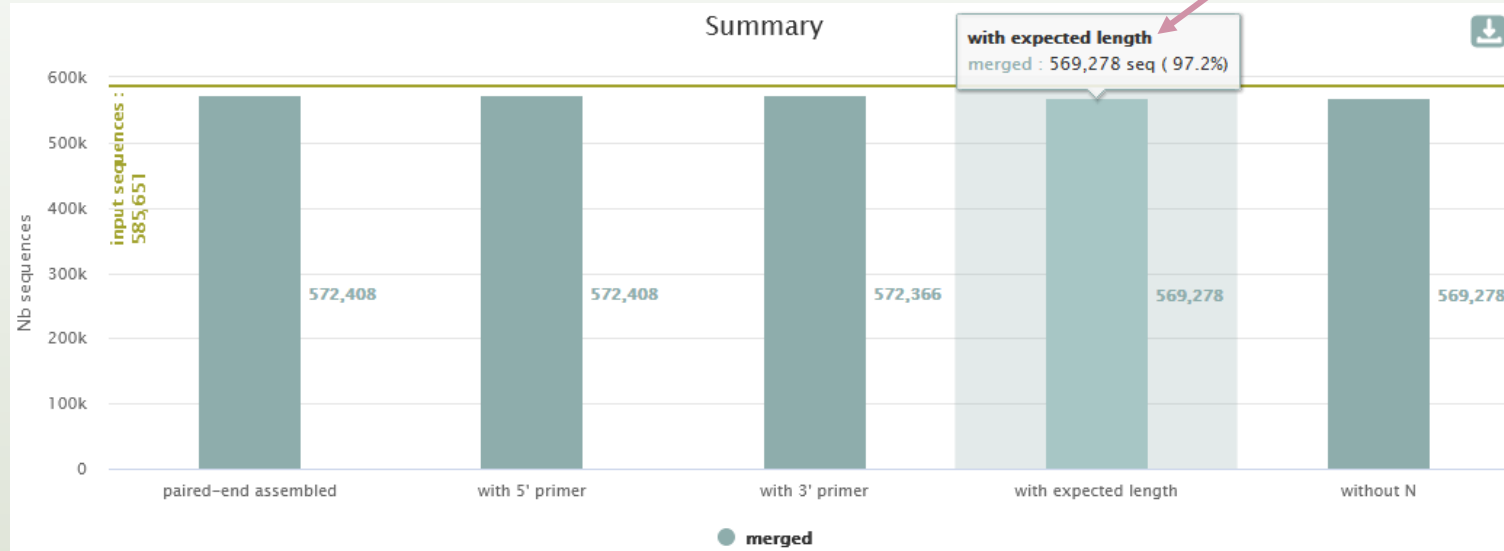Fasta sequence of all clean and dereplicated sequence *i.e.* only one copy of each sequence is kept

count table for each sequence in each sample

**Answer 4**

Q4: Explore the file « FROGS Pre-process: report.html »

By moving the mouse over the graphic, new information appears

You can download graphics or table in different formats

You can sort data in the table by clicking on the column headers

*Q5: Who loose a lot of sequences ?*

**53: FROGS Pre-process: report.html**

error
An error occurred with this dataset:

```
## Application
Software: preprocess.py (version: 3.2.2)
Command: /galaxydata/galaxy-preprod/my_tools/FROG
```

**52: FROGS Pre-process: count.tsv**

**51: FROGS Pre-process: dereplicated.
fasta**

## Dataset generation errors

**Dataset 53: FROGS Pre-process: report.html**

Tool execution generated the following error message:

```
Fatal error: Exit code 1 ()
Traceback (most recent call last):
  File "/galaxydata/galaxy-preprod/my_tools/FROGS_dev/app/preprocess.py", line 1290, in <module>
    process( args )
  File "/galaxydata/galaxy-preprod/my_tools/FROGS_dev/app/preprocess.py", line 1141, in process
    raise_exception( Exception( "\n\n#ERROR : The filters have eliminated all sequences (see summary for more details).\n\n" ))
  File "/galaxydata/galaxy-preprod/my_tools/FROGS_dev/lib/frogsUtils.py", line 45, in raise_exception
    raise exception
Exception:

#ERROR : The filters have eliminated all sequences (see summary for more details).
```

If your outputs are red, click on the bug to read the error message

it is likely that you did not enter the 3' primer in the right direction



Summary

All outputs are green but check the report.html

**65: FROGS Pre-process: report.html**

**64: FROGS Pre-process: count.tsv**

**63: FROGS Pre-process: dereplicated.fasta**



paired–end assembled
merged : 572,451 seq ( 97.75%)

Summary

Nb sequences

input sequences : 585,651

600k

400k

200k

0

572,451

572,451

977

976

976

paired–end assembled          with 5' primer          with 3' primer          with expected length          without N

● merged

Error in 3' primer sequence.
Primers must be similar with 10% of errors (~1 or 2 bases per primer)

**FROGS Pre-process** merging, denoising and dereplication. (Galaxy Version 3.2.1)          ▼ Options

**Sequencer**

Illumina                                                                                                          ▼

Select the sequencing technology used to produce the sequences.

**Input type**

TAR Archive                                                                                                      ▼

Samples files can be provided in a single TAR archive or sample by sample (with one or two files each).

> **TAR archive file**
>
> [📄] [🗐] [📂]  1: http://genoweb.toulouse.inra.fr/~formation/15_FROGS/Webinar_data/chaillou_withprimers_64renamedsam...  ▼
>
> The TAR file containing the sequences file(s) for each sample.

**Are reads already merged ?**

No                                                                                                               ▼

The archive contains 1 file by sample : R1 and R2 pair are already merged in one sequence.

> **Reads 1 size**
>
> 300
>
> The maximum read1 size.
>
> **Reads 2 size**
>
> 300
>
> The maximum read2 size.
>
> **Mismatch rate.**
>
> 0.1
>
> The maximum rate of mismatch in the overlap region
>
> **Merge software**
>
> Vsearch                                                                                                        ▼
>
> Select the software to merge paired-end reads.
>
> **Would you like to keep unmerged reads?**
>
> [ Yes ] [ No ]
>
> No : Unmerged reads will be excluded; Yes : unmerged reads will be artificially combined with 100 N. (default No)

To check the sequence quality use FASTQC (present in galaxy tools)

**FastQC: fastq/sam/bam**
    FastQC:Read QC reports using
    FastQC



Quality scores across all bases (Sanger / Illumina 1.9 encoding)

Q6: How many sequences are there in the input file ?
Q7: How many sequences did not have the 5' primer?
Q8: How many sequences still are after pre-processing the data?



Total number of sequences before preprocessing: 585 651

All sequences have the 5' primer
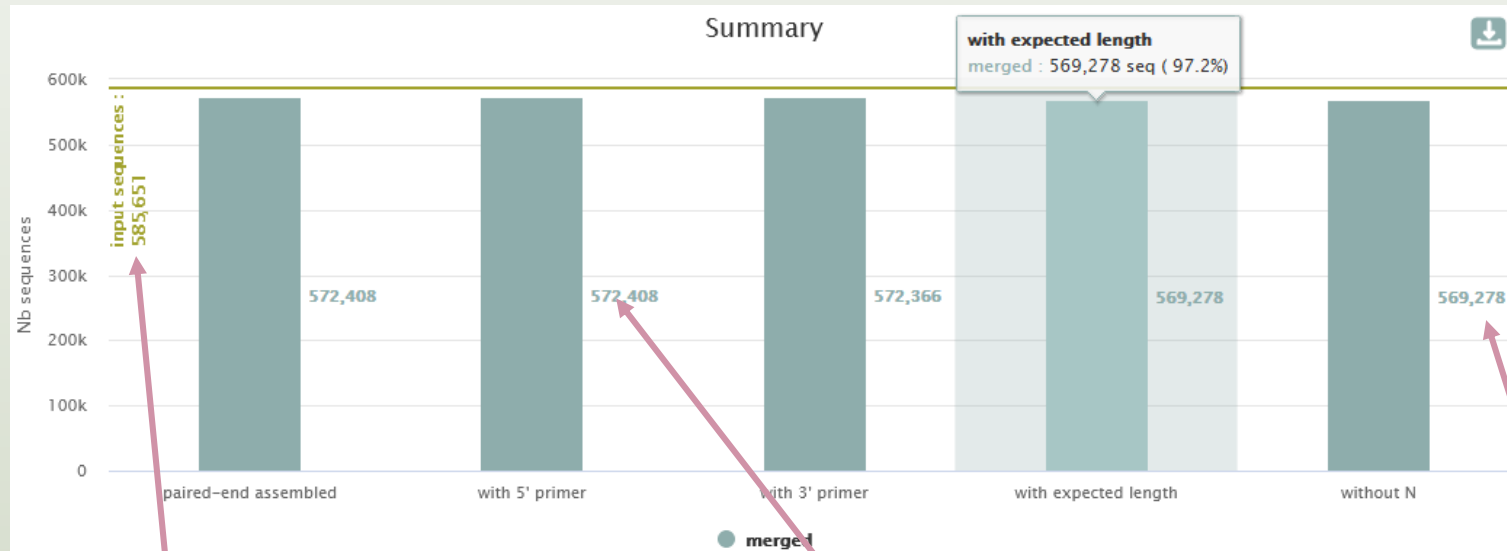
569 321 sequences are still after preprocessing

# Q9: How much time did it take to pre-process the data ?

**3: FROGS Pre-process: dereplicated.fasta**

287,252 sequences
format: **fasta**, database: **?**

## Application
Software: preprocess.py (version: 3.2.2)
Command: /galaxydata/galaxy-preprod
/my_tools/FROGS_dev/app/preprocess.py
illumina --output-dereplicated /galaxydata
/galaxy-prod/my_job_working_directory
/000/380/380454
/galaxy_dataset_731997.dat --ou

Click on « i »

**Tool: FROGS Pre-process**

| | |
|---|---|
| Name: | FROGS Pre-process: dereplicated.fasta |
| Created: | Tue 09 Mar 2021 05:00:30 PM (UTC) |
| Filesize: | 152.4 MB |
| Dbkey: | ? |
| Format: | fasta |
| Galaxy Tool ID: | FROGS_preprocess_3_2_2 |
| Galaxy Tool Version: | 3.2.2 |
| Tool Version: | |
| Tool Standard Output: | stdout |
| Tool Standard Error: | stderr |
| Tool Exit Code: | 0 |
| History Content API ID: | d7ff127129900fa8 |
| Job API ID: | 45c5decf7bd96ae1 |
| History API ID: | 96f266d5ffa0ae03 |
| UUID: | 58d5bf75-595e-422b-8c08-a16dbbe9110a |

| Input Parameter | Value | No |
|---|---|---|
| Sequencer | illumina | |
| Input type | archive | |
| TAR archive file | 1: http://genoweb.toulouse.inra.fr/~formation/15_FROGS/Webinar_data/chaillou_withprimers_64renamedsamples_V1V3_10000seq_R1R2.tar.gz | |
| Are reads already merged ? | paired | |
| Reads 1 size | 300 | |
| Reads 2 size | 300 | |
| Mismatch rate. | 0.1 | |
| Merge software | vsearch | |
| Would you like to keep unmerged reads? | False | |
| Minimum amplicon size | 400 | |
| Maximum amplicon size | 580 | |
| Sequencing protocol | standard | |
| 5' primer | AGAGTTTGATCCTGGCTCAG | |
| 3' primer | CCAGCAGCCGCGGTAAT | |

Retrieve the tool parameters

Stdout contains FROGS command lines and time execution

# Q10: What is the length of your merged reads before preprocessing ?

Show [All ▼] entries

Search: [_____]

**Select all samples**

| ✓ Samples ⇅ | process ⇅ | kept ⇅ | paired-end assembled ⇅ | with 5' primer ⇅ | with 3' primer ⇅ | with expected length ⇅ | without N ⇅ |
|---|---|---|---|---|---|---|---|
| ✓ BHT0.LOT01 | 9,282 | 97.90 | 9,087 | 9,087 | 9,087 | 9,087 | 9,087 |
| ✓ BHT0.LOT03 | 9,173 | 97.83 | 8,984 | 8,984 | 8,984 | 8,974 | 8,974 |
| ✓ BHT0.LOT04 | 9,171 | 97.79 | 8,969 | 8,969 | 8,968 | 8,968 | 8,968 |
| ✓ BHT0.LOT05 | 9,109 | 97.56 | 8,890 | 8,890 | 8,888 | 8,887 | 8,887 |
| ✓ BHT0.LOT06 | 9,193 | 97.86 | 8,996 | 8,996 | 8,996 | 8,996 | 8,996 |

Q10: What is the length of your merged reads before preprocessing ?

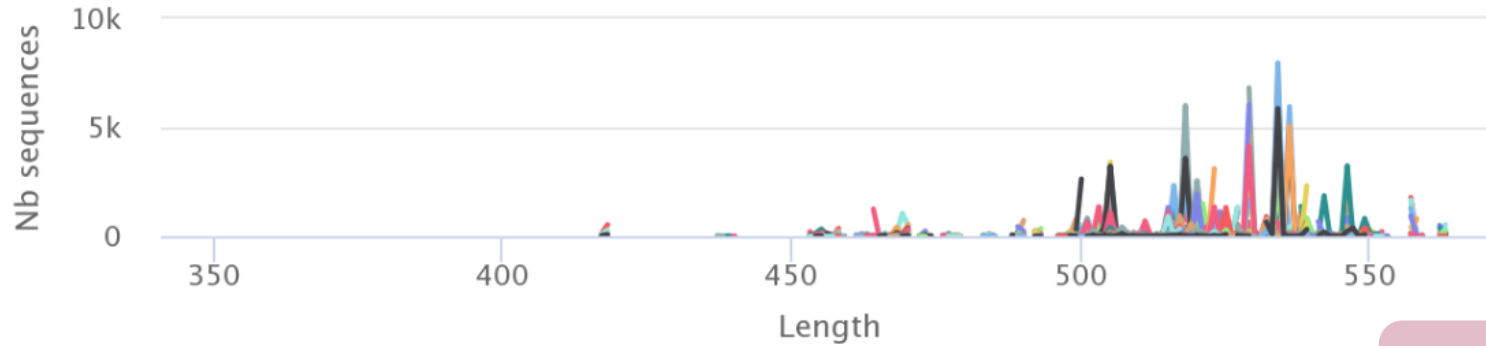| | | | | | | | |
|---|---|---|---|---|---|---|---|
| VHT0.LOT07 | 9,337 | 97.03 | 9,064 | 9,064 | 9,064 | 9,060 | 9,060 |
| VHT0.LOT08 | 9,436 | 97.33 | 9,192 | 9,192 | 9,192 | 9,184 | 9,184 |
| VHT0.LOT10 | 9,165 | 97.64 | 8,983 | 8,983 | 8,982 | 8,949 | 8,949 |

**With selection:** 📈 Display amplicon lengths    📈 Display preprocessed amplicon lengths

at the bottom of the table

Q10: What is the length of your merged reads before preprocessing ?



