



HAL
open science

Functional inference integrated in the FROGS suite

Moussa Samb, Maria Bernard, Géraldine Pascal

► **To cite this version:**

Moussa Samb, Maria Bernard, Géraldine Pascal. Functional inference integrated in the FROGS suite. JOBIM2020, Jun 2020, Montpellier, France. , 10.1038/s41587- . hal-03176828

HAL Id: hal-03176828

<https://hal.science/hal-03176828>

Submitted on 22 Mar 2021

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



Moussa SAMB¹, Maria BERNARD² and Géraldine PASCAL¹

¹ GenPhySE, Université de Toulouse, INRAE, ENVT, F-31326, Castanet-Tolosan, France

² Univ. Paris-Saclay, INRAE, AgroParisTech, GABI, SIGENAE, F-78352, Jouy-en-Josas, France

Corresponding author: geraldine.pascal@inrae.fr



Functional inference integrated in the FROGS suite

Concepts

Metabarcoding principle

Determine the diversity of an environment by amplification and sequencing of a genetic marker.



Bioinformatics analysis with FROGS

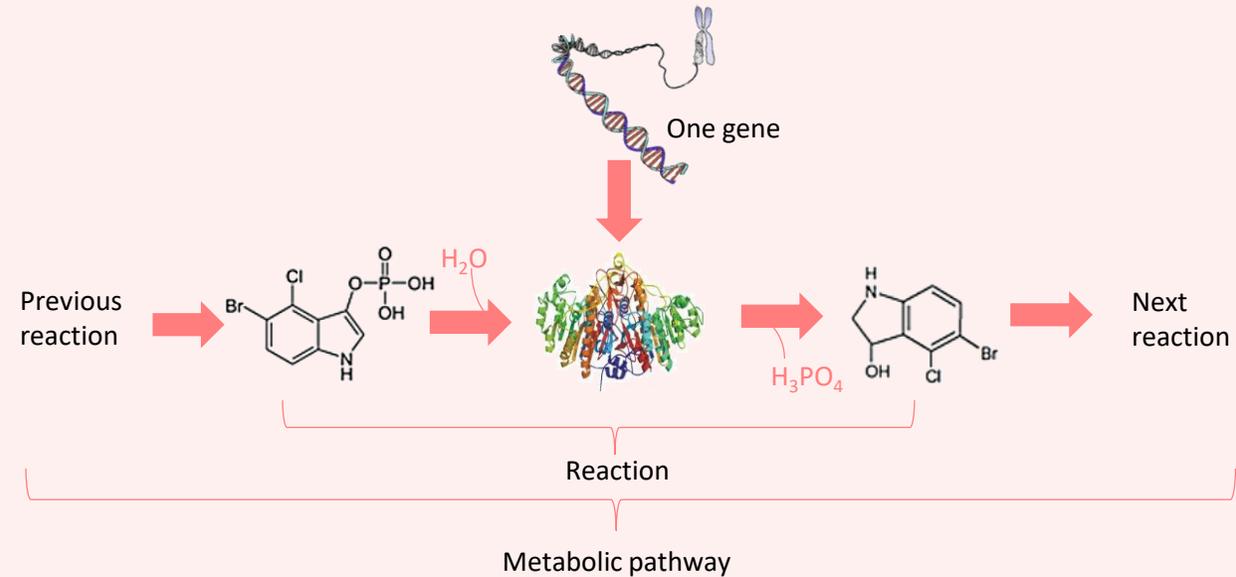


Abundance table displaying microbial diversity per samples

	Taxonomic affiliation	Sample1	Sample2	Sample3
OTU1	Species A	3500	6300	210
OTU2	Species B	0	460	36
OTU3	Species C	400	700	500

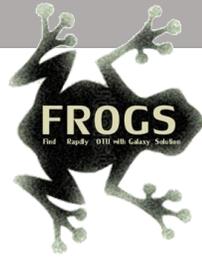
Functional inference

Assuming that an organism has a metabolic function by the presence in the organism's genome of a known sequence having that function.



