Gladstone Scientific Retreat 2024

# Unzipping the metagenome:

## strain-level discovery in the gut microbiome

**Byron J. Smith**

Bioinformatics Fellow

# First Thing: Thank You!

Gladstone Scientific Retreat 2024

# Introduction:

# The gut microbiome and shotgun metagenomics

The Gut Microbiome is Challenging

**The Gut Microbiome is Challenging**

- Enormous number of species

**The Gut Microbiome is Challenging**

- Enormous number of species

- Highly dynamic across people and time

**The Gut Microbiome is Challenging**

- Enormous number of species

- Highly dynamic across people and time

- Very hard to study in the lab

**The Gut Microbiome is Challenging**

- Enormous number of species

- Highly dynamic across people and time

- Very hard to study in the lab

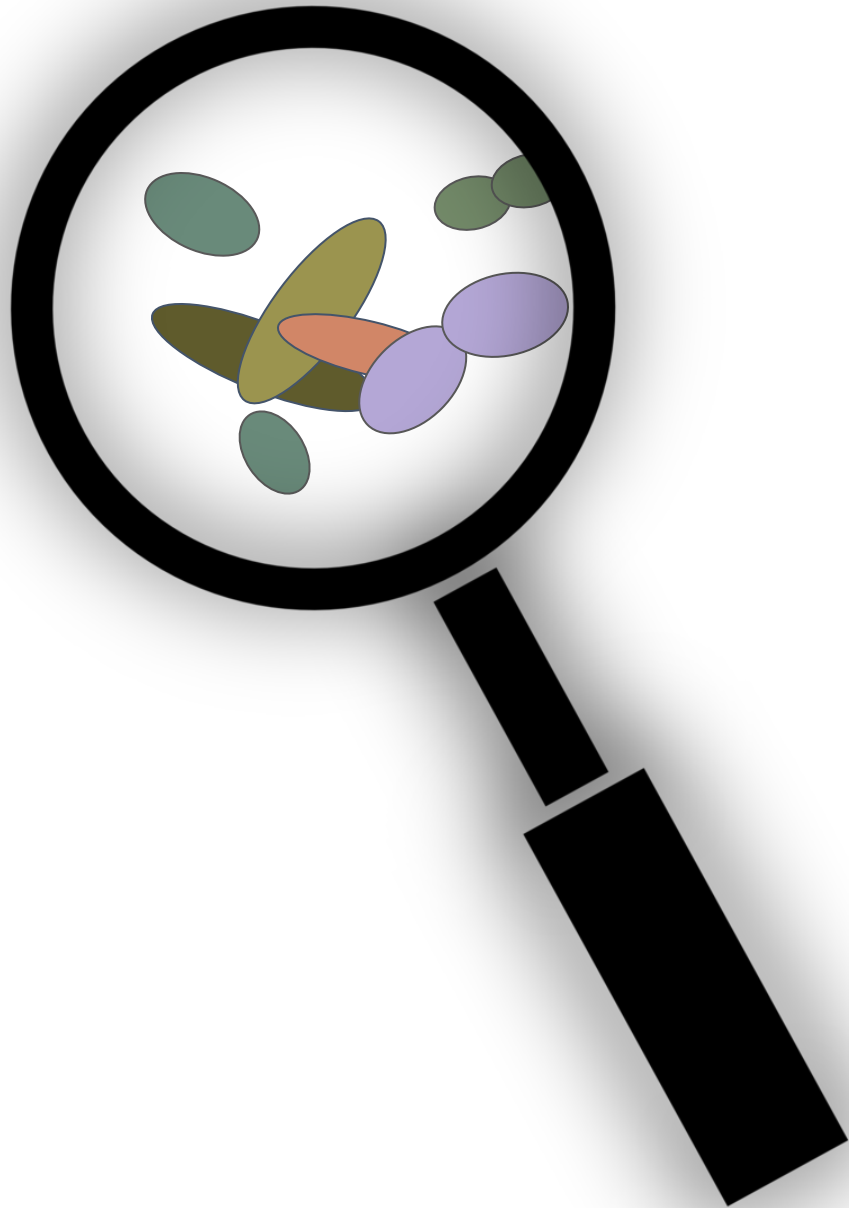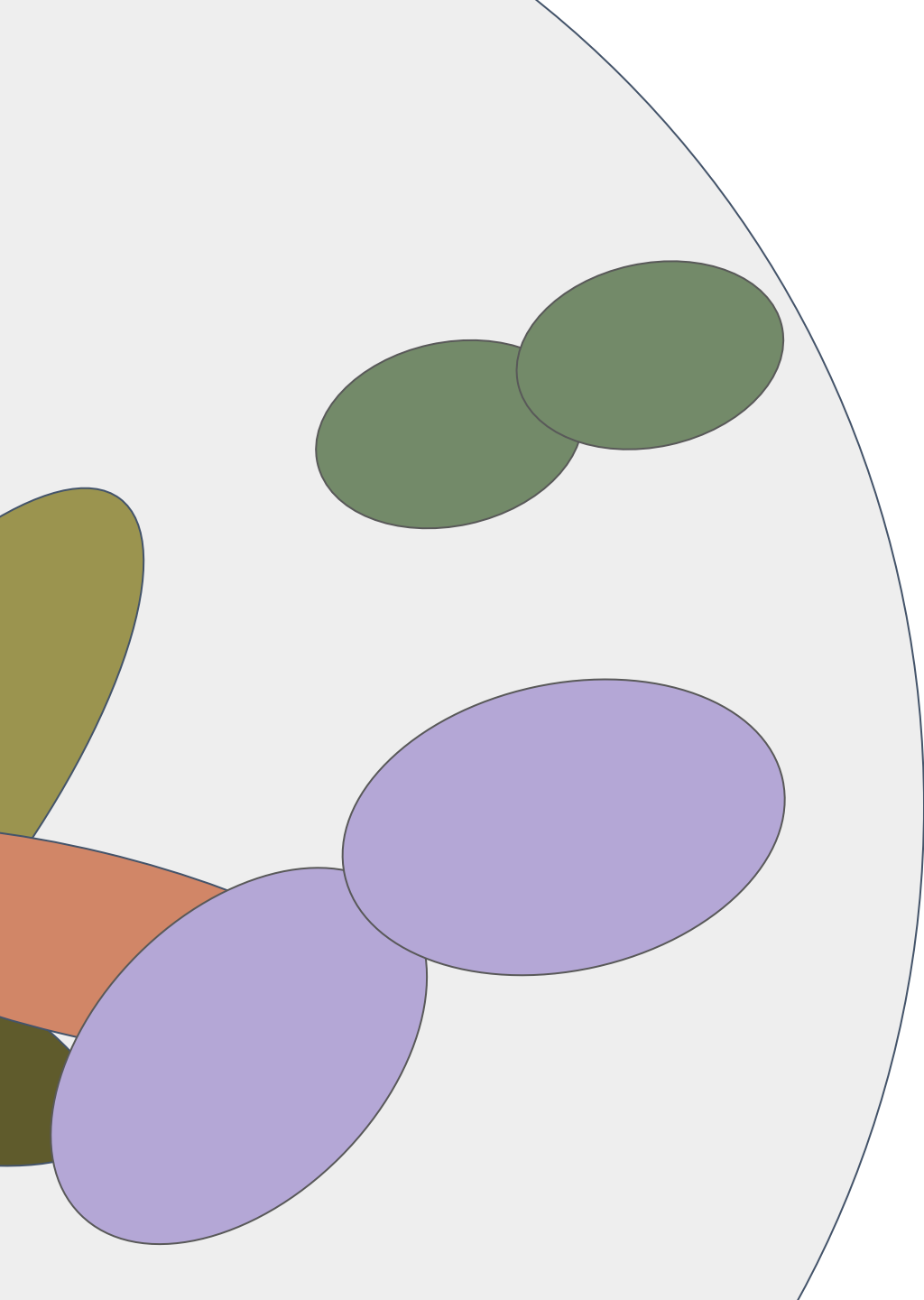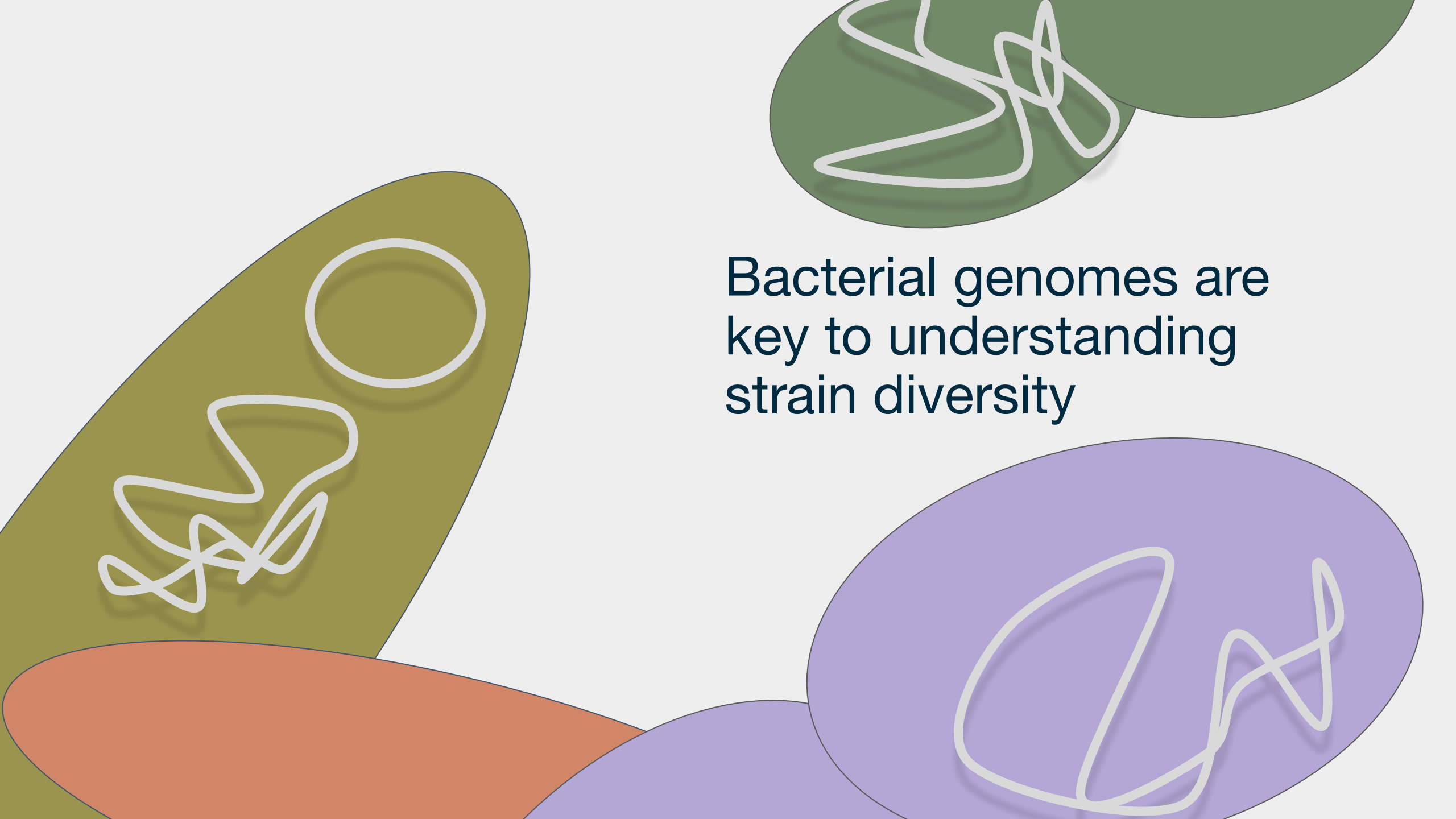- **Strains within species have different gene content and functional potential**

# The Gut Microbiome is Challenging

- Enormous number of species

- Highly dynamic across people and time

- Very hard to study in the lab

- **Strains within species have different gene content and functional potential**

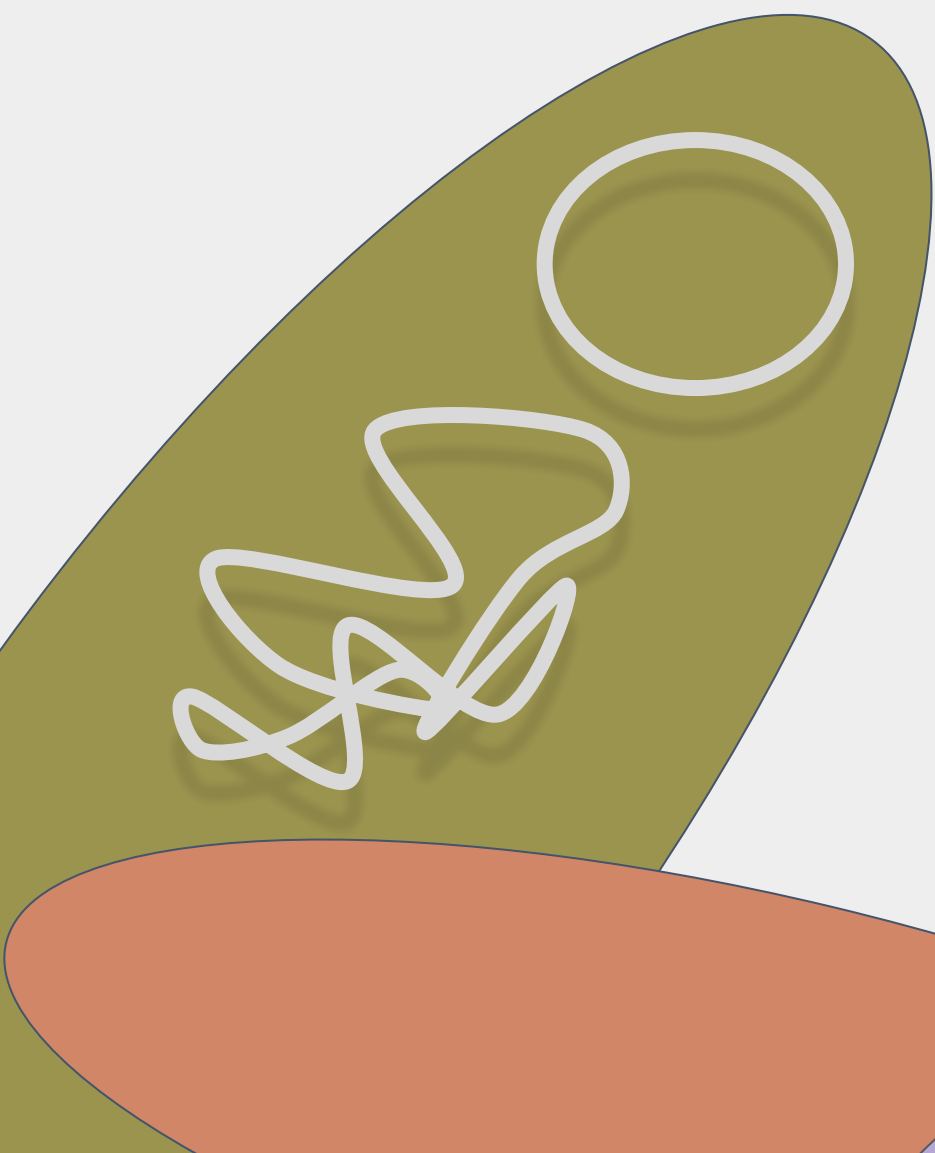Bacterial genomes are key to understanding strain diversity