Amplicon length distribution before trimming and filtering

Before preprocessing:
343 < sequence length < 570

# Q11: What can you tell about the samples, based on amplicon size distributions ?



Preprocessed Amplicon Length distribution

« Filet Cabillaud » samples

« Saumon Fumé » samples

« Bœuf Haché » samples

For each EnvType, we can observe different amplicon sizes. They correspond to different species.
*N.B.* amplicons with same size can represent different species.

# Preprocess tool in brief

| | Take in charge |
|---|---|
| Illumina | ✔ |
| 454 | ✔ |
| Merged data | ✔ |
| Not merged data | ✔ |
| Without primers | ✔ |
| Only R1 or only R2 | 🚫 |
| Too distant R1 and R2 to be merged | ✔ |
| Over-overlapping R1 R2 | ✔ |

| | Take in charge |
|---|---|
| Archive .tar.gz | ✔ |
| Fastq | ✔ |
| Fasta | 🚫 |
| With only 1 primer | 🚫 |
| Multiplexed data | 🚫 |
| Demultiplexed data | ✔ |
| | |
| | |
| | |

# Processed data by FROGS in brief

# Clustering tool

# Why do we need clustering ?

Amplication and sequencing and are not perfect processes

A

B

Expected

A

B

Natural variability ?
Technical noise?
Contaminant?
Chimeras?

Results

A

B

Natural variability ?
Technical noise?
Contaminant?
Chimeras?

Expected

Results

16S variability
*Cf.* RRNDB (ribosomal RNA operons database)
max. 21 copies of 16S in bacteria (*Photobacterium damselae*)
ex. *E. coli* 7 copies

# To have the best accuracy:

## Method: All against all

- Very accurate

- Requires a lot of memory and/or time

=> Impossible on very large datasets without strong filtering or sampling

# How traditional clustering works ?

# Input order dependent results

# Single a priori clustering threshold



compromise threshold
unadapted threshold

natural limits of clusters

# Swarm clustering method



This sequences is the seed of the cluster.
Only the seed is kept for next processes.

The abundances of each sequence in the cluster are added together.
And the total abundance is given to the seed.

# Comparison Swarm and 3% clusterings



radius (97%)

Radius expressed as a percentage of identity with the central amplicon (97% is by far the most widely used clustering threshold)

# Comparison Swarm and 3% clusterings



TARA V9 (264 samples)

TARA V9 (908 samples)

crown size (numbers of amplicons in the OTU)

crown size (numbers or amplicons in the OTU)

seed abundance (numbers of copies)

seed abundance (numbers of copies)

identity (%)
97
90

clusters produced with swarm using d = 1

More there is sequences, more abundant clusters are enlarged (more amplicon in the cluster).
More there are sequences, more there are artefacts

# SWARM

A robust and fast clustering method for amplicon-based studies.

The purpose of **swarm** is to provide a novel clustering algorithm to handle large sets of amplicons.

**swarm** results are resilient to input-order changes and rely on a small **local** linking threshold $d$, the maximum number of differences between two amplicons.

**swarm** forms stable high-resolution clusters, with a high yield of biological information.

longer but more accurate

small OTU (made of 2 rare amplicons)

virtual amplicon

Mahé, Frédéric et al. "Swarm v2: highly-scalable and high-resolution amplicon clustering." *PeerJ* vol. 3 e1420. 10 Dec. 2015, doi:10.7717/peerj.1420

# Cluster stat tool

**FROGS Clusters stat** Process some metrics on clusters. (Galaxy Version 3.2.1)   ▼ Options

**Abundance file**

[ 🗎 | ⎘ | 🗁 ]   6: FROGS Clustering swarm: abundance.biom                    ▼

Clusters abundance (format: BIOM).

✔ Execute

# Practice:

LAUNCH CLUSTERING AND CLUSTERSTAT TOOLS

# Exercise

Go to « 16S » history

Launch the Clustering SWARM tool on that data set with guideline 3.2 *i.e. aggregation distance =1*

→ objectives :
- understand the outputs from clustering
- understand the ClusterStat utility

# Exercise

1. How many clusters do you get ?

Launch FROGS **Cluster Stat tools** on the previous abundance biom file

FROGS Clusters stat Process
some metrics on clusters.

# Exercise

2. Interpret the boxplot: **Clusters size summary**

3. Interpret the table: **Clusters size details - How many single singletons do you find?**

4. What can we say by observing the **sequence distribution**?

5. How many clusters share "BHT0.LOT08" with at least one other sample?

6. How many clusters could we expect to be shared ?

7. How many sequences represent the 106 specific clusters of "CDT0.LOT06"?

8. This represents what proportion of "CDT0.LOT06"?

9. What do you think about it?

10. How do you interpret the « Hierarchical clustering » ?

Q1: How many clusters do you get ?

Q2: Interpret the boxplot: **Clusters size summary**

Q3: Interpret the table: **Clusters size details - How many single singletons do you find?**



| Clusters distribution | Sequences distribution | Samples distribution |

Clusters
19,847

Sequences
569,278

Most of clusters are single singletons

Clusters size summary

Clusters size distribution

| Decile | Value |
|--------|-------|
| Min | 1 |
| 1 | 1 |
| 2 | 1 |
| 3 | 1 |
| 4 | 1 |
| Median | 1 |
| 6 | 1 |
| 7 | 1 |
| 8 | 1 |
| 9 | 1 |
| Max | 84,850 |

**Q4: What can we say by observing the sequence distribution?**



Cumulative sequences proportion by cluster size

Most of sequences are contained in big clusters

N.B.: Select area to zoom in.

The small clusters represent few sequences

| | Total clusters | Shared clusters | Own clusters | Total sequences | Shared sequences | Own sequences |
|---|---|---|---|---|---|---|
| BHT0.LOT01 | 491 | 114 | 377 | 9,087 | 8,709 | 378 |
| BHT0.LOT03 | 433 | 140 | 293 | | | |
| BHT0.LOT04 | 474 | 152 | 322 | | | |
| BHT0.LOT05 | 475 | 153 | 322 | | | |
| BHT0.LOT06 | 490 | 156 | 334 | 8,996 | 8,662 | 334 |
| BHT0.LOT07 | 531 | 165 | 366 | 9,059 | 8,690 | 369 |
| BHT0.LOT08 | 430 | 201 | 229 | 8,715 | 8,486 | 229 |
| BHT0.LOT10 | 4?? | ??? | 308 | 8,938 | 8,630 | 308 |
| CDT0.LOT02 | | | 490 | 9,259 | 8,767 | 492 |
| CDT0.LOT04 | | | 302 | 8,917 | 8,609 | 308 |
| CDT0.LOT05 | 380 | 241 | 139 | 8,516 | 8,377 | 139 |
| CDT0.LOT06 | 362 | 256 | 106 | 8,370 | 8,264 | 106 |
| CDT0.LOT07 | 489 | 100 | 389 | | | ?89 |
| CDT0.LOT08 | 556 | 162 | 394 | | | ?8 |
| CDT0.LOT09 | 456 | 150 | 306 | | | ?8 |
| CDT0.LOT10 | 465 | 157 | 308 | | | ?8 |

Q5: How many clusters share "BHT0.LOT08" with at least one other sample?
Q6: How many clusters could we expect to be shared ?
Q7: How many sequences represent the 106 specific clusters of "CDT0.LOT06"?
Q8: This represents what proportion of "CDT0.LOT06"?
Q9: What do you think about it?

201 clusters of BHT0.LOT08 are common at least once with another sample

~30 % of the specific clusters of CDT0.LOT06 represent around ~1% of sequences
Could be interesting to remove if individual variability is not the concern of user

## Q10: How do you interpret the « Hierarchical clustering » ?



The « Hierachical clustering » is established with a Bray Curtis distance particularly well adapted to abundance table of very heterogenous values (very big and very small figures).

Newick tree available too, can be copied and pasted an tree viewer

**Newick**

(((((CDT0.LOT02,CDT0.LOT08):0.312,(CDT0.LOT04,((CDT0.LOT05,CDT0.LOT06):0.518,(CDT0.LOT09,(CDT0.LOT07,CDT0.LOT10):0.533):0.582):0.757):0.816):0.840,(((FCT0.LOT07,(FCT0.LOT03,FCT0.LOT05):0.257):0.262,((FCT0.LOT01,FCT0.LOT08):0.352,(FCT0.LOT06,(FCT0.LOT02,FCT0.LOT10):0.427):0.631):0.805):0.832,(((MVT0.LOT07,SFT0.LOT03):0.493,(FST0.LOT06,(SFT0.LOT06,(SFT0.LOT08,(SFT0.LOT01,SFT0.LOT07):0.132):0.345):0.354):0.570):0.655,(((MVT0.LOT06,(MVT0.LOT05,MVT0.LOT08):0.439):0.511,((FST0.LOT02,(FST0.LOT03,FST0.LOT05):0.147):0.179,((SFT0.LOT02,(SFT0.LOT04,SFT0.LOT05):0.211):0.227,((MVT0.LOT01,MVT0.LOT03):0.161,(MVT0.LOT09,MVT0.LOT10):0.341):0.466):0.526):0.661):0.681,(DLT0.LOT04,((((DLT0.LOT05,DLT0.LOT06):0.173,(DLT0.LOT08,((VHT0.LOT07,(VHT0.LOT01,VHT0.LOT08):0.095):0.184,(DLT0.LOT01,DLT0.LOT03):0.231):0.267):0.325):0.411,((BHT0.LOT04,(BHT0.LOT08,((BHT0.LOT01,BHT0.LOT07):0.224,(BHT0.LOT05,BHT0.LOT06):0.231):0.309):0.352):0.462,((VHT0.LOT03,VHT0.LOT06):0.387,(VHT0.LOT02,(BHT0.LOT10,(VHT0.LOT04,VHT0.LOT10):0.272):0.336):0.401):0.463):0.590):0.711,(BHT0.LOT03,((FST0.LOT07,(FST0.LOT01,(FST0.LOT08,FST0.LOT10):0.254):0.388):0.408,(DLT0.LOT07,DLT0.LOT10):0.440):0.666):0.734):0.745):0.827):0.856):0.875):0.911):0.938);

Samples distribution tab

Answer 10

N.B.: Hierarchical clustering is not all a phylogenetic tree ! Please consult with caution.

Open in FigTree v1.4.4

Not very clear

# Chimera removal tool

# What is chimera ?

PCR-generated chimeras are typically created when an aborted amplicon acts as a primer for a heterologous template. Subsequent chimeras are about the same length as the non-chimeric amplicon and contain the forward (for.) and reverse (rev.) primer sequence at each end of the amplicon.

**Chimera: from 5 to 45% of reads** (Schloss 2011)



*aborted amplification*

*next cycle's "primer"*

*chimeric sequence*

for.        rev.

# A smart removal chimera to be accurate

**We use a sample cross-validation**

### Sample A

a ▬▬▬▬ x1000
b ▬▬▬▬ x500
c ▬▬▬▬ x100
d ▬▬▬▬ x50
e ▬▬▬▬ x10
f ▬▬▬▬ x10
g ▬▬▬▬ x5

### Sample B

b ▬▬▬▬ x1000
d ▬▬▬▬ x500
h ▬▬▬▬ x100
i ▬▬▬▬ x50
f ▬▬▬▬ x10
e ▬▬▬▬ x10
g ▬▬▬▬ x5

" **d** " is view as chimera by Vsearch
Its " parents " are presents

" **d** " is view as normal sequence by Vsearch
Its " parents " are absents

⇒ For FROGS "d" is not a chimera
⇒ For FROGS "g" is a chimera, "g" is removed
⇒ FROGS increases the detection specificity

# Practice:

LAUNCH THE REMOVE CHIMERA TOOL

# Exercise

Go to «  16S » history

Launch the « FROGS Remove Chimera » tool

Follow by the « FROGS ClusterStat » tool

→ objectives :
- understand the efficiency of the chimera removal
- make links between small abundant OTUs and chimeras

# Exercise

1. Understand the « FROGS remove chimera : report.html»
   a. How many clusters are kept after chimera removal?
   b. How many sequences that represent ? So what abundance?
   c. What do you conclude ?

2. What is the size of the largest removed cluster of chimeras?

# Exercise

3. Rename html output in Chimera_report.html

Launch « FROGS ClusterStat » tool on non_chimera_abundance.biom

4. Compare the HTML files
   a. Of what are mainly composed singleton ? (compare with previous report.html)
   b. What are their abundance?
   c. What do you conclude ?

Answer 1

Remove summary

Clusters

Kept : 5,922

Removed : 13,925

Abundance

Removed : 14,304

Kept : 554,974

5922 clusters are kept.
The 13925 removed clusters
represent ~2.5 % of sequences

Here, chimera clusters
represent many clusters ~70%
but very few sequences.

Removed clusters are low
abundance clusters.

