# Computational Challenges in a Densely Sequenced Tree of Life

**Katie Pollard**

Gladstone Institutes

Chan Zuckerberg Biohub

UC San Francisco

Computational Challenges in Very Large-Scale 'Omics Workshop

July 21, 2022 • Simons Institute for the Theory of Computing

# A very large-scale 'omics problem

ACTGATG
CATCGAT
ATGCTAC
GATCGAT
CGATCTT
ATCGAAG

50 million sequences
300 bp each
from 100s of species mixed

Code to search for matches

ATGCATC
GATCTAC
GATCGAT
TTCGATC
AAATCGA

~300K genomes
~5 million bp each

## Problems we solved

- **50% of species have no genome: <10% now**

- **Code takes years to run or costs $10K/month in cloud: runs on laptop**
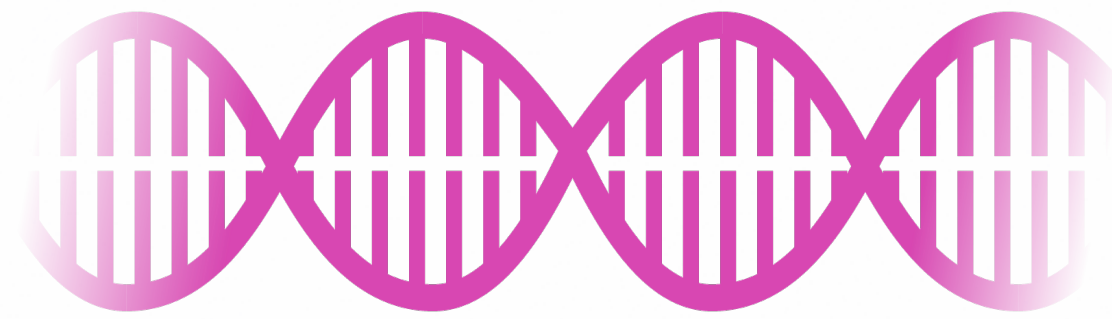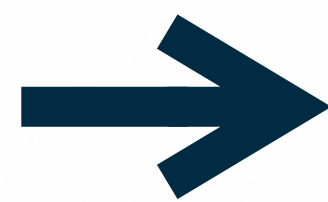
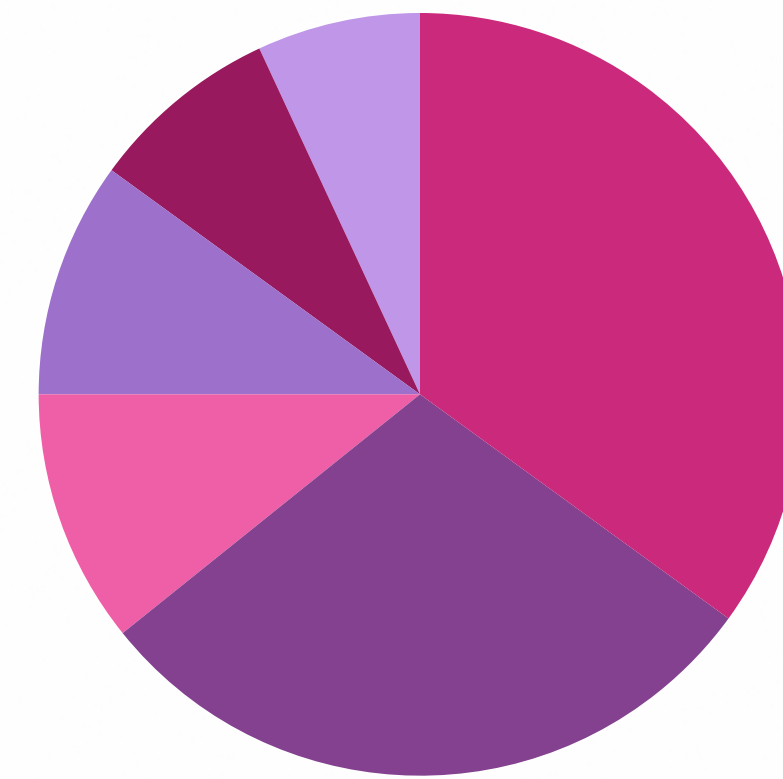- **Closely related species confound alignment and bias downstream statistics: mitigated**

Re: our AWS bill
"Let's call up
Jeff Bezos and
talk about this"
Oct 2010

Getty Images

# Metagenotyping single nucleotide variants (SNVs)



(A) COMMUNITY 1

ACGCTTC 70%

ACG**G**TTC 30%

**Similar approach for gene copy number variants (CNVs)**

## Using Genetic Variation

- **Phenotype associations**
    - human traits
    - microbe traits

- **Microbiome evolution**
    - mutation
    - selection
    - recombination
    - demography / ancestry

- **Strain / gene tracking**
- **Human evolution**
- **Genomic technologies**
- **Precision therapies**
- **Clinical decision making**

Zhao et al. (2022)                    Garud & Pollard (2019)

# Challenge 0:
# Species without a genome in the database are invisible

# Most species had no genome



Data: ~8K metagenomes from SRA, EBI, JGI
Analysis: MicrobeCensus (github/snayfach/MicrobeCensus/)

Nayfach et. al (2016)

# But this is changing
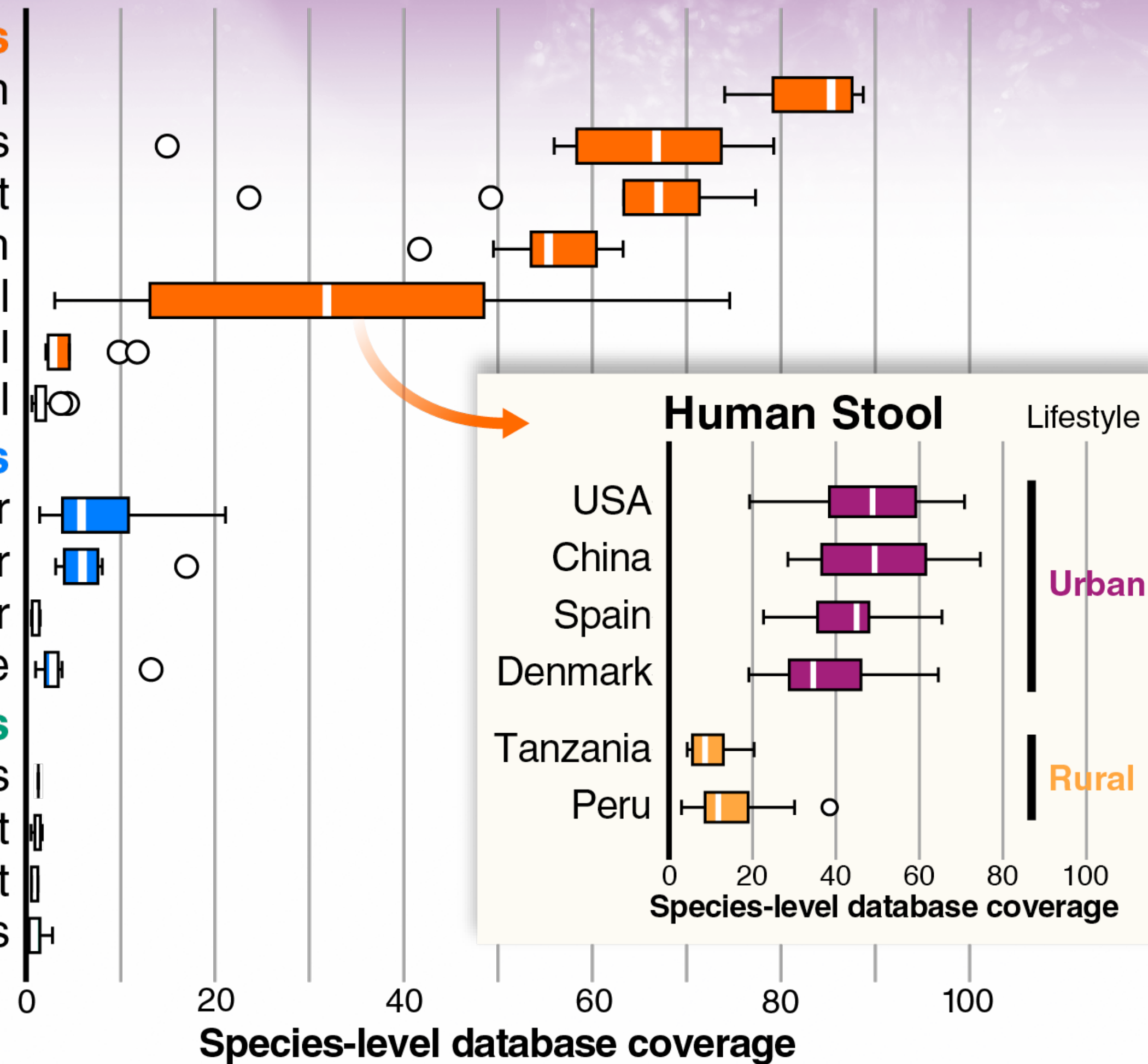
## Host-associated Metagenomes
- Human Skin
- Human Airways
- Human Urogenital Tract
- Human Mouth
- Human Stool
- Laboratory Mouse Stool
- Wild Baboon Stool

