

Genome contamination has minimal impact on the delineation of species

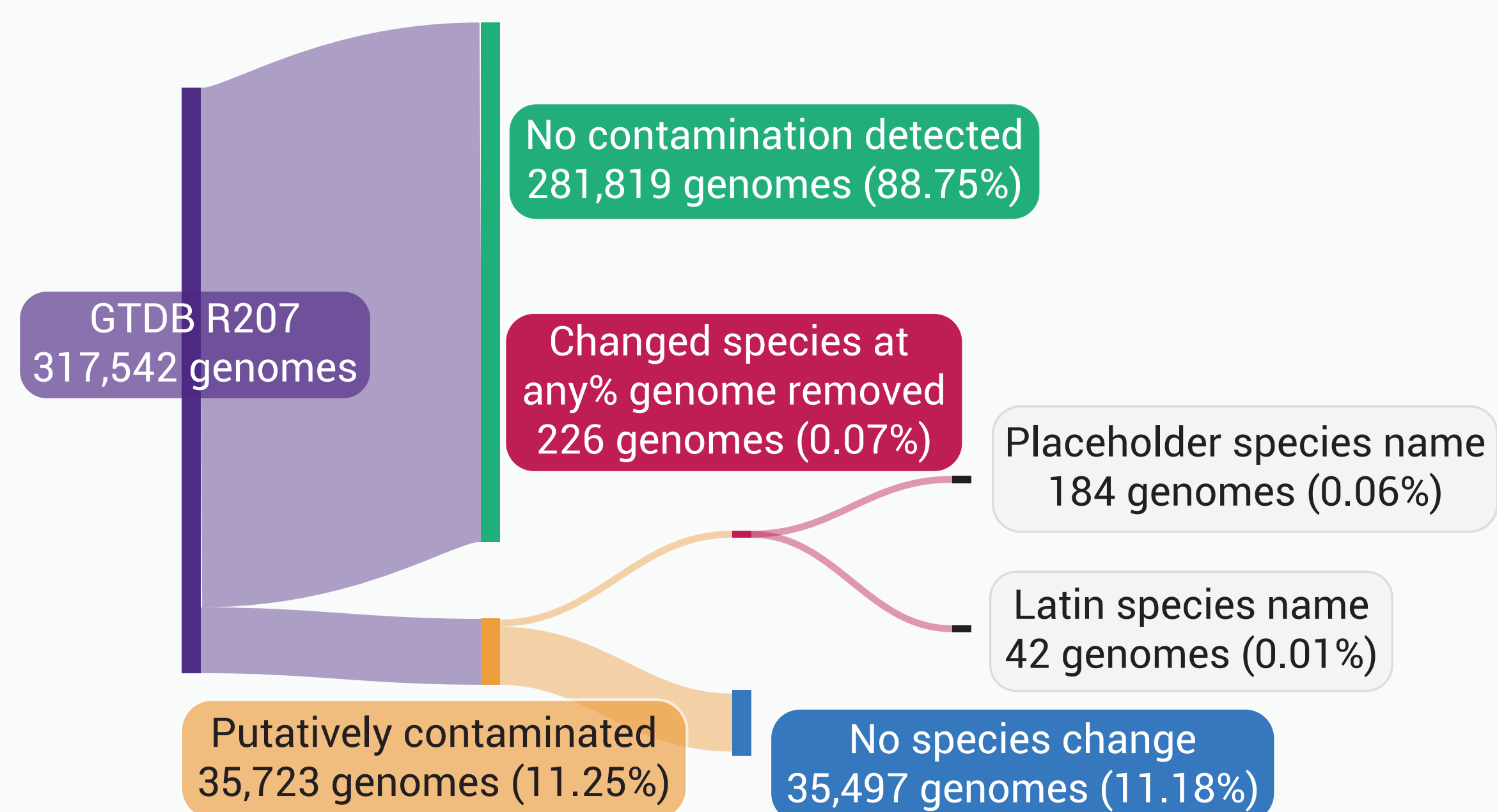
Aaron J. Mussig¹, Pierre-Alain Chaumeil¹, Maria Chuvochina¹, Christian Rinke¹, Donovan H. Parks¹, Philip Hugenholtz¹

¹ The University of Queensland, School of Chemistry and Molecular Biosciences, Australian Centre for Ecogenomics, St Lucia, QLD, Australia

Contamination is a known issue, but to what extent does this impact taxonomy?

