

# Session de formation 2018



- 12 Mars**      Guide de survie à Linux : les commandes de base pour débuter sur un serveur linux
- 13 Mars**      Linux avancé : manipuler et filtrer des fichiers sans connaissance de programmation
- 15 Mars**      Initiation à l'utilisation du cluster bioinformatique itrop
- 22 Mars**      Initiation à git
- 23 Mars**      Initiation aux gestionnaires de workflow South Green: Galaxy ou TOGGLE
- 26 Mars**      Initiation aux analyses de données transcriptomiques
- 23 Avril**      **Initiation aux analyses de données metabarcoding**





IRD

Institut de Recherche  
pour le Développement

# South Green

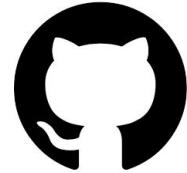
bioinformatics platform



plateau i-trop



[www.southgreen.fr](http://www.southgreen.fr)



<https://github.com/SouthGreenPlatform>



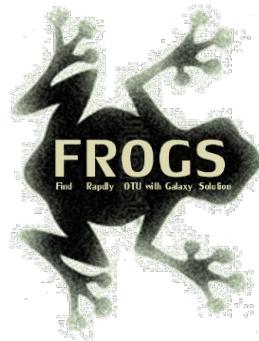
*The South Green portal: a comprehensive resource for tropical and Mediterranean crop genomics*, Current Plant Biology, 2016

# Session de formation 2018



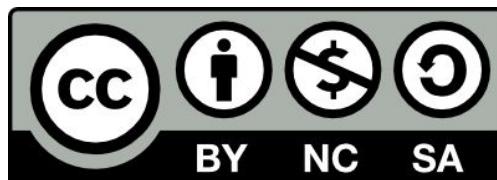
- Toutes nos formations :  
<https://southgreenplatform.github.io/trainings/>
- Topo & TP : [Metagenomics](#)
- Environnement de travail : [Logiciels à installer](#)

# Initiation aux analyses de données metabarcoding



[www.southgreen.fr](http://www.southgreen.fr)

<https://southgreenplatform.github.io/trainings>



# Planning

## 1. Introduction générale

## 2. Partie pratique

*Practice 1:* Obtaining an OTU table with FROGS in Galaxy

*Practice 2:* Visualizing and plotting all sample results with Phinch

*Practice 3:* Handling and visualisation of OTU table using PhyloSeq R package

## 3. Conclusions

## What metagenomics is ?

Metagenomics ( Environmental Genomics or Community Genomics) is the study of genomes recovered from environmental samples without the need for culturing them

Metagenomics processes data using bioinformatics tools

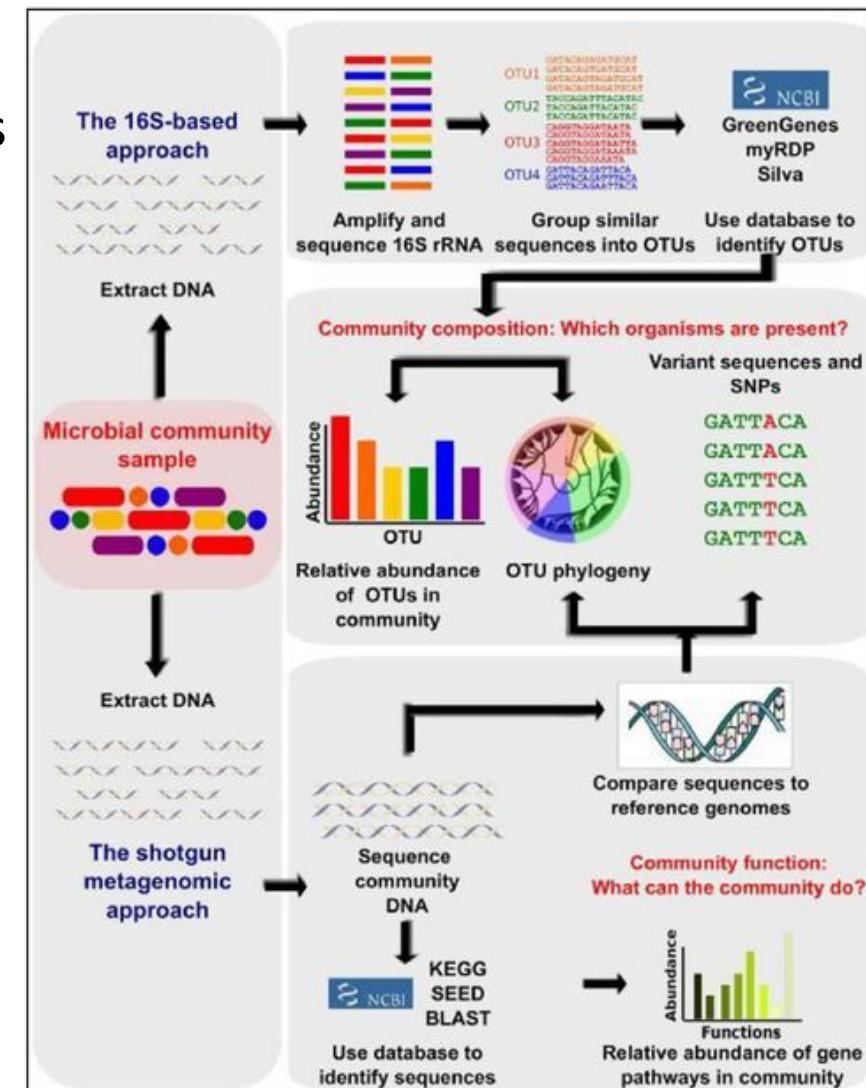
=> Organisms can be studied directly in their environments bypassing the need to isolate each species

=> There are significant advantages for viral metagenomics, because of difficulties cultivating the appropriate host

# Two main strategies in metagenomics

We can distinguish targeted metagenomics or shot-gun metagenomics :

- 16S rRNA metabarcoding is used to characterize the bacterial communities of an environment
- Whole-genome sequencing when the goal is to identify gene functions and pathways, or reconstruct microbial genomes.



# Markers genes vs Shotgun metagenomics

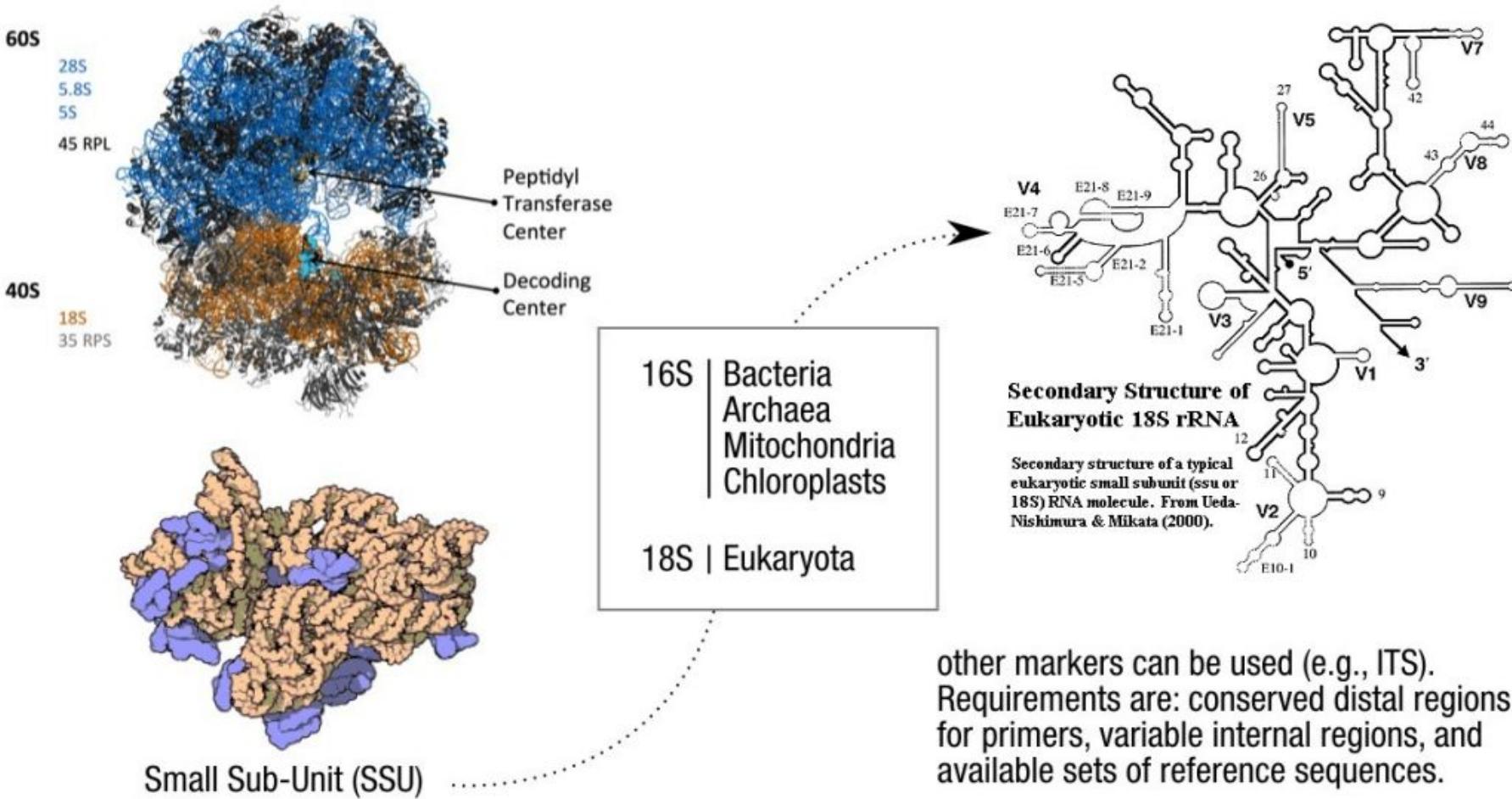
Marker Gene Profiling	Shotgun Metagenomics Profiling
Less expensive (~\$100 per sample)	Still very expensive (~\$1000 per sample)
Computational needs can be met by desktop / small server computers	Usually requires huge computational resources (cluster of computers)
Provides mainly taxonomic profiling	Provides both taxonomic and functional profiling
For 16S, majority of genes can be assigned at least to phylum level	Many more unassigned gene fragments ("wasted" data)
Relatively free of host DNA contamination	Prone to host DNA contamination

# Strategies in Diversity Characterisation

Technique	Advantages and challenges	Main applications
Metataxonomics using amplicon sequencing of the 16S or 18S rRNA gene or ITS	<ul style="list-style-type: none"> <li>+ Fast and cost-effective identification of a wide variety of bacteria and eukaryotes</li> <li>- Does not capture gene content other than the targeted genes</li> <li>- Amplification bias</li> <li>- Viruses cannot be captured</li> </ul>	<ul style="list-style-type: none"> <li>* Profiling of what is present</li> <li>* Microbial ecology</li> <li>* rRNA-based phylogeny</li> </ul>
Metagenomics using random shotgun sequencing of DNA or RNA	<ul style="list-style-type: none"> <li>+ No amplification bias</li> <li>+ Detects bacteria, archaea, viruses and eukaryotes</li> <li>+ Enables <i>de novo</i> assembly of genomes</li> <li>- Requires high read count</li> <li>- Many reads may be from host</li> <li>- Requires reference genomes for classification</li> </ul>	<ul style="list-style-type: none"> <li>* Profiling of what is present across all domains</li> <li>* Functional genome analyses</li> <li>* Phylogeny</li> <li>* Detection of pathogens</li> </ul>
Meta-transcriptomics using sequencing of mRNA	<ul style="list-style-type: none"> <li>+ Identifies active genes and pathways</li> <li>- mRNA is unstable</li> <li>- Multiple purification and amplification steps can lead to more noise</li> </ul>	* Transcriptional profiling of what is active

# Metabarcoding strategie

# A universal gene: ribosomal RNA

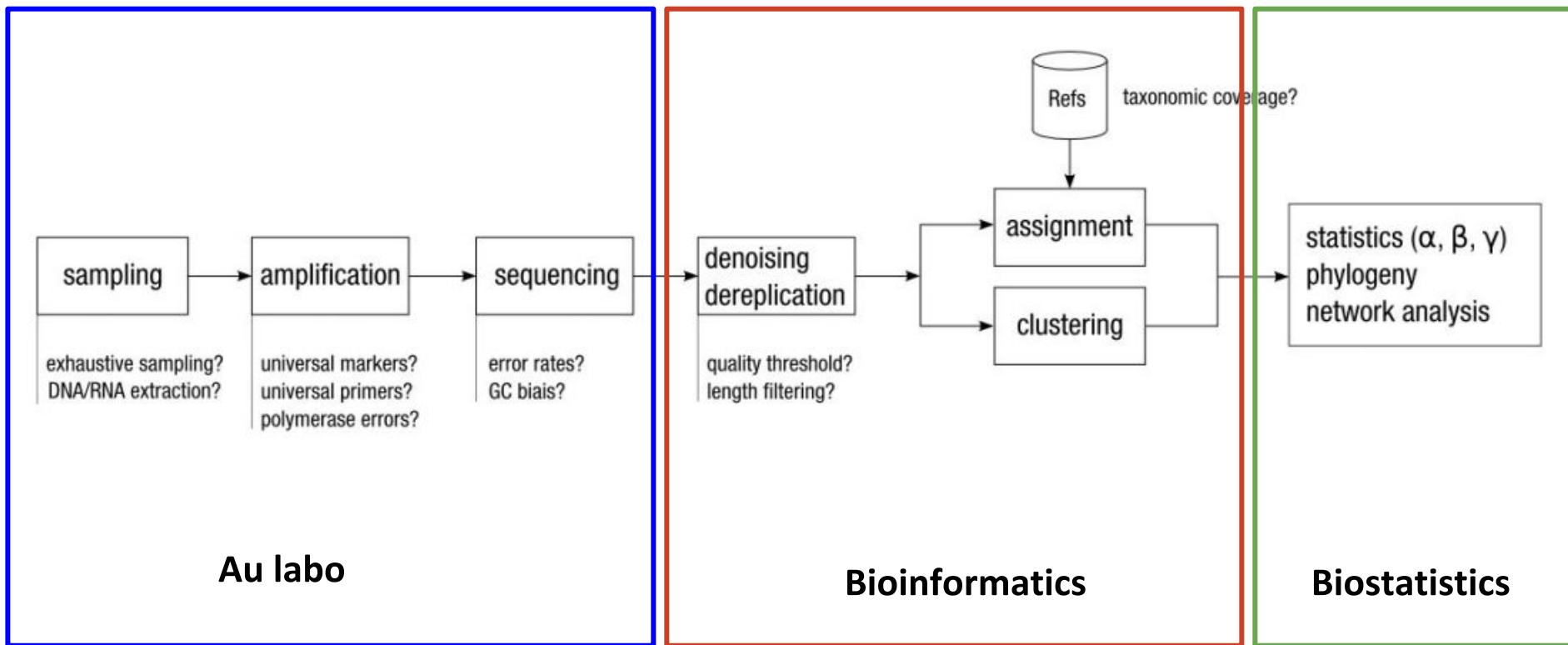


## Projets métagénomiques

4900 [projets sur NCBI](#) (avril 2018)

