

# Comparative genomics

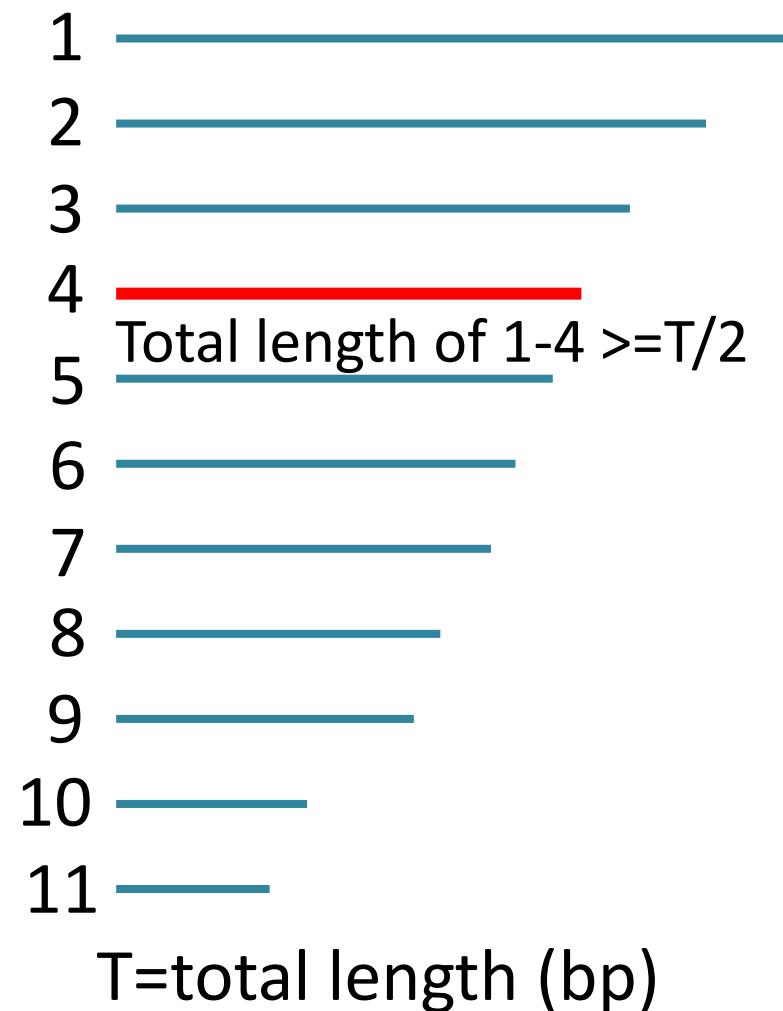
Bioinformatics Applications (PLPTH813)

Sanzhen Liu

4/11/2019

## Assembly statistics – N50

- N50: A statistic used for assessing the contiguity of a genome assembly.
- The contigs in an assembly are **sorted** by size and added, starting with the largest. The contig N50 is **the size of the contig that makes the total greater than or equal to 50% of the total contig size.**



# Long-read assemblies

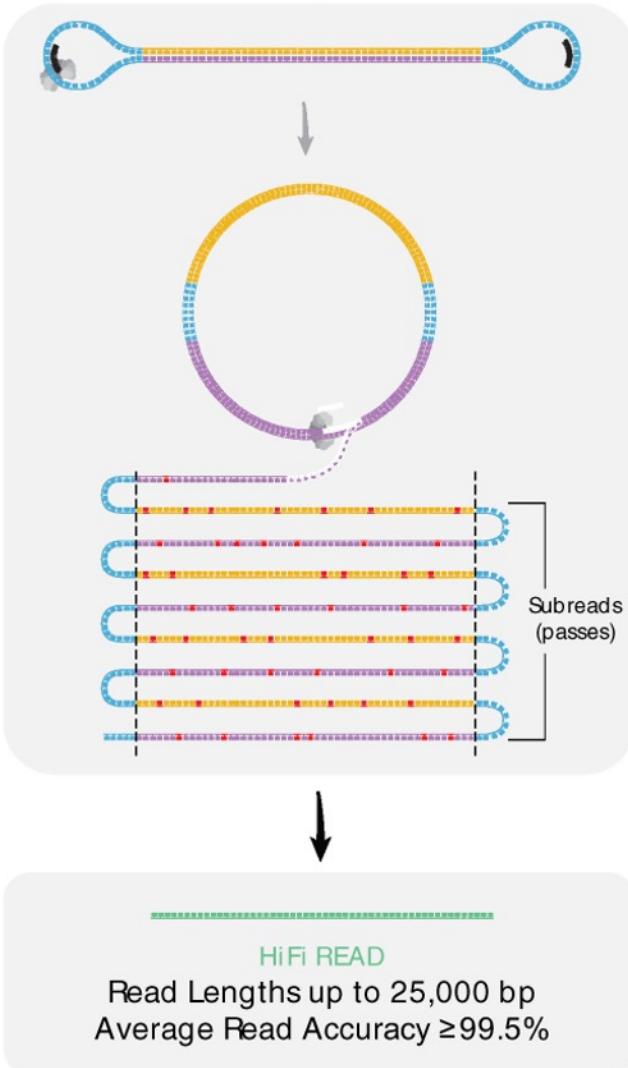
Software	Description
Canu	K-mer-based overlap computation
FALCON	Assembles phased diploid genomes
Flye	Uses A-Brujin graph

# Contig polishing

Software	Description
Arrow	Assembly error correction with <b>PacBio</b> long-read alignments
Nanopolish	Assembly error correction using voltage (raw) <b>Nanopore</b> data

Pilon	Assembly error correction using Illumina data
-------	---

# Assembly meets long HiFi PacBio reads

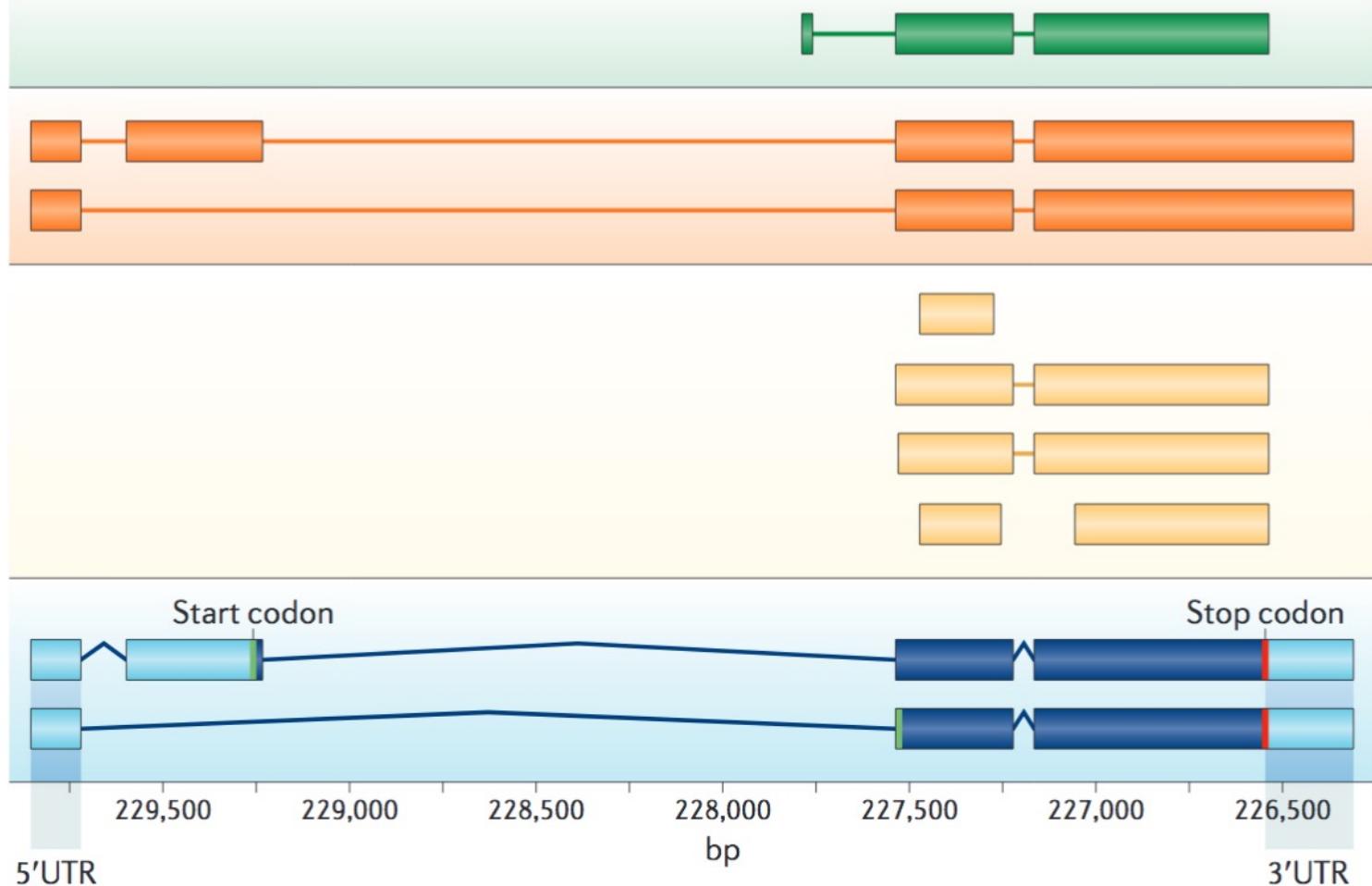


hifiasm:  
High-contiguity assembly  
fast

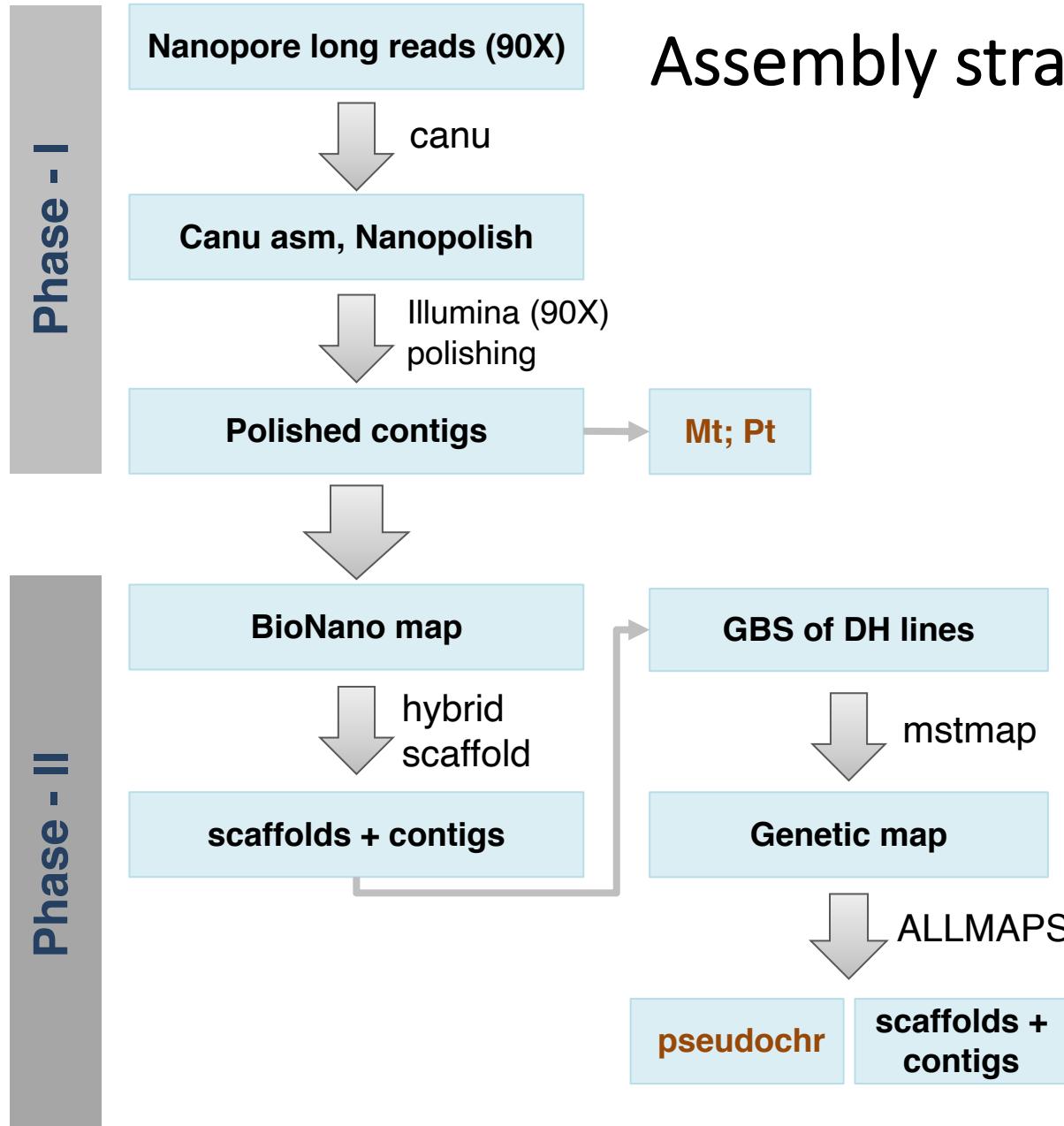
<https://github.com/chhylp123/hifiasm>

# Gene annotation: *ab initio* prediction + evidence

Gene prediction  
(SNAP)