We assess the impact of putative contamination on species delineated using average nucleotide identity (ANI).

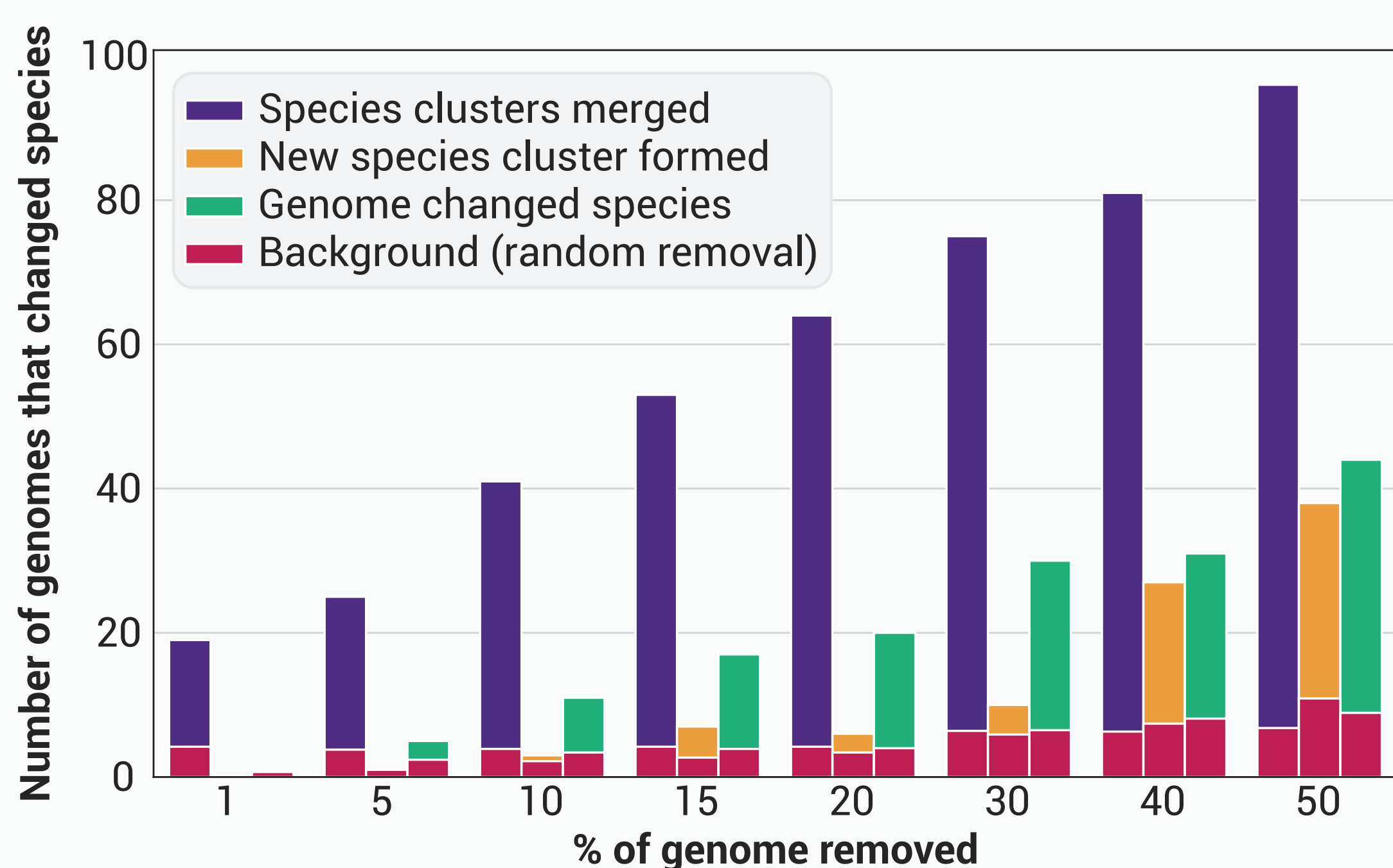
Species clustering is robust to contamination



• 35,723 of 317,542 (11.25%) GTDB R207 genomes were putatively contaminated.

