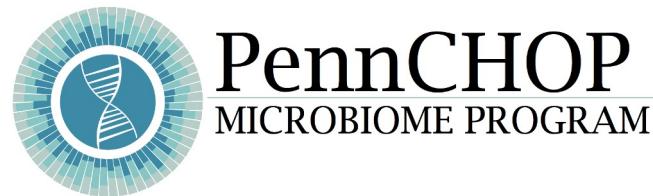


Fear and Trembling in Bacterial Genome Assembly

2023-03-15

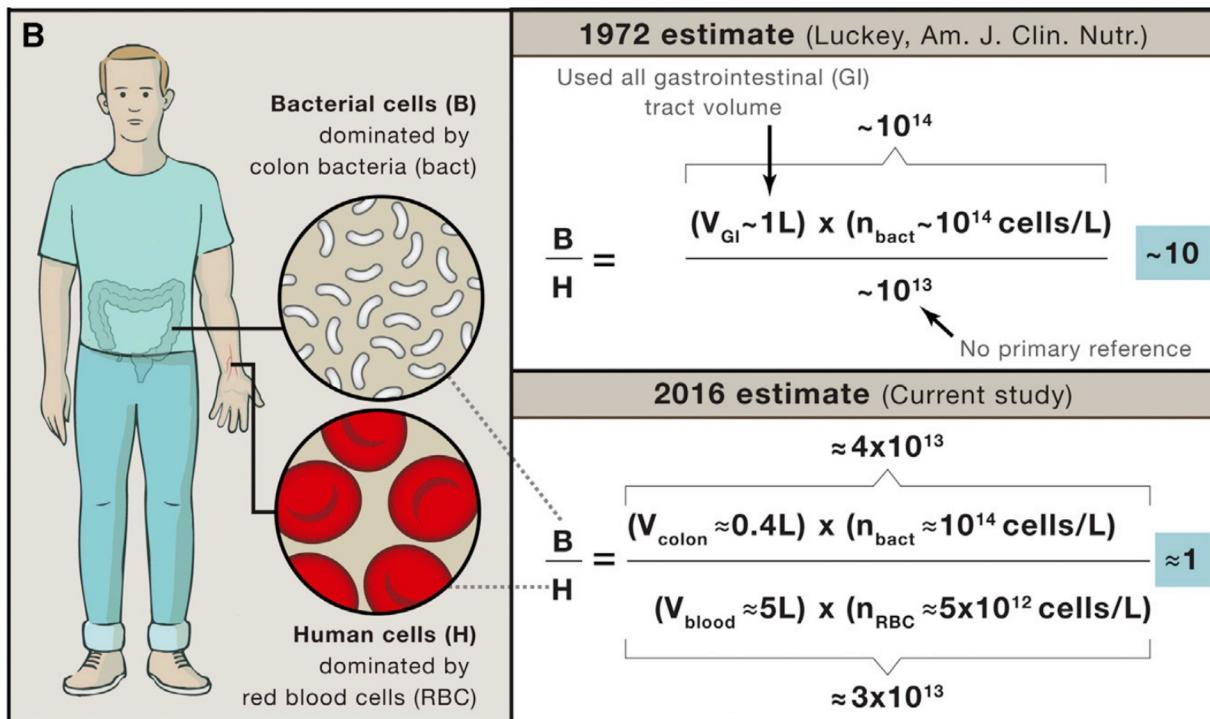


Kyle Bittinger

*Division of Gastroenterology, Hepatology, and Nutrition
CHOP Microbiome Center*

Illustration: Arwa Abbas

We co-exist with about as many bacterial cells as human cells

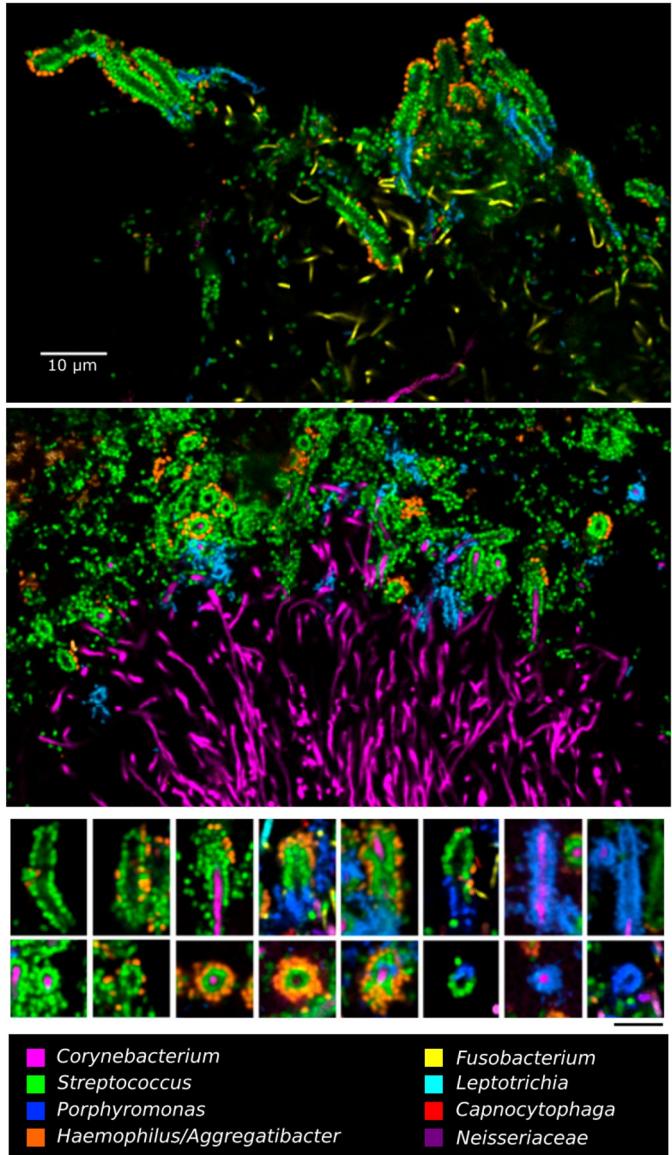


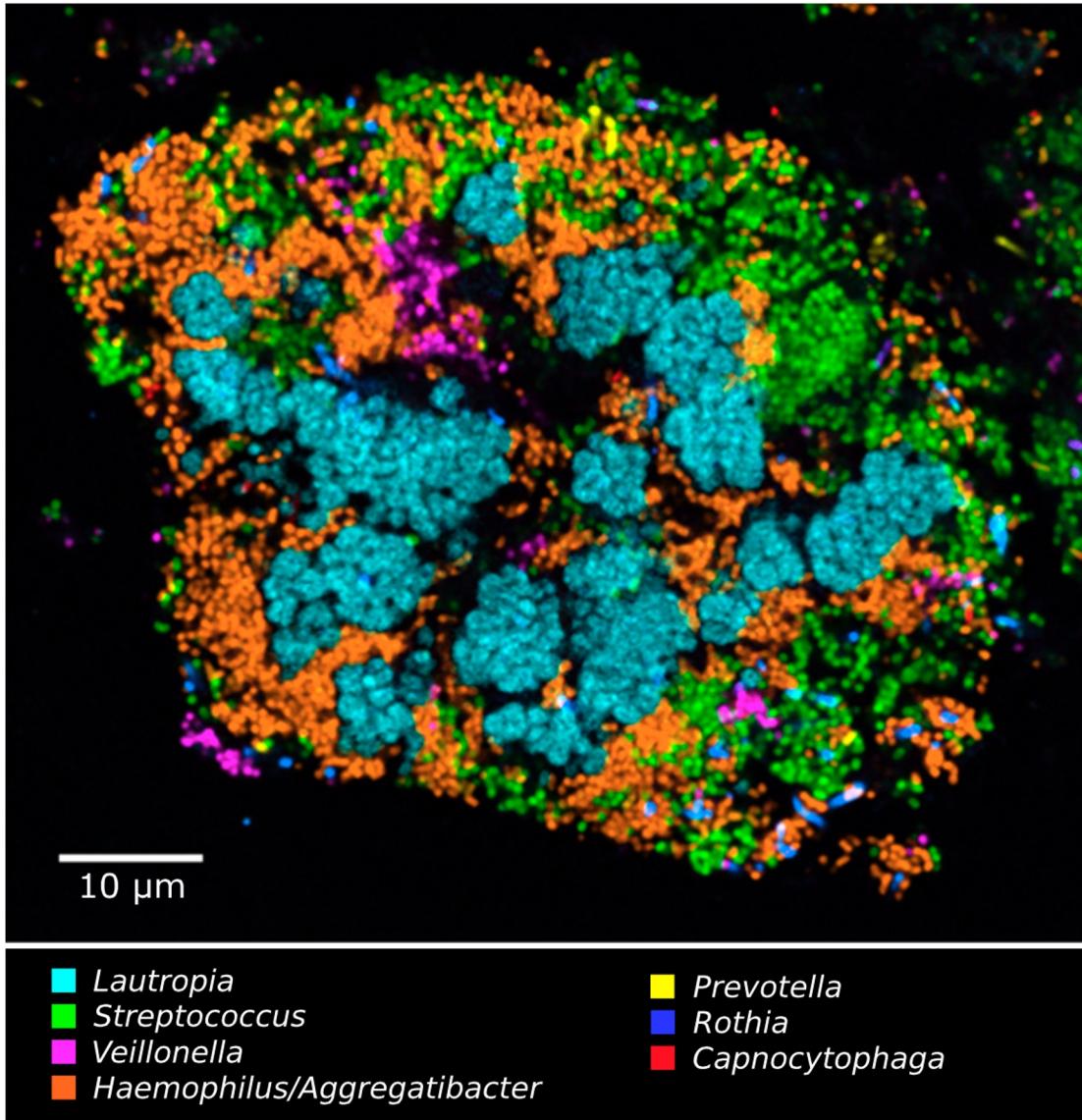
Sender R. Cell 164, 337 (2016).

