

### Core workflow

**Pre-process**  
 Multiplex sequences (23S, 18S, 16S, 454) → MiSeq Fastq R1, MiSeq Fastq R2 → Already contiged → FROGS Pre-process (dereplicated\_file, count\_file, summary\_file)

**Clustering**  
 Paired-end reads merging: Flash  
 Trimming: Cutadapt → FROGS Clustering swarm (abundance\_biom, seed\_file, swarms\_composition)

**Chimera**  
 Denoising and Clustering with a local threshold: SWARM [2] → FROGS Remove chimera (non\_chimera\_fasta, out\_abundance\_biom, out\_abundance\_count, summary\_file)

**Filters**  
 Efficient Chimera Removal: VSEARCH [3] + homemade cross-validation → FROGS Filters (Abundance file, Sequences file, output\_biom, output\_excluded, output\_summary)

**Affiliation**  
 Wide choice of filters. We advise to filter OTU abundances at 0.005% [4].  
 Double taxonomic affiliation RDP classifier [5] until species. Blast+ [6] with equal multi-hits. Databanks: Silva, Greengenes → FROGS Affiliation OTU (OTU seed sequence, abundance file, biom\_affiliation, summary)

### Optional tools for a complete analysis

Access: <http://galaxy-workbench.toulouse.inra.fr> Interconnectable with BIOM format

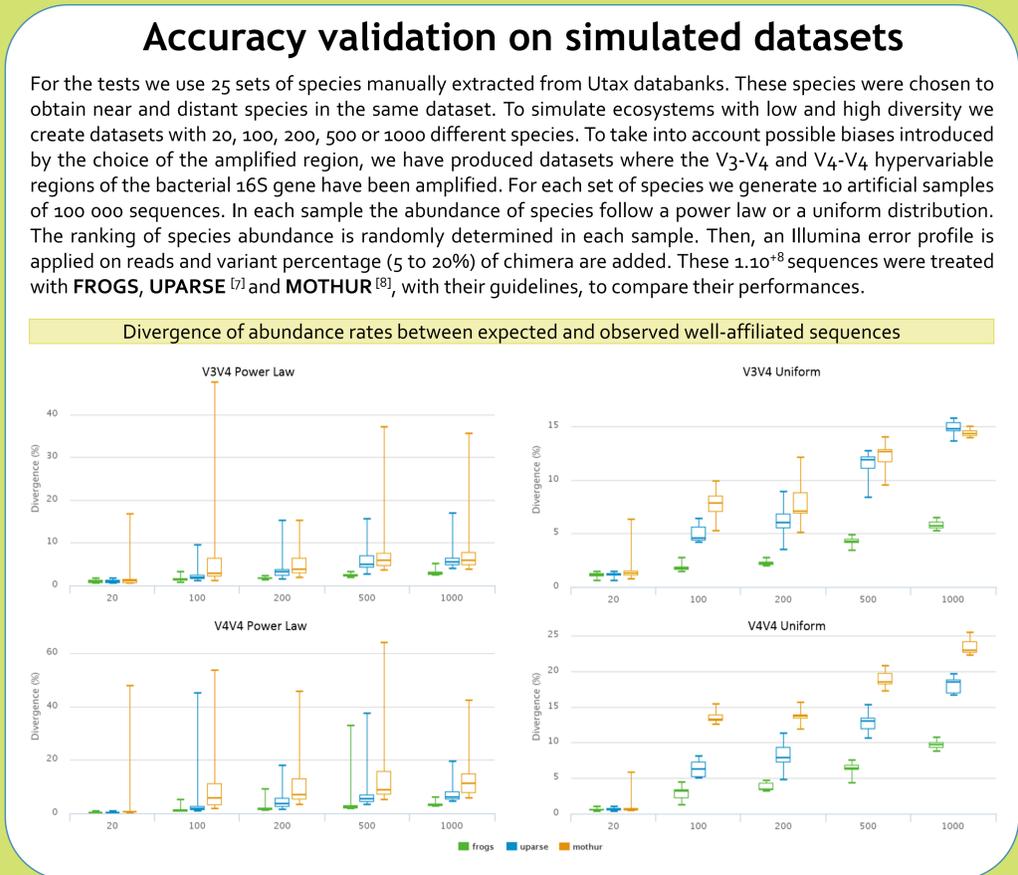
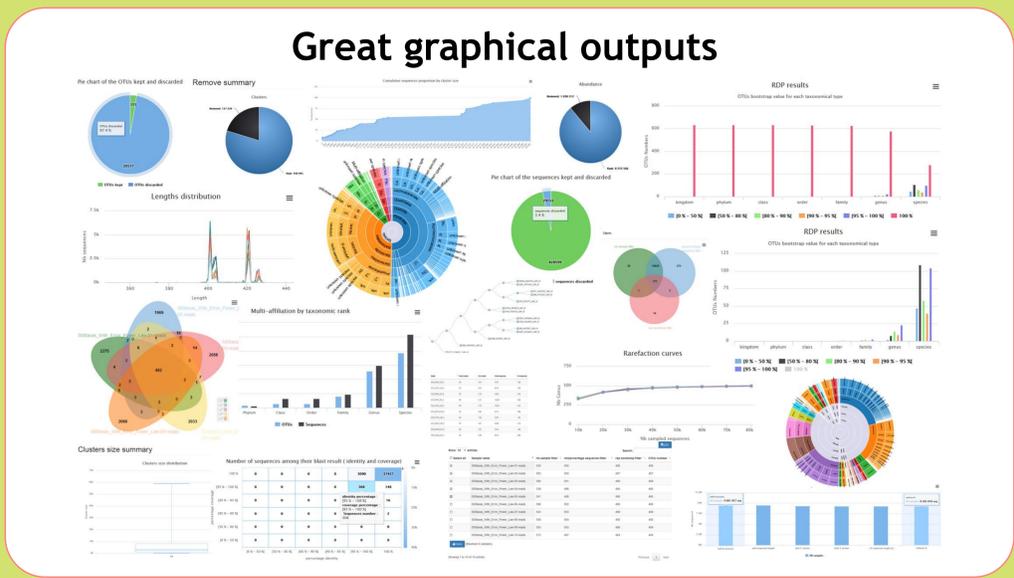
**Pre-process**: FROGS Pre-process, FROGS Clusters stat, FROGS Affiliations stat

