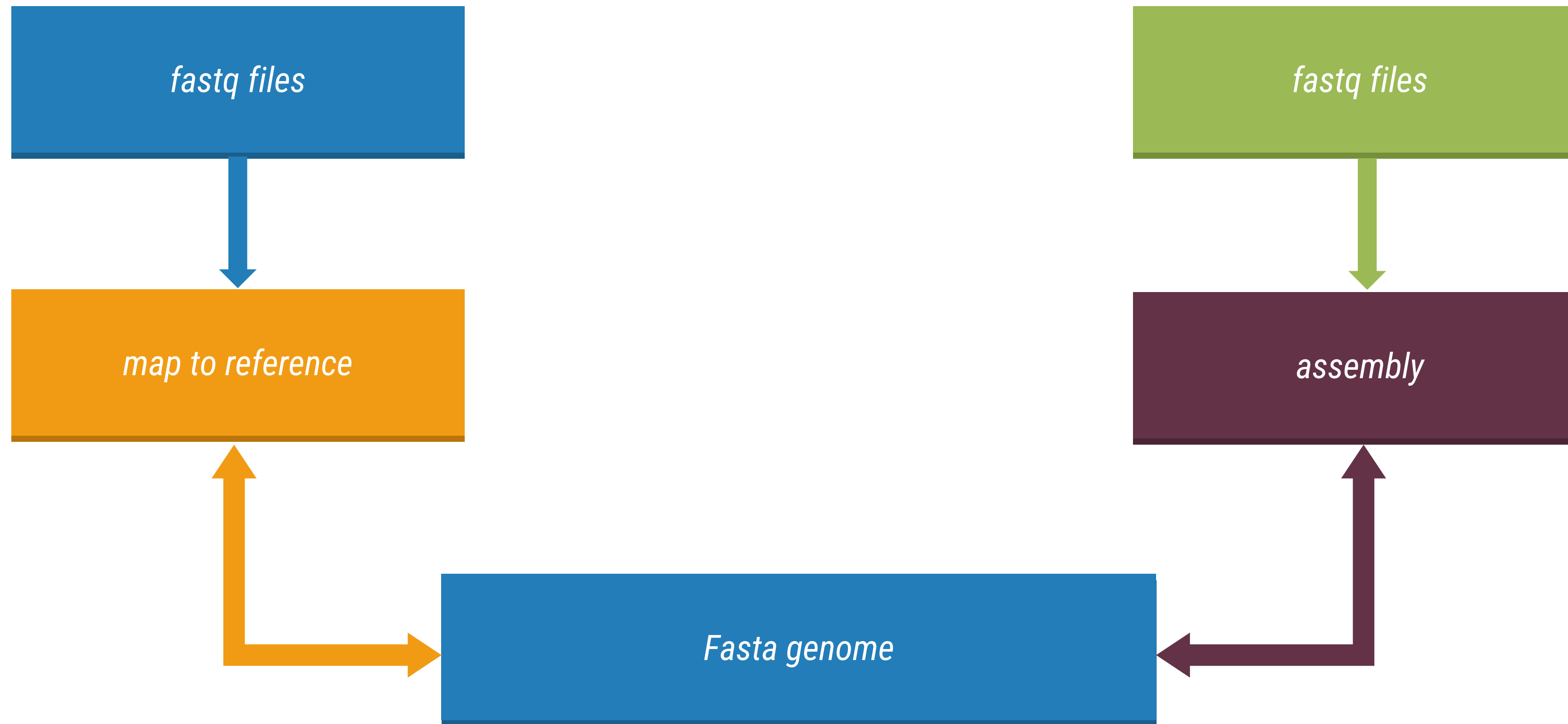


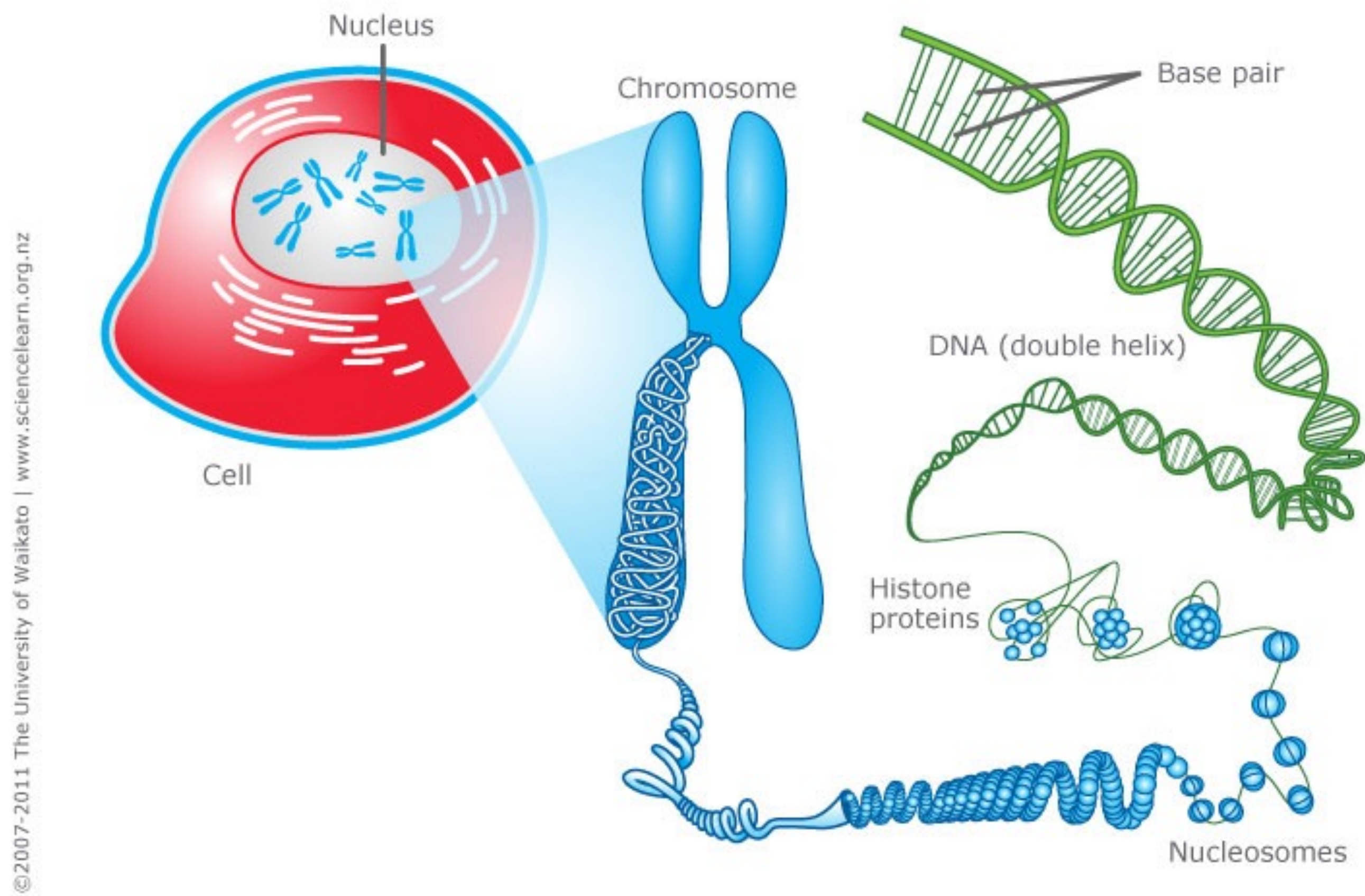
GENOME ASSEMBLY AND ANNOTATION



Having good quality fastq data is important!!

