

# **Plant Genomics**

**Alex Harkess**

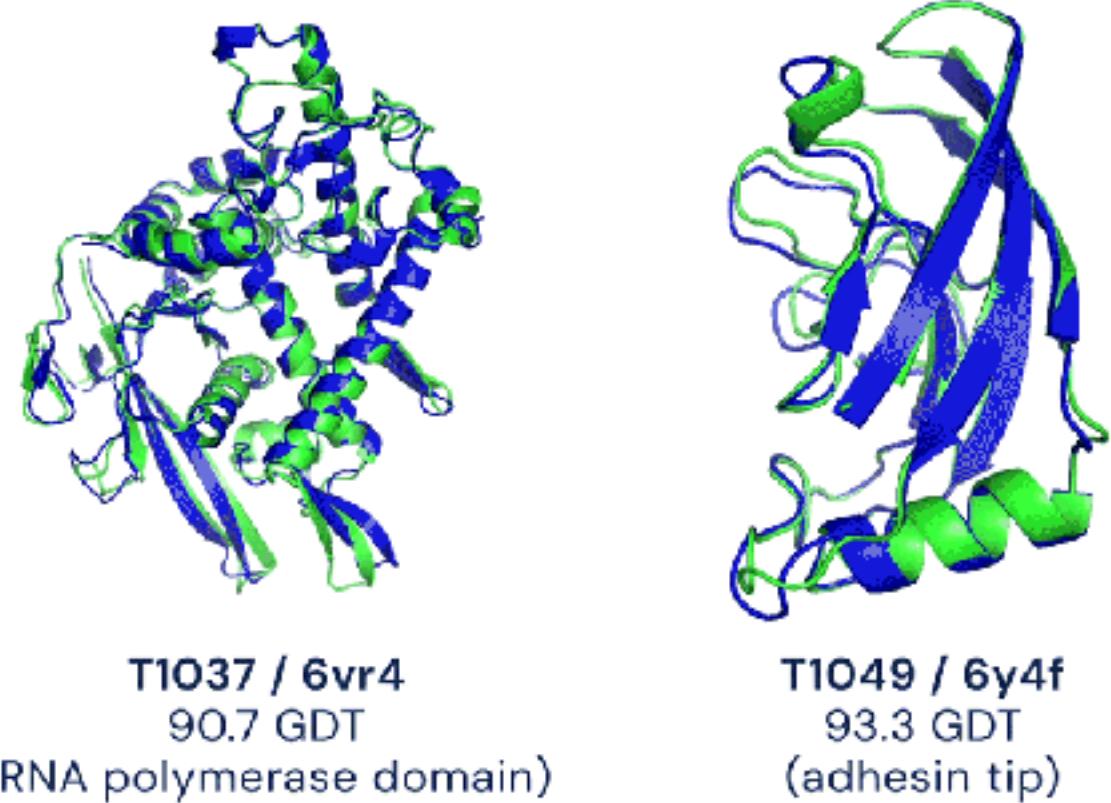
**Omics, genomics, transcriptomics, so many -omics**



# What is computational biology?

The development and application of large-scale datasets to study biological questions.

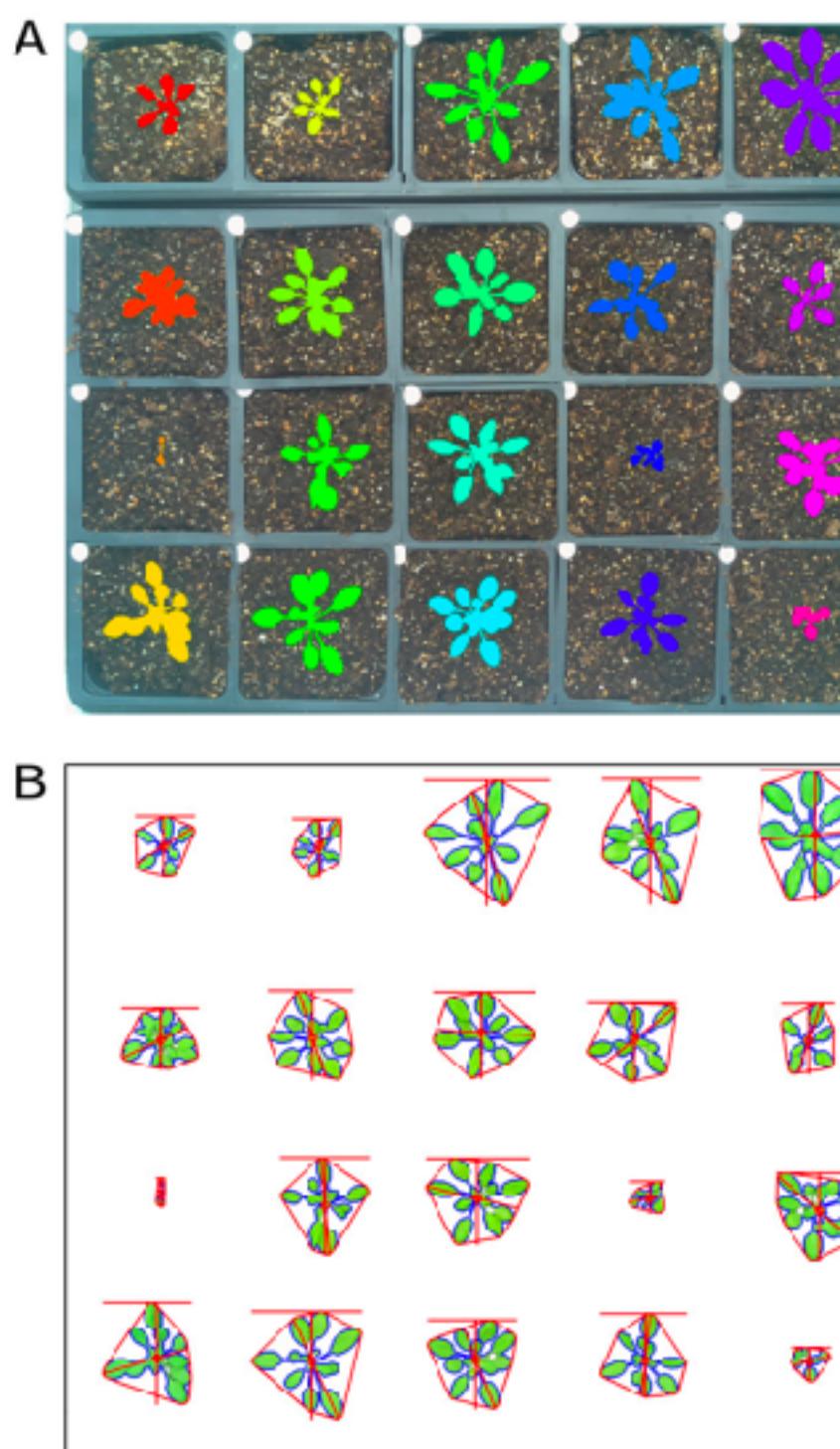
Protein modeling  
(AlphaFold)



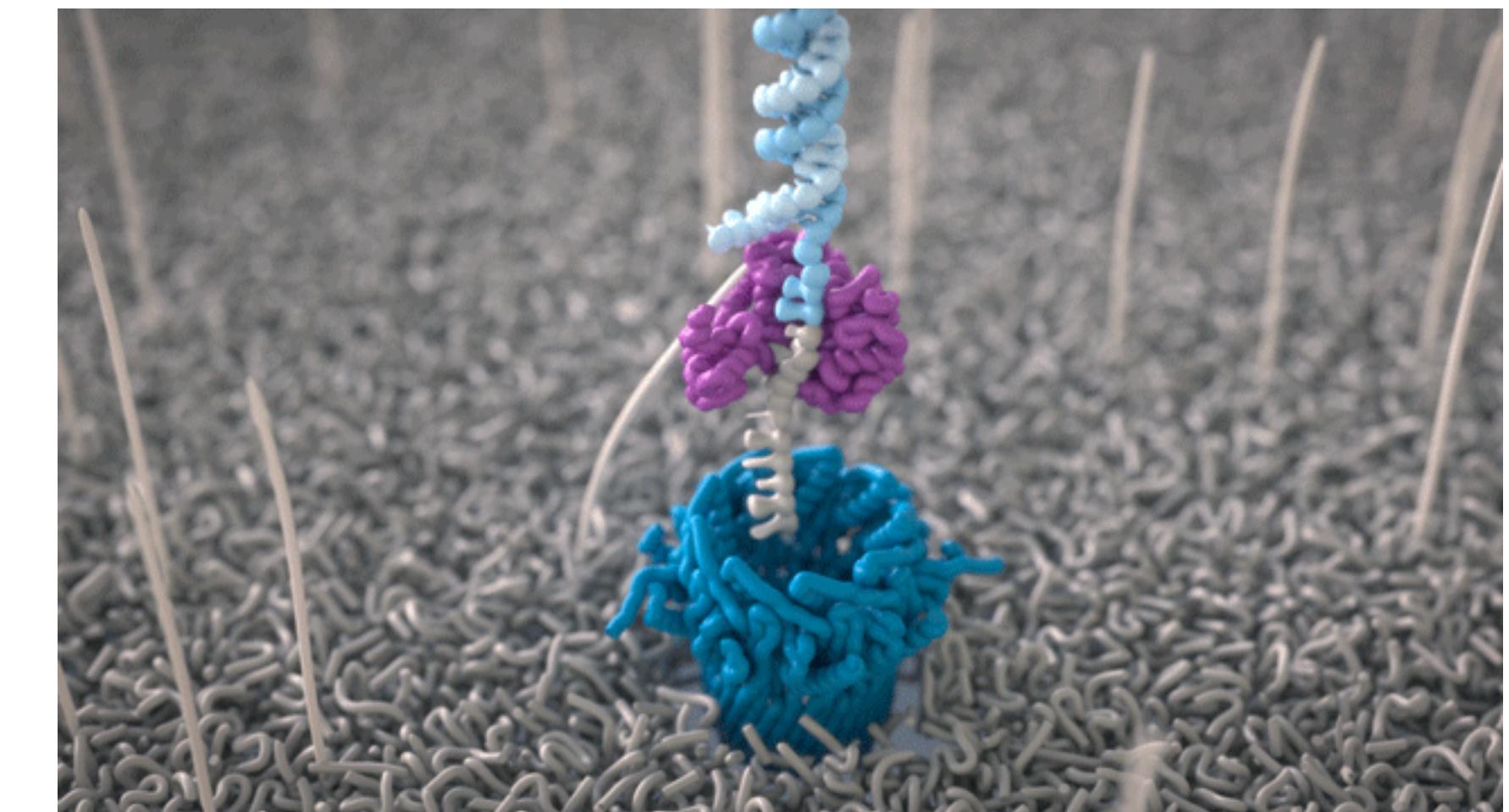
- Experimental result
- Computational prediction

[deepmind.com](https://deepmind.com)

Image analysis  
(PlantCV2)



Nucleic acid sequencing  
(DNA/RNA)



Gehan et al. 2017

You will be a stronger computational biologist  
if you **understand the biology of plants**

Plant genomes are complicated  
but **the technology is getting better**

How do we approach genome projects  
because **plants are absurd organisms**

How do we apply computational genomics  
to **assemble plant genomes**

You will be a stronger computational biologist  
if you **understand the biology of plants**

