

Leveraging public genomes: a case study in *Lactobacillus*



Stijn Wittouck

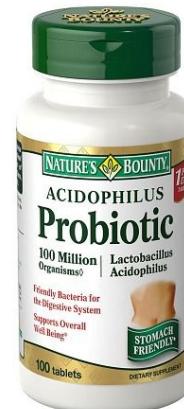
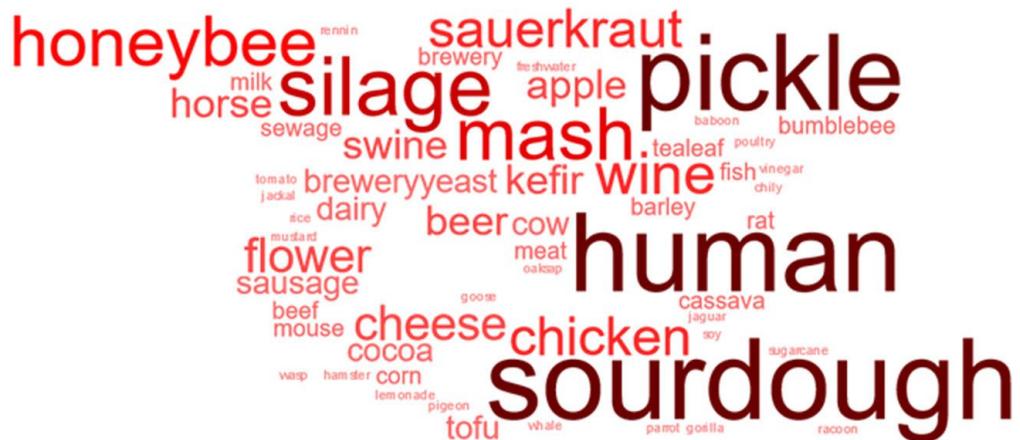
promoter: Sarah Lebeer (UA)

copromoter: Vera van Noort (KUL)



Background

Goal of PhD: study evolution of the *Lactobacillus* Genus Complex using public genomes





Challenges

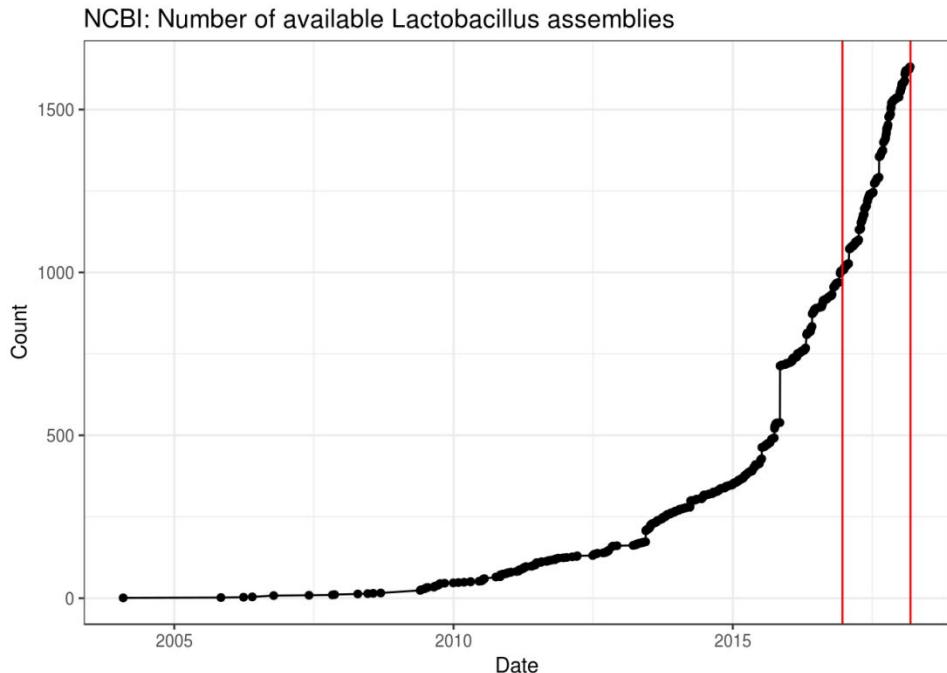
1) Data size

2,110 LGC genomes

=> 4,452,100 comparisons

5,438,039 LGC genes

=> 29,572,268,165,521 comparisons



source: blog of Sander Wuyts (swuyts.wordpress.com)

