Guest lecturer: Jillian Cieslik

# Lecture 19
# Genome assembly



*Spigelia marilandica*
(c) Matilda Adams/
Missouri Botanical Garden

Course:     Practical Bioinformatics (BIOL 4220)
Instructor:  Michael Landis
Email:       michael.landis@wustl.edu

# Lecture 19 outline

Last time: jupyter, matplotlib

This time: genome assembly

- genome sequences
- genome sizes
- genome assembly

# Sequencing

true sequence      ACGGTATATATACCGA

sequence
copies

ACGGTATATATACCGA
ACGGTATATATACCGA
ACGGTATATATACCGA

sequence
fragments
(reads)

ACGGTATA TATACCGA
ACGGTATAT ATACCGA
AC GGTATATA TACCGA
ACGGTA TATATACC GA

# Assembly

unordered reads
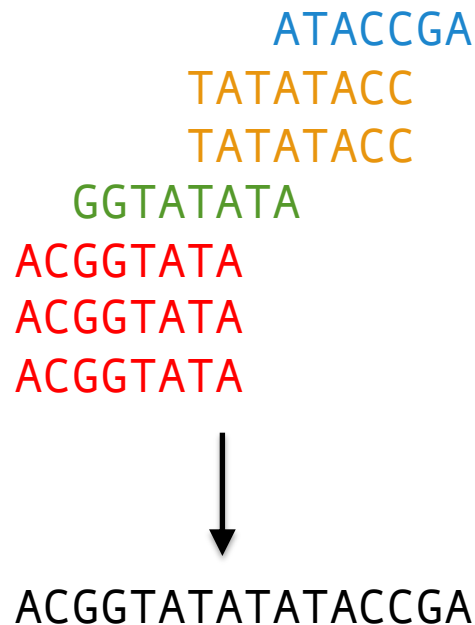
ATACCGA
ACGGTATA
GGTATATA
TATATACC

aligned reads

ATACCGA
TATATACC
GGTATATA
ACGGTATA

assembled sequence

ACGGTATATATACCGA

# Assembly

ATACCGA
TATATACC
TATATACC
GGTATATA
ACGGTATA
ACGGTATA
ACGGTATA

↓

ACGGTATATATACCGA

Would be easy if we knew how reads were aligned

We would retrieve the original genome sequence with no effort

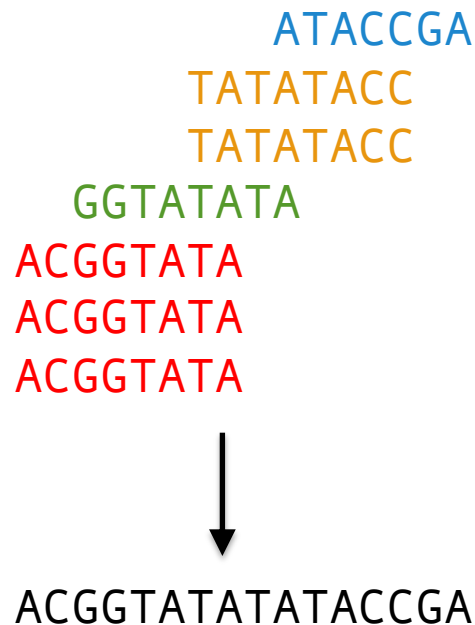Instead, we have an unordered and unaligned bag of reads