# Assembly strategy



# Outline

- Introduction of comparative genomics and structural variation
- Approaches
  1. Comparative genome hybridization
  2. Paired-end reads
  3. Read depth
  4. Whole genome assembly
- Pangenomics
- Case study: CGRD

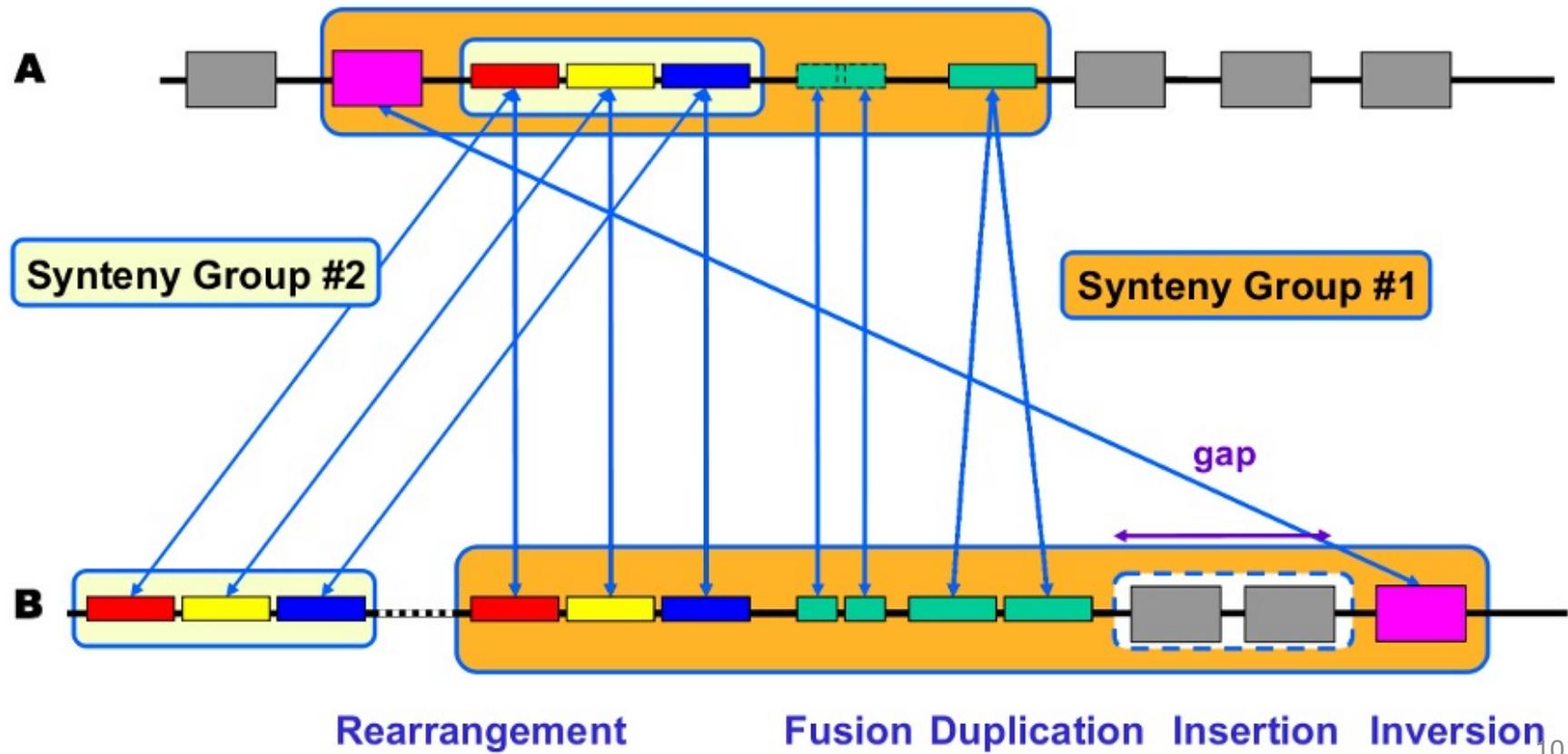
# Comparative genomics among genomes

## Conservation and variability

- Homologs: Genomic regions derived from a common ancestral gene.
- Orthologs: Homologs from the divergence of lineages.
- Paralogs: Homologs derived from their duplication within a lineage.
- Homeologs: The subset of paralogs created by WGD.  
(synonyms: ohnolog; syntenic paralog)

# Synteny

**Synteny** is usually referred to as the conservation of blocks of order within two sets of chromosomes that are being compared with each other. Syntenic regions are evidence by homologous genes arranged in a collinear order.



# Intra-species genome rearrangements and structural variation (SV)

Variation of *hundreds of bp* to several Mb in size

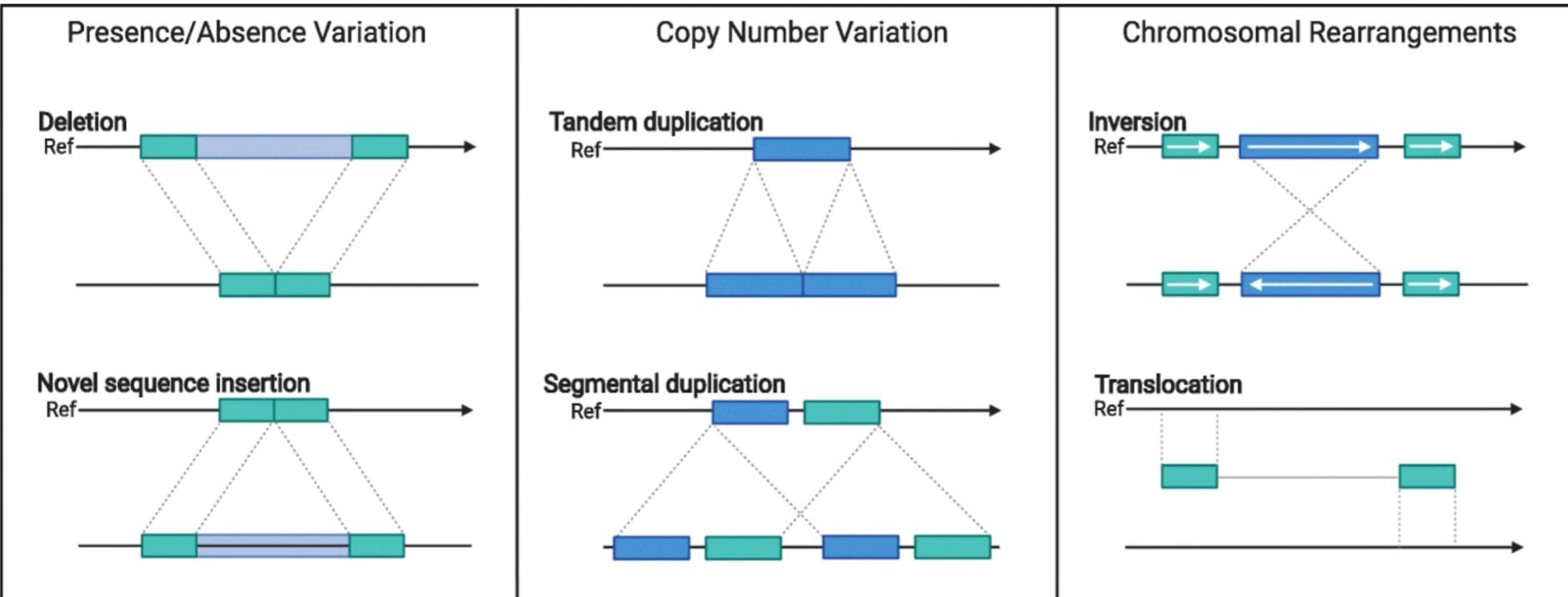
## Balanced variation

- Inversion
- Translocation

## Unbalanced variation:

- Copy number variation (CNV)  
(Presence/Absence variation, PAV)

# Structure variation types within species

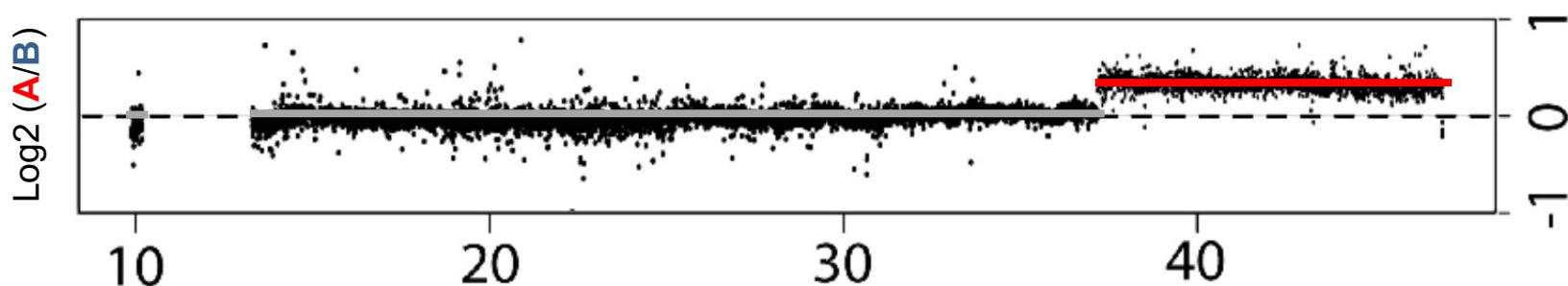
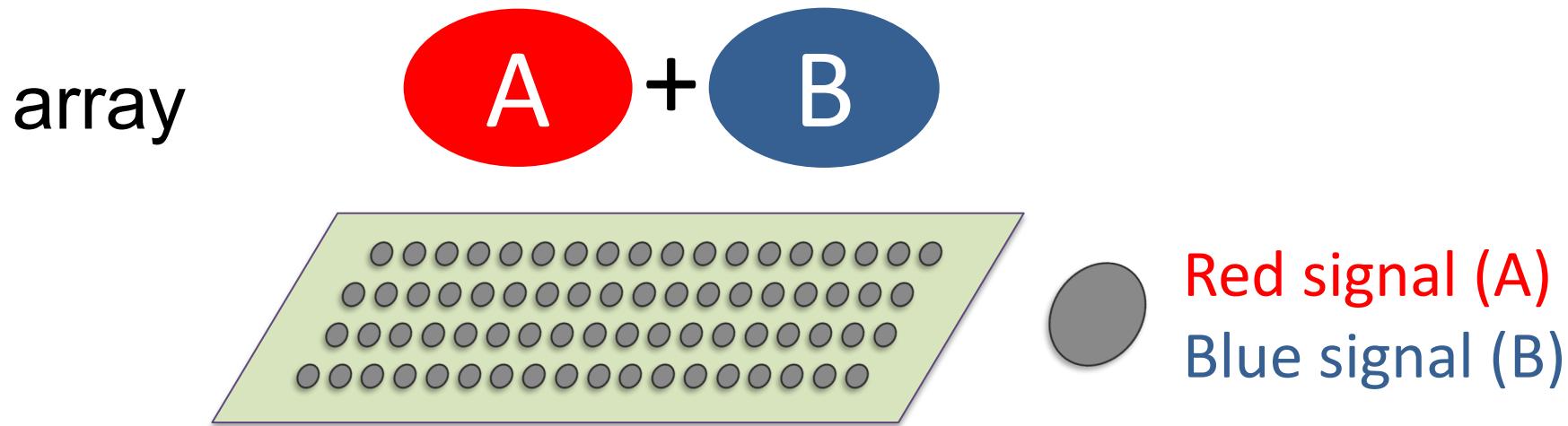


