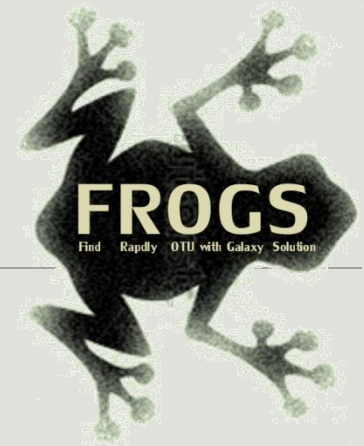


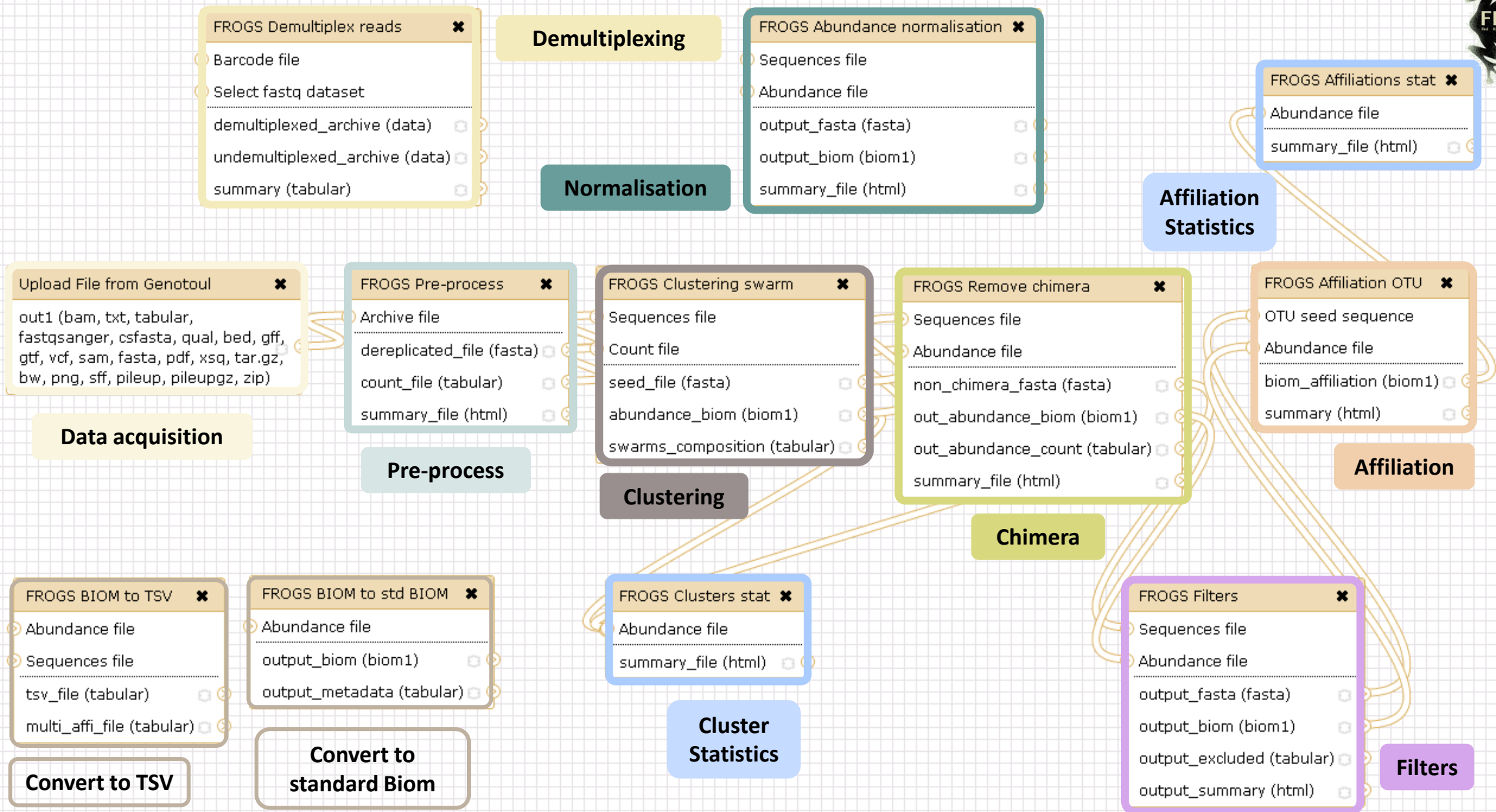
Exercices: Metagenomics

Find Rapidly OTU with Galaxy Solution



FRÉDÉRIC ESCUDIÉ* and LUCAS AUER*, MARIA BERNARD, LAURENT CAUQUIL, KATIA VIDAL, SARAH MAMAN, MAHENDRA MARIADASSOU, GUILLERMINA HERNANDEZ-RAQUET, GÉRALDINE PASCAL

* THESE AUTHORS HAVE CONTRIBUTED EQUALLY TO THE PRESENT WORK.



Exercice 1

LAUNCH UPLOAD TOOLS

Upload data

Your turn: exo 1



Create the 1st history **multiplexed**

Import files « **multiplex.fastq** » and « **barcode.tabular** » present in the **Genotoul** folder /work/formation/FROGS/



Create the 2nd history **454**

Import file « **454.fastq.gz** » present in the **Genotoul** folder /work/formation/FROGS/
(datatype fastq or fastq.gz is the same !)



Create the 3rd history **MiSeq R1 R2**

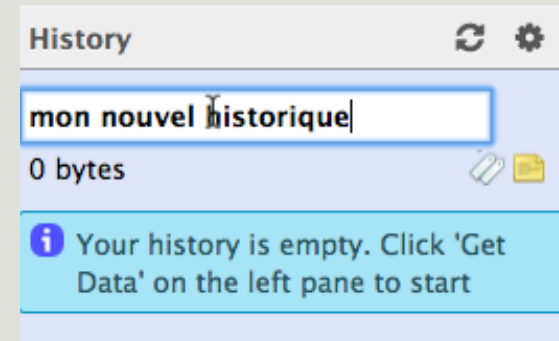
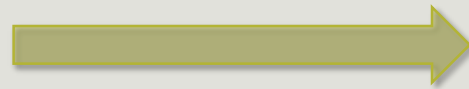
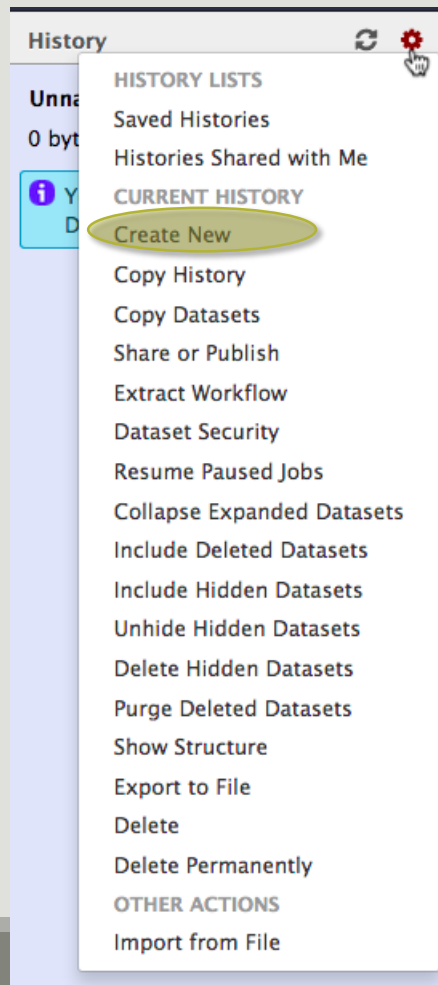
Import files « **sampleA_R1.fastq** » and « **sampleA_R2.fastq** » present in the **Genotoul** folder /work/formation/FROGS/



Create the 4th history **MiSeq contiged**

Import archive file « **100spec_90000seq_9samples.tar.gz** » present in the **Genotoul** folder /work/formation/FROGS/

History creation



Upload data: different methods

Tools

search tools

YOUR DATA

Upload Data

Upload File

Upload File from genotoul

EBI SRA ENA SRA

UCSC Main table browser

UCSC Test table browser

UCSC Archaea table browser

Get Microbial Data

BioMart Central server

Compress zip or tar file

Download Data

Upload File (version 1.1.3)

File Format:

Auto-detect

Which format? See help below

File:

Choisissez un fichier | Aucun fichier choisi

TIP: Due to browser limitations, uploading files larger than 2GB is guaranteed to fail (below) or FTP (if enabled by the site administrator).

URL/Text:

Here you may specify a list of URLs (one per line) or paste the contents of a file.

Convert spaces to tabs:

Yes

Use this option if you are entering intervals by hand.

Default method, your files are on **your** computer or accessible on the internet, they are copied on your Galaxy account

You can only upload one local file at a time
→ 10 samples ≥ 10 uploads
You can upload multiple files using URLs
but only smaller than 2Go



Each uploaded file will consume your Galaxy's quota!

Upload data: different methods

Tools

search tools

YOUR DATA

Upload Data

Upload File

Upload File from genotoul

EBI SRA ENA SRA

Upload File (version 1.0.0)

Path to file:

/work/frogs/Donnees_simulees/100WEPL_setA.tar.gz

Path must be like : /work/USERNAME/somewhere/afile

File type: Do not forget to precise the input file type

tar.gz

Execute

Specific SIGENAE GENOTOUL method. It allows you to access to your files in **your work account** on the Genotoul **without consuming your Galaxy quota**.

And if you have multiple samples ?

See [How to create an archiveTAR.ppt](#)



How to transfer files on /work of Genotoul?

See [How to transfert to genotoul.ppt](#)

Upload data: different methods

Tools

FROGS - FIND RAPIDLY OTU WITH GALAXY SOLUTION

FROGS pipeline

[Upload archive from your computer](#)

[Demultiplex reads](#) Split by samples the reads in function of inner barcode.

[FROGS Pre-process](#) Step 1 in metagenomics analysis (16S/18S): denoising and dereplication.

Upload archive (version 1.0.0)

File:
 Aucun fichier choisi
TIP: Due to browser limitations, uploading files larger than 2GB is guaranteed to fail. To upload large files, use the URL method.

