

# Analyse de données métagénomiques 16S - FROGS

*Module 20*

Olivier Rué 

MaIAGE - Migale

September 12, 2023



This work is licensed under a [Creative Commons Attribution-ShareAlike 2.0 Generic License](https://creativecommons.org/licenses/by-sa/2.0/)

# FROGS



This work is licensed under a [Creative Commons Attribution-ShareAlike 2.0 Generic License](https://creativecommons.org/licenses/by-sa/2.0/)

# FROGS team



- FROGS is a INRAE development project



Vincent DARBOT



Maria BERNARD



Olivier RUÉ



Lucas AUER



Laurent CAUQUIL



Patrice DÉHAIS

Developers

Biology experts

Galaxy  
support



Mahendra  
MARIADASSOU

Statistical expert



Géraldine  
PASCAL

Coordinator



This work is licensed under a [Creative Commons Attribution-ShareAlike 2.0 Generic License](https://creativecommons.org/licenses/by-sa/2.0/)

# FROGS articles

Bioinformatics, 34(8), 2018, 1287–1294  
doi:10.1093/bioinformatics/btz791  
Advance Access Publication Date: 7 December 2017  
Original Paper

OXFORD

Sequence analysis

## FROGS: Find, Rapidly, OTUs with Galaxy Solution

Frédéric Escudié<sup>1,\*</sup>, Lucas Auer<sup>2,\*</sup>, Maria Bernard<sup>3</sup>,  
Mahendra Mariadassou<sup>4</sup>, Laurent Cauquil<sup>5</sup>, Katia Vidal<sup>5</sup>, Sarah Maman<sup>5</sup>,  
Guillermina Hernandez-Raquet<sup>6</sup>, Sylvie Combes<sup>5</sup> and  
Géraldine Pascal<sup>5,\*</sup>

<sup>1</sup>Bioinformatics platform Toulouse Midi-Pyrenees, MIAT, INRA Auzeville CS 52627 31326 Castanet Tolosan cedex, France, <sup>2</sup>INRA, UMR 1136, Université de Lorraine, INRA-Nancy, 54280, Champenoux, France, <sup>3</sup>GABI, INRA, AgroParisTech, Université Paris-Saclay, Jouy-en-Josas, France, <sup>4</sup>MaIAGE, INRA, Université Paris-Saclay, 78350 Jouy-en-Josas, France, <sup>5</sup>GenPhySE, Université de Toulouse, INRA, INPT, ENVT, Castanet Tolosan, France and <sup>6</sup>Laboratoire d'ingénierie des Systèmes Biologiques et des Procédés-LISBP, Université de Toulouse, INSA, INRA, CNRS, Toulouse, France

\*To whom correspondence should be addressed.

The authors wish it to be known that, in their opinion, the first two authors should be regarded as Joint First Authors. Associate Editor: Bonnie Berger

Received on May 10, 2017; revised on December 1, 2017; editorial decision on December 4, 2017; accepted on December 5, 2017

### Abstract

**Motivation:** Metagenomics leads to major advances in microbial ecology and biologists need user friendly tools to analyze their data on their own.

**Results:** This Galaxy-supported pipeline, called FROGS, is designed to analyze large sets of amplicon sequences and produce abundance tables of Operational Taxonomic Units (OTUs) and their taxonomic affiliation. The clustering uses Swarm. The chimera removal uses VSEARCH, combined with original cross-sample validation. The taxonomic affiliation returns an innovative multi-affiliation output to highlight databases conflicts and uncertainties. Statistical results and numerical graphical illustrations are produced along the way to monitor the pipeline. FROGS was tested for the detection and quantification of OTUs on real and *in silico* datasets and proved to be rapid, robust and highly sensitive. It compares favorably with the widespread mothur, UPARSE and QIIME.

**Availability and implementation:** Source code and instructions for installation: <https://github.com/geraldinepascal/FROGS.git>. A companion website: <http://frogs.toulouse.inra.fr>.  
**Contact:** geraldine.pascal@inra.fr

**Supplementary information:** Supplementary data are available at *Bioinformatics* online.

### 1 Introduction

The expansion of high-throughput sequencing of rRNA amplicons has opened new horizons for the study of microbial communities. By making it possible to study all micro-organisms from a given environment without the need to cultivate them, metagenomics has led to major advances in many fields of microbial ecology, from the study of the impact of microbiota on human and animal pathologies

(Hess *et al.*, 2011; Hooper *et al.*, 2012; Jovel *et al.*, 2016) to the study of biodiversity in environmental ecosystems and the search for biomarkers of pollution (Andres and Bertin, 2016; de Vargas *et al.*, 2015). Determining the composition of a microbial ecosystem, at low cost and great depth, is still largely based on the amplification and sequencing of biodiversity marker genes, also called amplicons, such as rRNA genes and ITS. The clustering of sequences into

© The Author 2017. Published by Oxford University Press. All rights reserved. For Permissions, please e-mail: [journals.permissions@oup.com](mailto:journals.permissions@oup.com)

1287

Downloaded from <https://academic.oup.com/bioinformatics/advance-article-abstract/doi/10.1093/bioinformatics/btz791/7170232> by guest on 28 August 2023

OXFORD

Briefings in Bioinformatics, 22(6), 2021, 1–6

<https://doi.org/10.1093/bib/bbab118>  
Problem Solving Protocol

## FROGS: a powerful tool to analyse the diversity of fungi with special management of internal transcribed spacers

Maria Bernard<sup>1</sup>, Olivier Rué<sup>1</sup>, Mahendra Mariadassou<sup>2</sup> and  
Géraldine Pascal<sup>3</sup>

Corresponding author: Géraldine Pascal, GenPhySE, Université de Toulouse, INRAE, ENVT, F-31326, Castanet Tolosan, France. Tel: +33 (0)5 61 28 51 05.  
E-mail: [geraldine.pascal@inrae.fr](mailto:geraldine.pascal@inrae.fr)  
<sup>1</sup>Maria Bernard and Olivier Rué are joint first authors.

### Abstract

Fungi are present in all environments. They fulfil important ecological functions and play a crucial role in the food industry. Their accurate characterization is thus indispensable, particularly through metabarcoding. The most frequently used markers to monitor fungi are ITSs. These markers are the best documented in public databases but have one main weakness: polymerase chain reaction amplification may produce non-overlapping reads in a significant fraction of the fungi. When these reads are filtered out, traditional metabarcoding pipelines lose part of the information and consequently produce biased pictures of the composition and structure of the environment under study. We developed a solution that enables processing of the entire set of reads including both overlapping and non-overlapping, thus providing a more accurate picture of fungal communities. Our comparative tests using simulated and real data demonstrated the effectiveness of our solution, which can be used by both experts and non-specialists on a command line or through the Galaxy-based web interface.

**Key words:** fungi; ITS; metabarcoding; workflow; amplicon; metagenomics

### Introduction

Using amplicon sequencing to describe the microbial composition of an environment is a time saving and cost-effective strategy and can be used even for very large-scale surveys [1]. Most studies currently focus on the bacterial fraction of microbial communities but the fungal fraction is equally important, as fungi are ubiquitous and provide several ecosystem services [2]. Unfortunately, studying the fungal fraction using metabarcoding has its own challenges. Indeed, in fungi, there is no equivalent of the 16S rRNA gene, which is widely used and highly suitable

for bacteria. The best candidates are internal transcribed spacers (ITS), but these are more difficult to manipulate. The main problem with ITS is size polymorphism, with a size range of 361–1475 bases in UNITE 7.1 [3] (unlike 16S where 95% of the sequences have a length between 1205 and 1556 bases). Most studies describing ITS data analyses process either (i) paired-end reads but filter out non-overlapping, non-mergeable reads, thus systematically discarding taxa with longer ITS, or (ii) single-end reads, thus limiting taxonomic resolution and losing the benefit of information contained in longer sequences [4, 5].

Maria Bernard is a bioinformatics engineer. She is a member of a platform team conducting NGS sequence analysis and designing software. She specializes in workflow development in particular for metabarcoding analysis.

Olivier Rué is a bioinformatics engineer. He is in charge of data analysis at the Migale bioinformatics facility. He specializes in the analysis of metabarcoding and metagenomics data.

Mahendra Mariadassou has a PhD in statistics. He is involved in the development of new statistical methods and tools for metabarcoding analysis.

Géraldine Pascal has a PhD in bioinformatics and coordinates the FROGS project. She is currently involved in designing solutions for long read problems, workflow development and metagenomics analysis.

Submitted: 19 April 2021. Received (in revised form): 19 July 2021.

© The Author(s) 2021. Published by Oxford University Press. All rights reserved. For Permissions, please email: [journals.permissions@oup.com](mailto:journals.permissions@oup.com)

1

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



This work is licensed under a [Creative Commons Attribution-ShareAlike 2.0 Generic License](https://creativecommons.org/licenses/by-sa/2.0/)



RÉPUBLIQUE  
FRANÇAISE  
Liberté  
Égalité  
Fraternité

INRAE

mis@le

# How to use FROGS



- Command line `</>`

```
remove_chimera.py
--input-biom clustering.biom \
--input-fasta clustering.fasta \
--non-chimera remove_chimera.fasta \
--out-abundance remove_chimera.biom \
--summary remove_chimera.html
```

- Galaxy instances via web 

**FROGS\_3 Remove chimera** Remove PCR chimera in each sample (Galaxy Version 4.1.0+galaxy1)

**Sequences file (format: FASTA)**

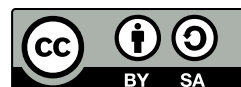
The sequences file

**Abundance type**

Select the type of file where the abundance of each sequence by sample is stored.

**Abundance file (format: BIOM)**

It contains the count by sample for each sequence.



- Migale - <https://galaxy.migale.inrae.fr/>

- Genotoul - <http://www.genotoul.fr/>

- IFB - <https://metabarcoding.usegalaxy.fr/>

This work is licensed under a [Creative Commons Attribution-ShareAlike 2.0 Generic License](https://creativecommons.org/licenses/by-sa/2.0/)

# FROGS docs and help



- 🌐 Website: <https://frogs.toulouse.inrae.fr>
- 🐙 Github: <https://github.com/geraldinepascal/FROGS.git>
- 📧 Newsletter: subscription request at [frogs-support@inrae.fr](mailto:frogs-support@inrae.fr)
- ? Need help
  - [frogs-support@inrae.fr](mailto:frogs-support@inrae.fr) for generic questions
  - [help-migale@inrae.fr](mailto:help-migale@inrae.fr) for bugs/quotas/errors with Galaxy Migale instance



This work is licensed under a [Creative Commons Attribution-ShareAlike 2.0 Generic License](https://creativecommons.org/licenses/by-sa/2.0/)

# TP1: Introduction to Galaxy



This work is licensed under a [Creative Commons Attribution-ShareAlike 2.0 Generic License](https://creativecommons.org/licenses/by-sa/2.0/)



RÉPUBLIQUE  
FRANÇAISE  
*Liberté  
Égalité  
Fraternité*

INRAE **mission**

# Sequencing data



This work is licensed under a [Creative Commons Attribution-ShareAlike 2.0 Generic License](https://creativecommons.org/licenses/by-sa/2.0/)



# FASTQ format

```
@ST-E00114:1342:HHMGVCCX2:1:1101:3123:2012 1:N:0:TCCGGAGA+TCAGAGCC
CTTGGTCATTTAGAG
```

+

```
***<<*AEF???* **
```

```
@ST-E00114:1342:HHMGVCCX2:1:1101:11556:2030 1:N:0:TCCGGAGA+TCAGAGCC
CATTGGCCATATCAT
```

+

```
AAAE??<<*???* **
```

## Meaning

```
@Identif1er1 (comment)
```

```
XXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXX
```

+

```
QQQQQQQQQQQQQQQQQQQQQQQQQQQQQQQQQQQQQQQQQQQQQQQQQQQQQQQQQQQQ
```

```
@Identif1er2 (comment)
```

```
XXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXX
```

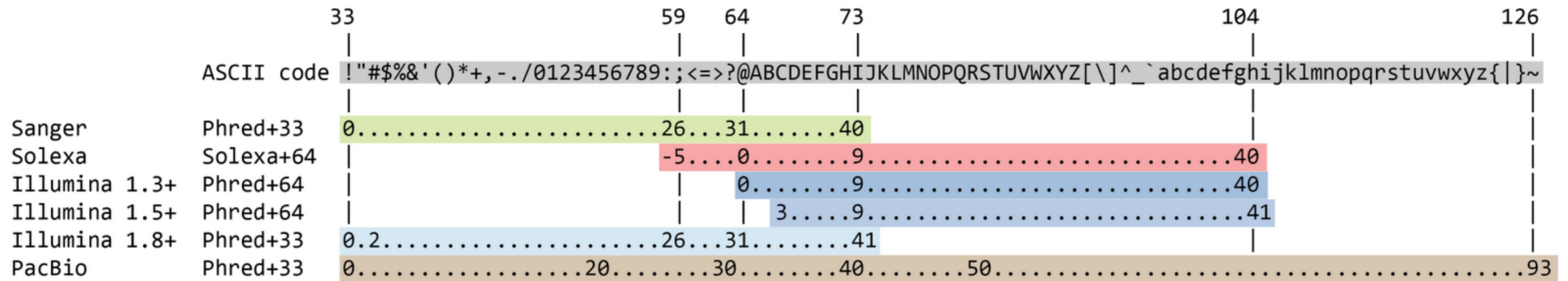
+

```
QQQQQQQQQQQQQQQQQQQQQQQQQQQQQQQQQQQQQQQQQQQQQQQQQQQQQQQQQQQQ
```



This work is licensed under a [Creative Commons Attribution-ShareAlike 2.0 Generic License](https://creativecommons.org/licenses/by-sa/2.0/)

# Quality score encoding



This work is licensed under a [Creative Commons Attribution-ShareAlike 2.0 Generic License](https://creativecommons.org/licenses/by-sa/2.0/)