Coletta et al., Genome Biol, 2021

# Outline

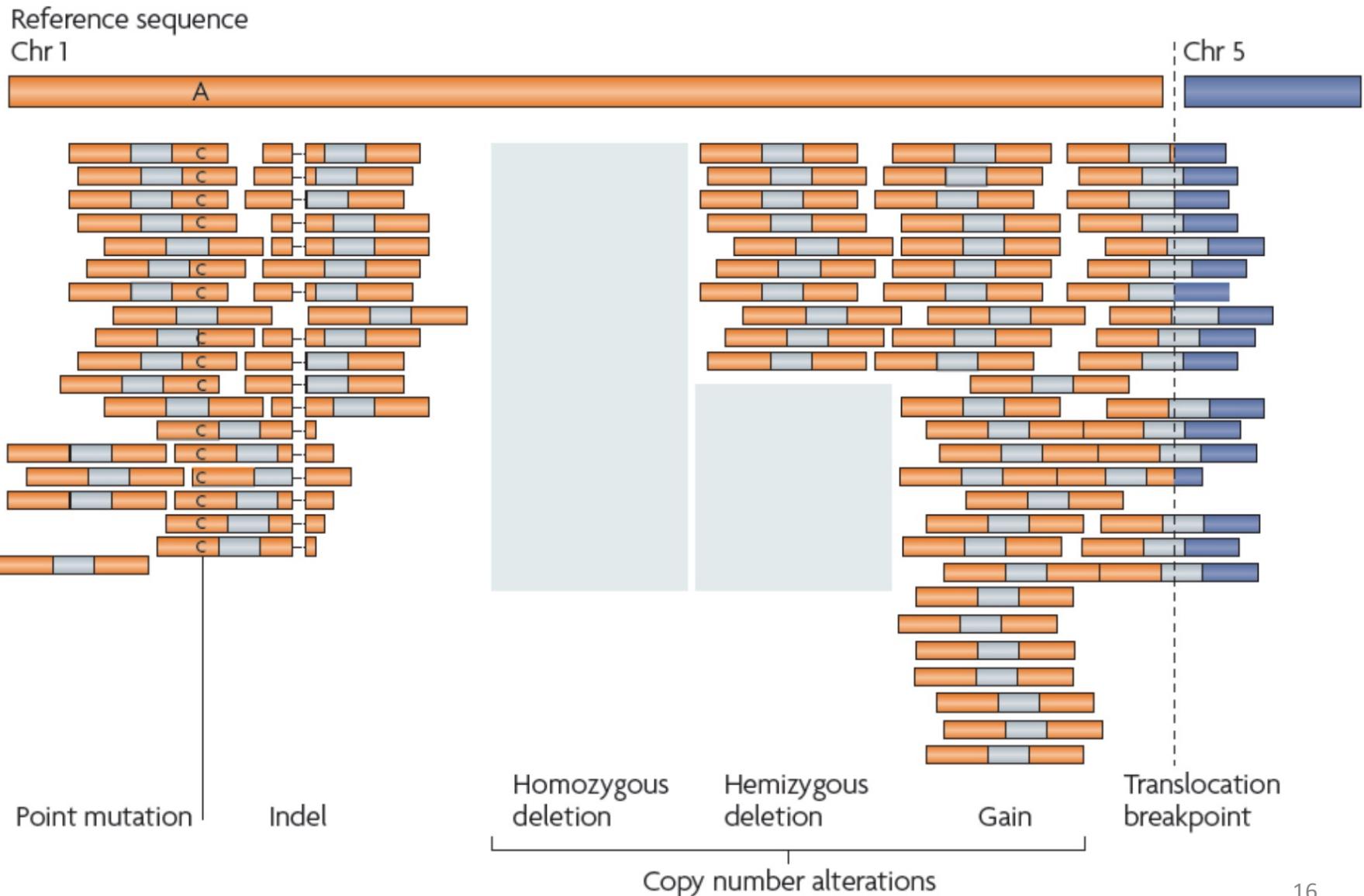
- Introduction of comparative genomics and structural variation
- Approaches
  1. Comparative genome hybridization
  2. Paired-end reads
  3. Read depth
  4. Whole genome assembly
- Pangenomics
- Case study: CGRD

# array Comparative Genomic Hybridization (aCGH)

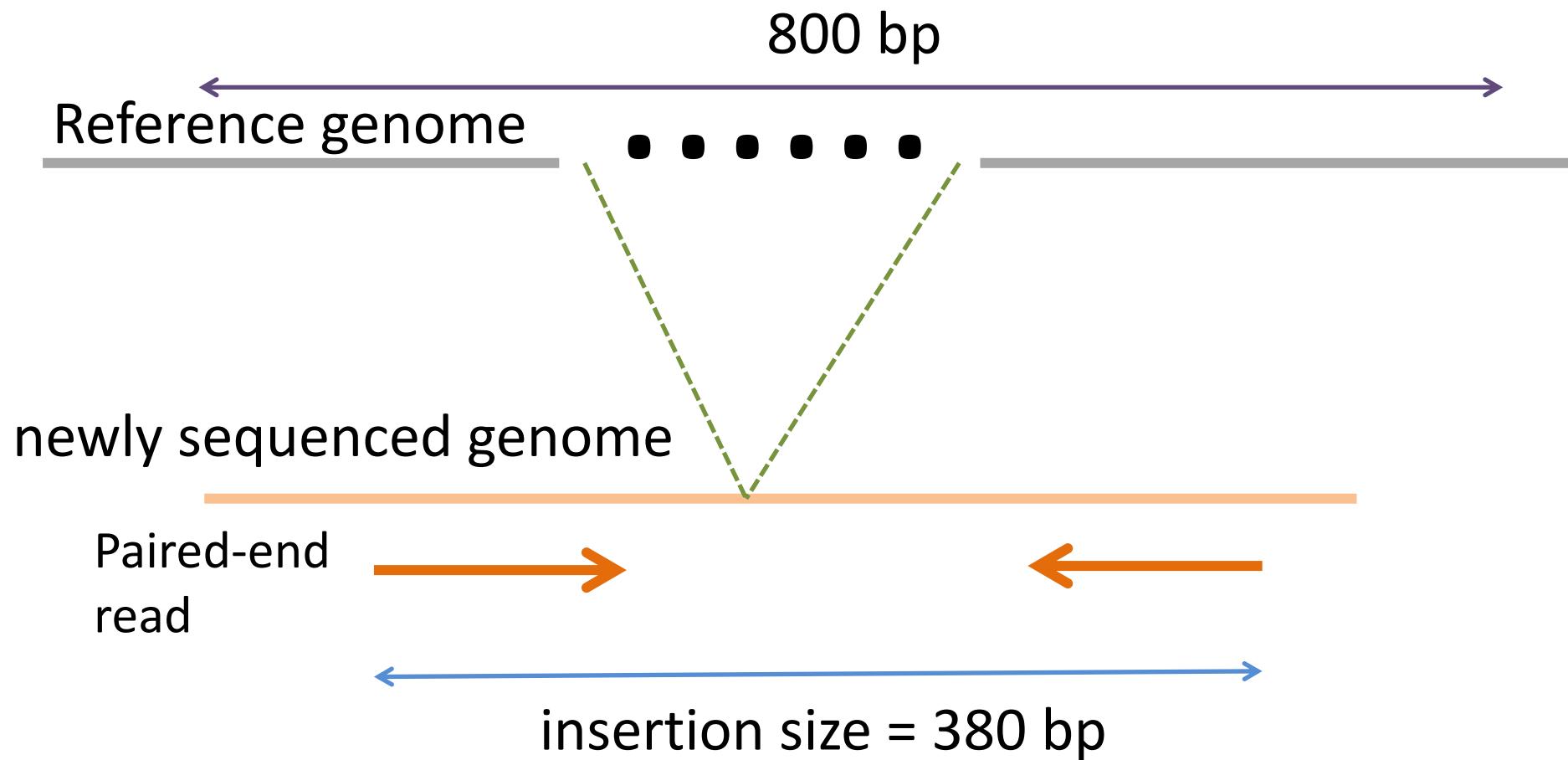


NGS provides information for the discovery of variants, including structural variation

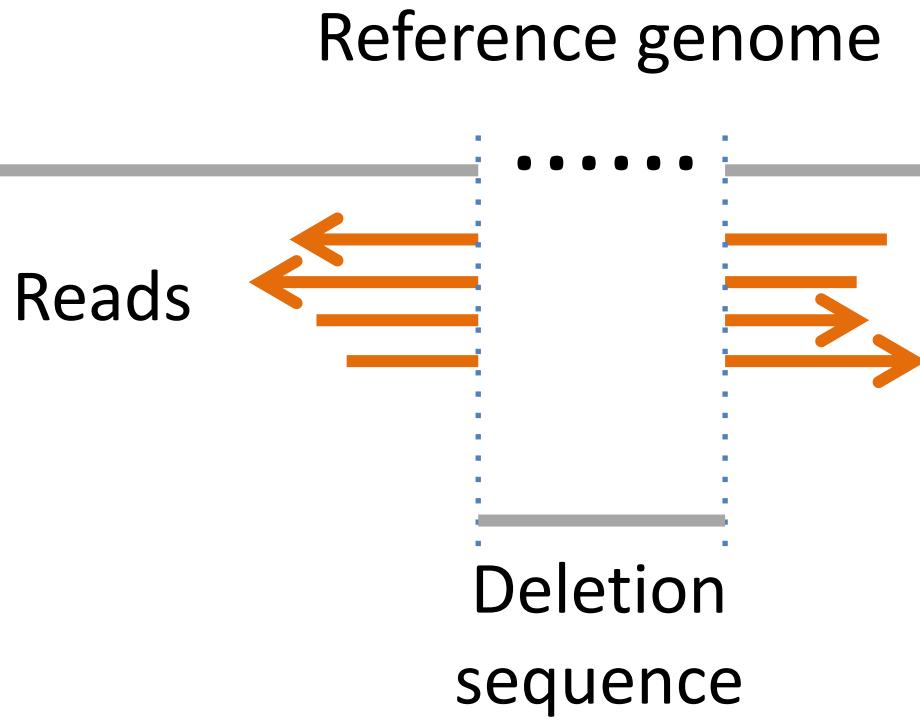
# Variants in sequencing reads



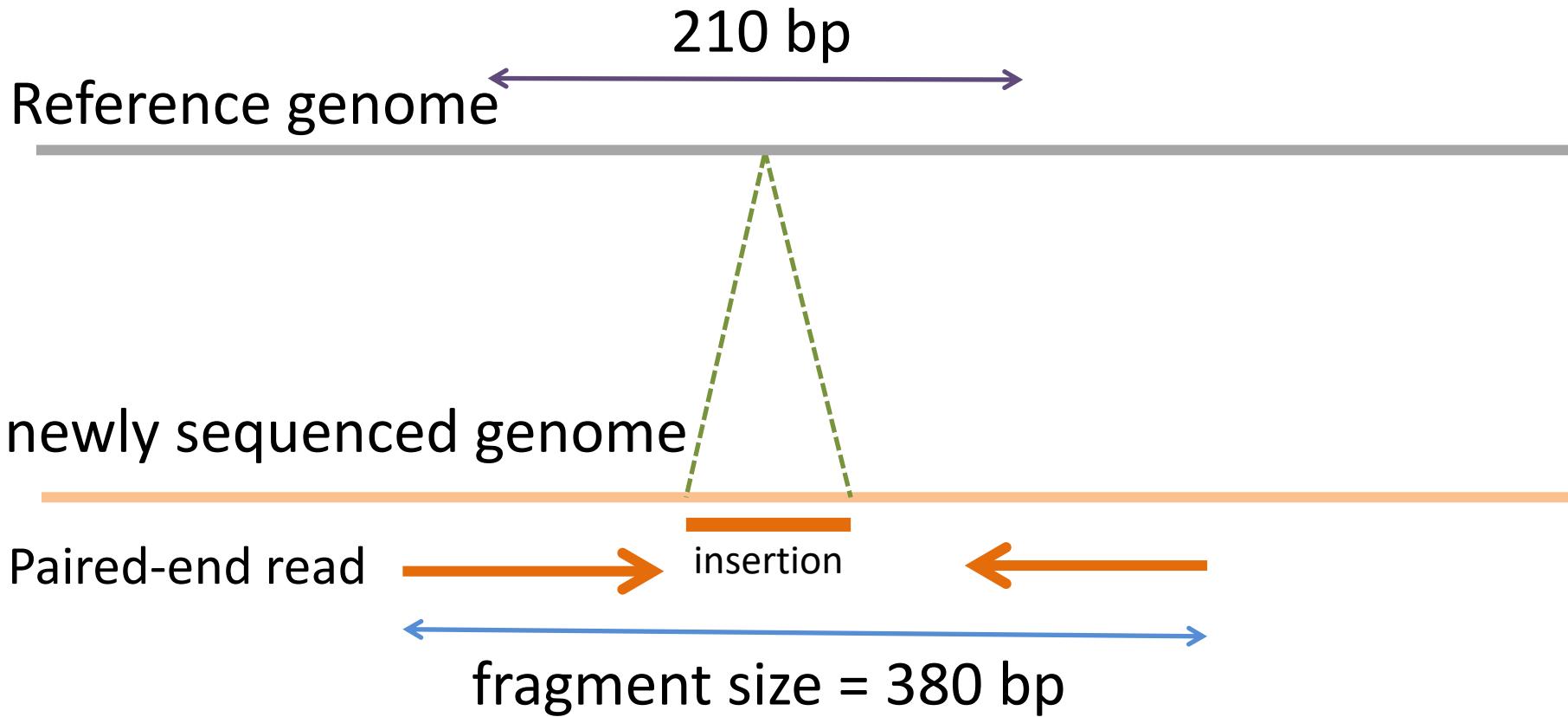
# Paired-end reads to find "deletion" relative to the reference



