# Pollen identification by ITS2 metabarcoding: curation of the sequences retrieved from GenBank to build a reference database

Quaresma, Andreia[1,2,3]; Keller, Alexander[4]; Rufino, José[5]; van der Steen, Jozef[6]; Pinto, M. Alice[1,2]

[1]Centro de Investigação de Montanha (CIMO), Instituto Politécnico de Bragança, Campus de Santa Apolónia, 5300-253 Bragança, Portugal; [2]Laboratório Associado para a Sustentabilidade e Tecnologia em Regiões de Montanha (SusTEC), Instituto Politécnico de Bragança, Campus de Santa Apolónia, 5300-253 Bragança, Portugal; [3]Departamento de Biologia, Faculdade de Ciências da Universidade do Porto, Rua do Campo Alegre, S/N, Edifício FC4, 4169-007, Porto, Portugal; [4]Organismic and Cellular Interactions, Biocenter, Ludwig-Maximilians-Universität München, Planegg, Germany; [5]Research Centre in Digitalization and Intelligent Robotics (CeDRI), Campus de Santa Apolónia, 5300-253 Bragança, Portugal; [6]Alveus AB Consultancy, Oisterwijk, Netherlands

## Framework

A powerful way of studying the quality of the environment is by examining pollen collected by honey bees as it contains information on the availability of plant sources, and spatial and temporal floral diversity. The botanical origin of pollen is typically addressed by classical palynology, a costly approach that provides low taxonomic resolution, is time-consuming and labour intensive, and requires plant taxonomy expertise. However, as high-throughput sequencing is becoming increasingly affordable, pollen metabarcoding is gaining momentum and it is a promising alternative to classical palynology. Given that one of the main drawbacks of pollen metabarcoding is the lack of good quality reference databases for the barcode of choice, developing tools to help curating reference databases is of uttermost importance.

## BCdatabaser

BCdatabaser was used to automatically generate the standardized database, for arbitrary barcodes and taxonomic groups, from GenBank (Keller et al. 2020)

## GenBank sequences

*Viridiplantae* ITS2 sequences were downloaded from GenBank

## Problems

Misidentified sequences were detected in GenBank downloads

### Plant sequences assigned to the wrong taxa



### Fungi sequences identified as plants species



## Script

A script in bash and R was developed to curate the ITS2 global reference database for the EU space

## ITS2 database curation steps

1. Removal of Fungi sequences using the RDP Classifier Fungal ITS database

2. Pairwise alignment using vsearch v2.14.1 (Rognes et al. 2016)

3. Development of a R script for the removal of low identity alignments among plant species

4. Development of single databases for the 27 EU countries from the curated global database using Euro+Med PlantBase (https://www.emplantbase.org/home.html), and GBIF (https://www.gbif.org/) list

Some crop taxa were undetected (e.g. *Malus sp.* and *Pyrus sp.*), as revealed by comparing metabarcoding against palynological data

## Crops reference database

A reference database was developed for the European Union crops by adding more sequences/species from GenBank to overcome the identity percentage removal after the pairwise alignment

## Developed pipeline

✓ Allows easy and fast development of reference databases

✓ Facilitates regular update of reference databases

✓ Can be applied to other barcodes and organisms

### References

Keller, A., Hohlfeld, S., Kolter, A., Schultz, J., Gemeinholzer, B., Ankenbrand, M. J. (2020). BCdatabaser: on-the-fly reference database creation for (meta-) barcoding. Bioinformatics, 36(8), 2630-2631.

Rognes, T., Flouri, T., Nichols, B., Quince, C., Mahé, F. (2016). VSEARCH: a versatile open source tool for metagenomics. PeerJ, 4, e2584.