## Q2: What is the size of the largest removed cluster of chimeras?

| Sample | Clusters kept | Cluster abundance kept | Chimeric clusters removed | Chimeric abundance removed | Abundance of the most abundant chimera removed | Individual chimera detected | Individual chimera abundance detected | Abundance of the most abundant individual chimera detected |
|--------|---------------|------------------------|---------------------------|----------------------------|-----------------------------------------------|----------------------------|---------------------------------------|-----------------------------------------------------------|
| VHT0.LOT02 | 205 | 8,862 | 366 | 410 | 19 | 372 | 446 | 19 |
| MVT0.LOT10 | 253 | 9,312 | | | 10 | 169 | 304 | 92 |
| VHT0.LOT08 | 261 | 8,852 | | | 10 | 310 | 344 | 1 |
| VHT0.LOT01 | 197 | 8,831 | | | 8 | 365 | 382 | 8 |

The largest cluster of chimeras contained 19 sequences.

92 chimeras are detected but only 10 are removed because 82 have been invalidated by the cross validation

## Q3: Rename html output in Chimera_report.html

**11: FROGS Remove chimera: report.html**

Attributes    Convert Format    Dat

**Edit Attributes**

**Name:**

Chimera_report.html

**Info:**

## Application
Software :/gal
/galaxy-preprod/my_tools

**11: Chimera_report.html**

## Answer 4

Q4a: Of what are mainly composed singleton ? (compare with previous report.html)
Q4b: What are their abundance?
Q4c: What do you conclude ?

| Cluster size ↑↓ | Number of cluster ↑↓ | % of all clusters ↑↓ |
|---|---|---|
| 1 | 19,078 | 96.13 |
| 2 | 148 | 0.75 |
| 3 | 22 | 0.11 |
| 4 | 10 | 0.05 |

Cluster_Stat report after clustering

Most small clusters are composed of chimeras

| Cluster size ↑↓ | Number of cluster ↑↓ | % of all clusters ↑↓ |
|---|---|---|
| 1 | 5,287 | 89.28 |
| 2 | 48 | 0.81 |
| 3 | 15 | 0.25 |
| 4 | 7 | 0.12 |

Cluster_Stat report after chimera removing

# OTU Filter tool

# OTU Filter

**Goal:** This tool deletes OTU among conditions enter by user. If an OTU reply to at least 1 criteria, the OTU is deleted.

**Criteria:**

**The OTU prevalence:** The number of times the OTU is present in the environment, *i.e.* the number of samples where the OTU must be present.

**OTU size:** An OTU that is not large enough for a given proportion or count will be removed.

**Biggest OTU:** Only the X biggest are conserved.

**Contaminant:** If OTU sequence matches with phiX, chloroplastic/mitochondrial 16S of A. Thaliana or your own contaminant sequence.

**One tool, 4 criteria**

1
2
3
4

**FROGS OTU Filters** Filters OTUs on several criteria. (Galaxy Version 3.2.2)  ▼ Options

**Sequences file**

9: FROGS Remove chimera: non_chimera.fasta

The sequence file to filter (format: fasta).

**Abundance file**

10: FROGS Remove chimera: non_chimera_abundance.biom

The abundance file to filter (format: BIOM).

**Minimum prevalence**

4

Fill the field only if you want this treatment. Keep OTU if it is present in at least this number of samples.

**Minimum OTU abundancy as proportion or count. We recommend to use a proportion of 0.00005.**

as proportion

> **Minimum proportion of sequences abundancy to keep OTU**
>
> 0.00005
>
> Fill the field only if you want this treatment. Example: 0.00005, recommended by Bokulich et al 2013, to keep OTU with at least 0,005% of all sequences).

**N biggest OTUs**

Fill the fields only if you want this treatment. Keep the N biggest OTU.

**Search for contaminant OTU.**

Use contaminant fasta file from the server

Either you use your own contaminant fasta file or you select one among available ones.

> **Contaminant databank**
>
> phiX
>
> For example the phiX databank (the phiX is a control added in Illumina sequencing technologies).

✔ Execute

**FROGS OTU Filters** Filters OTUs on several criteria. (Galaxy Version 3.2.2) ▾ Options

**Sequences file**

9: FROGS Remove chimera: non_chimera.fasta ▾

The sequence file to filter (format: fasta).

**Abundance file**

10: FROGS Remove chimera: non_chimera_abundance.biom ▾

The abundance file to filter (format: BIOM).

**Minimum prevalence**

Prevalence

4

Fill the field only if you want this treatment. Keep OTU if it is present in at least this number of samples.

Here, user wants that each OTU are present in at least 4 samples.

**Minimum OTU abundancy as proportion or count. We recommend to use a proportion of 0.00005.**

as proportion ▼

**Minimum proportion of sequences abundancy to keep OTU**

0.00005

Fill the field only if you want this treatment. Example: 0.00005, recommended by Bokulich et al 2013, to keep OTU with at least 0,005% of all sequences).

## OR

**Minimum OTU abundancy as proportion or count. We recommend to use a proportion of 0.00005.**

as count ▼

**Minimum number of sequences to keep OTU**

2

Fill the field only if you want this treatment. Ex: 2 to keep OTU with at least 2 sequences, so remove single singleton.

Here, user wants that each OTU has an abundance representing at least 0.005% of total number of sequences.

Here, user wants that each OTU has an abundance at least equals to 2 sequences -> single singleton will be removed.

**3**

**N biggest OTUs**

50

Fill the fields only if you want this treatment. Keep the N biggest OTU.

Here, user wants to keep the 50 biggest OTU.

Search for contaminant OTU.

Use contaminant fasta file from the server ▼

Either you use your own contaminant fasta file or you select one among available ones.

**Contaminant databank**

phiX ⟳

For example the phiX databank (the phiX is a control added in Illumina sequencing technologies).

**Remove phiX sequence (use as buffer while sequencing)**

**OR**

Search for contaminant OTU.

Use contaminant fasta file from the server ▼

Either you use your own contaminant fasta file or you select one among available ones.

**Contaminant databank**

Arabidopsis TAIR10 Chloroplast and mitochondrie ⟳

For example the phiX databank (the phiX is a control added in Illum

**Remove chloroplastic and mitochondrial 16S sequences of _A. Thaliana_**

**OR**

Search for contaminant OTU.

Use contaminant fasta file from the history ▼

Either you use your own contaminant fasta file or you select one among available ones.

**Select a contaminante reference from history**

📄 ⧉ 📁 | 31: contaminant.fasta ▼

**Place in your history (with getadata tool) your own file of contaminant sequences in fasta format.**

4

# Practice:

LAUNCH THE OTU FILTER TOOL

# Exercice:

Go to history «  16S » history

Launch « OTU Filter » tool with non_chimera_abundance.biom, non_chimera.fasta

Use 3 criteria to filter OTUs:

- OTU must be present at least in 4 samples
- Each OTU must represented a minimum of 0.005 % [1] of the totality of the sequences
- OTU of phiX [2] must be removed

→ objective : play with filters, understand their impacts on falses-positives OTUs

[1] *Nat Methods. 2013 Jan;10(1):57-9. doi: 10.1038/nmeth.2276. Epub 2012 Dec 2.*
**Quality-filtering vastly improves diversity estimates from Illumina amplicon sequencing.**
*Bokulich NA1, Subramanian S, Faith JJ, Gevers D, Gordon JI, Knight R, Mills DA, Caporaso JG.*

[2] https://www.illumina.com/products/by-type/sequencing-kits/cluster-gen-sequencing-reagents/phix-control-v3.html

# Exercice:

1. What are the output files of  "OTU Filter" ?

2. Explore "FROGS Filter : report.html" file. How many OTUs have you removed ? How many OTUs do they remain ? Which sample loses the most OTUs and for what reason?

3. Build the Venn diagram on the two filters. How many OTUs have you removed with each filter ?

4. How many own OTU remains in BHT0.LOT08  ? To retrieve this information, which tool do you need to launch previously ?

**FROGS OTU Filters** Filters OTUs on several criteria. (Galaxy Version 3.2.2)    ▼ Options

**Sequences file**

▢ ⧉ ▭ | 9: FROGS Remove chimera: non_chimera.fasta

The sequence file to filter (format: fasta).

**Abundance file**

▢ ⧉ ▭ | 10: FROGS Remove chimera: non_chimera_abundance.biom

The abundance file to filter (format: BIOM).

**16: FROGS OTU Filters: report.html**    👁 ✏ ✖

**15: FROGS OTU Filters: excluded.tsv**    👁 ✏ ✖

**Minimum prevalence**

**14: FROGS OTU Filters: abundance.biom**    👁 ✏ ✖

| 4 |

Fill the field only if you want this treatment. Keep OTU if it is present in at least this number of samples.

**13: FROGS OTU Filters: sequences.fasta**    👁 ✏ ✖

**Minimum OTU abundancy as proportion or count. We recommend to use a proportion of 0.00005.**

| as proportion                                                          ▼ |

**Minimum proportion of sequences abundancy to keep OTU**

| 0.00005 |

Fill the field only if you want this treatment. Example: 0.00005, recommended by Bokulich et al 2013, to keep OTU with at least 0,005%
~~l~~ sequences).

0.005% = 0.00005

~~B~~gest OTUs

|  |

Fill the fields only if you want this treatment. Keep the N biggest OTU.

**Search for contaminant OTU.**

| Use contaminant fasta file from the server                              ▼ |

Either you use your own contaminant fasta file or you select one among available ones.

**Contaminant databank**

| phiX                                                                    ▼ |

For example the phiX databank (the phiX is a control added in Illumina sequencing technologies).

✔ Execute

Filters by OTUs    Filters by samples

# Details by samples

Sort by **Kept** to find the answer

⬇ CSV

Show [ 10 ⬍ ] entries                                    Search: [                    ]

| Sample name ⇅ | Initial ⇅ | Kept ⇅ | Present in less than 4 samples ⇅ | Abundance < 0.005% (i.e 28 sequences ) ⇅ | Present in databank of contaminants ⇅ |
|---|---|---|---|---|---|
| SFT0.LOT06 | 433 | 34 | 376 | 398 | 0 |
| SFT0.LOT07 | 275 | 66 | 188 | 209 | 0 |
| SFT0.LOT01 | 308 | 70 | 216 | 238 | |
| SFT0.LOT08 | 324 | 88 | 215 | 236 | |
| CDT0.LOT02 | 234 | 92 | 141 | 142 | |
| MVT0.LOT10 | 253 | 96 | 155 | 157 | |
| SFT0.LOT03 | 196 | 97 | 92 | 98 | 0 |
| BHT0.LOT01 | 172 | 98 | 72 | 74 | 0 |
| CDT0.LOT07 | 187 | 99 | 87 | 88 | 0 |
| SFT0.LOT05 | 214 | 105 | 107 | 108 | 0 |

This sample have only very small clusters that are shared by very few other samples.

Showing 1 to 10 of 64 entries

Previous  **1**  2  3  4  5  6  7  Next

# Filters intersections

Draw a Venn to see which OTUs had been deleted by the filters chosen (Maximum 6 options):

- ✓ Present in less than 4 samples
- ✓ Abundance < 0.005% (i.e 28 sequences )
- ✓ Present in databank of contaminants

⚙ Venn

Present in less than 4 samples

Abundance < 0.005% (i.e 28 sequences )

5

5376

46

0

0

0

0

Present in databank of contaminants

- No phiX sequence.
- Most clusters are both small and not shared by 4 samples.

# Affiliation tool

**FROGS Affiliation OTU** Taxonomic affiliation of each OTU's seed (version 3.2.3)  ▾ Options

**Using reference database**

silva138.1 16S

Select reference from the list

**Also perform RDP assignation?**

Yes | No  — Optional

Taxonomy affiliation will be perform thanks to Blast. This option

**Taxonomic ranks**

Domain Phylum Class Order Family Genus Species

The ordered taxonomic ranks levels stored in the taxonomical ref

**OTU seed sequence**

31: FROGS Affiliation Filters: sequences.fasta

OTU sequences (format: fasta).

**Abundance file**

35: FROGS Affiliation Filters: abundance.biom

OTU abundances (format: BIOM).

✔ Execute

OR

For more details on FROGS databanks:
http://genoweb.toulouse.inra.fr/frogs_databanks/
assignation/readme.txt

Reference database dropdown list:

silva138.1 16S
silva138.1 pintail100 16S
silva138.1 pintail80 16S
silva138.1 pintail50 16S
silva138.1 18S
silva138.1 23S
silva138.1 28S
silva138 16S
silva138 pintail100 16S
silva138 pintail80 16S
silva138 pintail50 16S
silva138 18S
silva138 SSU
silva132 LSU
silva132 28S
silva132 16S
silva132_pintail100 16S
silva132_pintail80 16S
silva132_pintail50 16S
silva132 18S
silva132 23S
greengenes13_5
midas_S132_3.6
midas_S123_2.1.3
Psyringae CTS 20200131
pr2_4.12.0
rpoB_122017
Unite_Fungi_8.2_20200204
Unite_Euka_8.2_20200204
Unite_Fungi_8.0_18112018
Unite_Euka_8.0_18112018
RSyst_Diatom_7

Second dropdown list:

DAIRYdb_v1.1.2
EZBioCloud_052018
PHYMYCO-DB_2013
BOLD_COI-5P_022019
BOLD_COI-5P_1percentN_022019
MIDORI_UNIQUE_COI_20180221
MIDORI_UNIQUE_COI_MARINE_20180221
silva128 16S
silva128_pintail100 16S
silva128_pintail80 16S
silva128_pintail50 16S
silva128 18S
silva128 23S
silva123 16S
silva123 23S
silva123 18S
midas_S119_1.20
pr2_4.11.0
pr2_gb203_4.5
Unite_s_7.1_20112016

For ITS

132

# 1 Cluster = 2 affiliations

RDPClassifier*: one affiliation with bootstrap, on each taxonomic subdivision.

Bacteria;(1.0);Actinobacteriota;(1.0);Actinobacteria;(1.0);Propionibacteriales;(1.0);Propionibacteriaceae;(1.0);Cutibacterium;(1.0);Cutibacterium acnes;(0.57);

NCBI Blastn+** : one affiliation with identity %, coverage %, e-value,  alignment length and a special tag "**Multi-affiliation**".

