

# Floria: metagenome strain haplotyping with short/long reads

**Jim Shaw**<sup>1\*</sup>, Jean-Sébastien Gounot<sup>2\*</sup>,  
Hanrong Chen<sup>2</sup>, Niranjan Nagarajan<sup>2#</sup>, Yun William Yu<sup>1,3#</sup>

\* Equal contribution

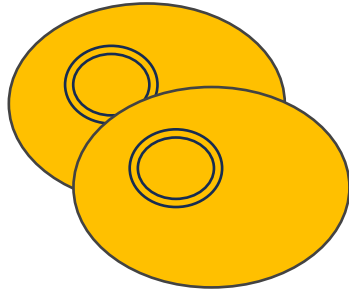
<sup>1</sup>University of Toronto

<sup>2</sup>Genome Institute of Singapore

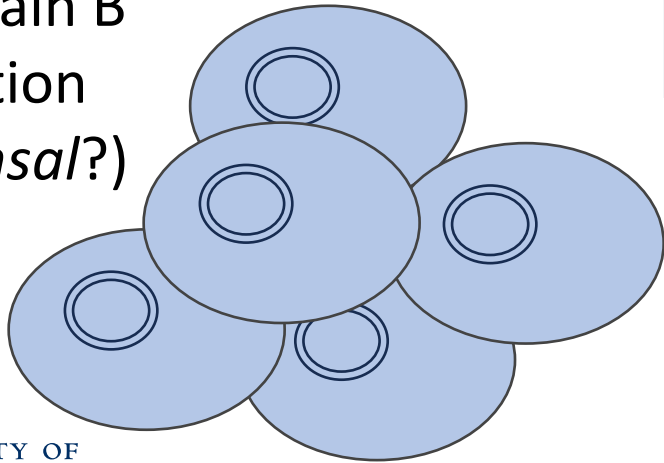
<sup>3</sup>Carnegie Mellon University

# Strain-level heterogeneity in metagenomes is important!

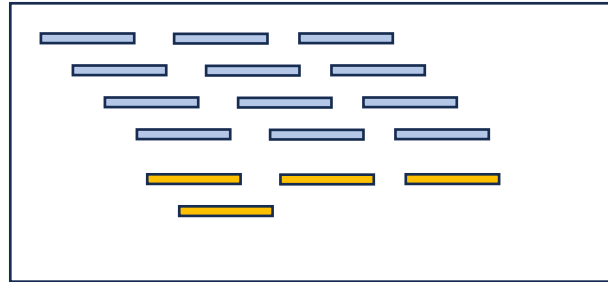
E. coli strain A  
population  
(*pathogenic?*)



E. coli strain B  
population  
(*commensal?*)



Strain-mixed  
metagenomic reads



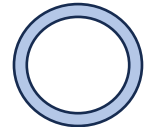
Assembly

What is assembled?

A) only



B) only



C) both

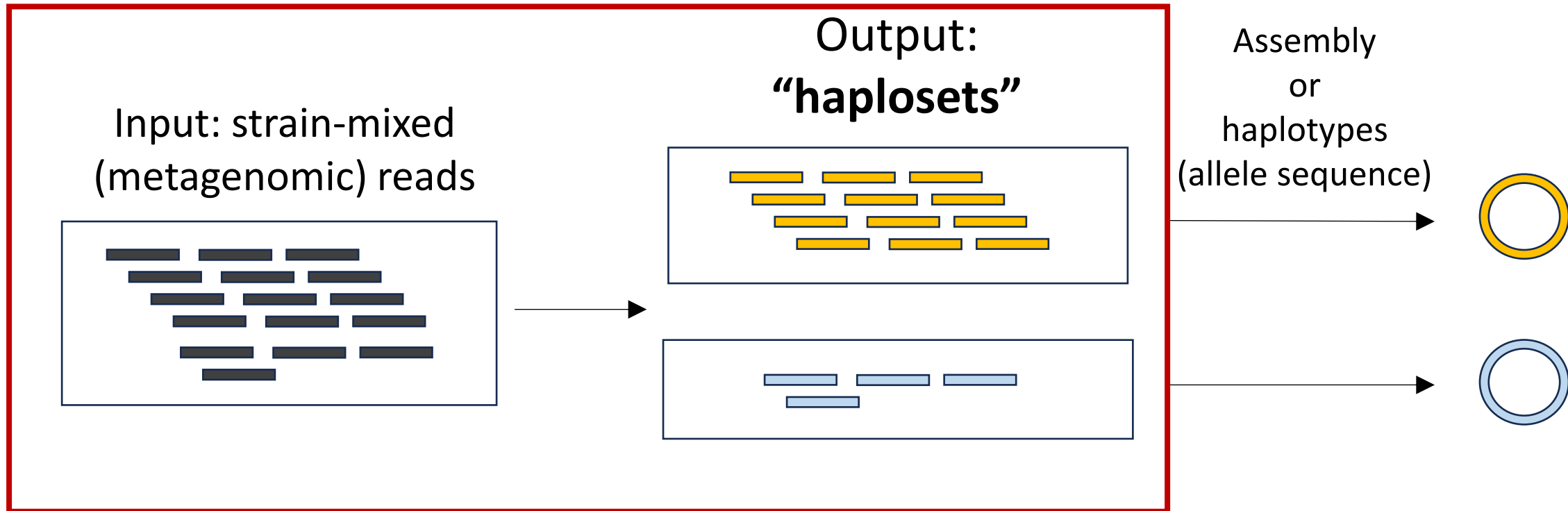


Answer: depends on technology + algorithms

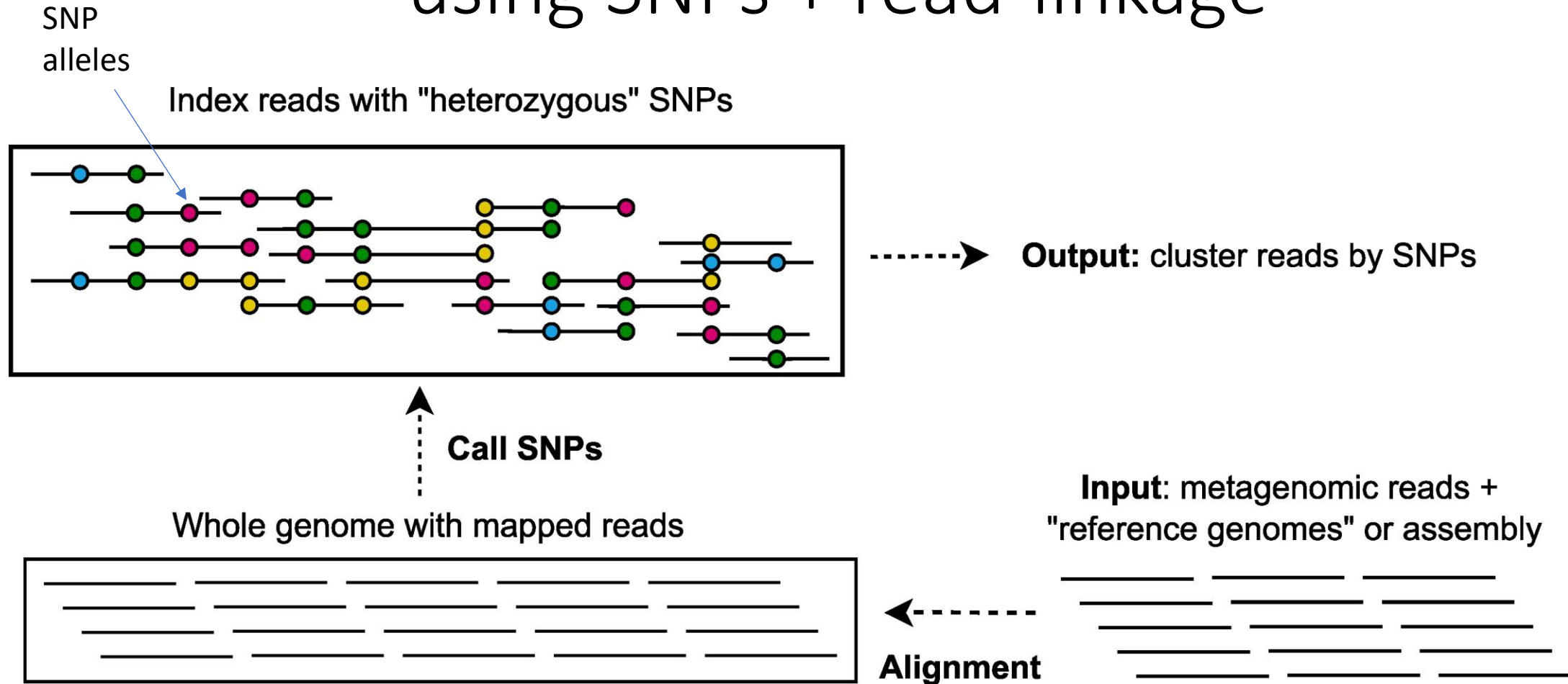
Technology	Resulting assembly
Short reads (algorithm: SPAdes)	<b>One</b> strain (lost information)
Low-fidelity long-reads (algorithm: metaFlye)	<b>One</b> strain (lost information)
High-fidelity long-reads (algorithm: hifiasm or metaMDBG)	<b>Both</b> strains