# Quality score

Measure of the quality of the identification of the nucleobases generated by automated DNA sequencing

Phred Quality Score	Probability of incorrect base call	Base call accuracy
10	1 in 10	90%
20	1 in 100	99%
30	1 in 1000	99.9%
40	1 in 10,000	99.99%
50	1 in 100,000	99.999%
60	1 in 1,000,000	99.9999%



This work is licensed under a [Creative Commons Attribution-ShareAlike 2.0 Generic License](https://creativecommons.org/licenses/by-sa/2.0/)

# FASTQ compression

- Compression is essential to deal with FASTQ files (reduce disk storage)
- *extension: file.fastq.gz*
- Tools are (almost all) able to deal with compressed files 👍



This work is licensed under a [Creative Commons Attribution-ShareAlike 2.0 Generic License](https://creativecommons.org/licenses/by-sa/2.0/)

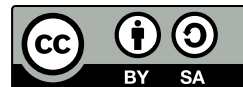
# Quality control



This work is licensed under a [Creative Commons Attribution-ShareAlike 2.0 Generic License](https://creativecommons.org/licenses/by-sa/2.0/)

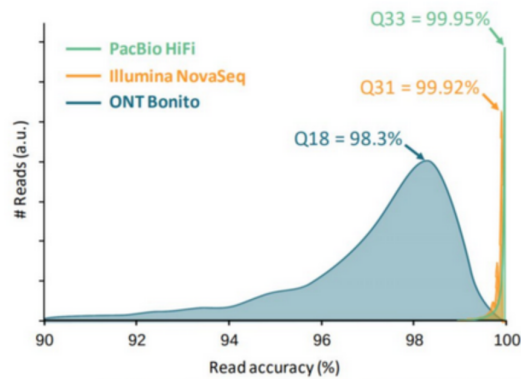
# Quality control

- One of the most easy step in bioinformatics ...
- ... but one of the most important
- check if everything is ok
- Indicates if/how to clean reads
- Shows possible sequencing problems
- The results must be interpreted in relation to what has been sequenced



This work is licensed under a [Creative Commons Attribution-ShareAlike 2.0 Generic License](https://creativecommons.org/licenses/by-sa/2.0/)

# Reads are not perfect



PacBio HiFi: HG003 18 kb library, Sequel II System Chemistry 2.0, [precisionFDA Truth Challenge V2](#)  
 Illumina: HG002 2x150 bp NovaSeq library, [precisionFDA Truth Challenge V2](#)  
 ONT: Bonito BCM Nanopore Tech, Lupton Dec. 2020 and Bonito Basecalling with P3.4.1

Step	Sample	DNA extraction	Fragmentation	PCR amplification	Flow cell hybridization	Cluster generation	Sequencing by Synthesis	
Error source	Mutagenesis	Oxidative damage	Oxidative damage	Polymerase mistakes		Cluster PCR errors	Phasing	Fluorophore crosstalk
Error category	Biological variants	DNA damage	DNA damage	PCR errors		Sequencing errors	Sequencing errors	
Detection						← Ignored	Called →	



This work is licensed under a [Creative Commons Attribution-ShareAlike 2.0 Generic License](#)

Sequencing error profiles of Illumina sequencing instruments [3]

# Why QC'ing your reads?

Try to answer to (not always) simple questions:

- Are data conform to the expected level of performance?
  - Size / Number of reads / Quality
- Residual presence of adapters or indexes?
- (Un)expected technical biases?
- (Un)expected biological biases?

 Warning

QC without context leads to misinterpretation!





# TP2: Quality control



This work is licensed under a [Creative Commons Attribution-ShareAlike 2.0 Generic License](https://creativecommons.org/licenses/by-sa/2.0/)

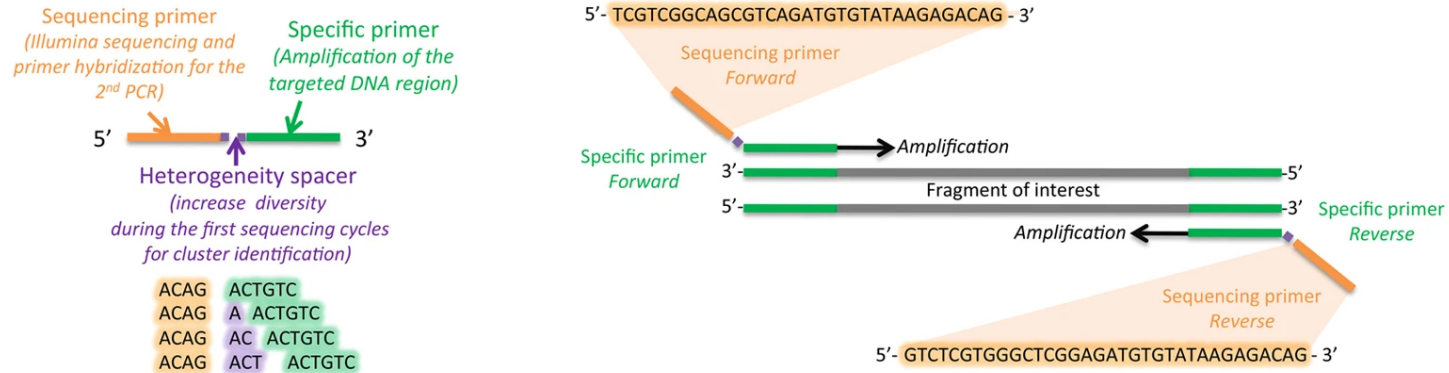
# Demultiplexing



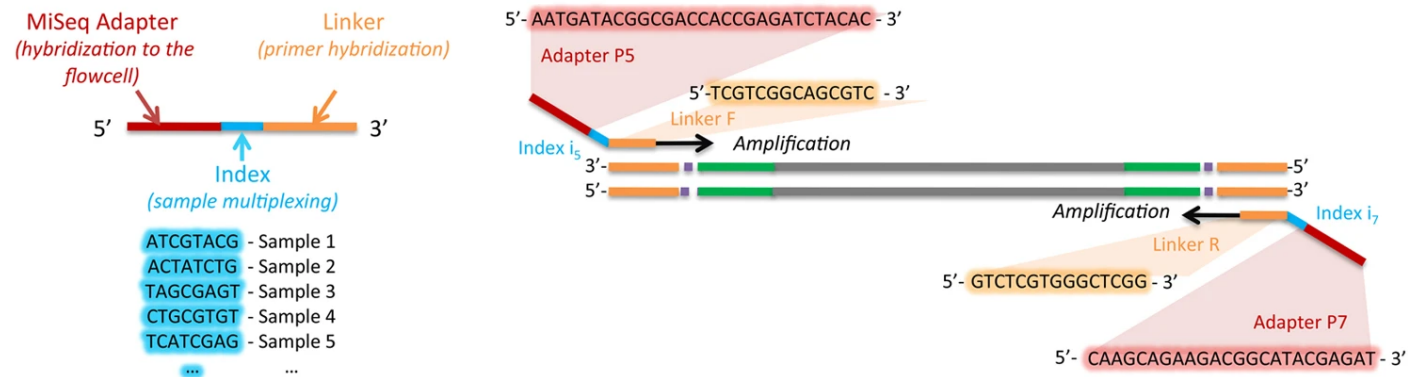
This work is licensed under a [Creative Commons Attribution-ShareAlike 2.0 Generic License](https://creativecommons.org/licenses/by-sa/2.0/)

# Multiplexing principle

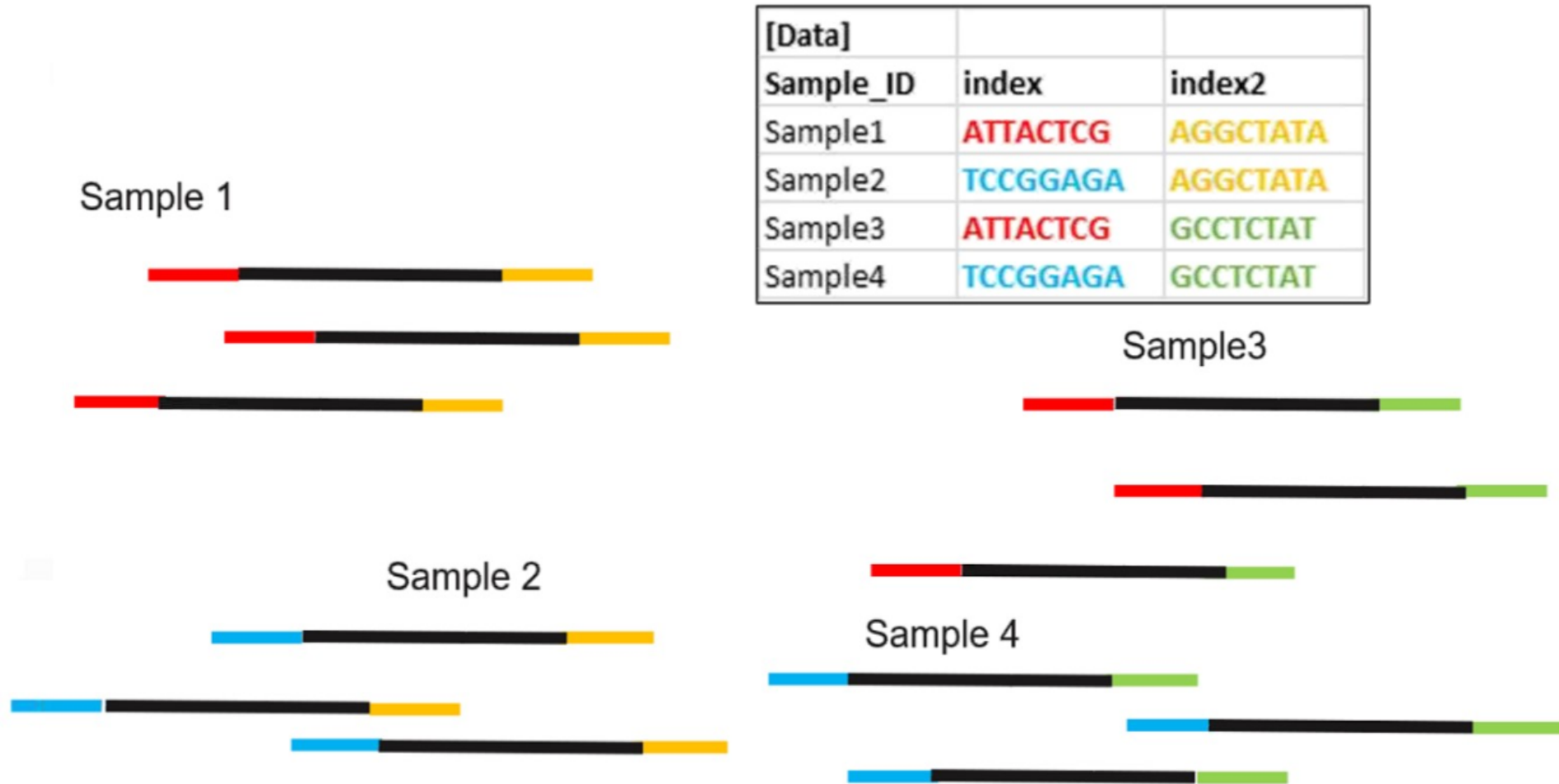
## Step 1 : Amplification of selected fragments



## Step 2 : Add Illumina adaptor sequences and indexes to multiplex samples



# Demultiplexing by bioinformatics



This work is licensed under a [Creative Commons Attribution-ShareAlike 2.0 Generic License](https://creativecommons.org/licenses/by-sa/2.0/)

Credits: [https://biocellgen-public.svi.edu.au/mig\\_2019\\_scrnaseq-workshop/processing-raw-scrna-seq-data.html](https://biocellgen-public.svi.edu.au/mig_2019_scrnaseq-workshop/processing-raw-scrna-seq-data.html)

# Demultiplexing tool

- Assign each read to FASTQ files depending on barcode found
- BARCODE FILE is expected to be tabular:
  - 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)



# TP DEMULTIPLEXING



This work is licensed under a [Creative Commons Attribution-ShareAlike 2.0 Generic License](https://creativecommons.org/licenses/by-sa/2.0/)

# Advices

- 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.



