

# Genome Assembly

Jie (Jessie) Li

PhD

Bioinformatics Core

Genome Center

UCDavis



# UC DAVIS

## GENOME CENTER

Worldclass Research & Core Facilities

Genomics & Transcriptomics

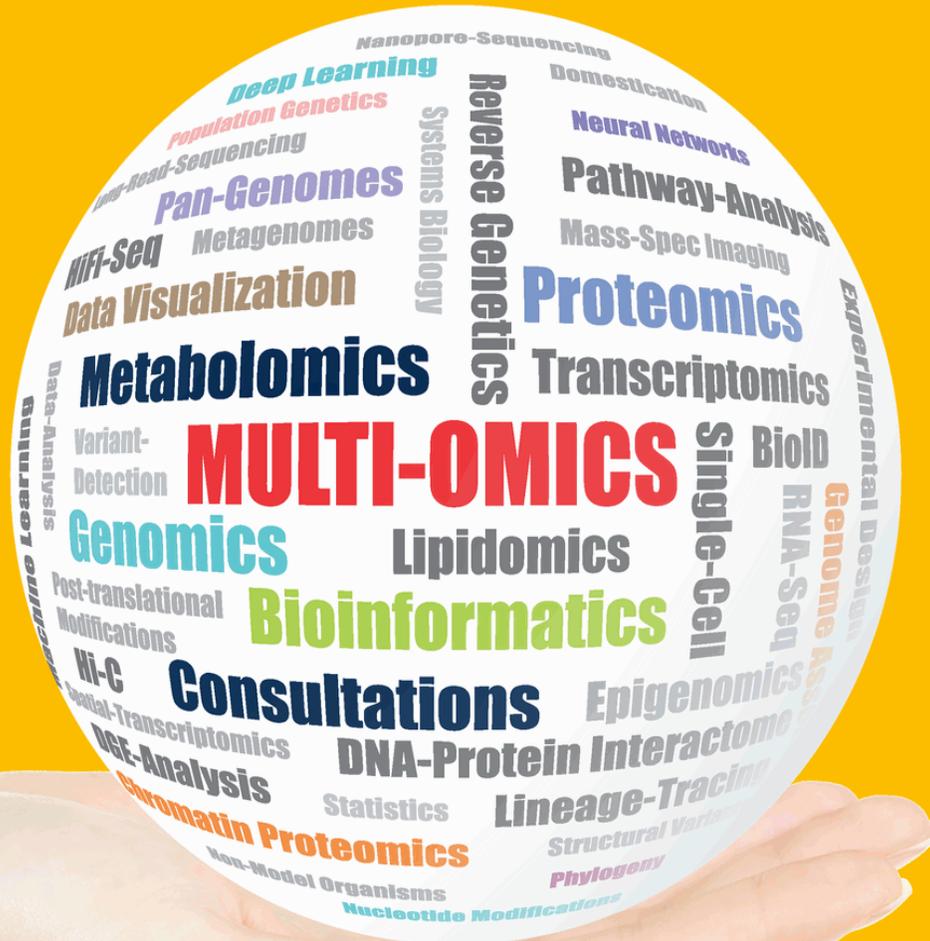
Proteomics

Metabolomics

Bioinformatics

Yeast One-Hybrid

Tilling



# UC Davis Bioinformatics Core in the Genome Center

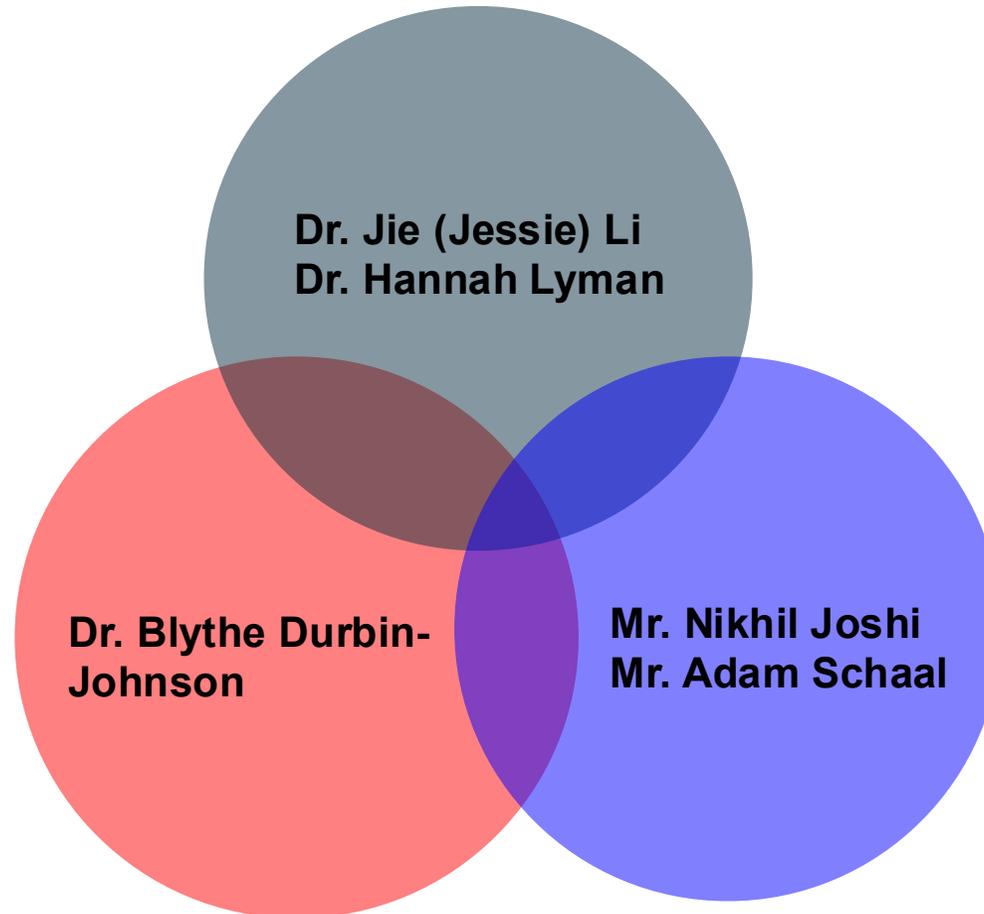
**Technical Director**  
Jean Challacombe