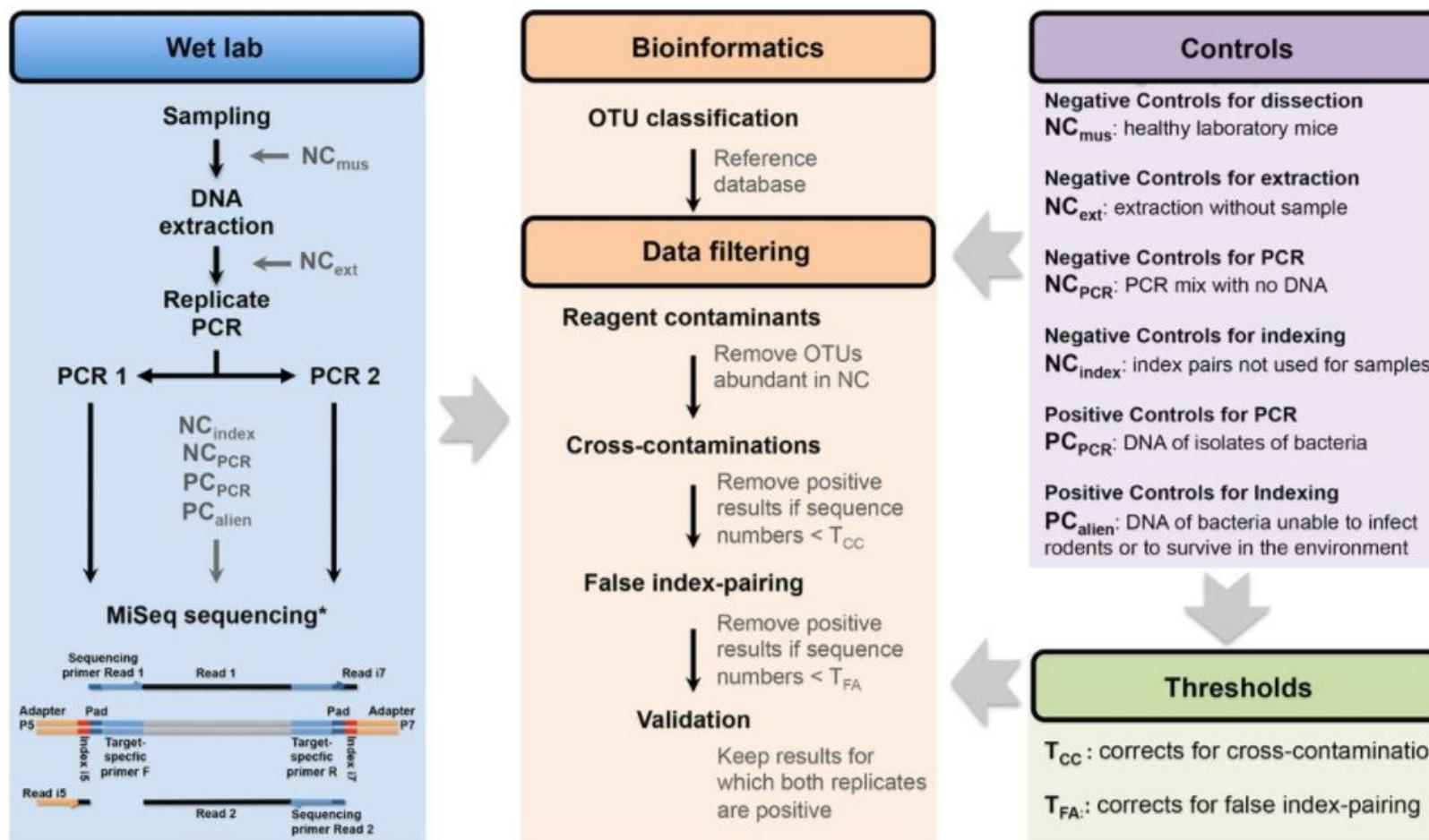
- Sable de plage
- Moustique
- Corail
- Glace
- Air de la ville de Singapour
- Surface de la cuvette des toilettes
- Fromages
- ...

# Amplicon-based studies general pipeline

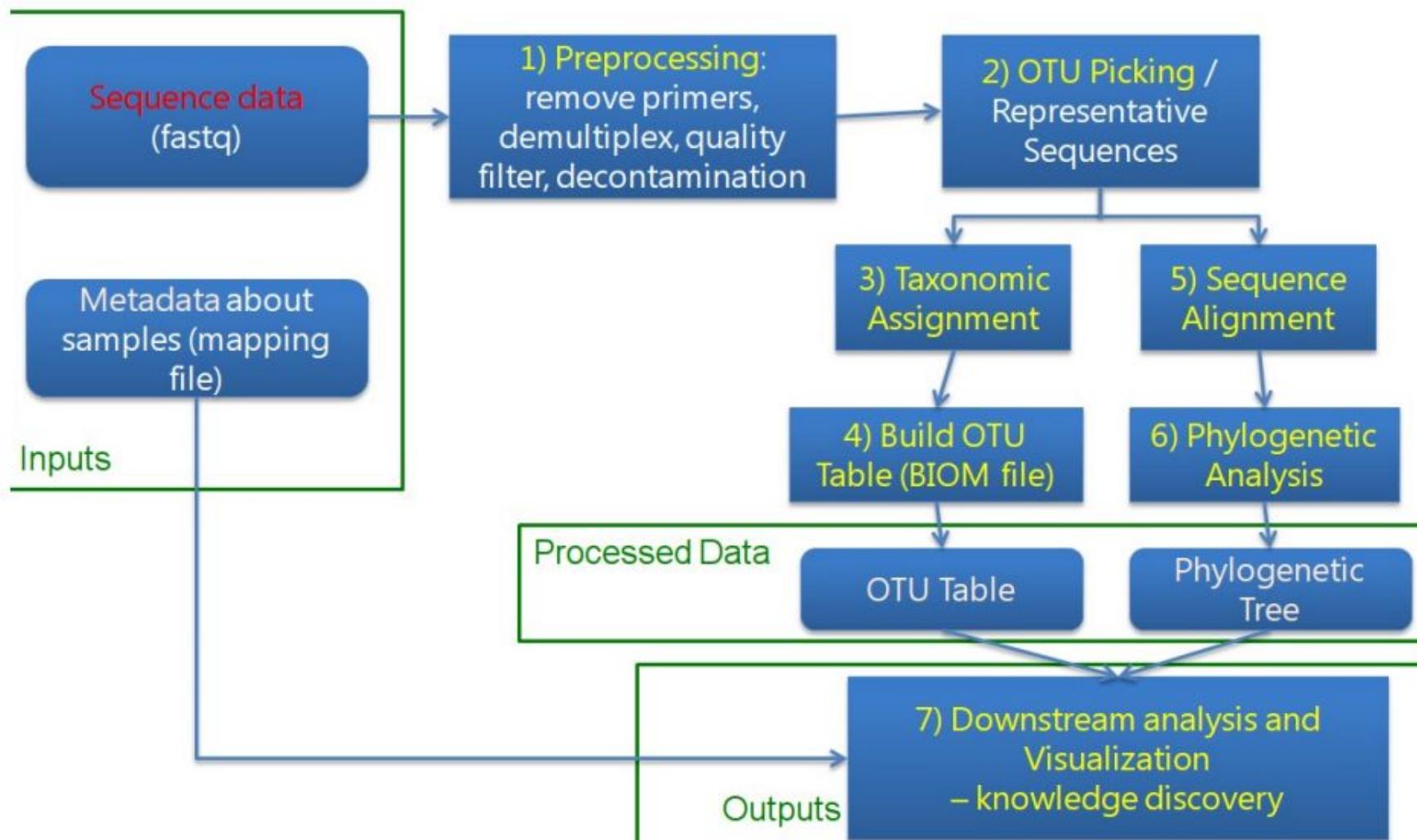


**TODAY !**

# Workflow of the wet laboratory, bioinformatics, and data filtering procedures in the process of data filtering for 16S rRNA amplicon sequencing.



## Overall bioinformatics workflow



## Major metagenomics pipelines

targeted amplification



Mothur (2009)  
Patrick Schloss  
open-source  
single piece  
most cited  
stats



Qiime (2010)  
Gregory Caporaso  
open-source  
python wrapper  
most used(?)  
stats



Uparse (2013)  
Robert Edgar  
closed-source  
usearch commands  
popular  
no stats

# Which bioinformatics solutions?

	Disadvantages
QIIME	Installation problem Command lines
UPARSE	Global clustering command lines
MOTHUR	Not MiSeq data without normalization Global hierarchical clustering Command lines
MG-RAST	No modularity No transparency



QIIME allows analysis of high-throughput community sequencing data  
 J Gregory Caporaso et al, *Nature Methods*, 2010; doi:10.1038/nmeth.f.303

Introducing mothur: Open-source, platform-independent, community-supported software for describing and comparing microbial communities.

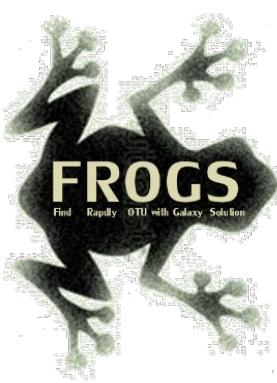
Schloss, P.D., et al., *Appl Environ Microbiol*, 2009, doi: 10.1128/AEM.01541-09

UPARSE: Highly accurate OTU sequences from microbial amplicon reads  
 Edgar, R.C. et al, *Nature Methods*, 2013, dx.doi.org/10.1038/nmeth.2604

The metagenomics RAST server – a public resource for the automatic phylogenetic and functional analysis of metagenomes

F Meyer et al, *BMC Bioinformatics*, 2008, doi:10.1186/1471-2105-9-386

# FROGS: Find, Rapidly, OTUs with Galaxy Solution



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

*Bioinformatics*, Volume 34, Issue 8, 15 April 2018, Pages 1287–1294, <https://doi.org/10.1093/bioinformatics/btx791>

<https://github.com/geraldinepascal/FROGS.git>

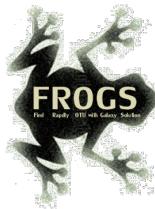


## *Practice 1:* **Obtaining an OTU table with FROGS in Galaxy**

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



- Use platform Galaxy
- Set of modules= Tools to analyze your “big” data
- Independent modules
- Run on Illumina/454 data 16S, 18S, and 23S
- New clustering method
- Many graphics for interpretation
- User friendly, hiding bioinformatics infrastructure/complexity



# FROGS Pipeline on Galaxy

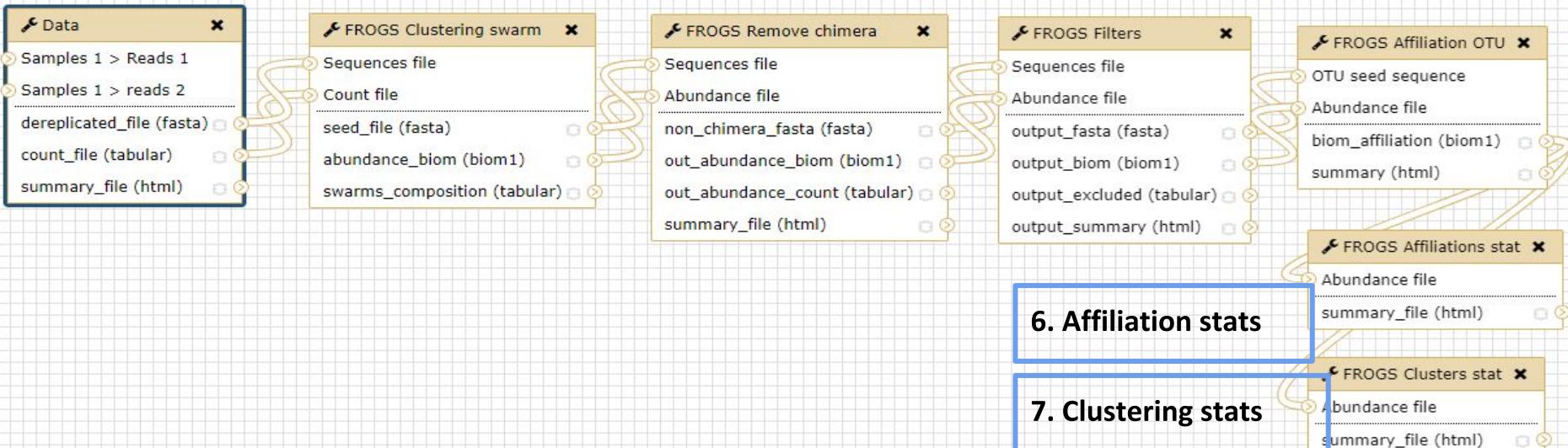
**1. Pre-process**

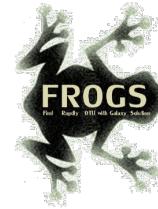
**2. Clustering**

**3. Remove Chimera**

**4. Filtering**

**5. Affiliation OTU**

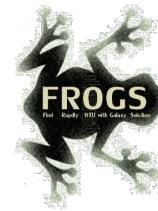




## 1.Pre-process

A preprocessing tool :

- merges paired sequences into contigs with flash,
- cleans the data with cutadapt,
- deletes the chimeras with VSEARCH combined with a cross-validation method and
- dereplicates sequences with a home-made python script.