## Marine Metagenomes
- Surface Water Layer
- Dcm Layer
- Mixed Layer
- Mesopelagic Zone

## Soil Metagenomes
- Temperate Grasslands
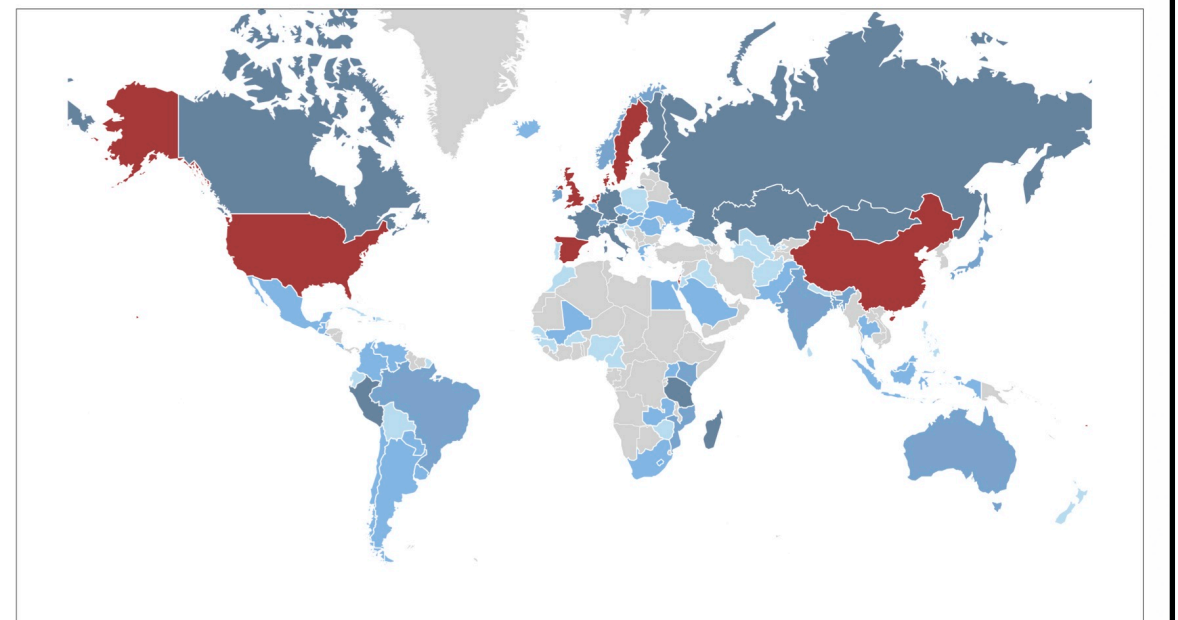- Temperate Forest
- Tropical Forest
- Deserts

**Species-level database coverage**

### Human Stool

Lifestyle

- USA
- China
- Spain
- Denmark

Urban

- Tanzania
- Peru

Rural

**Species-level database coverage**

Nayfach et. al (2016)

### UHGG Resource

Shotgun metagenomes
- 31 countries, 6 continents
- Different lifestyles & ages

Number of genomes: <10, 10–100, 101–1,000, 1,001–10,000, >10,000

286,799 gut genomes
4,644 species
81% of species MAG-only
50% increase in diversity
>2K disease associations