**Scientific Director**  
Ian Korf

**Computational Biology**

**Biostatistics**

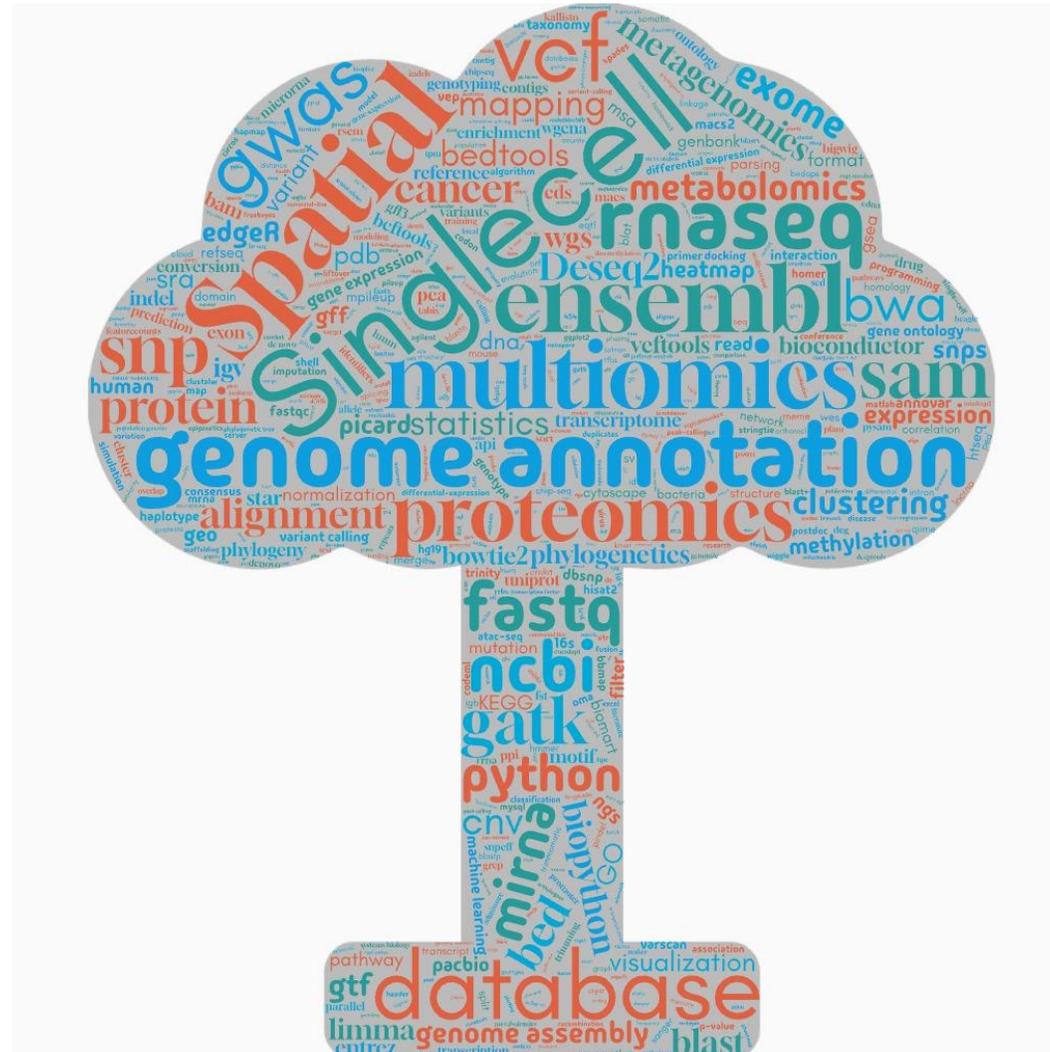
**Computer  
Science**



**Dr. Blythe Durbin-  
Johnson**

**Mr. Nikhil Joshi  
Mr. Adam Schaal**

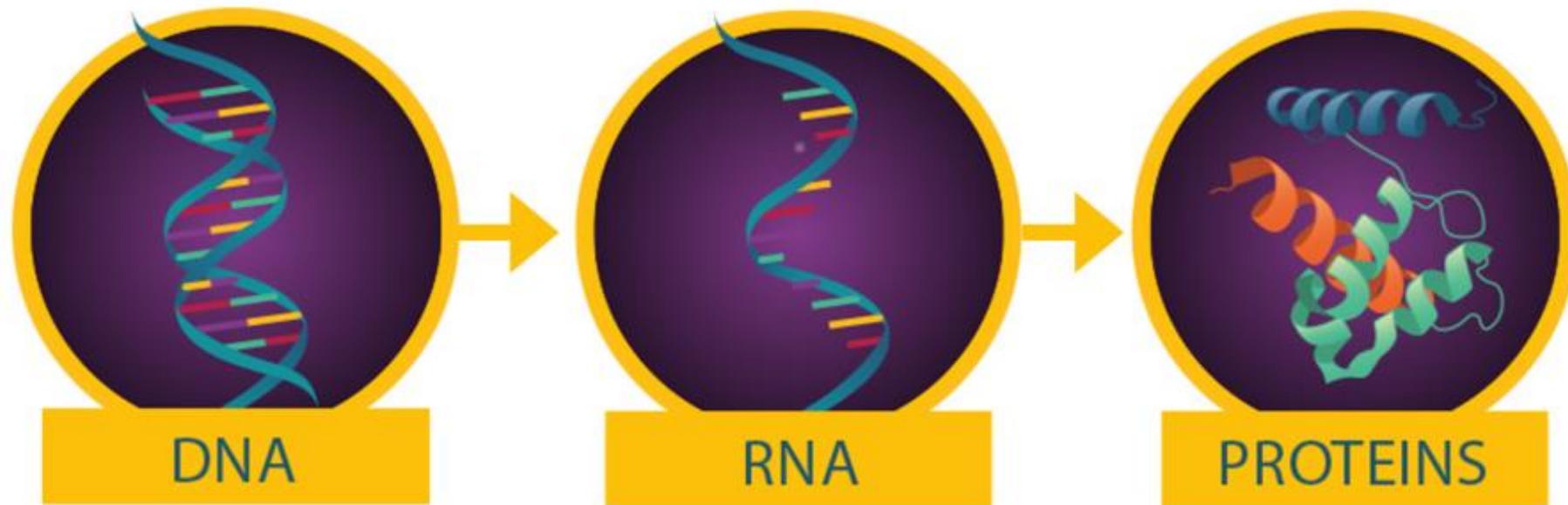
# UC Davis Bioinformatics Core in the Genome Center



# Outline

- Why assemble genomes and challenges
- Existing technologies for generating data
- Bioinformatics of genome assembly

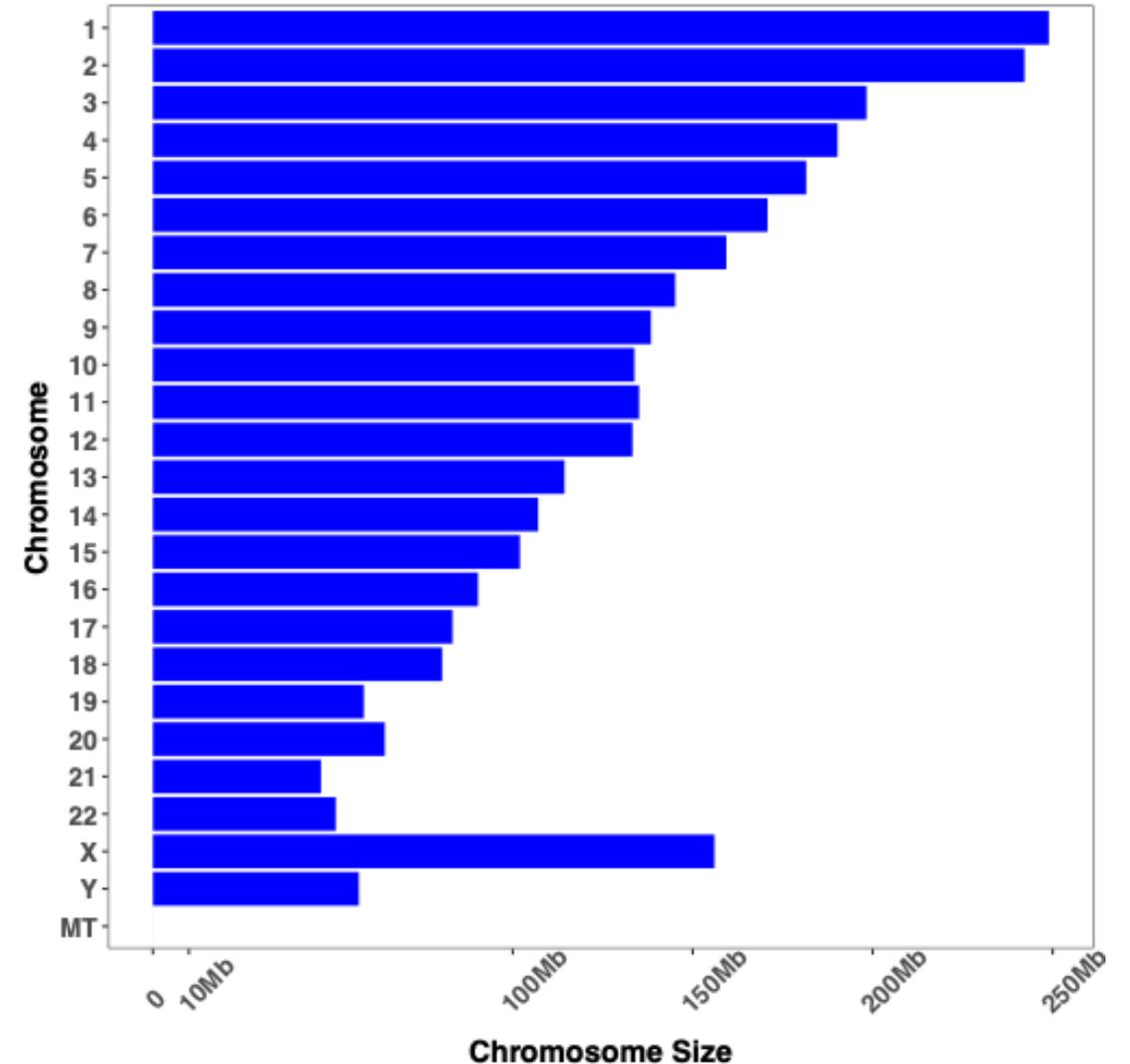
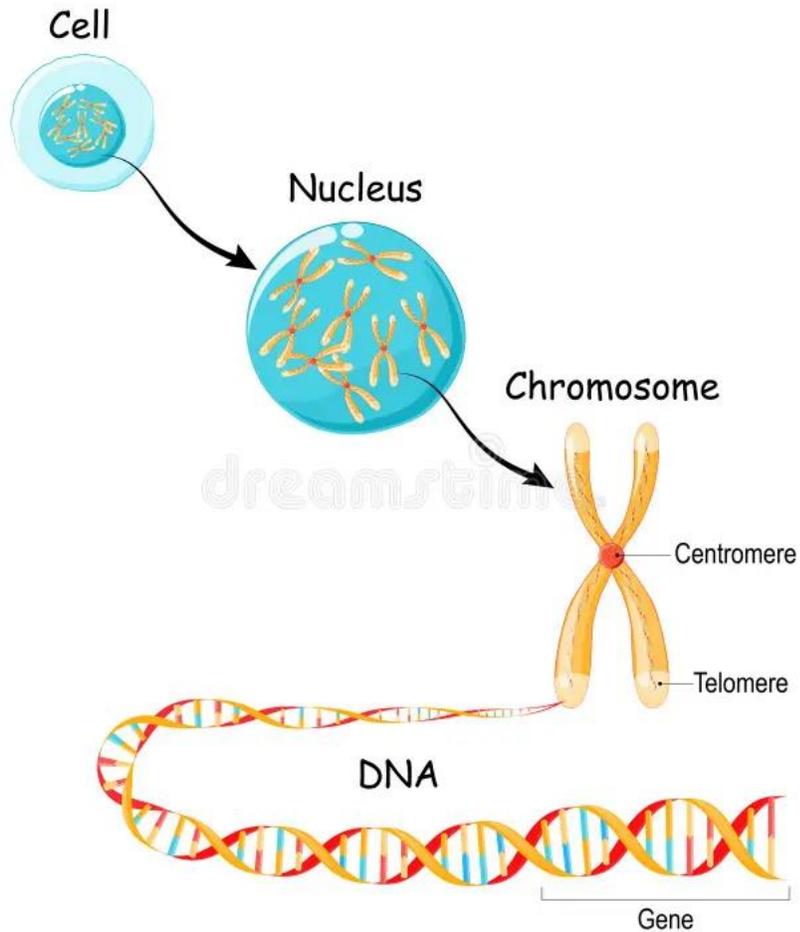
# Why genome assembly