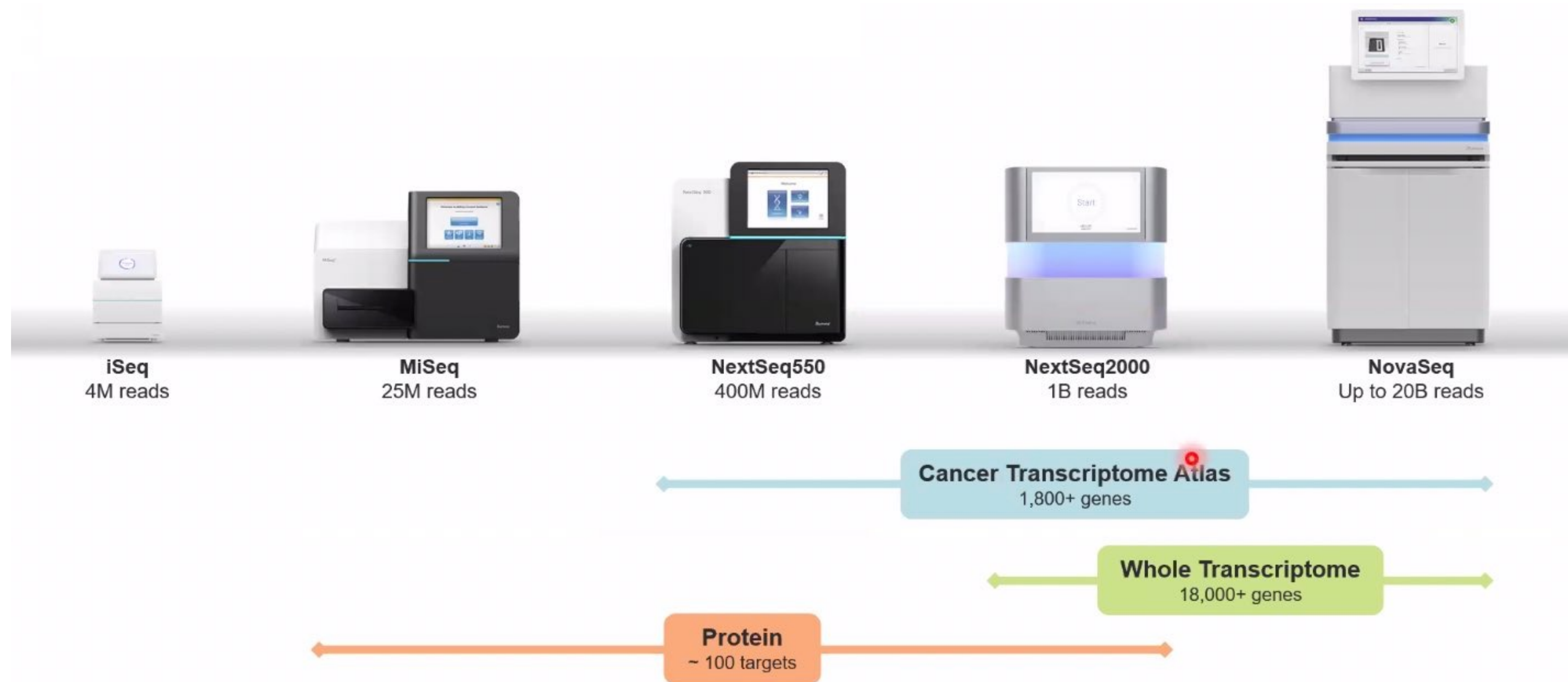


Genome sequencing is conceptually straight forward



```
CCACAAGTTCCTGACTGCCTGACTTCCCTTCACCGACTGGCACTTTCCACTCGGATCGCC
AGCGACCGTTACTAAAAAACAACATCGAATACTGTCTGCAAGACAGTGCAATAAAGCA
AATGAAAATAATTAGAATAAGAATAATGTTAATAATGATACCAAAAATTCTTCGGCTGGA
ACTGATGTGACTCTATGCATAATGTGAAATTTCCATGACGAACGAACACGCATCCTACAC
CAGATTTTGAGTAATGTTCTCTATATATGCATCTAATTCCTAAGATAAATGGGTGTGAG
CAGCAACTAGAATTGGAAAGAACCACTGAACAGCTTGGTACCTTTTGCAGAGCTACACGC
CACGATTCTAAAGCGCTCGATTCTGCTGGATACGGCTGGATGTCACTGGCCTCCTCCGA
AGATGATCATACTTGAAGGATTTCTTCGACGCATCATCCACGTAGCGACCTCCTCCAATA
TCGTTGAATGGCGTCACAAGTGTCTAAGCGACATTTTTTGAAAGTGTAAGCAATACTTGA
TCAACTTCTTTTCCGAATCCTCGATGAATAATCATTAGAACAATACGTTTCCATTTTTTA
CTGCACCCTTTTCCATACAAATAGTTACAAGTTAGATTGGACAATGACAACAAATGAACA
CCGAGTTCATGTTGAGAAAACATGCACTAGGGAATCGACCGCCTTGTGCAGCAGCATTTC
ATGGTGAAAAGACAAGATCTGTCATAGATGAGTTTAAGTTTGAACAGTTCCCCTTCAAAT
TCCACACTATCACAGACCATTCCCCGAAGAATGTTCGACCTCTAGACCTGACCTCTACGA
```


Sequencing genomes is easy...