• Identifying contamination using a reference database has limitations:
i) Taxonomically novel genomes are insufficiently represented.
ii) Order that genomes appear in the GUNC reference database impacts results.

Up to 226 genomes (0.07%) change species assignment when removing putatively contaminated contigs



• Even with removing half of the genome, there were only 226 genome that changed species.

• 184 of the genome changes were in nomenclaturely unimportant species (i.e. placeholder species).

• The 226 genomes belonged to 206 unique species clusters from 134 unique genera.

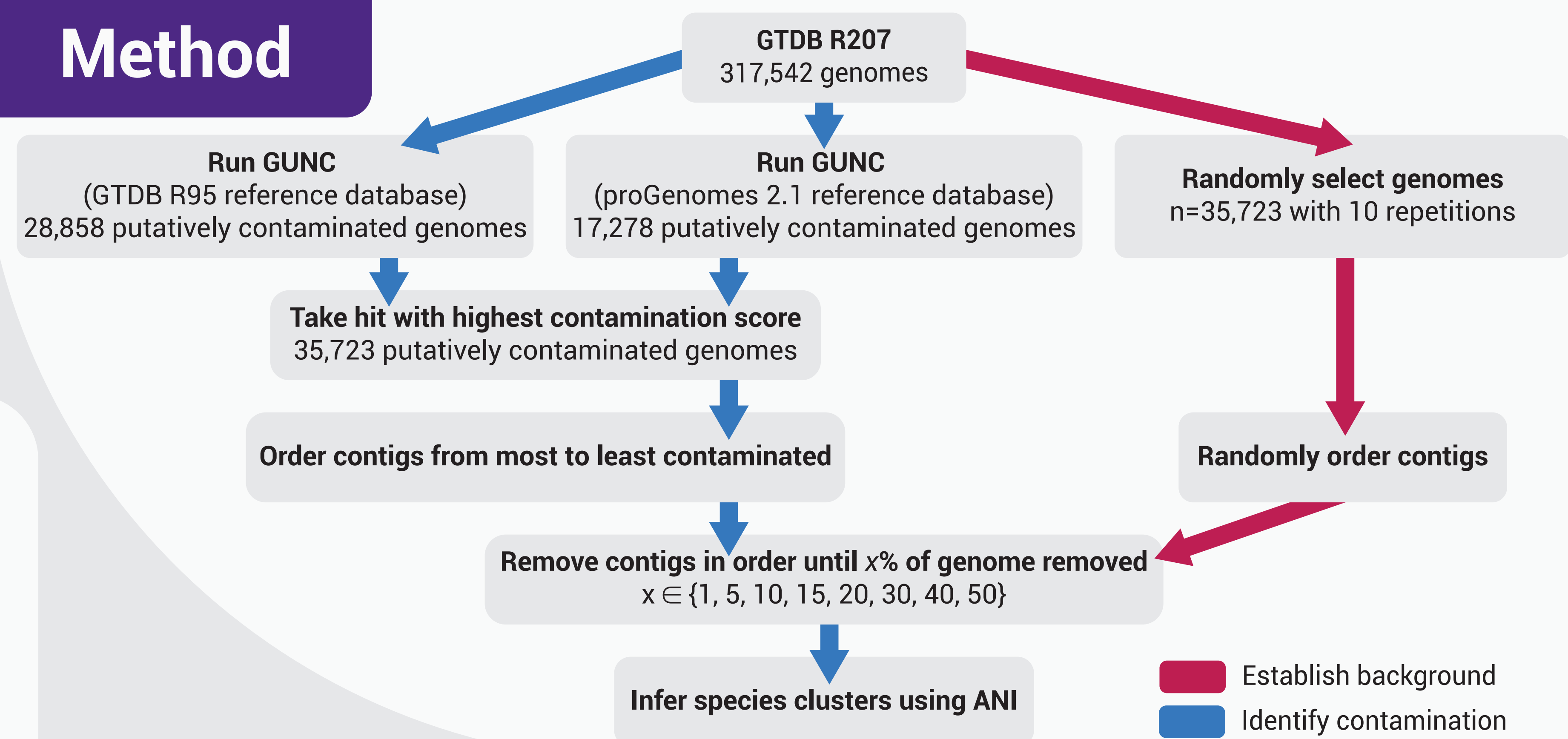
Closely related genomes blur interspecies boundaries

• 76.75% of the interspecies ANI comparisons were >92% ANI for the 226 genomes that changed species, compared to 16.58% for all genomes in GTDB R207.

