

Pan-Genomes Are The New Reference Genomes

35th TBI Winterseminar, Feb. 2020, Bled

What should a reference sequence be able to represent?

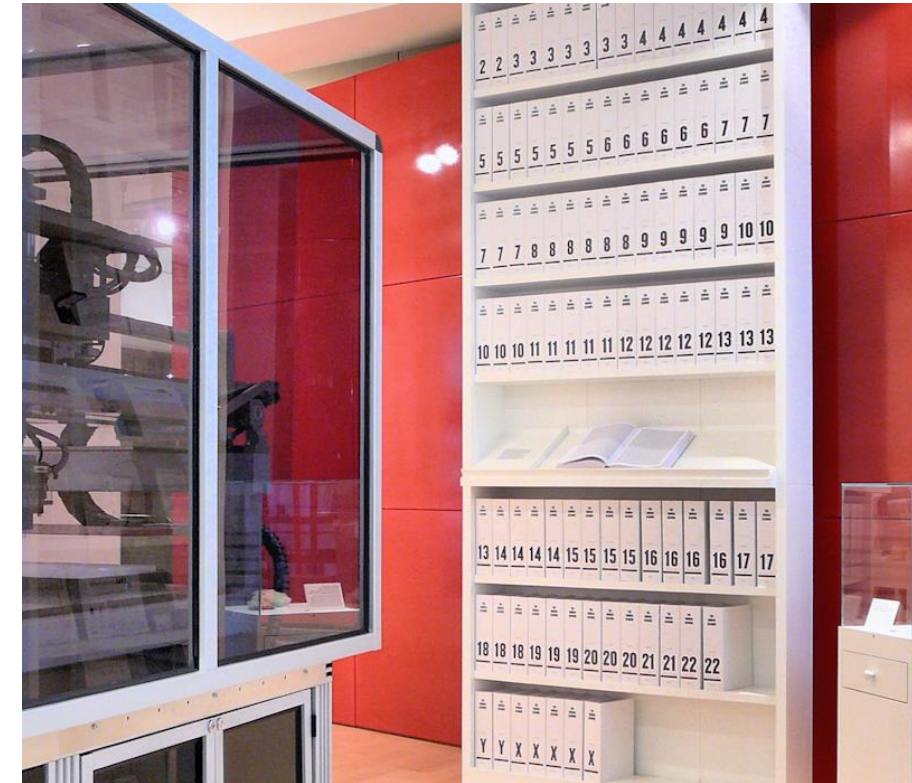
“Type strain/reference strain is usually the firstly isolated strain of the species and exhibits all of the relevant phenotypic and genotypic properties cited in the species circumscriptions.”

- Single genomes
- Functional genome
- Consensus from a population
- Maximal genome/Pan-genome

- → not one single genome as reference sequence, but rather the pan-genome

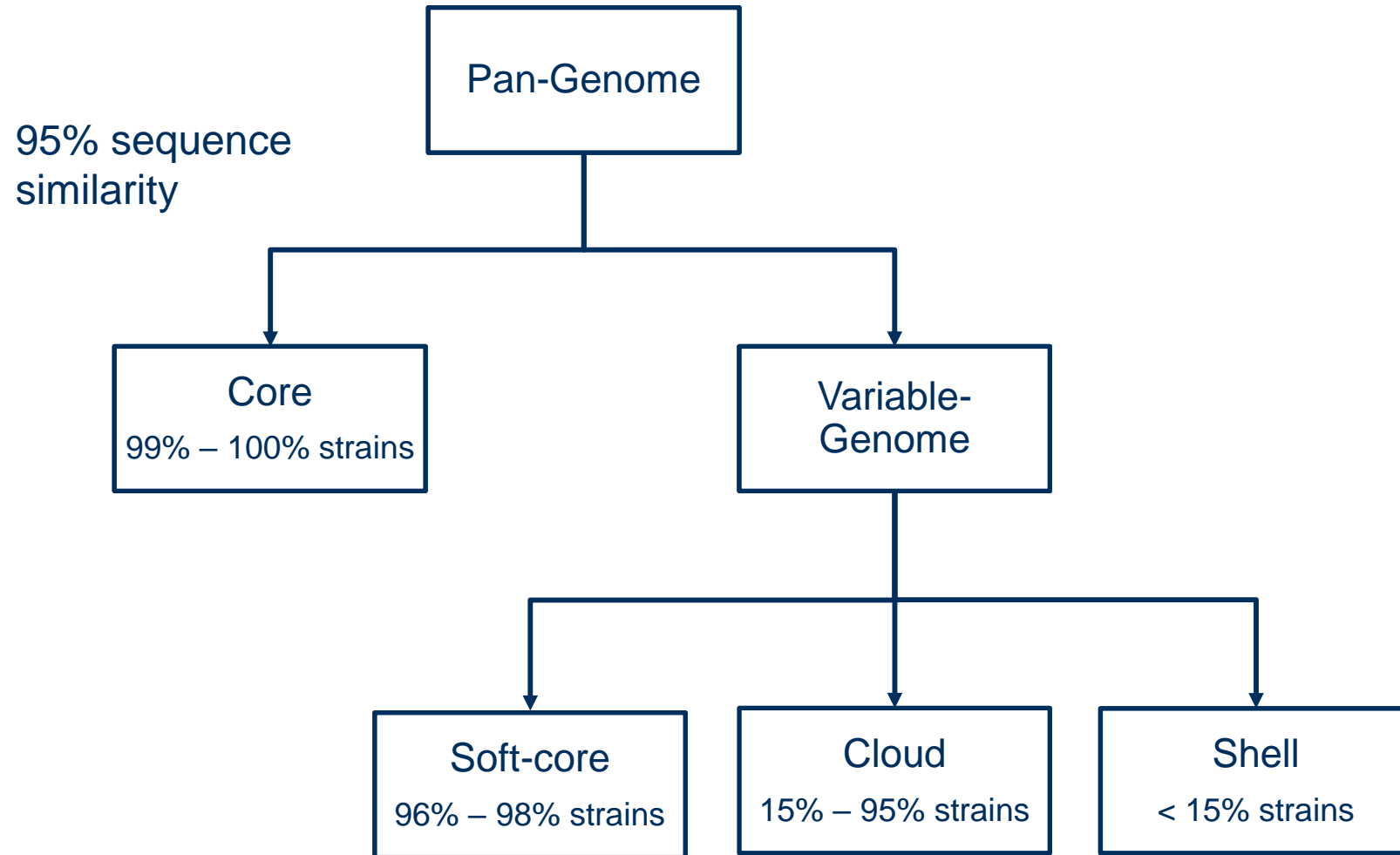
- Replace traditional linear reference genomes by richer data structures

