

# From Detection to Amplification: Metagenomics Perspectives on Antimicrobial Resistance

Christina Boucher

Dept of Computer and Information Science and Engineering

Herbert Wertheim College of Engineering

[www.christinaboucher.com](http://www.christinaboucher.com)



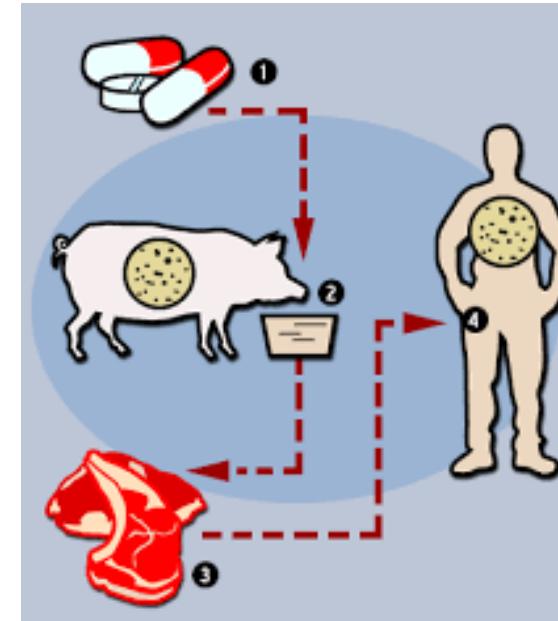


# THE TROUBLE WITH CHICKEN

MAY 12, 2015 // 54:41



**Pigs and humans are sharing antibiotic-resistance genes**

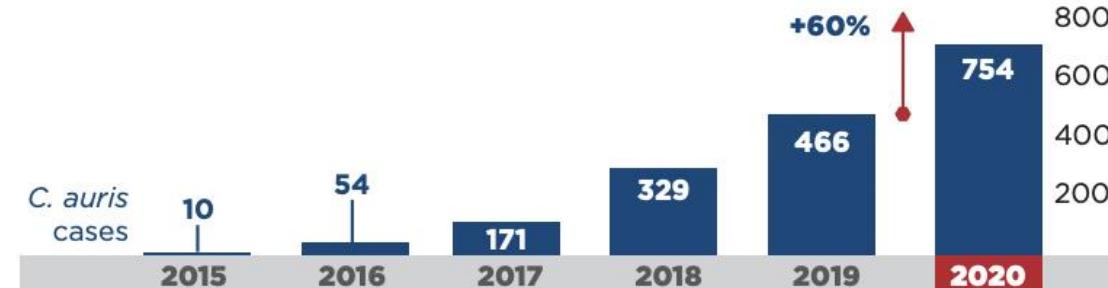
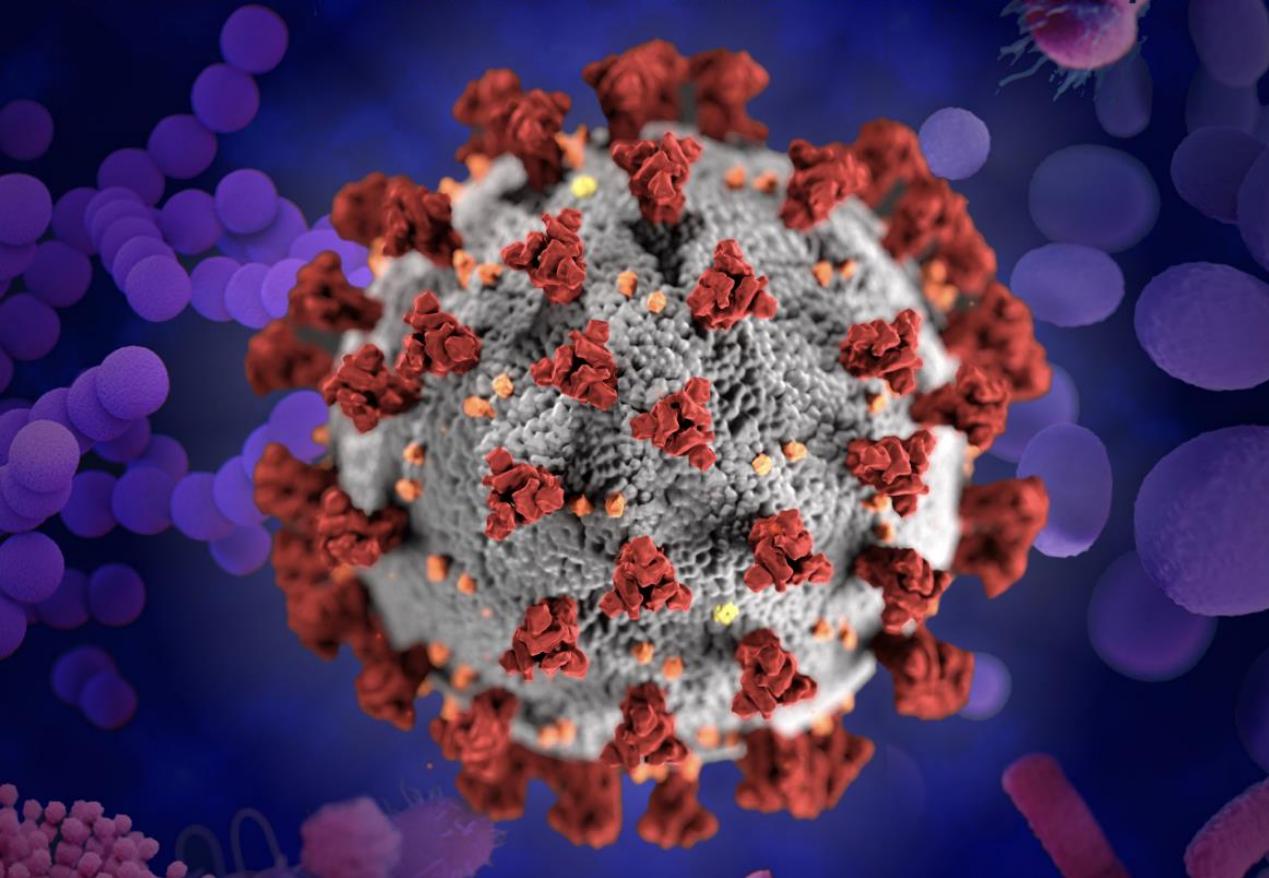


**EU figures show high levels of antibiotic resistance in foodborne bacteria**

By Joe Whitworth on April 6, 2022

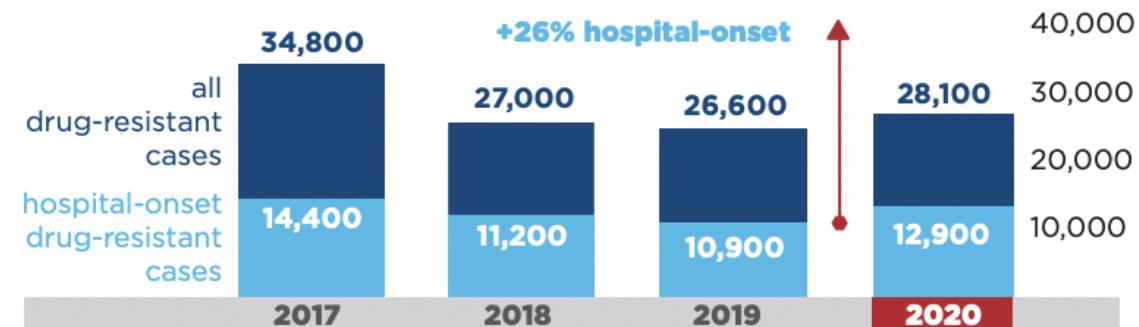
# COVID-19

## U.S. IMPACT ON ANTIMICROBIAL RESISTANCE

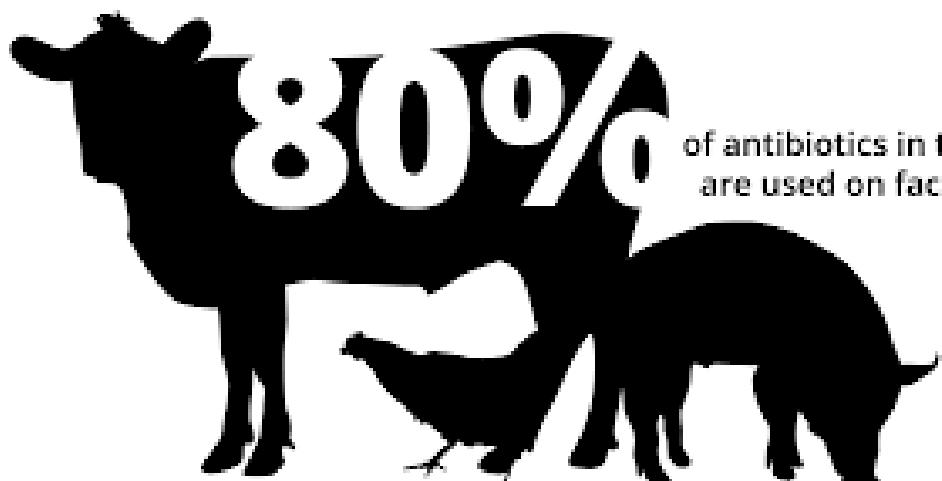
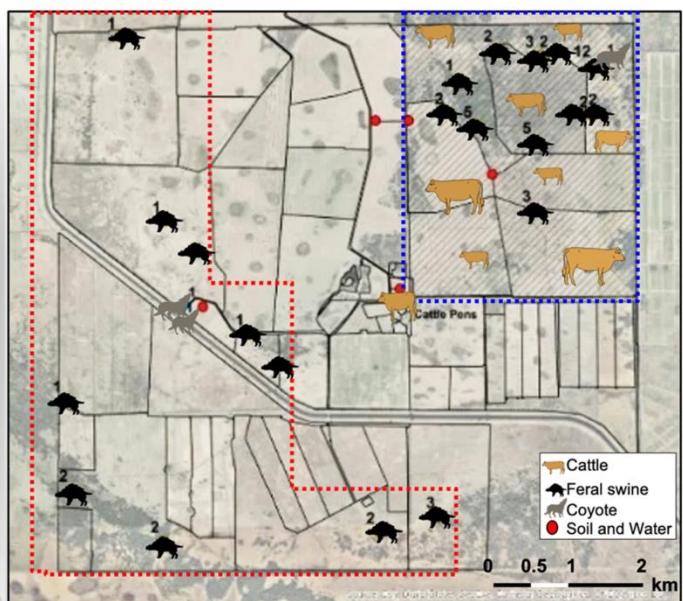
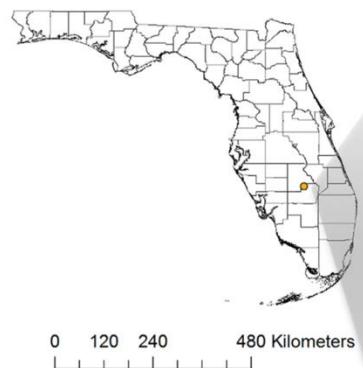


COVID-19: U.S. Impact on Antimicrobial Resistance, Special Report 2022

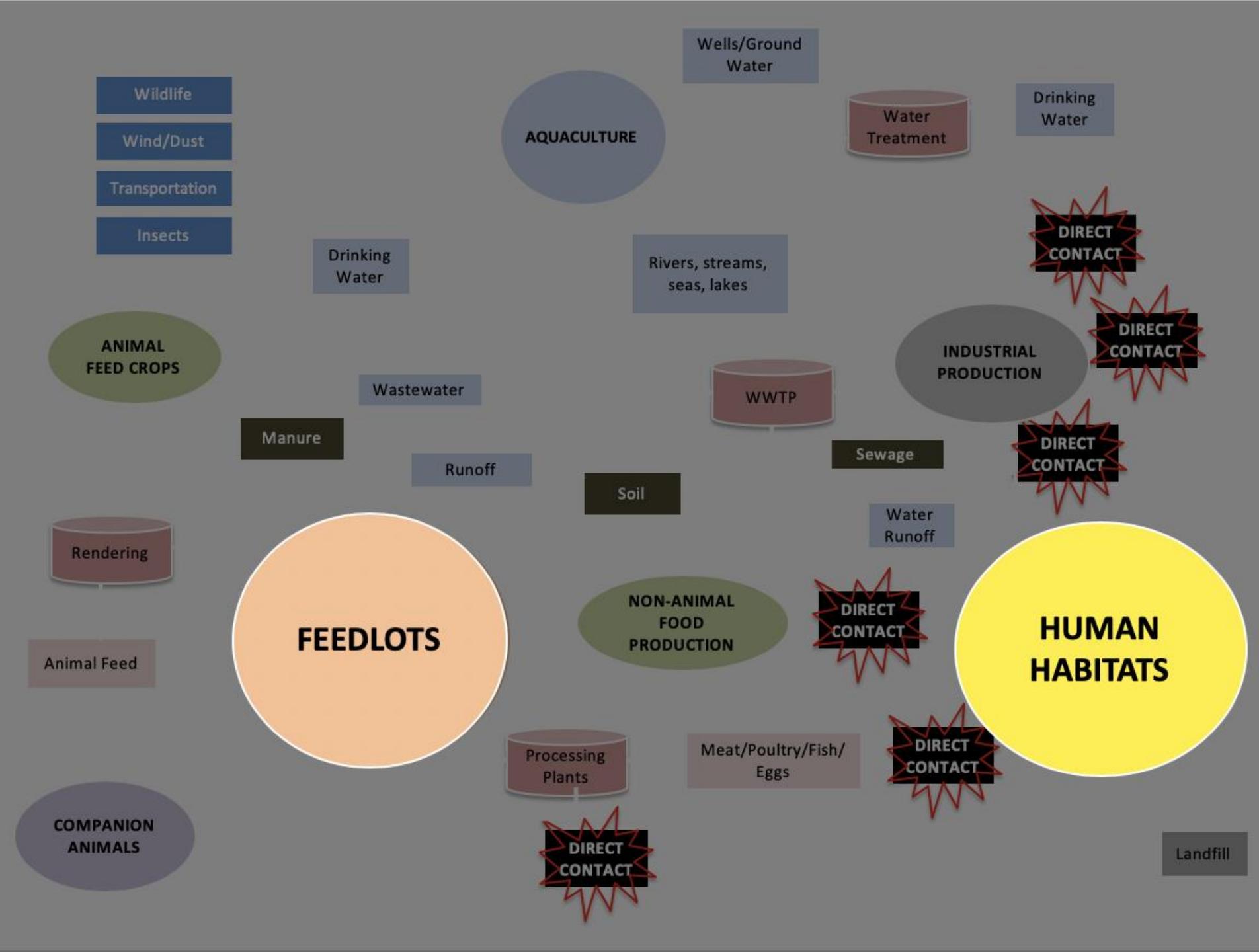
Drug-resistant *Candida* cases decreased until 2020, when hospital-onset cases increased during the COVID-19 pandemic.

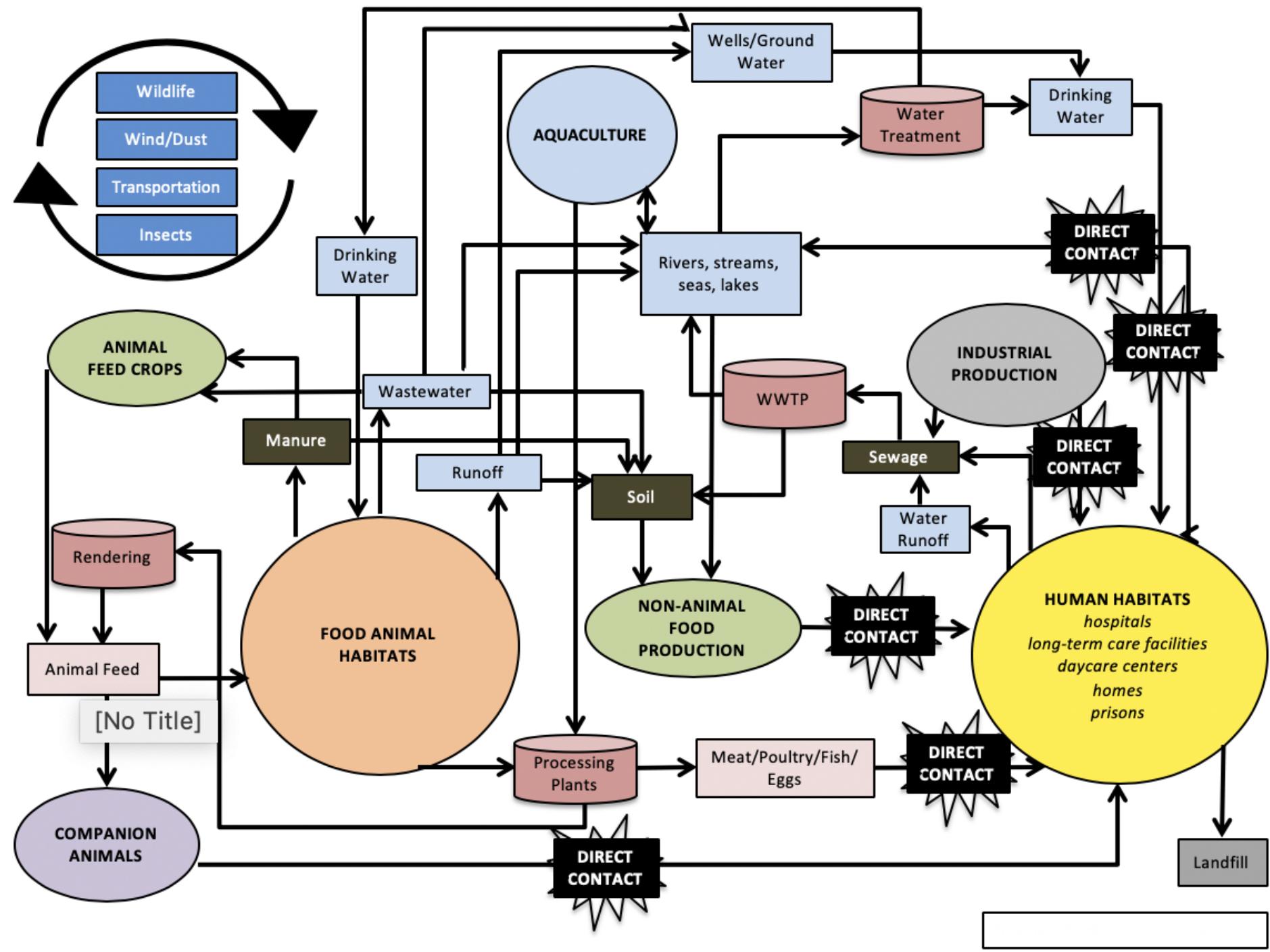


# AR Pathogens Cause Infections Across the One Health Spectrum



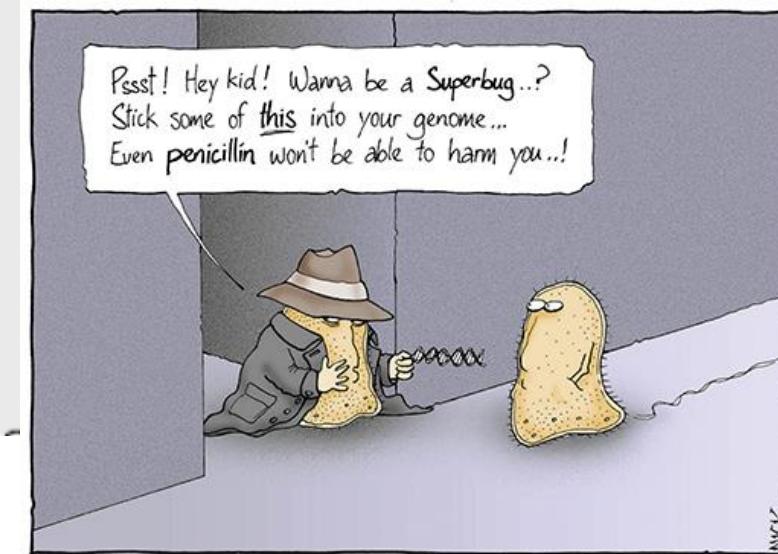
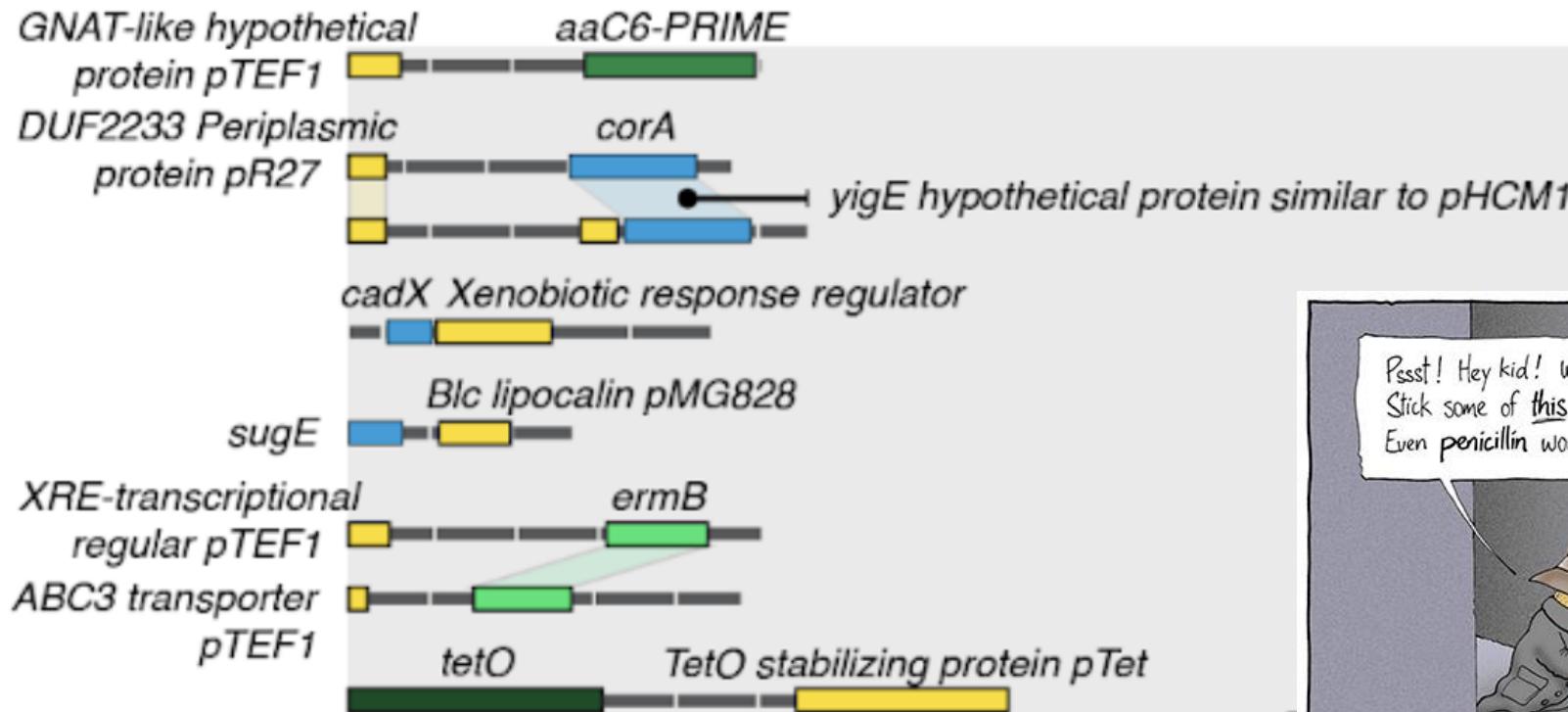
[Transmission of antibiotic resistance at the wildlife-livestock interface](#) S Lee et al.  
Communications biology 5 (1), 1-12





# Detection of Resistance

# Antimicrobial Resistant Genes (ARGs)



It was on a short-cut through the hospital kitchens that Albert was first approached by a member of the Antibiotic Resistance.

# Culture Based Methods

## Old Method



"indicator"  
bacterial  
species

# Microbiome is a Complex System

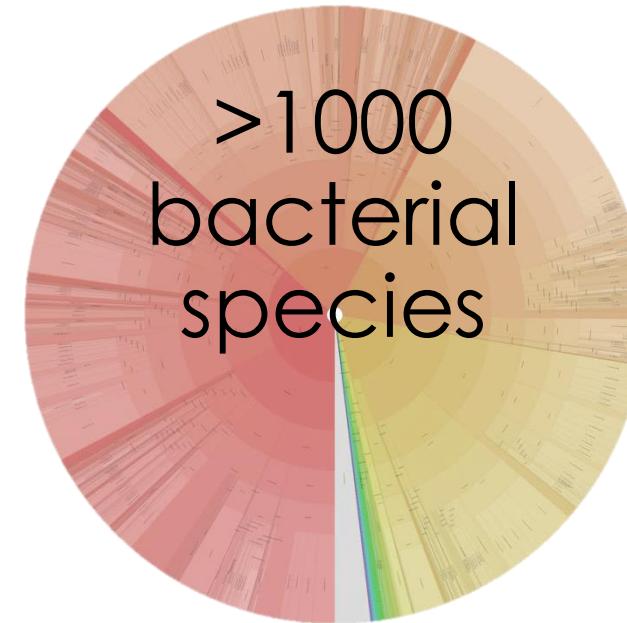
**Old Method**

vs.

**Reality**



“indicator”  
bacterial  
species



>1000  
bacterial  
species

# Microbiome is a Complex System

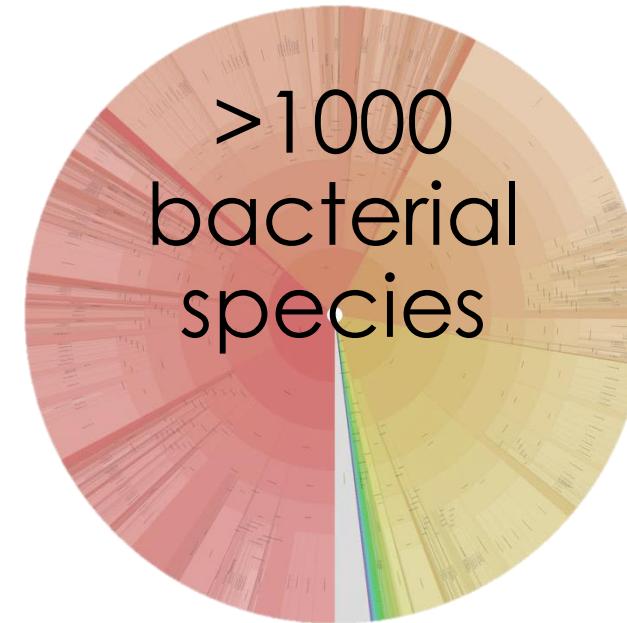
**Old Method**

vs.