Nayfach et al (2019)
Almeida et al (2019)
Also: Culturomics, single-cell

**Data: ~8K metagenomes from SRA, EBI, JGI**
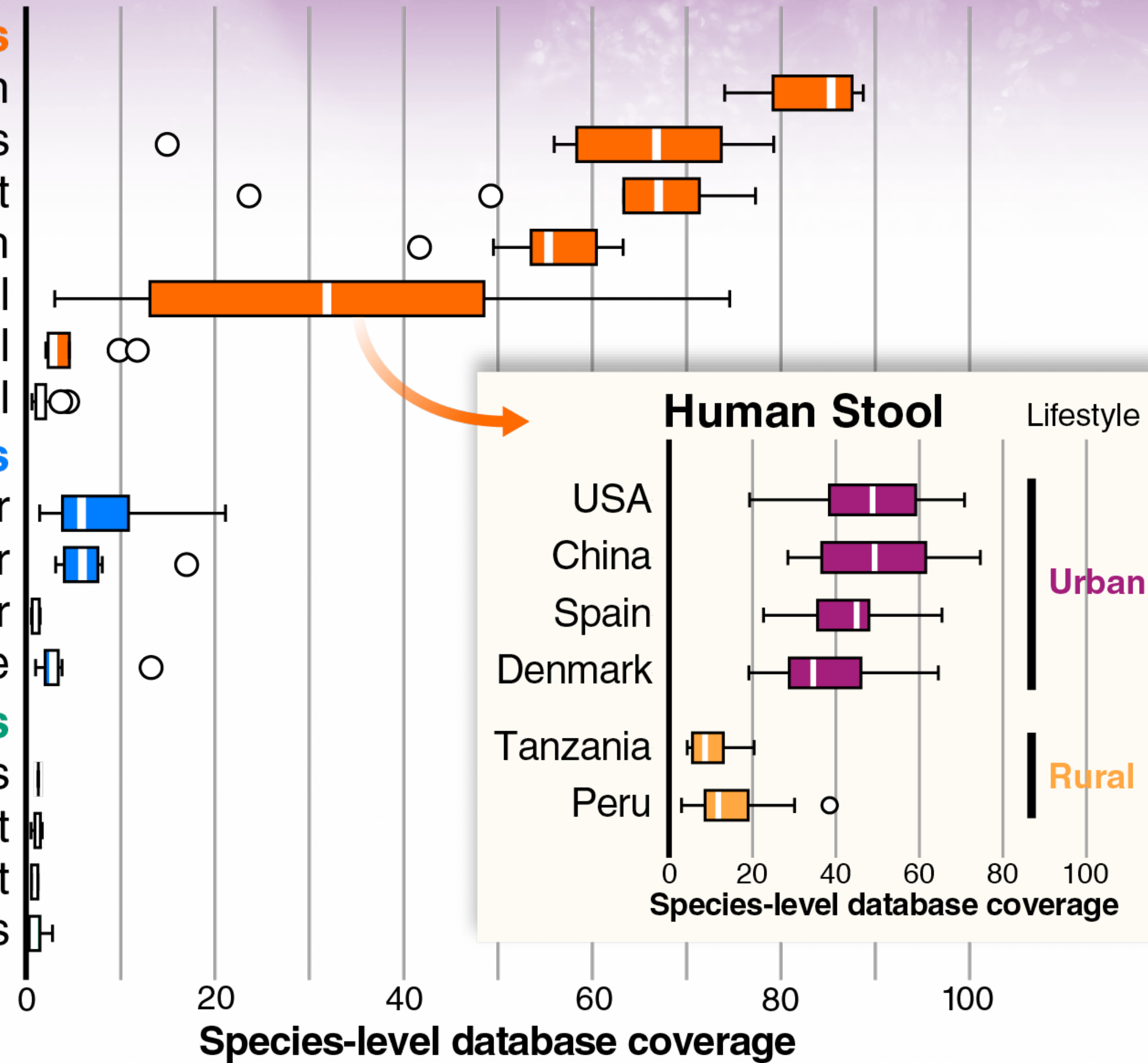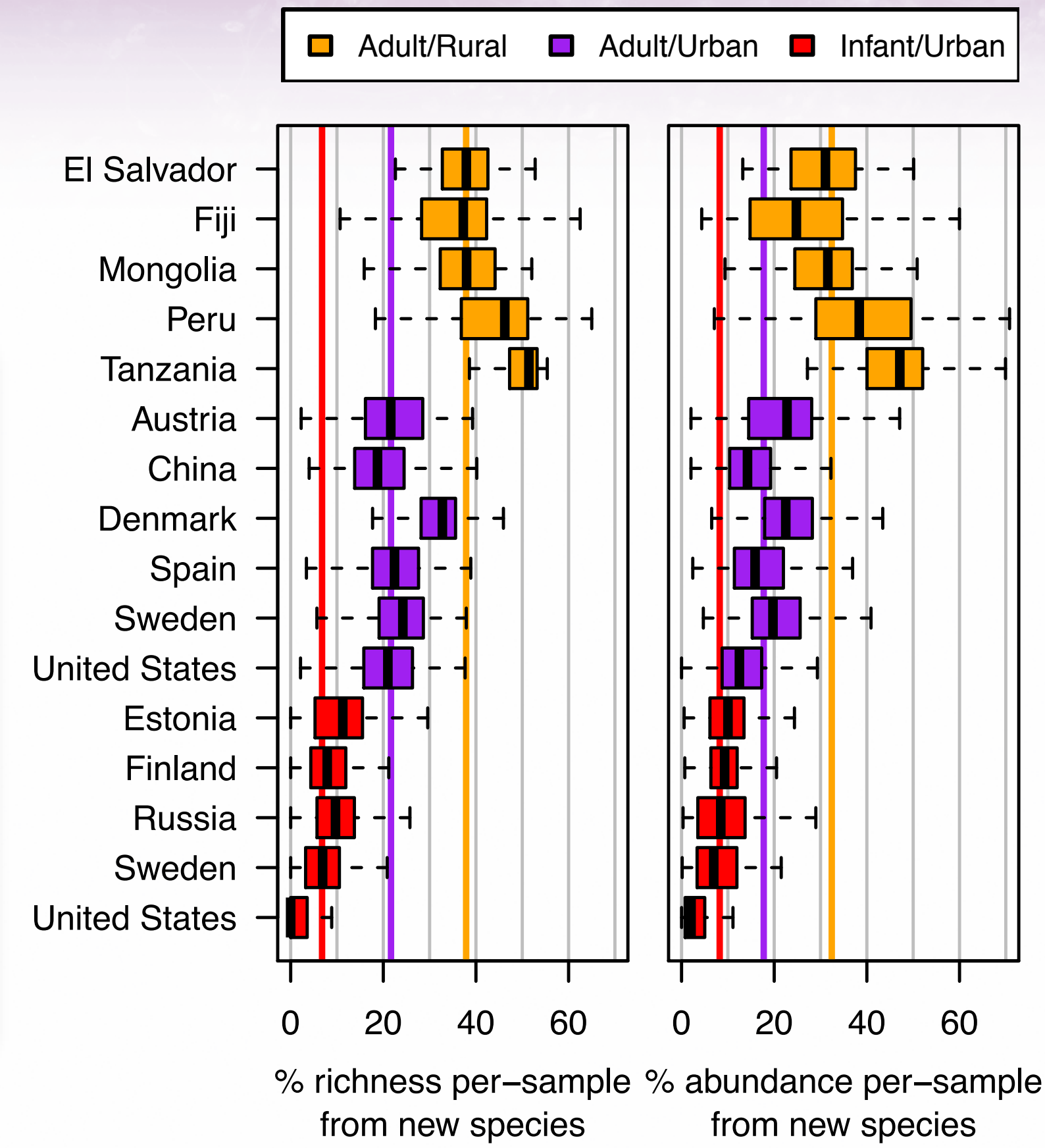**Analysis: MicrobeCensus (github/snayfach/MicrobeCensus/)**

# But this is changing



Host-associated Metagenomes, Marine Metagenomes, Soil Metagenomes box plots of Species-level database coverage, with Human Stool inset and "MAGs are closing the gap" panel.

Nayfach et. al (2016)

Nayfach et. al (2019)
Almeida et al (2019)

# Genome explosion

**NCBI Assembly**

**MAG DBs**

# More Genomes = Good News?

Human gut microbiome alignment rate now > 80%

But… new problems arise

# Challenge 1:
# Closely related species "compete" for reads and bias metagenotypes

# Closely related species are common

CRS = two species with at least one pair of genomes that have average nucleotide identity (ANI) 92%-95%



Zhao et. al (2022b)

# Read competition in dense lineages



Zhao et. al (2022b)

# Read competition in dense lineages



**Metagenomic Simulation**

Intra−species ANI
- 100
- 99
- 98
- 97
- 96
- 95

Cross−mapping rate

Average nucleotide identity of closest relative in database (%)

Zhao et. al (2022b)

# MIDAS2 mitigates low alignment uniqueness & cross-mapping



**Paired-end filtering, MAPQ<10**

**Dropping undetected species from db**

Zhao et. al (2022a)

