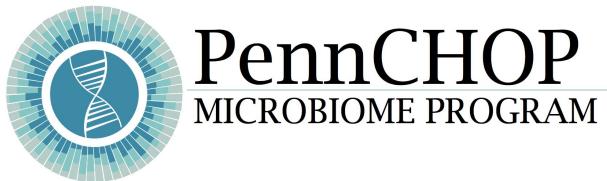


Fear and Loathing in Bacterial Genome Assembly

2023-03-15

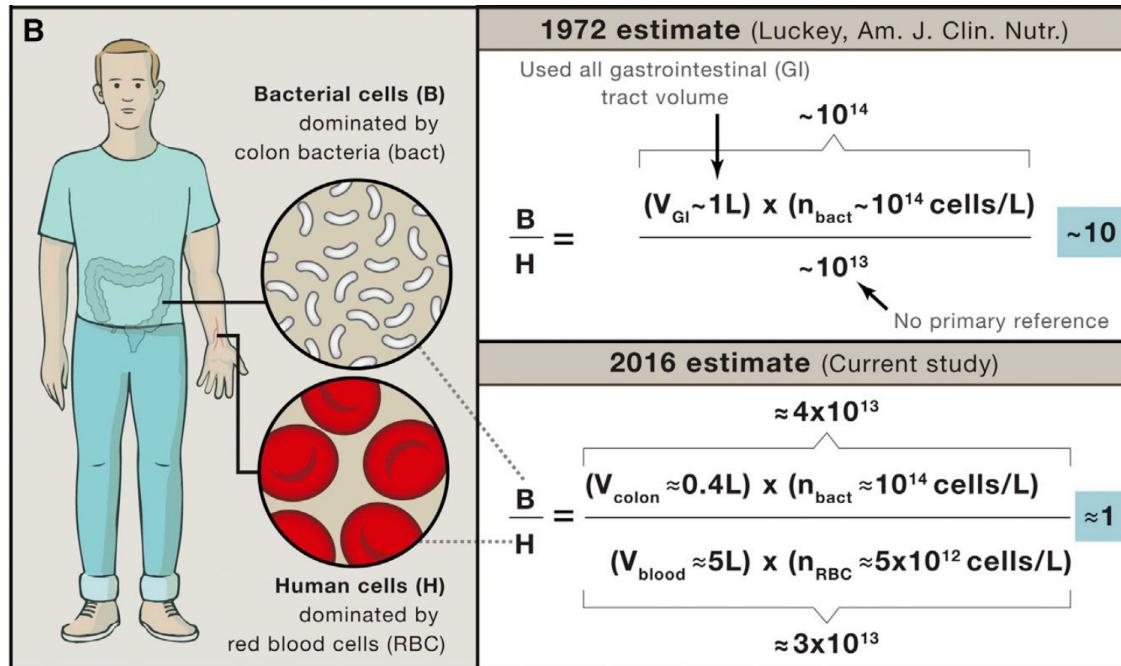


Kyle Bittinger

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

Illustration: Arwa Abbas

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

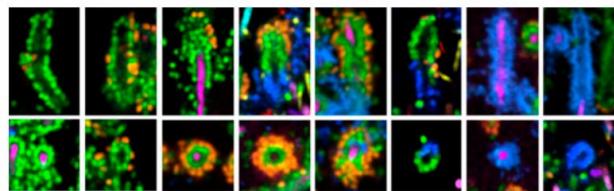
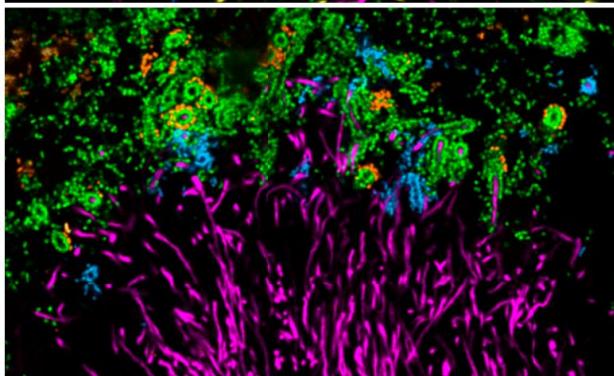
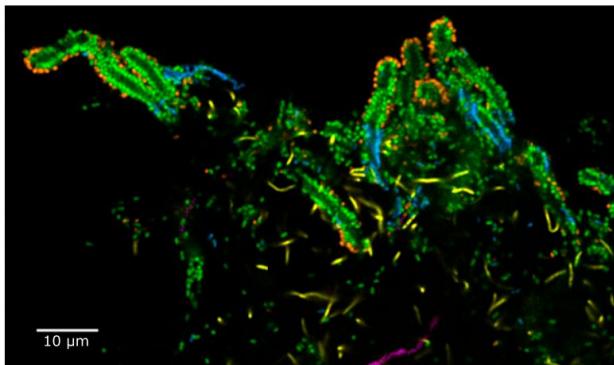


Sender R. Cell 164, 337 (2016).

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



Complex bacterial communities have spatial structure

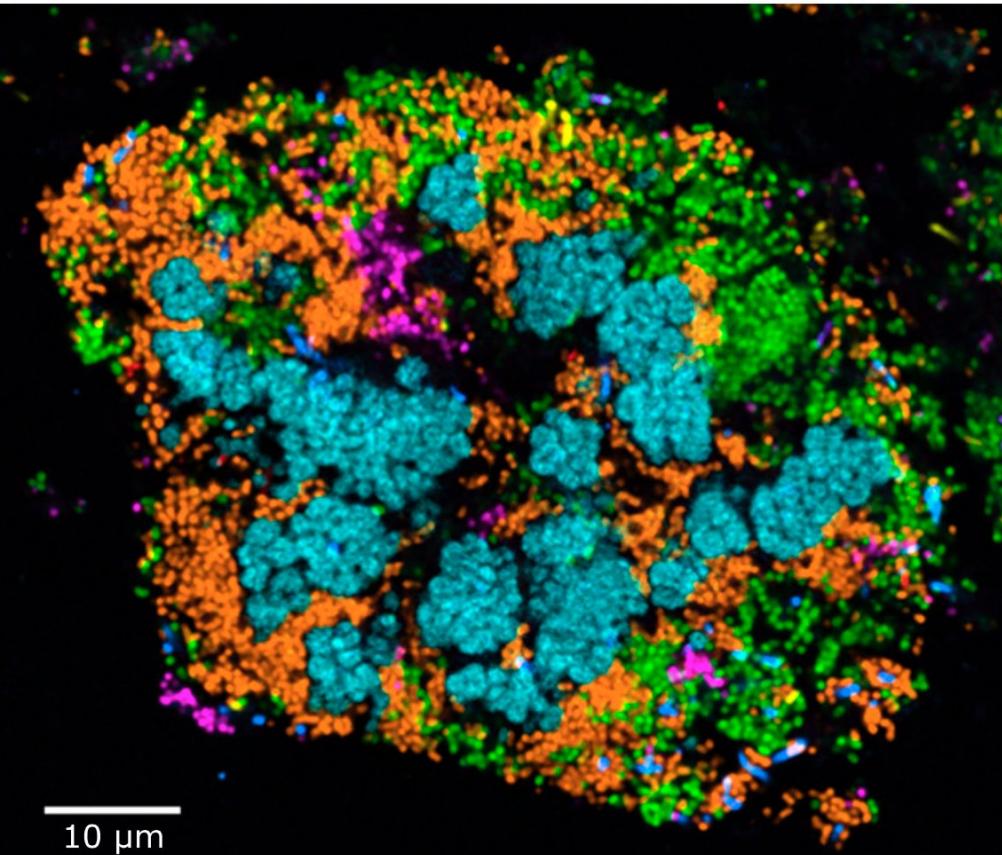
Example: oral bacteria in plaque

Base of "hedgehog structure" formed by *Corynebacterium*.

Streptococcus forms "corncob" near tip.

Haemophilus located in a third layer, growing on *Streptococcus*.

Mark Welch JL. *PNAS* E791 (2016).