Plant genomes are complicated  
but **the technology is getting better**

How do we approach genome projects  
because **plants are absurd organisms**

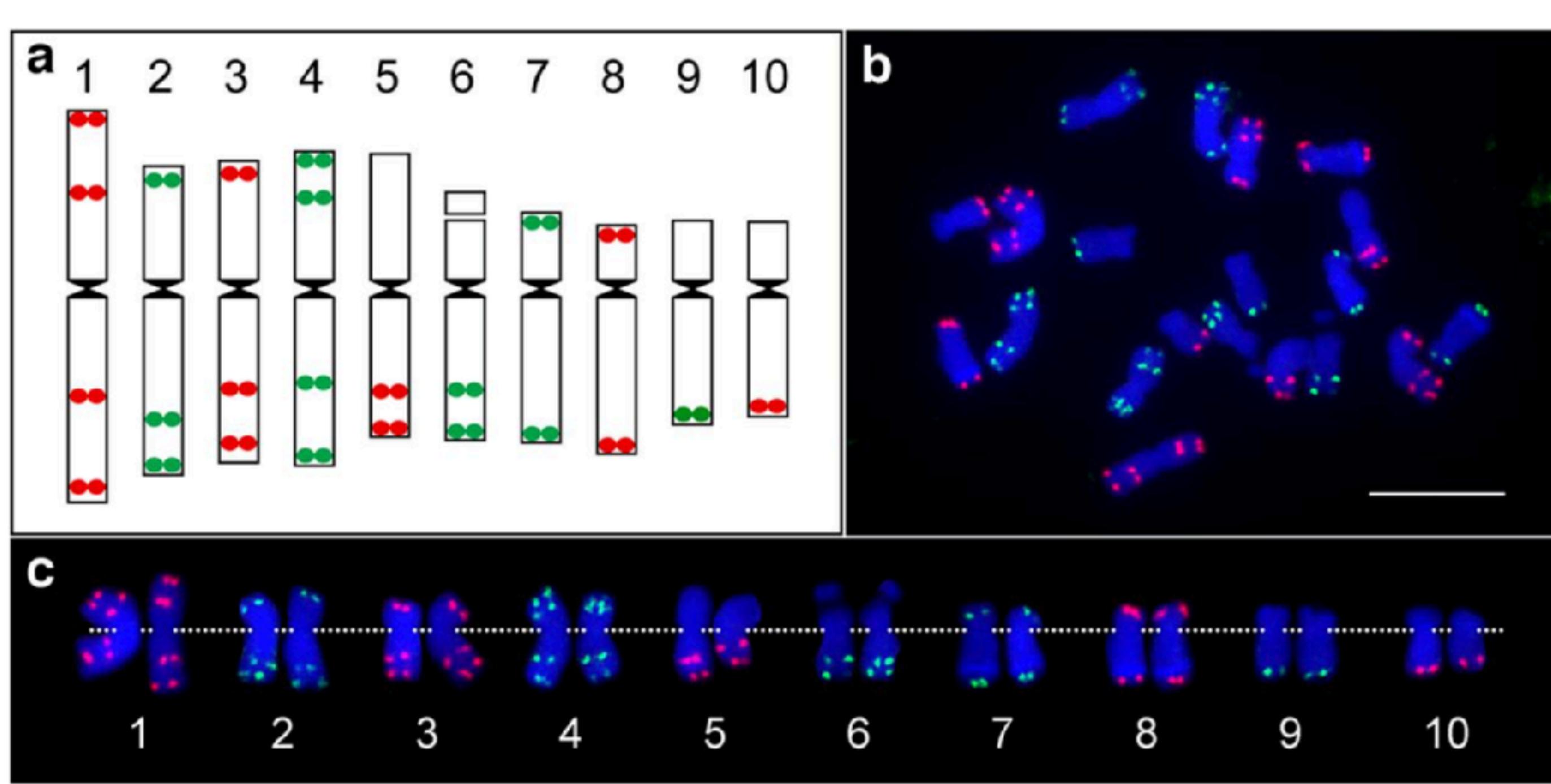
How do we apply computational genomics  
to **assemble plant genomes**

# What is genomics

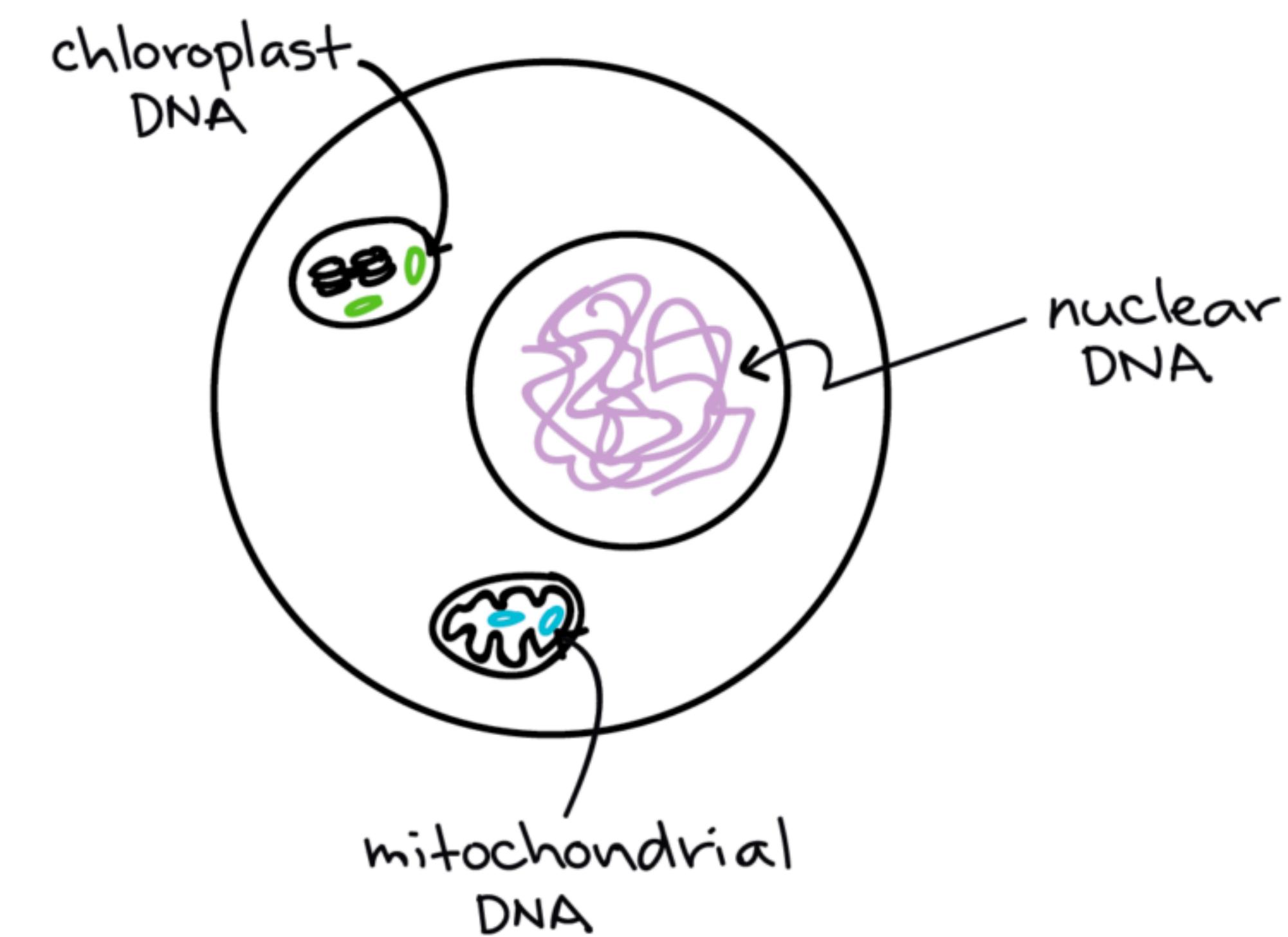
The diagram illustrates the hierarchical structure of genetic material. At the top, a double helix represents DNA (Deoxyribonucleic Acid). A segment of the DNA is labeled 'Gene'. An inset shows 'Base Pairs' with the sequence: T-A, C-G, G-C, T-A, C-G, A-T, G-C, T-A. Below the DNA, a single-stranded loop represents a chromosome. This chromosome is shown wrapped around protein complexes labeled 'Histones' and 'Nucleosomes'. The nucleosomes appear as small yellow circles along the DNA strand. A large, X-shaped structure at the bottom represents a cell during prophase, with chromosomes visible within it. The background features a grid of DNA sequence code.

# Plants have 3 distinct genomes

## Nuclear genome



Braz et al. 2020



# Nuclear genome sizes vary widely in plants



*Utricularia gibba*  
Creeping bladderwort  
80 Mb



*Zea mays*  
Corn  
2.7 Gb



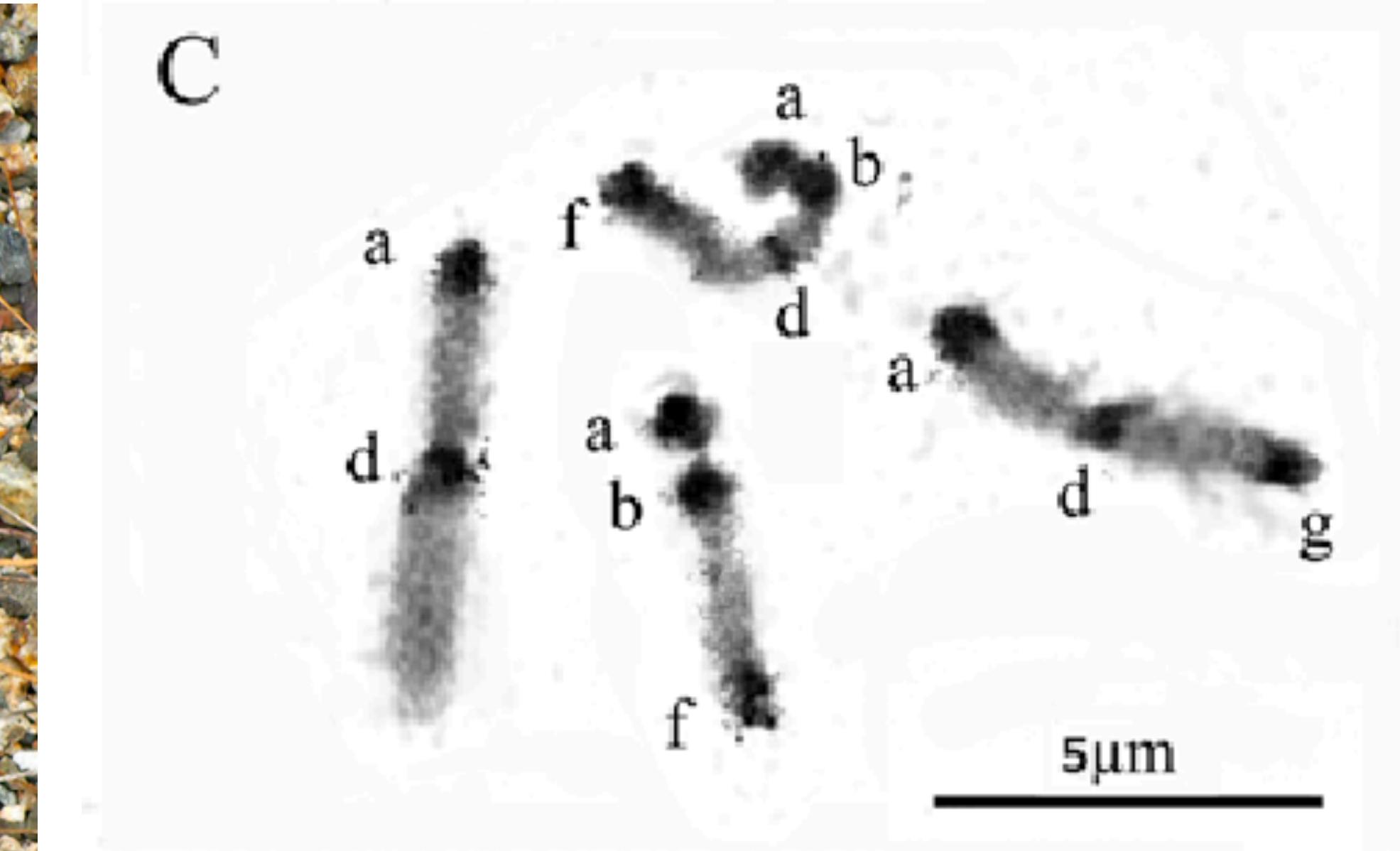
*Paris japonica*  
149 Gb

# ...and in chromosome numbers

**N = 2 chromosome pairs**



Credit: Wikipedia



*Haplopappus gracillis*

$2N = 4$

(2 pairs of chromosomes)

Castiglione et al. 2008

# ...and in chromosome numbers

**N = 630 chromosome pairs**



<http://cookislands.bishopmuseum.org/>



*Ophioglossum reticulatum*  
 $2N = 1260$   
(630 pairs of chromosomes)  
Abraham et al. 1954

