



**HAL**  
open science

## **UTOPIA: an automatically UpdaTed, cOmPlete and consistent ITS reference dAtabase**

Sébastien Theil, Etienne Rifa

### ► **To cite this version:**

Sébastien Theil, Etienne Rifa. UTOPIA: an automatically UpdaTed, cOmPlete and consistent ITS reference dAtabase. 5. Colloque de Génomique Environnementale, Oct 2019, La Rochelle, France. <hal-02340446>

**HAL Id: hal-02340446**

**<https://hal.science/hal-02340446v1>**

Submitted on 5 Jun 2020

**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.



Distributed under a Creative Commons CC BY 4.0 - Attribution - International License

# UTOPIA : AN AUTOMATICALLY UPDATED, COMPLETE AND CONSISTENT ITS REFERENCE DATABASE

Sebastien Theil<sup>1</sup>, Etienne Rifa<sup>1</sup>

<sup>1</sup>Université Clermont Auvergne, INRA, VetAgro Sup, UMR Unité Mixte de Recherche sur le Fromage, 20 côte de Reyne, 15000 Aurillac, France

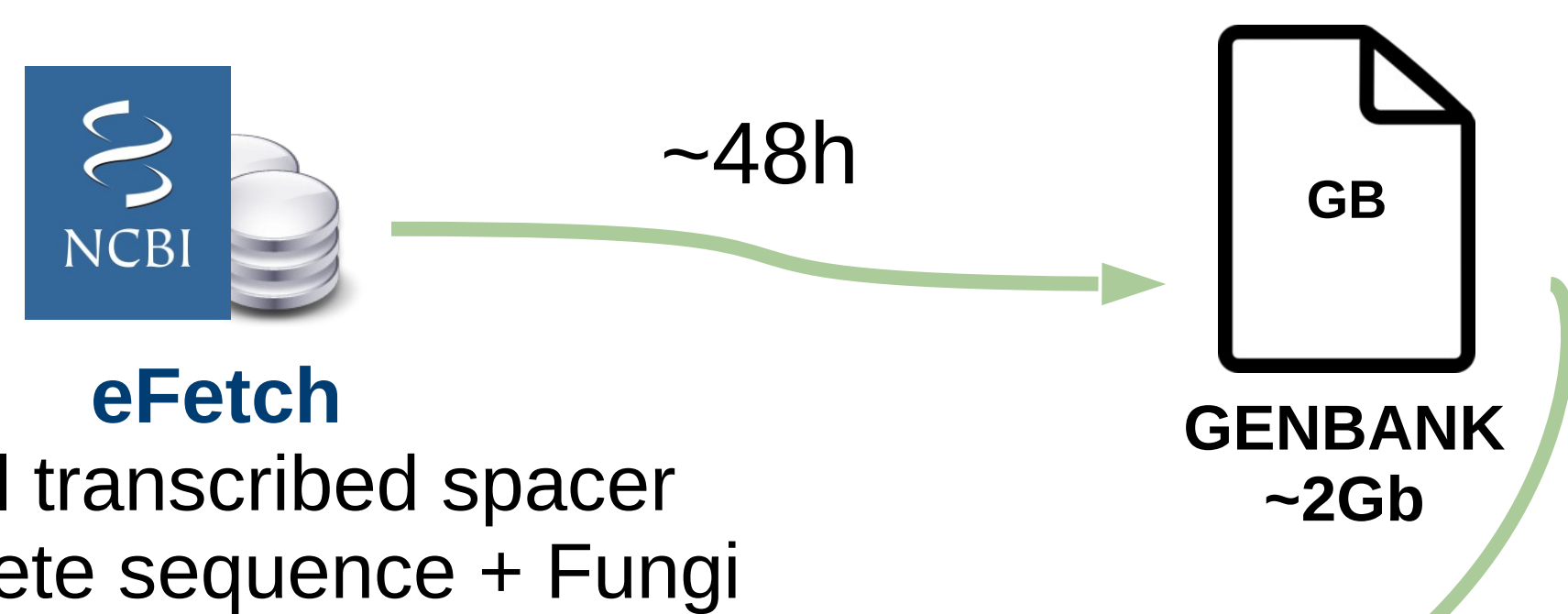
## Context

Taxonomic assignment in metabarcoding analysis is a critical and challenging step. As more organisms being sequenced, taxonomy is evolving fast with multiple taxa rearrangement and thousand of new sequences uploaded each year. The internal transcribed spacer (ITS) is an ubiquitous sequence used as a barcode to identify fungi species in complex environmental samples.