# Split reads to find "exact deletion sequence"

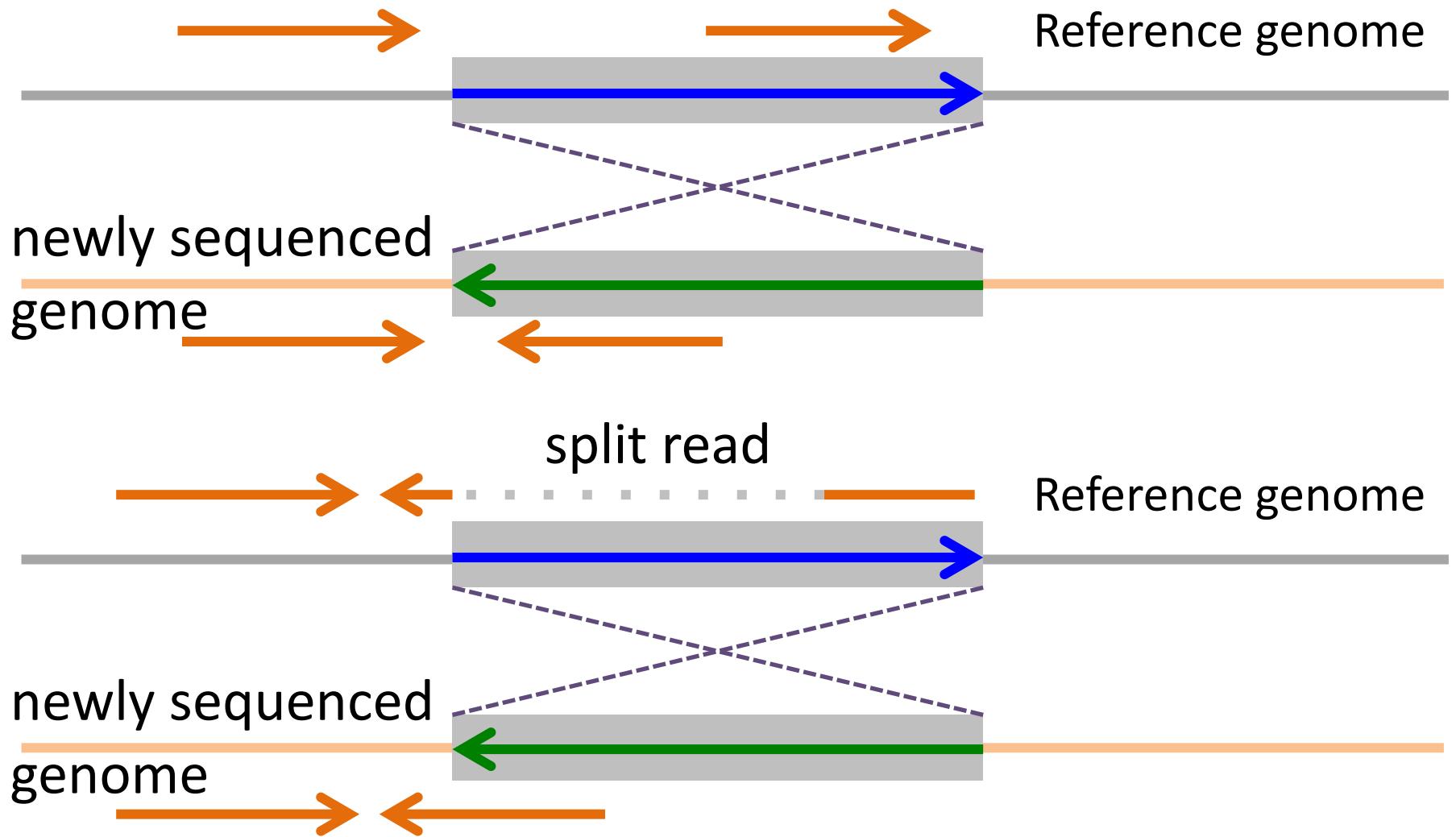


# Paired-end reads to find "insertion" relative to the reference

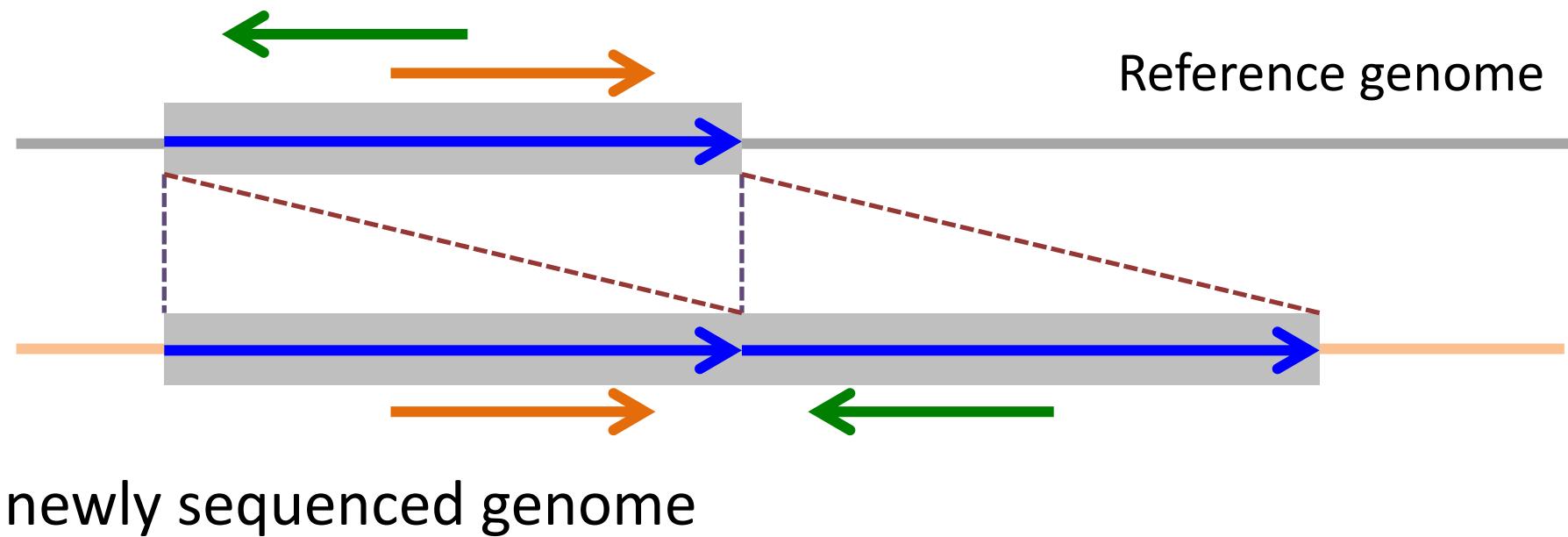


The size of insertions that can be identified by PE reads is determined by fragment sizes and read lengths.