## 1.Pre-process

FROGs takes the reads (R1 and R2) from multiple samples and performs the following steps:

- If the data is not in contigs, R1 et R2 will be overlapped
- Contigs that are too big or too small will be filtered out.
- Sequences that are too small or of poor quality will be filtered out.
- Sequences will be de-replicated: duplicates will be removed but the number of duplicates will be recorded.

FROGS was designed to support multiplexed and demultiplexed sequences ( Run FROGS Demultiplexing before Pre-process)



The goal of Flash (**Fast Length Adjustment of Short reads**) is to merge R1 and R2

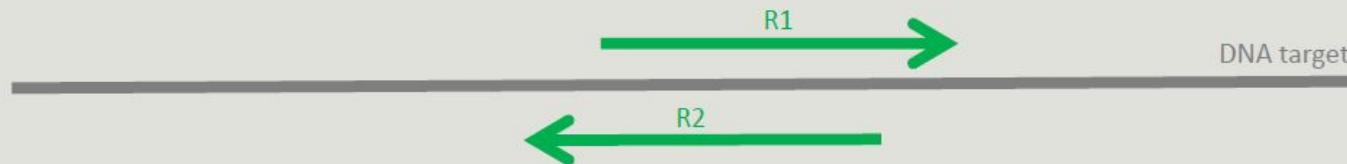
1<sup>st</sup> case: Impossible to merge



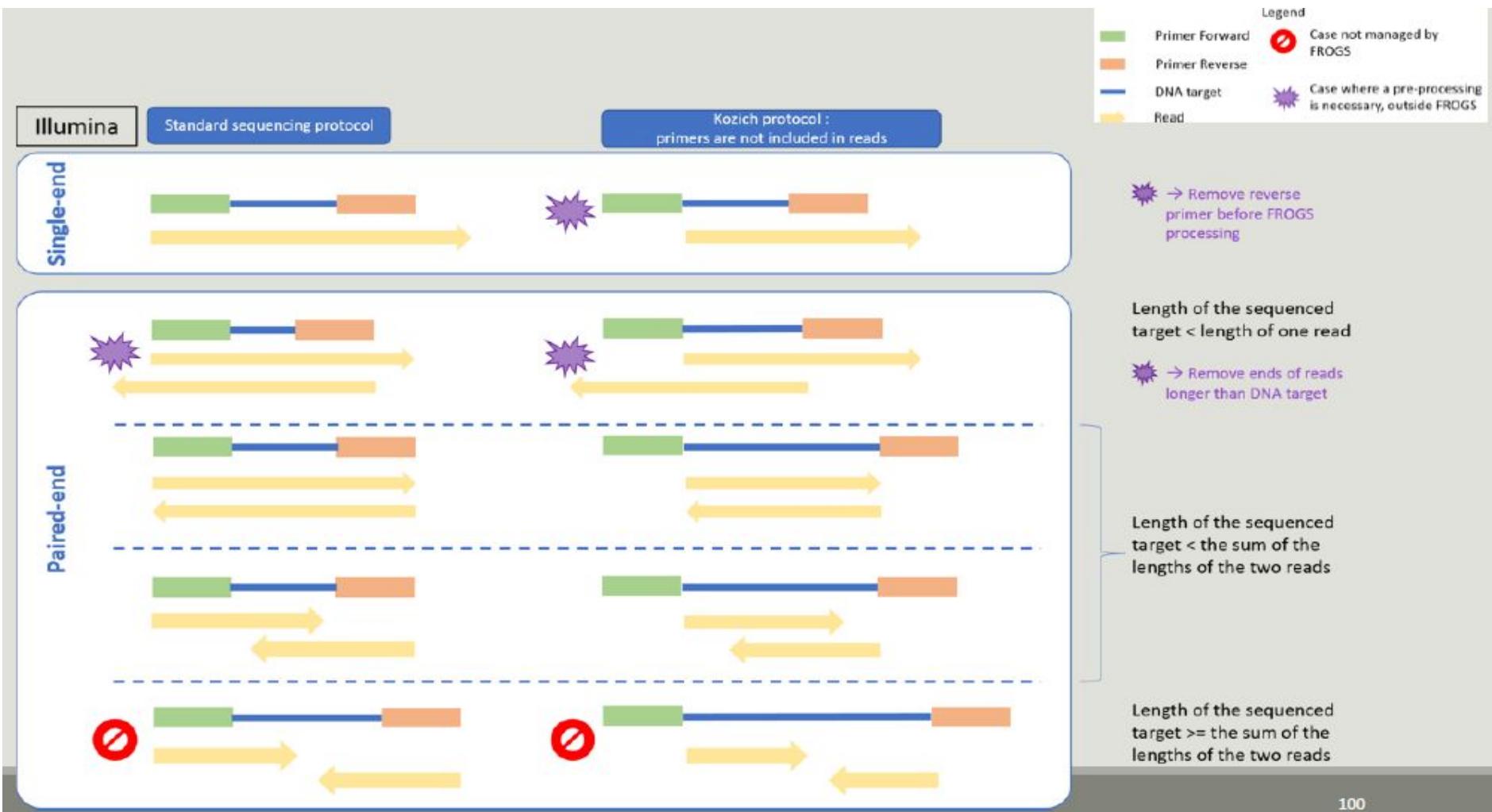
2<sup>nd</sup> case: flash have to find overlapping region between R1 and R2



3<sup>rd</sup> case: R1 and R2 cover entirely the target region



# Standard vs Kozich protocol

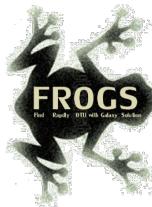


100

# Preprocess tool in bref

	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	soon
On-overlapping R1 R2	✗

	Take in charge
Archive .tar.gz	✓
Fastq	✓
Fasta	✗
With only 1 primer	✗
Multiplexed data	✗
Demultiplexed data	✓



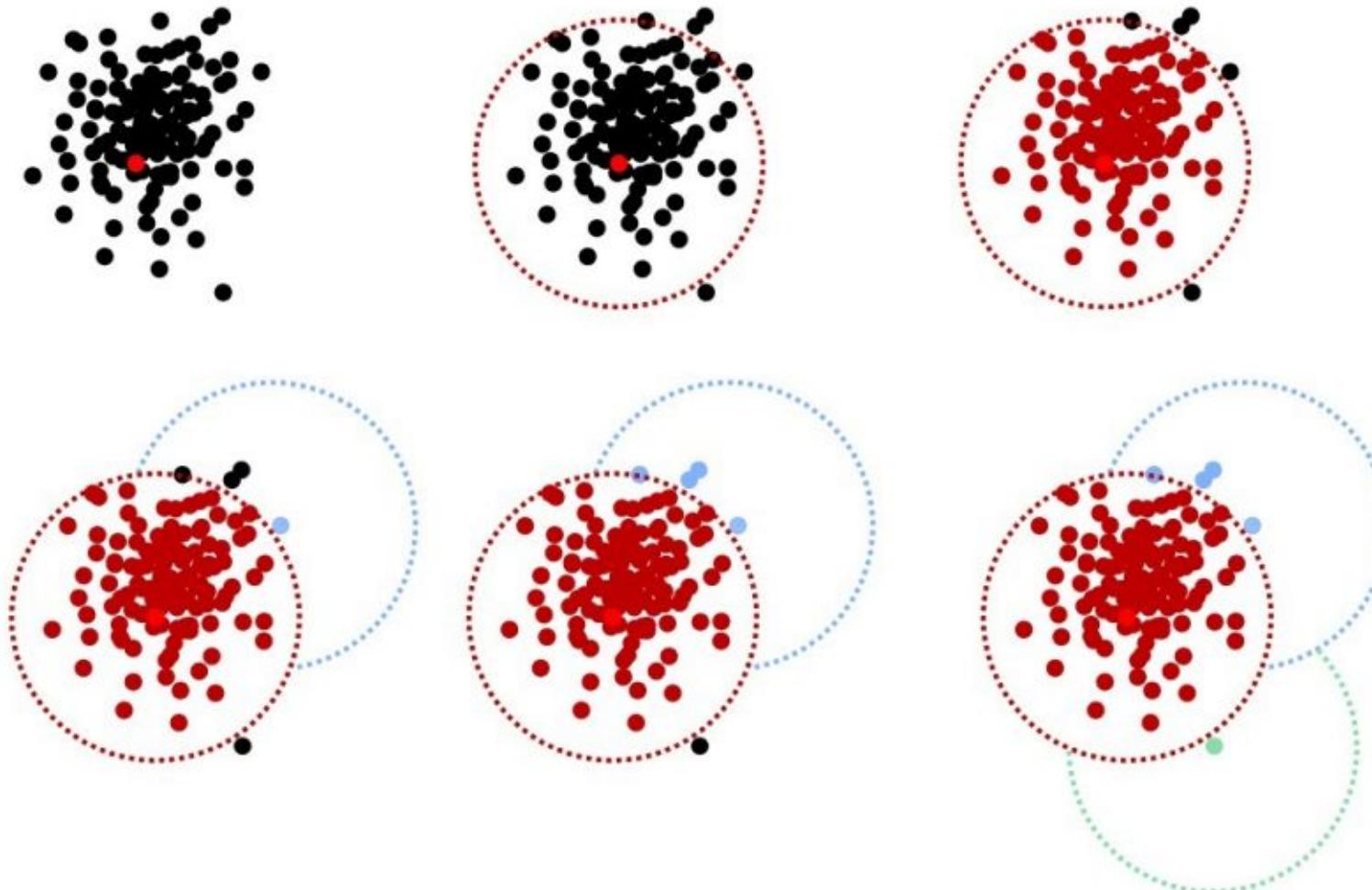
## 2. Clustering Swarm

In this step, sequences are clustered into groups using [Swarm](#). This takes the pre-processed fasta and counts files and does the following:

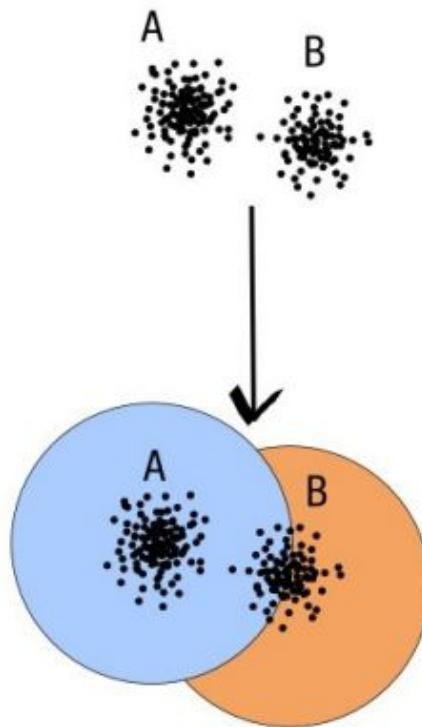
- Sorts reads by abundance.
- Clusters the reads into pre-clusters using Swarm and distance parameter of 1.
- Sorts these pre-clusters by abundance.
- Cluster the pre-clusters using Swarm and a user-specified distance.



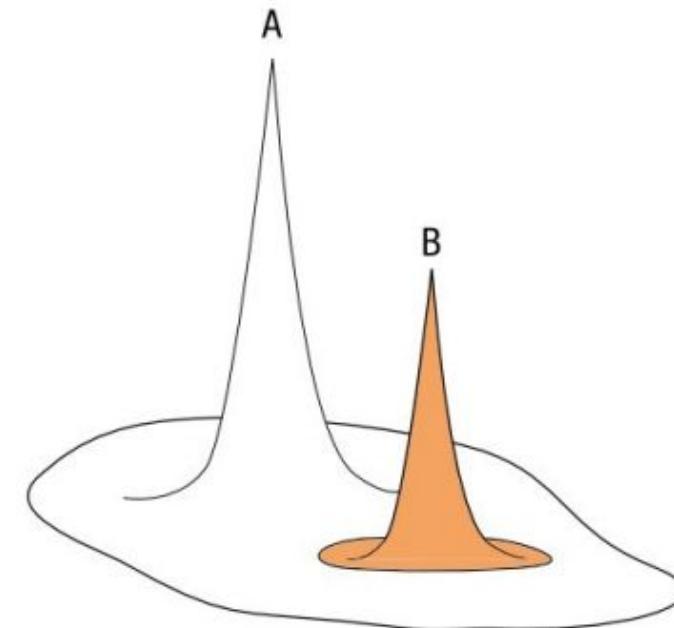
## How traditional clustering works?