# Human reference genome

- Human Genome Project
  - 1990 – 2003
  - Cost \$3Billion
- Genome Reference Consortium continued to improve the quality of the human reference. There were 79 “unresolved” gaps existing in May 2020 that accounts for ~5% of the genome.
- April 2022, Telomere-2-Telomere (T2T) consortium published a complete genome with non-Y chromosomes.
- August 2023, T2T reported the complete Y chromosome sequence

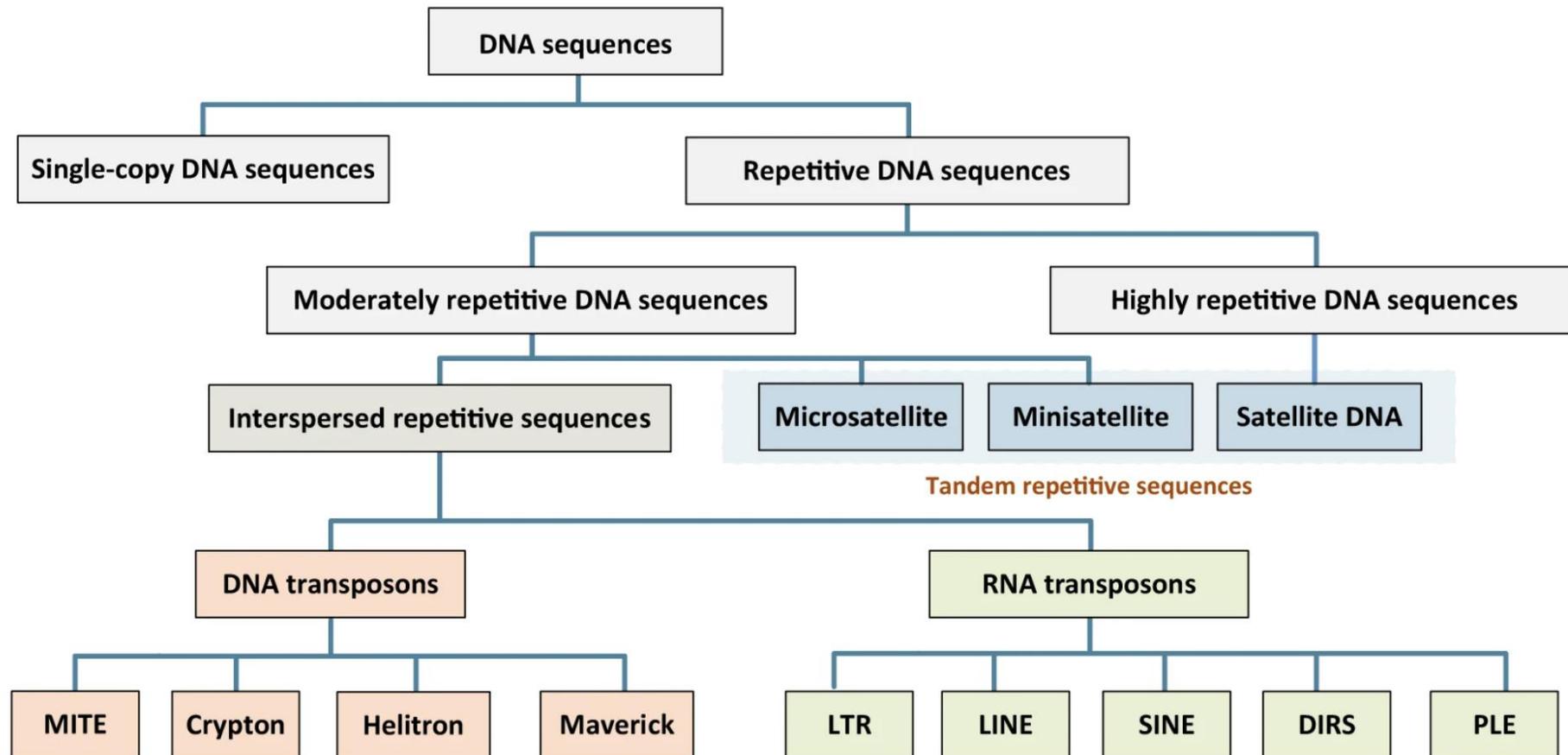
# First challenge in genome assembly



# Second challenge in genome assembly

- Repetitive sequences (repeats)
  - Patterns of nucleic acids that occur in multiple copies throughout a genome.
  - Content of repeats varies widely, ~50% in human and up to ~80% in some plant genomes.
  - Variable arrangements and sizes, with the largest repeat array in 100 Mbp.