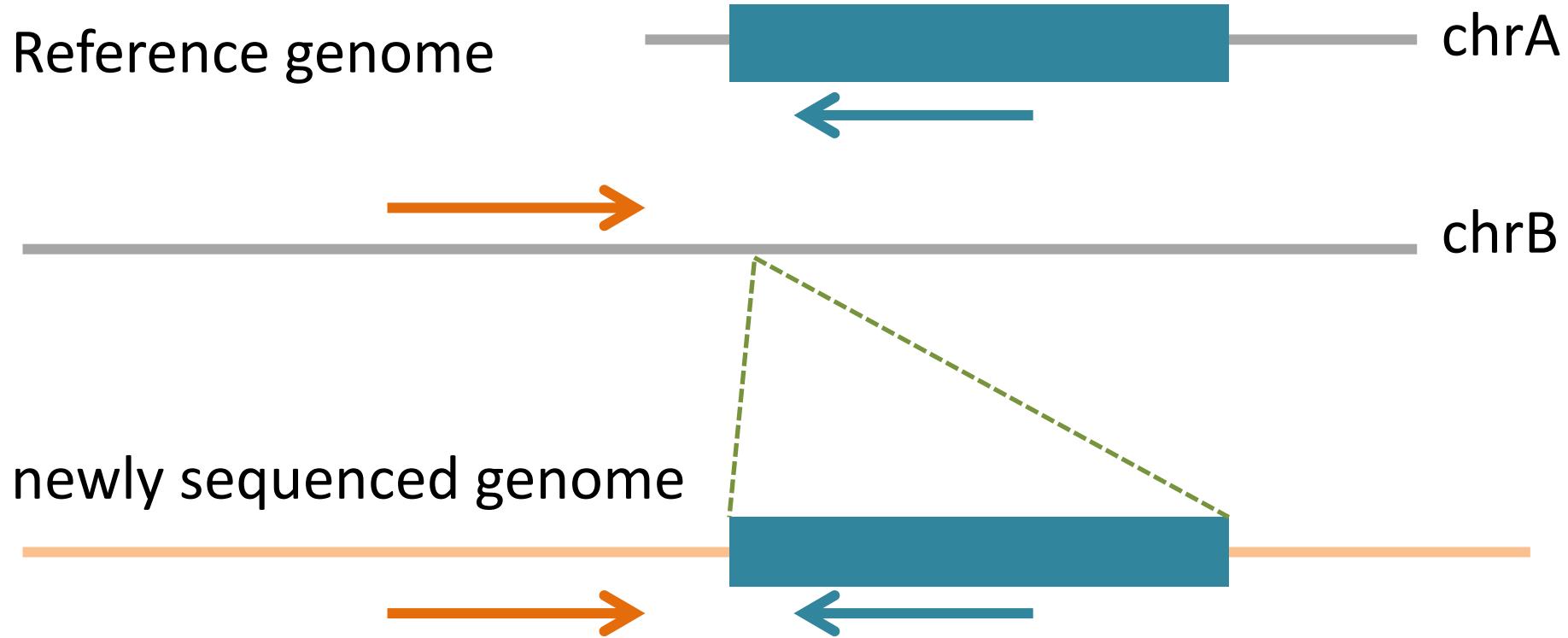
# inversion



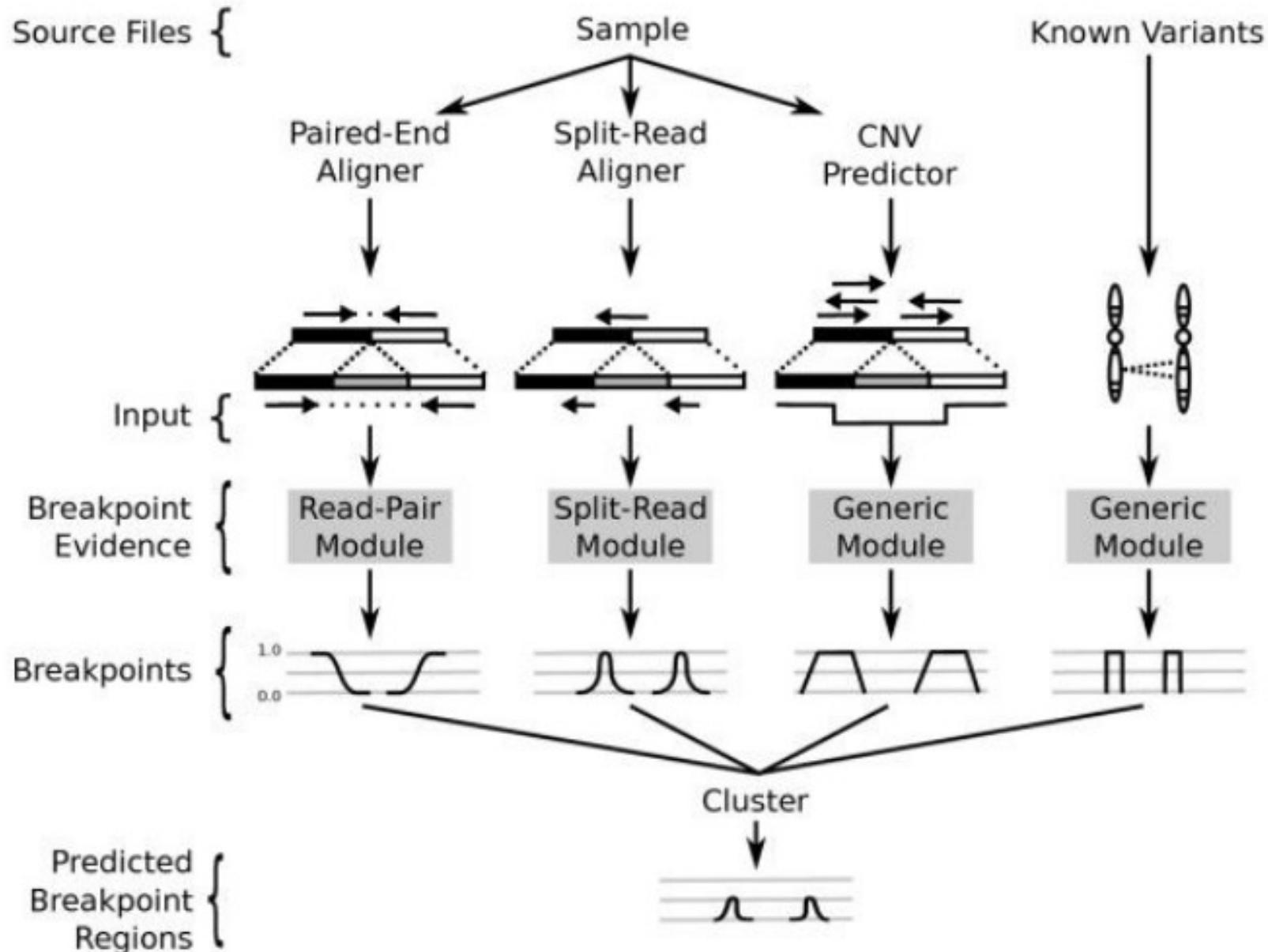
# Tandem duplication



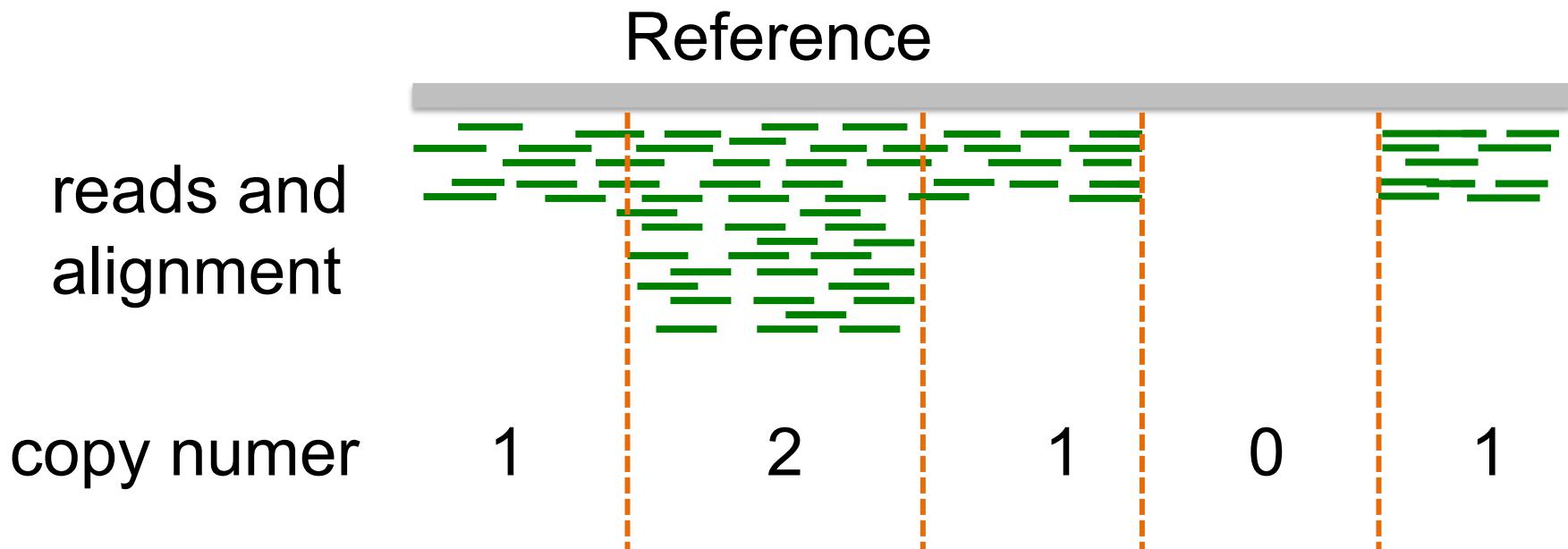
# Translocation



# LUMPY: an integrative framework for SV discovery

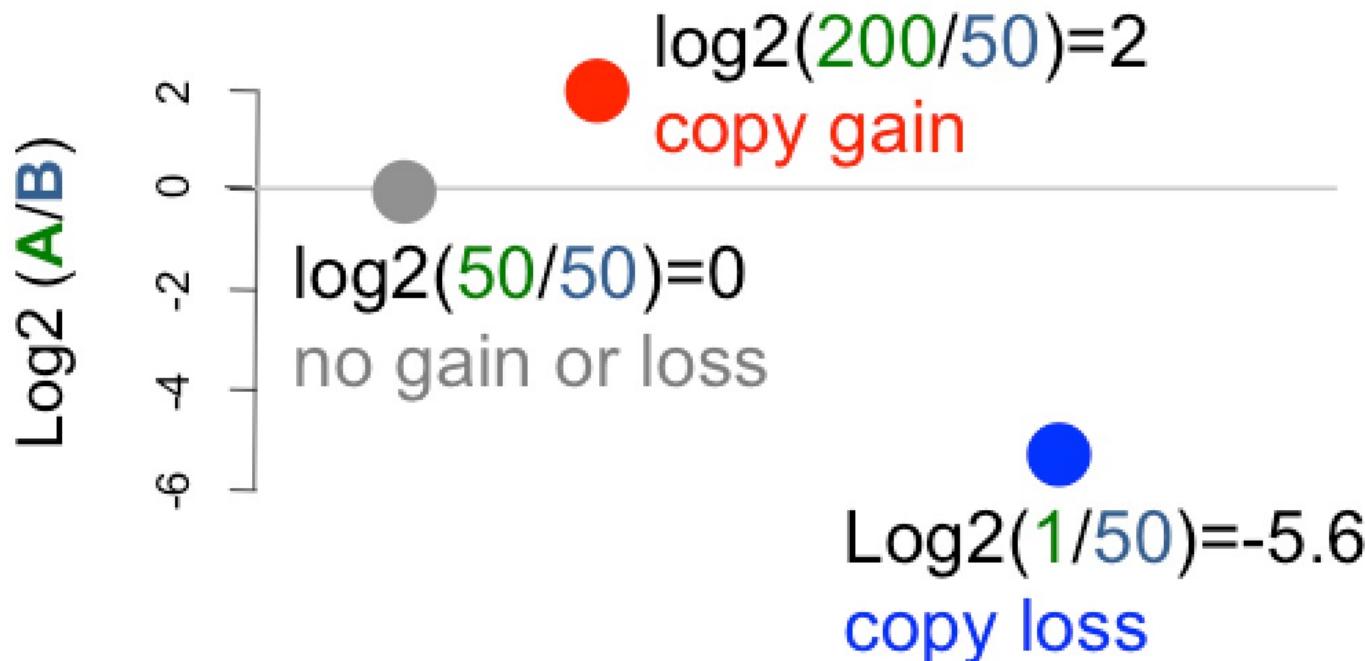
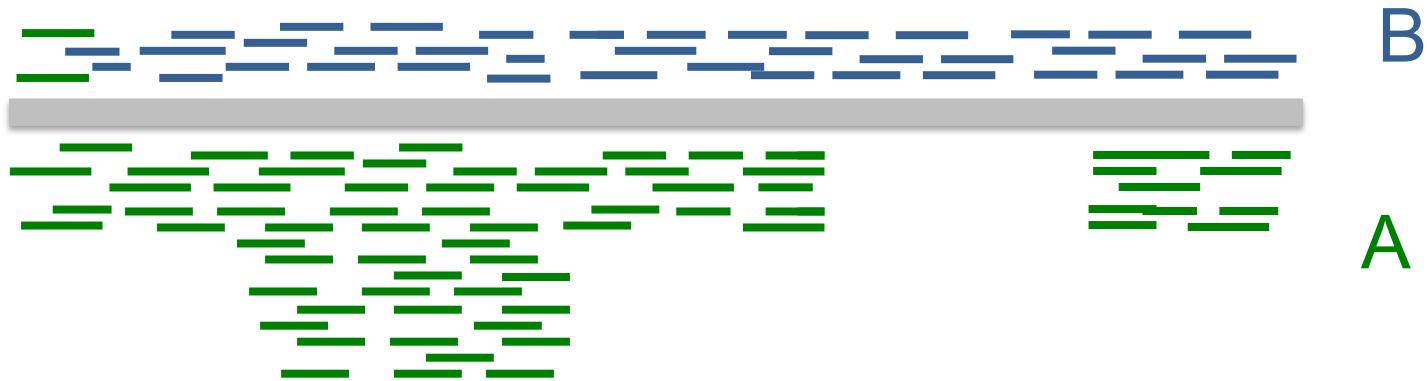


# Read depths



1. PCR bias (GC content, complicated structure)
2. Alignment bias (e.g., highly repetitive regions)
3. Nucleotide polymorphisms