**https://github.com/czbiohub/MIDAS2**

# Mitigation strategies help…

# But can we do better by avoiding alignment?

Gut bacteria species
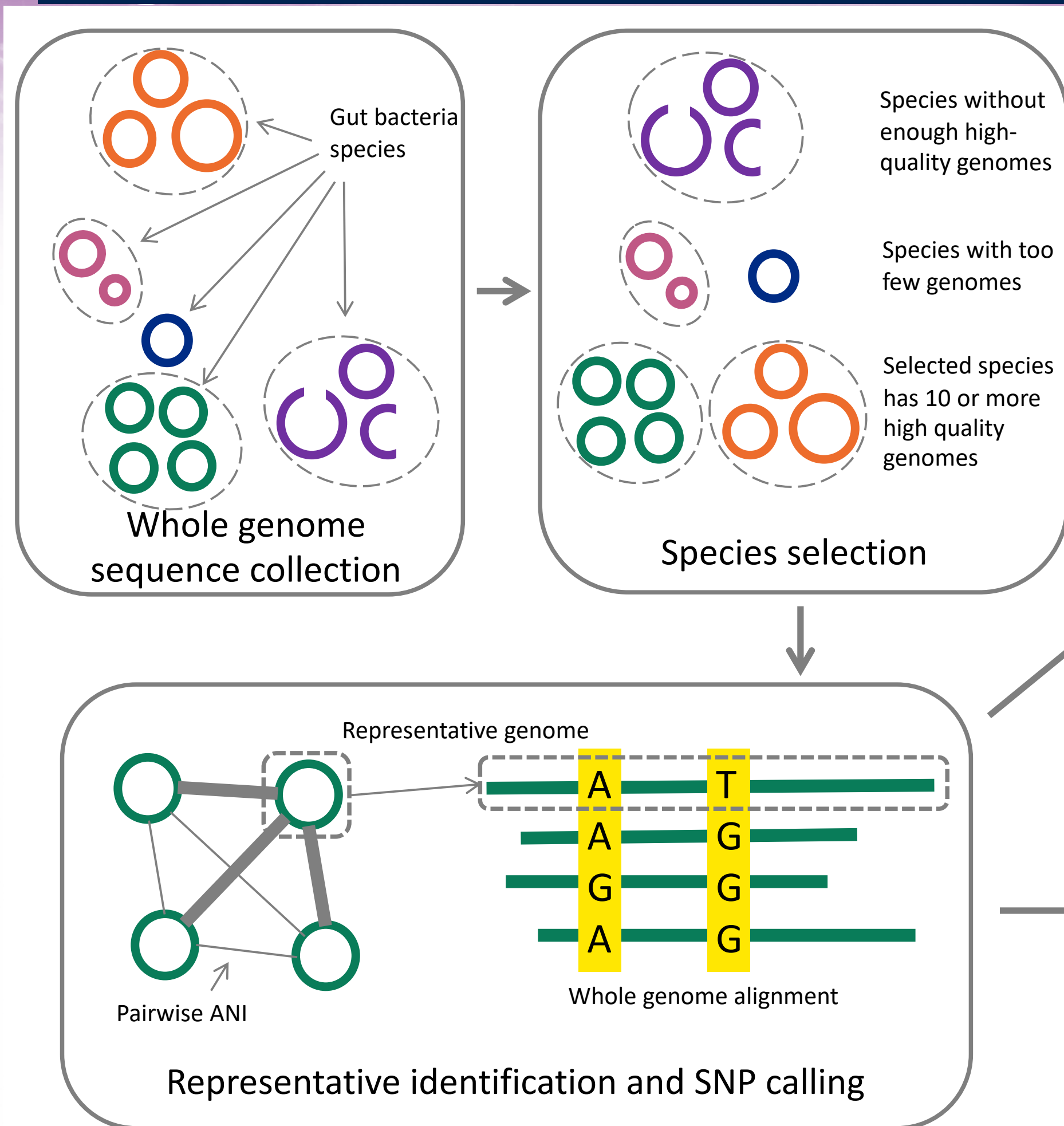
Whole genome sequence collection

Species without enough high-quality genomes

Species with too few genomes

Selected species has 10 or more high quality genomes

Species selection

Representative genome

Pairwise ANI

Whole genome alignment

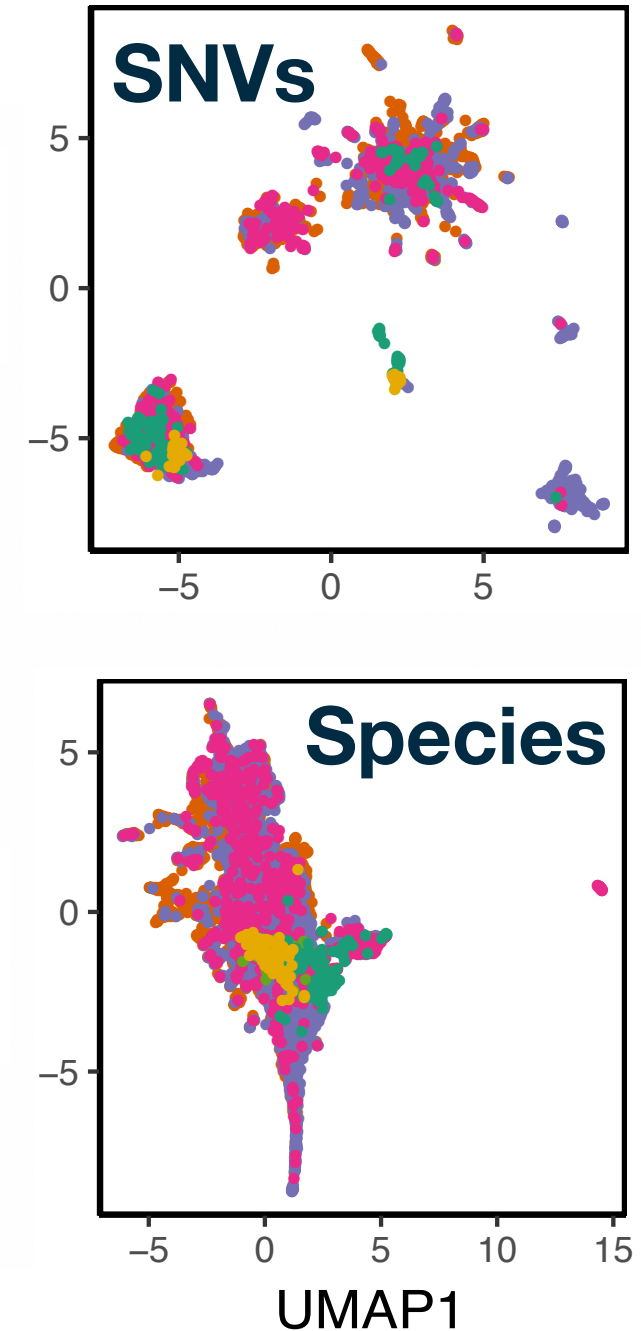Representative identification and SNP calling

**Maast**

Shi et al. (2021)

**Compression > bzip2, rapid exact matching**
**Prefix filter, Suffix array, Colex sort**

https://github.com/zjshi/gt-pro
https://github.com/zjshi/Maast

# GTPRO: 100x faster, more accurate



SERVER

LAPTOP

False discovery rate (FDR): GT-Pro, MIDAS, metaSNV

Sensitivity (%) vs Sequence coverage: x.001, x.01, x.1, x1, x2, x5, x10, x15

Allele sharing score: 0 0.02 0.04 0.06 0.08 0.1 0.12 0.14 0.16 0.18 0.2 0.22

SNVs

Species

Metagenome simulations with varying sequencing coverage

Global stool samples (N=7,459)

Shi et al. (2021)