True
genome



Sequenced
reads



Assembled
genome

# How do we assemble reads?

ATACCGA
TATATACC
TATATACC
GGTATATA
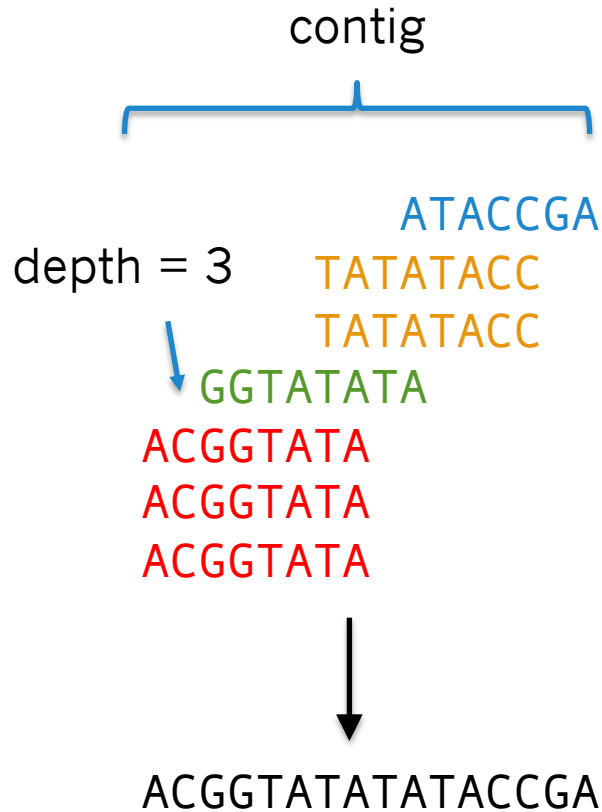ACGGTATA
ACGGTATA
ACGGTATA

↓

ACGGTATATATACCGA

Can we do global pairwise alignments for each pair of reads?

Let's make a block of contiguously mapped reads (contig)

Reads can align to any contig

Read mapped to contig with best score

# Basic unit of assembly

contig

We want high-coverage contigs

ATACCGA

depth = 3

TATATACC

TATATACC

GGTATATA

ACGGTATA

ACGGTATA

ACGGTATA

ACGGTATATATACCGA

depth = # reads mapped for one site

$$\text{avg. coverage} = \frac{\text{\# mapped sites}}{\text{contig size}}$$

$$\text{est. coverage} = \frac{\text{\# reads * read length}}{\text{genome size}}$$

$$\text{avg. coverage} = \frac{(8 + 8 + 8 + 8 + 8 + 8 + 7)}{16}$$

# Short read dataset sizes

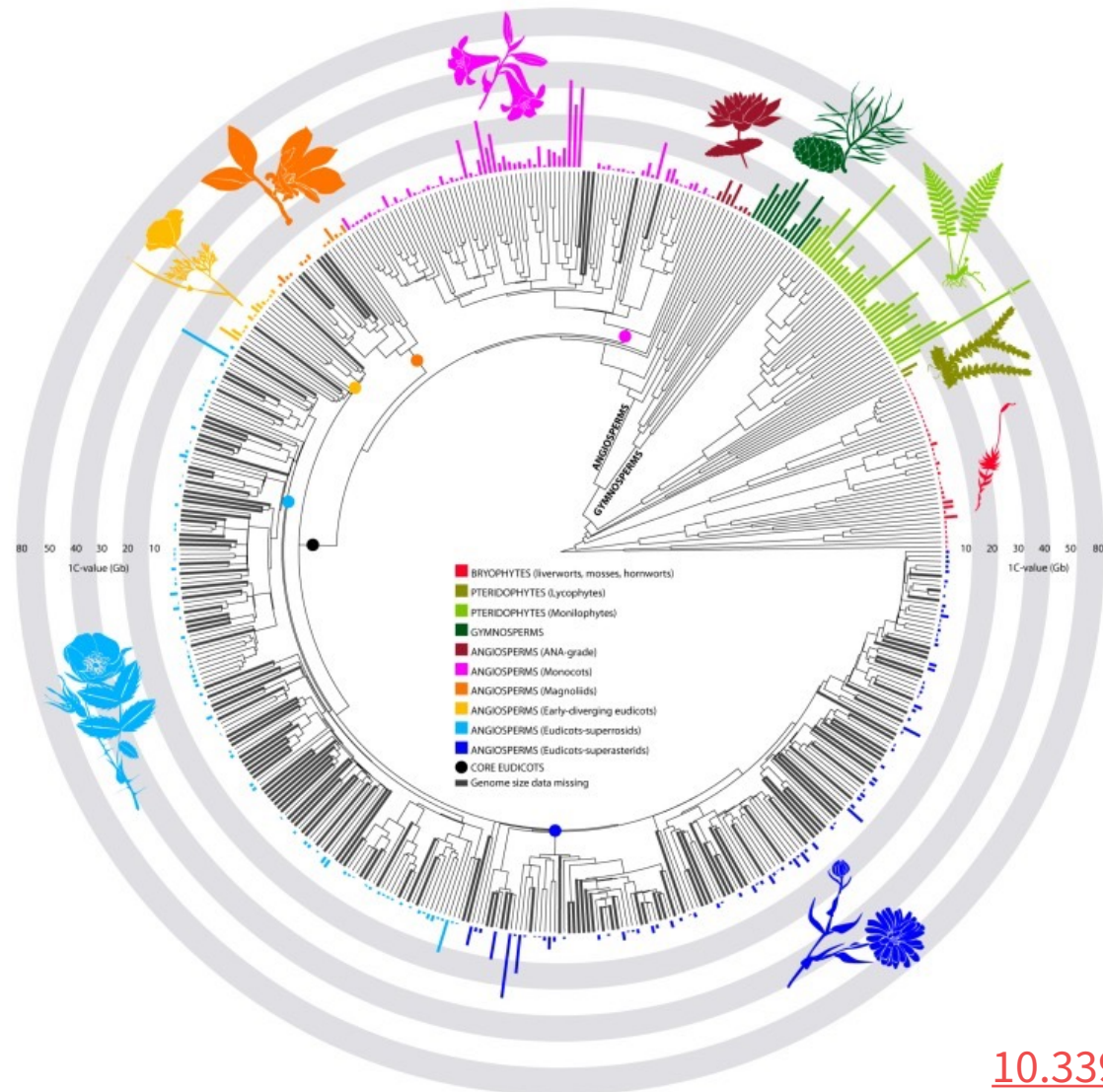How many 150 bp length reads needed
for 30x coverage?