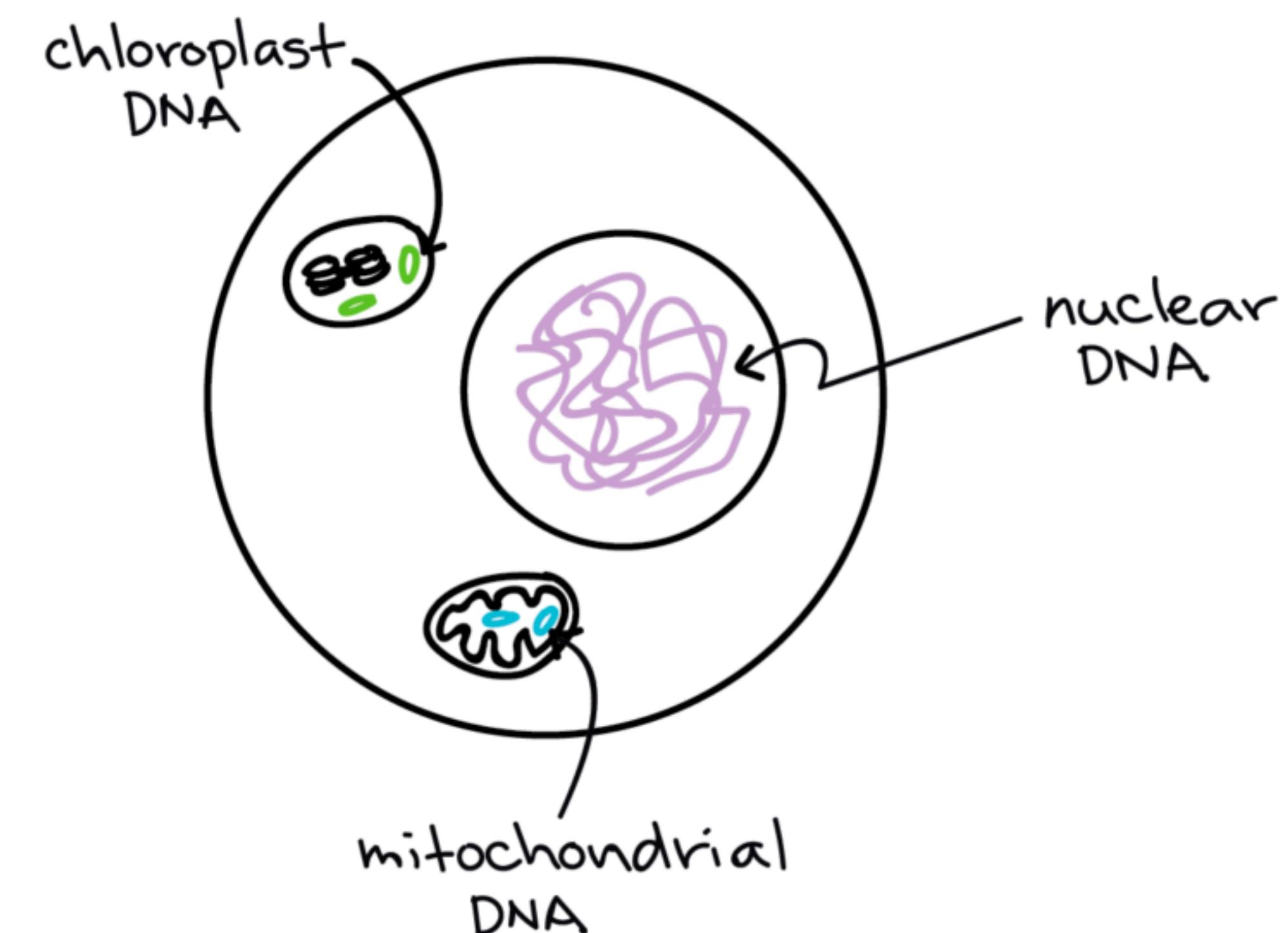
# Plants have 3 distinct genomes

## Chloroplast genome

- Typically ~150 kilobase circular genome
- 80-90 protein coding genes
- ~30 transfer RNA (tRNA) genes
- ~4 ribosomal RNA (rRNA) genes

### Useful for phylogenetics

- maternally inherited
- low substitution rate compared to nuclear genes
- and do not recombine

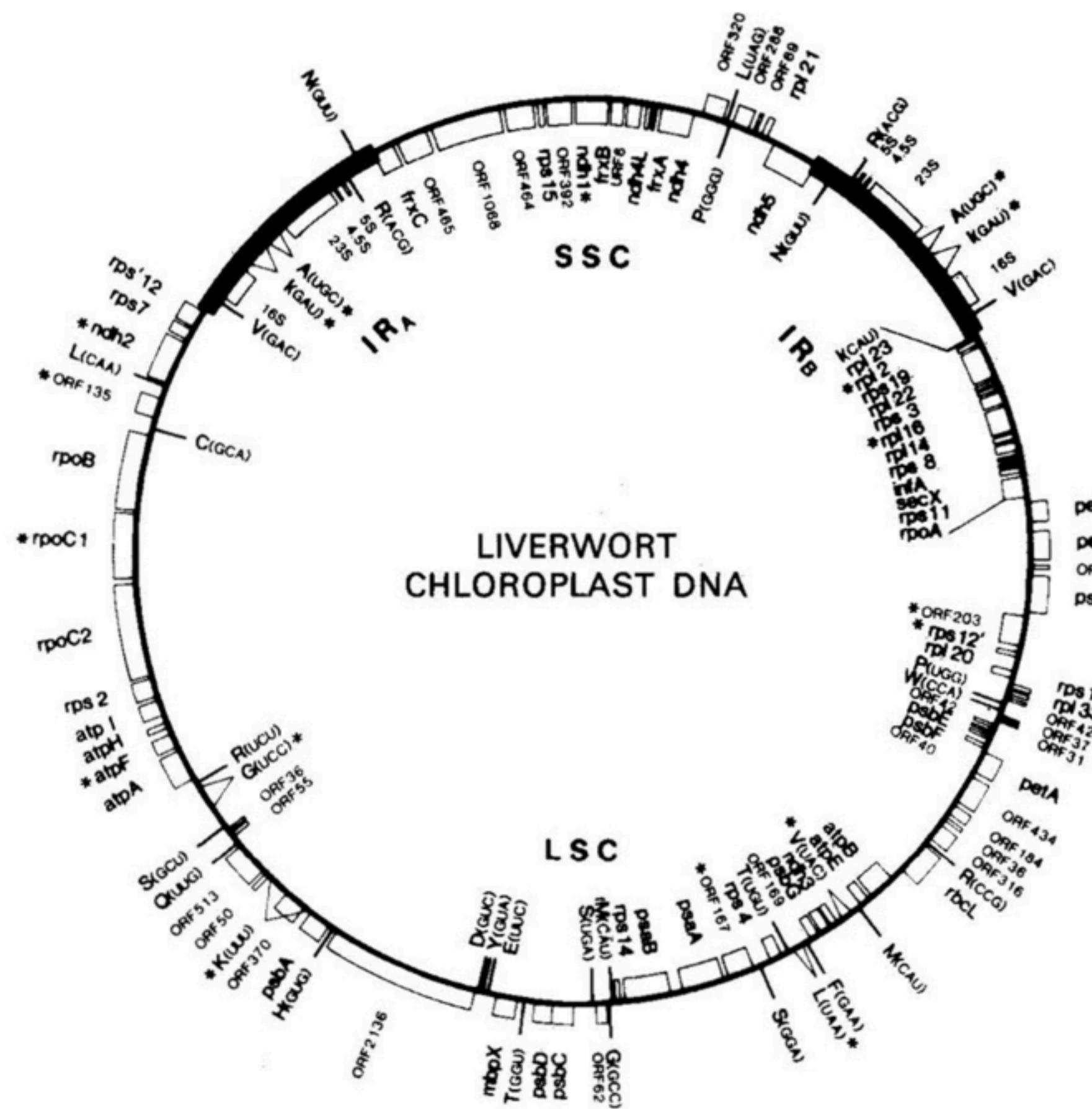


# The first chloroplast genome (1986)

## *Marchantia polymorpha* (a liverwort)

LETTERS TO NATURE

573



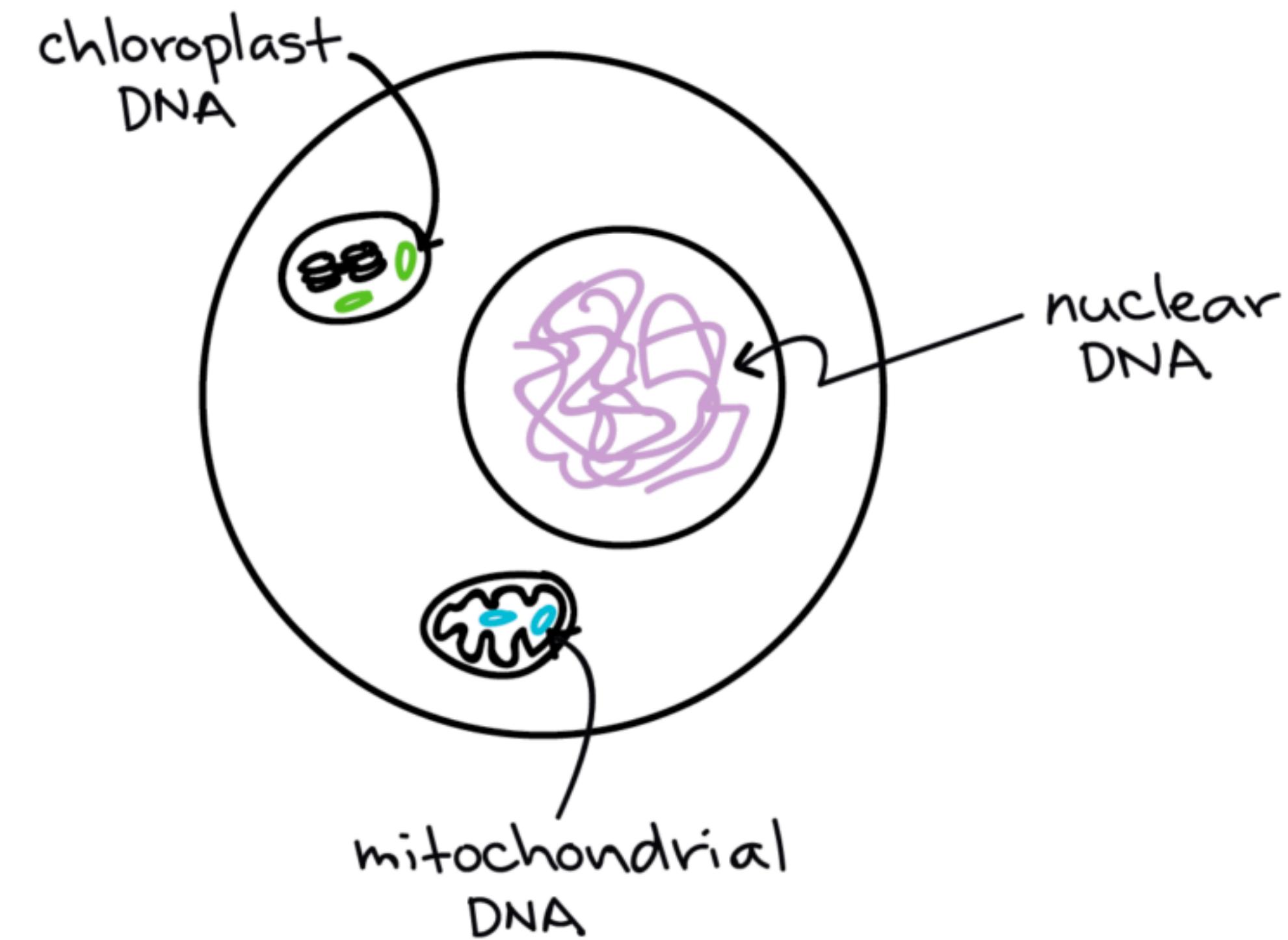
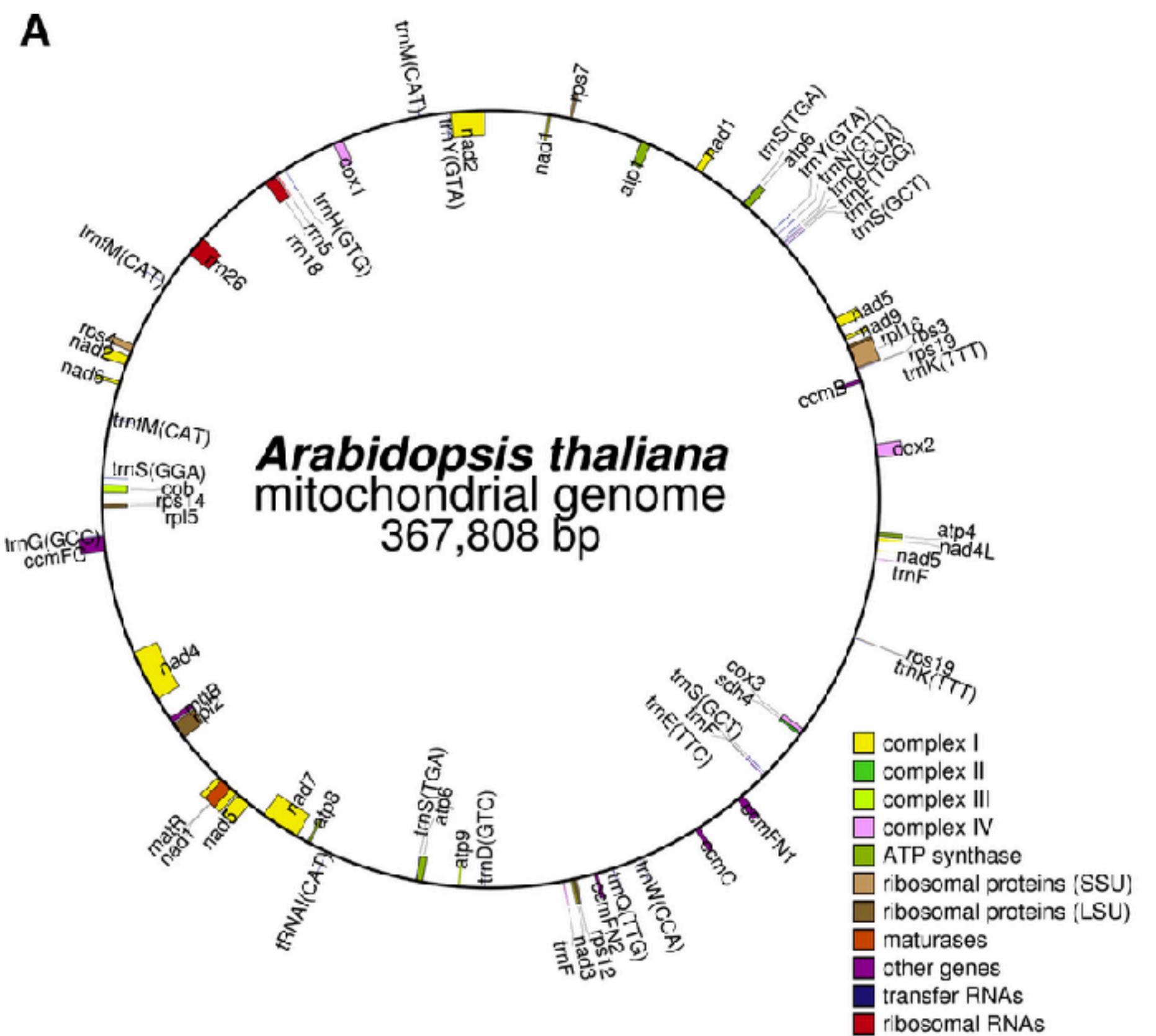
**Table 1** Genes encoded by the chloroplast DNA from a liverwort, *M. polymorpha*

