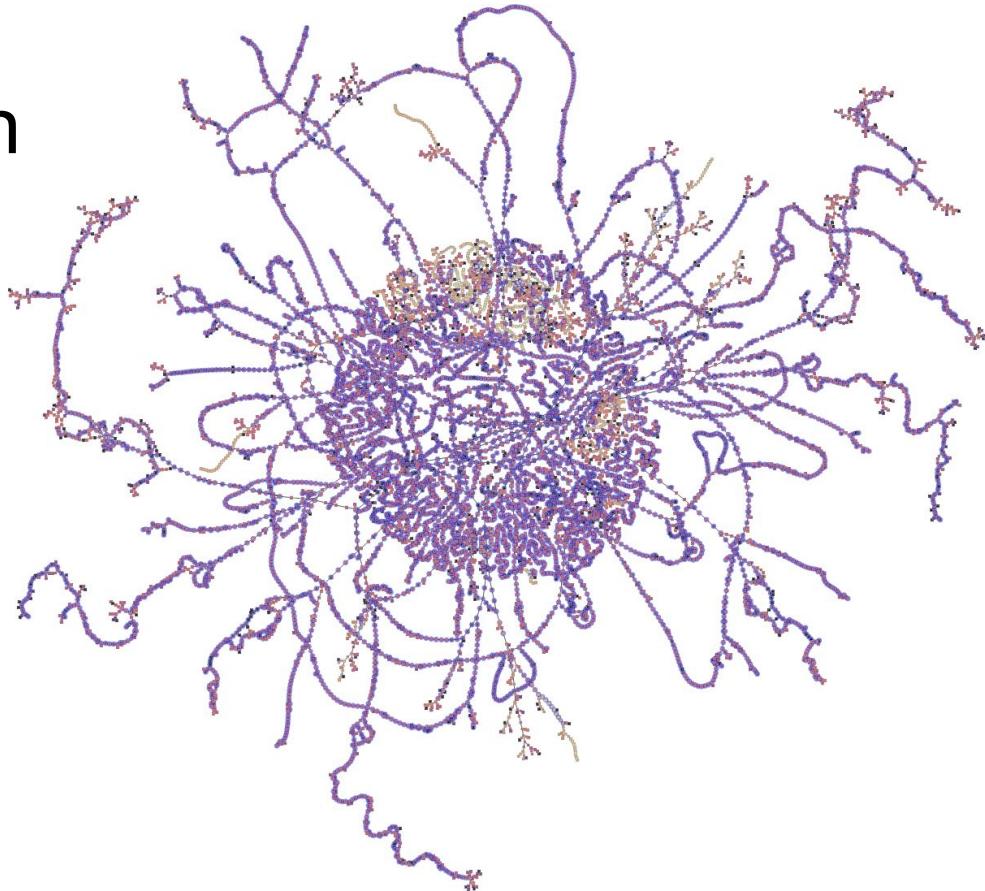


# Disentangling depth on the De Bruijn graph:

Combining quantification  
and *de novo* assembly

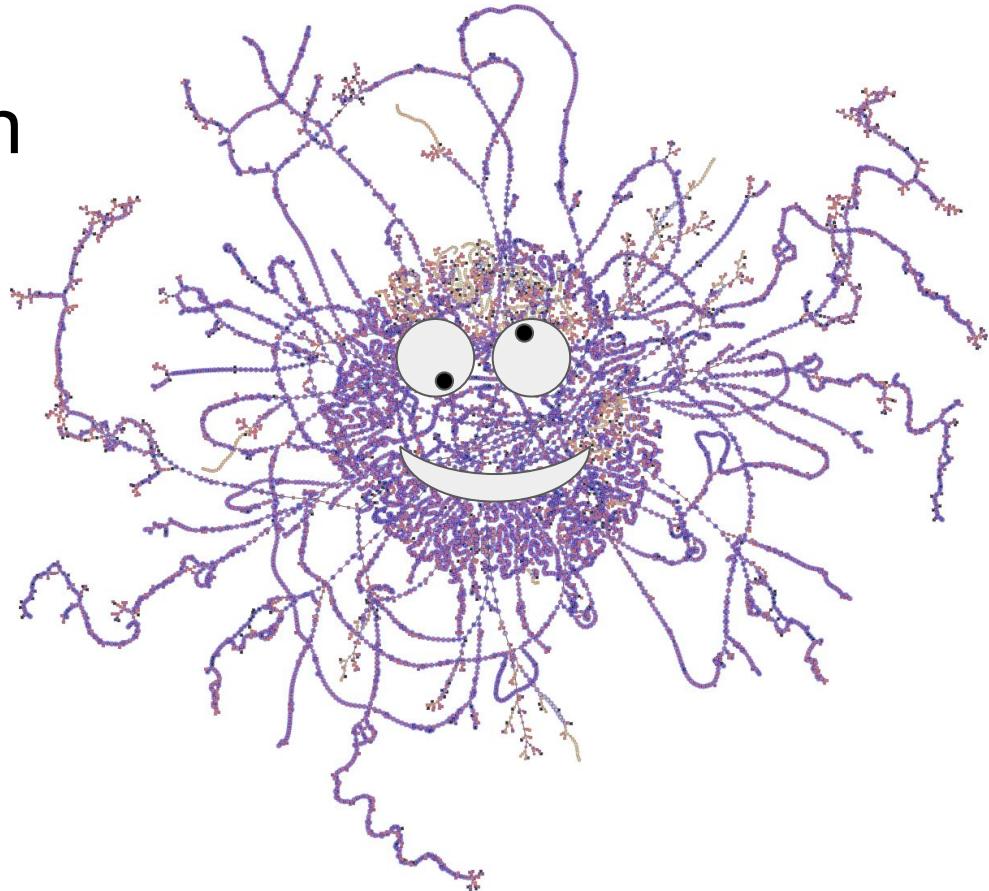
Byron J. Smith  
IGGSy  
2024-07-04



# Disentangling depth on the De Bruijn graph:

Combining quantification  
and *de novo* assembly

Byron J. Smith  
IGGSy  
2024-07-04



# Acknowledgments / Find Me



# Acknowledgments / Find Me

## Pollard Lab

- *Katie Pollard*
- Veronika Dubinkina
- et al.

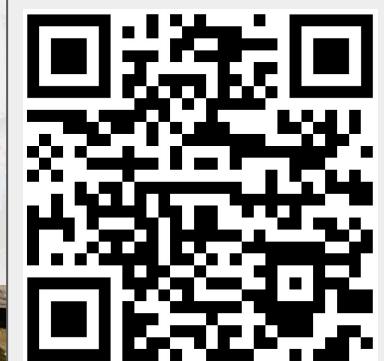
## Engelhardt Lab

- *Archit Verma*
- Dylan Cable

**GLADSTONE**  
INSTITUTES

**Twitter:** @ByronJSmith

**Email:** [byron.smith  
@gladstone.ucsf.edu](mailto:byron.smith@gladstone.ucsf.edu)



# Outline of the presentation

- **Motivation:**

Profiling the microbiome with metagenomics

- **Method:**

Assembly graph deconvolution with StrainZip

- **Demonstration:**

Application to a complex, ground-truthed dataset



# Motivation

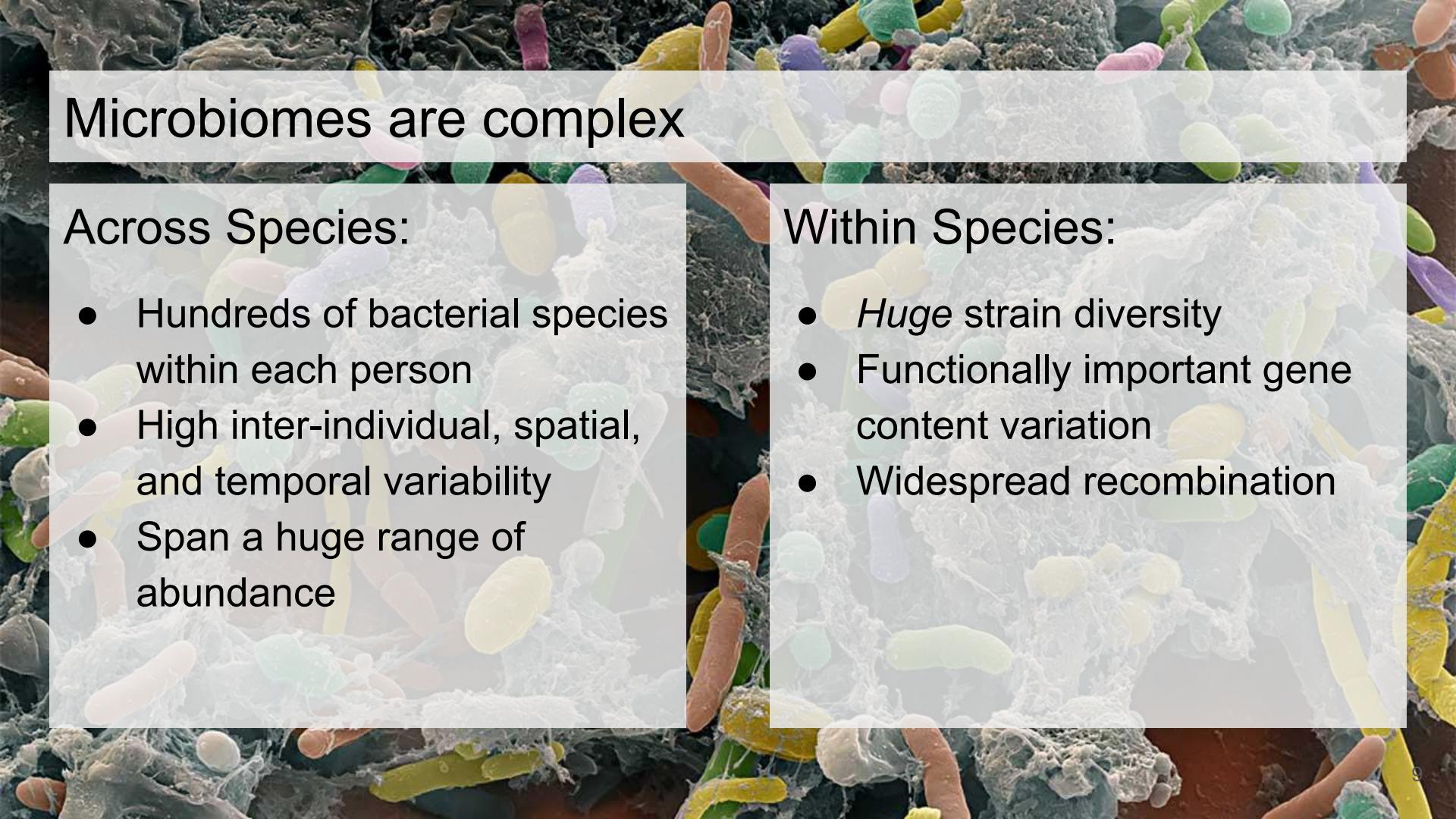


Microbiomes are complex

# Microbiomes are complex

## Across Species:

- Hundreds of bacterial species within each person
- High inter-individual, spatial, and temporal variability
- Span a huge range of abundance



# Microbiomes are complex

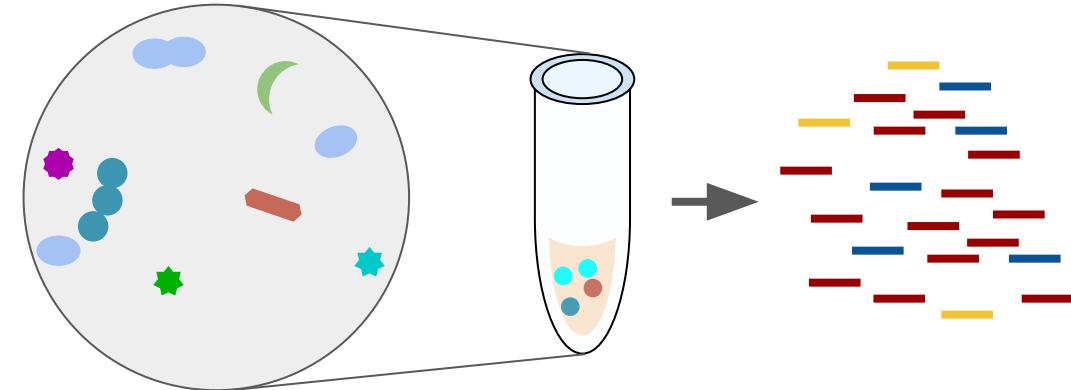
## Across Species:

- Hundreds of bacterial species within each person
- High inter-individual, spatial, and temporal variability
- Span a huge range of abundance

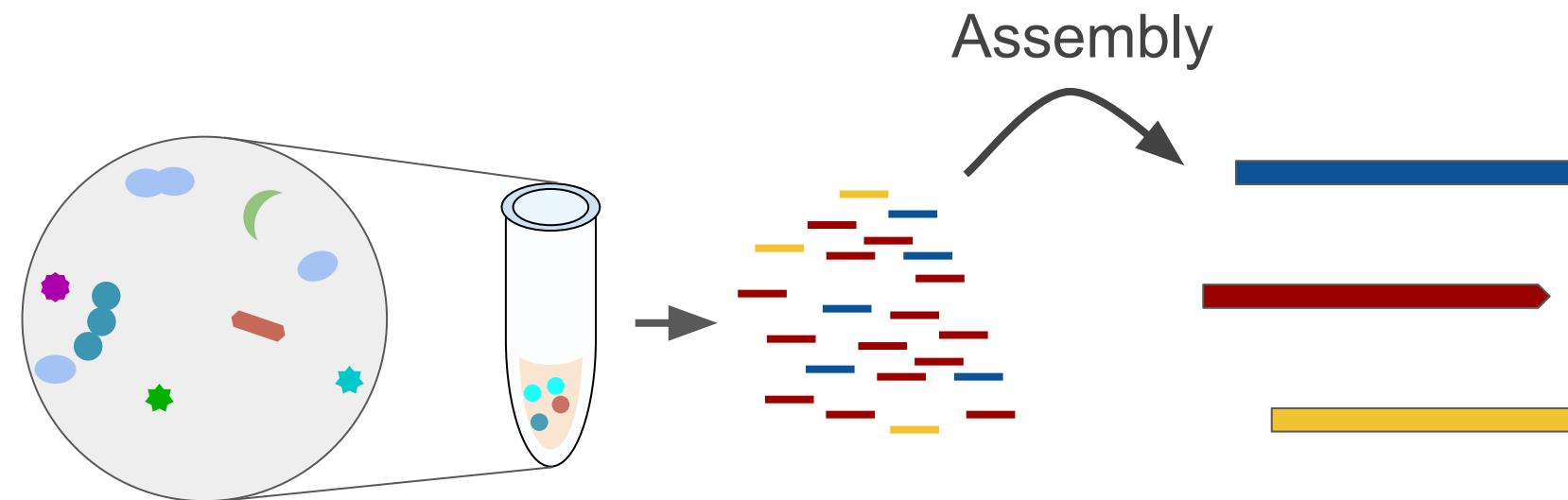
## Within Species:

- *Huge* strain diversity
- Functionally important gene content variation
- Widespread recombination

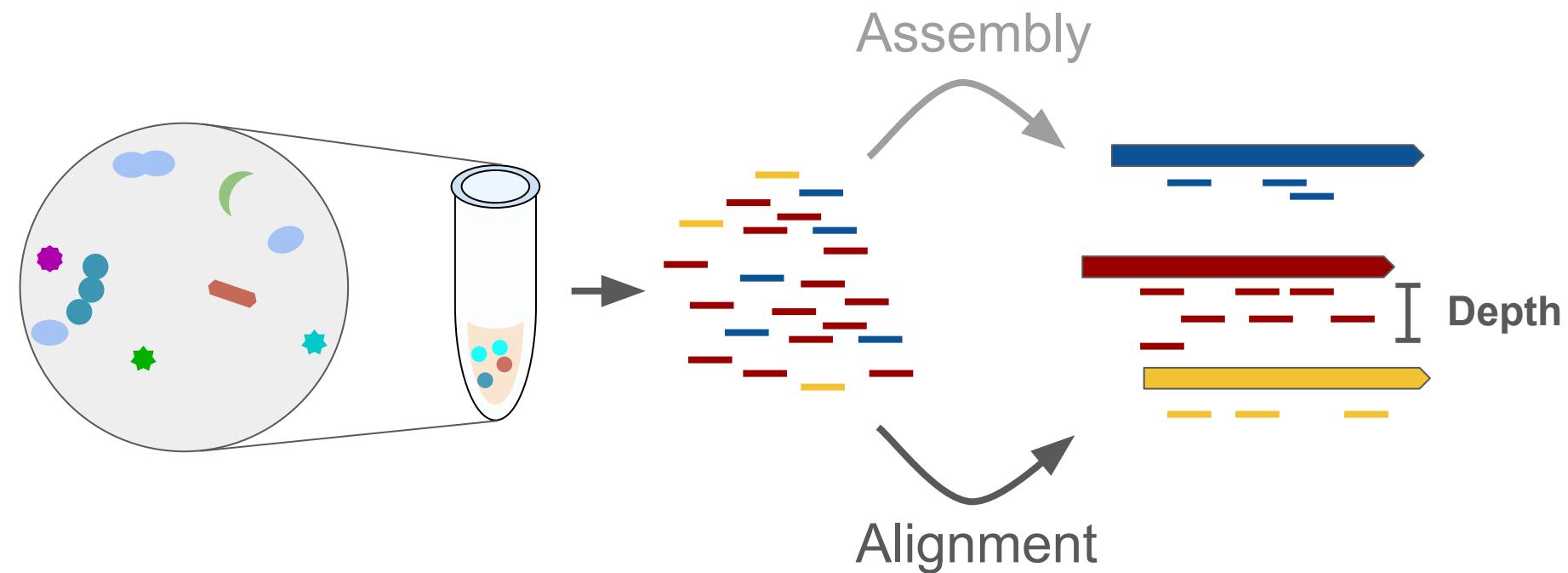
# Metagenomics enables modern microbiome science