| Species | #bp | #reads |
|---|---|---|
| SARS-CoV-2 | $2 \times 10^4$ | $4 \times 10^3$ |
| E. coli | $4.5 \times 10^6$ | $9 \times 10^5$ |
| Human | $3.2 \times 10^9$ | $6.4 \times 10^8$ |
| Fern | $1.6 \times 10^{11}$ | $3.2 \times 10^{10}$ |

# Genome sizes

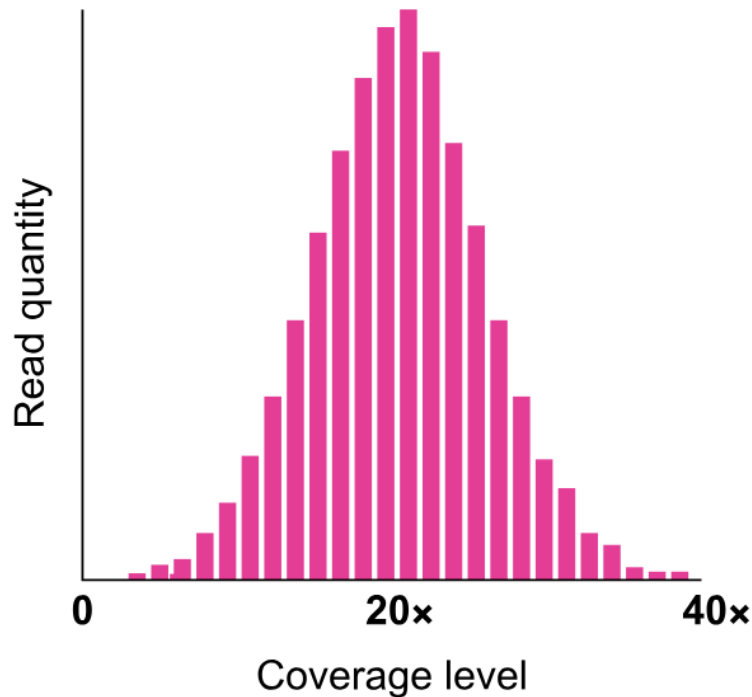What generates variation in genome size?