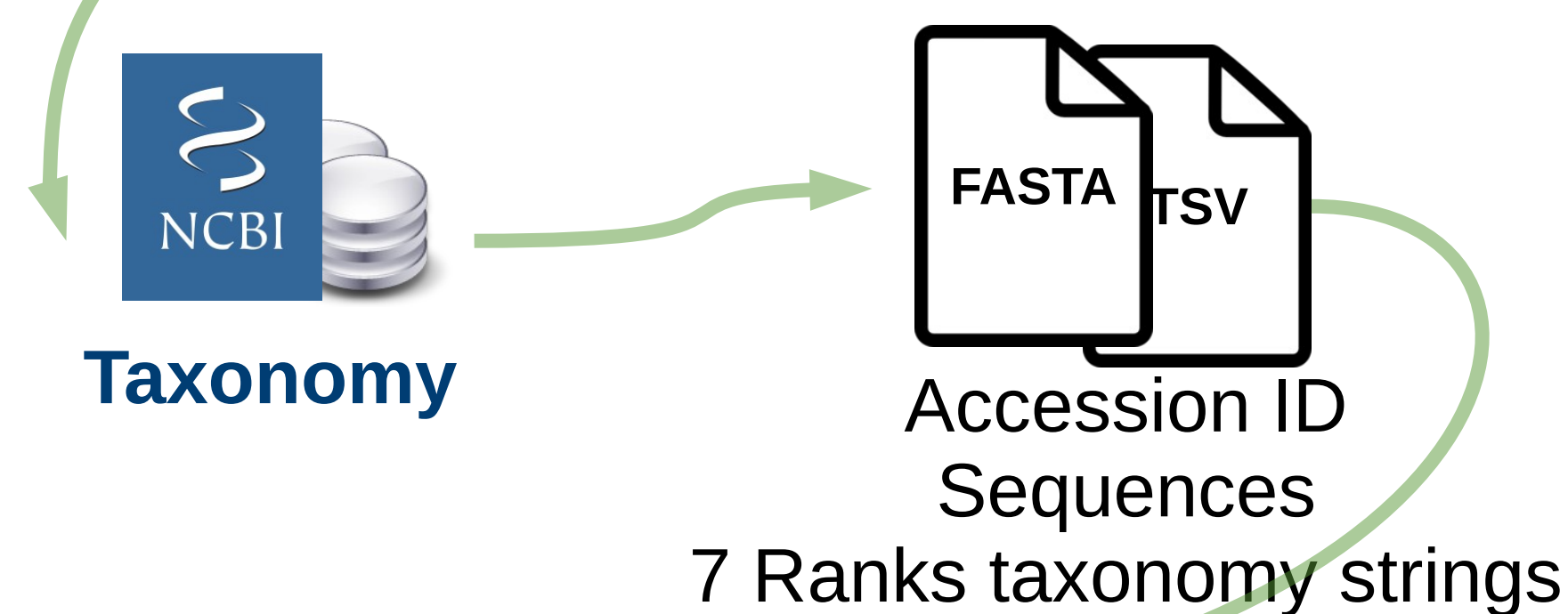
Currently used databases like UNITE, offer a good and reliable reference, but update frequency is generally low, and new strain sequences can take several years to be integrated. UTOPIA provides a workflow that produce an updated ITS reference database directly from the NCBI genbank and taxonomy database.