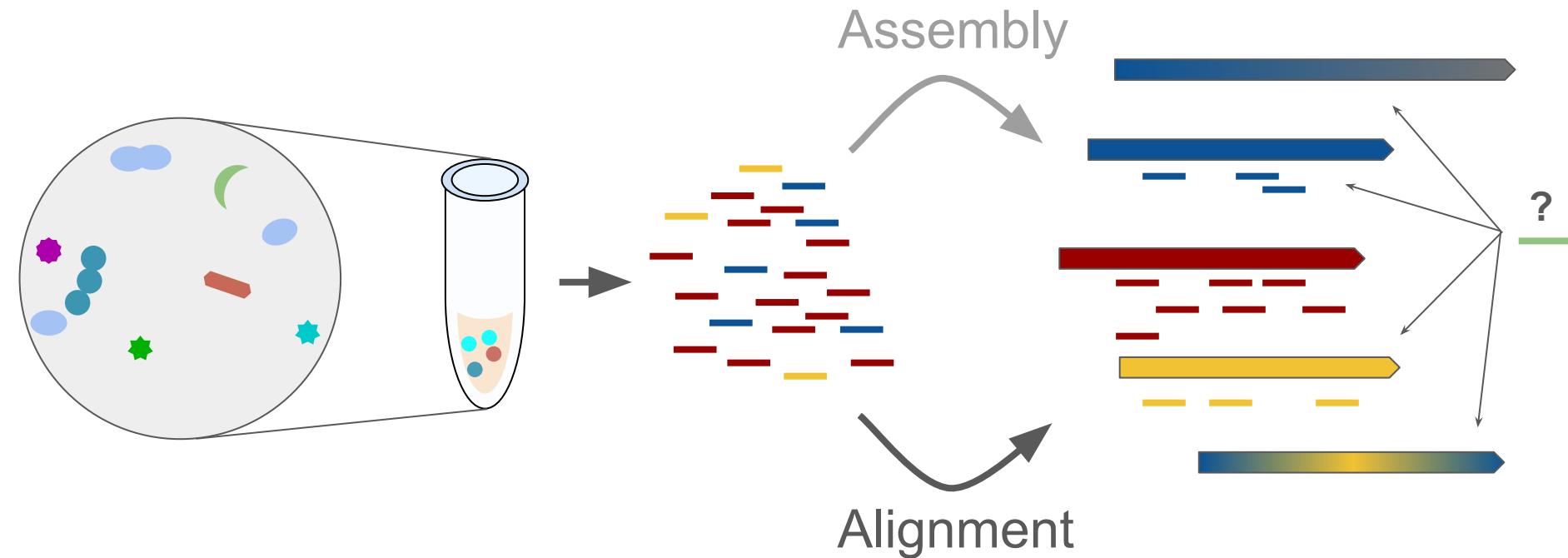
# Metagenomics enables modern microbiome science

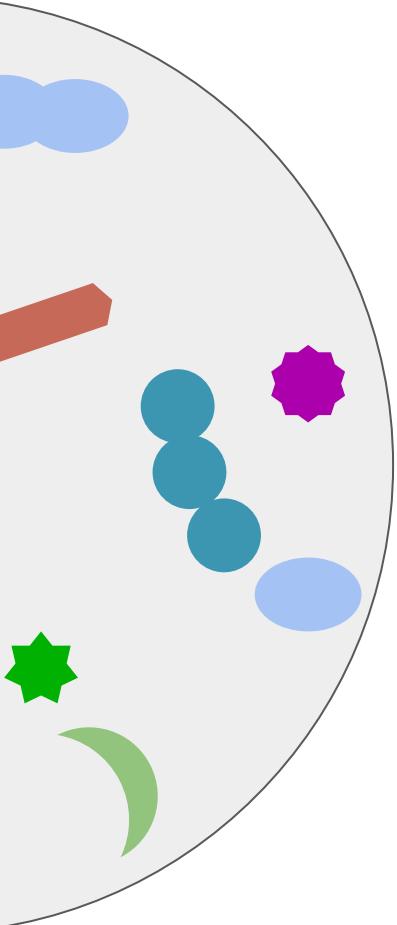


# Assembly and depth quantification are complementary



# Closely related sequences are a major challenge for alignment



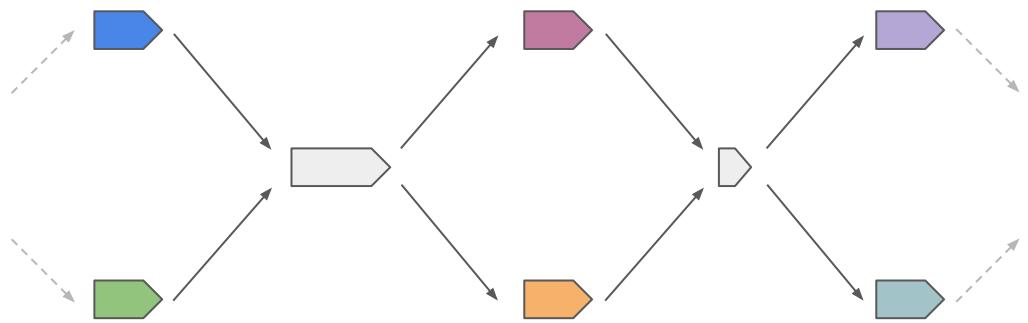


***Goal:***  
*Reference-free depth  
estimation among related  
sequences in metagenomes*

# Method

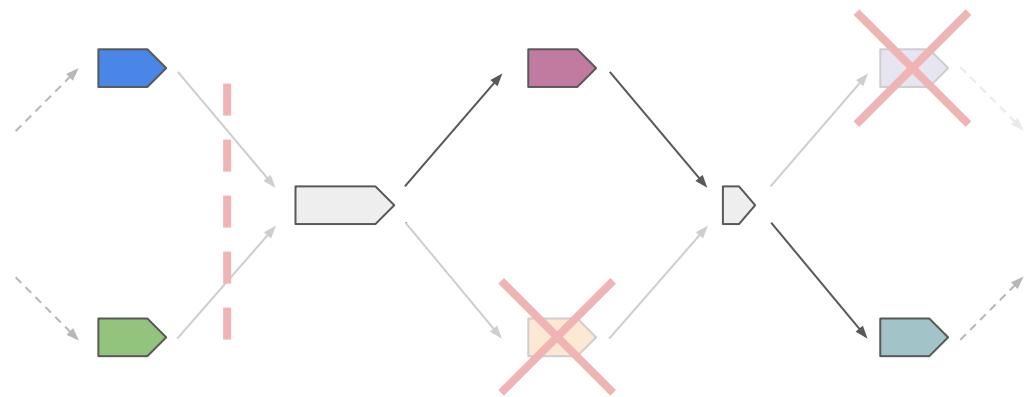
Closely related sequences are a major challenge for both metagenomic assembly and alignment

Complex graph structure

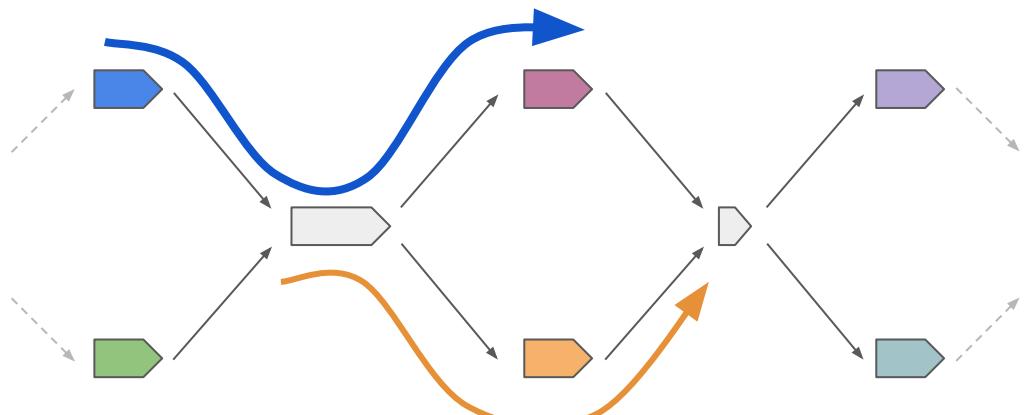


# Closely related sequences are a major challenge for alignment

Complex graph structure  
leads to low-quality assembly



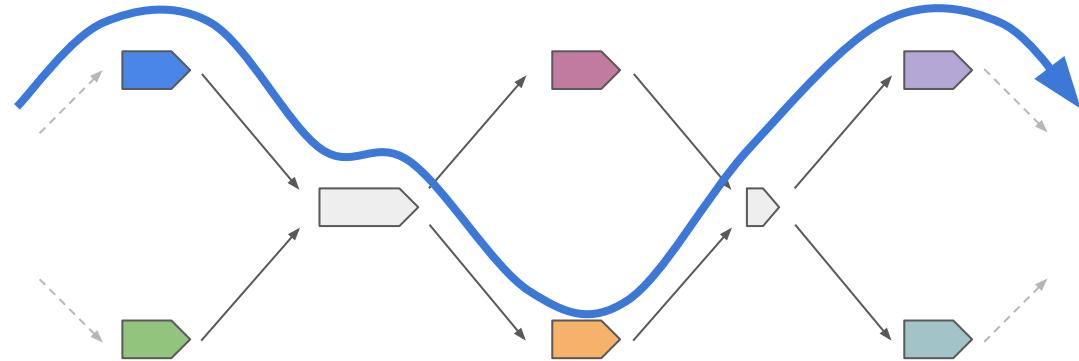
# Closely related sequences are a major challenge for alignment



Complex graph structure  
leads to low-quality assembly

Graph-pangenome  
approaches account for this  
variability better than linear  
references

# Closely related sequences are a major challenge for alignment

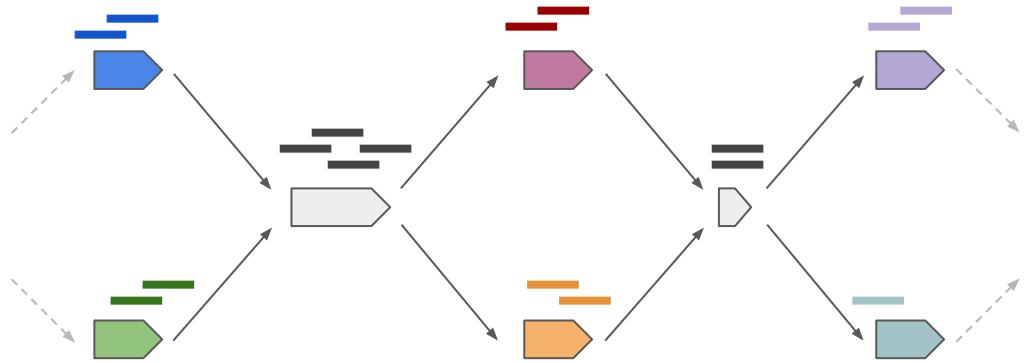


Complex graph structure  
leads to low-quality assembly

Graph-pangenome  
approaches account for this  
variability better than linear  
references

But long reads too expensive  
for profiling multiple samples

# Closely related sequences are a major challenge for alignment



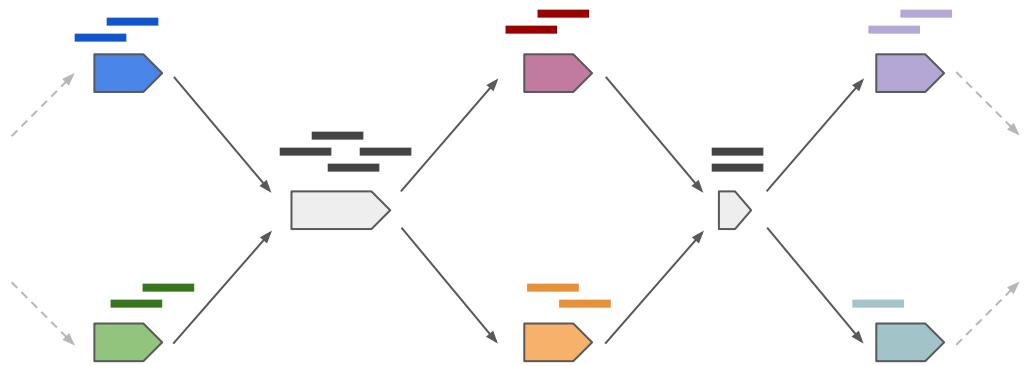
Complex graph structure  
leads to low-quality assembly

Graph-pangenome  
approaches account for this  
variability better than linear  
references

But long reads too expensive  
for profiling multiple samples

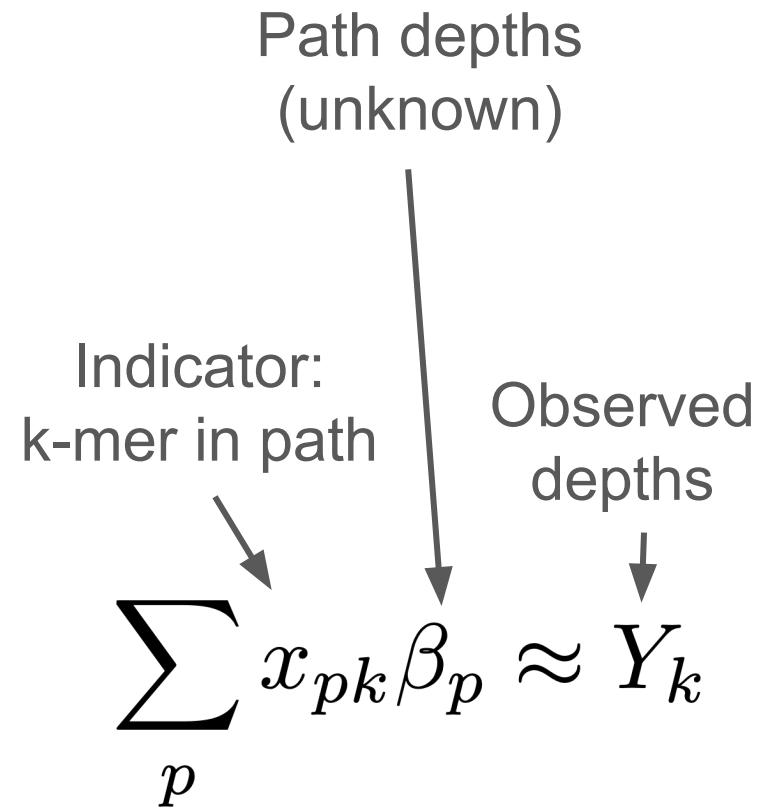
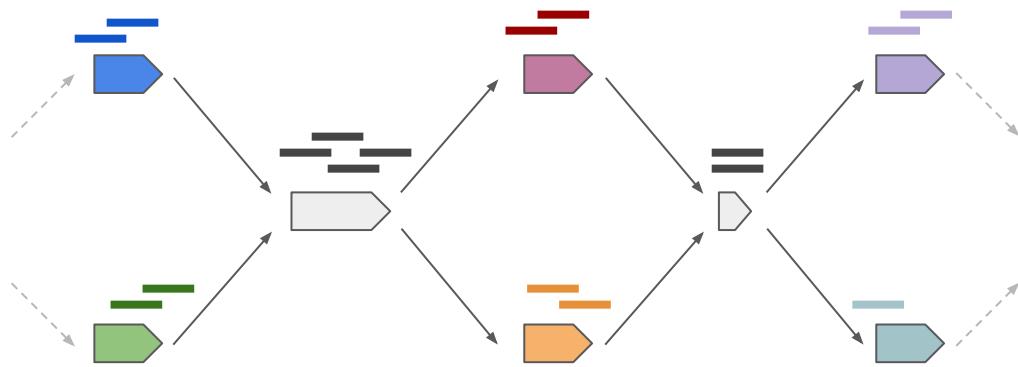
And short-reads are  
inherently ambiguous