# Swarm: fast, exact and high-resolution 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.

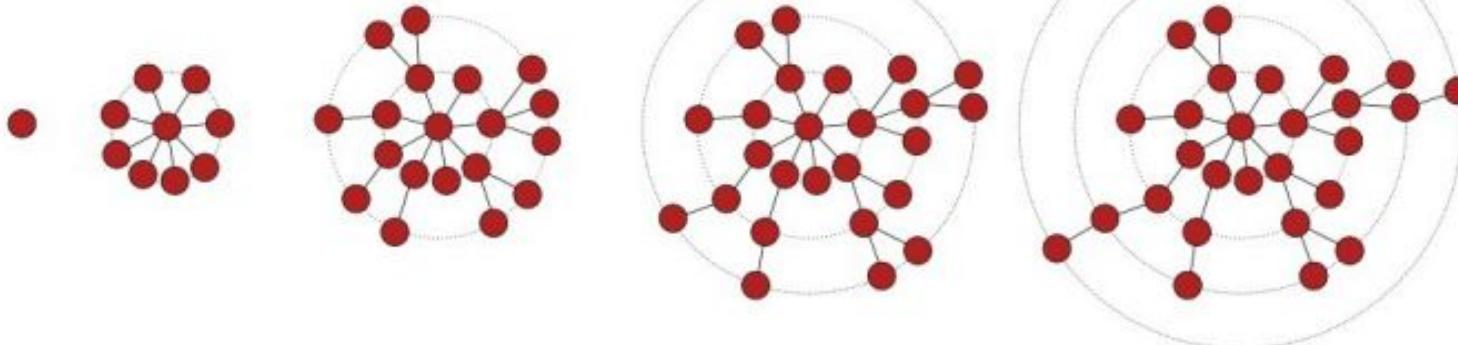
Swarm uses local clustering threshold, not a global clustering threshold

# Swarm clustering method

## growth phase

	ACGT	ACGT	ACGT	Avoid & speed-up comparisons
differences	1	1	2	- composition-based prefiltering - memoization - fast Needleman-Wunsch

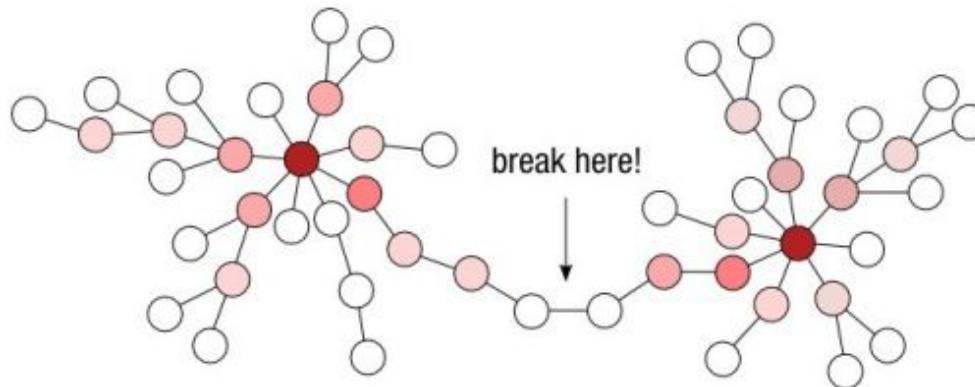
OTU grows iteratively



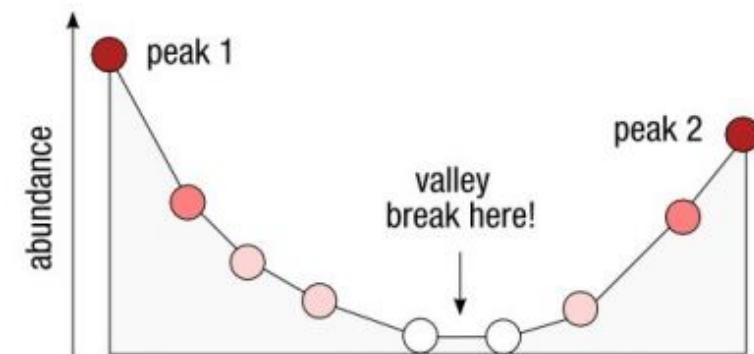
initial seed (randomly picked from amplicon dataset)

no more closely related amplicons,  
the process stops

## Swarm clustering method breaking phase



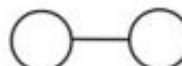
Take into account the abundance of amplicons  
to produce higher-resolution clusters.



Assuming that original sequences are more  
abundant than erroneous copies.

# Swarm clustering method grafting phase

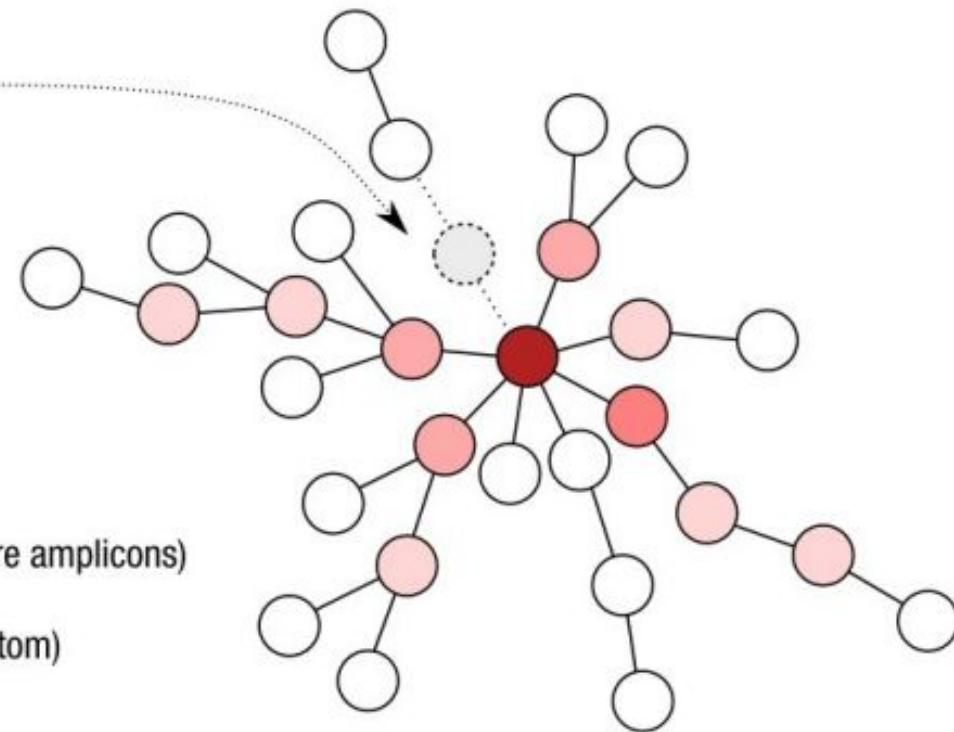
Postulate the existence of an intermediate amplicon to be able to graft a small OTU onto a bigger one.



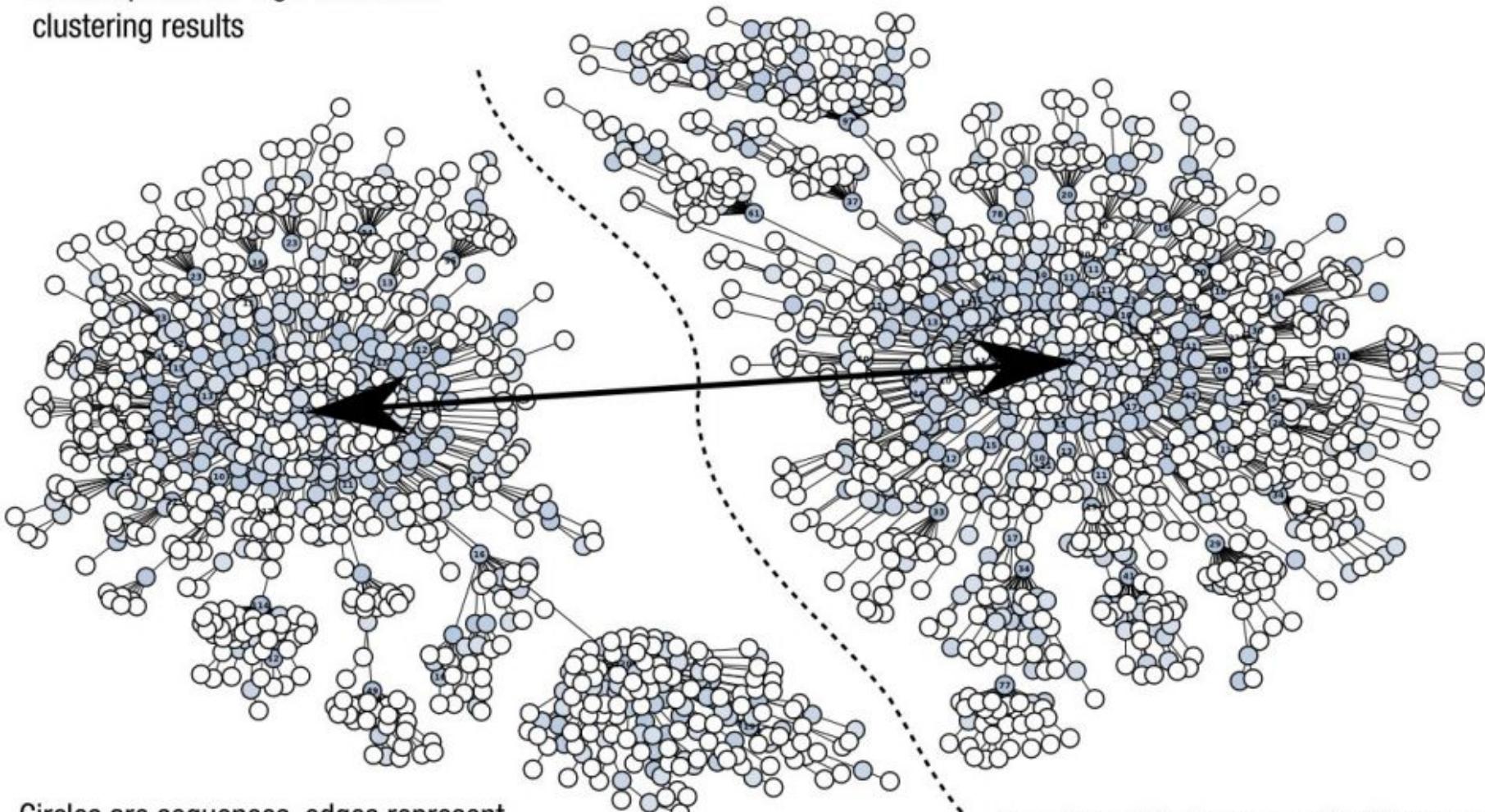
small OTU (made of 2 rare amplicons)



virtual amplicon (or phantom)



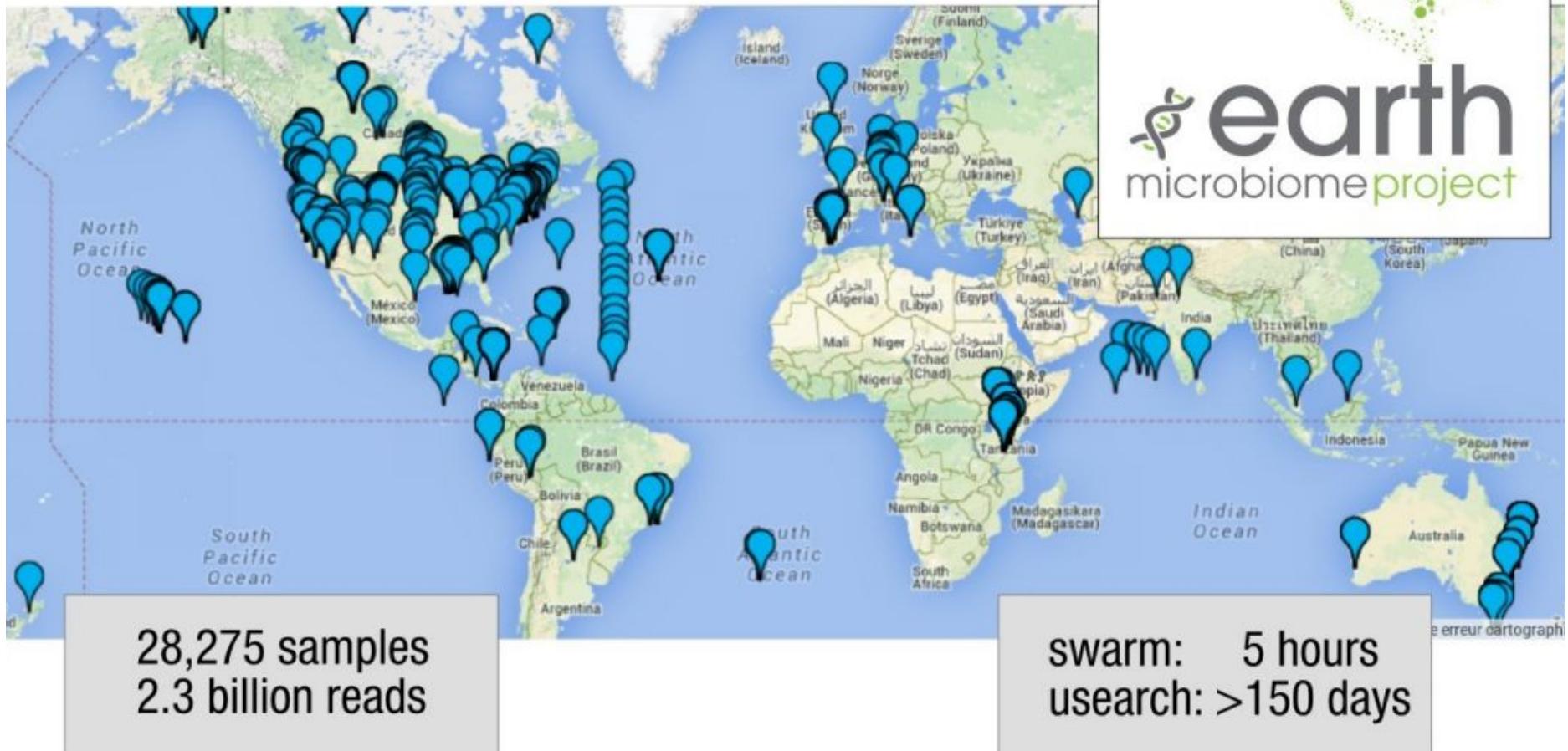
Swarm produces high-resolution clustering results