Lautropia

Streptococcus

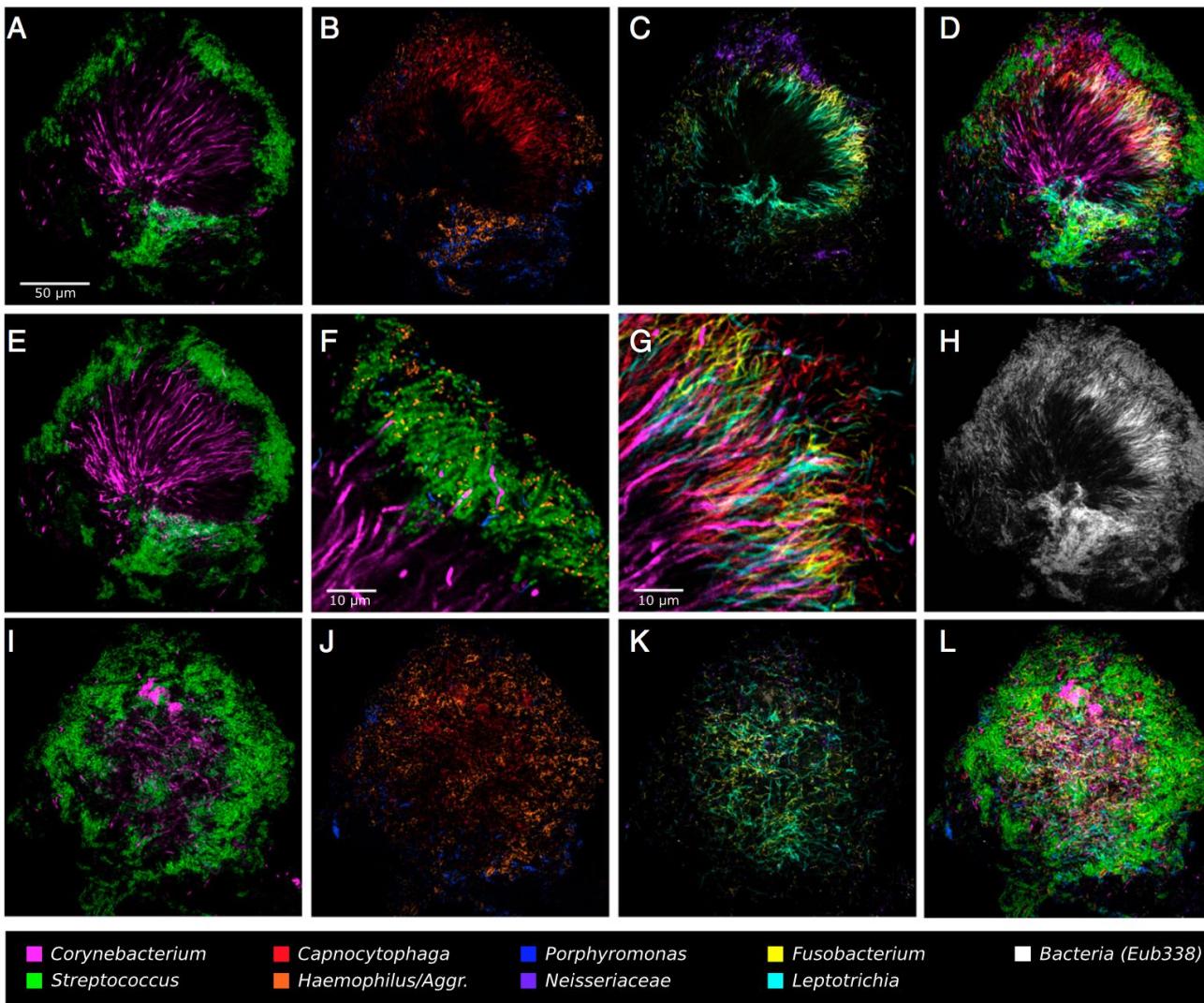
Veillonella

Haemophilus/Aggregatibacter

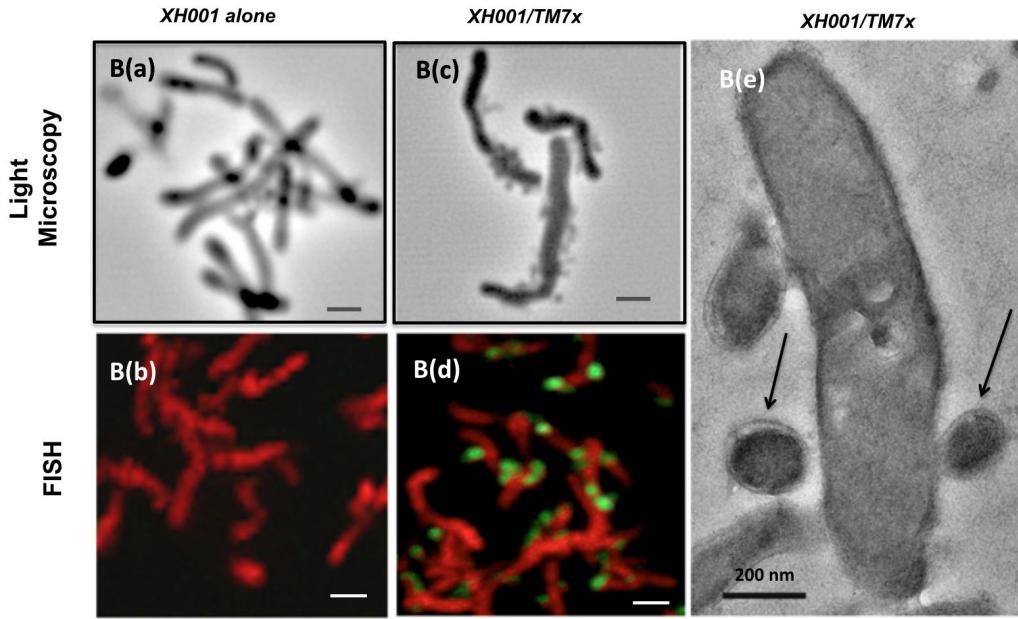
Prevotella

Rothia

Capnocytophaga



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



He PNAS 112, 244 (2015).

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

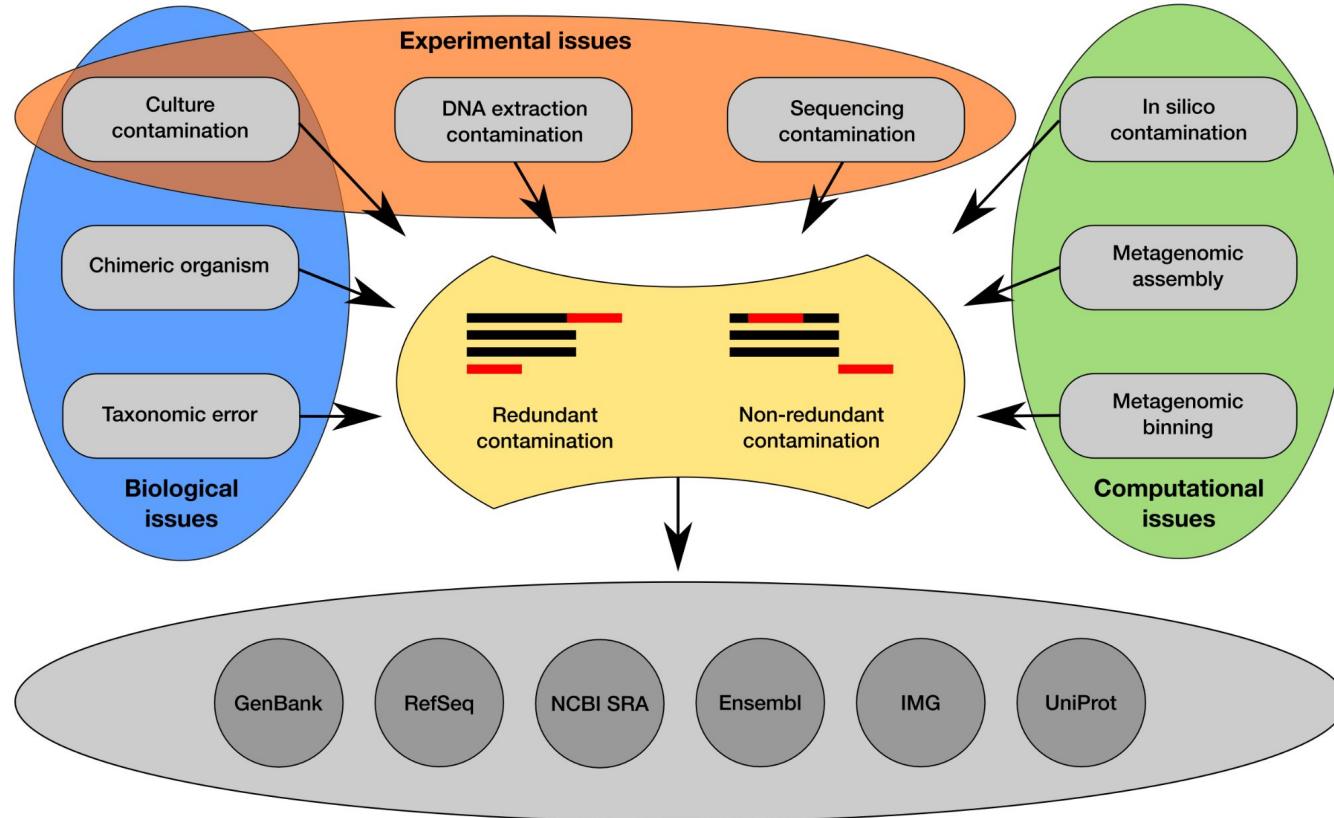
TM7 bacteria not cultured until 2014. (Soro V. *Appl Environ Microbiol.* 80, 6480)