**KEY IDEA:** The expected depth of a k-mer is the sum of the paths that include that k-mer

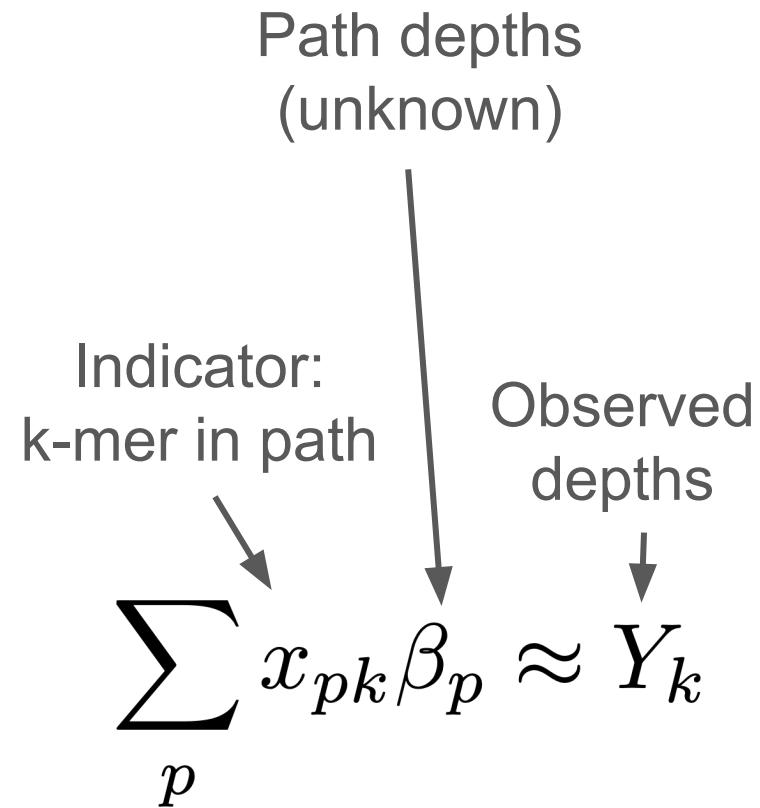
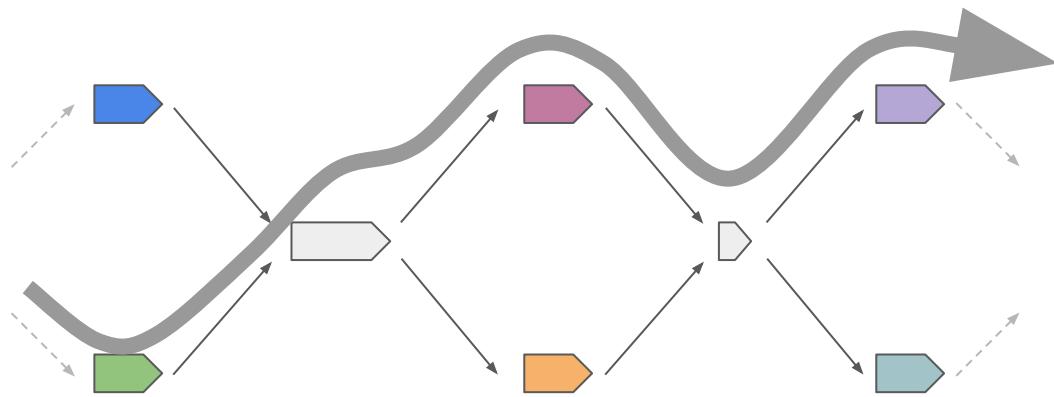


$$\sum_p x_{pk} \beta_p \approx Y_k$$

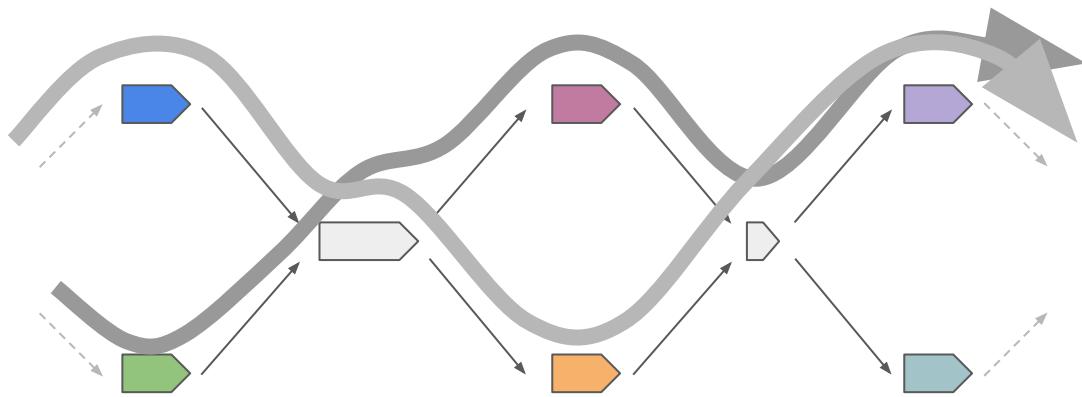
**KEY IDEA:** The expected depth of a k-mer is the sum of the paths that include that k-mer



**KEY IDEA:** The expected depth of a k-mer is the sum of the paths that include that k-mer



**KEY IDEA:** The expected depth of a k-mer is the sum of the paths that include that k-mer



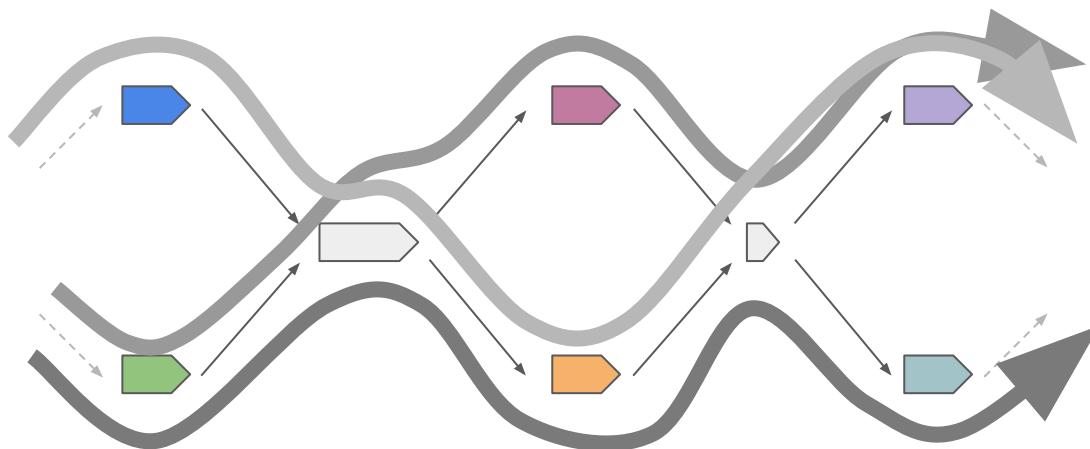
Path depths  
(unknown)

Indicator:  
k-mer in path

Observed  
depths

$$\sum_p x_{pk} \beta_p \approx Y_k$$

**KEY IDEA:** The expected depth of a k-mer is the sum of the paths that include that k-mer



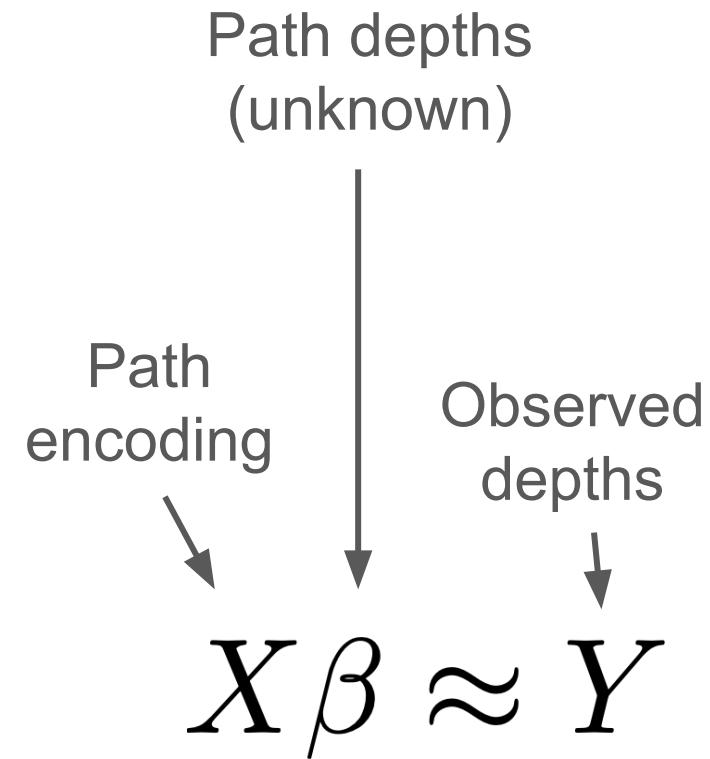
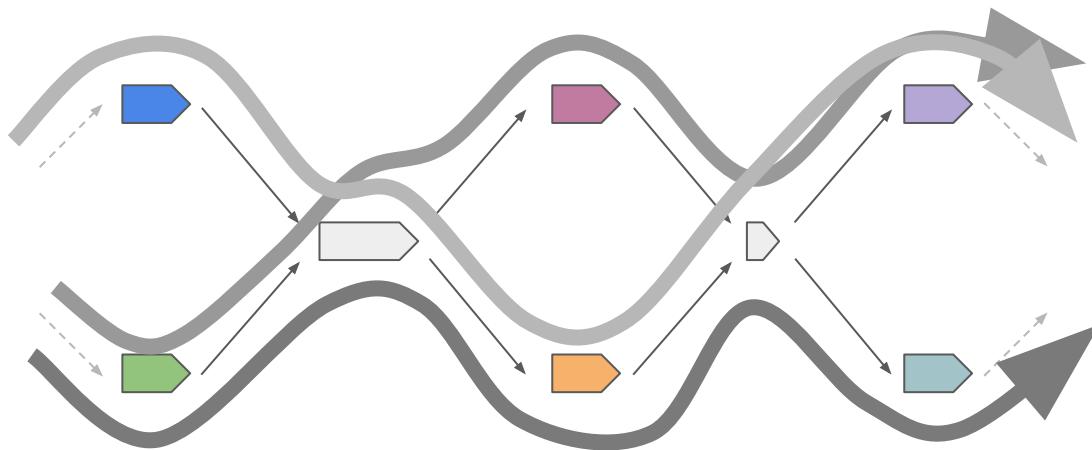
Path depths  
(unknown)

Indicator:  
k-mer in path

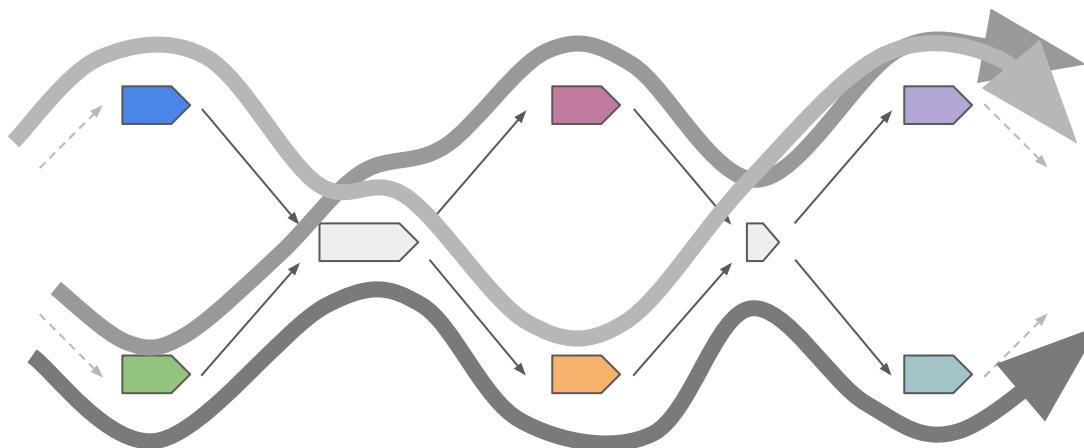
Observed depths

$$\sum_p x_{pk} \beta_p \approx Y_k$$

**KEY IDEA:** The expected depth of a k-mer is the sum of the paths that include that k-mer



**KEY IDEA:** The expected depth of a k-mer is the sum of the paths that include that k-mer

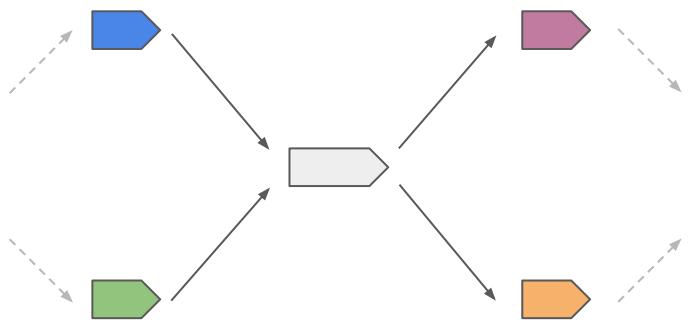


**Deconvolution:** Inferring the depth of these latent paths based on observed k-mer depths

Estimate this  
From these

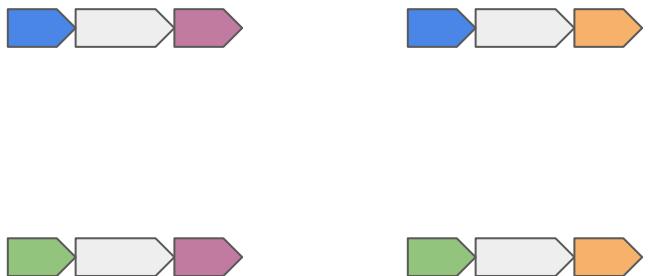
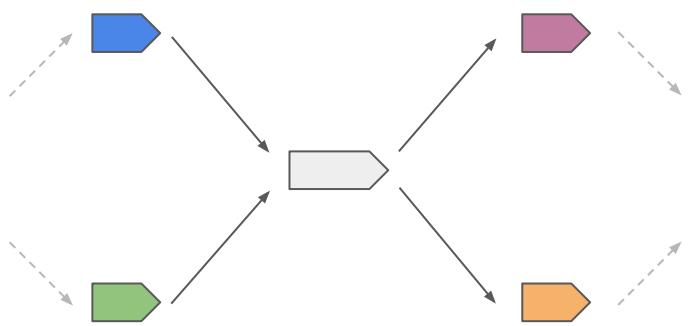
$$X\beta \approx Y$$

We can enumerate all possible paths on our assembly graph



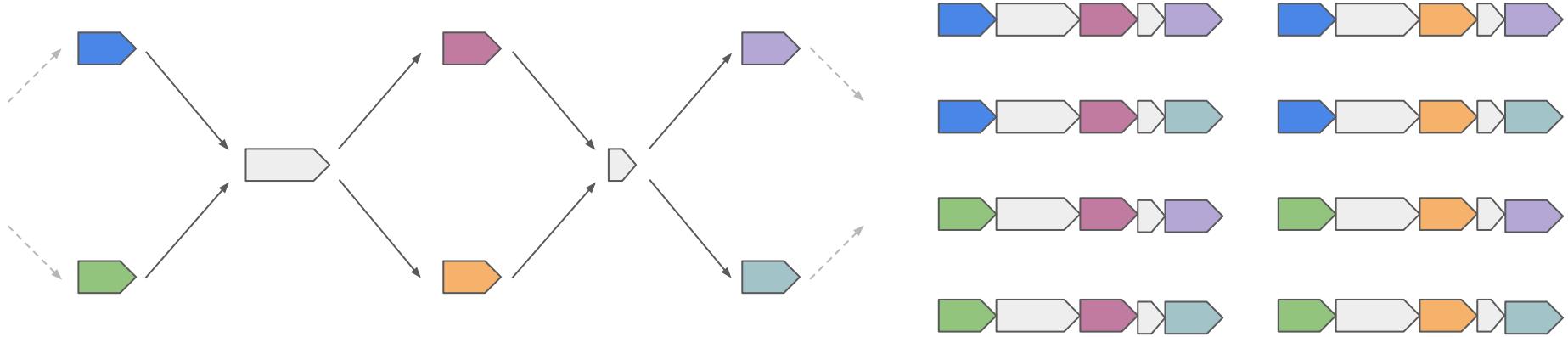
$$X\beta \approx Y$$

We can enumerate all possible paths on our assembly graph



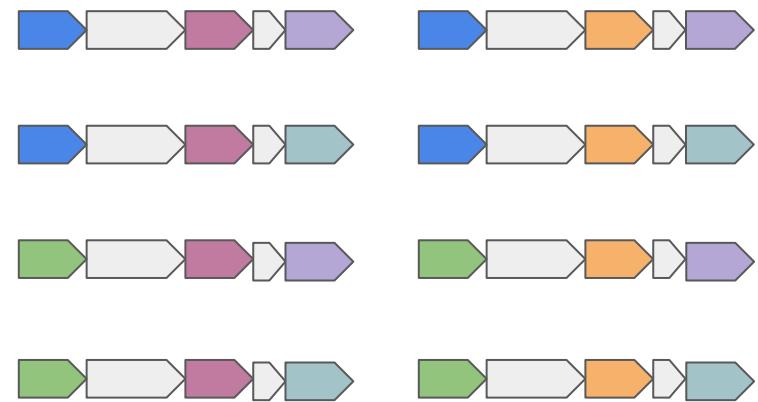
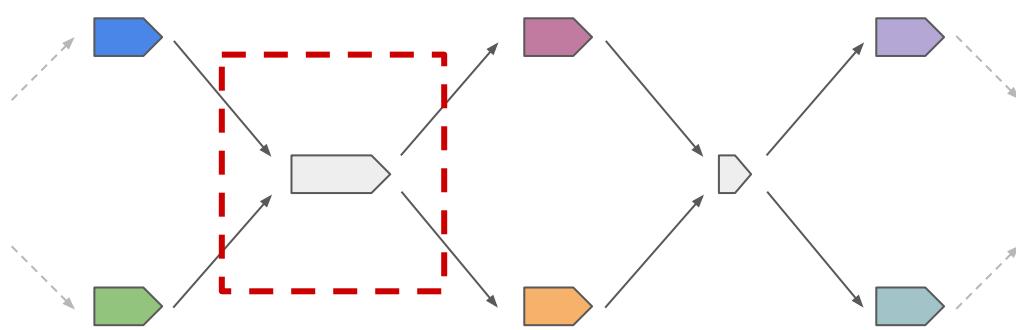
$$\boxed{X}^\beta \approx Y$$

We can enumerate all possible paths on our assembly graph

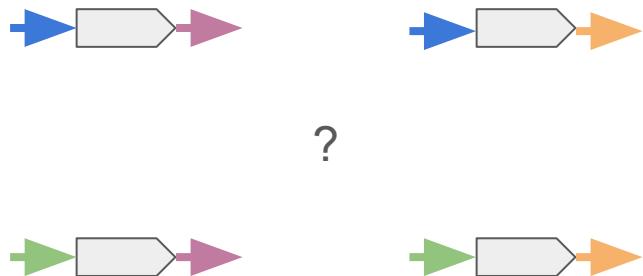
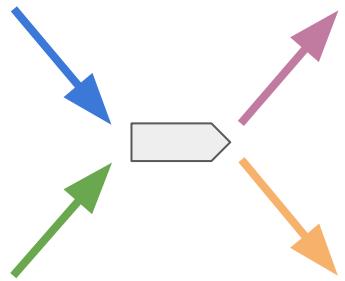


