



Requirements: Steps to process 16S data

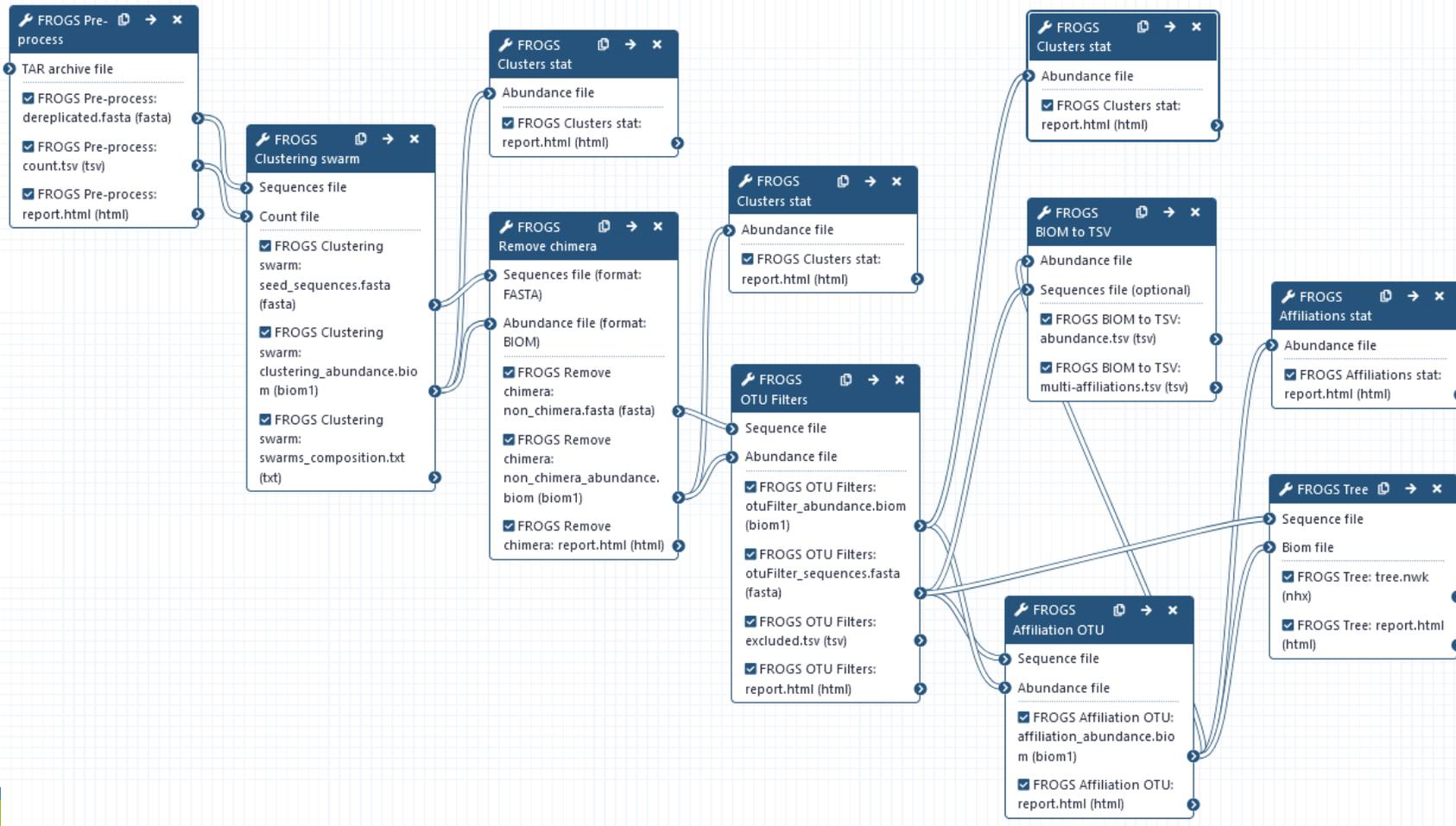
June 2022 - Webinar

FROGS Practice

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



Bioinformatics process



Acquisition des données

- Vous devez créer un nouvel historique nommé « 16S »
- Vous devez y mettre grâce à l'outil getdata de Galaxy le jeu de données nommé « chaillou_withprimers_64renamedsamples_V1V3_10000seq_R1R2.tar.gz »
- Dans getdata - **Upload File from your computer**

1. Dans la fenêtre cliquez sur 

2. Choisir « tar » dans le type de fichier



A small rectangular window titled "Type" with a dropdown menu. The dropdown menu is open, showing the word "tar" selected.

3. Et coller cette url:

http://genoweb.toulouse.inra.fr/~formation/15_FROGS/Webinar_data/chaillou_withprimers_64renamedsamples_V1V3_10000seq_R1R2.tar.gz

4. Puis 

Acquisition des données

Download from web or upload from disk

Regular Composite Collection Rule-based

You added 1 file(s) to the queue. Add more files or click 'Start' to proceed.