**Reality**

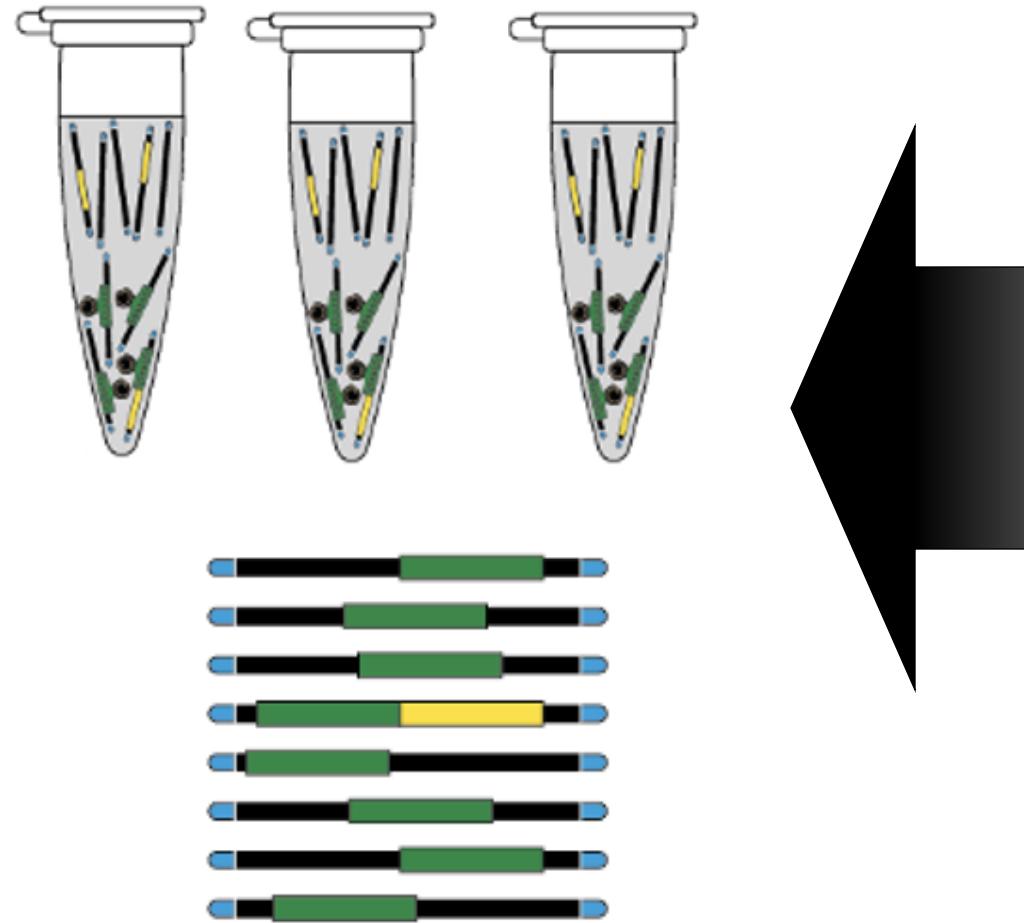


“indicator”  
bacterial  
species

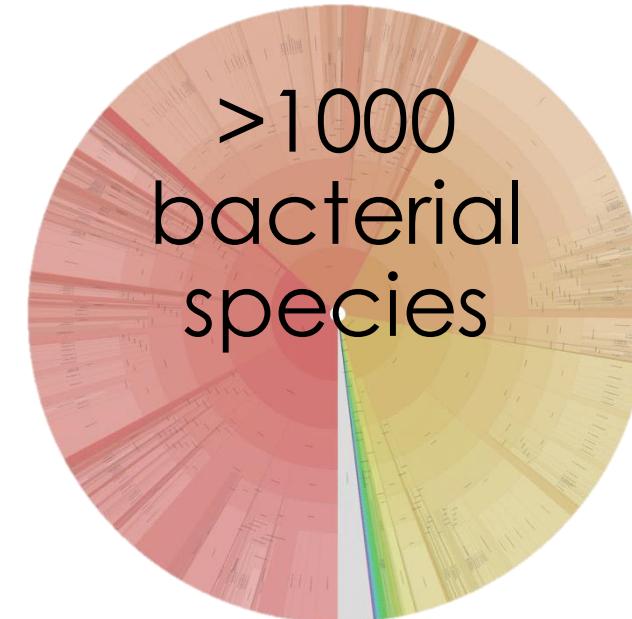


>1000  
bacterial  
species

# Microbiome is a Complex System



**Reality**



# Introduction of the Resistome

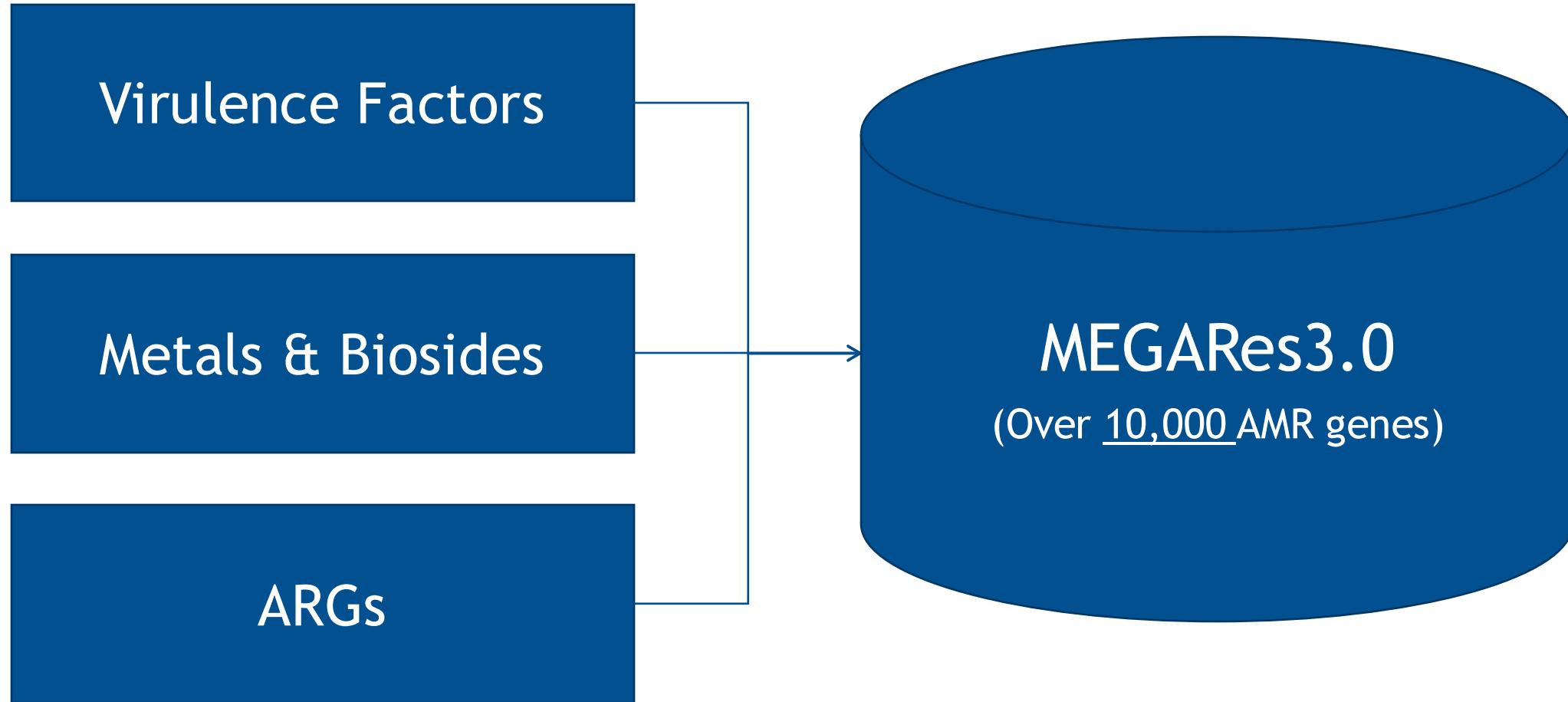
## Resistome:

All ARGs in both pathogenic and non-pathogenic bacteria

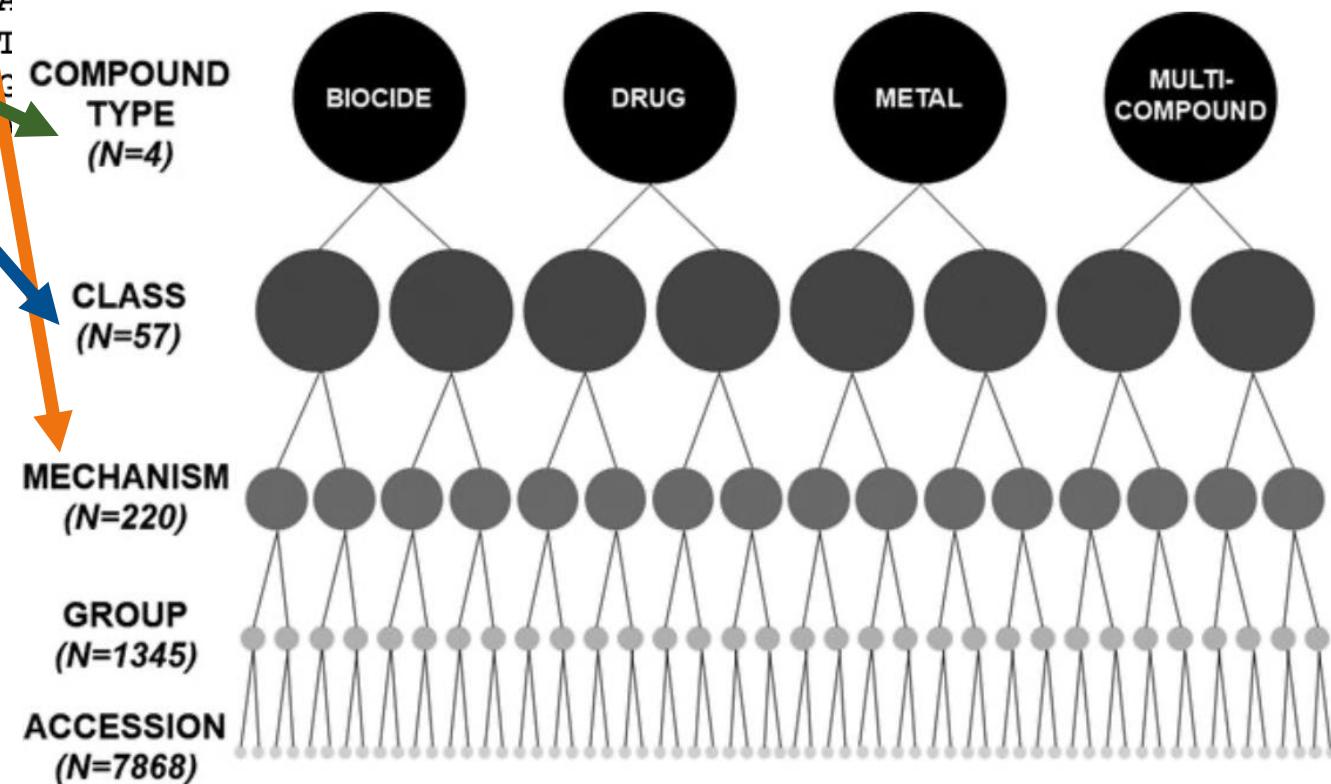
>1000  
bacterial species



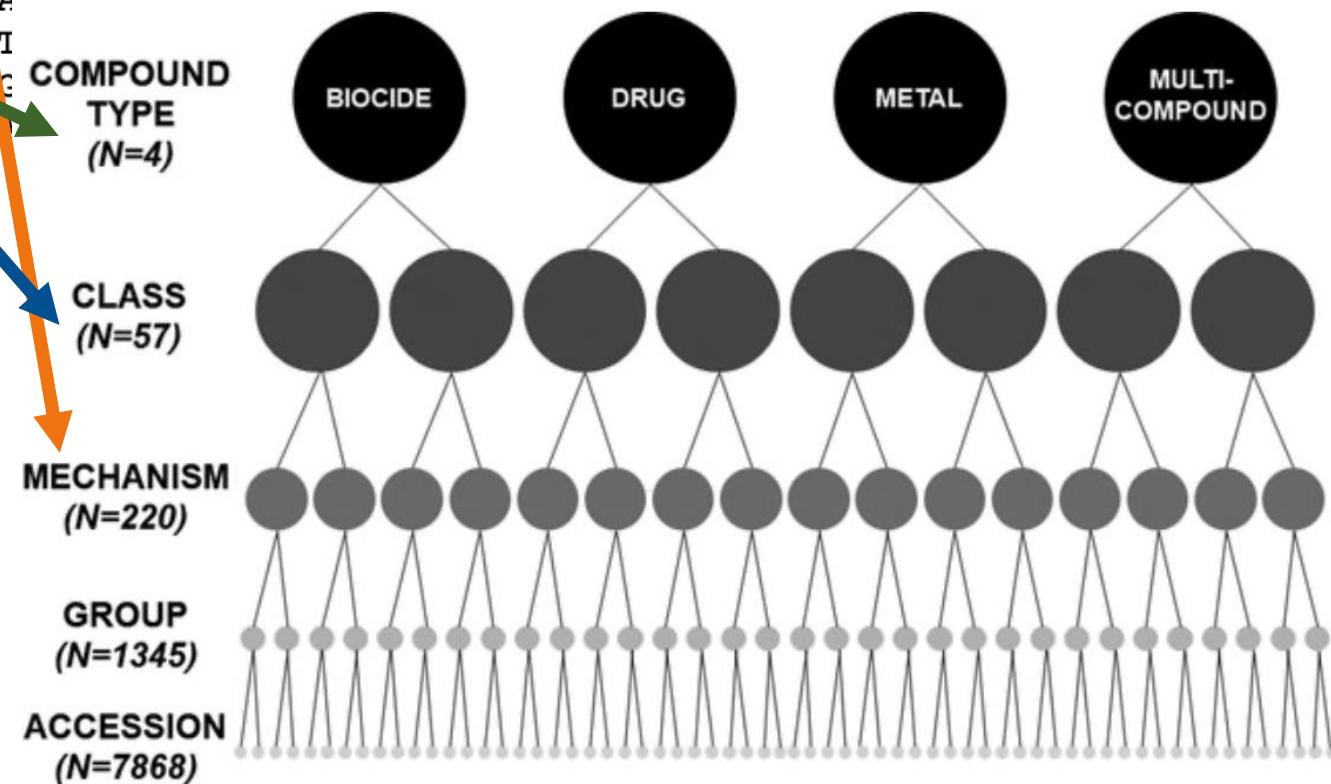
# Introduction of MEGARes



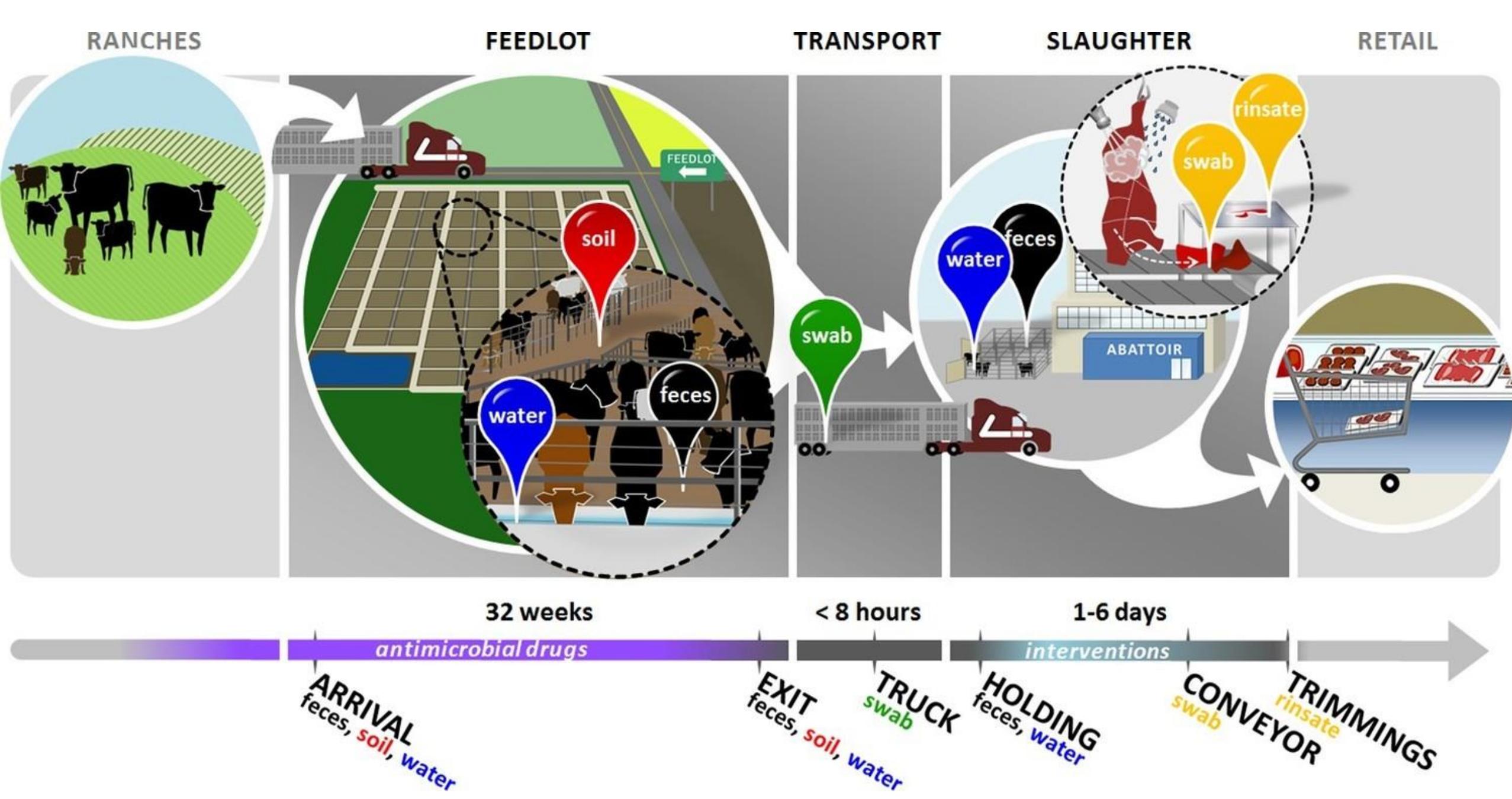
# Introduction of MEGARes

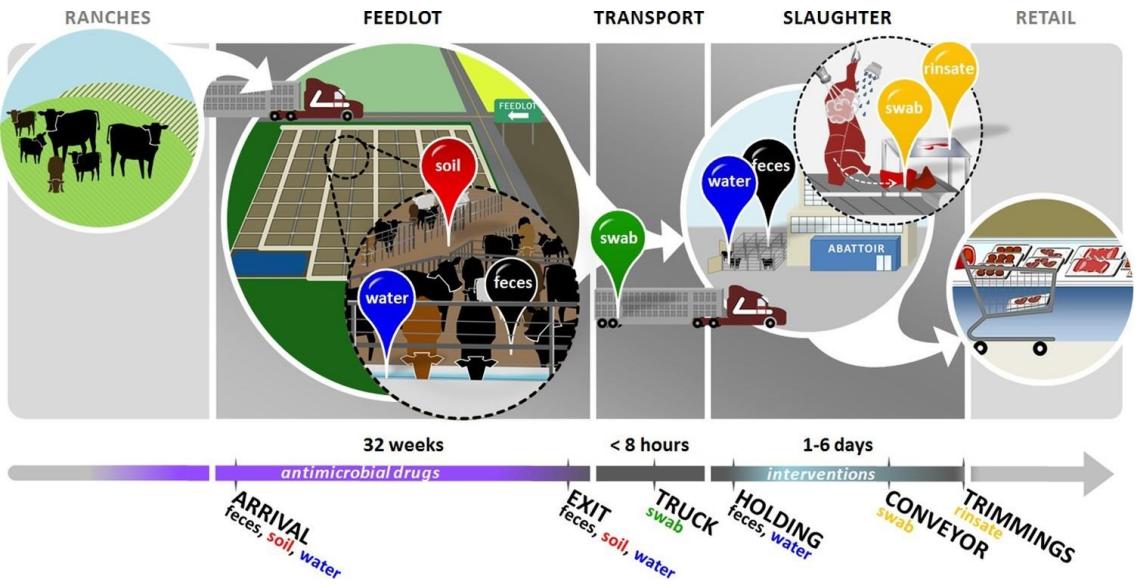


# Introduction of MEGARes



# A Case Study of Resistance





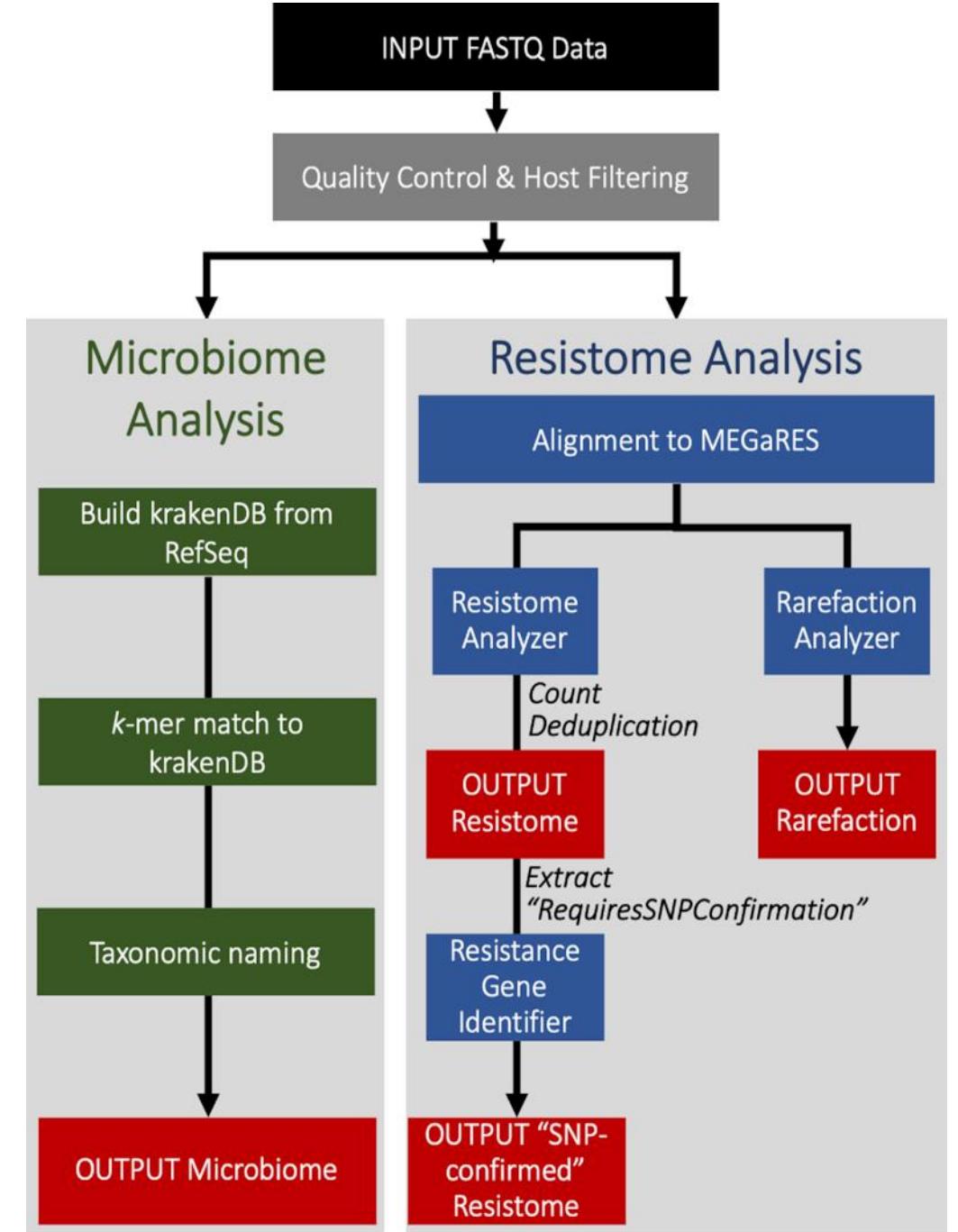
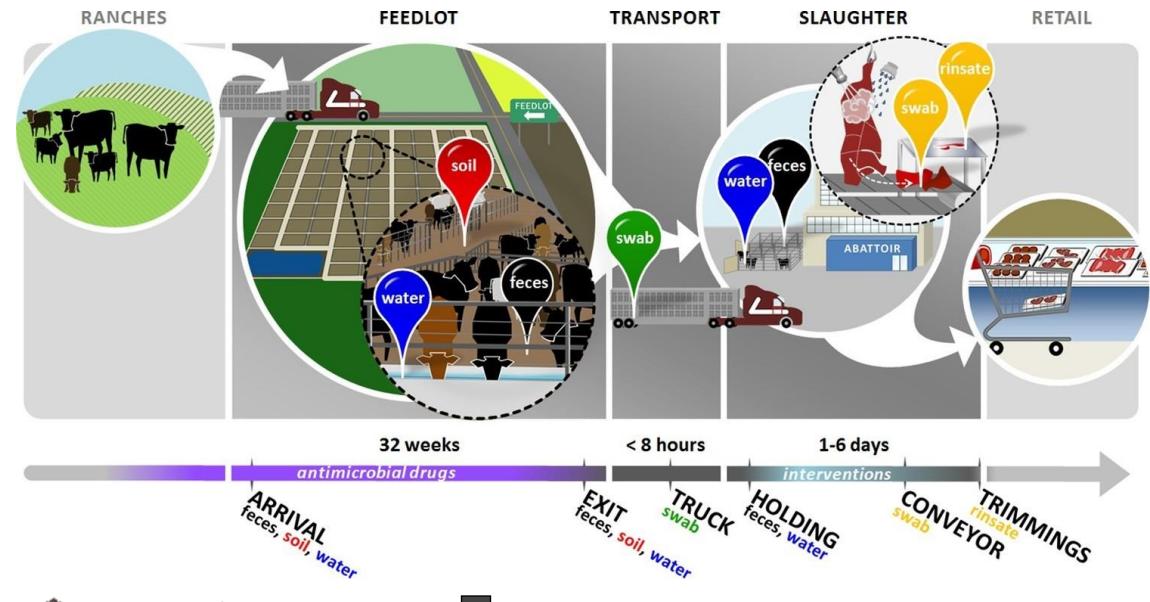
## Sample Preparation