# Introduction: Haplotyping and phasing

# Computational goal: reads → strain-level “haplosets”



# Main idea: computational phasing (haplotyping) using SNPs + read-linkage



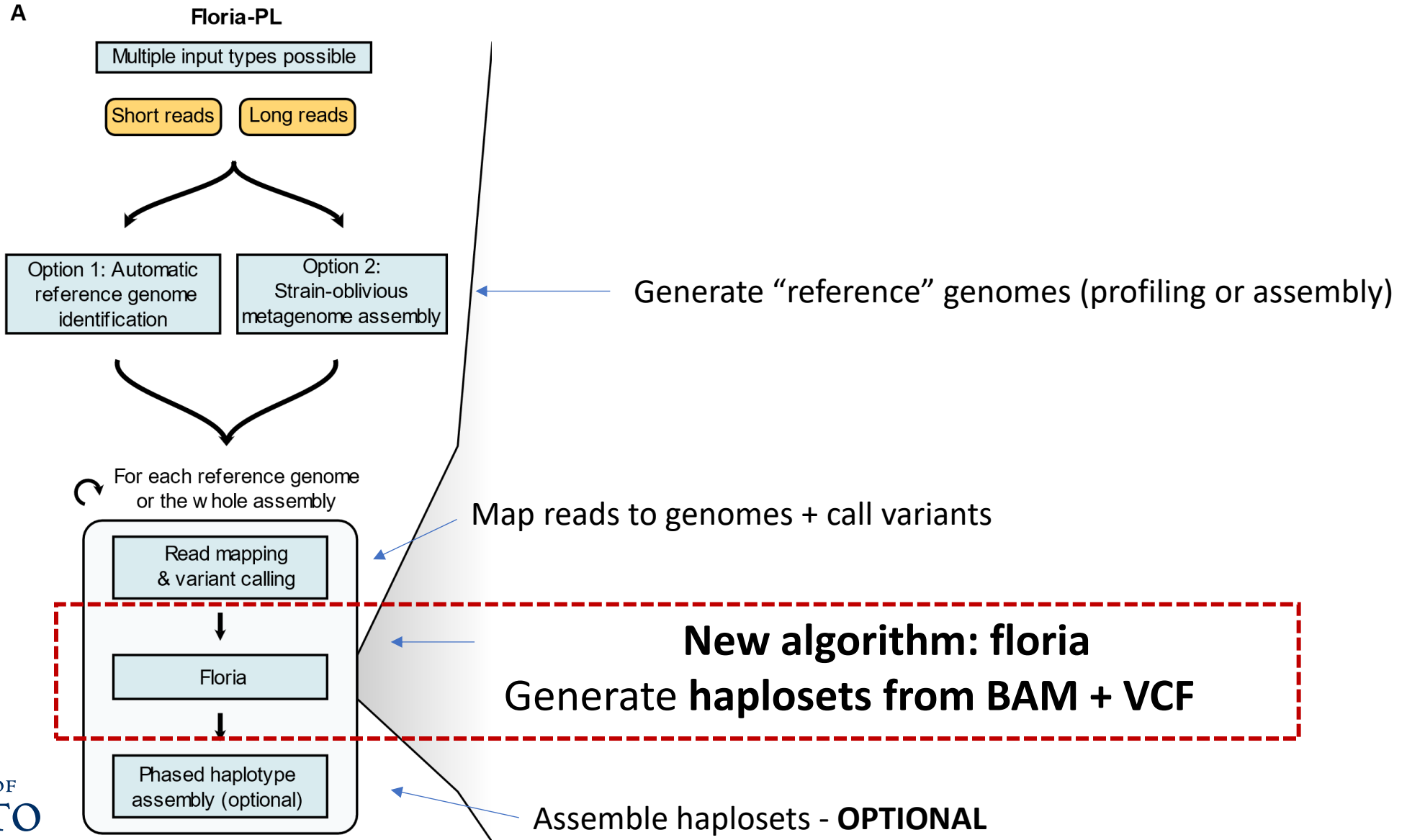
# Contributions (Shaw and Gounot et al. 2024)

- **(1) Floria**: read clustering (phasing) algorithm from **alignments + SNPs**
  - Written in Rust; documentation + conda install
- **(2) Floria-PL**: end-to-end pipeline (fastq -> assemblies)
  - Written in Snakemake - **integrating** floria



The image shows two side-by-side screenshots. The left screenshot is the Floria documentation website, featuring a blue header with the 'floria' logo and 'latest' version indicator. Below the header is a search bar labeled 'Search docs'. A dark sidebar on the left contains a list of navigation links: 'Introduction', 'Quick start', 'Tutorials', 'Usage: input and output information', 'Utility scripts and visualization', and 'How-to-guides'. The right screenshot is the GitHub repository page for 'floria - metagenomic read-based strain phasing'. It includes a 'Home' icon, the repository name, and a link to 'Edit on GitHub'. The main heading is 'floria - metagenomic read-based strain phasing'. The description states: 'floria is a software package for strain-level phasing of metagenomic sequencing samples. Given a BAM and a VCF file, floria clusters reads into strain-level clusters. Floria can:'. A bulleted list follows: '• Works with short or noisy long reads (preferably long)', '• Phase metagenomic reads or single genome reads', '• Multithreaded, takes minutes per contig/genome', and '• Minimal parameter tuning and automatic determination of strain number'.

# Floria-PL: pipeline integrating floria





# Floria: read clustering by optimization + network flows

A

Floria-PL

Multiple input types possible

Short reads

Long reads

Default

Option 1: Automatic  
reference genome  
identification

Option 2:  
Strain-oblivious  
metagenome assembly

For each reference genome  
or the whole assembly

Read mapping  
& variant calling

Floria

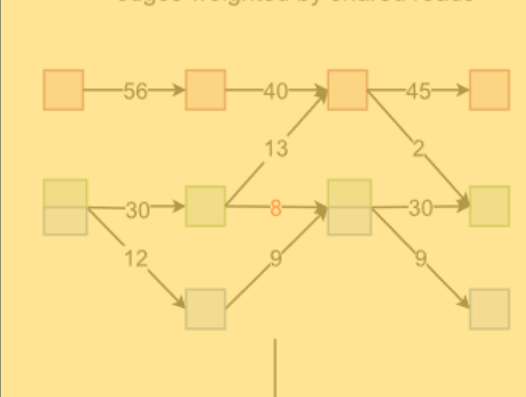
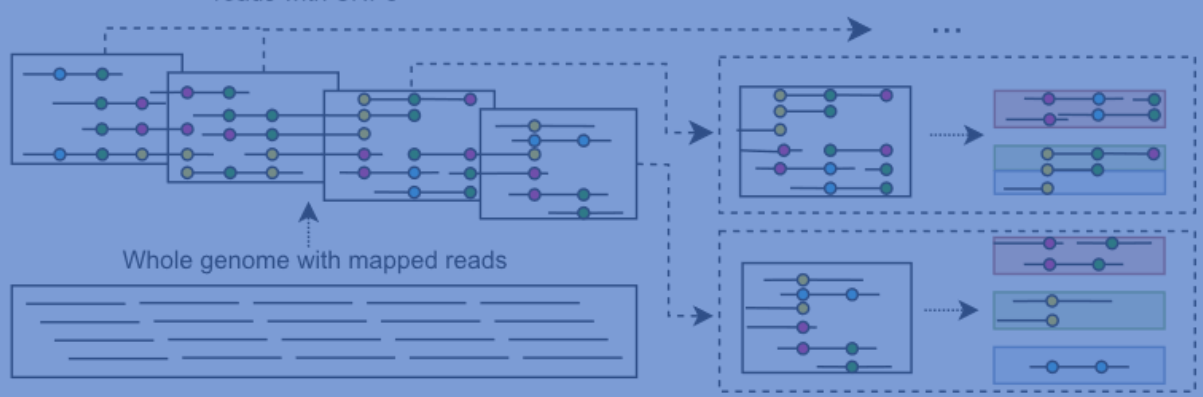
Phased haplotype  
assembly (optional)

B

(1) Segment contigs into blocks and index  
reads with SNPs

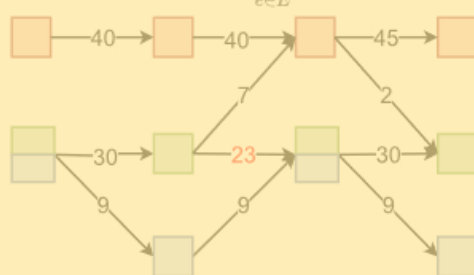
(2) Partition reads for each block by optimizing  
MEC objective; obtain "local" strain count

(3) Form a DAG with read sets as nodes and  
edges weighted by shared reads



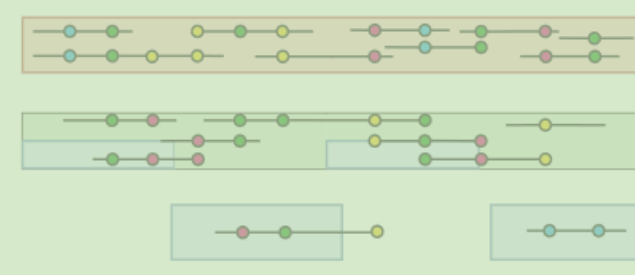
Get new coverage flow by minimizing flow "f"  
against shared reads "w" through LP

$$(4) \arg \min_{f \in Flow} \sum_{e \in E} |w(e) - f(e)|$$



(5) Obtain paths by removing paths with  
maximal minimum flow through DP on DAG

(6) Obtain haplosets (strain-specific  
read sets) by aggregating reads  
over paths



UNIVERSITY OF  
TORONTO

# Step 1: local clustering

A

Floria-PL

Multiple input types possible

Short reads Long reads

Default

Option 1: Automatic reference genome identification

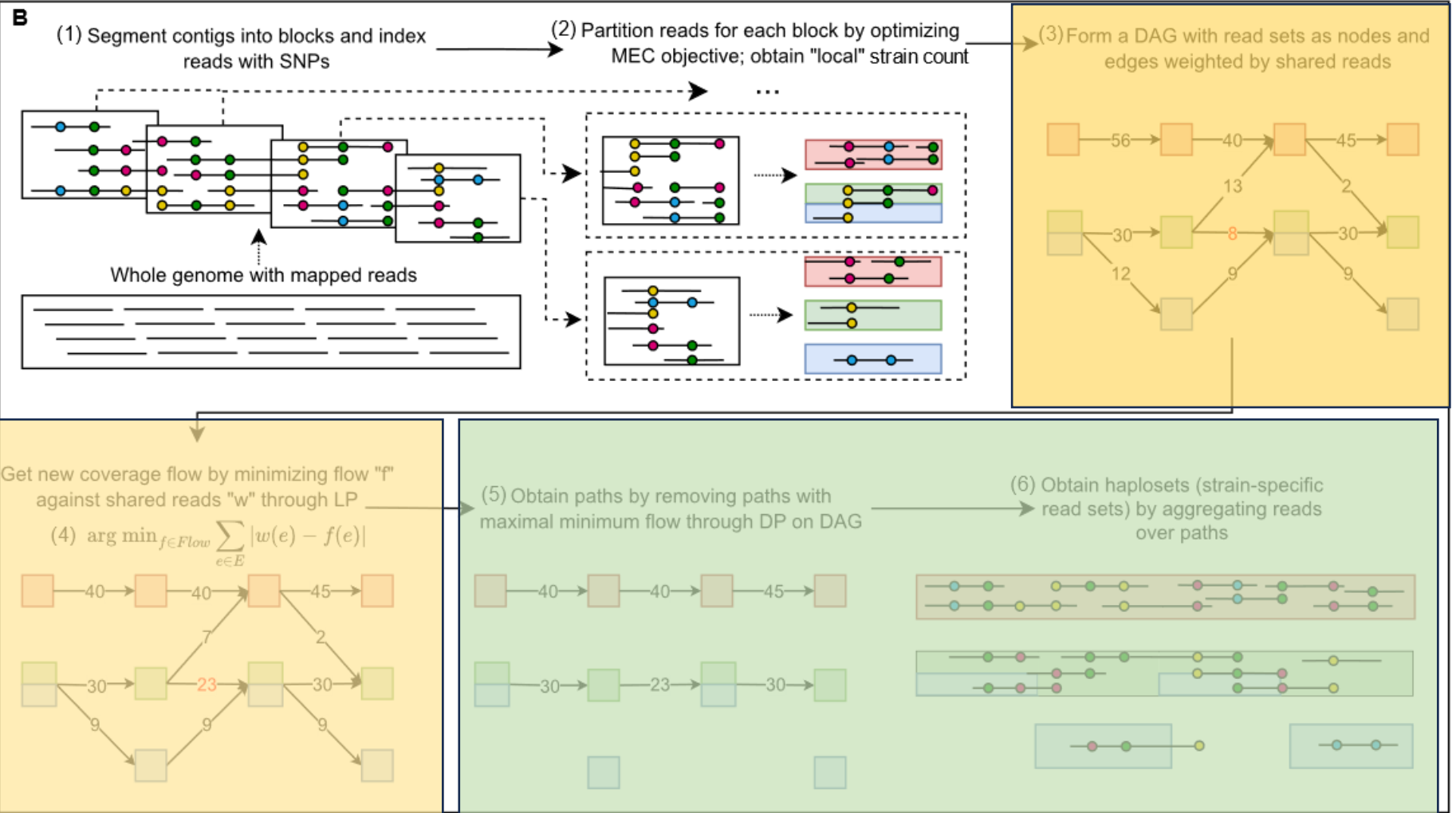
Option 2: Strain-oblivious metagenome assembly