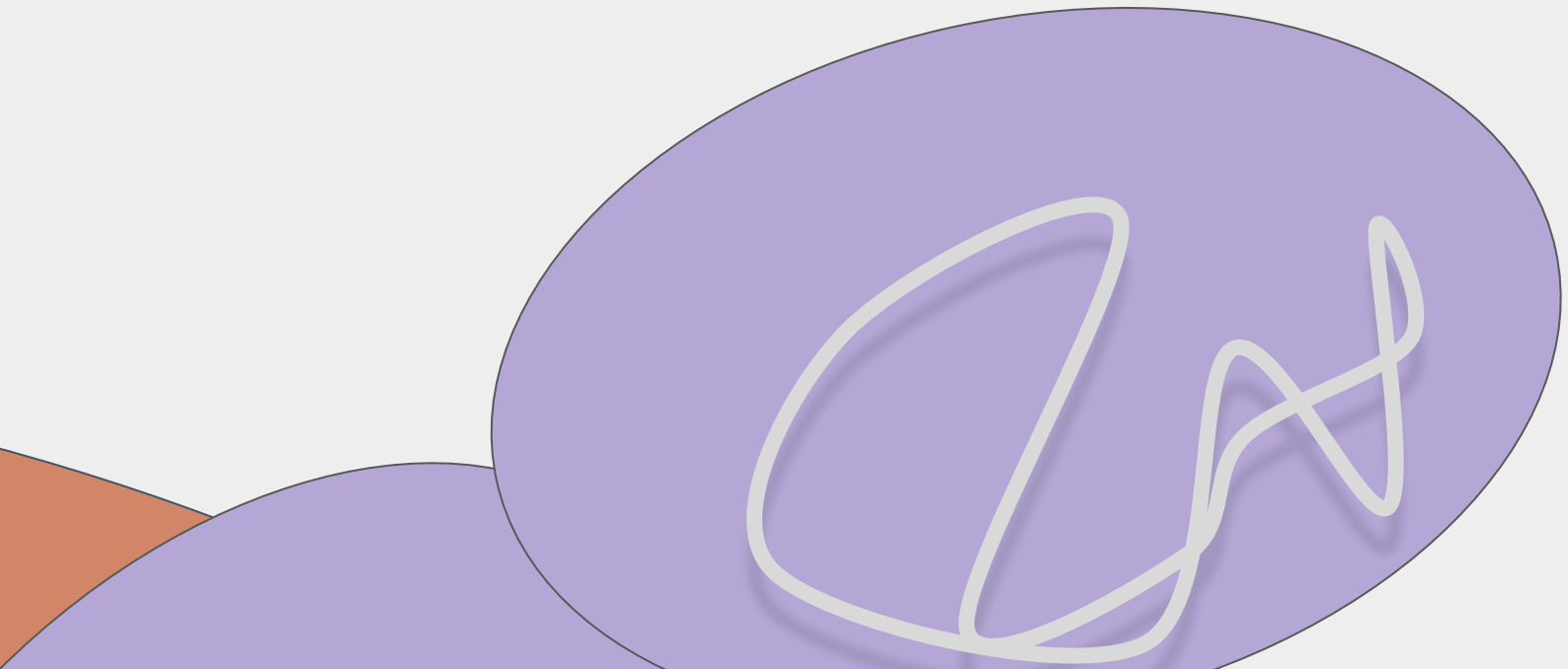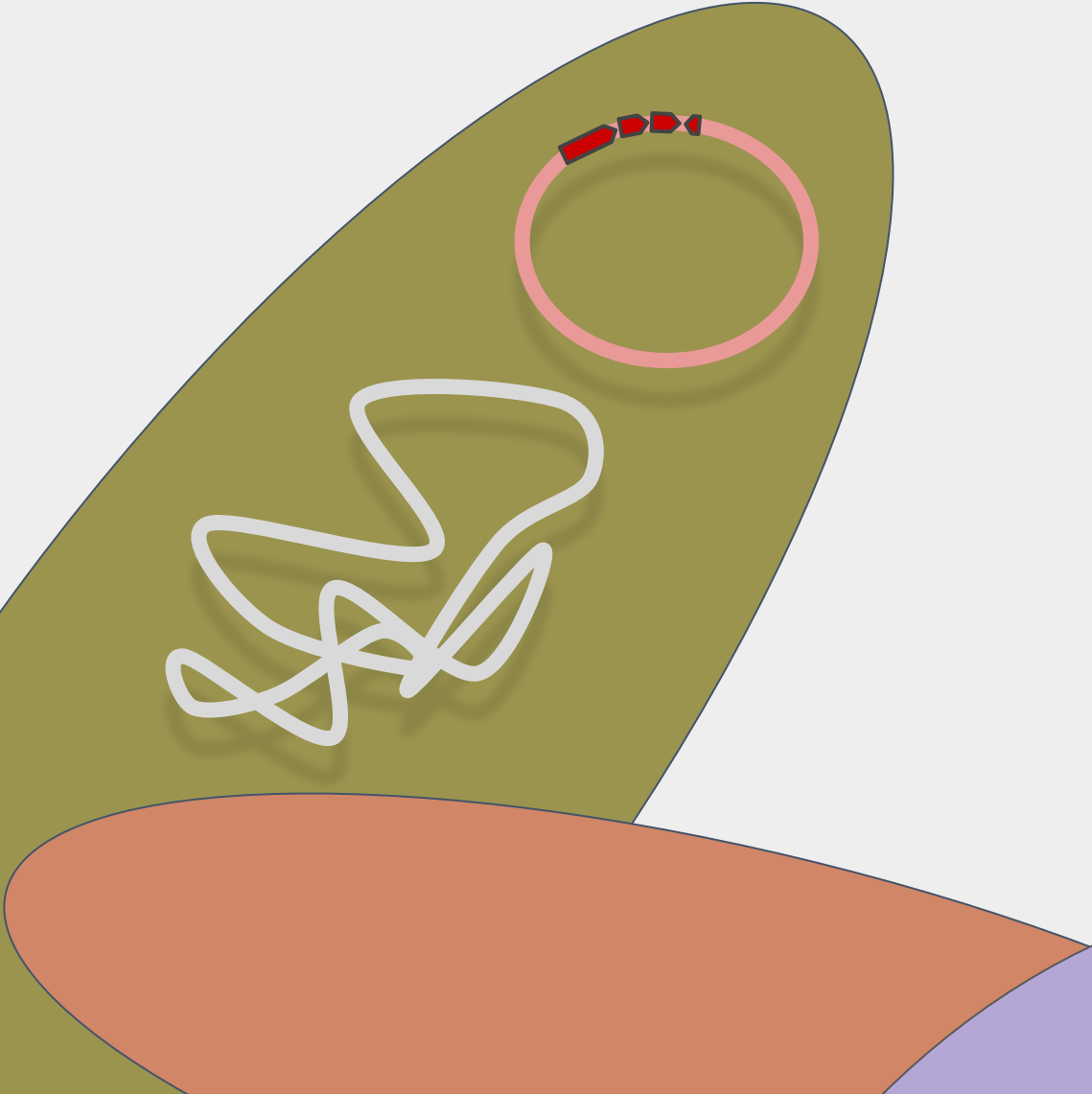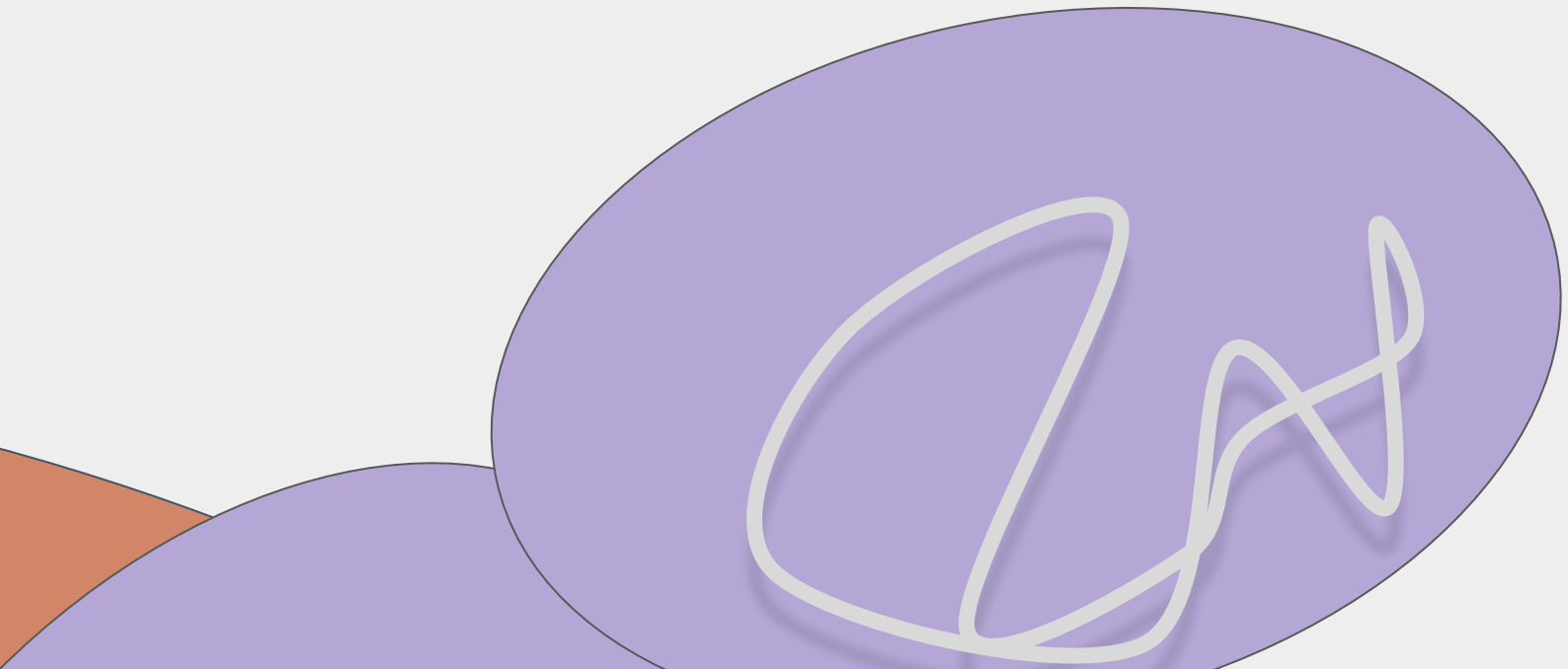Bacterial genomes are key to understanding strain diversity

Bacterial genomes are key to understanding strain diversity
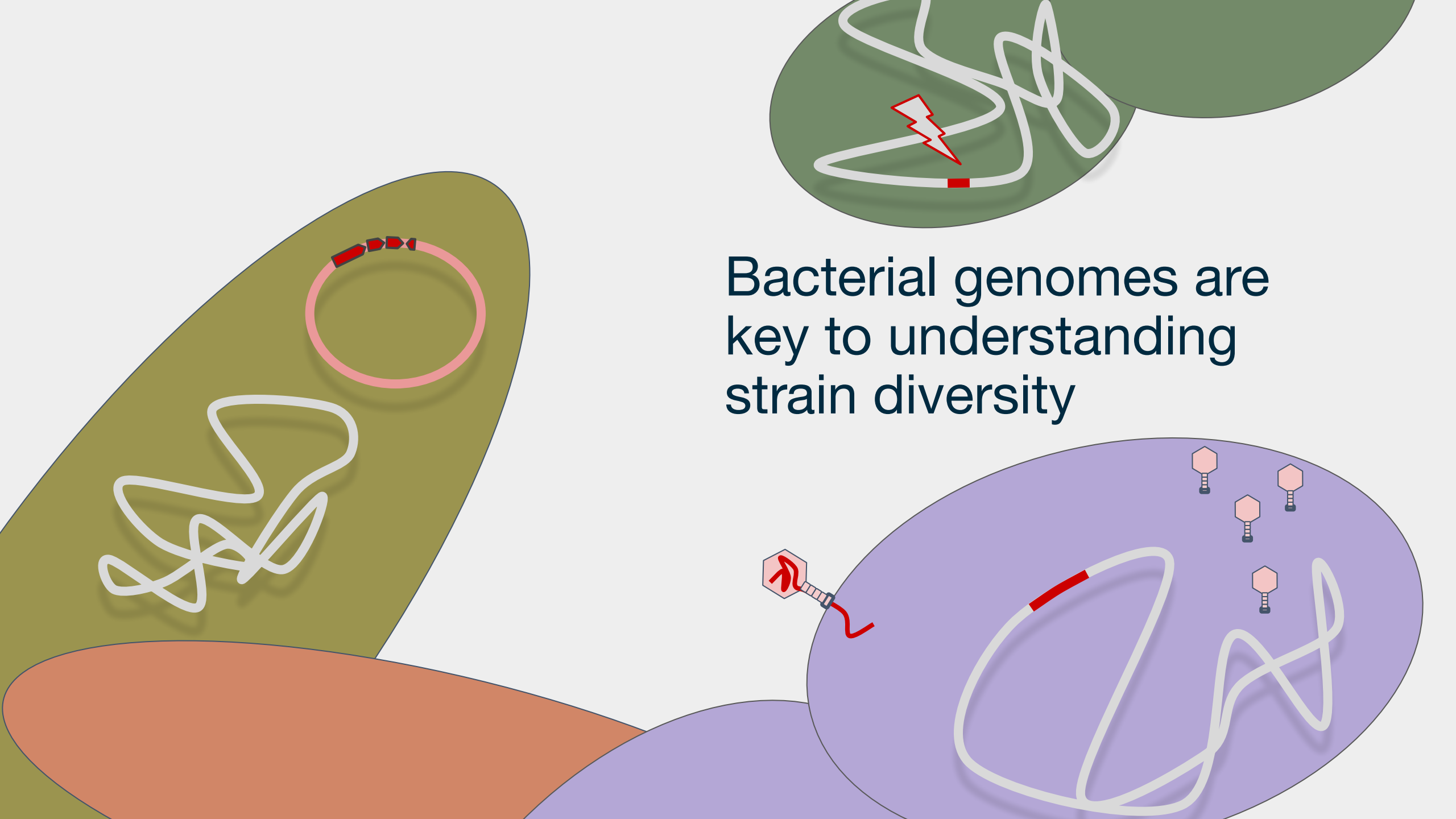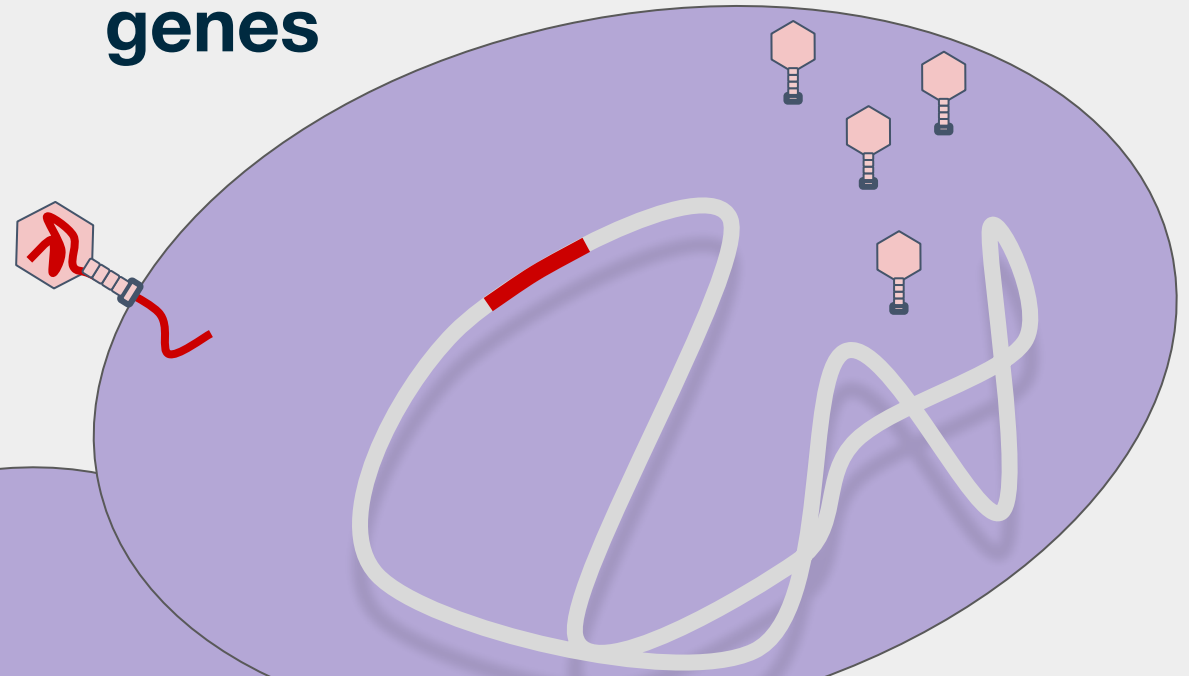
Bacterial genomes are key to understanding strain diversity