The first version of the human reference genome, 2001



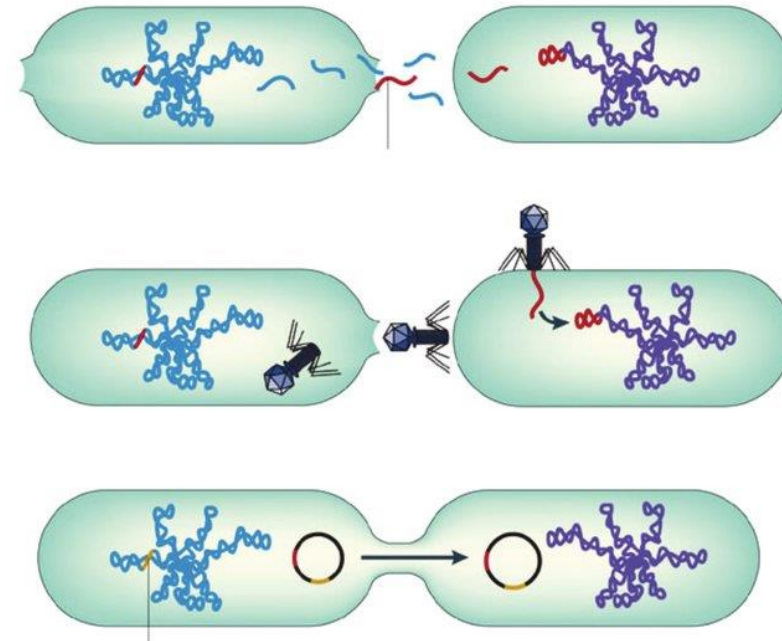
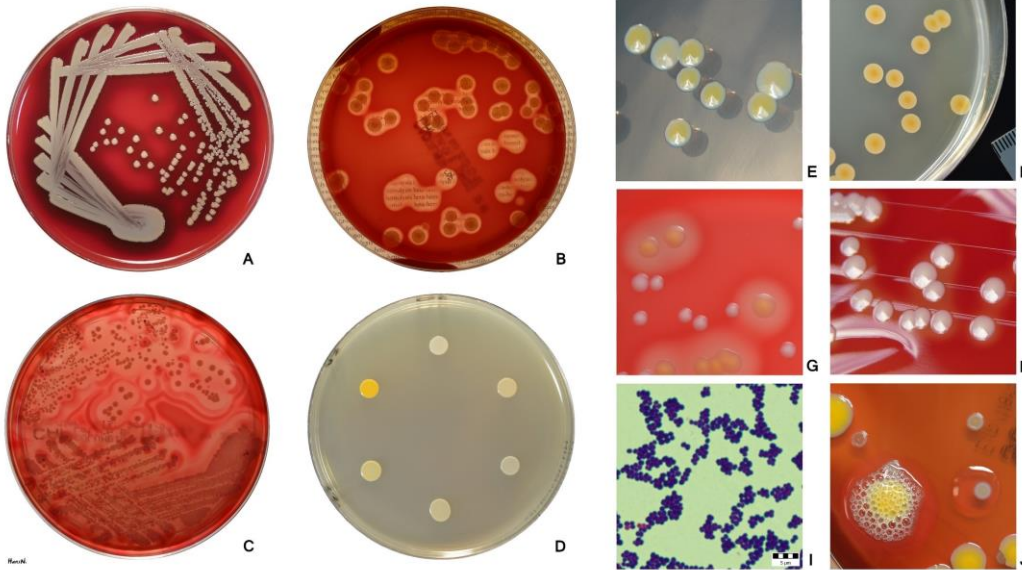
https://en.wikipedia.org/wiki/Reference_genome

The pan-genome is a representation of all genomic content in a certain species or phylogenetic clade



Staphylococcus aureus: a common bacterium in the upper respiratory tract and on the skin

- Human microbiota
- Horizontal gene transfer
→ pathogenic
- Antibiotic resistance



Adapted from Furuya and Lowy, 2006

The power of the pan-genome: phylogeny based on the top 10 “most relevant“ *Staphylococcus aureus* strains

1	>AP017922.1 <i>Staphylococcus aureus</i> DNA, strain: JP080
2	>CP013231.1 <i>Staphylococcus aureus</i> strain UTSW MRSA 55
3	>CP038021.1 <i>Staphylococcus aureus</i> strain 04-002
4	>CP038268.1 <i>Staphylococcus aureus</i> strain O55 isolate B118
5	>CP038819.1 <i>Staphylococcus aureus</i> strain O82
6	>CP039848.1 <i>Staphylococcus aureus</i> strain 2030RH1
7	>CP040623.1 <i>Staphylococcus aureus</i> strain D592-HR
8	>CP040801.1 <i>Staphylococcus aureus</i> strain S15
9	>LN626917.1 <i>Staphylococcus aureus</i> strain ILRI_Eymole1/1
10	>NC_002951.2 <i>Staphylococcus aureus</i> strain COL