...but this grows exponentially with graph complexity

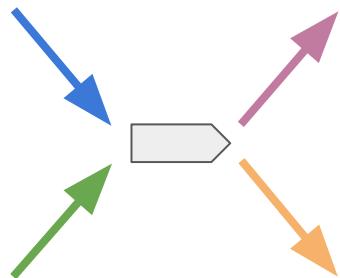
**KEY IDEA:** A single "junction" is the minimum unit of deconvolution



Divide and conquer: a single "junction" is the minimum unit of deconvolution

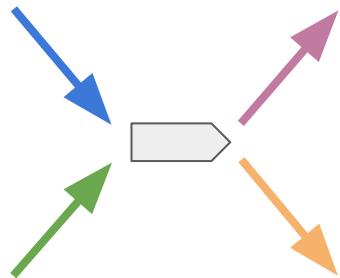


# Linear model of path depths



$$\begin{array}{c} \text{Legend:} \\ \text{Blue arrow: } \downarrow \\ \text{Green arrow: } \downarrow \\ \text{Purple arrow: } \downarrow \\ \text{Orange arrow: } \downarrow \end{array} \quad \begin{matrix} & \downarrow & \downarrow & \downarrow & \downarrow \\ & \text{p}_1 & \text{p}_2 & \text{p}_3 & \text{p}_4 \\ \text{X} = & \begin{array}{cccc} 1 & 1 & 0 & 0 \\ 0 & 0 & 1 & 1 \\ 1 & 0 & 1 & 0 \\ 0 & 1 & 0 & 1 \end{array} & \times & \begin{array}{c} \text{p}_1 \\ \text{p}_2 \\ \text{p}_3 \\ \text{p}_4 \end{array} \\ & \downarrow & \downarrow & \downarrow & \downarrow \\ & \text{e}_1 & \text{e}_2 & \text{e}_3 & \text{e}_4 \end{matrix} \quad \approx \quad \text{Y}$$

But not all paths exist: picking active paths is model selection

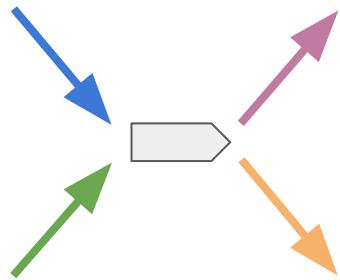


$$\begin{array}{c} \text{---} \\ \text{---} \\ \text{---} \\ \text{---} \end{array} \xrightarrow{\quad} \begin{array}{cccc} 1 & 1 & 0 & 0 \\ 0 & 0 & 1 & 1 \\ 1 & 0 & 1 & 0 \\ 0 & 1 & 0 & 1 \end{array} \times \begin{array}{c} p_1 \\ p_2 \\ p_3 \\ p_4 \end{array} \approx \begin{array}{c} e_1 \\ e_2 \\ e_3 \\ e_4 \end{array}$$

The matrix  $X$  is a 4x4 binary matrix. The vector  $\beta$  is a 4x1 column vector of probabilities  $p_1, p_2, p_3, p_4$ . The vector  $Y$  is a 4x1 column vector of error terms  $e_1, e_2, e_3, e_4$ .

$$\hat{\beta} = \operatorname{argmin} L(X\beta | Y)$$

But not all paths exist: picking active paths is model selection



$$\begin{array}{c} \text{Legend:} \\ \text{Blue arrow: } 1 \quad \text{Green arrow: } 0 \\ \text{Purple arrow: } 1 \quad \text{Orange arrow: } 0 \end{array}$$

1	1	0	0
0	0	1	1
1	0	1	0
0	1	0	1

$$X \times \beta \approx Y$$

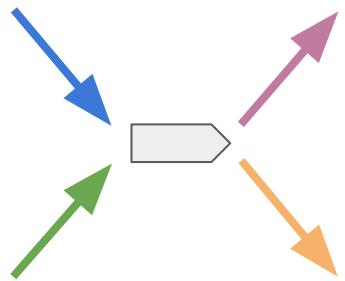
p <sub>1</sub>	0
p <sub>2</sub>	0
p <sub>3</sub>	0
p <sub>4</sub>	0

e <sub>1</sub>
e <sub>2</sub>
e <sub>3</sub>
e <sub>4</sub>

The matrix  $X$  has 4 rows and 4 columns. The first row has entries 1, 1, 0, 0. The second row has entries 0, 0, 1, 1. The third row has entries 1, 0, 1, 0. The fourth row has entries 0, 1, 0, 1. To the right of  $X$  is a multiplication sign  $\times$ . To the right of  $\beta$  is an approximation symbol  $\approx$ . To the right of  $Y$  is an equals sign  $=$ . Below  $X$  is the label  $X$ . Below  $\beta$  is the label  $\beta$ . Below  $Y$  is the label  $Y$ .

$$\hat{\beta} = \operatorname{argmin} L(X\beta | Y)$$

Last trick: To increase our power to pick paths, combine multiple samples

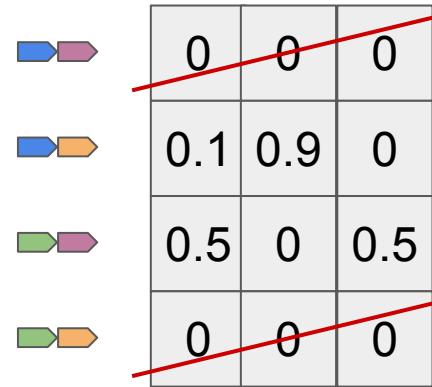
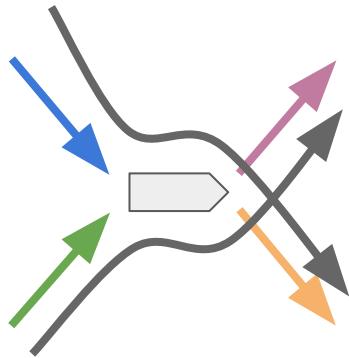


$$\begin{array}{c} \text{Samples} \\ \downarrow \\ \begin{matrix} & \textcolor{blue}{\downarrow} & \textcolor{blue}{\downarrow} & \textcolor{green}{\downarrow} & \textcolor{green}{\downarrow} \\ \textcolor{blue}{\rightarrow} & \begin{matrix} 1 & 1 & 0 & 0 \\ 0 & 0 & 1 & 1 \\ 1 & 0 & 1 & 0 \\ 0 & 1 & 0 & 1 \end{matrix} & \times & \begin{matrix} \text{Samples} \\ \downarrow \\ \begin{matrix} p_{1,1} & p_{1,2} & p_{1,3} \\ p_{2,1} & p_{2,2} & p_{2,3} \\ p_{3,1} & p_{3,2} & p_{3,3} \\ p_{4,1} & p_{4,2} & p_{4,3} \end{matrix} \xrightarrow{\quad 0 \quad} \\ \beta \end{matrix} & \approx & \begin{matrix} \text{Samples} \\ \downarrow \\ \begin{matrix} e_{1,1} & e_{1,2} & e_{1,3} \\ e_{2,1} & e_{2,2} & e_{2,3} \\ e_{3,1} & e_{3,2} & e_{3,3} \\ e_{4,1} & e_{4,2} & e_{4,3} \end{matrix} \xrightarrow{\quad 0 \quad} \\ Y \end{matrix} \end{array}$$

The diagram illustrates a matrix multiplication process. On the left, a matrix  $X$  is multiplied by a matrix  $\beta$ . The result is approximately equal to a matrix  $Y$ . Above the matrices, a bracket labeled "Samples" indicates that each row of  $X$  and  $\beta$  represents a sample. Red arrows point from the bottom-right corner of both  $\beta$  and  $Y$  to the value "0", indicating that the product of the last row of  $X$  and the last column of  $\beta$  is zero.

$$\hat{\beta} = \operatorname{argmin} L(X\beta | Y)$$

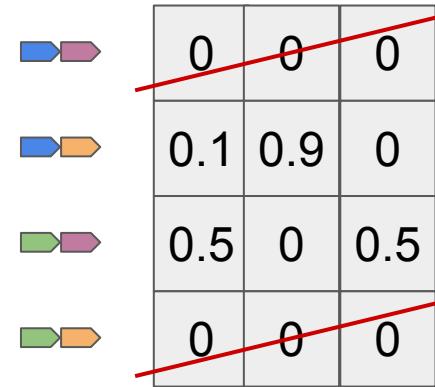
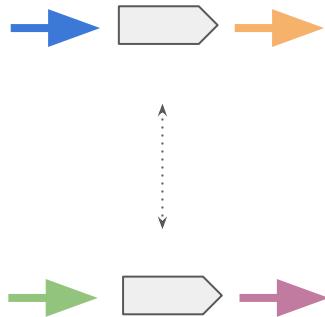
# Drop paths with no depth in any sample



$$\hat{\beta}$$

Used statistical linkage to resolve ambiguity about  
which of possible paths are "real"

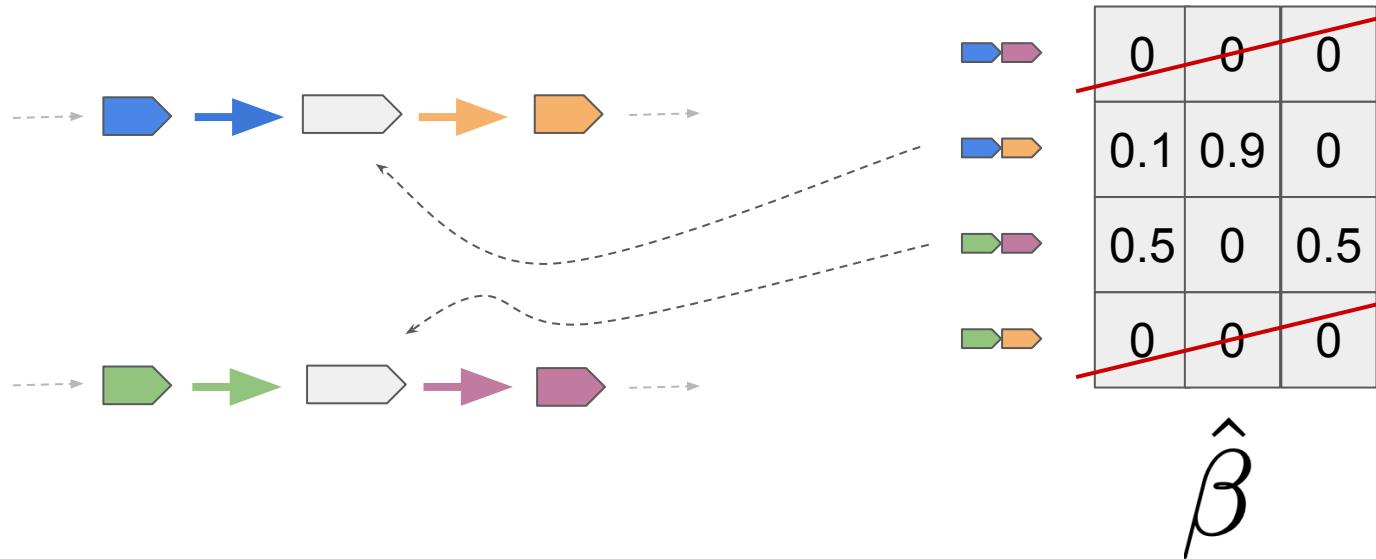
# Resolve ambiguity, longer linear sequences