# Second challenge in genome assembly

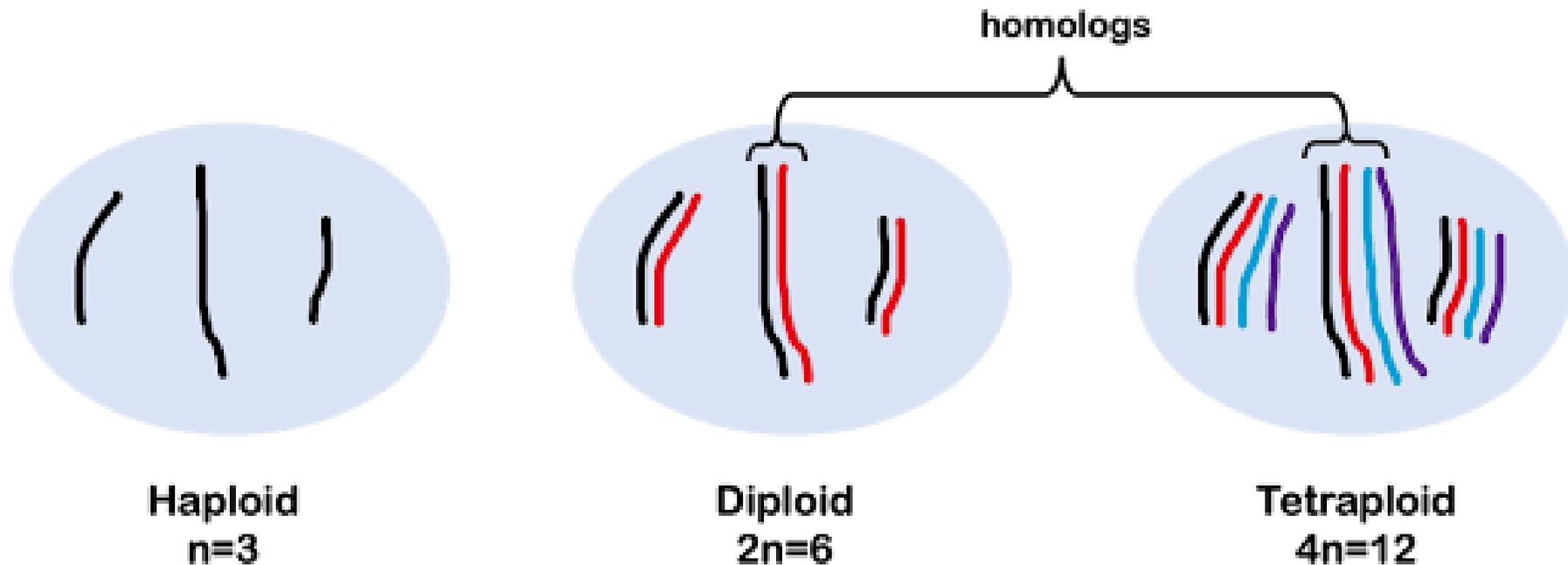


# Second challenge in genome assembly

- Repeats impact essential biological processes
  - Evolution of genome sizes
  - Regulation of gene expression
  - Variation induction
  - Disease

# Third challenge in genome assembly

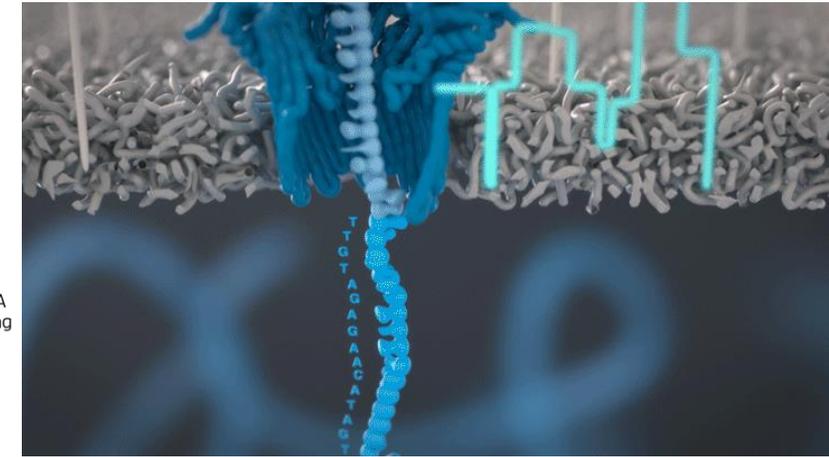
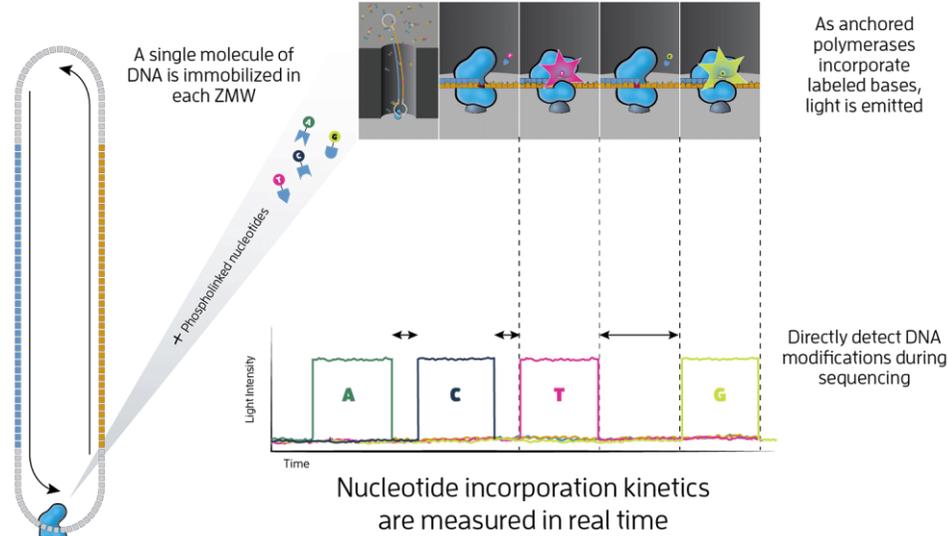
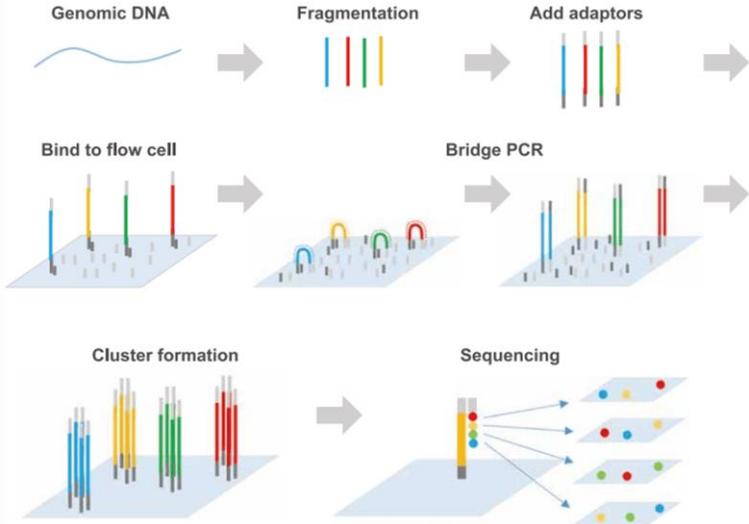
- Haplotype phasing



# Complementary approaches

Short-Read (Illumina, AVITI)	PacBio	PromethION Nanopore
Still-imaging of clusters (~1000 clonal molecules)	Movie recordings fluorescence of single molecules; HiFi: single-molecule circular consensus	Recording of electric current through a pore
Short reads - 2x300 bp AVITI, Miseq	HiFi reads: Up to 23 kb, N50 18 kb	Up to 100 kb, N50 25 kb Ultra-long protocol up to 1Mb
Repeats are mostly <u>not</u> analyzable	spans retro elements	spans retro elements
High output - up to 1 Tb per lane	up to 130 Gb HiFi data per SMRT-cell	Up to 100 Gb per flowcell Up to 80 million reads "cDNA-seq"
High accuracy (< 0.5 %)	Raw data error rate 15 % HiFi CCS data < 0.1%	Raw data error rate 2-10 %
Considerable base composition bias	No base composition bias	Some systematic errors
Very affordable	Costs 4 x higher	Costs 4 x higher
<i>De novo</i> assemblies result in thousands of scaffolds	"Near perfect" genome assemblies; lowest error rate	"Near perfect" genome assemblies; highest contiguity