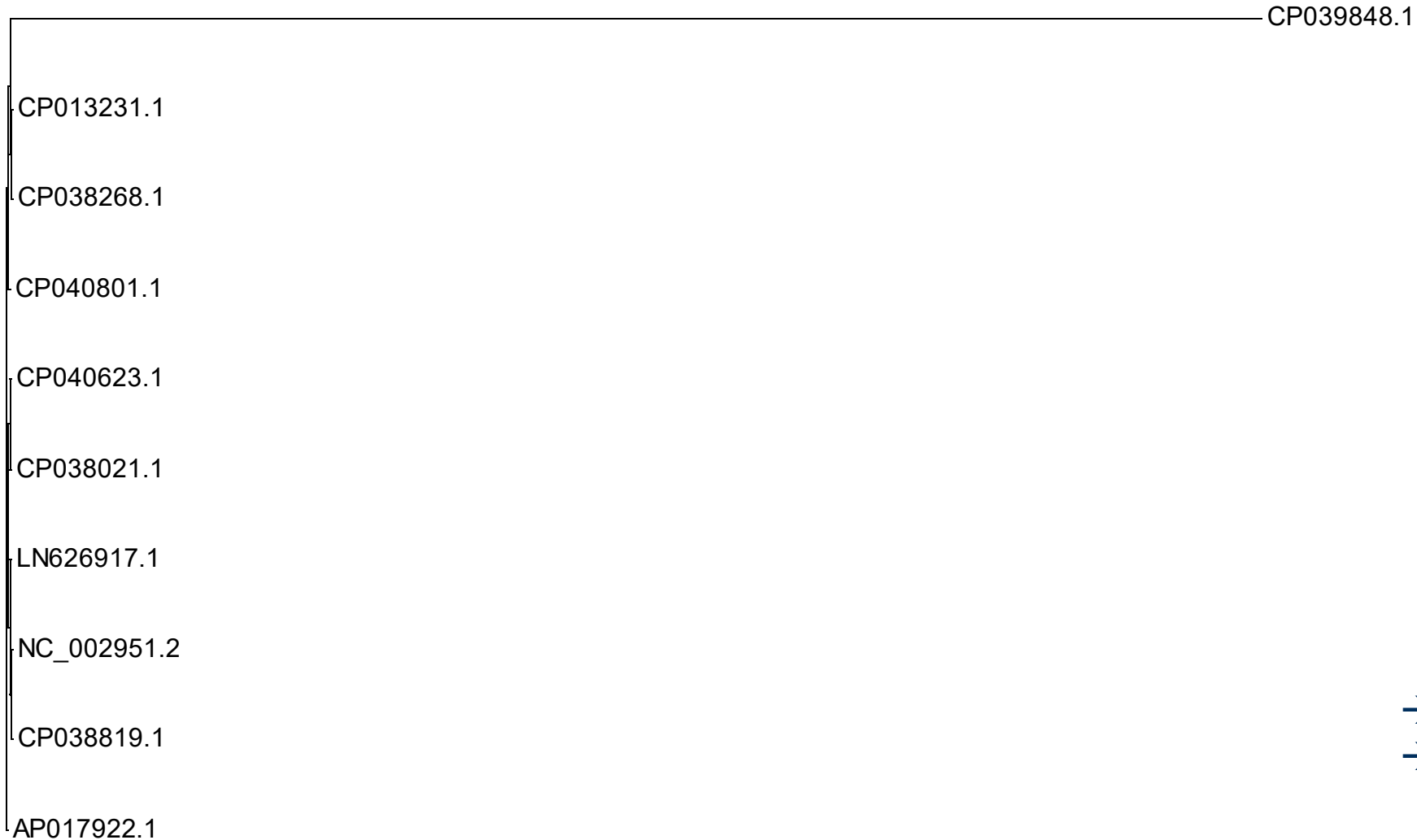
Genotyping approaches:

- 16S rRNA
→ one gene
- MLST (Multilocus sequence typing)
→ 7 housekeeping genes
- Core-genome
→ 1991 core genes

Staphylococcus aureus phylogeny

Based on: 16S rRNA (1 gene)