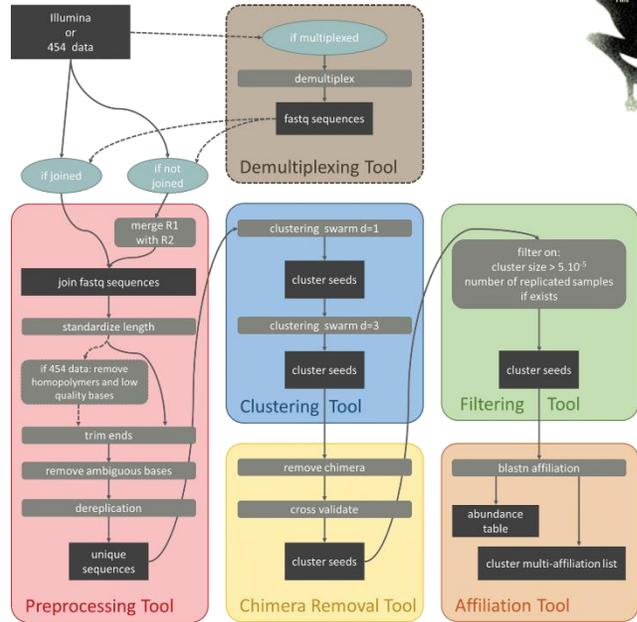
What metabolic functions are present in a microbial community ?

Our metabarcoding analysis tool

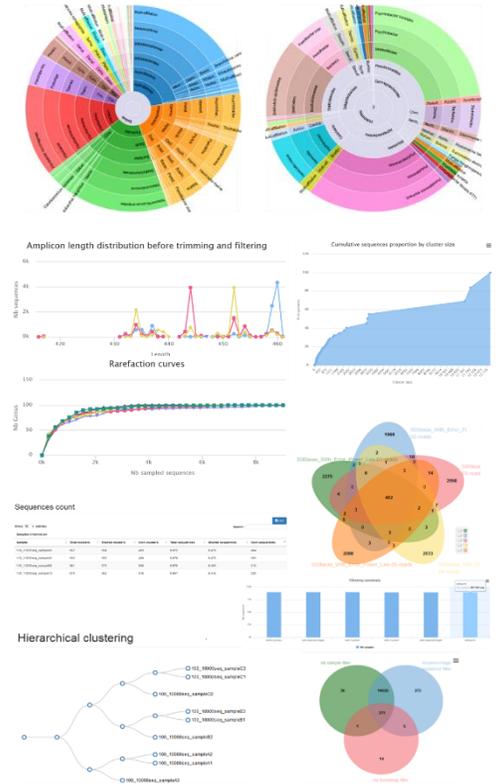


FROGS

Find, Rapidly, OTUs with Galaxy Solution
Software for analysis of metabarcoding data



- Analysis of large sets of DNA amplicon sequences (16S, ITS, rpoB...)
- Accuracy and speed
- Table of Abundance of OTUs (Operational Taxonomic Units) and Taxonomic Affiliations
- Graphical user interface under Galaxy platform (user friendly), or command line (python)
- 22 bioinformatics and statistics tools
- Treatment of ITS sequences (yeasts and fungi)
- Management of non-overlapping pairs of readings (ITS, RPOB, LSU or any other marker of size > 600 bases)
- + 3500 downloads of FROGS anaconda.org/bioconda/frogs
- At least 13 Galaxy platforms around the world offer FROGS
- + 25 000 visits on frogs.toulouse.inra.fr

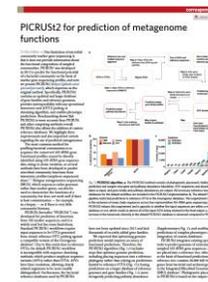


Inference functional tool

PICRUSt2

PICRUSt (Phylogenetic investigation of communities by reconstruction of unobserved states) is an open-source tool to predict metagenomic content from amplicon-based data.

The PICRUSt2 method consists of phylogenetic placement, hidden-state prediction and sample-wise gene and pathway abundance tabulation. OTU or ASV sequences and abundances are taken as input and gene family and pathway abundances are output. All necessary reference tree and trait databases for the default workflow are included in the PICRUSt2 implementation. PICRUSt2 is composed of 4 python applications. Users can run PICRUSt2 steps by steps or thanks to an inclusive command line. No graphic interface exists to run PICRUSt2 for non-expert users.