$\hat{\beta}$

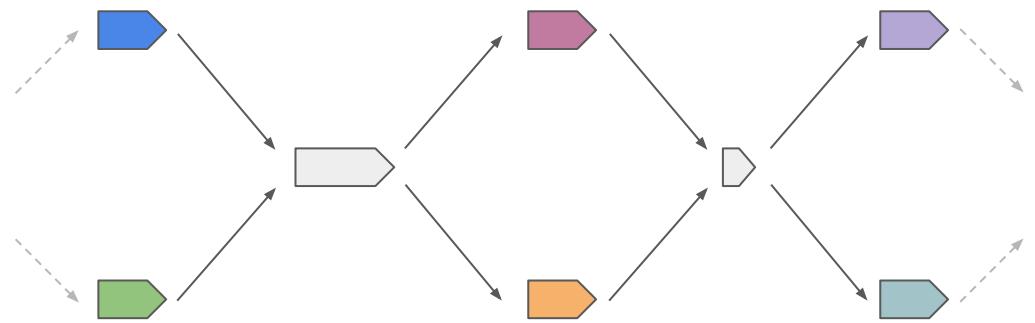
Can "unzip" this unitig into two paths

# Resolve ambiguity, longer linear sequences

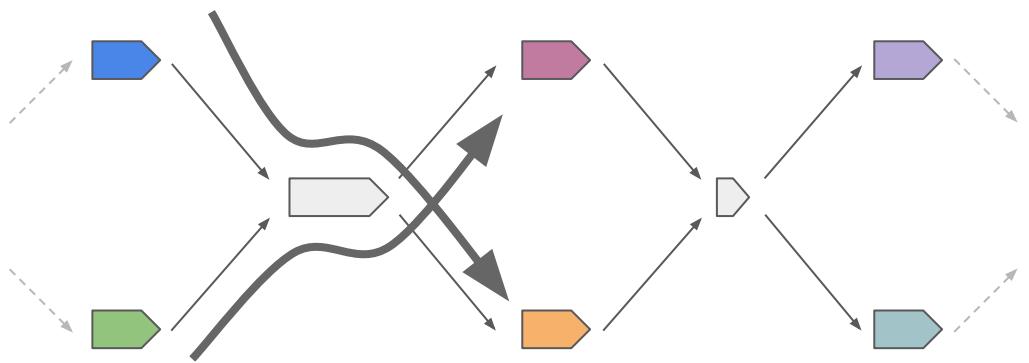


Newly split unitigs already have depths estimated across samples

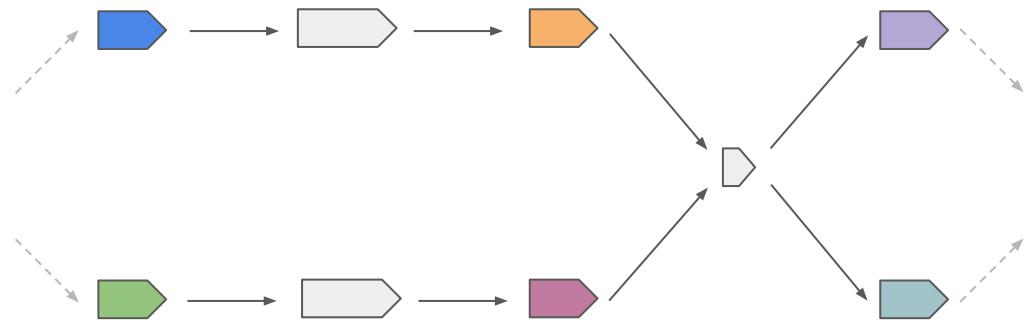
# Iteratively unzipping local junctions



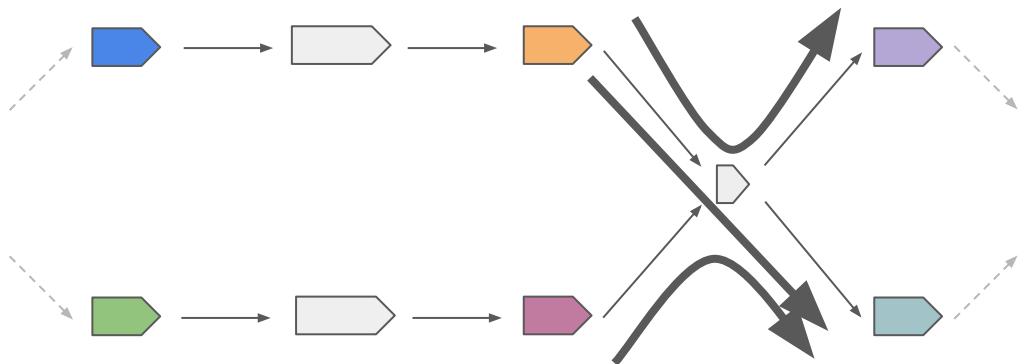
# Iteratively unzipping local junctions



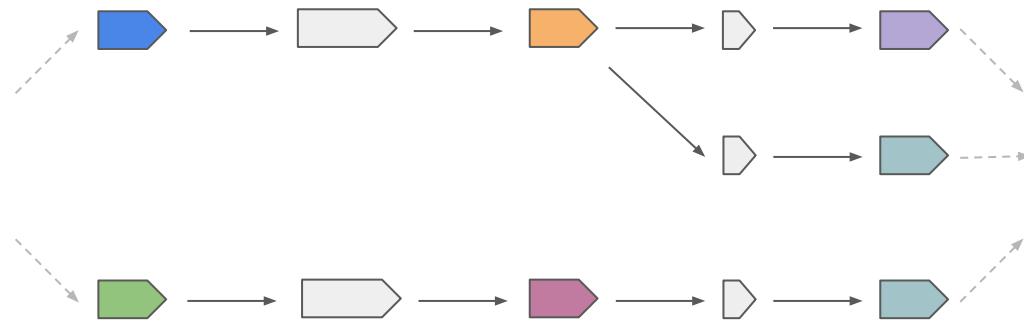
# Iteratively unzipping local junctions



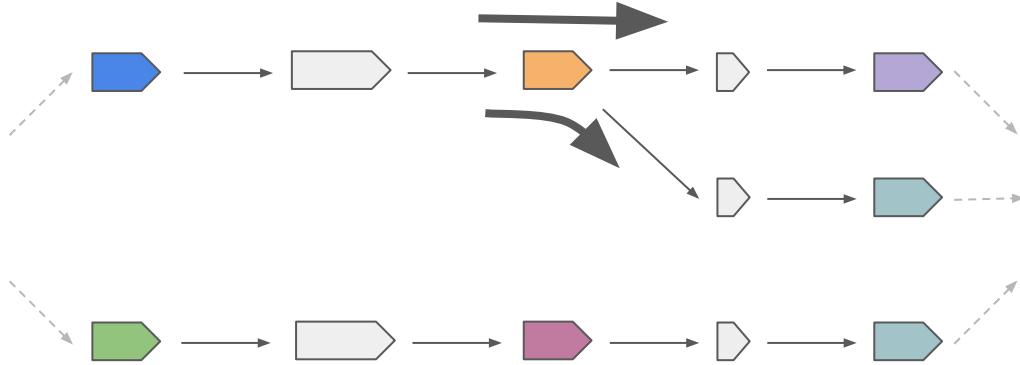
# Iteratively unzipping local junctions



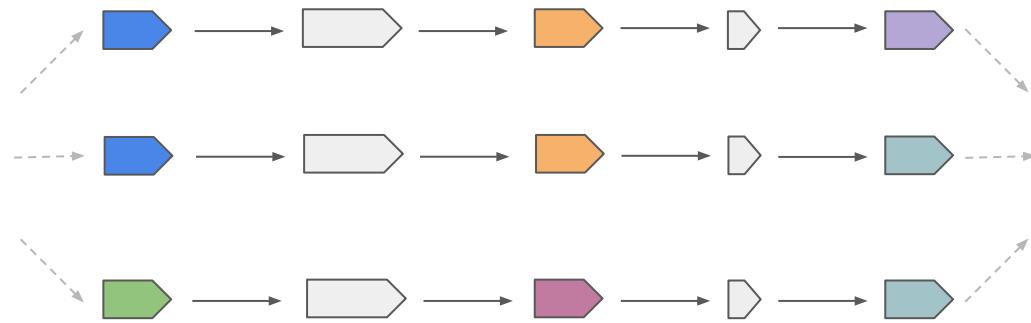
# Iteratively unzipping local junctions



# Iteratively unzipping local junctions



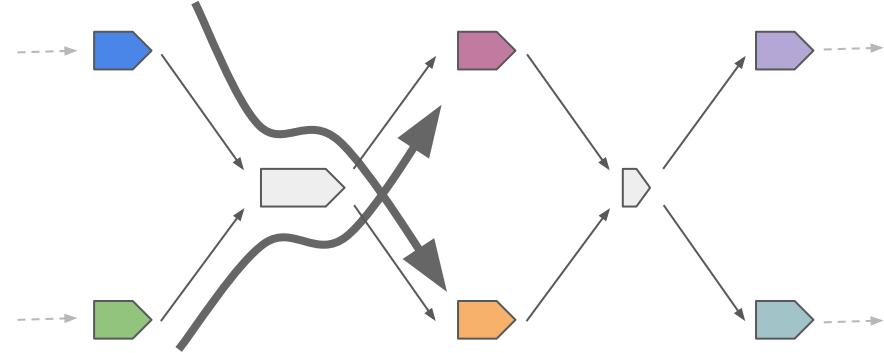
# Iteratively unzipping local junctions



# StrainZip

Assembly Graph Deconvolution for  
Quantification of Strain-Specific  
Sequences across Metagenomes

<https://github.com/bsmith89/StrainZip>



$$\begin{array}{c} \textcolor{blue}{\downarrow} \quad \textcolor{green}{\downarrow} \\ \begin{matrix} & 1 & 1 & 0 & 0 \\ 0 & 0 & 1 & 1 \\ 1 & 0 & 1 & 0 \\ 0 & 1 & 0 & 1 \end{matrix} \end{array}$$

$$\begin{array}{c} \textcolor{purple}{\downarrow} \quad \textcolor{orange}{\downarrow} \\ \begin{matrix} p_{1,1} & p_{1,2} & p_{1,3} \\ p_{2,1} & p_{2,2} & p_{2,3} \\ p_{3,1} & p_{3,2} & p_{3,3} \\ p_{4,1} & p_{4,2} & p_{4,3} \end{matrix} \end{array}$$

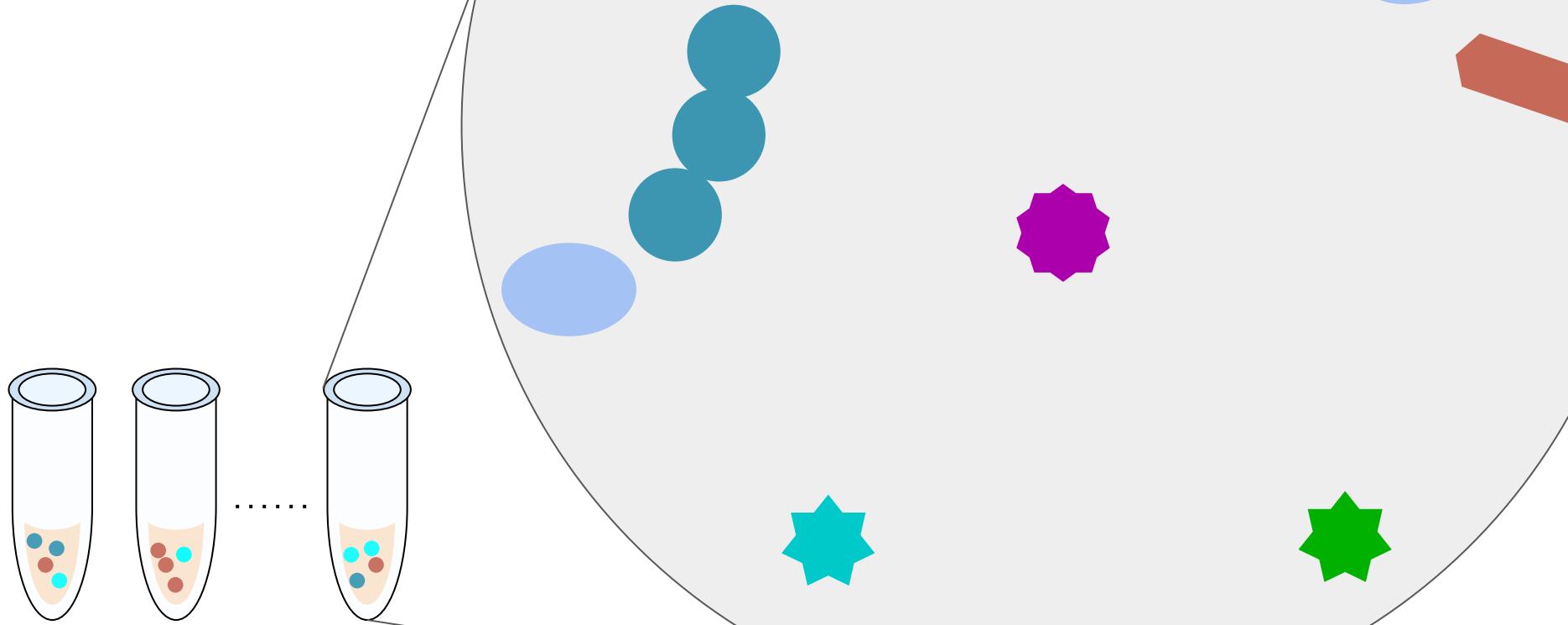
$\times$

$$\begin{array}{c} \approx \\ \begin{matrix} e_{1,1} & e_{1,2} & e_{1,3} \\ e_{2,1} & e_{2,2} & e_{2,3} \\ e_{3,1} & e_{3,2} & e_{3,3} \\ e_{4,1} & e_{4,2} & e_{4,3} \end{matrix} \end{array}$$

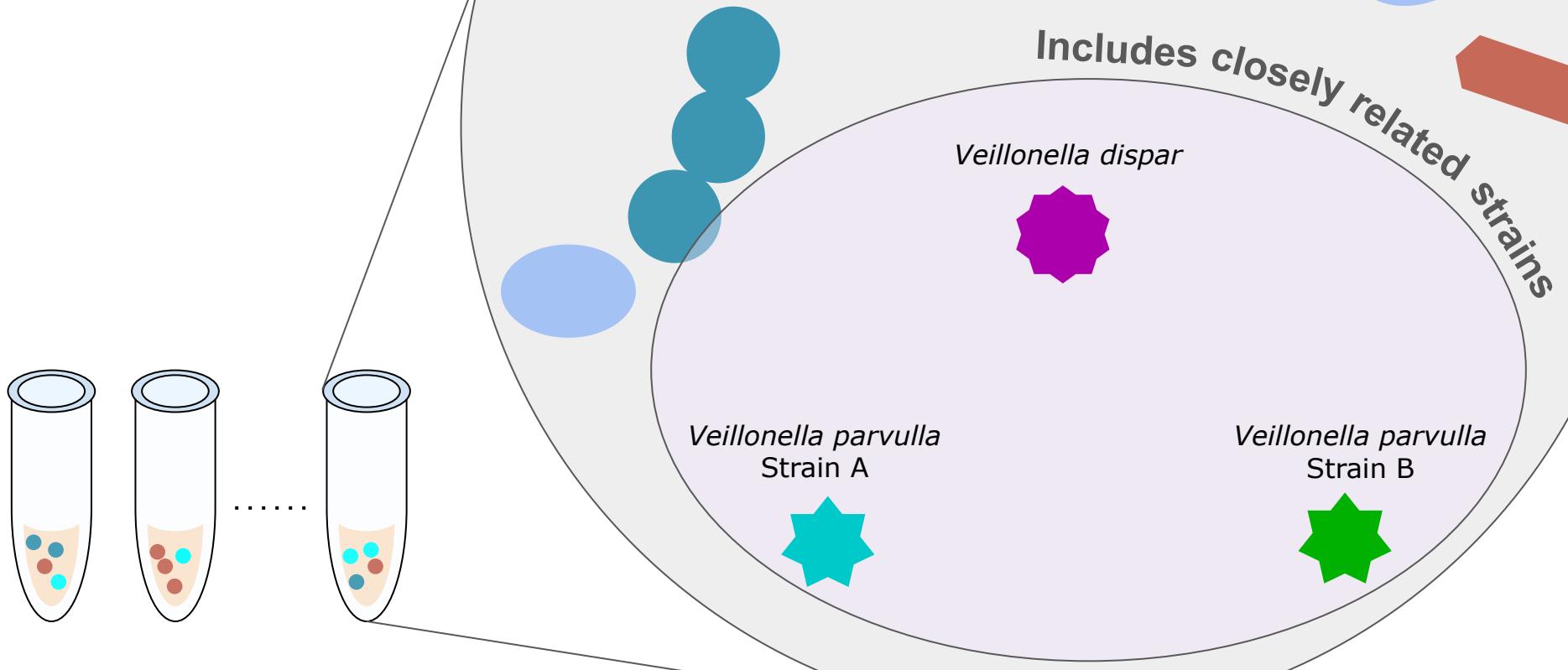
$$\begin{array}{c} \approx \\ \begin{matrix} e_{1,1} & e_{1,2} & e_{1,3} \\ e_{2,1} & e_{2,2} & e_{2,3} \\ e_{3,1} & e_{3,2} & e_{3,3} \\ e_{4,1} & e_{4,2} & e_{4,3} \end{matrix} \end{array}$$

# Demonstration

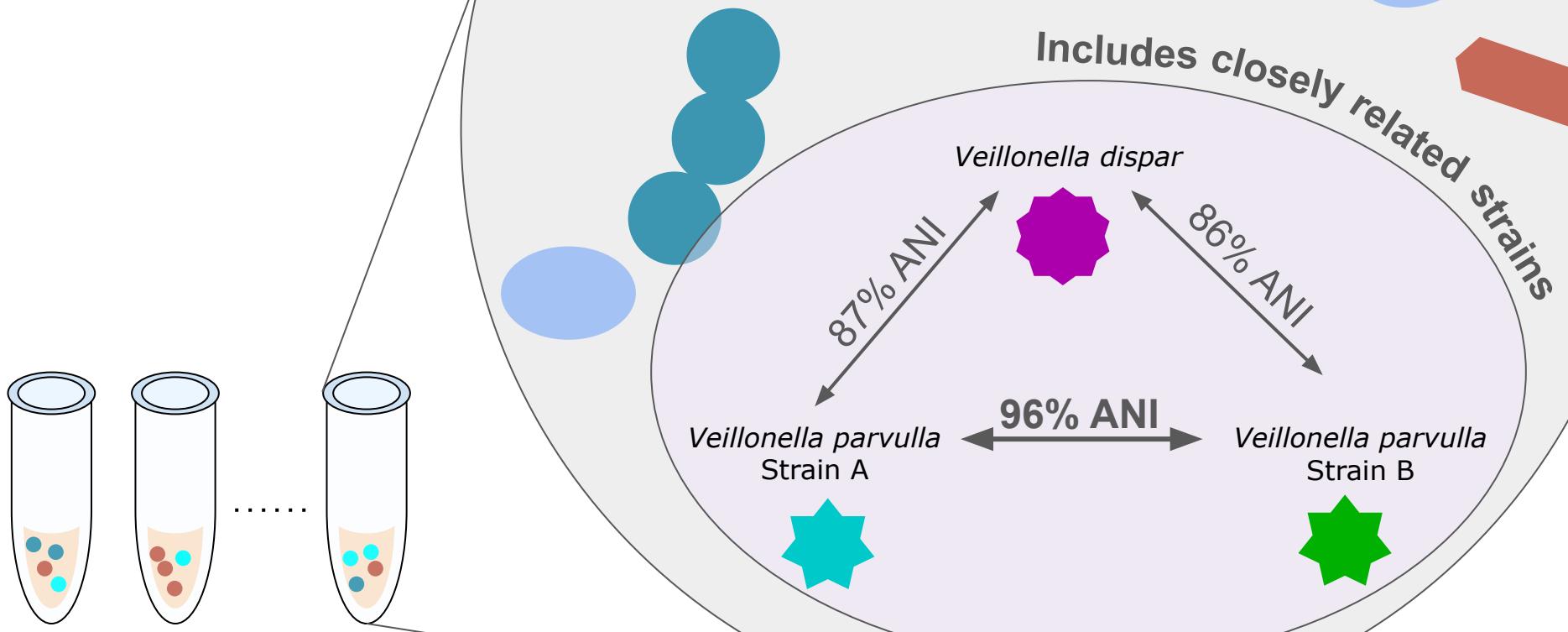
**hCOM2** is a complex (125 species), synthetic community  
with high-quality, reference genomes for all strains