# Unique k-mers beat alignment at known SNVs

# But current approach only works on SNVs discovered in genomes
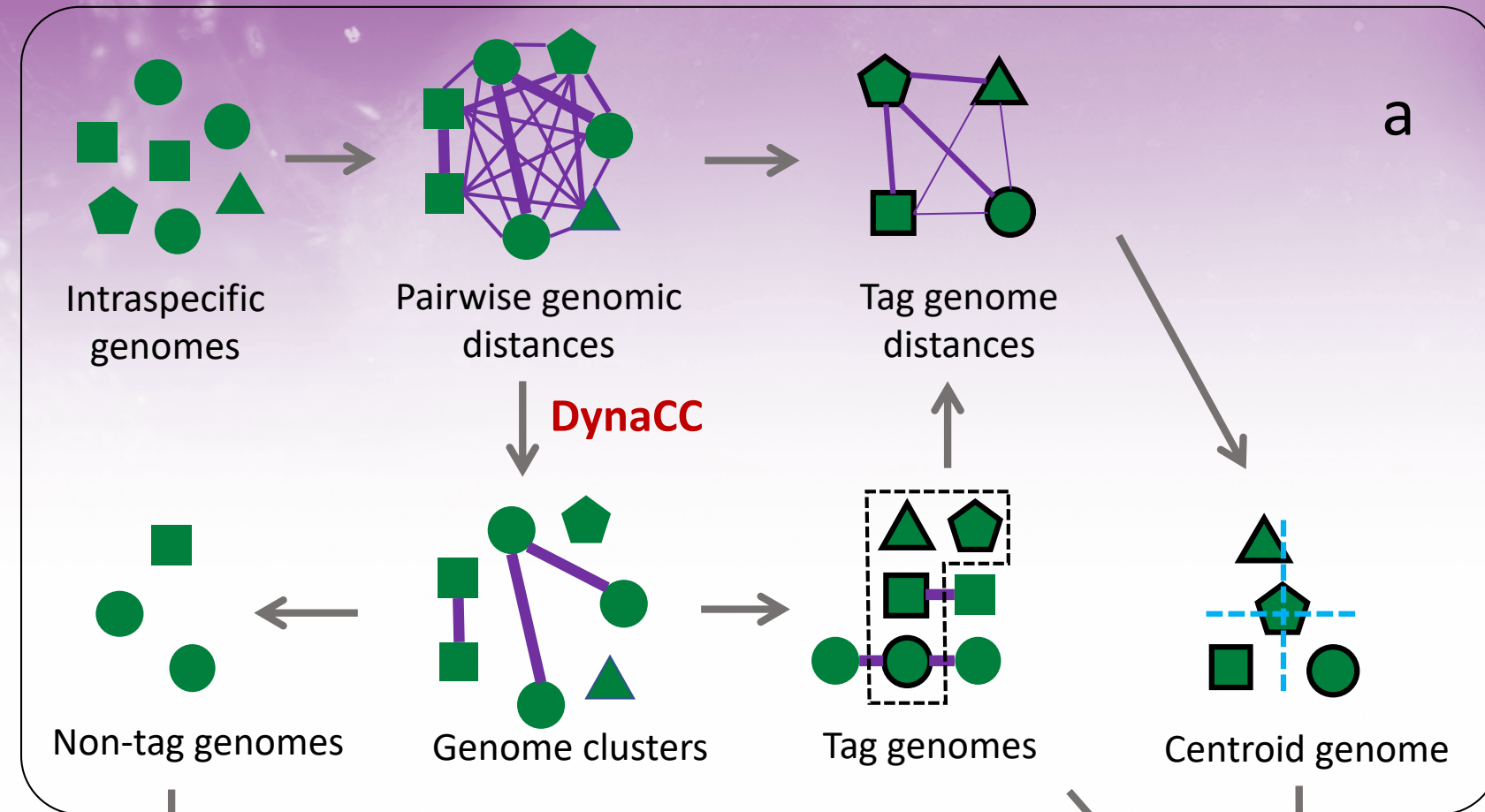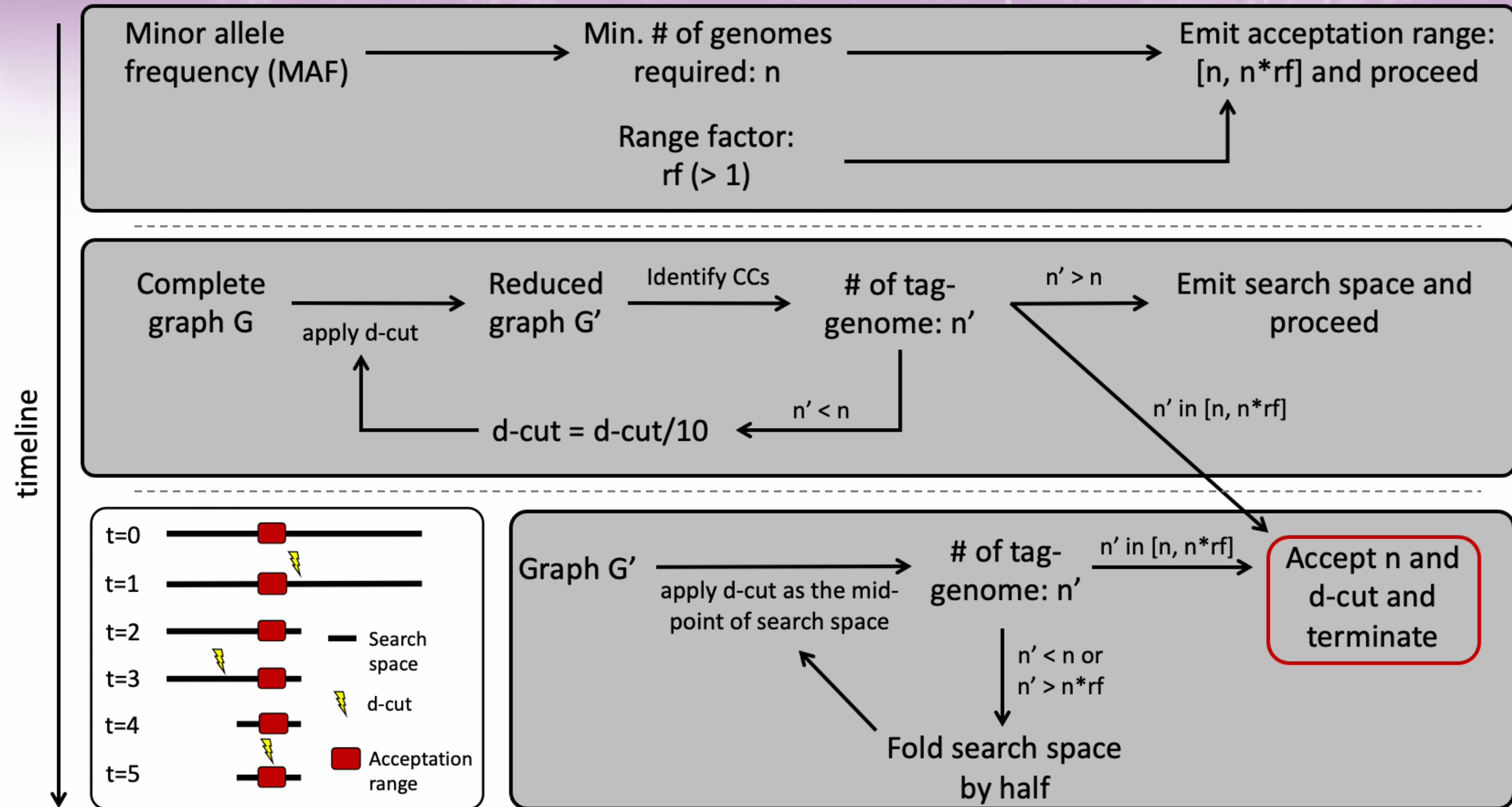