1.0E-4

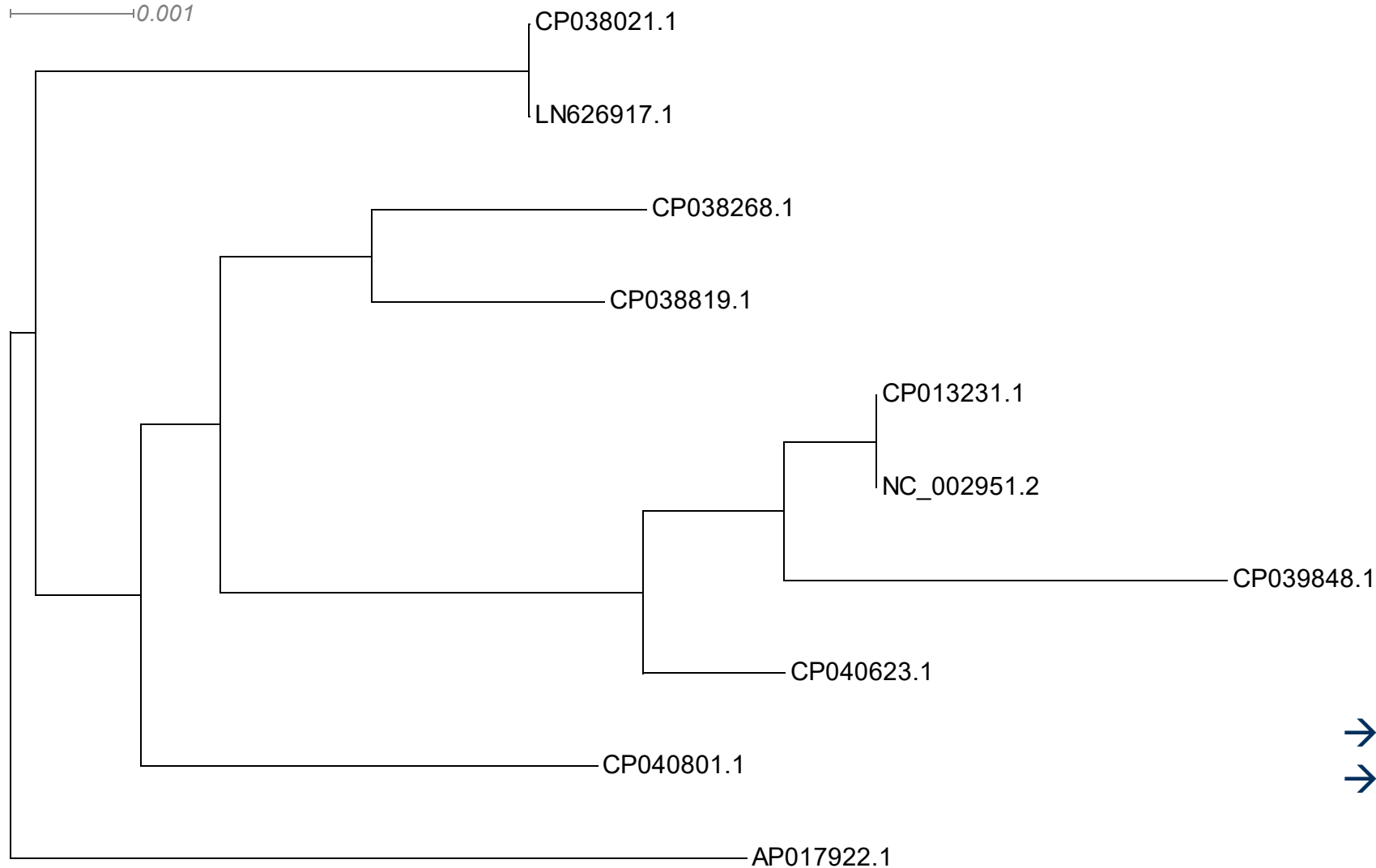


→ Closely related strains

→ One outlier

Staphylococcus aureus phylogeny

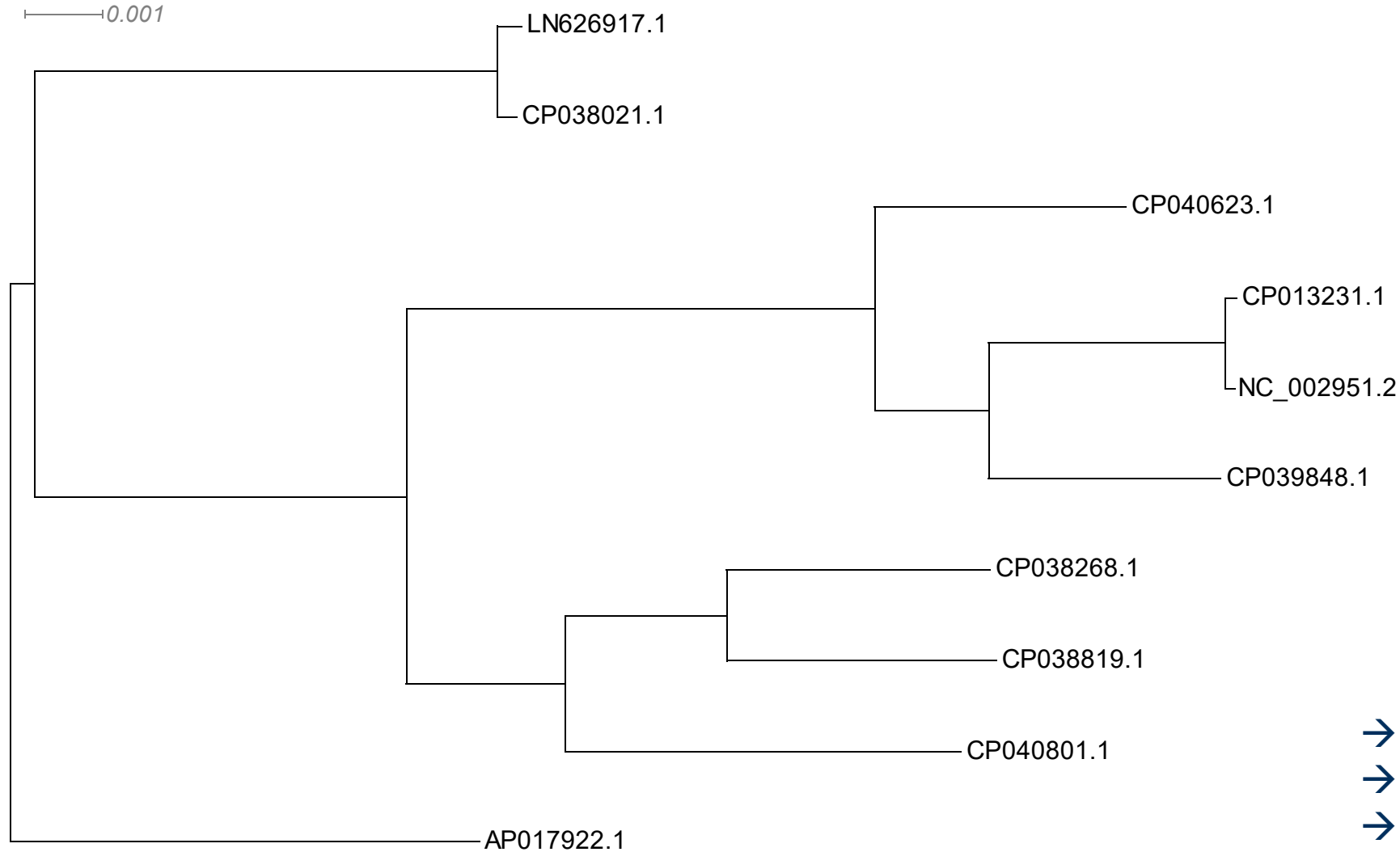
Based on: MLST (7 housekeeping genes)