For each reference genome or the whole assembly

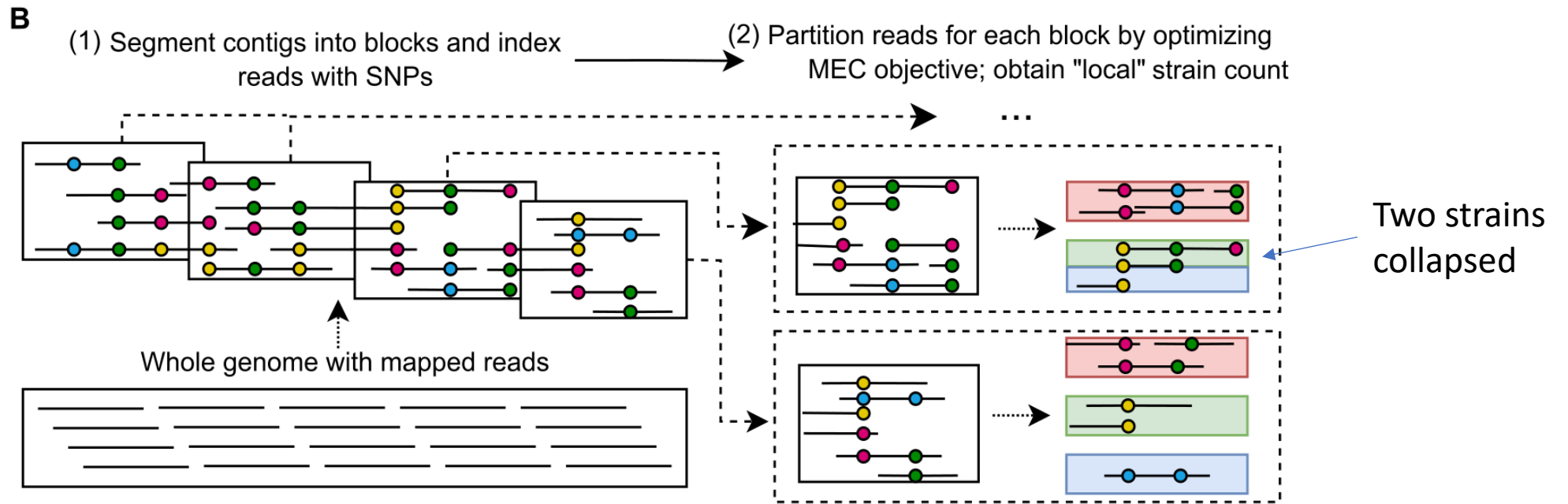
Read mapping & variant calling

Floria

Phased haplotype assembly (optional)



UNIVERSITY OF  
TORONTO



1. Clustering objective: **minimum error correction (MEC)** score
  - NP-Hard (Lancia et al., 2001)
  - **Floria**: beam search heuristic
    - often used in Natural Language Processing (NLP)
2. **# of strains?** → iteratively cluster until MEC score plateaus

# Step 1: local clustering

A

Floria-PL

Multiple input types possible

Short reads

Long reads

Default

Option 1: Automatic  
reference genome  
identification

Option 2:  
Strain-oblivious  
metagenome assembly

For each reference genome  
or the whole assembly

Read mapping  
& variant calling

Floria

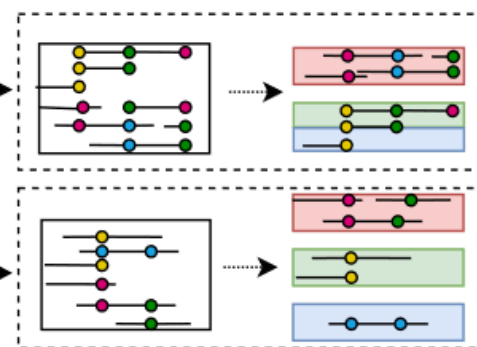
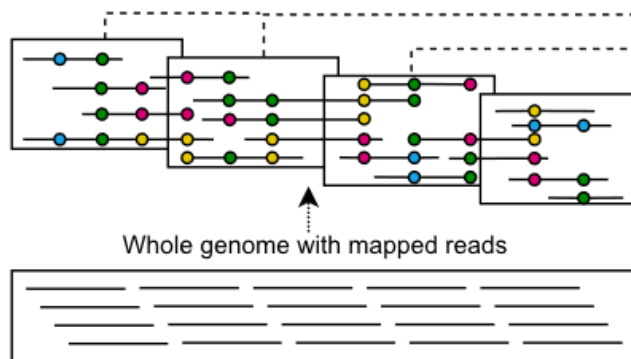
Phased haplotype  
assembly (optional)

B

(1) Segment contigs into blocks and index  
reads with SNPs

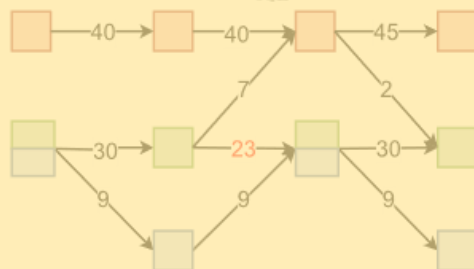
(2) Partition reads for each block by optimizing  
MEC objective; obtain "local" strain count

(3) Form a DAG with read sets as nodes and  
edges weighted by shared reads



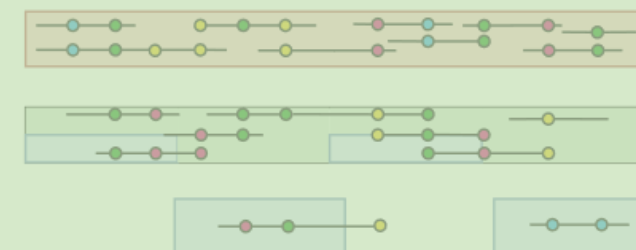
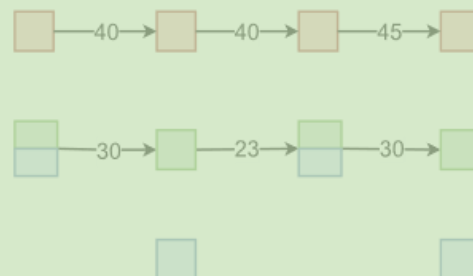
Get new coverage flow by minimizing flow "f"  
against shared reads "w" through LP

$$(4) \arg \min_{f \in Flow} \sum_{e \in E} |w(e) - f(e)|$$



(5) Obtain paths by removing paths with  
maximal minimum flow through DP on DAG

(6) Obtain haplosets (strain-specific  
read sets) by aggregating reads  
over paths



UNIVERSITY OF  
TORONTO

# Step 2: network flows

A

Floria-PL

Multiple input types possible

Short reads Long reads

Default

Option 1: Automatic reference genome identification

Option 2: Strain-oblivious metagenome assembly

For each reference genome or the whole assembly

Read mapping & variant calling

Floria

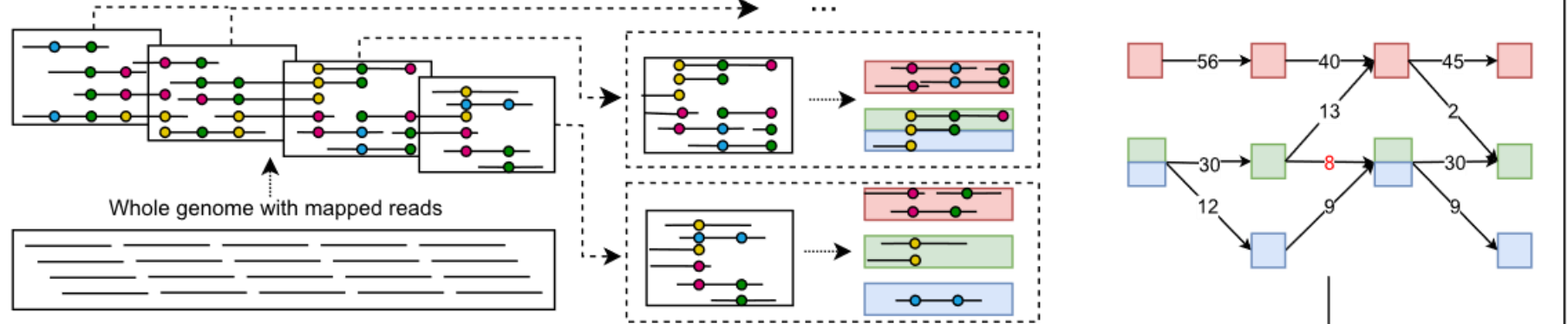
Phased haplotype assembly (optional)

B

(1) Segment contigs into blocks and index reads with SNPs

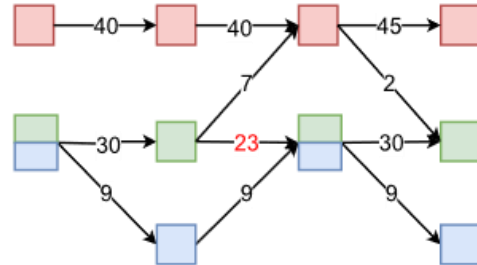
(2) Partition reads for each block by optimizing MEC objective; obtain "local" strain count