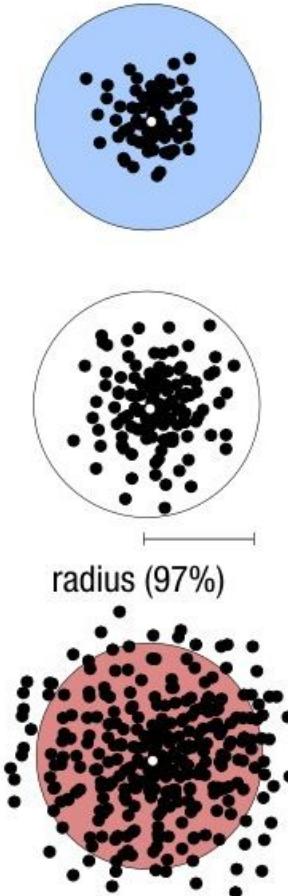
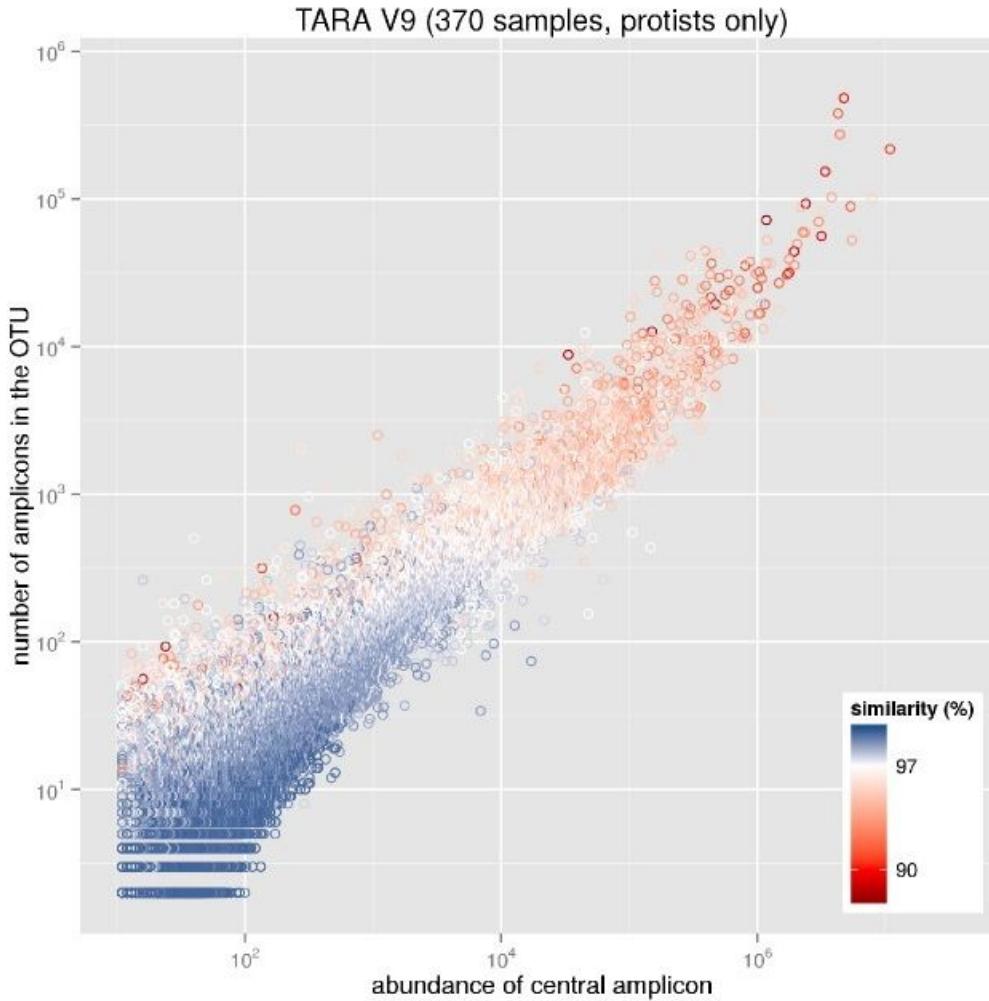
Circles are sequences, edges represent one difference (substitution or indel)

Less than 1% divergence (3 differences) between the two peaks of abundance

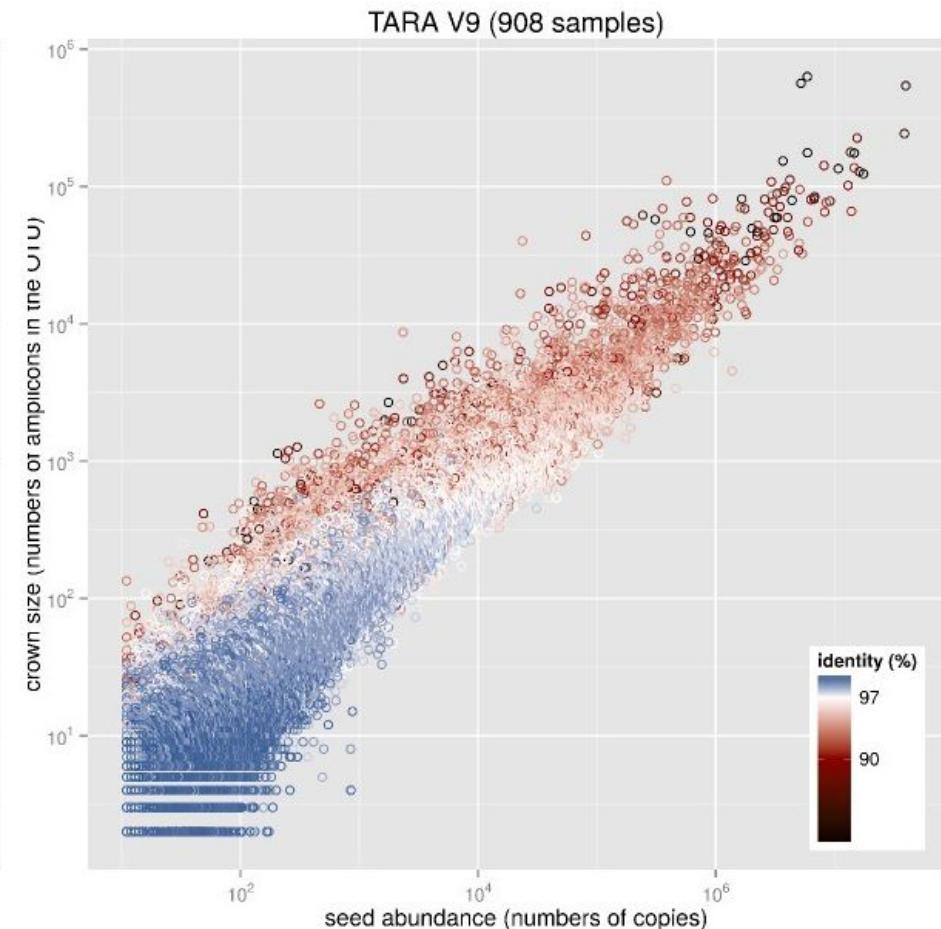
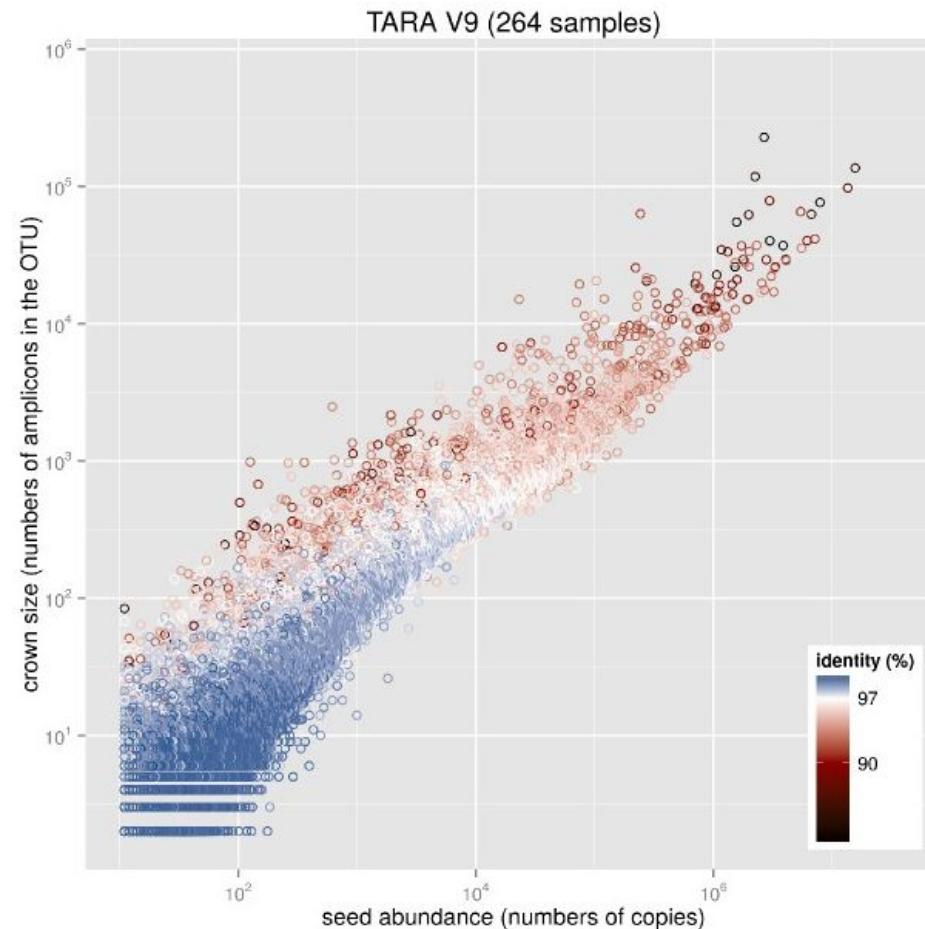
## Swarm 2.0 is a highly scalable denoising-clustering method



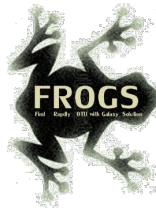
# What if we'd used a fix 97%-clustering threshold?



# Seed abundance vs cloud vs cluster radius shows 97%-threshold inadequacy



clusters produced with swarm using  $d = 1$



## 3. Remove 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 and reverse primer sequence at each end of the amplicon.

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

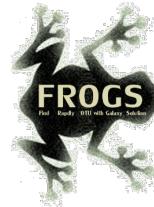
A: GTCGCTACTACCGATTGAACGTTTAGTGAGGTCTCGGACTGTGAGCCTGGCGGGTTG

|||||||||

B: TACTACCAAATGAGTTAGCGTTAGTGAGGT AAGACGACCAAACTGTAGCGTTAG

—————

C: GTCGCTACTACCGATTGAACGTTTAGTGAGGT AAGACGACCAAACTGTAGCGTTAG



## 3. Remove chimera

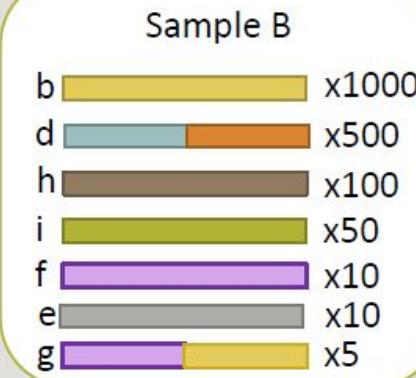
Closely-related sequences may form chimeras (mixed sequences) during PCR (library prep). This step removes these sequences by the following method:

- Splits input data into samples
- Uses **vsearch** to find chimeras in each sample
- Removes chimeras



### 3. Remove chimera

Chimera removal tool uses VSEARCH combined with an innovative chimera cross-validation.



“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

# vsearch: open-source alternative for usearch

clustering, chimera detection, dereplication, searching, sorting, masking and shuffling

**usearch** (Rob Edgar):

- very important for metagenomics,
- 1,000 citations,
- fundation for QIIME,
- closed-source & costly



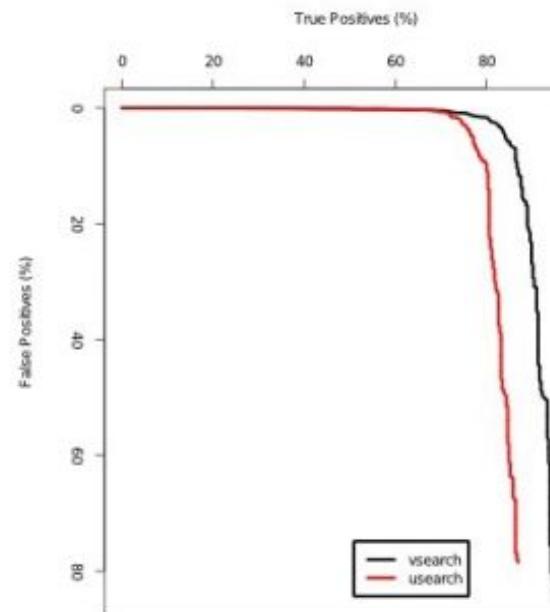
**growing success:**

- many happy users,
- faster and improved,
- fundation for QIIME 2.0

**vsearch:**

- free and open-source,
- fast,
- documented,
- revive the research field

Torbjørn Rognes  
Oslo University





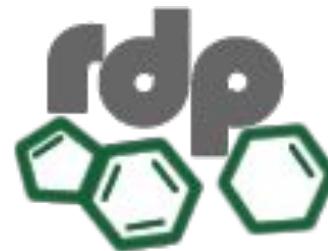
## 4. Filters

- The OTUs (Operational Taxonomic Units) have now been clustered. A filtering tool allows to remove noisy data. In this step, we will filter out some of the OTUs that are either not in at least 2 samples, and contain at least 2 sequences. Last allows eliminate singletons.
- Filters can be also done after affiliation taxonomy.



## 5. Affiliation OTU

- An OTU is a cluster of sequences. This step adds the taxonomy to the abundance file. It uses the SILVA database for rRNA.
- Affiliation tool returns taxonomic affiliation for each OTU using two methods with a unique multi-affiliation output





# Affiliation Strategy of FROGS

## Double Affiliation with :

1. RDPClassifiers
2. Blastn+ : all identical Best Hits with the tag “Multi-affiliation”.

V3 – V4	Bacteria   Firmicutes   Clostridia   Clostridiales   Lachnospiraceae   Pseudobutyribrio   16S unknown species
V3 – V4	Bacteria   Firmicutes   Clostridia   Clostridiales   Lachnospiraceae   Pseudobutyribrio   16S Butyribrio fibrisolvens
V3 – V4	Bacteria   Firmicutes   Clostridia   Clostridiales   Lachnospiraceae   Pseudobutyribrio   16S rumen bacterium 8   9293-9
V3 – V4	Bacteria   Firmicutes   Clostridia   Clostridiales   Lachnospiraceae   Pseudobutyribrio   16S Pseudobutyribrio xylovorans
V3 – V4	Bacteria   Firmicutes   Clostridia   Clostridiales   Lachnospiraceae   Pseudobutyribrio   16S Pseudobutyribrio ruminis



**FROGS Affiliation:** Bacteria | Firmicutes | Clostridia | Clostridiales | Lachnospiraceae | Pseudobutyribrio | **Multi-affiliation**