This work is licensed under a [Creative Commons Attribution-ShareAlike 2.0 Generic License](https://creativecommons.org/licenses/by-sa/2.0/)



RÉPUBLIQUE  
FRANÇAISE  
*Liberté  
Égalité  
Fraternité*

INRAE **mis@le**

# FROGS preprocess



This work is licensed under a [Creative Commons Attribution-ShareAlike 2.0 Generic License](https://creativecommons.org/licenses/by-sa/2.0/)



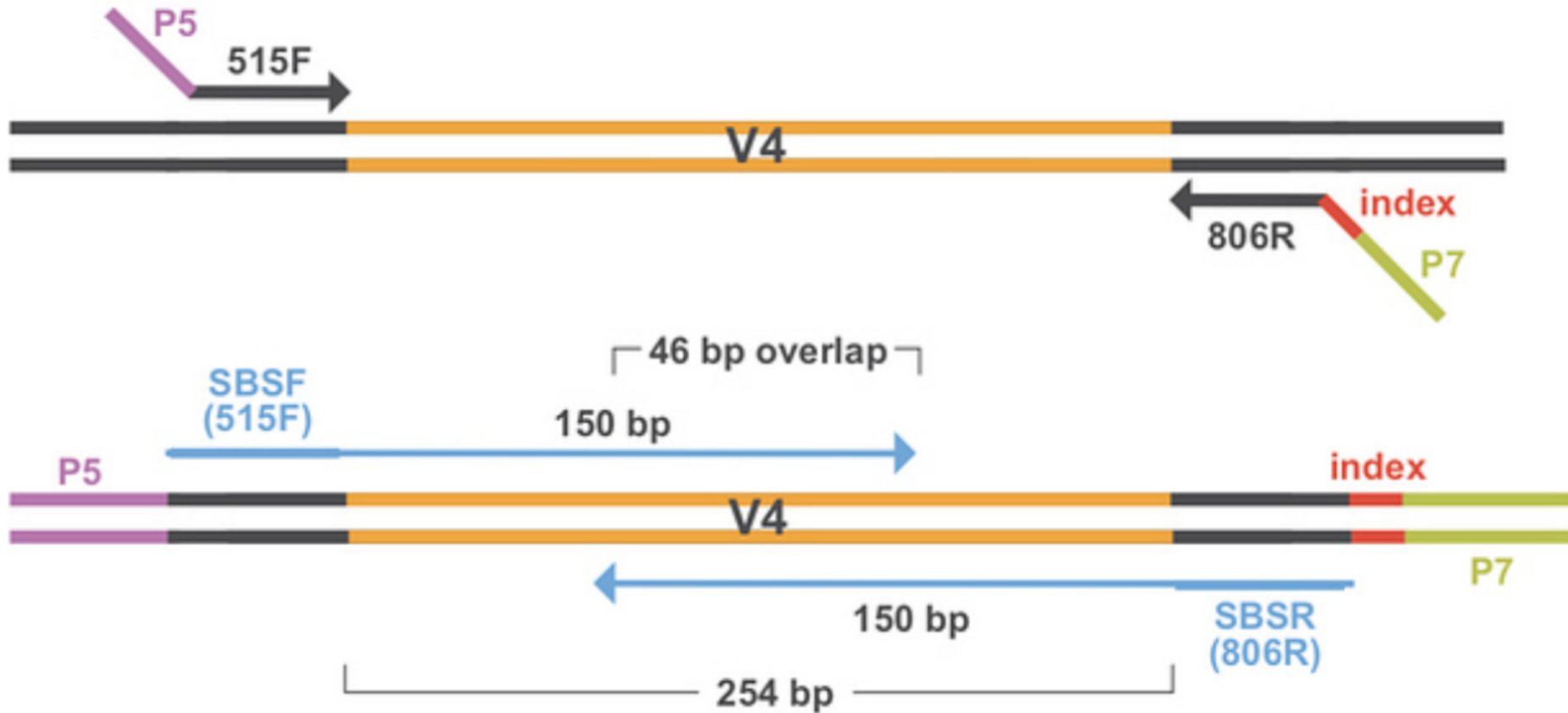
# What FROGS preprocess does?

- Merging of R1 and R2 reads with vsearch [4], flash [5] or pear [6] (only in command line)
- Deletes sequences without good primers
- Finds and removes adapter sequences with cutadapt
- Deletes sequence with not expected lengths
- Deletes sequences with ambiguous bases (N)
- Dereplication
- *removing homopolymers (size = 8) for 454 data*
- *quality filter for 454 data*



This work is licensed under a [Creative Commons Attribution-ShareAlike 2.0 Generic License](https://creativecommons.org/licenses/by-sa/2.0/)

# Merging of paired-end reads



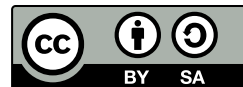
This work is licensed under a [Creative Commons Attribution-ShareAlike 2.0 Generic License](https://creativecommons.org/licenses/by-sa/2.0/)

# TP FROGS preprocess



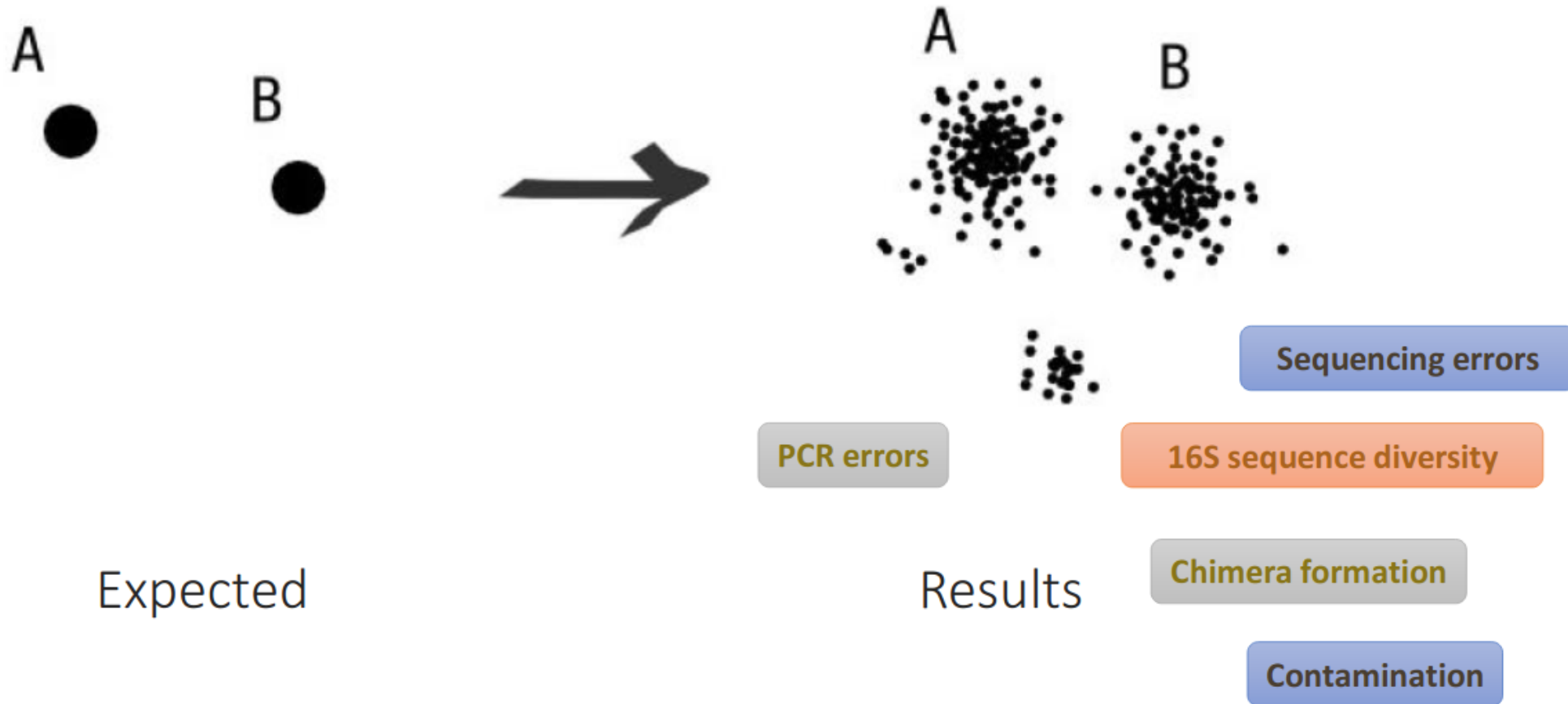
This work is licensed under a [Creative Commons Attribution-ShareAlike 2.0 Generic License](https://creativecommons.org/licenses/by-sa/2.0/)

# Clustering



This work is licensed under a [Creative Commons Attribution-ShareAlike 2.0 Generic License](https://creativecommons.org/licenses/by-sa/2.0/)

# Sequencing data are noised



This work is licensed under a [Creative Commons Attribution-ShareAlike 2.0 Generic License](https://creativecommons.org/licenses/by-sa/2.0/)

# How to deal with these noised sequences?

- Comparison all against all
  - Very accurate
  - Requires a lot of memory and/or time
- Clustering
  - closed-reference / open-reference
  - de novo
- Denoising



This work is licensed under a [Creative Commons Attribution-ShareAlike 2.0 Generic License](https://creativecommons.org/licenses/by-sa/2.0/)

# Vocabulary

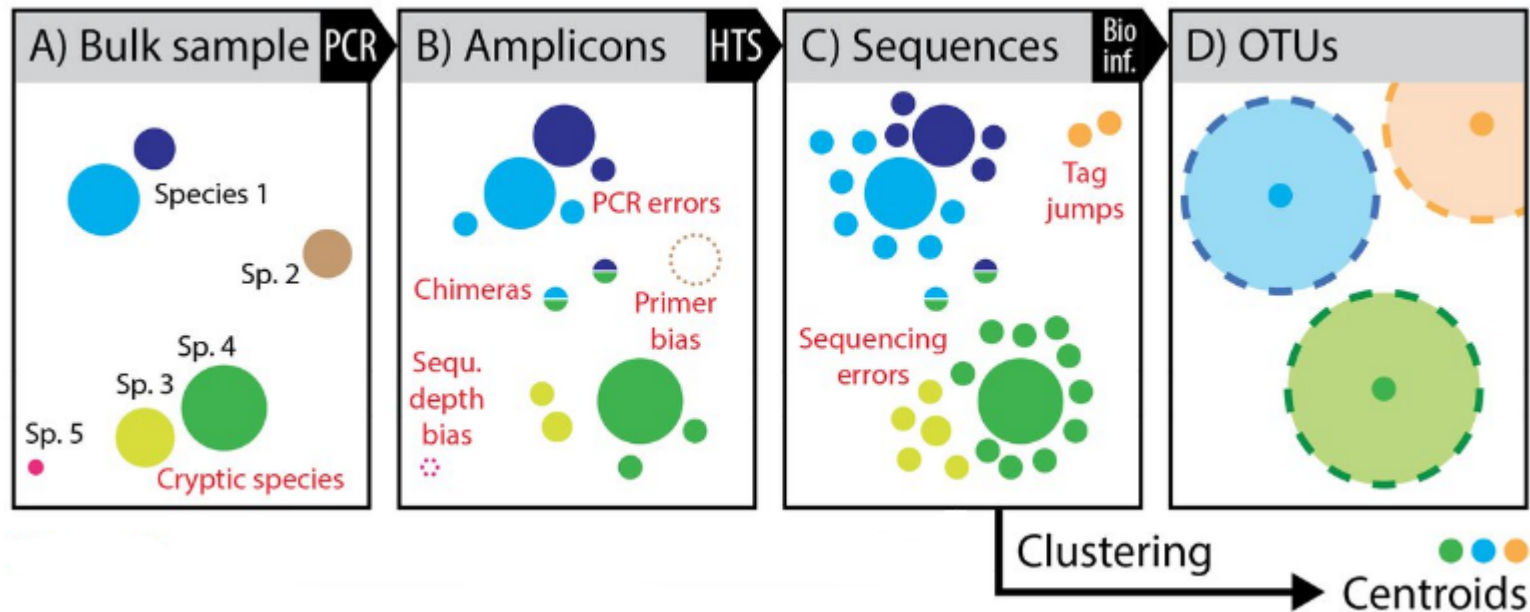
- A lot of terms for features built by softwares
  - OTUs, zOTUS, ASVs, ESVs...
- A recent review establishes the vocabulary [7]
  - OTUs / ASVs / swarm clusters
- *ASVs are identical denoised reads with as few as 1 base pair difference between variants, representing an inference of the biological sequences prior to amplification and sequencing errors*
- OTUs are formed with a % threshold clustering
- Swarm clusters are a third feature type



This work is licensed under a [Creative Commons Attribution-ShareAlike 2.0 Generic License](https://creativecommons.org/licenses/by-sa/2.0/)