Name	Size	Type	Genome	Settings	Status
<input type="text" value="New File"/>	128 b	tar	unspecified (?)		0%

Download data from the web by entering URLs (one per line) or directly paste content.

```
http://genoweb.toulouse.inra.fr/~formation/15_FROGS/Webinar_data/chaillou_withprimers_64renamedsamples_V1V3_10C
```

Type (set all): Auto-detect Genome (set all): unspecified (?)

16S dataset presentation:

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

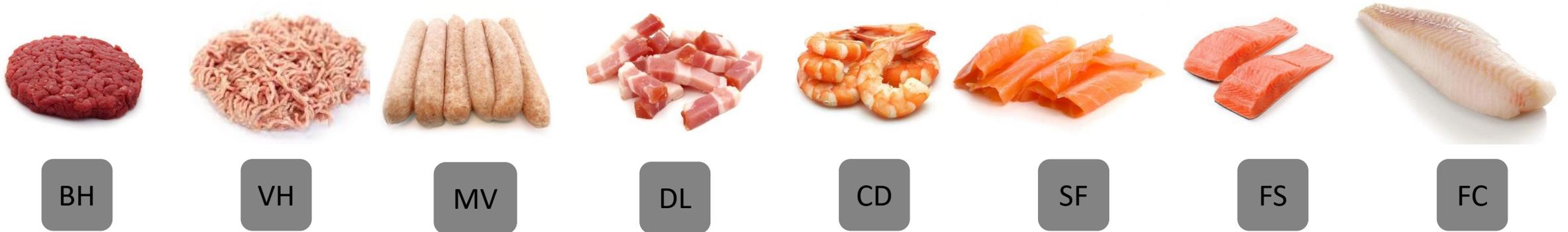
Comparison of meat and seafood bacterial communities.

8 environment types (EnvType) :

- Meat → Ground Beef, Ground veal, Poultry sausage, Diced bacon
- Seafood → Cooked schrimps, Smoked salmon, Salmon filet, Cod filet



16S dataset presentation:



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.

Les 11 étapes FROGS à lancer avant la formation

1. Preprocess (nettoyage des données)
2. Clustering (regroupement des séquences par similarité)
3. Cluster_stat (outil compagnon de visualisation des résultats)
4. Remove Chimera (suppression des chimères)
5. Cluster Stat
6. OTU Filter (suppression des cluster de trop faibles abondances ou pas assez fréquemment rencontrés)
7. Cluster Stat
8. Affiliation OTU (affiliation des OTUs à une taxonomie)
9. Affiliation Stat (outil compagnon de visualisation des résultats)
10. BIOM to TSV (transformation des données en une table d'abondance)
11. FROGS Tree (reconstruction de l'arbre phylogénétique des OTUs)

Pre-process tool

Les paramètres d'entrée

Input Parameter	Value
Sequencer	illumina
Input type	archive
TAR archive file	2 chaillou_withprimers_64renamedsamples_V1V3_10000seq_R1R2.tar.gz (Purged)
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

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

   2: chaillou_withprimers_64renamedsamples_V1V3_10000seq_R1R2.tar.gz

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 (--mismatch-rate)

Merge software

Vsearch

Select the software to merge paired-end reads (--merge-software)

Would you like to keep unmerged reads?

No

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

Minimum amplicon size

The minimum size for the amplicons (with primers) (--min-amplicon-size)

Maximum amplicon size

The maximum size for the amplicons (with primers) (--max-amplicon-size)

Sequencing protocol

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

5' primer

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

3' primer

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

Email notification

Send an email notification when the job completes.

✓ Execute

Ce que vous devez obtenir dans le report.html

Preprocess summary



Details on merged sequences

Show entries

Samples	before process	% kept	paired-end assembled	with 5' primer	with 3' primer	with expected length	without N
BHTO.LOT01	9,282	97.92	9,089	9,089	9,089	9,089	9,089
BHTO.LOT03	9,173	97.83	8,984	8,984	8,984	8,974	8,974
BHTO.LOT04	9,171	97.79	8,969	8,969	8,968	8,968	8,968
BHTO.LOT05	9,109	97.56	8,890	8,890	8,888	8,887	8,887
BHTO.LOT06	9,193	97.86	8,996	8,996	8,996	8,996	8,996

Si vous n'obtenez pas cela, vérifiez vos paramètres et en particulier faites bien un copier coller pour rentrer les primers.

Sinon, contactez-nous à frogs-training@inrae.fr

Si vous obtenez les mêmes chiffres, veuillez passer à l'étape suivante.