Bacteria;Actinobacteriota;Actinobacteria;Propionibacteriales;Propionibacteriaceae;Cutibacterium;Multi-affiliation

Identity: 100% and Coverage: 100%

* Appl. Environ. Microbiol. August 2007 vol. 73 no. 16 5261-5267. doi : 10.1128/AEM.00062-07
**Naïve Bayesian Classifier for Rapid Assignment of rRNA Sequences into the New Bacterial Taxonomy.**
Qiong Wang, George M.Garrity, James M. Tiedje and James R. Cole

** BMC Bioinformatics 2009, 10:421.  doi:10.1186/1471-2105-10-421
**BLAST+: architecture and applications**
Christiam Camacho, George Coulouris, Vahram Avagyan, Ning Ma, Jason Papadopoulos,Kevin Bealer and Thomas L Madden

# Affiliation Strategy of FROGS

Blastn+ with "**Multi-affiliation**" management

Bacteria;Firmicutes;Bacilli;Staphylococcales;Staphylococcaceae;Staphylococcus;Staphylococcus xylosus

Bacteria;Firmicutes;Bacilli;Staphylococcales;Staphylococcaceae;Staphylococcus;Staphylococcus saprophyticus

Strictly identical (V1-V3 amplification) on 499 nucleotides

Which one to choose?

# Affiliation Strategy of FROGS

Blastn+ with "**Multi-affiliation**" management

Bacteria;Firmicutes;Bacilli;Staphylococcales;Staphylococcaceae;Staphylococcus;Staphylococcus xylosus

Bacteria;Firmicutes;Bacilli;Staphylococcales;Staphylococcaceae;Staphylococcus;Staphylococcus saprophyticus

Strictly identical (V1-V3 amplification) on 499 nucleotides

Bacteria;Firmicutes;Bacilli;Staphylococcales;Staphylococcaceae;Staphylococcus;**Multi-affiliation**

We cannot choose without preconceived ideas.

# Practice:

LAUNCH THE FROGS AFFILIATION TOOL

# Exercice:

Go to history «  16S » history

Launch the « FROGS Affiliation » tool with

- SILVA 138.1 16S database pintail 100

→ objectives :
- understand abundance tables columns
- understand the BLAST affiliation

**FROGS Affiliation OTU** Taxonomic affiliation of each OTU's seed by RDPtools and BLAST (Galaxy Version 3.2.2)    ▼ Options

**Using reference database**

| silva138.1 pintail100 16S | ▼ |

Select reference from the list

**Also perform RDP assignation?**

| Yes | No |

Taxonomy affiliation will be perform thanks to Blast. This option allow you to perform it also with RDP classifier (default No)

**OTU seed sequence**

| 📄 | 🗇 | 📁 | 32: FROGS OTU Filters: sequences.fasta | ▼ |

OTU sequences (format: fasta).

**Abundance file**

| 📄 | 🗇 | 📁 | 33: FROGS OTU Filters: abundance.biom | ▼ |

OTU abundances (format: BIOM).

✔ Execute

# Exercise

1. What are the « **FROGS Affiliation tool** » output files ?

2. How many sequences are affiliated by BLAST ?

3. How many OTU have a "multiaffiliation" at Order ranks ?

4. Click on the « eye » button on the BIOM output file, what do you understand ?

# Exercise

Use the **Biom_to_TSV tool** on this last file and click again on the "eye" on the new output generated.

👁

FROGS BIOM to TSV Converts a BIOM file in TSV file. (Galaxy Version 3.2.2)    ▾ Options

**Abundance file**

| 📄 | 🗐 | 📁 | | 37: FROGS Affiliation OTU: affiliation.biom | ▾ |

The BIOM file to convert (format: BIOM).

Transform the biom file in tsv file (easy to manipulate on excel or R)

**Sequences file (optional)**

| 📄 | 🗐 | 📁 | | 32: FROGS OTU Filters: sequences.fasta | ▾ |

The sequences file (format: fasta). If you use this option the sequences will be add in TS...

Optional but very useful, insert sequence of OTU in the abundance table

**Extract multi-alignments**

| Yes | No |

If you have used FROGS ... ...ut multiple alignments in a second TS...

Build the multi_affiliations.tsv: the list of possible affiliations for each ambiguous OTU with multiaffiliation

✔ Execute

140

# Exercise

5. Click again on the "eye" on the new output generated.

Or open it in your favorite spreadsheet (Excel, google sheet, Calc...) !

Now, what do you think about the file format? What does it contain?

# Exercise

6.  Observe and describe

◦   In FROGS BIOM to TSV: abundance_silva.tsv,  the different columns of cluster 3

a.  how would you qualify the alignment between the OTU3 seed and the sequences of the silva database?

b.  What does it mean e-value = 0 ?

c.  What is the header of column that shows the sequence of OTU seed ?

d.  How many sequences have OTU3 in total ?

e.  How many sequences have OTU3 in MVT0.LOT10 ? What is the sample where OTU3 is absent ?

# Exercise

7. Observe and describe

◦ In FROGS BIOM to TSV: multi_affiliations.tsv, identifies the lines corresponding to cluster3

   a. Why cluster3 has a multiaffiliation for species ?

   b. Why "Bacteria;Firmicutes;Bacilli;Lactobacillales;Lactobacillaceae;Latilactobacillus;Lactobacillus sakei" is present 74 times ?

# Exercise

**Answer 1**

**19: FROGS Affiliation OTU: report.html**

**18: FROGS Affiliation OTU: affiliation.biom**

**Answer 2**

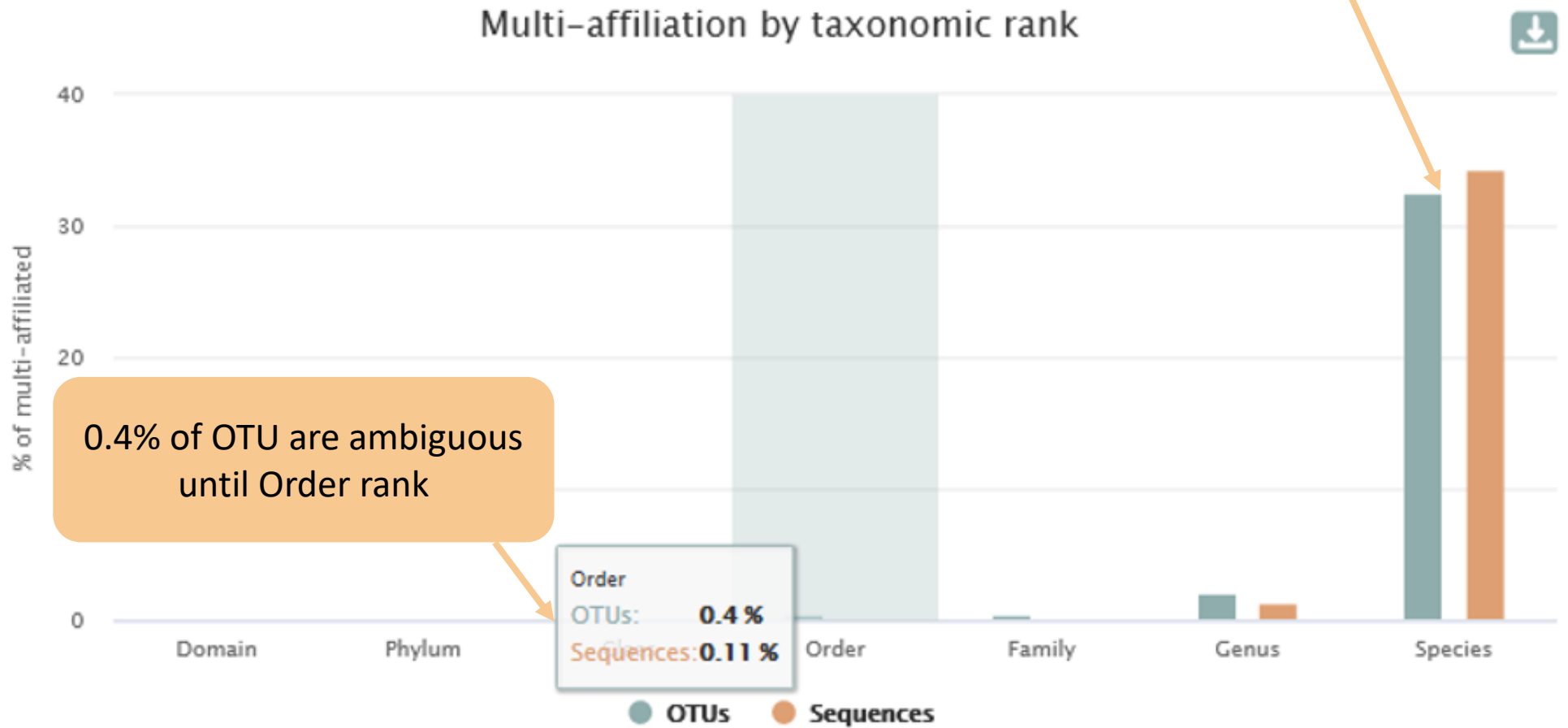OTUs affiliation

Without affiliation : 0

With affiliation : 495

Sequences affiliation

Without affiliation : 0

With affiliation : 547,520

All sequences have a blast affiliation

# Q4: Click on the « eye » button on the BIOM output file, what do you understand ?

{"matrix_type": "sparse", "shape": [495, 64], "date": "2021-03-03T11:57:55", "matr:
, 2, 23], [1, 3, 18], [1, 4, 19], [1, 5, 20], [1, 6, 29], [1, 7, 3], [1, 8, 1], [1,
9, 69], [2, 30, 98], [2, 31, 93], [2, 32, 38], [2, 33, 1682], [2, 34, 1598], [2, 3!
, 846], [3, 44, 210], [3, 45, 190], [3, 46, 122], [3, 47, 13], [3, 48, 3], [3, 49,
4, 61, 335], [4, 62, 540], [4, 63, 1943], [5, 0, 2408], [5, 1, 603], [5, 2, 1372],
, [7, 7, 24], [7, 9, 139], [7, 11, 7], [7, 12, 1], [7, 13, 37], [7, 14, 4], [7, 17.
46, 1], [9, 47, 4], [9, 51, 7], [9, 52, 4], [9, 56, 4], [9, 59, 4], [9, 60, 3], [9.
, [11, 47, 236], [11, 49, 24], [11, 50, 26], [11, 51, 44], [11, 52, 30], [11, 54, !
, 59, 71], [12, 60, 119], [12, 61, 16], [12, 62, 92], [12, 63, 272], [13, 0, 19],
27, 2], [14, 28, 3], [14, 29, 6], [14, 30, 8], [14, 31, 3], [14, 32, 10], [14, 33,
9], [17, 4, 17], [17, 5, 17], [17, 6, 20], [17, 7, 14], [17, 8, 3], [17, 9, 9], [1:
[18, 21, 34], [18, 22, 40], [18, 23, 105], [18, 25, 152], [18, 26, 2], [18, 27, 25
[20, 16, 16], [20, 17, 5], [20, 18, 1064], [20, 19, 12], [20, 20, 30], [20, 21, 33:
33, 43], [21, 34, 52], [21, 35, 59], [21, 36, 48], [21, 37, 44], [21, 38, 45], [21
, [23, 6, 16], [23, 7, 2], [23, 9, 2], [23, 10, 12], [23, 11, 27], [23, 12, 1], [2:
, [25, 30, 5], [25, 31, 23], [25, 36, 2], [25, 37, 16], [25, 38, 39], [25, 39, 4],
7, 16, 25], [27, 17, 7], [27, 18, 60], [27, 19, 40], [27, 20, 74], [27, 21, 41], [:
29, 23, 15], [29, 24, 4], [29, 25, 519], [29, 26, 1], [29, 27, 79], [29, 28, 1318].
31, 43, 16], [31, 44, 36], [31, 45, 91], [31, 46, 11], [31, 47, 2], [31, 56, 5], [:
76], [35, 12, 42], [35, 13, 2], [35, 14, 33], [35, 15, 78], [36, 0, 7], [36, 3, 1].
38, 28, 295], [38, 29, 45], [38, 30, 135], [38, 31, 566], [38, 32, 3], [38, 36, 3].
], [41, 17, 2], [41, 20, 5], [41, 21, 4], [41, 22, 1], [41, 23, 9], [41, 28, 1], [·
], [43, 38, 8], [43, 40, 2], [43, 42, 7], [43, 44, 3], [43, 46, 3], [43, 56, 2], [·
7, 11, 14], [47, 12, 1], [47, 13, 2], [47, 14, 1], [47, 15, 1], [47, 20, 2], [47, :
500], [50, 25, 21], [50, 26, 1], [50, 27, 1], [50, 28, 7], [50, 30, 6], [50, 31, 2
84], [52, 29, 3], [52, 30, 2], [52, 31, 21], [52, 32, 1], [52, 33, 6], [52, 34, 3].
, [54, 52, 1], [54, 55, 1], [54, 58, 3], [54, 60, 2], [55, 3, 8], [55, 4, 7], [55,
2 21 [57 6 21 [57 7 21 [57 9 11 [57 10 161 [57 11 921 [57 12 10

The biom file is not a human readable format. It is only very useful for bioinformaticians. To read the abundance table you have to transform the BIOM file in TSV file thanks to **BIOM_to_TSV tool**.

**Q5: what do you think about the TSV file format? What does it contain?**

The TSV format: tabular separated Value.
Universal format, ideal for different spreadsheets.