# OTU paradigm

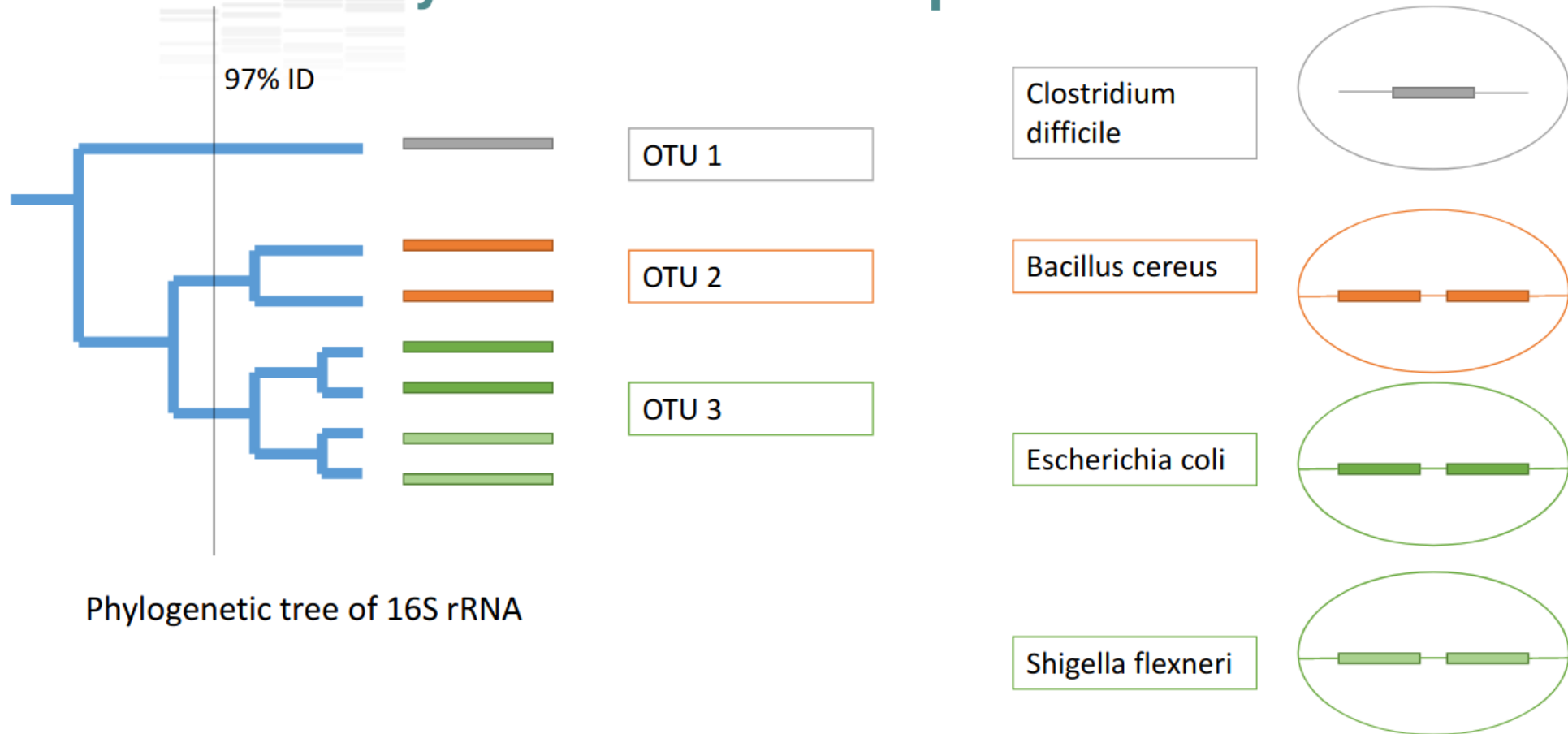
- Operational Taxonomic Unit





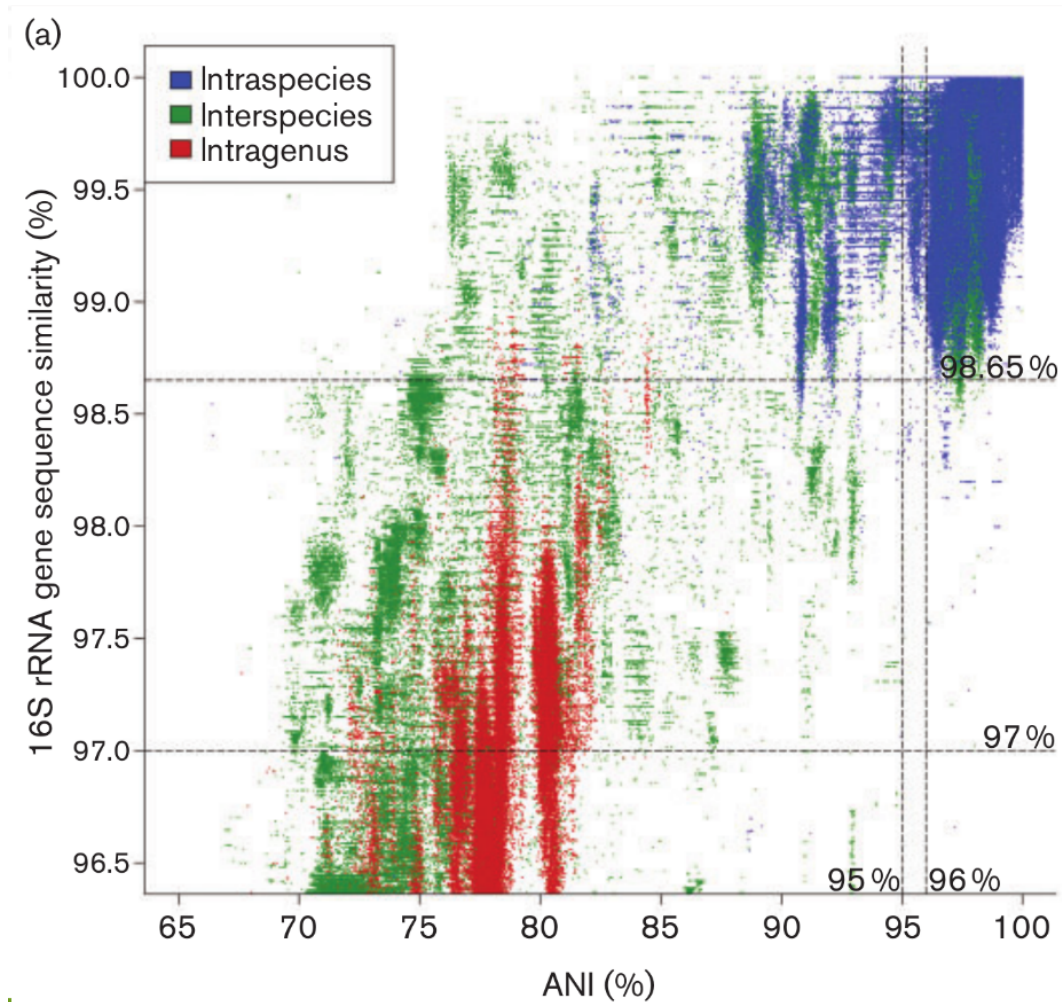
# Operational Taxonomic Units

## OTUs: a Proxy for « Bacterial Species »



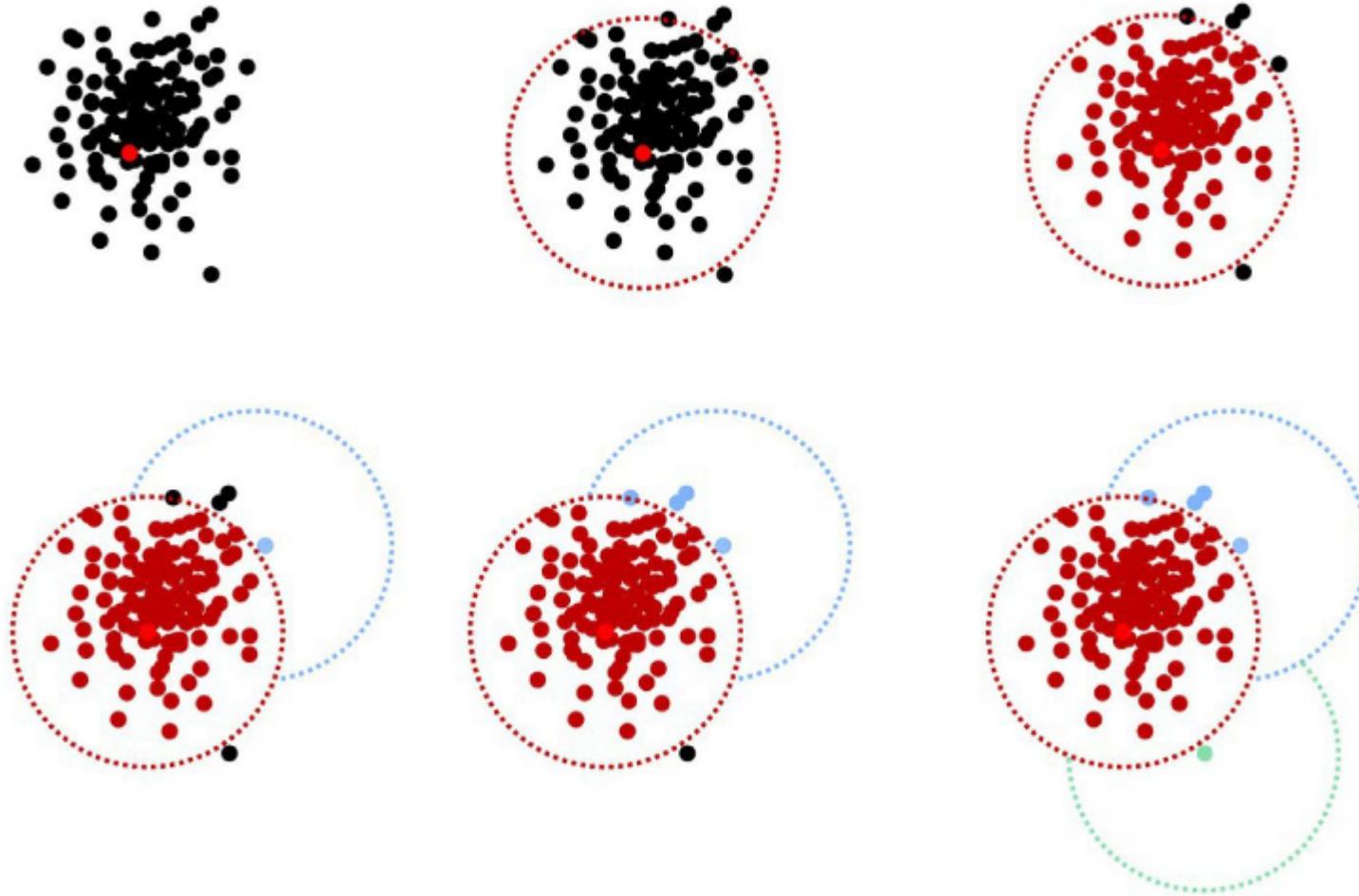
This work is licensed under a [Creative Commons Attribution-ShareAlike 2.0 Generic License](https://creativecommons.org/licenses/by-sa/2.0/)

# Operational Taxonomic Units



This work is licensed under a [Creative Commons Attribution-ShareAlike 2.0 Generic License](https://creativecommons.org/licenses/by-sa/2.0/)

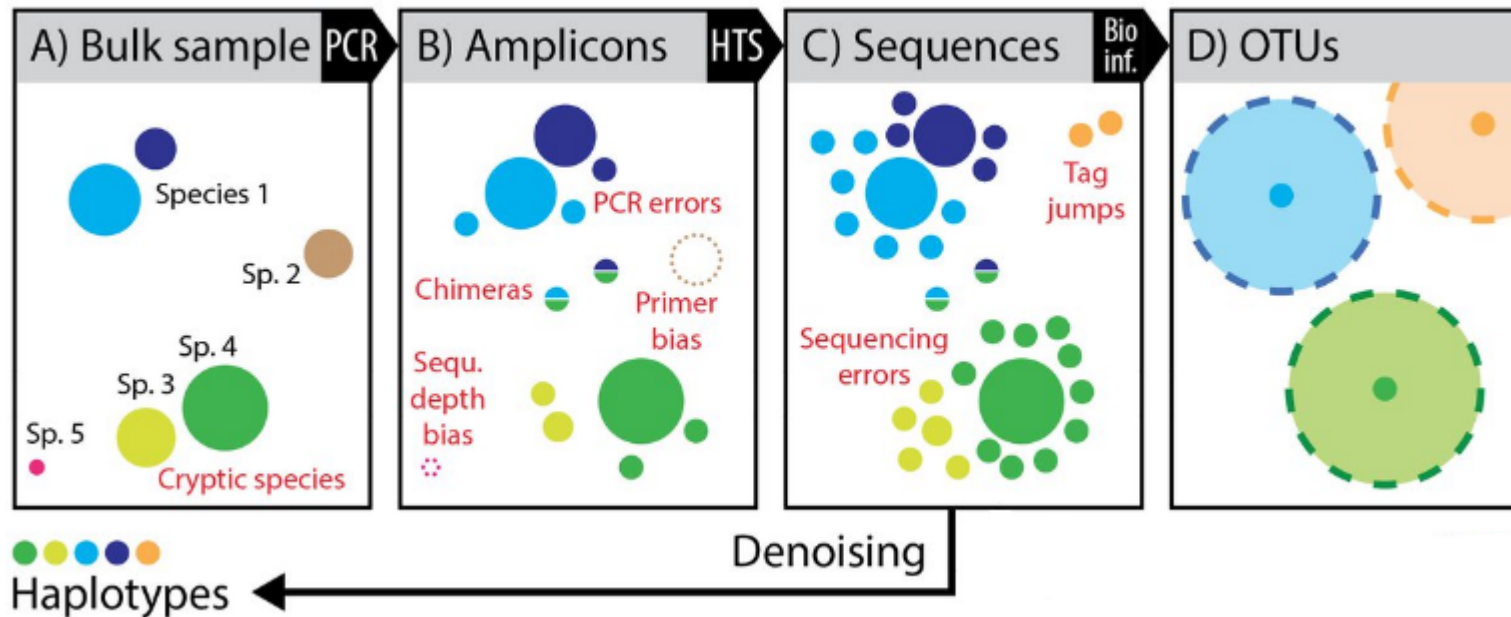
# Operational Taxonomic Units



This work is licensed under a [Creative Commons Attribution-ShareAlike 2.0 Generic License](https://creativecommons.org/licenses/by-sa/2.0/)

# ASV paradigm

- Amplicon Sequence Variants

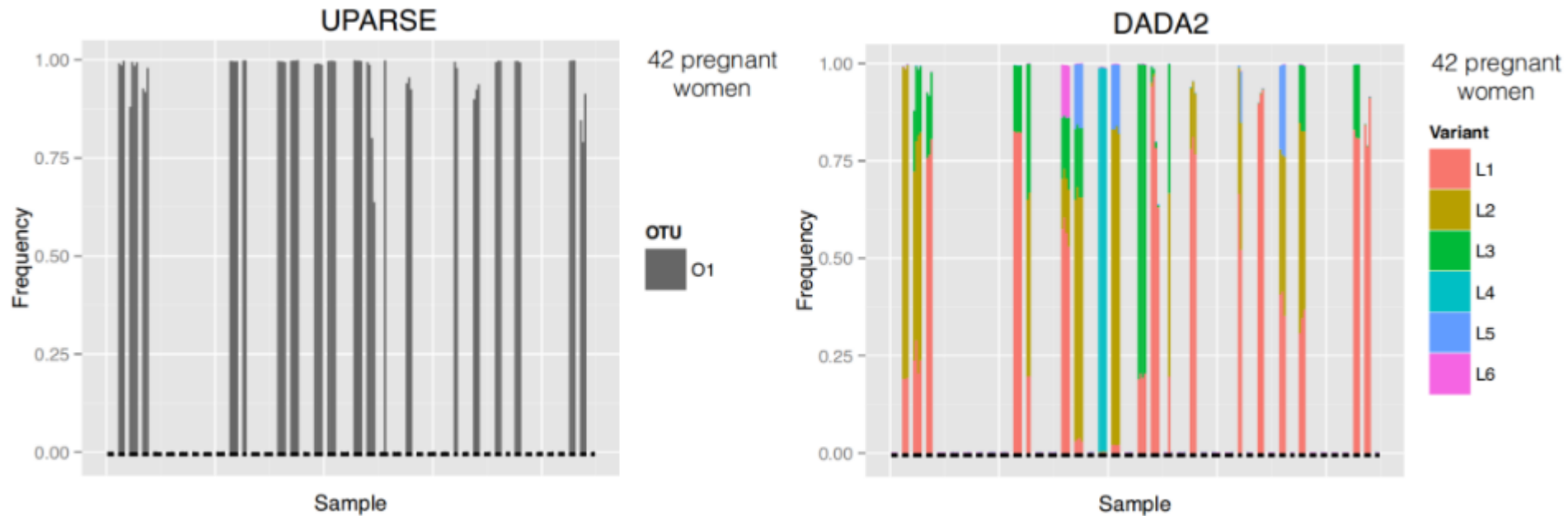


ASV are inferred by a de novo process in which biological sequences are discriminated from errors on the basis of the expectation that biological sequences are more likely to be repeatedly observed than are error-containing sequences



This work is licensed under a [Creative Commons Attribution-ShareAlike 2.0 Generic License](https://creativecommons.org/licenses/by-sa/2.0/)

# ASV resolution



- ASV resolution changes the composition for these samples



This work is licensed under a [Creative Commons Attribution-ShareAlike 2.0 Generic License](https://creativecommons.org/licenses/by-sa/2.0/)

Credits: [https://benjjneb.github.io/dada2/SMBS\\_DADA2.pdf](https://benjjneb.github.io/dada2/SMBS_DADA2.pdf)

# Swarm

Swarm [8] is a notably different sequence clustering approach, which, while technically a clustering algorithm, may also be considered a denoising method when using the fastidious method with  $d=1$ . It relies on the maximum number of differences between reads (local linking threshold) and forms clusters that are resilient to input-order changes, thus creating stable, high-resolution features (herein referred to as swarm-clusters). When using the fastidious method with  $d=1$ , swarm aims to produce clusters centered around real biological sequences, where clusters represent sequence variants.