URL:

Here you may specify the archive URL.

i What it does

If you have an **archive** on your **own computer** and **smaller than 2Go**, you may use this specific FROGS tool to upload your samples archive instead of the default « Upload File » of Galaxy.

Back main
documentation

Exercice 2

LAUNCH DEMULTIPLEX READS TOOL

The tool parameters depend on the input data type

FROGS Demultiplex reads (version 1.1.0)

Barcode file:

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

Single or Paired-end reads:

 Select between paired and single end data

Select fastq dataset:

 Specify dataset of your single end reads

barcode mismatches:

 Number of mismatches allowed in barcode

barcode on which end ?:

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

Forward
 Reverse
 Both ends
 Execute

FROGS Demultiplex reads (version 1.1.0)

Barcode file:

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

Single or Paired-end reads:

 Select between paired and single end data

Select first set of reads:

 Specify dataset of your forward reads

Select second set of reads:

 Specify dataset of your reverse reads

barcode mismatches:

 Number of mismatches allowed in barcode

barcode on which end ?:

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

Forward
 Reverse
 Both ends
 Execute

FROGS Demultiplex reads ✕

- Barcode file
- Select fastq dataset

demultiplexed_archive (data) 🗑

undemultiplexed_archive (data) 🗑

summary (tabular) 🗑

Demultiplexing

Exercise 2

In **multiplexed** history launch the demultiplex tool:

« The Patho-ID project, rodent and tick's pathobioms study, financed by the metaprogram INRA-MEM, studies zoonoses on rats and ticks from multiple places in the world, the co-infection systems and the interactions between pathogens. In this aim, they have extracted hundreds of or rats and ticks samples from which they have extracted 16S DNA and sequenced them first time on Roche 454 platform and in a second time on Illumina Miseq platform. For this courses, they authorized us to publicly shared some parts of these samples. »

Parasites & Vectors (2015) 8:172 DOI 10.1186/s13071-015-0784-7. **Detection of *Orientia* sp. DNA in rodents from Asia, West Africa and Europe.** Jean François Cosson, Maxime Galan, Emilie Bard, Maria Razzauti, Maria Bernard, Serge Morand, Carine Brouat, Ambroise Dalecky, Khalilou Bâ, Nathalie Charbonnel and Muriel Vayssier-Taussat

Exercise 2


In **multiplexed** history launch the demultiplex tool:




Data are single end reads

→ only 1 fastq file

Samples are characterized by an association of two barcodes in forward and reverse strands






→ multiplexing « both ends »

```
2: /work/frogs     
/Formation/multiplex.fastq
```

```
1: /work/frogs     
/Formation/barcode.txt
```

Exercise 2

Demultiplex tool asks for 2 files: one « fastq » and one « tabular »

1. Play with pictograms   
2. Observe how is built a fastq file. 
3. Look at the stdout, stderr when available (in the  pictogram)

FROGS Demultiplex reads (version 1.1.0)

Barcode file:

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

Single or Paired-end reads:

Select between paired and single end data

Select fastq dataset:

Specify dataset of your single end reads

barcode mismatches:

Number of mismatches allowed in barcode

barcode on which end ?:

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

What it does

Classify single or paired end reads in function of barcode forward or reverse in the first or both reads.

Command line:

```
demultiplex.py --input-R1 *FQ_INPUT1* [--input-R2 *FQ_INPUT2*] --input-barcode *TXT_I
```

History

FROGS multiplexed

2.1 MB

2: multiplex.fastq

2.1 MB

format: fastqsanger, database: ?

Epilog : job finished at Fri Nov 6 15:08:03 CET 2015

```
@HNHOSKD01ALD0H
ATCTAGTGATAAGTTCCTGTTTCATCCTAAGTCCATTATT
+
FFFFFFFFFDDA554444889422=<>40004444>>
@HNHOSKD01B8SLE
ATAGCTGATTGGTTTAAGCGGATAGGGATTAGATACCC
```

1: barcode.tabular

10 lines

format: tabular, database: ?




Epilog : job finished at Fri Nov 6 15:07:53 CET 2015




| 1 | 2 | 3 |
|-----------|---------|---------|
| MgArd0001 | ACAGCGT | TGTACGT |
| MgArd0009 | ACAGTAG | TGTACGT |
| MgArd0017 | ACGTCAG | TGTACGT |
| MgArd0029 | ACTCAGT | TGTACGT |
| MgArd0038 | ACTCGTC | TGTACGT |
| MgArd0046 | AGCAGTC | TGTACGT |




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

Results

9: FROGS Demultiplex   
reads: report

8: FROGS Demultiplex   
reads: undemultiplexed.tar.gz

7: FROGS Demultiplex   
reads: demultiplexed.tar.gz

Back to the main
documentation

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

| #sample | count |
|-----------|-------|
| ambiguous | 0 |
| MgArd0009 | 65 |
| MgArd0017 | 152 |
| MgArd0038 | 1185 |
| MgArd0029 | 172 |
| unmatched | 492 |
| MgArd0001 | 85 |
| MgArd0081 | 209 |
| MgArd0046 | 373 |
| MgArd0054 | 217 |
| MgArd0073 | 454 |
| MgArd0062 | 1109 |

With barcode mismatches >1 sequence can correspond to several samples. So these sequences are non-affected to a sample.

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

Exercises 3

LAUNCH THE PRE-PROCESS READS TOOL

Exercise 3.1

Go to « 454 » history

Launch the pre-process tool on that data set

→ objective : understand the parameters

1- Test different parameters for « minimum and maximum amplicon size »

2- Enter these primers: Forward: ACGGGAGGCAGCAG Reverse: AGGATTAGATACCCTGGTA

454

Size range of 16S V3-V4:
[380 – 500]

FROGS Pre-process (version 1.4.2)

Sequencer:

Illumina

Select the sequencer family used to produce the sequences.

Input type:

Files by samples

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

Reads already contiged ?:

Yes

The inputs contain 1 file by sample : Reads 1 and Reads 2 are already contiged by pair.

Samples

Samples 1

Name:

mon_echantillc

The sample name.

Sequence file:

6: /work/formation/FROGS/454.fastq.gz

FASTQ file of contiged reads.

Add new Samples

Minimum amplicon size:

380

The minimum size for the amplicons.

Maximum amplicon size:

500

The maximum size for the amplicons.

Sequencing protocol:

Illumina standard

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

5' primer:

ACGGGAGGCAGCAG

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

3' primer:

AGGATTAGATACCCTGGTA

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

Execute

Primers used for sequencing V3-V4:
Forward: ACGGGAGGCAGCAG
Reverse: AGGATTAGATACCCTGGTA

Exercise 3.1

What do you understand about amplicon size, which file can help you ?

Do you understand how enter your primers ?

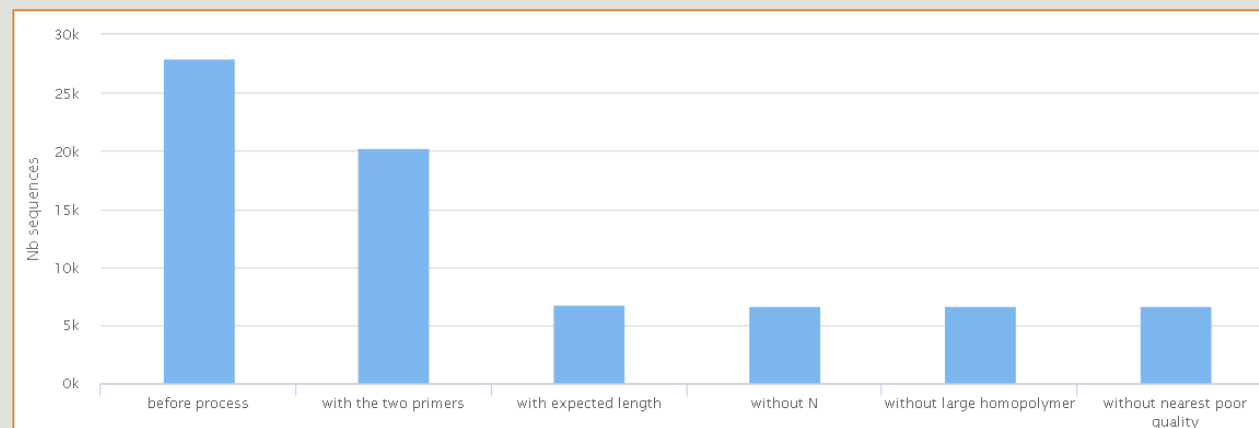
What is the « FROGS Pre-process: dereplicated.fasta » file ? 

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

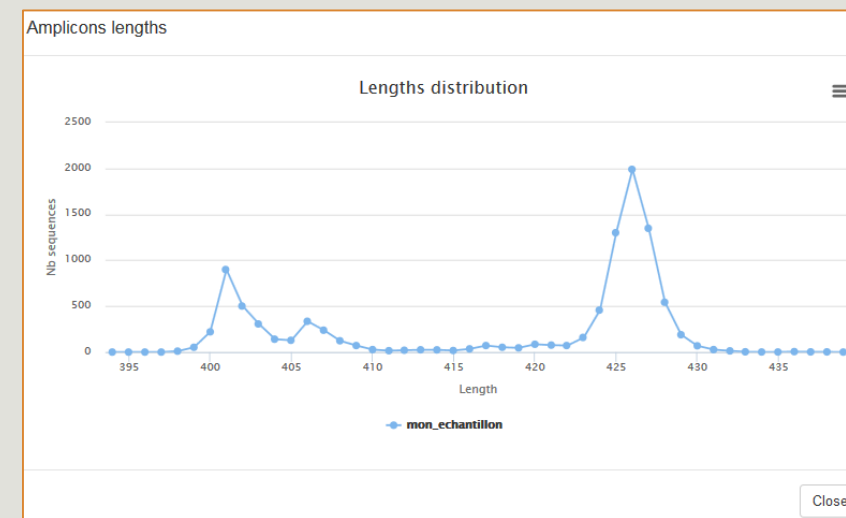
Explore the file « FROGS Pre-process: report.html » 

Who loose a lot of sequences ?

| <input type="checkbox"/> Samples | before process | with the two primers | with expected length | without N | without large homopolymer | without nearest poor quality |
|-------------------------------------|----------------|----------------------|----------------------|-----------|---------------------------|------------------------------|
| <input type="checkbox"/> sample_454 | 28,009 | 20,227 | 6,806 | 6,677 | 6,675 | 6,672 |



To be kept, sequences must have the 2 primers



To adjust your filtering, check the distribution of sequence lengths.