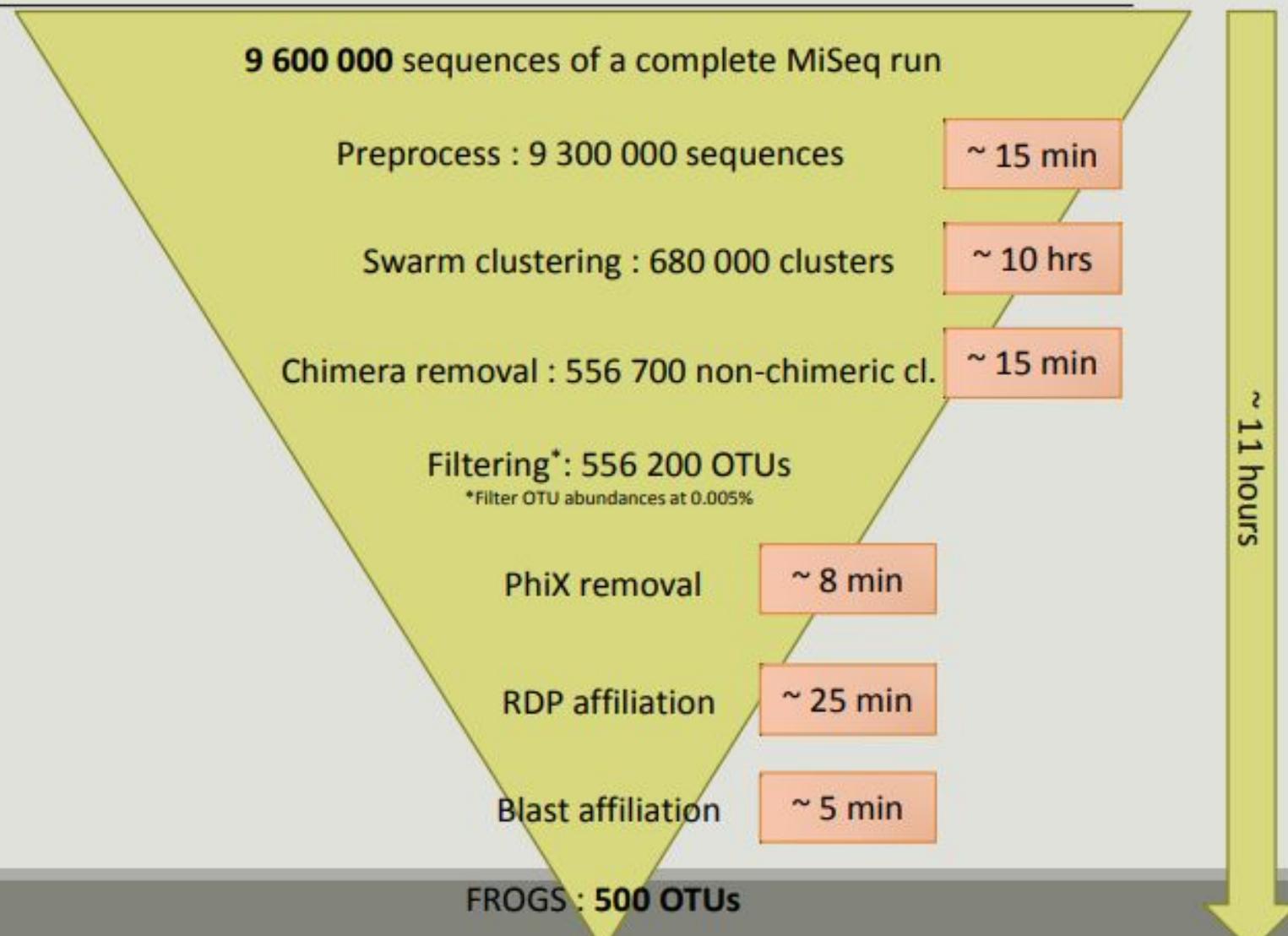
Steps	Description
1	<u>RDPClassifier</u> is used with database to associate to each OTU a taxonomy and a bootstrap (example: <i>Bacteria</i> ; (1.0); <i>Firmicutes</i> ; (1.0); <i>Clostridia</i> ; (1.0); <i>Clostridiales</i> ; (1.0); <i>Clostridiaceae</i> 1; (1.0); <i>Clostridium sensu stricto</i> ; (1.0);).
2	<u>blastn+</u> is used to find alignment between each OTU and the database. Only the best hits with the same score has reported.
3	For each OTU with several <u>blastn+</u> results a consensus is determined on each taxonomic level. If all the taxa in a taxonomic rank are identical the taxon name is reported otherwise <i>Multi-affiliation</i> is reported. By example, if you have an OTU with two corresponding sequences, the first is a <i>Bacteria</i> ; <i>Proteobacteria</i> ; <i>Gamma</i> <i>Proteobacteria</i> ; <i>Enterobacteriales</i> , the second is a <i>Bacteria</i> ; <i>Proteobacteria</i> ; <i>Beta Proteobacteria</i> ; <i>Methylophilales</i> , the consensus will be <i>Bacteria</i> ; <i>Proteobacteria</i> ; <i>Multi-affiliation</i> ; <i>Multi-affiliation</i> .



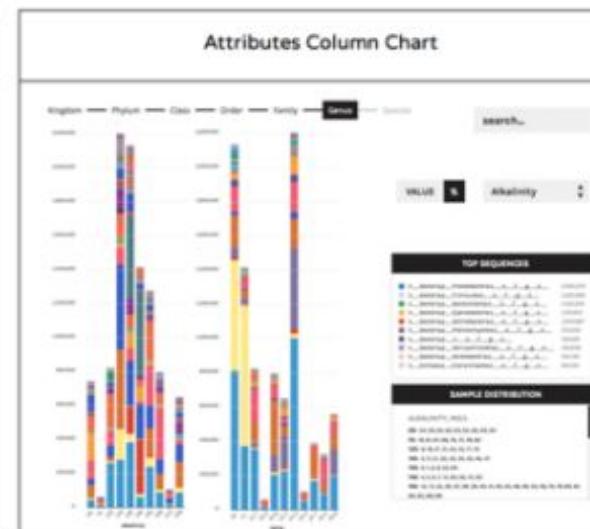
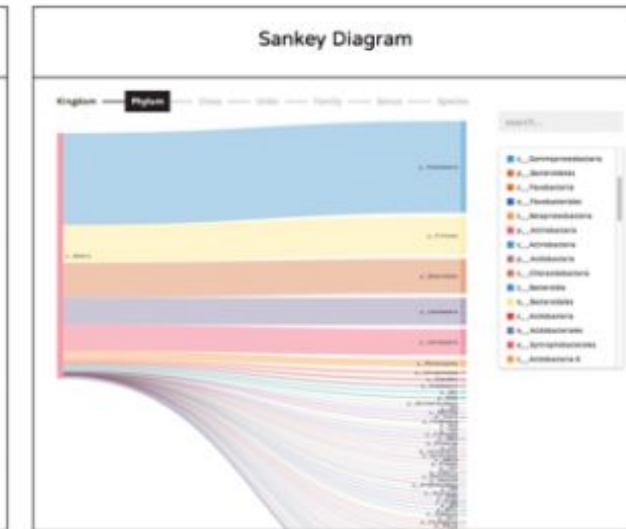
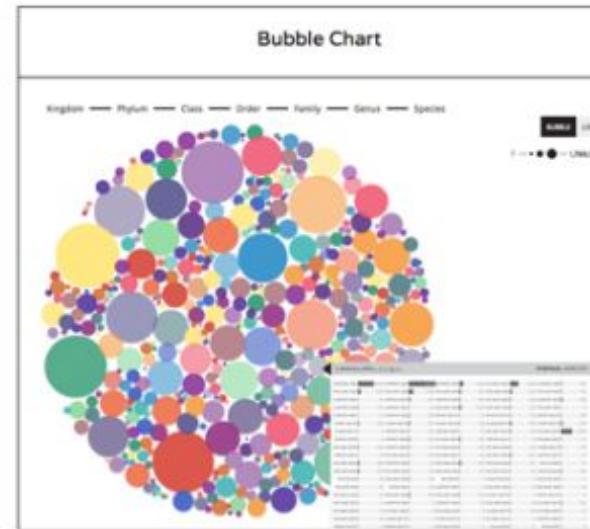
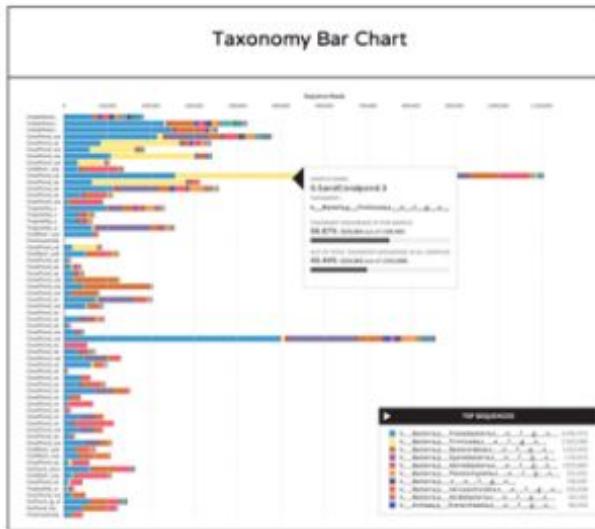
## 6. Affiliation Stats

This step computes some statistics from the analysis and generates a report of the OTUs/taxonomy found.

# Speed on real datasets



# Practice 2: Visualizing and plotting all sample results with Phinch



# Community analysis



- diversity indices and metrics
  - alpha-diversity
  - beta-diversity
- ordination techniques to visualize beta-diversity patterns
  - PCoA, NMDS
- common plots in microbiome papers