Since FROGS uses swarm (with the fastidious method with  $d=1$ ) and strongly promotes denoising by chimera removal and cluster filtering, FROGS produces ASVs.



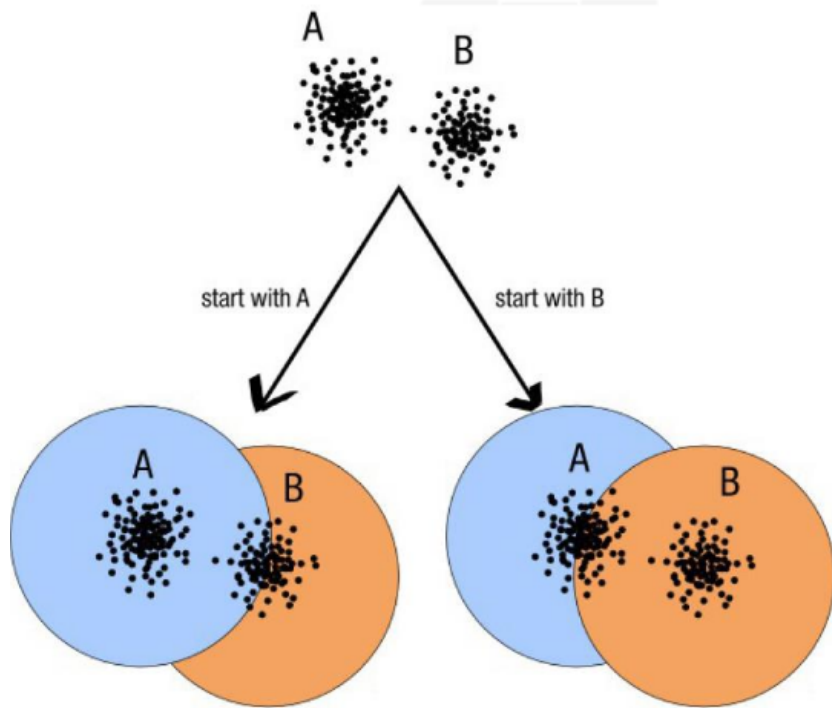
This work is licensed under a [Creative Commons Attribution-ShareAlike 2.0 Generic License](https://creativecommons.org/licenses/by-sa/2.0/)

# Why Swarm?

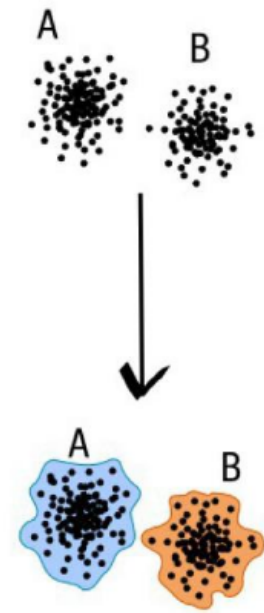
- Fixed clustering threshold is a real problem
- OTUs construction is input-order dependent



This work is licensed under a [Creative Commons Attribution-ShareAlike 2.0 Generic License](https://creativecommons.org/licenses/by-sa/2.0/)



decreasing length,  
decreasing abundance,  
external references



natural limits of clusters

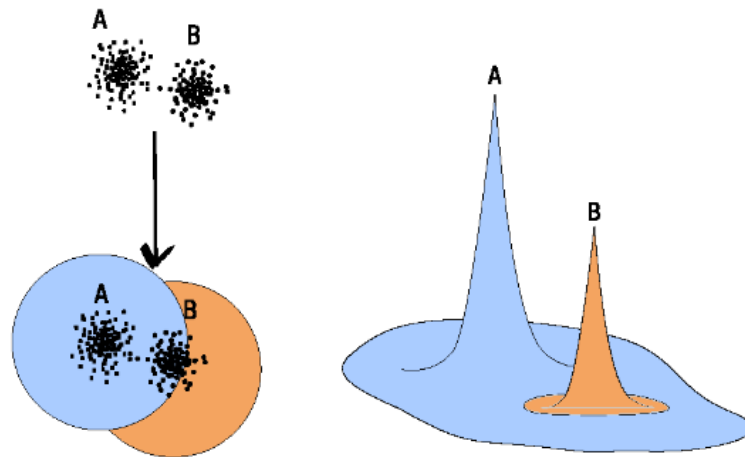


This work is licensed under a [Creative Commons Attribution-ShareAlike 2.0 Generic License](https://creativecommons.org/licenses/by-sa/2.0/)



# Swarm: A smart idea

## swarm large-scale clustering



clustering threshold (often 97%)  
is most of the time unadapted and  
can mask diversity.

swarm uses abundance values and a  
new clustering strategy to delineate  
natural high-quality OTUs.

*agglomeration rather than division*



Torbjørn Rognes  
Oslo University



Submitted 4 August 2015  
Accepted 31 October 2015  
Published 10 December 2015

Corresponding authors  
Frédéric Mahé,  
mahé@bioc.uzh.ch  
Torbjørn Rognes,  
torognes@ifi.uio.no

Academic editor  
Gilles van Wessel

Additional Information and  
Declarations can be found on  
page 10

DOI 10.7717/peerj.1420

© Copyright  
2015 Mahé et al.

Distributed under  
Creative Commons CC-BY 4.0

OPEN ACCESS

### Swarm v2: highly-scalable and high-resolution amplicon clustering

Frédéric Mahé<sup>1</sup>, Torbjørn Rognes<sup>2,3</sup>, Christopher Quince<sup>4</sup>,  
Colomban de Vargas<sup>5,6</sup> and Micah Dunthorn<sup>1</sup>

<sup>1</sup> Department of Ecology, Technische Universität Kaiserslautern, Kaiserslautern, Germany

<sup>2</sup> Department of Informatics, University of Oslo, Oslo, Norway

<sup>3</sup> Department of Microbiology, Oslo University Hospital, Rikshospitalet, Oslo, Norway

<sup>4</sup> Warwick Medical School, University of Warwick, Warwick, United Kingdom

<sup>5</sup> UMR 7144, EPEP-Evolution des Protistes et des Écosystèmes Pélagiques, Station Biologique de Roscoff, CNRS, Roscoff, France

<sup>6</sup> UMR7144 Station Biologique de Roscoff, Sorbonne Université, UPMC Univ Paris 06, Roscoff, France

#### ABSTRACT

Previously we presented Swarm v1, a novel and open source amplicon clustering program that produced fine-scale molecular operational taxonomic units (OTUs), free of arbitrary global clustering thresholds and input-order dependency. Swarm v1 worked with an initial phase that used iterative single-linkage with a local clustering threshold ( $d$ ), followed by a phase that used the internal abundance structures of clusters to break chained OTUs. Here we present Swarm v2, which has two important novel features: (1) a new algorithm for  $d = 1$  that allows the computation time of the program to scale linearly with increasing amounts of data; and (2) the new fastidious option that reduces under-grouping by grafting low abundant OTUs (e.g., singletons and doubletons) onto larger ones. Swarm v2 also directly integrates the clustering and breaking phases, dereplicates sequencing reads with  $d = 0$ , outputs OTU representatives in fasta format, and plots individual OTUs as two-dimensional networks.

**Subjects:** Biodiversity, Bioinformatics, Environmental Sciences, Microbiology, Molecular Biology  
**Keywords:** Environmental diversity, Barcoding, Molecular operational taxonomic units

#### INTRODUCTION

Traditional *de novo* amplicon clustering methods that can handle large high-throughput sequencing datasets (e.g., Edgar, 2010; Ghods, Liu & Pop, 2011; Fu et al., 2012) suffer from two fundamental problems. First, they rely on an arbitrary fixed global clustering threshold to group amplicons into molecular operational taxonomic units (OTUs). Global clustering thresholds have rarely been justified and are not applicable to all taxa and marker lengths (e.g., Caron et al., 2009; Nebel et al., 2011; Dunthorn et al., 2012; Brown et al., 2013). Second, there is variability in the clustering results due to amplicon input order (Koeppel & Wu, 2013; Mahé et al., 2014).

To solve these problems, we previously introduced the open source Swarm v1 program that implemented an initial clustering phase written in C++, then a breaking phase written in Python (Mahé et al., 2014). Swarm's clustering phase (Fig. 1A) was novel in its approach to single linkage clustering in that, instead of using a global clustering (e.g., Hartmann et al., 2012; Huse et al., 2010), amplicons were iteratively added together using a

How to cite this article: Mahé et al. (2015), Swarm v2: highly-scalable and high-resolution amplicon clustering. PeerJ 2:e1420.  
DOI 10.7717/peerj.1420



This work is licensed under a [Creative Commons Attribution-ShareAlike 2.0 Generic License](https://creativecommons.org/licenses/by-sa/2.0/)

# Swarm

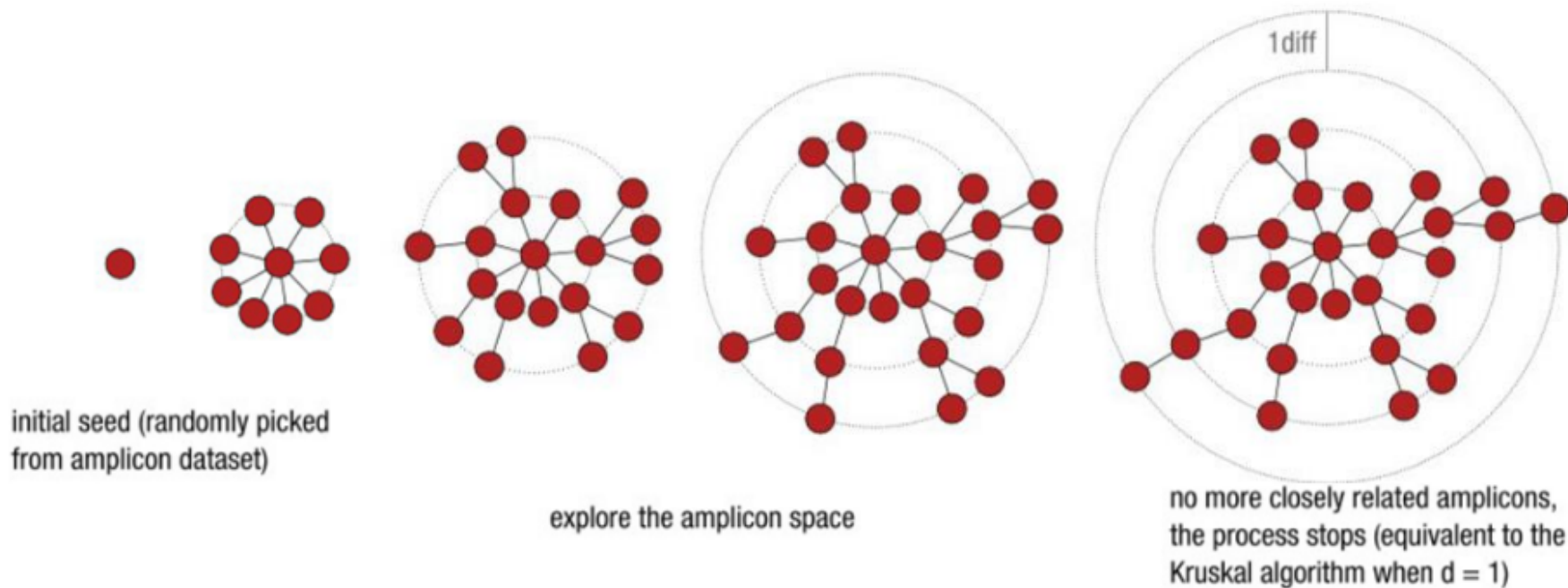
- A robust and fast clustering method for amplicon-based studies
- The purpose of swarm is to provide a novel clustering algorithm to handle large sets of amplicons
- swarm results are resilient to input-order changes and rely on a small local linking threshold  $d$ , the maximum number of differences between two amplicons
- swarm forms stable high-resolution clusters, with a high yield of biological information
- Default: forms a lot of low-abundant OTUs that are in fact artifacts and need to be removed
- **Swarm (fastidious method +  $d=1$ ) clusters + filters → ASVs**



This work is licensed under a [Creative Commons Attribution-ShareAlike 2.0 Generic License](https://creativecommons.org/licenses/by-sa/2.0/)