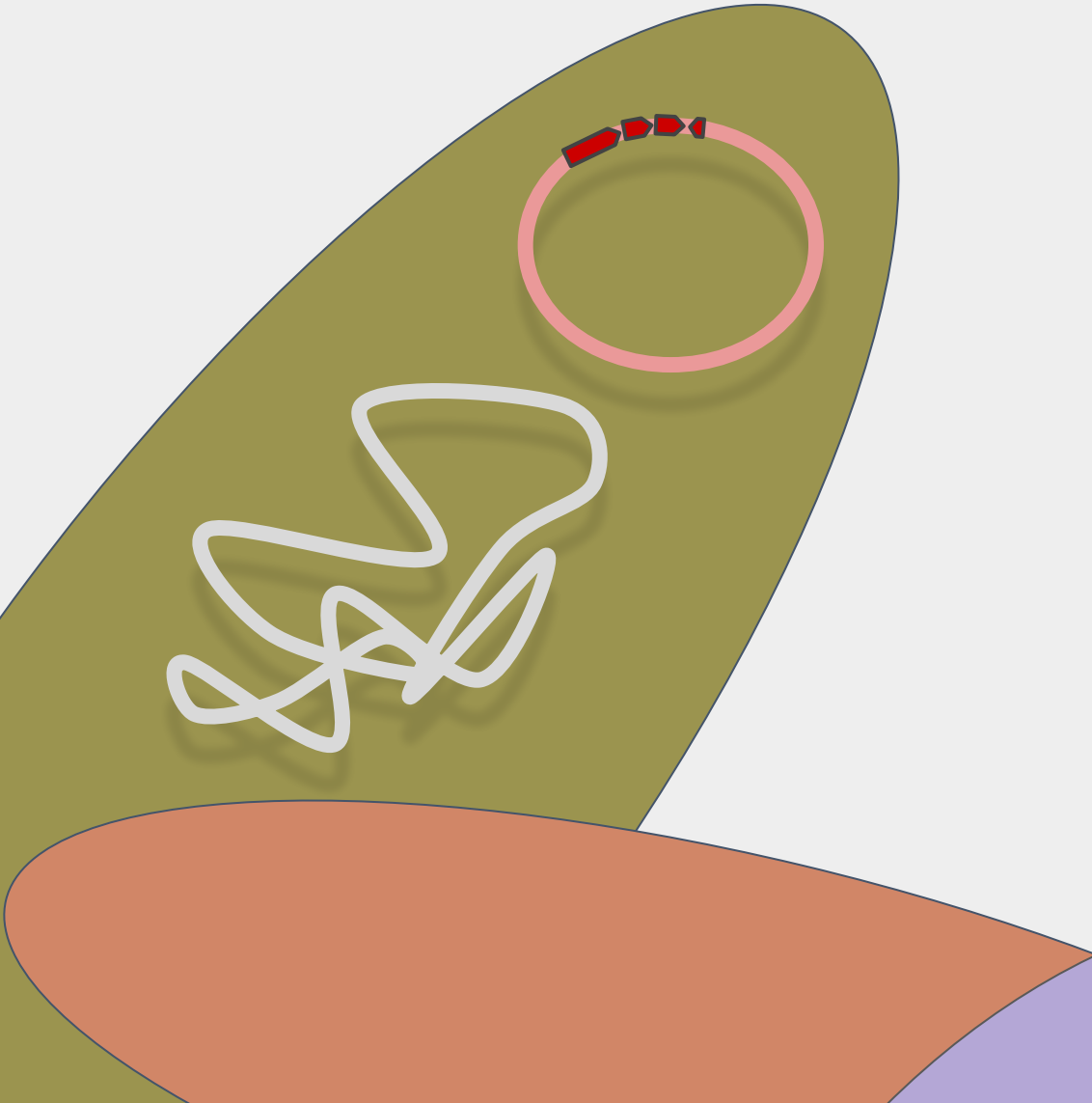Bacterial genomes are key to understanding strain diversity

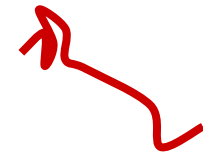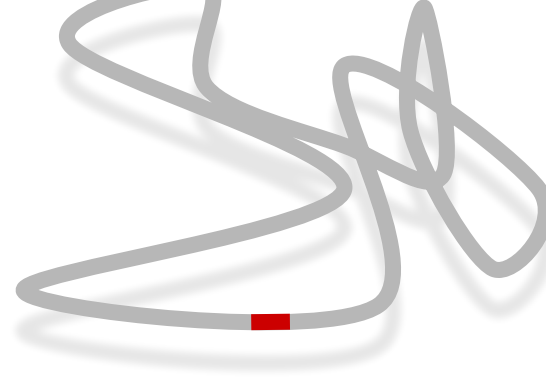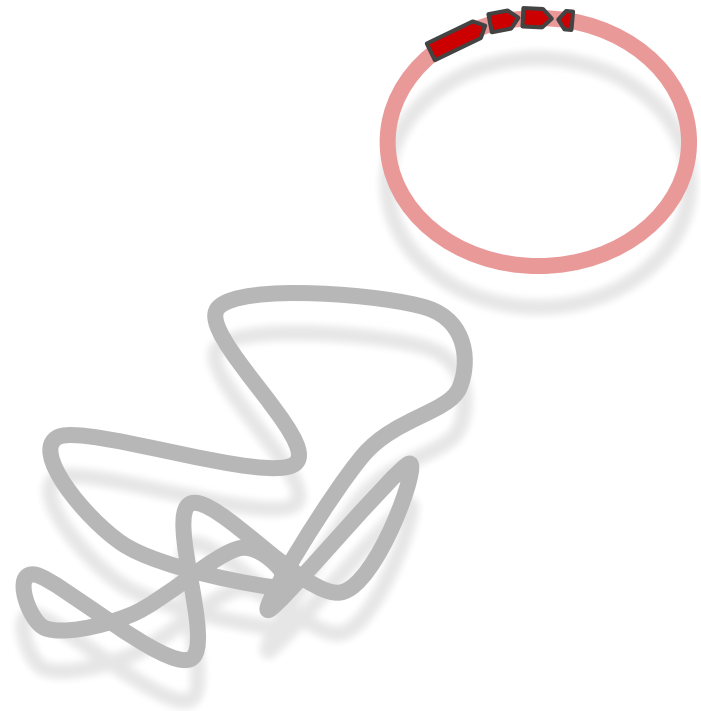Bacterial genomes are key to understanding strain diversity
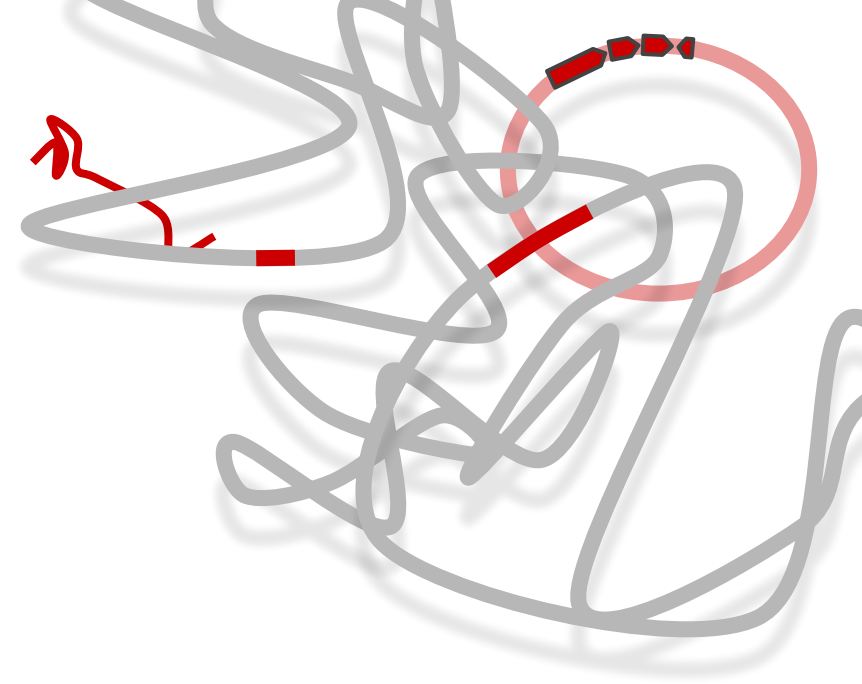
**Phage encoded antibiotic resistance genes**

# **Metagenomic** sequencing surveys all genomes

# Short-read, shotgun metagenomes enable modern microbiome science

Requirements:

# Short-read, shotgun metagenomes enable modern microbiome science

Requirements:
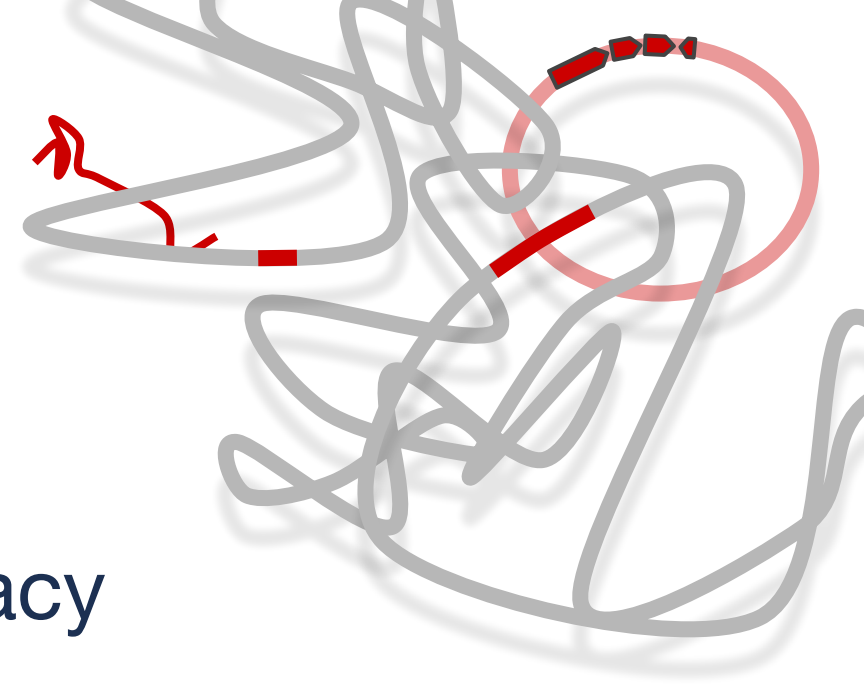
- strain-resolved genome sequences

  ➢ high accuracy

# Short-read, shotgun metagenomes enable modern microbiome science

Requirements:

- strain-resolved genome sequences
- capture low-abundance organisms

➢ high accuracy

➢ very deep sequencing

# Short-read, shotgun metagenomes enable modern microbiome science

Requirements:

- strain-resolved genome sequences
- capture low-abundance organisms
- lots of samples and longitudinal designs

➢ high accuracy

➢ very deep sequencing

➢ cheap

# Short-read, shotgun metagenomes enable modern microbiome science

Requirements:

- strain-resolved genome sequences ➢ high accuracy

- capture low-abundance organisms ➢ very deep sequencing

- lots of samples and longitudinal designs ➢ cheap

- long sequences ➢ ...

# Turning short reads into long sequences

ACGT ▶

# Turning short reads into long sequences

**ACGT**

**ACGT**

# Turning short reads into long sequences

ACGT ▶

ACGT **CCTT** ▶

**CCTT** ➤

# Turning short reads into long sequences

# Turning short reads into long sequences



Reconstruct linear genomes from overlapping sequences

...GGTAGAGCGTGGGACGTAGGGTTAACCTTAGAAAGCTAGAAAACCGCGCGCCCT...