Clustering tool

Les paramètres d'entrée

Input Parameter	Value
Sequences file	3 FROGS Pre-process: dereplicated.fasta
Count file	4 FROGS Pre-process: count.tsv
FROGS guidelines version	3.2
Aggregation distance clustering	1
Refine OTU clustering	True

 **FROGS Clustering swarm** Single-linkage clustering on sequences (Galaxy Version 4.0.0+galaxy1)

Sequences file

   15: FROGS OTU Filters: otuFilter_sequences.fasta

The dereplicated sequences file (format: FASTA)

Count file

   22: FROGS BIOM to TSV: multi-affiliations.tsv

It contains the count by sample for each sequence (format: TSV)

FROGS guidelines version

New guidelines from version 3.2

The denoising step before a d3 clustering is no longer recommended since FROGS 3.2, but you can still choose it.

Aggregation distance clustering



Maximum number of differences between sequences in each aggregation swarm step. (recommended d=1) (--distance)

Refine OTU clustering

Yes

Clustering will be performed with the swarm --fastidious option. It is recommended and only usable in association with a distance of 1 (default and recommended: Yes) (--fastidious)

Email notification

Send an email notification when the job completes.

 **Execute**

Cluster stat tool

Les paramètres d'entrée

Input Parameter	Value
Abundance file	7 FROGS Clustering swarm: clustering_abundance.biom

 **FROGS Clusters stat** Process some metrics on clusters (Galaxy Version 4.0.0+galaxy1)

Abundance file

   18: FROGS Affiliation OTU: affiliation_abundance.biom

Clusters abundance (format: BIOM)

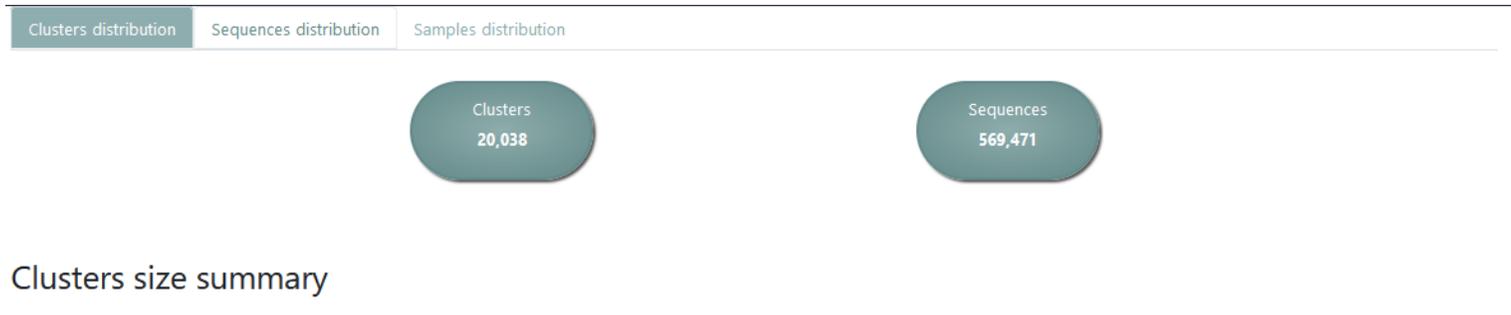
Email notification



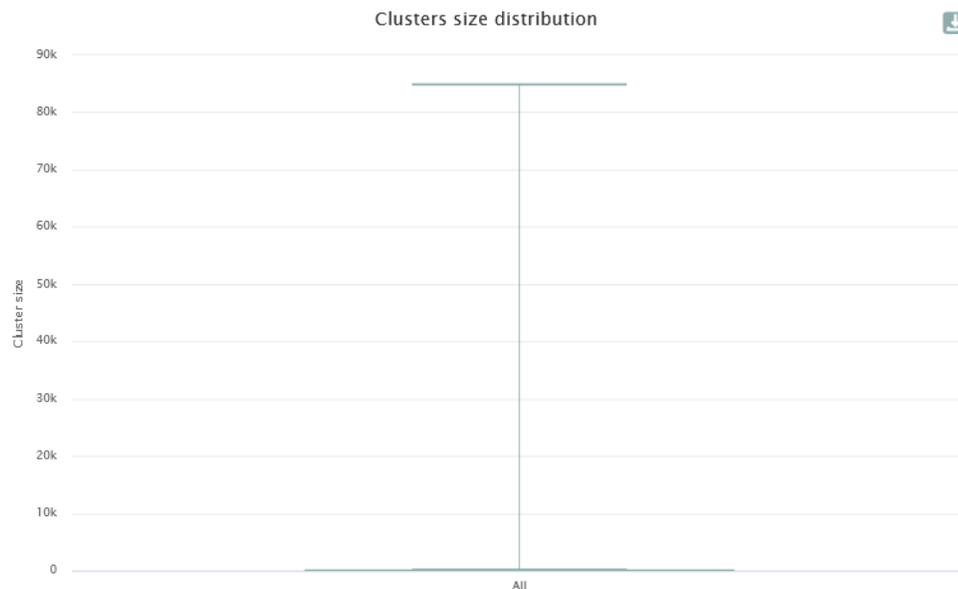
Send an email notification when the job completes.