Back to the main documentation

Exercise 3.2

Go to « [MiSeq R1 R2](#) » history

Launch the pre-process tool on that data set

→ objective: understand flash software

Sequencer:

Select the sequencer family used to produce the sequences.

Input type:

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

Reads already contiged ?:

The inputs contain 1 file by sample : Reads 1 and Reads 2 are already contiged by pair.

Samples

Samples 1

Name:

The sample name.

Reads 1:

R1 FASTQ file of paired-end reads.

reads 2:

R2 FASTQ file of paired-end reads.

Reads 1 size:

The read1 size.

Reads 2 size:

The read2 size.

Primers used for this sequencing :

Forward: CCGTCAATTC

Reverse: CCGCNGCTGCT

Lecture 5' → 3'

>ERR619083.M00704

CCGTCAATTCATTGAGTTTCAACCTTGC GGCCG TACTTCCCAGGC GG TACGTT
 TATCGCGTTAGCTTCGCCAAGCACAGCATCCTGCGCTTAGCCAACGTACATCG
 TTTAGGGTGTGGACTACCCGGGTATCTAATCCTGTTTCGCTACCCACGCTTTCG
 AGCCTCAGCGTCAGTGACAGACCAGAGGCCGCTTTCGCCACTGGTGTTCCTC
 CATATATCTACGCATTTACCCGCTACACATGGAATTCCACTCTCCCCTTCTGC
 ACTCAAGTCAGACAGTTTCCAGAGCACTCTATGGTTGAGCCATAGCCTTTTAC
 TCCAGACTTTCCTGACCGACTGCACTCGCTTTACGCCAATAAATCCGGACAA
 CGCTTGCCACCTACGTATTA**CCGCNGCTGCT**

MiSeq
R1 R2

Real 16S sequenced
fragment

Expected amplicon size:

Maximum amplicon length expected in approximately 90% of the amplicons.

Minimum amplicon size:

The minimum size for the amplicons.

Maximum amplicon size:

The maximum size for the amplicons.

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'.

3' primer:

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

Flash ?

Exercise 3.2

Interpret « FROGS Pre-process: report.html » file.

Exercise 3.3

Go to« [MiSeq contiged](#) » history

Launch the pre-process tool on that data set

→ objective: understand output files

Exercise 3.3

3 samples are **technically replicated** 3 times : 9 samples of 10 000 sequences each.

| | | |
|-----------------------------|-----------------------------|-----------------------------|
| 100_10000seq_sampleA1.fastq | 100_10000seq_sampleB1.fastq | 100_10000seq_sampleC1.fastq |
| 100_10000seq_sampleA2.fastq | 100_10000seq_sampleB2.fastq | 100_10000seq_sampleC2.fastq |
| 100_10000seq_sampleA3.fastq | 100_10000seq_sampleB3.fastq | 100_10000seq_sampleC3.fastq |

Exercise 3.3

- 100 species, covering all bacterial phyla
- Power Law distribution of the species abundances
- Error rate calibrated with real sequencing runs
- 20% chimeras
- 9 samples of 10 000 sequences each (90 000 sequences)



Exercise 3.3

“Grinder (v 0.5.3) (Angly et al., 2012) was used to simulate the PCR amplification of full-length (V3-V4) sequences from reference databases. The reference database of size 100 were generated from the LTP SSU bank (version 115) (Yarza et al., 2008) by

- (1) filtering out sequences with a N,
- (2) keeping only type species
- (3) with a match for the forward (ACGGRAGGCAGCAG) and reverse (TACCAGGGTATCTAATCCT) primers in the V3-V4 region and
- (4) maximizing the phylogenetic diversity (PD) for a given database size. The PD was computed from the NJ tree distributed with the LTP.”

MiSeq
contiged

process (version 1.4.2)

sequencer family used to produce the sequences.

Input type:

Archive

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

Archive file:

1: /work/formation/FROGS/100spec_90000seq_9samples.tar.gz

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

Reads already contiged ?:

Yes

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

Minimum amplicon size:

380

The minimum size for the amplicons.

Maximum amplicon size:

500

The maximum size for the amplicons.

Sequencing protocol:

Illumina standard

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

5' primer:

ACGGGAGGCAGCAG

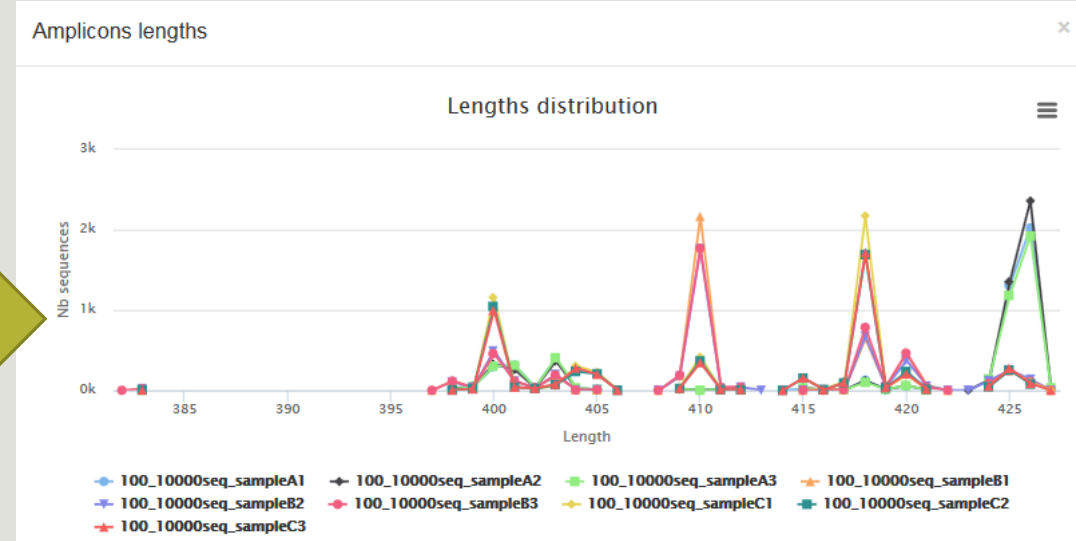
The 5' primer sequence (wildcards are accepted). The ori

3' primer:

TAGGATTAGATACCCTGGT

The 3' primer sequence (wildcards are accepted). The ori

Execute



Primers used for this sequencing :
5' primer: ACGGGAGGCAGCAG
3' primer: TAGGATTAGATACCCTGGTA
Lecture 5' → 3'

Exercise 3.3 - Questions

1. How many sequences are there in the input file ?
2. How many sequences did not have the 5' primer?
3. How many sequences still are after pre-processing the data?
4. How much time did it take to pre-process the data ?
5. What can you tell about the sample based on sequence length distributions ?

Exercise 4

LAUNCH CLUSTERING AND CLUSTERSTAT TOOLS

Clustering tools

Exercise 4

Go to « [MiSeq contiged](#) » history

Launch the Clustering SWARM tool on that data set with aggregation distance = 3 and the denoising

→ objectives :

- understand the denoising efficiency
- understand the ClusterStat utility

Exercise 4

1. How much time does it take to finish?
2. How many clusters do you get ?

Exercise 4

3. Edit the biom and fasta output dataset by adding **d1d3**



Attributes Convert Format Datatype Permissions

Edit Attributes

Name:
warm: seed_sequencesd1d3.fasta

Info:
Application
Software :/usr/local/bioinfo
/src/galaxy-test/galaxy-

Annotation / Notes:

FROGS Clusters stat Process
some metrics on clusters.

4. Launch FROGS Cluster Stat tools on the previous abundance biom file

Exercise 4

5. Interpret the boxplot: **Clusters size summary**
6. Interpret the table: **Clusters size details**
7. What can we say by observing the **sequence distribution**?
8. How many clusters share “sampleB3” with at least one other sample?
9. How many clusters could we expect to be shared ?
10. How many sequences represent the 550 specific clusters of “sampleC2”?
11. This represents what proportion of “sampleC2”?
12. What do you think about it?
13. How do you interpret the « Hierarchical clustering » ?

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

Exercise 5

LAUNCH THE REMOVE CHIMERA TOOL

Exercise 5

Go to « [MiSeq contiged](#) » history

Launch the « FROGS Remove Chimera » tool

Follow by the « FROGS ClusterStat » tool on the swarm d1d3 non chimera abundance biom

→ objectives :

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

FROGS Remove chimera ✕

- Sequences file
- Abundance file

- non_chimera_fasta (fasta)
- out_abundance_biom (biom1)
- out_abundance_count (tabular)
- summary_file (html)

Chimera

FROGS Remove chimera (version 1.3.0)

Sequences file:
5: FROGS Clustering swarm: d1d3_seed_sequences.fasta

The sequences file (format: fasta).

Abundance type:
BIOM file

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

Abundance file:
6: FROGS Clustering swarm: d1d3_abundance.biom

It contains the count by sample for each sequence.

Exercise 5

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

Exercise 5

2. Launch « FROGS ClusterStat » tool
on non_chimera_abundanced1d3.biom
3. Rename outputs in summary_nonchimera_d1d3.html
4. Compare the HTML files
 - a. Of what are mainly composed singleton ? (compare with precedent summary.html)
 - b. What are their abundance?
 - c. What do you conclude ?

The weakly abundant OTUs are mainly false positives, our data would be much more exact if we remove them

Exercise 6

LAUNCH DE LA TOOL FILTERS

Your turn: exo 6

Go to history« **MiSeq contiged** »

Launch « Filters » tool with non_chimera_abundanced1d3.biom, non_chimerad1d3.fasta

Apply 2 filters :

- **proportion/number of sequences threshold to remove an OTU: 0.00005***
- **Remove OTUs that are not present at least in XX samples;
how many samples do you choose? : 3**

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

FROGS Filters

- Sequences file
- Abundance file

- output_fasta (fasta)
- output_biom (biom1)
- output_excluded (tabular)
- output_summary (html)

Input

Filters

FROGS Filters (version 1.1.0)

Sequences file:

 The sequence file to filter (format: fasta).

Abundance file:

 The abundance file to filter (format: BIOM).