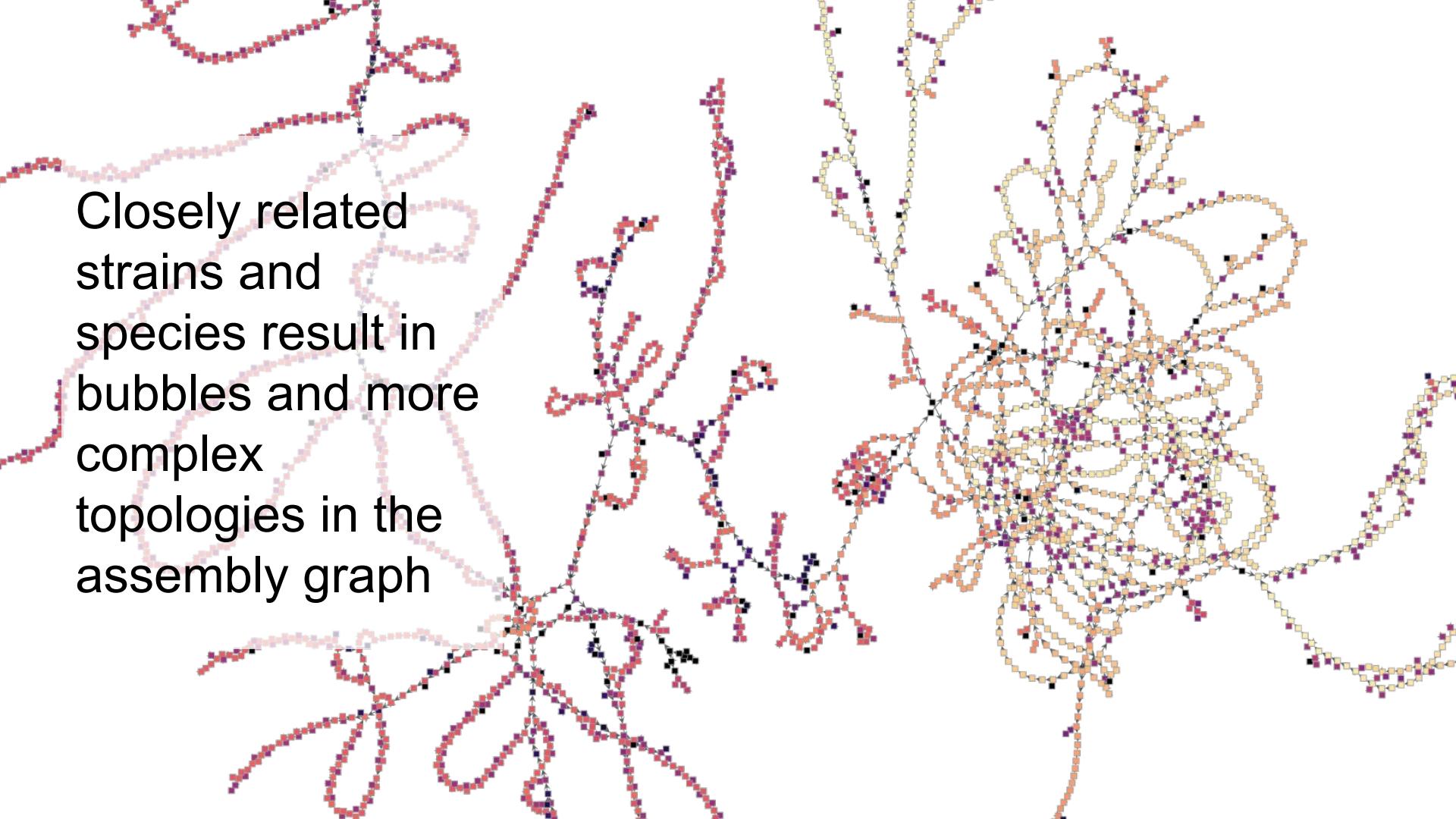


**hCOM2** is a complex (125 species), synthetic community with high-quality, reference genomes for all strains

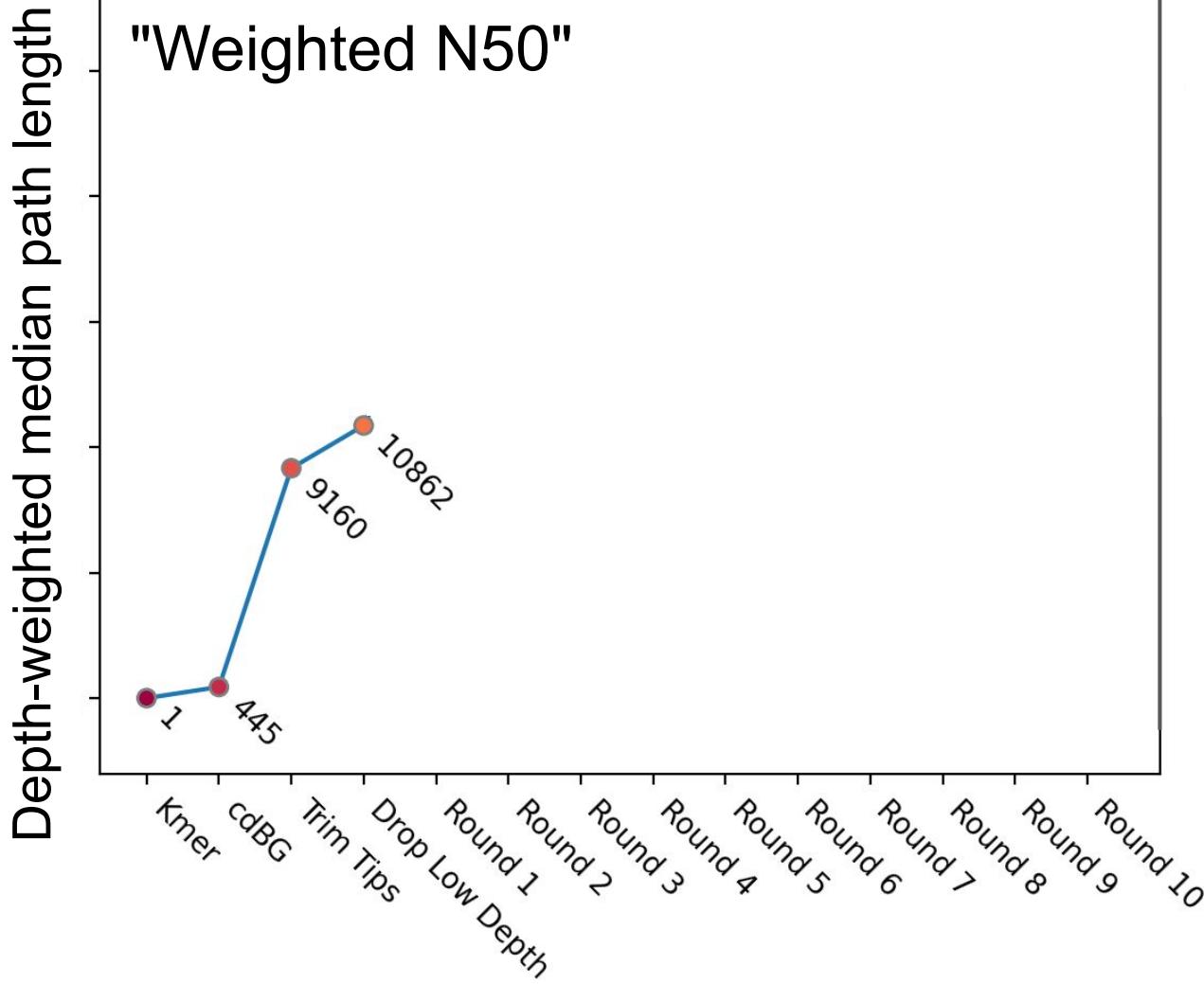


**hCOM2** is a complex (125 species), synthetic community with high-quality, reference genomes for all strains

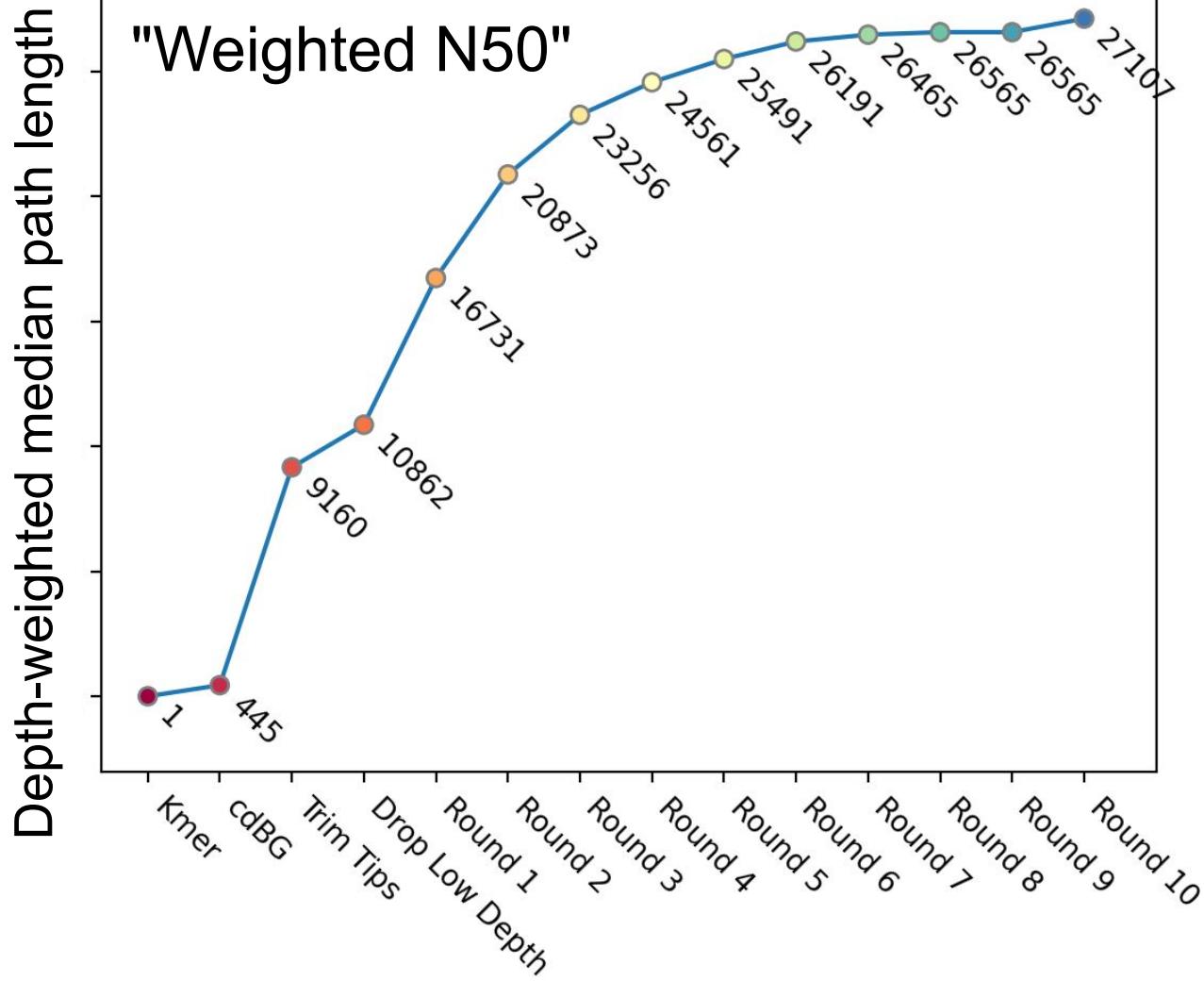




Closely related strains and species result in bubbles and more complex topologies in the assembly graph



Complex  
assembly graph  
results in short  
path lengths

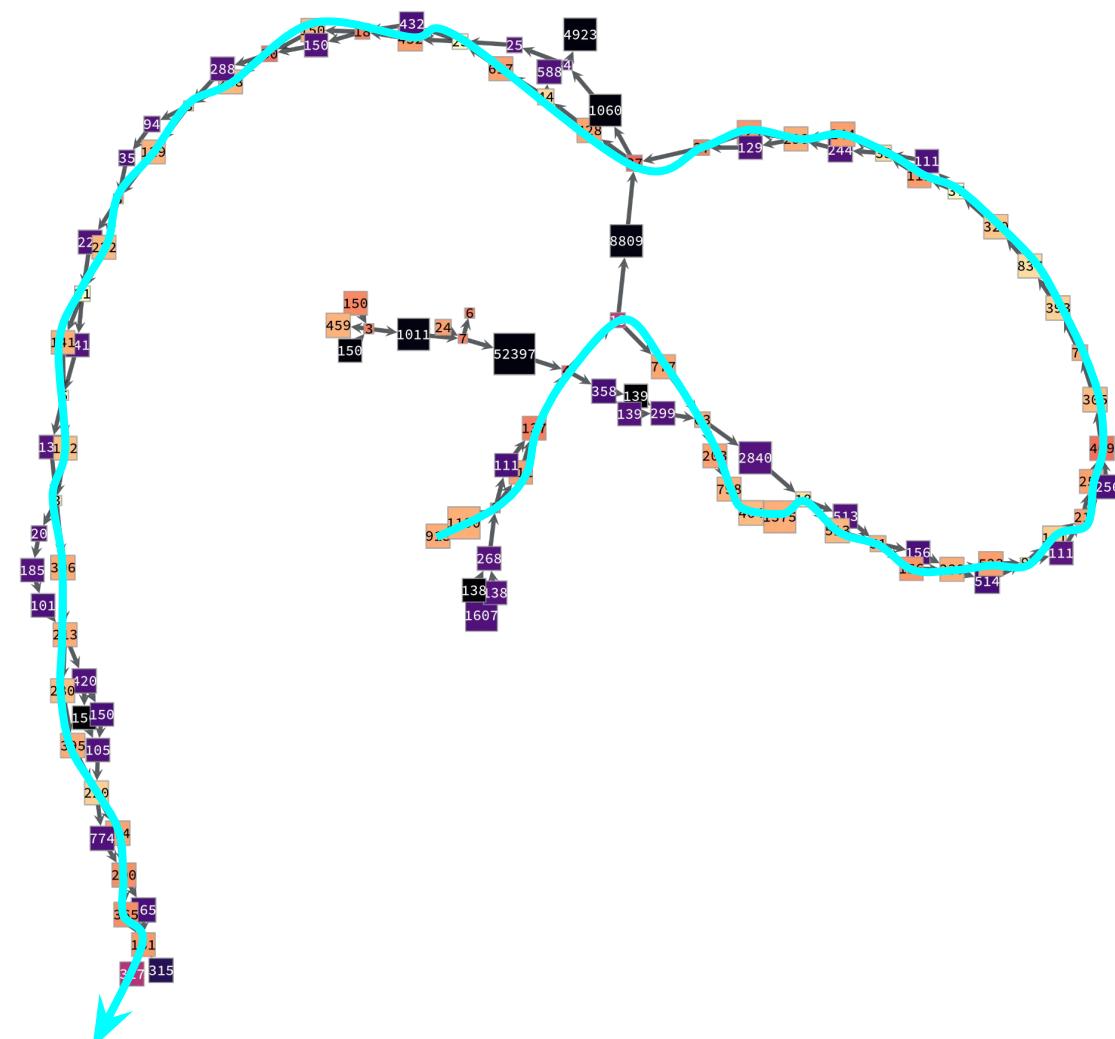


Path lengths increase over successive rounds of deconvolution

Closely related strains are interspersed in the assembly graph

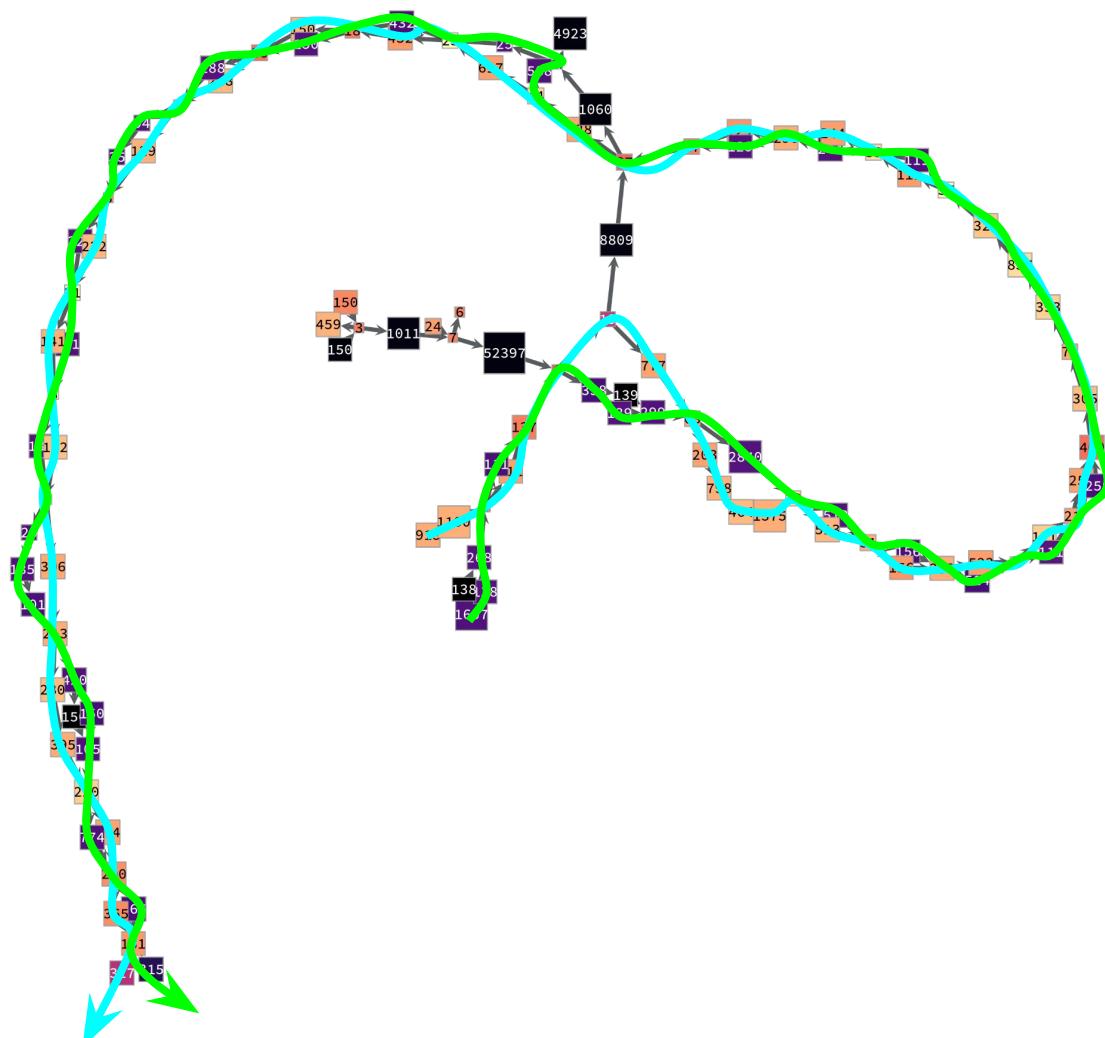


Deconvolution  
recovers longer,  
strain-specific  
sequences



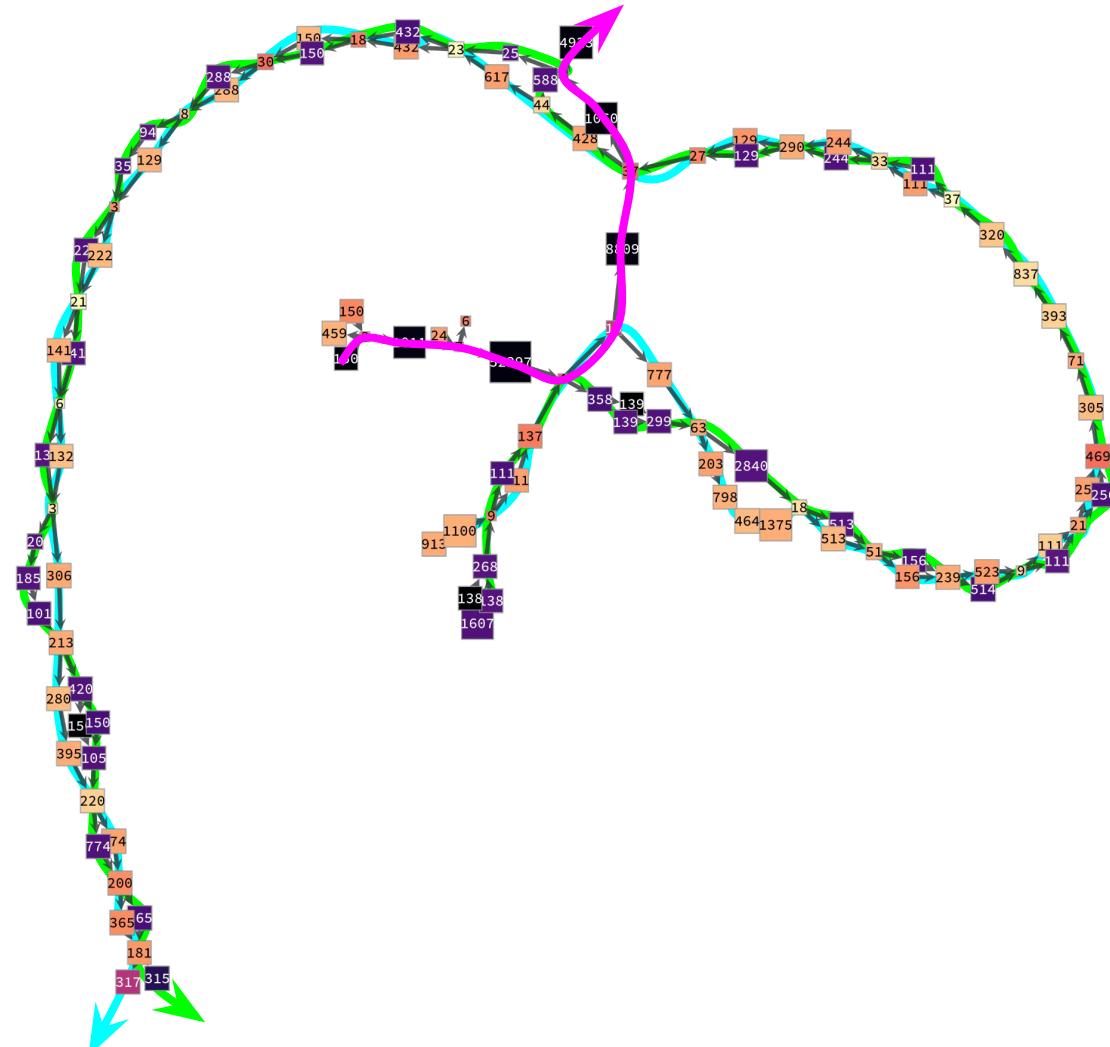
Deconvolution  
recovers longer,  
strain-specific  
sequences