# Challenge 2:
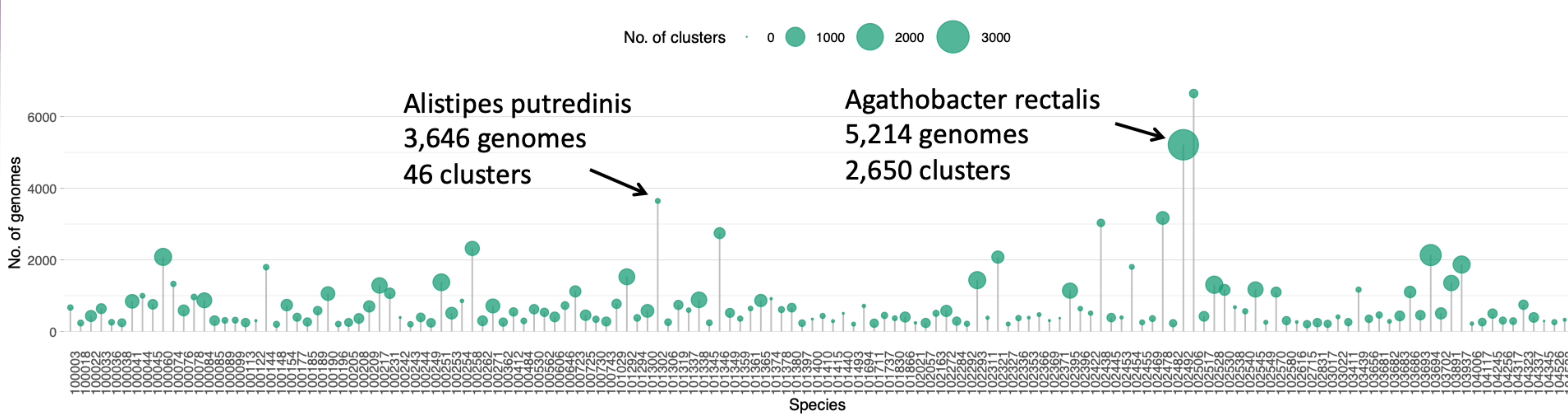# How to align and call variants in so many genomes?
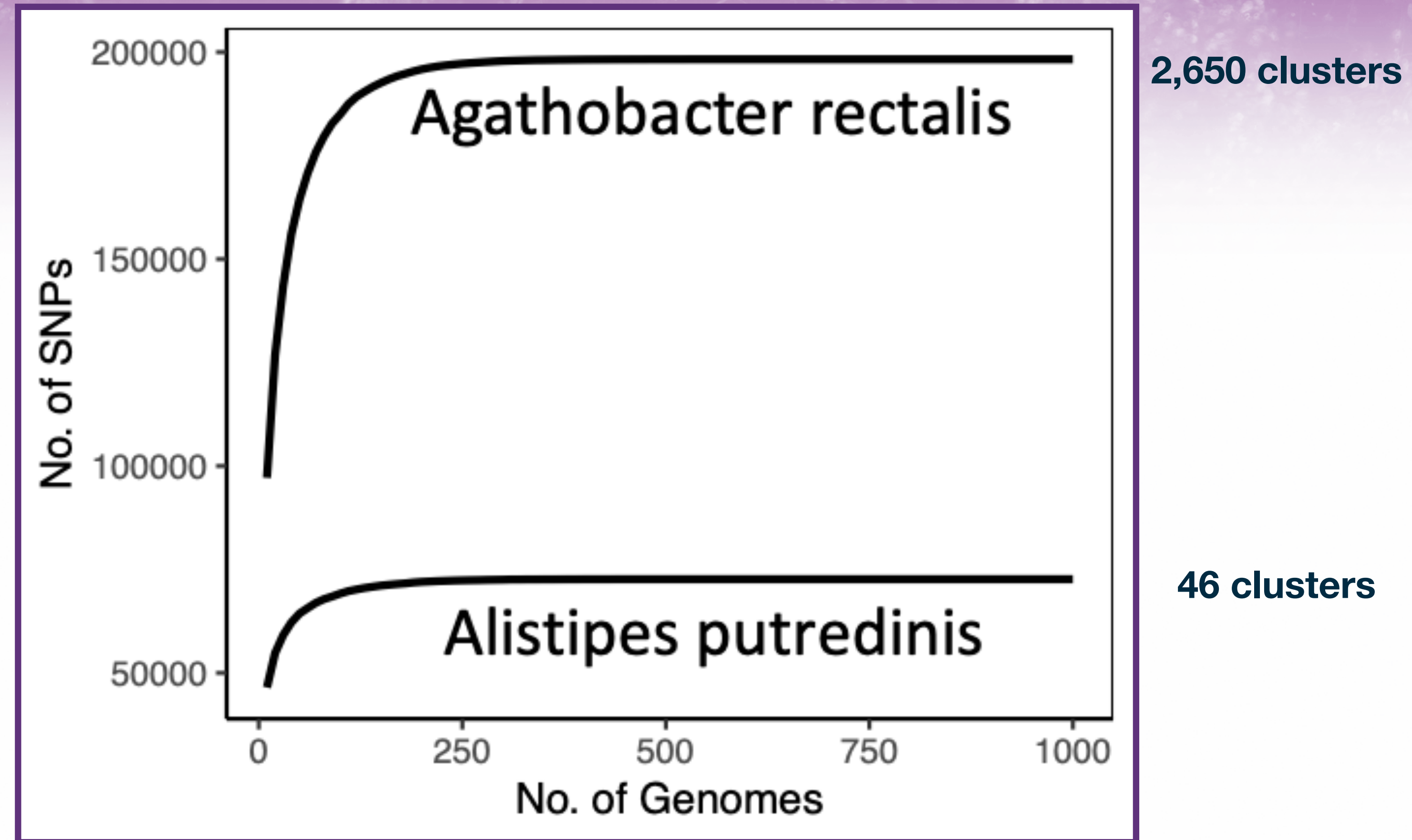
# Maast: fast variant discovery from genomes