- Better resolution
- Still identical strains

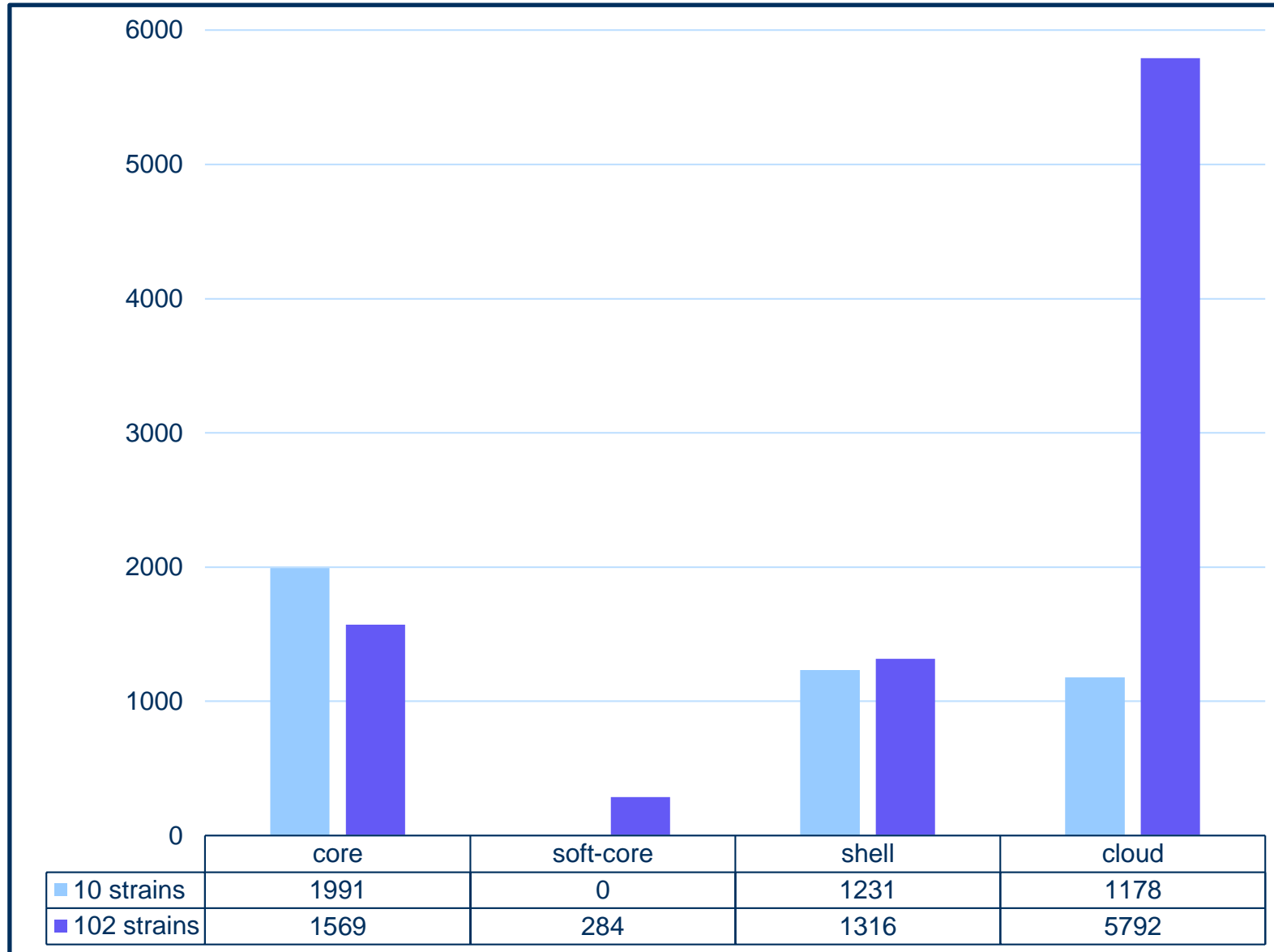
Staphylococcus aureus phylogeny

Based on: Core-genome (1991 genes)



- No identical
- Changed topology
- Robust?

The *Staphylococcus aureus* core-genome is robust: 10 vs 102 strains

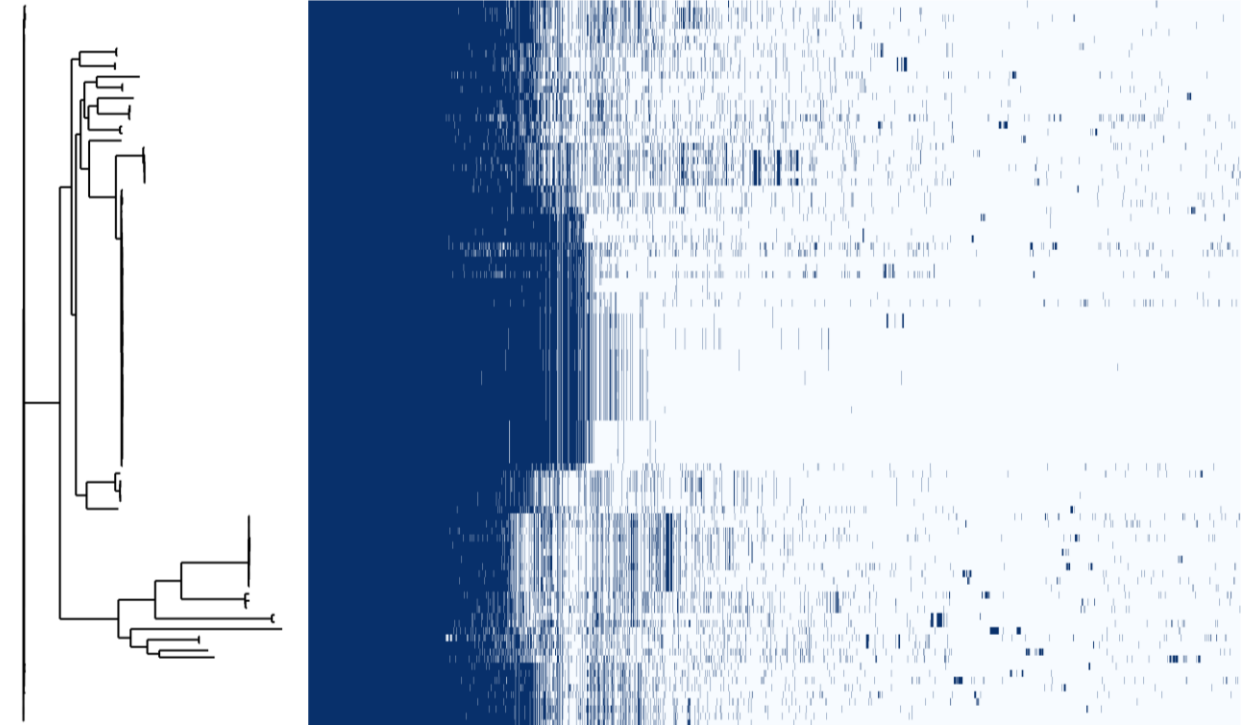
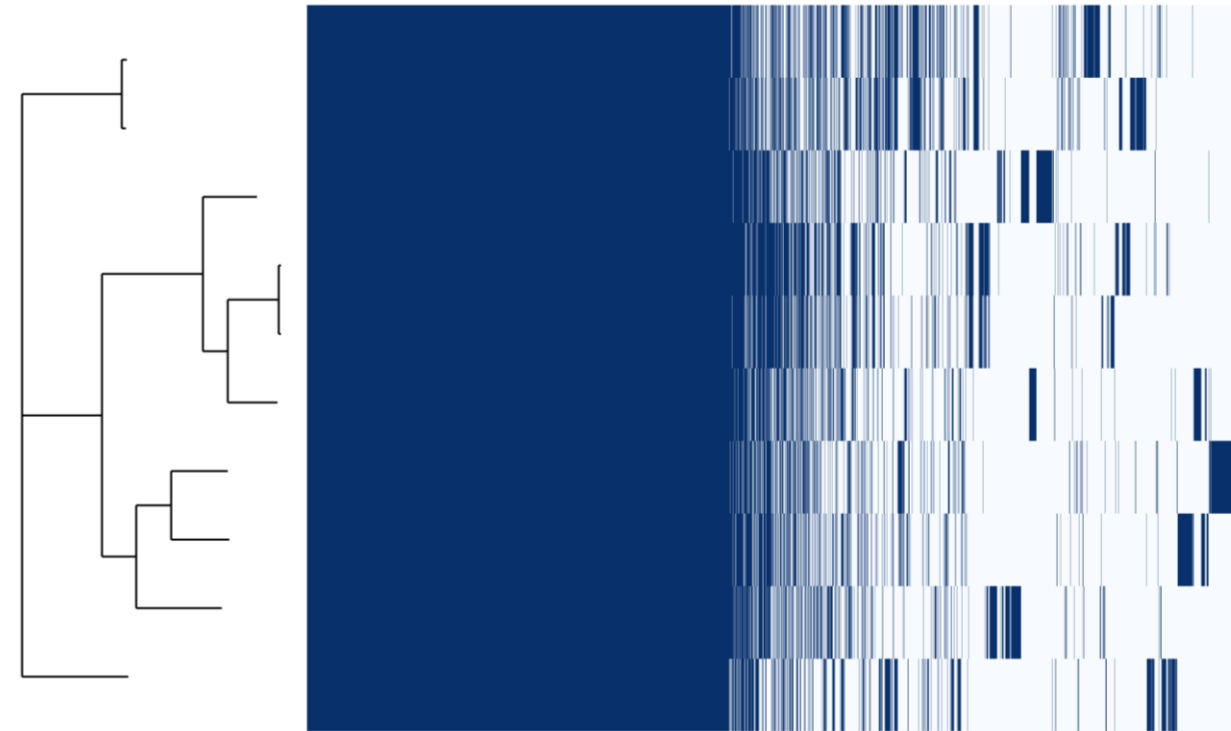