*** THE FILTERS ON OTUS IN SAMPLES, OTUS SIZE and SEQUENCE PERCENTAGE:

If you want to filter OTUs on their abundance and occurrence.

Remove OTUs that are not present at least in XX samples; how many samples do you choose? :

 Fill the field only if you want this treatment.

Proportion/number of sequences threshold to remove an OTU:

 Fill the field only if you want this treatment. Use decimal to express proportion (0.01 for 1%) integer to express number of sequence (1 for singleton).

When sorted by abundance, how many OTU do you want to keep ?:

 Fill the fields only if you want this treatment.

*** THE FILTERS ON RDP:

If you want to filter OTUs on their taxonomic affiliation produced by RDP.

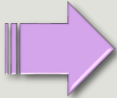
*** THE FILTERS ON BLAST:

If you want to filter OTUs on their taxonomic affiliation produced by Blast.

*** THE FILTERS ON CONTAMINATIONS:

If you want to filter OTUs on classical contaminations.

Output



- 92: **FROGS Filters: report.html**   
- 91: **FROGS Filters: excluded.tsv**   
- 90: **FROGS Filters: abundance.biom**   
- 89: **FROGS Filters: sequences.fasta**   

Your turn: exo 6

1. What are the output files of “Filters” ?
2. Explore “FROGS Filter : summary.html” file.
3. How many OTUs have you removed ?
4. Build the Venn diagram on the two filters.
5. How many OTUs have you removed with each filter “**abundance > 0.005%**”, “**Remove OTUs that are not present at least in 3 samples**”?
6. How many OTUs do they remain ?
7. Is there a sample more impacted than the others ?
8. To characterize these new OTUs, do not forget to launch “FROGS Cluster Stat” tool, and rename the output HTML file.

Exercises 7

LAUNCH THE « FROGS AFFILIATION » TOOL

Exercise 7.1

Go to « [MiSeq contiged](#) » history

Launch the « FROGS Affiliation » tool with

- SILVA 123 16S database
- FROGS Filters abundance biom and fasta files (after swarm d1d3, remove chimera and filter low abundances)

→ objectives :

- understand abundance tables columns
- understand the RDP and BLAST affiliation complementarity

FROGS Affiliation OTU ✕

- OTU seed sequence
- Abundance file

biom_affiliation (biom1) 🗑️ 🔍

summary (html) 🗑️ 🔍

Affiliation



FROGS Affiliation OTU (version 0.7.0)

Using reference database:
 ▼
Select reference from the list

OTU seed sequence:
 ▼
OTU sequences (format: fasta).

Abundance file:
 ▼
OTU abundances (format: BIOM).

Exercise 7.1

1. What are the « FROGS Affiliation » output files ?
2. How many sequences are affiliated by BLAST ?
3. Click on the « eye » button on the BIOM output file, what do you understand ? 
4. Use the Biom_to_TSV tool on this last file and click again on the "eye" on the new output generated. 
 What do the columns ?
 What is the difference if we click on case or not ? What consequence about weight of your file ?

FROGS BIOM to TSV (version 2.1.0)

Abundance file:
 17: FROGS Affiliation OTU: affiliation.biom
 The BIOM file to convert (format: BIOM).

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

Extract multi-alignments:

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

Execute

Tools

[FROGS - FIND RAPIDLY OTU WITH GALAXY SOLUTION](#)

FROGS pipeline

[FROGS Upload archive](#) from your computer

[FROGS Demultiplex reads](#)
 Split by samples the reads in function of inner barcode.

[FROGS Pre-process Step 1](#) in metagenomics analysis: denoising and dereplication.

[FROGS Clustering swarm](#)
 Step 2 in metagenomics analysis : clustering.

[FROGS Remove chimera](#) Step 3 in metagenomics analysis : Remove PCR chimera in each sample.

[FROGS Filters](#) Filters OTUs on several criteria.

[FROGS Affiliation OTU](#) Step 4 in metagenomics analysis : Taxonomic affiliation of each OTU's seed by RDPTools and BLAST

[FROGS BIOM to TSV](#) Converts a BIOM file in TSV file.

[FROGS Clusters stat](#) Process some metrics on clusters.

[FROGS Affiliations stat](#)
 Process some metrics on taxonomies.

[FROGS BIOM to std BIOM](#)
 Converts a FROGS BIOM in fully compatible BIOM.

[FROGS Abundance normalisation](#)

Exercise 7.1

5. Compare RDP and Blast affiliations - Cluster_2388

| #rdp_tax_and_bootstrap | blast_subject | blast_evalue | blast_len | blast_perc_query_coverage | blast_perc_identity | blast_taxonomy |
|---|---------------------------------|--------------|-----------|---------------------------|---------------------|--|
| Bacteria;(1.0);Planctomycetes;(1.0);Planctomycetacia;(1.0);Planctomycetales;(1.0);Planctomycetaceae;(1.0);Telmatocola;(1.0);Telmatocola sphagniphila;(1.0); | JN880417.1.1422 | 0.0 | 360 | 88.88 | 99.44 | Bacteria;Planctomycetes;Planctomycetacia;Planctomycetales;Planctomycetaceae;Telmatocola;Telmatocola sphagniphila |

Blast JN880417.1.1422 vs our OTU

OTU length : 405

Excellent blast but no matches at the beginning of OTU. Chimera ?

Telmatocola sphagniphila strain SP2 16S ribosomal RNA gene, partial sequence
Sequence ID: [ref|NR_118328.1](#) Length: 1422 Number of Matches: 1

Range 1: 375 to 734 [GenBank](#) [Graphics](#) [Next Match](#) [Previous Match](#)

| Score | Expect | Identities | Gaps | Strand |
|---------------|---|--------------|-----------|-----------|
| 654 bits(354) | 0.0 | 358/360(99%) | 0/360(0%) | Plus/Plus |
| Query 46 | CGCGTGCGCGATGAAGGCCTTCGGGTTGTAAGCGCGAAAGAGGTAATAAAGGGAAACT | 105 | | |
| Sbjct 375 | CGCGTGCGCGATGAAGGCCTTCGGGTTGTAAGCGCGAAAGAGGSAATAAAGGGAAACT | 434 | | |
| Query 106 | GATTGAACCTCAGTAAGCTCGGGCTAAGTTTGTGCCAGCAGCCGCGGTAAGACGAACCGA | 165 | | |
| Sbjct 435 | GATTGAACCTCAGTAAGCTCGGGCTAAGTTTGTGCCAGCAGCCGCGGTAAGACGAACCGA | 494 | | |
| Query 166 | GCGAACGTTGTTTCGGAATCACTGGGCATAAAGGGCGCGTAGGCGGGTTTCTAAGTCCGTG | 225 | | |
| Sbjct 495 | GCGAACGTTGTTTCGGAATCACTGGGCATAAAGGGCGCGTAGGCGGGTTTCTAAGTCCGTG | 554 | | |
| Query 226 | GTGAAATACTTCAGCTCAACTGGAGAACTGCCTCGGATACTGGGAATCTCGAGTAATGTA | 285 | | |
| Sbjct 555 | GTGAAATACTTCAGCTCAACTGGAGAACTGCCTCGGATACTGGGAATCTCGAGTAATGTA | 614 | | |
| Query 286 | GGGGCACGTGGAACGGCTGGTGGAGCGGTGAAATGCGTTGATATCAGTCGGAACCTCCGGT | 345 | | |
| Sbjct 615 | GGGGCACGTGGAACGGCTGGTGGAGCGGTGAAATGCGTTGATATCAGTCGGAACCTCCGGT | 674 | | |
| Query 346 | GGCGAAGGCGATGTGCTGGACATTTACTGACGCTGAGGCCGCGAAAGCCAGGGGAGCAAAC | 405 | | |
| Sbjct 675 | GGCGAAGGCGATGTGCTGGACATTTACTGACGCTGAGGCCGCGAAAGCCAGGGGAGCAAAC | 734 | | |

Telmatocola sphagniphila strain SP2 16S ribosomal RNA gene, partial sequence
NCBI Reference Sequence: NR_118328.1
[FASTA](#) [Graphics](#)

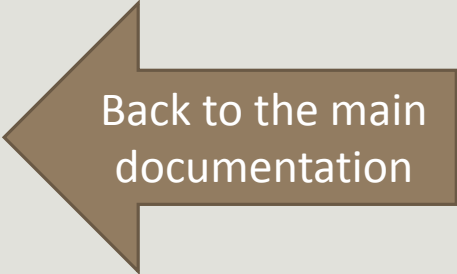
[Go to:](#)

LOCUS NR_118328 1422 bp rRNA linear BCT 03-FEB-2015
DEFINITION Telmatocola sphagniphila strain SP2 16S ribosomal RNA gene, partial sequence
ACCESSION [NR_118328](#)
VERSION [NR_118328.1](#) GI:645321338
DBLINK Project: [33175](#)
BioProject: [PRJNA33175](#)
KEYWORDS RefSeq.
SOURCE Telmatocola sphagniphila
ORGANISM [Telmatocola sphagniphila](#)
Bacteria; Planctomycetes; Planctomycetia; Planctomycetales;
Planctomycetaceae.
REFERENCE 1 (bases 1 to 1422)
AUTHORS Kulichevskaya,I.S., Serkebaeva,Y.M., Kim,Y., Rijpstra,W.I.,
Damste,J.S., Liesack,W. and Dedysh,S.N.
TITLE Telmatocola sphagniphila gen. nov., sp. nov., a novel dendriform
planctomycete from northern wetlands
JOURNAL Front Microbiol 3, 146 (2012)
PUBMED [22529844](#)
REMARK Publication Status: Online-Only
REFERENCE 2 (bases 1 to 1422)
CONSRM NCBI RefSeq Targeted Loci Project
TITLE Direct Submission
JOURNAL Submitted (28-APR-2014) National Center for Biotechnology
Information, NIH, Bethesda, MD 20894, USA
REFERENCE 3 (bases 1 to 1422)
AUTHORS Dedysh,S.N.
TITLE Direct Submission
JOURNAL Submitted (20-OCT-2011) Winogradsky Institute of Microbiology RAS,
Prospect 60-Letya Otyabrya 7/2, Moscow 117312, Russia
COMMENT REVIEWED [REFSEQ](#): This record has been curated by NCBI staff. The
reference sequence is identical to [JN880417:1-1422](#).