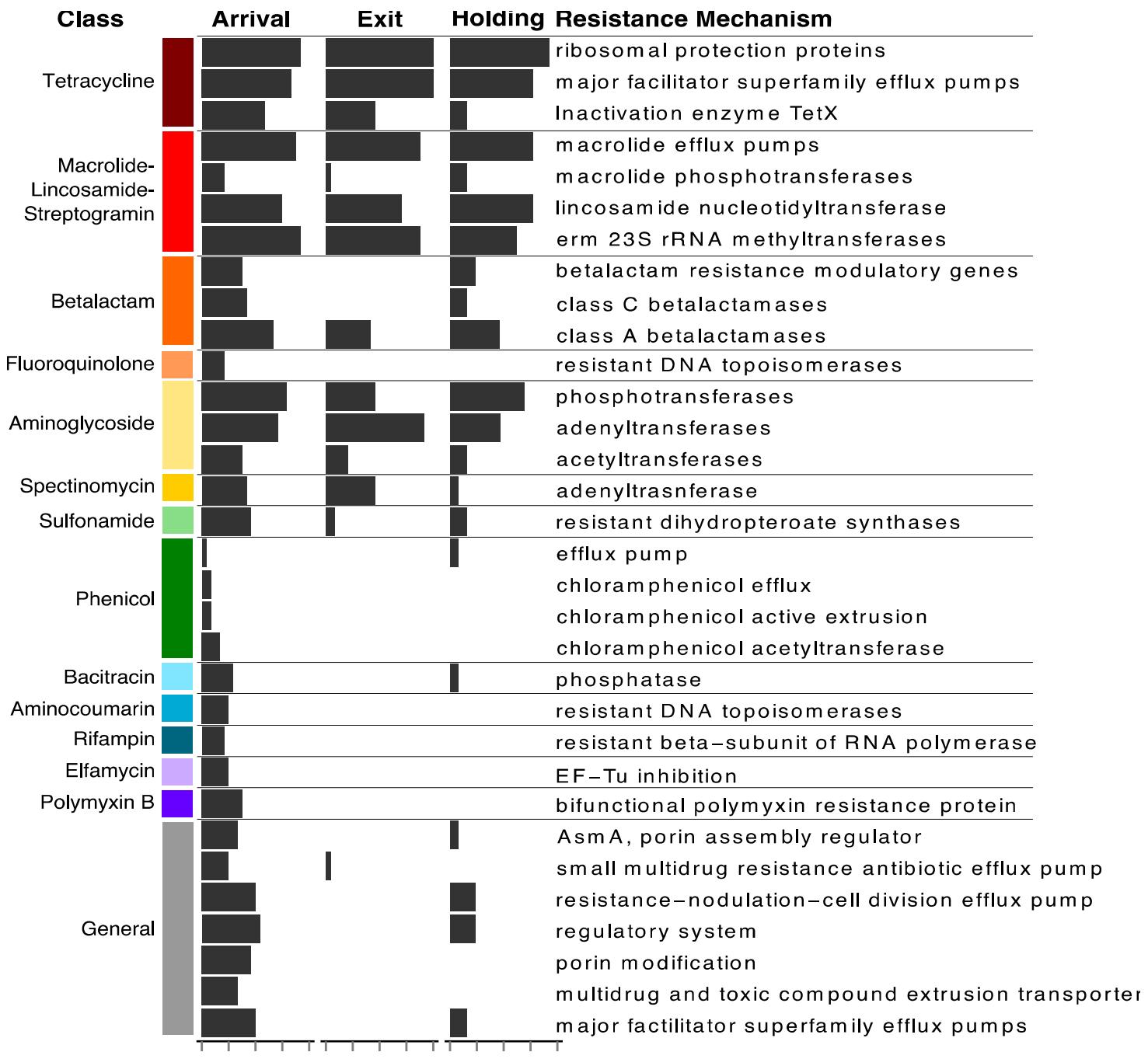


## Sequencing

# AMR+ Pipeline

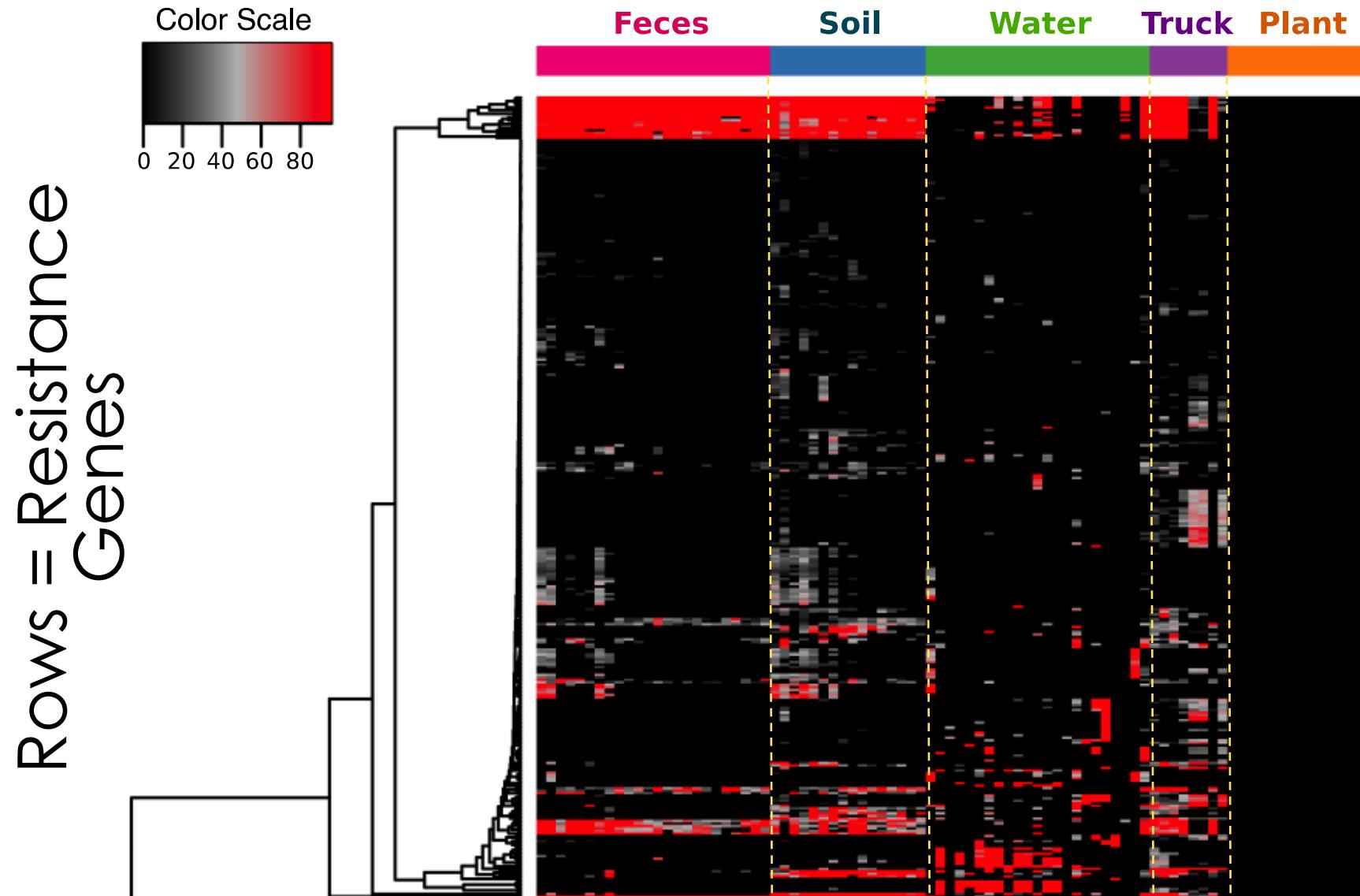
20



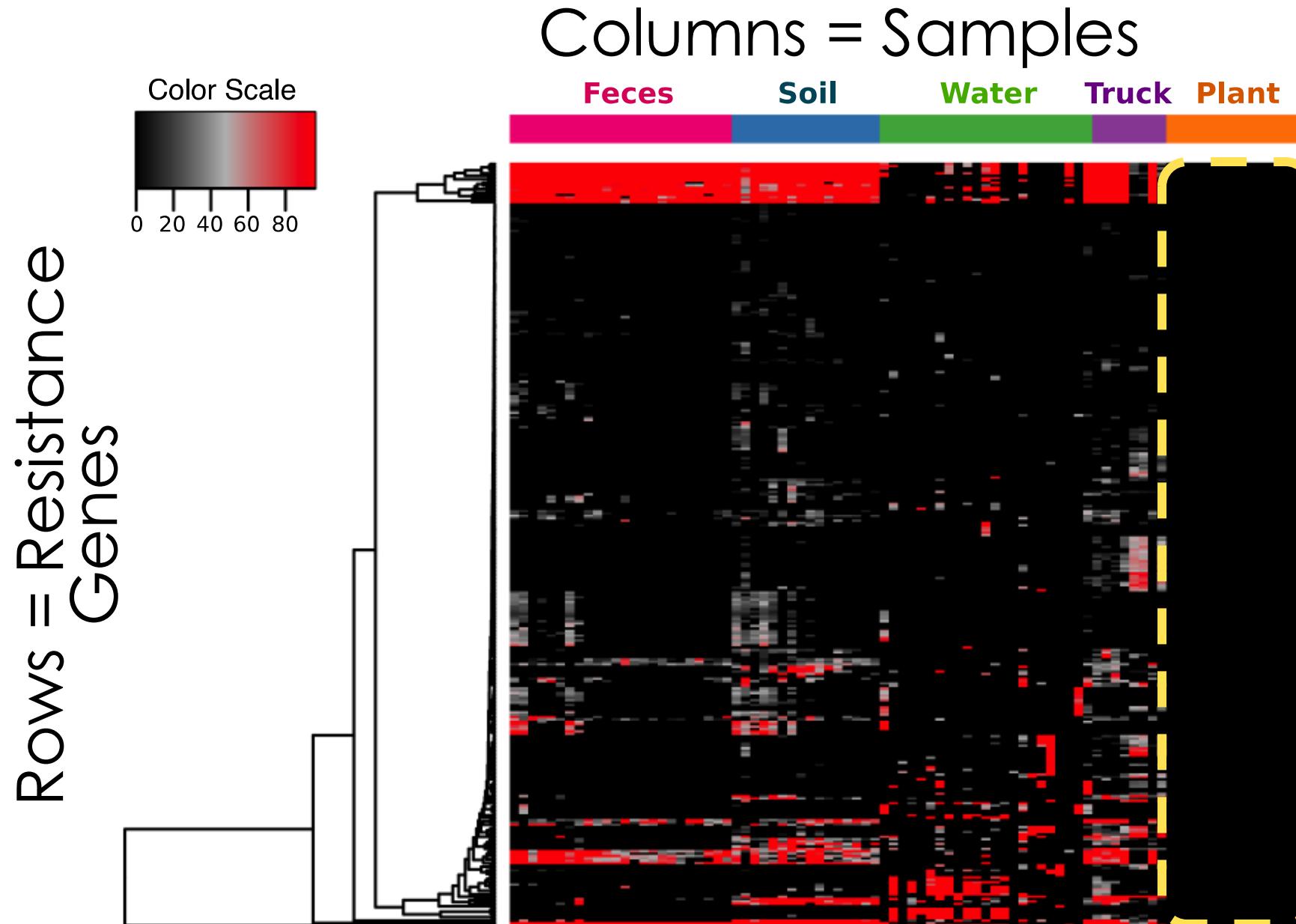


# The resistome clustered by matrix

Columns = Samples



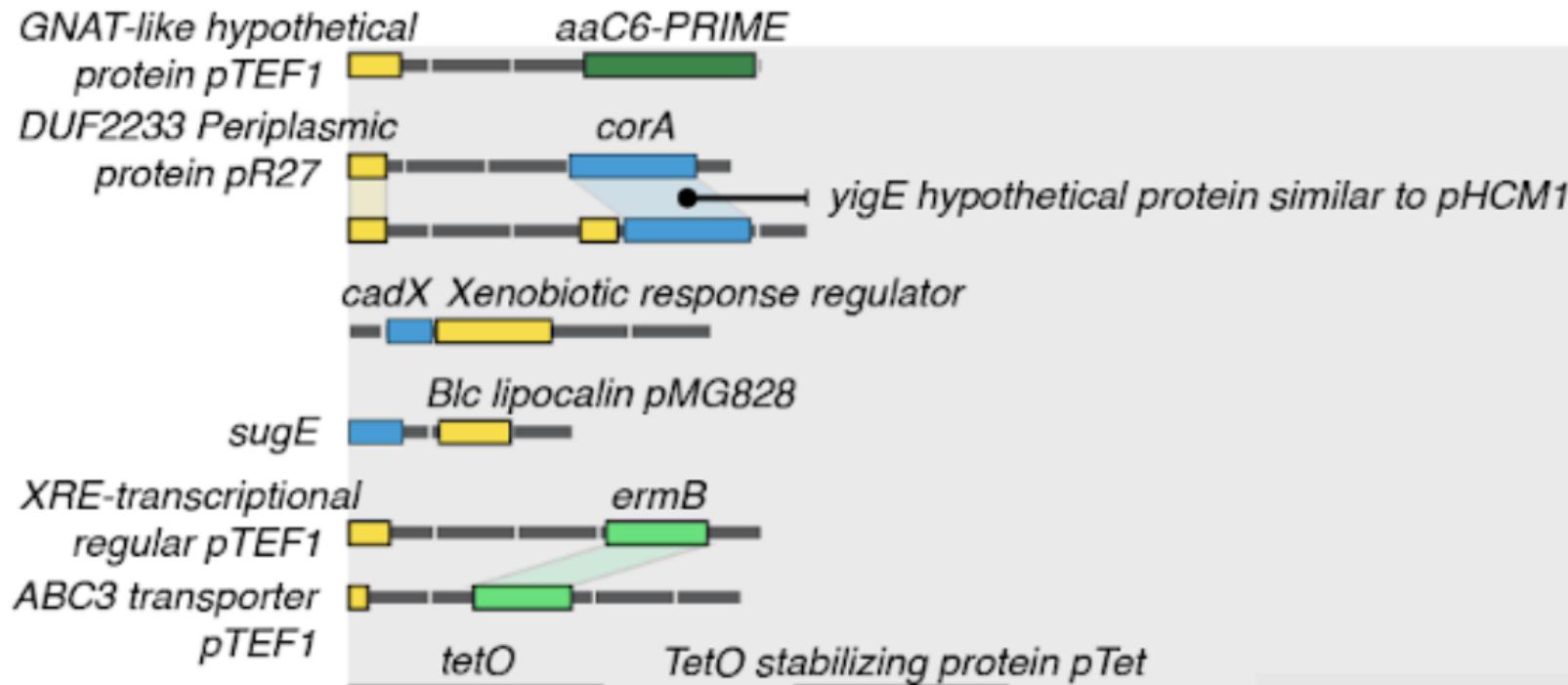
# No ARDs identified in plant samples



# Amplification of Resistance

# Challenges Related to Metagenomics

Human - TELSeq

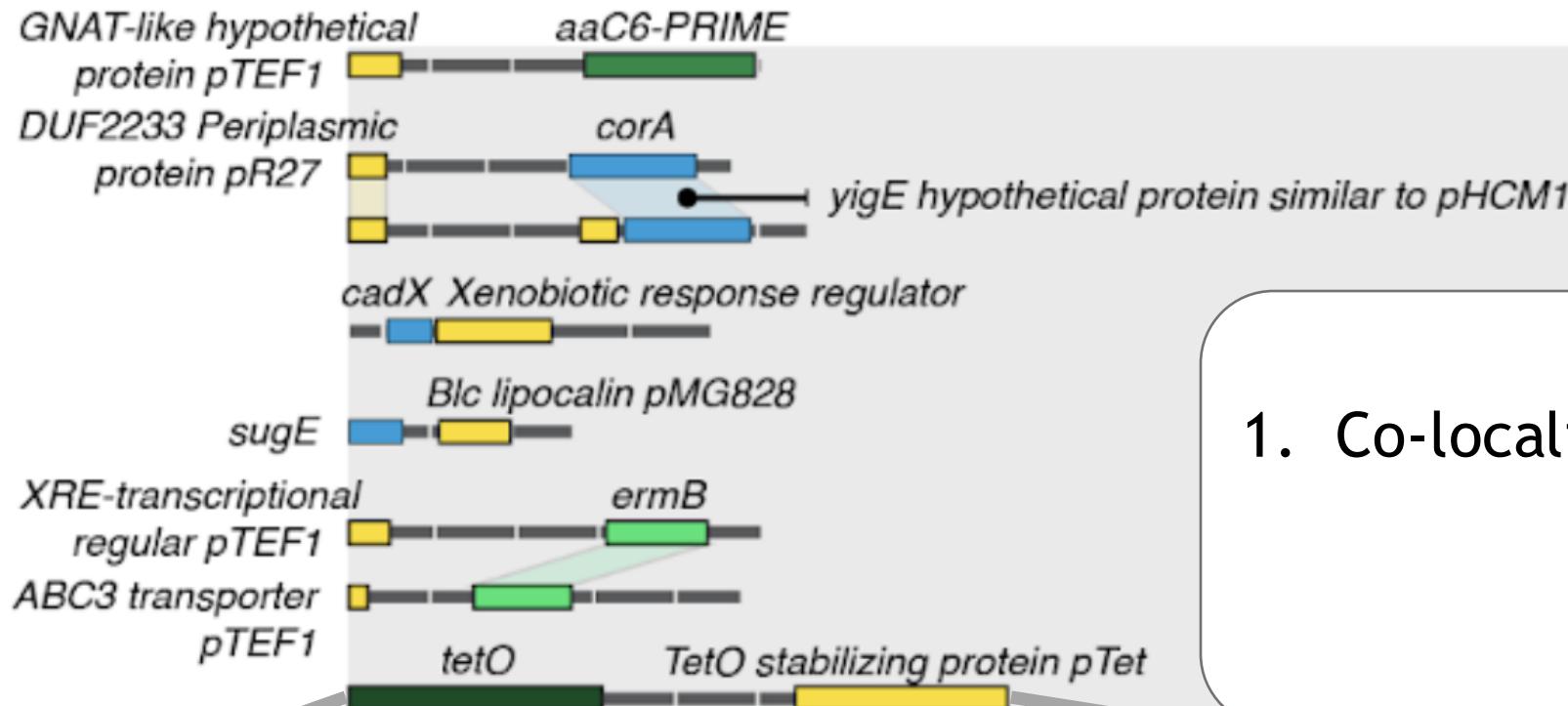


tetO

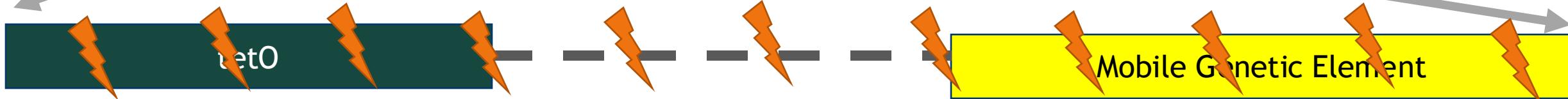
Mobile Genetic Element

# Challenges Related to Metagenomics

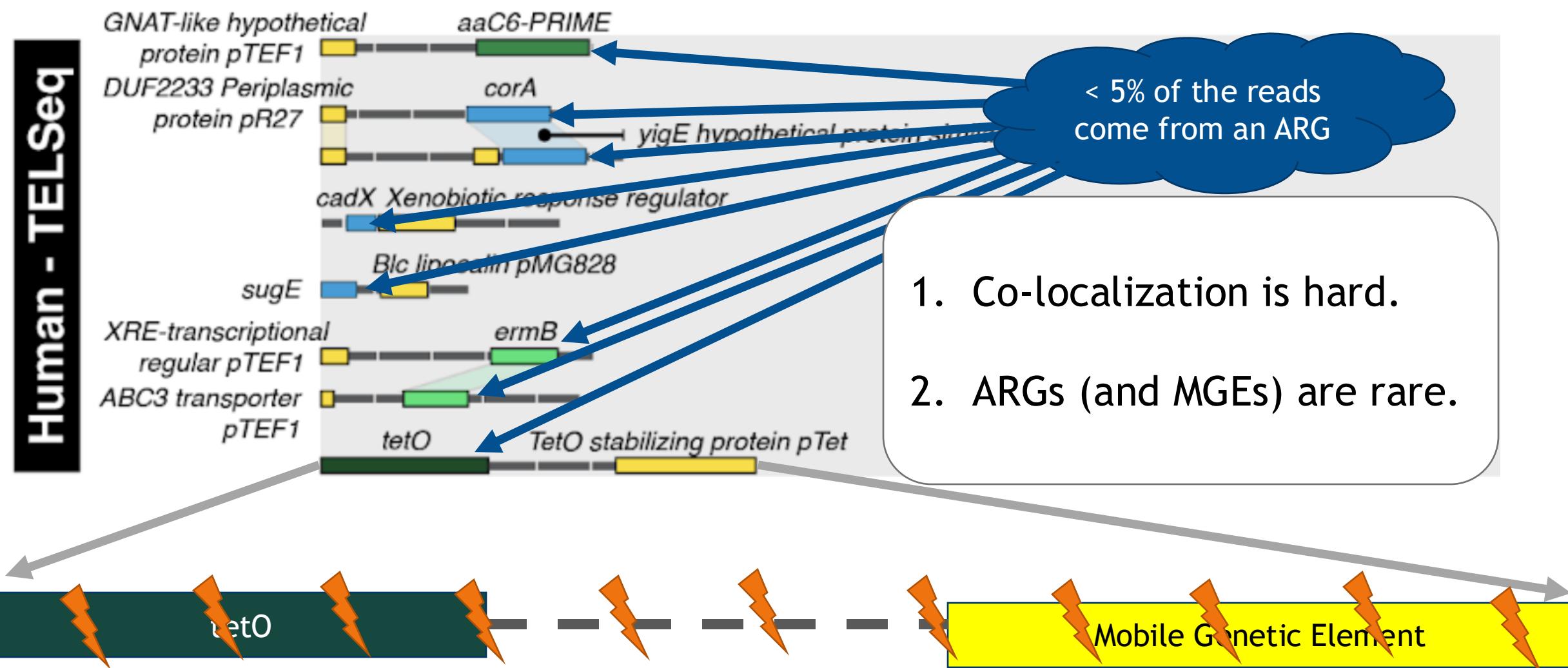
## Human - TELSeq



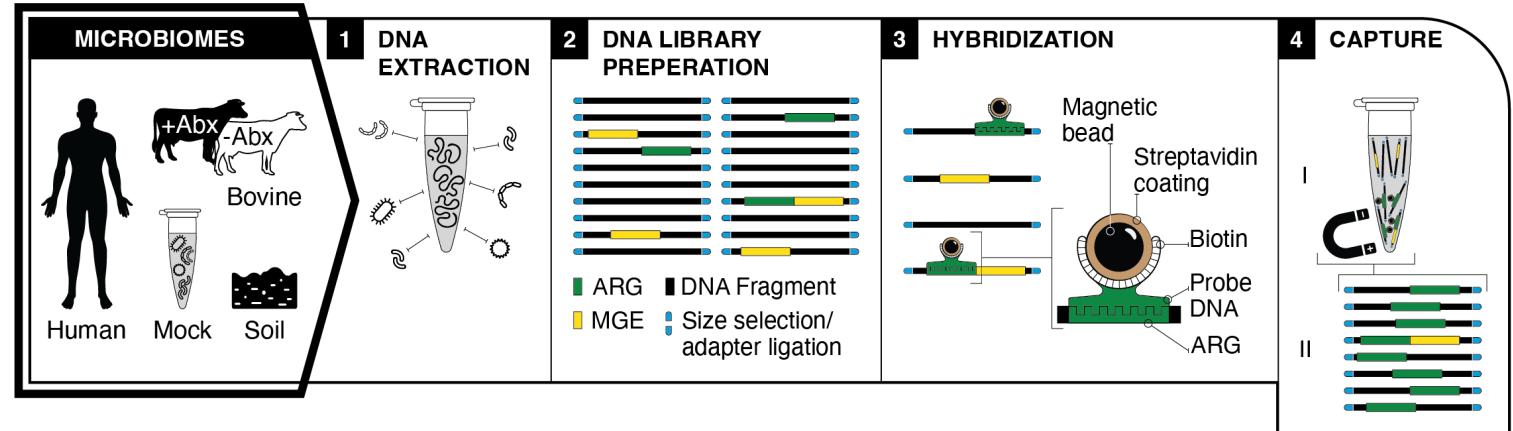
1. Co-localization is hard.



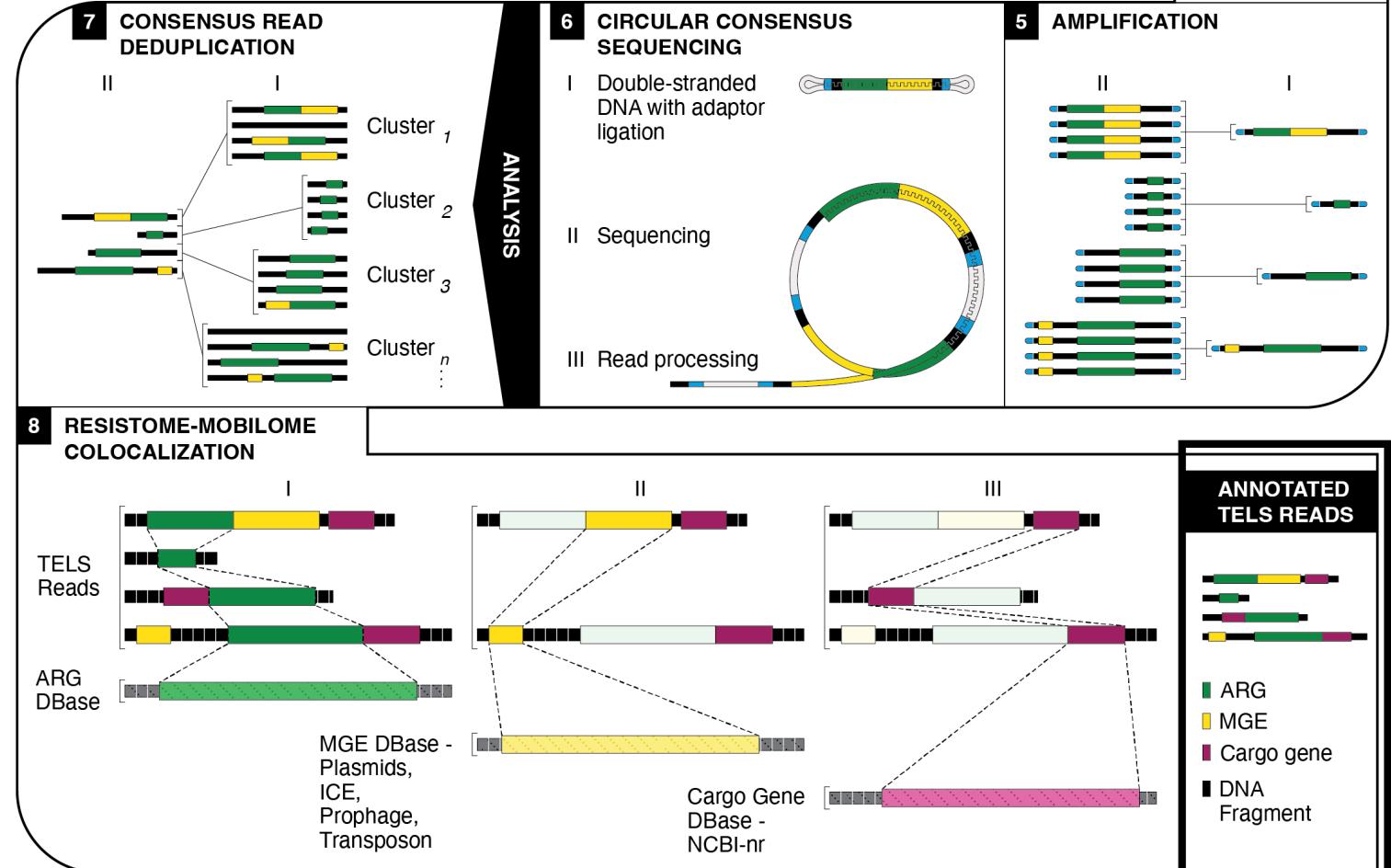
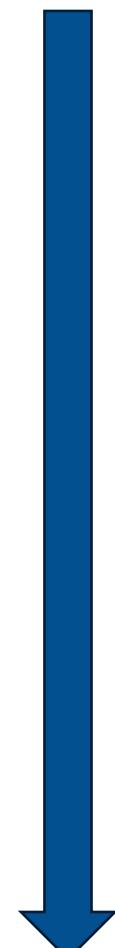
# Challenges Related to Metagenomics