TM7x, pictured, is an obligate epibiont of Actinomyces.

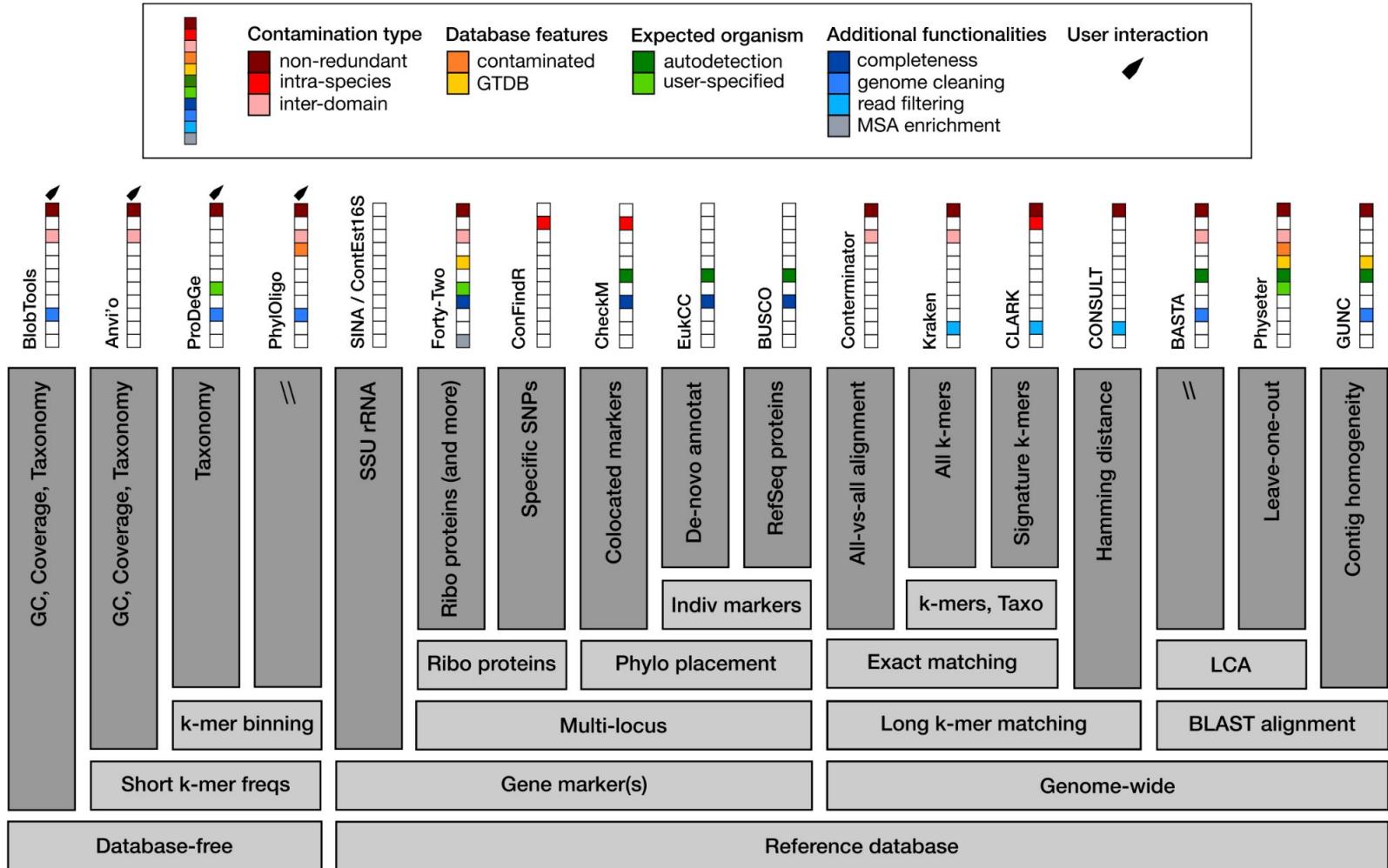
Complete genome: 705kb

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

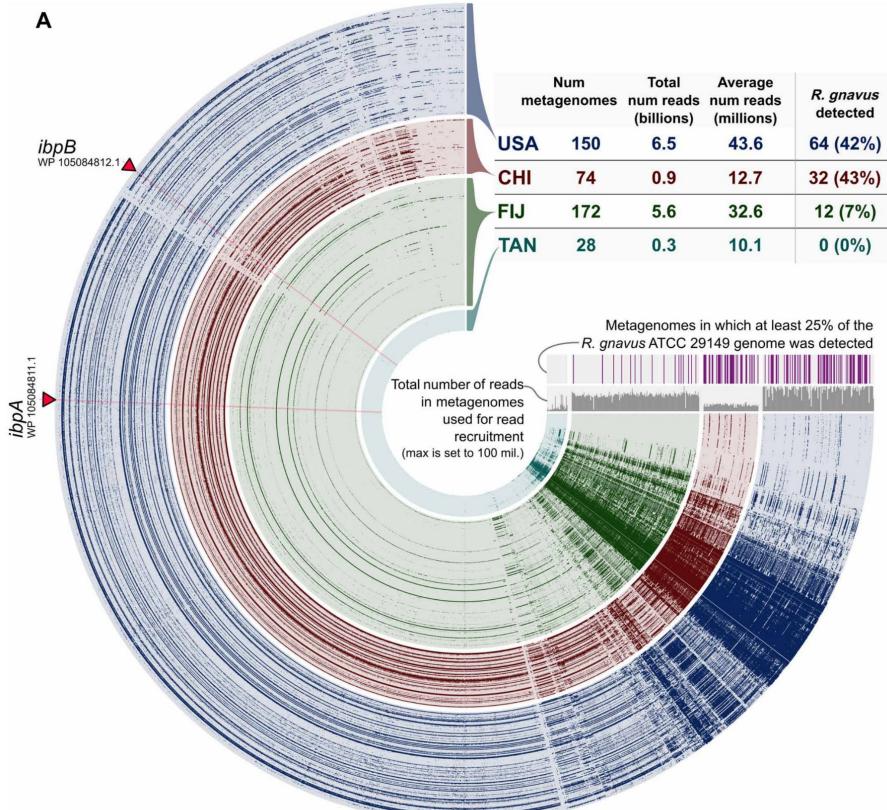
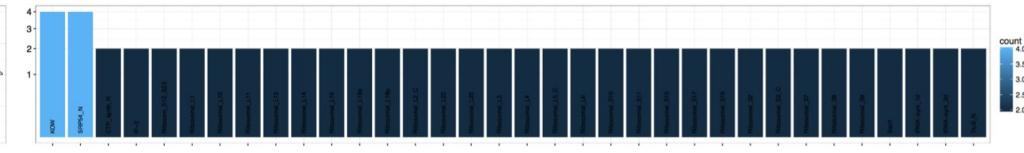
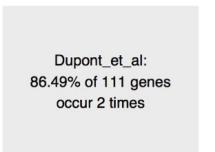
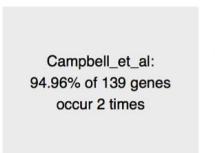
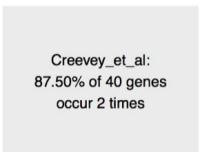
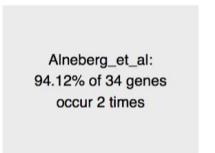
A