# Problem:

# Problem: Closely related strains make read-chaining ambiguous

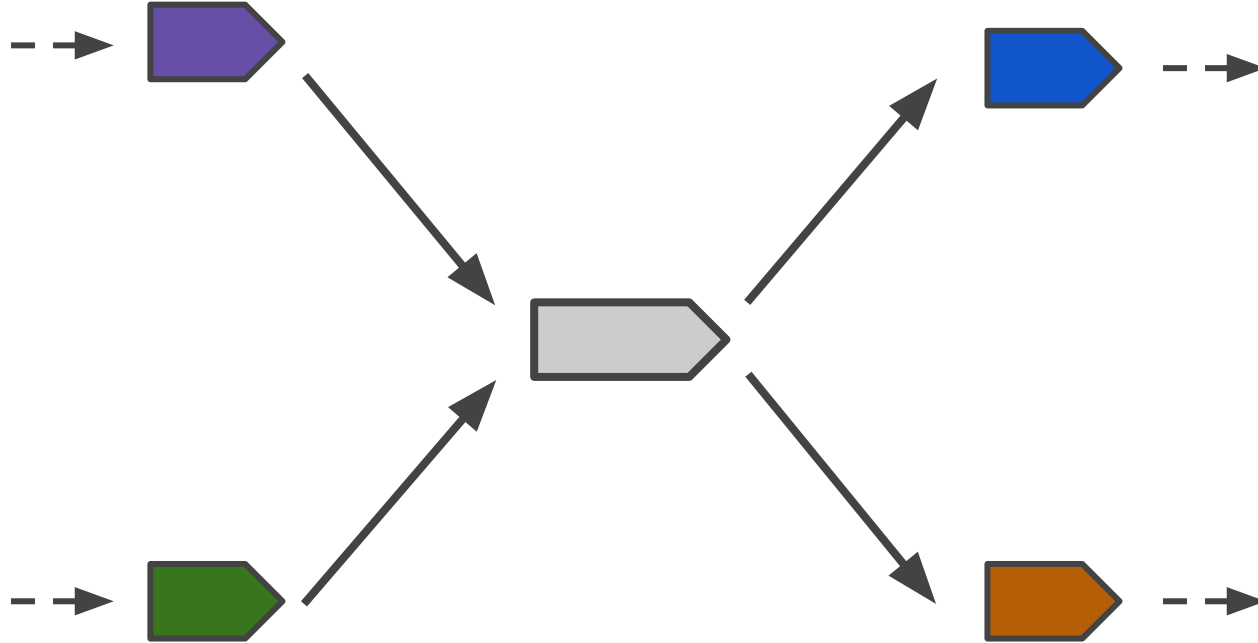# Problem: Closely related strains make read-chaining ambiguous

# Problem: Closely related strains make read-chaining ambiguous

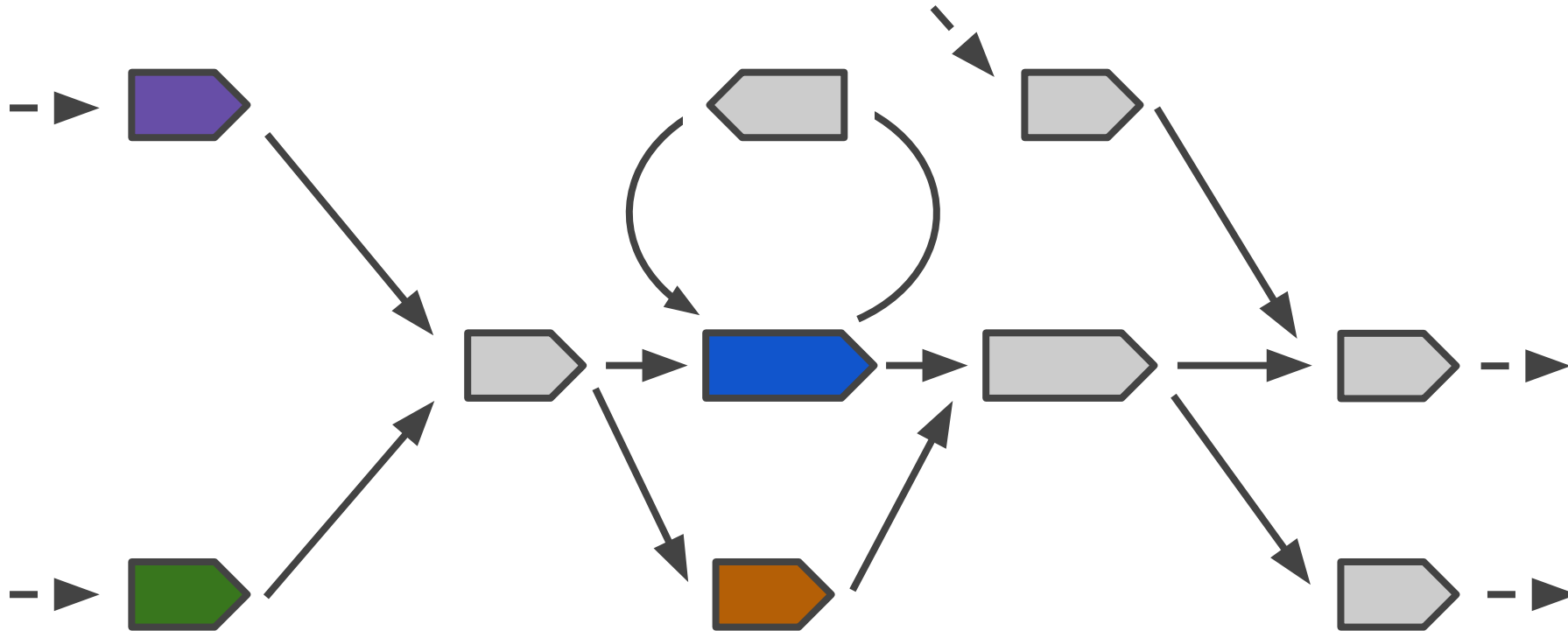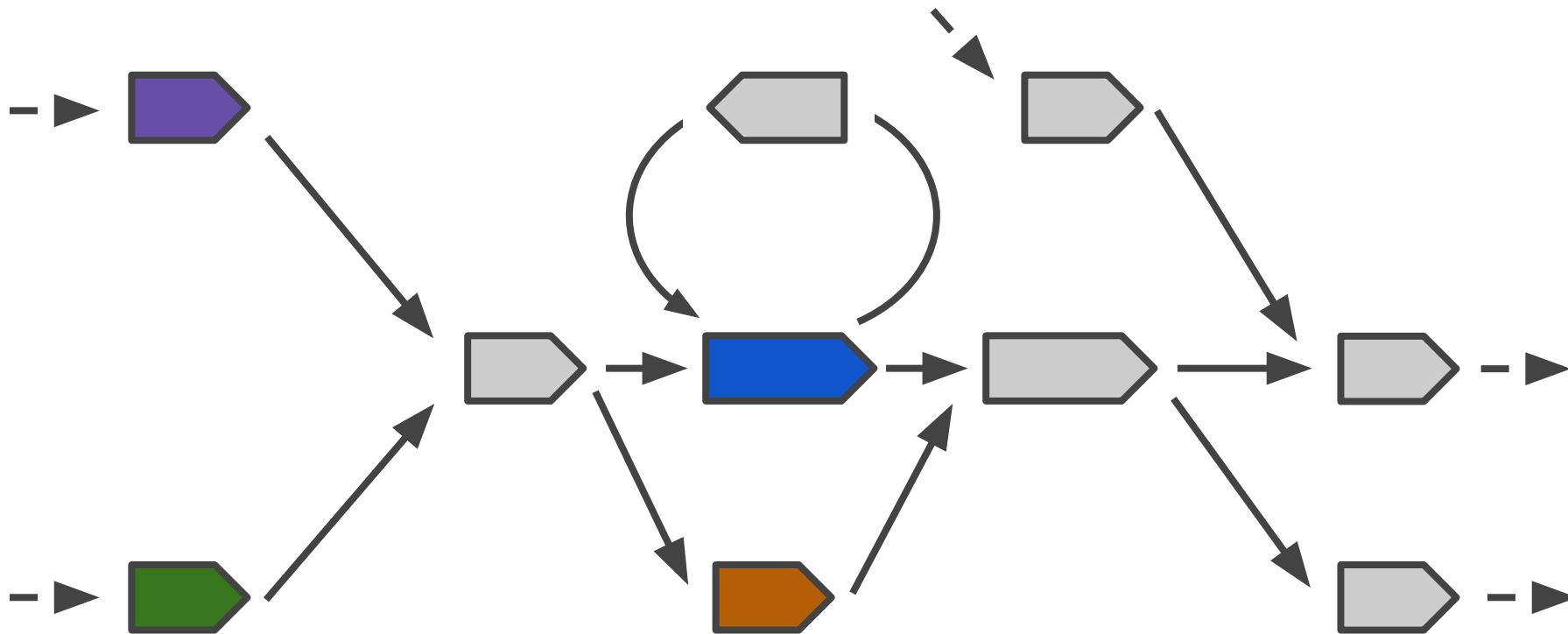# Can be represented as a graph of sequences linked by their overlaps

# Can be represented as a graph of sequences linked by their overlaps

# Can be represented as a graph of sequences linked by their overlaps

# Can be represented as a graph of sequences linked by their overlaps



(This problem also comes up for mRNA alternative splicing)

And real metagenomes
are **very** complex

And real metagenomes are **very** complex
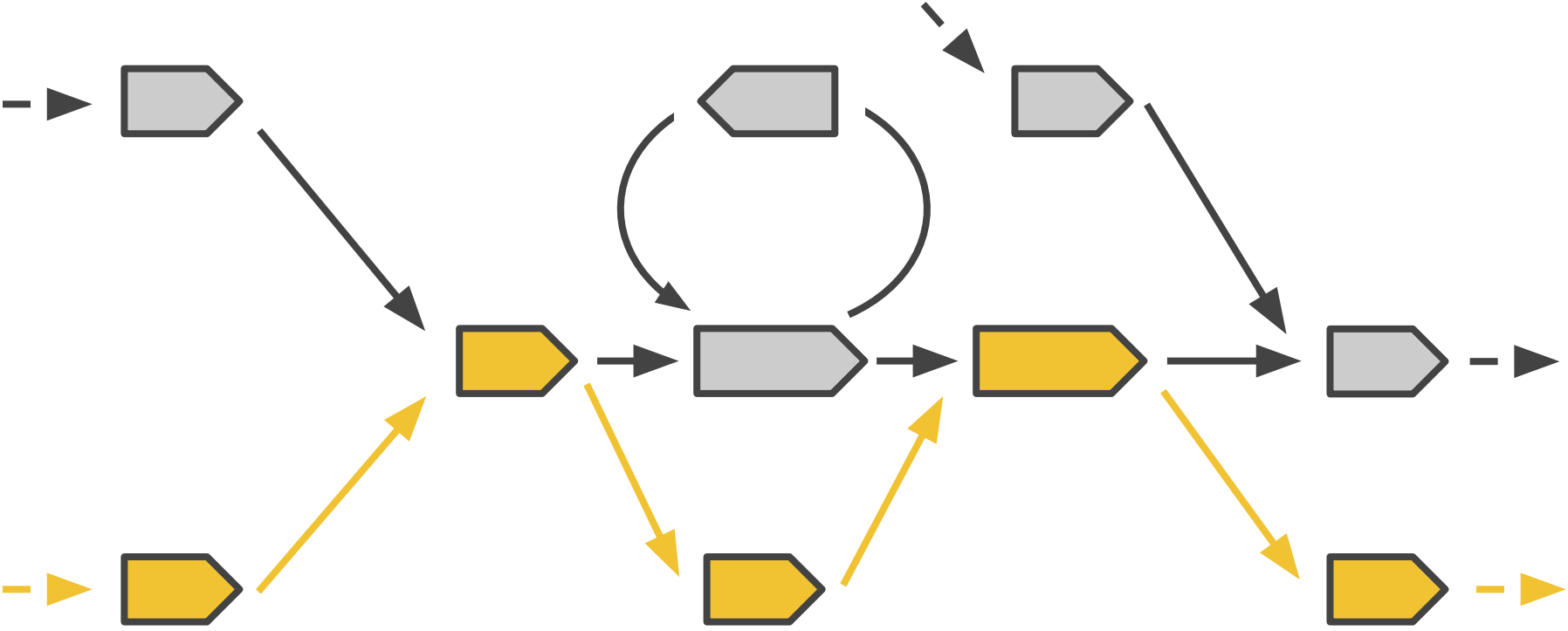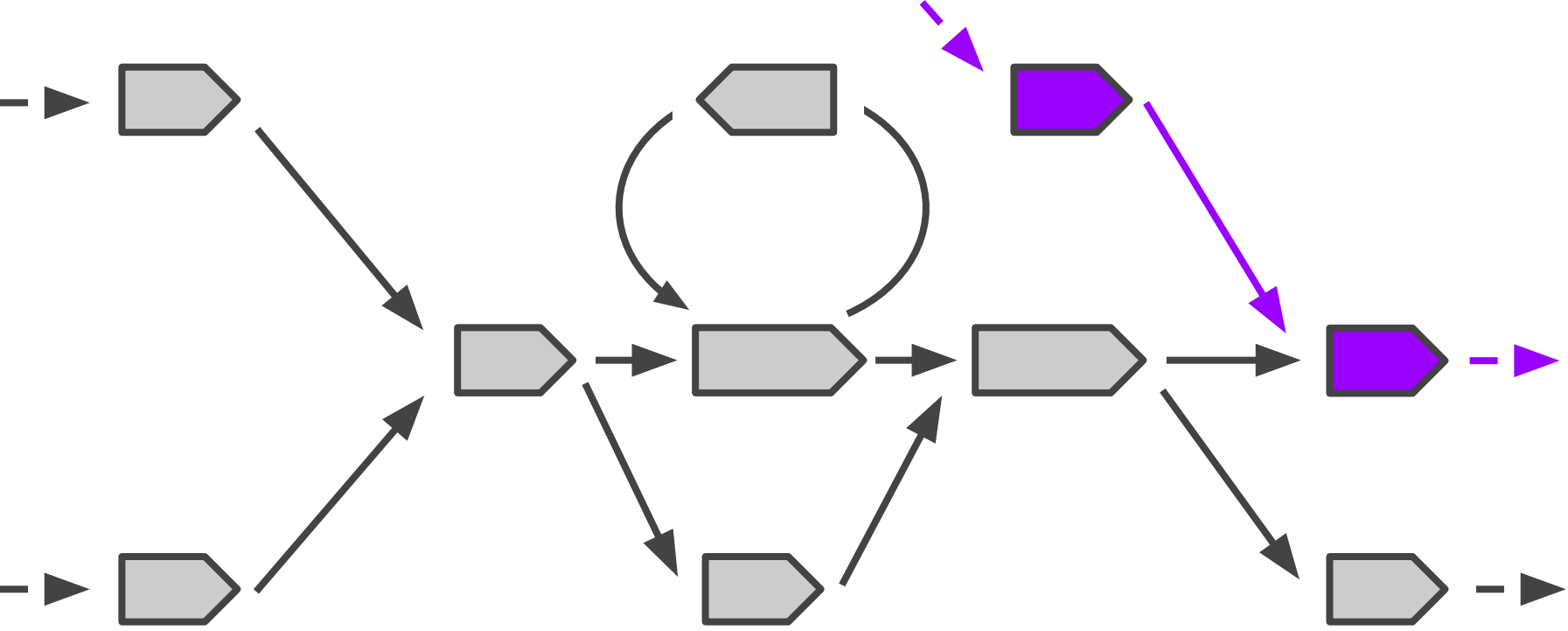
And real metagenomes are **very** complex