→ Robust
Core-genome

The *Staphylococcus aureus* core-genome is robust: 10 vs 102 strains

Pan-genome matrix
10 strains

Pan-genome matrix
102 strains

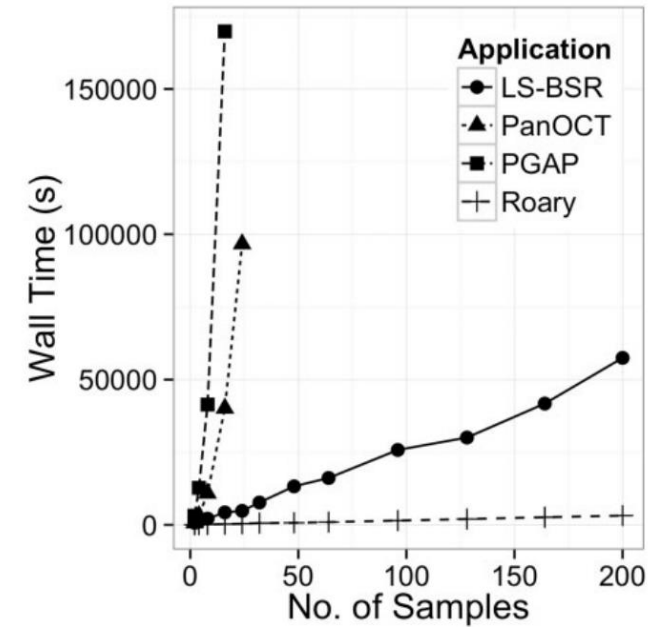


Next step: The *Staphylococcus aureus* pan-genome



- 40,000 strains
- A single reference strain is not representative for a whole species
- Use pan-genome to characterize species

- Runtime
- Storage
- Computational solution



Outlook - The *Staphylococcus aureus* pan-genome

Computational challenges

Data structures

Design goals:

- Construction and maintenance
- Coordinate system
- Biological features and computational layers
- Data retrieval
- Searching
- Comparing



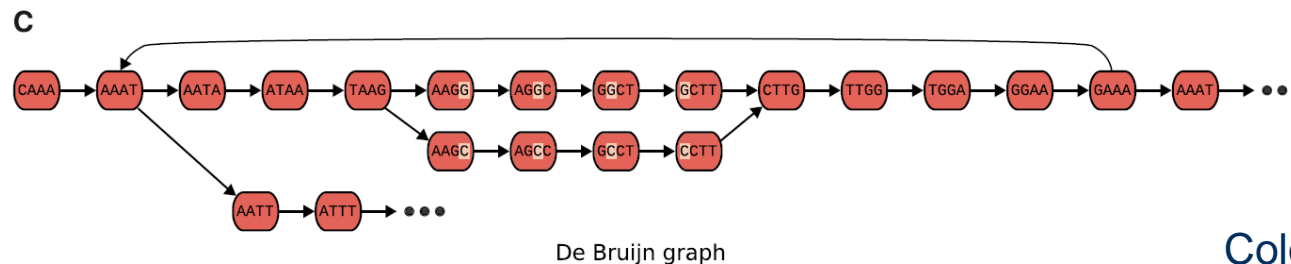
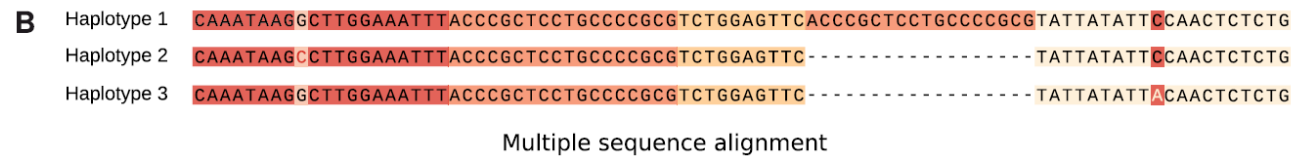
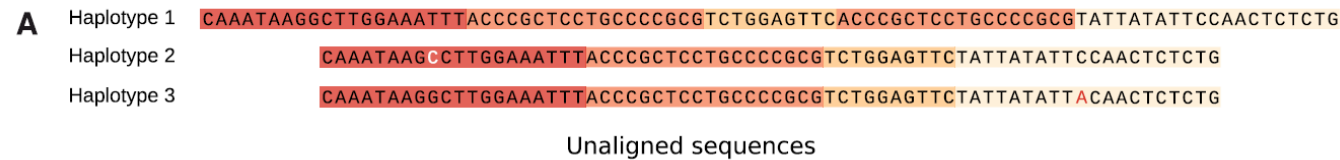
Outlook - The *Staphylococcus aureus* pan-genome

Computational challenges

Computational pan-genomics: status, promises and challenges

The Computational Pan-Genomics Consortium*

Data structure
 Approaches



Outlook - The *Staphylococcus aureus* pan-genome

Computational challenges

Variant calling and genotyping approaches

- Difference between newly sequenced genome and a reference

Visualization

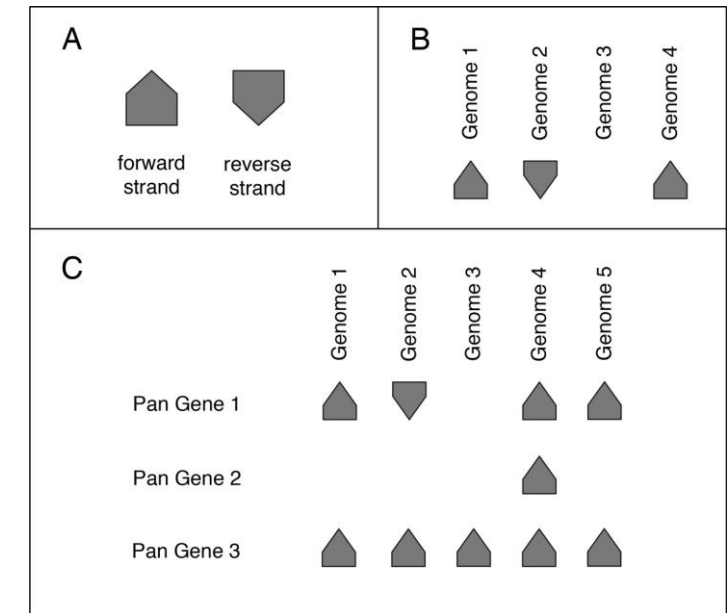
- View large genome sets
- Homology relationships

OXFORD

Briefings in Bioinformatics, 19(1), 2018, 118–135
doi: 10.1093/bib/bbw089
Advance Access Publication Date: 21 October 2016
Paper

Computational pan-genomics: status, promises and challenges

The Computational Pan-Genomics Consortium*



Pan-Tetris, Henning et al., 2015

Conclusion

- A single reference genome is not feasible to represent a whole species, but rather the pan-genome.
- Move away from linear reference genomes towards reference systems (graph based).
- Solve computational challenges in terms of storage and visualization.