# Read depths

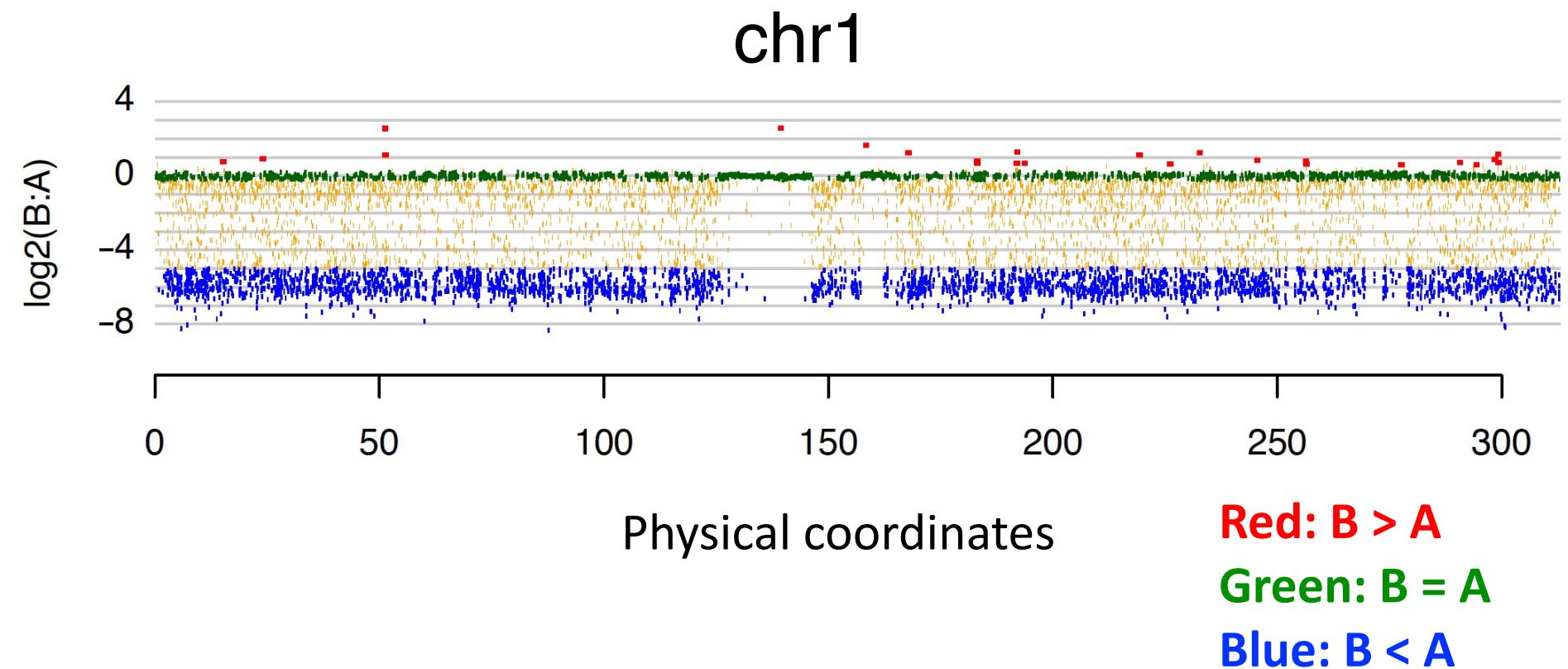


# Comparative Genomic Read Depth (CGRD)

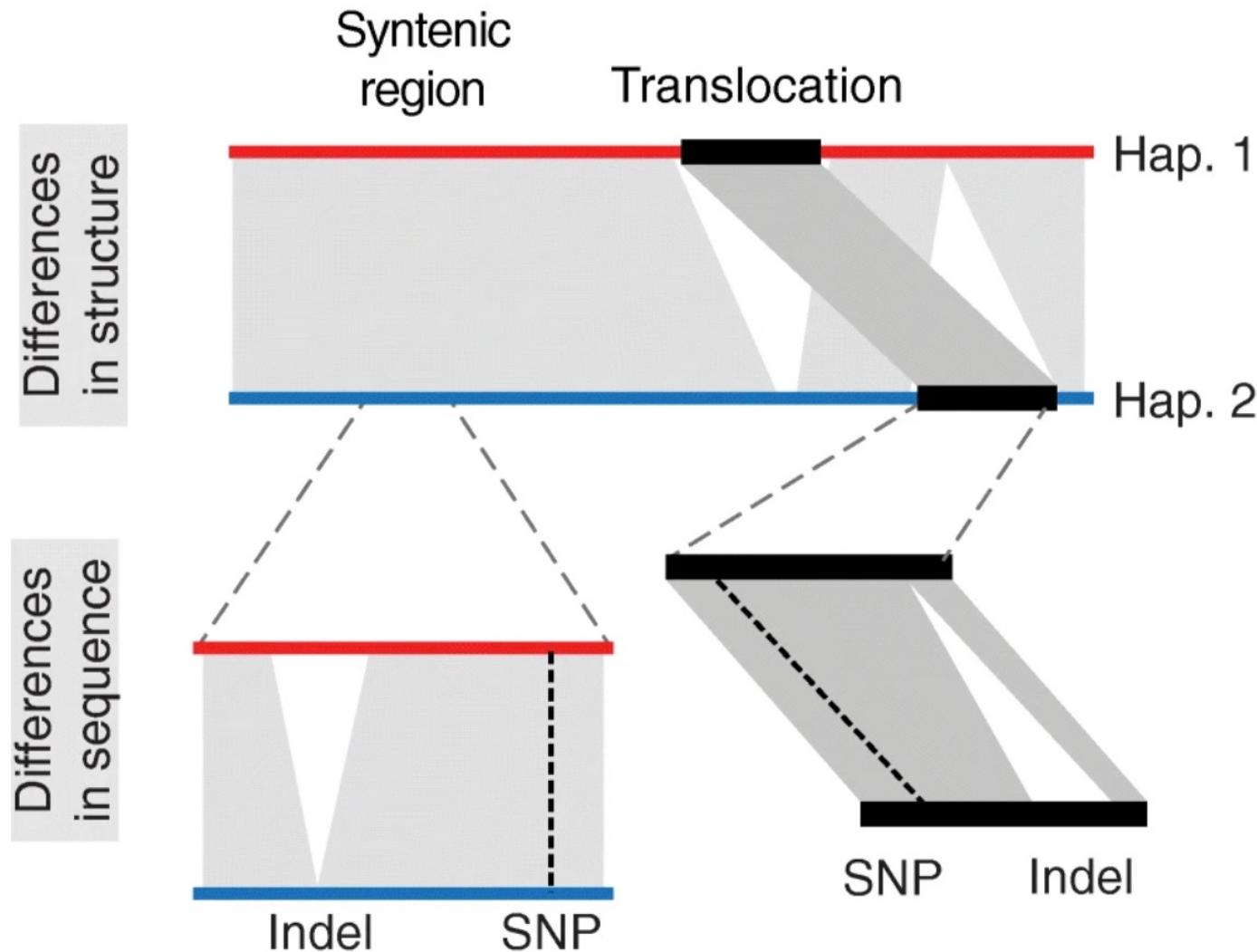
**CGRD** is a pipeline to compare sequencing read depths from two samples along a reference genome.

1. define effective genomic bins each of which harbors certain non-repetitive sequences
2. align reads and count read depths per bin for both samples
3. combine neighbor bins with similar fold changes in read depth between the two samples  
(segmentation)

# Comparative Genomic Read Depths (CGRD)



# SyRI: finding genomic variation using whole-genome assemblies



# SyRI - algorithm

**Input:**  
whole genome  
alignment

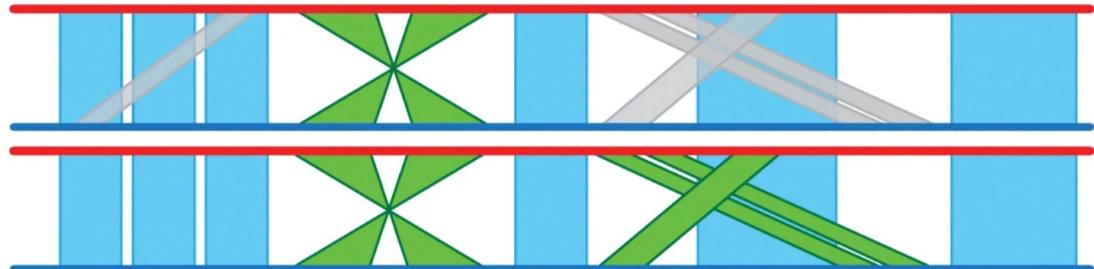
**Aligner:**  
NUCmer  
minimap2

**Step 1:**  
annotate  
syntenic regions



**Step 2:**  
annotate  
structural  
rearrangements

**2a:** annotate inversions



**2b:** annotate transpositions,  
duplications, and remove  
redundant alignments

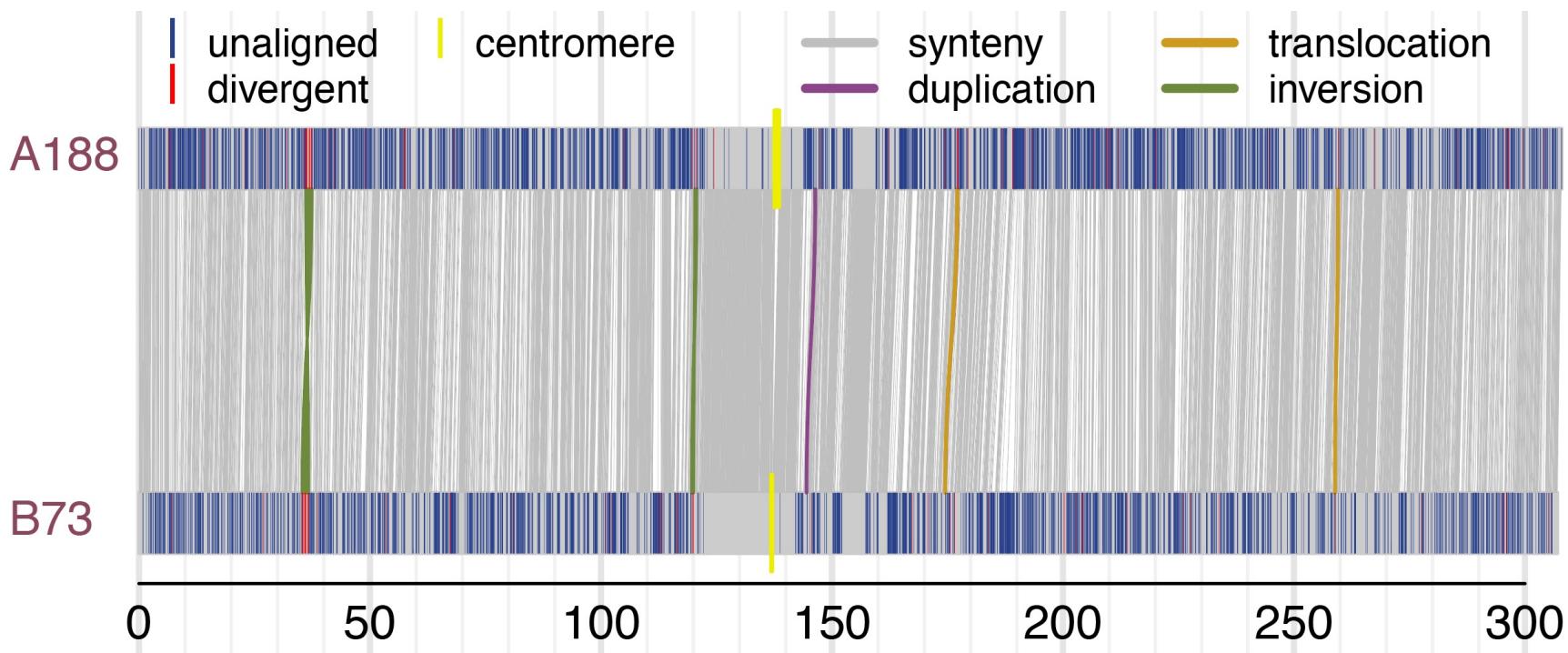


**2c:** annotate translocations and  
duplications between chromosomes



# Structure variation with SyRI

- SyRI is a tool using the **genome assemblies** to discover genomic rearrangement and local sequence differences (Goel et al., 2019)

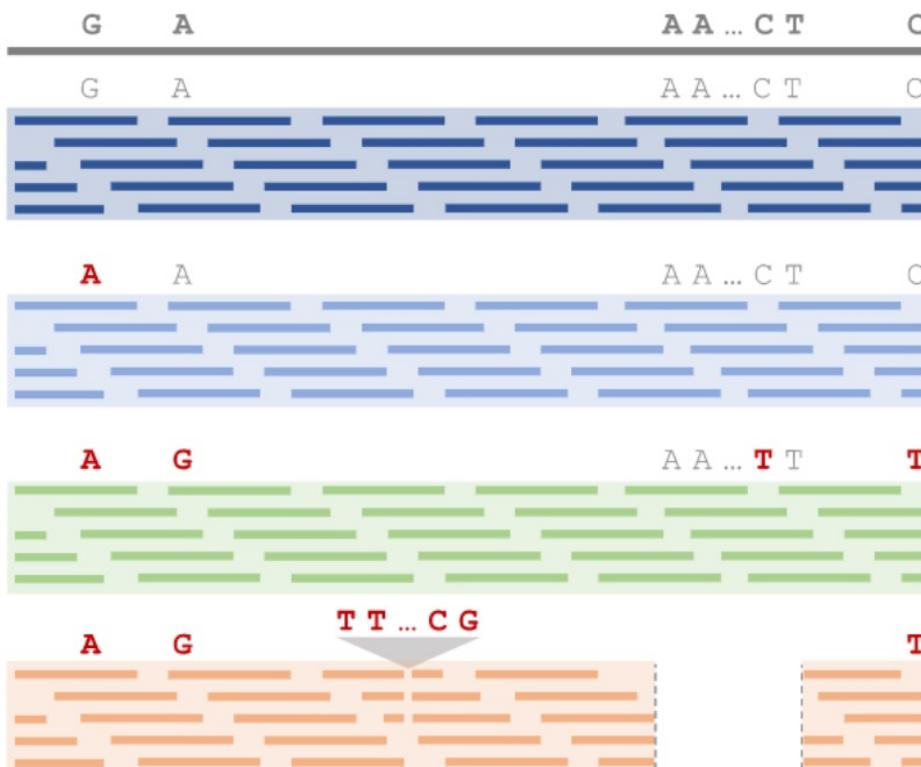


# Outline

- Introduction of comparative genomics and structural variation
- Approaches
  1. Comparative genome hybridization
  2. Paired-end reads
  3. Read depth
  4. Whole genome assembly
- **Pangenomics**
- Case study: CGRD

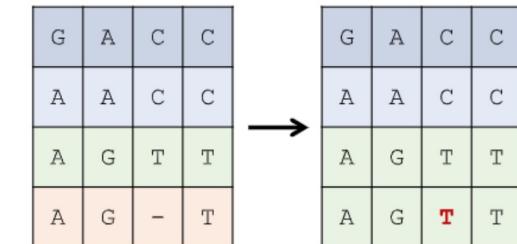
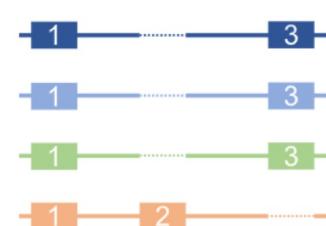
# Issues associated with single reference genomes