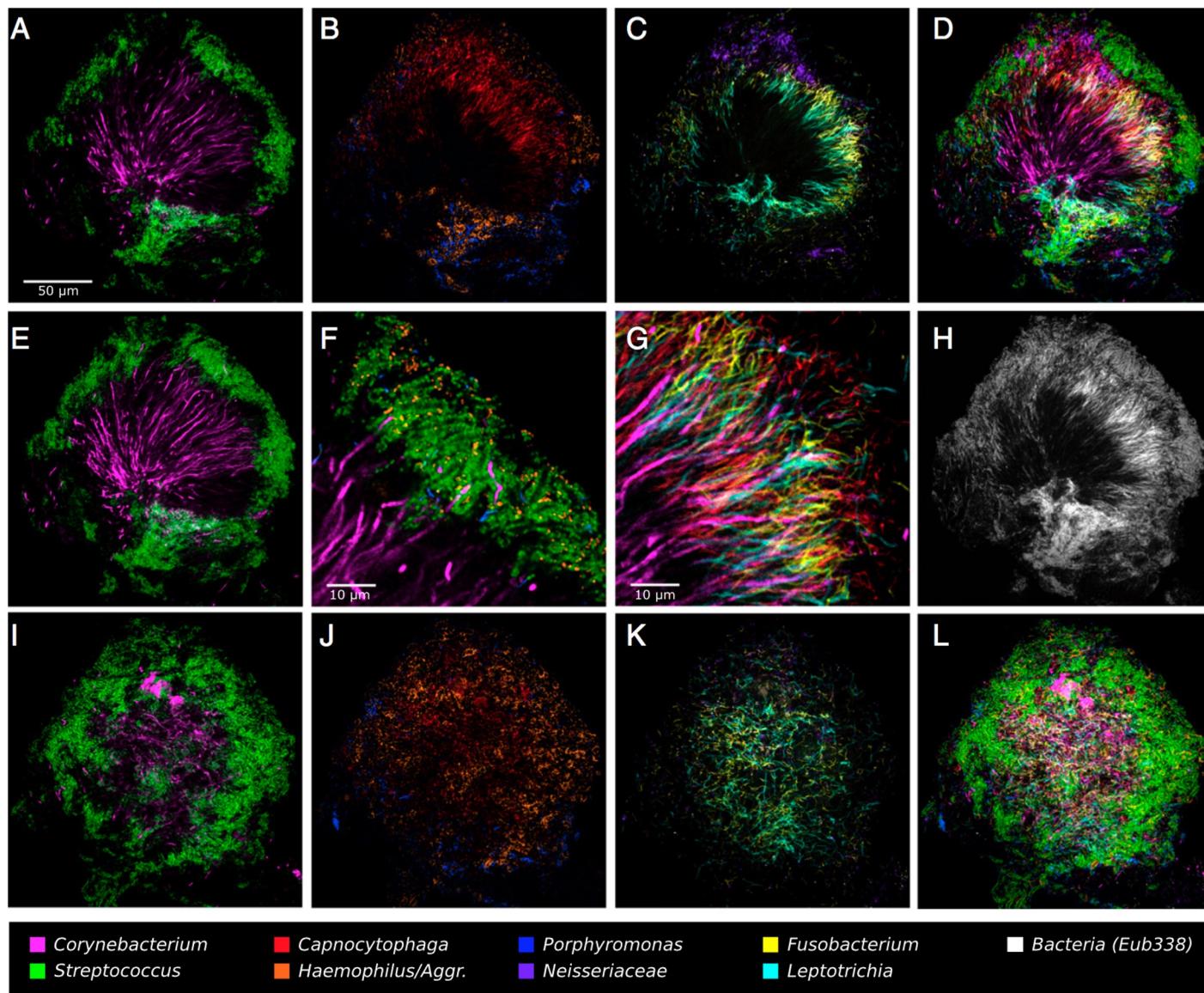
Roughly 10 trillion bacteria per human

Number of bacteria in colon about 10-50x that in rest of body

Red blood cells account for ~84% of human cells







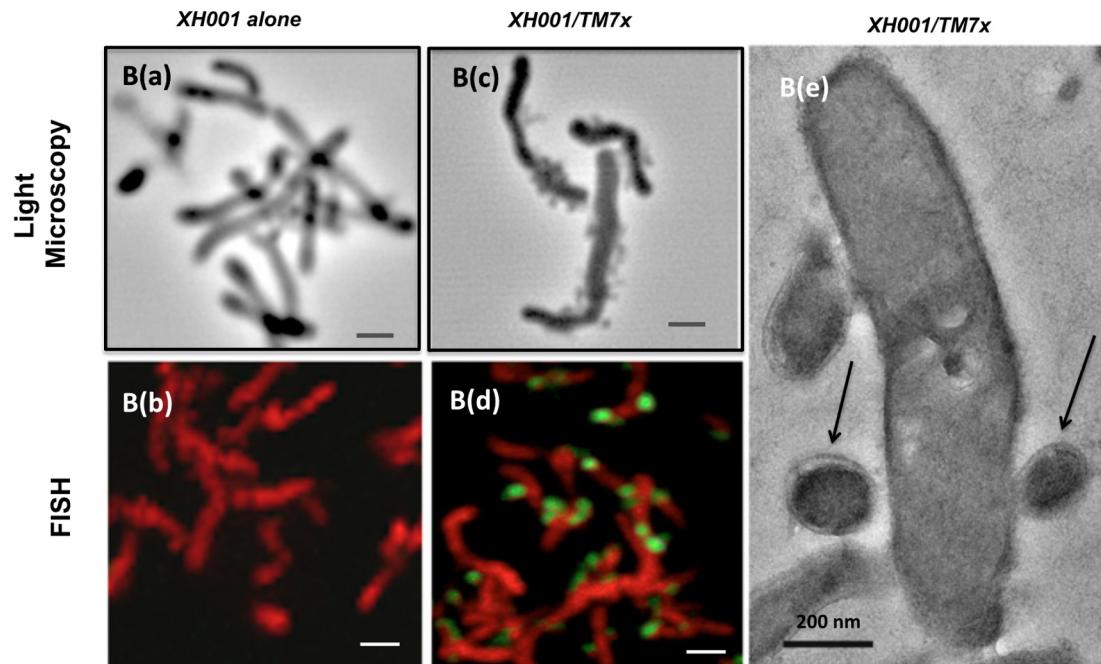
TM7 first observed in DNA-based studies, 18 years later in culture

Candidate division TM7 was established using 16S sequences recovered from a peat sample in 1996. (Rheims H. *J Ind Microbiol.* **17**, 159)

TM7 bacteria not cultured until 2014. (Soro V. *Appl Environ Microbiol.* **80**, 6480)

TM7x, pictured, is an obligate epibiont of Actinomyces.