<b>rRNA genes (IR<sub>A</sub> and IR<sub>B</sub>)</b>	<b>Genes for photosynthesis</b>	<b>Genes predicted by amino-acid sequence homology</b>
16S, 23S, 4.5S and 5S	<i>rbcL</i> , large subunit of ribulose bisphosphate carboxylase	<i>ndh1</i> *, homologous to mammalian mitochondrial URF1
<b>tRNA genes</b>	<i>psaA</i> , photosystem I P <sub>700</sub> chlorophyll a apoprotein	<i>ndh2</i> * , URF2 (same as above)
37 tRNAs (see Fig. 1)	<i>psaB</i> (same as above)	<i>ndh3</i> , URF3 (same as above)
<b>RNA polymerase genes</b>	<i>psbA</i> , photosystem II 32K protein	<i>ndh4</i> , URF4 (same as above)
<i>rpoA</i> : homologous to <i>E. coli</i> $\alpha$ -subunit	<i>psbB</i> , photosystem II P <sub>680</sub> chlorophyll a apoprotein	<i>ndh4L</i> , URF4L (same as above)
<i>rpoB</i> : homologous to <i>E. coli</i> $\beta$ -subunit	<i>psbC</i> (same as above)	<i>ndh5</i> , URF5 (same as above)
<i>rpoC1</i> *: homologous to <i>E. coli</i> $\beta'$ -subunit	<i>psbD</i> , photosystem II D2 protein	URF6, homologous to <i>Aspergillus nidulans</i>
<i>rpoC2</i> : homologous to <i>E. coli</i> $\beta'$ -subunit	<i>psbE</i> , cytochrome b <sub>559</sub>	mitochondrial URFC (human)
<b>Ribosomal protein genes and related genes</b>	<i>psbF</i> (same as above)	mitochondrial URF6)
50S subunit	<i>psbG</i> , photosystem II G protein	<i>frxA</i> , homologous to 4Fe-4S type ferredoxin
<i>rpl2</i> *, <i>rpl14</i> , <i>rpl16</i> *,	<i>atpA</i> , ATPase F <sub>1</sub> subunit $\alpha$	<i>frxB</i> (same as above)
<i>rpl20</i> , <i>rpl21</i> , <i>rpl22</i> ,	<i>atpB</i> , ATPase F <sub>1</sub> subunit $\beta$	<i>frxC</i> , homologous to 4Fe-4S
<i>rpl23</i> , <i>rpl33</i>	<i>atpE</i> , ATPase F <sub>1</sub> subunit $\epsilon$	protein found in <i>R. capsulata</i>
	<i>atpF</i> *, ATPase F <sub>0</sub> subunit I	<i>mbpX</i> , homologous to ATP-binding
	<i>atpH</i> , ATPase F <sub>0</sub> subunit III	subunit of inner membrane
	<i>atpI</i> , ATPase F <sub>0</sub> subunit IV	permease in <i>E. coli</i> ( <i>malK</i> ) or
<b>Others</b>	<i>petA</i> , cytochrome <i>f</i>	<i>Salmonella typhimurium</i> ( <i>hisP</i> ) <sup>29</sup>
<i>infA</i> , <i>secX</i>	<i>petB</i> *, cytochrome b <sub>6</sub>	
	<i>petD</i> *, subunit 4 of cytochrome b <sub>6</sub> / <i>f</i> complex	
		<b>Unidentified genes</b>
		More than 28 ORFs

Ohyama et al. 1986 *Nature*

# Plants have 3 distinct genomes

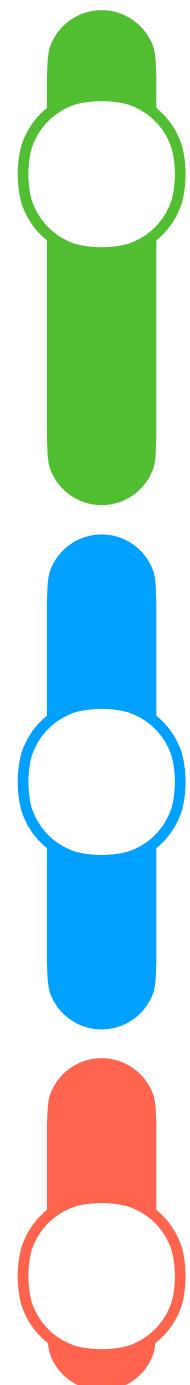
# Mitochondrial genome



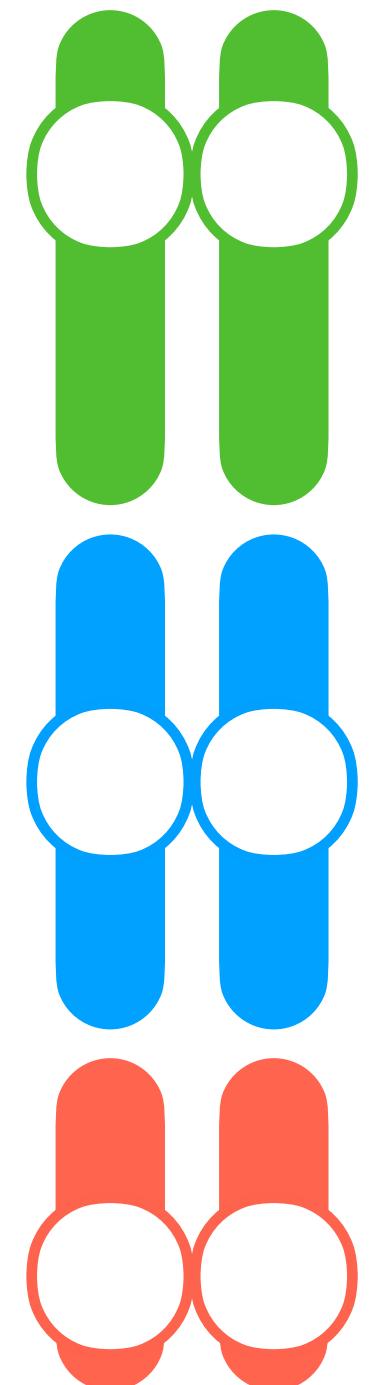
Johnston 2019

# Plants vary in ploidy

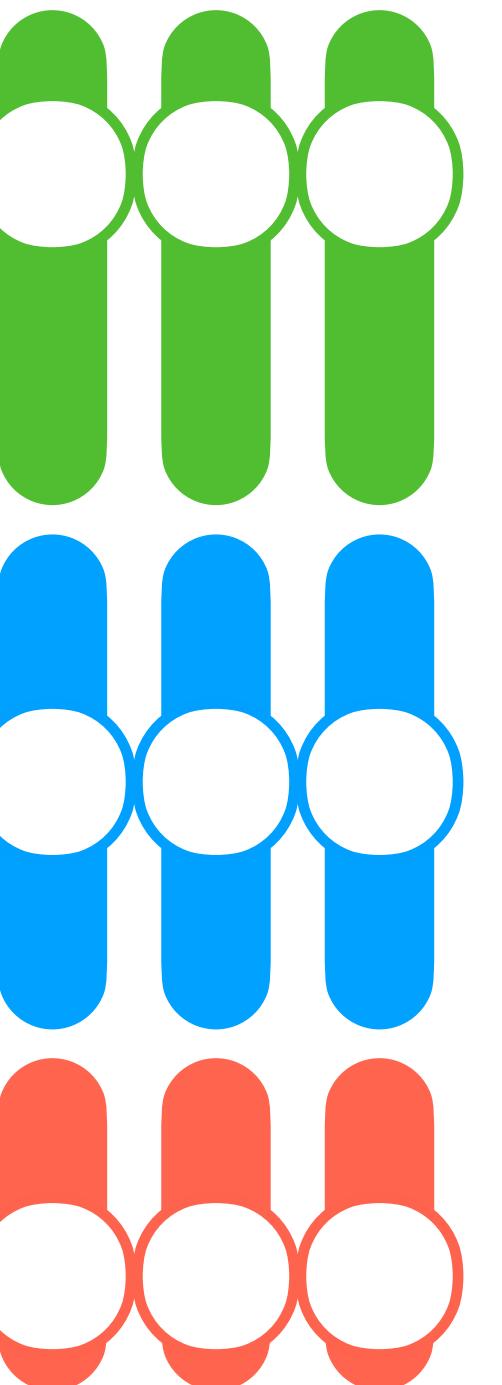
Haploid ( $n$ )