(3) Form a DAG with read sets as nodes and edges weighted by shared reads



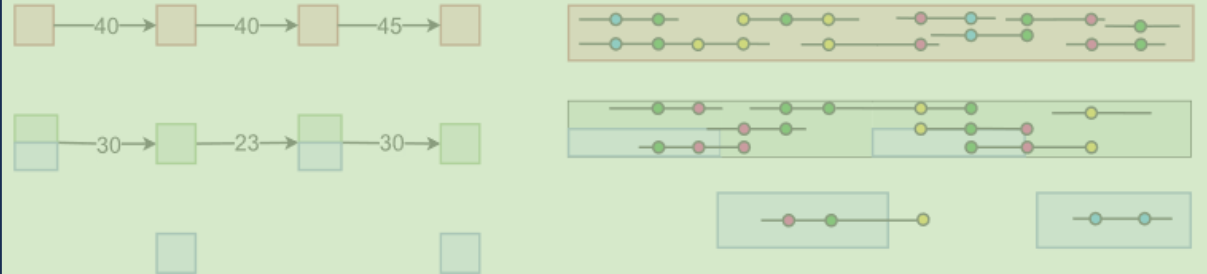
Get new coverage flow by minimizing flow "f" against shared reads "w" through LP

$$(4) \arg \min_{f \in Flow} \sum_{e \in E} |w(e) - f(e)|$$



(5) Obtain paths by removing paths with maximal minimum flow through DP on DAG

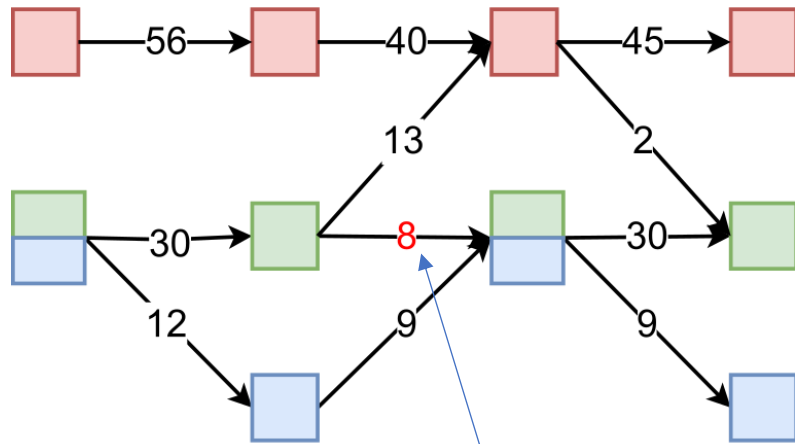
(6) Obtain haplosets (strain-specific read sets) by aggregating reads over paths



UNIVERSITY OF  
TORONTO

(DAG = Directed acyclic graph)

(3) Form a DAG with read sets as nodes  
and edges weighted by shared reads →



Error!  
bad edge

By linear programming:

$$\arg \min_{f \in Flow} \sum_{e \in E} |w(e) - f(e)|$$

$f$  – network flow

$w$  – original weights

# Step 2: network flows

A

Floria-PL

Multiple input types possible

Short reads Long reads

Default

Option 1: Automatic reference genome identification

Option 2: Strain-oblivious metagenome assembly

For each reference genome or the whole assembly

Read mapping & variant calling

Floria

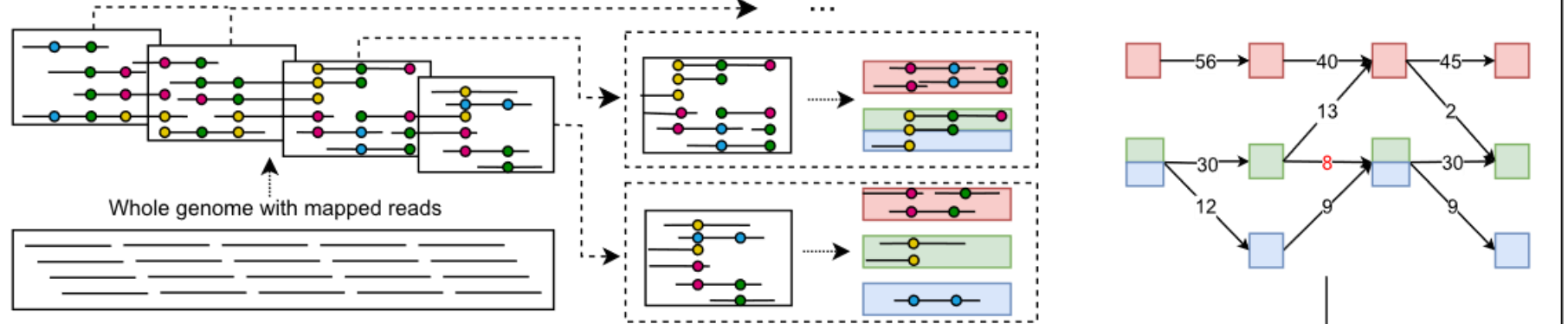
Phased haplotype assembly (optional)

B

(1) Segment contigs into blocks and index reads with SNPs

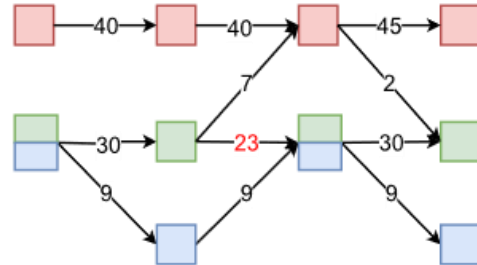
(2) Partition reads for each block by optimizing MEC objective; obtain "local" strain count

(3) Form a DAG with read sets as nodes and edges weighted by shared reads



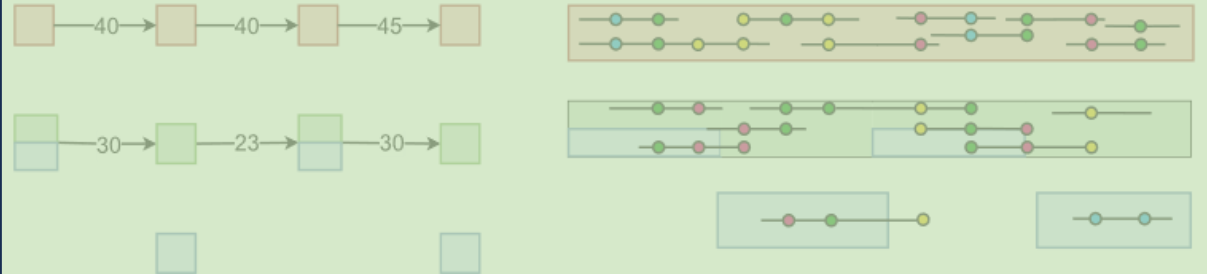
Get new coverage flow by minimizing flow "f" against shared reads "w" through LP

$$(4) \arg \min_{f \in Flow} \sum_{e \in E} |w(e) - f(e)|$$



(5) Obtain paths by removing paths with maximal minimum flow through DP on DAG

(6) Obtain haplosets (strain-specific read sets) by aggregating reads over paths



UNIVERSITY OF  
TORONTO



# Step 3: Obtaining haploset paths

A

Floria-PL

Multiple input types possible

Short reads Long reads

Default

Option 1: Automatic reference genome identification

Option 2: Strain-oblivious metagenome assembly

For each reference genome or the whole assembly

Read mapping & variant calling

Floria

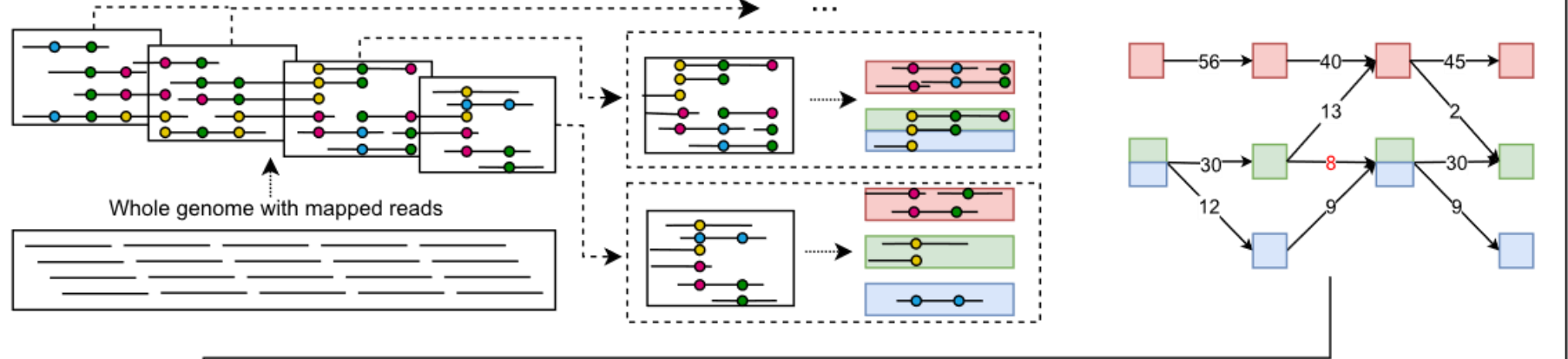
Phased haplotype assembly (optional)

B

(1) Segment contigs into blocks and index reads with SNPs

(2) Partition reads for each block by optimizing MEC objective; obtain "local" strain count