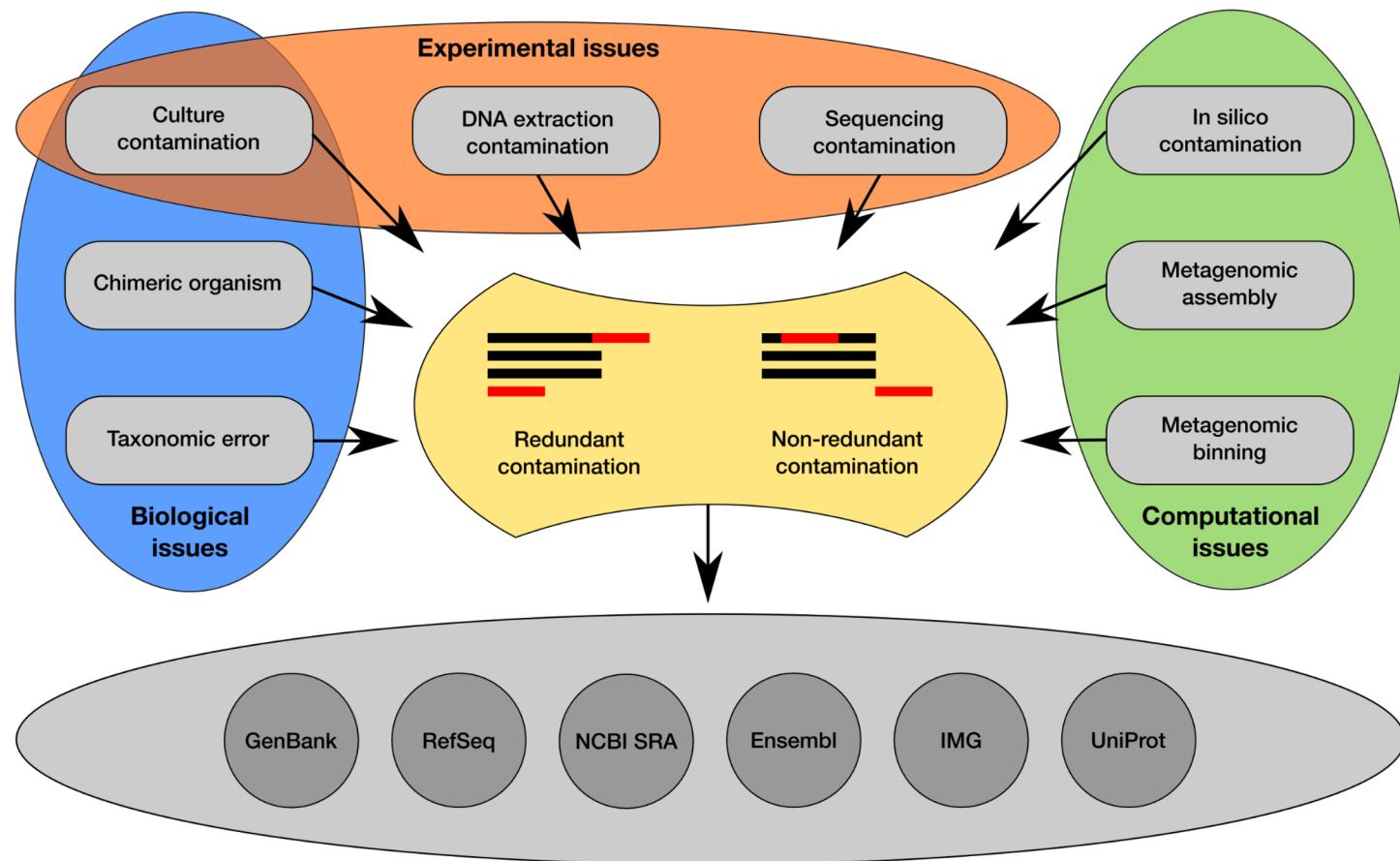
Complete genome: 705kb



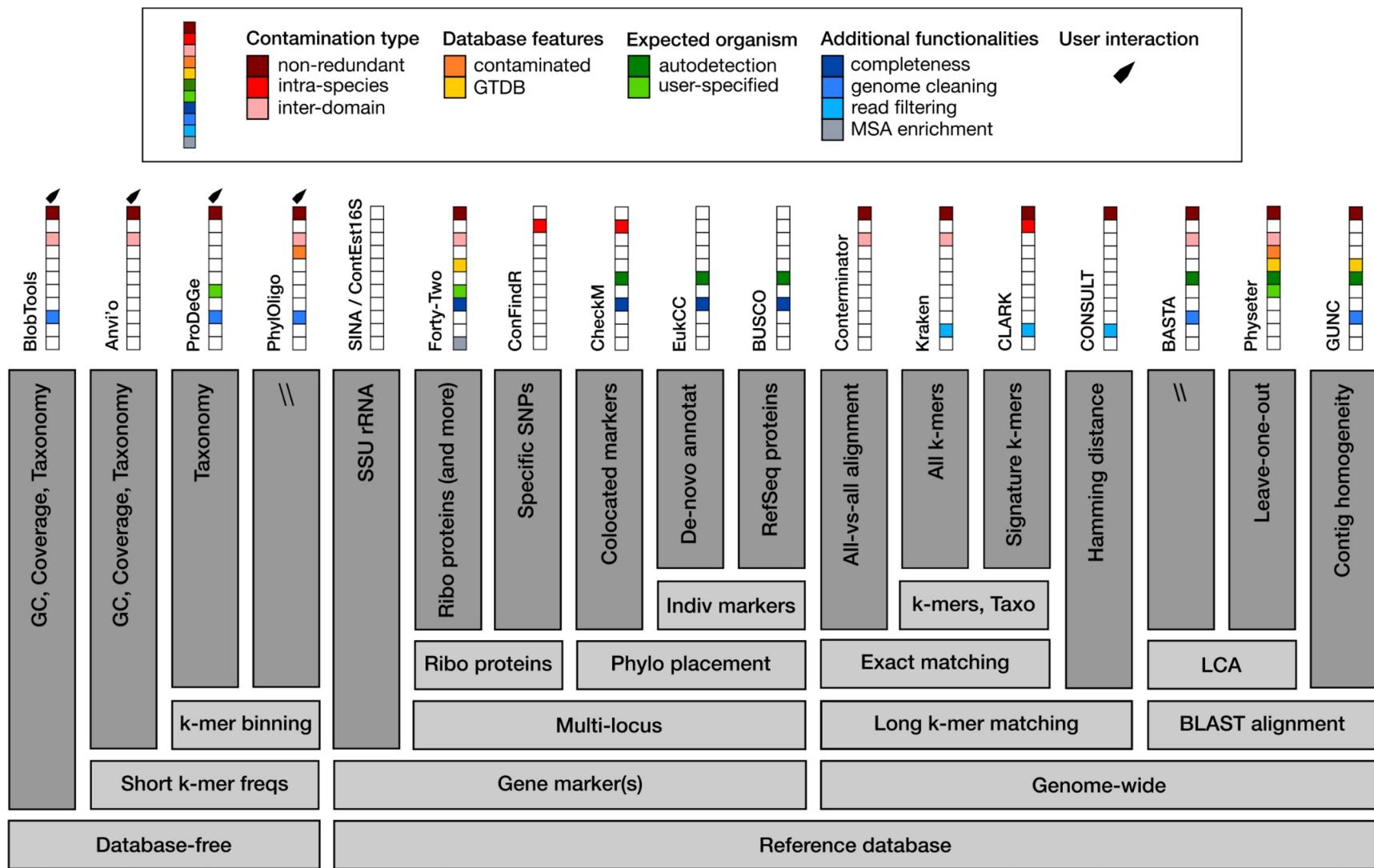
He PNAS **112**, 244 (2015).

1

Contamination detection in bacterial genomes



Cornet L, Baurain D. Contamination detection in genomic data: more is not enough. *Genome Biol.* 2022 Feb 21;23(1):60.



2

Genomes as metagenomes

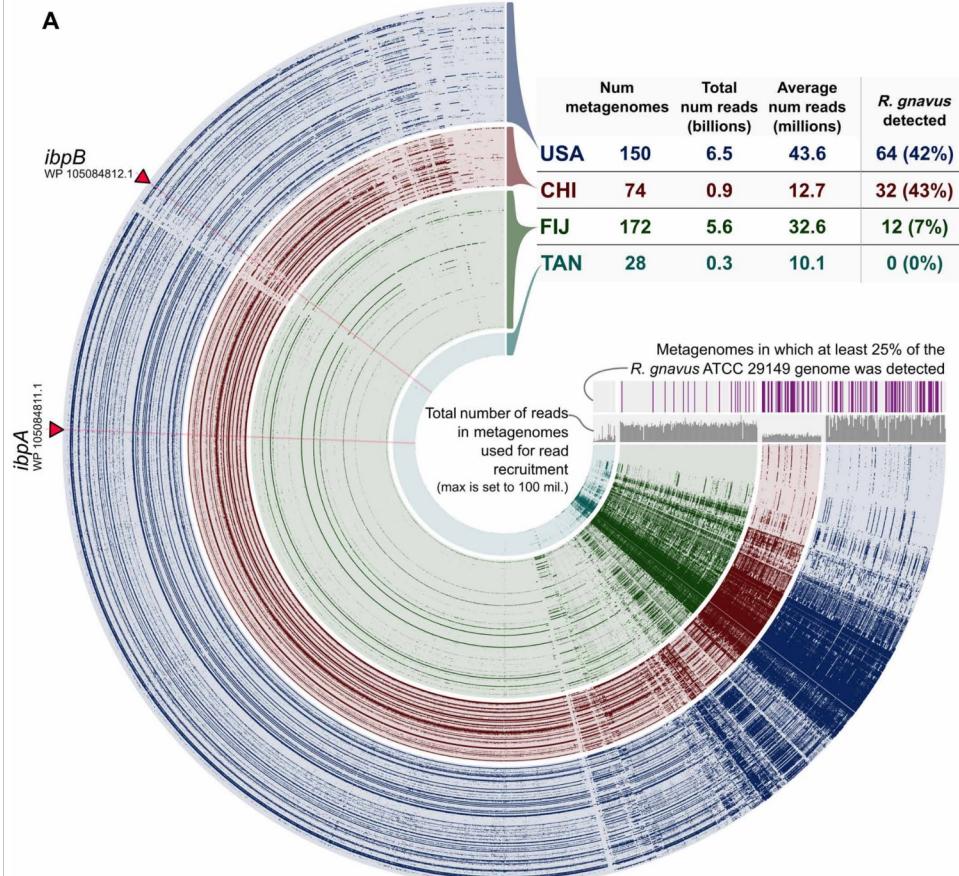
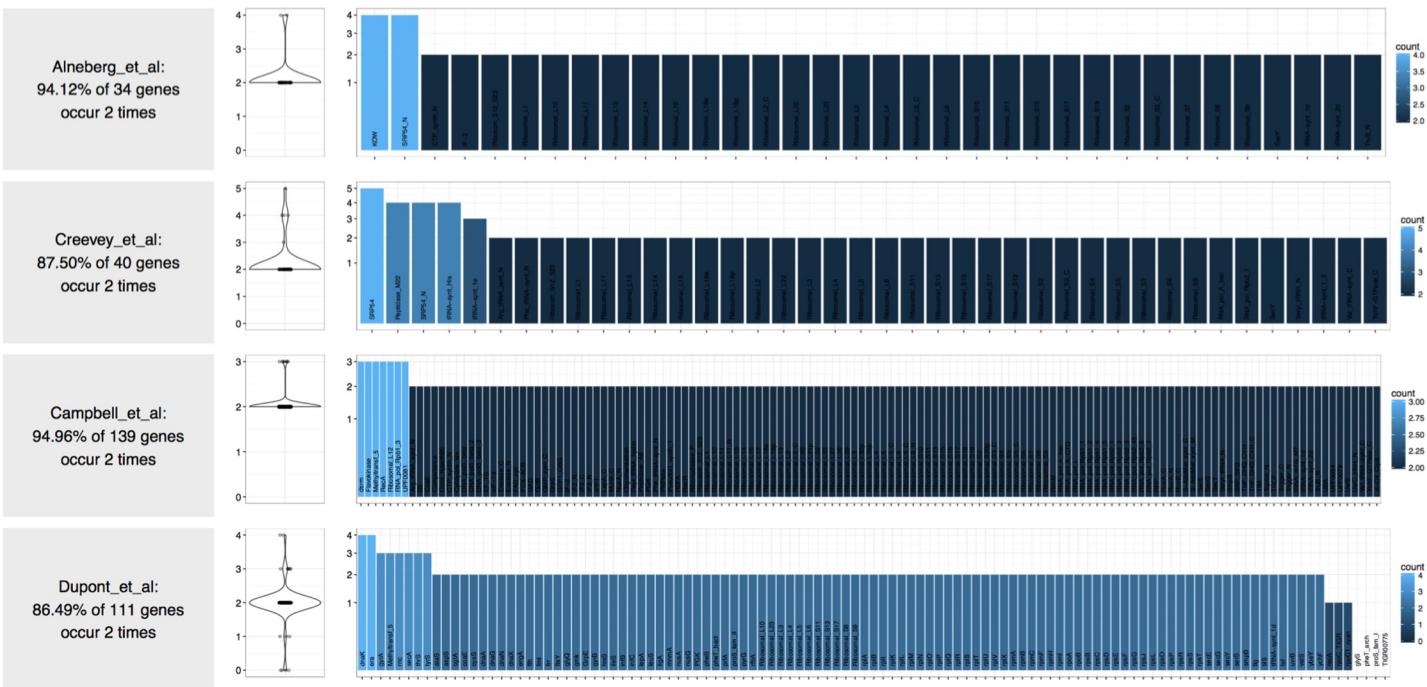


Fig. 5. Distribution of *R. gnavus* and its superantigens across human metagenomes. Dendrogram alignment of the *R. gnavus* ATCC 29149 genome to 424 human metagenomes (data files S3 and S4). Each spoke represents one gene in the *R. gnavus* genome, and each layer represents an individual human metagenome. The two superantigen genes are labeled. Intensity represents coverage of the open reading frame in the metagenome.