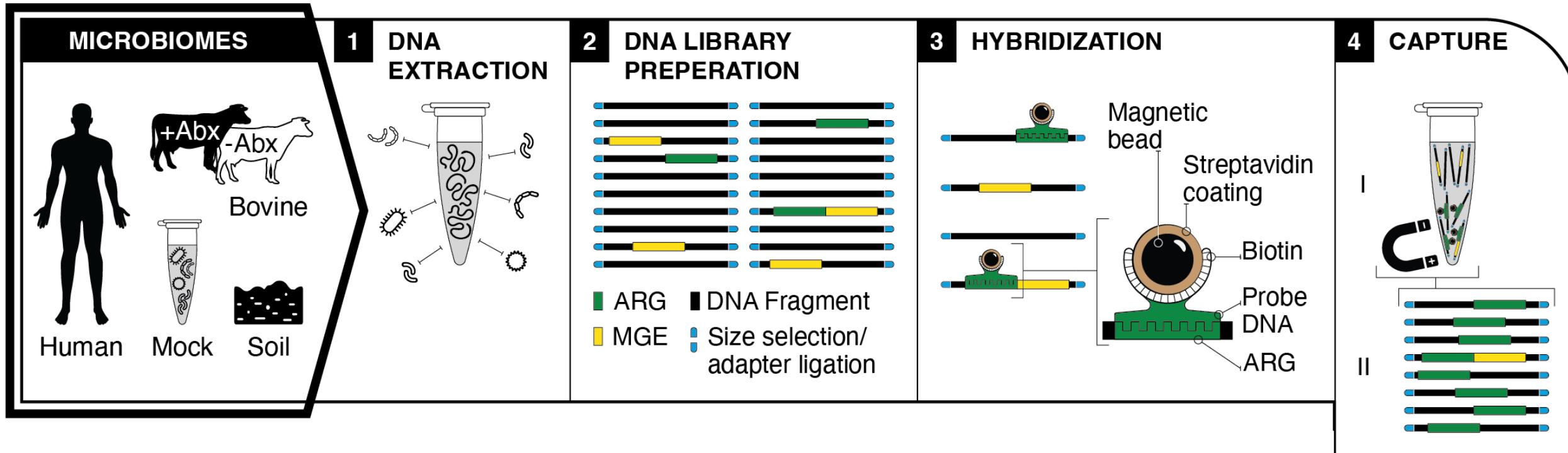
# Presequencing



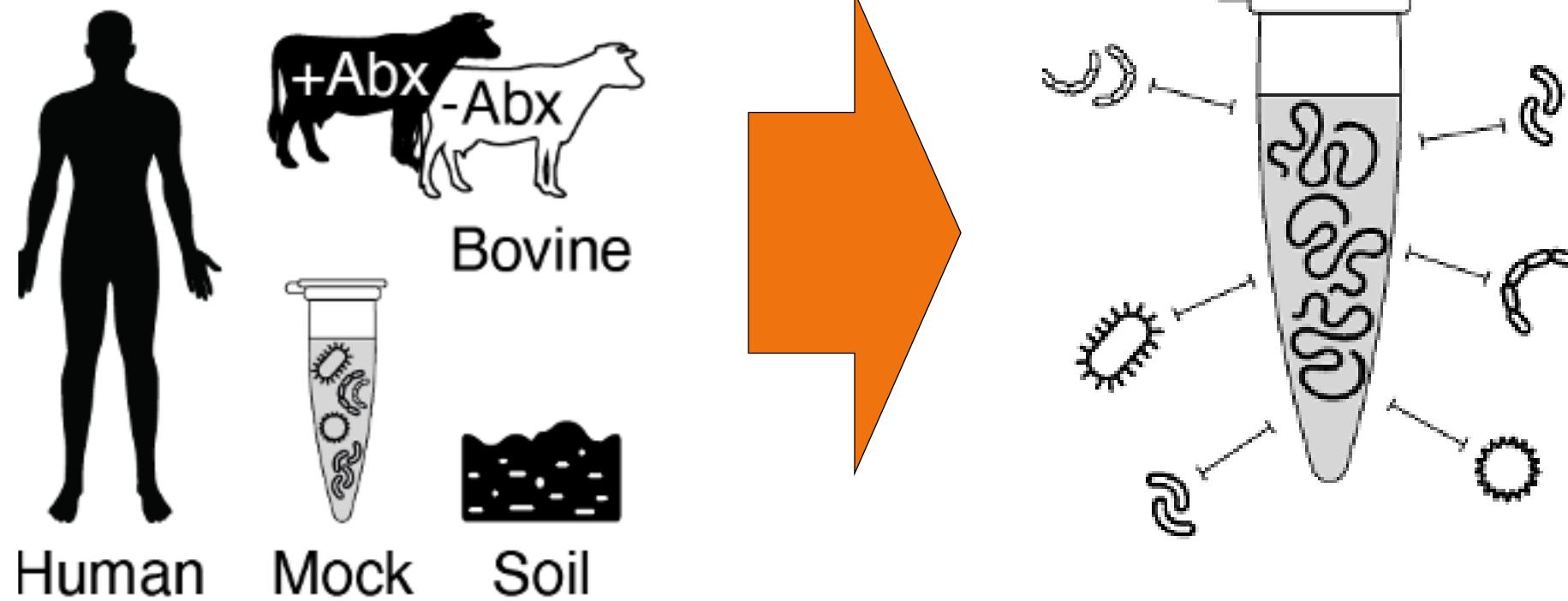
# Postsequencing



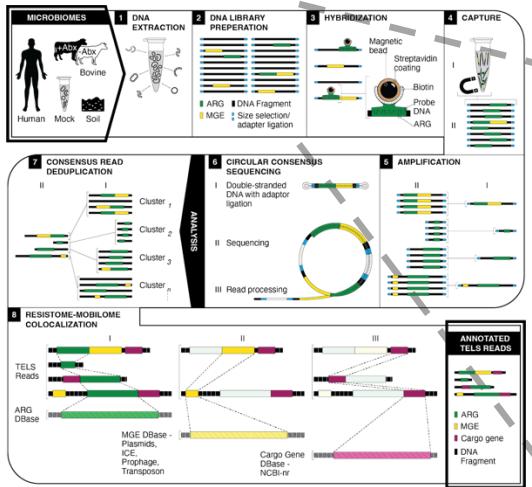
# Introduction to Amplification



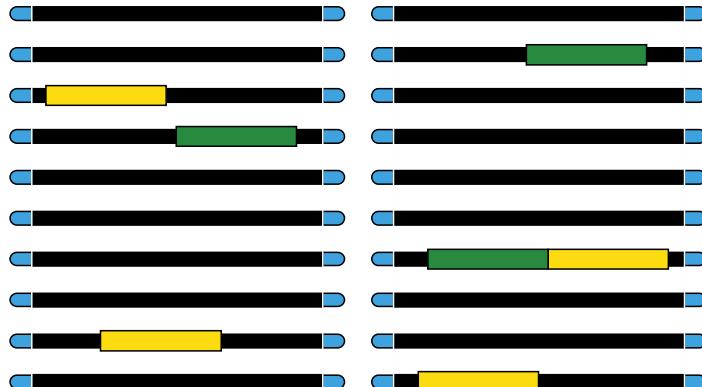
# Step 1: DNA Extraction



# Application of Baits

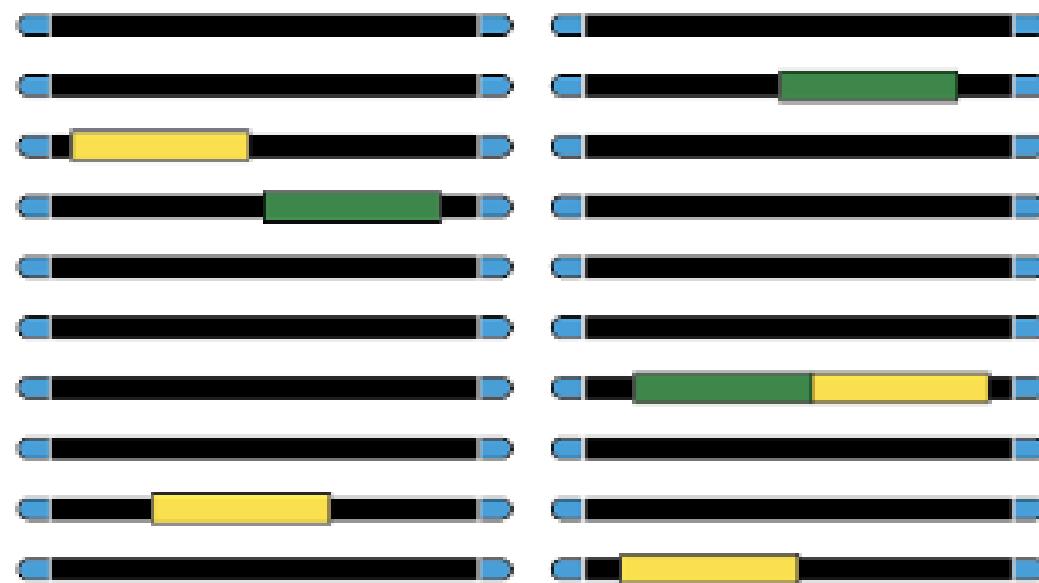


## 2 DNA LIBRARY PREPARATION



■ ARG ■ DNA Fragment  
■ MGE ■ Size selection/  
adapter ligation

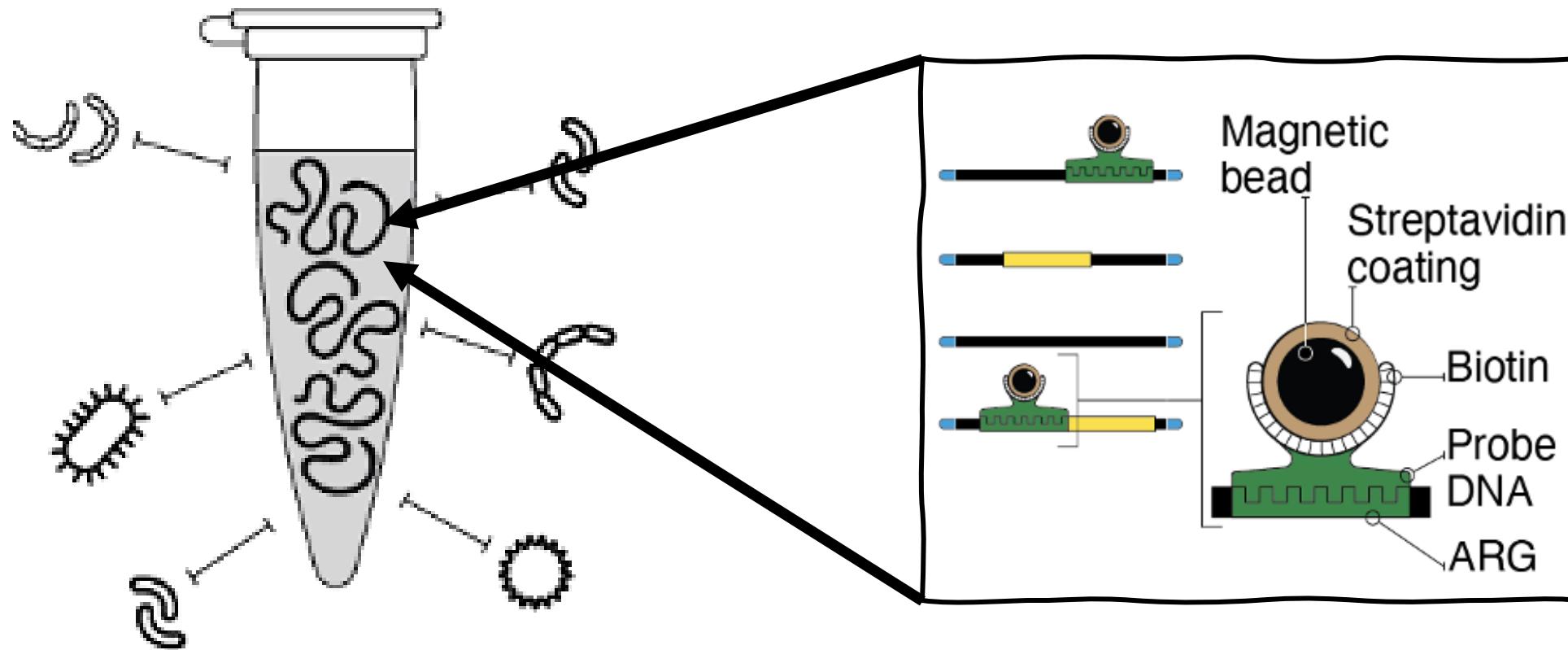
# Amplification via Probes



■ ARG ■ DNA Fragment  
■ MGE : Size selection/  
adapter ligation



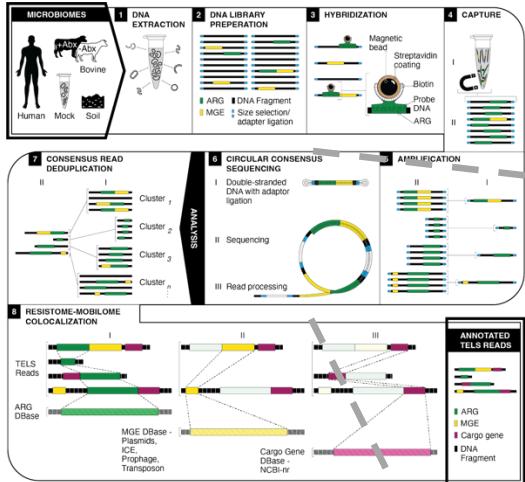
# Amplification via Probes



# 4: Sequencing



# Pre-analysis

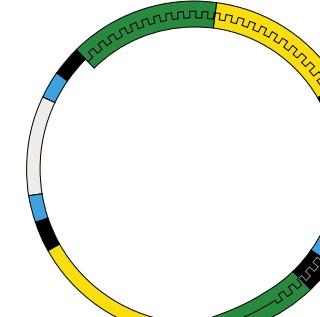


## 6 CIRCULAR CONSENSUS SEQUENCING

- I Double-stranded DNA with adaptor ligation



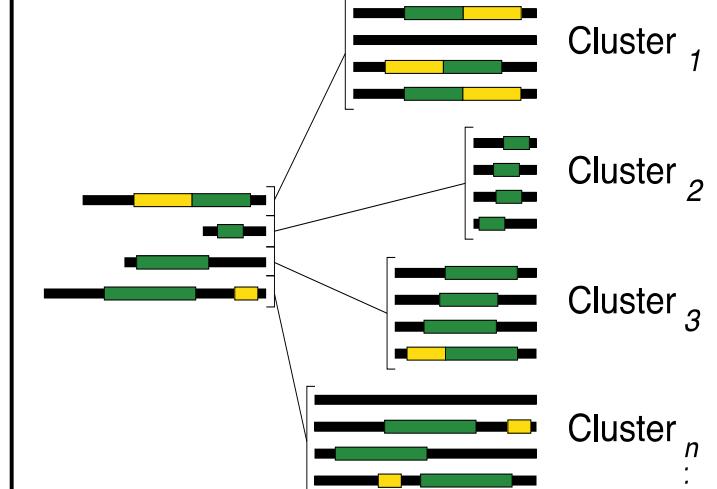
- II Sequencing



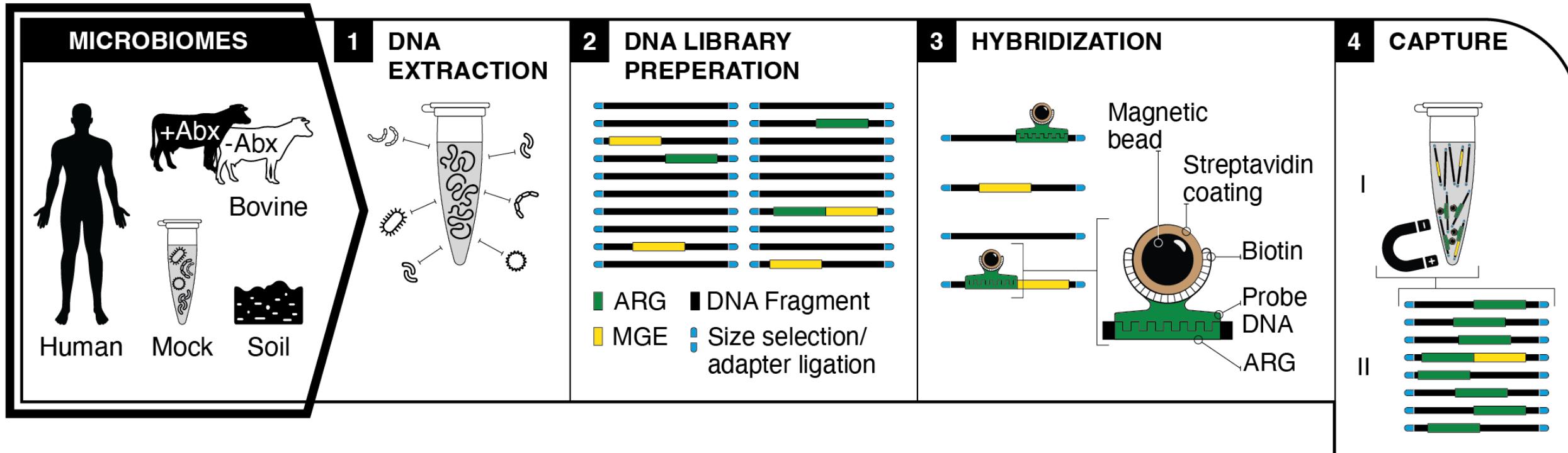
- III Read processing

## 7 CONSENSUS READ DEDUPLICATION

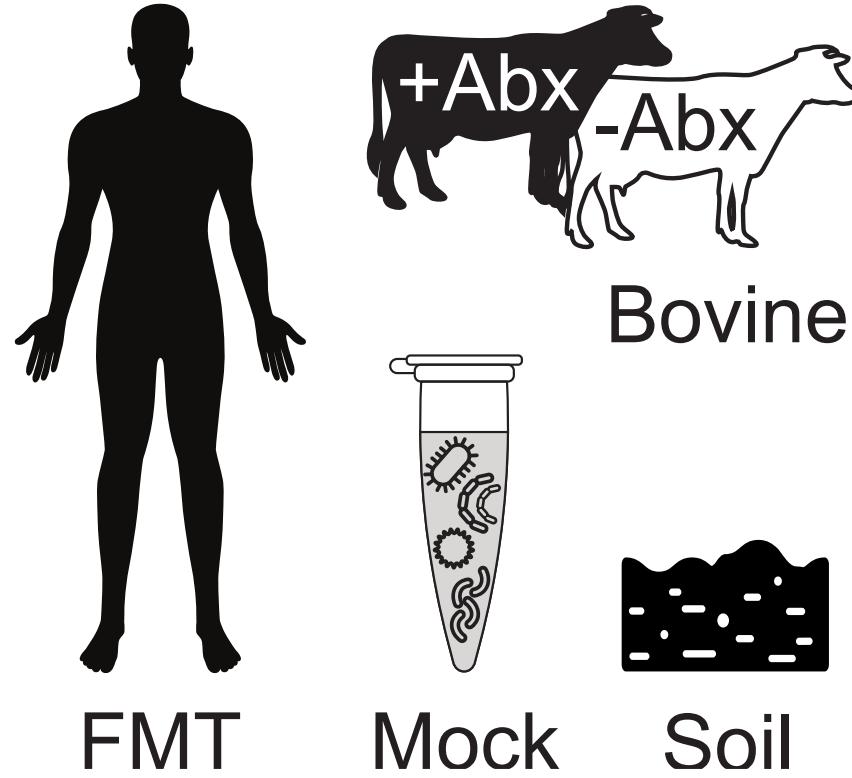
- II



# Introduction to Amplification

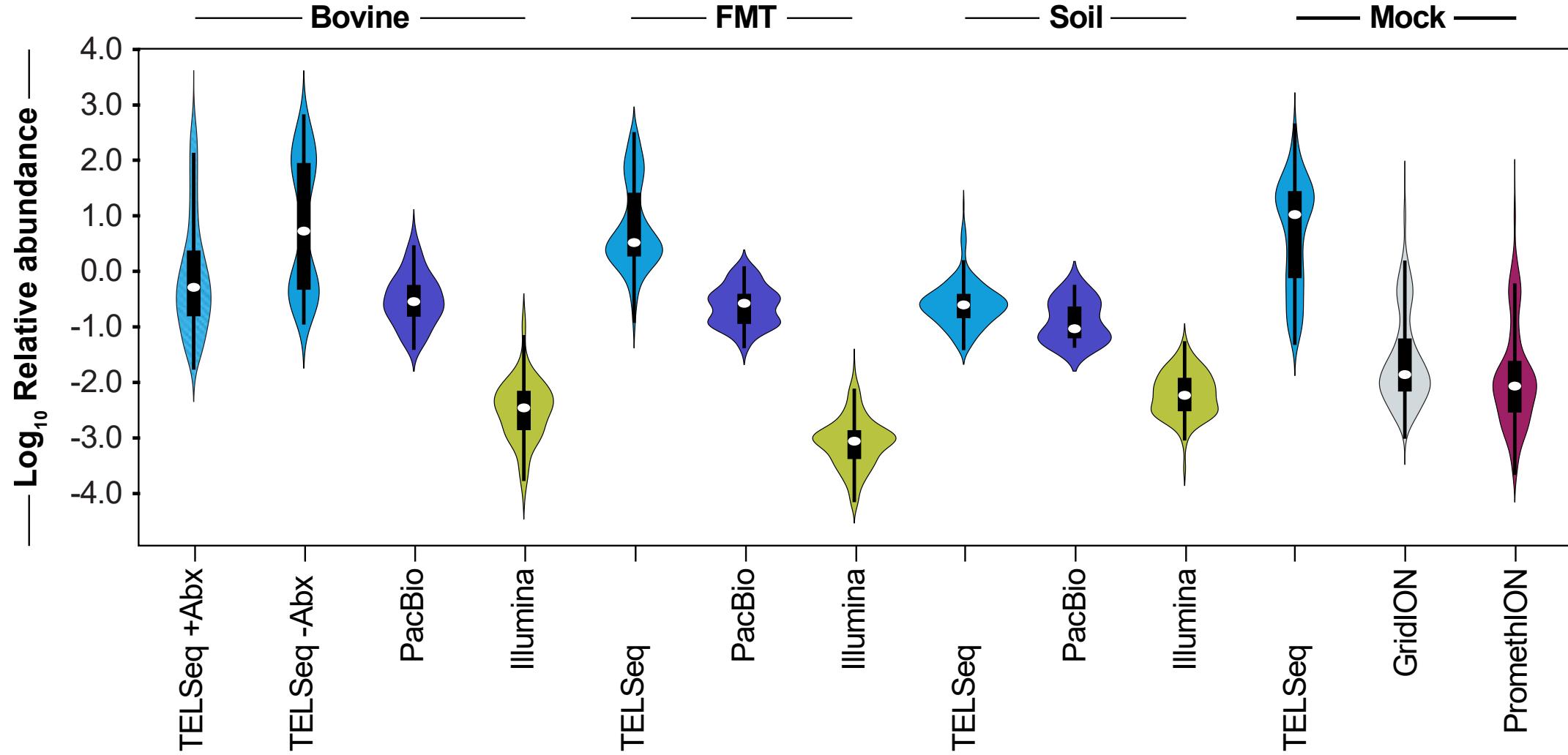


# TELS I: Datasets



- Fecal Microbiota Transplant from healthy human donor (**FMT**).
- Holstein Friesian dairy cow with recent exposure to antibiotic (**+ABX**).
- Holstein Friesian dairy cow deemed systematically healthy (**-ABX**).
- Soil sample from strip of un-utilized prairie in Mower County, Minnesota (**SOIL**).
- ZymoBiomics™ microbial community (**MOCK**).

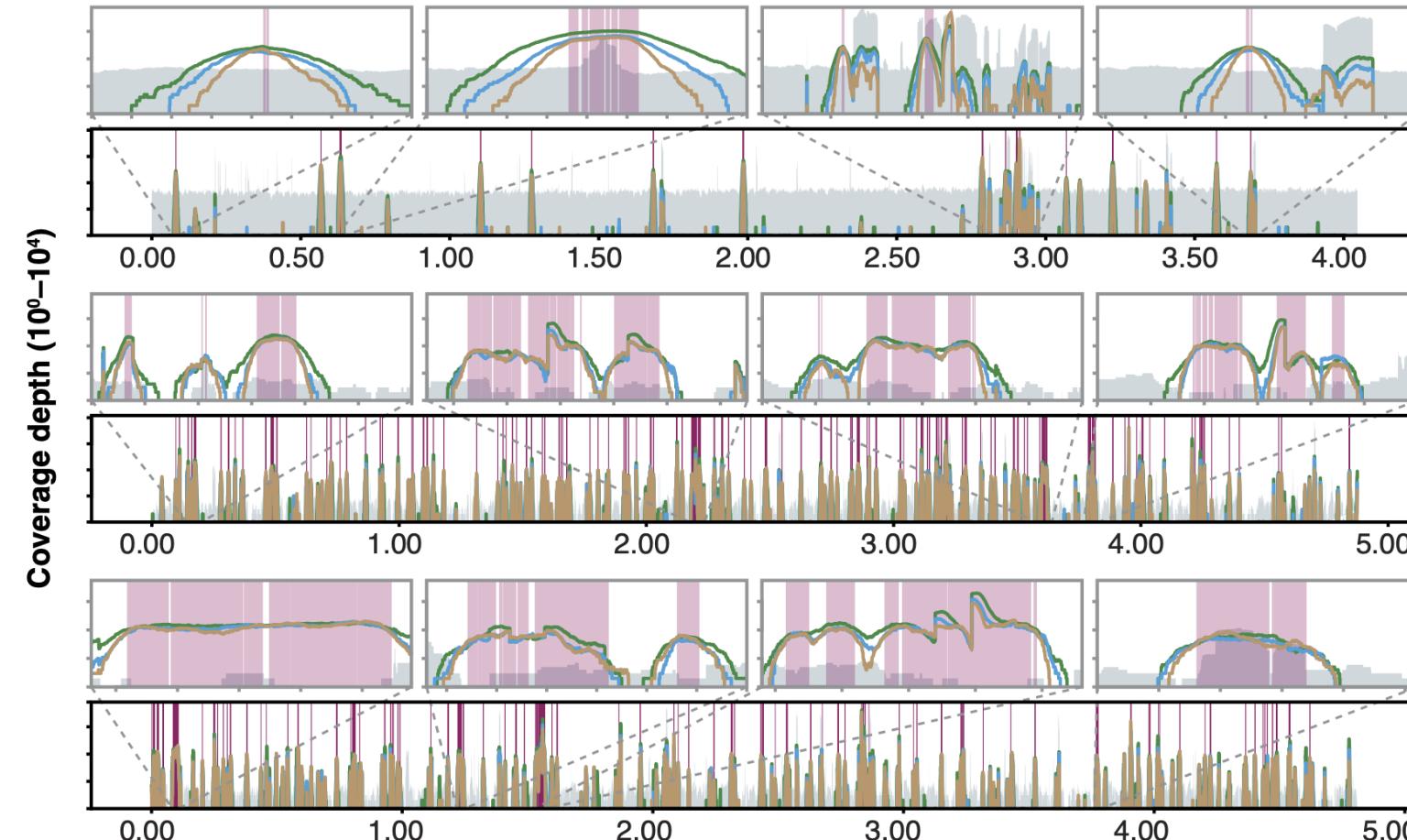
# Log Relative Abundance of ARGs



# Coverage Plots

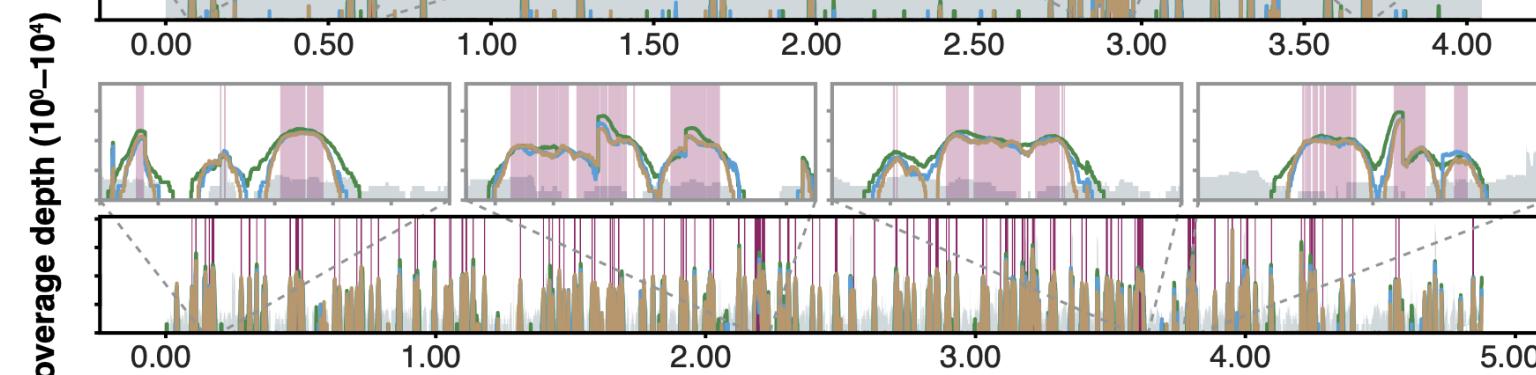
c. *Bacillus subtilis*

00.890000%



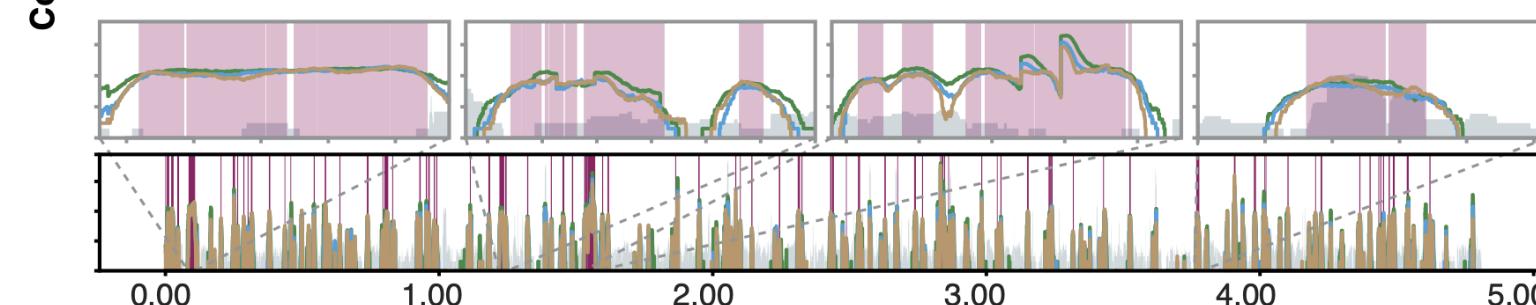
d. *Escherichia coli*

00.089000%

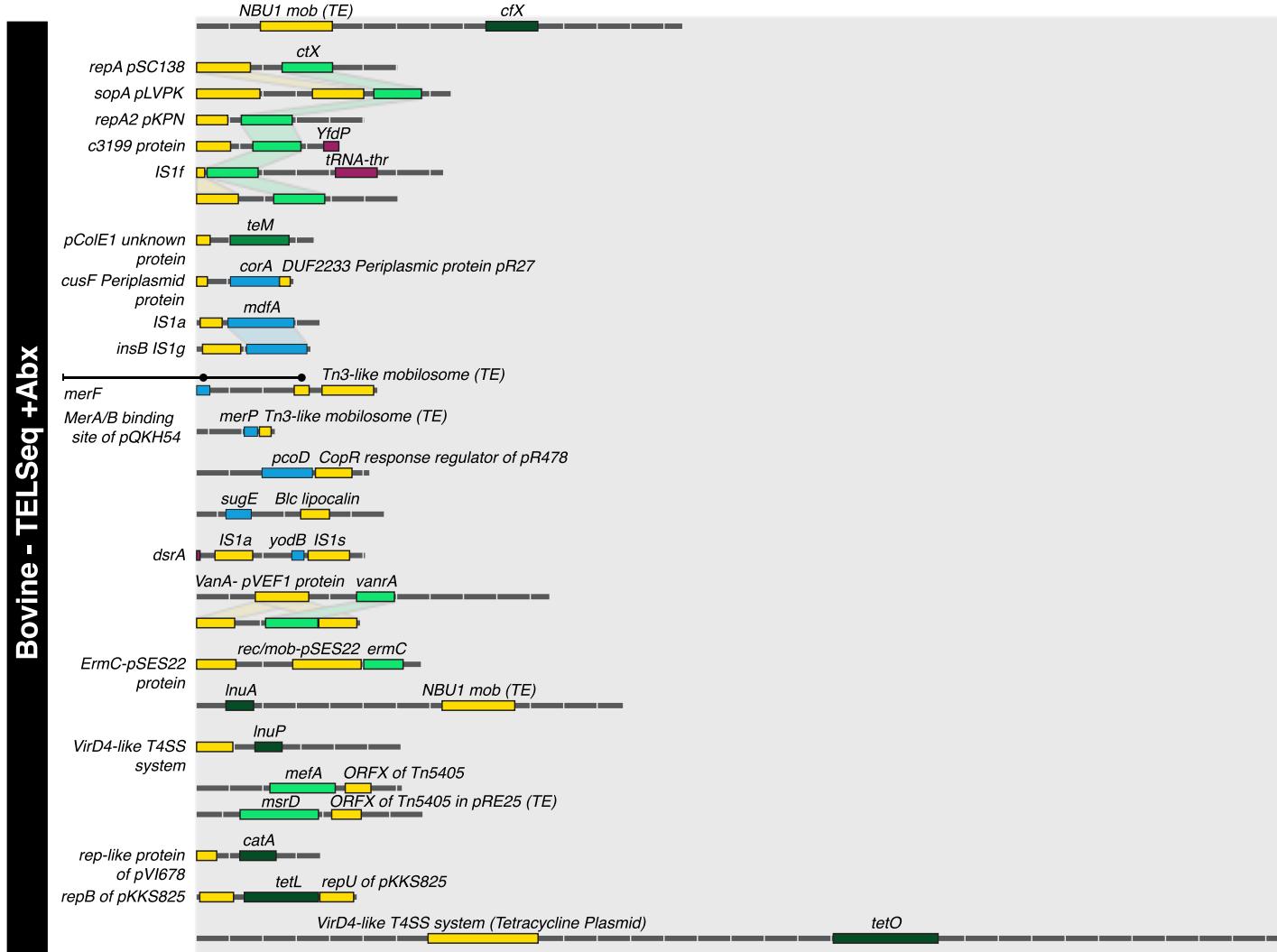


e. *Salmonella enterica*

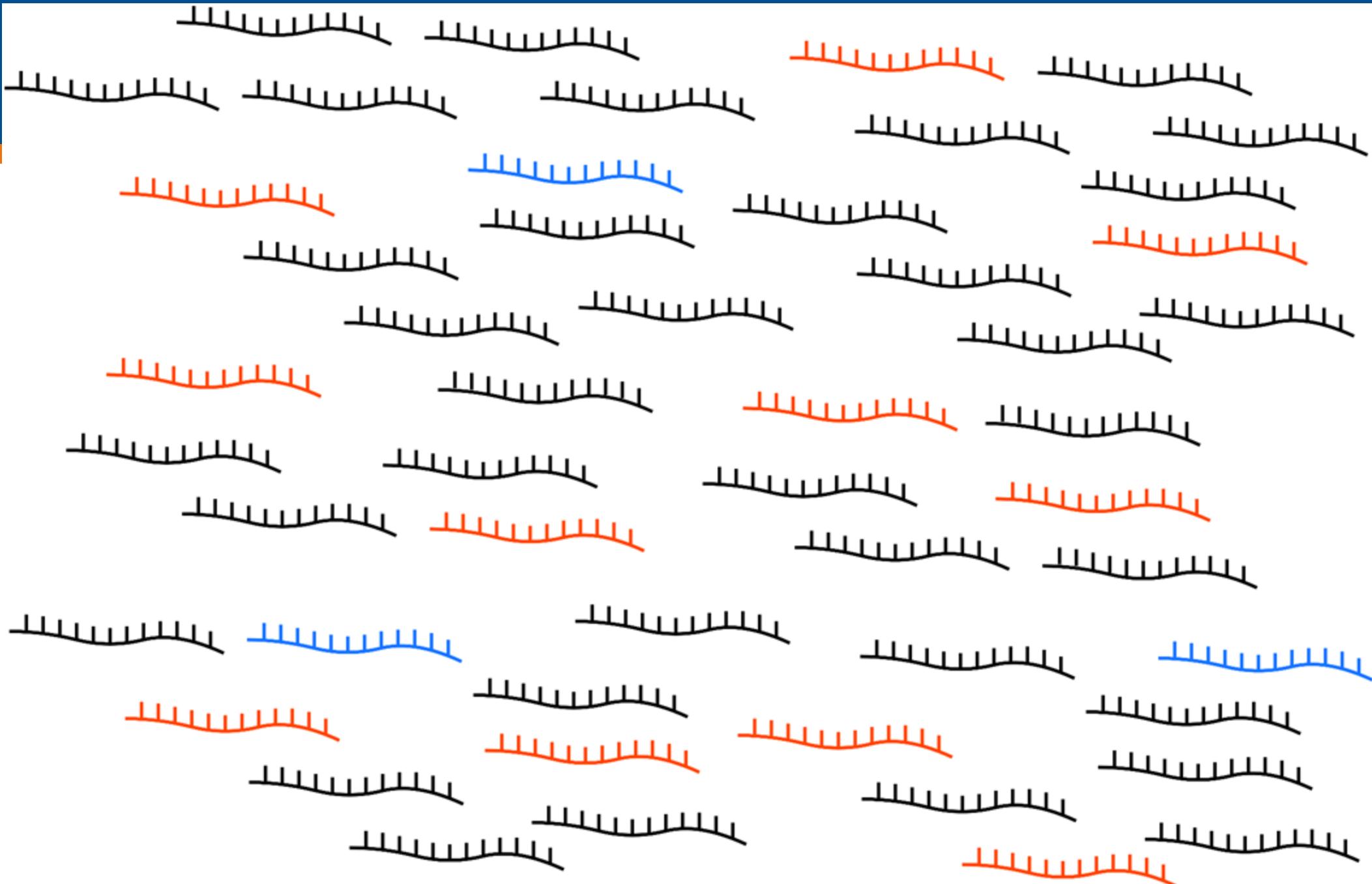
00.089000%

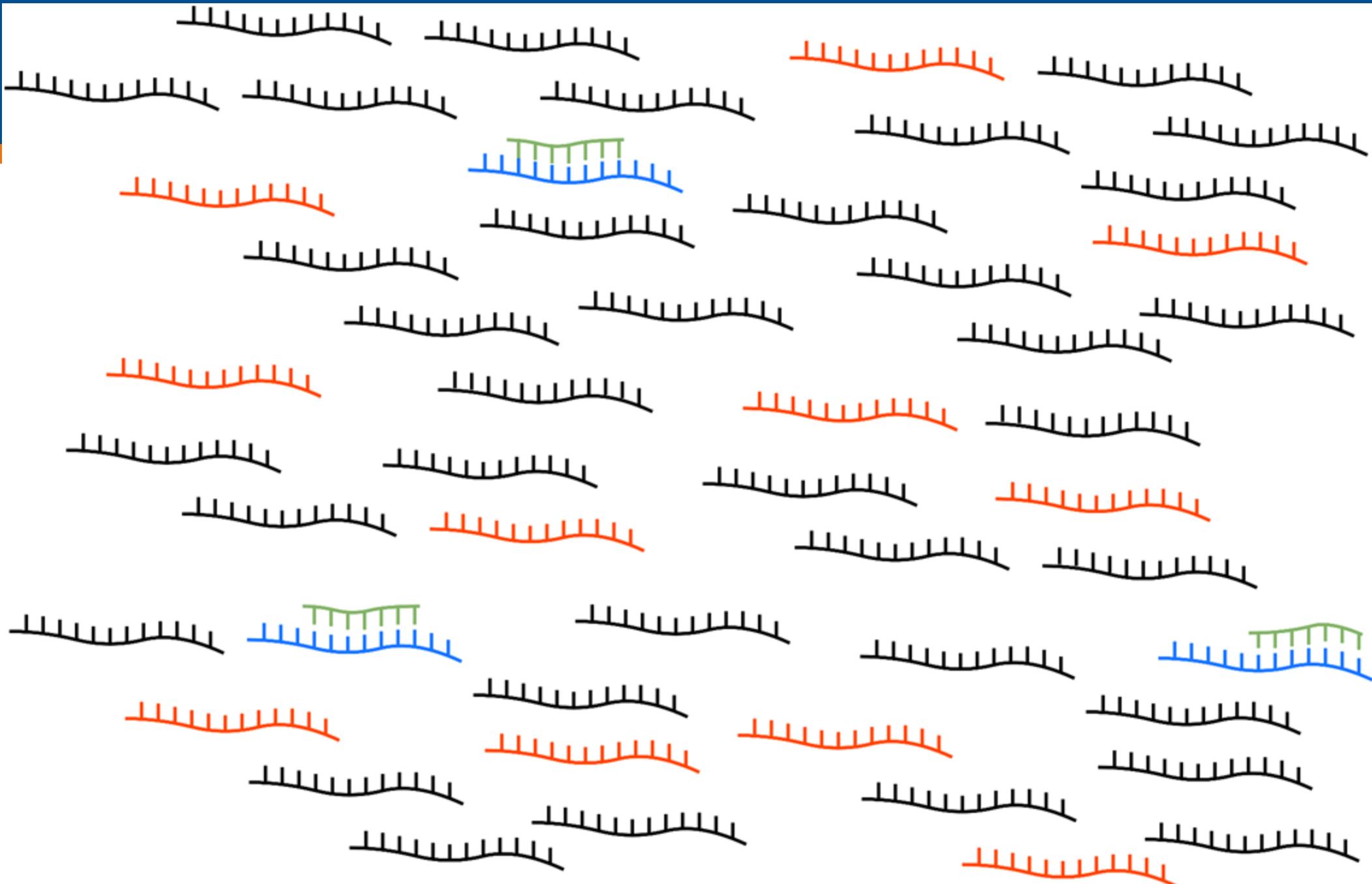


# Read Annotation



# Amplification of Resistance







# Manufacturing Costs

Tier 1: 1- 499 kbp (up to 57,606 baits) **\$12,796.00**

Tier 2: 0.5 - 2.999 Mbp (up to 57,606 baits) **\$16,560.00**

Tier 3: 3.0 - 5.999 Mbp (up to 57,606 baits) **\$21,578.00**

Tier 4: 6.0 - 11.999 Mbp (up to 115,212 baits) **\$32,517.50**

Tier 5: 12.0 - 24.0 Mbp (up to 230,424 baits) **\$40,546.50**

# Manufacturing Costs

Tier 1: 1- 499 kbps (up to 57,606 baits)	\$12,796.00
Tier 2: 0.5	Other Considerations: \$16,560.00
Tier 3: 3.0	Coverage, Specificity, GC-content, \$21,578.00 Blacklisting contaminants,
Tier 4: 6.0	) \$32,517.50
Tier 5: 12.0 - 24.0 Mbps (up to 230,424 baits)	\$40,546.50



I can create  
synthetic  
probes and  
apply them to  
enrich for  
target DNA!!