 **Execute**

Ce que vous devez obtenir dans le report.html



The cluster size is the sum of the abundances of the sequences grouped in a cluster.



Decile	Value
Min	1
1	1
2	1
3	1
4	1
Median	1
6	1
7	1
8	1
9	1
Max	84,849

Si vous n'obtenez pas cela, vérifiez vos paramètres et seulement après contactez-nous à

frogs-training@inrae.fr

Si vous obtenez les mêmes chiffres, veuillez passer à l'étape suivante.

Ce que vous devez obtenir dans le report.html

Clusters size details

Show entries Search: [CSV](#)

Cluster size	Number of cluster	% of all clusters
1	19,267	96.15
2	150	0.75
3	22	0.11
4	10	0.05
5	8	0.04
6	15	0.07
7	14	0.07
8	8	0.04
9	10	0.05
10	5	0.02

Ce que vous devez obtenir dans le report.html

Clusters distribution Sequences distribution **Samples distribution**

Sequences count

Show entries

 CSV

Search:

Sample	Total clusters	Shared clusters	Own clusters	Total sequences	Shared sequences	Own sequences
BHT0.LOT01	493	114	379	9,089	8,709	380
BHT0.LOT03	433	140	293	8,974	8,679	295
BHT0.LOT04	474	152	322	8,968	8,646	322
BHT0.LOT05	475	152	323	8,887	8,564	323
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	493	186	307	8,937	8,630	307
CDT0.LOT02	596	95	501	9,270	8,767	503
CDT0.LOT04	465	162	303	8,918	8,609	309

Showing 1 to 10 of 64 entries

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

Chimera removal tool

Les paramètres d'entrée

Input Parameter	Value
Sequences file (format: FASTA)	6 FROGS Clustering swarm: seed_sequences.fasta
Abundance type	biom
Abundance file (format: BIOM)	7 FROGS Clustering swarm: clustering_abundance.biom

 **FROGS Remove chimera** Remove PCR chimera in each sample (Galaxy Version 4.0.0+galaxy1)

Sequences file (format: FASTA)

   15: FROGS OTU Filters: otuFilter_sequences.fasta

The sequences file

Abundance type

BIOM file

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

Abundance file (format: BIOM)

   18: FROGS Affiliation OTU: affiliation_abundance.biom

It contains the count by sample for each sequence.

Email notification

Send an email notification when the job completes.

 **Execute**

Ce que vous devez obtenir dans le report.html

Remove summary



Si vous n'obtenez pas cela, vérifiez vos paramètres et seulement après contactez-nous à

frogs-training@inrae.fr

Si vous obtenez les même chiffres, veuillez passer à l'étape suivante.

Chimera detection by sample

Chimera are first detected by sample, and finally only clusters always detected as chimera in all samples including thoses clusters are removed.

Show entries Search: [CSV](#)

Sample	Clusters kept	% Clusters kept	Cluster abundance 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
BHT0.LOT01	173	35.09	8,767	96.46	320	322	2	325	420	73
BHT0.LOT03	199	45.96	8,738	97.37	234	236	2	235	240	4

Pensez a relancer un cluster_Stat

Comparer avec le report.html du précédent cluster stat.

Que pensez-vous des chiffres contenus dans la colonne « own clusters » de la table de l'onglet « Samples distribution » ?

Est-ce attendu ?



Clusters distribution Sequences distribution **Samples distribution**

Sequences count

Show entries Search: [CSV](#)

Sample	Total clusters	Shared clusters	Own clusters	Total sequences	Shared sequences	Own sequences
BHT0.LOT01	173	103	70	8,767	8,697	70
BHT0.LOT03	199	139	60	8,738	8,678	60
BHT0.LOT04	230	151	79	8,724	8,645	79
BHT0.LOT05	221	145	76	8,630	8,554	76
BHT0.LOT06	213	150	63	8,717	8,654	63
BHT0.LOT07	225	152	73	8,746	8,673	73
BHT0.LOT08	254	198	56	8,539	8,483	56
BHT0.LOT10	253	173	80	8,694	8,614	80
CDT0.LOT02	240	94	146	8,912	8,766	146
CDT0.LOT04	246	162	84	8,698	8,609	89

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.