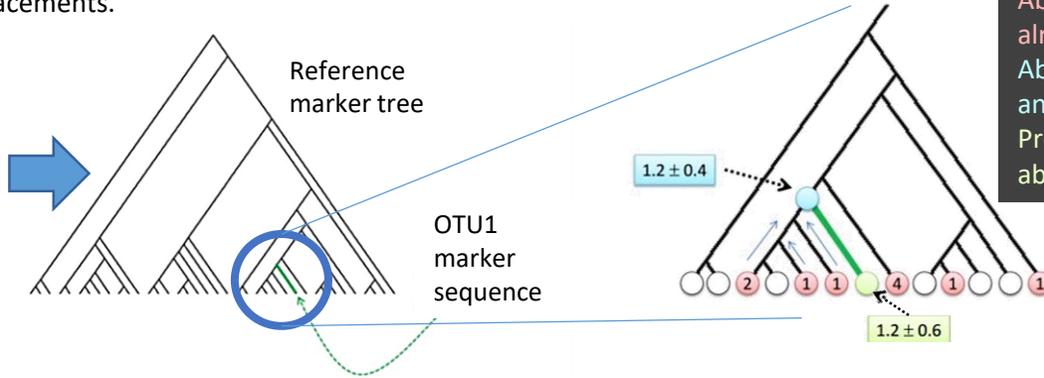


Douglas, G.M., Maffei, V.J., Zaneveld, J.R. *et al.* PICRUSt2 for prediction of metagenome functions. **Nat Biotechnol** 38, 685–688 (2020).
<https://doi.org/10.1038/s41587-020-0548-6>

1. Placement of each sequence OTUs on a reference phylogenetic tree:

From a table of OTUs, PICRUSt2 aligns the representative sequence of each OTUs on a phylogenetic tree reconstructed from the sequences of the 16S rRNA gene (or ITS or 18S). The 16S rRNA gene sequence tree is reconstructed from 41,926 16S rRNA gene sequences from IMG database genomes. Phylogenetic placement in PICRUSt2 is based on the output of three tools: HMMER (<http://www.hmmer.org>) to place OTUs, EPA-ng to determine the optimal position of these placed OTUs in a reference phylogeny, and GAPPA to output a new tree incorporating the OTU placements.

	S1	S2	S3
OTU1	3500	6300	210
OTU2	0	460	36
OTU3	400	700	500



Abundance of genes already known
Abundance of ancestral state genes
Predicted gene abundance

Repeat these steps for each gene
Repeat these steps for each OTU

2. Prediction of gene family copy numbers of OTUs and 3. Determine gene family abundance per sample.

Gene families can be inferred among KEGG Orthologs, EC number, COGs, Pfam and TIGRFAM database for 16S marker and EC number only for ITS and 18S marker. Hidden-state prediction approaches are used in PICRUSt2 to infer the genomic content of sampled sequences. The cast-

-or R package is used for core hidden-state prediction functions. OTUs are corrected by their 16S rRNA gene copy number and then multiplied by their functional predictions to produce a predicted metagenome. PICRUSt2 also provides the OTU contribution of each predicted function, allowing for taxonomy-informed statistical analyses to be conducted.

OTU abundance table

	S1	S2	S3
OTU1	3500	6300	210
OTU2	0	460	36
OTU3	400	700	500

16S copy number prediction table

	16S copy number
OTU1	7
OTU2	4
OTU3	1

Normalized OTU table

	S1	S2	S3
OTU1	3500/7	6300/7	210/7
OTU2	0/4	460/4	36/4
OTU3	400/1	700/1	500/1

4. Infer pathway abundance

Lastly, pathway abundances are inferred on the basis of structured pathway mappings. Metabolic pathway are built with MinPath tool (Ye and Doak, 2009). MetaCyc pathway abundances are calculated in PICRUSt2 through structured mappings of EC gene families to pathways.