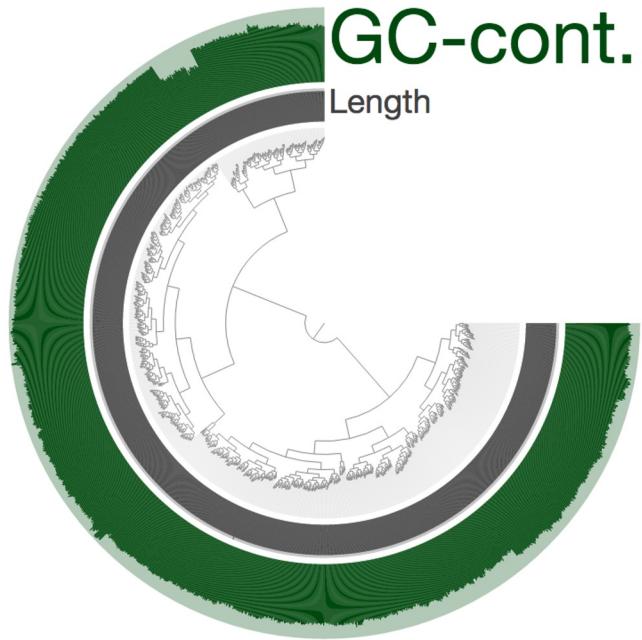
Eren AM, et al. Community-led, integrated, reproducible multi-omics with anvi'o. *Nat Microbiol.* 2021 Jan;6(1):3-6.

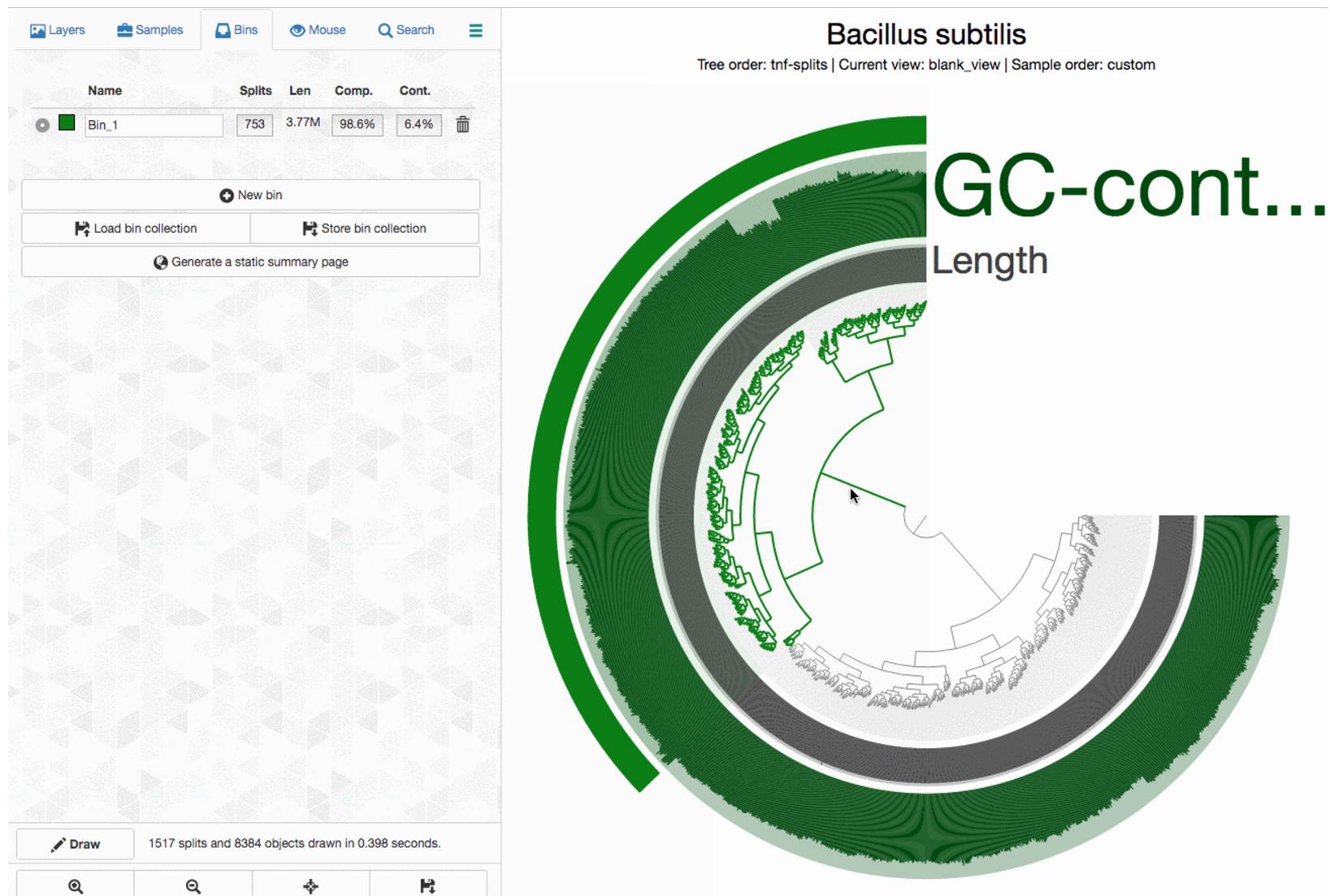


Example Combining two genomes in one cell: Stable cloning of the Synechocystis PCC6803 genome in the *Bacillus subtilis* 168 genome

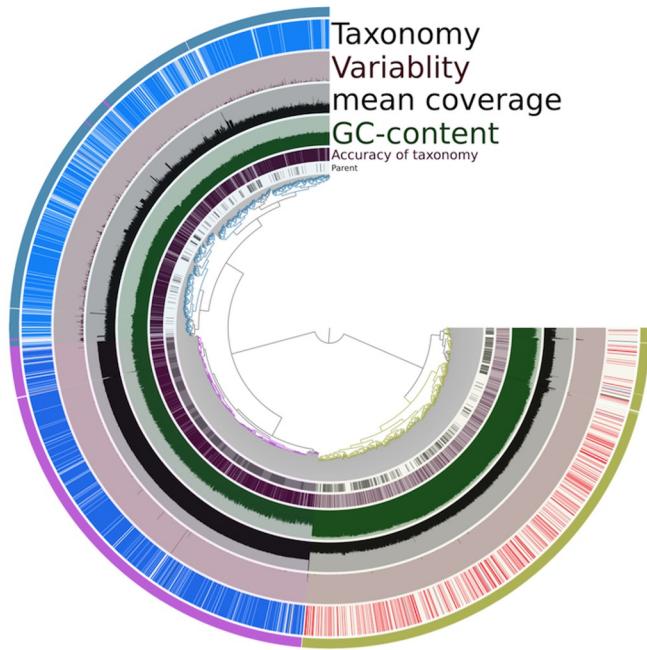


Bacillus subtilis
Tree order: tnf-splits | Current view: blank_view | Sample order: custom

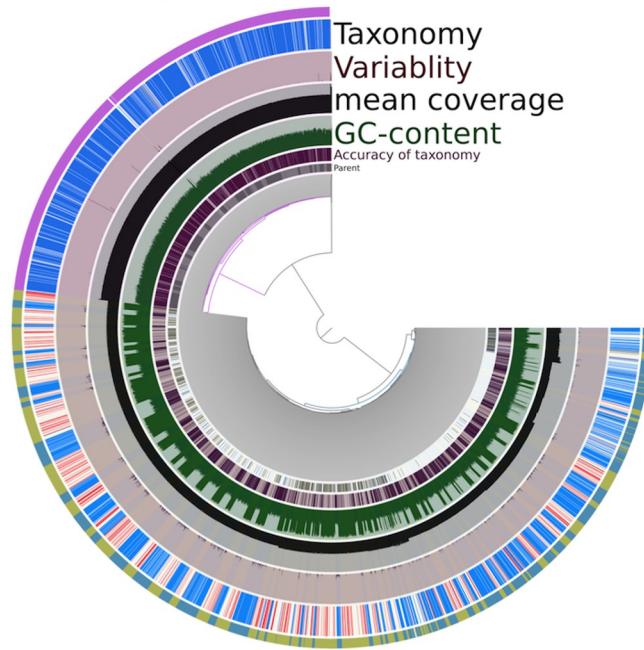




Clustering based on sequence composition



Clustering based on sequence composition and abundance



Bin	Taxonomy	Total Size	Num Contigs	N50	GC Content	Compl.	Contam.
selection_2	<i>Bacillus subtilis</i>	3.70 Mb	54	159,263	41.56%	98.47%	8.02%
selection_3	Unknown	4.20 Mb	1,394	3,773	69.34%	78.37%	4.19%
selection_1	<i>Bacillus anthracis</i>	3.29 Mb	1,794	1,858	34.99%	54.45%	10.60%

NCBI Taxonomy Browser

Entrez PubMed Nucleotide Protein Genome Structure PMC Taxonomy BioCollections

Search for as lock

Display 3 levels using filter: none

[\[Ruminococcus\].gnavus](#) ¹⁾

Taxonomy ID: 33038 (for references in articles please use NCBI:txid33038)

current name

Ruminococcus gnavus Moore et al. 1976 (Approved Lists 1980) Moore et al. 1976 in [Skerman VBD et al. (1980)]

type strain of *Ruminococcus gnavus* Moore et al. 1976 (Approved Lists 1980): [ATCC:29149](#), [VPI C7-9](#), [JCM:6515](#)

homotypic synonym:

"**Mediterraneibacter gnavus**" (Moore et al. 1976) Togo et al. 2018, effective name ²⁾

NCBI BLAST name: firmicutes

Rank: species

Genetic code: [Translation table 11 \(Bacterial, Archaeal and Plant Plastid\)](#)

[Lineage](#)(full)

[cellular organisms](#); [Bacteria](#); [Terrabacteria group](#); [Bacillota](#); [Clostridia](#); [Eubacteriales](#); [Lachnospiraceae](#); [Mediterraneibacter](#)

Entrez records			
Database name	Subtree links	Direct links	Links from type
Nucleotide	19,448	19,207	129
Protein	197,843	177,578	-
Structure	48	15	-
Genome	1	1	-
Popset	4	4	-
GEO Datasets	12	9	-
PubMed Central	10	4	-
Gene	7,503	1	-
SRA Experiments	144	122	-
Protein Clusters	2,504	2,504	-
Identical Protein Groups	78,137	76,776	-
BioProject	53	42	-
BioSample	1,408	1,382	11
Assembly	165	158	7
PubChem BioAssay	7	7	-
Taxonomy	6	1	-



Q Search NCBI ...