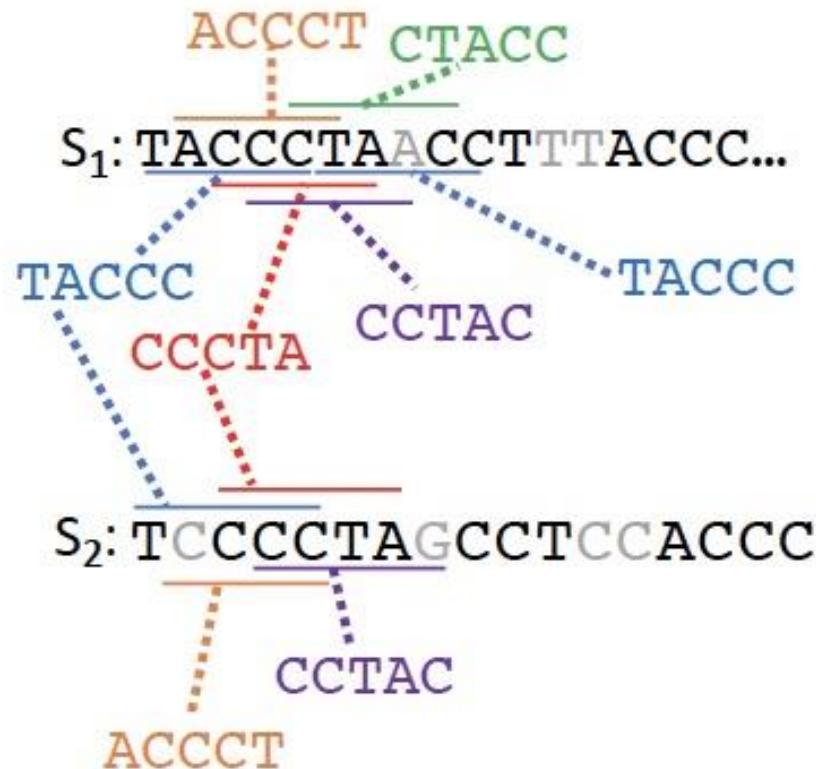
Great! I can  
run a k-mer  
counter on the  
database!  
That's easy.

Well, the probes bind  
to anything with 70%  
sequence homology and  
I need the smallest set  
possible...

Ugh.



# The Minimum Bait Cover Problem



b<sub>1</sub>: TACCC  
b<sub>2</sub>: ACCCT  
b<sub>3</sub>: CCCTA  
b<sub>4</sub>: CCTAC  
b<sub>5</sub>: CTACC

A set  $T = \{T_1, \dots, T_m\}$   $\theta$ -covers  $S = \{S_1, \dots, S_n\}$  if for every  $S_i$  there exists at least one in  $T$  that has Hamming distance at most  $\theta$  from  $S_i$

# The Minimum Bait Cover Problem

Input: two integers  $\theta \geq 0$  and  $L > 0$  and a set  $S$  of  $n$  strings  $S_1, \dots, S_n$  over a finite alphabet  $\Sigma$ .

Question: What is the smallest possible set  $T$  of  $L$  –length strings such that  $S$  is  $\theta$ -covered by  $T$ .

# Minimum Bait Cover in Hard

Proposition 1: Minimum Bait Cover is NP-hard, even for  $n = 1$ ,  $\theta = 1$ , and  $L = 2$ .

- Rules out an FPT or Slicewise Polynomial algorithm w.r.t.  $n, \theta, L$ .
- However, this result requires that  $|\Sigma|$  is large.

# Minimum Bait Cover in Hard... Very Hard

Theorem 2. For every  $k \geq 2$ , Minimum Bait Cover is NP-hard even for  $|\Sigma| = k$ ,  $S = \{S\}$ ,  $\theta = 0$  and  $L = O(\log |S|)$ .

- The problem remains intractable for small problem instances, e.g., DNA alphabet, a single sequence, Hamming distance of 0, and  $\log |S|$
- Solves an existing open problem in the literature.

# Heuristic Solution

- Build the de Bruijn graph, find all the unitigs, find all  $L$  length sequences on unitigs.
- Slide a window of length  $L$ . If the starting position is not yet covered, take it; if not go the next one.
- Dynamic programming algorithms to bootstrap a solution
- Expectation-Maximization algorithms.

# Heuristic Solution

1. Construct the FM-index and GSA of  $S_1, \dots, S_n$
2. Store a bit vector for each position of  $S_1, \dots, S_n$
3. Consider each position of every string in  $S_1, \dots, S_n$ , say  $S_i[j]$  if it not covered then we add the L-mer starting at  $S_i[j]$  to the bait set.
4. Use the FM-index and GSA to find all positions covered by the newly added bait.

# Heuristic Solution

1. Construct the FM-index and GSA of  $S_1, \dots, S_n$
2. Store a bit vector for each position of  $S_1, \dots, S_n$
3. Consider each position of every string in  $S_1, \dots, S_n$ , say  $S_i[j]$  if it not covered then we add the L-mer starting at  $S_i[j]$  to the bait set.
4. Use the FM-index and GSA to find all positions covered by the newly added bait.

# Heuristic Solution

```
For i = 1..|B| - k + 1 do:
```

```
    x ← B[i..i+k-1]
```

```
    [r1, r2] ← Fmindex.Search(x)
```



Find all k-mers

```
For j = r1..r2 do
```

```
    t, p ← GSA[j]
```

```
    B' ← St[p-i+1..p-i+L]
```

```
    if d(B', B) ≤ θ
```

```
        Mark St[p-i+1..p-i+L] as covered
```



Extend to  
k-mers to  
length L

# Competing Methods

- **CATCH** (Nature Biotech. 2019) has quadratic running time in practice.
- **MrBait** requires all-pairs alignment of baits as post-processing.
- **BaitFisher** requires a multiple sequence alignment.

GATACCGATTACCGATGGAAAGA

TGGGGAAAGATAACCGATTTCAG

ATGGAAAGATGGAATTACTCAGA

GACAGGGAAATAAGTCCTG

GATACCGA

TTACCGATGGAAAGA

TGGGGAAAGATAACCGATTTC

ATGGAAGATGGAATTACTCAGA

GACAGGGAAATAAGTCCTG

GATACCGATTACCGATGGAAAGA

TGGGGAAAGATACCGATTTCG

ATGGAAAGATGGAATTACTCAGA

GACAGGGAAATAAGTCCTG

GATACCGATTACCGATGGAAAGA

TGGGGAAAGATACCGATTTCGG

ATGGAAAGATGGAAATTACTCAGA

GACAGGGAAATAAGTCCTG

GATACCGATTACCGATGGAAAGA

TGGGGAAAAGATACCGATTTCGG

ATGGAAAGATGGAAATTACTCAGA

GACAGGGAAATAAGTCCTG

GATACCGATTACCGATGGAAAGA

TGGGGAAAGATACCGATTTCGG

ATGGAAAGATGGAAATTACTCAGA

GACAGGGAAATAAGTCCTG

GATACCGATTACCGATGGAAAGA

TGGGGAAAGATACCGATTTCG

ATGGAAAGATGGAAATTACTCAGA

GACAGGGAAATAAGTCCTG

GATACCGATTACCGATGGAAAGA

TGGGGAAAGATACCGATTTCG

ATGGAAAGATGGAAATTACTCAGA

GACAGGGAAATAAGTCCTG

GATACCGATTACCGATGGAAAGA

TGGGGAAAGATACCGATTTCG

ATGGAAAGATGGAAATTACTCAGA

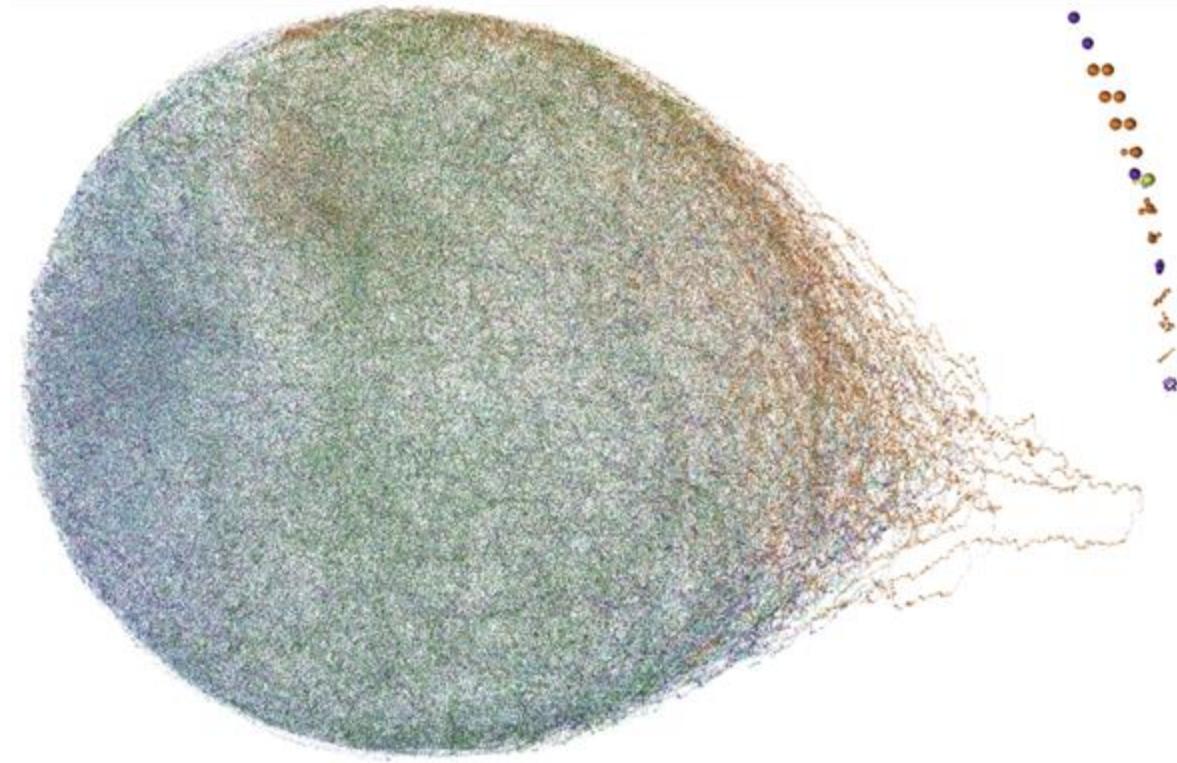
GACAGGGAAAATAAGTCCTG

# Datasets

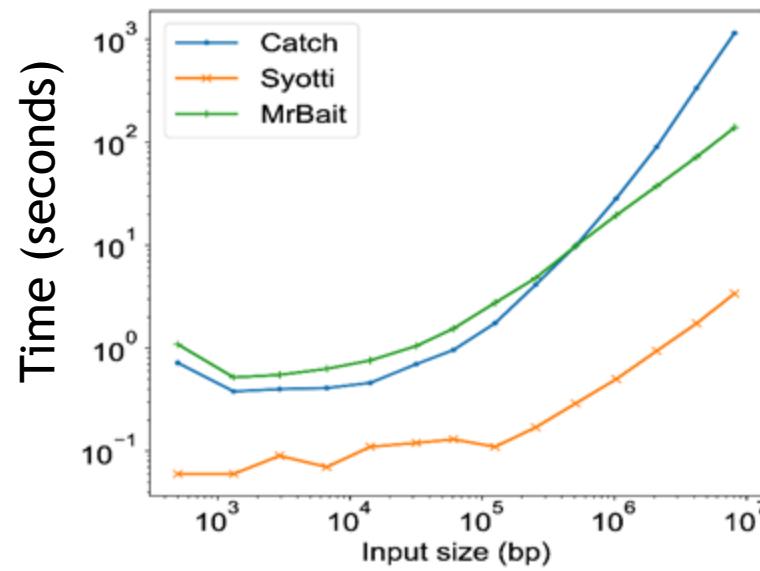
**MEGARes**: 7868 ARGs (**8.1 MB total**)

**Viral**: 420,000 viral genomes from 608 species (**1.3 GB total**)

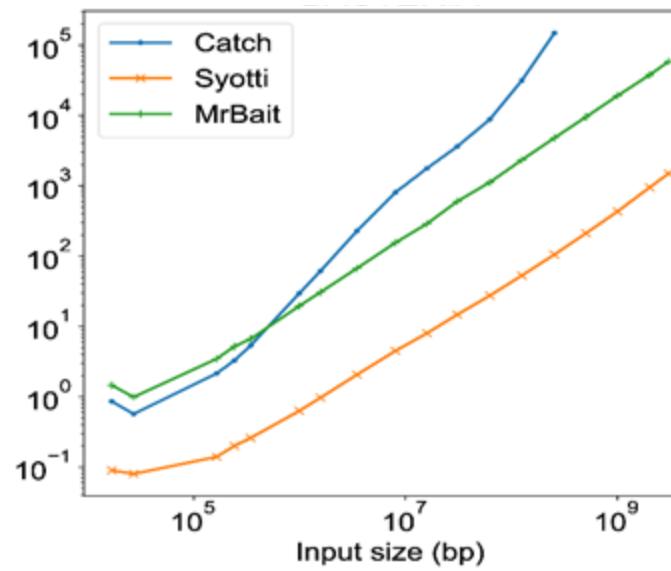
**Bacteria**: 1,000 bacterial genomes from 4 foodborne pathogen species (**3.0 GB total**)



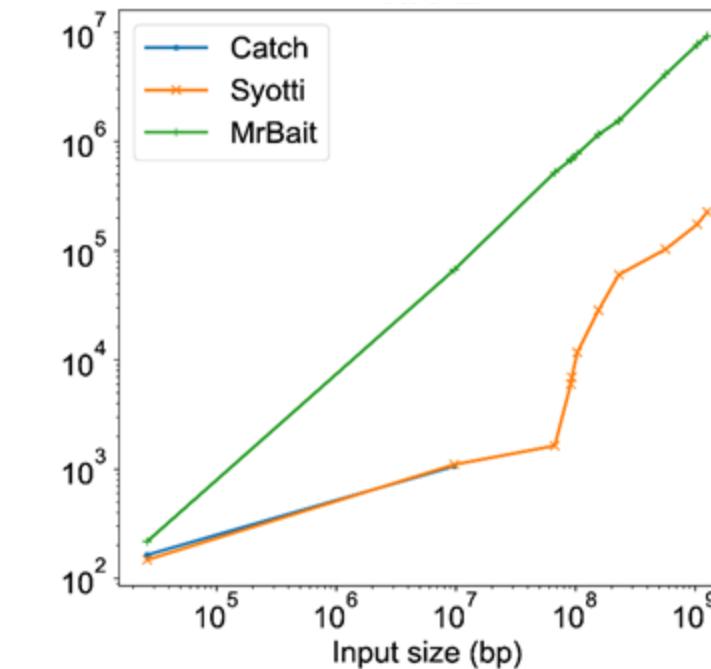
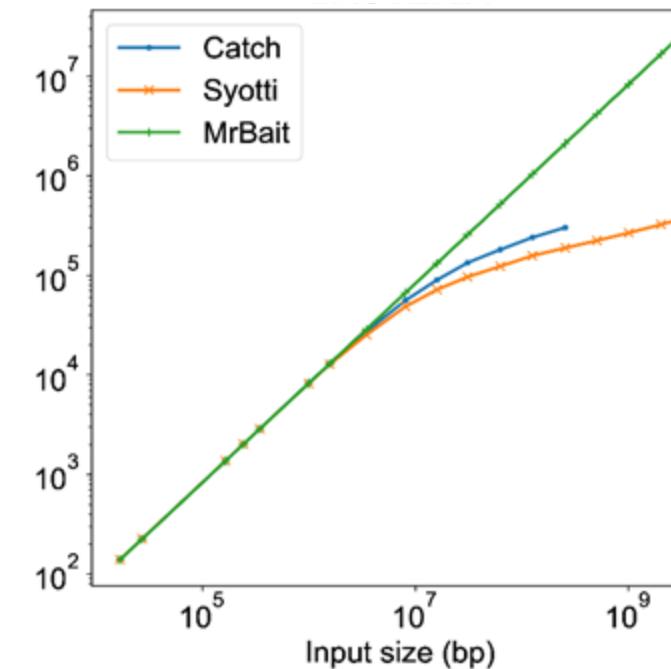
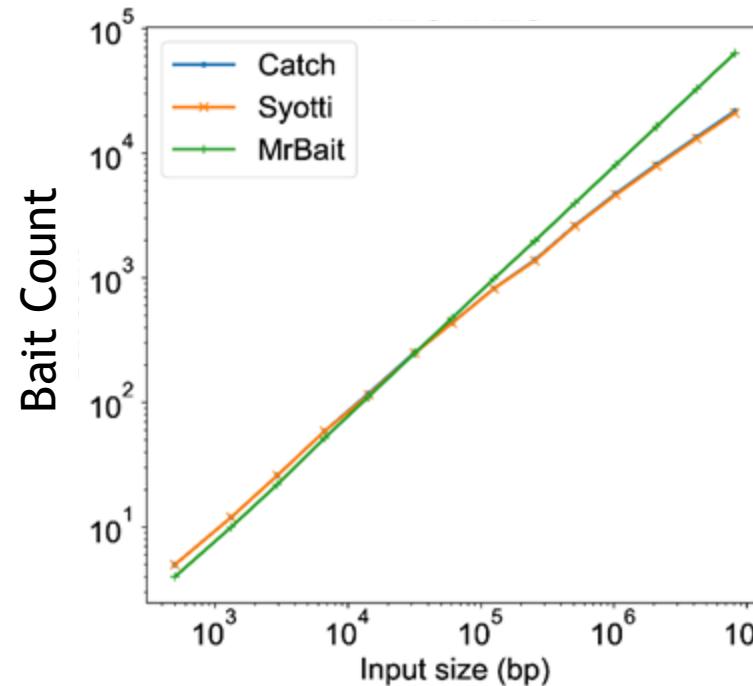
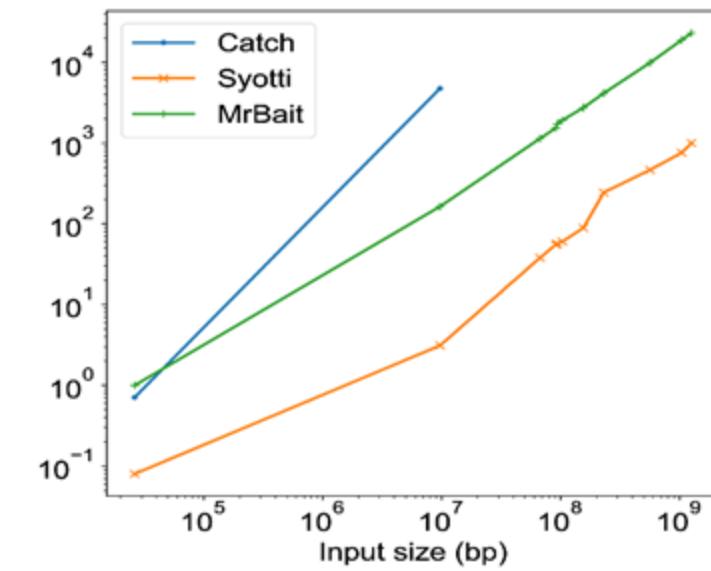
### MEGARes



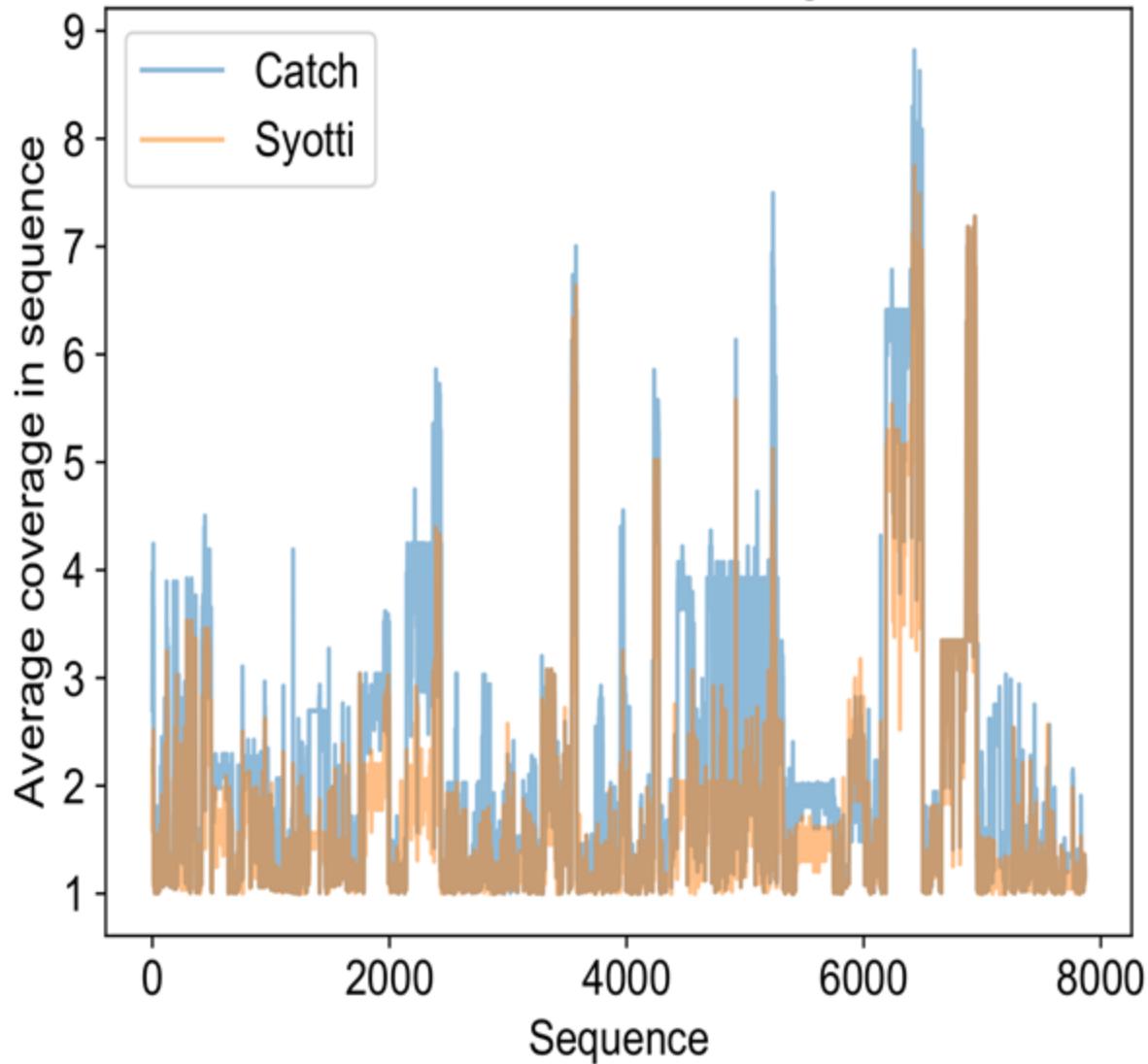
### Bacteria



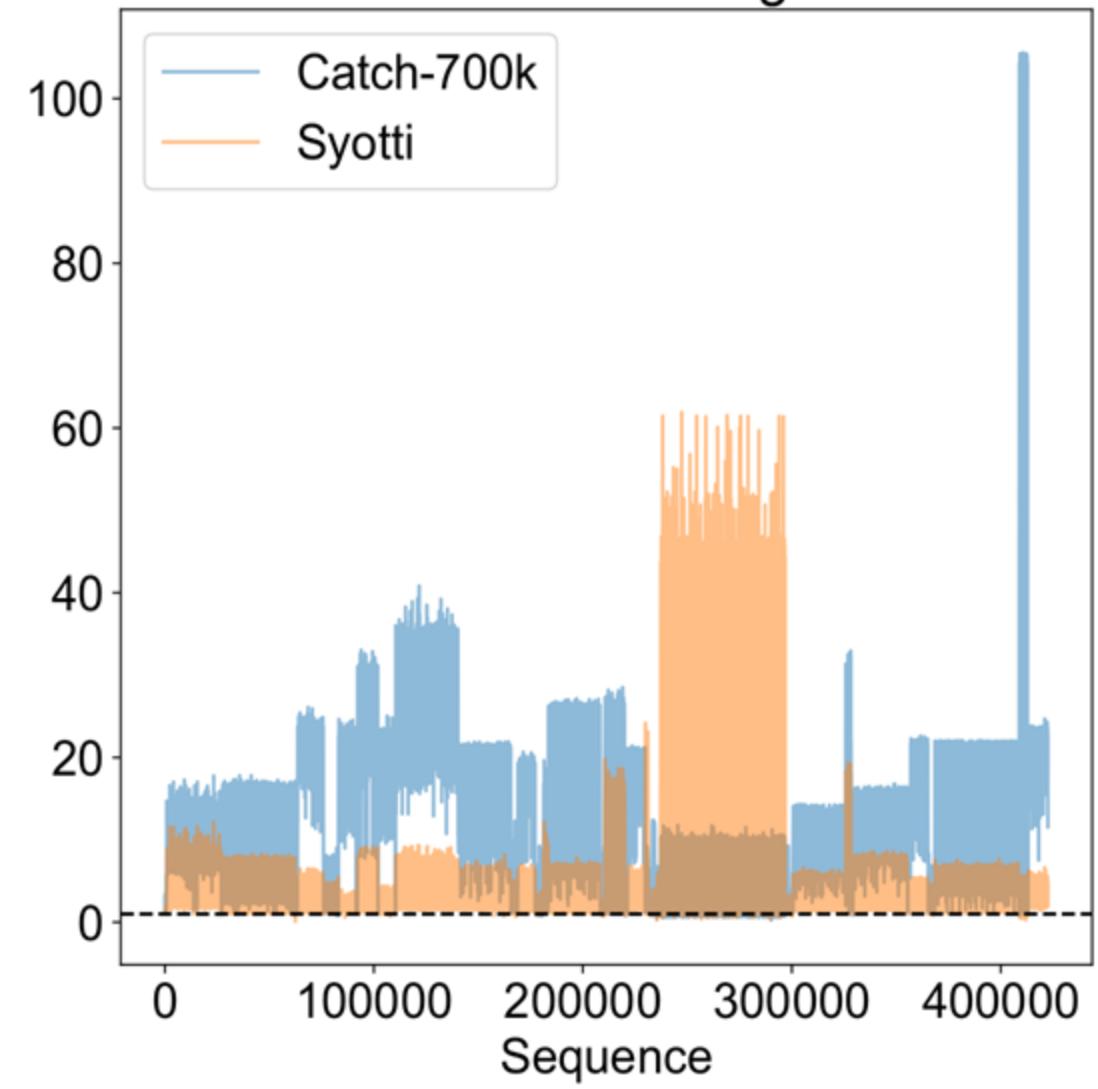
### Viral



### MEGARES Coverage

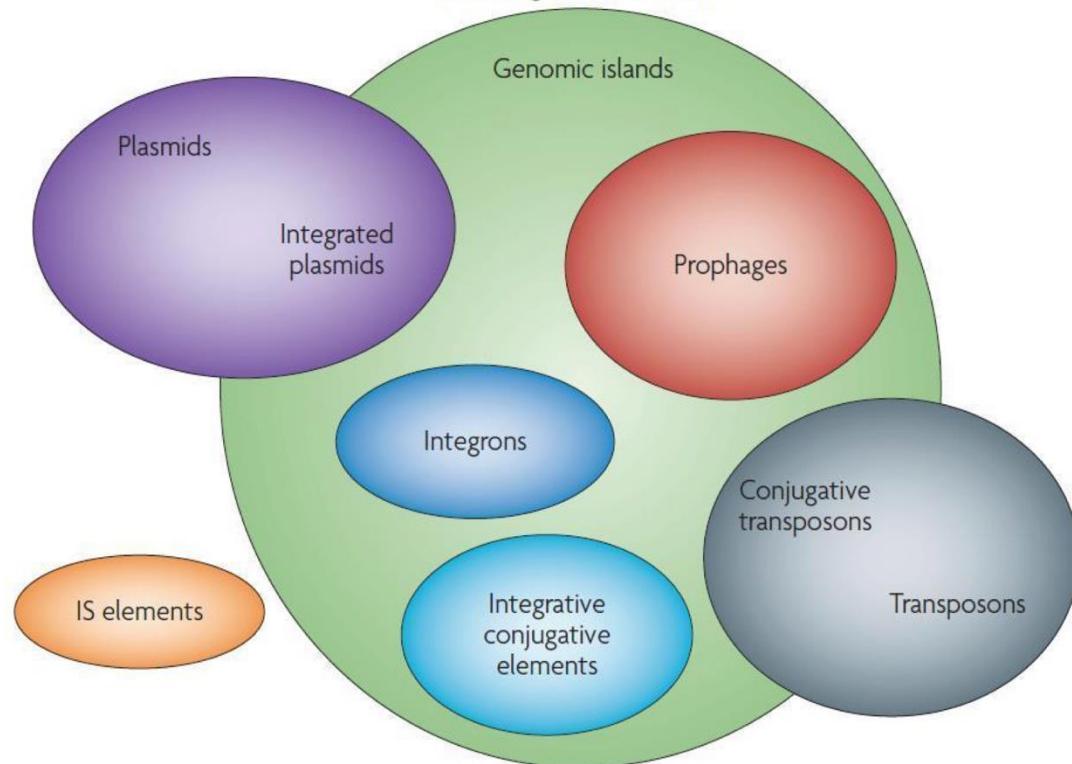


### VIRAL Coverage



# Open Areas for Discovery

# Detection of MGEs is Hard



Databases are overlapping and unorganized (i.e., no hierarchy).