Real genomic
sequences are
paths on the graph

Real genomic sequences are paths on the graph
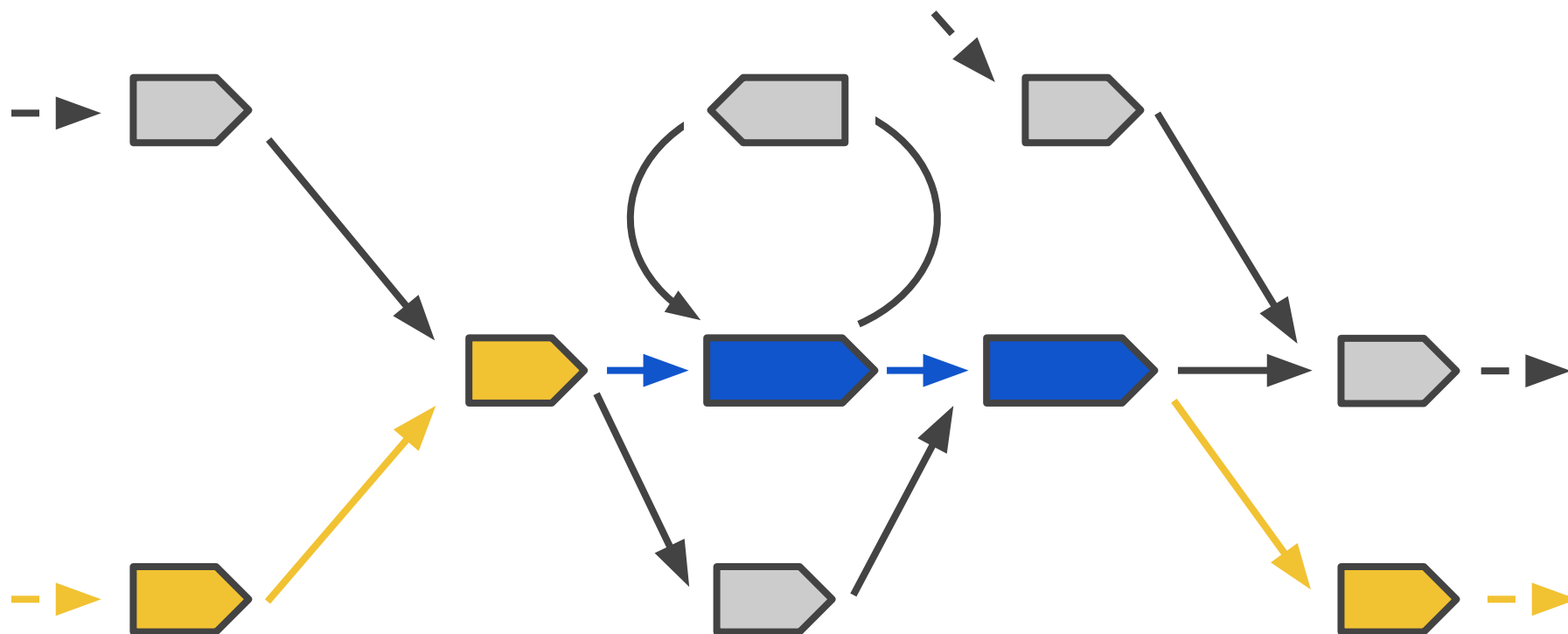
Real genomic sequences are paths on the graph
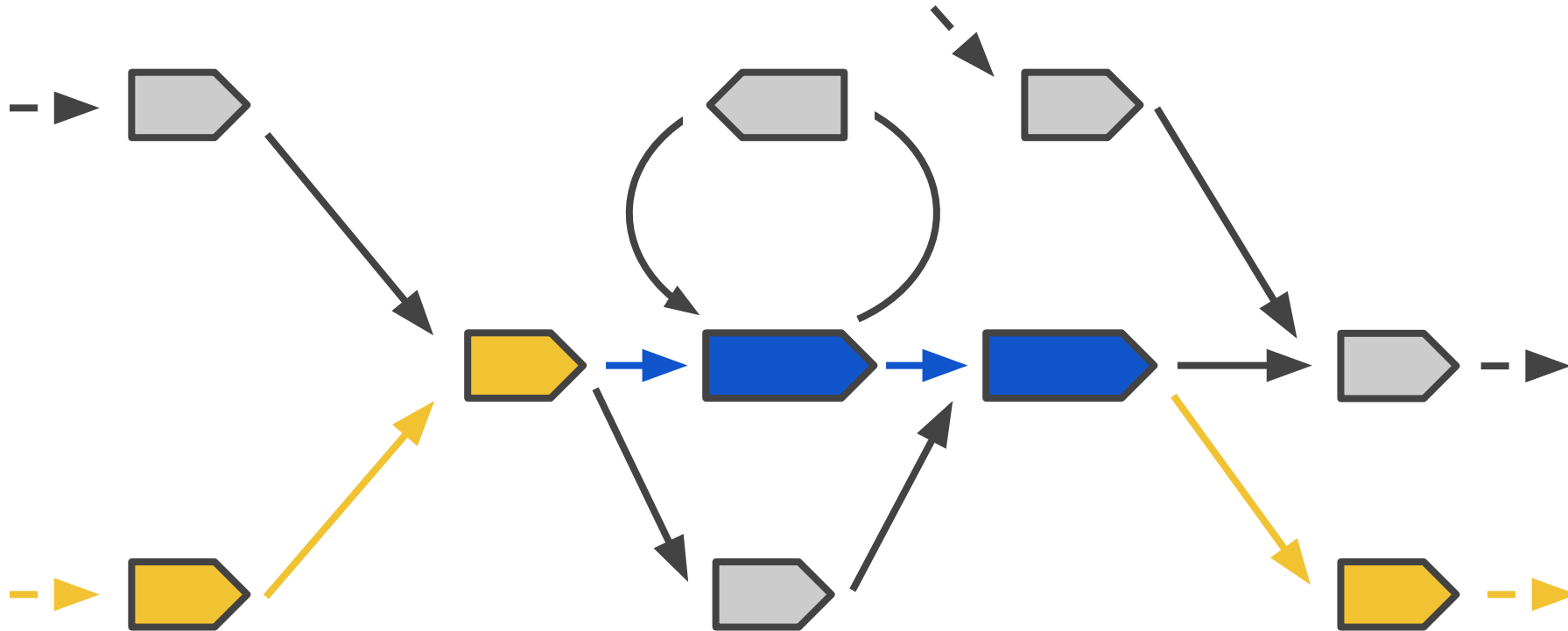
Real genomic sequences are paths on the graph

# Lots of incorrect paths also exist…

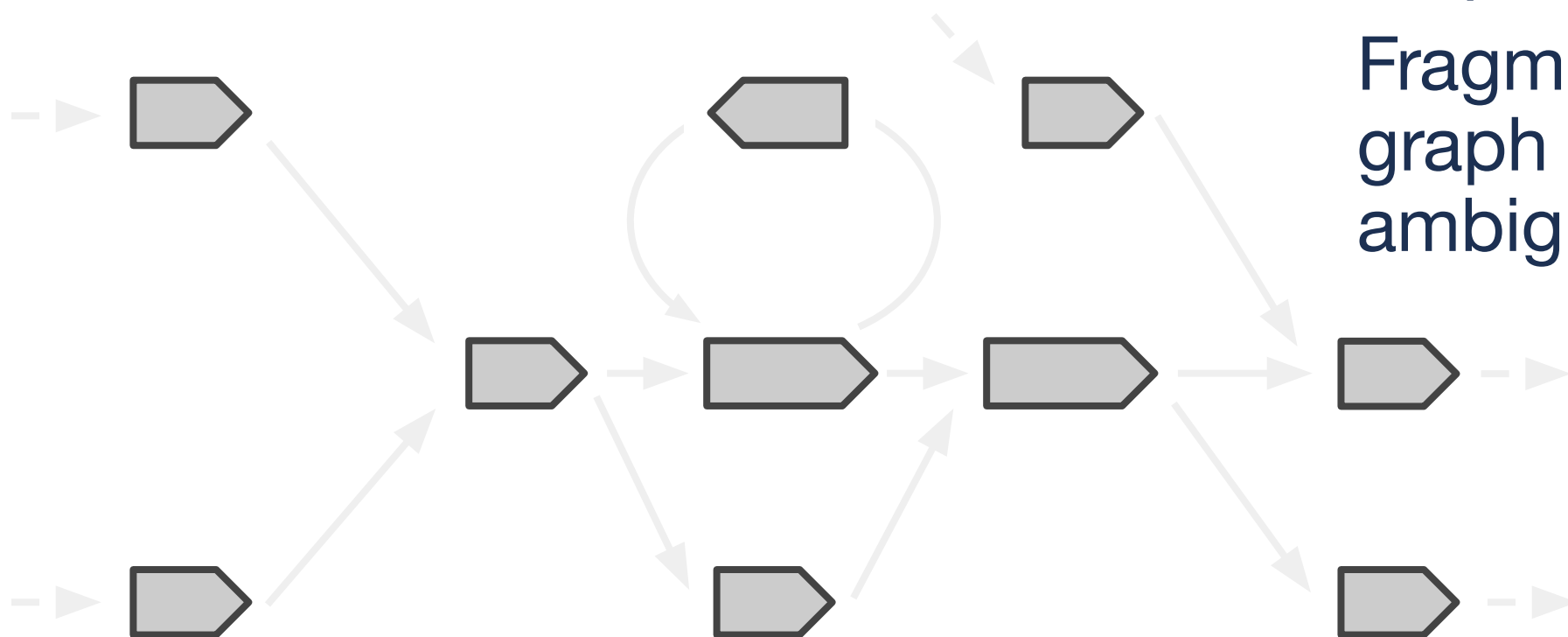# Lots of incorrect paths also exist…
*How do we avoid these?*

# Lots of incorrect paths also exist…
## *How do we avoid these?*

**Standard Tools:**
Filter out low-abundance sequences

# Lots of incorrect paths also exist...
*How do we avoid these?*

**Standard Tools:**

Filter out low-abundance sequences

Fragment the graph when it's ambiguous

Untangling the hairball

Untangling the hairball

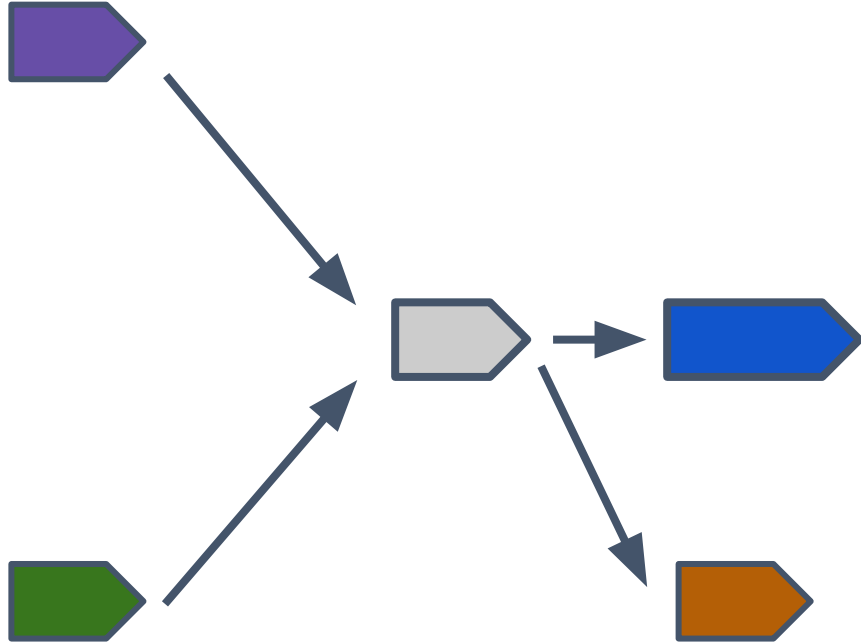Gladstone Scientific Retreat 2024
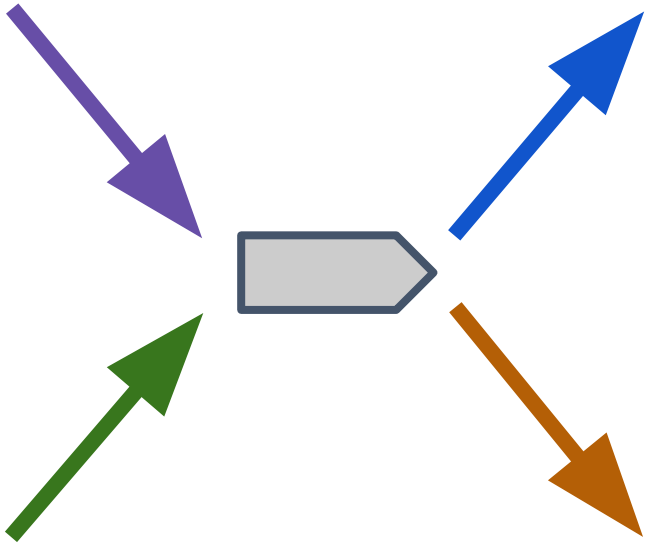
**StrainZip:**