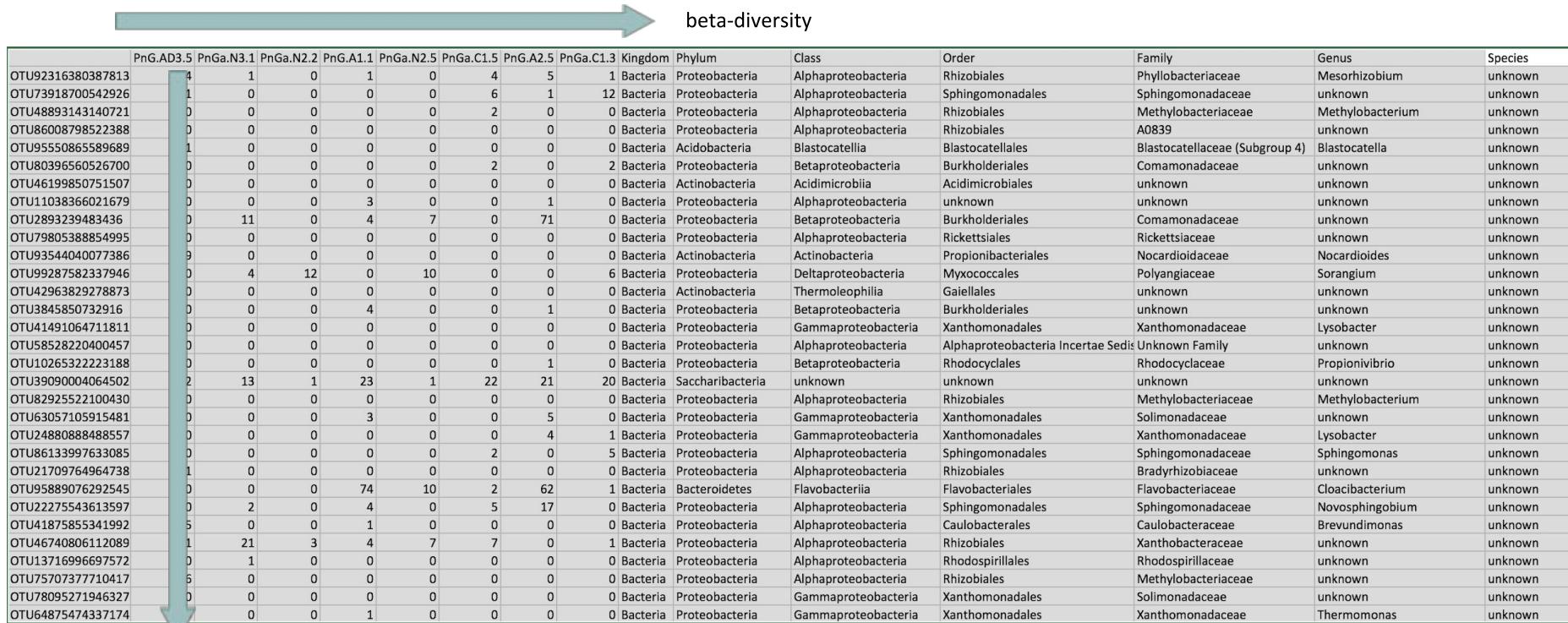
# Diversity indices and metrics

- alpha diversity vs beta-diversity

	PnG.AD3.5	PnGa.N3.1	PnGa.N2.2	PnG.A1.1	PnGa.N2.5	PnG.C1.5	PnG.A2.5	PnGa.C1.3	Kingdom	Phylum	Class	Order	Family	Genus	Species
OTU92316380387813	24	1	0	1	0	4	5	1	Bacteria	Proteobacteria	Alphaproteobacteria	Rhizobiales	Phyllobacteriaceae	Mesorhizobium	unknown
OTU73918700542926	101	0	0	0	0	6	1	12	Bacteria	Proteobacteria	Alphaproteobacteria	Sphingomonadales	Sphingomonadaceae	unknown	unknown
OTU48893143140721	0	0	0	0	0	2	0	0	Bacteria	Proteobacteria	Alphaproteobacteria	Rhizobiales	Methylobacteriaceae	Methylobacterium	unknown
OTU86008798522388	0	0	0	0	0	0	0	0	Bacteria	Proteobacteria	Alphaproteobacteria	Rhizobiales	A0839	unknown	unknown
OTU95550865589689	11	0	0	0	0	0	0	0	Bacteria	Acidobacteria	Blastocatellia	Blastocatellales	Blastocellaceae (Subgroup 4)	Blastocatella	unknown
OTU803965650526700	0	0	0	0	0	2	0	2	Bacteria	Proteobacteria	Betaproteobacteria	Burkholderiales	Comamonadaceae	unknown	unknown
OTU46199850751507	0	0	0	0	0	0	0	0	Bacteria	Actinobacteria	Acidimicrobia	Acidimicrobiales	unknown	unknown	unknown
OTU11038366021679	0	0	0	3	0	0	0	1	Bacteria	Proteobacteria	Alphaproteobacteria	unknown	unknown	unknown	unknown
OTU2893239483436	0	11	0	4	7	0	71	0	Bacteria	Proteobacteria	Betaproteobacteria	Burkholderiales	Comamonadaceae	unknown	unknown
OTU79805388854995	0	0	0	0	0	0	0	0	Bacteria	Proteobacteria	Alphaproteobacteria	Rickettsiales	Rickettsiaceae	unknown	unknown
OTU93544040077386	9	0	0	0	0	0	0	0	Bacteria	Actinobacteria	Actinobacteria	Propionibacterales	Nocardioidaceae	Nocardioides	unknown
OTU99287582337946	0	4	12	0	10	0	0	6	Bacteria	Proteobacteria	Deltaproteobacteria	Myxococcales	Polyangiaeae	Sorangium	unknown
OTU42963829278873	0	0	0	0	0	0	0	0	Bacteria	Actinobacteria	Thermoleophilia	Gaiellales	unknown	unknown	unknown
OTU3845850732916	0	0	0	4	0	0	0	1	Bacteria	Proteobacteria	Betaproteobacteria	Burkholderiales	unknown	unknown	unknown
OTU41491064711811	0	0	0	0	0	0	0	0	Bacteria	Proteobacteria	Gammaproteobacteria	Xanthomonadales	Xanthomonadaceae	Lysobacter	unknown
OTU58528220400457	0	0	0	0	0	0	0	0	Bacteria	Proteobacteria	Alphaproteobacteria	Alphaproteobacteria Incertae Sedis	Unknown Family	unknown	unknown
OTU10265322223188	0	0	0	0	0	0	0	1	Bacteria	Proteobacteria	Betaproteobacteria	Rhodocyclales	Rhodocyclaceae	Propionivibrio	unknown
OTU39090004064502	2	13	1	23	1	22	21	20	Bacteria	Saccharibacteria	unknown	unknown	unknown	unknown	unknown
OTU82925522100430	0	0	0	0	0	0	0	0	Bacteria	Proteobacteria	Alphaproteobacteria	Rhizobiales	Methylobacteriaceae	Methylobacterium	unknown
OTU63057105915481	0	0	0	3	0	0	5	5	Bacteria	Proteobacteria	Gammaproteobacteria	Xanthomonadales	Solimonadaceae	unknown	unknown
OTU24880888488557	0	0	0	0	0	0	4	1	Bacteria	Proteobacteria	Gammaproteobacteria	Xanthomonadales	Xanthomonadaceae	Lysobacter	unknown
OTU86133997633085	0	0	0	0	0	2	0	5	Bacteria	Proteobacteria	Alphaproteobacteria	Sphingomonadales	Sphingomonadaceae	Sphingomonas	unknown
OTU21709764964738	11	0	0	0	0	0	0	0	Bacteria	Proteobacteria	Alphaproteobacteria	Rhizobiales	Bradyrhizobiaceae	unknown	unknown
OTU95889076292545	0	0	0	74	10	2	62	1	Bacteria	Bacteroidetes	Flavobacteriia	Flavobacteriales	Flavobacteriaceae	Cloacibacterium	unknown
OTU22275543613597	0	2	0	4	0	5	17	0	Bacteria	Proteobacteria	Alphaproteobacteria	Sphingomonadales	Sphingomonadaceae	Novosphingobium	unknown
OTU41875855341992	105	0	0	1	0	0	0	0	Bacteria	Proteobacteria	Alphaproteobacteria	Caulobacterales	Caulobacteraceae	Brevundimonas	unknown
OTU46740806112089	61	21	3	4	7	7	0	1	Bacteria	Proteobacteria	Alphaproteobacteria	Rhizobiales	Xanthobacteraceae	unknown	unknown
OTU13716996697572	0	1	0	0	0	0	0	0	Bacteria	Proteobacteria	Alphaproteobacteria	Rhodospirillales	Rhodospirillaceae	unknown	unknown
OTU75707377710417	6	0	0	0	0	0	0	0	Bacteria	Proteobacteria	Alphaproteobacteria	Rhizobiales	Methylobacteriaceae	unknown	unknown
OTU78095271946327	0	0	0	0	0	0	0	0	Bacteria	Proteobacteria	Gammaproteobacteria	Xanthomonadales	Solimonadaceae	unknown	unknown
OTU64875474337174	0	0	0	1	0	0	0	0	Bacteria	Proteobacteria	Gammaproteobacteria	Xanthomonadales	Xanthomonadaceae	Thermomonas	unknown

# Diversity indices and metrics

- alpha diversity vs beta-diversity



The diagram illustrates the relationship between alpha diversity and beta diversity. A large green arrow points from the bottom left (labeled "alpha- diversity") towards the top right (labeled "beta-diversity"). The table below provides detailed data for various bacterial OTUs across different taxonomic levels.

	PnG.AD3.5	PnGa.N3.1	PnGa.N2.2	PnG.A1.1	PnGa.N2.5	PnG.C1.5	PnG.A2.5	PnG.C1.3	Kingdom	Phylum	Class	Order	Family	Genus	Species
OTU92316380387813	4	1	0	1	0	4	5	1	Bacteria	Proteobacteria	Alphaproteobacteria	Rhizobiales	Phyllobacteriaceae	Mesorhizobium	unknown
OTU73918700542926	1	0	0	0	0	6	1	12	Bacteria	Proteobacteria	Alphaproteobacteria	Sphingomonadales	Sphingomonadaceae	unknown	unknown
OTU48893143140721	0	0	0	0	0	2	0	0	Bacteria	Proteobacteria	Alphaproteobacteria	Rhizobiales	Methylobacteriaceae	Methylobacterium	unknown
OTU86008798522388	0	0	0	0	0	0	0	0	Bacteria	Proteobacteria	Alphaproteobacteria	Rhizobiales	A0839	unknown	unknown
OTU95550865589689	1	0	0	0	0	0	0	0	Bacteria	Acidobacteria	Blastocatellia	Blastocatellales	Blastocellaceae (Subgroup 4)	Blastocatella	unknown
OTU803965650526700	0	0	0	0	0	2	0	2	Bacteria	Proteobacteria	Betaproteobacteria	Burkholderiales	Comamonadaceae	unknown	unknown
OTU46199850751507	0	0	0	0	0	0	0	0	Bacteria	Actinobacteria	Acidimicrobia	Acidimicrobiales	unknown	unknown	unknown
OTU1038366021679	0	0	0	3	0	0	0	1	Bacteria	Proteobacteria	Alphaproteobacteria	unknown	unknown	unknown	unknown
OTU2893239483436	0	11	0	4	7	0	71	0	Bacteria	Proteobacteria	Betaproteobacteria	Burkholderiales	Comamonadaceae	unknown	unknown
OTU79805388854995	0	0	0	0	0	0	0	0	Bacteria	Proteobacteria	Alphaproteobacteria	Rickettsiales	Rickettsiaceae	unknown	unknown
OTU93544040077386	9	0	0	0	0	0	0	0	Bacteria	Actinobacteria	Actinobacteria	Propionibacterales	Nocardioidaceae	Nocardioides	unknown
OTU99287582337946	0	4	12	0	10	0	0	6	Bacteria	Proteobacteria	Deltaproteobacteria	Myxococcales	Polyangiaeae	Sorangium	unknown
OTU42963829278873	0	0	0	0	0	0	0	0	Bacteria	Actinobacteria	Thermoleophilia	Gaiellales	unknown	unknown	unknown
OTU3845850732916	0	0	0	4	0	0	0	1	Bacteria	Proteobacteria	Betaproteobacteria	Burkholderiales	unknown	unknown	unknown
OTU4491064711811	0	0	0	0	0	0	0	0	Bacteria	Proteobacteria	Gammaproteobacteria	Xanthomonadales	Xanthomonadaceae	Lysobacter	unknown
OTU58528220400457	0	0	0	0	0	0	0	0	Bacteria	Proteobacteria	Alphaproteobacteria	Alphaproteobacteria Incertae Sedis	Unknown Family	unknown	unknown
OTU10265322223188	0	0	0	0	0	0	0	1	Bacteria	Proteobacteria	Betaproteobacteria	Rhodocyclales	Rhodocyclaceae	Propionivibrio	unknown
OTU390900004064502	2	13	1	23	1	22	21	20	Bacteria	Saccharibacteria	unknown	unknown	unknown	unknown	unknown
OTU82925522100430	0	0	0	0	0	0	0	0	Bacteria	Proteobacteria	Alphaproteobacteria	Rhizobiales	Methylobacteriaceae	Methylobacterium	unknown
OTU63057105915481	0	0	0	3	0	0	5	0	Bacteria	Proteobacteria	Gammaproteobacteria	Xanthomonadales	Solimonadaceae	unknown	unknown
OTU24880888488557	0	0	0	0	0	0	4	1	Bacteria	Proteobacteria	Gammaproteobacteria	Xanthomonadales	Xanthomonadaceae	Lysobacter	unknown
OTU86133997633085	0	0	0	0	0	2	0	5	Bacteria	Proteobacteria	Alphaproteobacteria	Sphingomonadales	Sphingomonadaceae	Sphingomonas	unknown
OTU21709764964738	1	0	0	0	0	0	0	0	Bacteria	Proteobacteria	Alphaproteobacteria	Rhizobiales	Bradyrhizobiaceae	unknown	unknown
OTU95889076292545	0	0	0	74	10	2	62	1	Bacteria	Bacteroidetes	Flavobacteriales	Flavobacteriales	Cloacibacterium	unknown	unknown
OTU22275543613597	0	2	0	4	0	5	17	0	Bacteria	Proteobacteria	Alphaproteobacteria	Sphingomonadales	Sphingomonadaceae	Novosphingobium	unknown
OTU41875855341992	5	0	0	1	0	0	0	0	Bacteria	Proteobacteria	Alphaproteobacteria	Caulobacterales	Caulobacteraceae	Brevundimonas	unknown
OTU46740806112089	1	21	3	4	7	7	0	1	Bacteria	Proteobacteria	Alphaproteobacteria	Rhizobiales	Xanthobacteraceae	unknown	unknown
OTU13716996697572	0	1	0	0	0	0	0	0	Bacteria	Proteobacteria	Alphaproteobacteria	Rhodospirillales	Rhodospirillaceae	unknown	unknown
OTU75707377710417	6	0	0	0	0	0	0	0	Bacteria	Proteobacteria	Alphaproteobacteria	Rhizobiales	Methylobacteriaceae	unknown	unknown
OTU78095271946327	0	0	0	0	0	0	0	0	Bacteria	Proteobacteria	Gammaproteobacteria	Xanthomonadales	Solimonadaceae	unknown	unknown
OTU64875474337174	0	0	1	0	0	0	0	0	Bacteria	Proteobacteria	Gammaproteobacteria	Xanthomonadales	Xanthomonadaceae	Thermomonas	unknown

alpha- diversity

## Diversity indices and metrics

- alpha-diversity
  - within sample metric
  - different indices

Richness	Number of observed taxa
Chao	Richness + (estimated) number of unobserved species
Shannon	Evenness of the species abundance distribution
Inv-Simpson	Inverse probability that two sequences sampled at random come from the same species

## Diversity indices and metrics

- beta-diversity
  - among sample
  - pairwise distance matrix

	PnG.AD3.5	PnGa.N3.1	PnGa.N2.2	PnG.A1.1	PnGa.N2.5	PnGa.C1.5	PnG.A2.5	PnGa.C1.3
OTU92316380387813	24	1	0	1	0	4	5	1
OTU73918700542926	101	0	0	0	0	6	1	12
OTU48893143140721	0	0	0	0	0	2	0	0
OTU86008798522388	0	0	0	0	0	0	0	0
OTU95550865589689	11	0	0	0	0	0	0	0
OTU80396560526700	0	0	0	0	0	2	0	2
OTU46199850751507	0	0	0	0	0	0	0	0
OTU11038366021679	0	0	0	3	0	0	1	0
OTU2893239483436	0	11	0	4	7	0	71	0
OTU79805388854995	0	0	0	0	0	0	0	0
OTU93544004077386	9	0	0	0	0	0	0	0
OTU9287582337946	0	4	12	0	10	0	0	6
OTU42963829278873	0	0	0	0	0	0	0	0
OTU3845850732916	0	0	0	4	0			
OTU4149106471811	0	0	0	0	0			
OTUS58528220400457	0	0	0	0	0	0	0	0
OTU10265322223188	0	0	0	0	0	0	1	0
OTU39090004064502	2	13	1	23	1	22	21	20
OTU82925522100430	0	0	0	0	0	0	0	0
OTU63057105915481	0	0	0	3	0	0	5	0
OTU24880888488557	0	0	0	0	0	0	4	1
OTU86133997633085	0	0	0	0	0	2	0	5
OTU21709764964738	11	0	0	0	0	0	0	0
OTU95889076292545	0	0	0	74	10	2	62	1
OTU22275543613597	0	2	0	4	0	5	17	0
OTU41875855341992	105	0	0	1	0	0	0	0
OTU46740806112089	61	21	3	4	7	7	0	1
OTU13716996697572	0	1	0	0	0	0	0	0
OTU75707377710417	6	0	0	0	0	0	0	0
OTU78095271946327	0	0	0	0	0	0	0	0
OTU64875474337174	0	0	0	1	0	0	0	0

	PnG.AD3.5	PnGa.N3.1	PnGa.N2.2	PnG.A1.1	PnGa.N2.5	PnGa.C1.5	PnG.A2.5	PnGa.C1.3
PnG.AD3.5	0							
PnGa.N3.1	0.9498313	0						
PnGa.N2.2	0.8794906	0.0694657	0					
PnG.A1.1	0.6865317	0.9647475	0.9854593	0				
PnGa.N2.5	0.411009	0.0343057	0.7786044	0.7753815	0			
PnGa.C1.5	0.0540924	0.2975522	0.3981714	0.5222625	0.4122164	0		
PnG.A2.5	0.9867694	0.6644837	0.812524	0.1160494	0.0431276	0.6362839	0	
PnGa.C1.3	0.9568336	0.9218136	0.4782994	0.107668	0.3606226	0.4325389	0.5148707	0



Dimension reduction approaches to summarize the multidimensional matrix

## Diversity indices and metrics

- beta-diversity
  - among sample
  - pairwise distance matrix
  - different metrics
    - Bray-Curtis
    - Sorensen (presence/absence of OTU)
  - PCoA, NMDS or dendrogram to visualize beta-diversity
    - PCoA : like PCA but consider any distance matrix
    - NMDS : similar as PCoA but constrain all dimension in only 2d