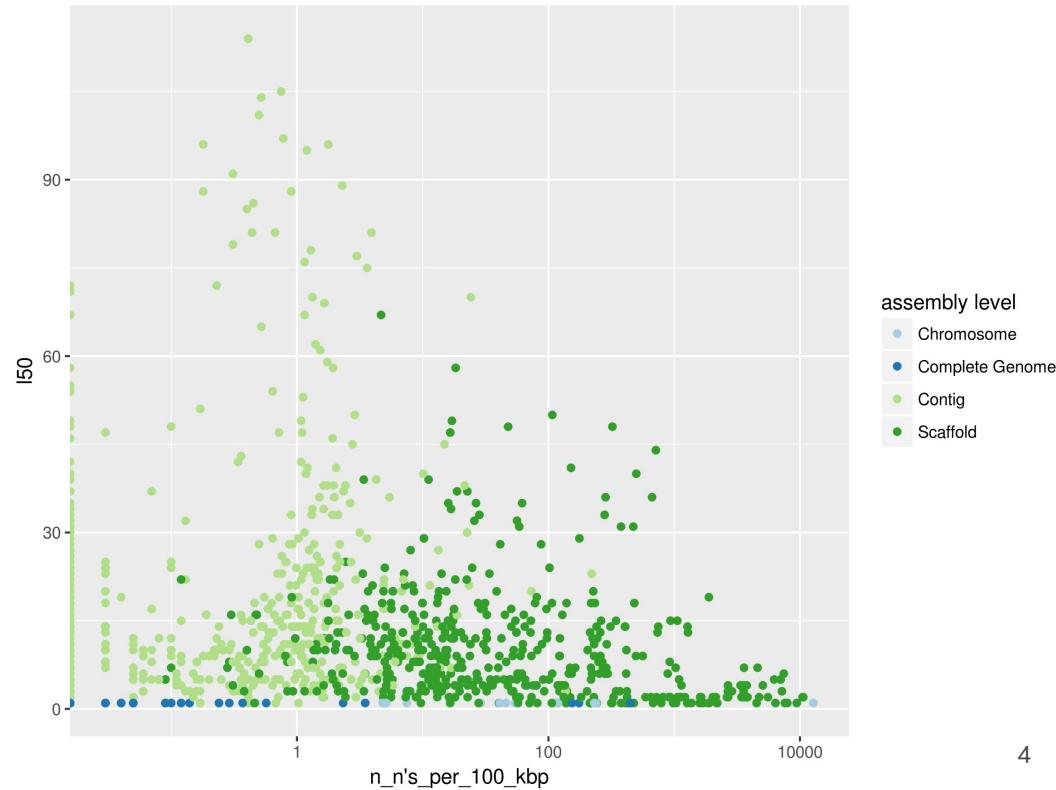


Challenges

2) Uncertain quality

Quality is very variable

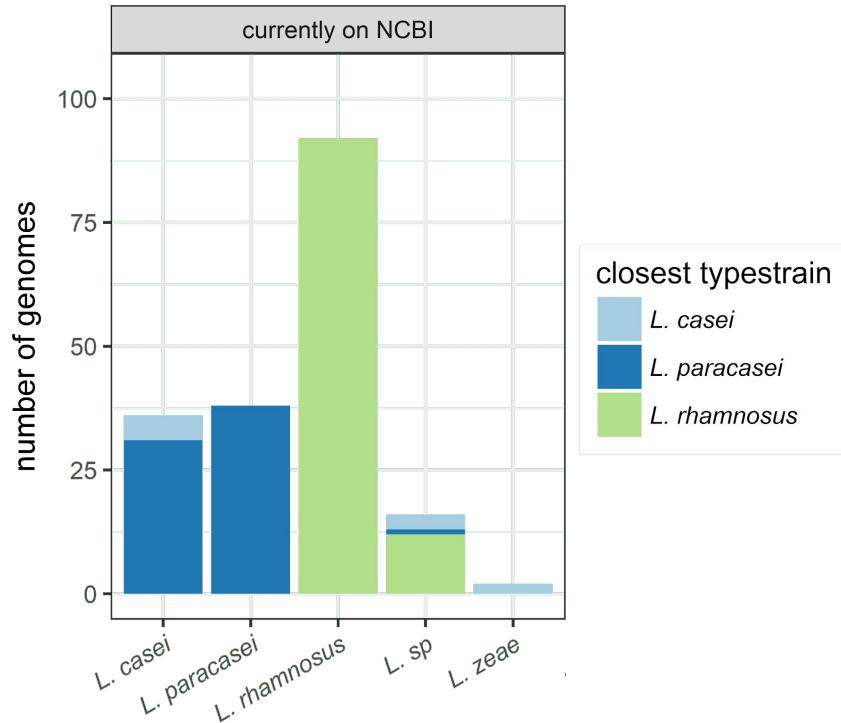
How to interpret QC measures?





Challenges

3) Uncertain taxonomy





Rapid single-copy core genes

Why?

- Intuitive genome quality control
 - Interpretation in terms of completeness and contamination
- Fast species delimitation (single-copy core nucleotide identity or SCNI)
 - Alternative to slow ANI calculations
- Phylogenetic tree inference



Rapid single-copy core genes

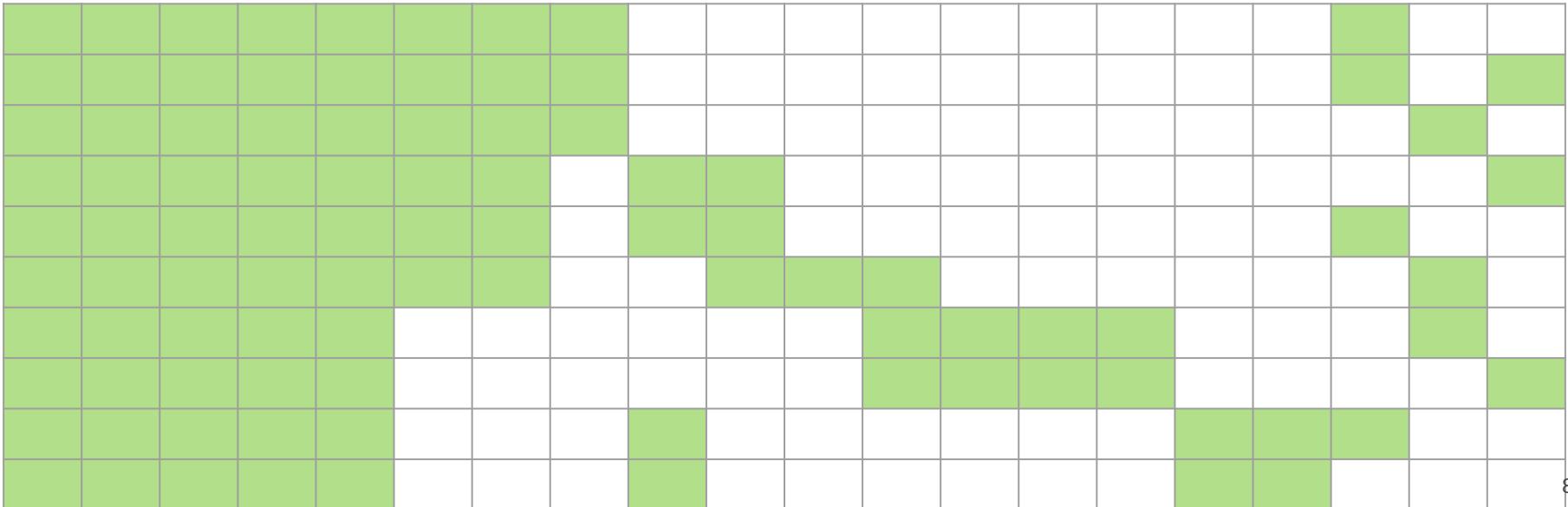
How?

Principle: **every core gene of a set of genomes is, by definition, also a core gene of a subset of those genomes**



Rapid single-copy core genes

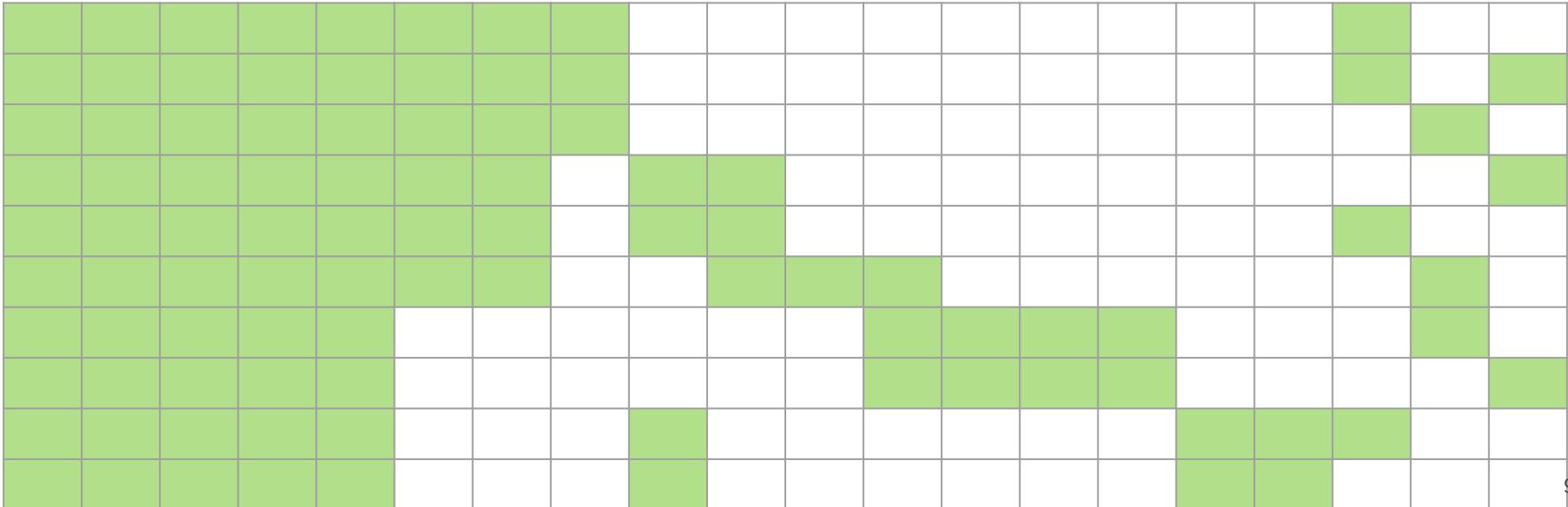
How to get all core genes?