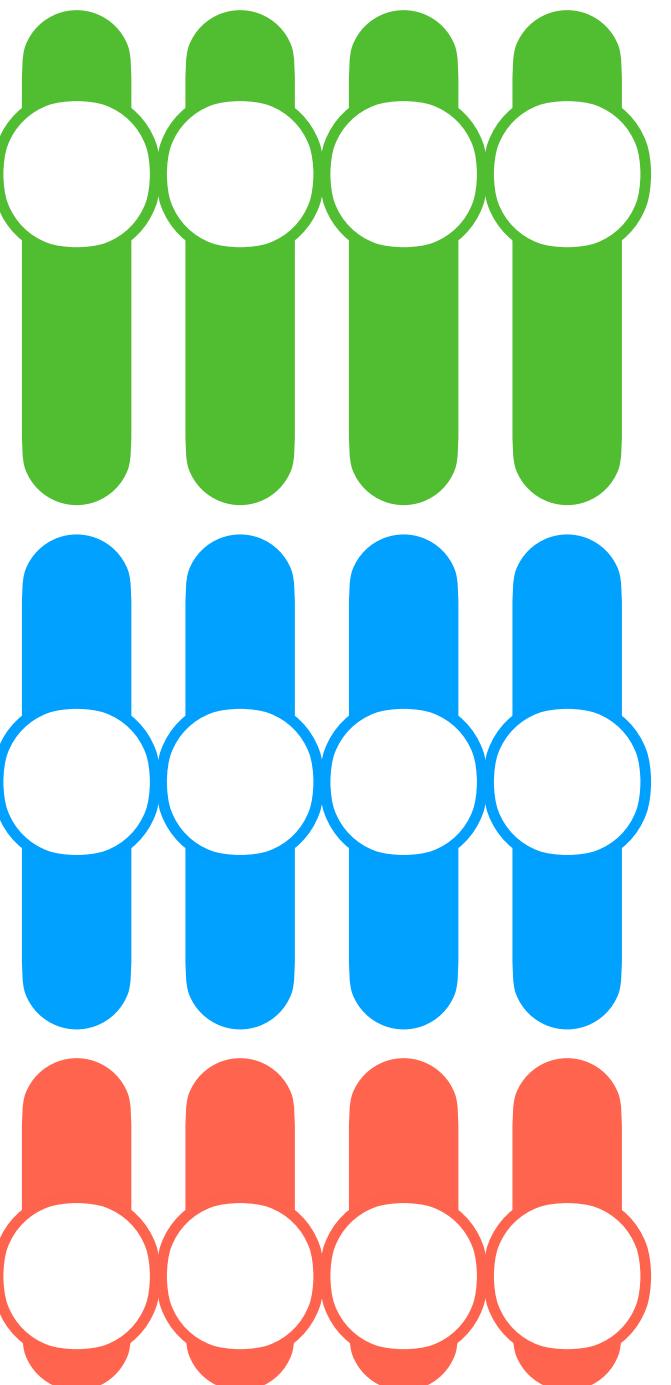
Diploid ( $2n$ )



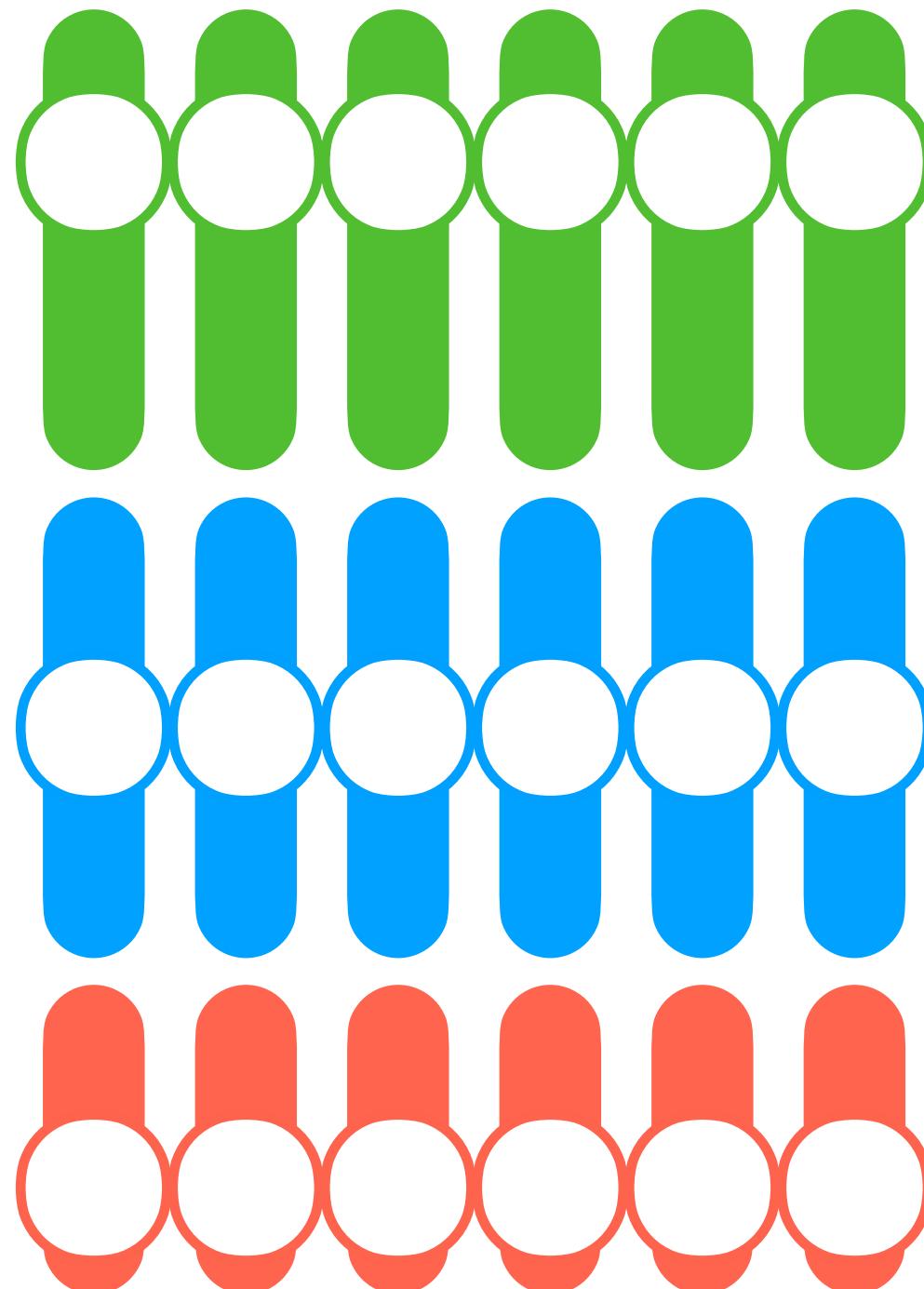
Triploid ( $3n$ )



Tetraploid ( $4n$ )

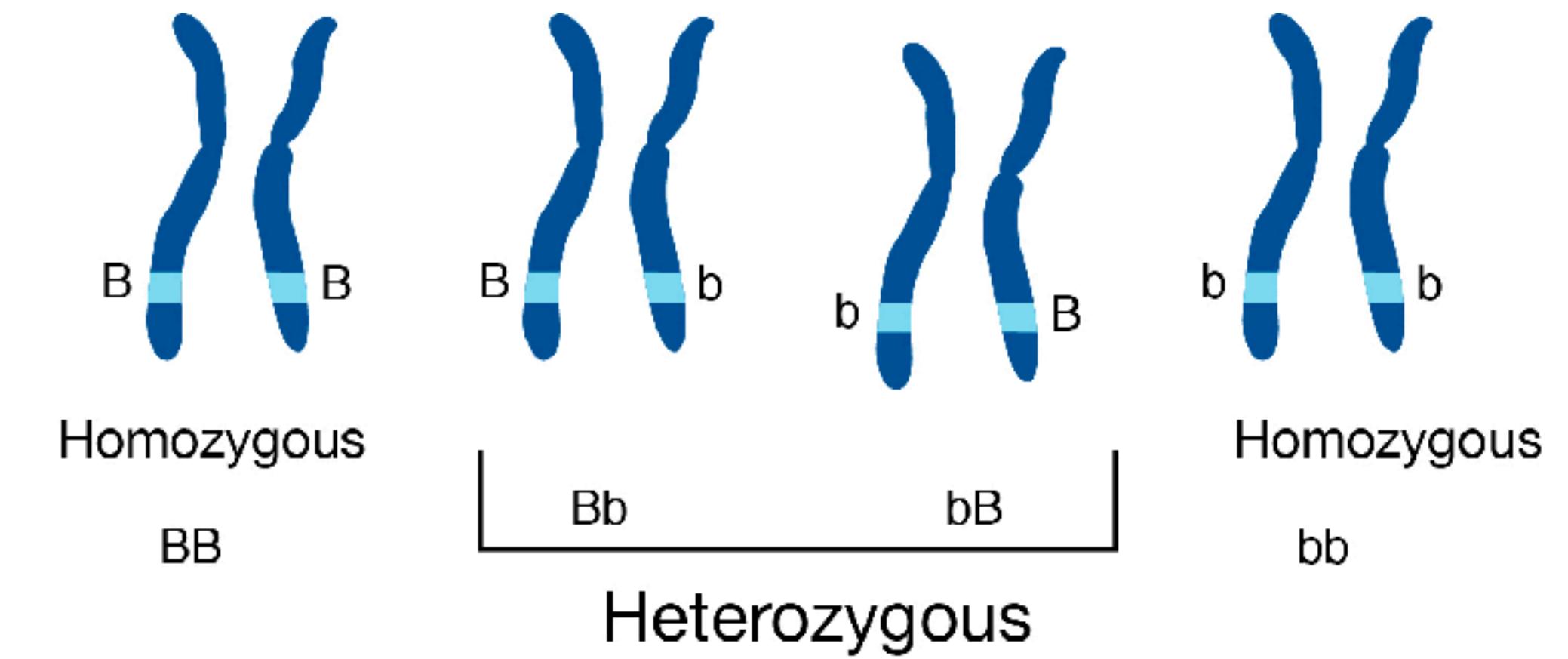


Hexaploid ( $6n$ )



# Plants vary in heterozygosity

- **Zygosity** = the degree to which both copies of a chromosome have the same allele
- **Homozygous** = one locus has the same allele on both chromosomes
- **Heterozygous** = one locus has different alleles
- **Hemizygous** = one locus has presence/absence variation of an allele

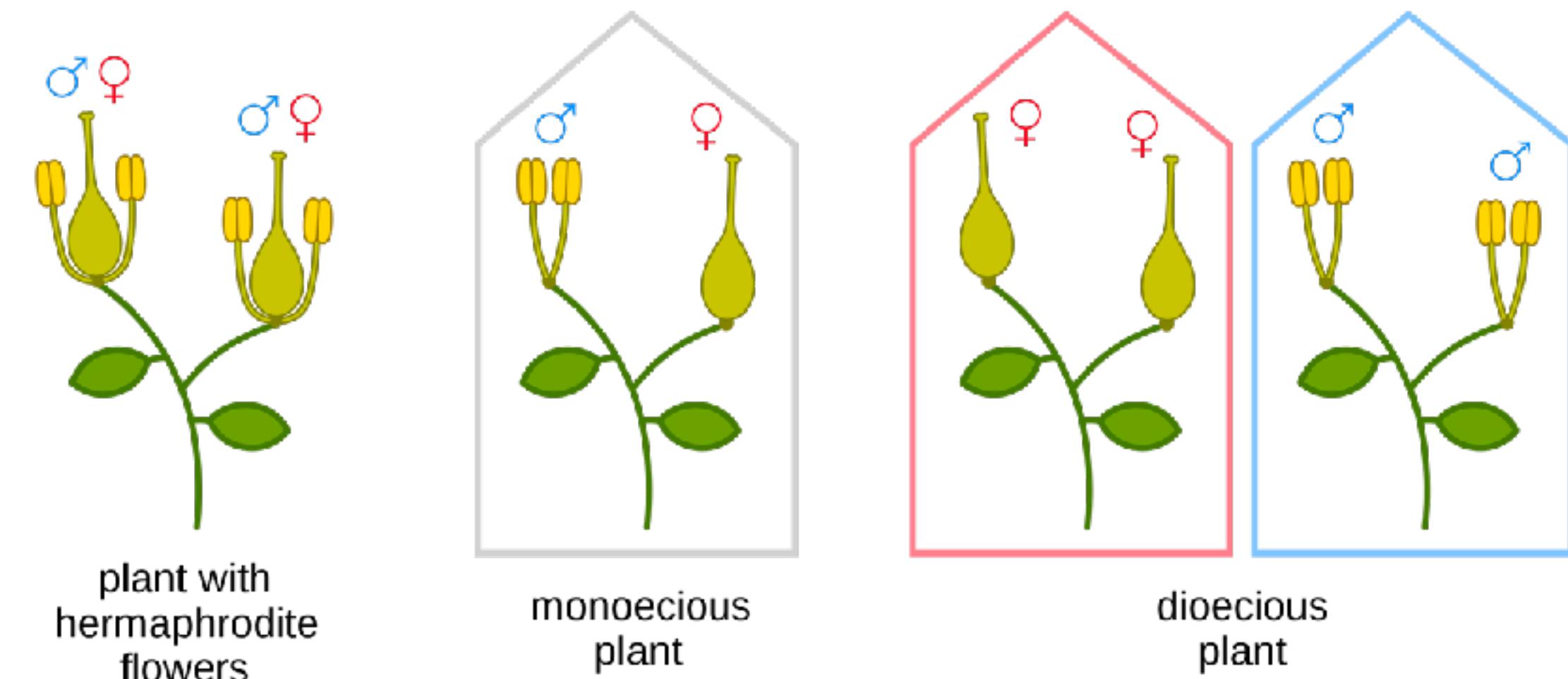


Credit: Wikipedia

- **B** and **b** are two alleles of the same gene on homologous chromosomes
- **Allele** = a form of a gene; a variant

# Heterozygosity relates to reproductive strategies

- Flowering plants (angiosperms) reproduce in a huge diversity of ways
- Hermaphrodite flowers
- Self-compatible, self-incompatible
- Obligate outcrossers (dioecious)
- And all combinations within!