• Significantly more non-representative genomes fall within multiple species ANI radii.

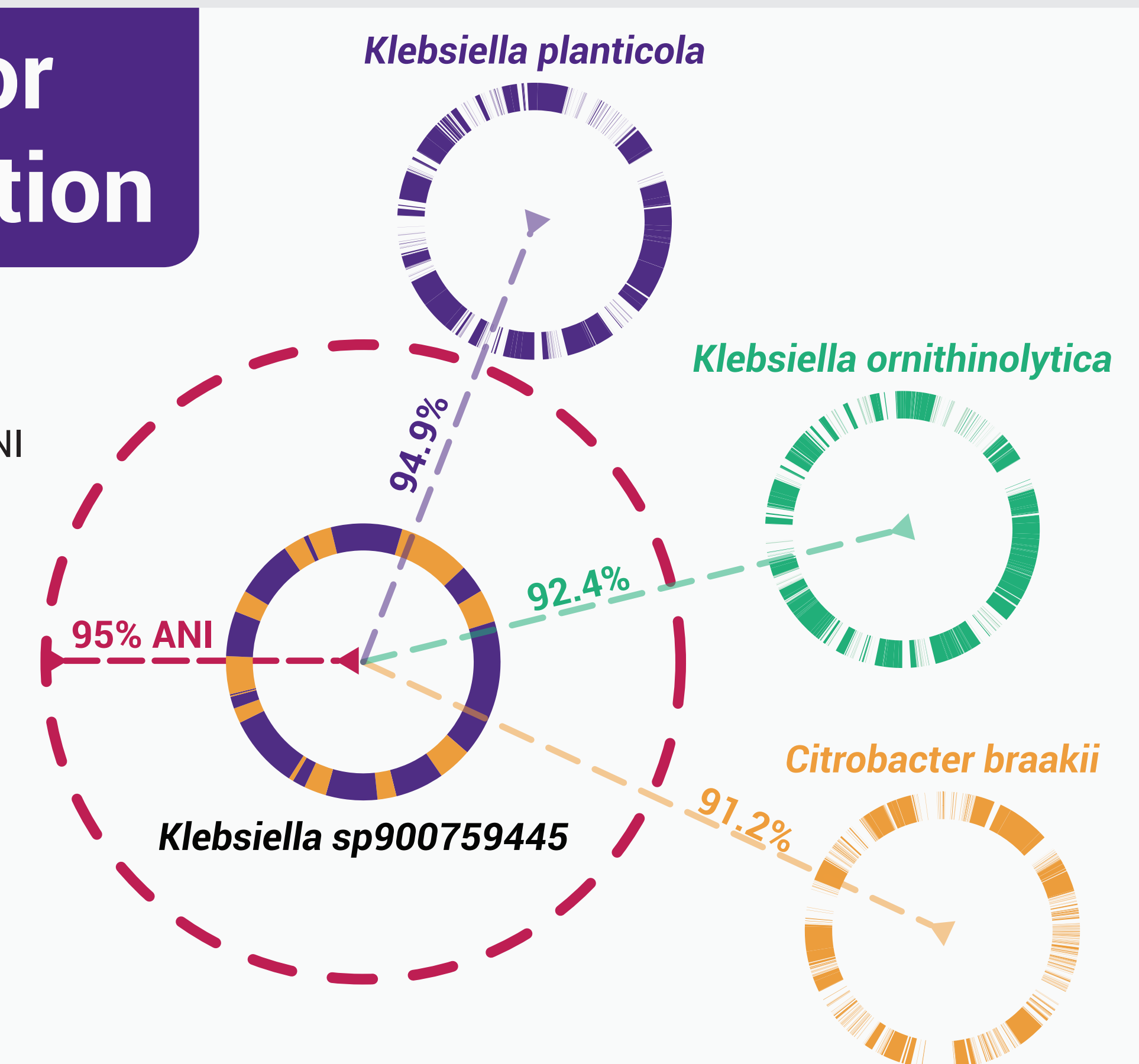
• Genomes are more likely to have a species reassignment if they are close to multiple representatives.