# d: the small local linking threshold

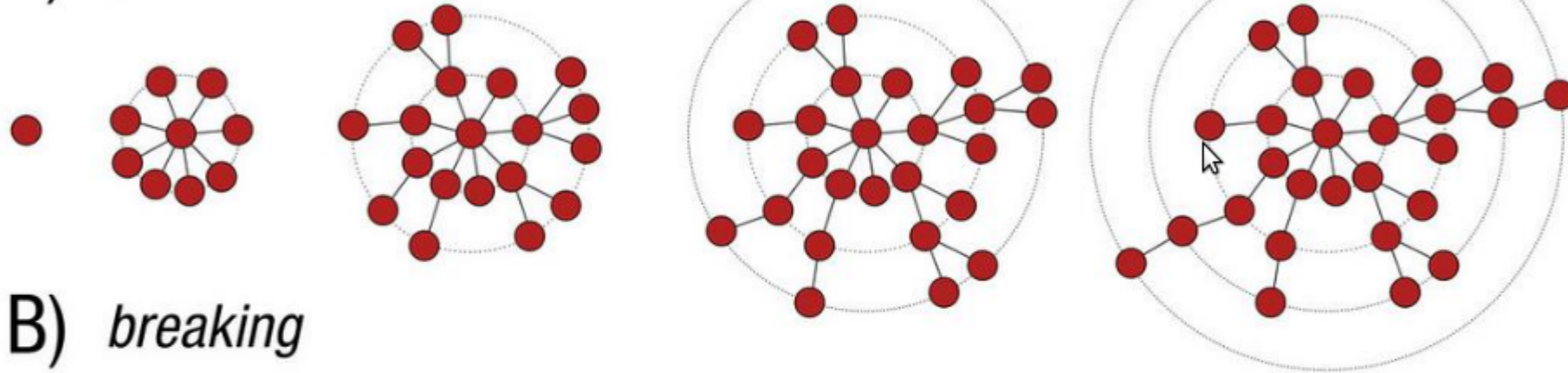
	ACGT	ACGT	ACGT
	AGGT	A - GT	A - - T
differences	1	1	2



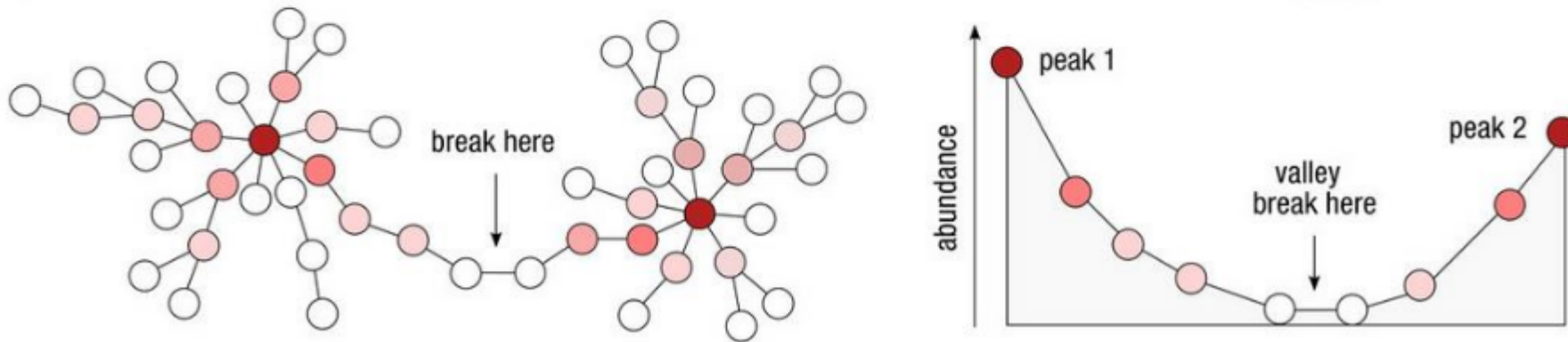
This work is licensed under a [Creative Commons Attribution-ShareAlike 2.0 Generic License](https://creativecommons.org/licenses/by-sa/2.0/)

# Swarm steps

## A) growth



## B) breaking



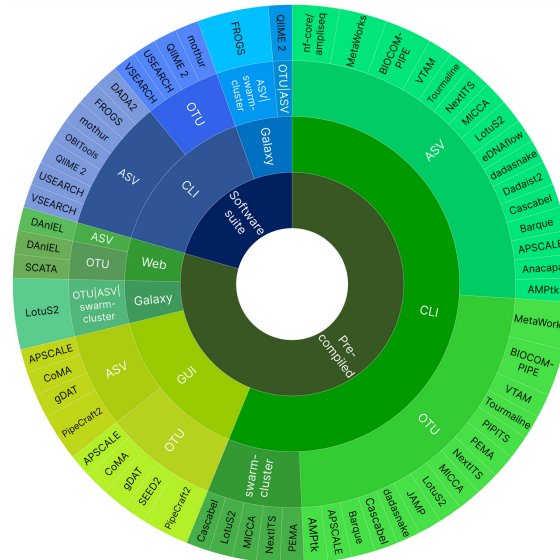
# Which method to choose?

**Single-end data**

AMPTk	Anacapa	BIOCOM-PIPE	Cascabel	CoMA	DADA2	dadasnake	
eDNAflow	FROGS	gDAT	JAMP	LotuS2	MetaWorks	MICCA	mothur
NextITS	nf-core/ampliseq	OBITools	PipeCraft2	QIIME 2	SCATA		
SEED2	USEARCH	VSEARCH	Tourmaline	VTAM			

**Paired-end data**

APSCALE	Barque	Dadaist2	DAnIEL	PIPITS	PEMA
---------	--------	----------	--------	--------	------



<b>Linux</b>	<b>macOS</b>	<b>Windows</b>
eDNAflow dadasnake MetaWorks NextITS	AMPTk Barque BIOCOM-PIPE Cascabel Dadaist2 JAMP OBITools PIPITS VSEARCH	Anacapa APSCALE CoMA DADA2 gDAT MICCA mothur nf-core/ampliseq PEMA PipeCraft2 USEARCH
DAnIEL	FROGS LotuS2 QIIME 2	Tourmaline VTAM
SCATA		SEED2

**Web-based (including Galaxy)**



This work is licensed under a [Creative Commons Attribution-ShareAlike 2.0 Generic License](https://creativecommons.org/licenses/by-sa/2.0/)

# Advantages and inconvenients

- *Which type of features to prefer may be context-dependent, and both may even be used in the same study*
- ASV demonstrate a biologically informative fine-scale resolution [9]
- But difficult to separate noise from a real signal in low abundant reads [10]
- ASVs represent stable and reproducible units across studies whereas OTUs are dataset-specific features (swarm clusters are not ⚠)
  - problematic for longitudinal and very big studies

📢 FROGS will soon offer the choice between swarm and dada2 for ASV creation



This work is licensed under a [Creative Commons Attribution-ShareAlike 2.0 Generic License](https://creativecommons.org/licenses/by-sa/2.0/)



RÉPUBLIQUE  
FRANÇAISE  
Liberté  
Égalité  
Fraternité

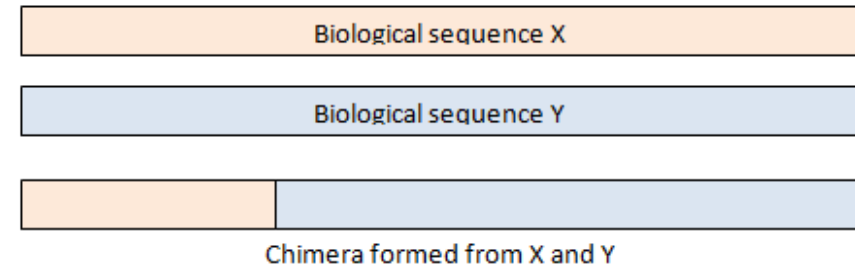
INRAE **mis@le**

# TP FROGS clustering



This work is licensed under a [Creative Commons Attribution-ShareAlike 2.0 Generic License](https://creativecommons.org/licenses/by-sa/2.0/)

# Chimera removal

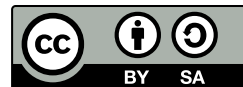


This work is licensed under a [Creative Commons Attribution-ShareAlike 2.0 Generic License](https://creativecommons.org/licenses/by-sa/2.0/)



# Chimera detection strategies

- Reference based: against a database of «genuine» sequences
  - dependant of the references used
- De novo: against abundant sequences in the samples 👍
- FROGS uses vsearch [4] as chimera removal tool



This work is licensed under a [Creative Commons Attribution-ShareAlike 2.0 Generic License](https://creativecommons.org/licenses/by-sa/2.0/)

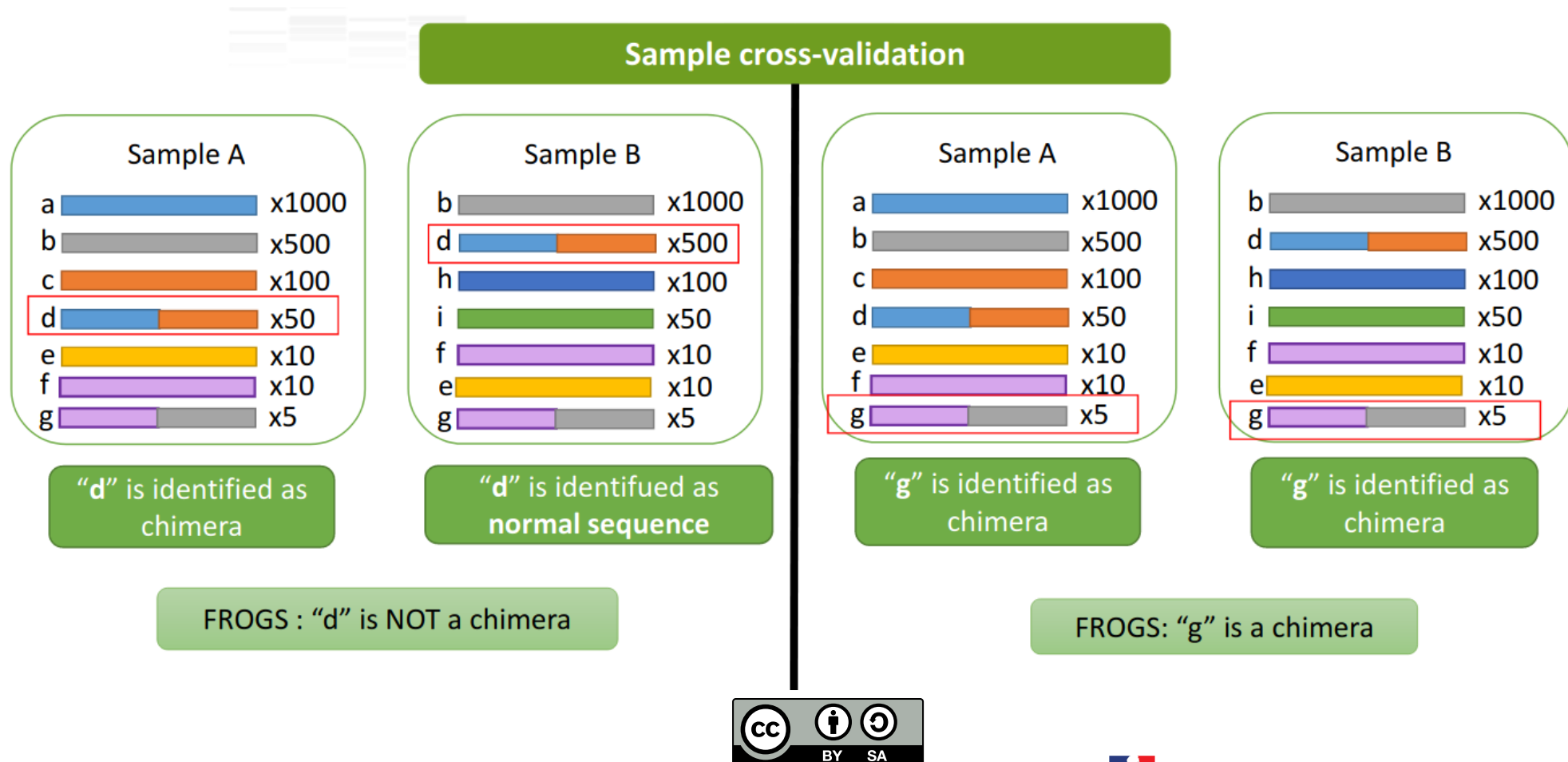


RÉPUBLIQUE  
FRANÇAISE  
Liberté  
Égalité  
Fraternité

INRAE **mission**

# A little extra: the sample-cross validation

- FROGS adds a sample-cross validation



# Chimera rates in samples

- From 5 to 40% in 16S data
- Few with ITS (<10%)

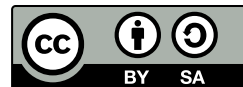
Samples	% Observed Chimera content			
	ABI3730	454 FLX Titanium		
	V1-V9	V1-V3	V3-V5	V6-V9
MC	5.99±3.07	14.26±10.34	14.75±9.45	13.49±8.52
gut	7.71±6.46	22.90±8.56	16.03±2.86	17.76±3.76
oral	7.22±6.35	20.55±11.73	10.98±4.01	9.10±5.02
skin	3.49±5.77	11.15±1.36	7.51±2.49	5.73±1.69
vaginal	6.31±6.64	12.60±6.70	6.62±3.51	3.00±1.65

\*Values are averages ± STDEV calculated from multiple replicates of MC, and from replicates of multiple clinical samples originating from different body sites.  
doi:10.1371/journal.pone.0039315.t001



This work is licensed under a [Creative Commons Attribution-ShareAlike 2.0 Generic License](https://creativecommons.org/licenses/by-sa/2.0/)

# TP Frogs remove chimera



This work is licensed under a [Creative Commons Attribution-ShareAlike 2.0 Generic License](https://creativecommons.org/licenses/by-sa/2.0/)