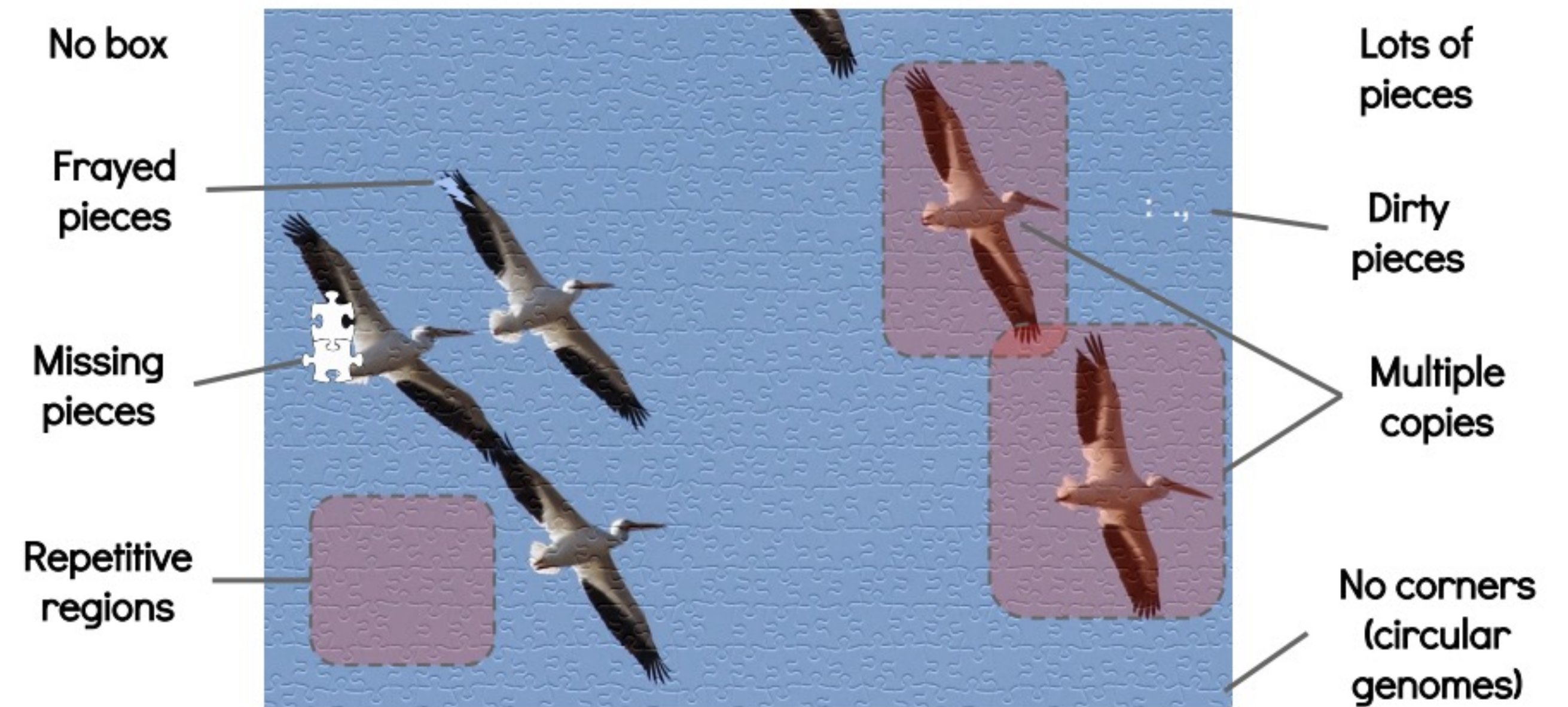


Sequencing is easy, constructing good genomes is not

- **Genome: *biologically***
 - “the haploid set of chromosomes in a gamete or microorganism, or in each cell of a multicellular organism”
 - “the complete set of genes or genetic material present in a cell or organism”

Sequencing genomes is easy, constructing good genomes is not

- **Genome: *bioinformatically***
 - Best guess, but often:
 - highly fragmented
 - misassembled to some degree
 - contaminated
 - duplicated or missing



Draft genomes
“manageable”(?)