Rapid single-copy core genes

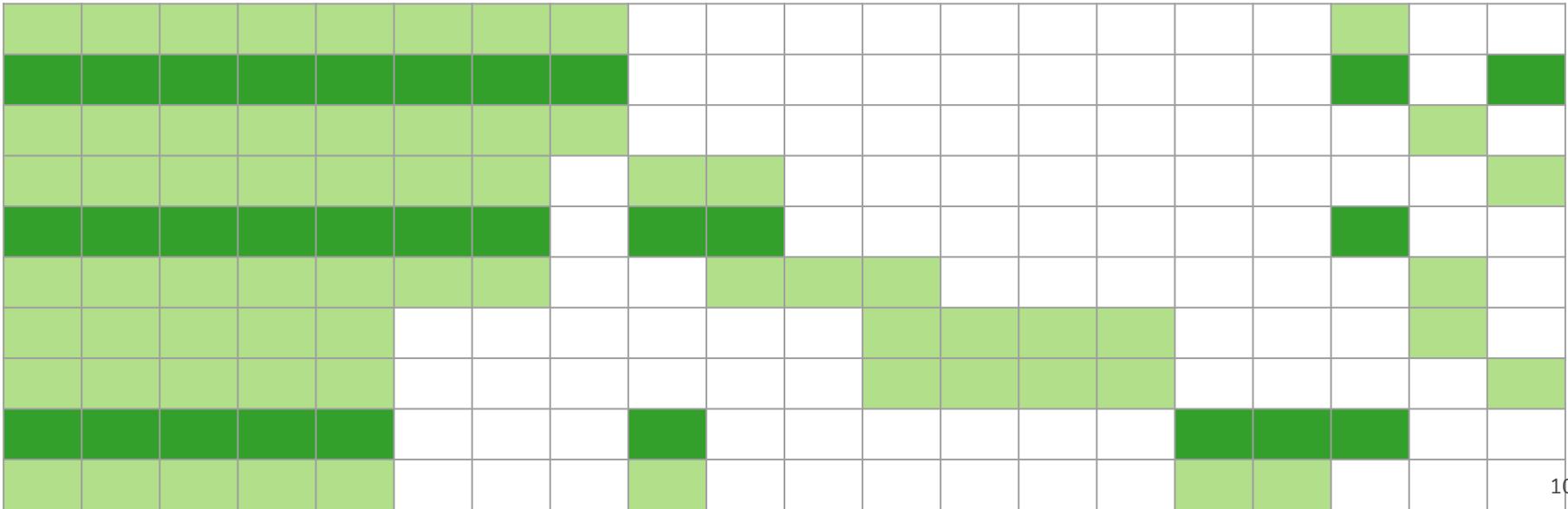
- 1) Select random “seed genomes” + make gene families





Rapid single-copy core genes

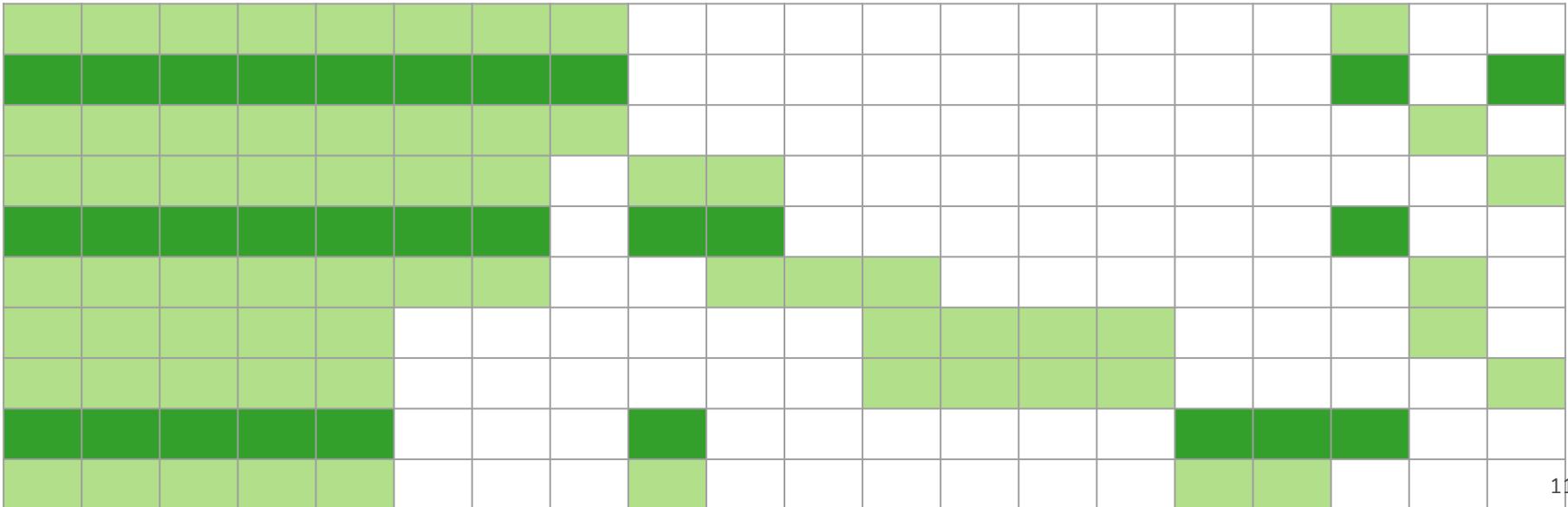
- 1) Select random “seed genomes” + make gene families