## UTOPIA

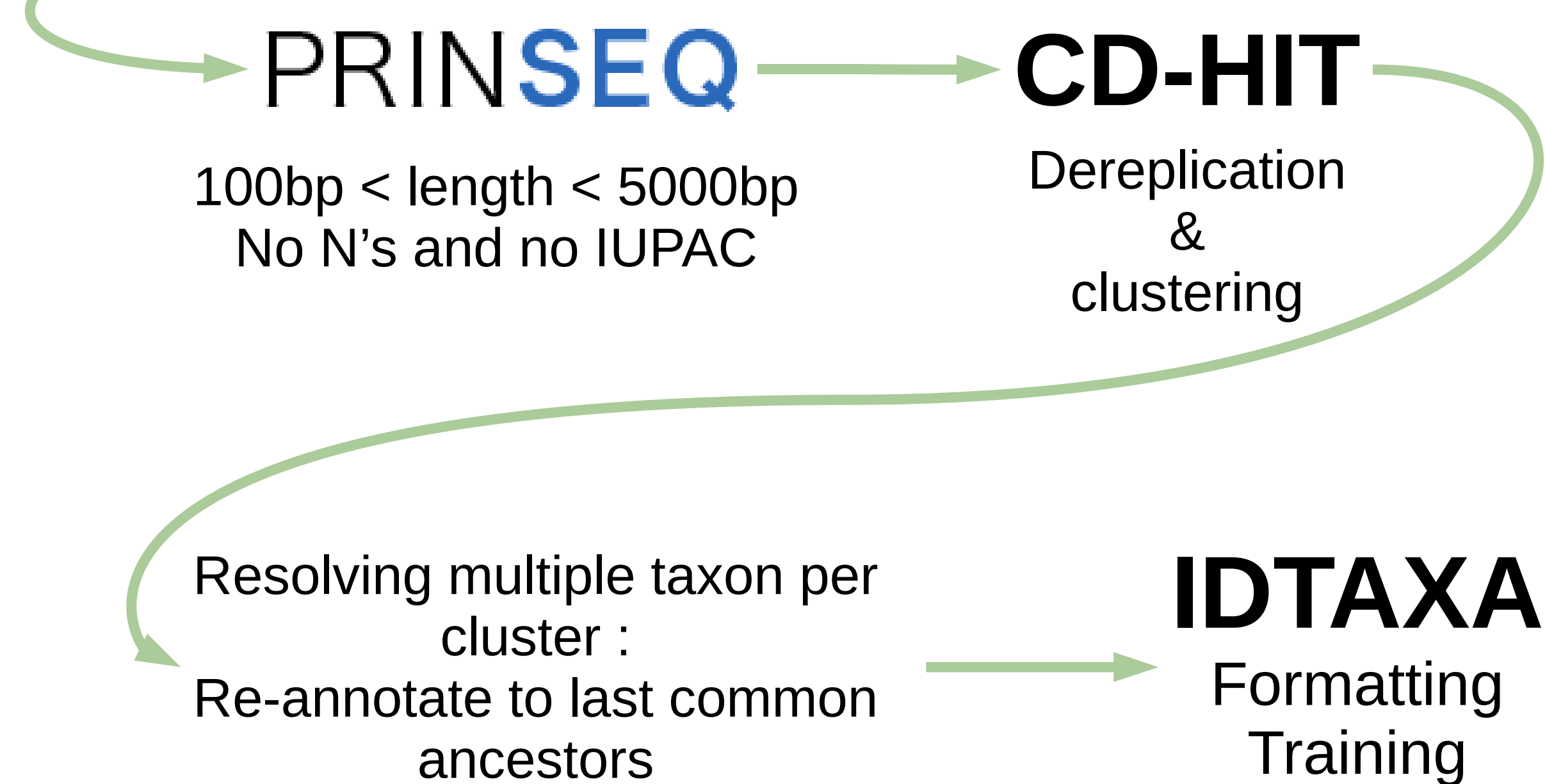
### Fetch



### Extract



### Filter



## Database statistics

### Taxon overview

Clustering threshold collapses species and few genus. Number of sequences and taxons is comparable to UNITE.

	utopia Id90, cov90	utopia Id99, cov99	Unite Dynamic	Unite Dynamic_s	unite INSDC
kindom	1	2	1	1	1
phylum	9	10	20	21	21
class	63	64	93	102	105
order	233	236	273	301	321
family	809	856	804	918	1007
genus	3739	4978	3578	4512	5256
specie	13941	42709	17278	22253	33024
nb_seq	28538	126317	35667	71042	887397

### Test on mock from Bakker *et al.* 2018

ITS1 (26 sequences)	Family				Genus			
	TPR	MCR	UCR	Acc	TPR	MCR	UCR	Acc
utopia id90, cov90	0,8	0,1	0,1	0,8	0,4	0,1	0,5	0,4
utopia id99, cov99	0,9	0,1	0	0,8	0,7	0,2	0,1	0,7
Unite Dynamic	0,7	0,1	0,2	0,7	0,6	0,1	0,3	0,5
Unite Dynamic_s	0,6	0,1	0,3	0,6	0,5	0,1	0,5	0,5
Unite INSDC	0,8	0	0,2	0,7	0,5	0,1	0,3	0,5

Metrics from **Edgar et al. 2018**:

TPR: True positive rate  
MCR: Misclassification rate  
UCR: under-classification rate  
Acc: Accuracy

Utopia has best accuracy at Genus level, especially for ITS2.