Chromosome-scale
genomes
HARD

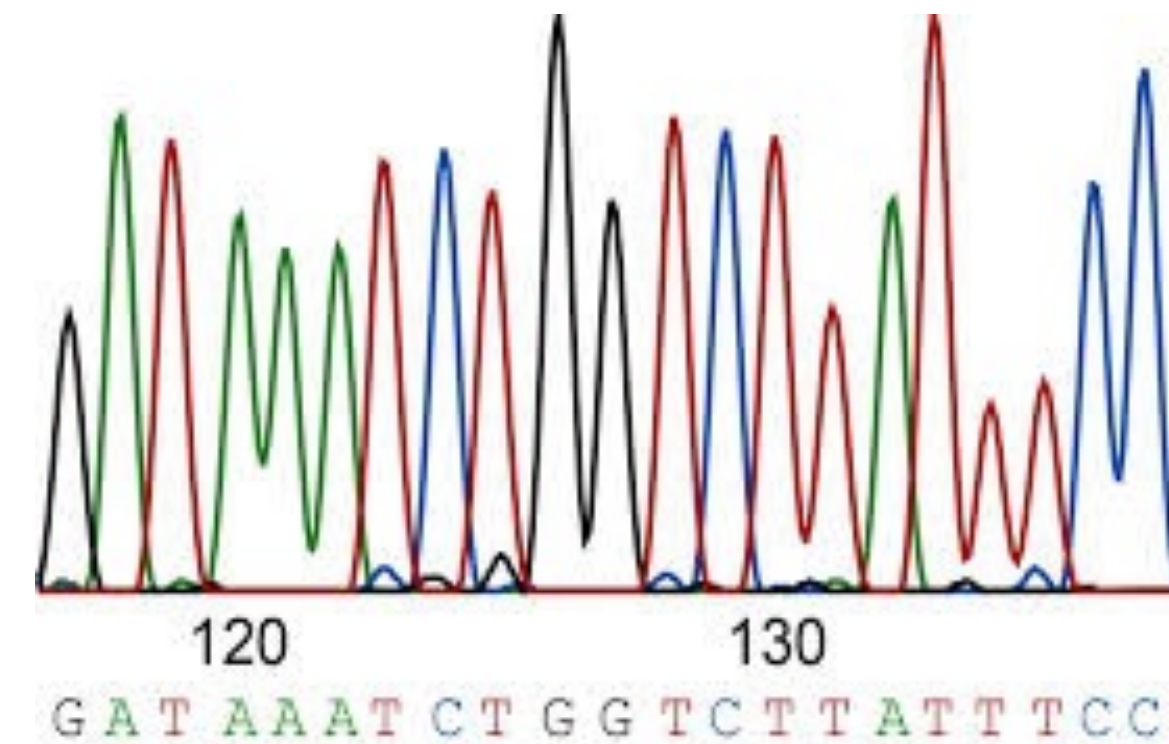
Time, money, expertise

New technologies are making genome assembly easier

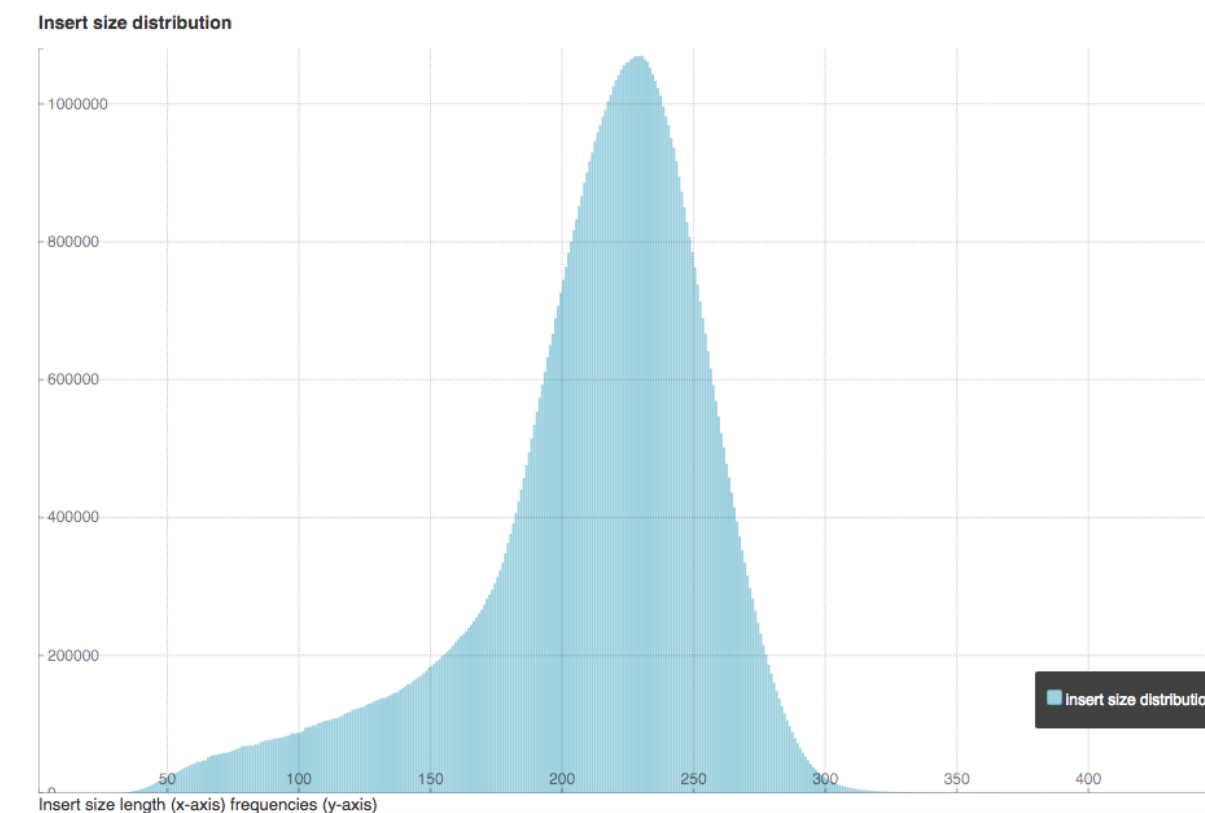
Sanger Sequencing: ABI

High Throughput Sequencing: Illumina

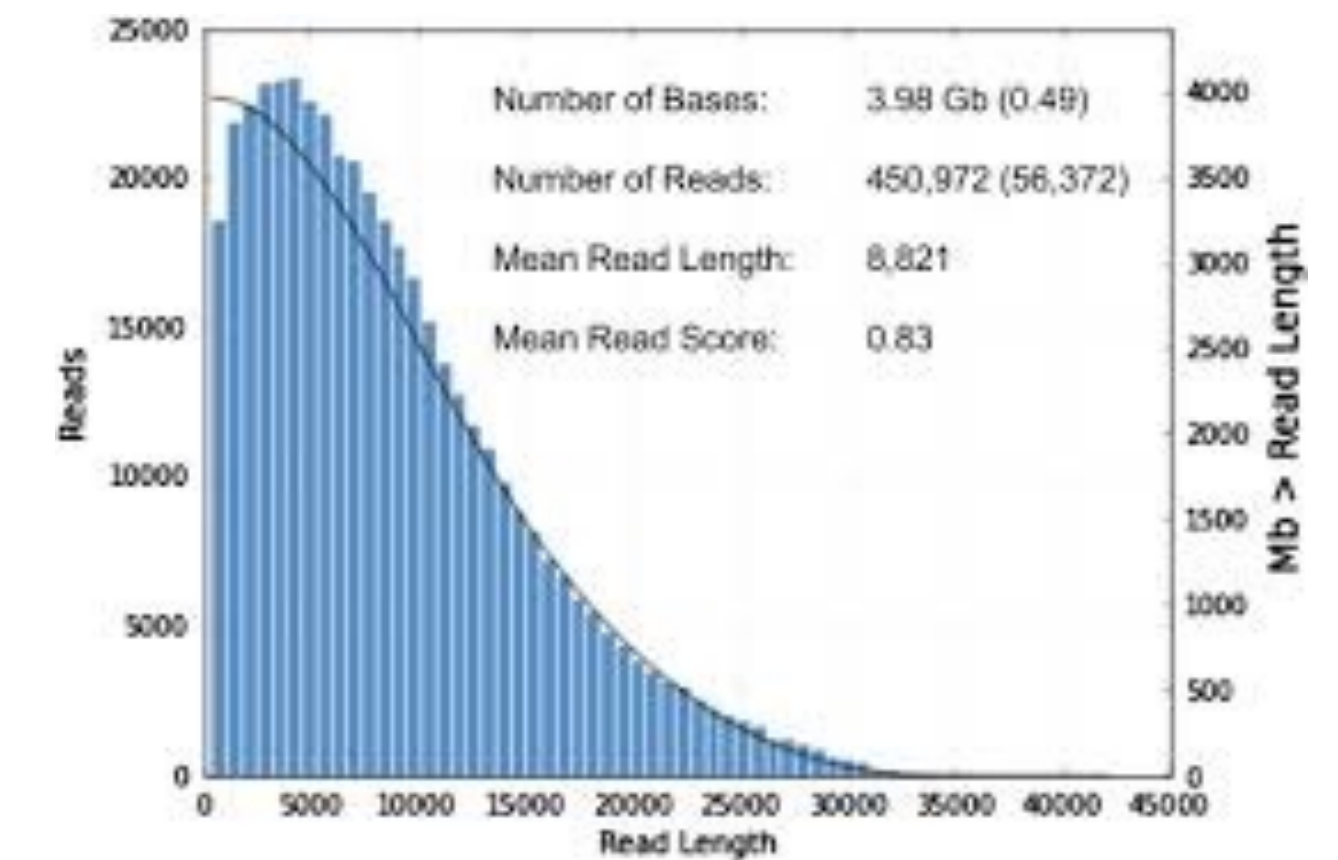
Long read sequencing: Pacbio & Nanopore