Log in

NCBI Datasets **Taxonomy** Genome Gene Command-line tools Documentation

Bacteria / Bacillota / Clostridia / Eubacteriales / Lachnospiraceae / Mediterraneibacter / [Ruminococcus] gnavus

[Ruminococcus] gnavus ATCC 29149

[ruminococcus] gnavus atcc 29149 is a strain of [ruminococcus] gnavus.

[Browse taxonomy](#)

Current scientific name [Ruminococcus] gnavus ATCC 29149

Taxonomic rank strain

NCBI Taxonomy ID 411470

For more details see [NCBI Taxonomy](#)

View the legacy [Genome page](#)

Genome

[Browse all 7 genomes](#)

Reference genome

ASM983137v1

Kyungpook National University (2020). Strain: ATCC 29149.

RefSeq: GCF_009831375.1



Search NCBI ...

Log in

NCBI Datasets Taxonomy Genome Gene Command-line tools Documentation

Genome

Download a genome data package including genome, transcript and protein sequence, annotation and a data report

Selected taxa

[Ruminococcus] gnavus

Enter one or more taxonomic names

Filters

Download

Select columns

189 genomes

Rows per page

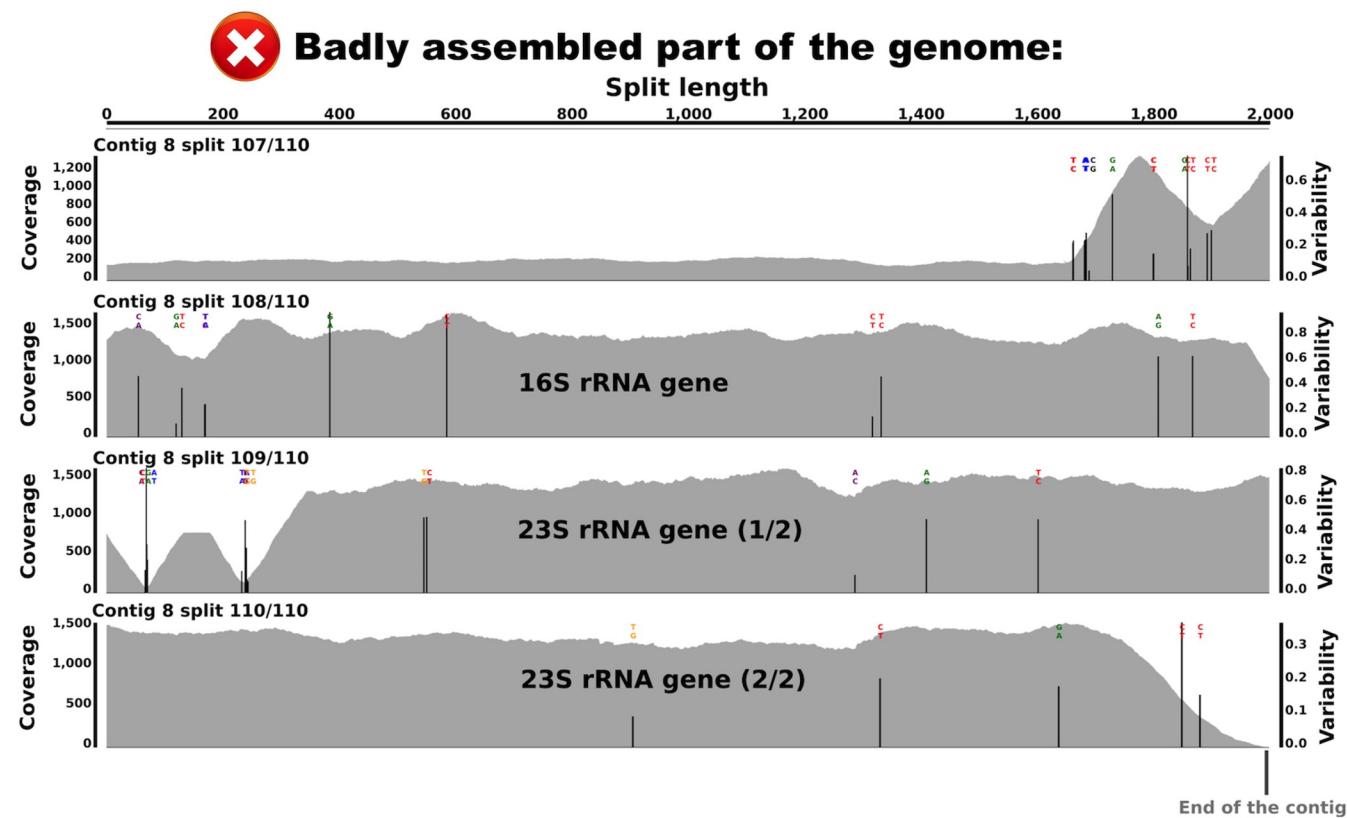
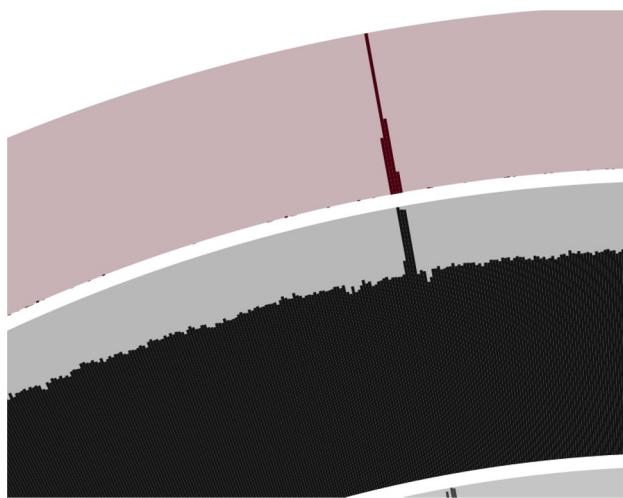
20

1-20 of 189



<input type="checkbox"/> Assembly	GenBank	RefSeq	Scientific name	Modifier	Annotation	Action
<input type="checkbox"/> ASM983137v1	GCA_009831375.1	GCF_009831375.1	[Ruminococcus] gnavus ATCC ...	ATCC 29149 (strain)		
<input type="checkbox"/> ASM2627830v1	GCA_026278305.1	GCF_026278305.1	[Ruminococcus] gnavus	JCM6515 (strain)		
<input type="checkbox"/> ASM2627832v1	GCA_026278325.1	GCF_026278325.1	[Ruminococcus] gnavus	JCM6515 (strain)		