Les paramètres d'entrée

Input Parameter	Value
Sequence file	10 FROGS Remove chimera: non_chimera.fasta
Abundance file	11 FROGS Remove chimera: non_chimera_abundance.biom
Minimum prevalence method	all
Minimum prevalence	4
Minimum OTU abundance as proportion or count. We recommend to use a proportion of 0.00005.	proportion
Minimum proportion of sequences abundance to keep OTU	5e-05
N biggest OTUs	Not available.
Search for contaminant OTU.	no

Sequence file

   15: FROGS OTU Filters: otuFilter_sequences.fasta

The sequence file to filter (format: FASTA)

Abundance file

   18: FROGS Affiliation OTU: affiliation_abundance.biom

The abundance file to filter (format: BIOM)

Minimum prevalence method

all samples

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 abundance as proportion or count. We recommend to use a proportion of 0.00005.

as proportion

Minimum proportion of sequences abundance to keep OTU

5e-05

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 (--min_abundance)

N biggest OTUs

Fill the fields only if you want this treatment. Keep the N biggest OTU (--nb-biggest-otu)

Search for contaminant OTU.

No contaminant filter

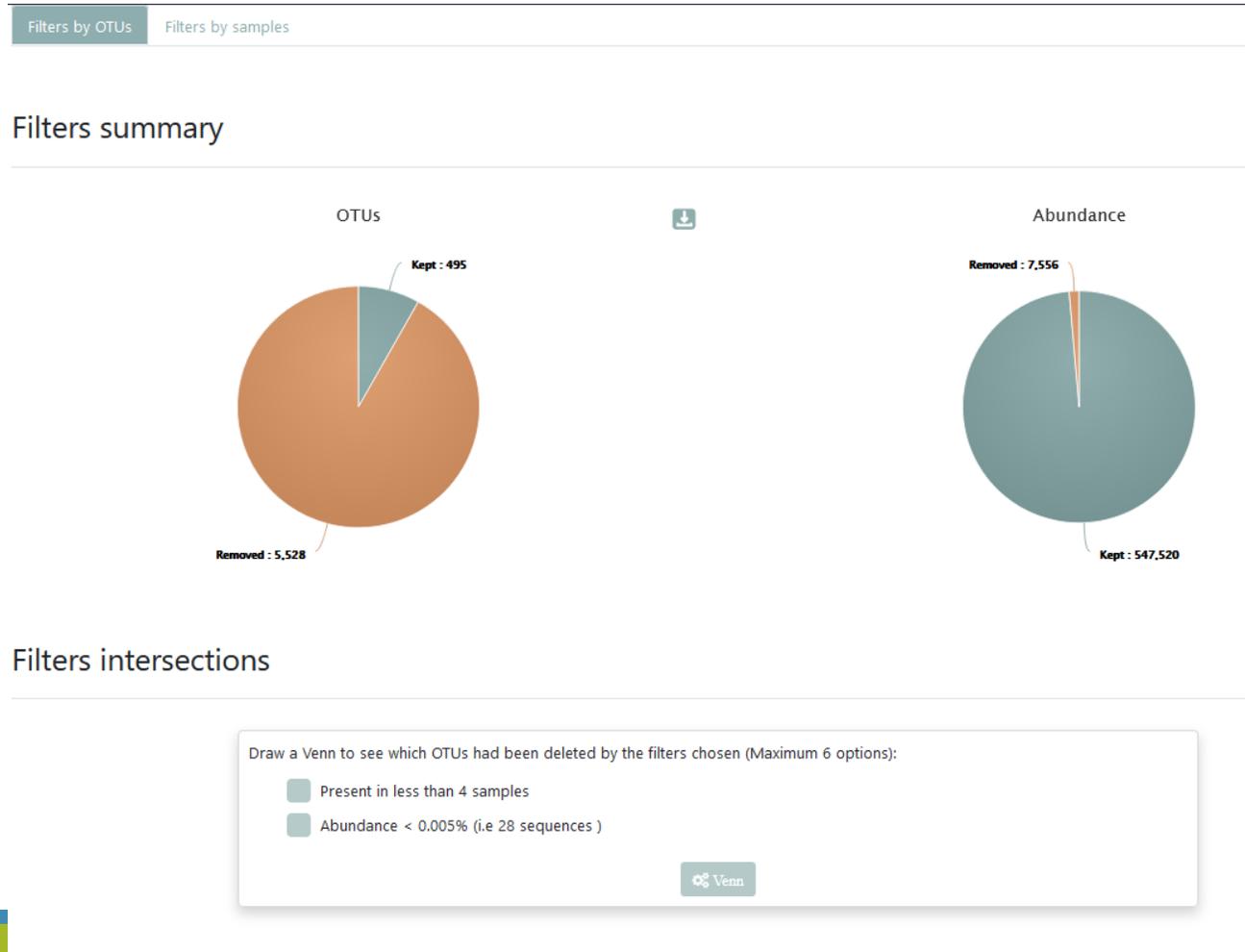
Either you use your own contaminant fasta file or you select one among available ones. (--contaminant)

Email notification

Send an email notification when the job completes.