Untangling the metagenome graph

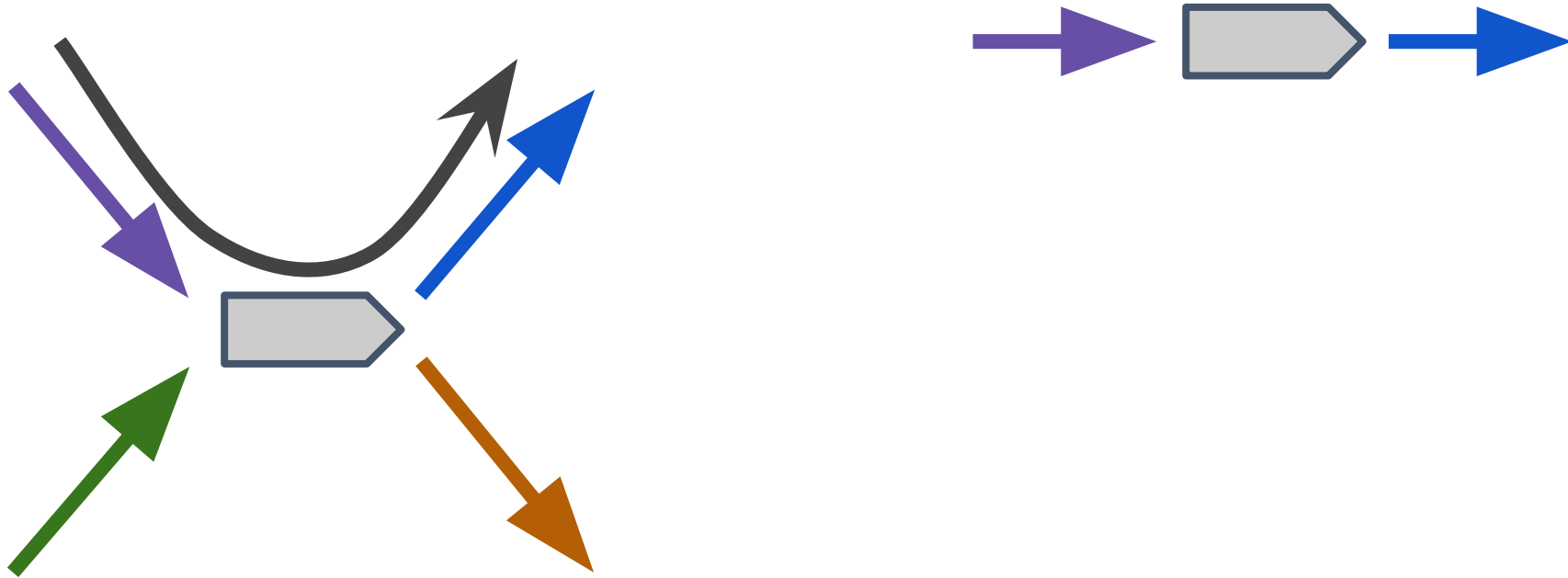# How can we recover long, accurate genome sequences from short reads?
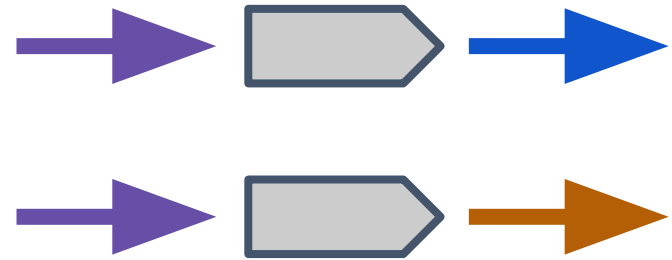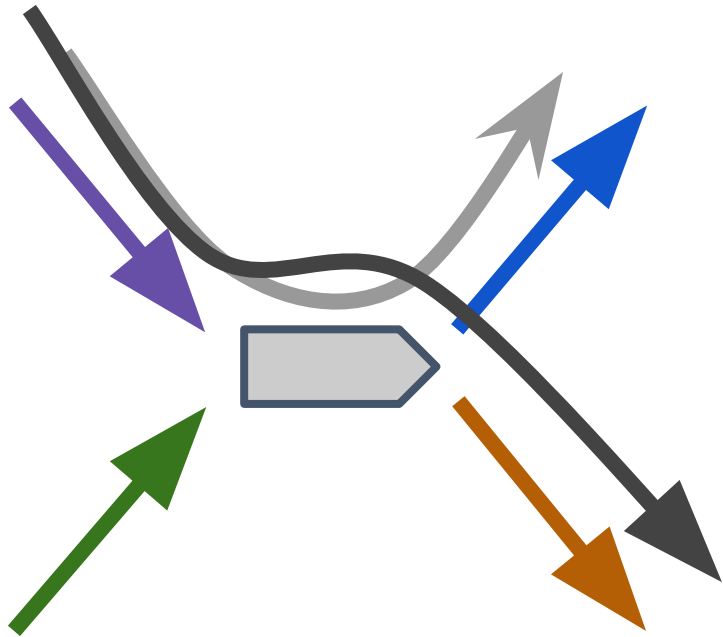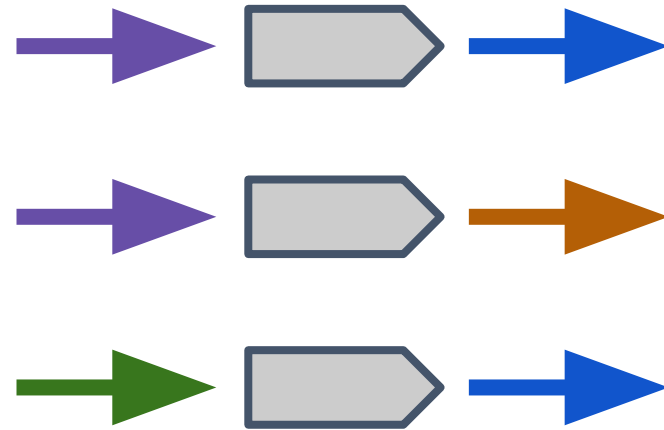
# How can we recover long, accurate genome sequences from short reads?

# Focus on just one junction at a time

# Focus on just one junction at a time
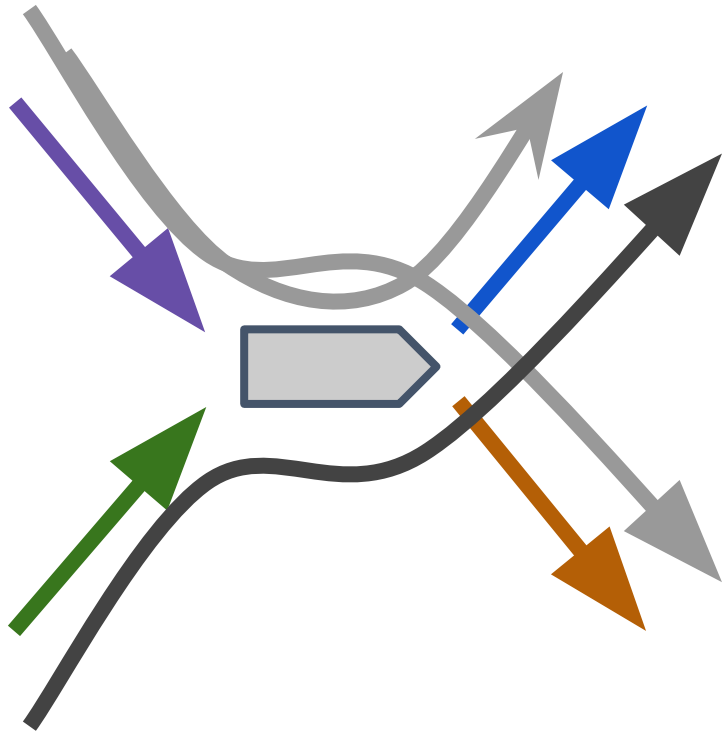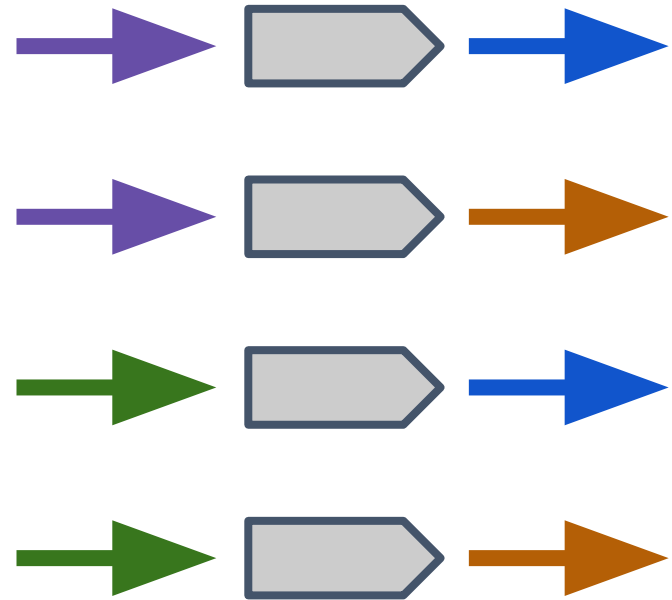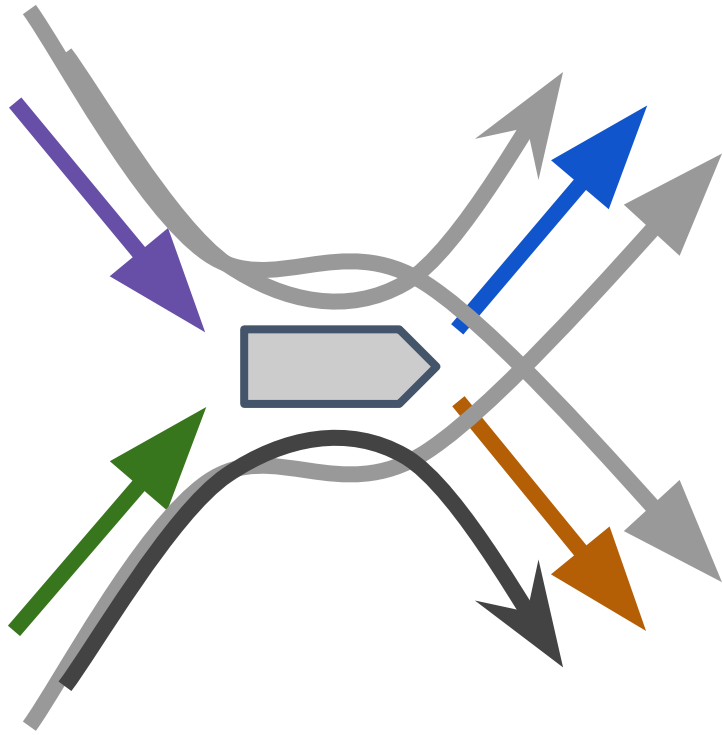
# Focus on just one junction at a time

# Focus on just one junction at a time

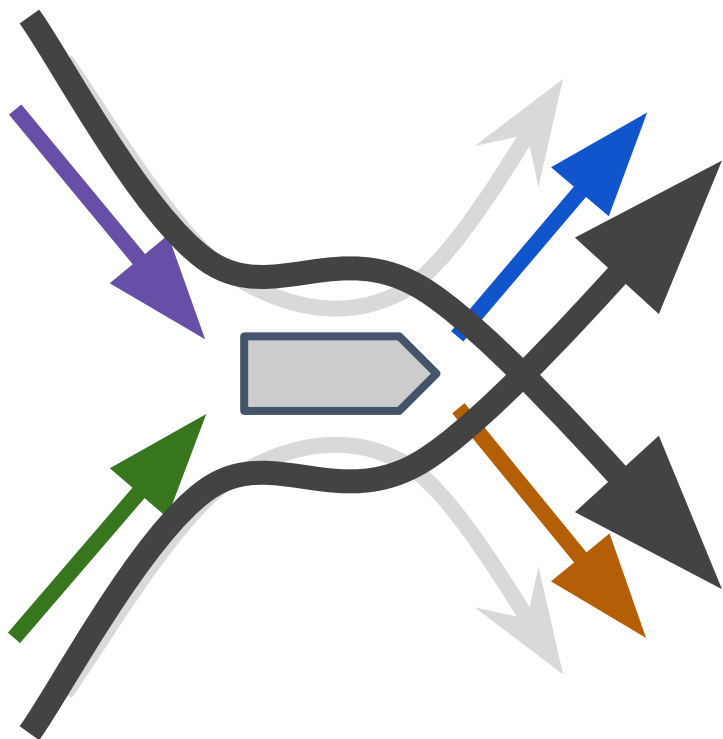# Focus on just one junction at a time

# Focus on just one junction at a time

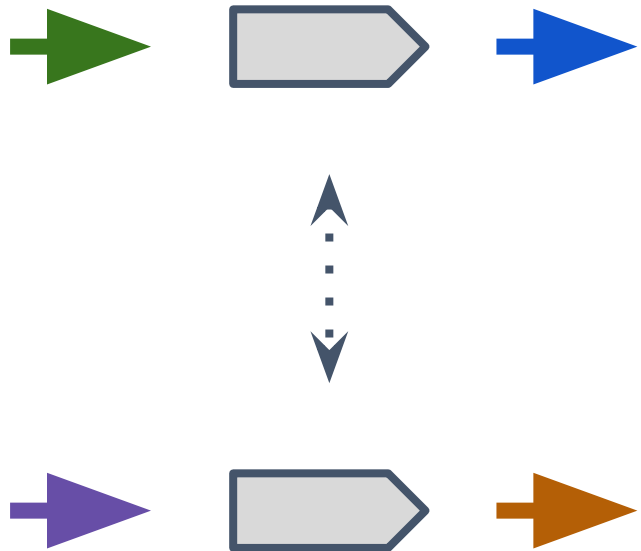# Focus on just one junction at a time
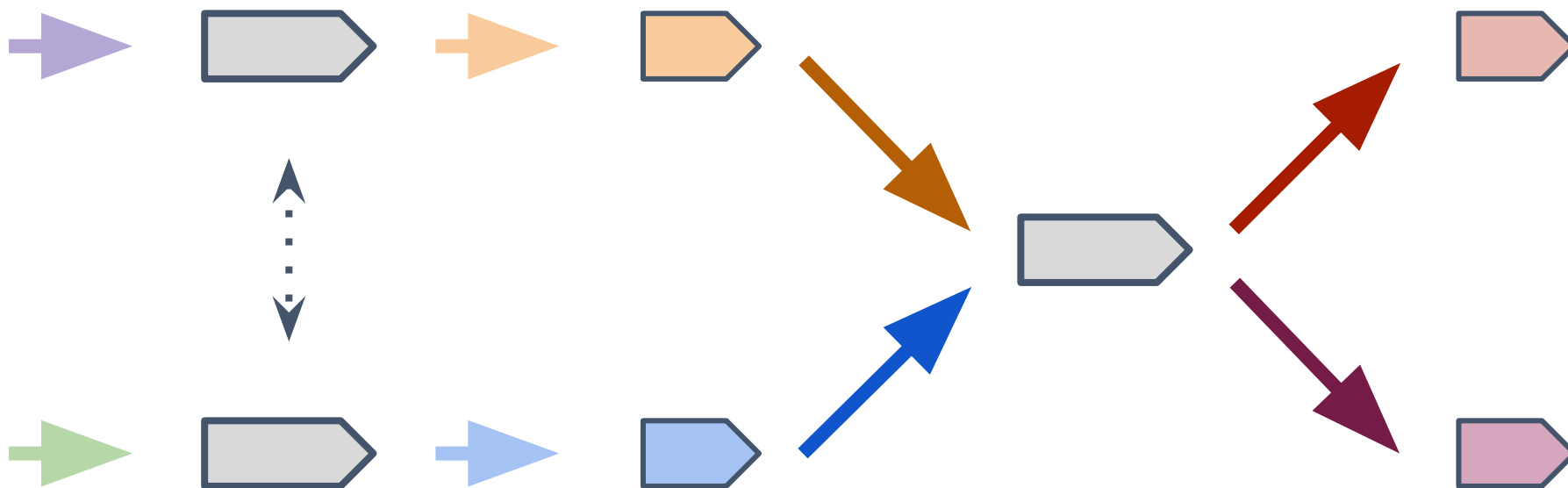# Select local paths