a

- Intraspecific genomes
- Pairwise genomic distances
- **DynaCC**
- Tag genome distances
- Non-tag genomes
- Genome clusters
- Tag genomes
- Centroid genome

## DynaCC algorithm flowchart

Minor allele frequency (MAF) → Min. # of genomes required: n → Emit acceptation range: [n, n*rf] and proceed

Range factor: rf (> 1)

Complete graph G — apply d-cut → Reduced graph G' — Identify CCs → # of tag-genome: n' — n' > n → Emit search space and proceed

d-cut = d-cut/10 ← n' < n

n' in [n, n*rf]

t=0
t=1
t=2
t=3
t=4
t=5

— Search space

⚡ d-cut

🟥 Acceptation range

Graph G' → apply d-cut as the mid-point of search space → # of tag-genome: n' — n' in [n, n*rf] → Accept n and d-cut and terminate

n' < n or n' > n*rf

Fold search space by half

timeline

Shi et al. (2022)

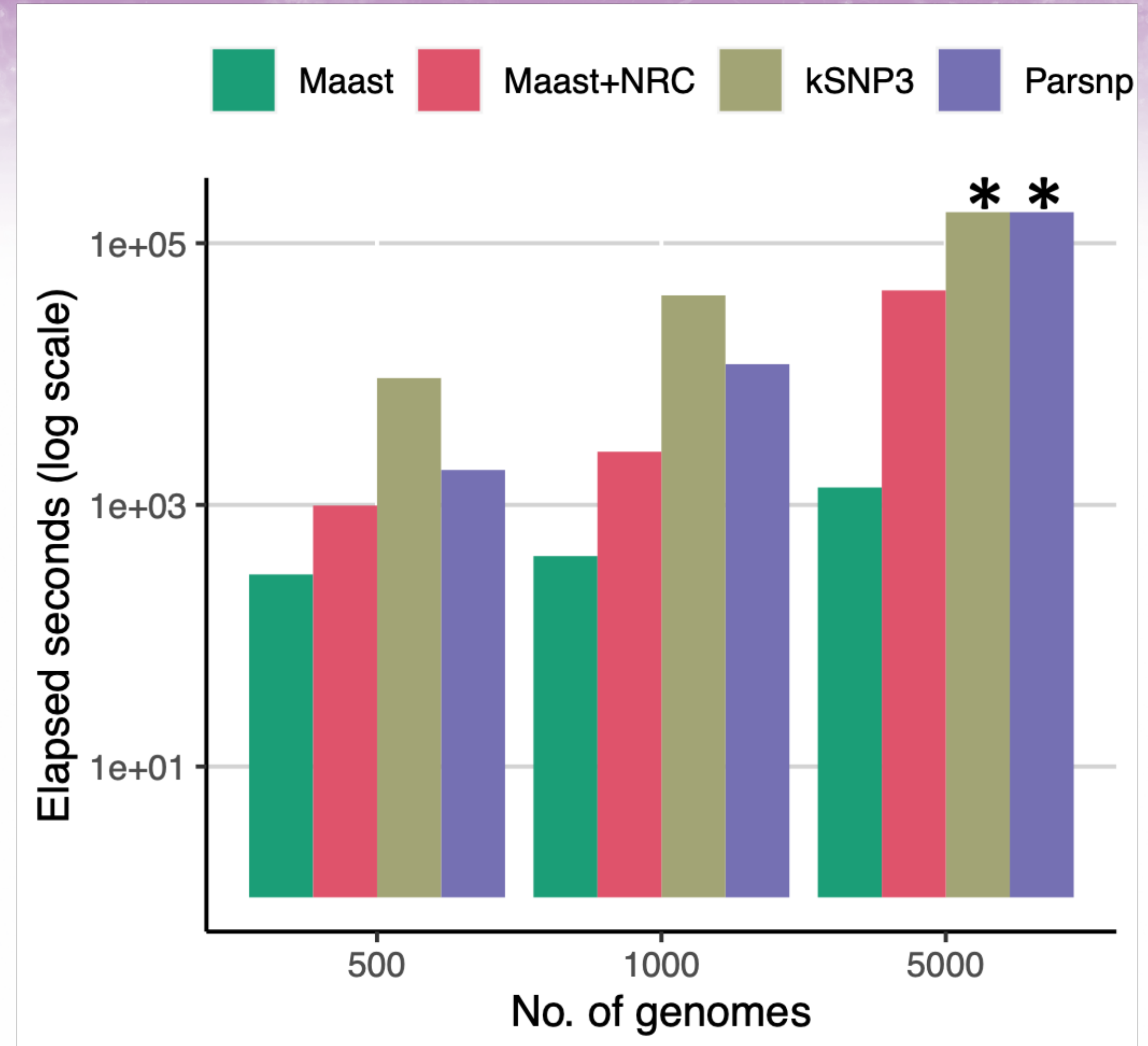# Genome redundancy offers solution
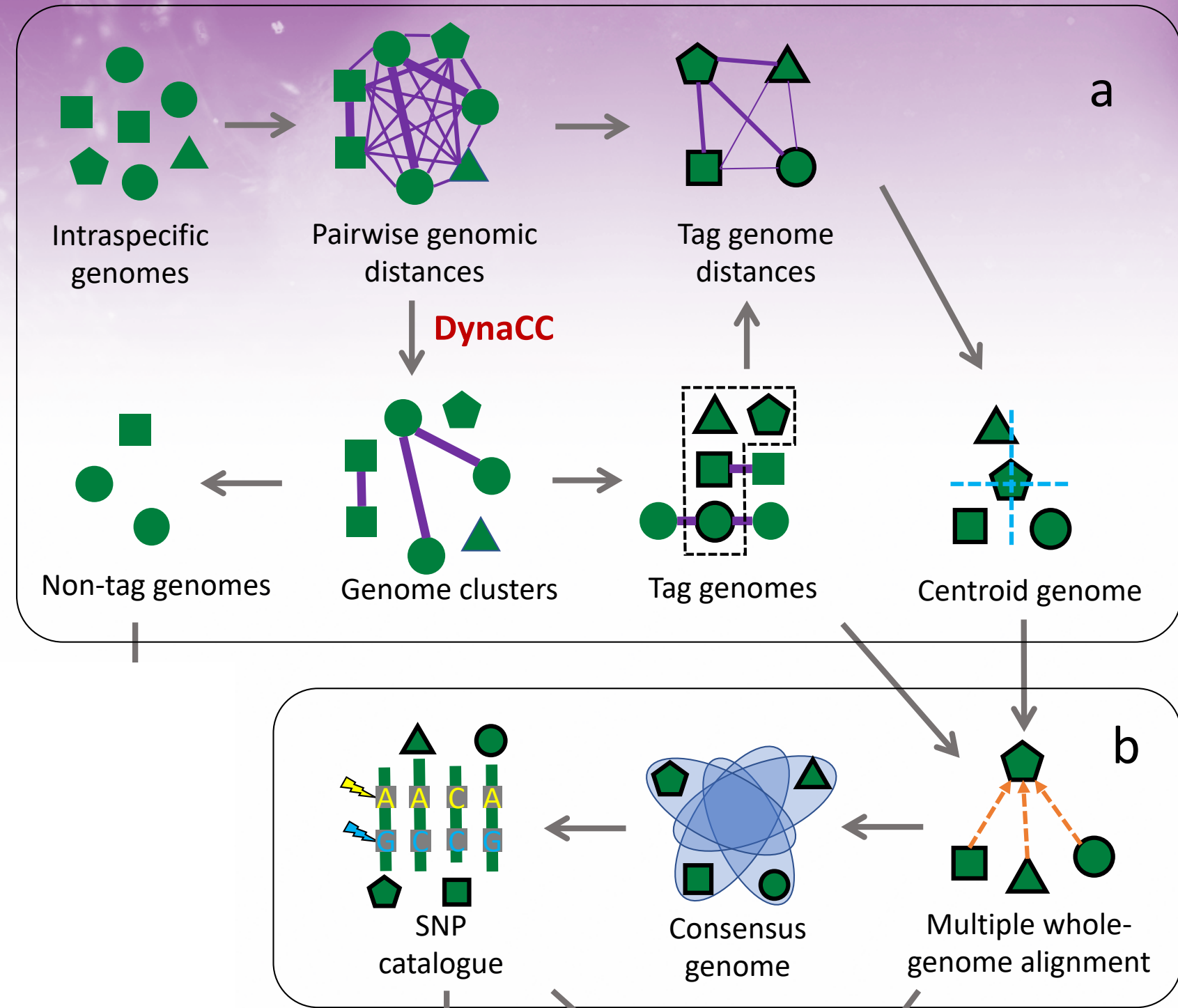


Shi et al. (2022)