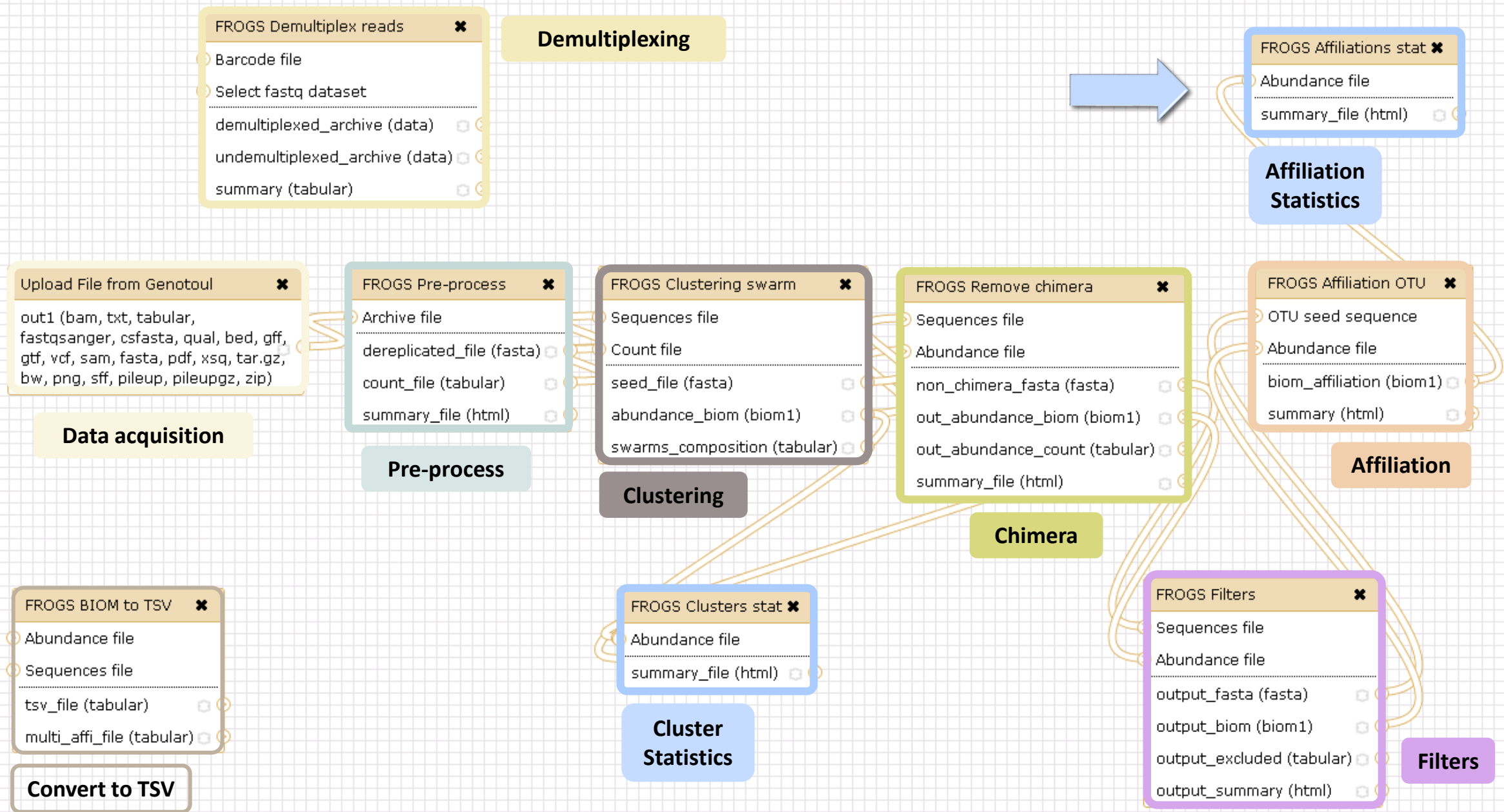
What do you think about this case ?

Observe the Cluster_4 line (big abundance!):

| #rdp_tax_and_bootstrap | blast_subject | blast_evalue | blast_len | blast_perc_query_coverage | blast_perc_identity | blast_taxonomy |
|---|-----------------|--------------|-----------|---------------------------|---------------------|---|
| Bacteria;(1.0);Thermotogae;(1.0);Thermotogae;(1.0);Thermotogales;(1.0);Thermotogaceae;(1.0);Petrotoga;(1.0);Petrotoga miotherma;(0.62); | FR733705.1.1499 | 0.0 | 419 | 100.0 | 100.0 | Bacteria;Thermotogae;Thermotogae;Thermotogales;Thermotogaceae;Petrotoga;Petrotoga miotherma |



Back to the main documentation



Exercise 7.2

Launch the « FROGS Affiliation Stat» tool

- On FROGS blast affiliation
- On FROGS RDP affiliation

→ objectives :

understand rarefaction curve and sunburst

understand the RDP and BLAST affiliation complementarity

Exercise 7.2

FROGS Affiliations stat (version 1.1.0)

Abundance file:

17: FROGS Affiliation OTU: affiliation.biom

OTUs abundances and affiliations (format: BIOM).

Rarefaction ranks:

Class Order Family Genus Species

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

Affiliation processed:

FROGS blast

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

Execute

FROGS Affiliations stat (version 1.1.0)

Abundance file:

17: FROGS Affiliation OTU: affiliation.biom

OTUs abundances and affiliations (format: BIOM).

Rarefaction ranks:

Class Order Family Genus Species

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

Affiliation processed:

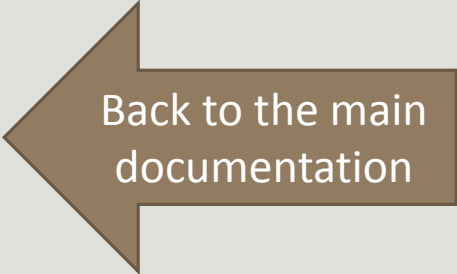
FROGS rdp

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

Execute

Exercise 7.2

1. Explore the Affiliation stat results.
2. What kind of graphs can you generate? What do they mean?
 - a) Common to Blast and RDP affiliation results
 - b) On Blast results
 - c) On RDP results



Back to the main
documentation

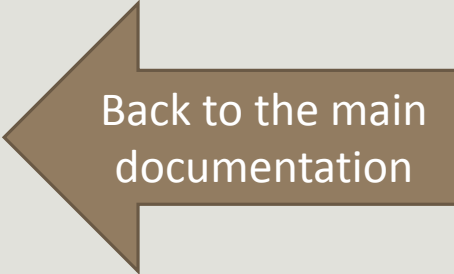
Exercise 8

LAUNCH NORMALIZATION TOOL

Exercise 8

Launch Normalization Tool

1. What is the smallest sequenced samples ?
2. Normalize your data from Affiliation based on that number of sequence
3. Explore the report HTML result.



Back to the main
documentation

Exercise 9

CREATE YOUR OWN WORKFLOW !

MiSeq
contiged

Exercice 9

Galaxy Sigenae - Welcome gpascal Analyze Data Workflow Shared Data Visualization Help User Using 18.3 GB

Your workflows

[+ Create new workflow](#) [↑ Upload or import workflow](#)

| Name | # of Steps |
|-------------------------|------------|
| formation workflow ▾ | 9 |
| demoNEM2015 workflow ▾ | 9 |
| FROGS_v1.0_06_05_2015 ▾ | 10 |

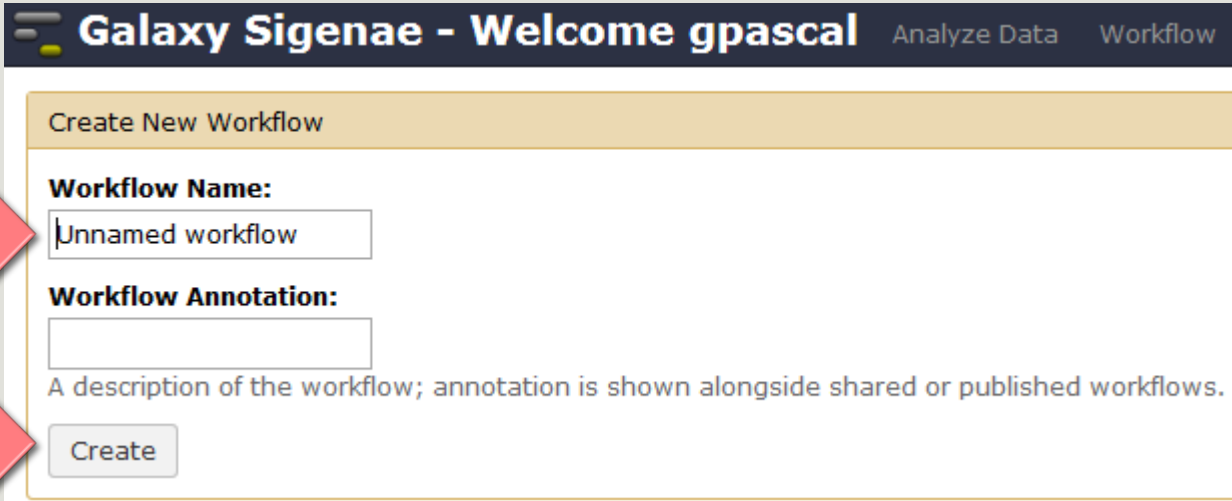
Workflows shared with you by others

No workflows have been shared with you.

Other options

[Configure your workflow menu](#)

Exercice 9



Galaxy Sigenae - Welcome gpascal Analyze Data Workflow

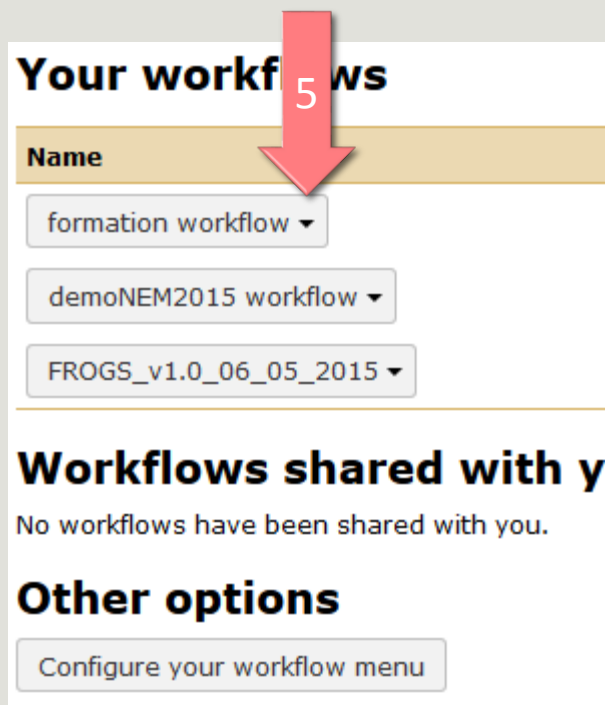
Create New Workflow

Workflow Name:

Workflow Annotation:

A description of the workflow; annotation is shown alongside shared or published workflows.