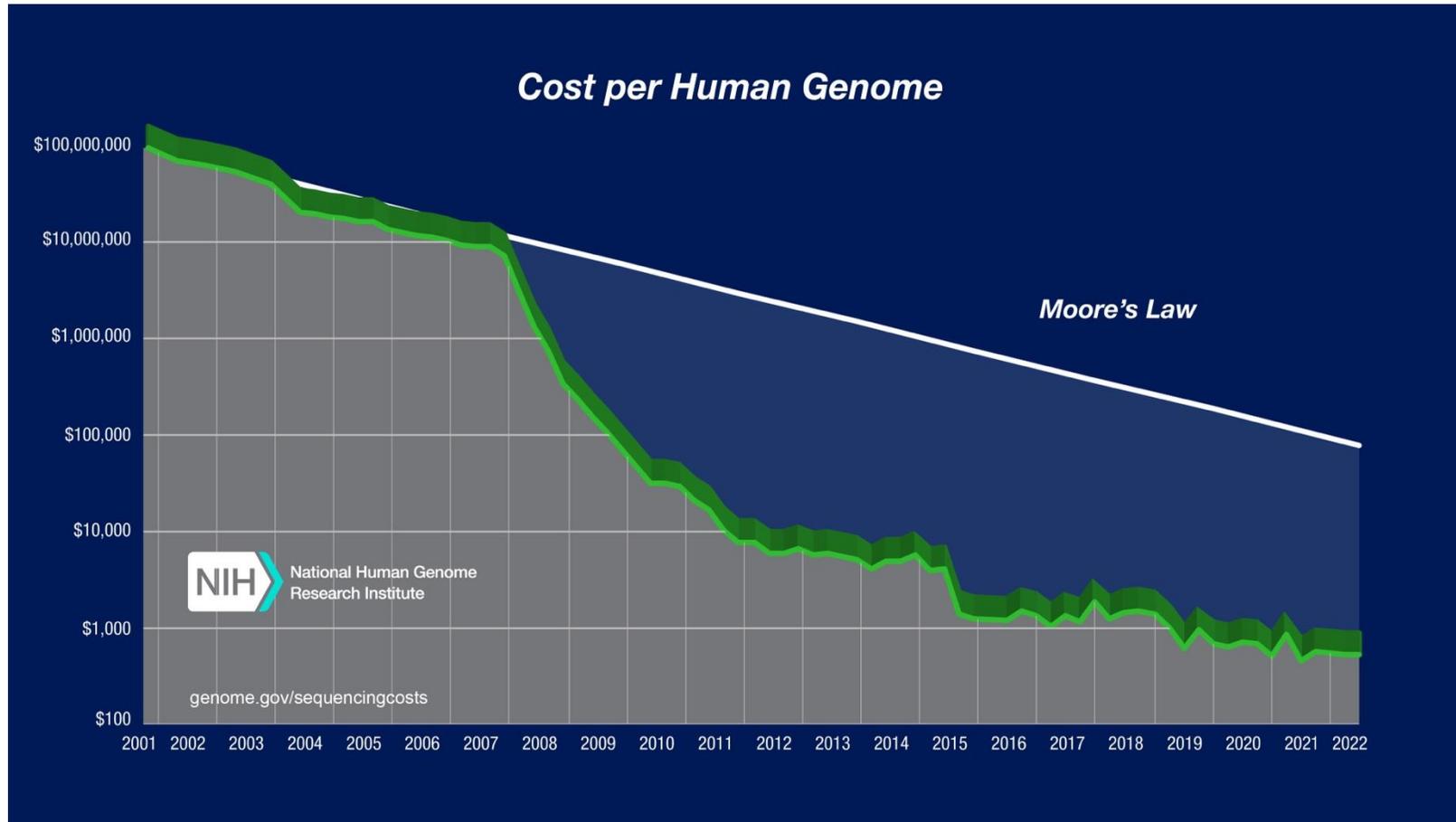
# Complementary approaches



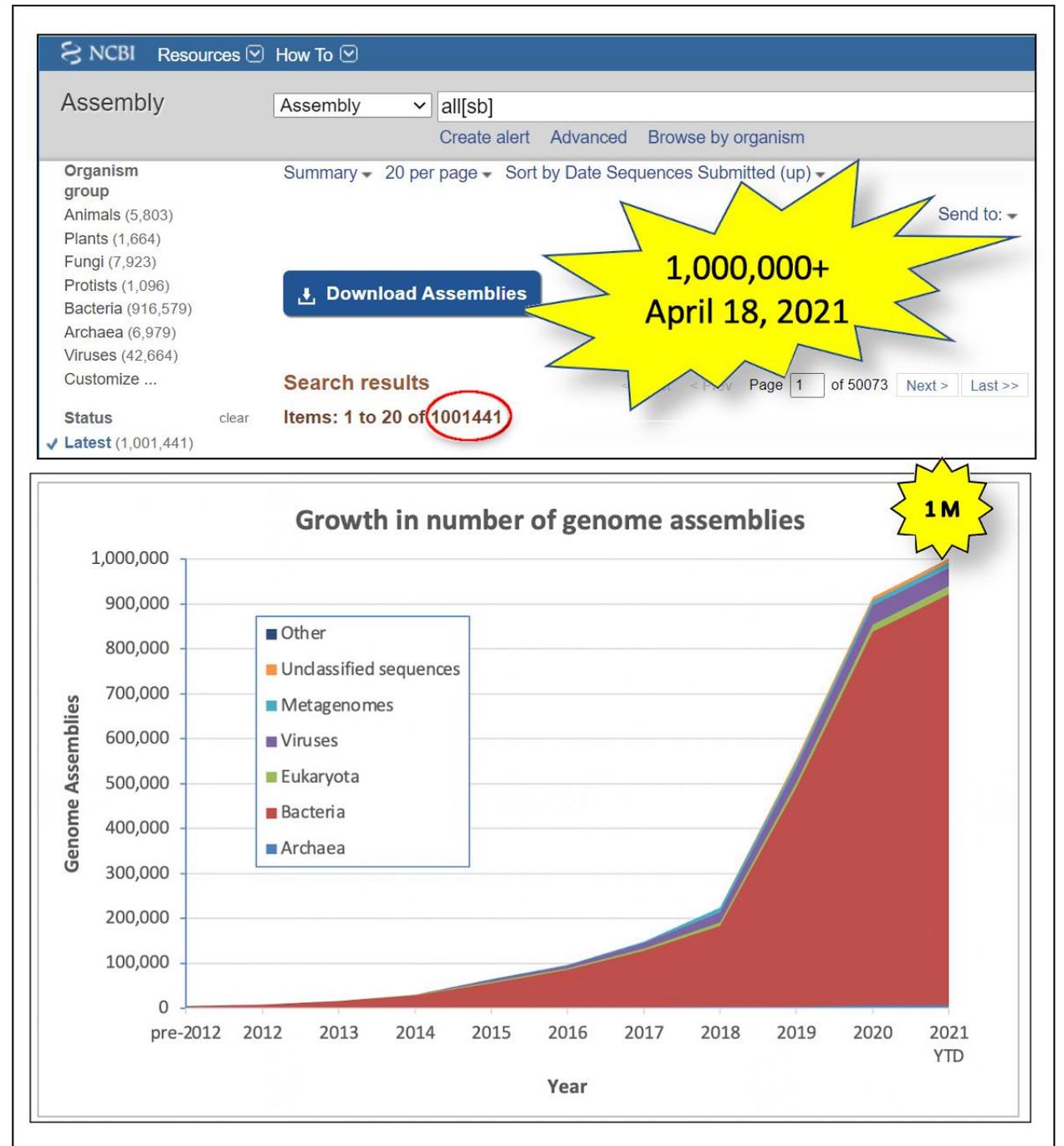
# Stages involved in assembling a genome

- High molecular weight DNA extract
  - Some tissues are easier to work with than others
  - Some species are easier to work with than others
- Sequencing library prep
- Sequencing
  - PacBio HiFi + ultra-long ONT/HiC
  - ONT
  - Illumina (very small genomes)
- Bioinformatically assemble the genome

# Sequencing cost keeps decreasing



# Number of genomes assembled increasing



# Number of genomes assembled increasing

The screenshot shows the NCBI Genome page. At the top, there is the NIH logo and the text "National Library of Medicine National Center for Biotechnology Information". Below this is a search bar with the placeholder "Search NCBI ..." and a "Log in" link. The main navigation menu includes "NCBI Datasets", "Taxonomy", "Genome" (which is highlighted), "Gene", "Command-line tools", and "Documentation".