This file contain the abundance table and information about affiliation of OTUs.

| #comment | blast_taxonomy | blast_subject | blast_perc_identity | blast_perc_query_coverage |
|---|---|---|---|---|
| no data | Bacteria;Firmicutes;Bacilli;Lactobacillales;Listeriaceae;Brochothrix;Brochothrix thermosphacta | multi-subject | 100 | 100 |
| no data | Bacteria;Proteobacteria;Gammaproteobacteria;Enterobacterales;Vibrionaceae;Photobacterium;unknown species | FJ456662.1.1555 | 100 | 100 |
| no data | Bacteria;Firmicutes;Bacilli;Lactobacillales;Lactobacillaceae;Latilactobacillus;Multi-affiliation | multi-subject | 100 | 100 |
| no data | Bacteria;Actinobacteriota;Actinobacteria;Propionibacteriales;Propionibacteriaceae;Cutibacterium;Multi-affiliation | multi-subject | 100 | 100 |
| no data | Bacteria;Firmicutes;Bacilli;Lactobacillales;Lactobacillaceae;Leuconostoc;Multi-affiliation | multi-subject | 100 | 100 |
| no data | Bacteria;Firmicutes;Bacilli;Lactobacillales;Streptococcaceae;Lactococcus;Lactococcus piscium | AM943029.1.1242 | 99.799 | 100 |
| no data | Bacteria;Firmicutes;Bacilli;Erysipelotrichales;Erysipelotrichaceae;ZOR0006;unknown species | HG792212.1.1536 | 94.203 | 100 |
| no data | Bacteria;Firmicutes;Bacilli;Lactobacillales;Streptococcaceae;Lactococcus;Multi-affiliation | multi-subject | 100 | 100 |
| no data | Bacteria;Firmicutes;Bacilli;Lactobacillales;Lactobacillaceae;Weissella;Weissella ceti | FN813251.1.1461 | 99.799 | 100 |

| blast_evalue | blast_aln_length | seed_id | seed_sequence | observation_name | observation_sum | BHT0.LOT01 | BHT0.LOT03 | BHT0.LOT04 | BHT0.LOT05 | BHT0.LOT06 | BHT0.LOT07 | BHT0.LOT08 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 497 | 17_41 | GACGAACGCTGGCGGC... | Cluster_1 | 84849 | 791 | 402 | 433 | 911 | 1232 | 653 | 441 |
| 0 | 492 | 17_611 | ATTGAACGCTGGCGGC... | Cluster_2 | 31333 | 22 | 4 | 23 | 18 | 19 | 20 | 29 |
| 0 | 520 | 17_595 | GACGAACGCTGGCGGC... | Cluster_3 | 40711 | 342 | 70 | 71 | 218 | 81 | 199 | 114 |
| 0 | 468 | 17_257 | GACGAACGCTGGCGGC... | Cluster_4 | 22275 | 146 | 1251 | 263 | 327 | 180 | 118 | 293 |
| 0 | 497 | 17_4 | GATGAACGCTGGCGGC... | Cluster_5 | 29355 | 1842 | 217 | 1243 | 1799 | 1623 | 1374 | 954 |
| 0 | 497 | 17_23 | GACGAACGCTGGCGGC... | Cluster_6 | 21301 | 2408 | 603 | 1372 | 2231 | 2597 | 2218 | 1981 |
| 0 | 483 | 57_5 | GATGAACGCTGGCGGC... | Cluster_7 | 15272 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 0 | 499 | 17_420 | GACGAACGCTGGCGGC... | Cluster_8 | 16252 | 54 | 33 | 51 | 10 | 72 | 1 | 50 |
| 0 | 497 | 57_3 | TGCAAGTCGAACGCAC... | Cluster_9 | 11525 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |

a.  how would you qualify the alignment between the OTU3 seed and the sequences of the silva database?

Alignment is perfect ! 100% indentity and 100% coverage between OTU3 seed and the 520 nucleotides of sequence from silva database

b.  What does it mean e-value = 0 ?

The expect value is a parameter that describes the number of hits one can "expect" to see by chance when searching a database of a particular size. The lower the e-value, or the closer it is to zero, the more "significant" the match is.

c.  What is the header of column that shows the sequence of OTU seed ?

Seed_sequence

d.  How many sequences have OTU3 in total ?

40711 found in column " observation_sum"

e.  How many sequences have OTU3 in MVT0.LOT10 ? What is the sample where OTU3 is absent ?

| MVT0.LOT10 |
| --- |
| 4 |
| 0 |
| 6722 |
| 13 |
| 20 |

| CDT0.LOT02 |
| --- |
| 64 |
| 1 |
| 0 |
| 0 |
| 3 |

We can remark that OTU3 is particularly present in MV samples and rare in CD samples

a.  Why cluster3 has a multiaffiliation for species ?

In multi-affiliations.tsv  file, for cluster_3, we observe that 75 affiliations are possible for this OTU at species rank.

All strictly equivalent 100% identity and 100% coverage with 75 different sequences of silva database.

| | | | | | |
|---|---|---|---|---|---|
| ctobacillus;Lactobacillus sakei | CP025206.1448122.1449699 | 100 | 100 | 0 | 520 |
| ctobacillus;Lactobacillus sakei | CP020806.1000690.1002267 | 100 | 100 | 0 | 520 |
| ctobacillus;Lactobacillus sakei | CP025839.1959094.1960671 | 100 | 100 | 0 | 520 |
| ctobacillus;unknown species | KF601977.1.1550 | 100 | 100 | 0 | 520 |
| ctobacillus;Lactobacillus sakei | CP020806.811637.813214 | 100 | 100 | 0 | 520 |
| ctobacillus;Lactobacillus sakei | CP020806.1103805.1105382 | 100 | 100 | 0 | 520 |
| ctobacillus;Lactobacillus sakei | CP020806.1109220.1110797 | 100 | 100 | 0 | 520 |

b.  Why "Bacteria;Firmicutes;Bacilli;Lactobacillales;Lactobacillaceae;Latilactobacillus;Lactobacillus sakei" is present 74 times ?

Because these are 74 different strains of *L. sakei*. They have blast ID different.

# Silva pintail or not pintail ?

Pintail* represents the probability that the rRNA sequence contains anomalies or is a chimera, where 100 means that the probability for being anomalous or chimeric is low.

4 ranks of available databases in FROGS: 50 pintail, 80 pintail or 100 pintail or no pintail filter.

silva138.1 16S

silva138.1 pintail100 16S

silva138.1 pintail80 16S

silva138.1 pintail50 16S

silva138.1 18S

silva138.1 23S

silva138.1 28S

Only for 16S !

* http://aem.asm.org/content/71/12/7724.abstract

# Exemple between silva 138.1 and silva 138.1 pintail 100

130 identical blast best hits on **SILVA 138.1 pintail 100** databank

Cluster_4   Bacteria;Actinobacteriota;Actinobacteria;Propionibacteriales;Propionibacteriaceae;Cutibacterium;Cutibacterium acnes

Cluster_4   Bacteria;Actinobacteriota;Actinobacteria;Propionibacteriales;Propionibacteriaceae;Cutibacterium;Cutibacterium acnes 6609

Cluster_4   Bacteria;Actinobacteriota;Actinobacteria;Propionibacteriales;Propionibacteriaceae;Cutibacterium;Cutibacterium acnes C1

Cluster_4   Bacteria;Actinobacteriota;Actinobacteria;Propionibacteriales;Propionibacteriaceae;Cutibacterium;Cutibacterium acnes KPA171202

Cluster_4   Bacteria;Actinobacteriota;Actinobacteria;Propionibacteriales;Propionibacteriaceae;Cutibacterium;Cutibacterium acnes SK137

Cluster_4   Bacteria;Actinobacteriota;Actinobacteria;Propionibacteriales;Propionibacteriaceae;Cutibacterium;Cutibacterium acnes TypeIA2 P.acn17

Cluster_4   Bacteria;Actinobacteriota;Actinobacteria;Propionibacteriales;Propionibacteriaceae;Cutibacterium;Cutibacterium acnes TypeIA2 P.acn31

Cluster_4   Bacteria;Actinobacteriota;Actinobacteria;Propionibacteriales;Propionibacteriaceae;Cutibacterium;Cutibacterium acnes TypeIA2 P.acn33

# Exemple between silva 138.1 and silva 138.1 pintail 100

267 identical blast best hits on **SILVA 138.1 full** databank

? Cluster_4  Bacteria;Actinobacteriota;Actinobacteria;Corynebacteriales;Corynebacteriaceae;Corynebacterium;unknown species
? Cluster_4  Bacteria;Actinobacteriota;Actinobacteria;Propionibacteriales;Propionibacteriaceae;Cutibacterium;Aureobasidium melanogenum
  Cluster_4  Bacteria;Actinobacteriota;Actinobacteria;Propionibacteriales;Propionibacteriaceae;Cutibacterium;Cutibacterium acnes
  Cluster_4  Bacteria;Actinobacteriota;Actinobacteria;Propionibacteriales;Propionibacteriaceae;Cutibacterium;Cutibacterium acnes 266
  Cluster_4  Bacteria;Actinobacteriota;Actinobacteria;Propionibacteriales;Propionibacteriaceae;Cutibacterium;Cutibacterium acnes 6609
  Cluster_4  Bacteria;Actinobacteriota;Actinobacteria;Propionibacteriales;Propionibacteriaceae;Cutibacterium;Cutibacterium acnes C1
  Cluster_4  Bacteria;Actinobacteriota;Actinobacteria;Propionibacteriales;Propionibacteriaceae;Cutibacterium;Cutibacterium acnes hdn-1
  Cluster_4  Bacteria;Actinobacteriota;Actinobacteria;Propionibacteriales;Propionibacteriaceae;Cutibacterium;Cutibacterium acnes HL096PA1
  Cluster_4  Bacteria;Actinobacteriota;Actinobacteria;Propionibacteriales;Propionibacteriaceae;Cutibacterium;Cutibacterium acnes KPA171202
  Cluster_4  Bacteria;Actinobacte       ctinobacteria;Propionibacteriales;Propionibacteriaceae;Cutibacterium;Cutibacterium acnes SK137
  Cluster_4  Bacteria;Actinobacte       ctinobacteria;Propionibacteriales;Propionibacteriaceae;Cutibacterium;unknown species
  Cluster_4  Bacteria;Actinobacte       Induces a multi-affiliation up to phylum rank       terium;Cutibacterium acnes TypeIA2 P.acn17
  Cluster_4  Bacteria;Actinobacteriota;Actinobacteria;Propionibacteriales;Propionibacteriaceae;Cutibacterium;Cutibacterium acnes TypeIA2 P.acn31
  Cluster_4  Bacteria;Actinobacteriota;Actinobacteria;Propionibacteriales;Propionibacteriaceae;Cutibacterium;Cutibacterium acnes TypeIA2 P.acn33
? Cluster_4  Bacteria;Firmicutes;Bacilli;Lactobacillales;Carnobacteriaceae;Dolosigranulum;unknown species

| | accession number | organism name | sequence length | sequence quality | alignment quality | pintail quality | SILVA | taxonomy |
|---|---|---|---|---|---|---|---|---|
| ☐ | KF100699 | *uncultured bacterium* | 1341 | | | | | Bacteria▸Firmicutes▸Bacilli... |

# How choose the good affiliation ?

| | | | | | |
|---|---|---|---|---|---|
| Cluster_64 | Bacteria;Firmicutes;Bacilli;Staphylococcales;Staphylococcaceae;Staphylococcus;Staphylococcus xylosus | D83374.1.1477 | 100 | 100 | 0 | 499 |
| Cluster_64 | Bacteria;Firmicutes;Bacilli;Staphylococcales;Staphylococcaceae;Staphylococcus;Staphylococcus xylosus | CP007208.2831760.2833315 | 100 | 100 | 0 | 499 |
| Cluster_64 | Bacteria;Firmicutes;Bacilli;Staphylococcales;Staphylococcaceae;Staphylococcus;Staphylococcus xylosus | CP007208.1649831.1651386 | 100 | 100 | 0 | 499 |
| Cluster_64 | Bacteria;Firmicutes;Bacilli;Staphylococcales;Staphylococcaceae;Staphylococcus;Staphylococcus xylosus | CP007208.1426849.1428404 | 100 | 100 | 0 | 499 |
| Cluster_64 | Bacteria;Firmicutes;Bacilli;Staphylococcales;Staphylococcaceae;Staphylococcus;Staphylococcus xylosus | CP007208.1544187.1545742 | 100 | 100 | 0 | 499 |
| Cluster_64 | Bacteria;Firmicutes;Bacilli;Staphylococcales;Staphylococcaceae;Staphylococcus;Staphylococcus xylosus | LT963439.723352. | | | | |
| Cluster_64 | Bacteria;Firmicutes;Bacilli;Staphylococcales;Staphylococcaceae;Staphylococcus;Staphylococcus xylosus | CP013922.15879 | | | | |
| Cluster_64 | Bacteria;Firmicutes;Bacilli;Staphylococcales;Staphylococcaceae;Staphylococcus;Staphylococcus xylosus | CP013922.2856345.2857902 | 100 | 100 | 0 | 499 |
| Cluster_64 | Bacteria;Firmicutes;Bacilli;Staphylococcales;Staphylococcaceae;Staphylococcus;Staphylococcus xylosus | CP013922.2851139.2852696 | 100 | 100 | 0 | 499 |
| Cluster_64 | Bacteria;Firmicutes;Bacilli;Staphylococcales;Staphylococcaceae;Staphylococcus;Staphylococcus xylosus | CP013922.2904966.2906523 | 100 | 100 | 0 | 499 |
| Cluster_64 | Bacteria;Firmicutes;Bacilli;Staphylococcales;Staphylococcaceae;Staphylococcus;Staphylococcus xylosus | CP013922.2899760.2901317 | 100 | 100 | 0 | 499 |
| Cluster_64 | Bacteria;Firmicutes;Bacilli;Staphylococcales;Staphylococcaceae;Staphylococcus;Staphylococcus xylosus | CP013922.1470936.1472493 | 100 | 100 | 0 | 499 |
| Cluster_64 | Bacteria;Firmicutes;Bacilli;Staphylococcales;Staphylococcaceae;Staphylococcus;Staphylococcus xylosus | CP013922.1685669.1687226 | 100 | 100 | 0 | 499 |
| Cluster_64 | Bacteria;Firmicutes;Bacilli;Staphylococcales;Staphylococcaceae;Staphylococcus;Staphylococcus saprophyticus | EU855225.1.1531 | 100 | 100 | 0 | 499 |