(3) Form a DAG with read sets as nodes and edges weighted by shared reads

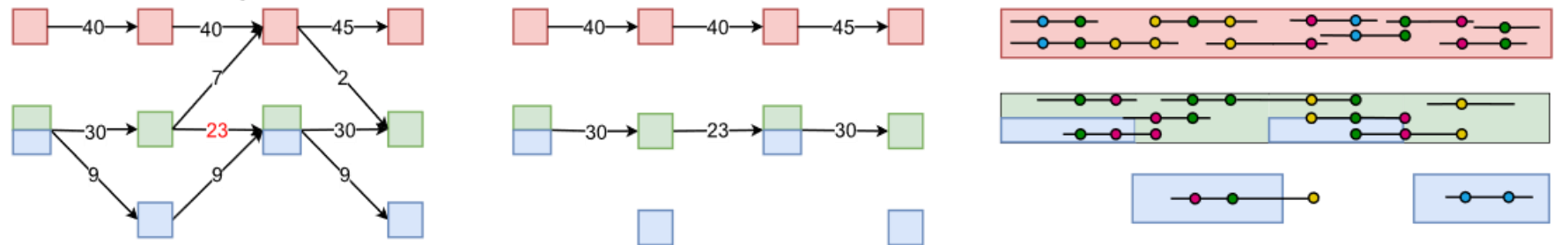


Get new coverage flow by minimizing flow "f" against shared reads "w" through LP

$$(4) \arg \min_{f \in Flow} \sum_{e \in E} |w(e) - f(e)|$$

(5) Obtain paths by removing paths with maximal minimum flow through DP on DAG

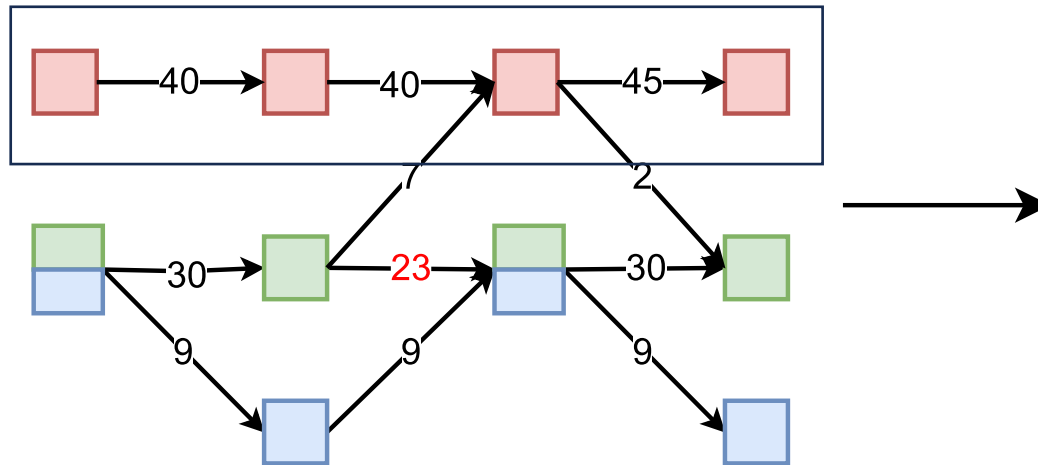
(6) Obtain haplosets (strain-specific read sets) by aggregating reads over paths



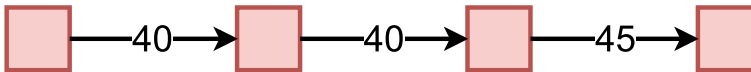
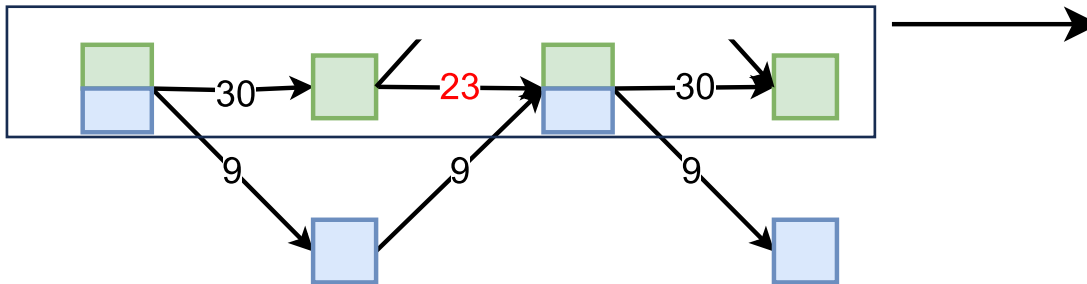
UNIVERSITY OF  
TORONTO

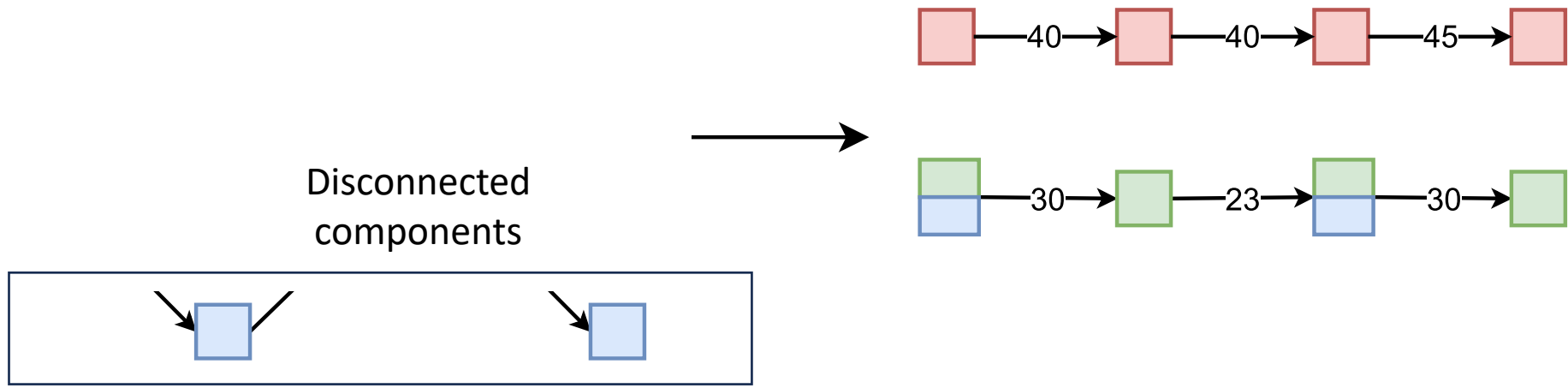


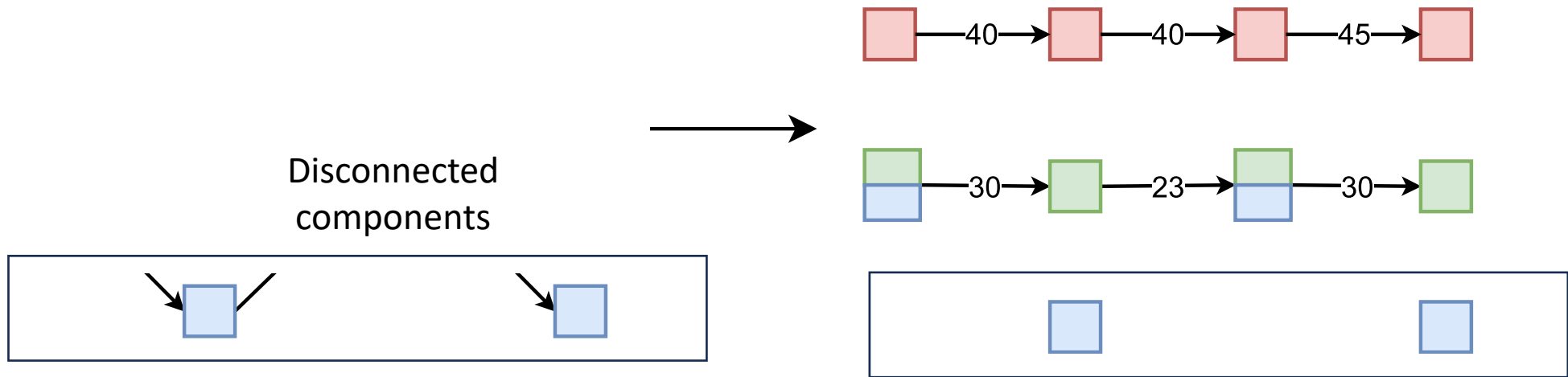
## Largest *minimum* flow path (via dynamic programming - DP)



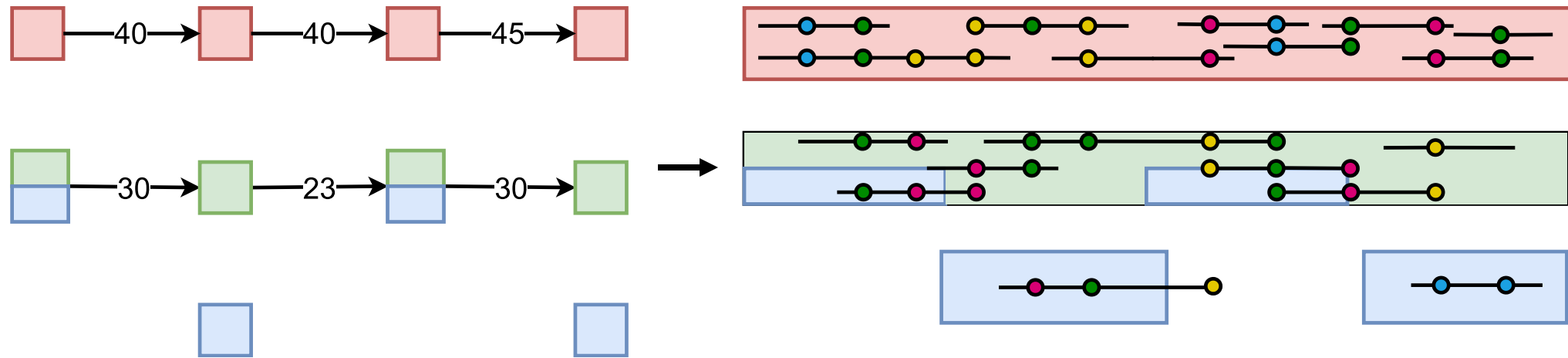
## Largest *minimum* flow path (via dynamic programming - DP)







(6) Obtain haplosets (strain-specific read sets)  
by aggregating reads over paths



# Step 3: trim flow graph to obtain haplosets

A

Floria-PL

Multiple input types possible

Short reads

Long reads

Default

Option 1: Automatic  
reference genome  
identification

Option 2: Strain-oblivious  
metagenome assembly

For each reference genome  
or the whole assembly

Read mapping  
& variant calling

Floria

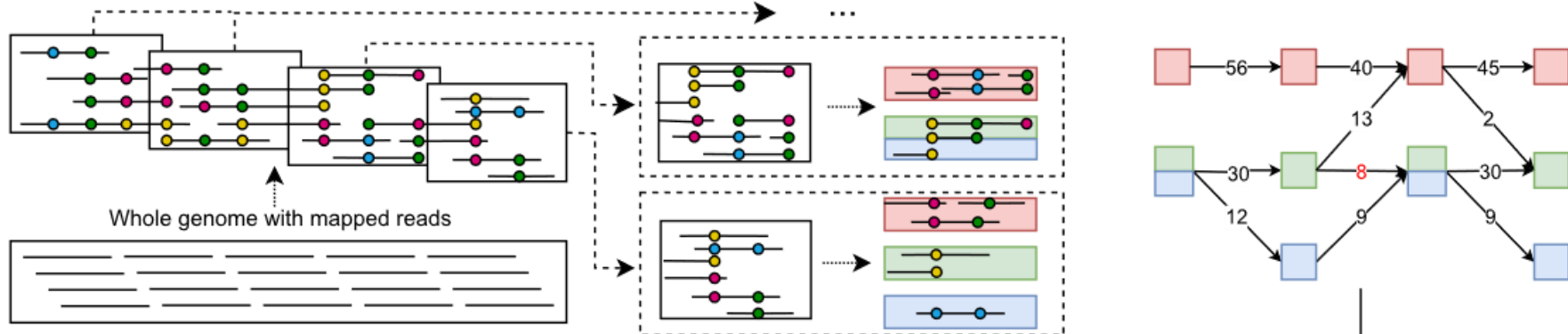
Phased haplotype  
assembly (optional)