## Genome

Search by taxonomic name or ID, Assembly name, BioProject, BioSample, WGS or Nucleotide accession

Search term

Try examples: [Homo sapiens](#) [GCF\\_000001405.40](#) [PRJNA489243](#) [SAMN15960293](#) [WFKY01](#) [GRCh38.p14](#) [NC\\_000913.3](#)

### Genomic data available from NCBI Datasets

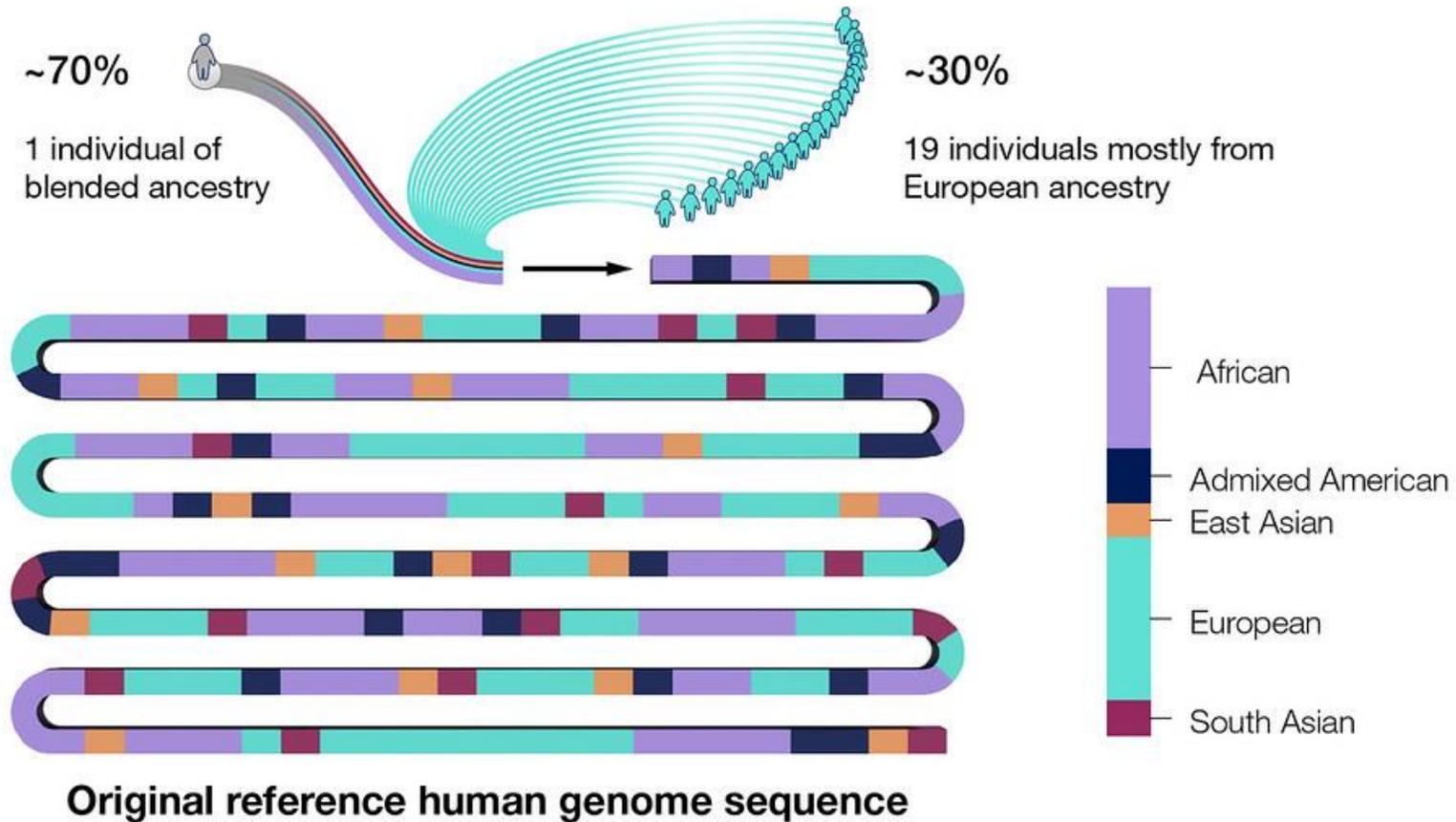
Click below to learn more about the genomic data available from NCBI Datasets.

Category	Count
Eukaryota	3.25M
Archaea	43.23K
Bacteria	2.61M
Viruses	-

### All Genomes

Category	Count
Total	3.25M
Reference	43.23K
Annotated	2.61M

# Human reference genome

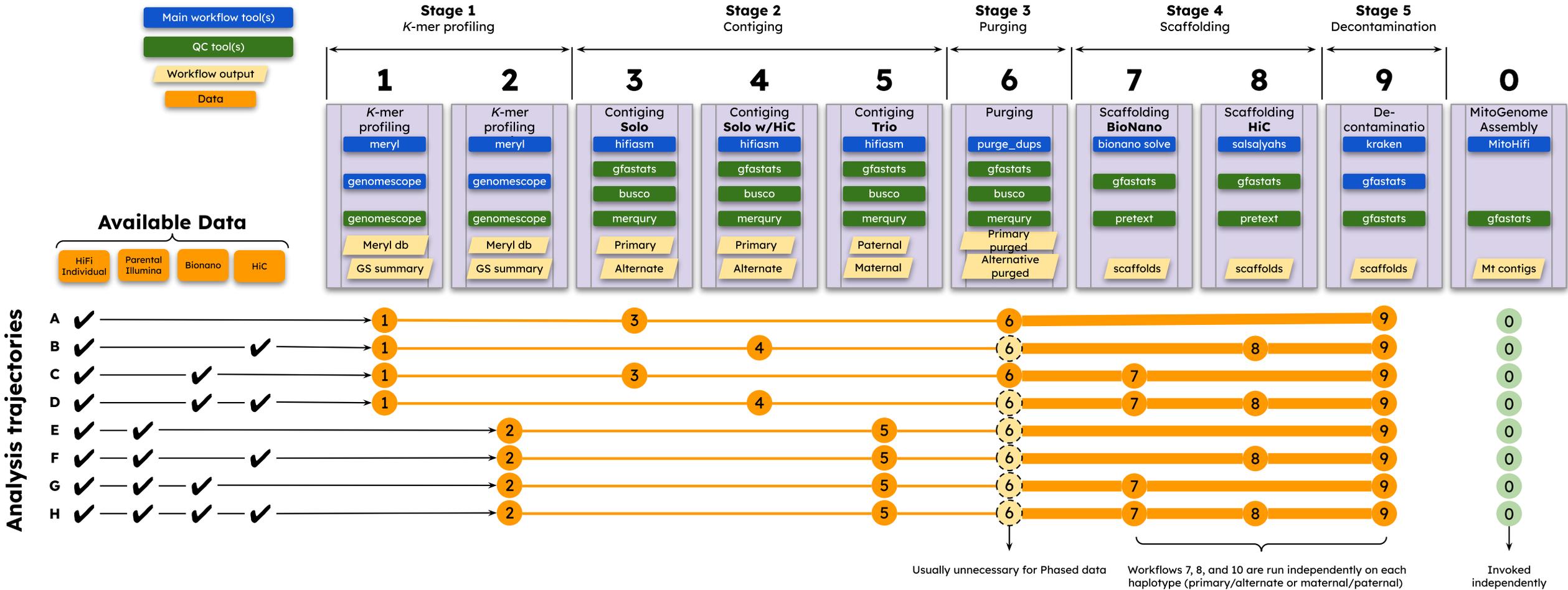


# Genome assembly efforts

- Earth BioGenome Project (EBP)
  - An initiative aiming to sequence and annotate all 1.5-1.8 Million eukaryotic species
  - A global network includes over 61 affiliated projects
  - Vertebrate Genomes Project (VGP)
- Telomere-to-Telomere (T2T) consortium
- Genome Reference Consortium (GRC)

**Generate haplotype phased telomere-to-telomere high quality genomes.**

# VGP assembly pipeline



# VGP assembly pipeline

- Stage 1: Data QC (Kmer profiling)
- Stage 2: Assemble reads to contigs
- Stage 3: Purge duplicates
- Stage 4: Scaffold
- Stage 5: Decontamination