Rapid single-copy core genes

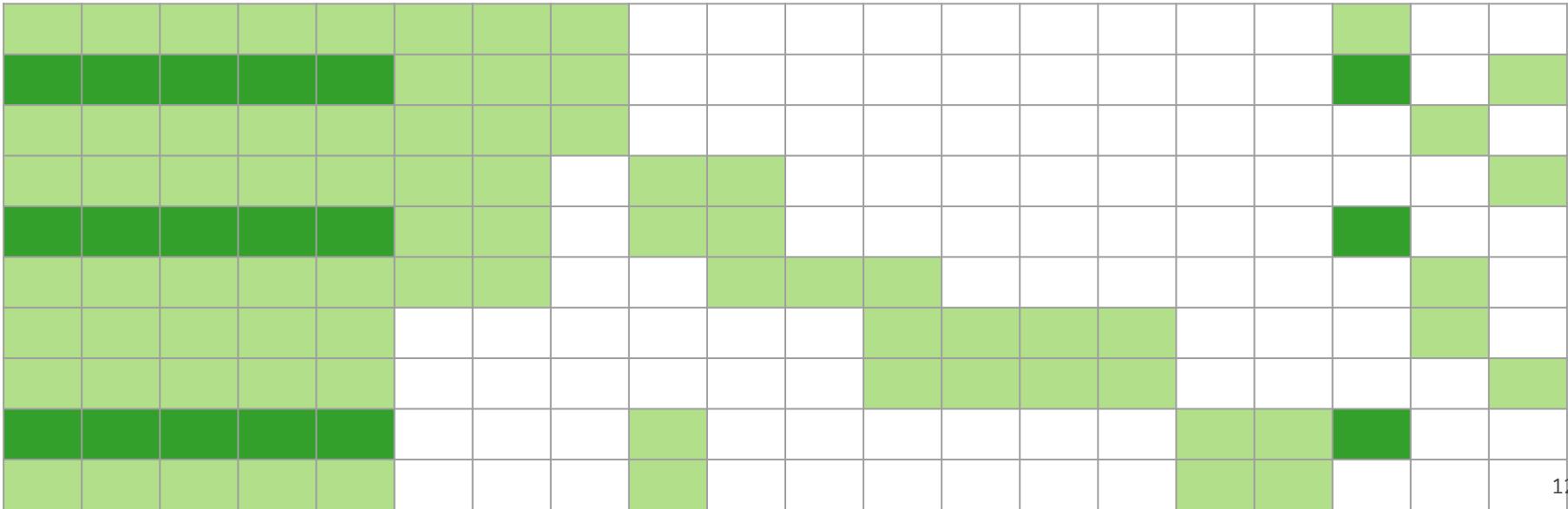
2) Keep only the “seed core genes”





Rapid single-copy core genes

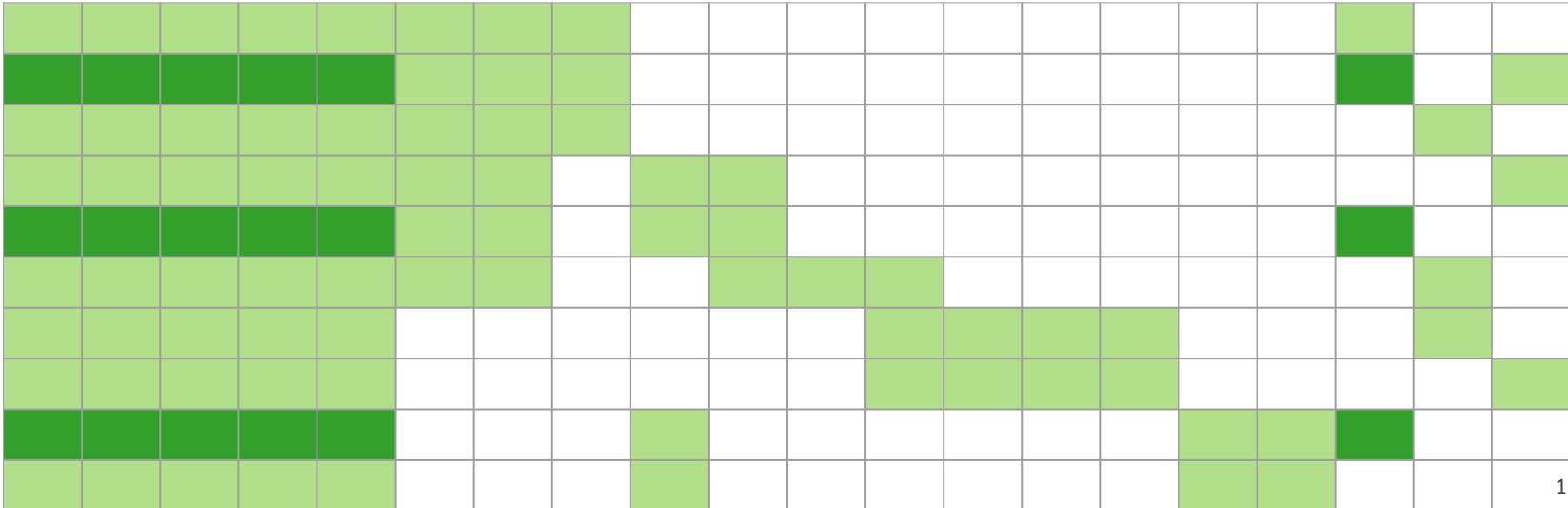
2) Keep only the “seed core genes”





Rapid single-copy core genes

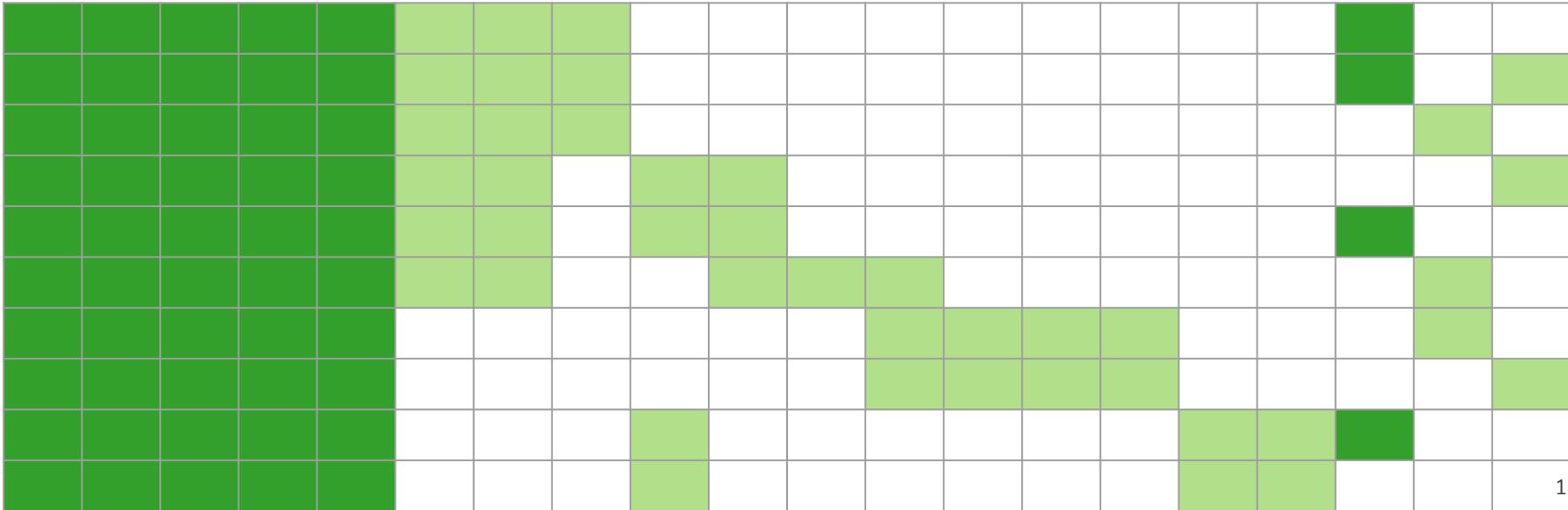
- 3) Find the seed core genes in all genomes (using profile HMMs)