Alignment does not yield a correct answer.

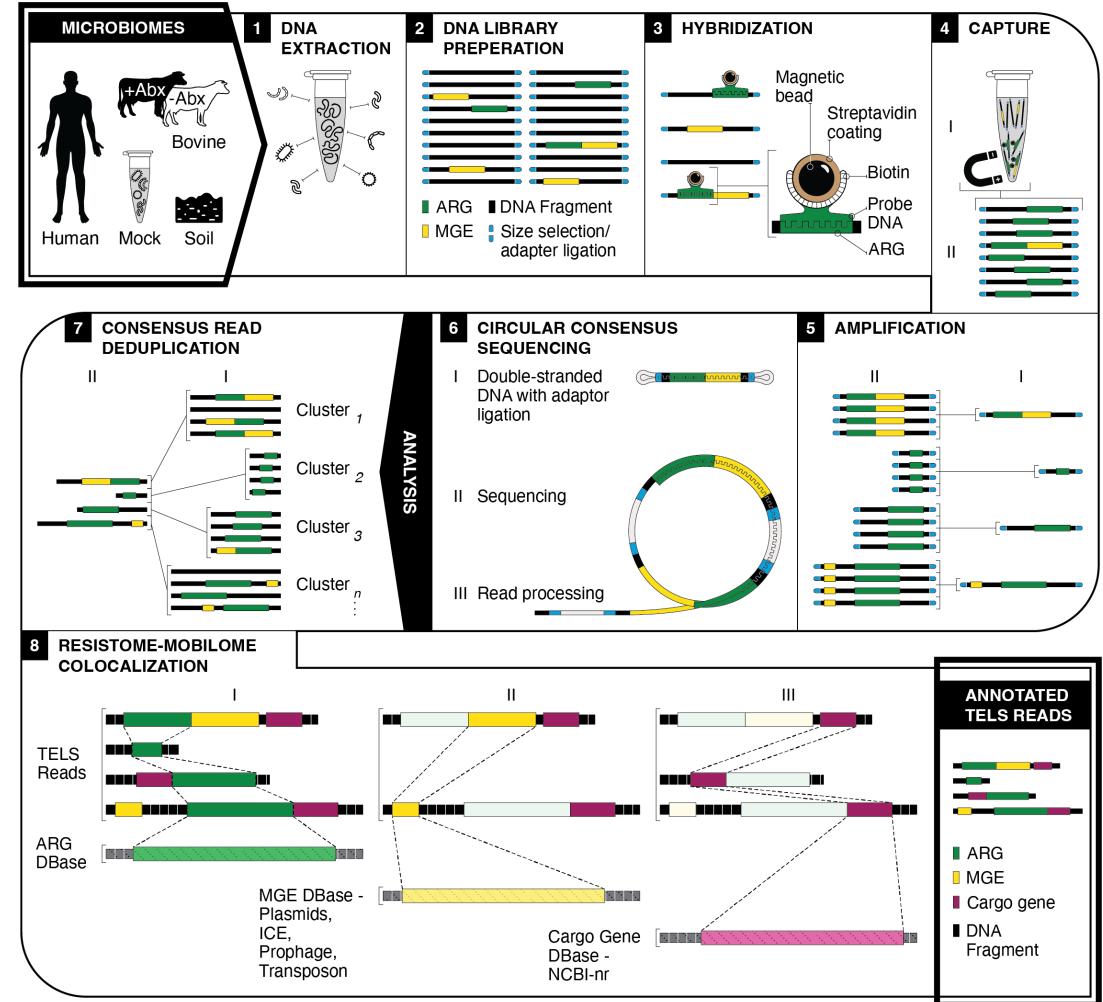
Only knowing broad information is not useful.

# Bait Design Remains Challenging

Can the problem be reformulated?

Does there exist approximate solutions?

Can off-target binding be controlled?



# Thank you

Jarno Alanko



Ilya Slizovskiy



Marco Oliva



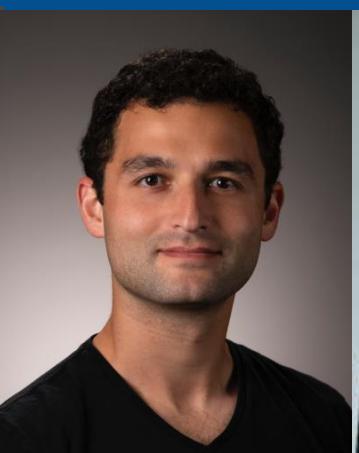
Mattia Prosperi



Travis Gagie



Dan Lokshantov



Noelle Noyes

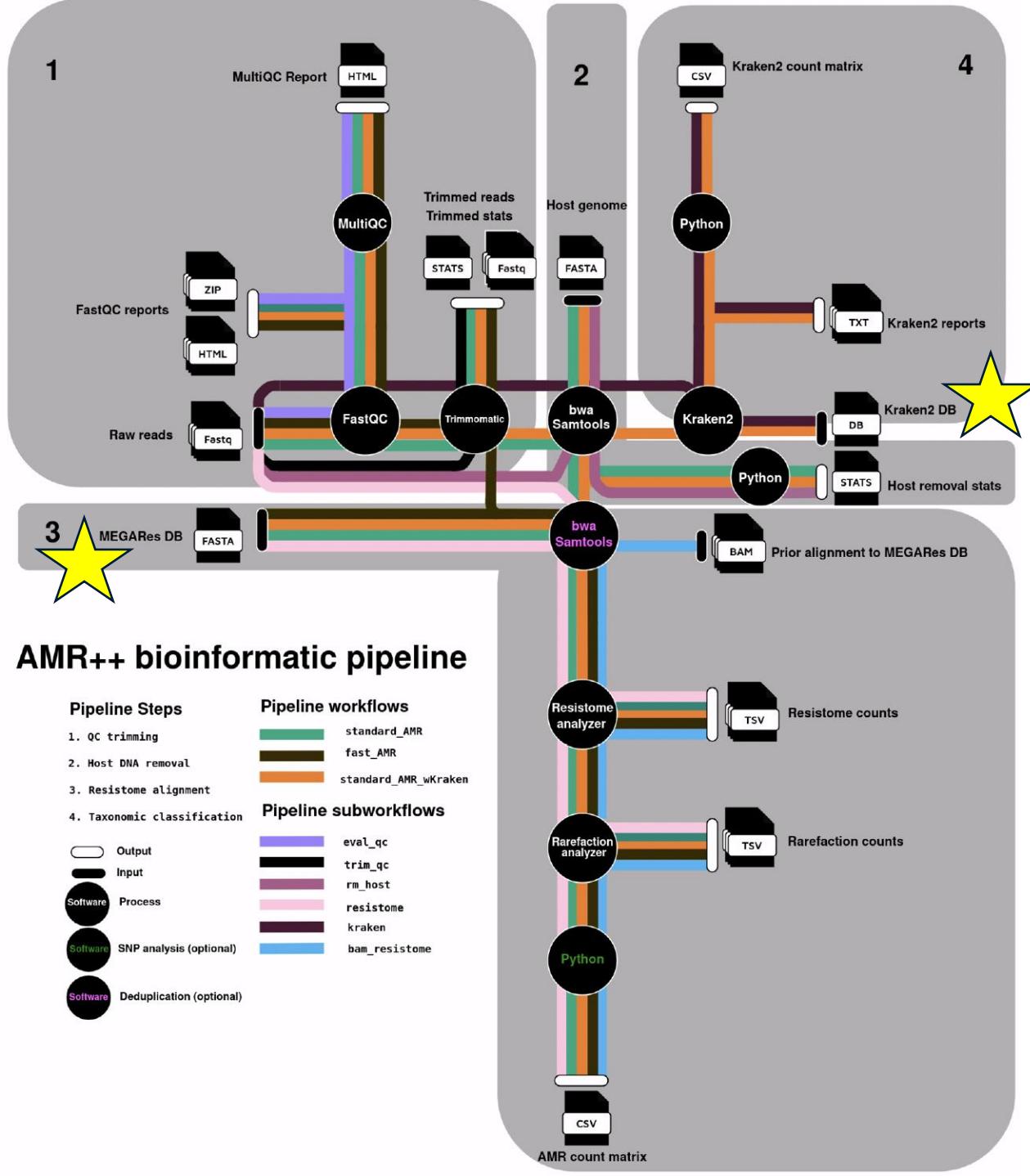


NSF EAGER, NSF SCH, and NSF IIBR

NIH BDMA (grants 308030, 314170, and 323233)



# AMR++ Pipeline



# Part 1: Databases for Resistome Analysis

# It's important to understand your database

- ✓ How was the database developed?
- ✓ How “good” are its sources?
- ✓ How closely has it been checked for inaccuracies or inconsistencies?
- ✓ How comprehensive is it?

# Two main databases in AMR

 Microbial Ecology Group

Home MEGAREs AMR++ Pipeline

## Microbial Ecology Group (MEG)

Research addressing the issues of microbial ecology in animal, public, and environmental health.

MEGAREs AMR++ Pipeline

 COLORADO STATE UNIVERSITY

 ATM | WT  
VERO  
VETERINARY EDUCATION,  
RESEARCH, & OUTREACH

 UF UNIVERSITY of FLORIDA

 UNIVERSITY OF MINNESOTA



 Search

## CARD

Use or Download Copyright & Disclaimer  
Help Us Curate #AMRCuration #WorkTogether

Browse Analyze Download About

Search

### The Comprehensive Antibiotic Resistance Database

A bioinformatic database of resistance genes, their products and associated phenotypes.  
8582 Ontology Terms, 6442 Reference Sequences, 4480 SNPs, 3354 Publications, 6480 AMR Detection Models  
Resistome predictions: 414 pathogens, 24291 chromosomes, 2662 genomic islands, 48212 plasmids, 172216 WGS assemblies, 279120 alleles  
YouTube: Canadian Bioinformatics Workshops 2024: Antimicrobial Resistant Gene (AMR) Analysis

**Browse**  
The CARD is a rigorously curated collection of characterized, peer-reviewed resistance determinants and associated antibiotics, organized by the Antibiotic Resistance Ontology (ARO) and AMR gene detection models.

**Analyze**  
The CARD includes tools for analysis of molecular sequences, including BLAST and the Resistance Gene Identifier (RGI) software for prediction of resistome based on homology and SNP models.

**Download**  
CARD data and ontologies can be downloaded in a number of formats, including lists of mutations and molecules with corresponding metadata and citations. RGI software is available as a command-line tool.

# *So how is MEGARes different then?*

The real difference lies in the annotation structure!

*(Plus, there's a lot of manual curation)*

# Ontology of Microbiome Data

PHYLUM

Proteobacteria

CLASS

Gammaproteobacteria

ORDER

Enterobacteriales

FAMILY

Enterobacteriaceae

GENUS

*Escherichia*

SPECIES

*Escherichia coli*

# MEGARes 3.0



Type: the type of antimicrobial compound, e.g. drugs, biocides, multi-compound



Class: the major antimicrobial chemical class, e.g. betalactams, aminoglycosides



Mechanism: the biological mechanism of resistance, e.g. penicillin binding protein



Group: the gene- or operon-level group for that sequence, e.g. SHV betalactamase, MCR-1

# Ontology of Resistome Data

TYPE

Drug

CLASS

Beta-lactam

MECHANISM

Class A beta-lactamases

GROUP

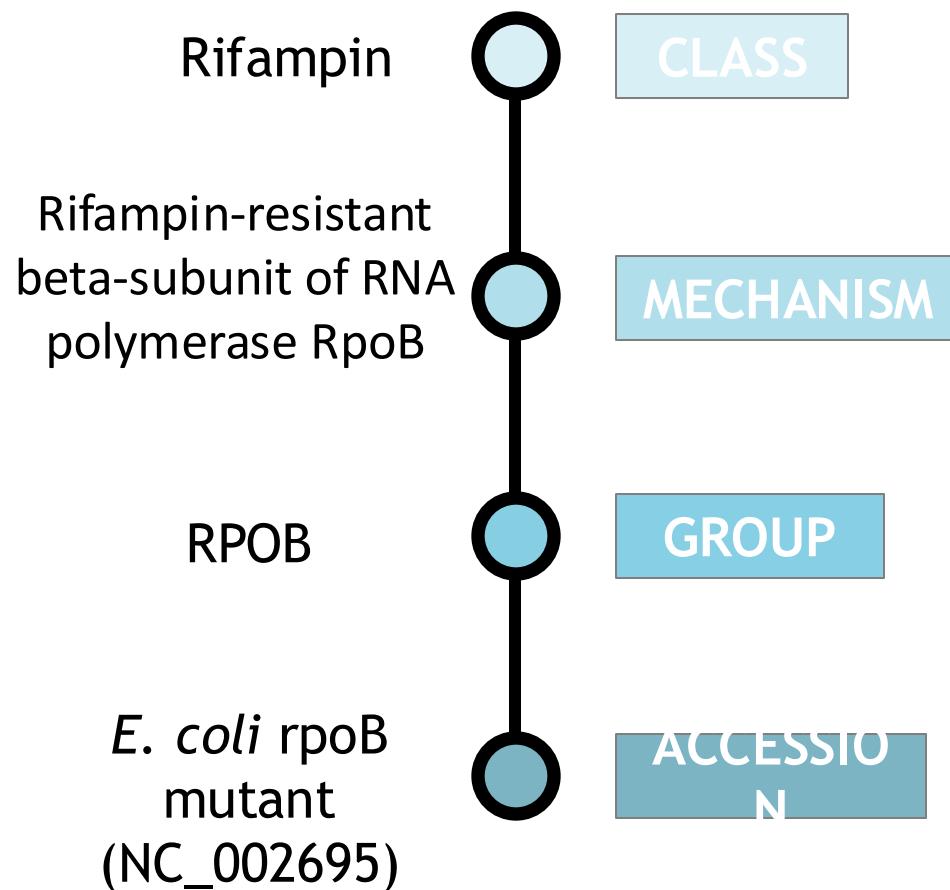
TEM

ACCESSION

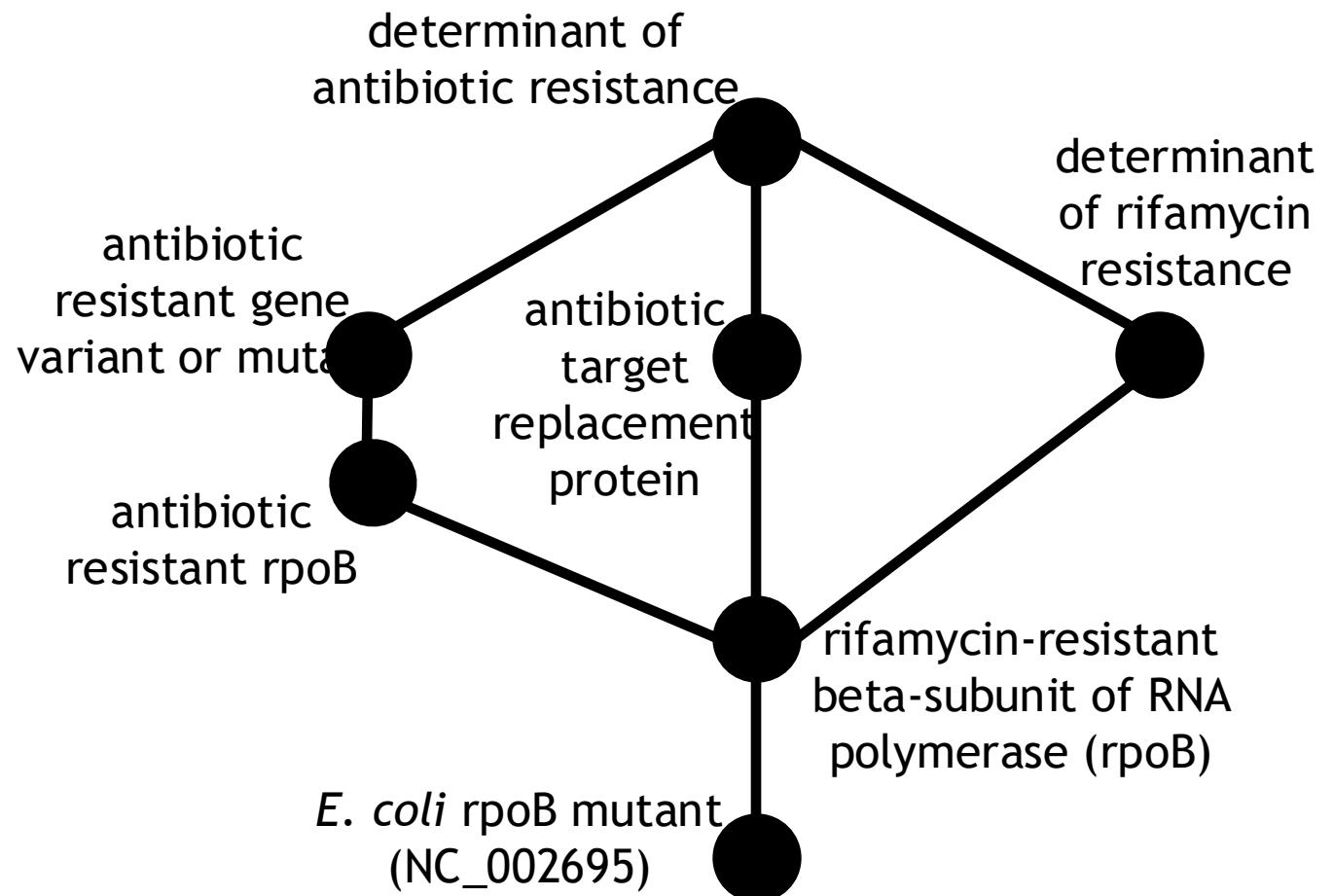
*bla*<sub>TEM-144</sub> (DQ256080)

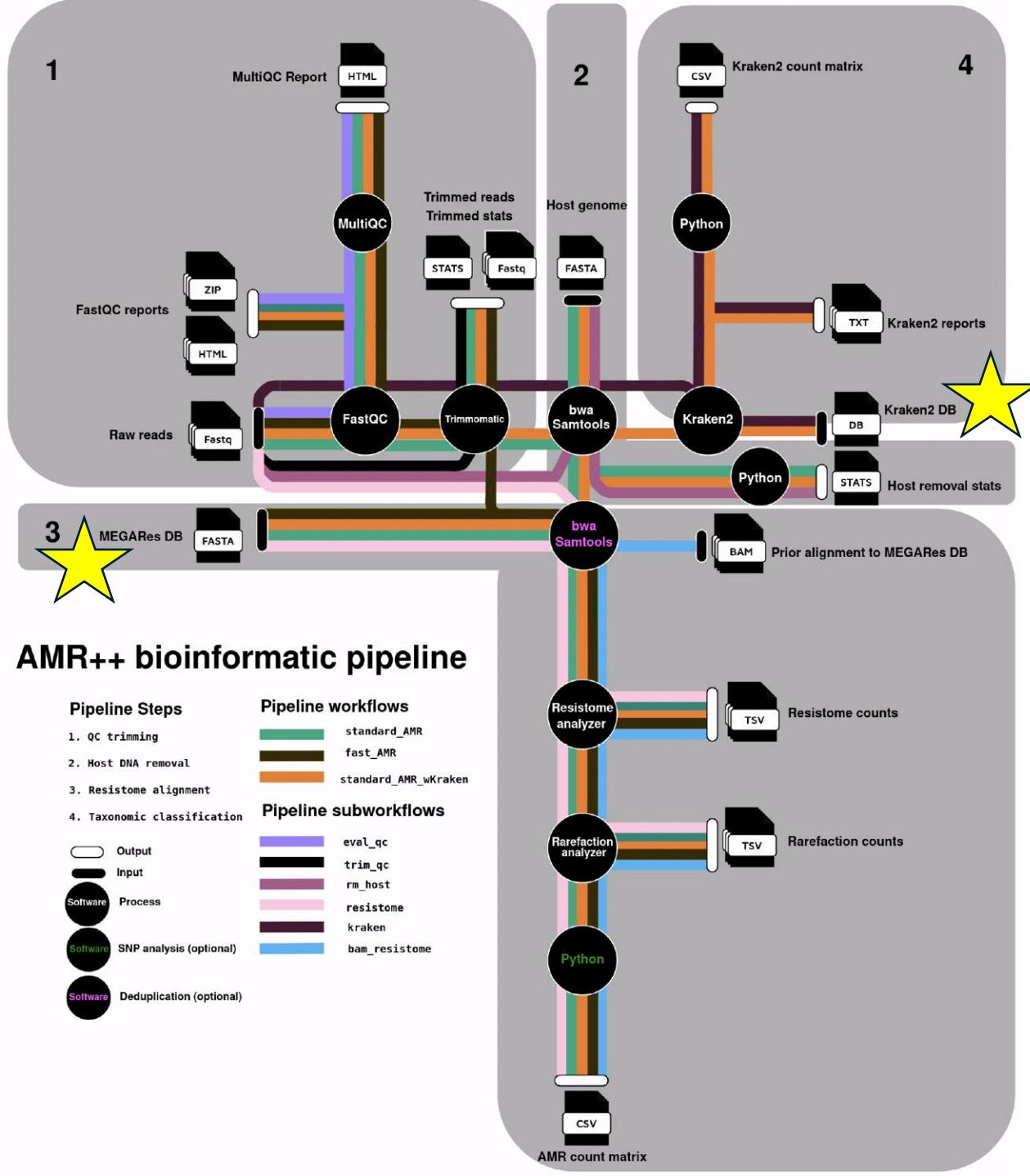
# MEGARes ontology is acyclical

## MEGARes



## CARD







Search MEGARes by keyword or annotation term



- [Acetate resistance](#)
- [Acid resistance](#)
- [Aluminum resistance](#)
- [Aminocoumarins](#)
- [Aminoglycosides](#)
- [Arsenic resistance](#)
- [Bacitracin](#)
- [betalactams](#)
- [Biguanide resistance](#)
- [Biocide and metal resistance](#)
- [Cadmium resistance](#)
- [Cationic antimicrobial peptides](#)
- [Chromium resistance](#)
- [Cobalt resistance](#)
- [Copper resistance](#)



[Browse](#) / Multi-drug resistance

## Multi-drug resistance

We define multi-drug resistance as genes and mechanisms that cause resistance to two or more different antibiotic classes. Typically, such mechanisms involve active extrusion of antibiotic molecules from the bacterial cell or mechanisms that prevent the drug from reaching its target.

### References

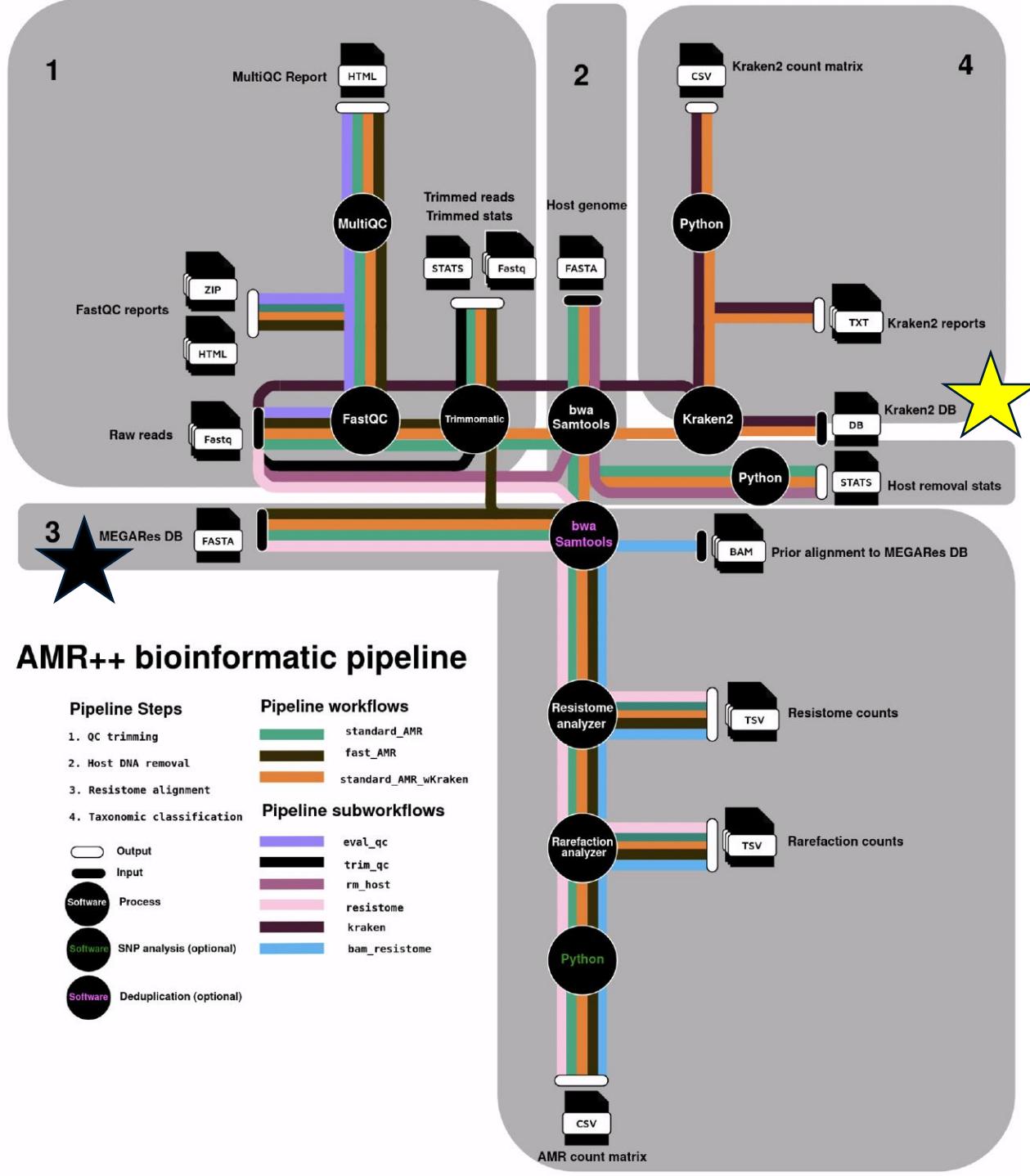
- <http://jnci.oxfordjournals.org/content/92/16/1295.full>

### Mechanisms

- [MDR 23S ribosomal RNA methyltransferase](#)
- [MDR 23S rRNA mutation](#)
- [MDR acetyltransferase](#)
- [MDR mutant porin proteins](#)
- [MDR regulator](#)
- [Multi-drug ABC efflux pumps](#)

>>MEG\_2694|Drugs|Multi-drug\_resistance|Multi-drug\_ABC\_efflux\_pumps|EATAV|RequiresSNPConfirmation

ATGTCTAAAATCGAAATAAAAATCTGACATTGGCTACGACAGCCAAGGCACATTATTATTGAACAAGCAAATCTAAATTGACACACAATGGAAACTAG  
GACTTATCGGACGAAACGGTCGAGGAAAGACAACCTTACTGAATATTCTACAAAACAAACTACCTTACAGGGCAAGTAATCCATCAGCAAGAATTGCCTA  
TTTCCCACAGACAAAAGATAAAAGAACGTTAACCTATTACGTGTTAAATGATATTACGGATTGAGATATGGGAAATCGAAAGAGAGCTCCAATTGATG  
CAAACAGATCCTGAAATCTTATGGAGAGAATTGAGCACACTATCGGGGGAGAGAAGACAAAAGTCCTACTGGCATTATTGTGGATGACACTCATTCC  
CGTTAATCGATGAACCAACGAATCATTGGATACTCTGGTAGAAAACAAGTAGCGGCTATTGAAAAAGAAAAACAGGCTCATCGTGGTCAGCCATGA  
CCGGGGATTATCGATGAAGTAGTGGACCATGTTAGCAATCGAAAAAGTCAACTGGAACTTTATCAAGGGATTCTATCTATGAAAGAACAGAAAAAA  
CTTCGTGATGAATTGAAATGGCTAAAATGAAAAAGAAGTCAGTAGGCTAAAGAAAACAGCAGCTGAAAAGCCGAATGGCTCGTCCCGAG  
AAGGAGATAAAACAAAGAAACAGTCGGATTCACTGATACTGAATCTAGACGAGTGAATAAGGAGCAGTGGTGCTGATGCTGCACGGACGATGAAACGATC  
CAAAGCAATCGTAATCGGATGGAGACCCAGATCAGCGAGAAAGAAAAACTATTAAAAGATATCGAATATATCGATTGACGATGAATAGCCAAGCGTCT  
CACCATAAGCGACTTTAAGCGTAGAAGATCTCAATTAGGGTATGAAAATCTGTTATTGAGCCAATTCACTTACAATCGAGCCTCATCAGCGGGTGGCGA  
TTTCAGGTCTAACGGTGCAGGAAAGTCATCCATTATCCATTATCTCTGGGGCATTCAACGGCAAGGTATAGGAGAAAATACCAGCAAAACATCTGAG  
CATTAGTTATGCAAGCCAAATTATGAAGACAATCGAGGAACGTTGGCGGAATTGCAGAGAAAACCAAGTAGACTACCAAGCATTGAAACAACCTCCGA  
AAGCTTGGATGGAAAGAGATGTTTCATAACAAGATCGAGCAGATGAGTATGGCCAACGGAAAAAGTGGATTGGCTAAATCTTATCACAGCCAGCTG  
AACTATATATGGGATGAACCATTGAATTATTGGATGTCTTCAATCAAGAACAAATTAGAACAACTGATCTGAACGTGAAACCTGCCATGTTACTAGTGG  
ACATGATCAAACCTTCTGGATAAAGTATCTACTGAGATTATTCTTCTTGAGAGAAATCTAA





## RefSeq: NCBI Reference Sequence Database

A comprehensive, integrated, non-redundant, well-annotated set of reference sequences including genomic, transcript, and protein.