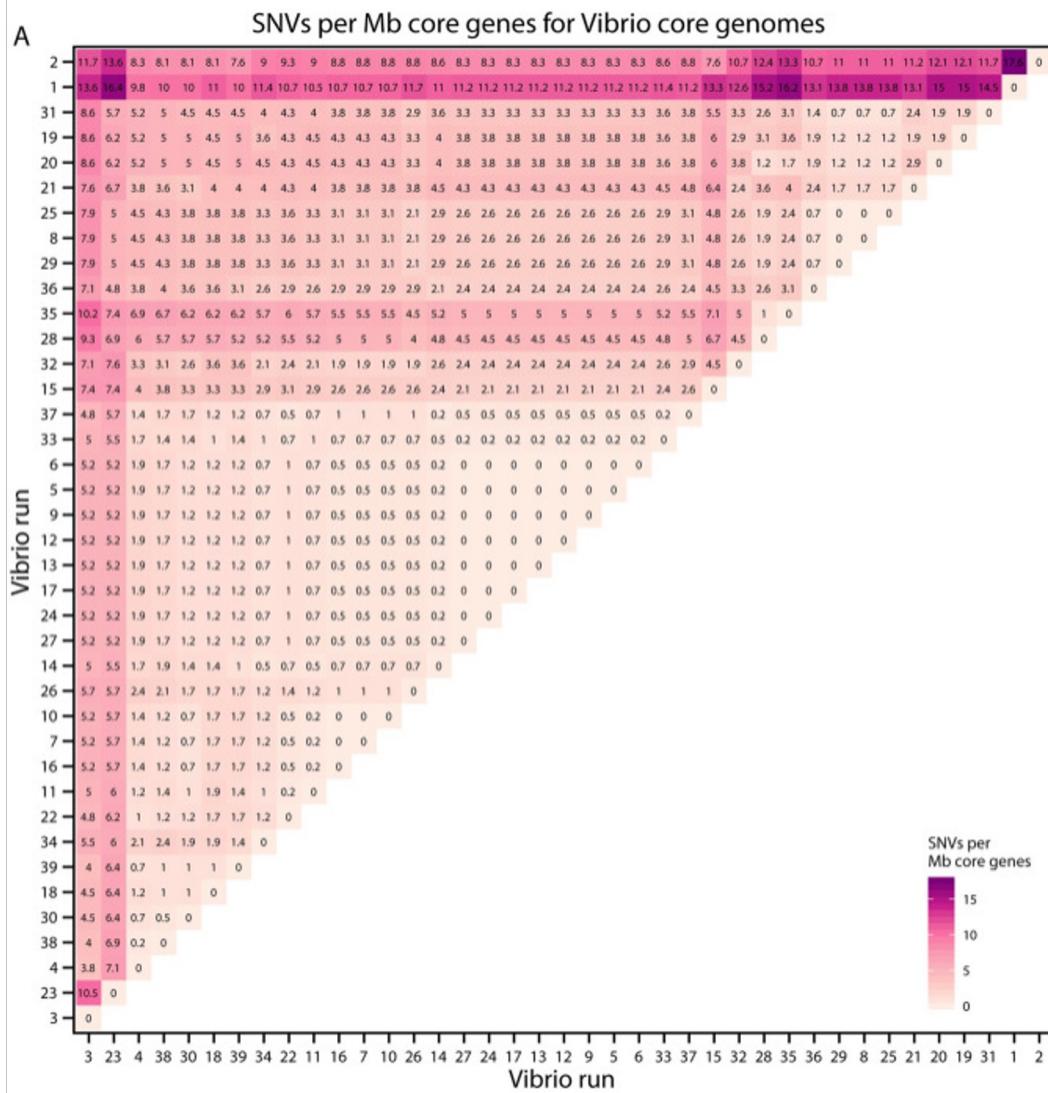




3

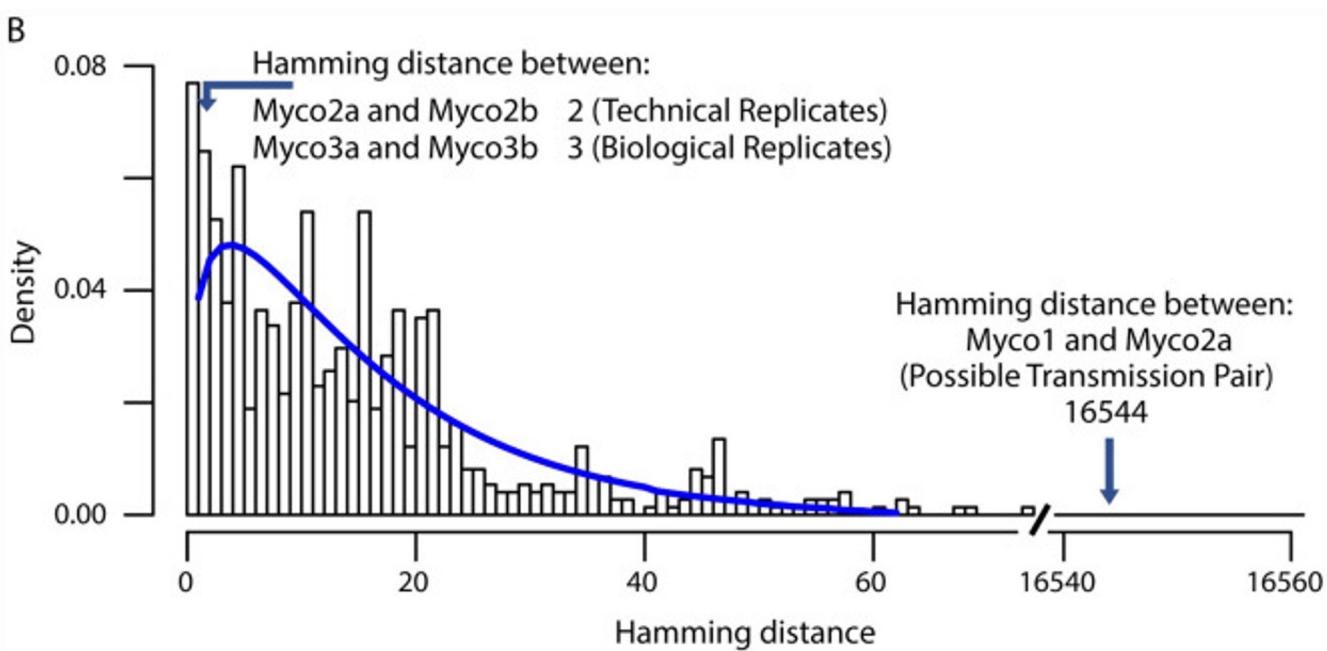
Technical replicates and strain-level differences

A



Genome assembly of a single *Vibrio campbellii* genome, 39 technical replicates

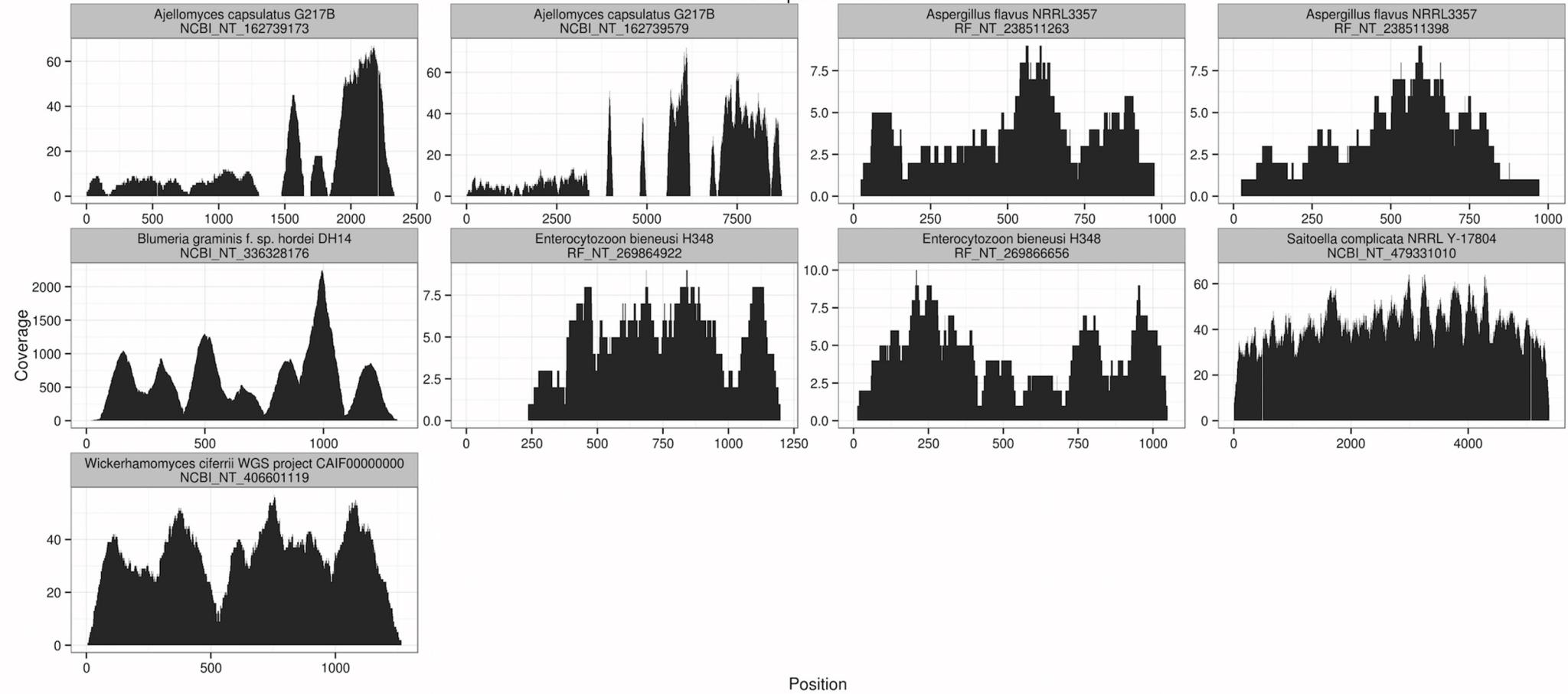
Gu CH, Zhao C, Hofstaedter C, Tebas P, Glaser L, Baldassano R, Bittinger K, Mattei LM, Bushman FD. Investigating hospital *Mycobacterium chelonae* infection using whole genome sequencing and hybrid assembly. *PLoS One*. 2020 Nov 9;15(11):e0236533.