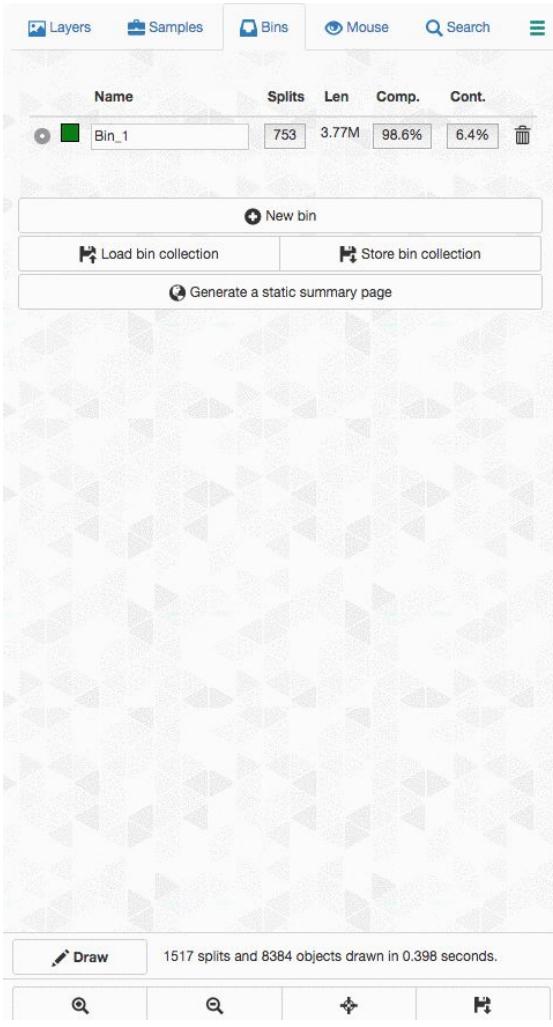
Fig. 5. Distribution of *R. gnavus* and its superantigens across human metagenomes. Dendrogram alignment of the *R. gnavus* ATCC 29149 genome to 424 human metagenomes (data files S3 and S4). Each spoke represents one gene in the *R. gnavus* genome, and each layer represents an individual human metagenome.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



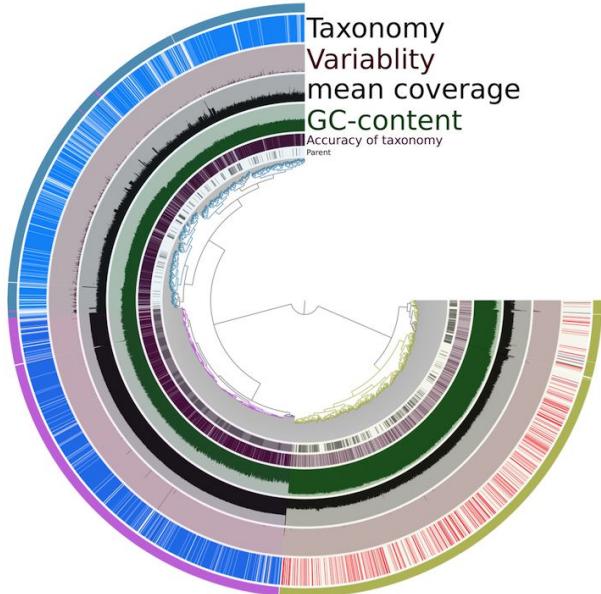
Eren AM, et al. Community-led, integrated, reproducible multi-omics with anvi'o. *Nat Microbiol*. 2021 Jan;6(1):3-6.

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

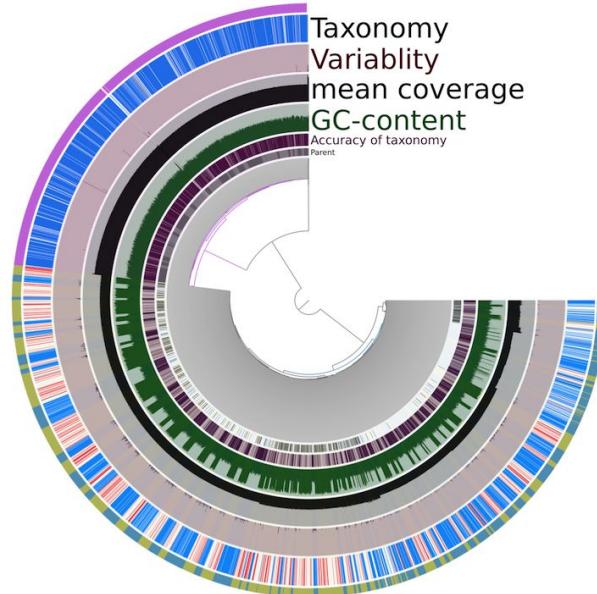




Clustering based on sequence composition



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%



Taxonomy Browser

Entrez

PubMed

Nucleotide

Protein

Genome

Structure

PMC

Taxonomy

BioCollections

Search for as lock

Display levels using filter:

[*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	1,503	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	-



Genome

Genome



txid33038[Organism:exp]



Search

Create alert Limits Advanced

Help



In June 2023, Genome record pages will be redirected to the new [Datasets Taxonomy page](#). [Learn more](#)

[*Ruminococcus*] *gnavus*

Representative genome: [*Ruminococcus*] *gnavus* ATCC 29149

Download sequences in FASTA format for [genome](#), [protein](#)

Download genome annotation in [GFF](#), [GenBank](#) or [tabular](#) format

BLAST against [*Ruminococcus*] *gnavus* [genome](#), [protein](#)

All 163 genomes for species:

Browse the [list](#)

Download sequence and annotation from [RefSeq](#) or [GenBank](#)

Display Settings: [▼](#) Overview

Send to: [▼](#)

[Organism Overview](#) ; [Genome Assembly and Annotation report \[163\]](#) ; [Genome Tree report \[107\]](#) ; [Plasmid Annotation Report](#)

ID: 979

[1]

[*Ruminococcus*] *gnavus*

Normal gastrointestinal bacterium

Lineage: Bacteria[35070]; Bacillota[5812]; Clostridia[2561]; Eubacteriales[2323]; Lachnospiraceae[811]; *Mediterraneibacter*[40]; [*Ruminococcus*] *gnavus*[1]

Ruminococcus gnavus. *Ruminococcus gnavus* represents 0.085% of the organisms identified from the human gut. This organism has also been shown to produce an antibacterial compound (Ruminococcin A) which may play a role in its colonization and persistence in the gut.

Summary

Sequence data: genome assemblies: 163 (See [Genome Assembly and Annotation report](#))

Statistics: median total length (Mb): 3.40826

median protein count: 3228

median GC%: 42.7

Tools

[BLAST Genome](#)

Related information

Assembly

BioProject

Gene

Components

Protein

PubMed

Taxonomy

Search details

txid33038[Organism:exp]