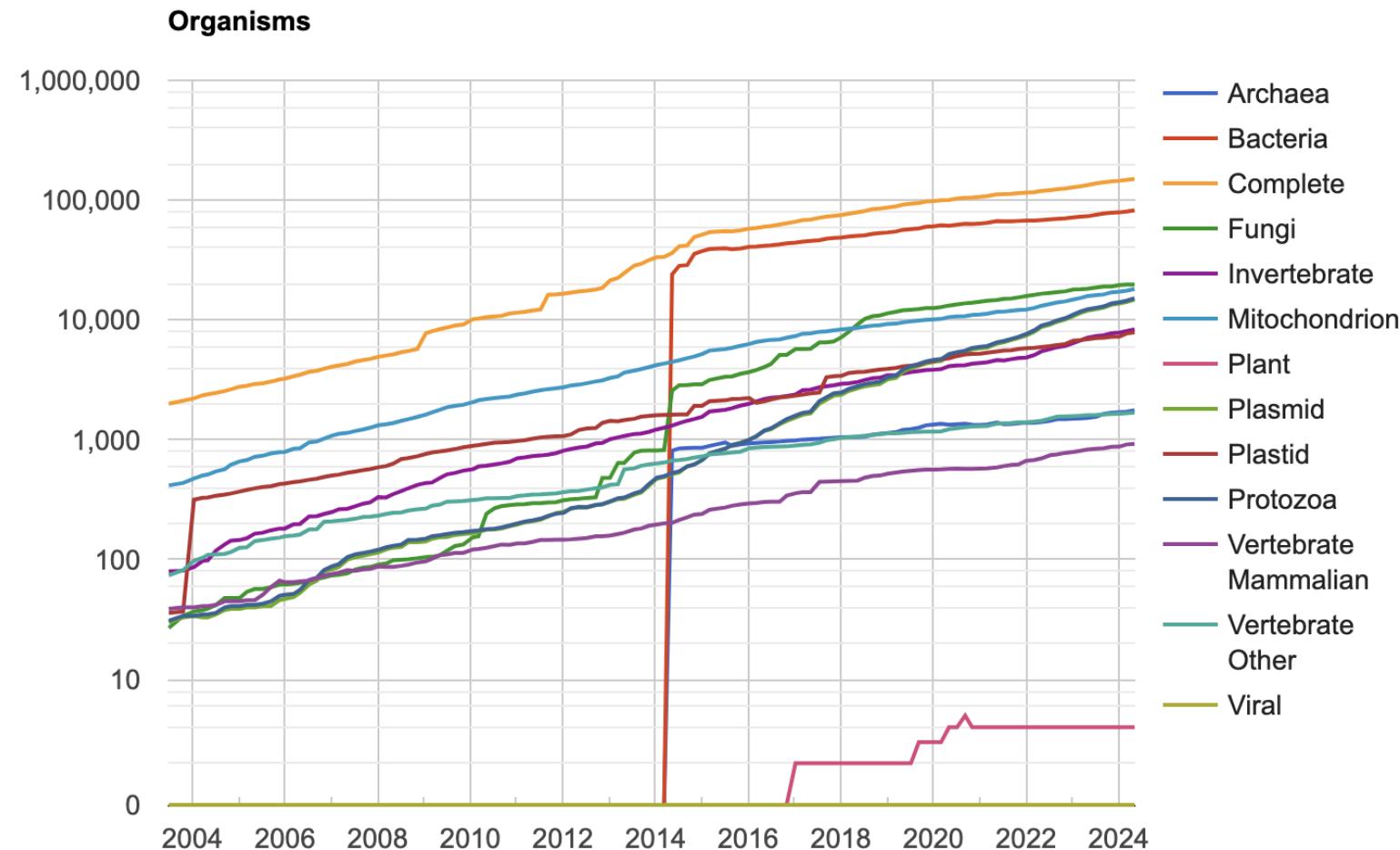
“Comprehensive” - it’s relative

“Integrated” -- ??

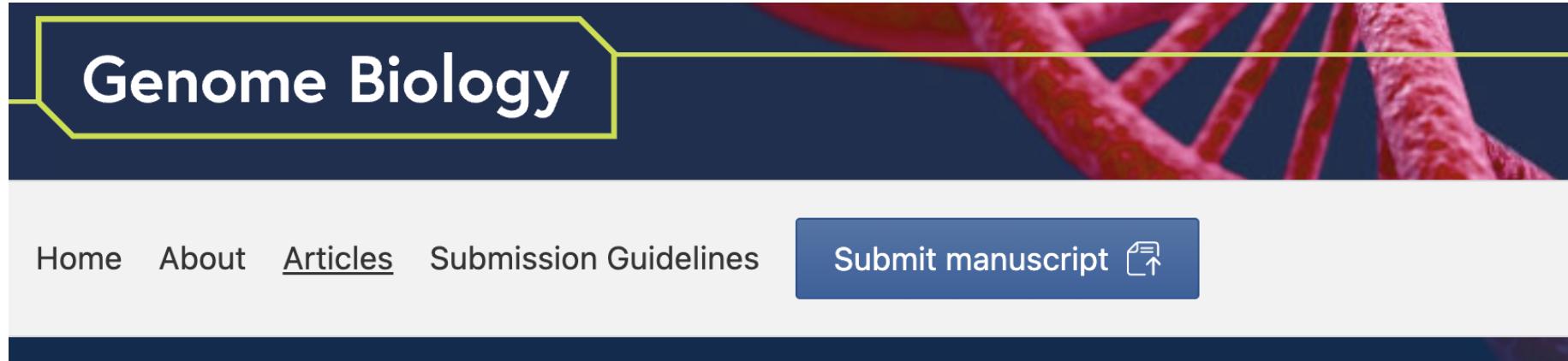
“Non-redundant” -- ??

“Well-annotated” -- not really

# RefSeq Bacteria Growth Statistics



Database size does matter, and there's good research on this:



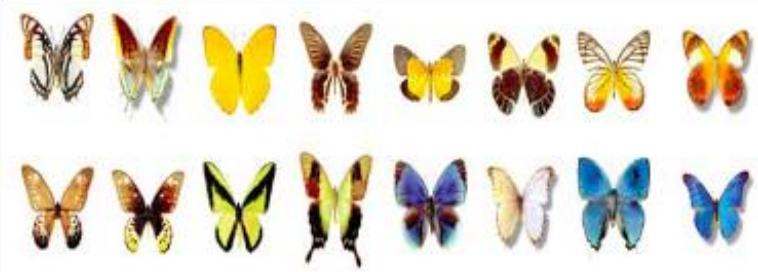
Open Letter | [Open access](#) | Published: 30 October 2018

## RefSeq database growth influences the accuracy of $k$ -mer-based lowest common ancestor species identification

[Daniel J. Nasko](#), [Sergey Koren](#), [Adam M. Phillippy](#) & [Todd J. Treangen](#)

[Genome Biology](#) 19, Article number: 165 (2018) | [Cite this article](#)

# The NCBI Taxonomy Database



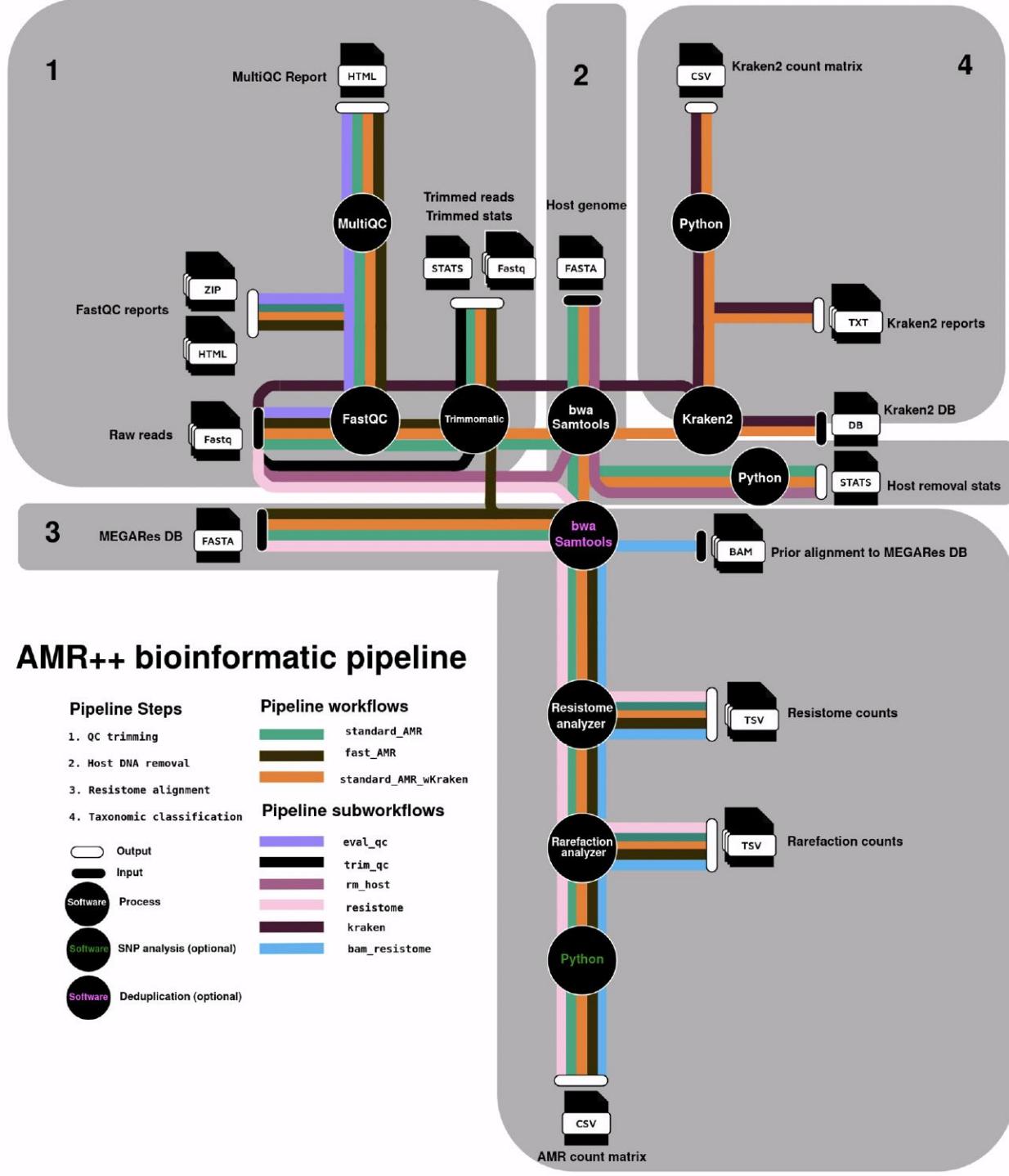
## Taxonomy

The Taxonomy Database is a curated classification and nomenclature for all of the organisms in the public sequence databases. This currently represents about 10% of the described species of life on the planet.

Also has its limitations.

In the end, keep in mind that taxonomies are attempting to bin DNA sequences into containers with “hard boundaries”.

# Part 2: Quality assessment of shotgun metagenomic data



# Why do we need to do quality assessment?

1. Because sequence data can have errors – both expected and unexpected
2. Because sequence data usually contain both technical and biological sequence

# Example of adapter read-through

## Adapter Content

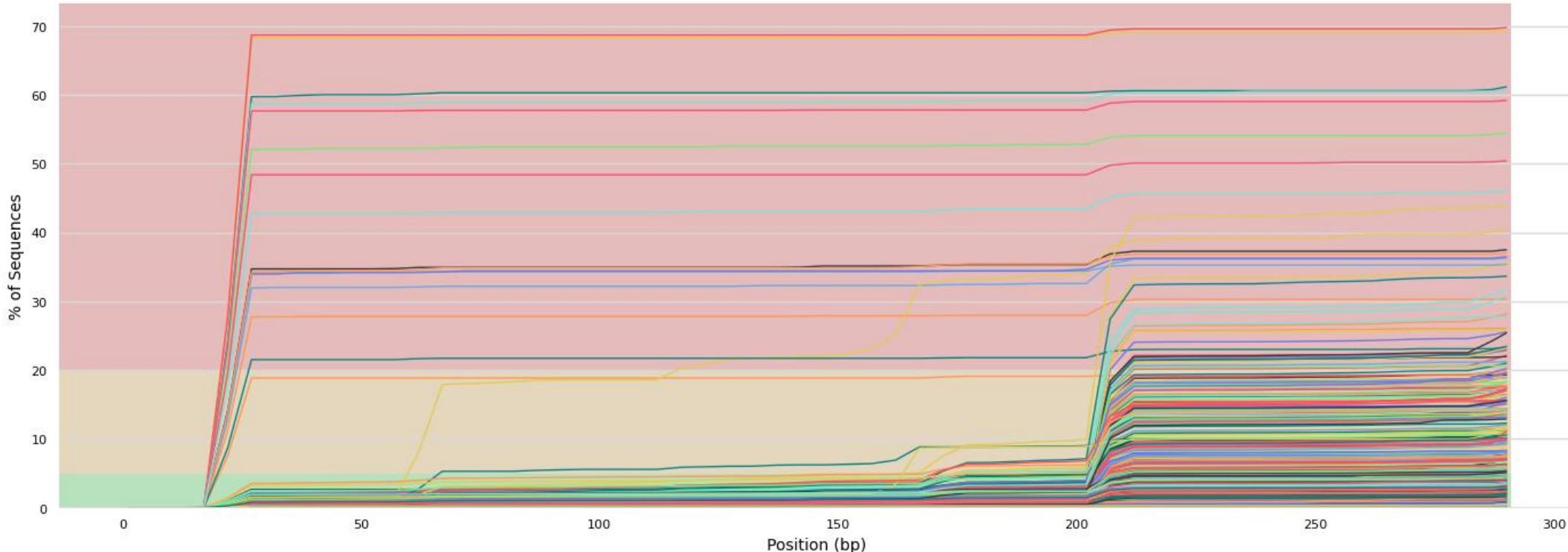
72 47 115

Help

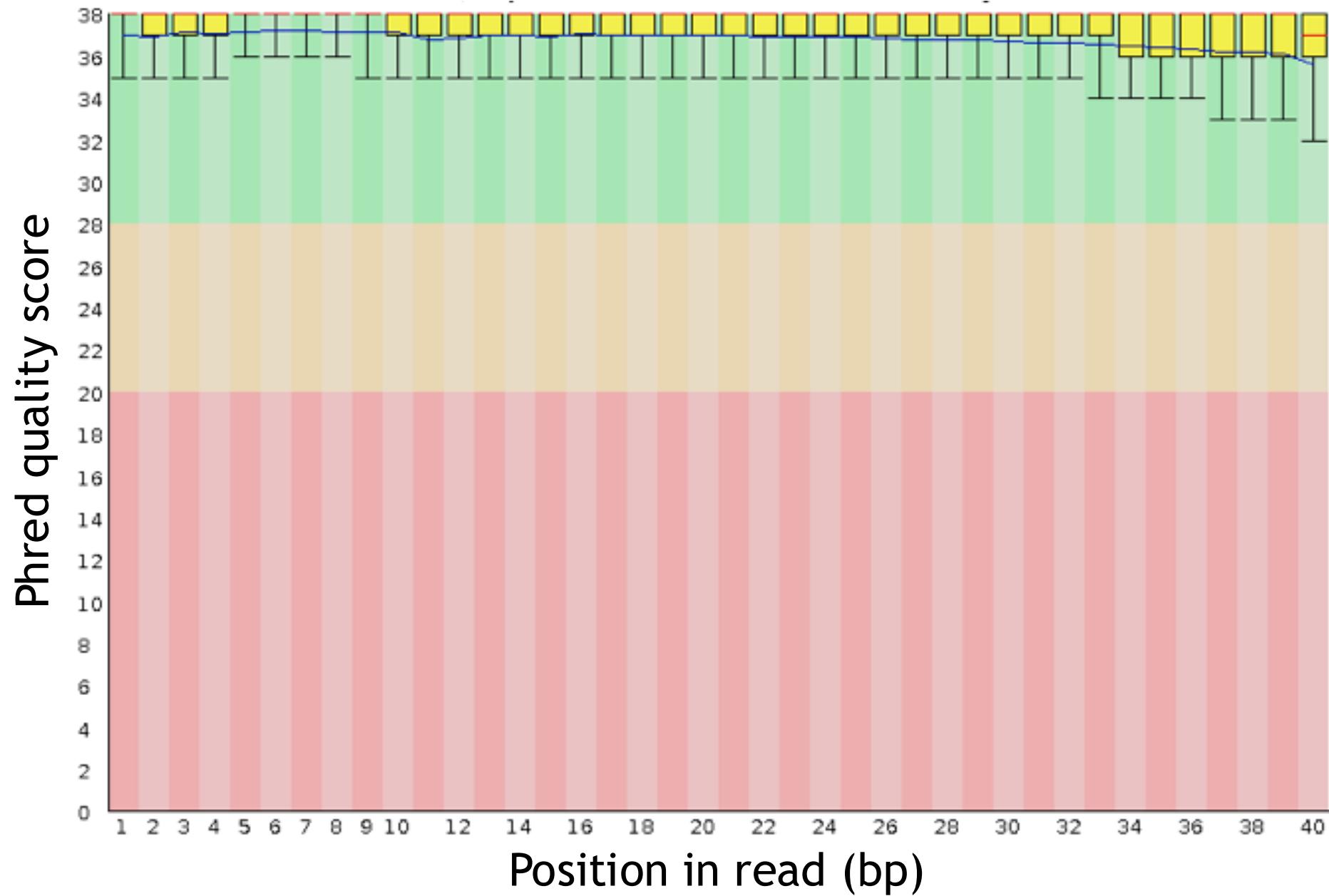
The cumulative percentage count of the proportion of your library which has seen each of the adapter sequences at each position.

Flat image plot. Toolbox functions such as highlighting / hiding samples will not work (see the docs).

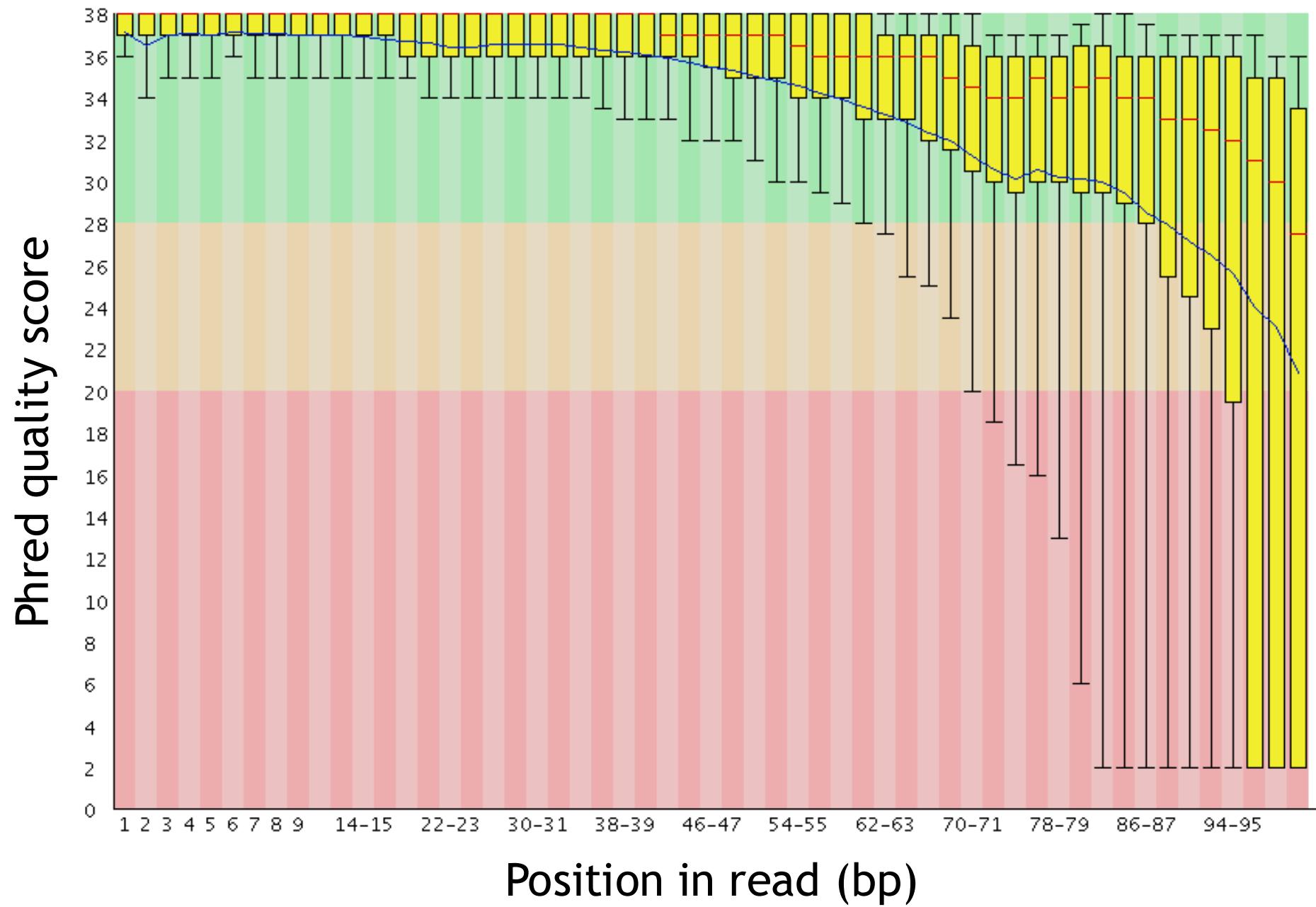
FastQC: Adapter Content



# FastQC Output: Good Illumina Data



# FastQC Output: Bad Illumina Data



# Trimming of FASTQ Sequences

Trimming corrects several error types:

- Illumina systematic bias (end of read base substitutions)
- Adapter sequences
- Missing mate-pairs
- Low quality reads



# Filtering of Off-Target Reads

- Samples collected from hosts contain host DNA
- Known contaminants might be present due to library preparation techniques

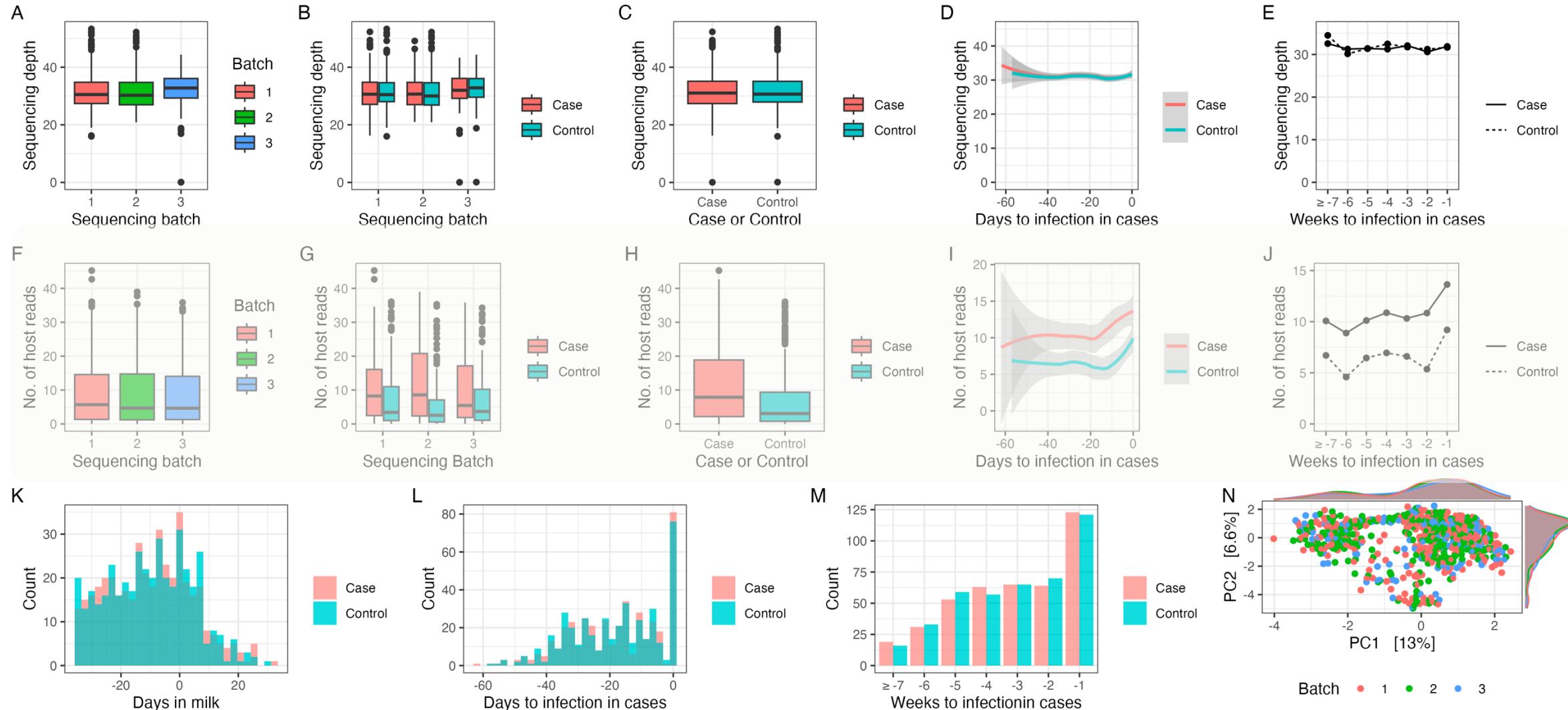


# Filtering of Off-Target Reads

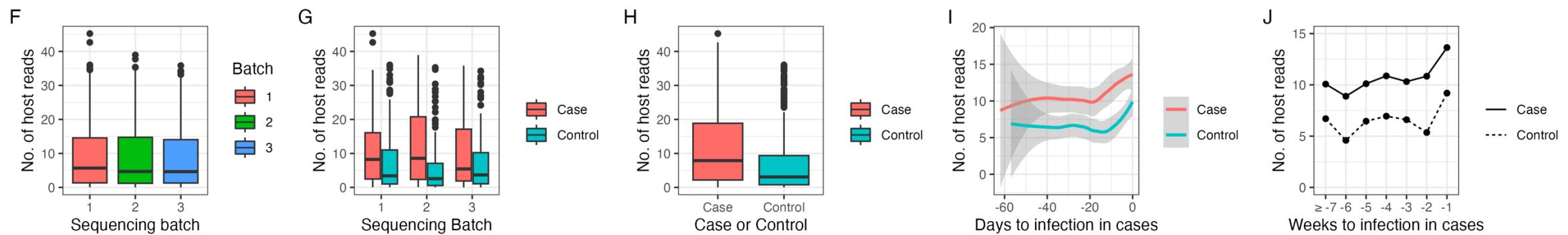
- Filtering host DNA removes the chance of false positive classification
- AMR++ filtering:
  - Single filtering genome
  - Multiple filtering genome



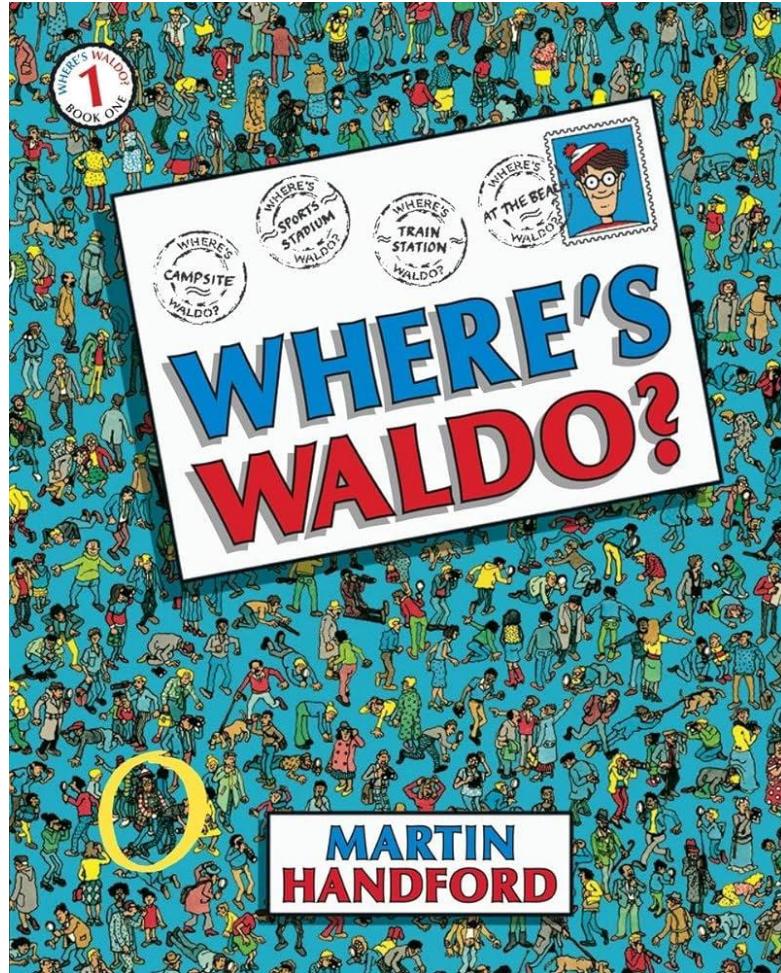
# This can uncover some interesting findings!