Normalized OTU table

	S1	S2	S3
OTU1	500	900	30
OTU2	0	115	9
OTU3	400	700	500

PICRUSt EC number prediction table

	EC:1.1.1.1	EC:1.1.1.2	EC:1.1.1.3
OTU1	2	0	2
OTU2	1	0	0
OTU3	2	4	2

EC prediction per sample table

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

Metabolic pathway prediction per sample table

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

FROGS PICRUSt Step1 Tool

Analyses de données Workflow Visualize Données partagées Admin Aide Utilisateur

FROGS_PICRUSt2_Step1 phylogenetic placement of reads:
HMMER, EPA-NG, GAPPA (Galaxy Version 1.0) ☆ Favorite ▾ Options

Fasta file
2: test.fasta
The sequence file to filter (format: fasta).

Biom file
1: test.biom1
The abundance table of OTUs (format: biom).

Reference Tree
Prokaryote

Total length
0.8
Proportion of the total length of an input query, minimum coverage between input query and sequences of reference tree.

Execute

Inputs from FROGS metabarcoding treatment:

1. OTU sequences in fasta
2. Abundance table from FROGS

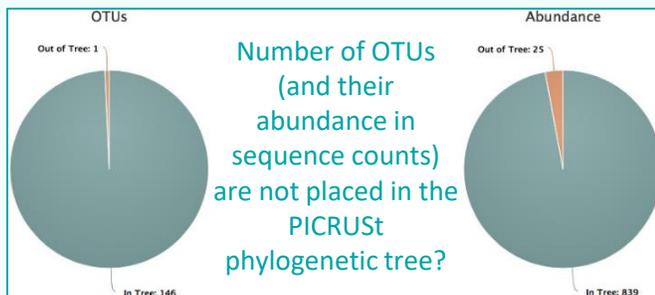
Outputs:

1. The phylogenetic tree file (Nwk format) including OTUs in the reference tree
2. A HTML report file with graphics
3. An excluded tsv file to show what OTUs could not be placed in the tree

The Newick phylogenetic tree

```
(( (2508501111:0.00005, Cluster_7:0.009689):0.469984, (2540341237:0.569658, ((2506520040:0.082301, (643228485:0.000001, 643692015:0.000001):0.05544):0.336813, (640753047:0.283299, 2734482011:0.314233):0.116214):0.075516):0.122207):0.419053, (2708742478-cluster:0.322535, 2619619054:0.567318):0.036112):0.051347):0.049542, (2737471717:0.458874, (2602041904:0.308354, (((2517572072-
```

Here, OTU7, is inserted into the reference PICRUSt tree built from IMG databank.



FROGS PICRUSt Step2 Tool

Analyses de données Workflow Visualize Données partagées Admin Aide Utilisateur

FPStep2 Predict gene copy: OTU_name, OTU_count, NSTI (Galaxy Version 1.0) ☆ Favorite ▾ Options

Tree file
158: FPStep1: tree.nwk
The tree file of OTU (format: newick).

Marker type
Prokaryote_16S

Phylogenetic distance methods
mp

Function databases
 Select/Unselect all

EC
 KO
 PFAM
 COG
 TIGRFAM
 PHENO

Execute

Users can choose between: mp, emp_prob, pic, scp, subtree_average

Nearest-Sequenced Taxon Index
The smaller the NSTI, the greater the confidence in the prediction.

Inputs :

1. The phylogenetic tree file from FPStep1
2. Abundance table from FROGS

Outputs:

1. The marker copy number prediction table (TSV).
2. Abundance table of PICRUSt function prediction (TSV). The more references selected, the wider the table will be. Here the output table will contain results for EC and KO.

16S copy number prediction table

sequence	16S_rRNA_Count	metadata_NSTI
Cluster_1	7	0.003336999999
Cluster_10	1	0.007553
Cluster_100	1	0.231634999999
Cluster_101	1	0.231634999999

PICRUSt EC number prediction table

sequence	EC:1.1.1.1	EC:1.1.1.10	EC:1.1.1
Cluster_1	2	0	0
Cluster_10	0	0	2
Cluster_100	0	0	3
Cluster_101	0	0	3