**Clustering**: FROGS Clustering swarm

**Chimera**: FROGS Remove chimera, FROGS Demultiplex reads

**Filters**: FROGS Filters, FROGS Abundance normalisation

**Affiliation**: FROGS Affiliation OTU, FROGS BIOM to std BIOM, FROGS BIOM to TSV, Format converter



### Speed on real datasets

9 600 000 sequences of a complete MiSeq run

- Preprocess : 9 300 000 sequences → 15 min
- Swarm clustering : 680 000 clusters → 10 hours
- Chimera removal : 556 700 non-chimeric cl. → 15 min
- Filtering : 556 200 OTUs (filter OTU abundances at 0.005%)
- PhiX removal → ~8 min
- RDP affiliation → ~25 min
- Blast affiliation → ~5 min
- 500 OTUs

~ 11 hours

### Announcements

- Evaluate FROGS on others metrics and datasets (mock community, real already known community).
- Add other databases for affiliation (ITS, Midas...).
- Add FROGS in the toolshed.
- Github repository available <https://github.com/geraldinepascal/FROGS.git>
- News letter: mail to [sympa@listes.inra.fr](mailto:sympa@listes.inra.fr) with object: sub frogs-newsletter

### References

- [1] Goecks (2010) Galaxy: a comprehensive approach for supporting accessible, [...] computational research in the life sciences.
- [2] Mahé (2014) Swarm: robust and fast clustering method for amplicon-based studies.
- [3] VSEARCH GitHub repository, doi 10.5281/zenodo.15524.
- [4] Bokulich (2013) Quality-filtering vastly improves diversity estimates from Illumina amplicon sequencing
- [5] Wang (2007) Naive Bayesian classifier for rapid diversity estimates of rRNA sequences into the new bacterial taxonomy.
- [6] Altschul (1990) Basic local alignment search tool.
- [7] Edgar. (2013) UPARSE: Highly accurate OTU sequences from microbial amplicon reads.
- [8] Schloss (2009) Introducing mothur: Open-source, [...] for describing and comparing microbial communities.