✓ Execute

Ce que vous devez obtenir dans le report.html



Si vous n'obtenez pas cela, vérifiez vos paramètres et en particulier le paramètre du filtre sur l'abondance

Et seulement après contactez-nous à

frogs-training@inrae.fr

Si vous obtenez les mêmes chiffres, veuillez passer à l'étape suivante.



Pensez à lancer à nouveau l'outil Cluster_Stat pour observer le nombre d'OTU partager ou pas entre les différents échantillons.

report.html of
ClusterStat tool

Because of the "prevalence = 4" criterion, there is no longer an "own cluster" for any sample.

Clusters distribution Sequences distribution **Samples distribution**

Sequences count

Show entries Search: [CSV](#)

Sample	Total clusters	Shared clusters	Own clusters	Total sequences	Shared sequences	Own sequences
BHT0.LOT01	98	98	0	8,690	8,690	0
BHT0.LOT03	135	135	0	8,377	8,377	0
BHT0.LOT04	150	150	0	8,643	8,643	0
BHT0.LOT05	140	140	0	8,544	8,544	0
BHT0.LOT06	145	145	0	8,646	8,646	0
BHT0.LOT07	150	150	0	8,671	8,671	0
BHT0.LOT08	195	195	0	8,479	8,479	0
BHT0.LOT10	165	165	0	8,606	8,606	0
CDT0.LOT02	92	92	0	8,750	8,750	0
CDT0.LOT04	161	161	0	8,605	8,605	0

Affiliation tool

Les paramètres d'entrée

Input Parameter	Value
Using reference database	16S_SILVA_Pintail100_138.1
Also perform RDP assignation?	False
Taxonomic ranks	Domain Phylum Class Order Family Genus Species
Sequence file	15 FROGS OTU Filters: otuFilter_sequences.fasta
Abundance file	14 FROGS OTU Filters: otuFilter_abundance.biom

 **FROGS Affiliation OTU** Taxonomic affiliation of each OTU's seed by RDPtools and BLAST (Galaxy Version 4.0.0+galaxy1)

Using reference database

16S SILVA Pintail100 138.1

Select reference from the list

Also perform RDP assignation?

No

Taxonomy affiliation will be perform thanks to Blast. This option allows to perform it also with RDP classifier tool (default No) (--rdp)

Taxonomic ranks

Domain Phylum Class Order Family Genus Species

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

Sequence file

   15: FROGS OTU Filters: otuFilter_sequences.fasta

The sequences to affiliated (format: FASTA)

Abundance file

   14: FROGS OTU Filters: otuFilter_abundance.biom

The abundance file (format: BIOM)

Email notification

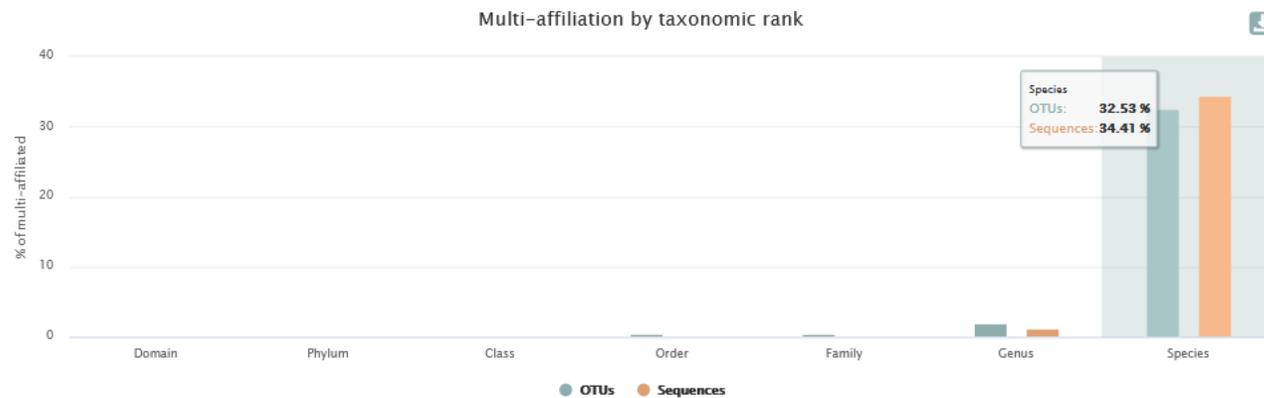
Send an email notification when the job completes.

 **Execute**

Ce que vous devez obtenir dans le report.html



Blast multi-affiliation summary



Si vous n'obtenez pas cela, vérifiez vos paramètres et en particulier avez-vous choisi la bonne banque ?

Et seulement après contactez-nous à

frogs-training@inrae.fr

Si vous obtenez les mêmes chiffres, veuillez passer à l'étape suivante.