B

(1) Segment contigs into blocks and index  
reads with SNPs

(2) Partition reads for each block by optimizing  
MEC objective; obtain "local" strain count

(3) Form a DAG with read sets as nodes and  
edges weighted by shared reads

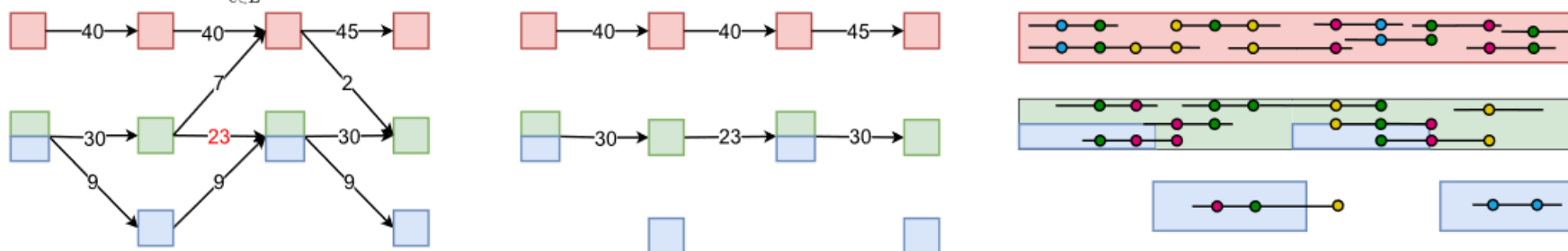


Get new coverage flow by minimizing flow "f"  
against shared reads "w" through LP

$$(4) \arg \min_{f \in Flow} \sum_{e \in E} |w(e) - f(e)|$$

(5) Obtain paths by removing paths with  
maximal minimum flow through DP on DAG

(6) Obtain haplosets (strain-specific  
read sets) by aggregating reads  
over paths



UNIVERSITY OF  
TORONTO

# Benchmarking results

# Benchmarking: simulated metagenome

## **Synthetic metagenomes:**

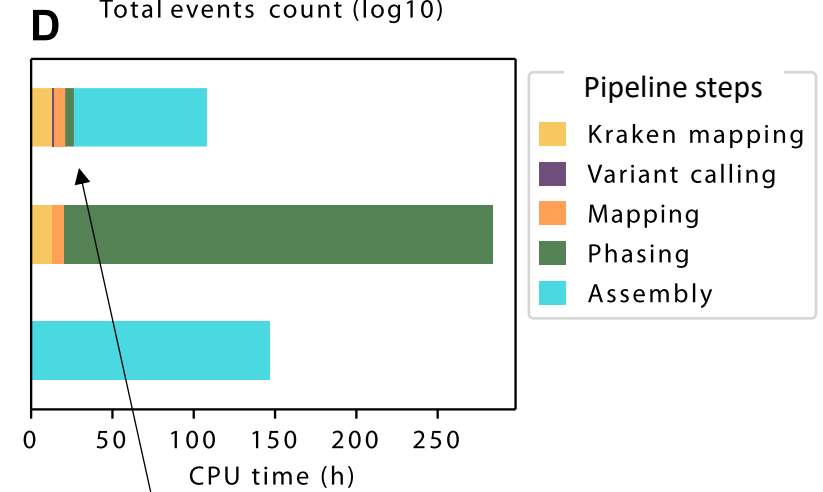
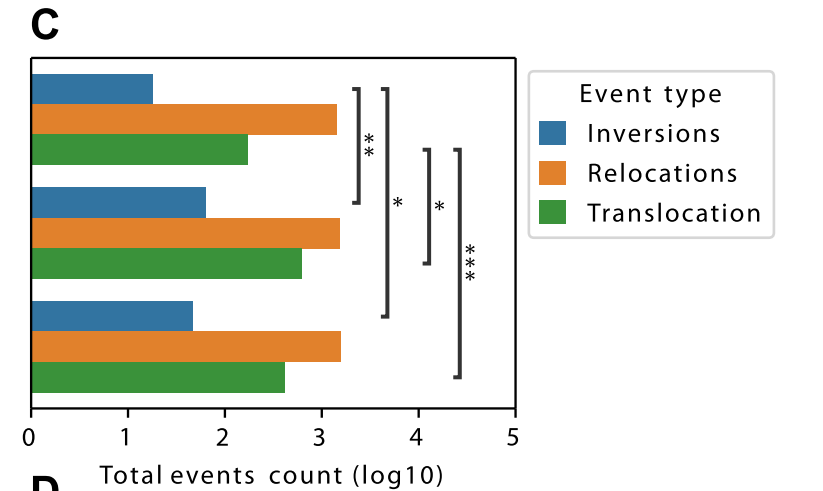
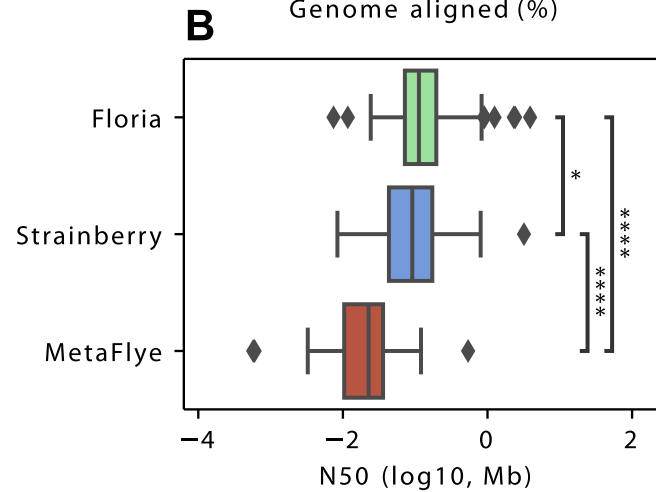
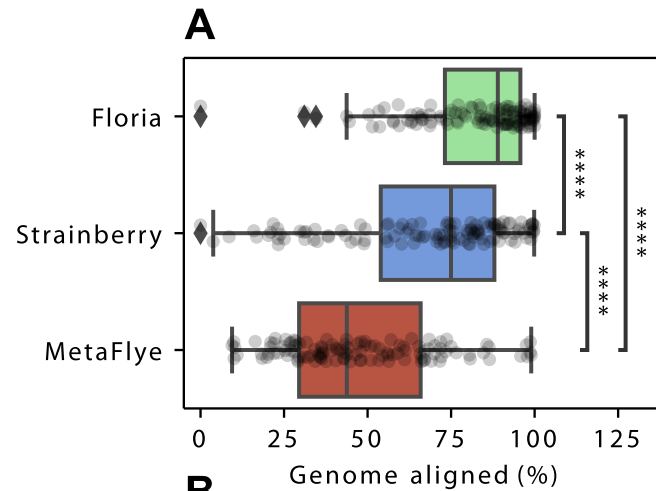
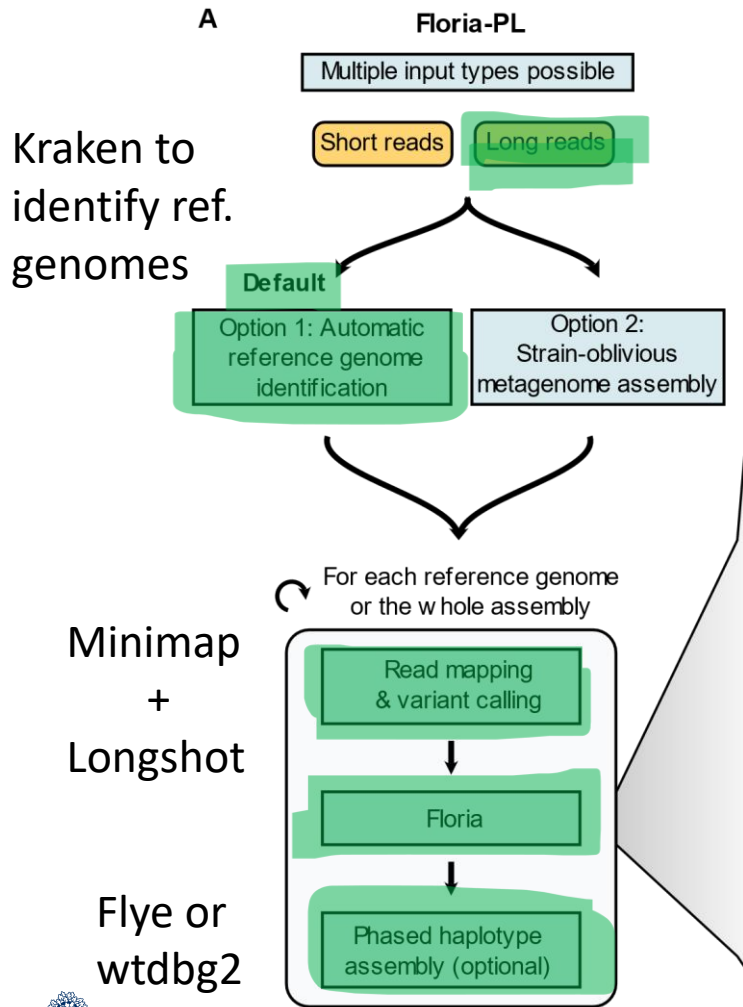
- 40 common gut species
- 1-5 strain genomes per species
- Synthetic noisy nanopore reads (**88% identity**)
  - Random strain coverages between [5,25]

## **Comparison against:**

1. metaFlye – metagenome assembler (Kolmogorov et al. 2020)
2. Strainberry – strain assembler (Vicedomini et al. 2021)



# Assembly benchmarking



Floria **phasing** time << **assembly** time!