Exercice 9



Your workflows

Name

- formation workflow ▾
- demoNEM2015 workflow ▾
- FROGS_v1.0_06_05_2015 ▾

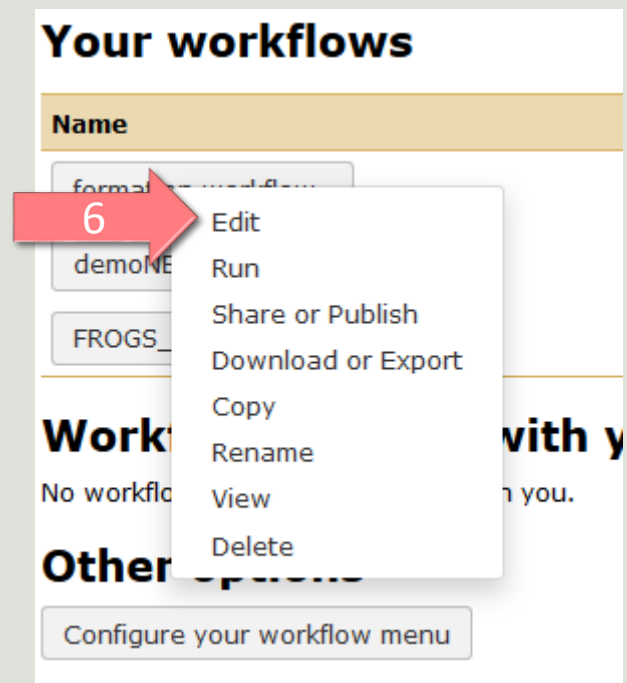
Workflows shared with you

No workflows have been shared with you.

Other options

[Configure your workflow menu](#)

A red arrow with the number 5 points to the first workflow name 'formation workflow'.



Your workflows

Name

- formation workflow ▾
- demoNEM2015 workflow ▾
- FROGS_v1.0_06_05_2015 ▾

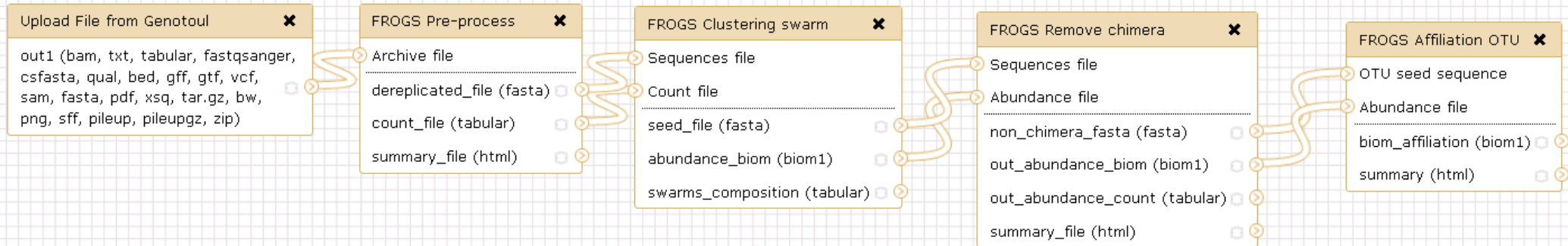
Workflows shared with you

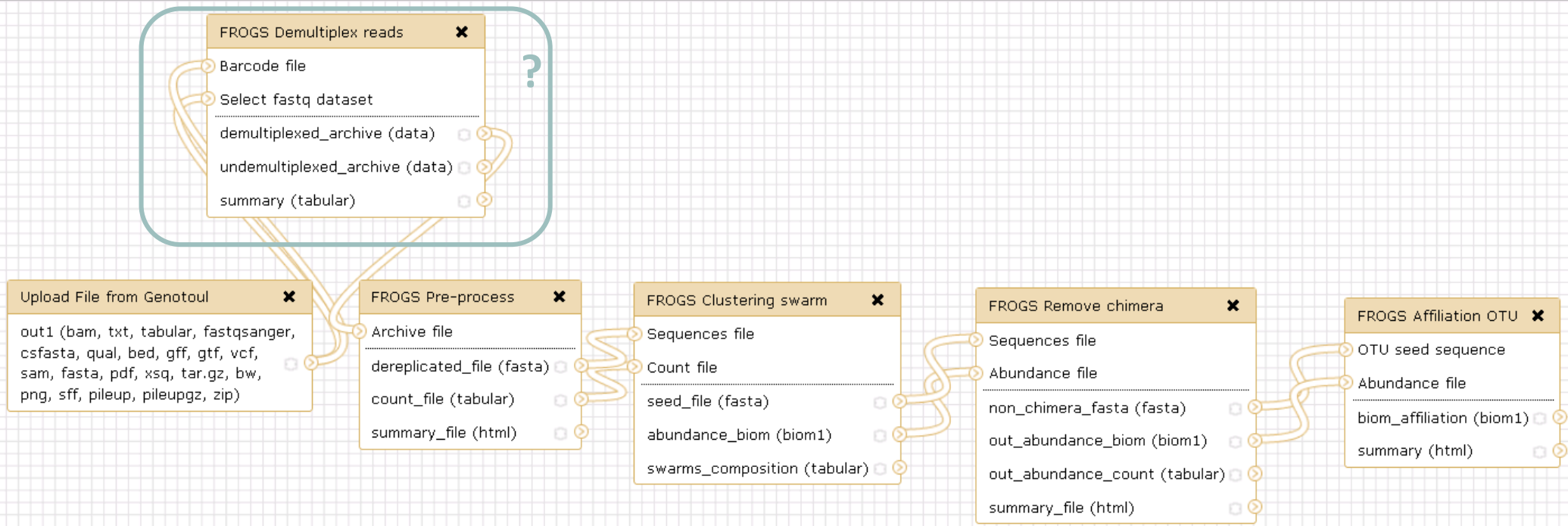
No workflows have been shared with you.