Rapid single-copy core genes

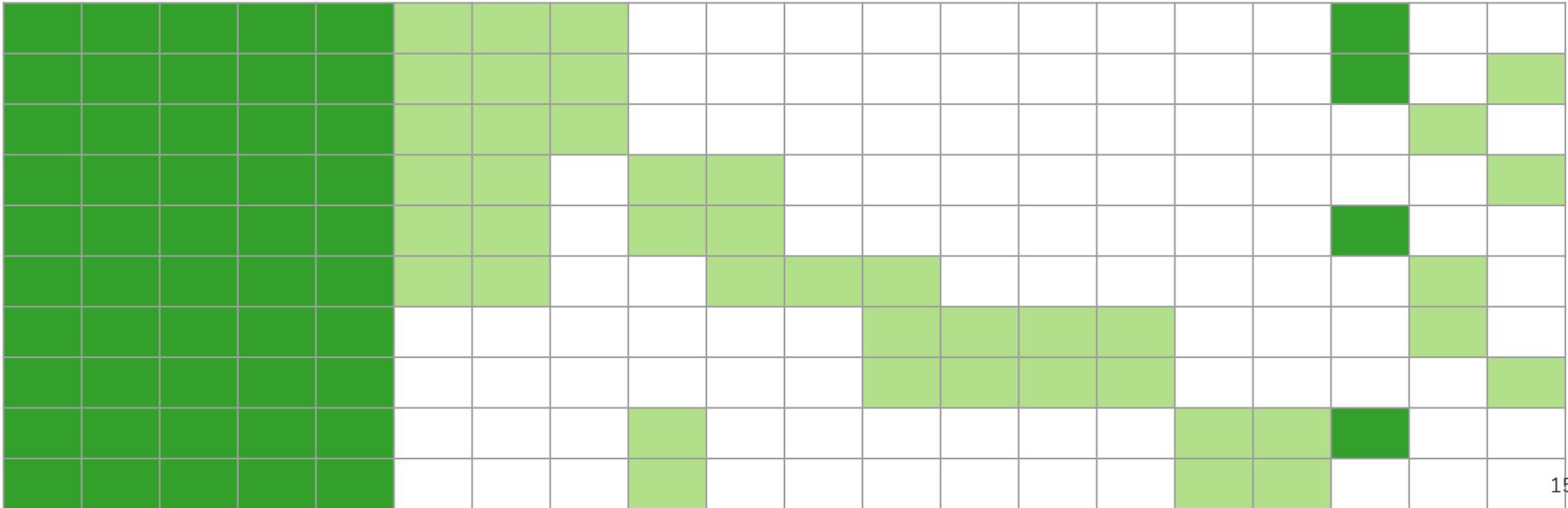
- 3) Find the seed core genes in all genomes (using profile HMMs)





Rapid single-copy core genes

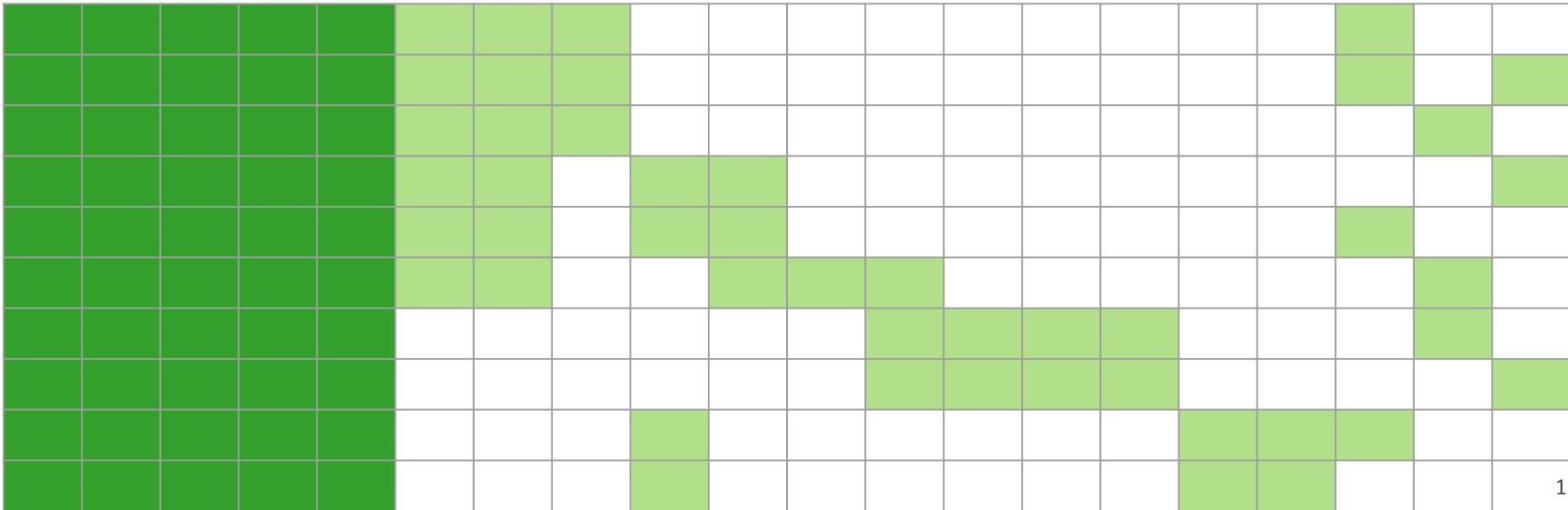
- 4) Keep only the core genes in all genomes