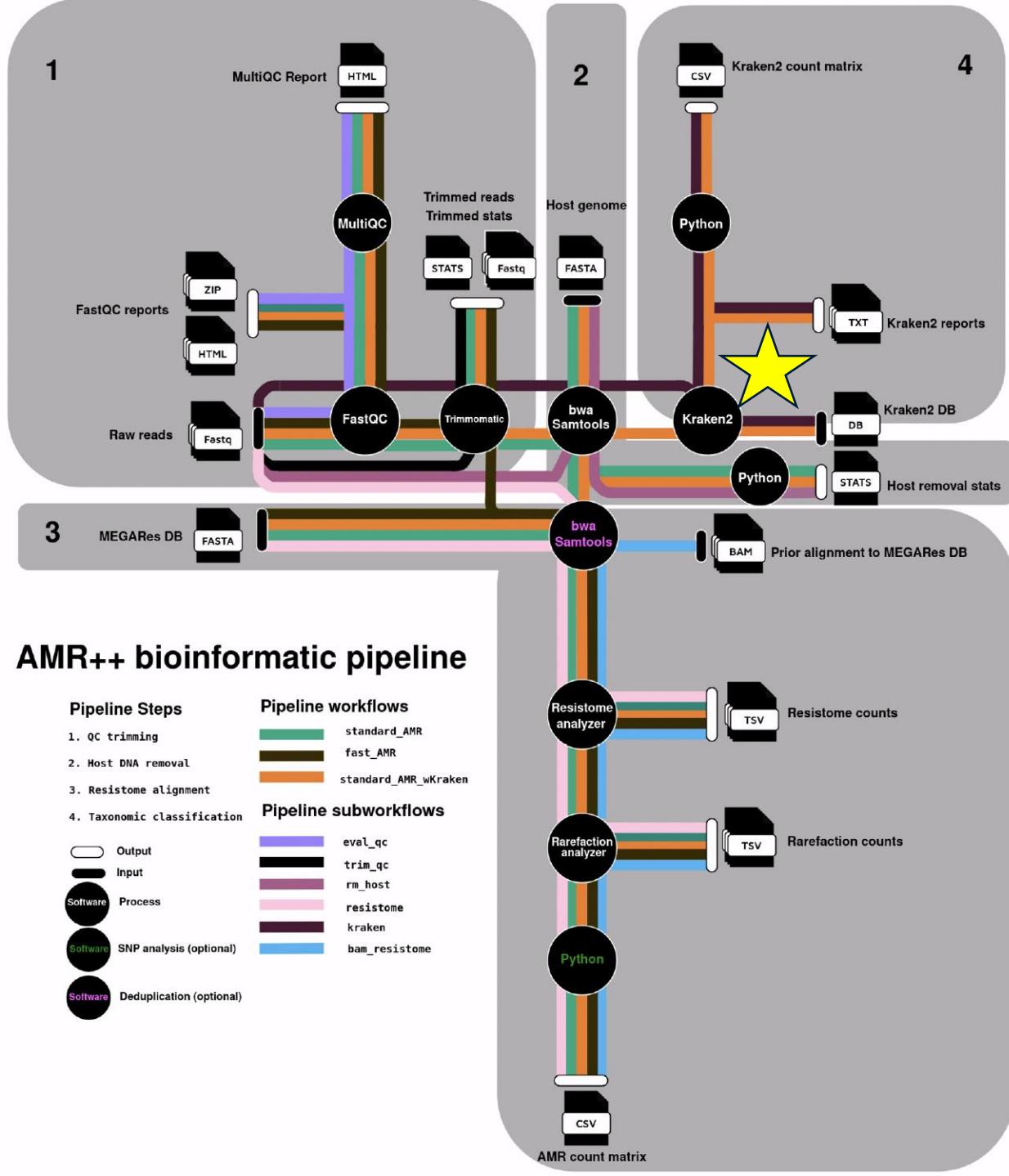


# This can uncover some interesting findings!



# Part 3: WHO is there?





Let's assume that we want to  
create a simple microbiome  
profile for 1 sample

Let's assume that we want to create a simple microbiome profile for 1 sample

Basically, we want to know:

***WHO*** is there?

What is their ***RELATIVE ABUNDANCE***?

# Let's start with WHO is there...

# Two approaches for determining who:

1. Marker-based
2. Genome-based

# Marker-based methods = “16S on steroids”

Examples:

[MetaPhlAn](#)

[PhyloSift](#)

And many others...

# Two approaches for determining who:

1. Marker-based
2. Genome-based

# How could the metagenome composition be used to bin and classify sequence data?

- GC-content, nucleotide frequencies, codon patterns
- $k$ -mer based
  - a. exact-match
  - b. profile

# Each Sequence has a k-mer Profile

For every sequence, there exists a profile of k-mers  
e.g., Sequence: ACGTGTG

3-mer profile:

ACG : 1

CGT : 1

GTG : 2

TGT : 1

# $k$ -mer Profiles Contain Information

- $k$ -mer profiles can help to identify a sequence
- Match the  $k$ -mer profile of unknown sequences to known  $k$ -mer profiles
- This is a classification problem:
  - We have an unknown observed sequence
  - We have known genomes
  - We can use matching of  $k$ -mer profiles to label the unknown sequence

# *What is Kraken?*

Kraken is actually a program that tries to match your sequence data to all of the genomes in RefSeq (*more on this later*)...

**So really, it's the NCBI RefSeq database,**  
“manipulated” in a special way

# *What is Kraken?*

Kraken pulls all of the prokaryotic and viral genomes from RefSeq and does some algorithmic maneuvering to create the “Kraken-DB”

Then, it uses the NCBI taxonomy to name everything.

# *What is Kraken?*

- This is kraken’s “database”; it is a database of  $k$ -mer profiles created from NCBI RefSeq genomes, which are also linked to NCBI’s taxonomy
- This database is large, and it requires a lot of computational resources to build it; however, once it is built, it can “sit” on your computer to be used over and over again (can also be updated periodically)

# How is it created (high level)?

1. Chop up all of the genomes in RefSeq
2. Compare all of the pieces to one another
  - a. Pieces ***without*** a match represent the RefSeq genome "as-is" (strain, species, genus, etc....)
  - b. Pieces ***with*** match(es) get moved up the taxonomy until they are unique within a "level" (the "lowest common ancestor" or LCA approach)

# Let's do an example....

After chopping up all RefSeq genomes into 4-mers, we have identified that the 4-mers ACCT and GGT<sub>A</sub> are contained within these genomes, and we want to uniquely associate them with taxa in the NCBI taxonomy.

Kingdom

Phylum

Class

Order

Family

Genus

Species

# Which species is it?



ACCT

*Salmonella enterica* ssp.  
*arizona*

ACCT

*Salmonella enterica* ssp.  
*enterica*

ACCT

*Salmonella bongori*

GGTA

*Escherichia coli*

Kingdom

Phylum

Class

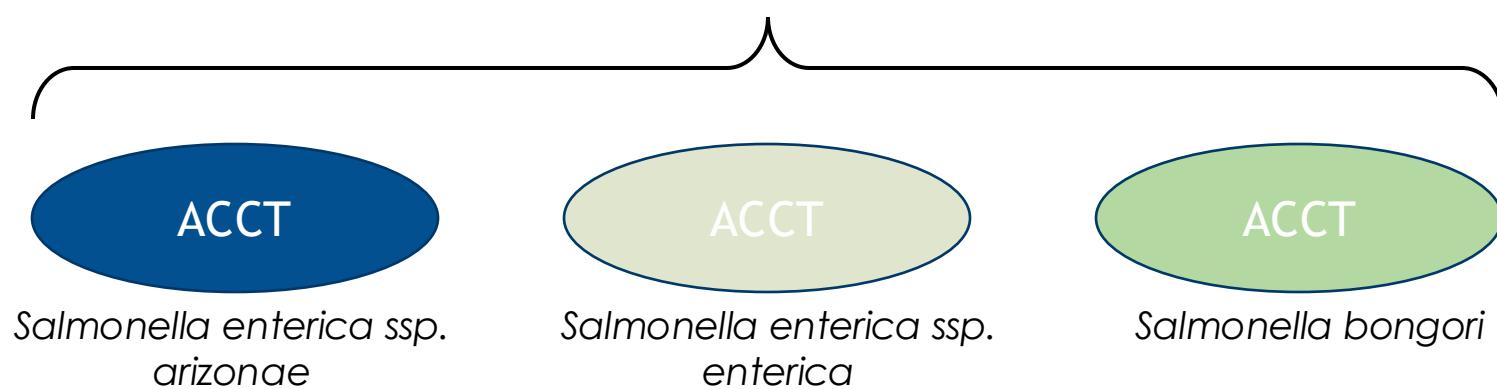
Order

Family

Genus

Species

## How could we resolve this?



Kingdom

Phylum

Class

Order

Family

Genus

Species



*Salmonella*



*Salmonella*



*Salmonella*

Kingdom

Phylum

Class

Order

Family

Genus

Species

**Uh-oh.....**



Kingdom

Phylum

Class

Order

Family

Genus

Species



*Enterobacteriaceae*



*Enterobacteriaceae*



*Enterobacteriaceae*



*Enterobacteriaceae*

# What do you end up with? The database...

<u>K-mer sequence</u>	<u>NCBI taxonomic ID</u>	<u>kraken ID</u>
ACCT	<i>Enterobacteriaceae</i>	1339265
GGTA	<i>Escherichia coli</i>	3927561

# How do we use this database to classify the reads in our metagenomic data?

<u>K-mer sequence</u>	<u>NCBI taxonomic ID</u>	<u>kraken ID</u>
ACCT	<i>Enterobacteriaceae</i>	1339265
GGTA	<i>Escherichia coli</i>	3927561

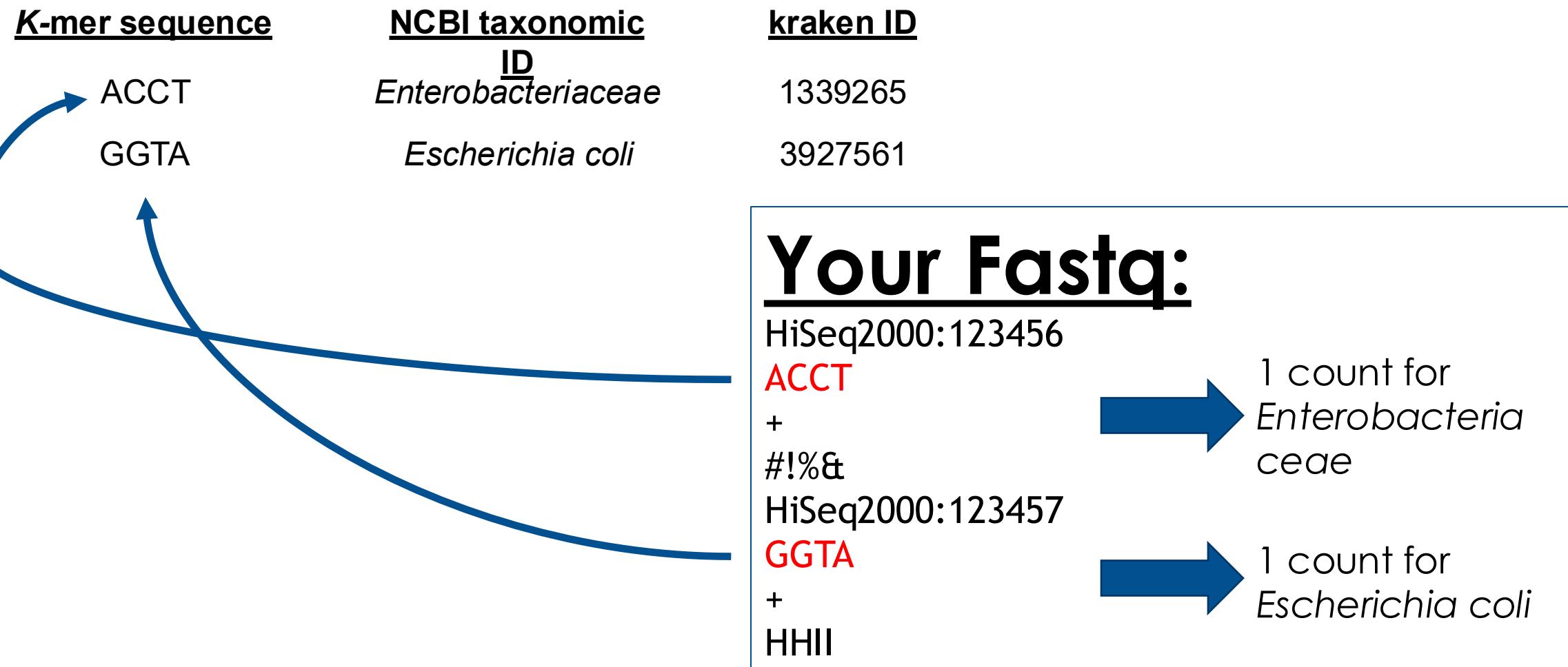
When you then run your sequence data through kraken, it looks for exact *k*-mer matches

<u>K-mer sequence</u>	<u>NCBI taxonomic ID</u>	<u>kraken ID</u>
ACCT GGTA	Enterobacteriaceae <i>Escherichia coli</i>	1339265
		3927561

**Your Fastq:**

```
HiSeq2000:123456
ACCT
+
#!%&
HiSeq2000:123457
GGTA
+
HHII
```

When you then run your sequence data through kraken, it looks for exact *k*-mer matches



# But sometimes, it is messy...

- Kraken's default k-mer size for the exact-matching process is 31
- Kraken takes every 31-mer from every read and attempts to match it (exactly!) to the kraken database
  - Some 31-mers don't match (exactly) to anything = unclassified
  - A single read is usually 150-300bp, and therefore it contains many 31-mers. Sometimes, those 31-mers match to *different* things. In that case...there is a confidence scoring system.

# Benefits and Limitations (some)

## The Good:

- Rapid
- High specificity (\*\*caveat: for those organisms in the RefSeq database)
- Pulls from a relatively well-curated database (\*\*caveat: plasmids)

## The Bad:

- High rate of “unclassified”
- Varying “degrees” of taxonomic resolution (and probably non-random)
- Pulls from a relatively small database

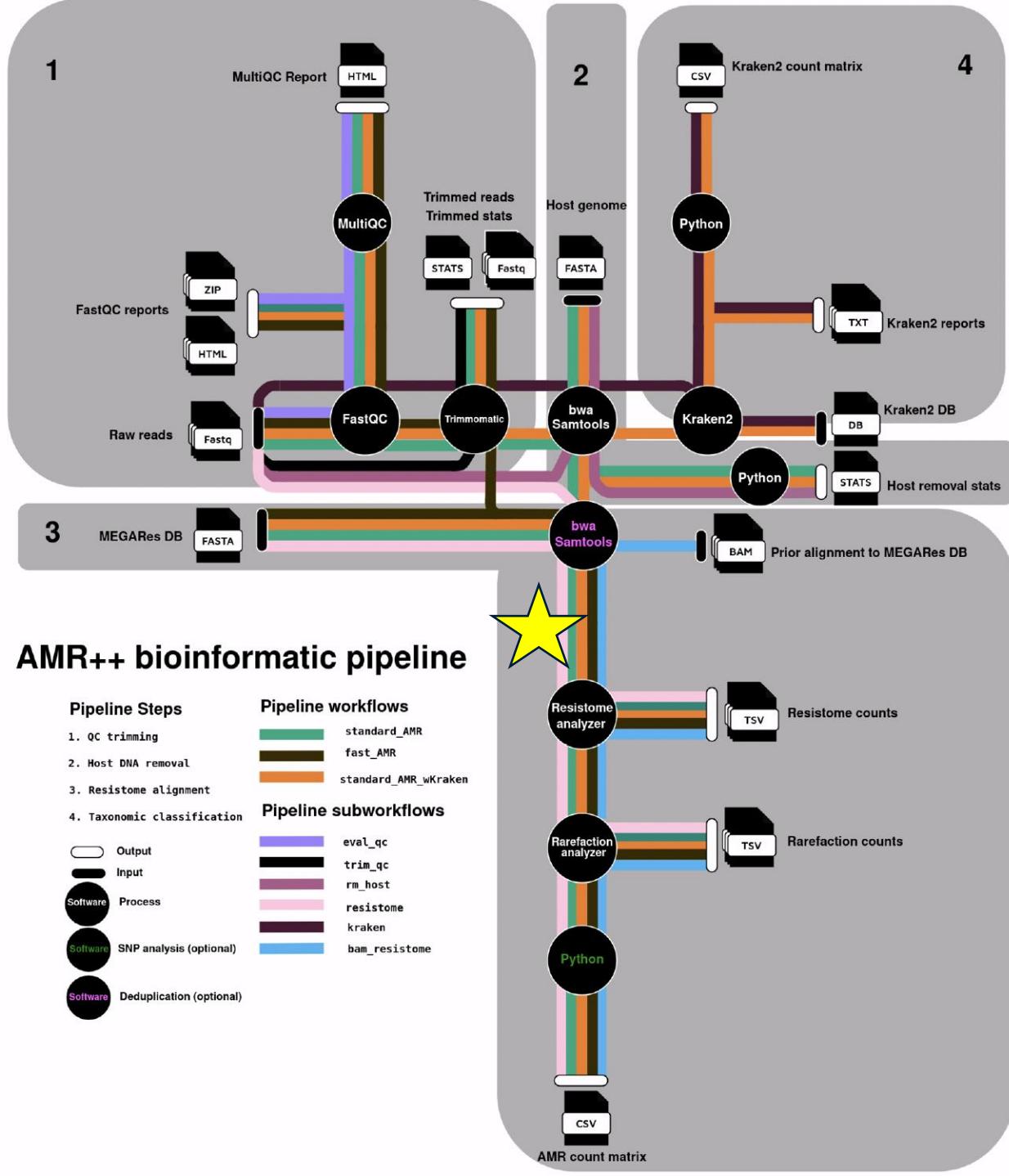
# More details (from user manual):

OUTPUT: 562:13 561:4 A:31 0:1 562:3 would indicate that:

- the first 13 *k*-mers mapped to taxonomy ID #562
- the next 4 *k*-mers mapped to taxonomy ID #561
- the next 31 *k*-mers contained an ambiguous nucleotide
- the next *k*-mer was not in the database
- the last 3 *k*-mers mapped to taxonomy ID #562

In this case, ID #561 is the parent node of #562. Here, a label of #562 for this sequence would have a score of  $C/Q = (13+3)/(13+4+1+3) = 16/21$ . A label of #561 would have a score of  $C/Q = (13+4+3)/(13+4+1+3) = 20/21$ . If a user specified a confidence threshold over 16/21, the classifier would adjust the original label from #562 to #561; if the threshold was greater than 20/21, the sequence would become unclassified.

# Part 4: AMR Detection



# What is our goal?

To describe and/or compare resistomes in a set of samples!!



To do that, we must identify and quantify AMR genes



To do that, we must find and count AMR genes in metagenomic data

Remember, that for all of the options we will cover today, it all comes down to matching...

*(And that's why the databases are so important!)*

# Basic Alignment Approach

# Sequence reads

TACGTGCAACCAAGACCAACCAGTCTTCCCCGCTTT  
 TAAGCTTACGTGCAACCAAGACCAACCATACTTCCCAG CTCGACACAAGA  
 GTAAAGCTTACGTGCAACCAAGACCAACCAGTCTTCCC CTCGACACAAGA  
 GTAGTAAGCTTACGTGCAACCAAGACCAACCAGAGATC CTCGACACAAGA  
 TTAGTAGTAAGCTTACGTGCAACCAAGACCAACCAGGC CTGCTCGACACAAGA  
 TTAGTAGTAAGCTTACGTGCAACCAAGACCAACCAGTC CAGCTTTCTCGACACAAGA  
 CTTAGTAGTAAGCTTACGTGCAACCAAGACCAACCAGT CCAGCTTTCTCGACACAAGA  
 CTTAGTAGTAAGCTTACGTGCAACCAAGACCAACCAG CAGCTTTCTCGACACAAGA  
 CCTTAGTAGTAAGCTTACGTGCAACCAAGACCAACCAC CCAGCTTTCTCGACACAAGA  
 CCTTAGTAGTAAGCTTACGTGCAACCAAGACCAACCAG TCCAGCTTTCTCGACACAAGA  
 ACCTTAGTAGTAAGCTTACGTGCAACCAAGACCAACCA TGCCAGCTTTCTCGACACAAGA  
 ACCTTAGTAGTAAGCTTACGTGCAACCAAGACCAACCA TTCCAGCTTTCTCGACACAAGA  
 TTACCTTAGTAGTAAGCTTACGTGCAACCAAGACCAAC CTTCCCAGCTTTCTCGACACAAGA  
 AAAACGTTACCTTAGTAGTAAGCTTACGTGCAACCAAA CTTCCCAGCTTTCTCGACACAAGA  
 AAAACGTTACCTTAGTAGTAAGCTTACGTGCAACCAAA GTCTTCCCAGCTTTCTCGACACAAGA  
 CGAAAAAACGTTACCTTAGTAGTAAGCTTACGTGCAAC AGTCTTCCCAGCTTTCTCGACACAAGA  
 ACGAAAAAACGTTACCTTAGTAGTAAGCTTACGTGCC CAGTCTTCCCAGCTTTCTCGACACAAGA  
 ACGAAAAAACGTTACCTTAGTAGTAAGCTTACGTGC CCACCTGTCTTCCCAGCTTTCTCGACACAAGA  
 AACGAAAAAACGTTACCTTAGTAGTAAGCTTACGTGC ACCACCAAGTCTTCCCAGCTTTCTCGACACAAGA  
 AACGAAAAAACGTTACCTTAGTAGTAAGCTTACGTGC GACCACCAAGTCTTCCCAGCTTTCTCGACACAAGA  
 AACGAAAAAACGTTACCTTAGTAGTAAGCTTACGTGC AAGACCAACCAAGTCTTCCCAGCTTTCTCGACACAAG  
 AACGAAAAAACGTTACCTTAGTAGTAAGCTTACGTGC AAGACCAACCAAGTCTTCCCAGCTTTCTCGACACAAG  
 TATAAACGAAAAAACGTTACCTTAGTAGTAAGATTAC CCAACACCAACCAAGTCTTCCCAGCTTTCTCGACACA  
 ATAATAAACGAAAAAACGTTACCTTAGTAGTAAGCTTA CCAACACCAACCAAGTCTTCCCAGCTTTCTCGACACA  
 ATATAAACGAAAAAACGTTACCTTAGTAGTAAGAT CACCAAGACCAACCAAGTCTTCCCAGCTTTCTCGA  
 AAATAAACGAAAAAACGTTACCTTAGTAGTAAGCT GCACCAAGACCAACCAAGTCTTCCCAGCTTTCTCGA  
 AAATAAACGAAAAAACGTTACCTTAGTAGTAAGCT GTGCACCAAGACCAACCAAGTCTTCCCAGCTTTCTCG  
 AAATAAACGAAAAAACGTTACCTTAGTAGTAAGCT ACGTGCACCAAGACCAACCAAGTCTTCCCAGCTTTCTCG  
 AAATAAACGAAAAAACGTTACCTTAGTAGTAAGCT ACGTGCACCAAGACCAACCAAGTCTTCCCAGCTTTCTCG  
 TTATAAACGAAAAAACGTTACCTTAGTAGTA TACGGCGCCAAGACCAACCAAGTCTTCCCAGCTTTCTCG  
 ATTCTAACGAAAAAACGTTACCTTAGTAGTA TTACGTGCAACCAAGACCAACCAAGCCTCCCAGCTTTCTCG  
 ATTCTAACGAAAAAACGTTACCTTAGTAGTAAGATTAC ATTCTAACGAAAAAACGTTACCTTAGTAGTAAGATTAC

# Important characteristics:

- The match does not have to be perfect
- Matching is performed at the unit of the read
- The alignment preserves a lot of information

# Benefits of alignment:

1. Utilizes all of the (relevant) sequence data
2. Relatively quick (computationally)
3. Preserves sequence variability
4. Preserves the original sequence (no “assembly error”)

# Disadvantages of alignment:

1. Does not provide genomic context
2. Can be prone to false-positive identification of, e.g., AMR genes
3. Must make arbitrary decisions when there are homologous stretches of DNA between reference sequences

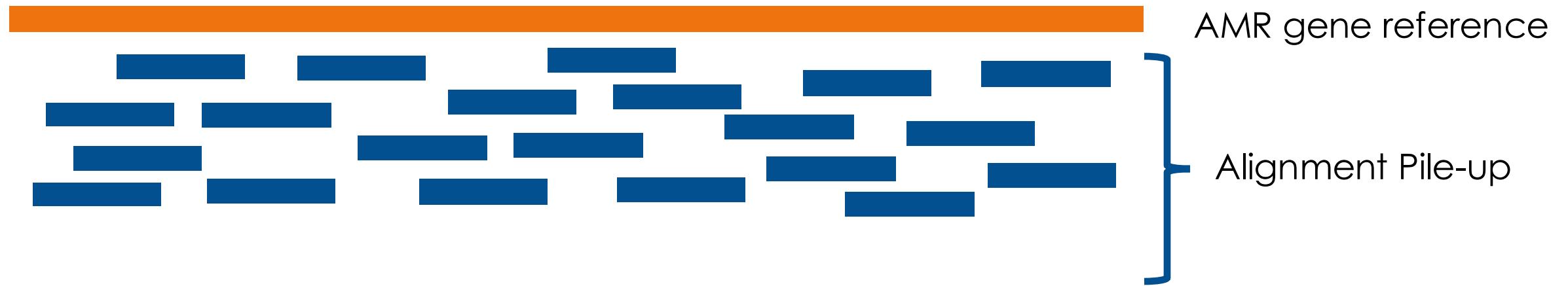
# Disadvantages of alignment:

1. Does not provide genomic context
2. Can be prone to false-positive identification of,  
e.g., AMR genes
3. Must make arbitrary decisions when there are  
homologous stretches of DNA between  
reference sequences

# Some examples:



# Some examples:



# Some examples:



# Some examples:



# So what do we do about this?

We use a “gene fraction” cutoff  
*(sometimes called gene “coverage”: AVOID THIS TERM)*

Gene Fraction: The % of nucleotides in a given reference sequence that have at least one mapped (aligned) read

# Disadvantages of alignment:

1. Does not provide genomic context
2. Can be prone to false-positive identification of, e.g., AMR genes
3. Must make arbitrary decisions when there are homologous stretches of DNA between reference sequences

AMR gene 1



AMR gene 2



} ***Reference***

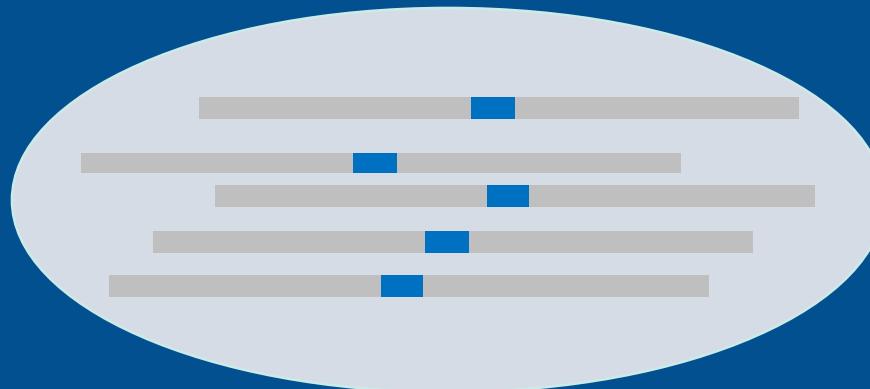
AMR gene 1



AMR gene 2



***Reference***



***Sample***

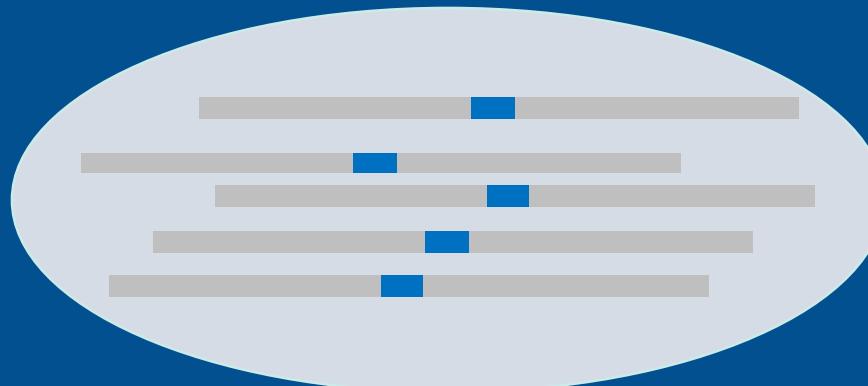
AMR gene 1

AMR gene 2

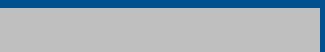
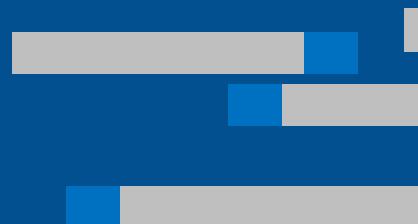
Reference



Sample



Data



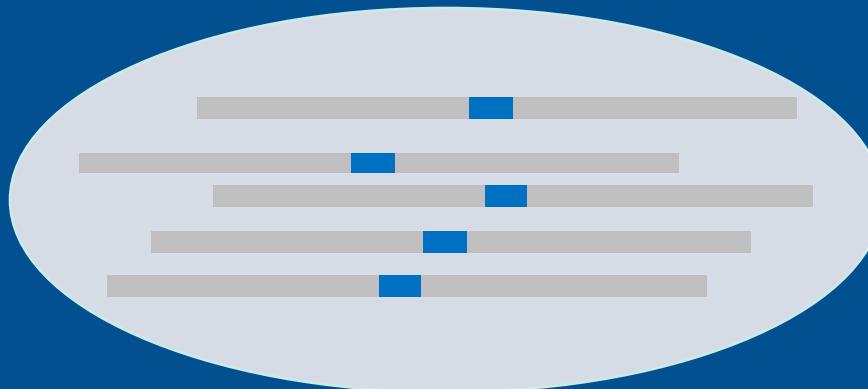
AMR gene 1



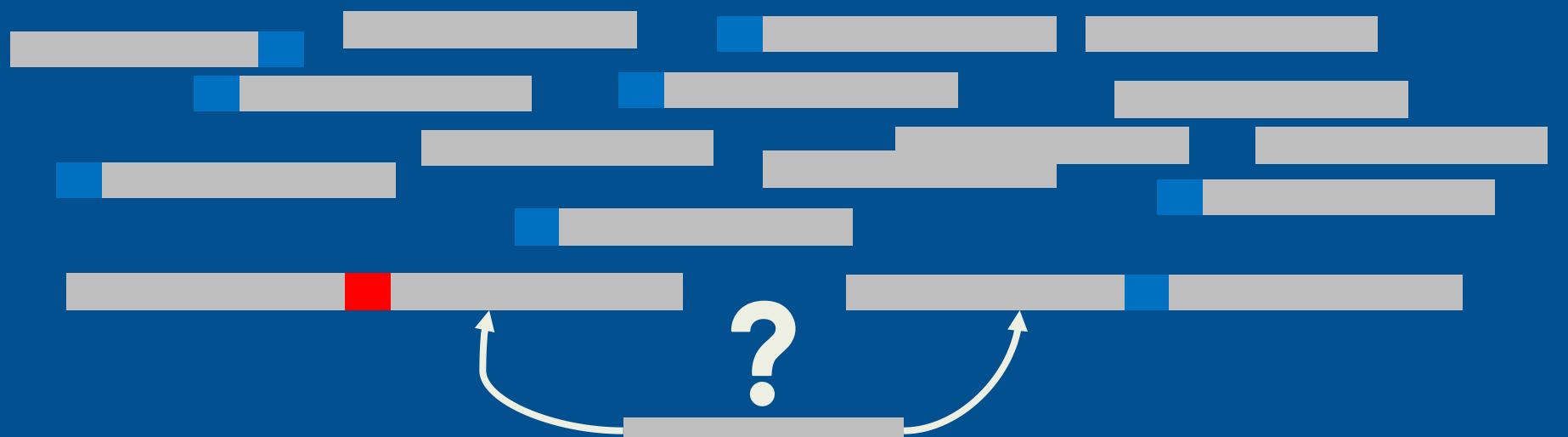
AMR gene 2



Reference



Sample



Data

?

Alignment

# What can we do?

Default of the aligner = random assignment

You can also force aligner to output ***all possible alignments*** (but this means you will “double count”)

You can “sidestep” the issue by conducting your analyses at a “higher level” – more on this later!