Affiliation Stat

Les paramètres d'entrée

Input Parameter	Value
Abundance file	18 FROGS Affiliation OTU: affiliation_abundance.biom
Taxonomic ranks	Domain Phylum Class Order Family Genus Species
Rarefaction ranks	Class Order Family Genus Species
Affiliation processed	FROGS_blast

 **FROGS Affiliations stat** Process some metrics on taxonomies (Galaxy Version 4.0.0+galaxy1)

Abundance file

   18: FROGS Affiliation OTU: affiliation_abundance.biom

Abundances and affiliations (format: BIOM)

Taxonomic ranks

Domain Phylum Class Order Family Genus Species

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

Rarefaction ranks

Class Order Family Genus Species

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

Affiliation processed

FROGS Blast

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

Email notification



Send an email notification when the job completes.

 **Execute**

Ce que vous devez obtenir dans le report.html

Taxonomy distribution Alignment distribution

Display global distribution

Show 10 entries Search:

CSV

Samples	Nb domain	Nb phylum	Nb class	Nb order	Nb family	Nb genus	Nb species	Nb otus	Nb sequences
BHT0.LOT01	1	7	9	20	35	54	77	98	8,690
BHT0.LOT03	1	5	8	25	46	88	120	135	8,377
BHT0.LOT04	1	7	10	27	51	89	126	150	8,643
BHT0.LOT05	1	5	7	22	40	69	116	140	8,544
BHT0.LOT06	1	6	10	28	47	91	125	145	8,646
BHT0.LOT07	1	6	9	28	51	90	124	150	8,671
BHT0.LOT08	1	6	9	27	53	109	166	195	8,479
BHT0.LOT10	1	4	7	26	50	106	144	165	8,606
CDT0.LOT02	1	6	8	22	36	58	85	92	8,750
CDT0.LOT04	1	5	7	22	41	74	138	161	8,605

With selection: Class Display rarefaction Display distribution

Si vous n'obtenez pas cela, vérifiez vos paramètres

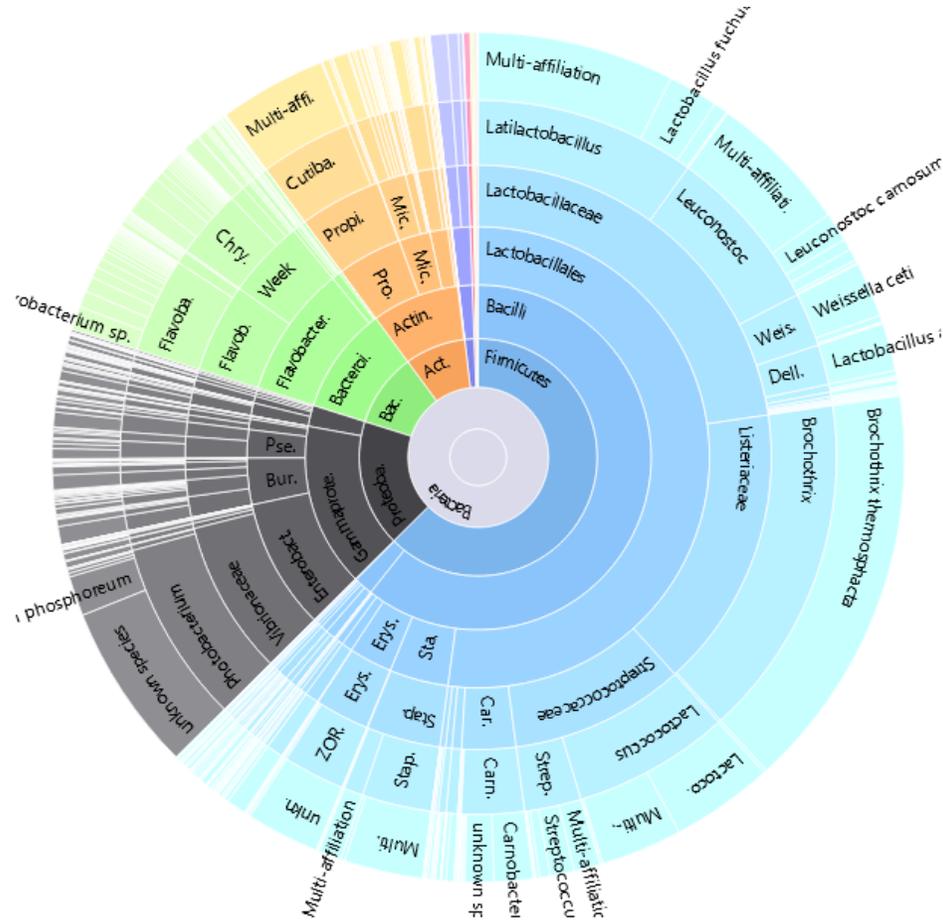
Et seulement après contactez-nous à

frogs-training@inrae.fr

Si vous obtenez les même chiffres, veuillez passer à l'étape suivante.