$X$  $\beta$  $Y$

Sparse linear regression across multiple samples

# Focus on just one junction at a time
# Select local paths
# Unzip

# Focus on just one junction at a time
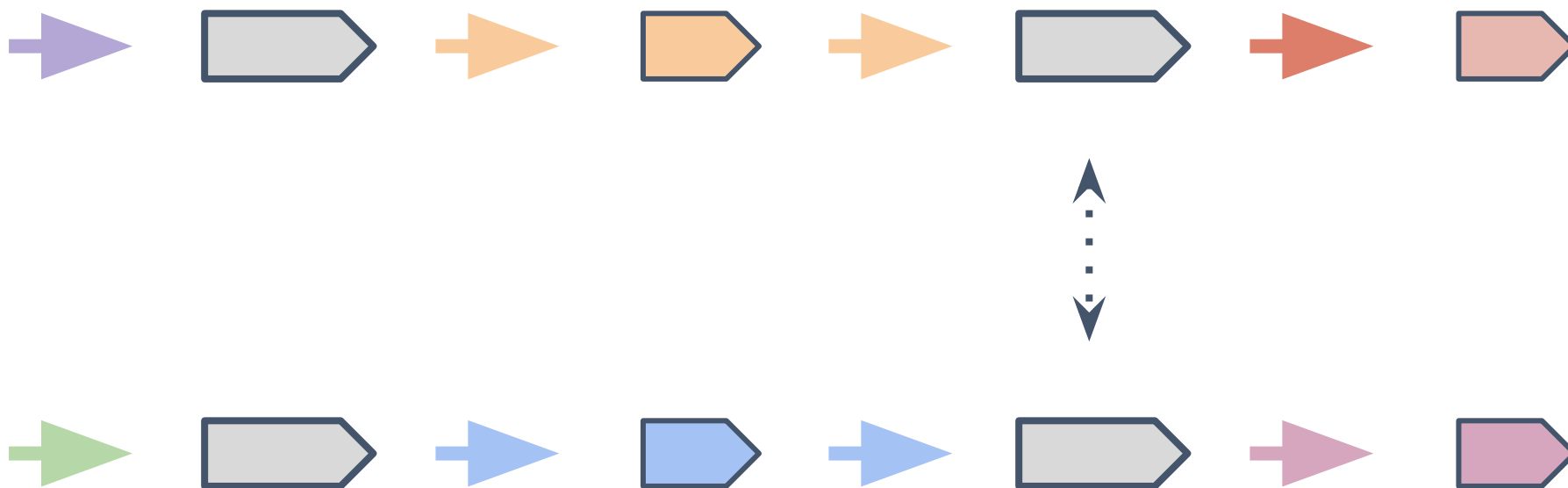# Select local paths
# Unzip
# Repeat

# Focus on just one junction at a time
## Select local paths
## Unzip
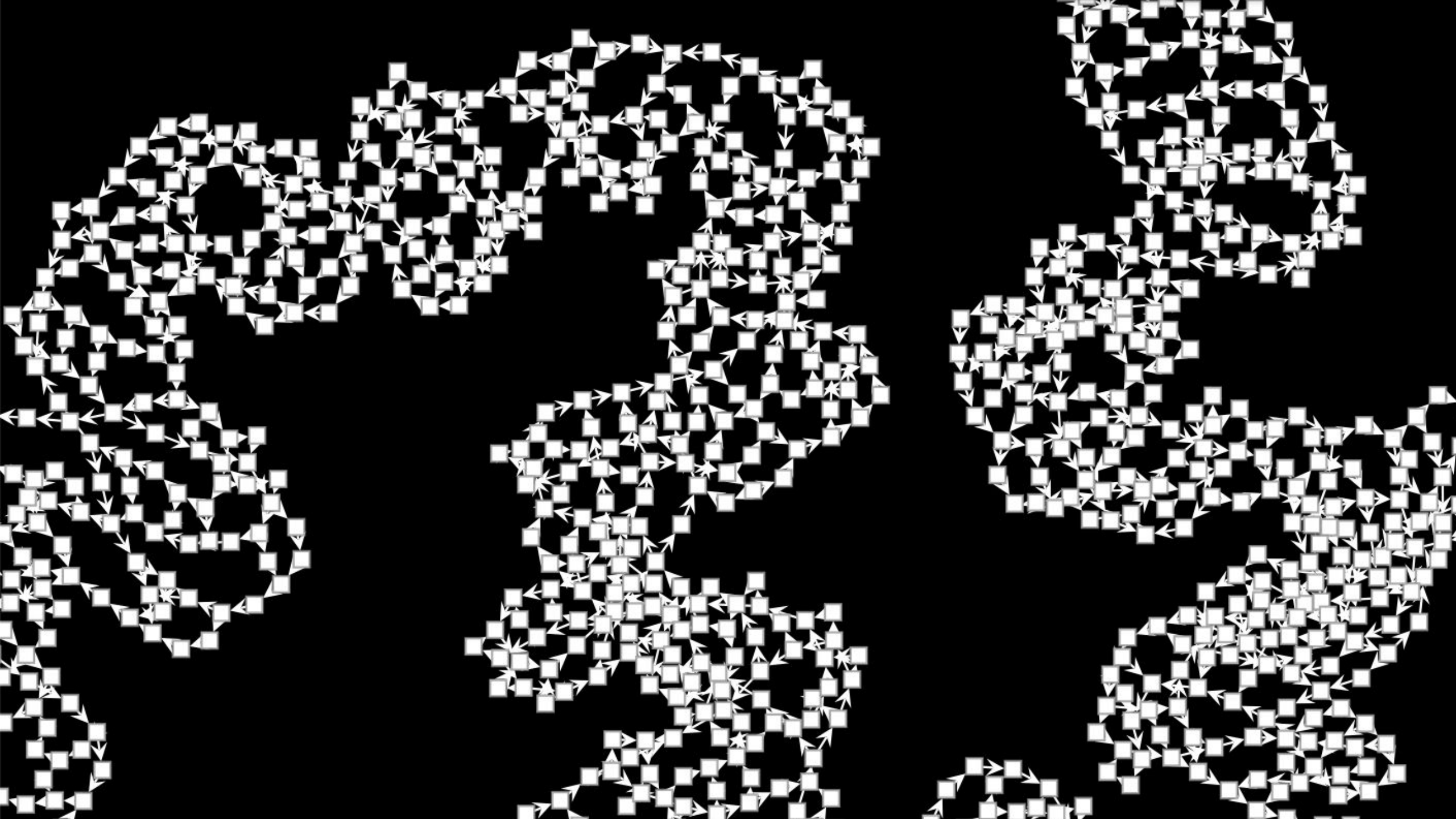## Repeat
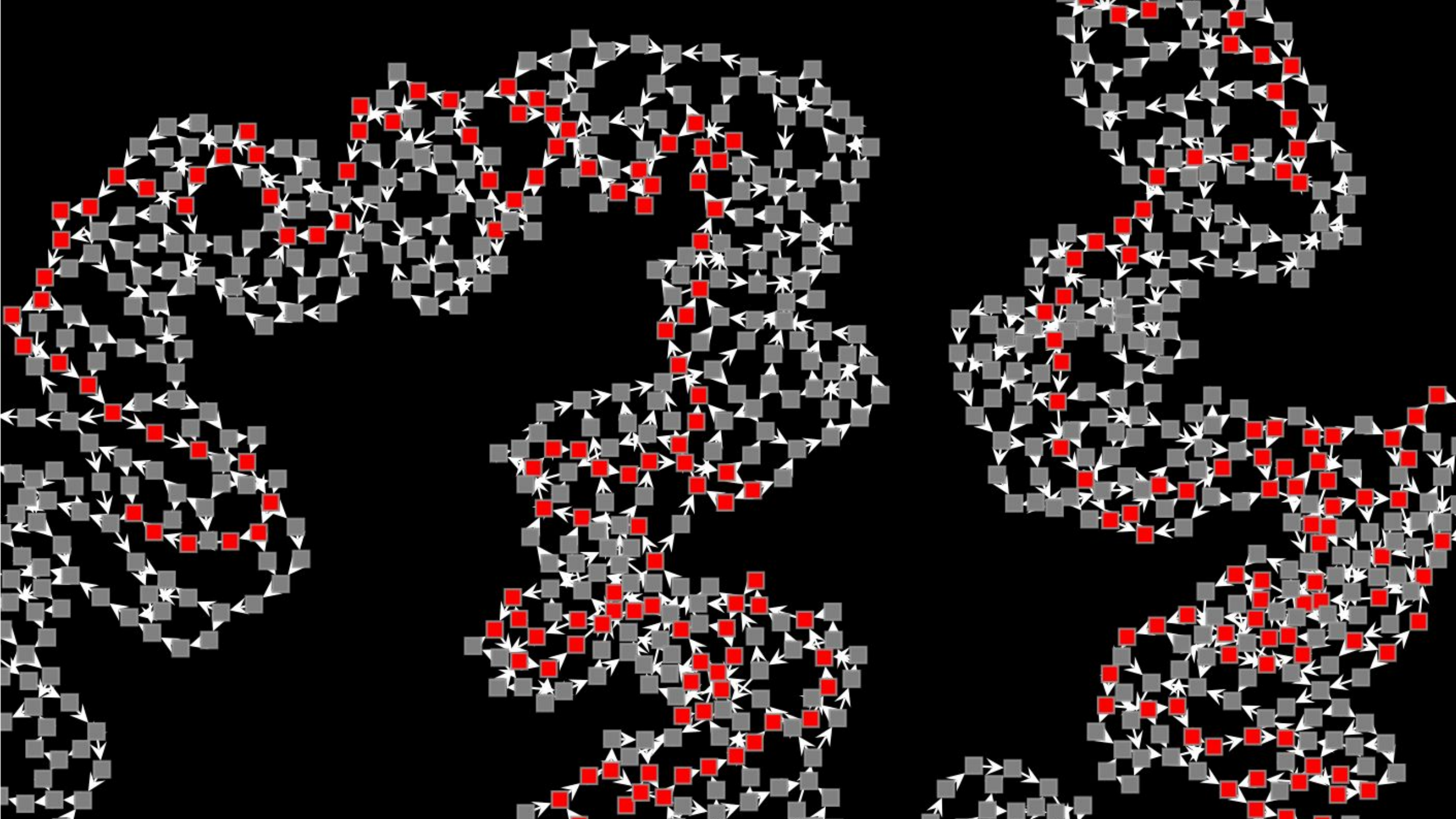
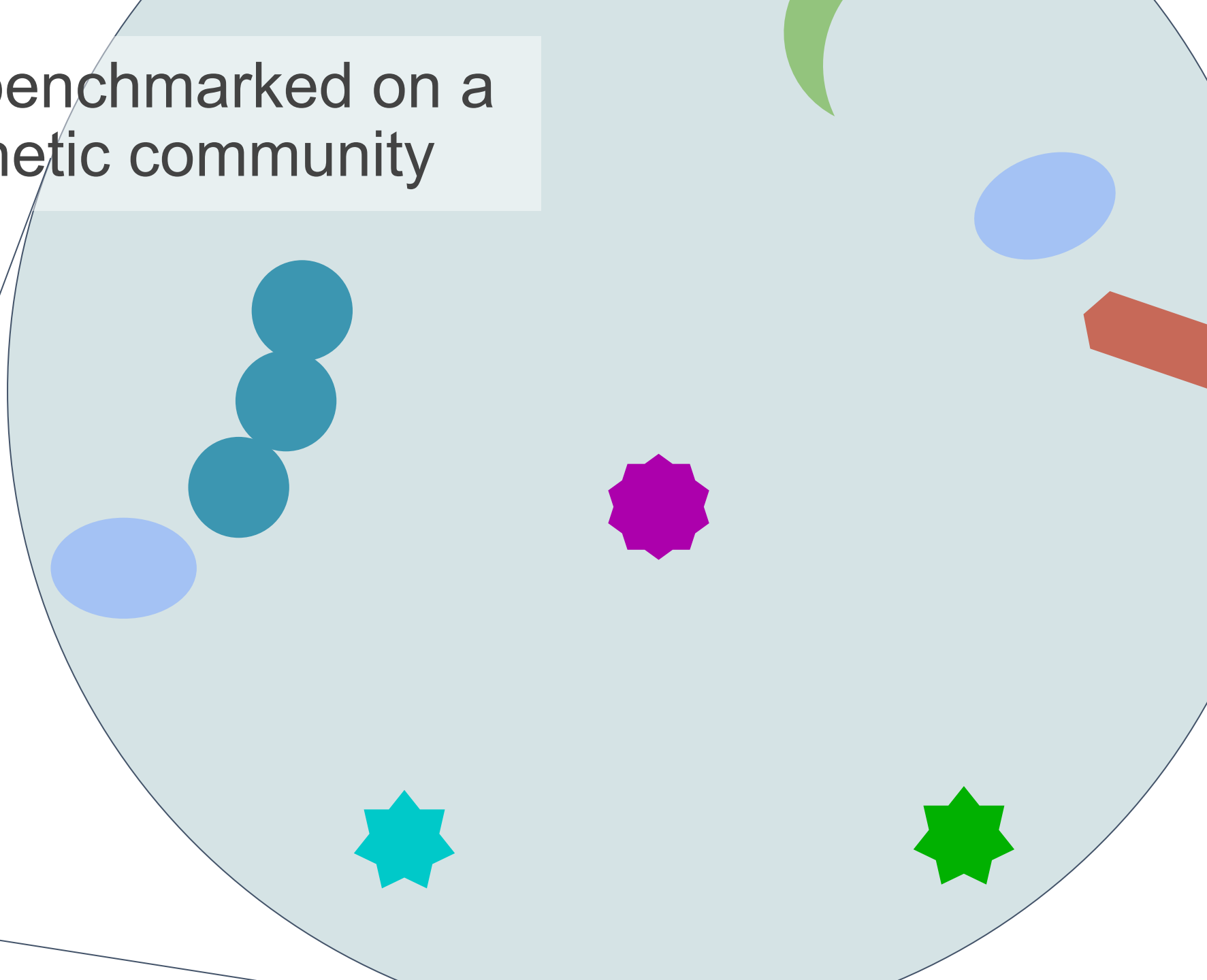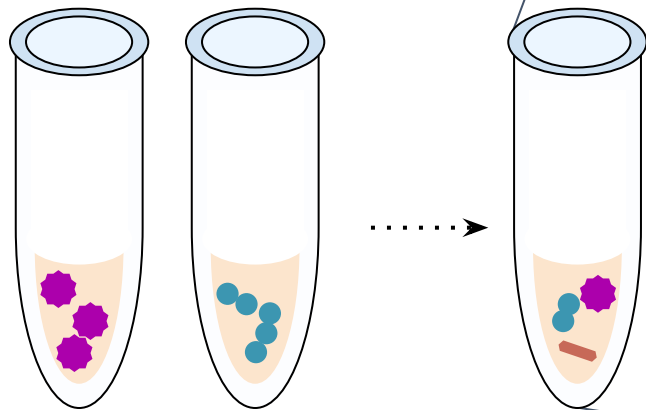Focus on just one junction at a time
Select local paths
Unzip
Repeat