Read length: 500-1000 bp

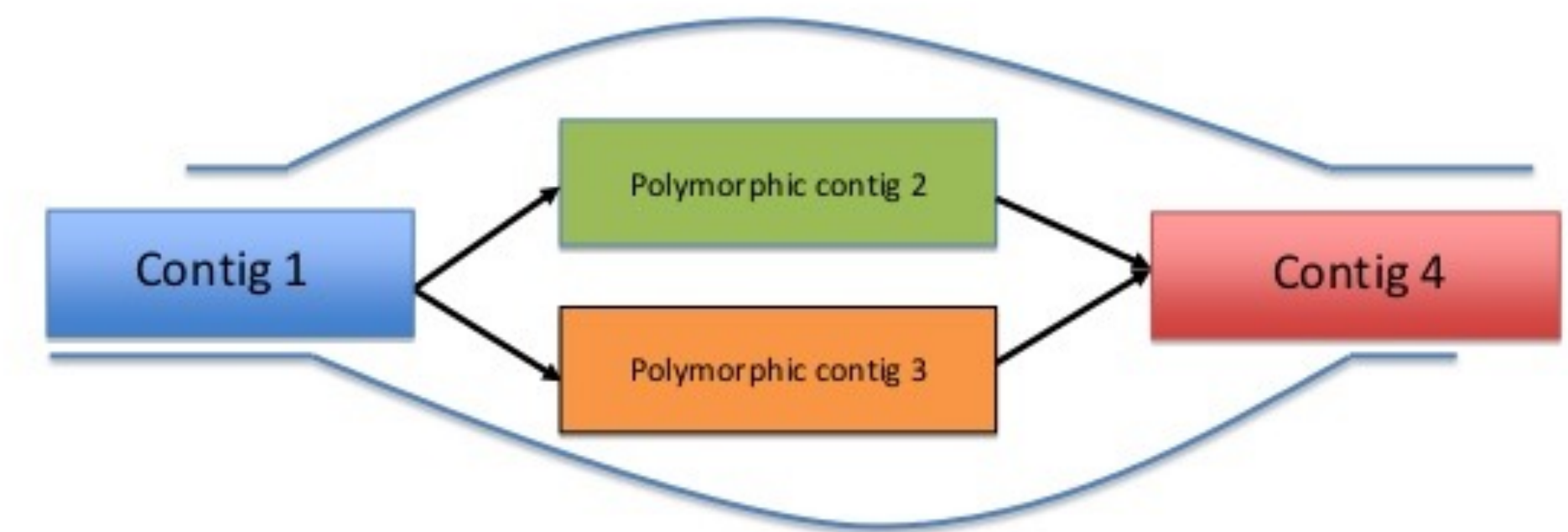
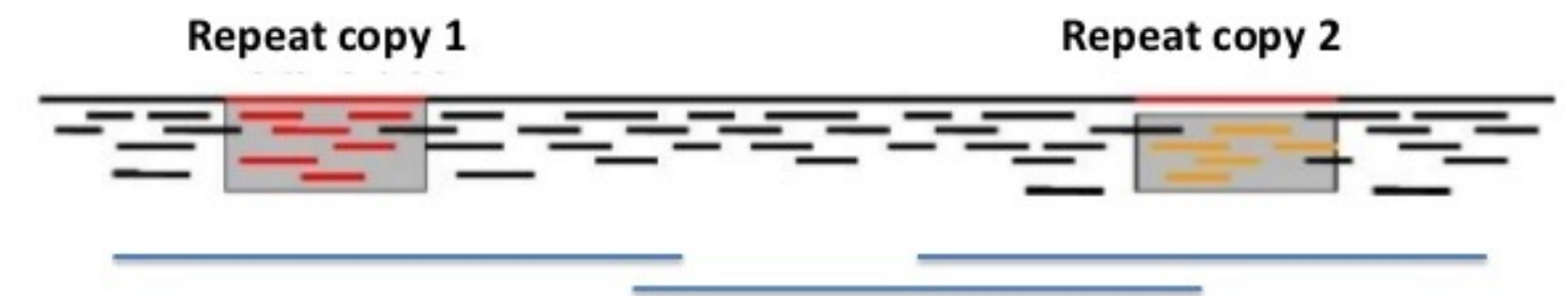
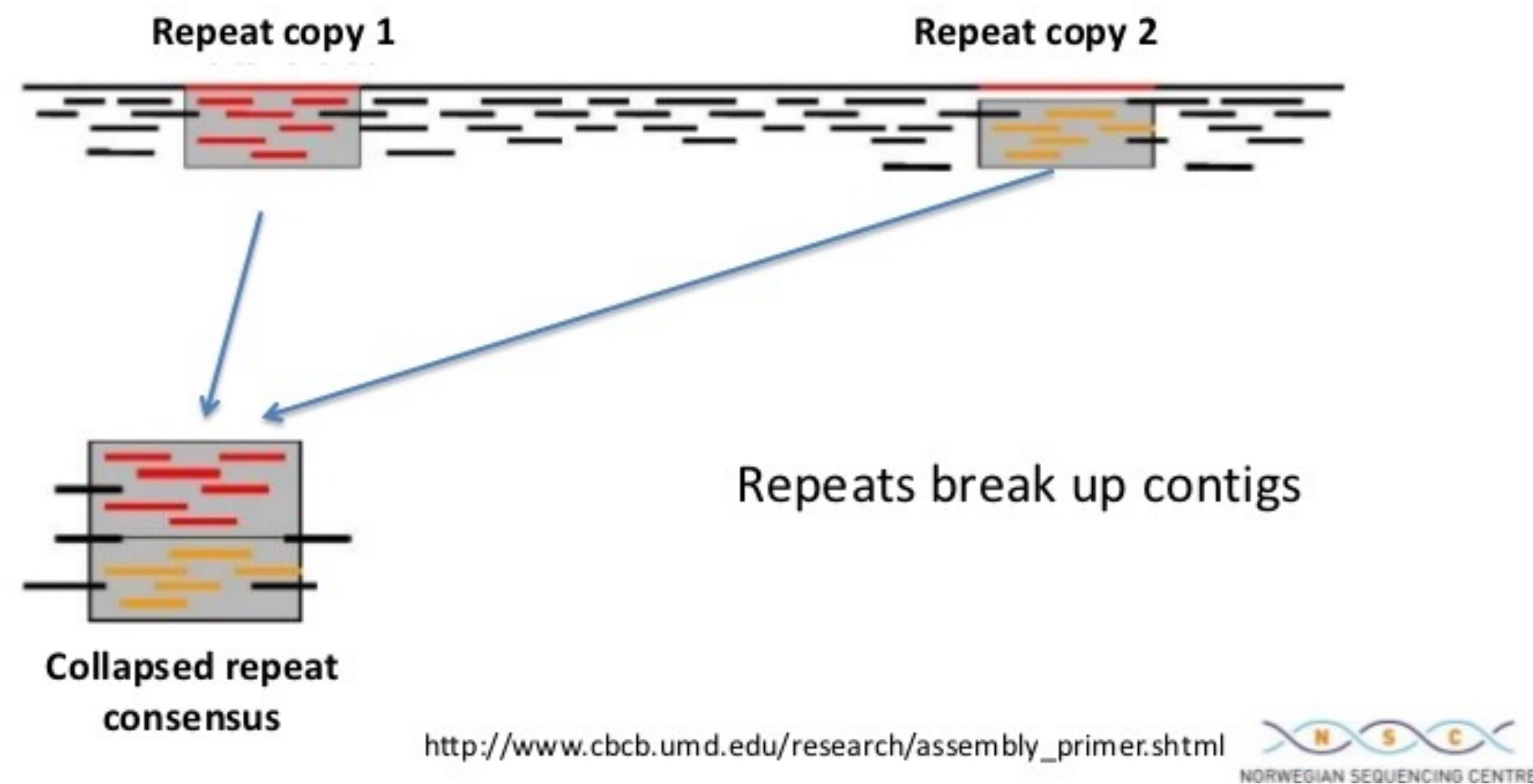


Read lengths: 100-300 bp
Insert lengths: ave 300-500 bp



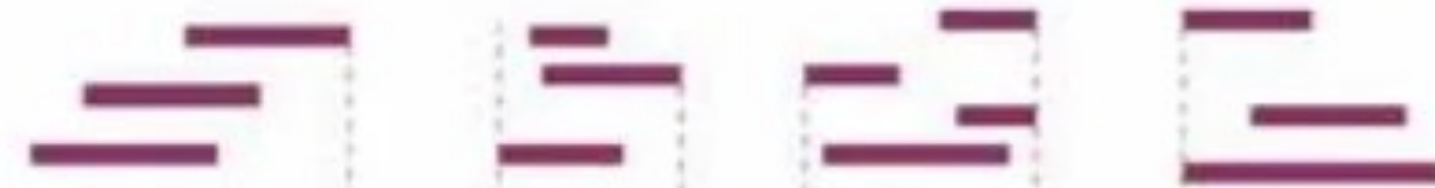
Read lengths: 5-10 kb
- Pacbio: up to 60 kb
- Nanopore: up to 1Mb