**2 choices for cluster 64**

# How choose the good affiliation ?

| | | | | | | |
|---|---|---|---|---|---|---|
| Cluster_64 | Bacteria;Firmicutes;Bacilli;Staphylococcales;Staphylococcaceae;Staphylococcus;Staphylococcus xylosus | D83374.1.1477 | 100 | 100 | 0 | 499 |
| Cluster_64 | Bacteria;Firmicutes;Bacilli;Staphylococcales;Staphylococcaceae;Staphylococcus;Staphylococcus xylosus | CP007208.2831760.2833315 | 100 | 100 | 0 | 499 |
| Cluster_64 | Bacteria;Firmicutes;Bacilli;Staphylococcales;Staphylococcaceae;Staphylococcus;Staphylococcus xylosus | CP007208.1649831.1651386 | 100 | 100 | 0 | 499 |
| Cluster_64 | Bacteria;Firmicutes;Bacilli;Staphylococcales;Staphylococcaceae;Staphylococcus;Staphylococcus xylosus | CP007208.1426849.1428404 | 100 | 100 | 0 | 499 |
| Cluster_64 | Bacteria;Firmicutes;Bacilli;Staphylococcales;Staphylococcaceae;Staphylococcus;Staphylococcus xylosus | CP007208.1544187.1545742 | 100 | 100 | 0 | 499 |
| Cluster_64 | Bacteria;Firmicutes;Bacilli;Staphylococcales;Staphylococcaceae;Staphylococcus;Staphylococcus xylosus | LT963439.723352.724884 | 100 | 100 | 0 | 499 |
| Cluster_64 | Bacteria;Firmicutes;Bacilli;Staphylococcales;Staphylococcaceae;Staphylococcus;Staphylococcus xylosus | CP013922.1587968.1589525 | 100 | 100 | 0 | 499 |
| Cluster_64 | Bacteria;Firmicutes;Bacilli;Staphylococcales;Staphylococcaceae;Staphylococcus;Staphylococcus xylosus | CP013922.2856345.2857902 | 100 | 100 | 0 | 499 |
| Cluster_64 | Bacteria;Firmicutes;Bacilli;Staphylococcales;Staphylococcaceae;Staphylococcus;Staphylococcus xylosus | CP013922.2851139.2852696 | 100 | 100 | 0 | 499 |
| Cluster_64 | Bacteria;Firmicutes;Bacilli;Staphylococcales;Staphylococcaceae;Staphylococcus;Staphylococcus xylosus | CP013922.2904966.2906523 | 100 | 100 | 0 | 499 |
| Cluster_64 | Bacteria;Firmicutes;Bacilli;Staphylococcales;Staphylococcaceae;Staphylococcus;Staphylococcus xylosus | CP013922.2899760.2901317 | 100 | 100 | 0 | 499 |
| Cluster_64 | Bacteria;Firmicutes;Bacilli;Staphylococcales;Staphylococcaceae;Staphylococcus;Staphylococcus xylosus | CP013922.1470936.1472493 | 100 | 100 | 0 | 499 |
| Cluster_64 | Bacteria;Firmicutes;Bacilli;Staphylococcales;Staphylococcaceae;Staphylococcus;Staphylococcus xylosus | CP013922.1685669.1687226 | 100 | 100 | 0 | 499 |
| Cluster_64 | Bacteria;Firmicutes;Bacilli;Staphylococcales;Staphylococcaceae;Staphylococcus;Staphylococcus saprophyticus | EU855225.1.1531 | 100 | 100 | 0 | 499 |

- you have a preconceived notion
- you are familiar with the environment being studied
- you are looking for specific organisms as pathogens
- you collect bibliographical information

Ex:
*Staphylococcus saprophyticus* is a bacterium that can cause urinary tract infections in young women
and
*Staphylococcus xylosus* exists as a commensal on the skin of humans and animals and in the environment. It appears to be <u>much more common in animals </u>than in humans. S. xylosus has very occasionally been identified as a cause of human infection.

Maybe, for this cluster, S. xylosus is better

# Affiliation explorer

https://shiny.migale.inrae.fr/app/affiliationexplorer



A very user-friendly tool, developed by Mahendra Mariadassou and his collaborators (Maiage unit - INRAE Jouy-en-Josas). It allows to modify very simply the affiliations of an abundance table from FROGS.

# Affiliation explorer

https://shiny.migale.inrae.fr/app/affiliationexplorer

Demo video

# Divergence on the composition of microbial communities at the different taxonomic ranks

With the first versions of FROGS where multi-affiliation did not yet exist.

Affiliations and abundances of FROGS OTUs are they reliable ?

| Taxonomic ranks | Average divergence of the affiliations of the 10 samples (%) 500setA | Average divergence of the affiliations of the 10 samples (%) 100setA |
|---|---|---|
| Kingdom | 0.00 | 0.00 |
| Phylum | 0.46 | 0.41 |
| Class | 0.64 | 0.50 |
| Order | 0.94 | 0.68 |
| Familly | 1.18 | 0.78 |
| Genus | 1.76 | 1.30 |
| Species | 23.87 | 34.80 |

Affiliation was chosen with arbitrary criterion among all strictly equivalent affiliation

solution

Report on abundance table, the multiple identical affiliations

## Only one best hit

| Taxonomic ranks | Average divergence of the affiliations of the 10 samples (%) 500setA | Average divergence of the affiliations of the 10 samples (%) 100setA |
|---|---|---|
| Kingdom | 0.00 | 0.00 |
| Phylum | 0.46 | 0.41 |
| Class | 0.64 | 0.50 |
| Order | 0.94 | 0.68 |
| Familly | 1.18 | 0.78 |
| Genus | 1.76 | 1.30 |
| Species | 23.87 | 34.80 |

## Multiple best hit

| Taxonomic ranks | Median divergence of the affiliations of the 10 samples (%) 500setA | Median divergence of the affiliations of the 10 samples (%) 100setA |
|---|---|---|
| Kingdom | 0.00 | 0.00 |
| Phylum | 0.46 | 0.41 |
| Class | 0.64 | 0.50 |
| Order | 0.93 | 0.68 |
| Familly | 1.17 | 0.78 |
| Genus | 1.60 | 1.00 |
| Species | 6.63 | 5.75 |

With the FROGS guideline
OTU filter on abundance < 0.005%

| Taxonomic ranks | Median divergence of the affiliations of the 10 samples (%) 500setA filter: 0.005% - 505 OTUs | Median divergence of the affiliations of the 10 samples (%) 100setA filter: 0.005% - 100 OTUs |
|---|---|---|
| Kingdom | 0.00 | 0.00 |
| Phylum | 0.38 | 0.38 |
| Class | 0.57 | 0.48 |
| Order | 0.81 | 0.64 |
| Familly | 1.08 | 0.74 |
| Genus | 1.43 | 0.76 |
| Species | 1.53 | 0.78 |

# Affiliation Stat

**FROGS Affiliations stat** Process some metrics on taxonomies. (Galaxy Version 3.2.2)     ▾ Options

**Abundance file**

📄 🗐 📁   18: FROGS Affiliation OTU: affiliation.biom   ▾

OTUs abundances and affiliations (format: BIOM).

**Taxonomic ranks**

Domain Phylum Class Order Family Genus Species

The ordered taxonomic ranks levels stored in BIOM. Each rank is separated

**Rarefaction ranks**

Class Order Family Genus Species

The ranks that will be evaluated in rarefaction. Each rank is separated by one space.

**Affiliation processed**

FROGS blast   ▾

Select the type of affiliation processed. If your affiliation has been processed with an external tool: use 'Custom'.

✔ Execute

If your OTU are affiliated with less taxonomic ranks (species is missing for example), change it.

160

# Practice:

LAUNCH THE FROGS AFFILIATION STAT TOOL

# Exercice:

Go to  history «  16S » history

Launch the « FROGS Affiliation Stat » tool on last affiliation_abundance.biom

→ objectives :

   understand rarefaction curves and the diversity diagram

# Exercice:

1. Build the **rarefaction** curve on genus rank with the 10 samples that contain the least number of different genus.

2. SFT0.LOT06 and MVT0.LOT10 have they been sequenced deeply enough?

3. Build the **distribution** on FC samples *i.e.* "Filet de Cabillaud"

4. How many sequences are some *Brochothrix thermosphacta* ?

5. On the total of sequences, what is the proportion affiliated to the Firmicutes?

6. Among Firmicutes, how many are Bacilli ?

7. But what is the proportion of Firmicutes in the total of sequence of all sample ?

8. How many OTUs are align perfectly with a database sequence ?

**Q1: Build the rarefaction curve on genus rank with the 10 samples that contain the least number of different genus.**

| | Samples | Nb domain | Nb phylum | Nb class | Nb order | Nb family | Nb genus | Nb species | Nb sequences |
|---|---|---|---|---|---|---|---|---|---|
| ☑ | SFT0.LOT06 | 1 | 4 | 5 | 9 | 14 | | | |
| ☑ | | | | 5 | 12 | 26 | 35 | 57 | 8,821 |
| ☑ | SFT0.LOT01 | 1 | 4 | 6 | 13 | 27 | 39 | 63 | 8,859 |
| ☑ | FCT0.LOT01 | 1 | 5 | 6 | 13 | 24 | 41 | 96 | 8,504 |
| ☑ | SFT0.LOT05 | 1 | 5 | 7 | 18 | 32 | 50 | 95 | 8,728 |
| ☑ | SFT0.LOT08 | 1 | 4 | 6 | 13 | 33 | 53 | 77 | 8,788 |
| ☑ | BHT0.LOT01 | 1 | 7 | 9 | 20 | 35 | 5 | | |
| ☑ | SFT0.LOT04 | 1 | 6 | 8 | 17 | 34 | 5 | | |
| ☑ | SFT0.LOT03 | 1 | 5 | 8 | 1 | | | | |
| ☑ | SFT0.LOT02 | 1 | 6 | 7 | 1 | | | | |
| ☐ | MVT0.LOT10 | 1 | 4 | 5 | 17 | 31 | 57 | 83 | 9,143 |
| ☐ | CDT0.LOT02 | 1 | 6 | 8 | 22 | 36 | 58 | 85 | 8,750 |

1. Sort the table by genus number

2. Select the 10 first samples

3. At the bottom of the table click on

**With selection:** Genus ▾ | 📈 Display rarefaction | ⬤ Display distribution

Q2: SFT0.LOT06 and MVT0.LOT10 have they been sequenced deeply enough?

For MVTO.LOT10, the plateau does not seem to have been reached. Perhaps should have been sequenced more deeply?

With ~8000 sequences, all genus for this species are represented



Rarefaction

**Rarefaction curves**

Nb sampled sequences

Nb Genus

MVT0.LOT10    SFT0.LOT06

**Q3: Build the distribution on FC samples *i.e.* "Filet de Cabillaud"**

Use search to find only FC samples

Select the 8 samples of FC

CSV

Show

Search: FC

| | Samples ↑↓ | Nb domain ↑↓ | Nb phylum ↑↓ | Nb class ↑↓ | Nb order ↑↓ | Nb family ↑↓ | Nb genus ↑↓ | Nb species ↑↓ | Nb sequences ↑↓ |
|---|---|---|---|---|---|---|---|---|---|
| ☑ | FCT0.LOT01 | 1 | 5 | 6 | 13 | 24 | 41 | 96 | 8,504 |
| ☑ | FCT0.LOT02 | 1 | 6 | 8 | 23 | 40 | 67 | 126 | 7,638 |
| ☑ | FCT0.LOT03 | 1 | 8 | 10 | 26 | 45 | 71 | 122 | 8,608 |
| ☑ | FCT0.LOT05 | 1 | 8 | 10 | 25 | 44 | 78 | 139 | 8,577 |
| ☑ | FCT0.LOT06 | 1 | 8 | 10 | 29 | 53 | 97 | | |
| ☑ | FCT0.LOT07 | 1 | 5 | 7 | 24 | 46 | 8 | | |
| ☑ | FCT0.LOT08 | 1 | 7 | 9 | 2 | | | | |
| ☑ | FCT0.LOT10 | 1 | 7 | 9 | 2 | | | | |

At the bottom of the table click on

**With selection:** Genus ⌄ 📈 Display rarefaction ● Display distribution

Answer 3 4 & 5

Q3: Build the **distribution** on FC samples *i.e.* "Filet de Cabillaud"

Click on to see *Brochothrix thermosphacta*

Q4: How many sequences are some *Brochothrix thermosphacta* ?
Q5: On the total of sequences, what is the proportion affiliated to the Firmicutes?
Q6: Among Firmicutes, how many are Bacilli ?



| Name | Size | Global % | Parent % |
|---|---|---|---|
| root | 67211 | | |
| Bacteria | 67211 | 100.000 | 100.000 |
| Firmicutes | 20741 | 30.860 | 30.860 |
| Bacilli | 20658 | 30.736 | 99.600 |
| Lactobacillales | 19871 | 29.565 | 96.190 |
| Listeriaceae | 2649 | 3.941 | 13.331 |
| Brochothrix | 2649 | 3.941 | 100.000 |
| Brochothrix thermosphacta | 2649 | 3.941 | 100.000 |

Brochothrix thermosphacta nb children: 0

Detail on selected:

| Name | Size | Global % | Parent % |
|---|---|---|---|
| root | 67211 | | |
| Bacteria | 67211 | 100.000 | 100.000 |
| Firmicutes | 20741 | 30.860 | 30.860 |
| Bacilli | 20658 | 30.736 | 99.600 |
| Lactobacillales | 19871 | 29.565 | 96.190 |
| Listeriaceae | 2649 | 3.941 | 13.331 |
| Brochothrix | 2649 | 3.941 | 100.000 |
| Brochothrix thermosphacta | 2649 | 3.941 | 100.000 |

Brochothrix thermosphacta nb children: 0

A table appears

- 2649 sequences are some *Brochothrix thermosphacta*
- Firmicutes represent ~30% of total of sequences of these samples
- 99.6% of Firmicutes are Bacilli

Q7: But what is the proportion of Firmicutes in the total of sequence of all sample ?