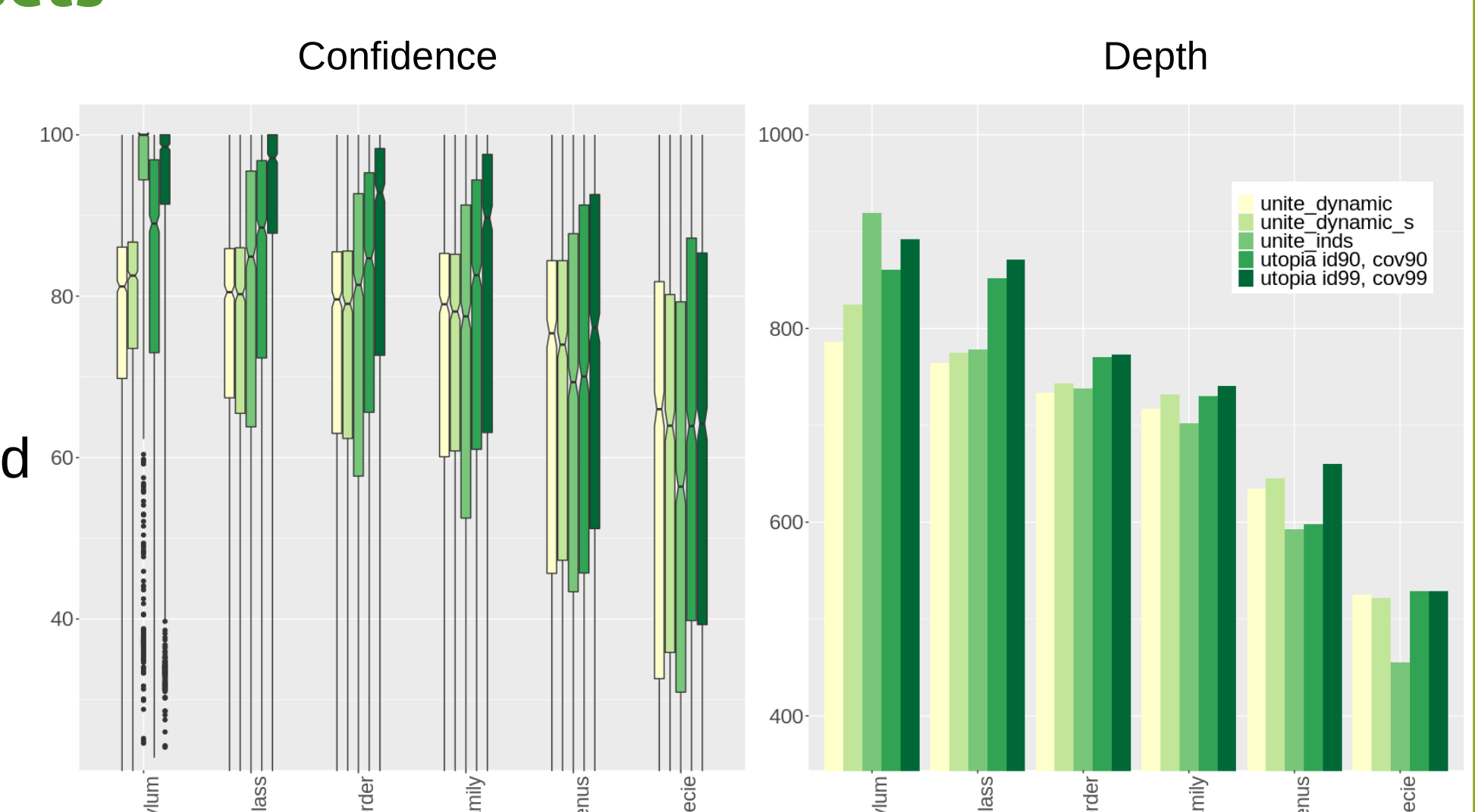
Taxonomies from NCBI match better those used by Bakker *et al.*

ITS2 (28 sequences)	Family				Genus			
	TPR	MCR	UCR	Acc	TPR	MCR	UCR	Acc
utopia id90, cov90	0,7	0,2	0	0,7	0,7	0,2	0	0,7
utopia id99, cov99	0,7	0,1	0,2	0,7	0,6	0,1	0,2	0,6
Unite Dynamic	0,4	0	0,6	0,4	0,3	0	0,7	0,3
Unite Dynamic_s	0,5	0	0,5	0,5	0,2	0	0,8	0,2
Unite INSDC	0,8	0,1	0	0,8	0,6	0,1	0,3	0,6

### Test on two real datasets

ITS1 dataset from raw milk:  
1348 sequences

Utopia has a better median confidence at family level, and performs the same for genus and species.  
Utopia99 and UNITE provide slightly the same annotation depth.



## Conclusions

We provide to the community a protocol and scripts to build an updated ITS database locally from the very last release of NCBI. Our database is created, filtered, and clustered fully automatically. Taxonomy is processed to obtain a 7 ranks assignment for each sequence and formatted to fit commonly used programs such as IDTAXA, or Phyloseq.

However we show here that UTOPIA is at least as good as UNITE in term of depth of annotation in real data and better in accuracy on the Bakker *et al.* mock community making our database a good alternative to assign fungus community with the latest sequences and taxonomies available.

Available on  repository at :  
<https://forgemia.inra.fr/umrf/utopia>