Many thanks to:

RNA
BIOINFORMATICS & HIGH-THROUGHPUT ANALYSIS

Leibniz | ipht 



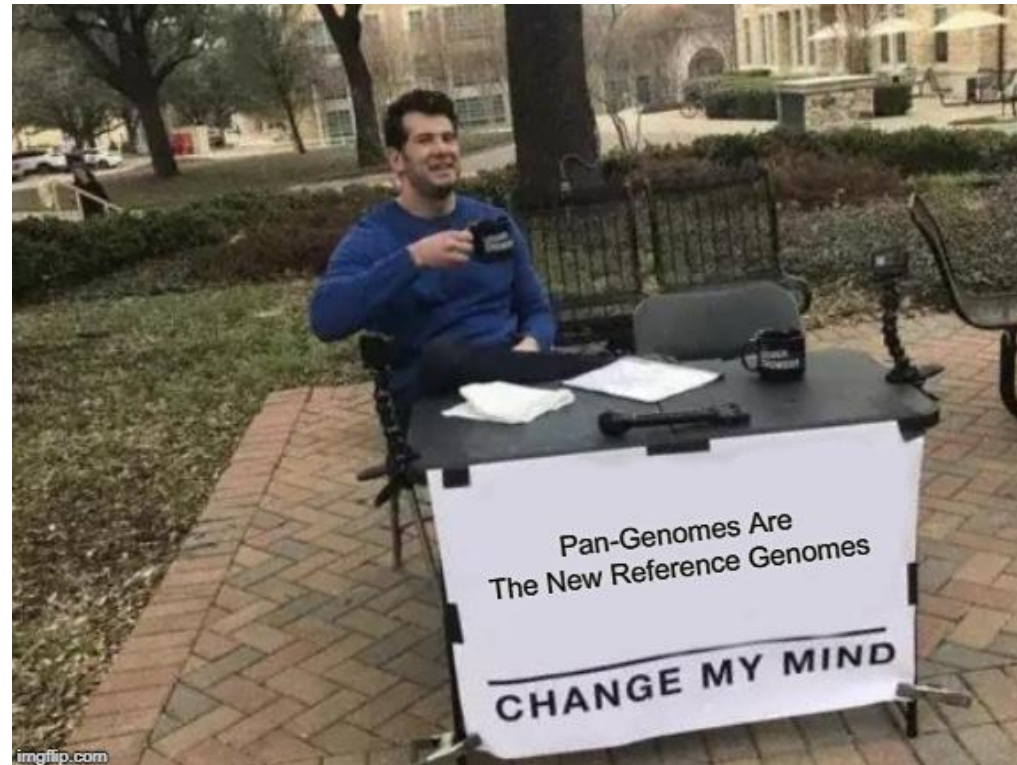
Max-Planck-Institut
für Menschheitsgeschichte