Taxonomy distribution    Alignment distribution

At the top of the page, click on

⬤ Display global distribution

⬇ CSV

Show  10 ⬍  entries                                                Search: [                    ]

| | Samples ⇅ | Nb domain ⇅ | Nb phylum ⇅ | Nb class ⇅ | Nb order ⇅ | Nb family ⇅ | Nb genus ⇅ | Nb species ⇅ | Nb sequences ⇅ |
|---|---|---|---|---|---|---|---|---|---|
| ☐ | BHT0.LOT01 | 1 | 7 | 9 | 20 | 35 | 54 | 77 | 8,690 |
| ☐ | BHT0.LOT03 | 1 | 5 | 8 | 25 | 46 | 88 | 120 | 8,377 |
| ☐ | BHT0.LOT04 | 1 | 7 | 10 | 27 | 51 | 89 | 126 | 8,643 |
| ☐ | BHT0.LOT05 | 1 | 5 | 7 | 22 | 40 | 69 | 116 | 8,544 |
| ☐ | BHT0.LOT06 | 1 | 6 | 10 | 28 | 47 | 91 | 125 | 8,646 |
| ☐ | BHT0.LOT07 | 1 | 6 | 9 | 28 | 51 | 90 | 124 | 8,671 |
| ☐ | BHT0.LOT08 | 1 | 6 | 9 | 27 | 53 | 109 | 166 | 8,479 |
| ☐ | BHT0.LOT10 | 1 | 4 | 7 | 26 | 50 | 106 | 144 | 8,606 |
| ☐ | CDT0.LOT02 | 1 | 6 | 8 | 22 | 36 | 58 | 85 | 8,750 |
| ☐ | CDT0.LOT04 | 1 | 5 | 7 | 22 | 41 | 74 | 138 | 8,605 |

**With selection:**   [ Class ▼ ]   📈 Display rarefaction    ⬤ Display distribution

Q7: But what is the proportion of Firmicutes in the <u>total</u> of sequence of all sample ?



Click on Firmicutes

| Name | Size | Global % | Parent % |
|------|------|----------|----------|
| root | 547518 | | |
| Bacteria | 547518 | 100.000 | 100.000 |
| Firmicutes | 342409 | 62.538 | 62.538 |

Firmicutes represent 62% of Bacteria

Q7: But what is the proportion of Firmicutes in the total of sequence of all sample ?



To focus on Firmicutes, **double click on**. After you can apply color among rank depth.

# Q8: How many OTUs are align perfectly with a database sequence ?

At the top of the page, click on this tab

210 sequences are aligned with 100% identity and 100% coverage with a sequence of silva.

Taxonomy distribution | **Alignment distribution**

## Number of OTUs among their alignment results

| Coverage | [0% – 1%[ | [1% – 50%[ | [50% – 80%[ | [80% – 90%[ | [90% – 95%[ | [95% – 99%[ | [99% – 100%[ | [100%] |
|---|---|---|---|---|---|---|---|---|
| [100%] | | | | 3 | 13 | 38 | 199 | 210 |
| [99% – 100%[ | | | | | | 4 | 6 | 4 |
| [95% – 99%[ | | | | | 1 | 7 | 3 | 1 |
| [90% – 95%[ | | | | | 1 | 4 | | |
| [80% – 90%[ | | | | | | 1 | | |
| [50% – 80%[ | | | | | | | | |
| [1% – 50%[ | | | | | | | | |
| [0% – 1%[ | | | | | | | | |

Identity

0
50
100
150
200
250

by OTUs

by sequences

# Filters on affiliations

**FROGS Affiliation Filters** Filters OTUs on several affiliation criteria. (Galaxy Version 3.2.2)    ▼ Options

**Sequences file**

🗋 📑 🗀 | 13: FROGS OTU Filters: sequences.fasta    ▼

The sequence file to filter (format: fasta).

**Abundance file**

🗋 📑 🗀 | 18: FROGS Affiliation OTU: affiliation.biom    ▼

The abundance file to filter (format: BIOM).

**Taxonomic ranks**

Domain Phylum Class Order Family Genus Species

The ordered taxonomic ranks levels stored in BIOM. Each rank is separated by one space.

**Filtering mode**

○ Hidding mode
⊙ Deleting mode

Do you want to delete OTU or hide affiliations

**Filter on Blast affiliations**    👁

**Maximum e-value (between 0 and 1)**

0

Fill the field only if you want this treatment

**Minimum identity % (between 0 and 1)**

0.99

Fill the field only if you want this treatment

**Minimum coverage % (between 0 and 1)**

0.99

Fill the field only if you want this treatment

**Minimum alignment length**

Fill the field only if you want this treatment

**Filter blast affiliations including these taxon / word**

1: Filter blast affiliations including these taxon / word    🗑

**Full or partial taxon name**

unknown species

ex: "unknown species" or "subsp."

2: Filter blast affiliations including these taxon / word

**Full or partial taxon name**

Firmicutes

ex: "unknown species" or "subsp."

3: Filter blast affiliations including these taxon / word

**Full or partial taxon name**

subsp.

ex: "unknown species" or "subsp."

➕ Insert Filter blast affiliations including these taxon / word
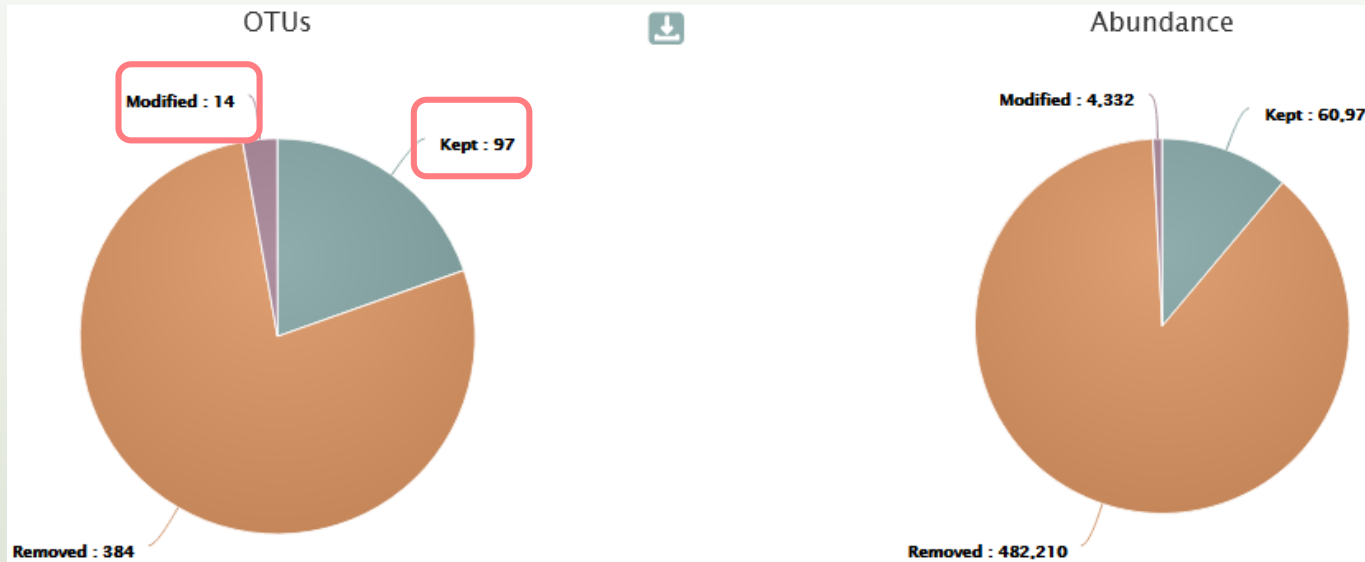
**Filter on RDP affiliations**    👁

**Taxonomical rank on which to apply bootstrap filter**

One of the available taxonomical rank name. Ex: Species

**Minimum bootstrap % (between 0 and 1)**

Fill these two fields if you want this treatment.

✔ Execute

Careful, it is case sensitive.
Firmicutes it's different of firmicutes !

Not open by default

2 modes: hidding or deleting mode.
All affiliations that enter in criteria of filter will be either hidden or deleted
• hidding: affiliation counting are not affected, affiliation are simply hidden
• deleting: all abundancies are computed again, affiliation have disappeared

# Practice:

LAUNCH THE FROGS AFFILIATION FILTER TOOL

# Exercice:

1. Apply filters to keep only sequences with perfect alignment with Silva sequences and affilliations without « unknown species » and « Firmicutes » terms. (deleting mode)

2. Apply filters to hide OTU affiliations that have not a perfect alignment with Silva sequences and the affilliations without « unknown species » and « Firmicutes » terms.

3. In deleting mode:

   - How many OTUs remain?

   - Among OTUs with multiaffiliation, How many were impacted/modified ?

4. In hidding mode:

   - What outputs change between deleted mode and hidding mode ?

**Answer 1**

**Answer 2**

FROGS Affiliation Filters Filters OTUs on several affiliation criteria. (Galaxy Version 3.2.2)    ▾ Options

**Sequences file**

13: FROGS OTU Filters: sequences.fasta
The sequence file to filter (format: fasta).

**Abundance file**

18: FROGS Affiliation OTU: affiliation.biom
The abundance file to filter (format: BIOM).

**Taxonomic ranks**

Domain Phylum Class Order Family Genus Species
The ordered taxonomic ranks levels stored in BIOM. Each rank is separated by one space.

**Filtering mode**

○ Hidding mode
◉ Deleting mode
Do you want to delete OTU or hide affiliations

Filter on Blast affiliations 👁

**Maximum e-value (between 0 and 1)**

Fill the field only if you want this treatment

**Minimum identity % (between 0 and 1)**

1
Fill the field only if you want this treatment

**Minimum coverage % (between 0 and 1)**

1
Fill the field only if you want this treatment

**Minimum alignment length**

Fill the field only if you want this treatment

**Filter blast affiliations including these taxon / word**

1: Filter blast affiliations including these taxon / word 🗑

**Full or partial taxon name**

unknown species
ex: "unknown species" or "subsp."

2: Filter blast affiliations including these taxon / word 🗑

**Full or partial taxon name**

Firmicutes
ex: "unknown species" or "subsp."

+ Insert Filter blast affiliations including these taxon / word

Filter on RDP affiliations 👁

✔ Execute

---

FROGS Affiliation Filters Filters OTUs on several affiliation criteria. (Galaxy Version 3.2.2)    ▾ Options

**Sequences file**

13: FROGS OTU Filters: sequences.fasta
The sequence file to filter (format: fasta).

**Abundance file**

18: FROGS Affiliation OTU: affiliation.biom
The abundance file to filter (format: BIOM).

**Taxonomic ranks**

Domain Phylum Class Order Family Genus Species
The ordered taxonomic ranks levels stored in BIOM. Each rank is separated by one space.

**Filtering mode**

◉ Hidding mode
○ Deleting mode
Do you want to delete OTU or hide affiliations

Filter on Blast affiliations 👁

**Maximum e-value (between 0 and 1)**

Fill the field only if you want this treatment

we want to keep the OTUs that have aligned perfectly with a sequencce of the silva bank *i.e.* 100% identity and 100% coverage

**Minimum alignment length**

Fill the field only if you want this treatment

**Filter blast affiliations including these taxon / word**

Enter key word

ex: "unknown species" or "subsp."

2: Filter blast affiliations including these taxon / word 🗑

**Full or partial taxon name**

Firmicutes
ex: "unknown species" or "subsp."

+ Insert Filter blast affiliations including these taxon / word

Filter on RDP affiliations 👁

✔ Execute

**Answer 3**

OTUs

Modified : 14

Kept : 97

Removed : 384

Abundance

Modified : 4,332

Kept : 60,978

Removed : 482,210

- Only 97 OTUs are kept <u>without modification</u>.
- 14 OTUs with multiaffliation were impacted/modified (all affiliations in the multi_affiliations with key words "unknown species" or "Firmicutes" were deleted).
  The consequences are either OTU have less multiaffiliations, or all multiaffiliations are impacted and OTU is deleted.
  The list of blast affiliations for multi-affiliated impacted OTUs are in **impacted_OTU.multiaffiliation.tsv**
- So, **111 OTUs** remains after filtering

: FROGS Affiliation Filters: report.html

FROGS Affiliation Filters: impacted_OTU.multi-affiliations.tsv

FROGS Affiliation Filters: impacted_OTU.tsv

FROGS Affiliation Filters: sequences.fasta

FROGS Affiliation Filters: abundance.biom

570 affiliations at species rank disappeared after filtering, including multi-affiliations of 54 OTUs.

## Taxon lost summary

Filtering criteria are applied by affiliation. So for blast, filters are not only applied on the blast consensus taxonomy but on each blast hit (cf multihit.tsv file).

For each OTU, none, part or all blast affiliations may be removed, resulting in unchanged / updated or deleted blast consensus taxonomy.