Search

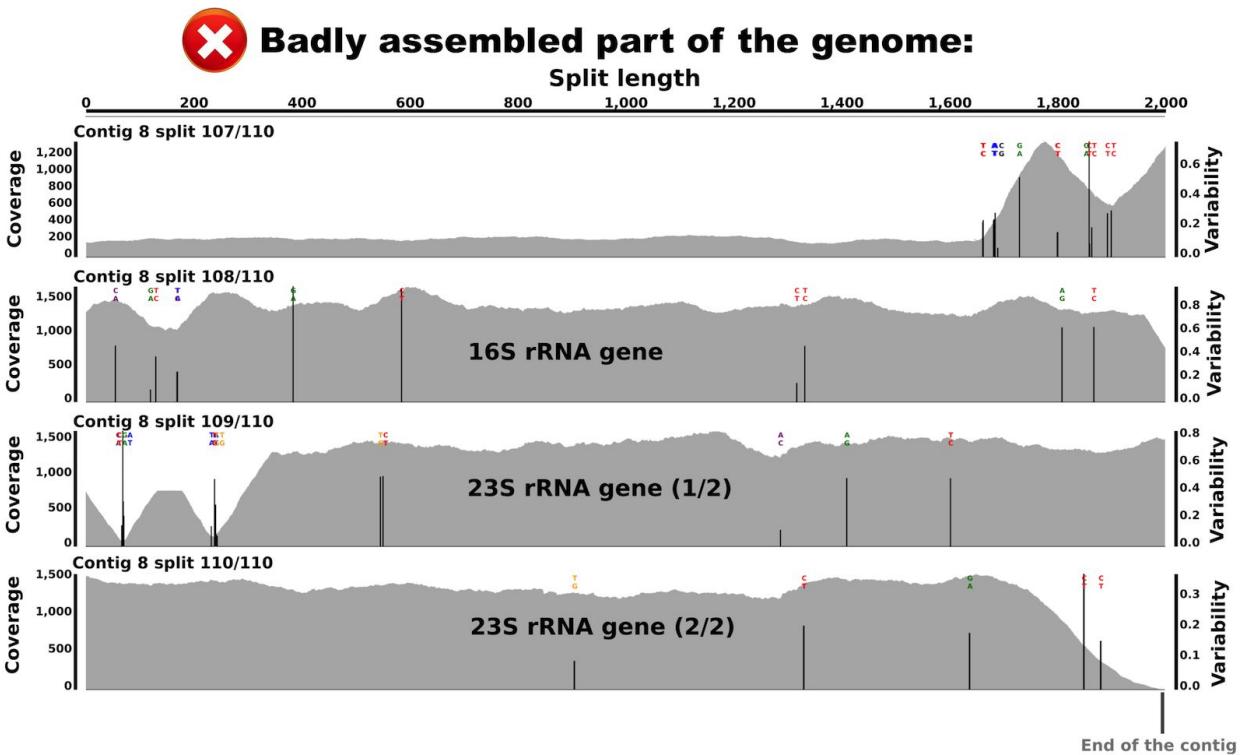
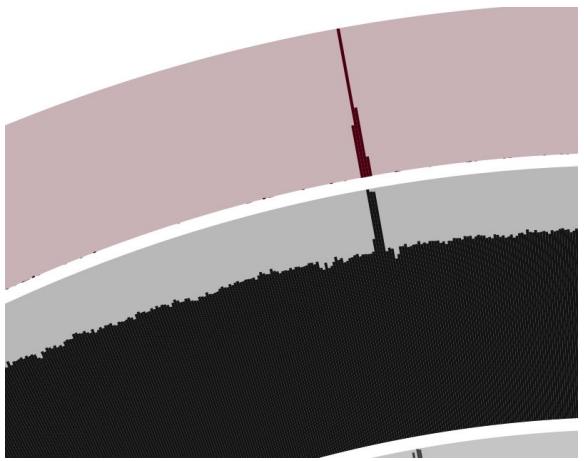
See more...

Sequence data: genome assemblies: 163 (See [Genome Assembly and Annotation report](#))

Statistics: median total length (Mb): 3.40826

median protein count: 3228

median GC%: 42.7

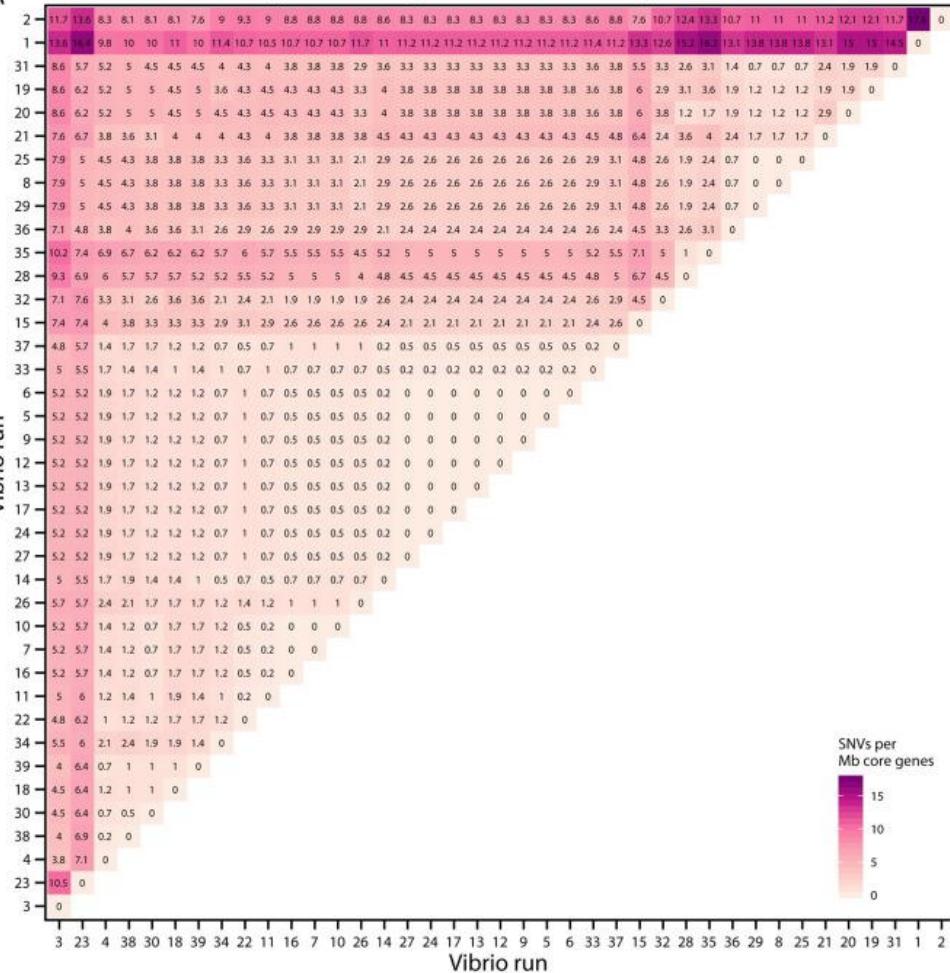


3

**Technical replicates and
strain-level differences**

A

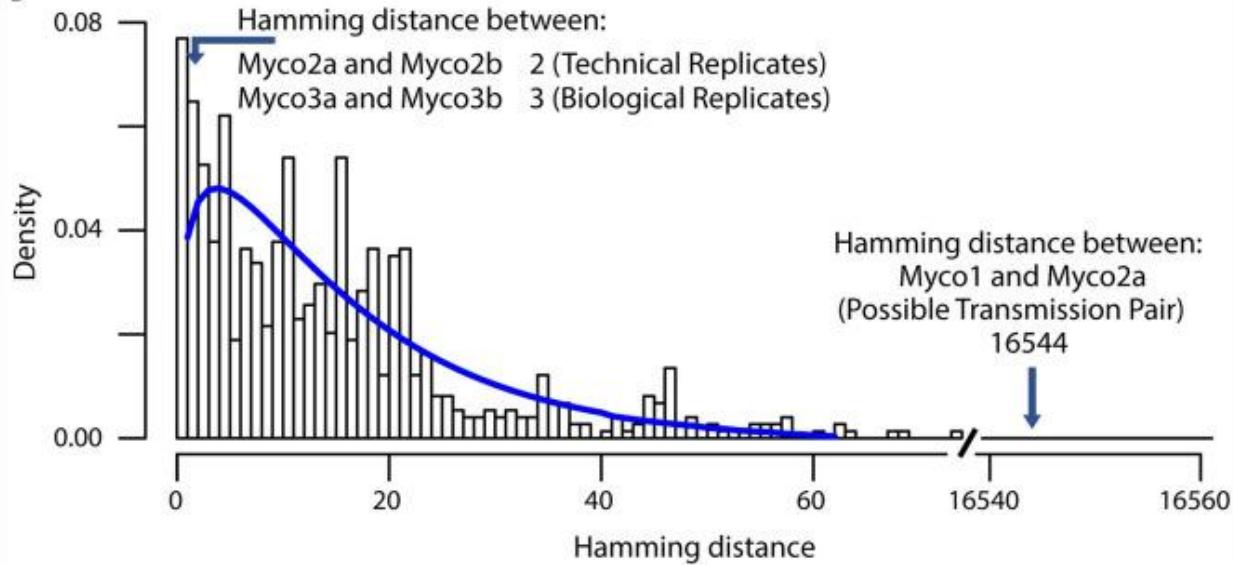
SNVs per Mb core genes for *Vibrio* core genomes



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.