## Linear reference genome



Region 2 is not present  
in the linear reference  
genome

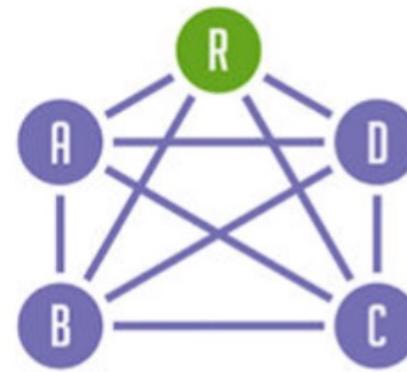
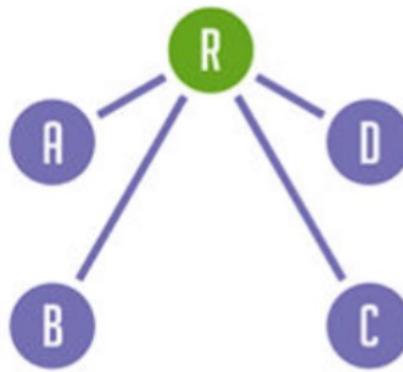
Wrong SNP imputation of  
region 3 leads to wrong  
haplotype



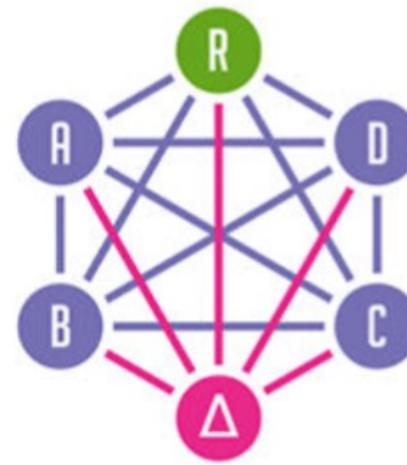
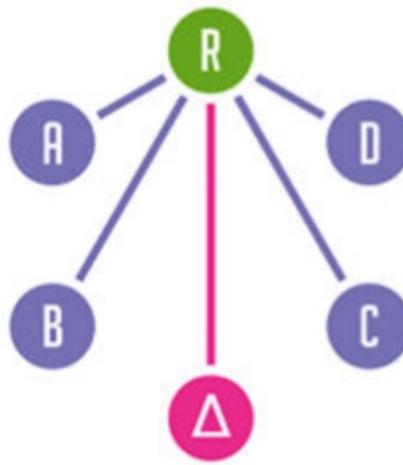
Coletta et al., Genome Biol, 2021

# Pangenomics

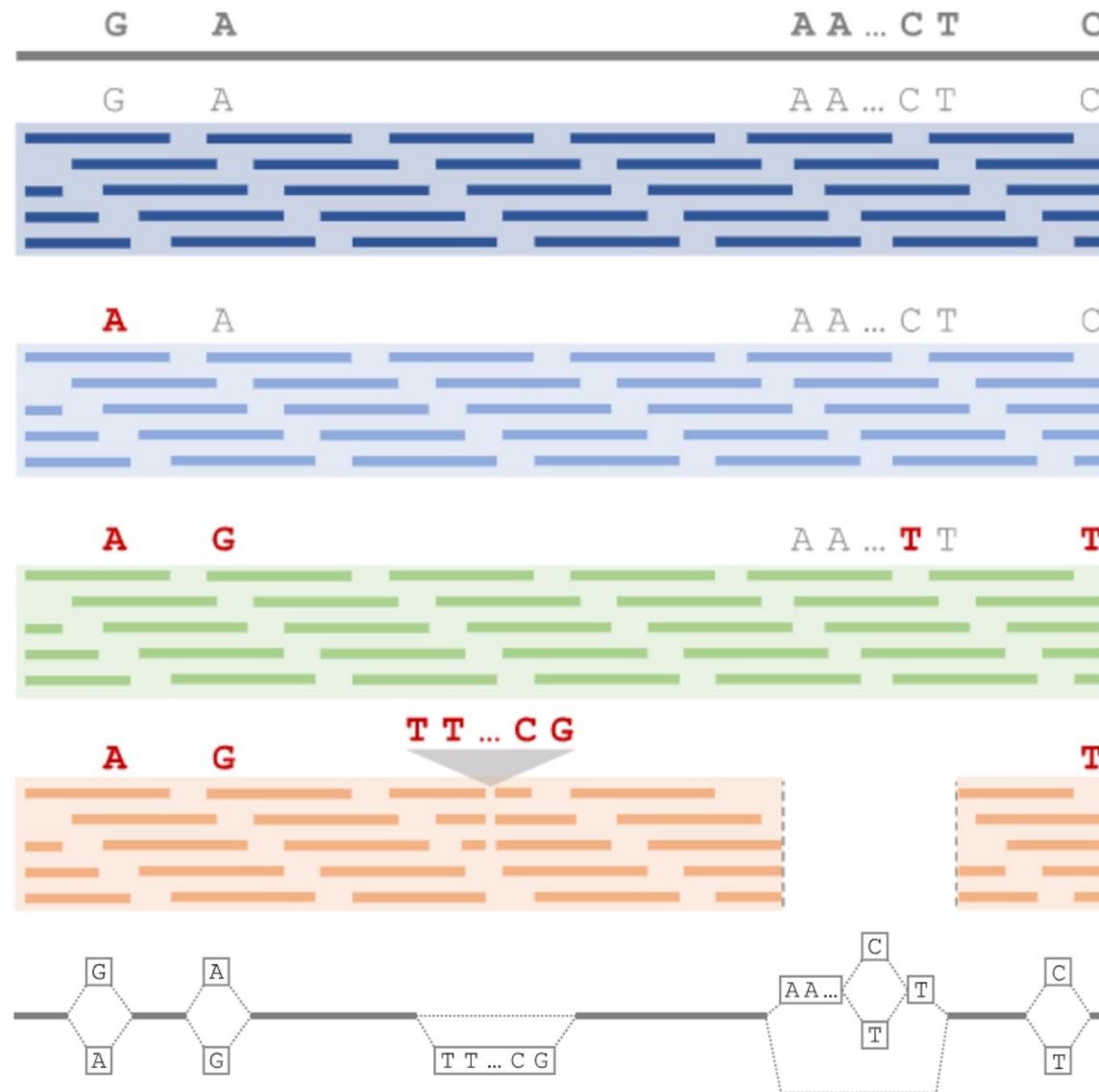
Reference model



Extending the model



# Construct a pangenome to facilitate genotyping

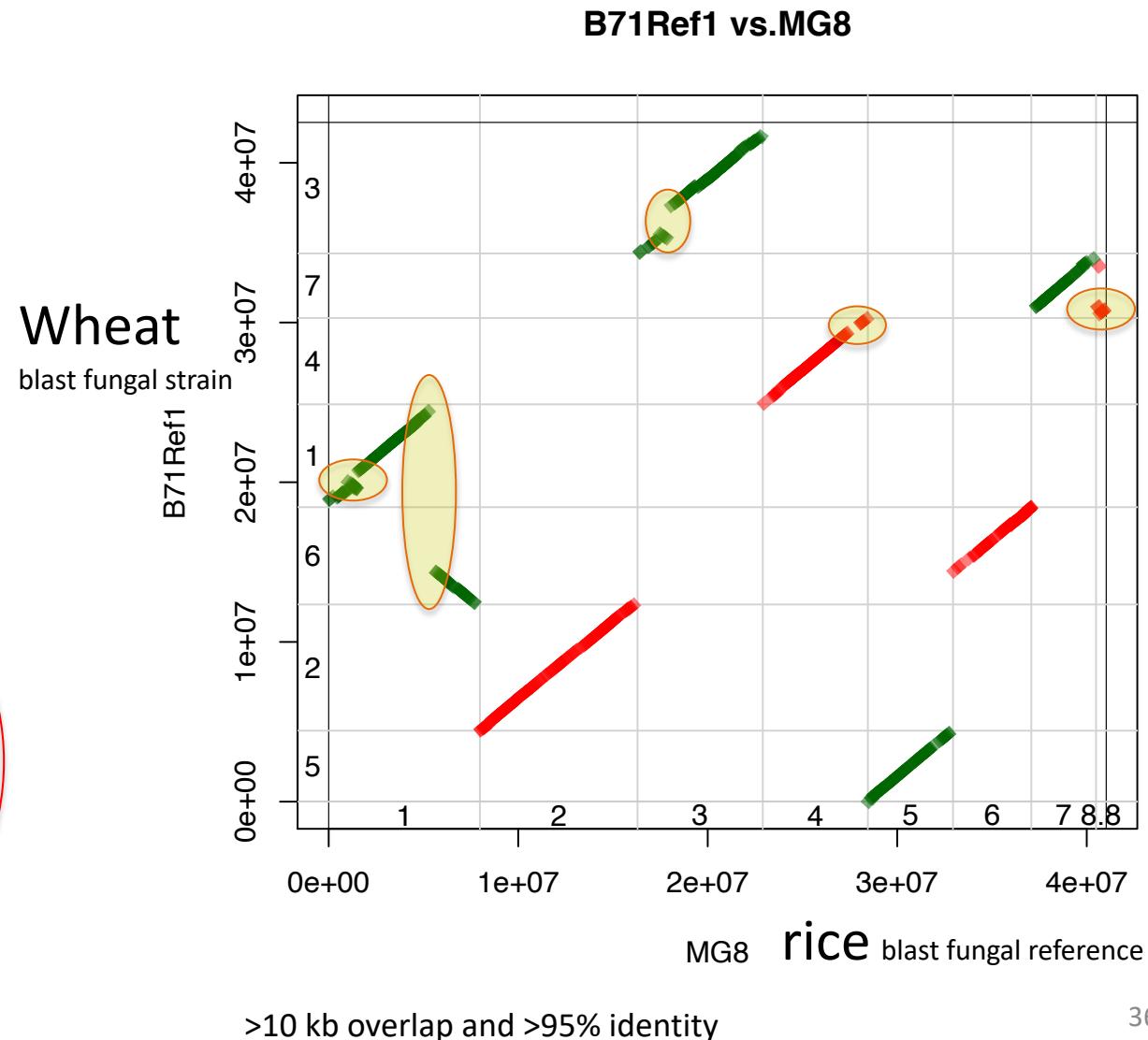


# Outline

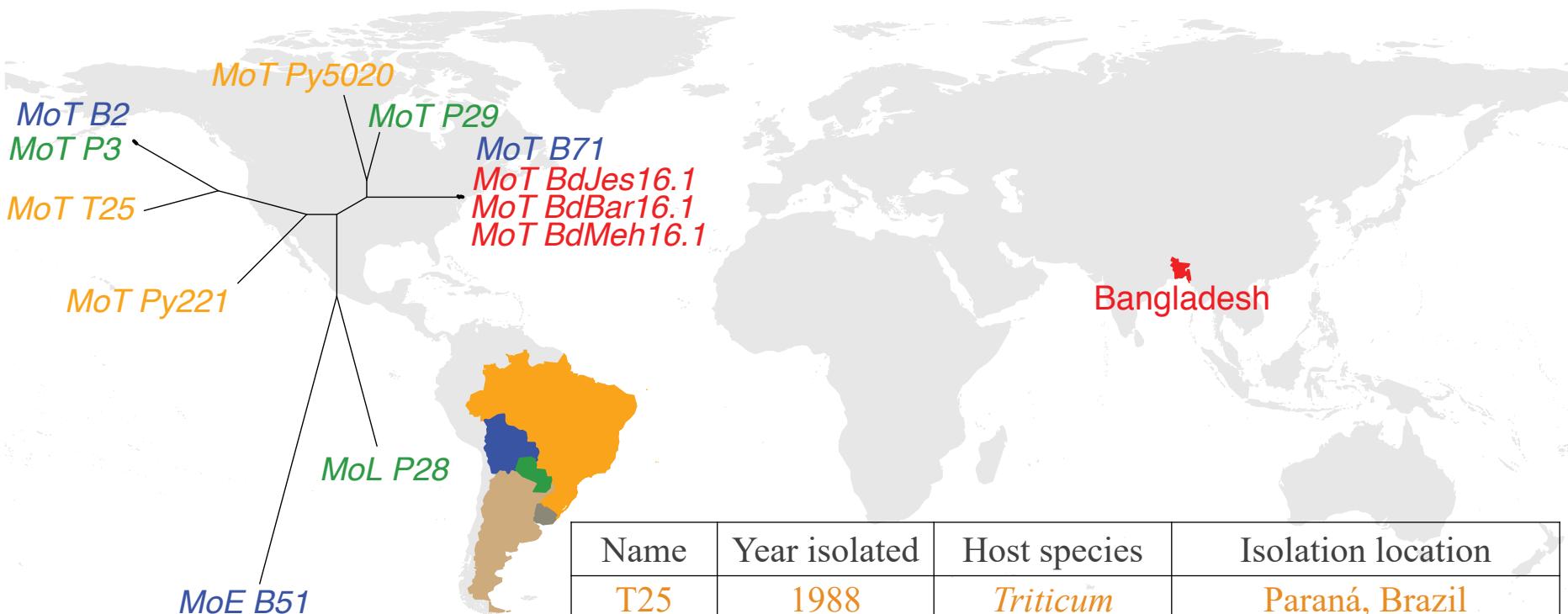
- Introduction of comparative genomics and structural variation
- Approaches
  1. Comparative genome hybridization
  2. Paired-end reads
  3. Read depth
  4. Whole genome assembly
- Pangenomics
- Case study: CGRD