# Abundance/Prevalence filters



This work is licensed under a [Creative Commons Attribution-ShareAlike 2.0 Generic License](https://creativecommons.org/licenses/by-sa/2.0/)

# How to filter clusters?

- Low abundant sequences
- Clusters not shown in few replicates
- Contamination



This work is licensed under a [Creative Commons Attribution-ShareAlike 2.0 Generic License](https://creativecommons.org/licenses/by-sa/2.0/)

# TP Frogs cluster filters



This work is licensed under a [Creative Commons Attribution-ShareAlike 2.0 Generic License](https://creativecommons.org/licenses/by-sa/2.0/)

# Taxonomic affiliation



This work is licensed under a [Creative Commons Attribution-ShareAlike 2.0 Generic License](https://creativecommons.org/licenses/by-sa/2.0/)



# Comparison of approaches



This work is licensed under a [Creative Commons Attribution-ShareAlike 2.0 Generic License](https://creativecommons.org/licenses/by-sa/2.0/)

# Table 1 Number of taxonomic groups identified by each classifier among Illumina 16S rRNA gene sequences (SRR3225706) from a mock microbiome sample [33]. Counts are provided with and without including any sequences in the RDP training set that are labeled as belonging to the 20 expected genera

From: [IDTAXA: a novel approach for accurate taxonomic classification of microbiome sequences](#)

		Classified to genus level <sup>a</sup> (%)	Groups present in the mock community							Absent from mock community <sup>b</sup>		
			Root	Domain	Phylum	Class	Order	Family	Genus	Order	Family	Genus
Using the RDP training set	BLAST	97.9	1	0	0	0	0	0	17	0	0	24
	IDTAXA	94.2	1	0	1	1	2	5	14	0	1	2
	MAPSeq	96.5	1	0	0	0	0	4	15	0	2	6
	QIIME	95.4	1	0	0	0	0	0	16	0	0	7
	RDP Classifier	93.3	1	1	2	3	6	8	15	0	2	6
	SINTAX	94.2	1	1	1	4	3	3	14	1	0	3
	SPINGO	96.5	1	0	0	0	0	0	17	0	0	3
With expected genera excluded from training data	BLAST	17.3	1	0	0	0	0	0	0	0	0	65
	IDTAXA	0.01	1	1	1	2	3	4	0	0	2	2
	MAPSeq	24.6	1	0	0	2	5	11	0	1	8	20
	QIIME	13.5	1	0	0	0	0	0	0	0	0	16
	RDP Classifier	3.83	1	1	2	3	6	9	0	0	3	12
	SINTAX	8.76	1	1	1	7	5	6	0	1	1	9
	SPINGO	26.7	1	0	0	0	0	0	0	0	0	15

<sup>a</sup>Percent of total sequences from the mock community that were classified to the genus rank

<sup>b</sup>Other rank levels (root, domain, phylum, and class) all had counts of zero



This work is licensed under a [Creative Commons Attribution-ShareAlike 2.0 Generic License](#)



# RDP problems

- Depends too much on the databank used!
- Gives one affiliation for each feature with bootstrap, on each subdivision

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



This work is licensed under a [Creative Commons Attribution-ShareAlike 2.0 Generic License](https://creativecommons.org/licenses/by-sa/2.0/)

# The FROGS recommendation


- Use Blast and not RDP
- Check Blast metrics to avoid concluding too fast
- Take care of the reference databank used!

Bacteria;Actinobacteriota;Actinobacteria;Propionibacteriales;Propionibacteriaceae



This work is licensed under a [Creative Commons Attribution-ShareAlike 2.0 Generic License](https://creativecommons.org/licenses/by-sa/2.0/)

# The FROGS databanks

- Command line: you can use your own databank
- Galaxy
  - You have access to **several databanks**
  - Admins have to add your databank
- The file must be well formatted, we can do it for you
-  For private databanks, contact us!



This work is licensed under a [Creative Commons Attribution-ShareAlike 2.0 Generic License](#)

# The FROGS extra: the multi-affiliations

- FROGS gives all identical hits

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

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

*Strictly identical (V1-V3 amplification) on 499 nucleotides*

- FROGS can't decide if it's one or another
- You have to check if you can choose between multi-affiliations



This work is licensed under a [Creative Commons Attribution-ShareAlike 2.0 Generic License](https://creativecommons.org/licenses/by-sa/2.0/)

# To help you

- <https://shiny.migale.inrae.fr/app/affiliationexplorer>
- a very user-friendly Shiny web app, allowing users to modify very simply the affiliations from a FROGS abundance file



This work is licensed under a [Creative Commons Attribution-ShareAlike 2.0 Generic License](https://creativecommons.org/licenses/by-sa/2.0/)

# Filter ASVs on their affiliation



This work is licensed under a [Creative Commons Attribution-ShareAlike 2.0 Generic License](https://creativecommons.org/licenses/by-sa/2.0/)



# Affiliation filters

- Remaining contamination?
- Want to analyse only the **Firmicutes**?
- 2 modes
  - *Deleting*: remove ASVs
  - *Hiding*: only the affiliation is modified, not the abundance



This work is licensed under a [Creative Commons Attribution-ShareAlike 2.0 Generic License](https://creativecommons.org/licenses/by-sa/2.0/)



RÉPUBLIQUE  
FRANÇAISE  
Liberté  
Égalité  
Fraternité

INRAE **mission**

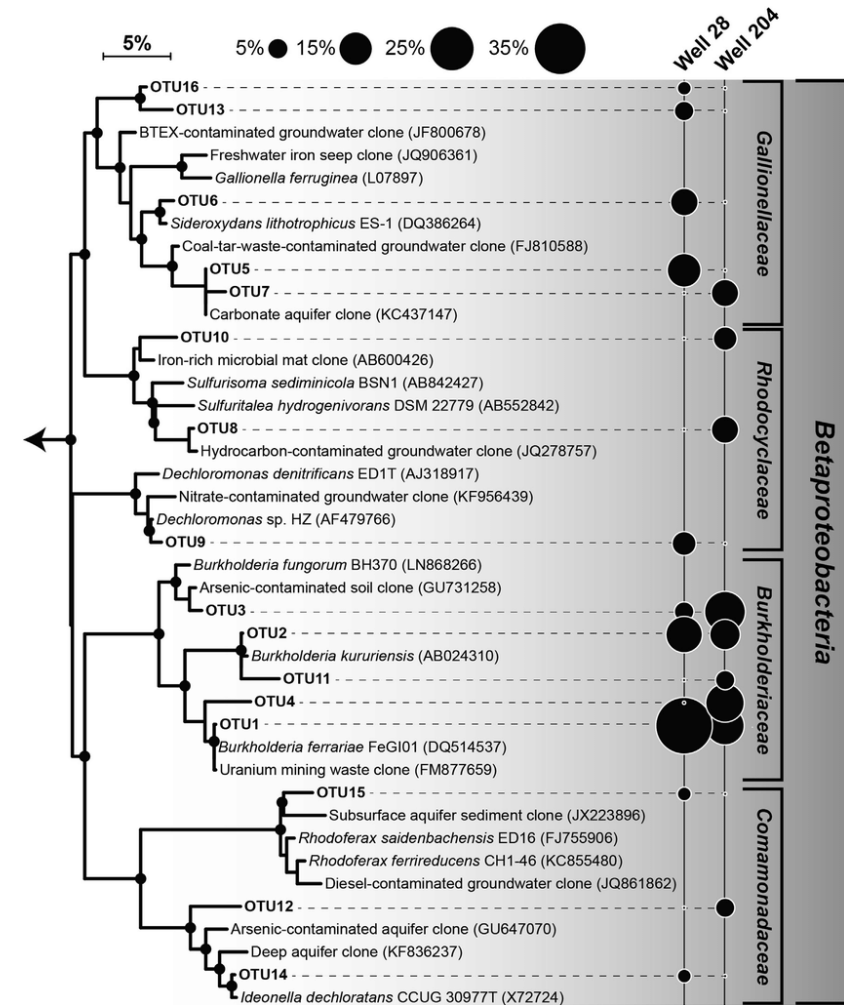
# Phylogenetic tree



This work is licensed under a [Creative Commons Attribution-ShareAlike 2.0 Generic License](https://creativecommons.org/licenses/by-sa/2.0/)

# FROGS tree

- This tool builds a phylogenetic tree thanks to affiliations of ASVs contained in the BIOM file
- Needed to compute beta-diversity indices based on phylogenetic distances
- Interesting to explore poor-characterized environments

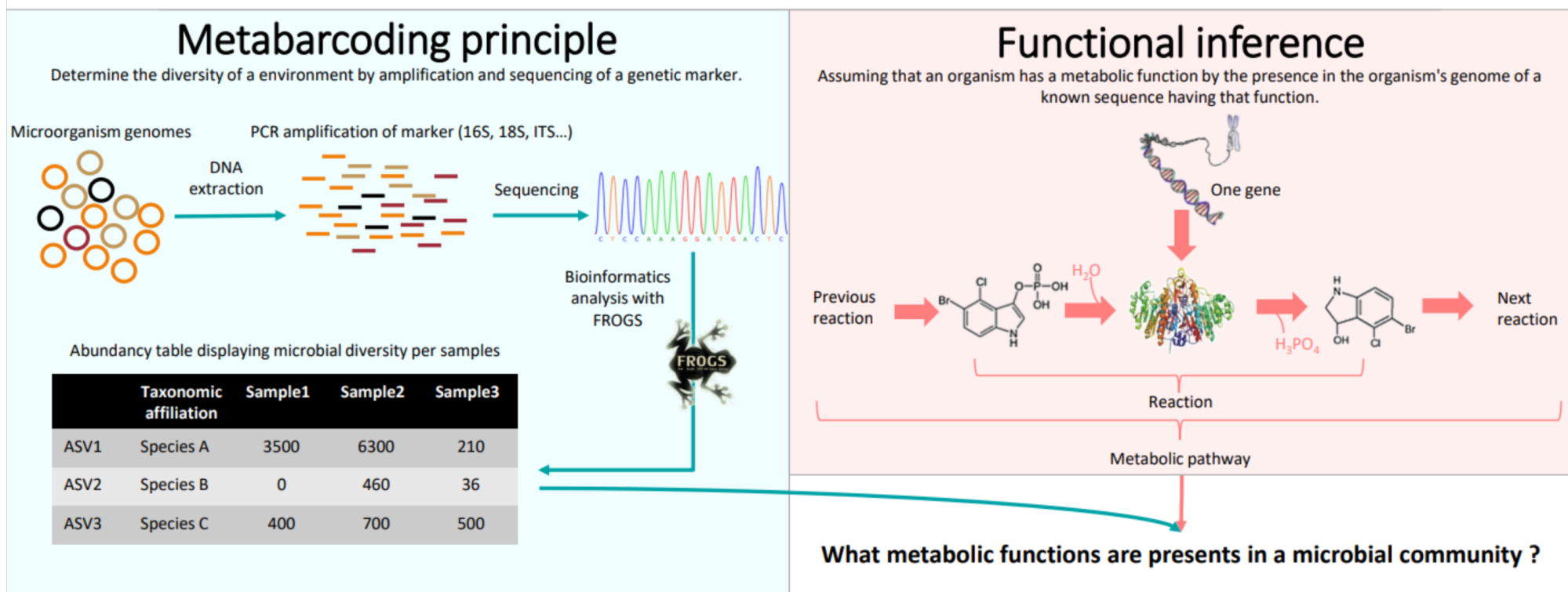


# FROGSfunc: function inference



This work is licensed under a [Creative Commons Attribution-ShareAlike 2.0 Generic License](https://creativecommons.org/licenses/by-sa/2.0/)

# Concepts



# Based on PICRUSt2

- PICRUSt [13] (Phylogenetic investigation of communities by reconstruction of unobserved states) is an open-source tool.
- It is a software for predicting functional abundances based only on marker gene sequences
- PICRUSt2 is composed of 4 python applications.
- No graphic interface exists to run PICRUSt2 for non-expert users.



This work is licensed under a [Creative Commons Attribution-ShareAlike 2.0 Generic License](https://creativecommons.org/licenses/by-sa/2.0/)

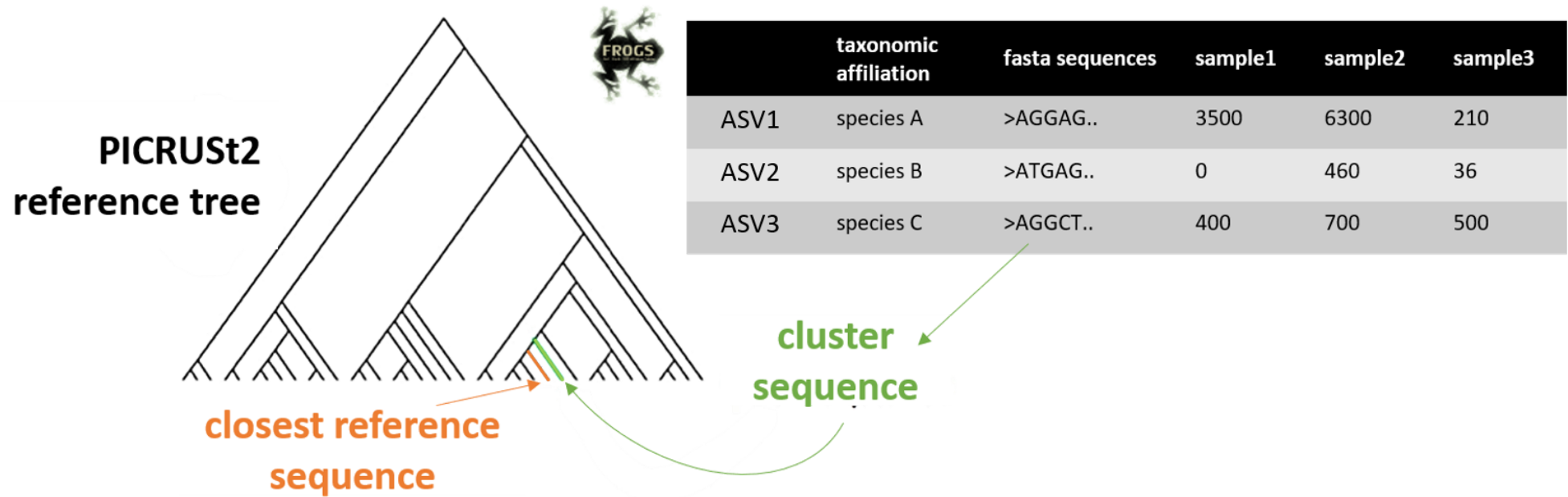
# How it works

1. Places the ASVs into a reference phylogenetic tree and predicts of marker copy number in each ASV.
2. Predicts number of function copy number in each ASV and calculates functions abundances in each sample and ASV abundances according to marker copy number.
3. Calculates pathway abundances in each sample.