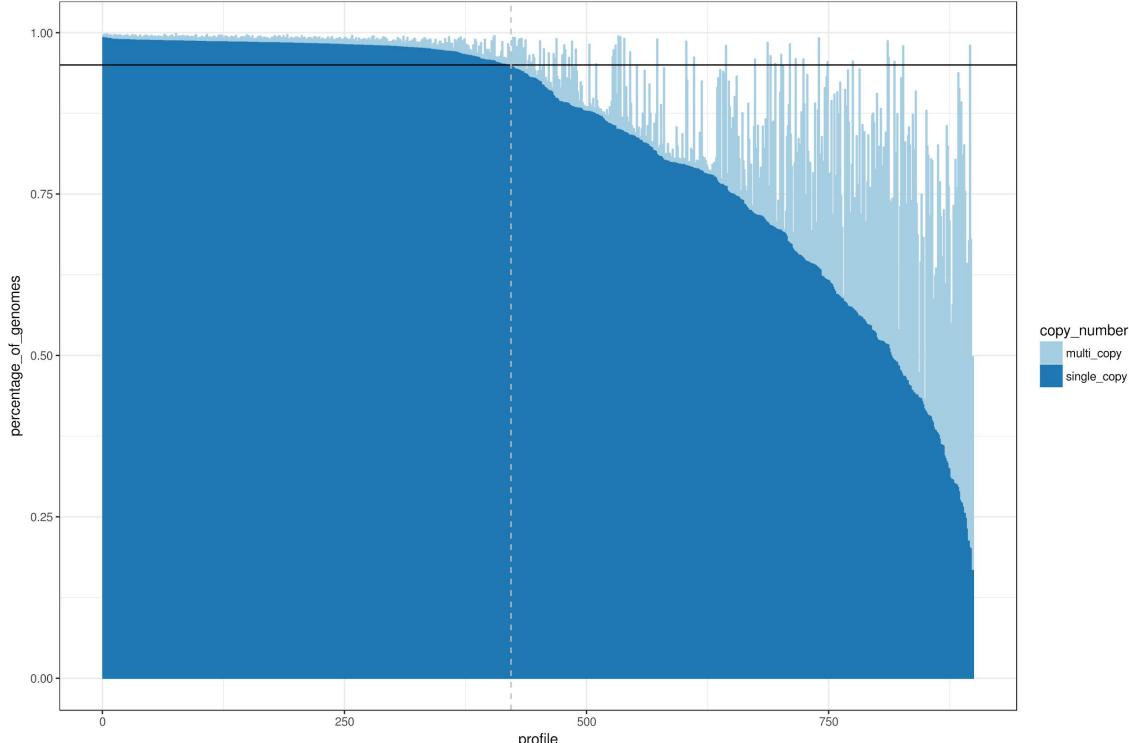
Rapid single-copy core genes

- 4) Keep only the core genes in all genomes



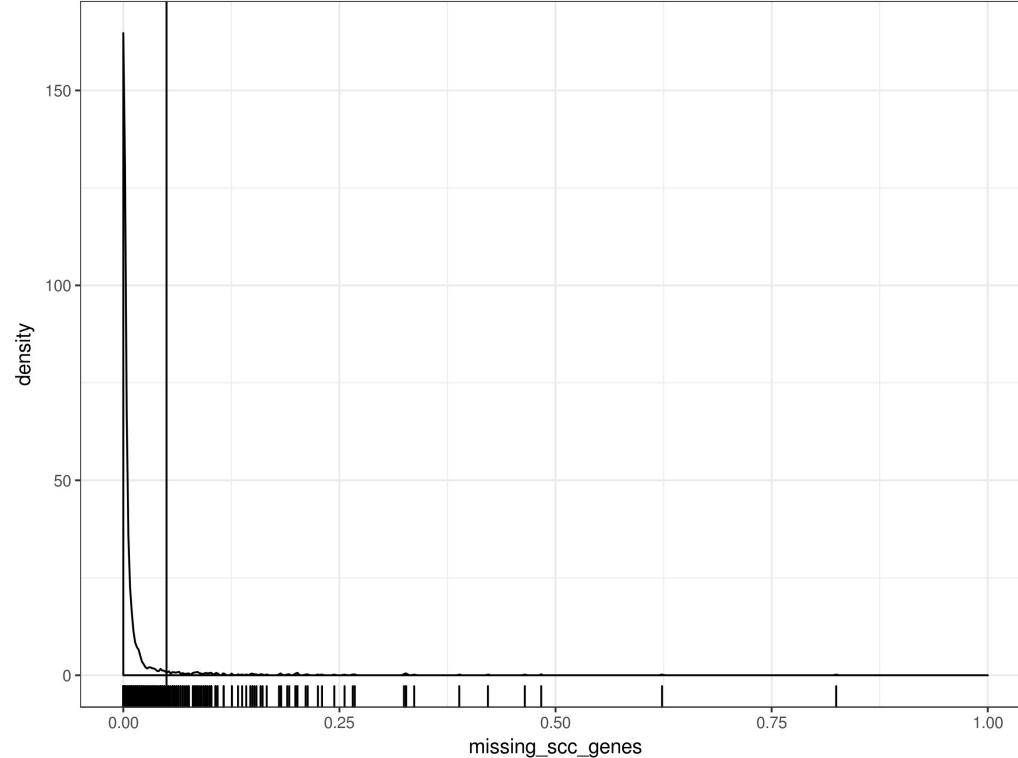


Results in the LGC



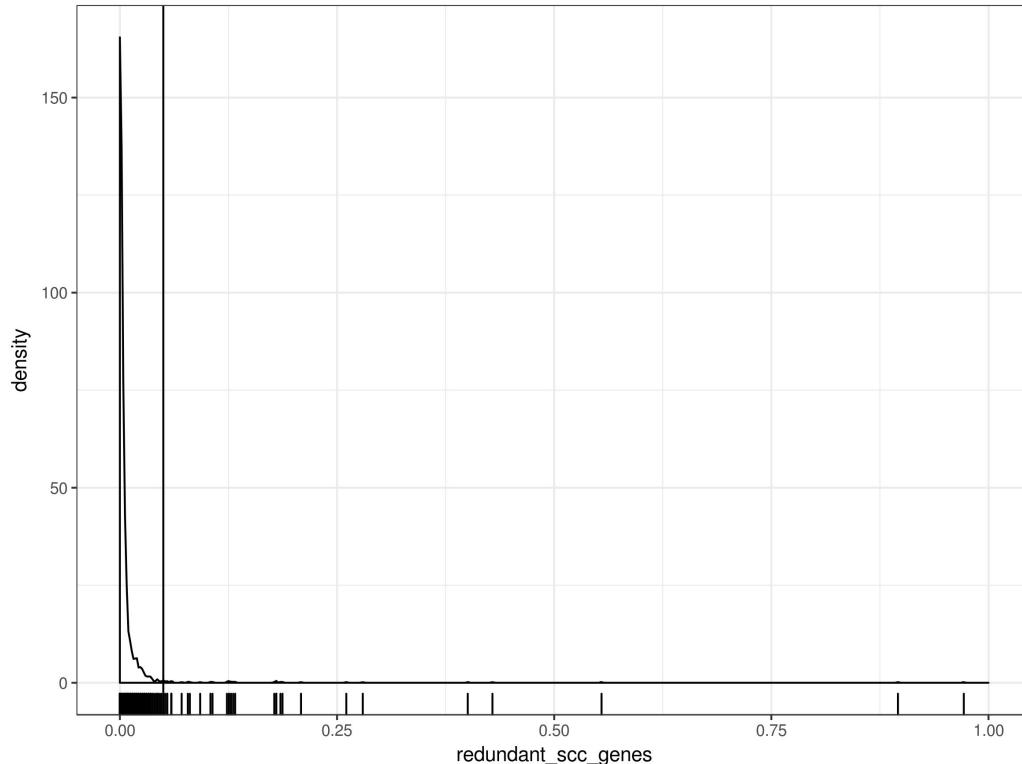


Results in the LGC



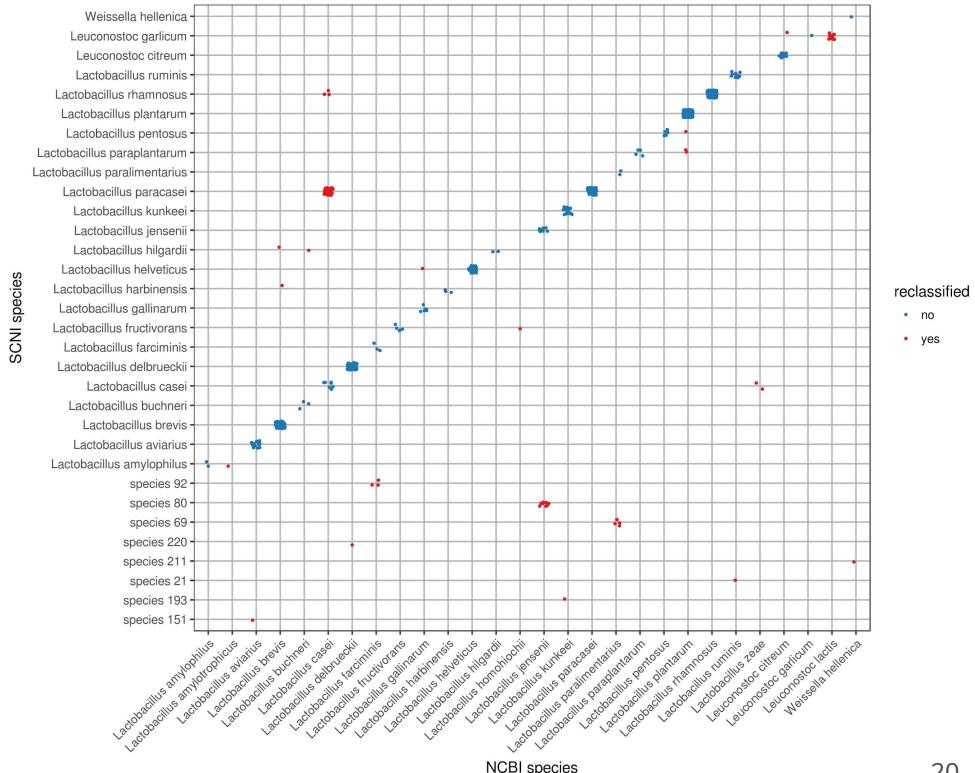
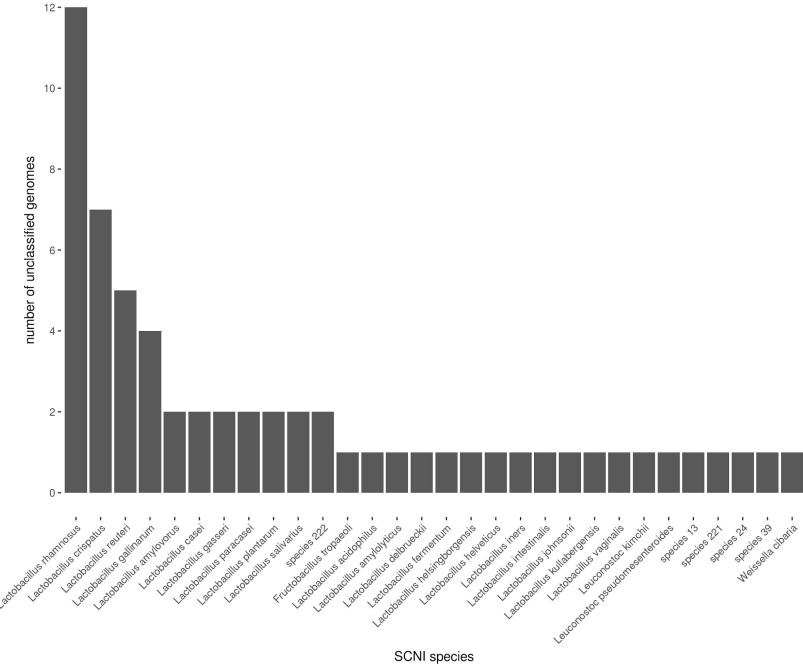


Results in the LGC





Results in the LGC





Implementation

Rapid bActerial Cross-species COmparative geNOmics

→ RACCOON