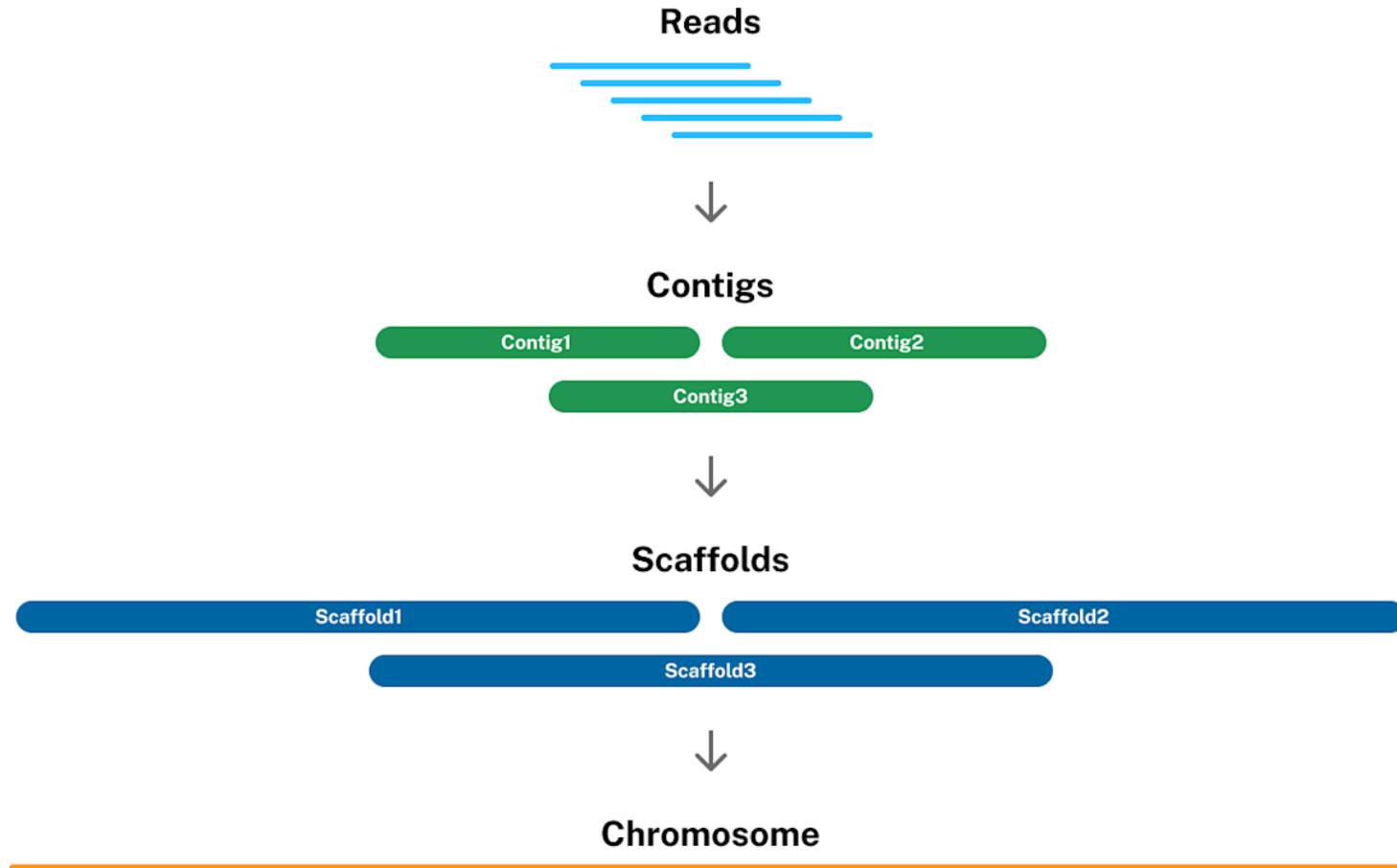
# Current approaches

- PacBio HiFi + HiC/Bionano
- PacBio HiFi + UL ONT
- ONT + UL ONT
- Trio data facilitate haplotype phasing

# VGP established assembly pipeline

- Stage 1: Data QC (Kmer profiling)
- Stage 2: Assemble reads to contigs
- Stage 3: Purge duplicates
- Stage 4: Scaffold
- Stage 5: Decontamination