- Intron length
- Repetitive regions
- Transposable elements
- Whole genome duplication/polyploidy
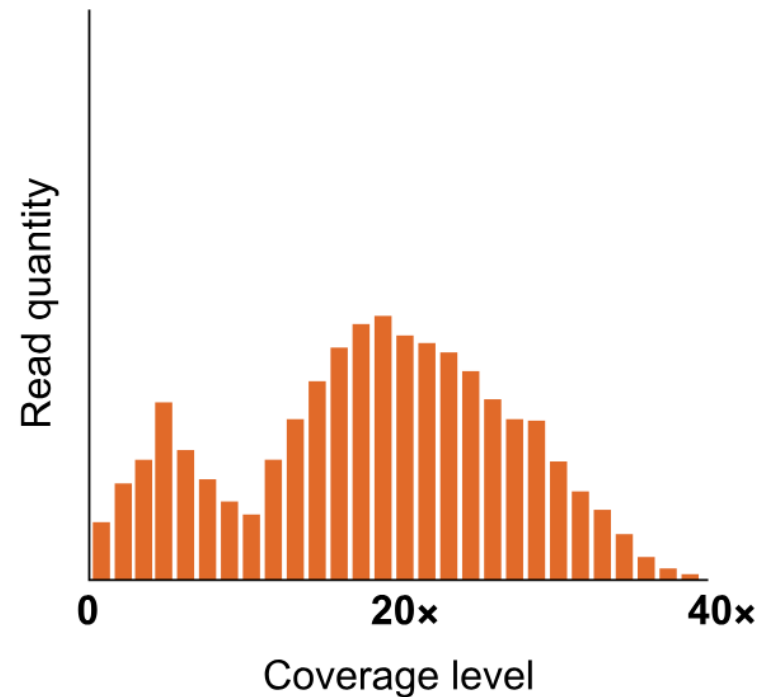- Number of genes (not always…)

# Genome sizes

# Coverage distributions

# Assembly problem

Naive assembly would require $N^2$ pairwise alignments.

TATATACC
| | | | |
ACGGTATA        ATACCGA
                | |   |
                ACGGTATA

        GGTATATA
        | | | | | | |
        ACGGTATA

                GGTATATA
                | | | | | | |
GGTATATA        ATACCGA        TATATACC
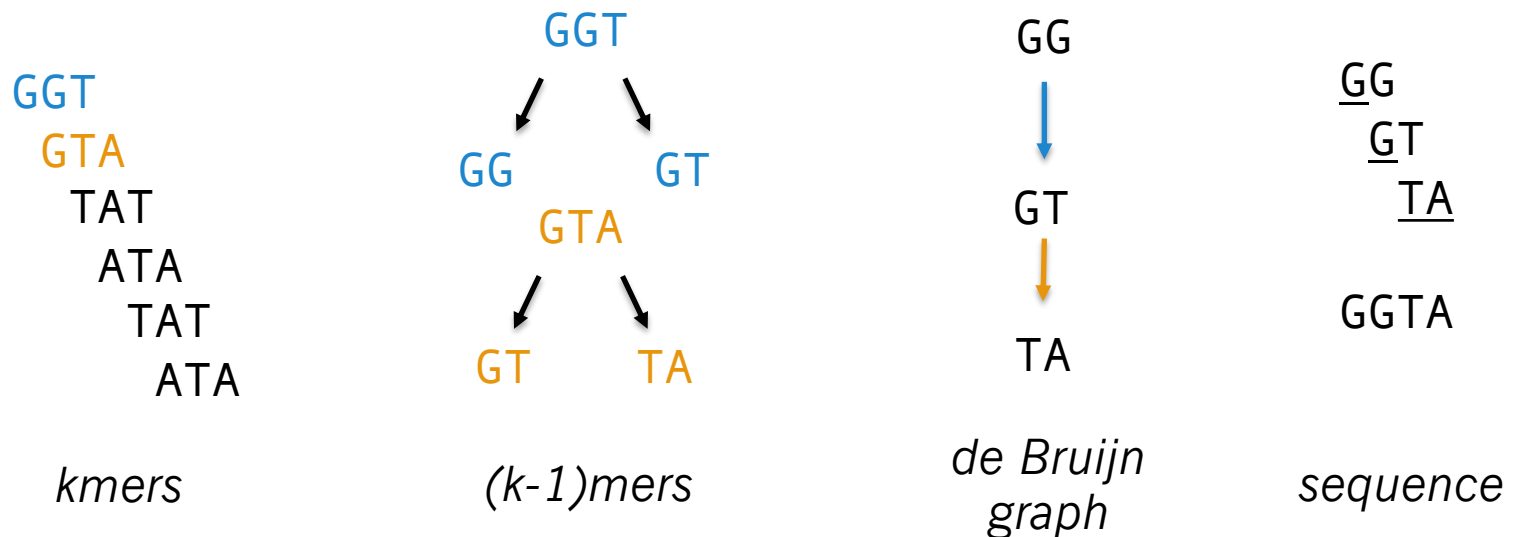| | | |         | | |          | | | |
TATATACC                       ACGGTATA