# Maast: fast variant discovery from genomes

# Maast: fast variant discovery from genomes
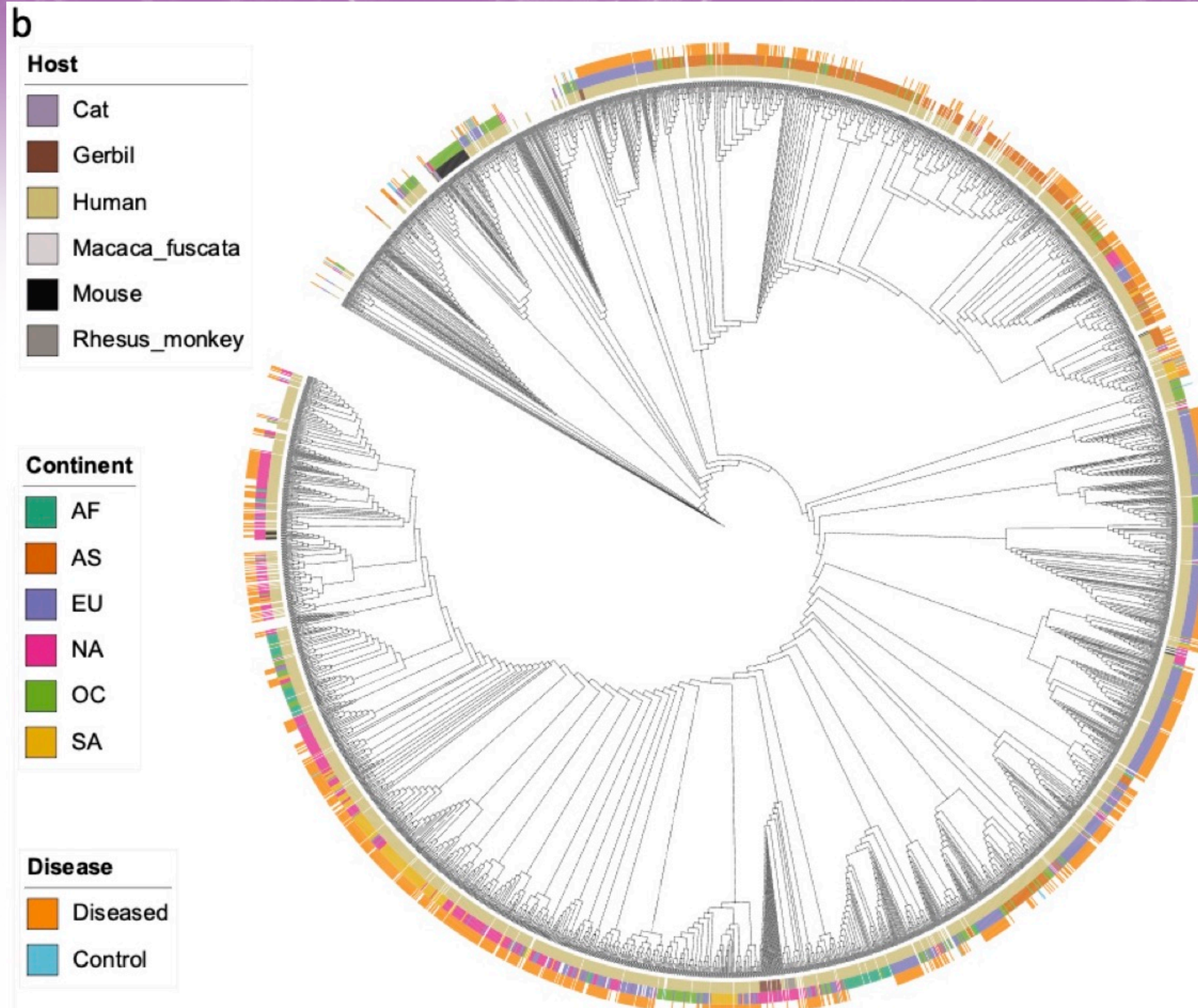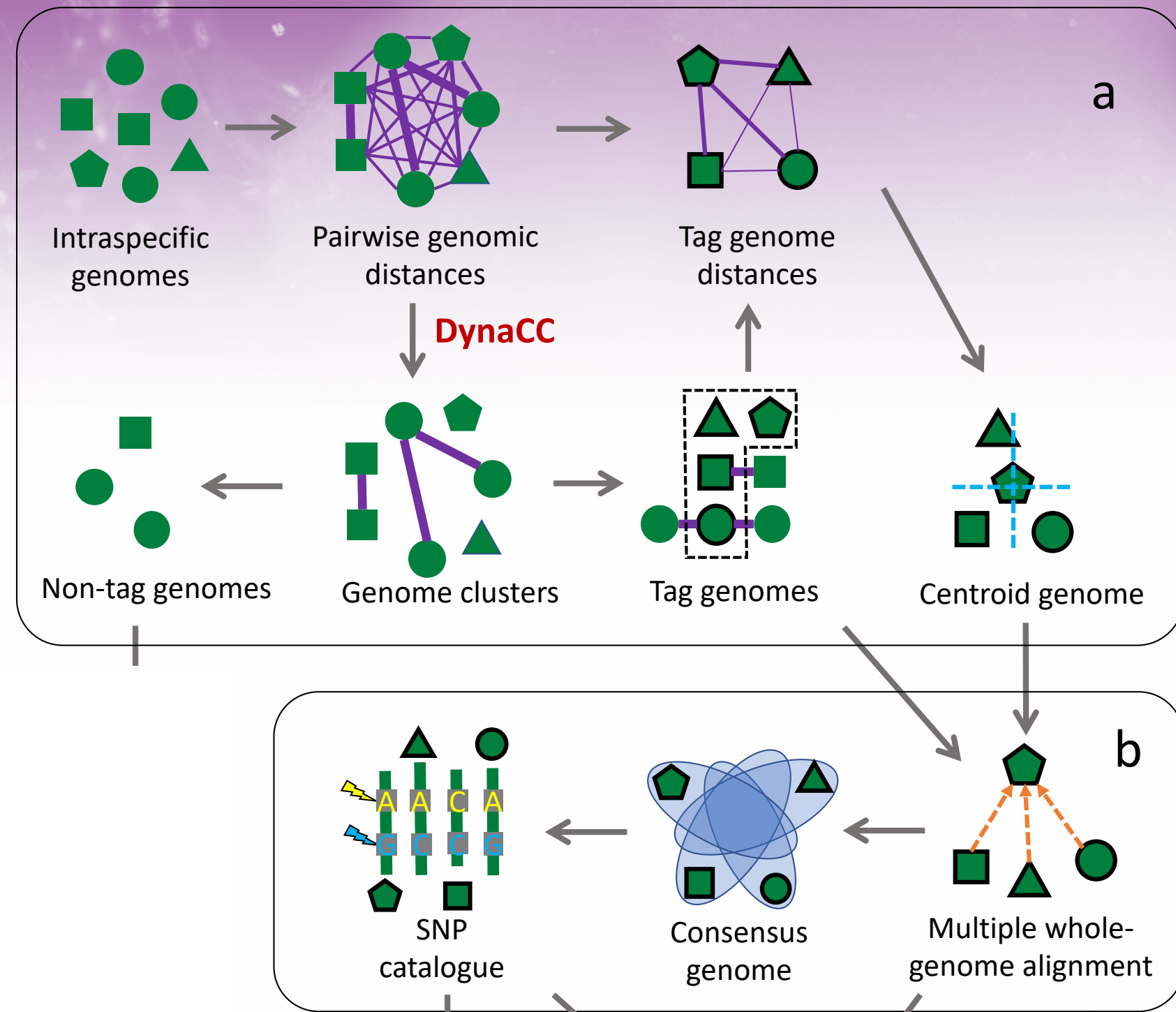


Shi et al. (2022)

# Maast: fast variant discovery from genomes

# Maast: fast variant discovery from genomes



3,068 *H. pylori* strains
Also: 37,096 SARS-CoV-2 strains

https://github.com/zjshi/Maast    Shi et al. (2022)

# Tag genomes speed up variant discovery and improve accuracy

# Sequencing effort should focus on new lineages not redundant ones

# Future Prospects

- Strategies beyond short-read aligners are needed, e.g.,
  - faster genome graph algorithms
  - probabilistic read mapping
  - read-to-read comparisons (reference databases for interpretation)
  - long reads / haplotypes
- Tools that use reference databases need to be flexibly implemented so that the algorithms and database can be tailored to the community

# Future Prospects

- Not just problems for bacterial communities.
  - CRS and redundant genomes in some lineages of archaea, eukaryotes, and viruses.
- These challenges affect all bioinformatics methods that compare reads to databases, not just metagenotyping.
- Democratizing large-scale bioinformatics is critical!

# Acknowledgements

Jason Shi

Chunyu Zhao

**pollard**lab