The detailed number of lost affiliations (not only the consensus) by rank are summarised. It may also precise if blast consensus multi-affiliation are lost.

| Affiliation method | Domain | Phylum | Class | Order | Family | Genus | Species |
|---|---|---|---|---|---|---|---|
| Blast | 0 | 7 | 12 | 31 ( including 1 multi-affiliation(s) ) | 56 ( including 1 multi-affiliation(s) ) | 164 ( including 4 multi-affiliation(s) ) | 570 ( including 54 multi-affiliation(s) ) |

In addition to the Firmicutes phylum that was deleted, there are 6 others that are deleted (unknow species or %id %cov)

31 Orders were deleted and 1 was a OTU with a multiaffiliation (-> Cluster_451)

Cluster_451 Bacteria;Firmicutes;Bacilli;Multi-affiliation;Multi-affiliation;Multi-affiliation;Multi-affiliation

| | | | | | |
|---|---|---|---|---|---|
| Cluster_451 | Bacteria;Firmicutes;Bacilli;Staphylococcales;Staphylococcaceae;Staphylococcus;Staphylococcus aureus | CP026068.13 | 100 | 100 | 0 | 497 |
| Cluster_451 | Bacteria;Firmicutes;Bacilli;Staphylococcales;Staphylococcaceae;Staphylococcus;Staphylococcus aureus | CP029082.38 | 100 | 100 | 0 | 497 |
| Cluster_451 | Bacteria;Firmicutes;Bacilli;Paenibacillales;Paenibacillaceae;Paenibacillus;Staphylococcus aureus | MIZO010000 | 100 | 100 | 0 | 497 |
| Cluster_451 | Bacteria;Firmicutes;Bacilli;Staphylococcales;Staphylococcaceae;Staphylococcus;Staphylococcus aureus | CP029030.22 | 100 | 100 | 0 | 497 |
| Cluster_451 | Bacteria;Firmicutes;Bacilli;Staphylococcales;Staphylococcaceae;Staphylococcus;Staphylococcus aureus | CP029671.97 | 100 | 100 | 0 | 497 |

FROGS Affiliation Filters: report.html

FROGS Affiliation Filters: impacted_OTU.multi-affiliations.tsv

FROGS Affiliation Filters: impacted_OTU.tsv

FROGS Affiliation Filters: sequences.fasta

FROGS Affiliation Filters: abundance.biom

*N.B.* The abundancy table (TSV format) of all deleted (or hidden according to the tool parameters) or modified OTUs are kept in **impacted_OTU.tsv**

| #comment | status | blast_taxonomy |
| --- | --- | --- |
| undesired_tax_in_blast | OTU_deleted | Bacteria;Firmicutes;Bacilli;Lactobacillales;Listeriaceae;Brochothrix;Brochothrix thermosphacta |
| undesired_tax_in_blast | OTU_deleted | Bacteria;Proteobacteria;Gammaproteobacteria;Enterobacterales;Vibrionaceae;Photobacterium;unknown species |
| undesired_tax_in_blast | OTU_deleted | Bacteria;Firmicutes;Bacilli;Lactobacillales;Lactobacillaceae;Latilactobacillus;Multi-affiliation |
| undesired_tax_in_blast | Blast_taxonomy_changed | Bacteria;Proteobacteria;Gammaproteobacteria;Pseudomonadales;Moraxellaceae;Psychrobacter;Multi-affiliation |
| blast_identity_lt_1.0;undesired_tax_in_blast | OTU_deleted | Bacteria;Firmicutes;Bacilli;Lactobacillales;Streptococcaceae;Lactococcus;Lactococcus piscium |
| blast_identity_lt_1.0;undesired_tax_in_blast | OTU_deleted | Bacteria;Firmicutes;Bacilli;Erysipelotrichales;Erysipelotrichaceae;ZOR0006;unknown species |
| undesired_tax_in_blast | OTU_deleted | Bacteria;Firmicutes;Bacilli;Lactobacillales;Streptococcaceae;Lactococcus;Multi-affiliation |
| blast_identity_lt_1.0;undesired_tax_in_blast | OTU_deleted | Bacteria;Firmicutes;Bacilli;Lactobacillales;Lactobacillaceae;Weissella;Weissella ceti |
| blast_identity_lt_1.0 | OTU_deleted | Bacteria;Bacteroidota;Bacteroidia;Flavobacteriales;Flavobacteriaceae;Flavobacterium;Flavobacterium sp. |
| blast_identity_lt_1.0 | OTU_deleted | Bacteria;Proteobacteria;Gammaproteobacteria;Enterobacterales;Vibrionaceae;Photobacterium;Photobacterium phosphoreum |
| blast_identity_lt_1.0;blast_coverage_lt_1.0;undesired_tax_in_blast | OTU_deleted | Bacteria;Firmicutes;Bacilli;Lactobacillales;Lactobacillaceae;Dellaglioa;Lactobacillus algidus |

In impacted_OTU.tsv
- #comment: the reason(s) why OTU was deleted (or hidden)
- #status: for deleted OTU (or hidden OTU), or for OTU with modified consensus taxonomy with affiliation (or multiaffiliation) was modified

**Q4: In hidding mode:** What outputs change between deleted mode and hidding mode ?

FROGS Affiliation Filters: report.html

FROGS Affiliation Filters: impacted_OTU.multi-affiliations.tsv

FROGS Affiliation Filters: impacted_OTU.tsv

FROGS Affiliation Filters: abundance.biom

In hidden mode: no **sequence.fasta** as output because none OTU was deleted

In hidden mode: **abundance.biom** contains all OTU but 111 have their affiliation that is hidden

| #comment | blast_taxonomy | blast_subjec | blast_perc_i | blast_perc_c | blast_evalue | blast_aln_le | seed_id | observation_ |
|---|---|---|---|---|---|---|---|---|
| undesired_tax_in_blast | no data | no data | no data | no data | no data | no data | 17_41 | Cluster_1 |
| undesired_tax_in_blast | no data | no data | no data | no data | no data | no data | 17_611 | Cluster_2 |
| undesired_tax_in_blast | no data | no data | no data | no data | no data | no data | 17_595 | Cluster_3 |
| undesired_tax_in_blast | Bacteria;Actinobacteriota;Actinobacteria;Propionibacteriales;Propionibacteriaceae;Cutibacterium;Multi-affiliation | multi-subjec | 100 | 100 | 0 | 468 | 17_257 | Cluster_4 |
| undesired_tax_in_blast | no data | no data | no data | no data | no data | no data | 17_4 | Cluster_5 |
| blast_identity_lt_1.0;undesired_tax_in_blast | no data | no data | no data | no data | no data | no data | 17_23 | Cluster_6 |
| blast_identity_lt_1.0;undesired_tax_in_blast | no data | no data | no data | no data | no data | no data | 57_5 | Cluster_7 |
| undesired_tax_in_blast | no data | no data | no data | no data | no data | no data | 17_420 | Cluster_8 |

« no data » appears in hidding mode

To see the content, think to transform the BIOM to TSV file with **BIOM_to_TSV tool**

# Normalization

SKIP →

# Normalization

Conserve a predefined number of sequence per sample:

- update Biom abundance file

- update seed fasta file

May be used when :

- Low sequencing sample

- Required for some statistical methods to compare the samples in pairs

# Exercise 8

1. What is the smallest sequenced samples ?

2. Normalize your data from Affiliation based on this number of sequence

3. Explore the report HTML result.

Q2: Normalize your data from Affiliation based on this number of sequence

**FROGS Abundance normalisation** Normalise OTUs abundance. (Galaxy Version 3.2.2)  ▼ Options

**Sequence file**

32: FROGS OTU Filters: sequences.fasta

Sequence file to normalise (format: fasta).

**Abundance file**

37: FROGS Affiliation OTU: affiliation.biom

Abundance file to normalise (format: BIOM).

**Number of reads**

7638

The final number of reads per sample.

✔ Execute

Q3: Explore the report HTML result.

Composition summary



*N.B.* if you normalize this datasets at 5000 or even 2000 sequences threshold, curiously you will not loose OTUs
But, **careful!** Generally, **more you normalize at low threshold, more you loose OTUs**

This reduction of data has not as consequence to loose OTUs

# FROGS Tree

CREATE A PHYLOGENETICS TREE OF OTUS

# FROGS Tree

This tool builds a phylogenetic tree thanks to affiliations of OTUs contained in the BIOM file

It uses MAFFT for the multiple alignment and FastTree for the phylogenetic tree.



2 outputs:

**FROGS Tree: report.html**

**FROGS Tree: tree.nwk**

OTUs

Abundance

Out of Tree : 0

Out of Tree : 0

In Tree : 495

In Tree : 547,518

# Tree View

Enabling zoom:

ON

Cluster_234 Bacteria Bacteroidota Bacteroidia Bacteroidales Prevotellaceae Alloprevotella ur

Cluster_325 Bacteria Bacteroidota Bacteroidia Bacteroidales Prevotellaceae Prevotella unkno

Cluster_351 Bacteria Bacteroidota Bacteroidia Bacteroidales Prevotellaceae Prevotella unkno

Cluster_356 Bacteria Bacteroidota Bacteroidia Bacteroidales Prevotellaceae Prevotella_7 unl

Cluster_251 Bacteria Bacteroidota Bacteroidia Bacteroidales Bacteroidaceae Bacteroides Ba

Cluster_330 Bacteria Bacteroidota Bacteroidia Bacteroidales Bacteroidaceae Bacteroides Pr

Cluster_131 Bacteria Bacteroidota Bacteroidia Bacteroidales Tannerellaceae Macellibacteroi

Cluster_450 Bacteria Bacteroidota Bacteroidia Bacteroidales Porphyromonadaceae Porphyro

Cluster_95 Bacteria Bacteroidota Bacteroidia Flavobacteriales Weeksellaceae Chryseobacte

Cluster_33 Bacteria Bacteroidota Bacteroidia Flavobacteriales Weeksellaceae Chryseobacte

Cluster_92 Bacteria Bacteroidota Bacteroidia Flavobacteriales Weeksellaceae Chryseobacte

Cluster_102 Bacteria Bacteroidota Bacteroidia Flavobacteriales Weeksellaceae Chryseobact

Cluster_101 Bacteria Bacteroidota Bacteroidia Flavobacteriales Weeksellaceae Chryseobact

The phylogentic tree in Newick format *i.e.* each mode is represented between brackets. This format is universal and can be used with all tree viewer



(dog:20, (elephant:30, horse:60):20):50

Our tree in nhx (= nwk) format

Exemple of visualization in FigTree from nhx file

```
(((((((((((((Cluster_234:0.25278,(Cluster_325:0.09784,Clu
67)0.972:0.02504,(Cluster_468:0.0269,(Cluster_138:0.0016
.782:0.00832,Cluster_277:0.01601)1.000:0.06764,Cluster_4
ter_47:0.13954,(Cluster_166:0.16129,(Cluster_403:0.22934
72:0.01332,(Cluster_400:0.00545,Cluster_473:0.01483)1.00
)0.829:0.01282,Cluster_240:0.12227)0.717:0.02027)0.981:0
uster_478:0.00249)0.000:0.00055,(Cluster_193:0.00055,Clu
359,Cluster_484:0.01913)0.880:0.03155)0.993:0.08088)0.45
0989)0.827:0.01144)0.870:0.01235,((Cluster_81:0.08926,Cl
05)0.862:0.00658,(Cluster_303:0.04337,Cluster_398:0.0311
237)0.953:0.01895,(Cluster_346:0.0235,((Cluster_369:0.01
Cluster_402:0.12402,(Cluster_309:0.02202,(Cluster_284:0.
.00054,(Cluster_427:0.00054,(Cluster_14:0.00402,Cluster_
0.791:0.02141,(Cluster_93:0.00054,Cluster_340:0.01463)0.
:0.03373)0.847:0.03692,Cluster_406:0.16125)0.831:0.03655
:0.04264)0.321:0.00907)0.487:0.01277,Cluster_129:0.06386
02802)0.763:0.02715,(Cluster_16:0.1183,(Cluster_63:0.062
```

# Practice:

# Exercice:

1. Create the phylogenetic tree that will be used for statistical analyses.



**FROGS Tree** Reconstruction of phylogenetic tree (Galaxy Version 3.2.2)

**OTUs sequence file**

14: FROGS OTU Filters sequences.fasta

OTUs sequence file (format: fasta). Warning: FROGS Tree does not work on more than 10000 sequences!

**Biom file**

19: FROGS Affiliation OTU: affiliation.biom

The abundance table of OTUs (format: biom).

✔ Execute

*For tutorial,* we ask you to create a phylogentic tree on affiliation.biom **before** "**affiliation filter**" process. Otherwise on your own data, create the phylogenetic tree on cleaned affiliation.biom

# Download your data

In order to share resources as well as possible, files that have not been accessed for more than 120 days are regularly purged. The backup of data generated using of Galaxy is <u>your responsibility</u>.

You have 2 backup possibilities:
1. Save your datasets one by one using the "floppy disk" icon.

2. Or export each history.
To export a history, from the "History" menu, click on the wheel, then "Export History to File":

**55: FROGS Affiliation OTU:**
**excluded_data_report.html**
11.4 KB
format: html, database: ?
## Application Software:
affiliation_OTU.py (version: 0.4.0)
Command: /usr/local/bioinfo
/src/galaxy-test/galaxy-dist/tools
/FROGS/affiliation_OTU.py
--reference /save/galaxy-
test/bank/FROGS/silva_119-1
/prokaryotes
/silva_119-1_prokaryotes.fasta
--abundance

HTML file

**istory**
HISTORY LISTS
Saved Histories
Histories Shared with Me
HISTORY ACTIONS
Create New
Copy History
Share or Publish
Show Structure
Extract Workflow
Delete
Delete Permanently
DATASET ACTIONS
Copy Datasets
Dataset Security
Resume Paused Jobs
Collapse Expanded Datasets
Unhide Hidden Datasets
Delete Hidden Datasets
Purge Deleted Datasets
DOWNLOADS
Export Tool Citations
Export History to File
OTHER ACTIONS
Import from File

To retrieve your history, click on the http link that appears automatically:

It is then possible to record the data :

Ouverture de Galaxy-History-historique-multiplex.tar.gz

Vous avez choisi d'ouvrir :

📦 Galaxy-History-historique-multiplex.tar.gz

qui est un fichier de type : archive gzip

à partir de : http://sigenae-workbench.toulouse.inra.fr

Que doit faire Firefox avec ce fichier ?

○ Ouvrir avec    Gestionnaire d'archives (défaut) ▾

○ Enregistrer le fichier

☐ Toujours effectuer cette action pour ce type de fichier.

Annuler    OK

This directory contains :

📁 datasets

📄 datasets_attrs.txt

📄 datasets_attrs.txt.provenance

📄 history_attrs.txt

📄 jobs_attrs.txt

1. in the "datasets" directory: Your Galaxy files.
2. in the files "-attrs.txt" : Metadata about your datasets, your jobs and your history.

# How to cite FROGS

Frédéric Escudié, Lucas Auer, Maria Bernard, Mahendra Mariadassou, Laurent Cauquil, Katia Vidal, Sarah Maman, Guillermina Hernandez-Raquet, Sylvie Combes, Géraldine Pascal.

"**FROGS**: Find, Rapidly, OTUs with Galaxy Solution." *Bioinformatics,* , Volume 34, Issue 8, 15 April 2018, Pages 1287–1294

# FROGS'docs

Website: http://frogs.toulouse.inrae.fr

Tuto: https://youtu.be/Kh6ZrlmKGoY







Pipeline FROGS on

http://sigenae-workbench.toulouse.inra.fr/galaxy/u/gpascal/w/to-test-frogs

All scripts on Github: https://github.com/geraldinepascal/FROGS.git

# To contact

FROGS:

frogs-support@inrae.fr


Galaxy:

support.sigenae@inrae.fr


Newsletter – subscription request:

frogs@inrae.fr

# Play list FROGS:

https://www.deezer.com/fr/playlist/5233843102?utm_source=deezer&utm_content=playlist-5233843102&utm_term=18632989_1545296531&utm_medium=web