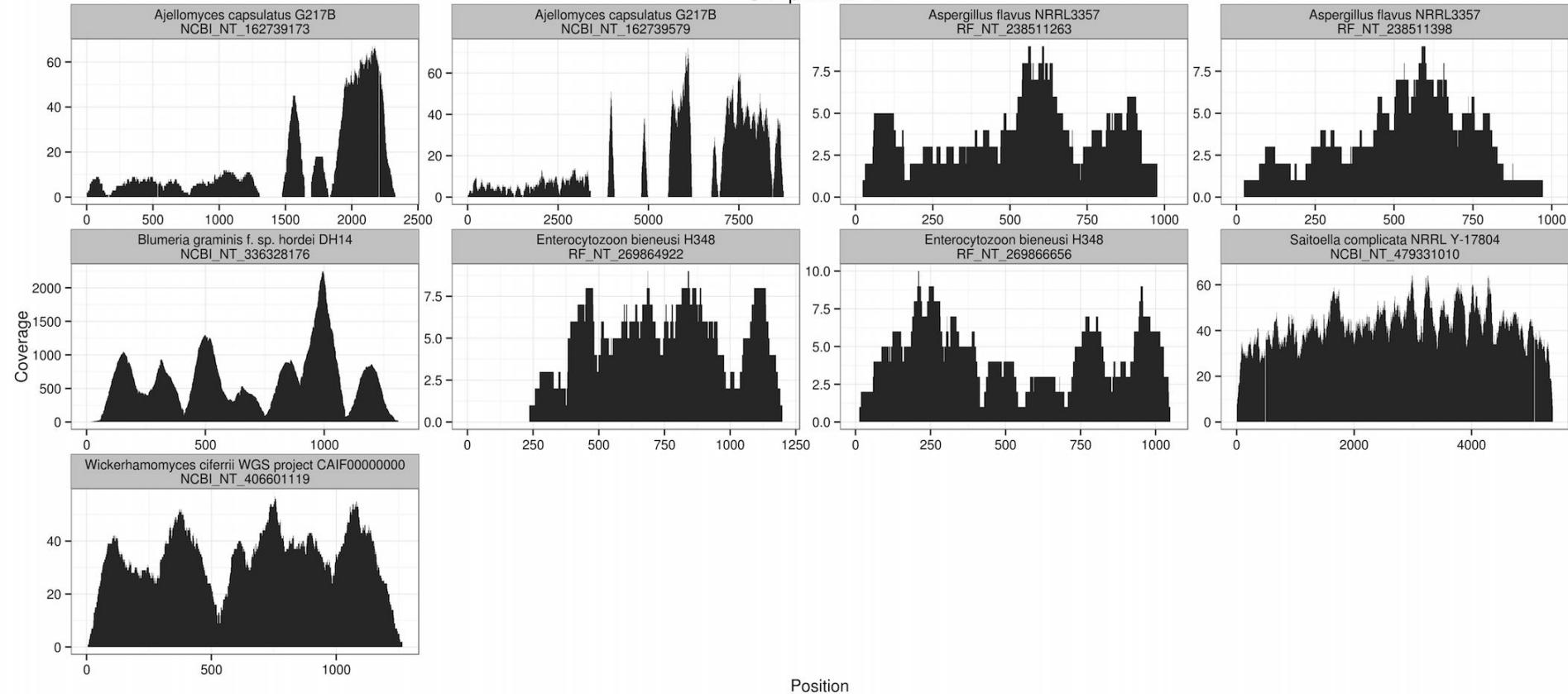
B



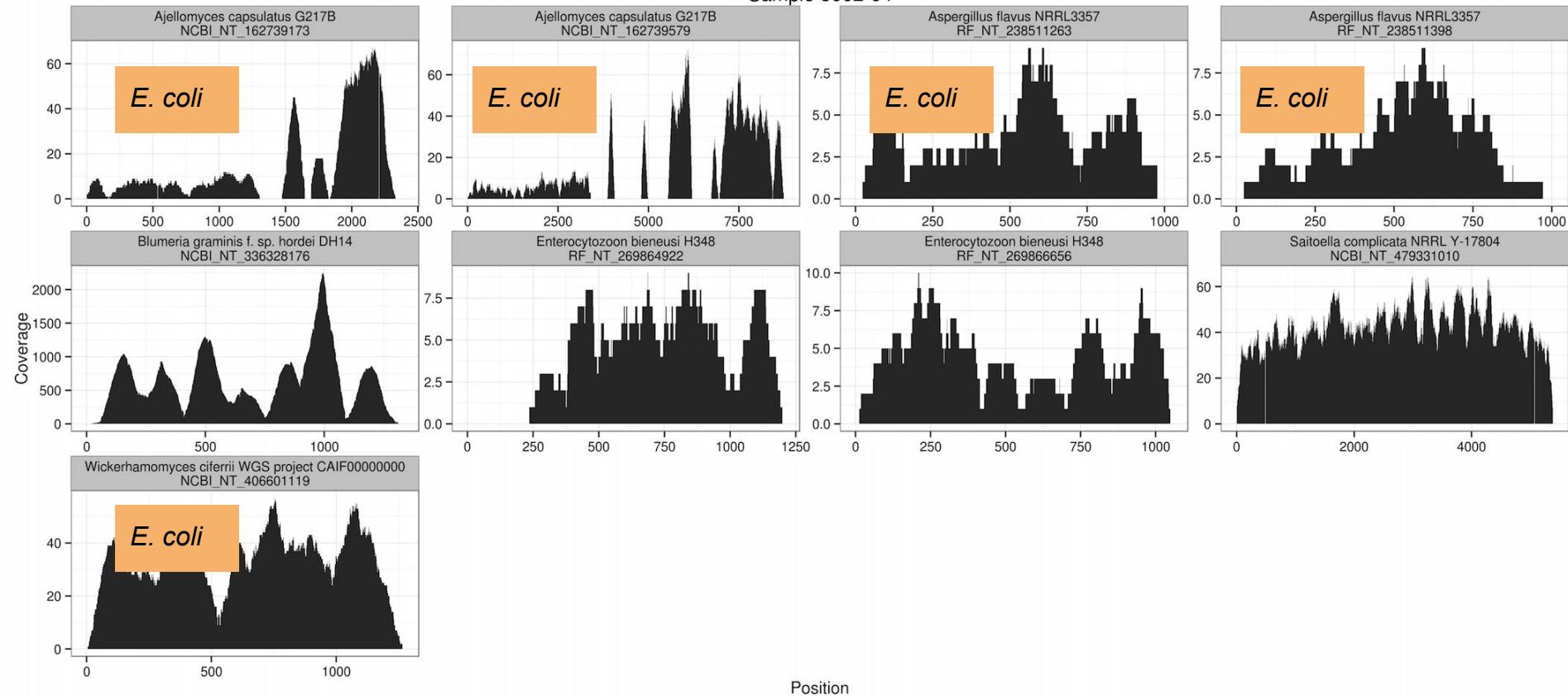
4

**Database contamination,
human DNA contamination**

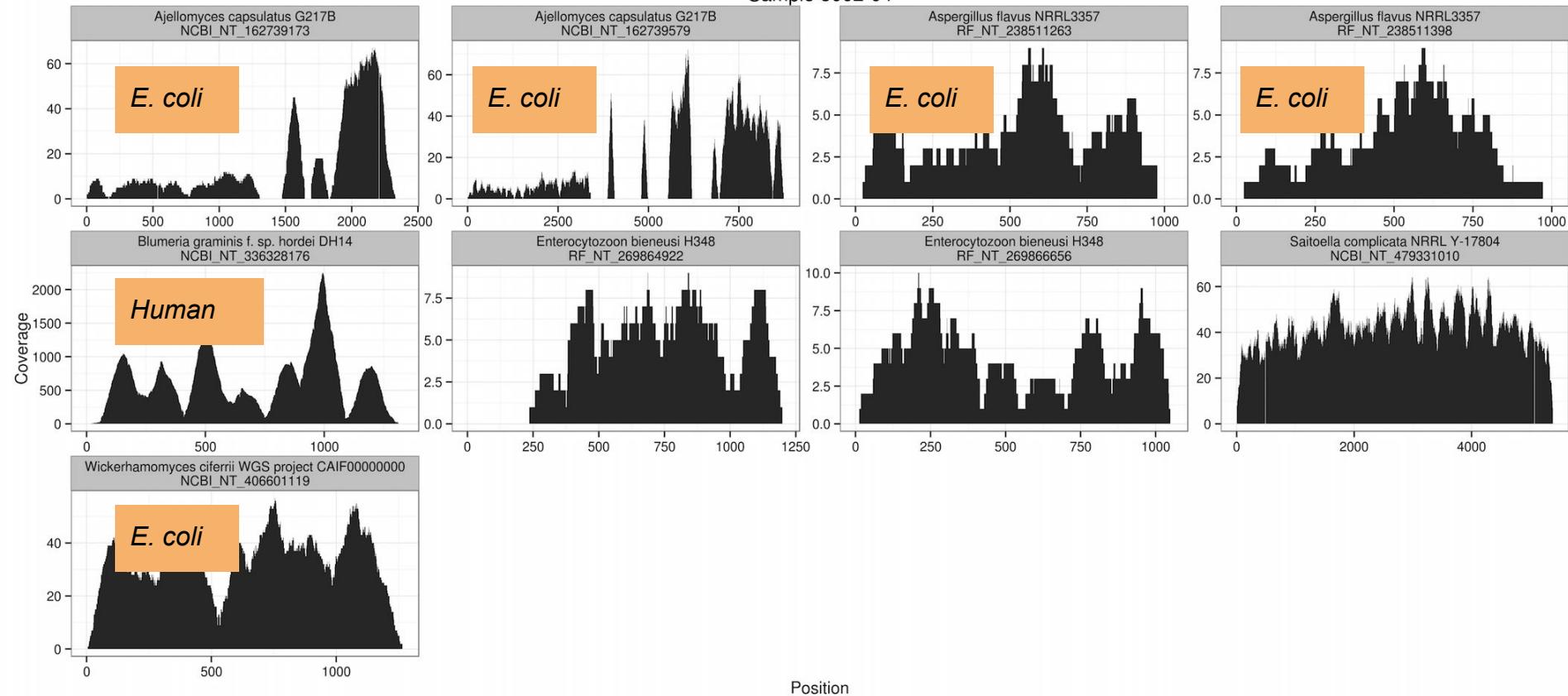
Sample 5002-04