Repeats / polymorphic loci can break genomes





Short reads



Little overlap results in gaps
Hard to assemble

Draft genome



Missing sequence leads to missed genes and limits biological interpretation

Long reads



Overlapping reads with no gaps
Easy to assemble

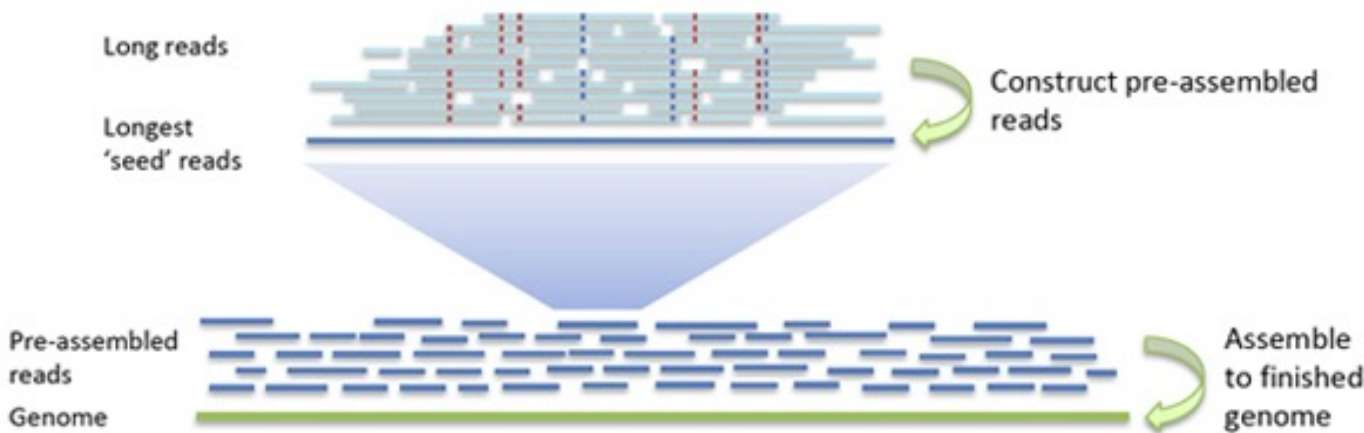
Complete genome



A comprehensive structural, functional, and organisational picture of the genome

Long read / range sequencing is key to good genomes

Pacific Biosciences (PacBio)



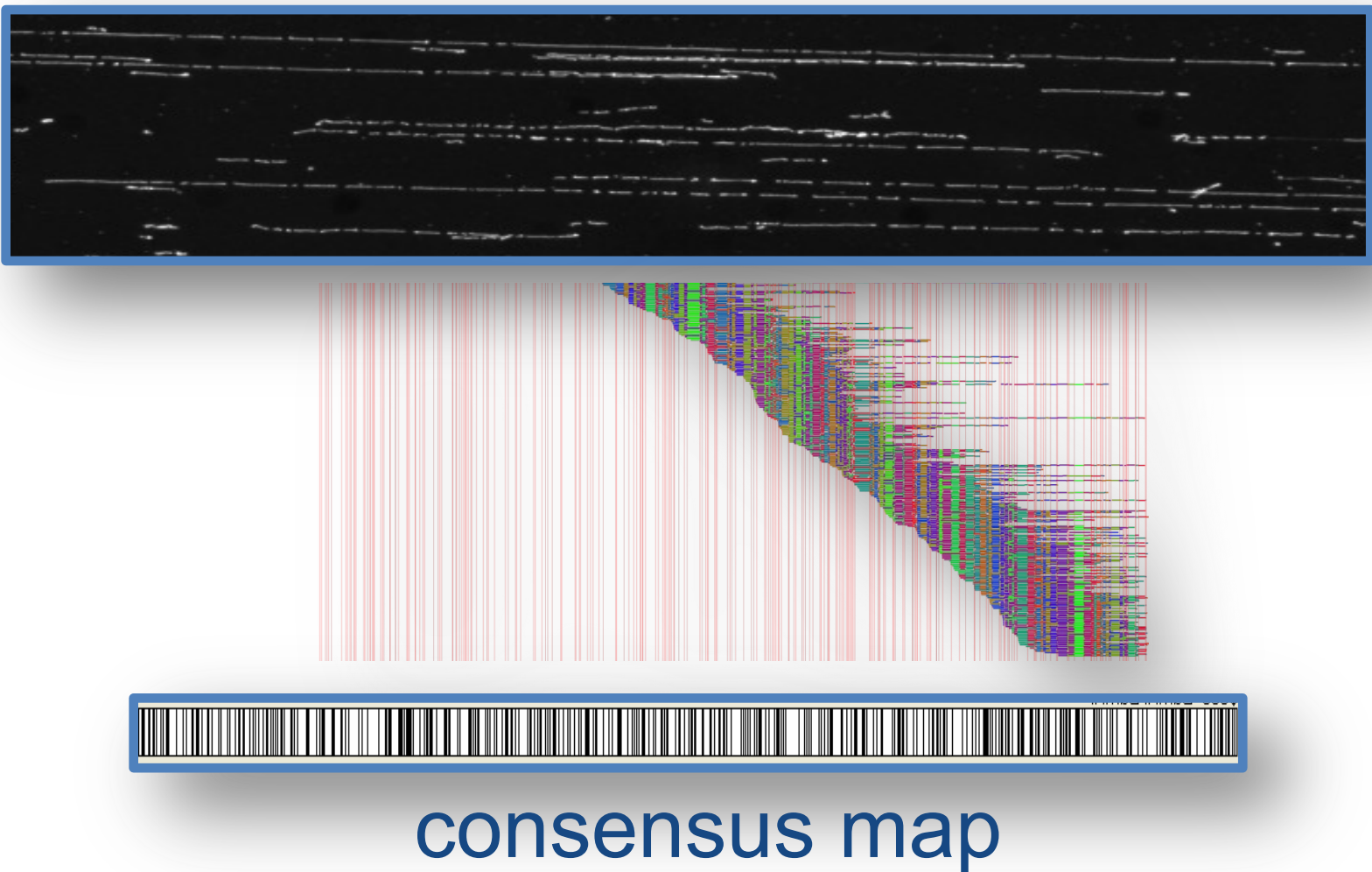
Linked reads (10X Genomics)