- In development
- Mostly python
- Temporarily some R
- Dependencies: orthofinder, hmmer, mafft

github.com/SWittouck/raccoon





Thank you!

prof. dr. ir. Sarah Lebeer

prof. dr. ir. Vera van Noort

dr. Conor Meehan

 ENdEMIC
Environmental Ecology and Applied Microbiology
University of Antwerp



**Research Foundation
Flanders**
Opening new horizons



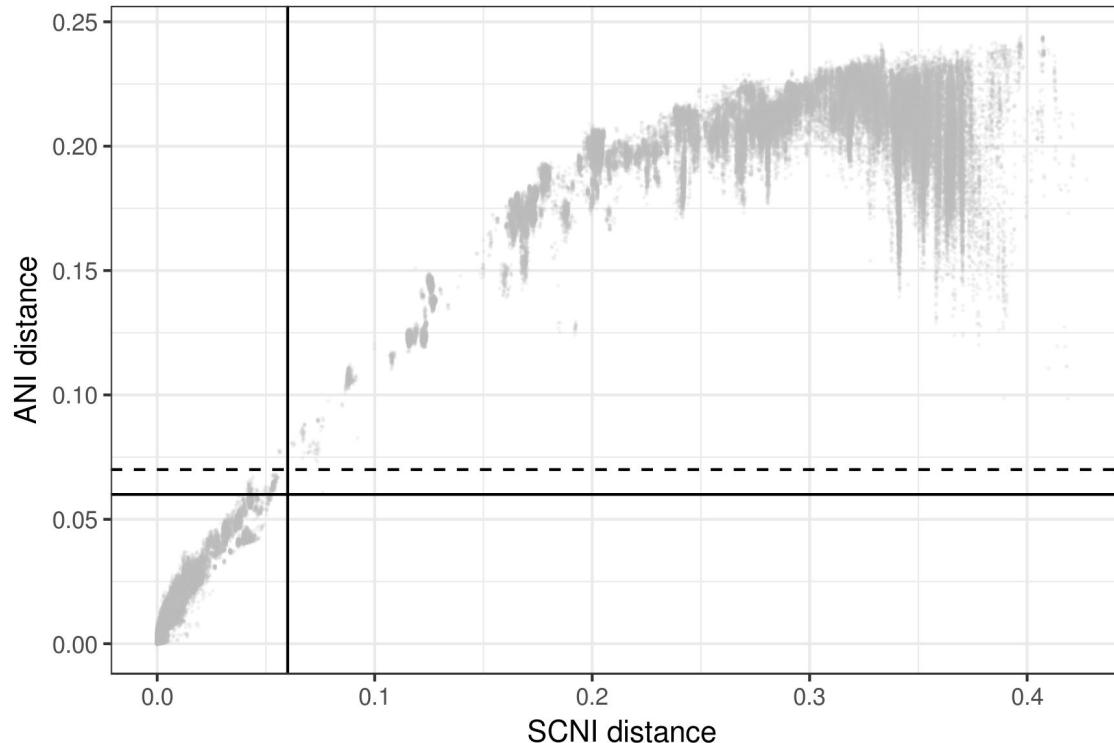
Background

Why study Lactobacillus genome evolution?