Max Planck Institute
for the Science of Human History

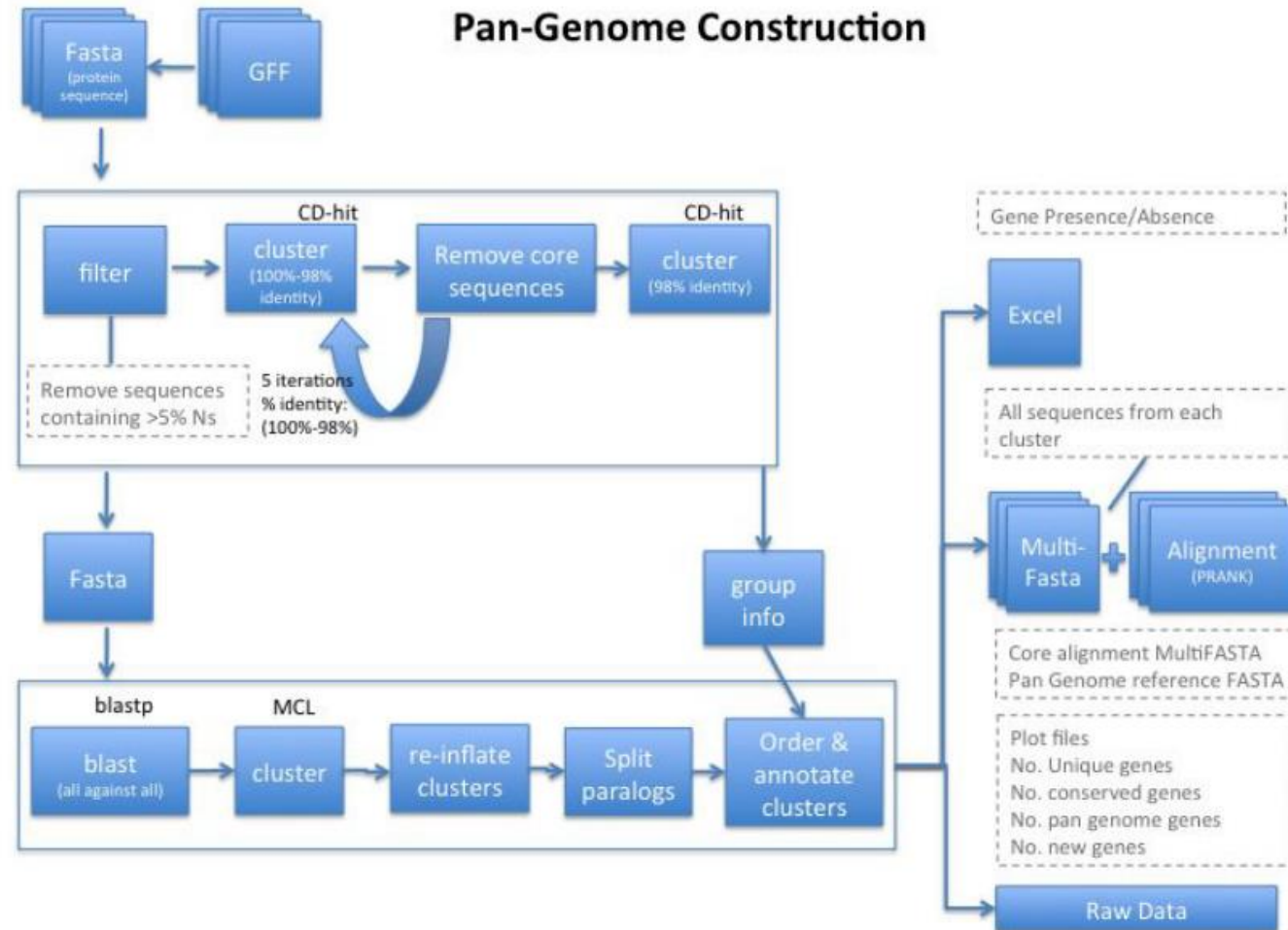
FRIEDRICH-SCHILLER-
UNIVERSITÄT
JENA



Thank you for your attention!



2.2 Method description



Sup. Fig. 13: A flowchart of the steps in the application.

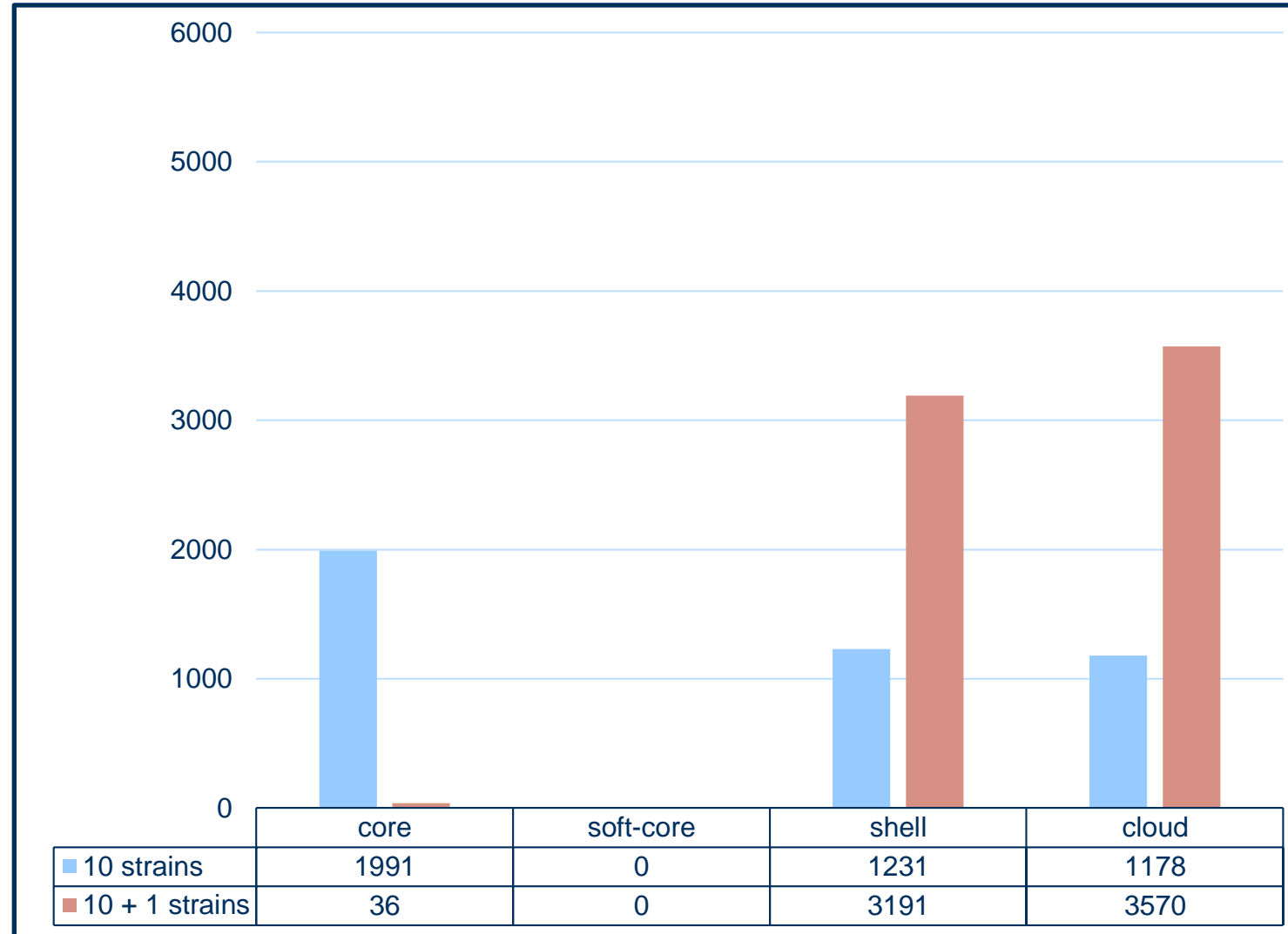
The problem of interspecies pan-genomes

Staphylococcus aureus and *Staphylococcus saprophyticus*

1	>AP017922.1 <i>Staphylococcus aureus</i> DNA, strain: JP080
2	>CP013231.1 <i>Staphylococcus aureus</i> strain UTSW MRSA 55
3	>CP038021.1 <i>Staphylococcus aureus</i> strain 04-002
4	>CP038268.1 <i>Staphylococcus aureus</i> strain O55 isolate B118
5	>CP038819.1 <i>Staphylococcus aureus</i> strain O82
6	>CP039848.1 <i>Staphylococcus aureus</i> strain 2030RH1
7	>CP040623.1 <i>Staphylococcus aureus</i> strain D592-HR
8	>CP040801.1 <i>Staphylococcus aureus</i> strain S15
9	>LN626917.1 <i>Staphylococcus aureus</i> strain ILRI_Eymole1/1
10	>NC_002951.2 <i>Staphylococcus aureus</i> strain COL
11	>AP008934.1 <i>Staphylococcus saprophyticus</i> strain ATCC 15305 DNA

The problem of interspecies pan-genomes

Staphylococcus aureus and *Staphylococcus saprophyticus*



Core: 99% – 100%; soft-core: 96% – 98%; shell: 15% – 95%; cloud: < 15%

The problem of interspecies pan-genomes

Staphylococcus aureus and *Staphylococcus saprophyticus*