Not possible for short read datasets!

*e.g.* $10^{18}$ alignments for $10^9$ reads

# de Bruijn graph

- Choose kmer length (often 40 < k < 100)
- Make left and right (k-1)mers for each kmer
- Add node for (k-1)mer if it doesn't exist
- Add edge from left (k-1)mer to right (k-1)mer



| kmers | (k-1)mers | de Bruijn graph | sequence |

# Graph construction

One read

AC → CG → GG → GT → TA → AT

ACGGTATA

ACG
 CGG
  GGT
   GTA
    TAT
     ATA

# Graph construction

Two reads



AC → CG → GG → GT → TA → AT

CC    GA

Read 1

ACGGTATA

ACG
CGG
GGT
GTA
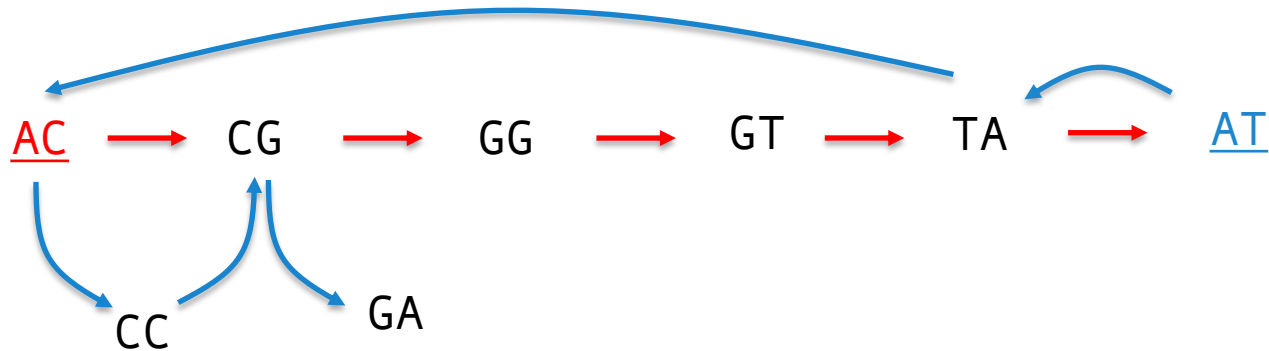TAT
ATA

Read 2

ATACCGA

ATA
TAC
ACC
CCG
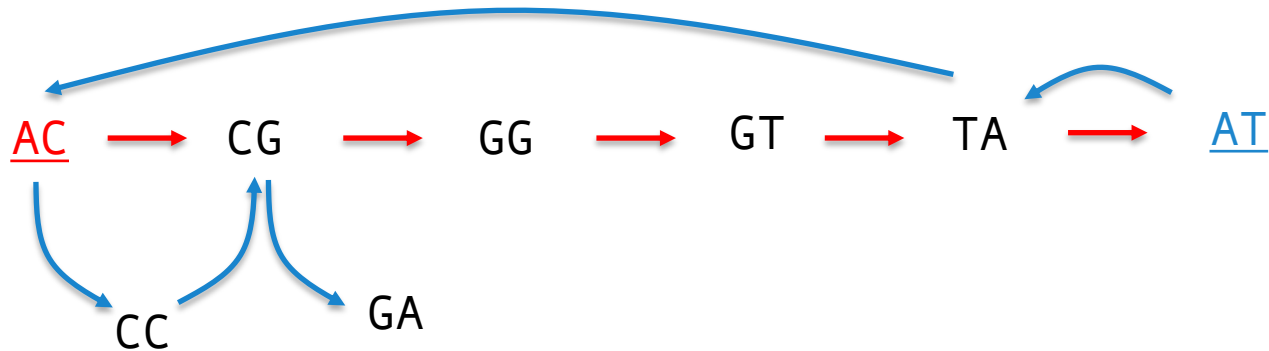CGA

# Graph traversal



Eulerian path: visit all nodes using each edge once

Starting at AC