- Characterize known probiotic genes
- Characterize niche-specificity → where to find new probiotic strains?
- Characterize niche-adaptation → where to find new probiotic genes?
- Use LGC as model system for microbial evolution in general

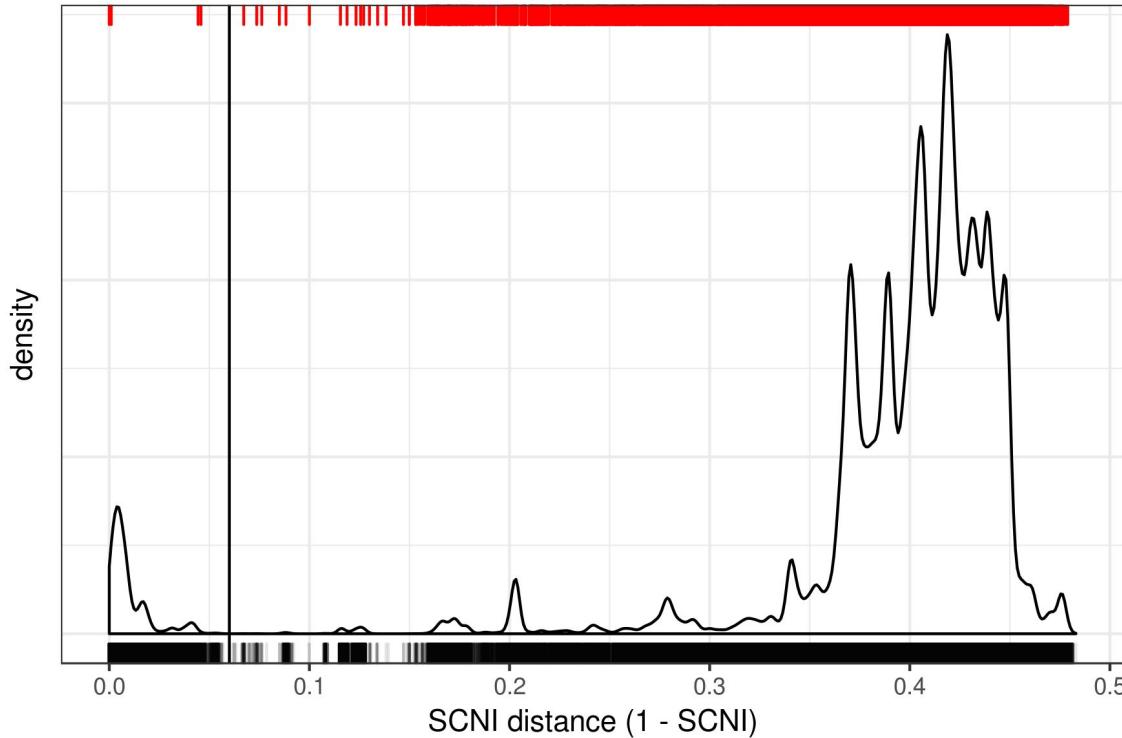


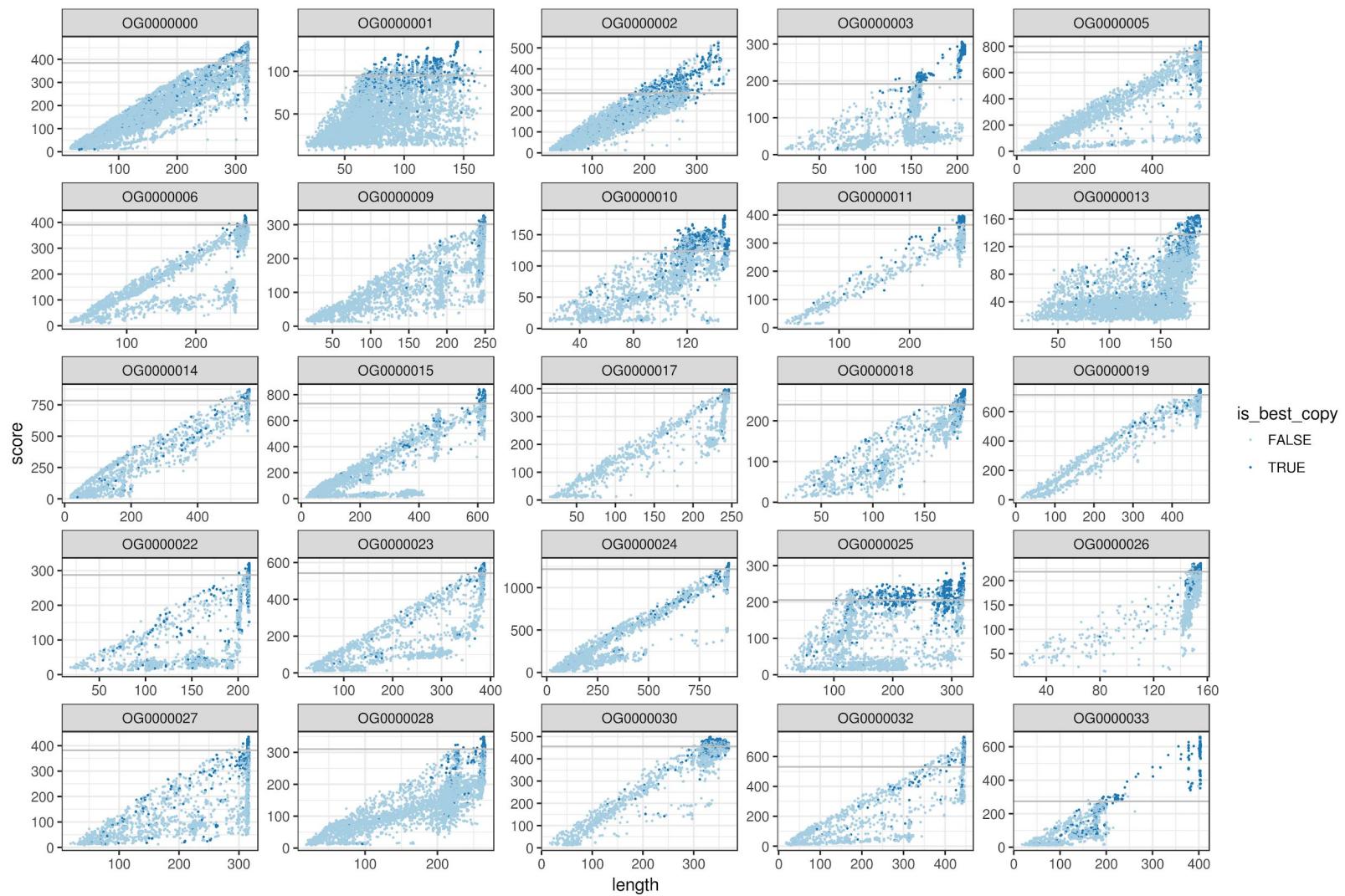
Rapid single-copy core genes





Results in the LGC







Hierarchical gene family clustering

Idea:

1. Cluster gene families independently for each species → pangenomes
2. Gather the pan-genome of each species
3. Cluster gene families on these pangenomes

Why?

- Higher resolution species-level gene families
- Faster