...including  
lower-abundance  
strains



Deconvolution  
recovers longer,  
strain-specific  
sequences

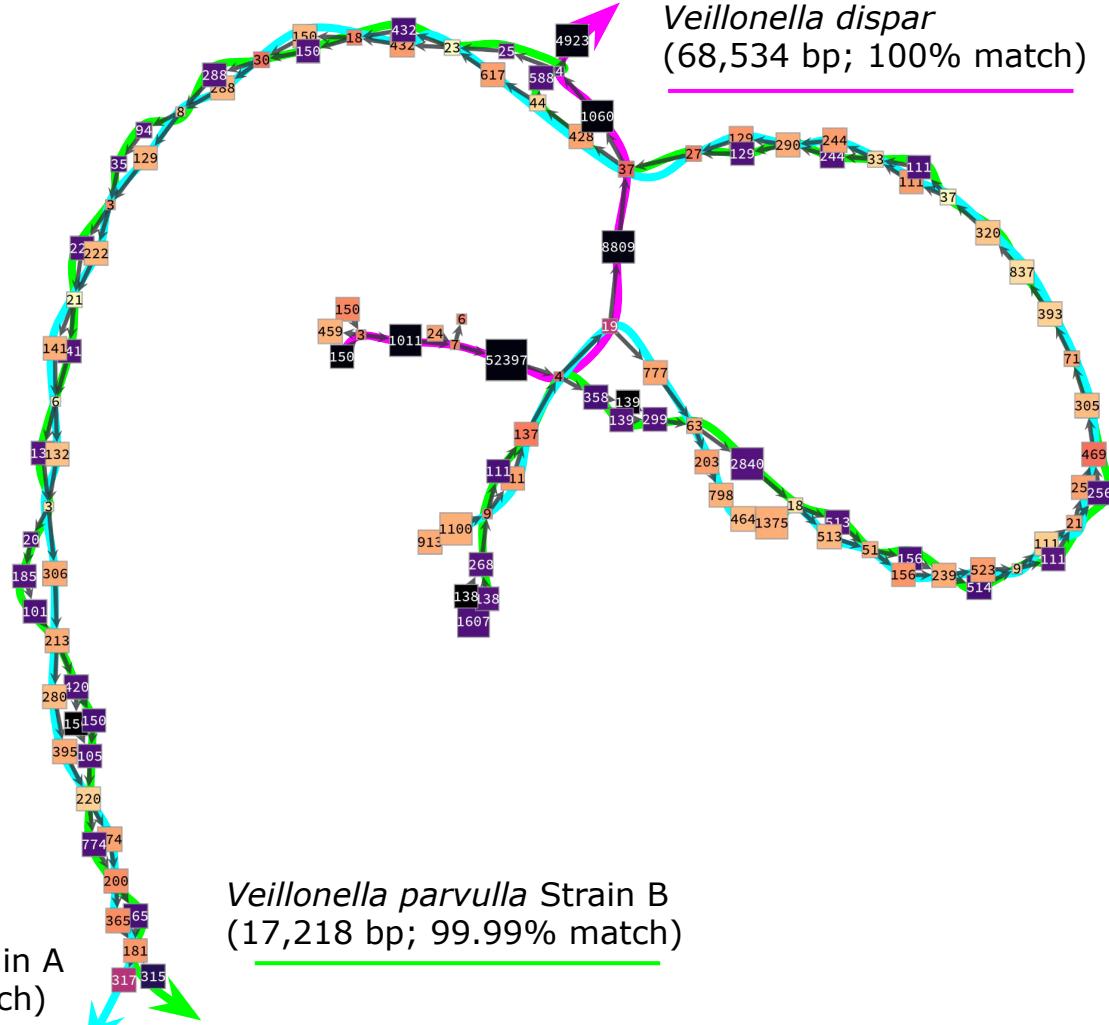
...including  
lower-abundance  
strains  
...and species

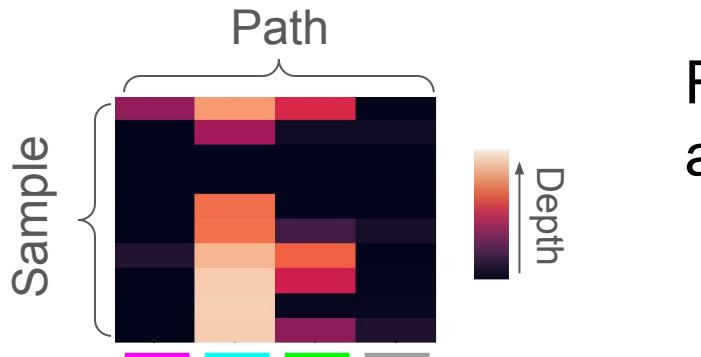


Deconvolution  
recovers longer,  
strain-specific  
sequences

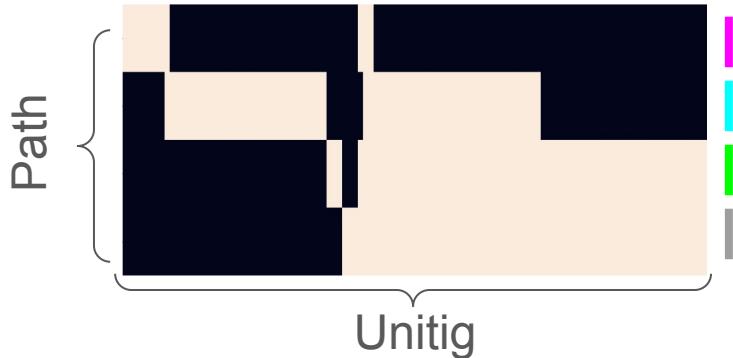
...including  
lower-abundance  
strains  
...and species  
...accurately

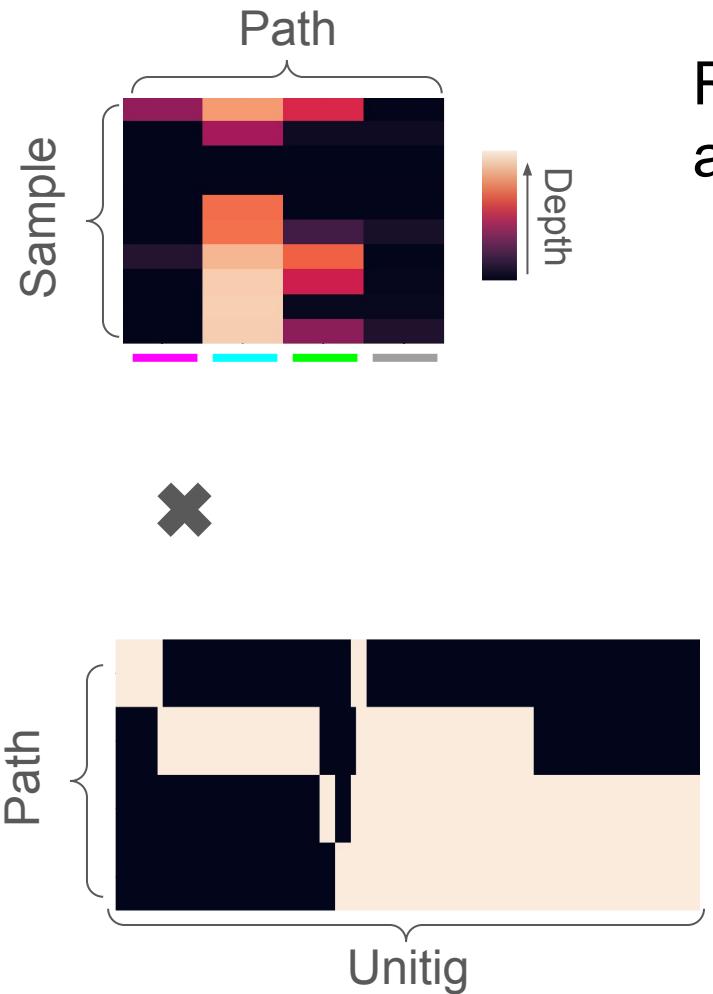
*Veillonella parvulla* Strain A  
(17,229 bp; 100% match)



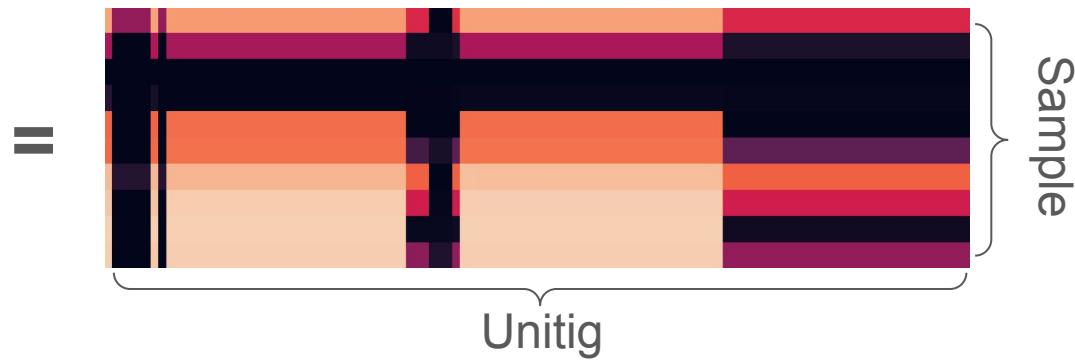


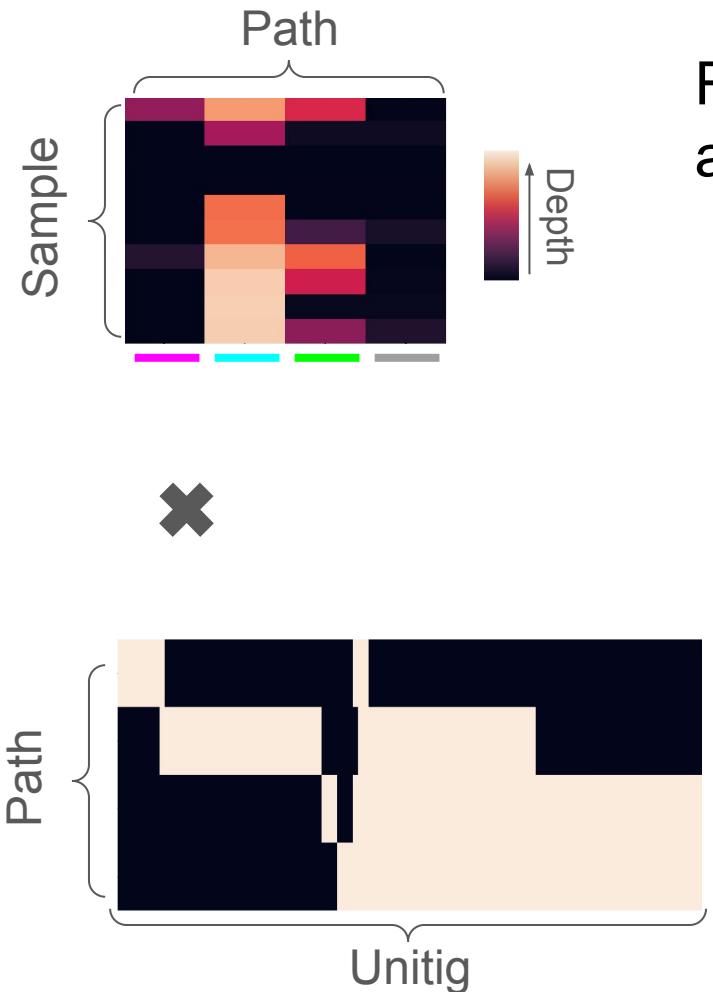
Result: both paths, and path depths across samples (without read mapping)



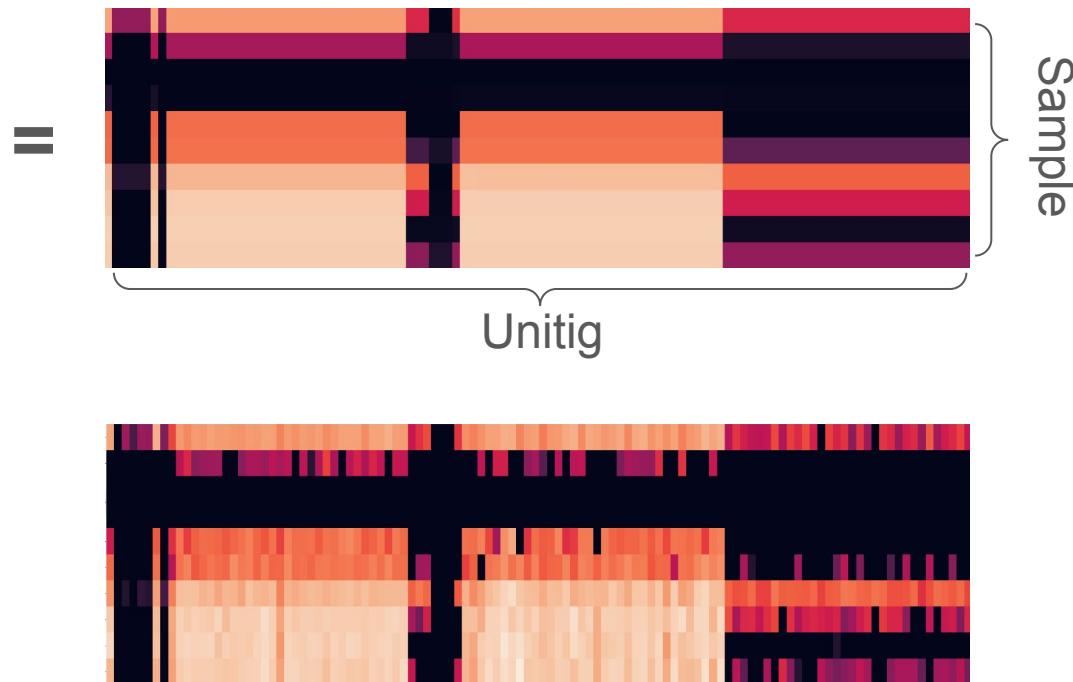


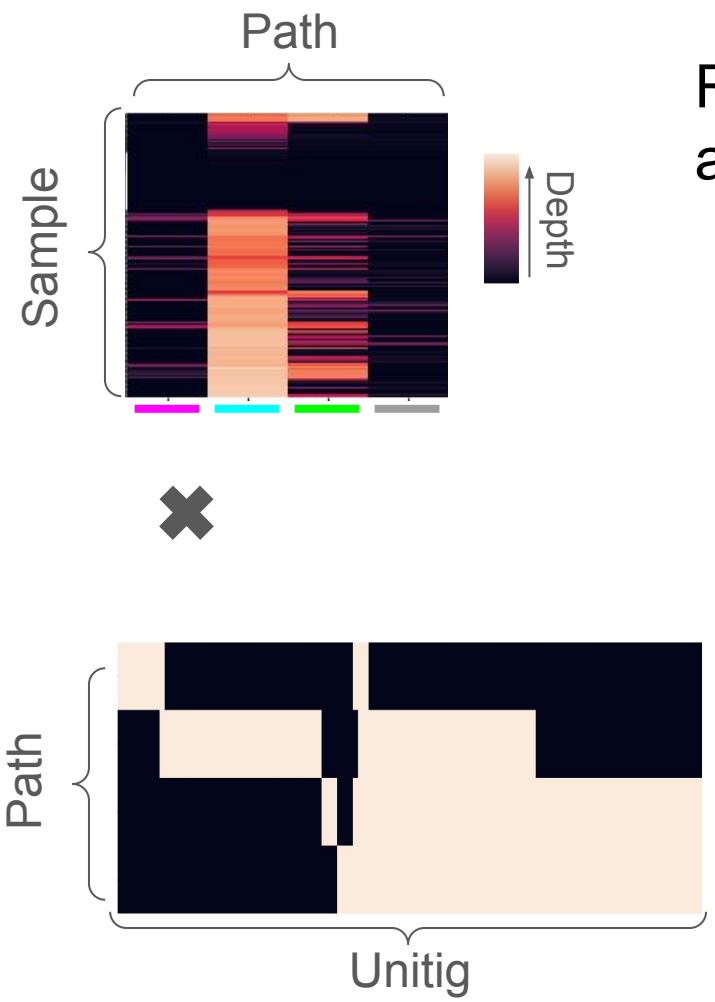
Result: both paths, and path depths  
across samples (without read mapping)