# MoT B71 final assembly vs. MG8 (70-15, rice strain)

Chr	Length (bp)
chr1	6,442,091
chr2	7,902,655
chr3	8,206,304
chr4	5,402,116
chr5	4,442,877
chr6	6,090,985
chr7	4,042,640
scaf1	941,816
scaf2	739,928
scaf3	104,670
scaf4	74,396
scaf5	69,131



# Illumina sequencing of additional eight strains

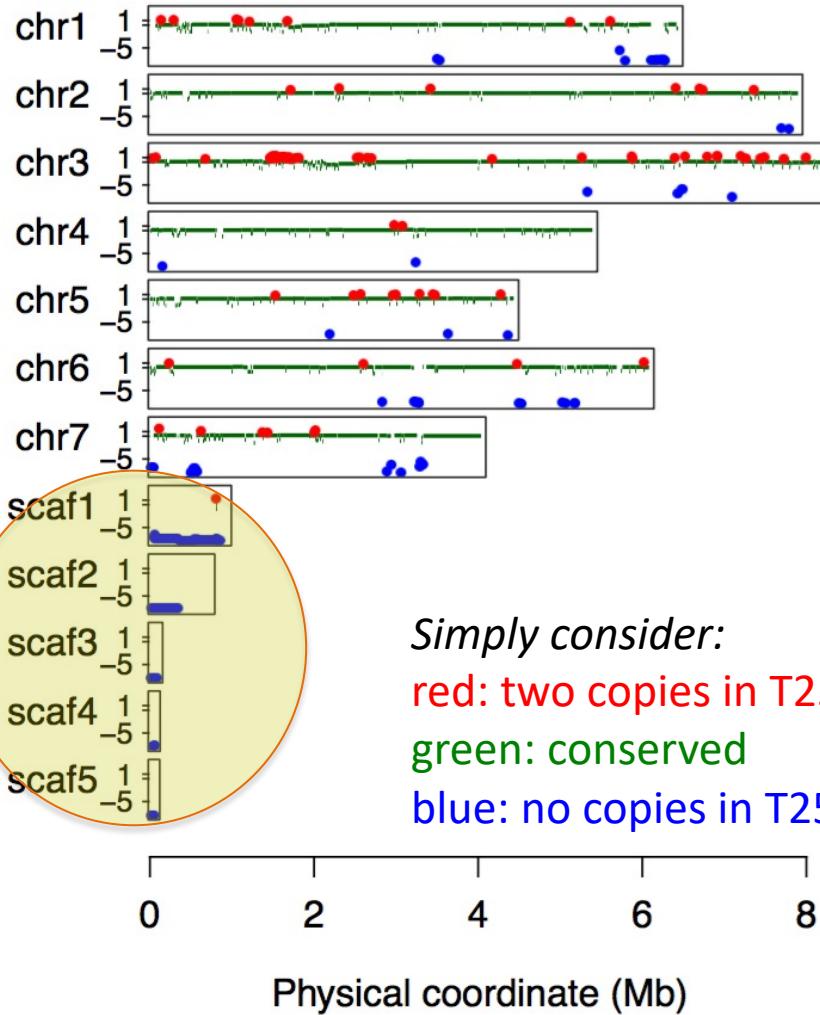


Name	Year isolated	Host species	Isolation location
T25	1988	<i>Triticum</i>	Paraná, Brazil
Py5020	2005	<i>Triticum</i>	Paraná, Brazil
Py22.1	2007	<i>Triticum</i>	Paraná, Brazil
B2	2011	<i>Triticum</i>	Quirusillas, Bolivia
B71	2012	<i>Triticum</i>	Okinawa, Bolivia
P3	2012	<i>Triticum</i>	Canindeyú, Paraguay
P28	2014	<i>Bromus</i>	Paraguay
P29	2014	<i>Bromus</i>	Paraguay
B51	2012	<i>Eleusine</i>	Quirusillas, Bolivia

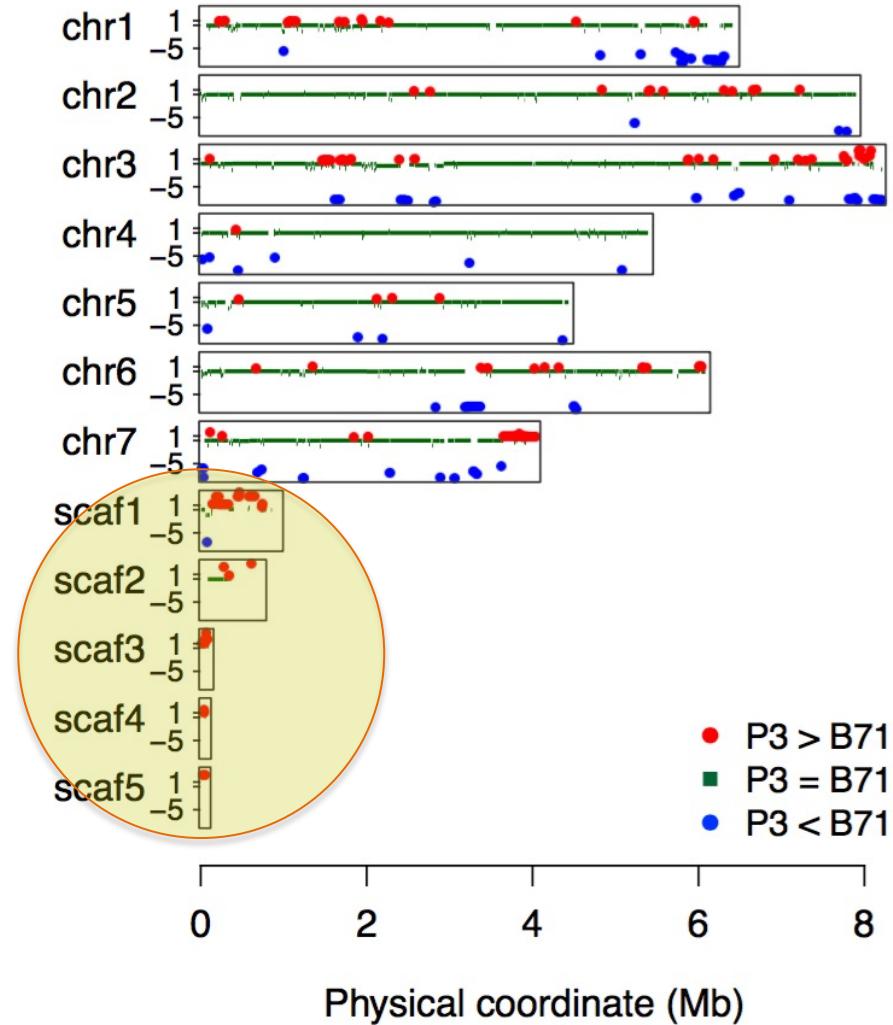
# Read depths to infer *copy number variation*

T25 vs B71

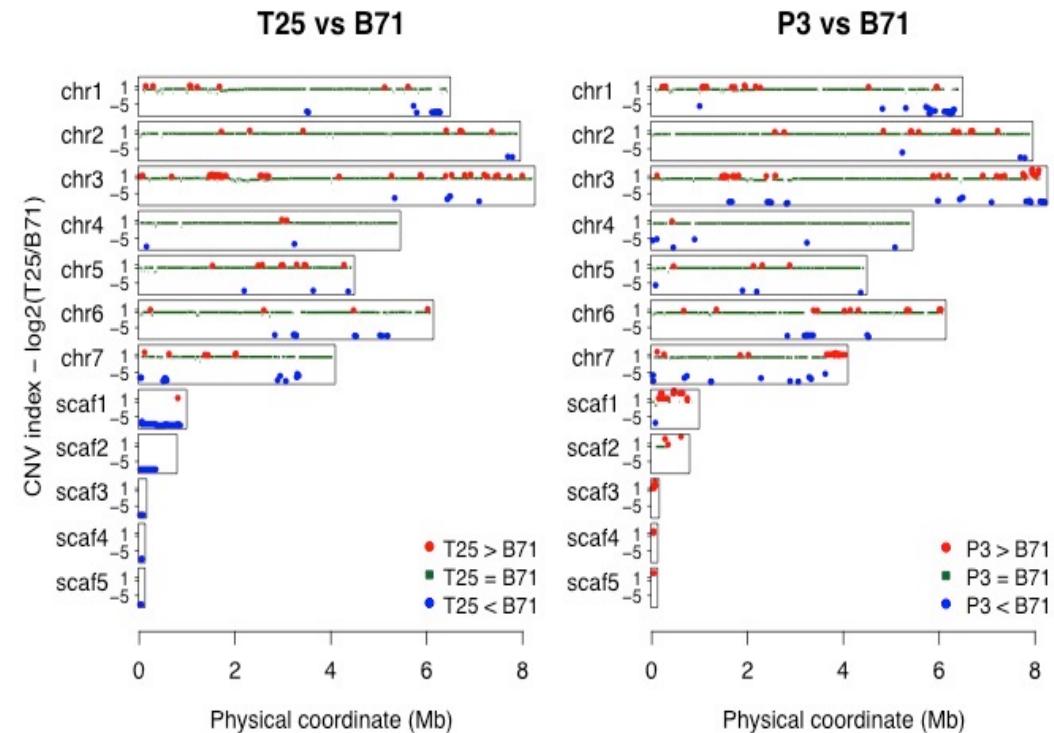
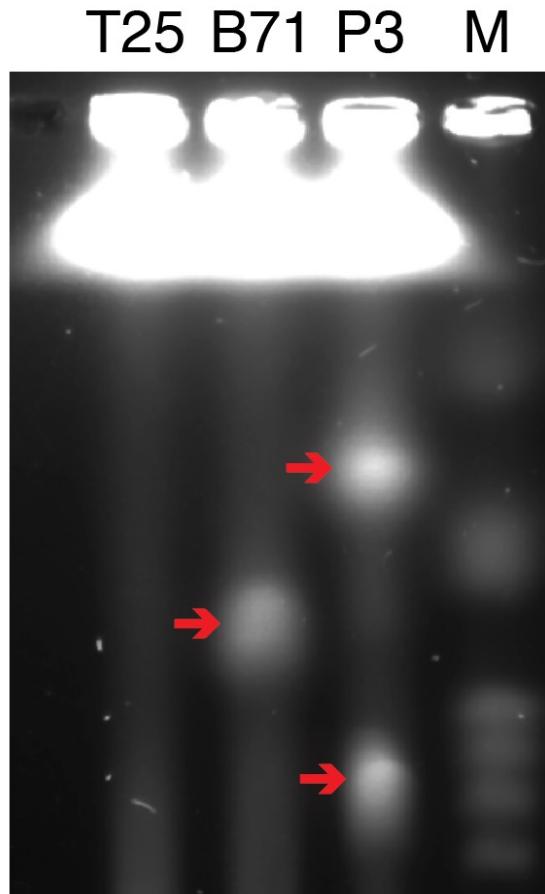
CNV index -  $\log_2(T25/B71)$



P3 vs B71



# *CHEF* gel to separate chromosome-sized DNA



T25 has no mini-chromosomes  
B71 has one and P3 has two

## Additional References

- Pinkel, D. & Albertson, D. G. Comparative genomic hybridization. *Annu. Rev. Genomics Hum. Genet.* 6, 331–354 (2005).
- Miller, W., Makova, K. D., Nekrutenko, A. & Hardison, R. C. Comparative genomics. *Annu. Rev. Genomics Hum. Genet.* 5, 15–56 (2004).
- Korbel, J. O. et al. Paired-end mapping reveals extensive structural variation in the human genome. *Science* 318, 420–426 (2007).
- Armstrong, J., Fiddes, I. T., Diekhans, M. & Paten, B. Whole-Genome Alignment and Comparative Annotation. *Annu Rev Anim Biosci* (2018). doi:10.1146/annurev-animal-020518-115005
- Kronenberg, Z. N. et al. High-resolution comparative analysis of great ape genomes. *Science* 360, (2018).