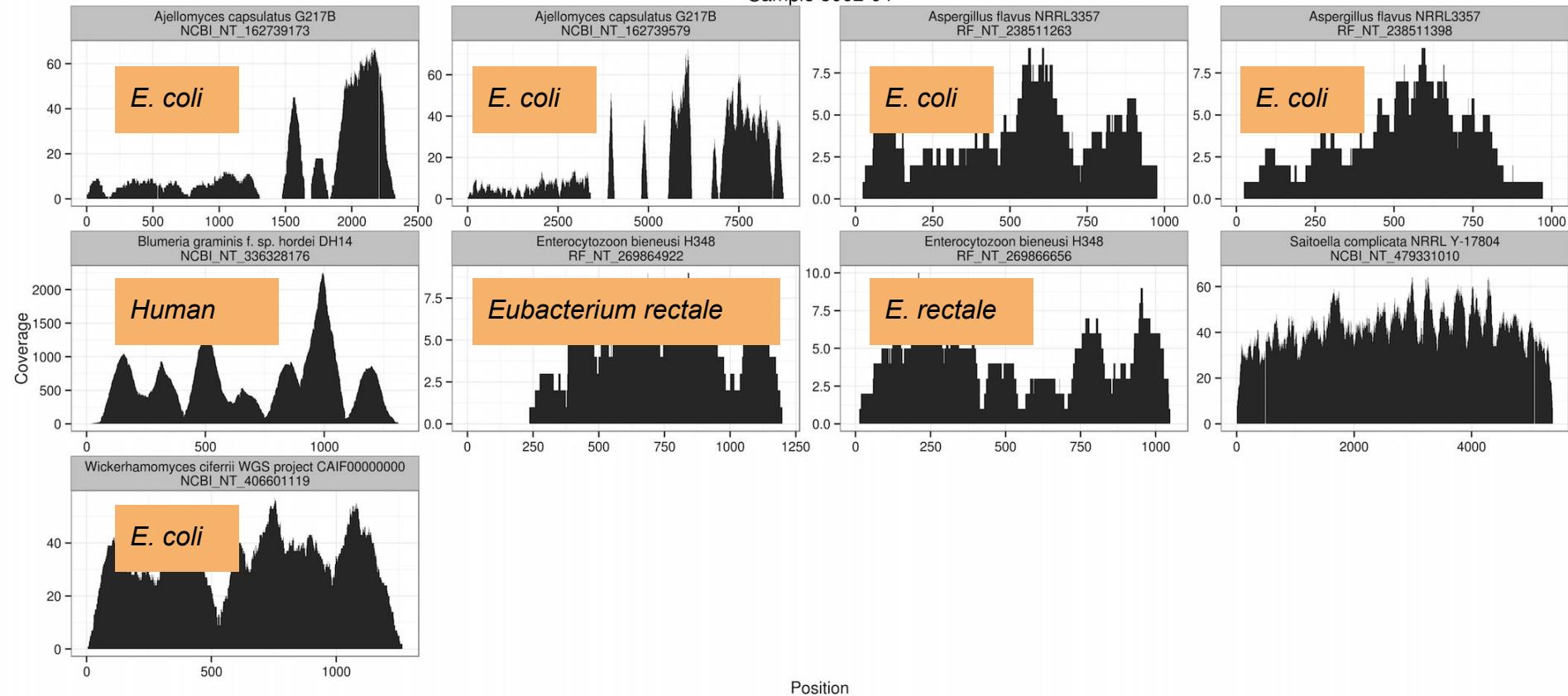
Sample 5002-04



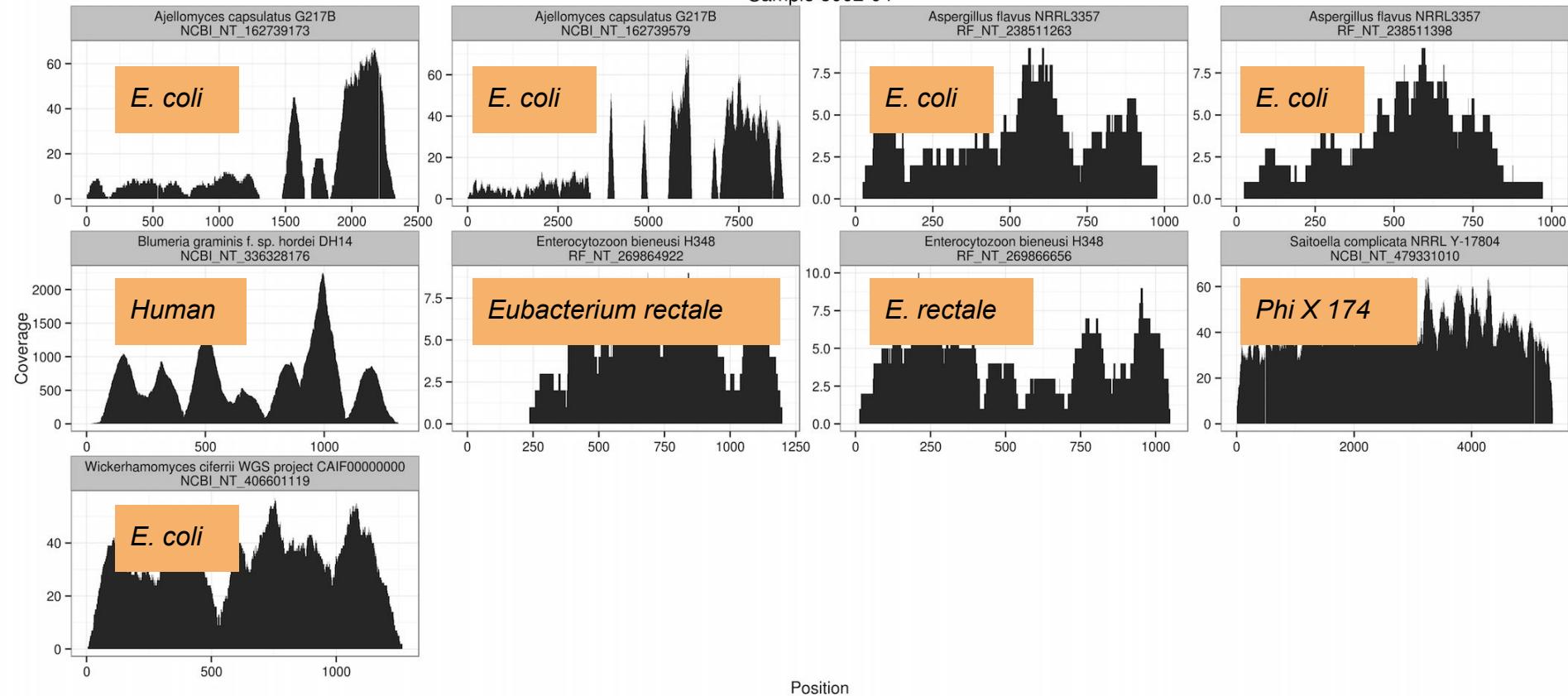
Sample 5002-04



Sample 5002-04



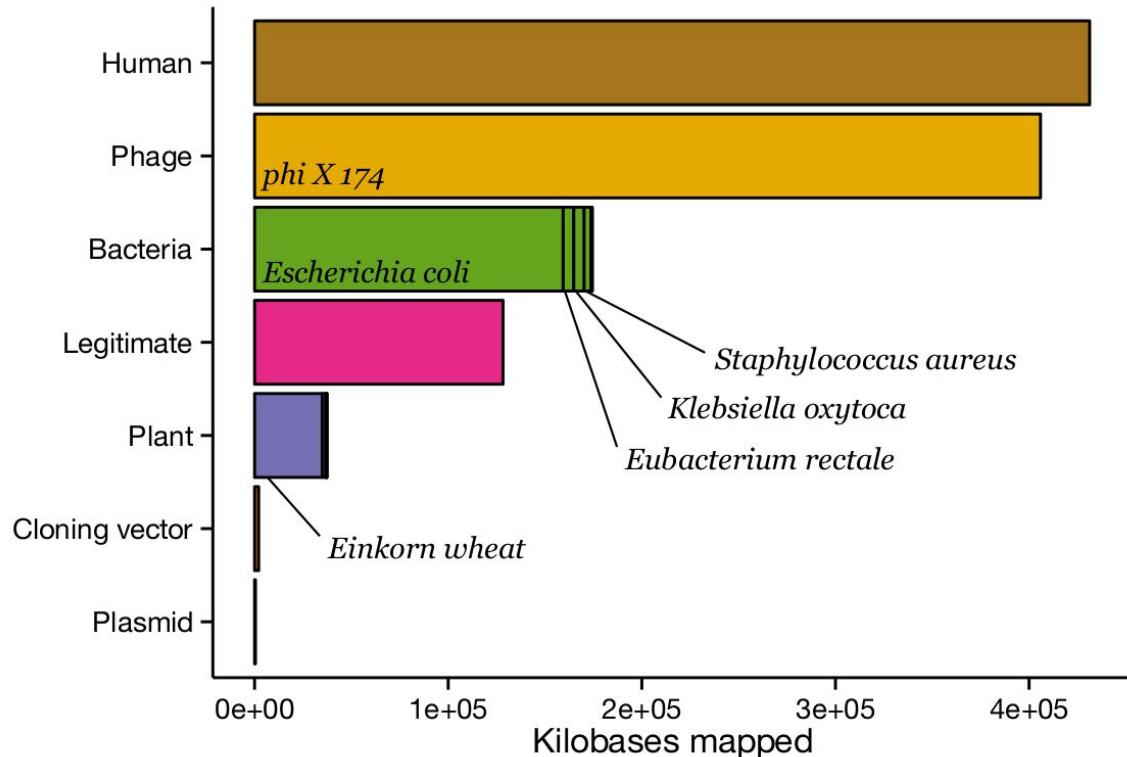
Sample 5002-04



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