Credit: Wikipedia

You will be a stronger computational biologist  
if you **understand the biology of plants**

Plant genomes are complicated  
but **the technology is getting better**

How do we approach genome projects  
because **plants are absurd organisms**

How do we apply computational genomics  
to **assemble plant genomes**

# ***De novo genome assembly***

The process of reconstructing the original DNA sequence with only the sequence **read** data



**Like a jigsaw puzzle...kind of.**

Lots of pieces

Easy to find the edges and corners

**But...**

Many pieces look the same

Pieces have 2 sides (forward/reverse strandedness)

Diploid organisms = 2 puzzles in same box

(But 95%+ of the pieces are identical)

There might not be a cover

Some of the pieces are wrong

Some of the pieces are from another puzzle

# Modern data generation is massive

As of 2021...



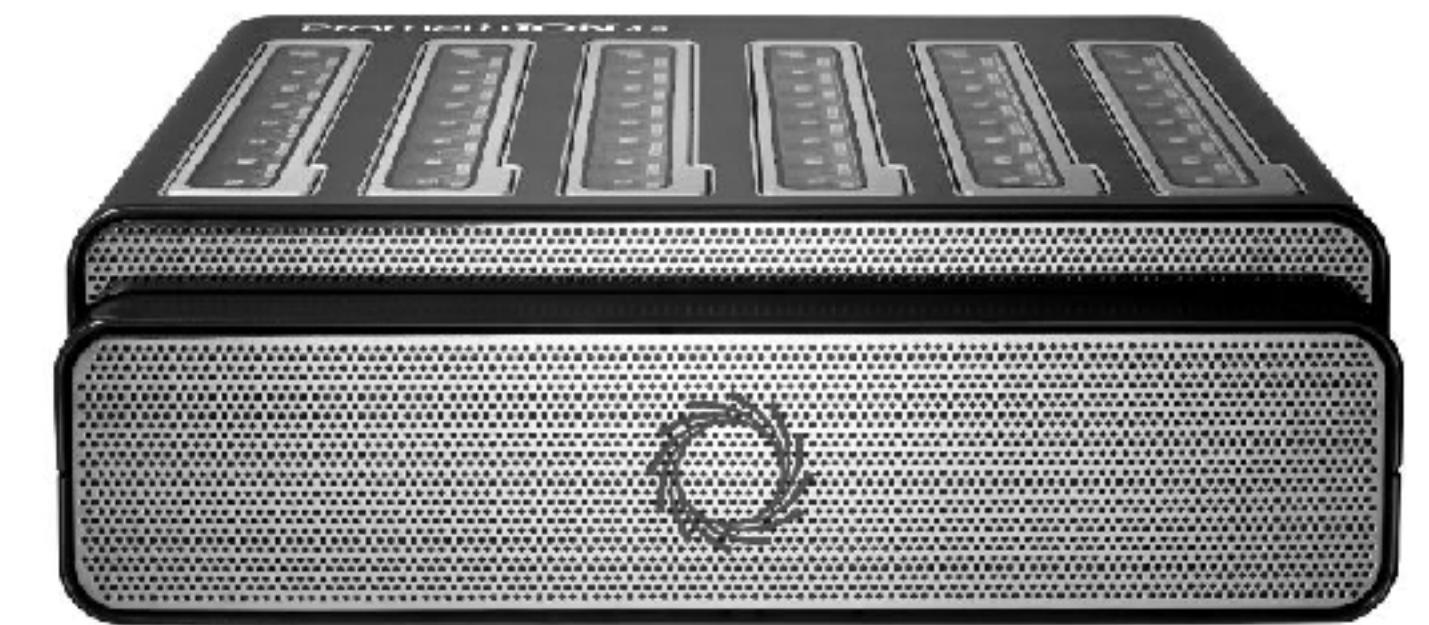
**Illumina NovaSeq6000**

Around \$1M instrument cost  
20 billion reads in 48 hours  
\$20,000 per flow cell run  
48 human genomes



**PacBio Sequel IIe**

Around \$500k instrument cost  
4 million reads in 30 hours  
\$3,500 per flow cell  
0.5 human genomes



**Oxford Nanopore PromethION**

Around \$500k instrument cost  
4 million reads in 30 hours  
\$2,000 per flow cell  
2 human genomes

# Modern data generation is massive

As of 2021...



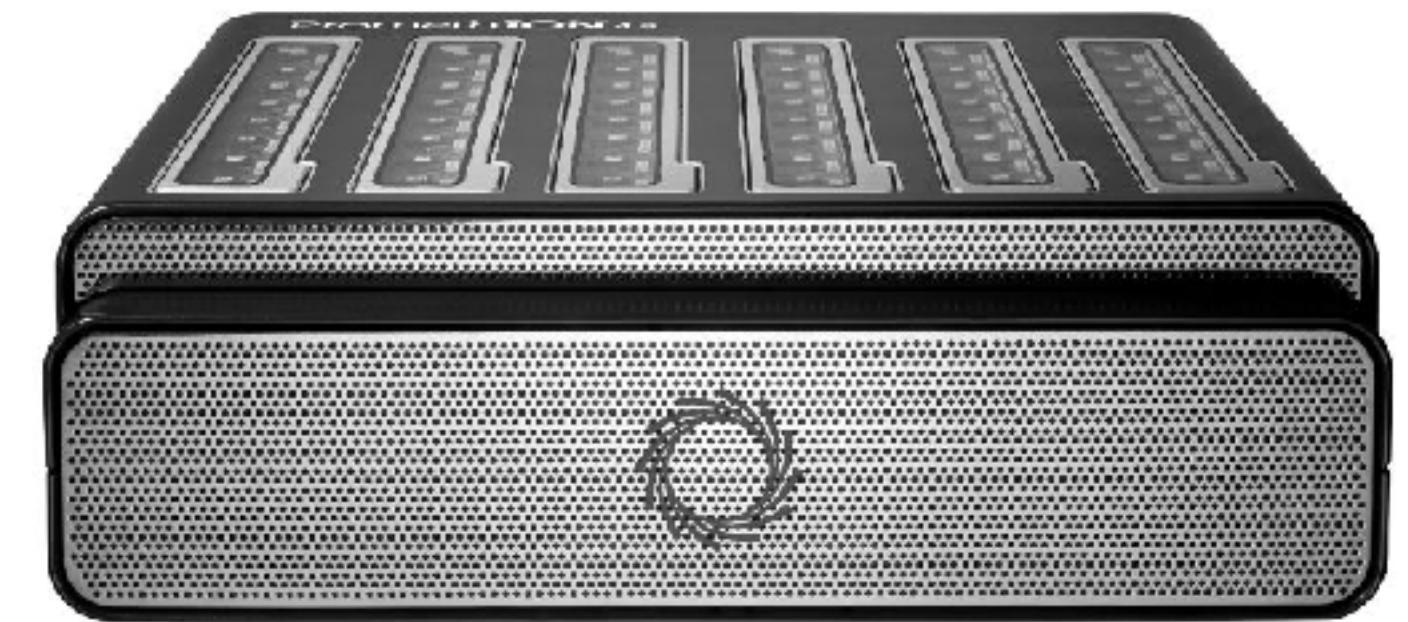
Illumina NovaSeq6000

Tiny reads  
accurate  
billions of reads  
super cheap (\$)



PacBio Sequel IIe

Huge reads  
super accurate  
millions of reads  
and expensive (\$\$\$)

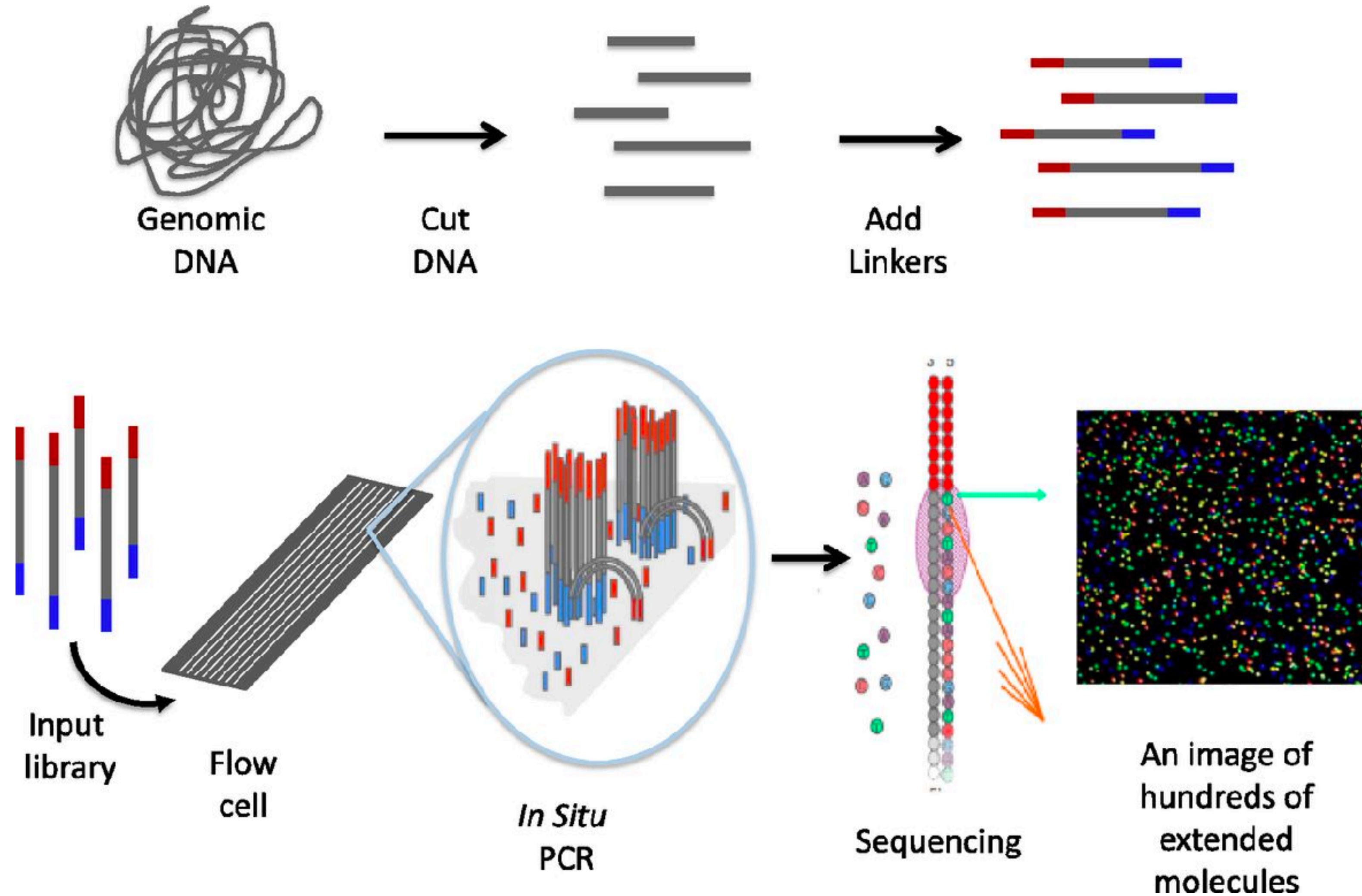


Oxford Nanopore  
PromethION

Huge reads  
accurate-ish...  
millions of reads  
inbetween (\$\$)