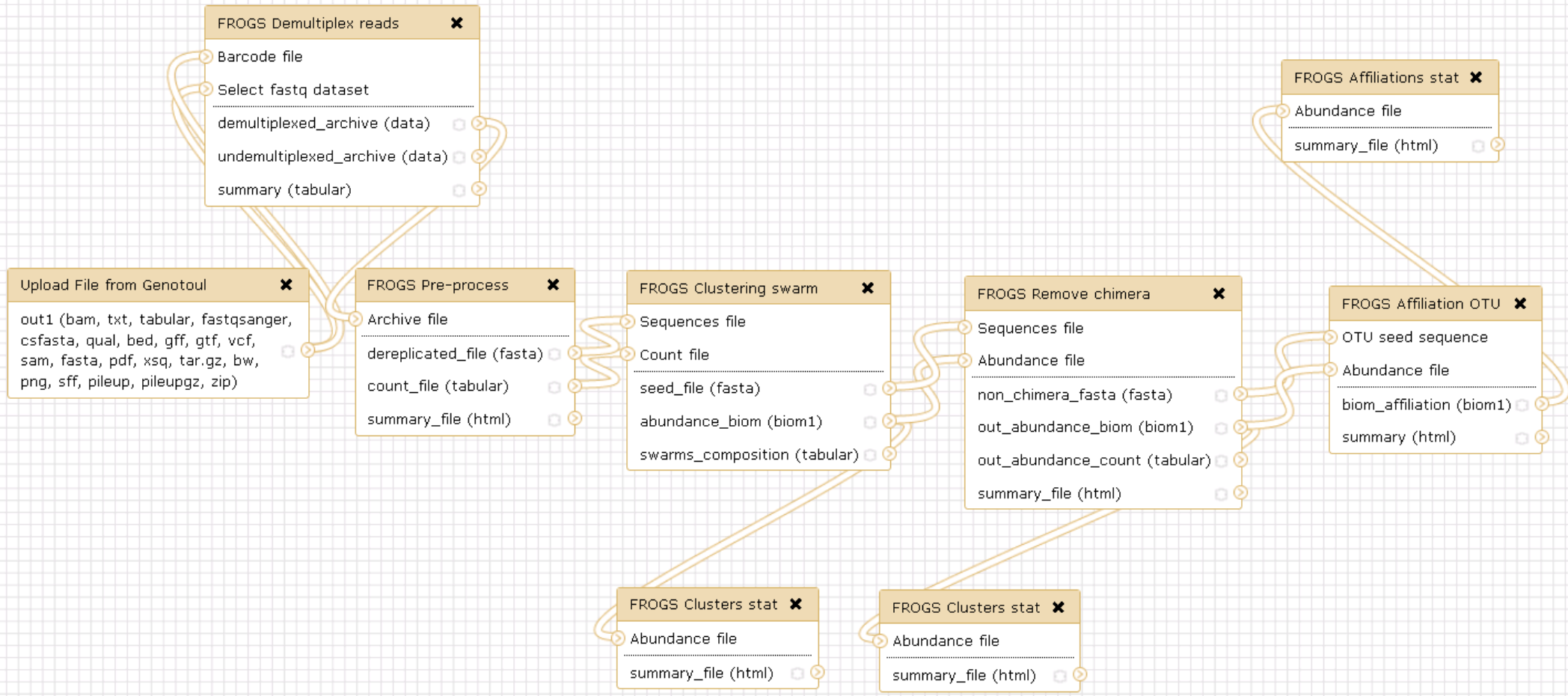
Other options

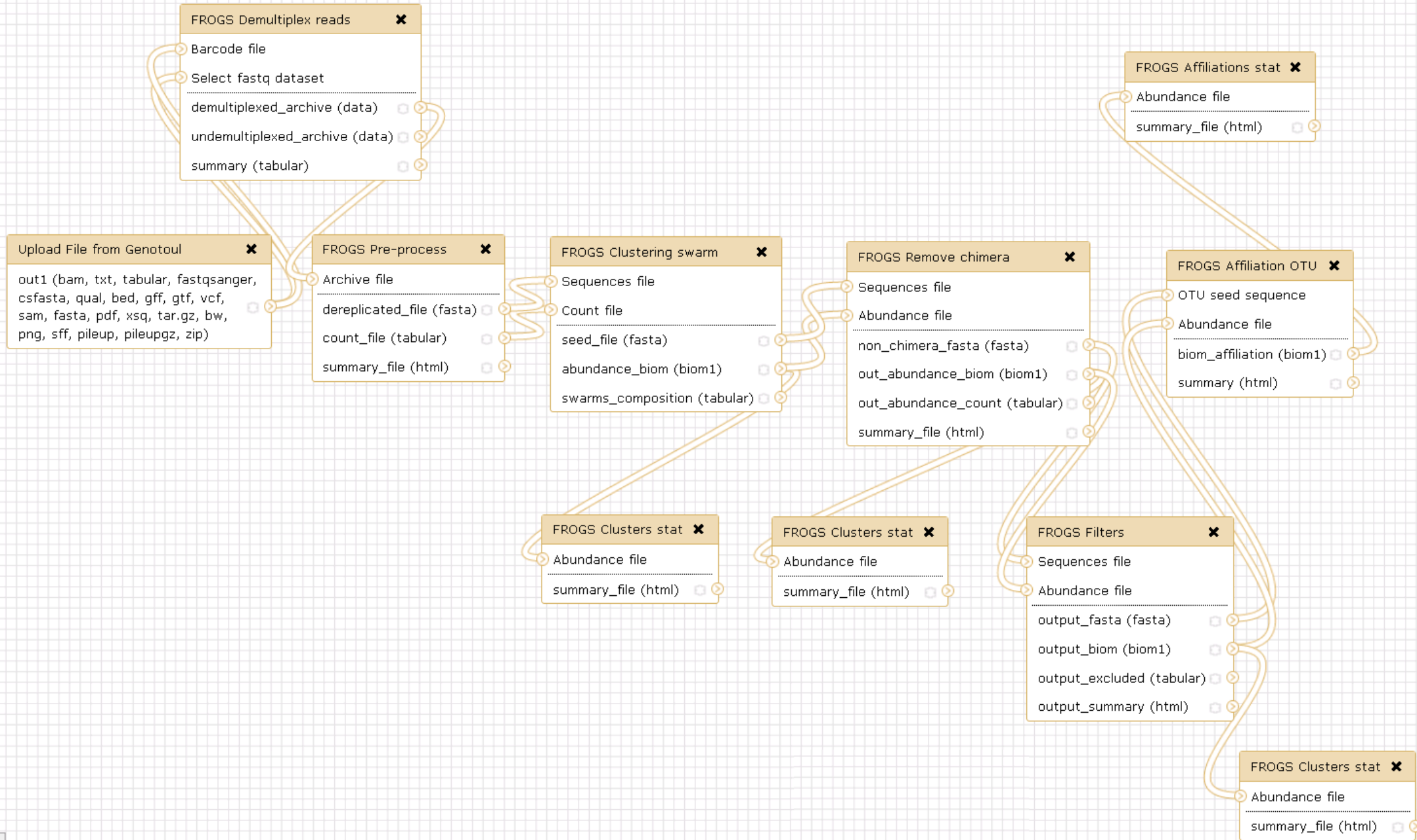
[Configure your workflow menu](#)

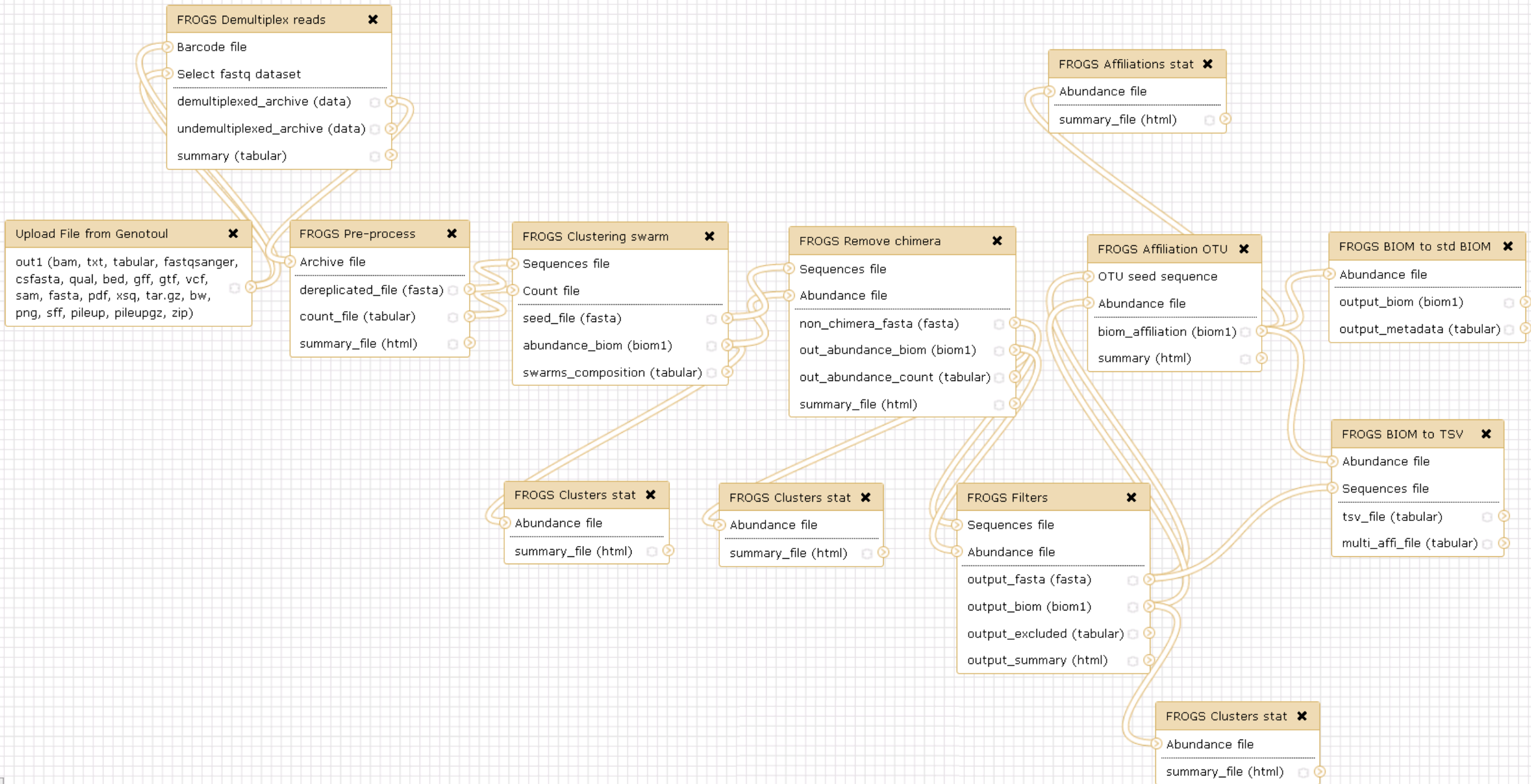
A context menu is open over the first workflow name 'formation workflow', listing the following actions: Edit, Run, Share or Publish, Download or Export, Copy, Rename, View, and Delete. A red arrow with the number 6 points to the 'Edit' option.