# Premise of genome assembly

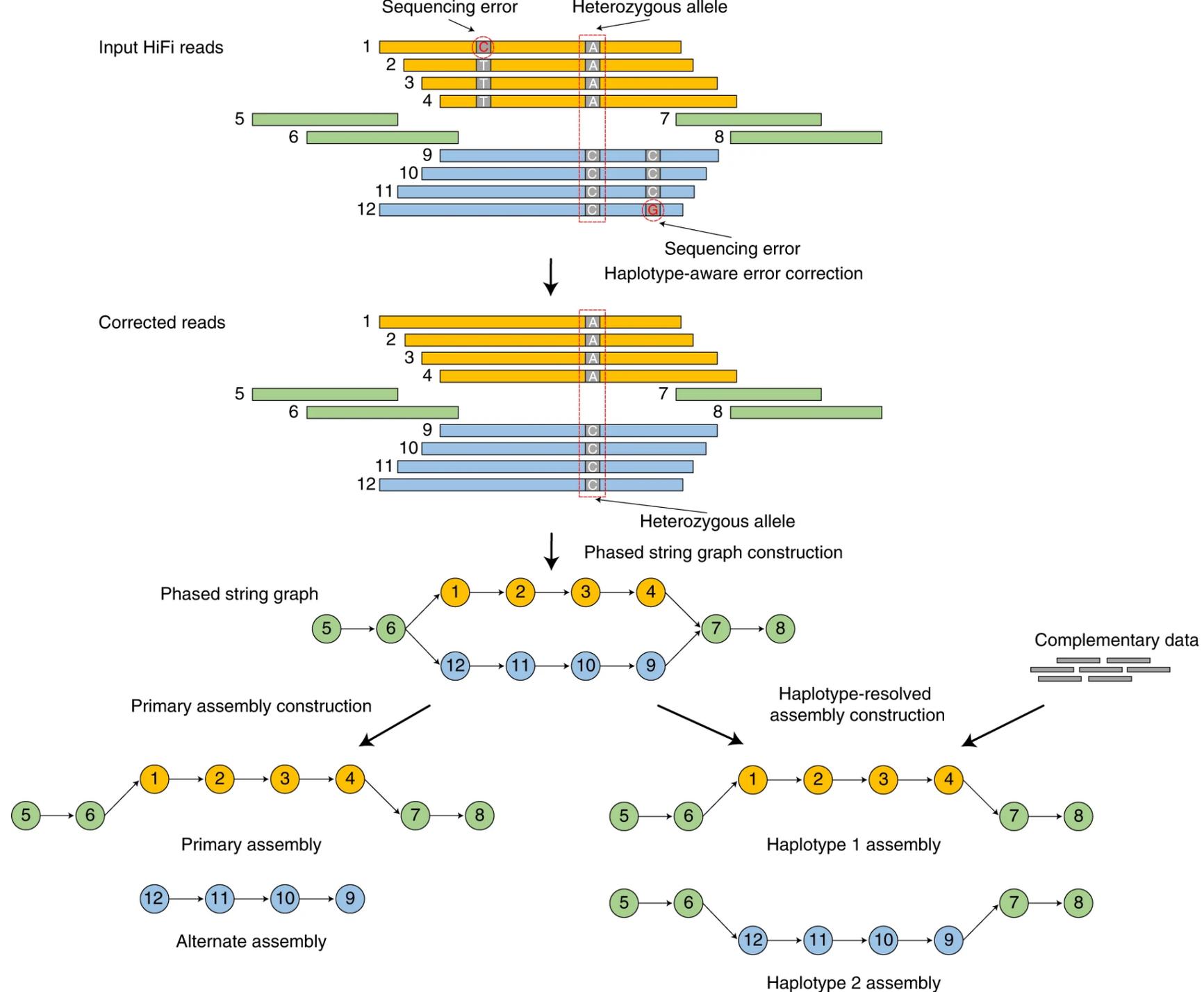


Premise of genome assembly

*GRAPH  
THEORY*

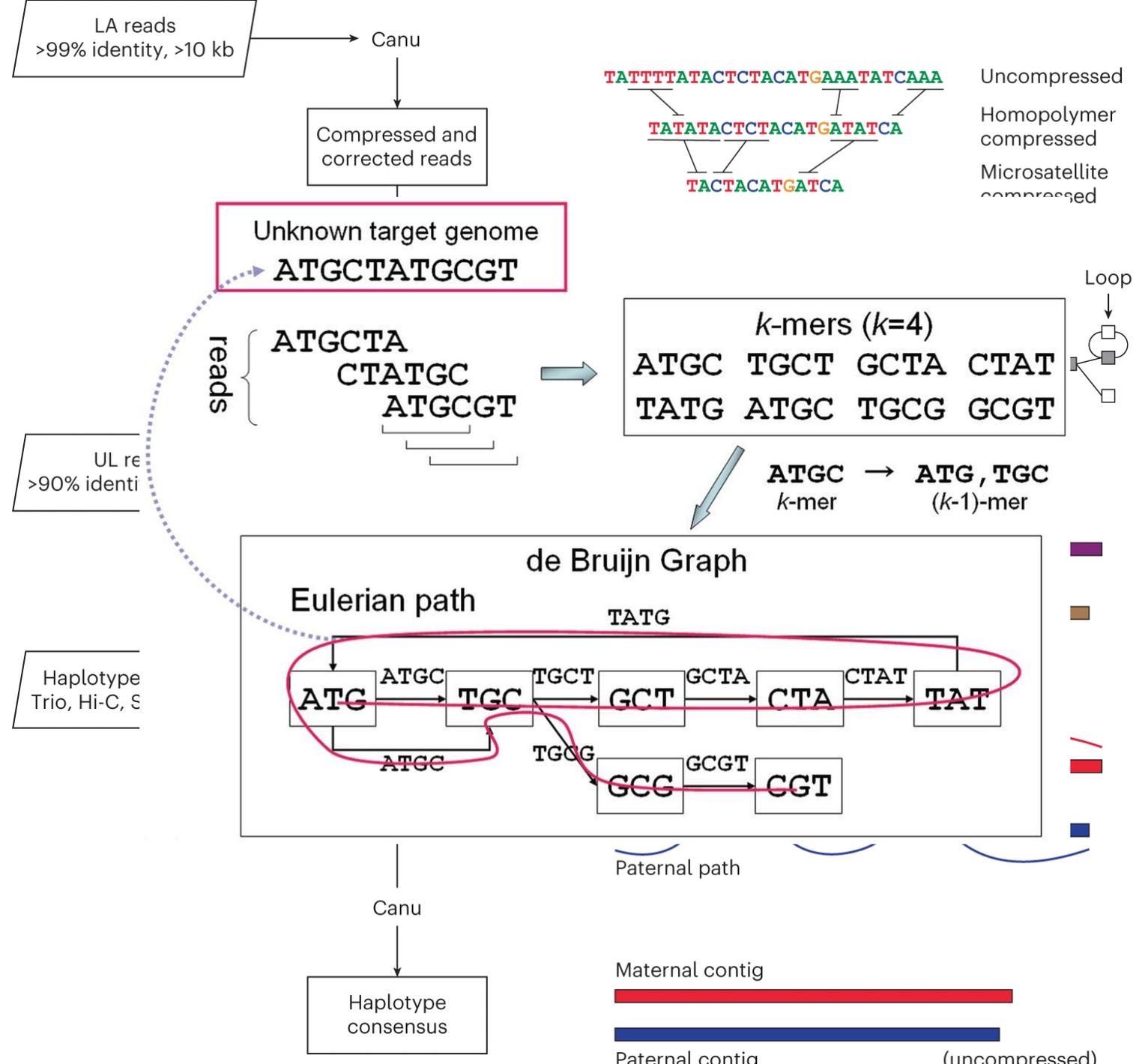
# HiFiasm

Cheng et. al, *Nat Methods*,  
2021



# Verkko

Rautiainen, M. *et al. Nat Biotechnol* 2023



# Quality assessment

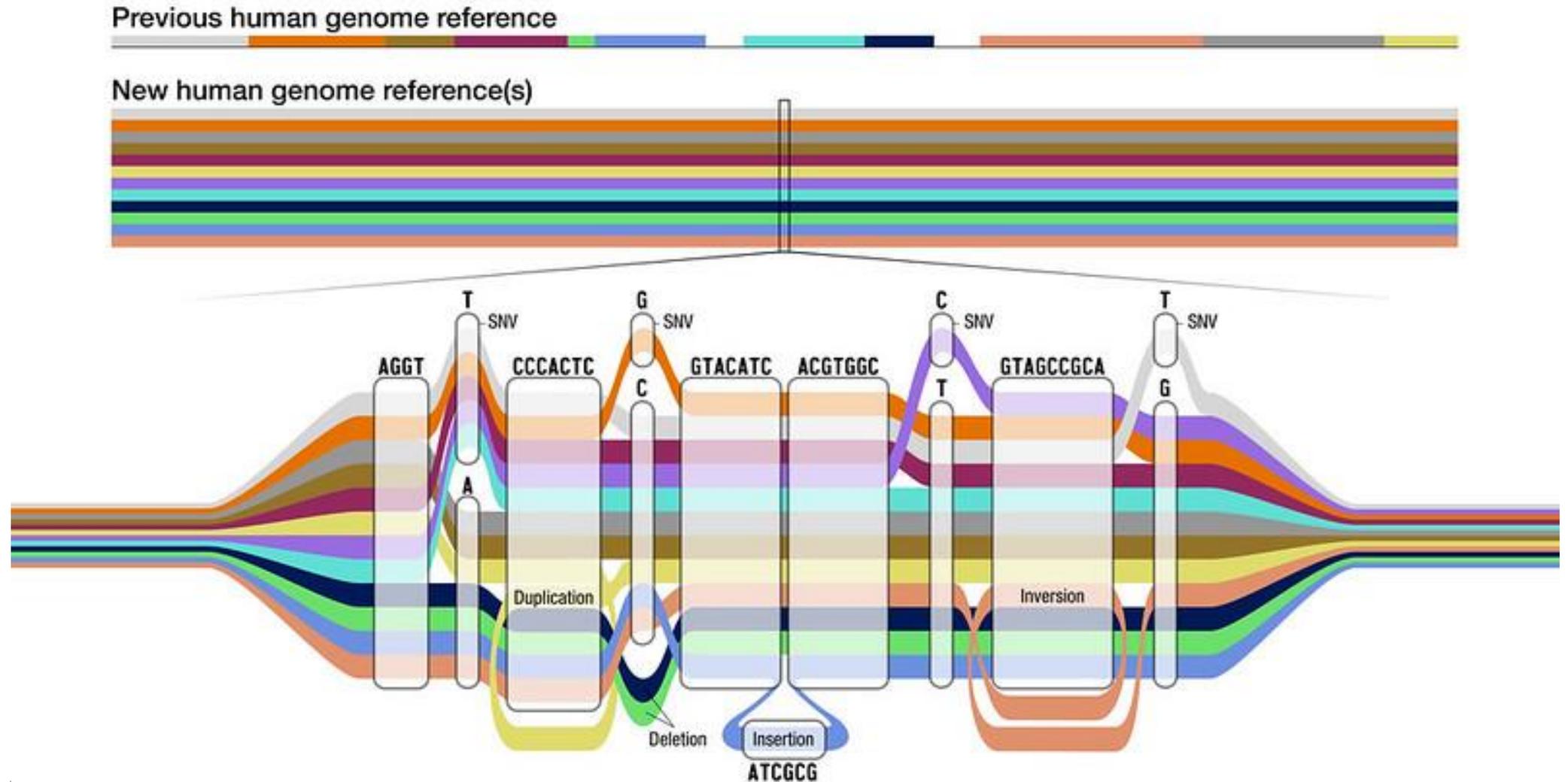
- Contiguity
  - N50
  - Gaps
- Completeness
  - BUSCO
  - Repeat structures
- Contamination

# Genome assembly

- Generate appropriate data
  - PacBio HiFi/ONT
  - HiC/ONT ultra long
  - Trio data facilitate better phasing
- Assemble
  - HiFiasm
  - Verkko
- Quality assessment



# Pangenome era



Thank you!