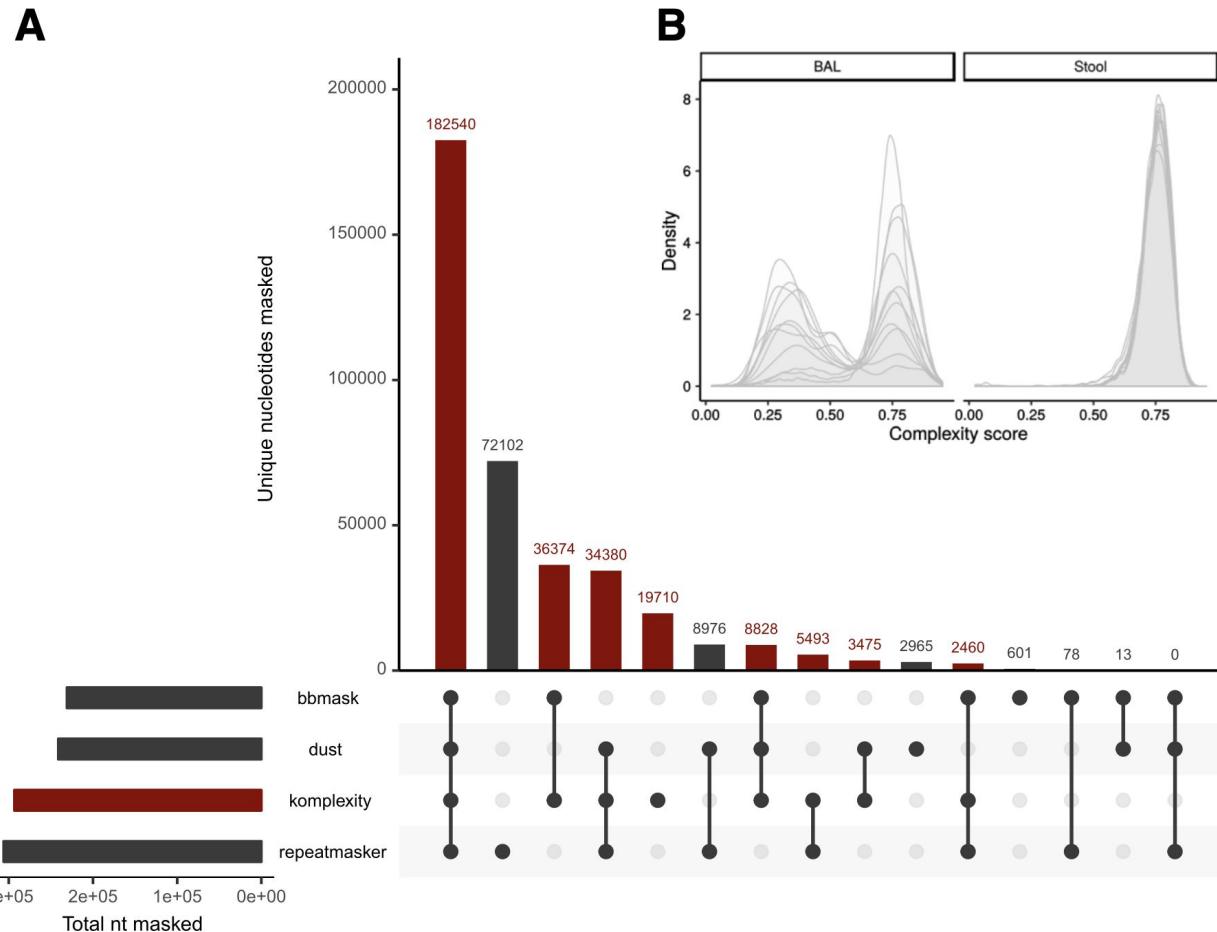
Total number of kilobases mapped to fungal genomes from NCBI.

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.



YES I'M
PARANOID



BUT AM I
PARANOID
ENOUGH?