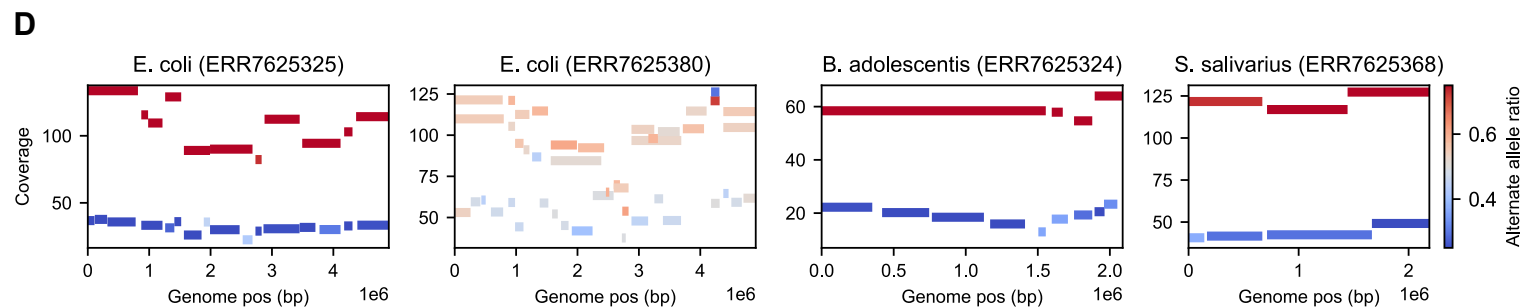


Results: **real** metagenomes  
Long read **AND** short read

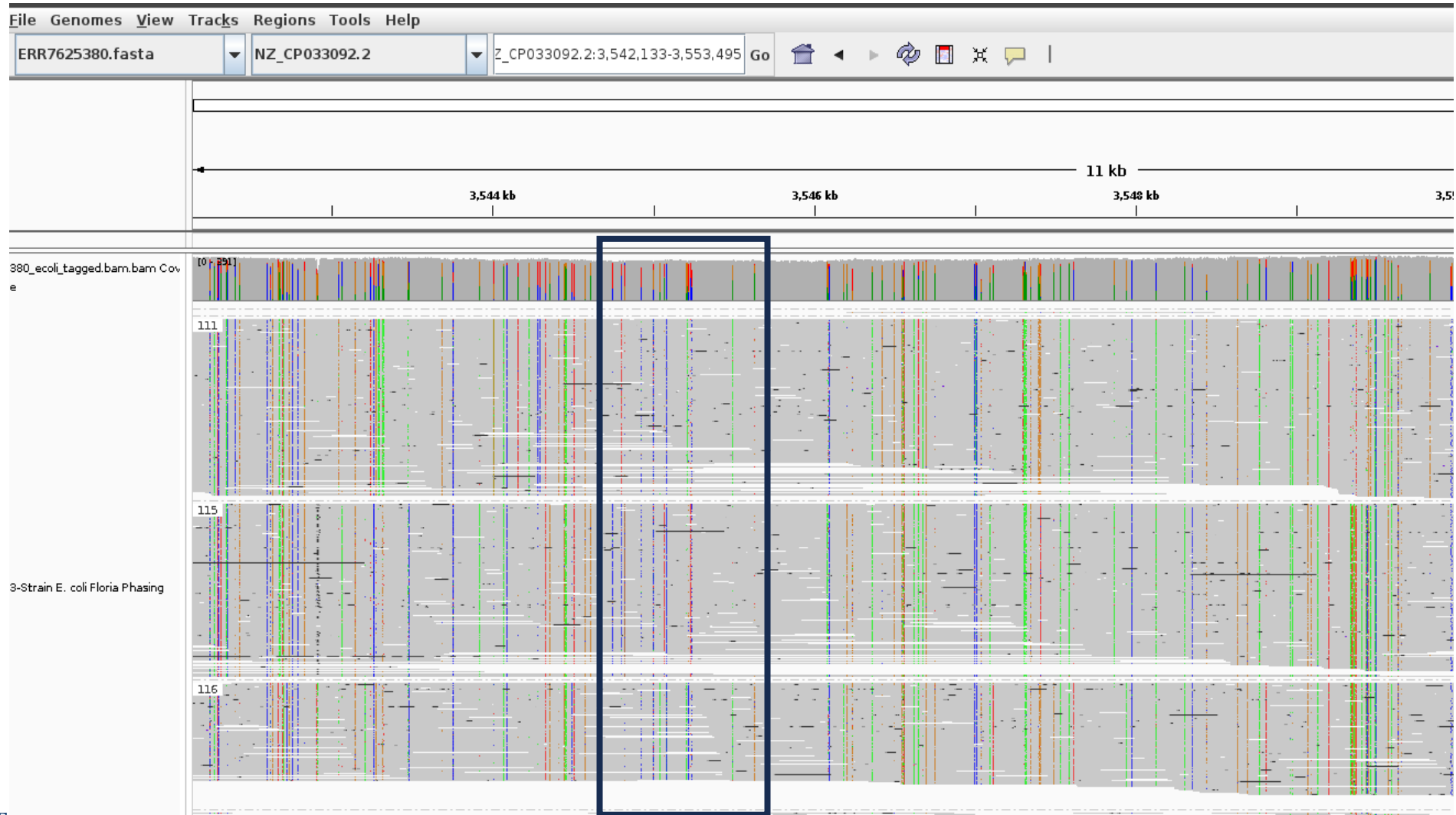
# Floria on 109 gut nanopore samples!

(dataset: Gounot et al. 2022, Nat. Comms.)

< 15 mins per sample  
for phasing!



# Floria allows for visualization – 3-strain *E. coli*



Consistently different alleles

# Longitudinal strain tracking with floria

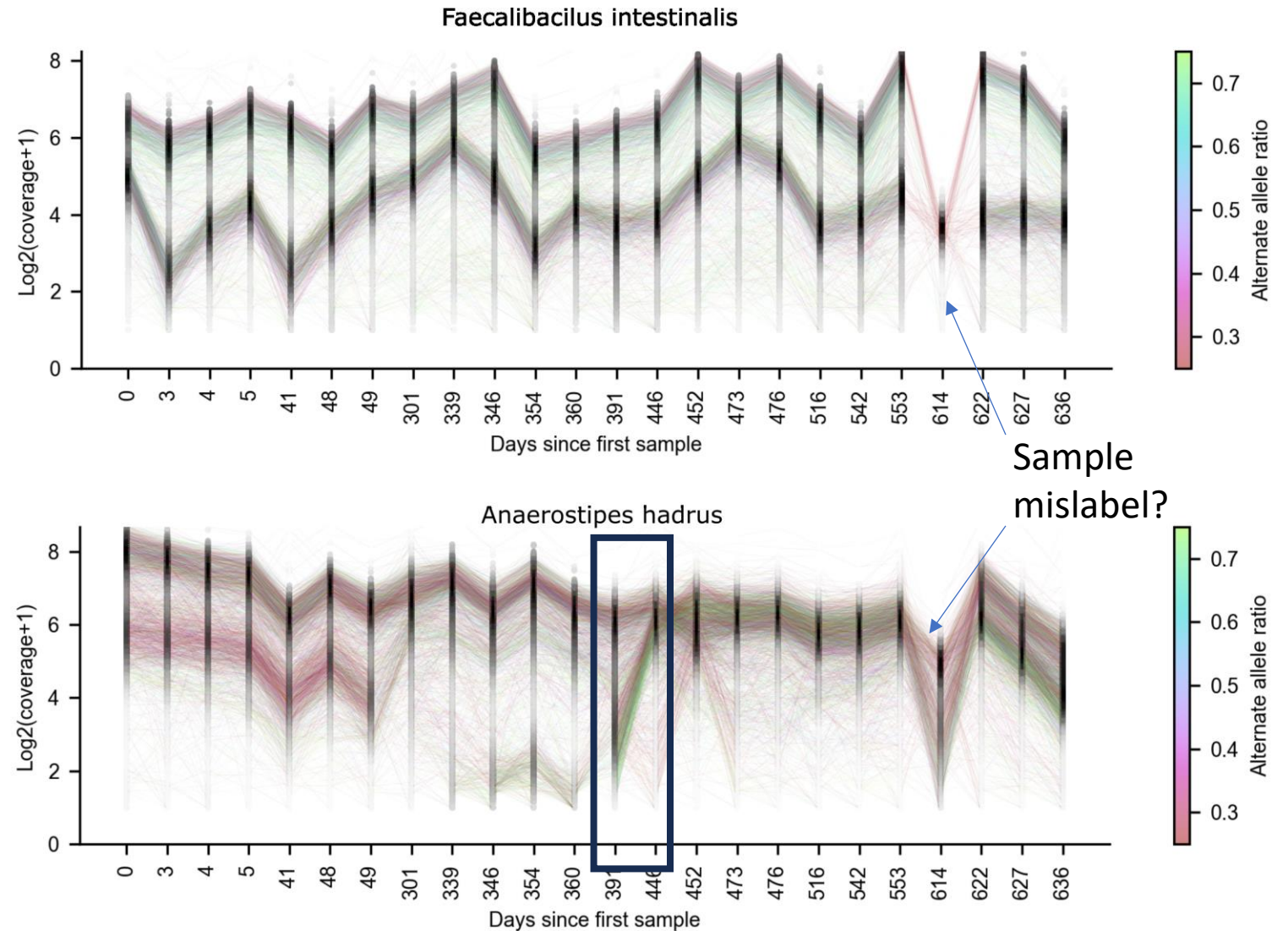
- 24 longitudinal **SHORT-READ** gut samples (636 days)
  - Dataset from “*Metabolic independence drives gut microbial colonization and resilience in health and disease*” by Watson et al. (Genome Biology 2023)
- Run floria → obtain haplosets → track across time

Two species with  
interesting patterns:

*F. intestinalis* - **stable**

*A. hadrus* – **transient**

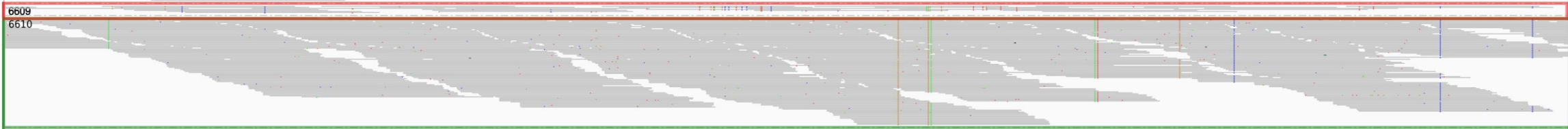
Line: haploset match across sample



Low-abundance strain to high-abundance strain  
emergence? **Visualize!**



## Day 391 - A. hadrus



Zoomed in on red  
haploset (“6609”):  
LOW coverage



NC\_021016.1:2,664,539-2,665,758

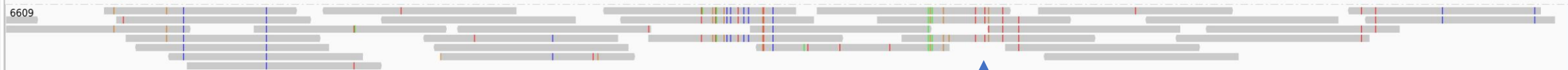


**Day 391**

Day 391 A. hadrus

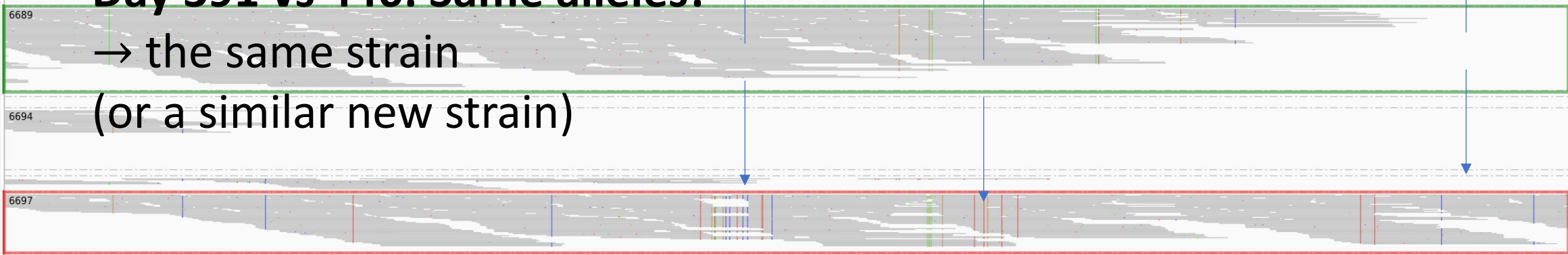


Day 391 A. hadrus magnified  
Haploset 6609



**Day 446**

Day 446 A. hadrus



**Day 391 vs 446: Same alleles!**  
→ the same strain  
(or a similar new strain)

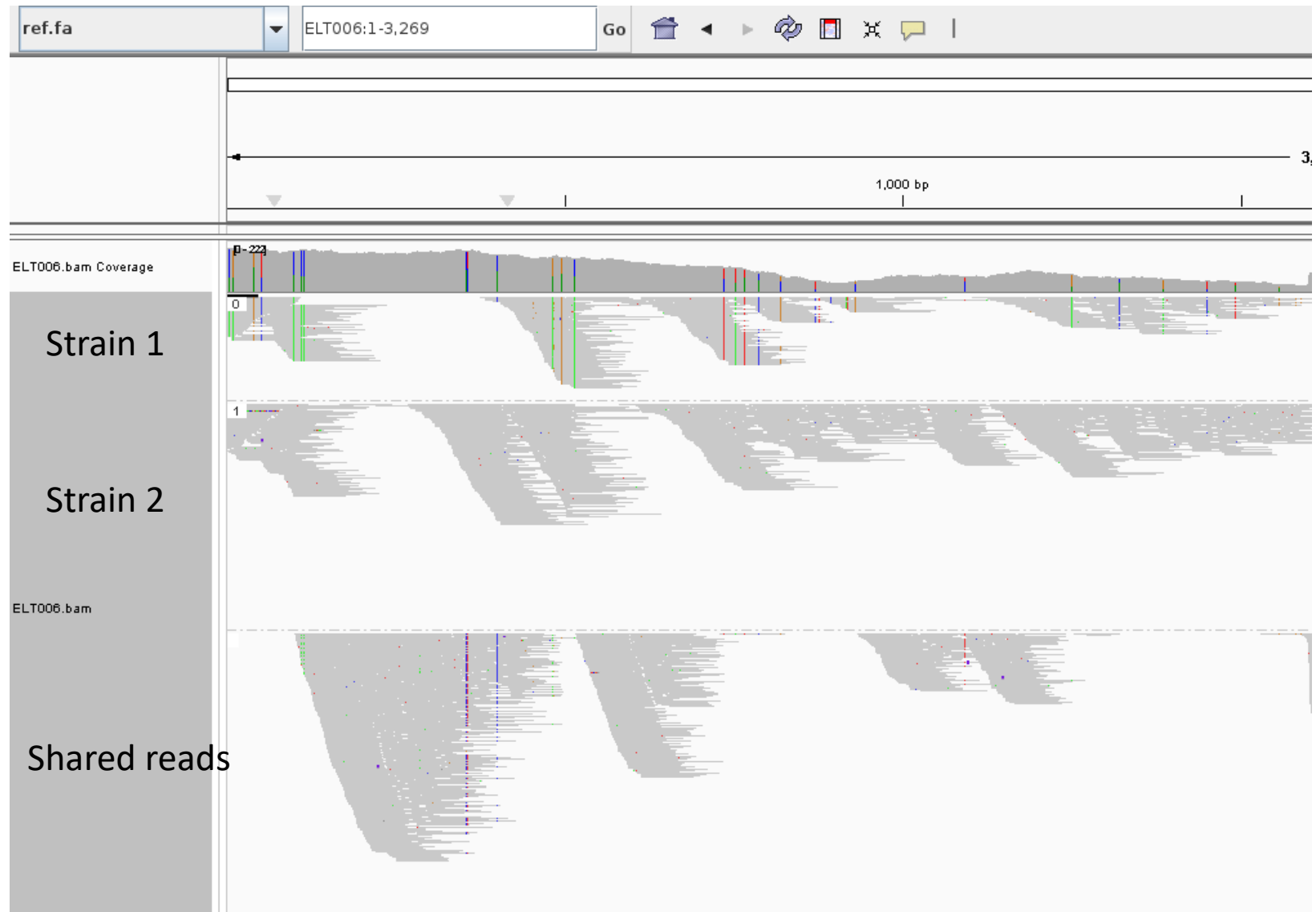
High coverage haploset



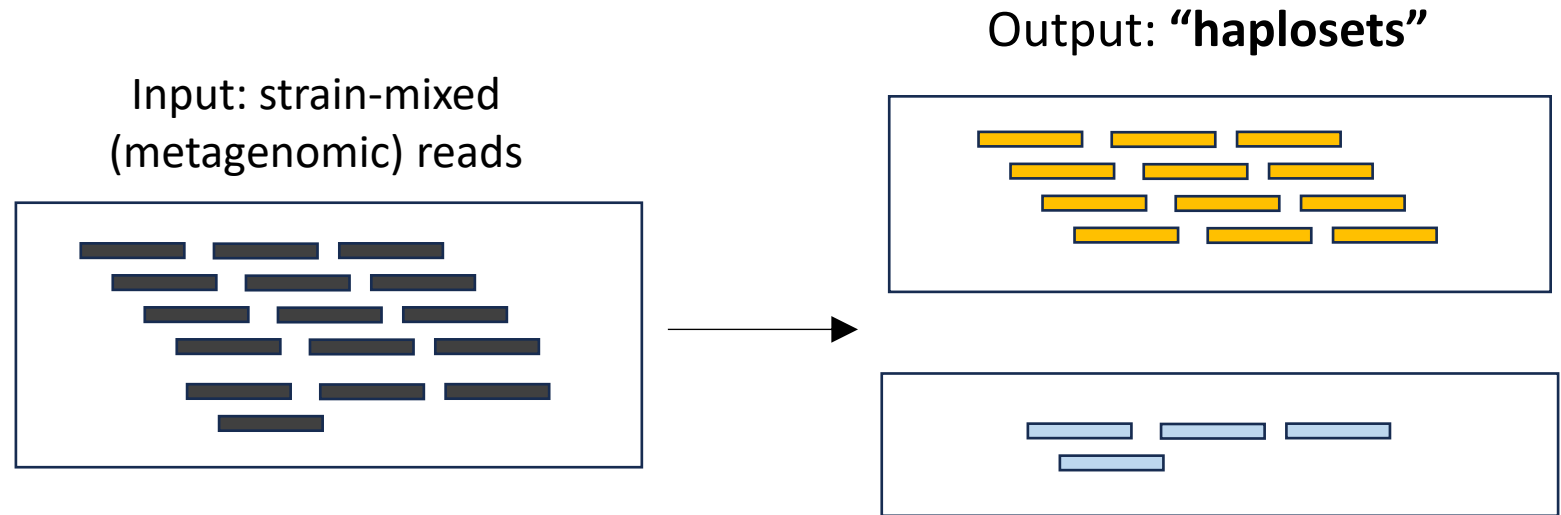
# Floria works on ancient viral metagenomes!?

**From Maxime Borry** (Postdoc at Max Planck for Evo Anthro):

**Ancient mixed infection of hepatitis B** (Kocher et al., 2021, Science)



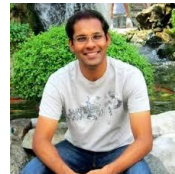
# Conclusion





- Developed **floria**: a strain-level read clustering (phasing) algorithm
  - Short OR long reads
  - < 20 min per metagenome
  - assembly optional
- Fast, informative, and versatile
  - can do strain-level metagenomic assembly...
  - but even more useful for **hypothesis generation and data sleuthing**

# Acknowledgements + authors

- Jim Shaw (**presenter**)
  - PhD at University of Toronto (2024)
  - incoming postdoc at Harvard Medical School / DFCI
- Jean-Sébastien Gounot
  - Research fellow at Genome Institute of Singapore
  - **co-lead author**
- Hanrong Chen
  - Postdoc at Genome Institute of Singapore
- Niranjan Nagarajan
  - Genome Institute of Singapore, **co-lead PI**
- Yun William Yu
  - Carnegie Mellon University, **co-lead PI**



## Floria: fast and accurate strain haplotyping in metagenomes

Jim Shaw, Jean-Sebastien Gounot, Hanrong Chen, Niranjan Nagarajan , Yun William Yu  [Author Notes](#)

*Bioinformatics*, Volume 40, Issue Supplement\_1, July 2024, Pages i30–i38,  
<https://doi.org/10.1093/bioinformatics/btae252>

**Published:** 28 June 2024



Natural Sciences and Engineering  
Research Council of Canada

Conseil de recherches en sciences  
naturelles et en génie du Canada



**Supported by NSERC CGS-D**

Thanks to Maxime Borry  
for ancient DNA results