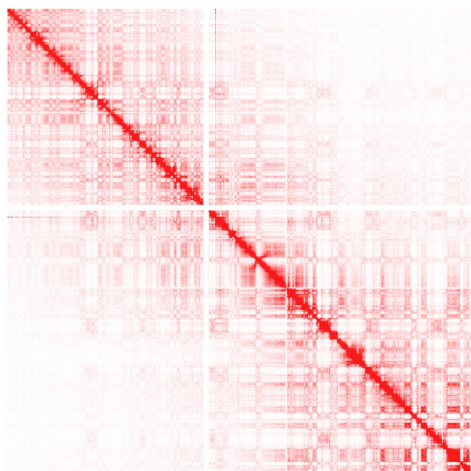
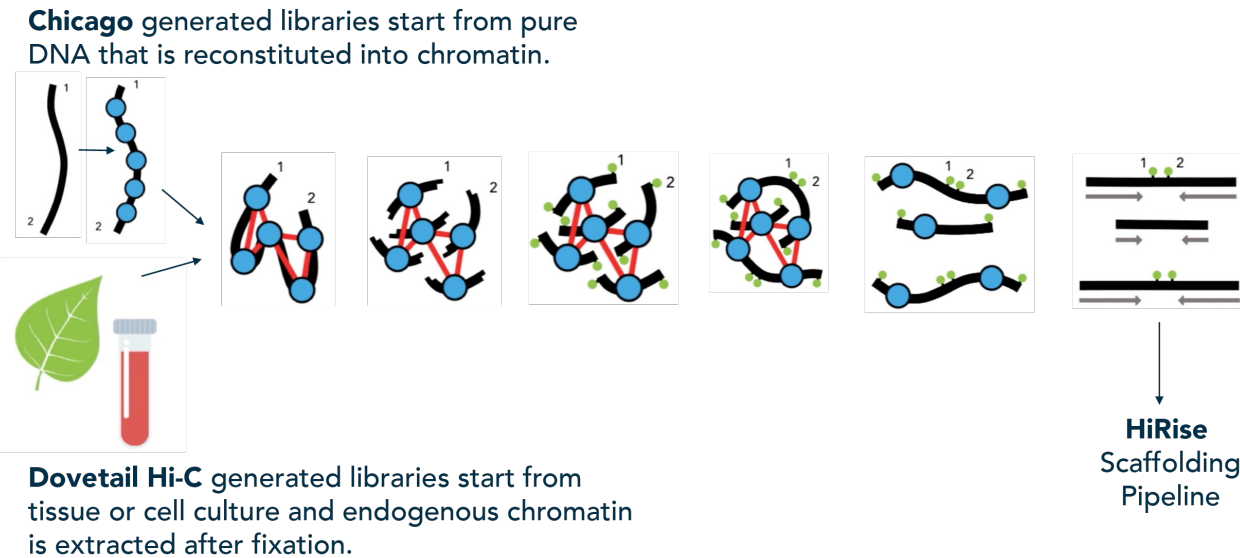
Oxford Nanopore



Optical Mapping (OpGen, Bionano genomics)



Chromosome conformation capture, ie Hi-C (Dovetail Genomics)

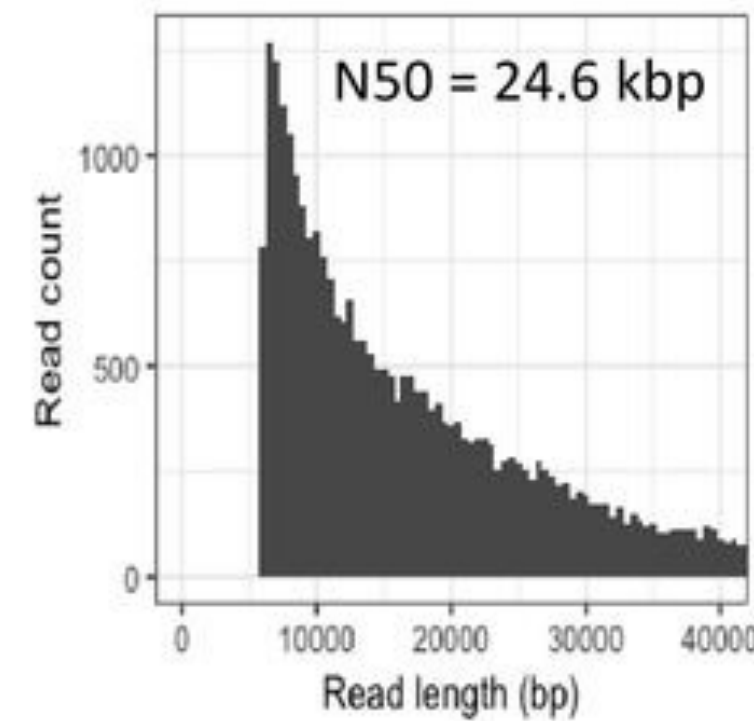







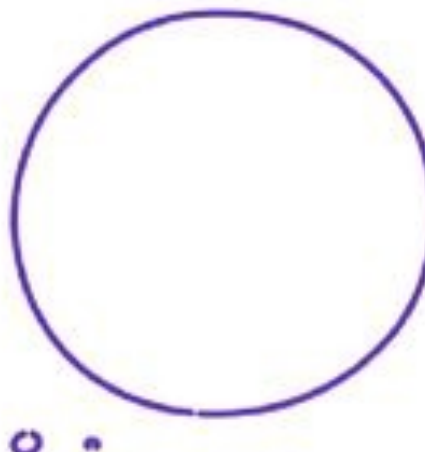
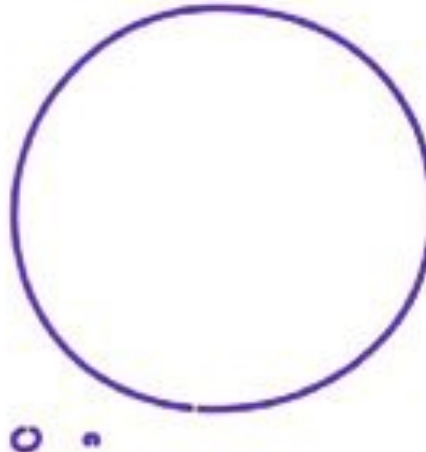