Is this an example of a Eulerian path?

ACGGTATACCGA

# Graph traversal



Eulerian path: visit all nodes using each edge once
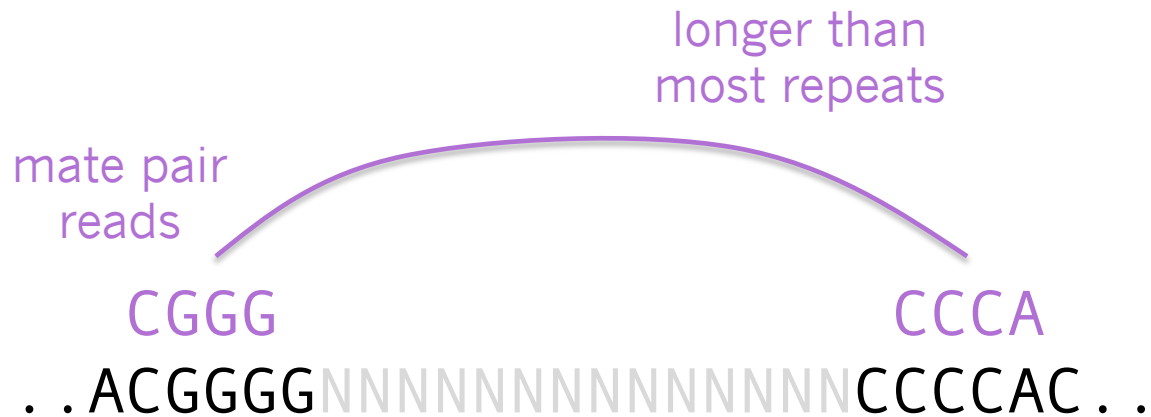
Starting at AT

How about these?

ATACCGA
or
ATACCGGTAT

# Repeat regions

*reads*                    ..ACGGGGTATATA

                                    TATATATATA

                                              TATATACCCCAC..

...ACGGGGTATATA

TATATATATATA          ????

TATATACCCCAC...                         TATATACCCCAC..

                                    TATATATATA

                              ..ACGGGGTATATA

                    ????

*contig #1*                    *contig #2*

..ACGGGGTATATA          TATATACCCCAC..
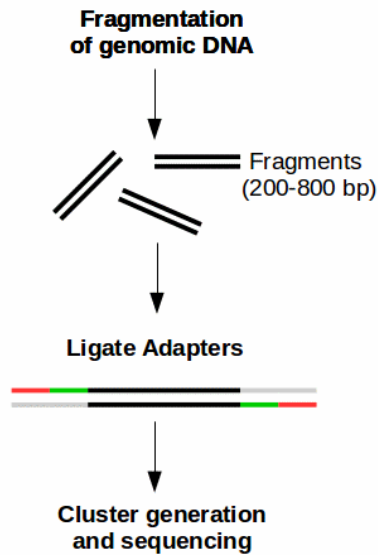
contigs form where assembly is ambiguous