Method

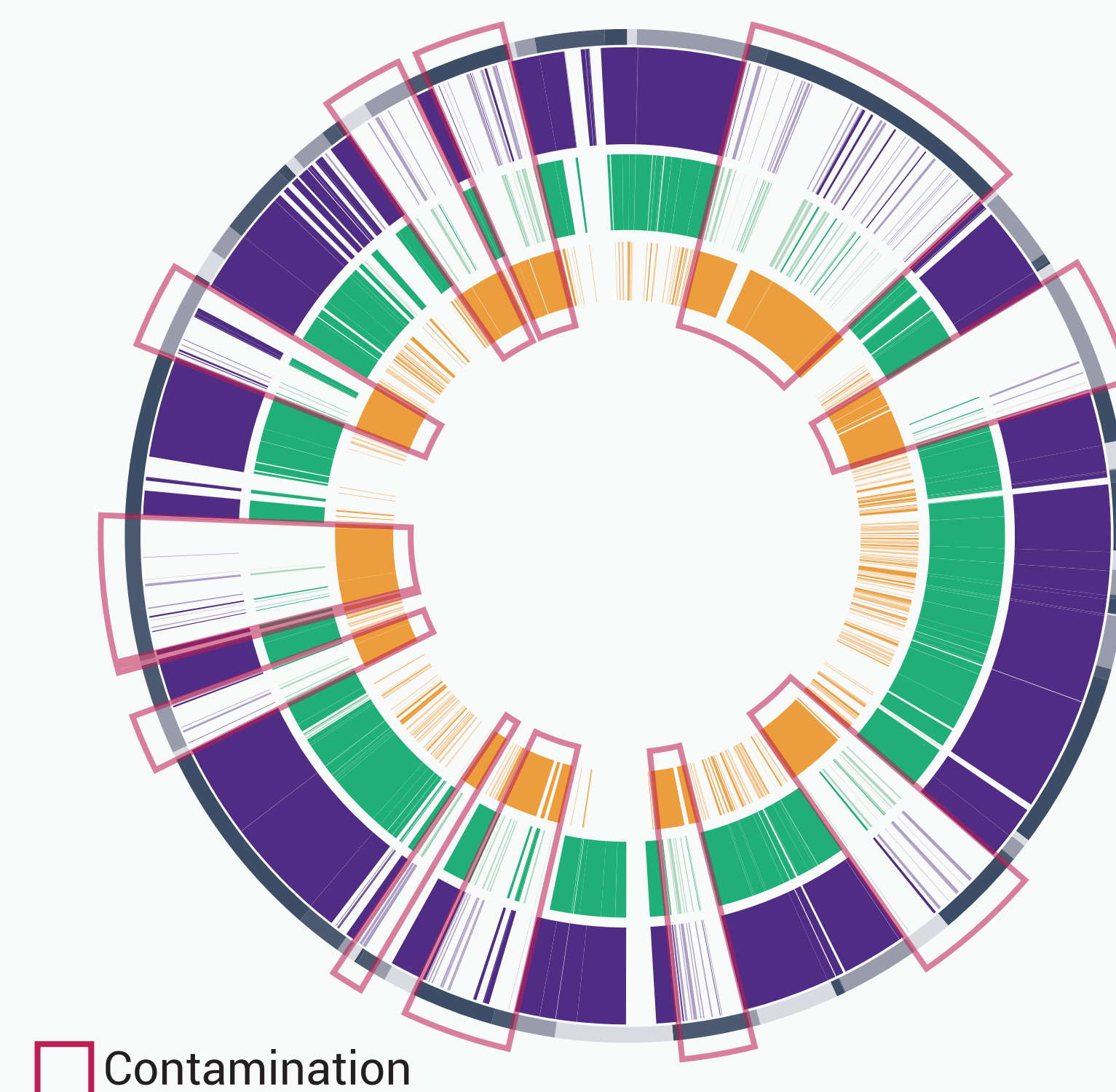


The perfect storm for species reclassification

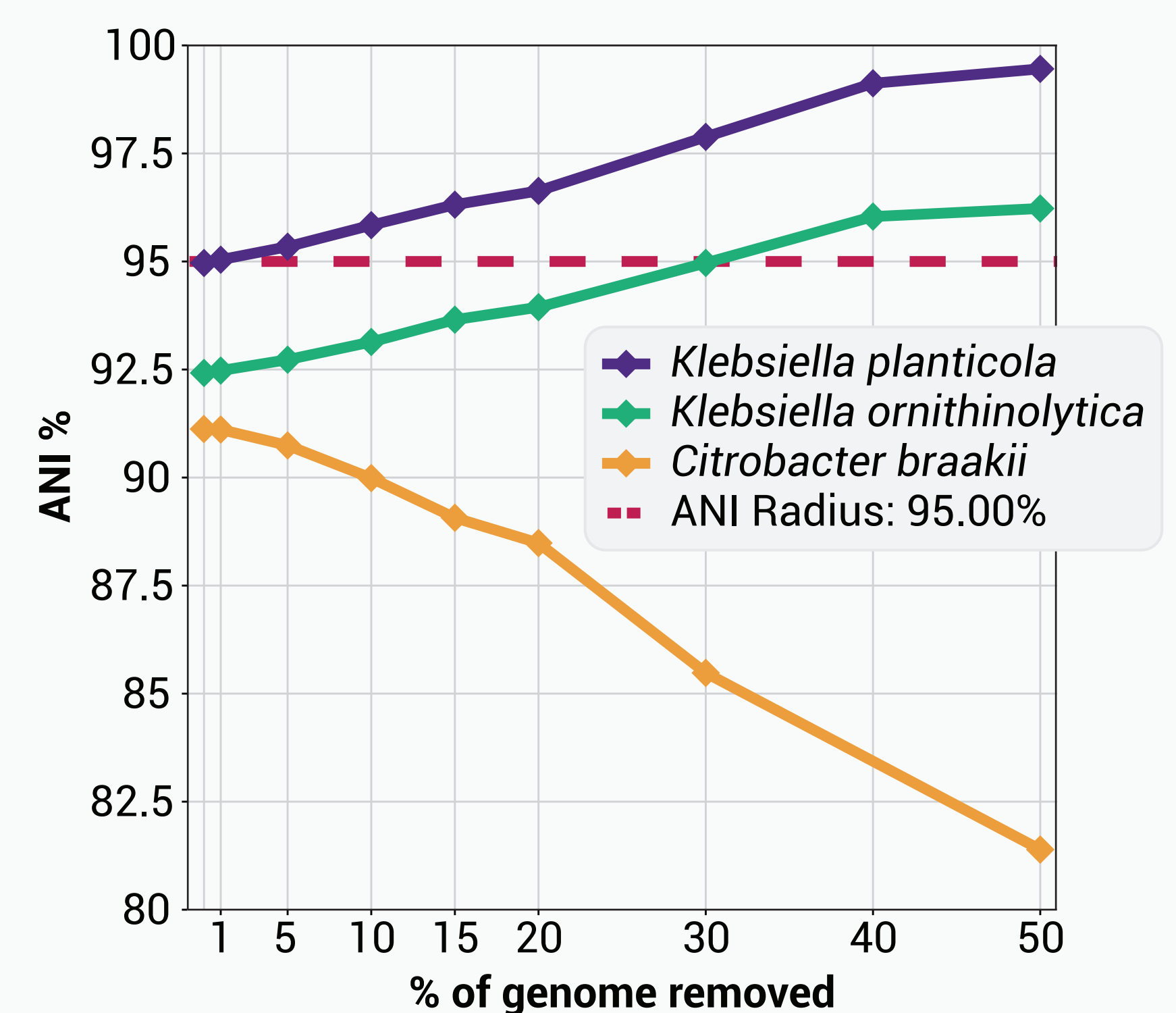
- Chimeric assemblies decrease interspecies ANI to closely related species.
- Removal of foreign contigs increases ANI to neighbouring species clusters.
- If interspecies ANI similarity is extremely close, a species change is more likely.
- Placeholder species are disproportionately affected.