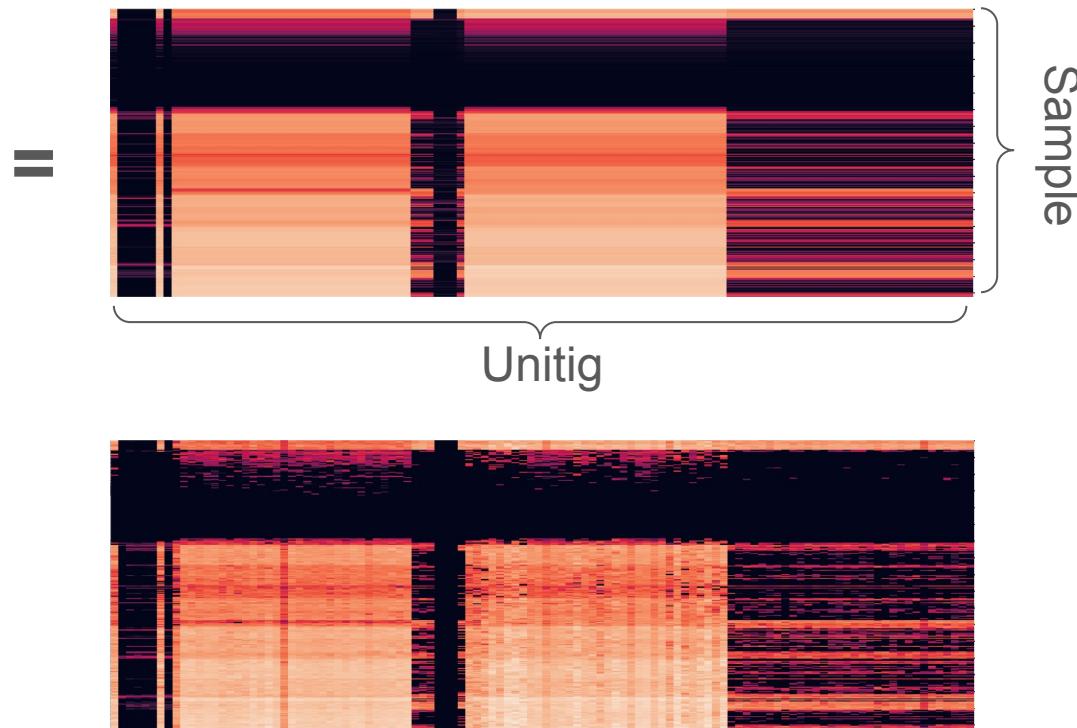


Result: both paths, and path depths  
across samples (without read mapping)



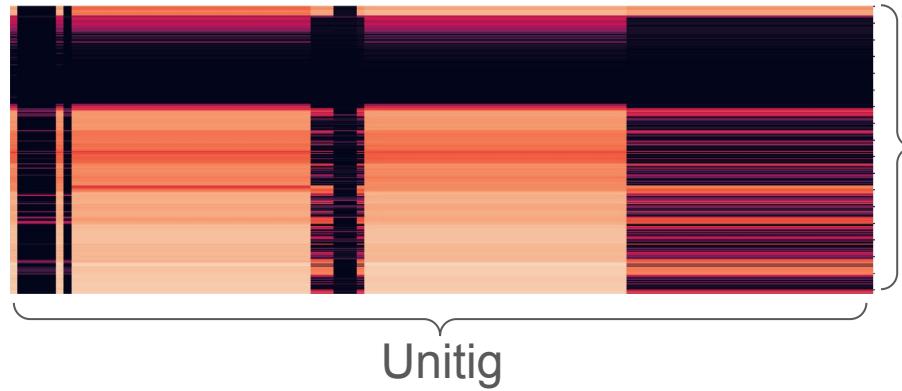


Result: both paths, and path depths  
across samples (without read mapping)

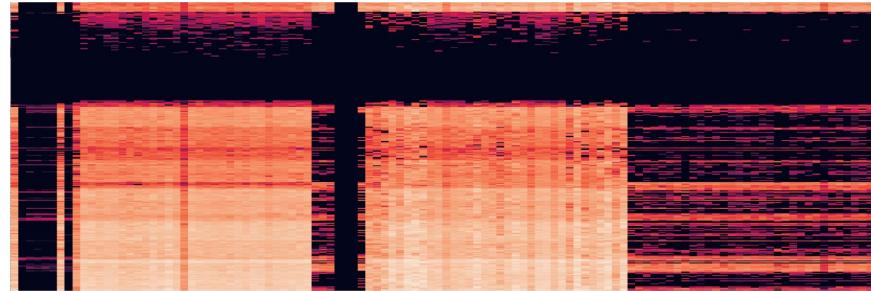


Estimated  
unitig  
depths  
closely  
match  
observed  
depths

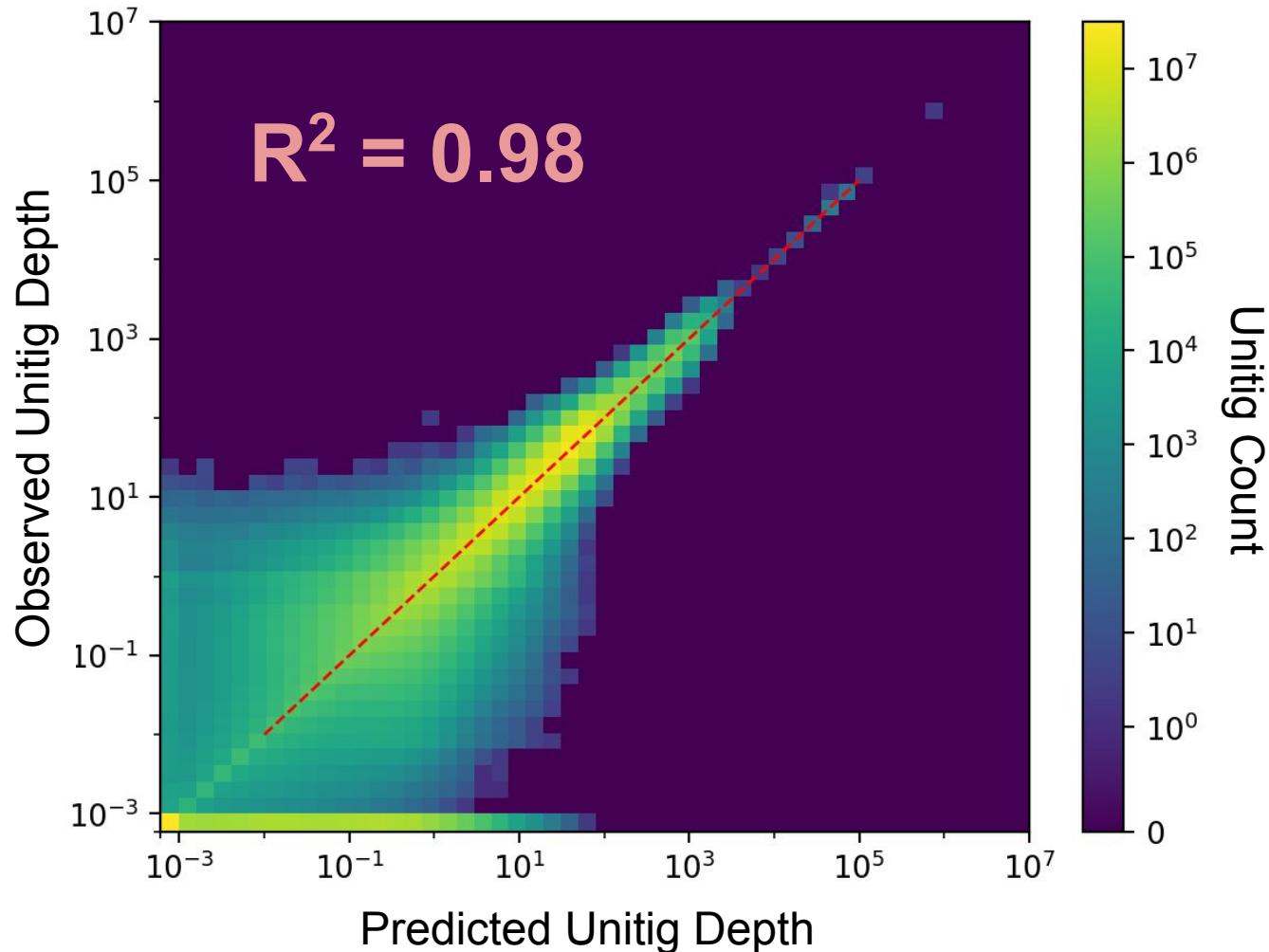
Predicted →



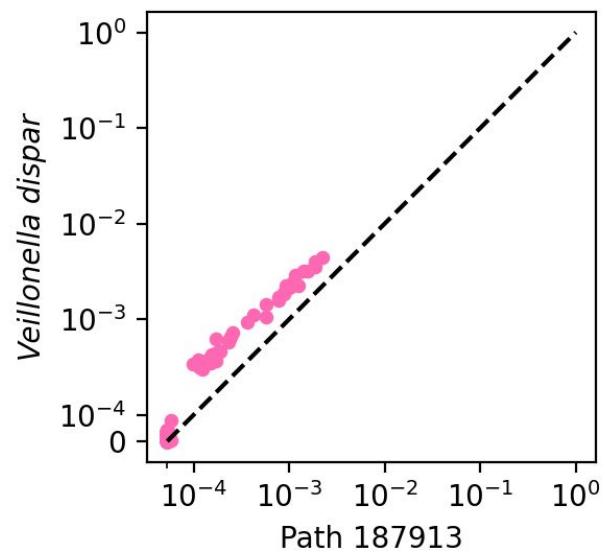
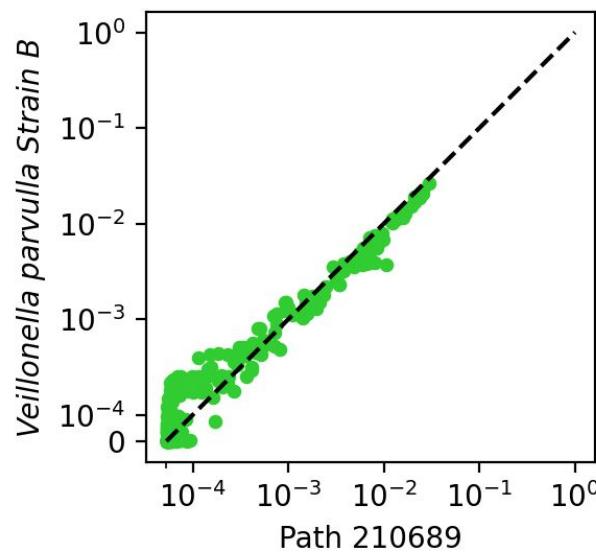
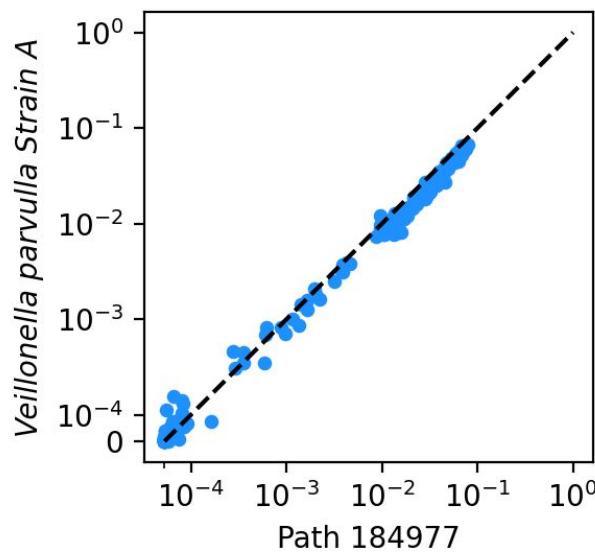
Observed →



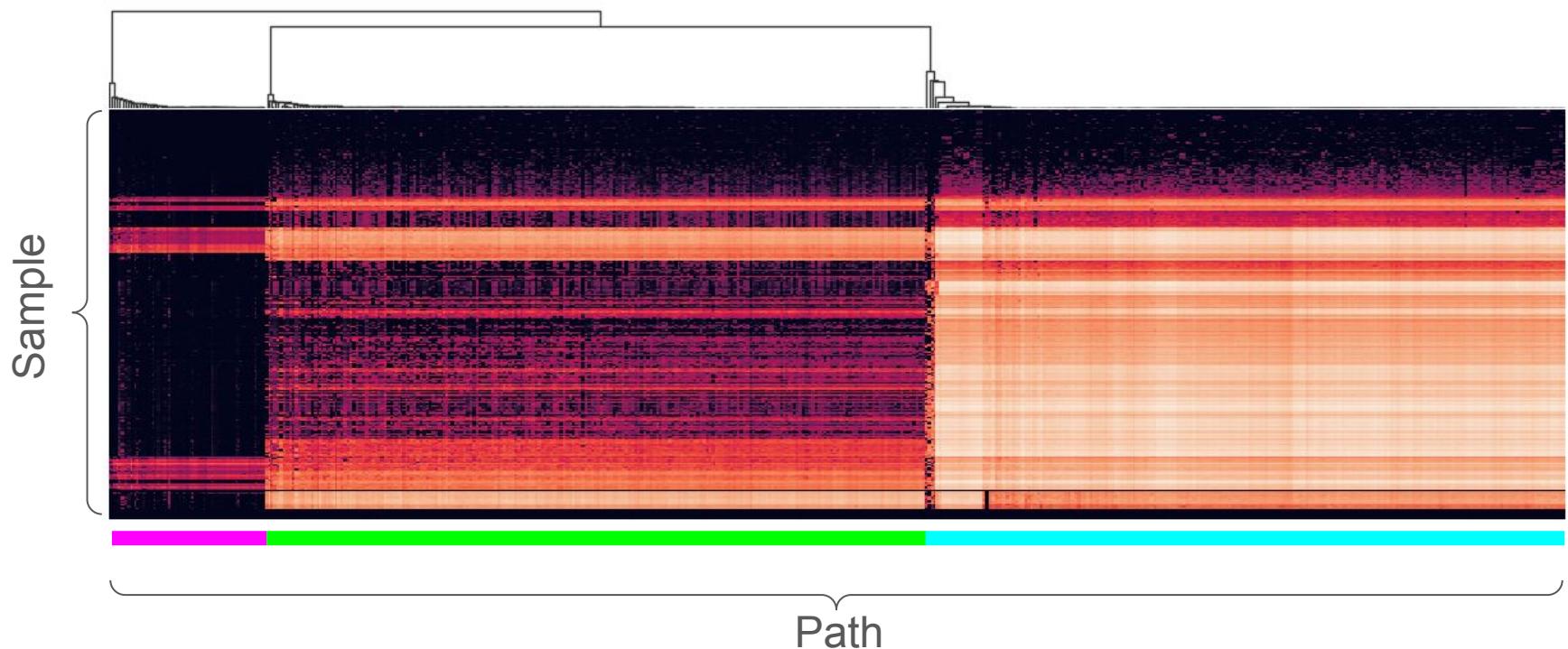
Estimated  
unitig  
depths  
closely  
match  
observed  
depths

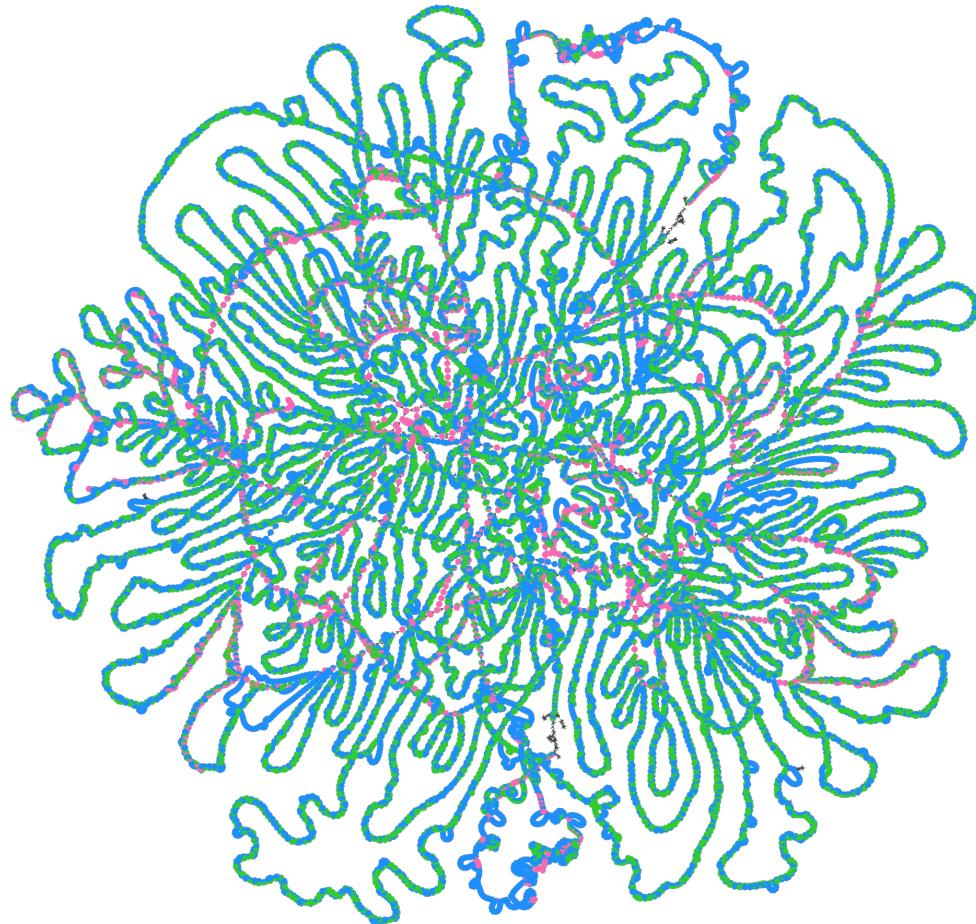


# Path depths match reference-based strain depth estimates



Clustering paths by depth combines multiple sequences from the same strain





Enables  
strain-resolved  
genome assembly  
from metagenomes





# Thank You!