# Scaffolds from contigs



Mate pair reads establish order and estimated
distance between pairs of contigs

# Scaffolds from contigs
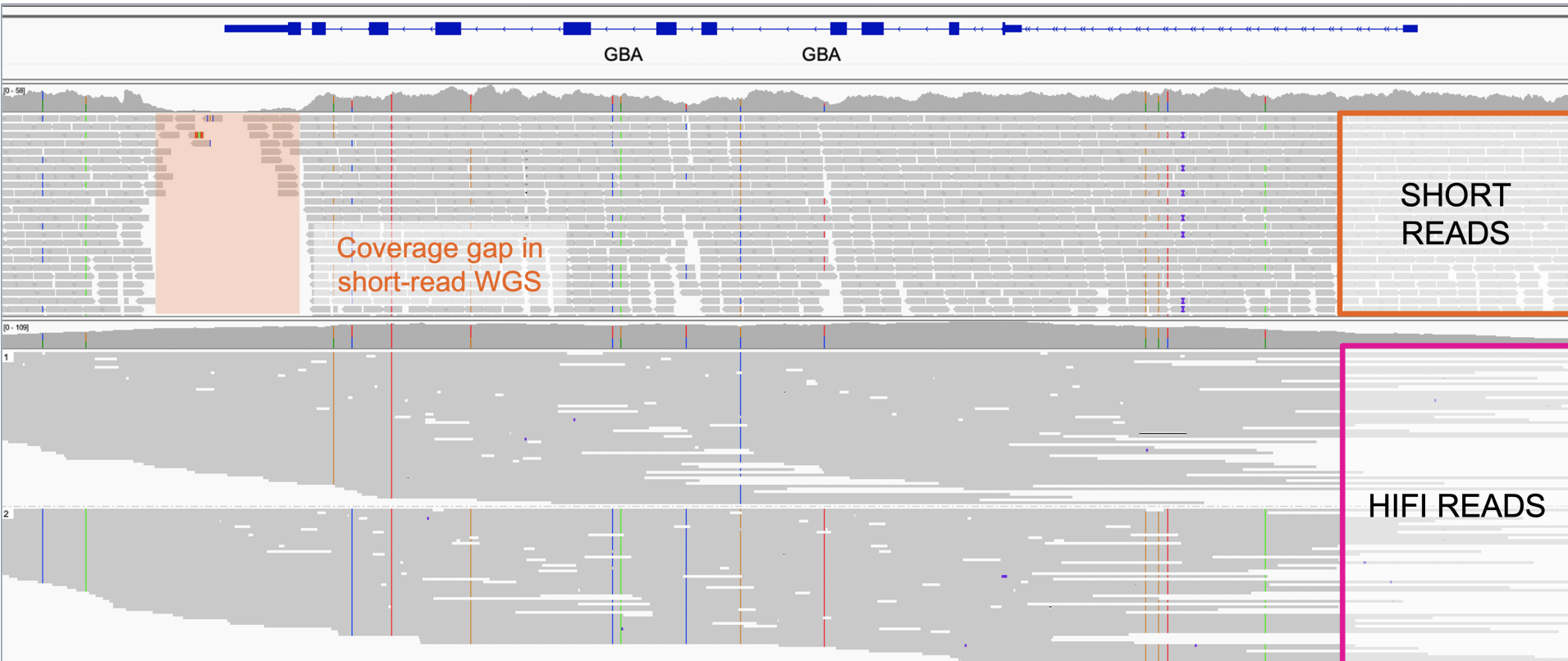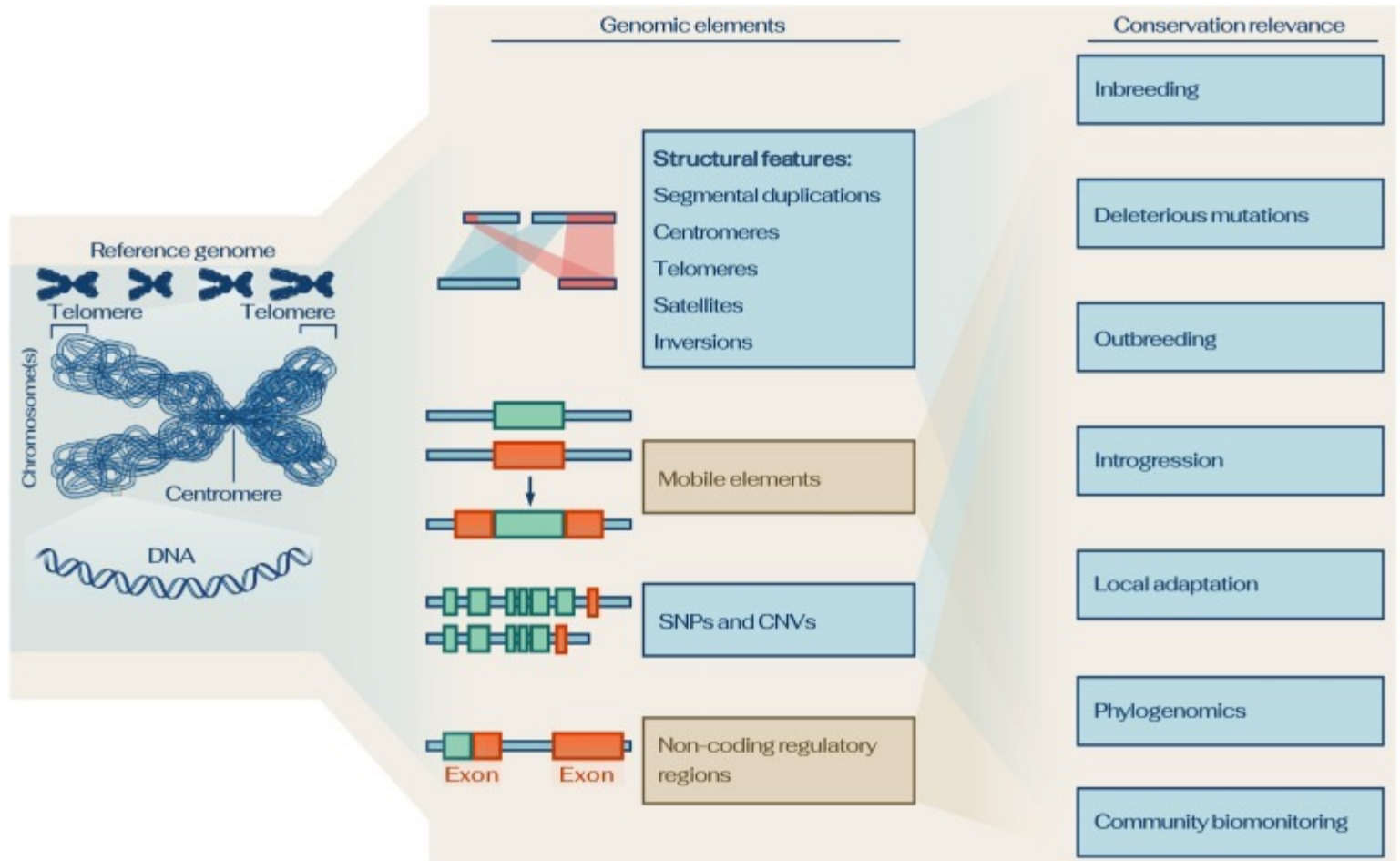
# What influences number of contigs?

- genome size
- repetitiveness of genome
- number of reads
- read length

# Reference Genomes



Trends in Ecology & Evolution

https://doi.org/10.1016/j.tree.2021.11.008

# Short read workflow

Lab focuses
on these steps

1. Assess quality of raw reads
2. Trim raw reads based on quality
3. Assemble trimmed reads into contigs
4. Assess quality of contigs
5. Scaffold contigs into genome
6. Assess/annotate genome

# Overview for Lab 19