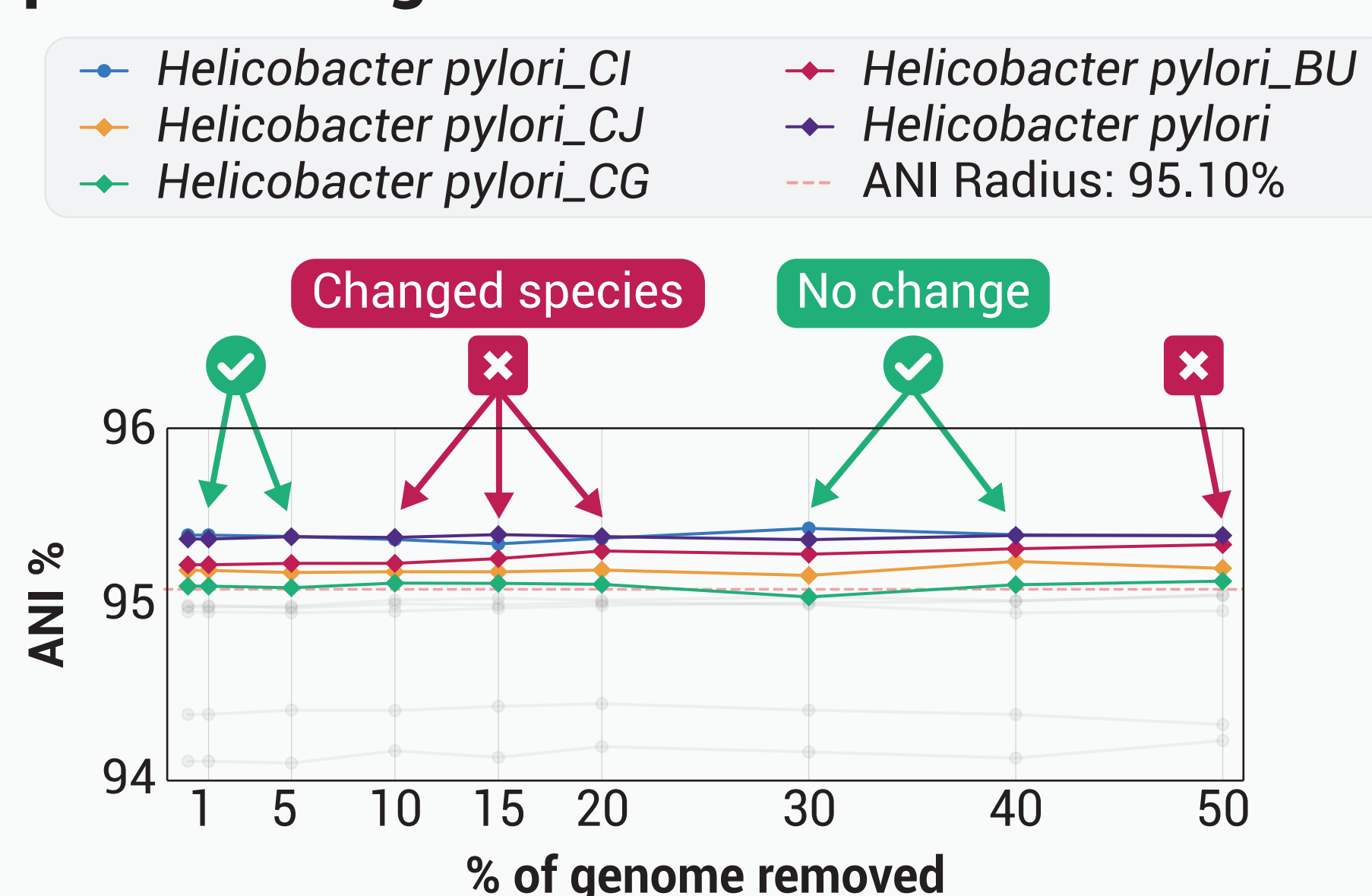
Whole genome alignment shows *C. braakii* contamination



Removing contamination merges *K. sp900759445* and *K. planticola*



Species assignments can be on a knives edge



Interspecies ANI is higher for genomes that changed species