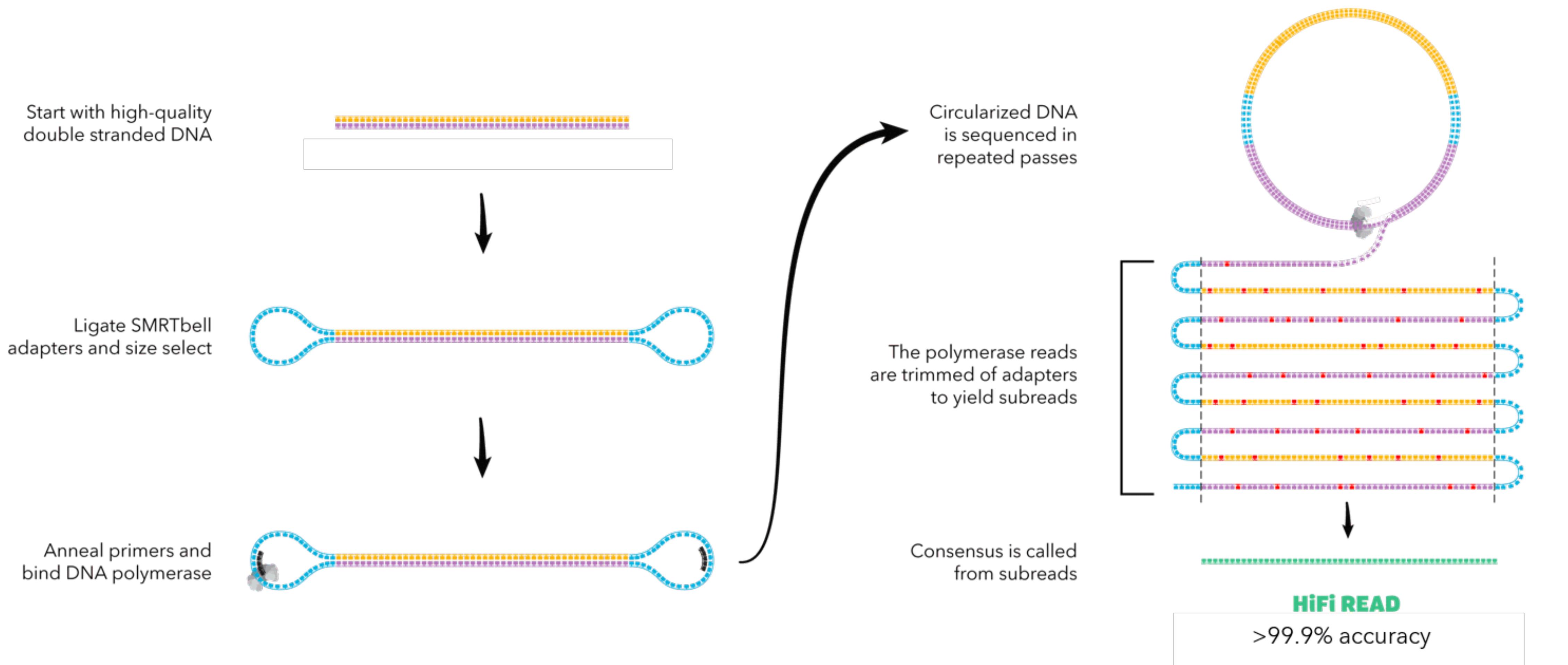
# Illumina sequencing

Short reads (~150 nucleotides)