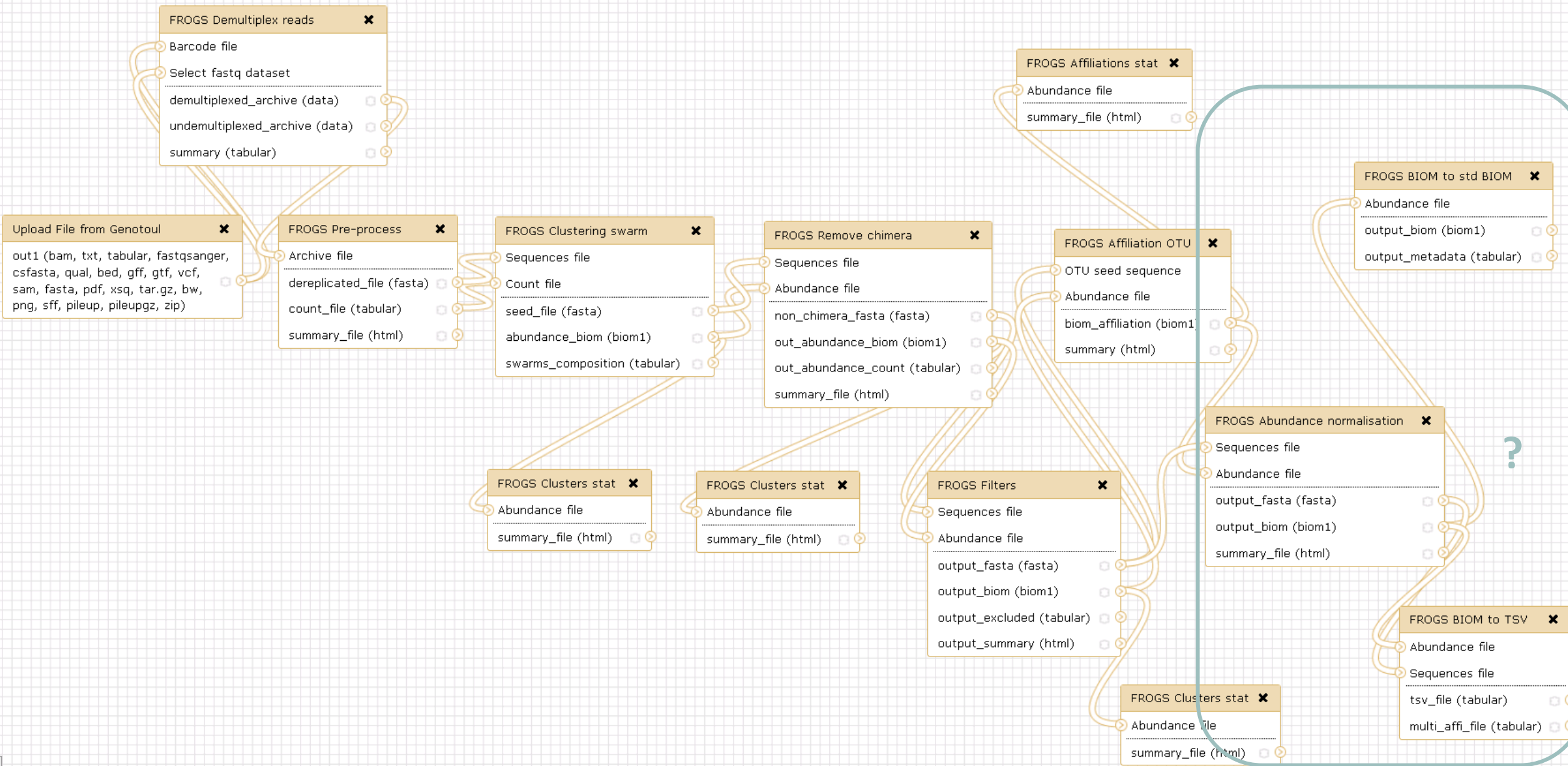


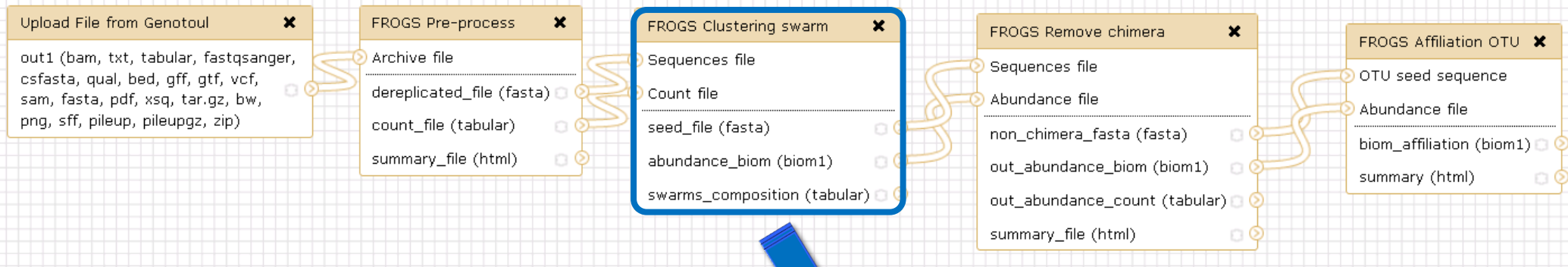












For each tool, think to:

- Fixe parameter ?

Tool: FROGS Clustering swarm

Version: 2.3.0

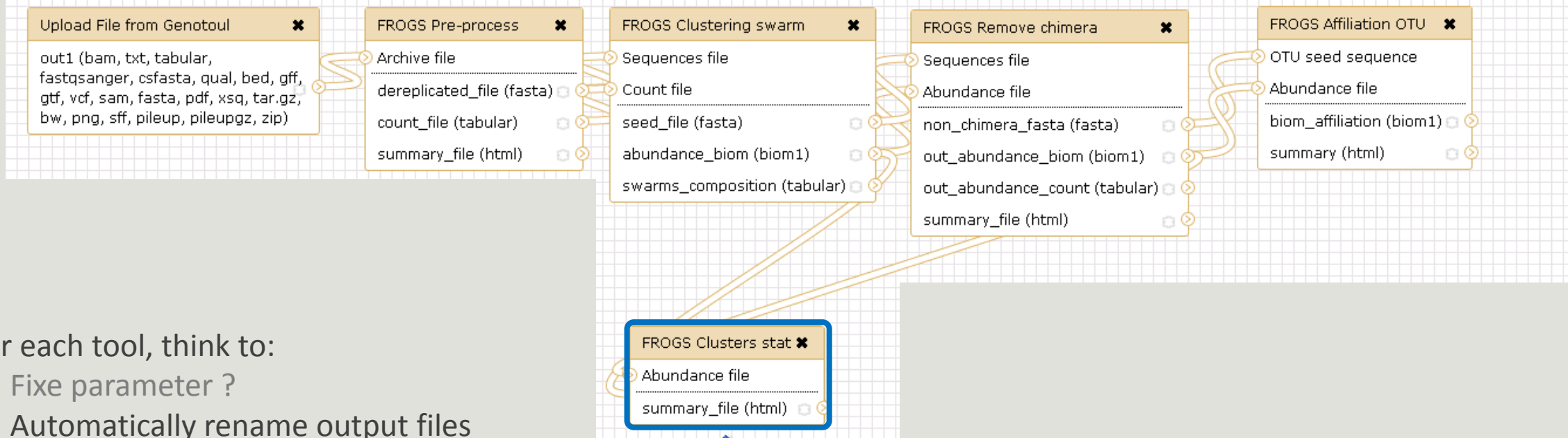
Sequences file
Data input 'sequence_file' (fasta)

Count file
Data input 'count_file' (tabular)

None: ▼

Aggregation distance: ▼

Performe denoising clustering step?: ▼ ?



For each tool, think to:

- Fixe parameter ?
- Automatically rename output files

Tool: FROGS Clusters stat

Version: 1.4.0

Abundance file
Data input 'biom' (biom1)

Edit Step Actions

Rename Dataset
summary_file

Add actions to this step; actions are applied when this workflow step completes.



Edit Step Actions

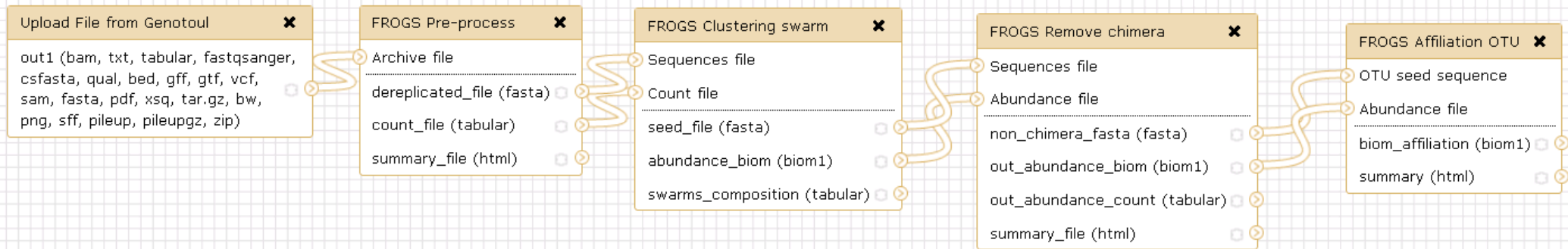
Rename Dataset
summary_file

Rename Dataset on summary_file

New output name:
swarm_cluster_stat.html

This action will rename the result dataset.

Add actions to this step; actions are applied when this workflow step completes.

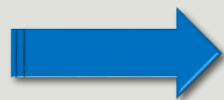


For each tool, think to:

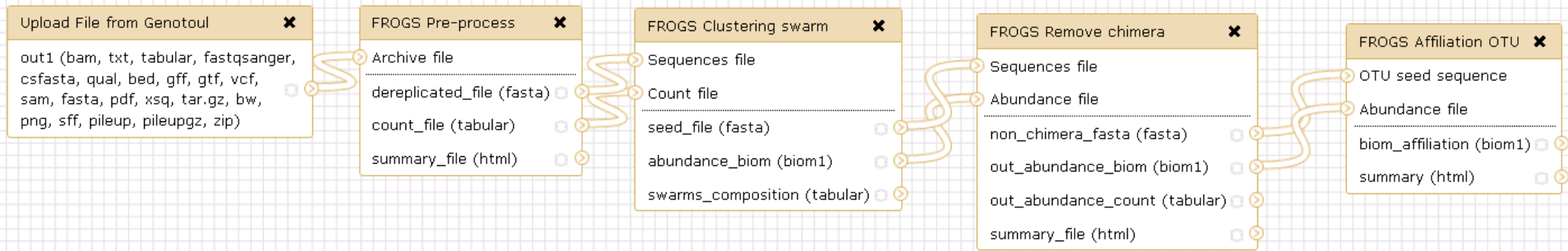
- Fixe parameter ?
- Automatically rename output files
- Hide intermediate files ?

FROGS Remove chimera ✕

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



- 11: FROGS Remove chimera: report.html** 👁️ ✎ ✕
- 10: FROGS Remove chimera: non chimera abundance.biom** 👁️ ✎ ✕
- 9: FROGS Remove chimera: non chimera.fasta** 👁️ ✎ ✕



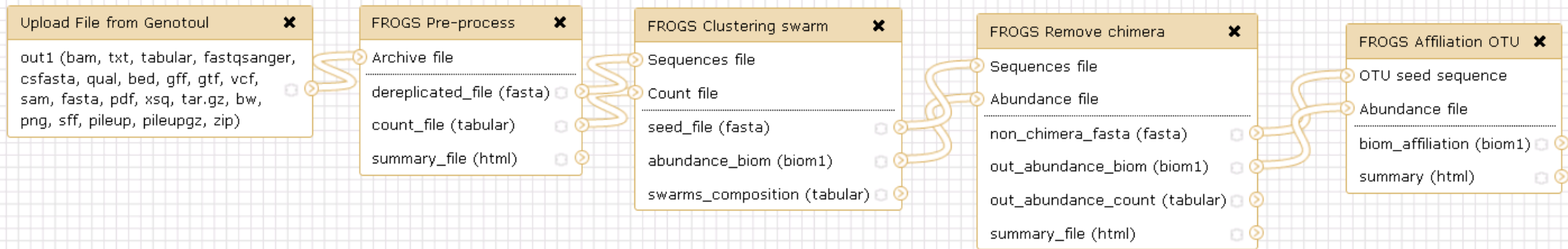
For each tool, think to:

- Fixe parameter ?
- Automatically rename output files
- Hide intermediate files ?

FROGS Remove chimera

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

Flag this as a workflow output. All non-flagged outputs will be hidden.

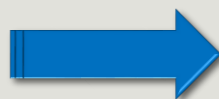





For each tool, think to:

- Fixe parameter ?
- Automatically rename output files
- Hide intermediate files ?

FROGS Remove chimera

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



11: FROGS Remove chimera: report.html   

Back to the main documentation