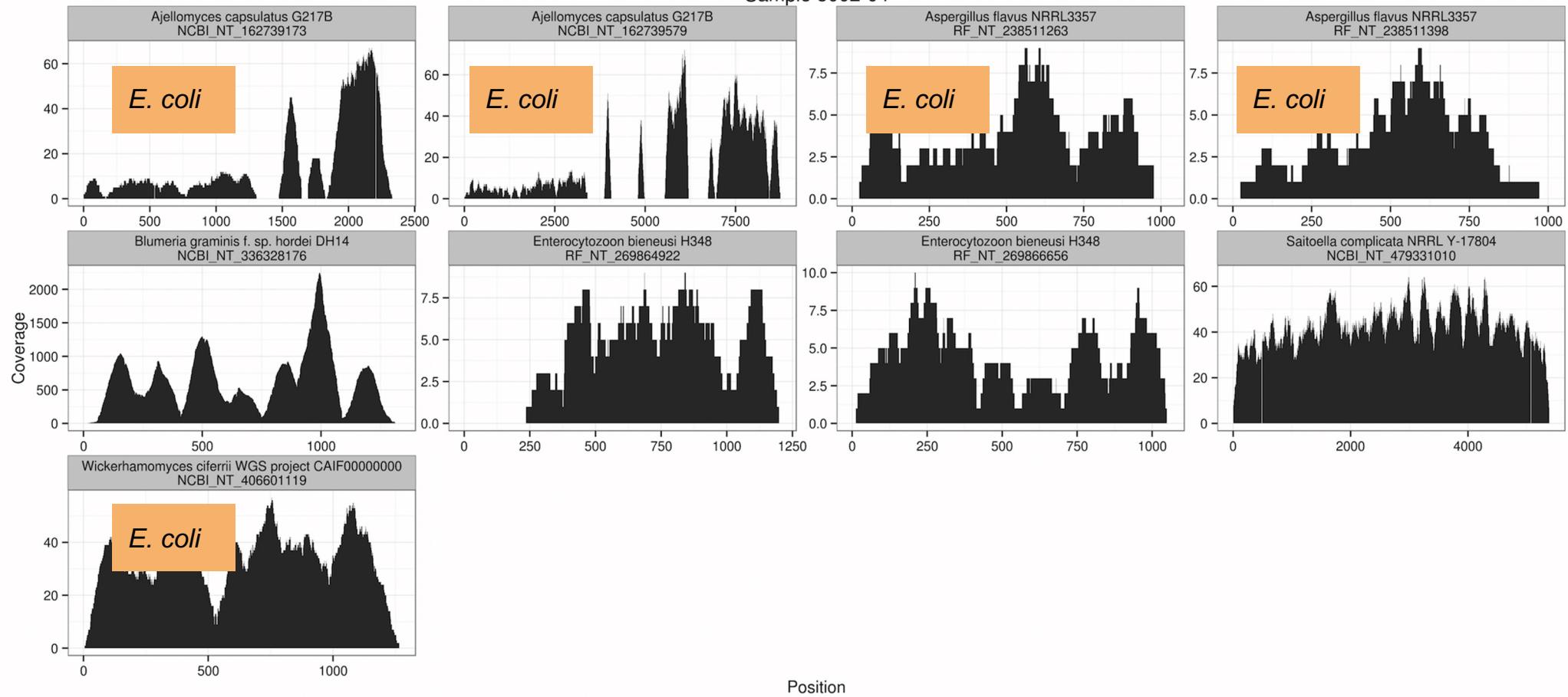
4

**Database contamination,
human DNA contamination**

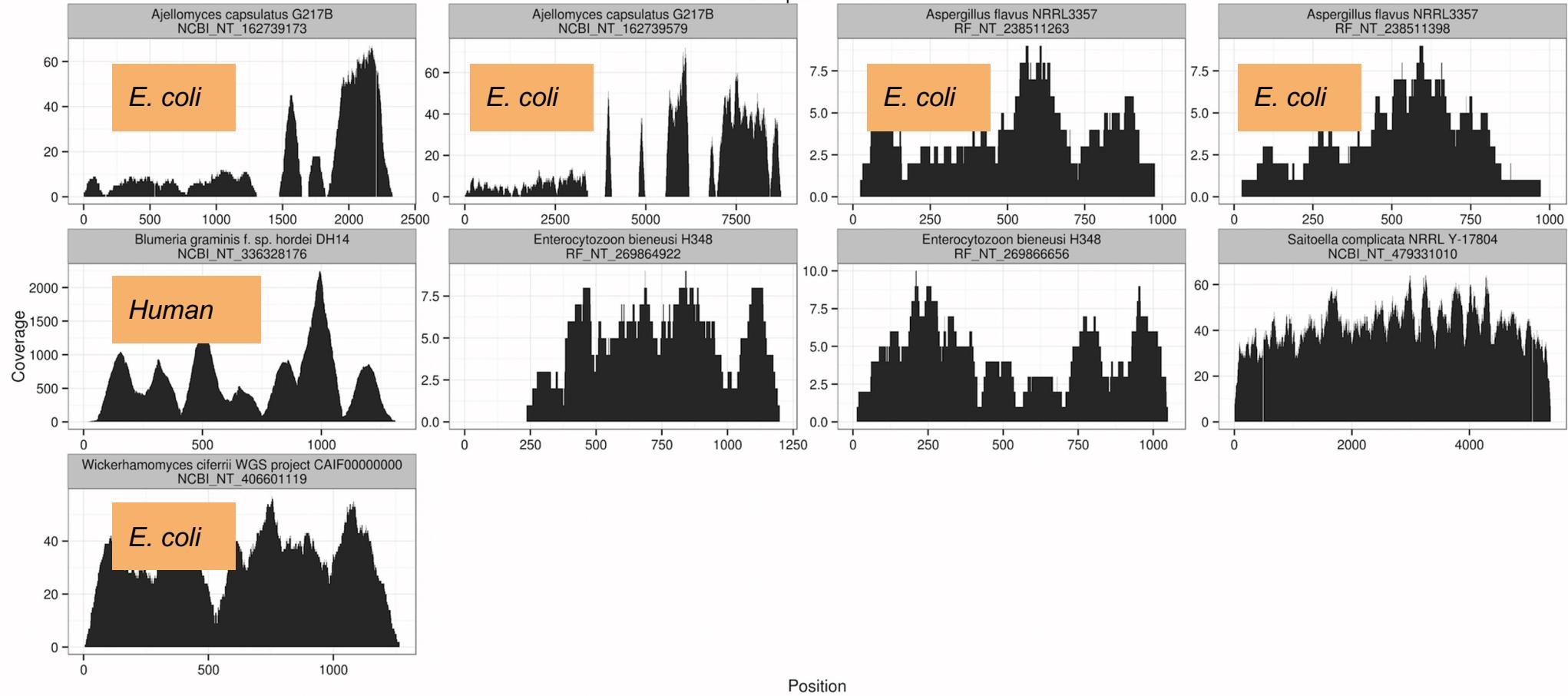
Sample 5002-04



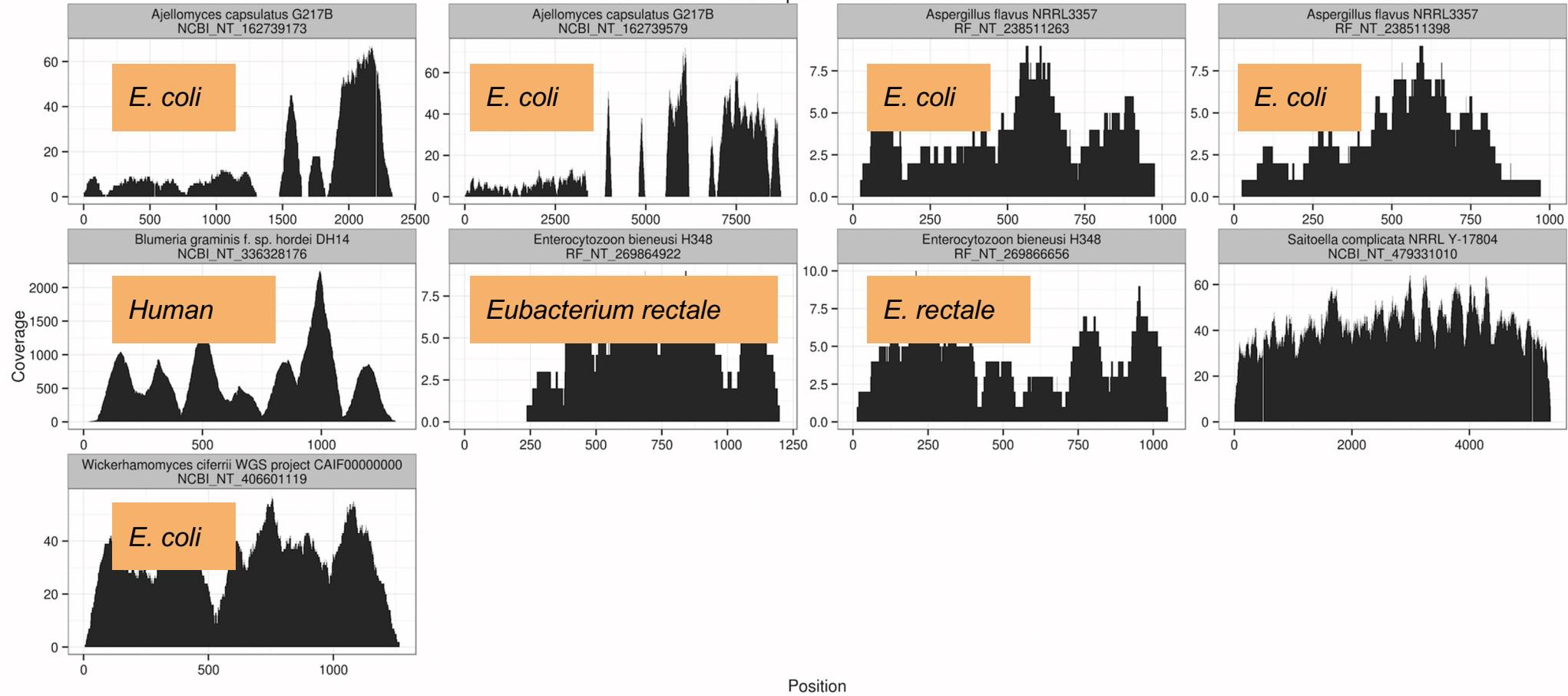
Sample 5002-04



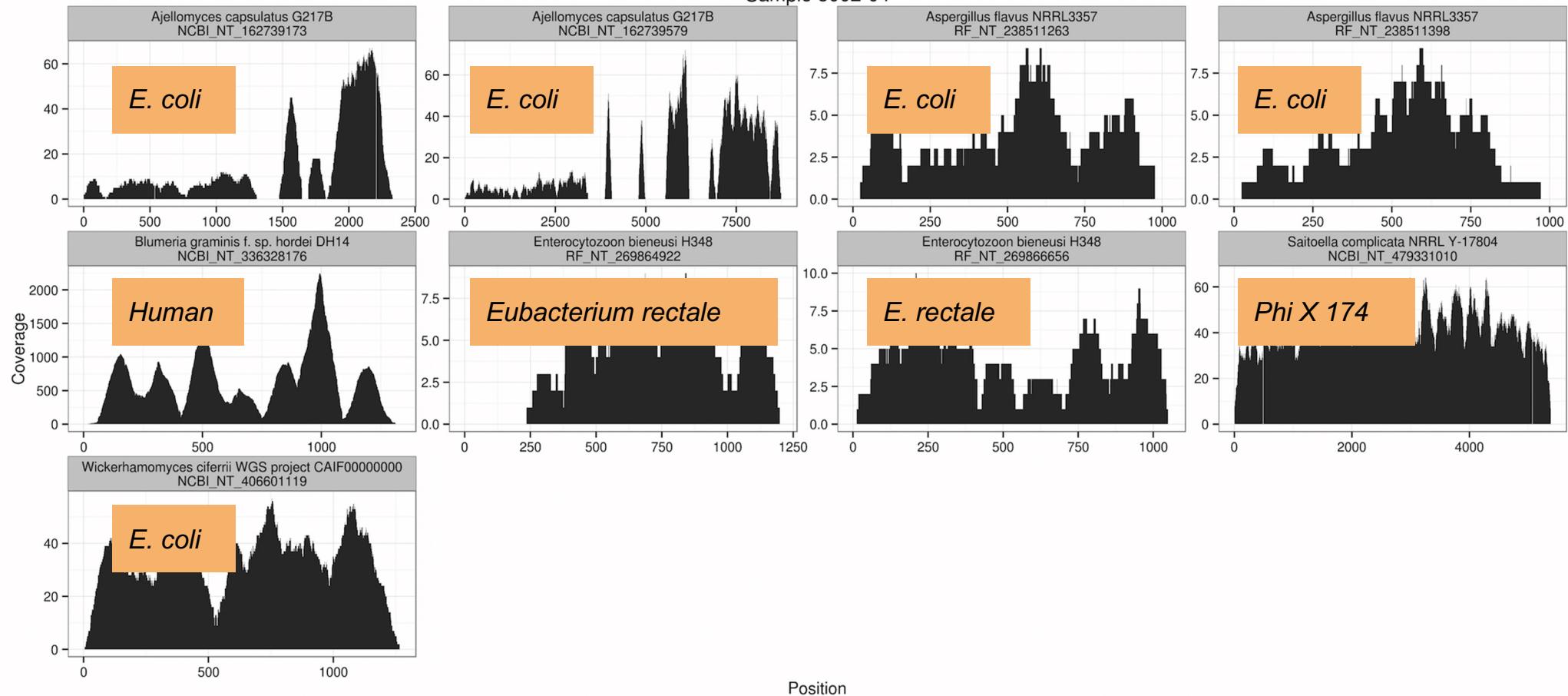
Sample 5002-04



Sample 5002-04



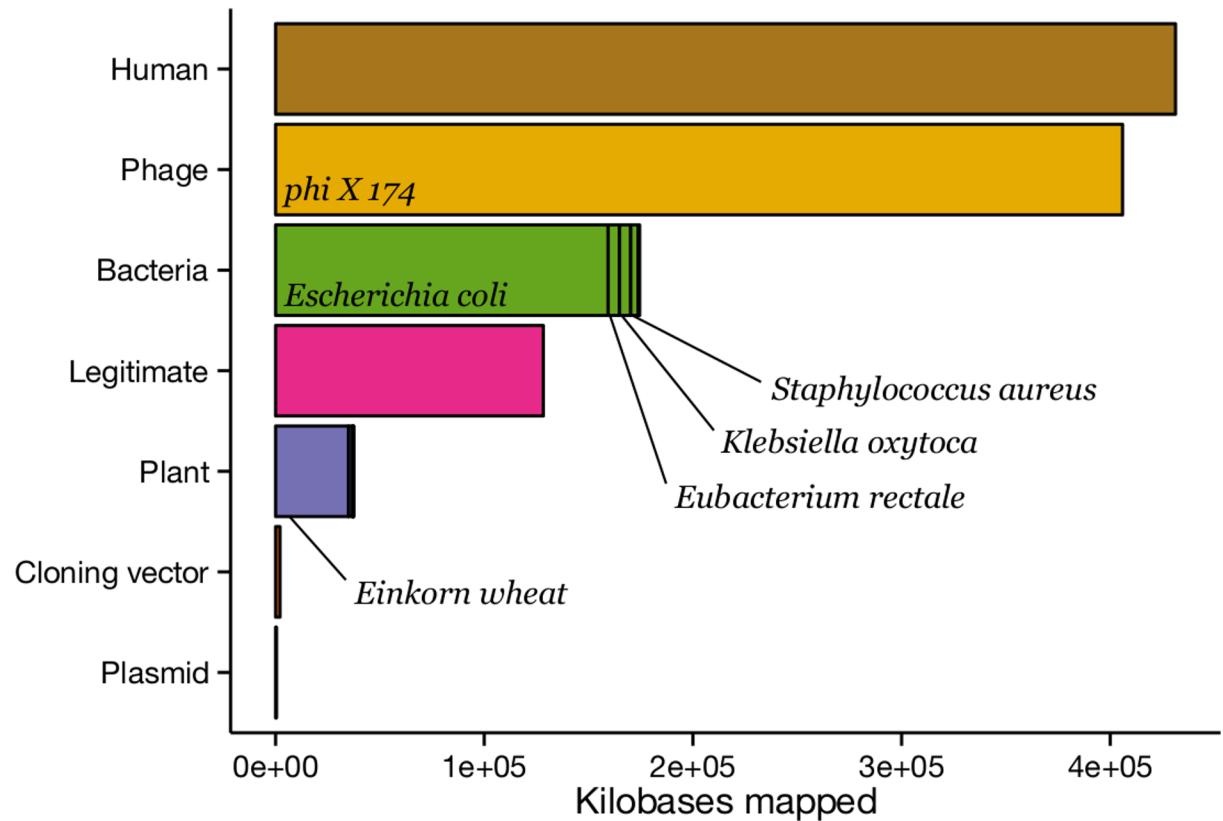
Sample 5002-04



Sources of misattribution in fungal genomes detected in study of pediatric Crohn's Disease

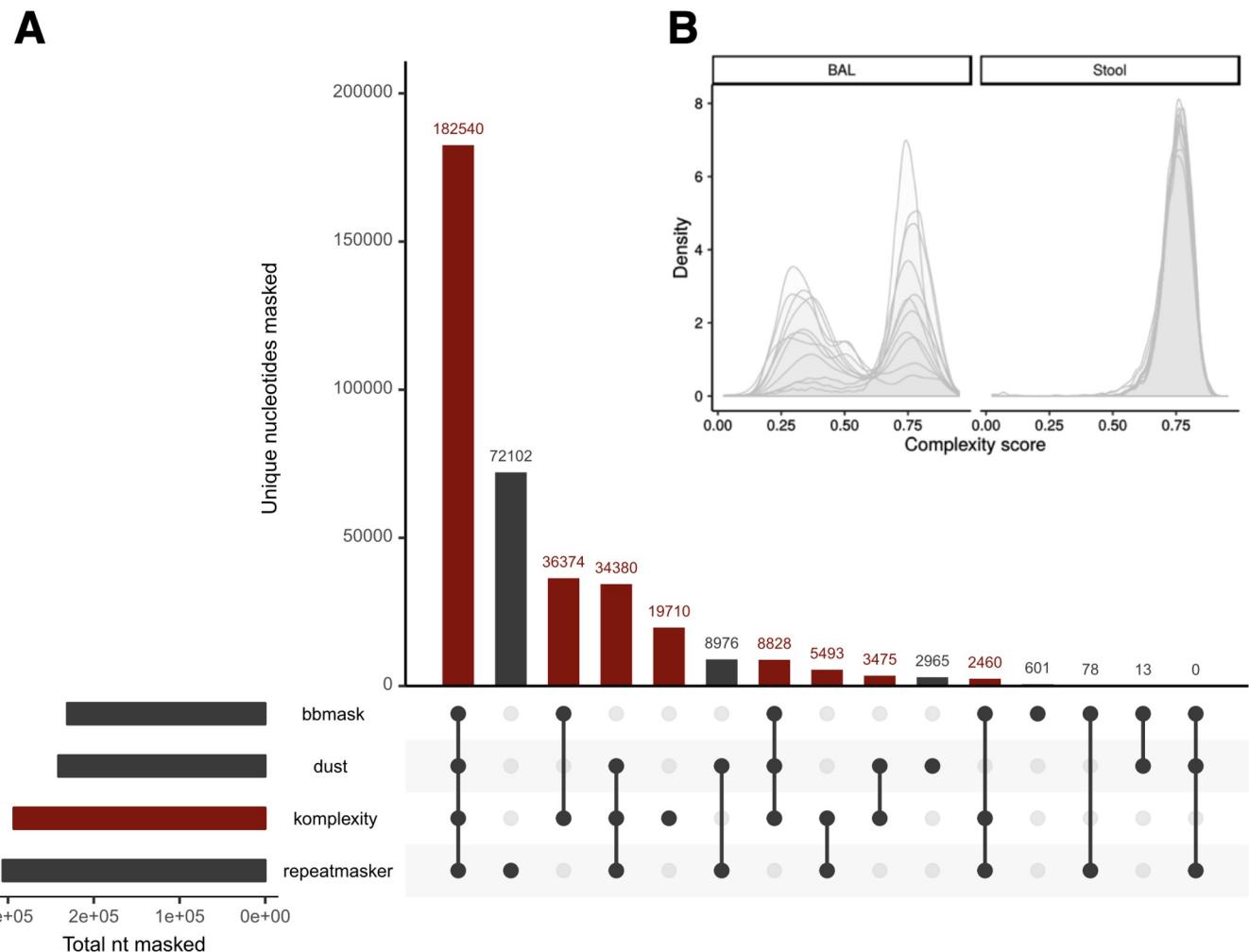
Total number of kilobases mapped to fungal genomes from NCBI.

Potential sources of misidentification evaluated by BLAST search to nt database followed by manual inspection.



Eukaryotic genomes contain **low-complexity DNA**, which is difficult to detect by sequence alignment.

DNA complexity can be quantified in order to remove low-complexity reads.



Clarke EL, Taylor LJ, Zhao C, Connell A, Lee JJ, Fett B, Bushman FD, Bittinger K. Sunbeam: an extensible pipeline for analyzing metagenomic sequencing experiments. *Microbiome*. 2019 Mar 22;7(1):46.

YES I'M
PARANOID



BUT AM I
PARANOID
ENOUGH?