This work is licensed under a [Creative Commons Attribution-ShareAlike 2.0 Generic License](https://creativecommons.org/licenses/by-sa/2.0/)

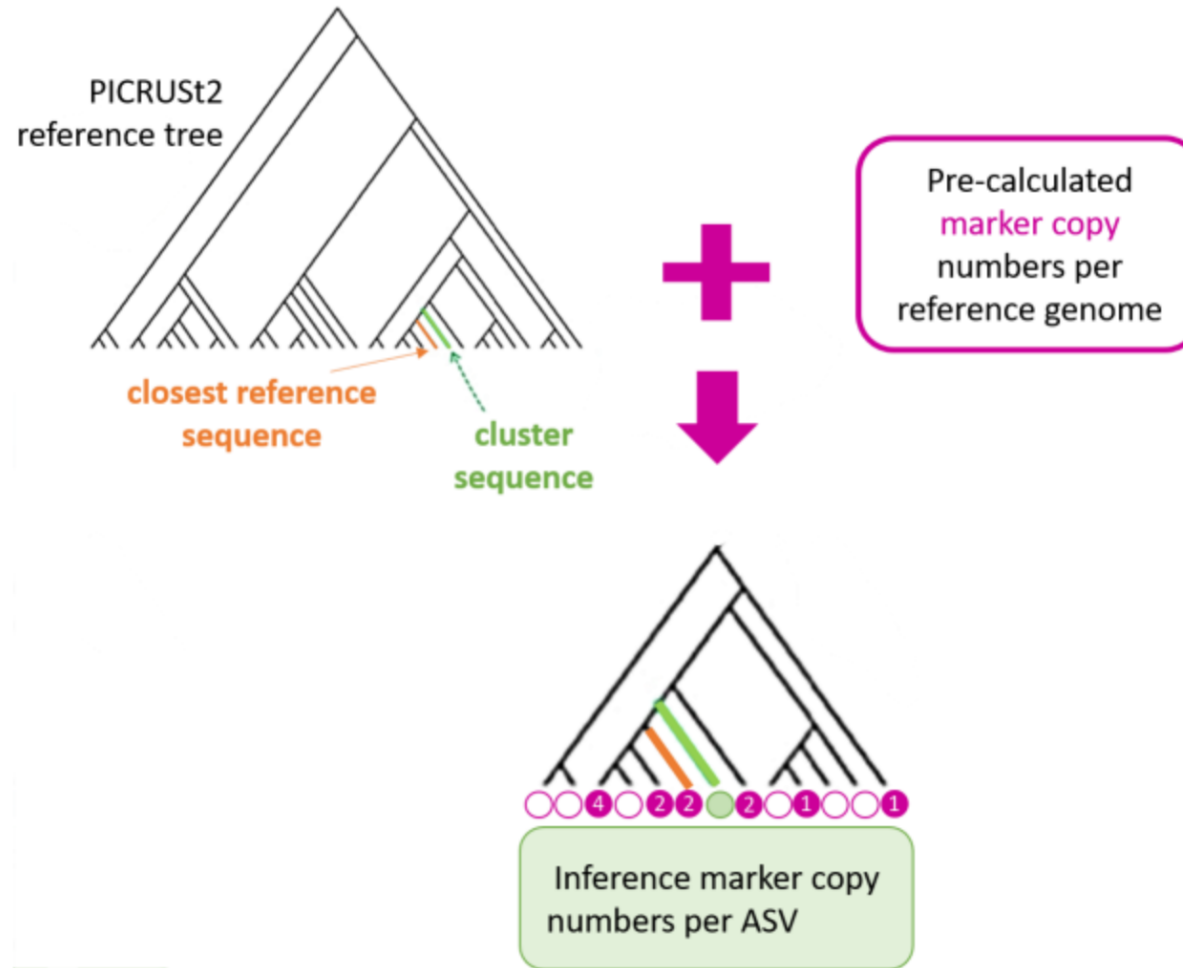
# FROGSfunc placeseqs and copynumber



This work is licensed under a [Creative Commons Attribution-ShareAlike 2.0 Generic License](https://creativecommons.org/licenses/by-sa/2.0/)



# FROGSfunc placeseqs and copynumber



This work is licensed under a [Creative Commons Attribution-ShareAlike 2.0 Generic License](https://creativecommons.org/licenses/by-sa/2.0/)

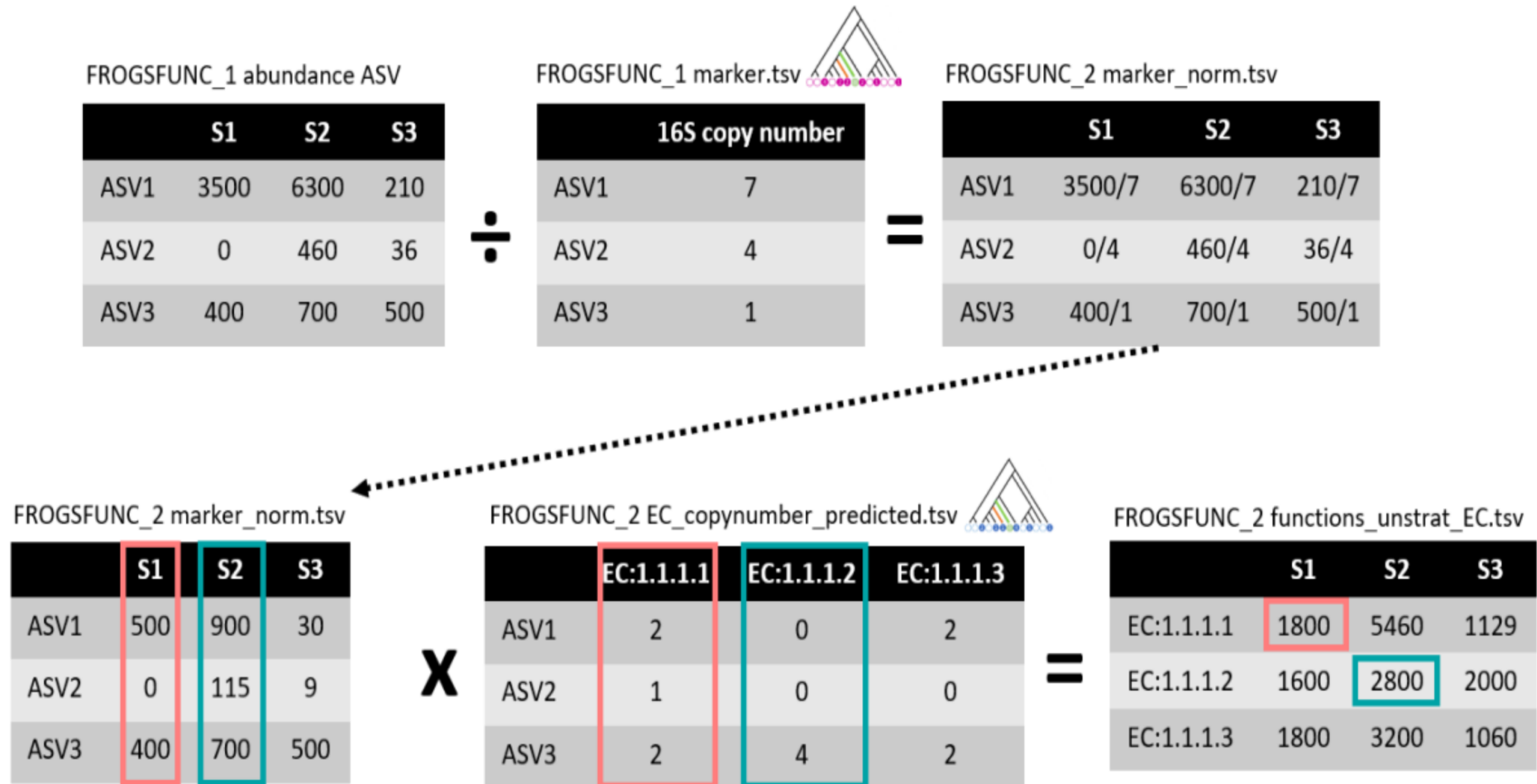
# NSTI

- NSTI scores are simply the average branch length that separates each ASV in your sample from a reference bacterial genome, weighted by the abundance of that ASV in the sample.
- PICRUSt2 sets NSTI threshold to 2 per default. Some studies have shown that this threshold is permissive. Thus, it is important to see if the taxonomies between PICRUSt2 and FROGS are quite similar or not, in order to potentially choose a more stringent threshold afterwards.
  - $0 < \text{Good} < 0.5$
  - $0.5 \leq \text{Medium} < 1$
  - $1 \leq \text{Bad} < 2$
  - To exclude  $\geq 2$



This work is licensed under a [Creative Commons Attribution-ShareAlike 2.0 Generic License](https://creativecommons.org/licenses/by-sa/2.0/)

# FROGSfunc functions



This work is licensed under a [Creative Commons Attribution-ShareAlike 2.0 Generic License](https://creativecommons.org/licenses/by-sa/2.0/)

# FROGSfunc pathways

FROGSFUNC\_2 functions\_unstrat\_EC.tsv

	S1	S2	S3
EC:1.1.1.1	1800	5460	1129
EC:1.1.1.2	1600	2800	2000
EC:1.1.1.3	1800	3200	1060

+

PICRUST2 map of gene families to pathways



FROGSFUNC\_3 pathways\_unstrat per sample and per reference

Pathways	S1	S2	S3
1CMET2-PWY	1289.7451	1485.2474	1233.5908
ANAEROFRUCAT-PWY	904.7455	1565.5453	1227.6231
ANAGLYCOLYSIS-PWY	1501.0804	1805.3271	1544.3206
ARG+POLYAMINE-SYN	0	49.3391	45.6559



This work is licensed under a [Creative Commons Attribution-ShareAlike 2.0 Generic License](https://creativecommons.org/licenses/by-sa/2.0/)

# References

1. Escudié F, Auer L, Bernard M, Mariadassou M, Cauquil L, Vidal K, et al. FROGS: Find, rapidly, OTUs with galaxy solution. *Bioinformatics*. 2017;34:1287–94.
2. Bernard M, Rué O, Mariadassou M, Pascal G. FROGS: a powerful tool to analyse the diversity of fungi with special management of internal transcribed spacers. *Briefings in Bioinformatics*. 2021;22. doi:[10.1093/bib/bbab318](https://doi.org/10.1093/bib/bbab318).
3. Stoler N, Nekrutenko A. Sequencing error profiles of Illumina sequencing instruments. *NAR Genomics and Bioinformatics*. 2021;3. doi:[10.1093/nargab/lqab019](https://doi.org/10.1093/nargab/lqab019).
4. Rognes T, Flouri T, Nichols B, Quince C, Mahé F. VSEARCH: A versatile open source tool for metagenomics. *PeerJ*. 2016;4:e2584.
5. Magoč T, Salzberg SL. FLASH: Fast length adjustment of short reads to improve genome assemblies. *Bioinformatics*. 2011;27:2957–63.
6. Zhang J, Kobert K, Flouri T, Stamatakis A. PEAR: A fast and accurate illumina paired-end reAd mergeR. *Bioinformatics*. 2013;30:614–20.
7. Hakimzadeh A, Abdala Asbun A, Albanese D, Bernard M, Buchner D, Callahan B, et al. A pile of pipelines: An overview of the bioinformatics software for metabarcoding data analyses. *Molecular Ecology Resources*. 2023.
8. Mahé F, Rognes T, Quince C, Vargas C de, Dunthorn M. Swarm v2: Highly-scalable and high-resolution amplicon clustering. *PeerJ*. 2015;3:e1420.
9. Couton M, Baud A, Daguin-Thiébaud C, Corre E, Comtet T, Viard F. High-throughput metabarcoding of ethanol is effective at jointly examining infraspecific and taxonomic diversity.



do not perform equally. Ecology and Evolution. 2021;11:5533–46.

10. De Santiago A, Pereira TJ, Mincks SL, Bik HM. Dataset complexity impacts both MOTU delimitation and biodiversity estimates in eukaryotic 18S rRNA metabarcoding studies. Environmental DNA. 2022;4:363–84.

11. Group JCHMPDGW. Evaluation of 16S rDNA-based community profiling for human microbiome research. PloS one. 2012;7:e39315.

12. Murali A, Bhargava A, Wright ES. IDTAXA: A novel approach for accurate taxonomic classification of microbiome sequences. Microbiome. 2018;6:1–14.

13. Douglas GM, Maffei VJ, Zaneveld JR, Yurgel SN, Brown JR, Taylor CM, et al. PICRUSt2 for prediction of metagenome functions. Nature biotechnology. 2020;38:685–8.



This work is licensed under a [Creative Commons Attribution-ShareAlike 2.0 Generic License](https://creativecommons.org/licenses/by-sa/2.0/)



This work is licensed under a [Creative Commons Attribution-ShareAlike 2.0 Generic License](https://creativecommons.org/licenses/by-sa/2.0/)