Klebsiella pneumoniae INF116

High Nanopore read depth

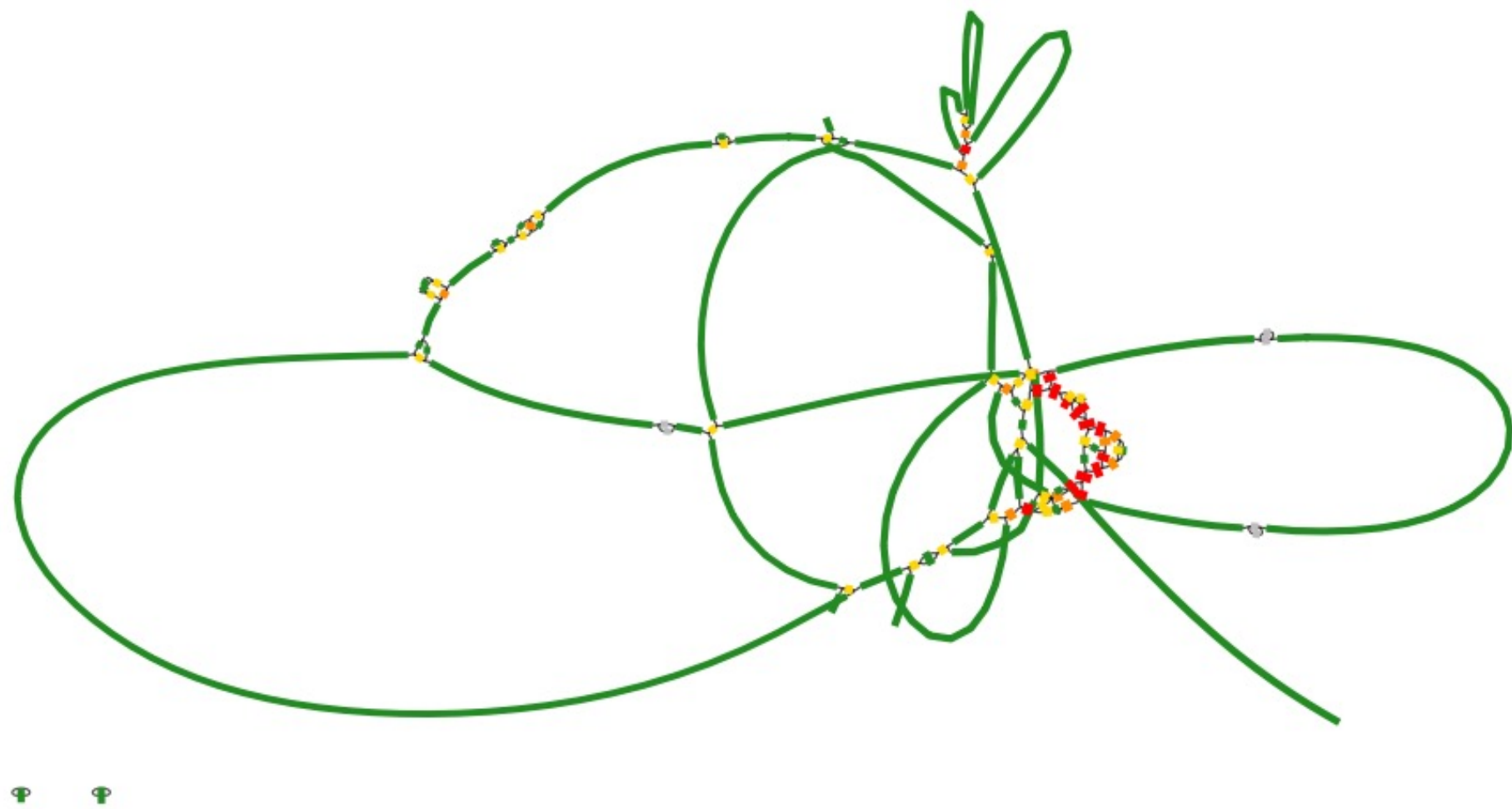
Nanopore reads:

- R9.4, whole flow cell
- Subsampled for length, quality
- 564 Mbp (>100x depth)
- Albacore 1.0.2
- Read length distribution:

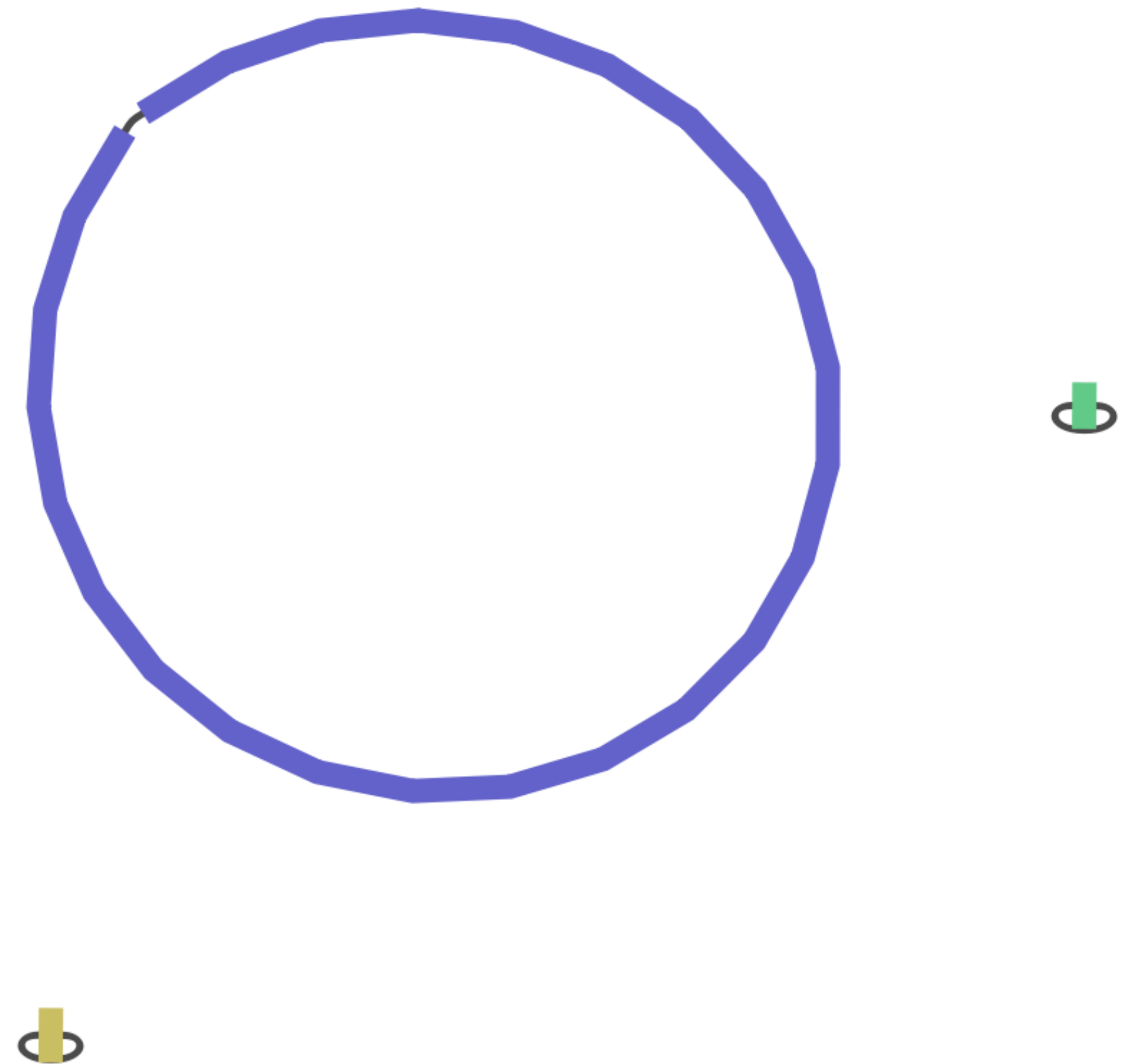


|     | | | | |
|---|--|--|--|----------|
| Illumina-only | Nanopore-only | + Nanopolish | Hybrid | |
|  |  |  |  | |
| 450 kbp | n/a | n/a | n/a | |
| { | Accuracy: | 99.15% | 99.54% | 100% |
| | Bases per error: | 118 | 218 | ∞ |

Klebsiella pneumoniae assemblies



Illumina – Spades



Illumina + Nanopore - Unicycler

Today's Agenda



Use git command to download assemblies



Use quast to assess various quality metrics



Annotate genomes using prokka



View genomes in Artemis

Questions?