# Diversity indices and metrics

- beta-diversity

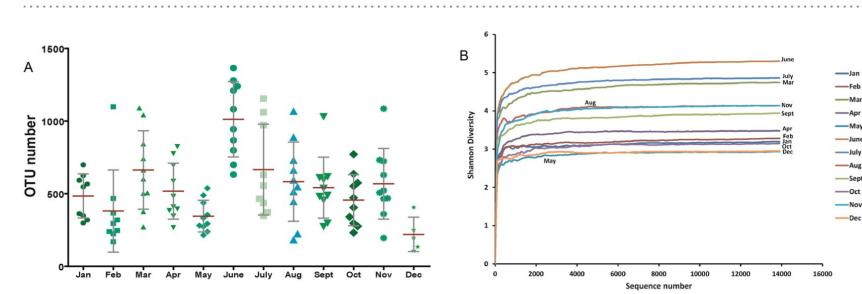
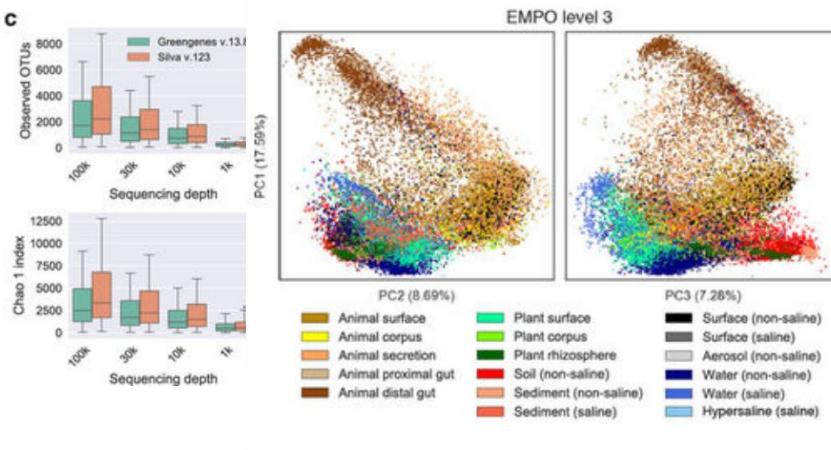
	PnG.AD3.5	PnGa.N3.1	PnGa.N2.2	PnG.A1.1	PnGa.N2.5	PnGa.C1.5	PnG.A2.5	PnGa.C1.3
OTU92316380387813	24	1	0	1	0	4	5	1
OTU73918700542926	101	0	0	0	0	6	1	12
OTU48893143140721	0	0	0	0	0	2	0	0
OTU8600879852388	0	0	0	0	0	0	0	0
OTU95550865589689	11	0	0	0	0	0	0	0
OTU80396560526700	0	0	0	0	0	2	0	2
OTU46199850751507	0	0	0	0	0	0	0	0
OTU11038366021679	0	0	0	3	0	0	1	0
OTU2893239483436	0	11	0	4	7	0	71	0
OTU79805388854995	0	0	0	0	0	0	0	0
OTU93544040077386	9	0	0	0	0	0	0	0
OTU9287582337946	0	4	12	0	10	0	0	6
OTU42963829278873	0	0	0	0	0	0	0	0
OTU3845850732916	0	0	0	4	0			
OTU41491064711811	0	0	0	0	0			
OTUS5852820400457	0	0	0	0	0	0	0	0
OTU10265322223188	0	0	0	0	0	0	1	0
OTU39090004064502	2	13	1	23	1	22	21	20
OTU82925522100430	0	0	0	0	0	0	0	0
OTU63057105915481	0	0	0	3	0	0	5	0
OTU2488088848857	0	0	0	0	0	0	4	1
OTU86133997633085	0	0	0	0	0	2	0	5
OTU21709764964738	11	0	0	0	0	0	0	0
OTU95889076292545	0	0	0	74	10	2	62	1
OTU22275543613597	0	2	0	4	0	5	17	0
OTU41875855341992	105	0	0	1	0	0	0	0
OTU46740806112089	61	21	3	4	7	7	0	1
OTU13716996697572	0	1	0	0	0	0	0	0
OTU75707377710417	6	0	0	0	0	0	0	0
OTU78095271946327	0	0	0	0	0	0	0	0
OTU64875474337174	0	0	0	1	0	0	0	0

	PnG.AD3.5	PnGa.N3.1	PnGa.N2.2	PnG.A1.1	PnGa.N2.5	PnGa.C1.5	PnG.A2.5	PnGa.C1.3
PnG.AD3.5	0							
PnGa.N3.1	0.9498313	0						
PnGa.N2.2	0.8794906	0.0694657	0					
PnG.A1.1	0.6865317	0.9647475	0.9854593	0				
PnGa.N2.5	0.411009	0.0343057	0.7786044	0.7753815	0			
PnGa.C1.5	0.0540924	0.2975522	0.3981714	0.5222625	0.4122164	0		
PnG.A2.5	0.9867694	0.6644837	0.812524	0.1160494	0.0431276	0.6362839	0	
PnGa.C1.3	0.9568336	0.9218136	0.4782994	0.107668	0.3606226	0.4325389	0.5148707	0



Dimension reduction approaches to summarize and visualize the multidimensional matrix

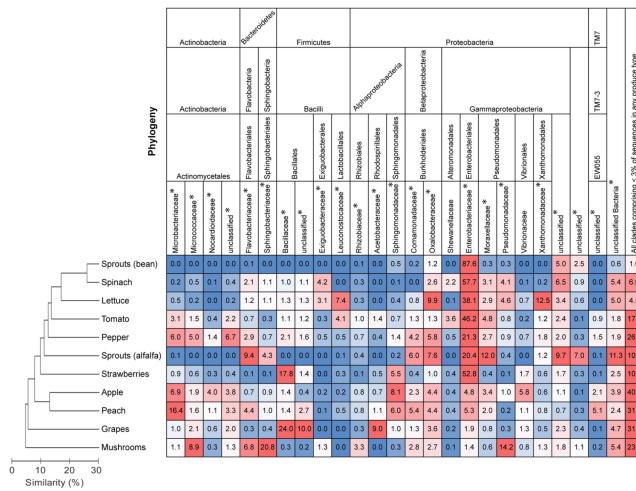
# Common plots in microbiome studies



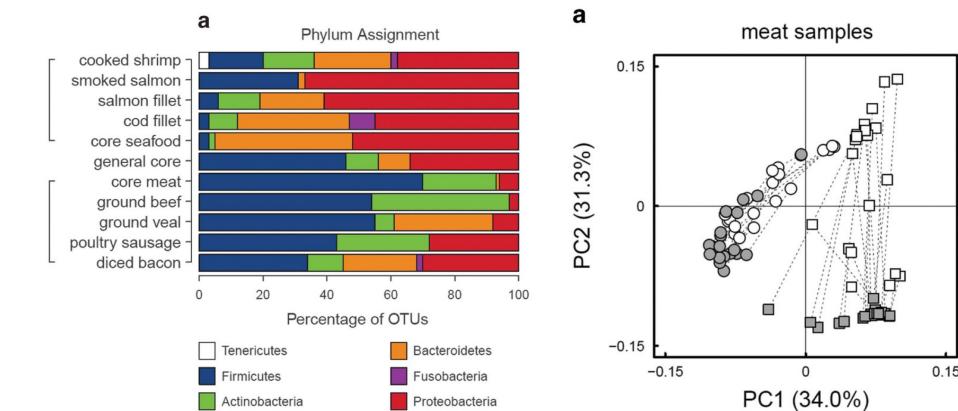
**Figure 1.** (A) Scatter plots showing OTU numbers with respect to total sequences obtained from 16S rRNA genes associated with months. The error bars were the mean with SD. (B) Shannon diversity curve of bacterial community derived from 12 months.

<http://dx.doi.org/10.1038/s41598-018-20862-8>

<https://www.nature.com/articles/nature24621>

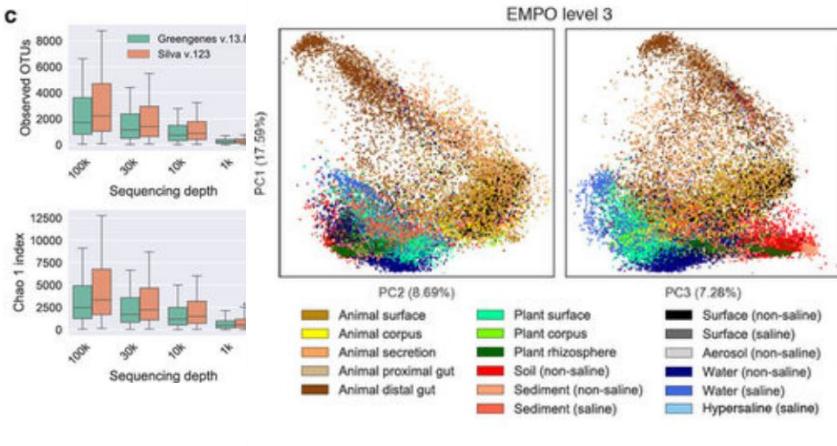


**Figure 2.** Relationships between bacterial communities on each produce type and relative abundances of bacterial families. The dendrogram is based on mean Bray-Curtis dissimilarities and shows differences among produce types in the overall composition of the bacterial communities. The heatmap shows mean relative abundances (%) of bacterial families on produce types. Only families and unclassified groupings representing at least three percent on any produce type are represented.

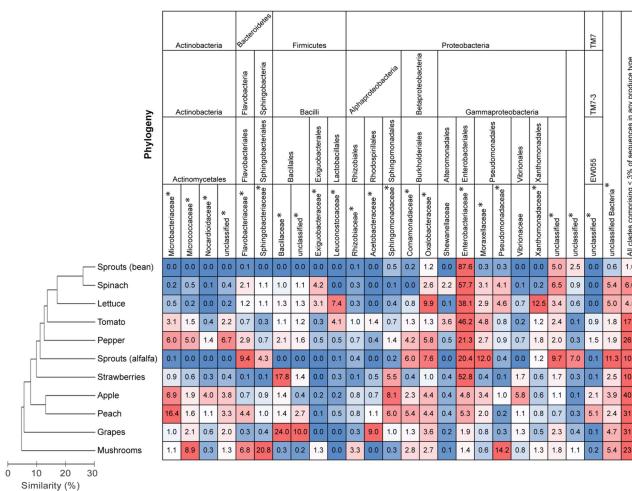


<http://www.nature.com/doifinder/10.1038/ismej.2014.202>

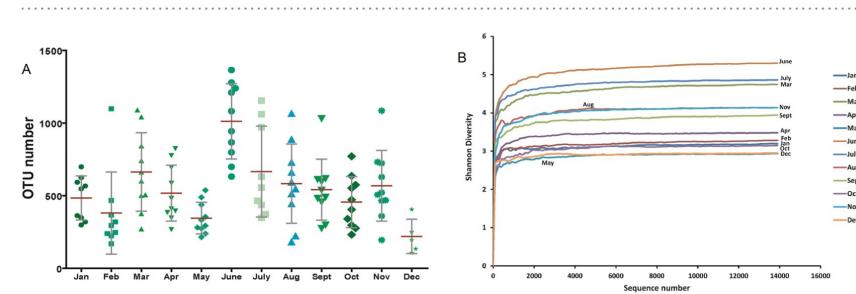
# Common plots in microbiome studies



<https://www.nature.com/articles/nature24621>

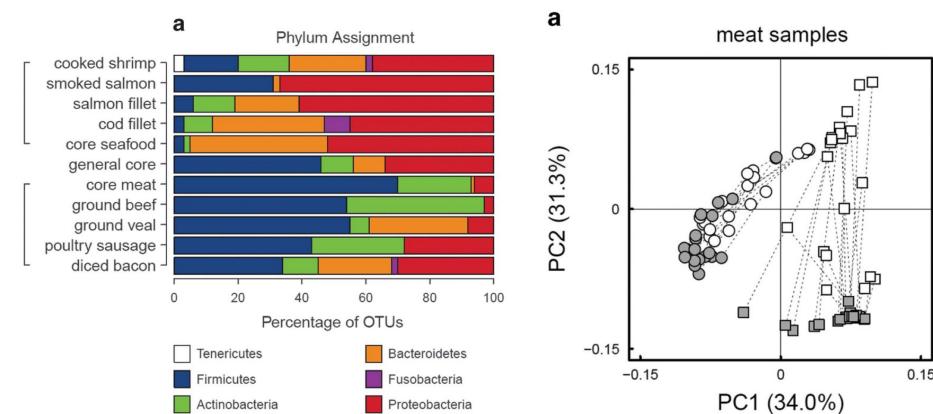


**Figure 2.** Relationships between bacterial communities on each produce type and relative abundances of bacterial families. The dendrogram is based on mean Bray-Curtis dissimilarities and shows differences among produce types in the overall composition of the bacterial communities. The heatmap shows mean relative abundances (%) of bacterial families on produce types. Only families and unclassified groupings representing at least three percent on any produce type are represented.  
doi:10.1371/journal.pone.0059310.g002



**Figure 1.** (A) Scatter plots showing OTU numbers with respect to total sequences obtained from 16S rRNA genes associated with months. The error bars were the mean with SD. (B) Shannon diversity curve of bacterial community derived from 12 months.

<http://dx.doi.org/10.1038/s41598-018-20862-8>



<http://www.nature.com/doifinder/10.1038/ismej.2014.202>

... to describe alpha/beta-diversity and taxonomy of the communities and link patterns to covariates

# *Practice 3: Handling and visualisation of OTU table using PhyloSeq*

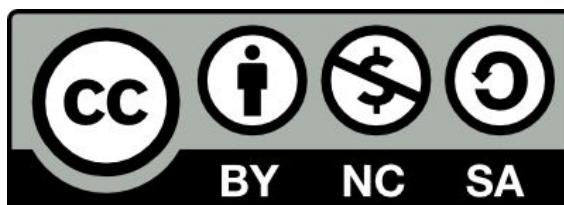
Hands-on session...

# Formateurs itrop / South Green + Collaborateurs UMR QualiSud CIRAD

- Alexis Dereeper
- Julie Orjuela-Bouniol
- Florentin Constancias



# Merci pour votre attention !



Le matériel pédagogique utilisé pour ces enseignements est mis à disposition selon les termes de la licence Creative Commons Attribution - Pas d'Utilisation Commerciale - Partage dans les Mêmes Conditions (BY-NC-SA) 4.0 International:

<http://creativecommons.org/licenses/by-nc-sa/4.0/>