Gladstone Scientific Retreat 2024

# Strain-resolved discovery

Performance benchmarked on a complex, synthetic community

Antibiotic resistance genes are widespread in the gut microbiome

**Antibiotic resistance genes are widespread in the gut microbiome**
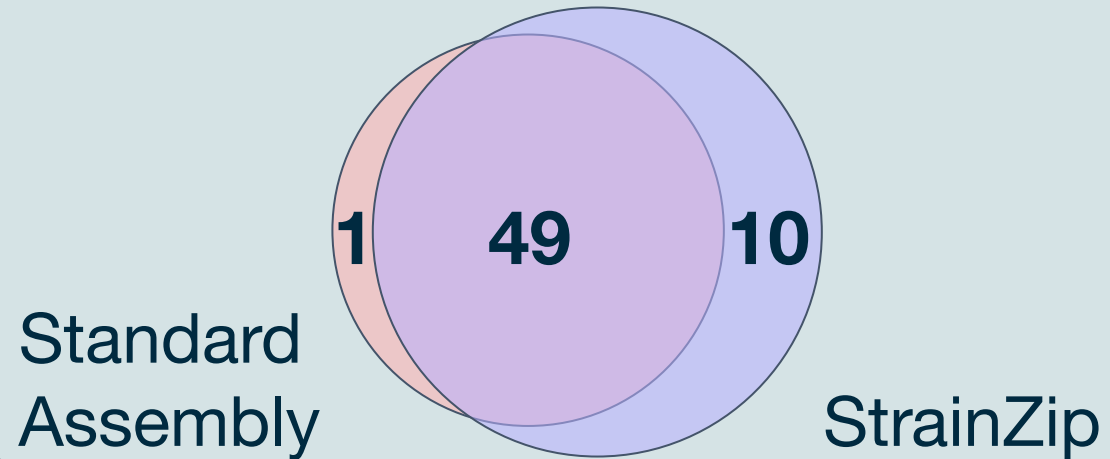
- Detection can inform treatment

**Antibiotic resistance genes are widespread in the gut microbiome**

- Detection can inform treatment

**No. of Unique Resistance Genes Found**

1   49   10

Standard Assembly   StrainZip

**Antibiotic resistance genes are widespread in the gut microbiome**
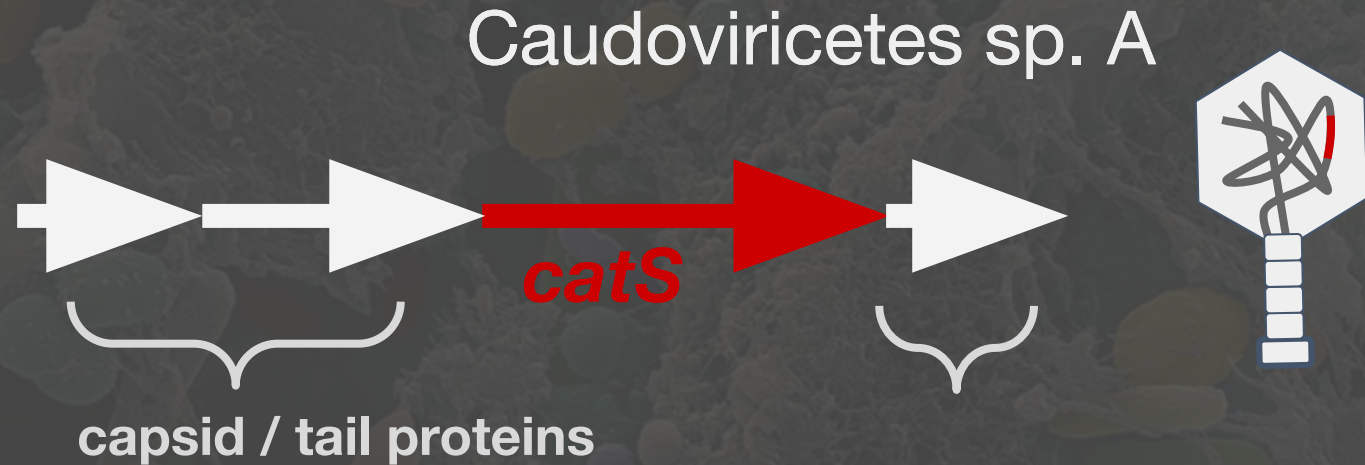
- Detection can inform treatment

*catS*

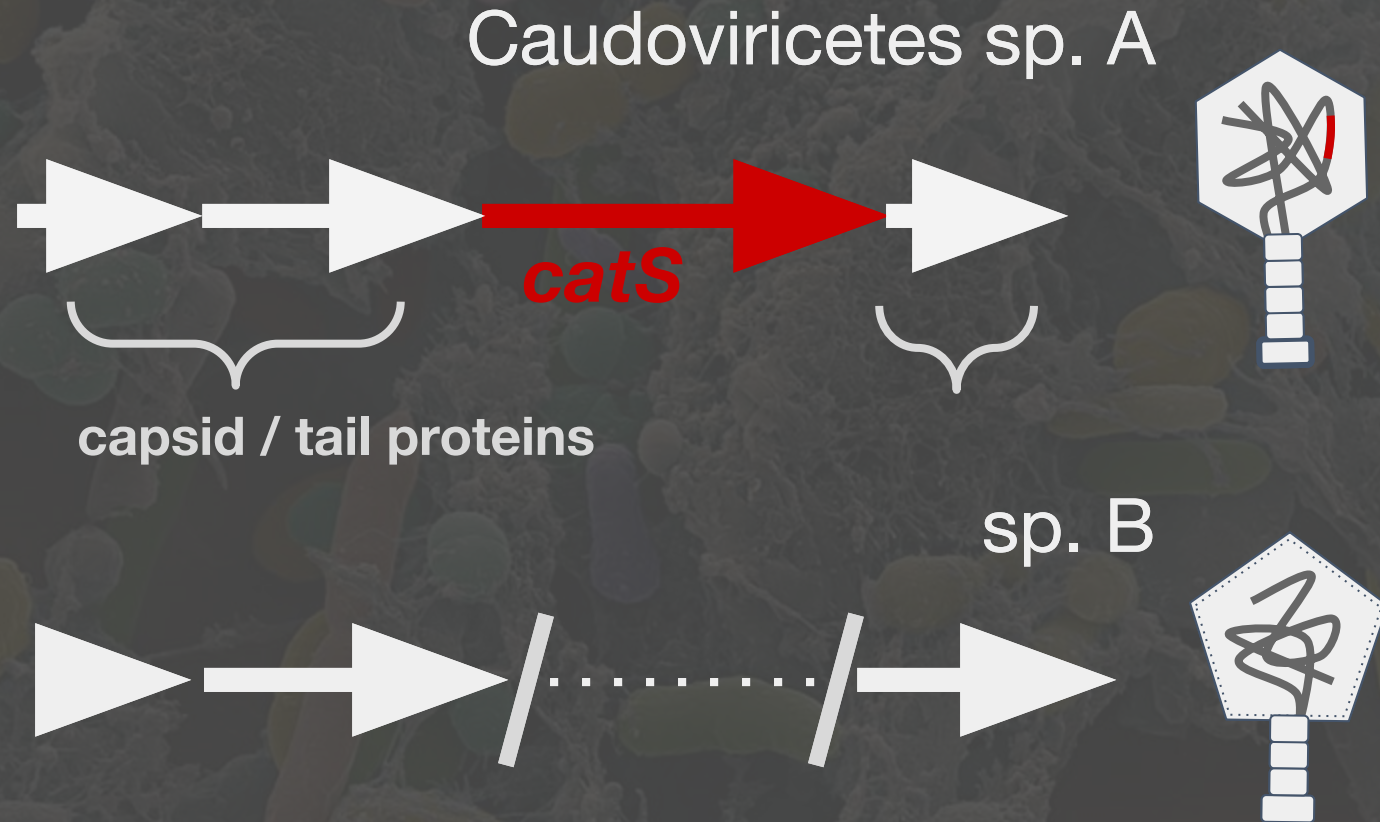**Antibiotic resistance genes are widespread in the gut microbiome**

- Detection can inform treatment

- Can be carried in phage genomes

- Long sequence fragments provide useful information

Caudoviricetes sp. A

*catS*

capsid / tail proteins

sp. B

# Antibiotic resistance genes are widespread in the gut microbiome
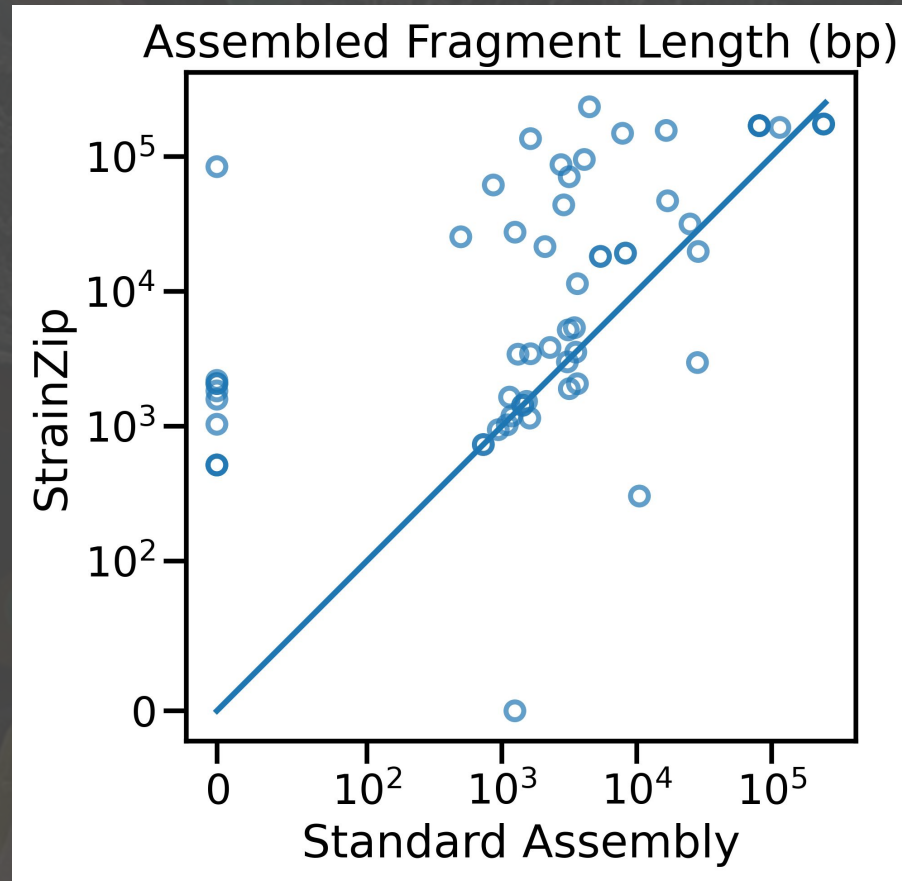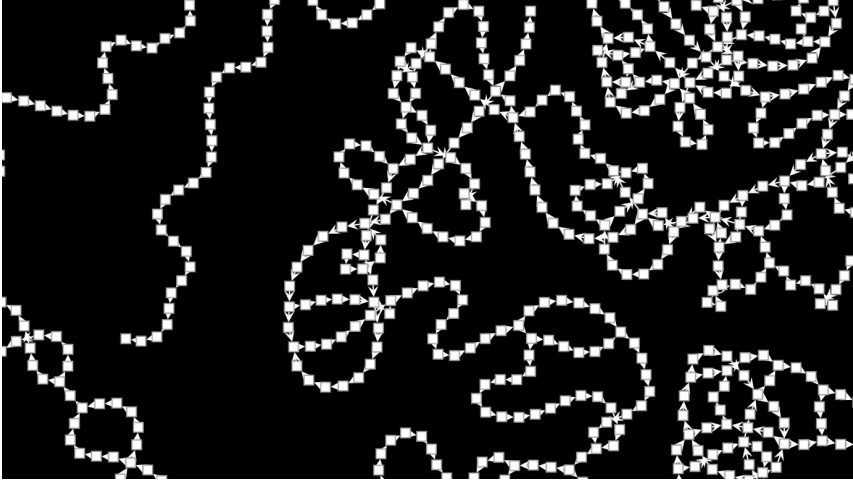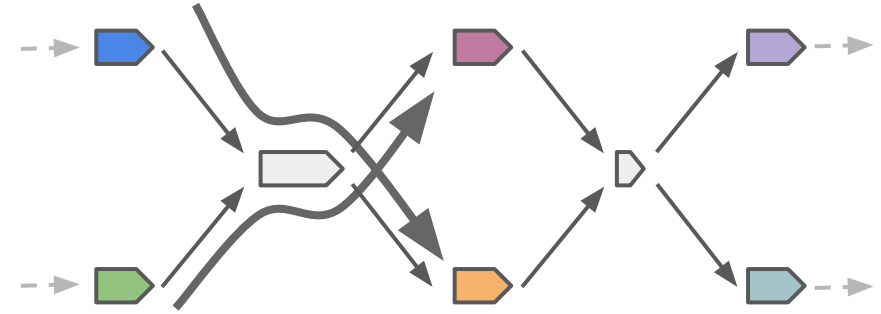
- Detection can inform treatment

- Can be carried in phage genomes

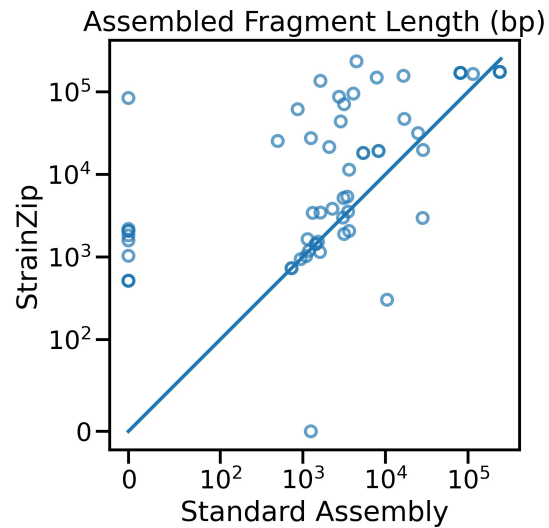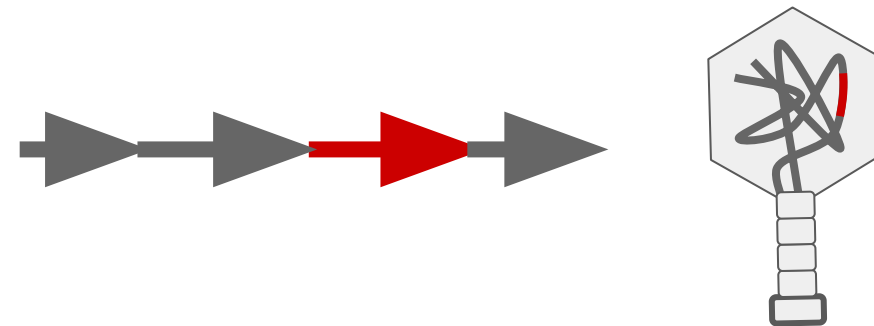- Long sequence fragments provide useful information

**Complex Metagenome Graphs**

**StrainZip Iteratively Unzips Junctions**

**Strain-Resolved Metagenomics**

Assembled Fragment Length (bp)

**Antibiotic Resistance Potential of Phage**

Gladstone Scientific Retreat 2024

# Thank You!
# Questions?