Ce que vous devez obtenir dans le report.html

Affichage de la distribution global des espèces



Pour faire cela, cliquez sur le bouton « Display global distribution »

Si vous obtenez les même chiffres, veuillez passer à l'étape suivante.

BIOM to TSV

CREATE A PHYLOGENETICS TREE OF OTUS

Les paramètres d'entrée

Input Parameter	Value
Abundance file	18 FROGS Affiliation OTU: affiliation_abundance.biom
Sequences file (optional)	15 FROGS OTU Filters: otuFilter_sequences.fasta
Extract multi-alignments	True

 **FROGS BIOM to TSV** Converts a BIOM file in TSV file (Galaxy Version 4.0.0+galaxy1)

Abundance file

   18: FROGS Affiliation OTU: affiliation_abundance.biom

The BIOM file to convert (format: BIOM)

Sequences file (optional)

   15: FROGS OTU Filters: otuFilter_sequences.fasta

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

Extract multi-alignments

Yes

If you have used FROGS affiliation on your data, you can extract information about multiple alignments in a second TSV.

Email notification

Send an email notification when the job completes.

 **Execute**

Ce que vous devez obtenir dans le report.html

#comment	blast_taxonomy	blast_subject	blast_perc_identity	blast_perc_query_coverage	blast_value
#comment	blast_taxonomy	blast_subject	blast_perc_identity	blast_perc_query_coverage	blast_value
no data	Bacteria;Firmicutes;Bacilli;Lactobacillales;Listeriaceae;Brochothrix;Brochothrix thermosphacta	multi-subject	100.0	100.0	0.0
no data	Bacteria;Proteobacteria;Gammaproteobacteria;Enterobacterales;Vibrionaceae;Photobacterium;unknown species	FJ456662.1.1555	100.0	100.0	0.0
no data	Bacteria;Firmicutes;Bacilli;Lactobacillales;Lactobacillaceae;Latilactobacillus;Multi-affiliation	multi-subject	100.0	100.0	0.0
no data	Bacteria;Actinobacteriota;Actinobacteria;Propionibacteriales;Propionibacteriaceae;Cutibacterium;Multi-affiliation	multi-subject	100.0	100.0	0.0
no data	Bacteria;Firmicutes;Bacilli;Lactobacillales;Lactobacillaceae;Leuconostoc;Multi-affiliation	multi-subject	100.0	100.0	0.0
no data	Bacteria;Firmicutes;Bacilli;Lactobacillales;Streptococcaceae;Lactococcus;Lactococcus piscium	AM943029.1.1242	99.799	100.0	0.0
no data	Bacteria;Firmicutes;Bacilli;Erysipelotrichales;Erysipelotrichaceae;ZOR0006;unknown species	HG792212.1.1536	94.203	100.0	0.0
no data	Bacteria;Firmicutes;Bacilli;Lactobacillales;Streptococcaceae;Lactococcus;Multi-affiliation	multi-subject	100.0	100.0	0.0
no data	Bacteria;Firmicutes;Bacilli;Lactobacillales;Lactobacillaceae;Weissella;Weissella ceti	FN813251.1.1461	99.799	100.0	0.0
no data	Bacteria;Bacteroidota;Bacteroidia;Flavobacteriales;Flavobacteriaceae;Flavobacterium;Flavobacterium sp.	AM934672.1.1486	98.775	100.0	0.0
no data	Bacteria;Proteobacteria;Gammaproteobacteria;Enterobacterales;Vibrionaceae;Photobacterium;Photobacterium phosphoreum	AB680520.1.1470	99.794	100.0	0.0
no data	Bacteria;Firmicutes;Bacilli;Lactobacillales;Lactobacillaceae;Dellaglistia;Lactobacillus algidus	AB033209.1.1505	99.803	99.60707269155206	0.0

Si vous n'obtenez pas cela, contactez-nous à frogs-training@inrae.fr

Si vous obtenez les mêmes chiffres, veuillez passer à l'étape suivante.

FROGS Tree

CREATE A PHYLOGENETICS TREE OF OTUS

Les paramètres d'entrée

Input Parameter	Value
Sequence file	15 FROGS OTU Filters: otuFilter_sequences.fasta
Biom file	18 FROGS Affiliation OTU: affiliation_abundance.biom

 FROGS Tree Reconstruction of phylogenetic tree (Galaxy Version 4.0.0+galaxy1)

Sequence file

   15: FROGS OTU Filters: otuFilter_sequences.fasta

Sequence file (format: FASTA). Warning: FROGS Tree does not work on more than 10000 sequences!

Biom file

   18: FROGS Affiliation OTU: affiliation_abundance.biom

The abundance file (format: BIOM)

Email notification



Send an email notification when the job completes.

 Execute