# PacBio HiFi sequencing

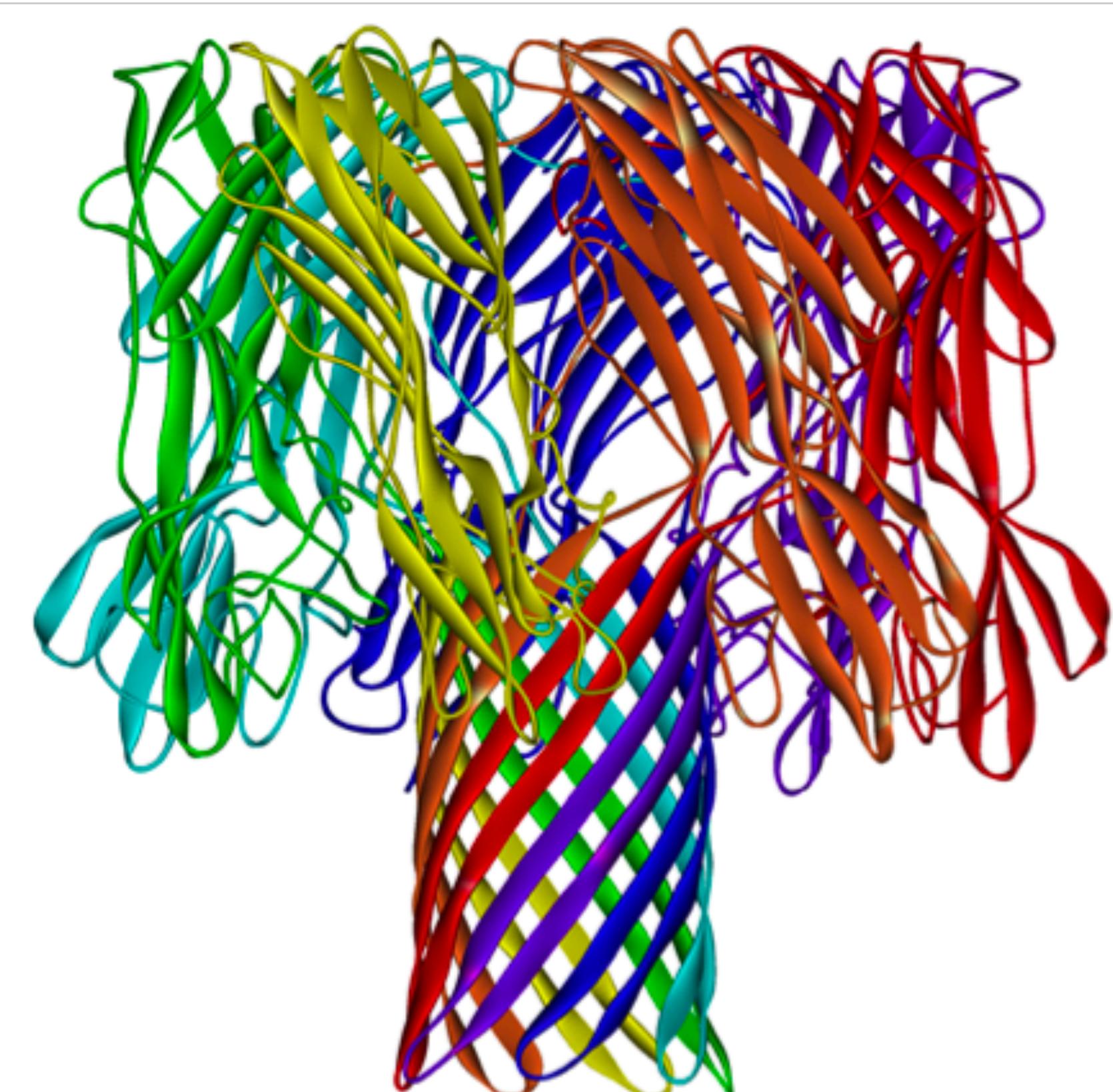
## Make bigger pieces



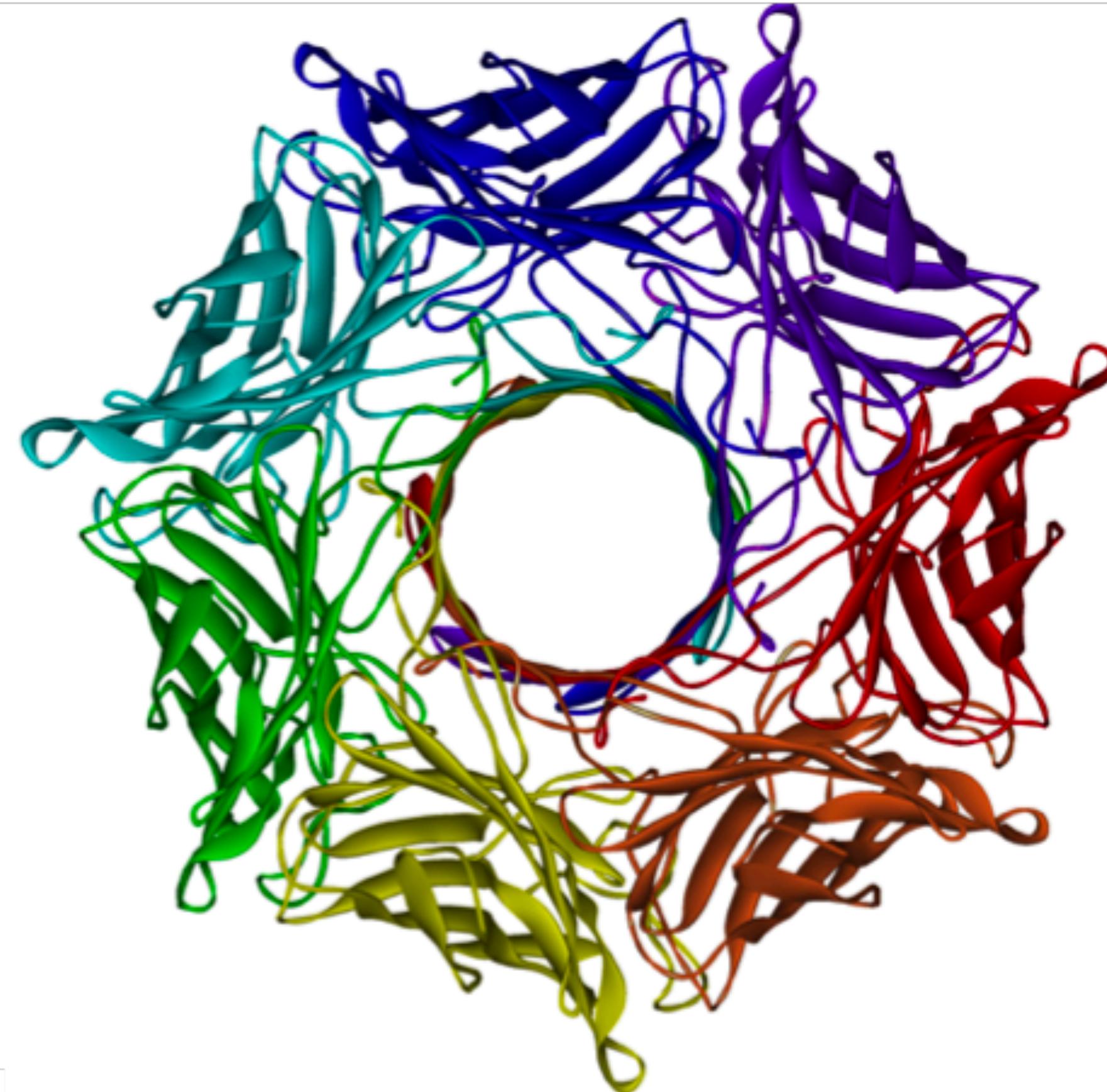
# Oxford Nanopore sequencing

**Make bigger pieces**

Side view



Top view

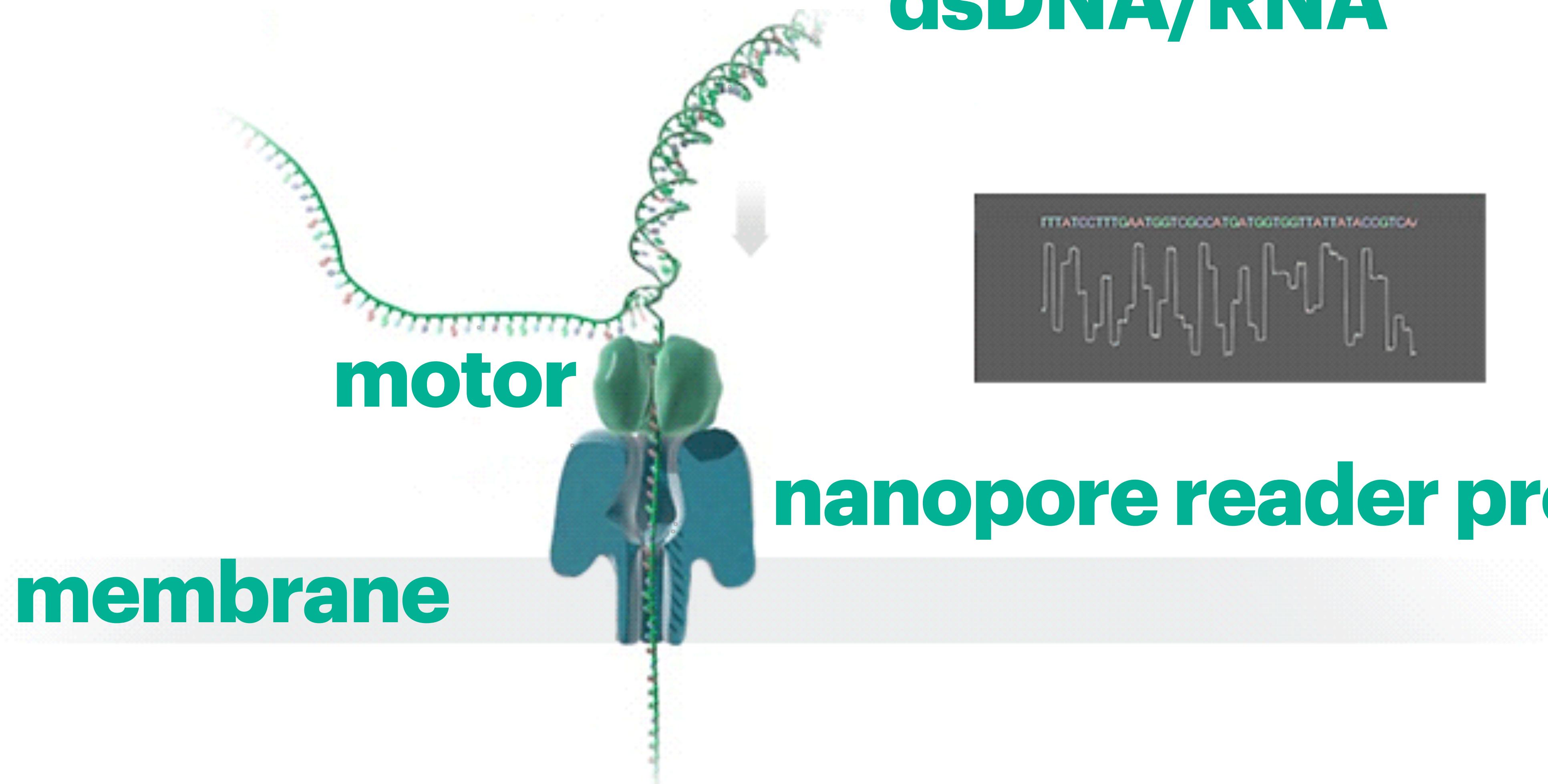


alpha hemolysin

# Oxford Nanopore sequencing

Make bigger pieces

dsDNA/RNA



motor

membrane

nanopore reader protein

You will be a stronger computational biologist  
if you **understand the biology of plants**

Plant genomes are complicated  
but **the technology is getting better**

How do we approach genome projects  
because **plants are absurd organisms**

How do we apply computational genomics  
to **assemble plant genomes**

# Major steps to start a genome project

**Identifying the accession you want to sequence**

Get “small DNA” for Illumina shotgun  
Genome size, ploidy, heterozygosity



**Isolate ultra HMW DNA**

Get “big DNA” for PacBio sequencing



**Isolate RNA from a variety of tissues and conditions**

Get RNA to annotate genes,  
or do an experiment  
(e.g. differential expression)

# Identifying the accession you want to sequence

## Before you start:

- **Questions you should ask first:**

- Do I know the genome size (flow cytometry or kmer-based)
- Is this a selfer or an outcrosser?
  - Has this specific individual been inbred for several generations?
  - Is it wild collected?
- Annual, or perennial? Can I return to this one plant?
- Do I have adequate plant material (e.g. young leaves)?
- Can I get a diversity of tissues for RNA sequencing?

You will be a stronger computational biologist  
if you **understand the biology of plants**

Plant genomes are complicated  
but **the technology is getting better**

How do we approach genome projects  
because **plants are absurd organisms**

How do we apply computational genomics  
to **assemble plant genomes**

# Long-read shotgun sequencing



- Produces DNA sequence **reads**
- Based on the haploid genome size, calculate coverage of reads
-

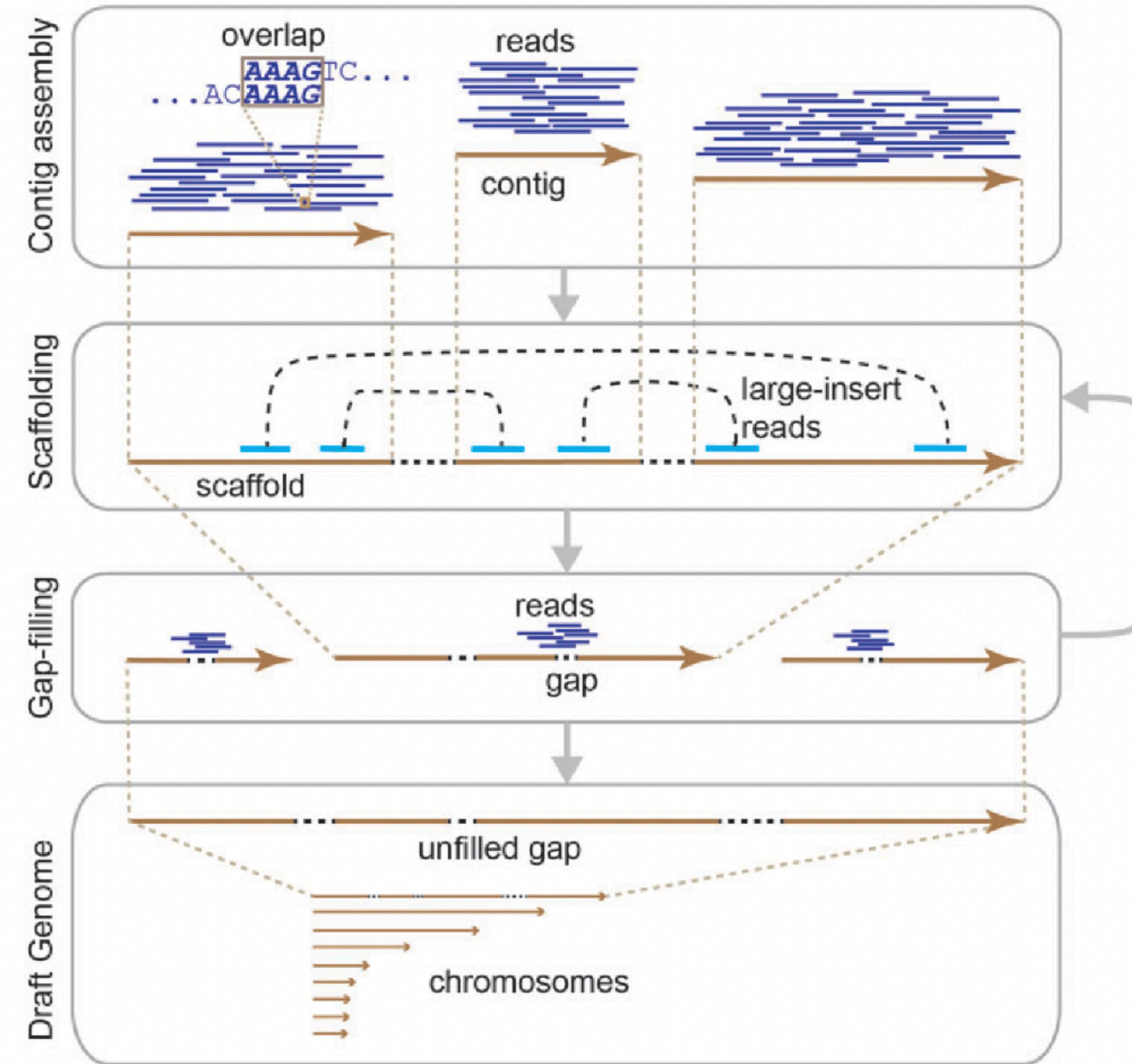
# Long-read shotgun sequencing



- Fragment DNA randomly into pieces
- Then sequence on long-read (e.g. PacBio, Nanopore)
  - or short-read platform (e.g. Illumina)

# Basic recipe for genome assembly

- **Find overlaps between the reads**
  - Depending how many reads we sequence, this adds up
- **Build a graph**
  - A picture of the connections between reads
- **Simplify the graph**
  - Sequencing errors complicate this
- **Traverse the graph**
  - Find the best path through the graph to build a consensus assembly



# de Bruijn vs OLC

## OLC Overlap-Layout-Consensus

**A**

ATATAT[ACTGGCGTATCGCAGTAAAC]GCGCCG

R1: ACTGGCGTAT

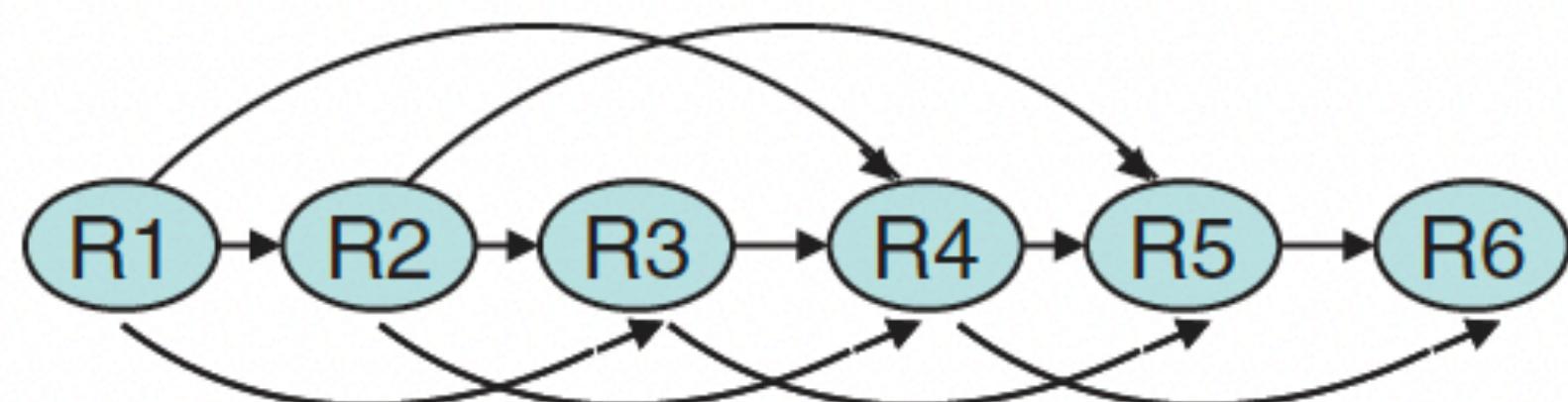
R2: TGGCGTATCG

R3: GGC GTATCGC

R4: CGTATCGCAG

R5: TATCGCAGTA

R6: CGCAGTAAAC



All-by-all pairwise read alignment to find overlaps

Construct a read graph where...

**Nodes** = reads

**Links** = two reads overlap larger than a cutoff length

**Number of nodes** = Number of reads

If requiring at least 5 nt overlaps....

R1 (read 1) overlaps with R2, R3, R4, but NOT R5

# DBG

## de Bruijn Graph

**B**

ATATAT[ACTGGCGTATCGCAGTAAAC]GCGCCG

K1: ACTGG

K2: CTGGC

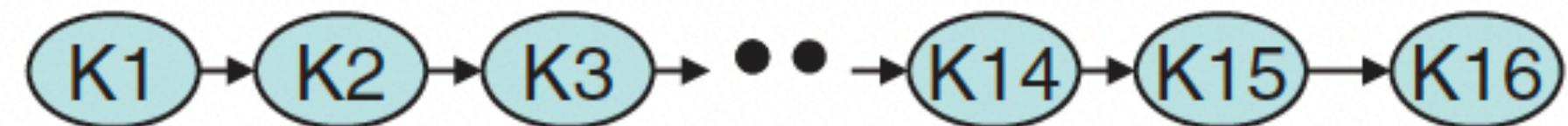
K3: TGGCG

K.: .....

K14: AGTAA

K15: GTAAA

K16: TAAAC



Break reads into **kmers** (e.g. k=5)  
**kmer** = a sequence with length k

Construct a read graph where...

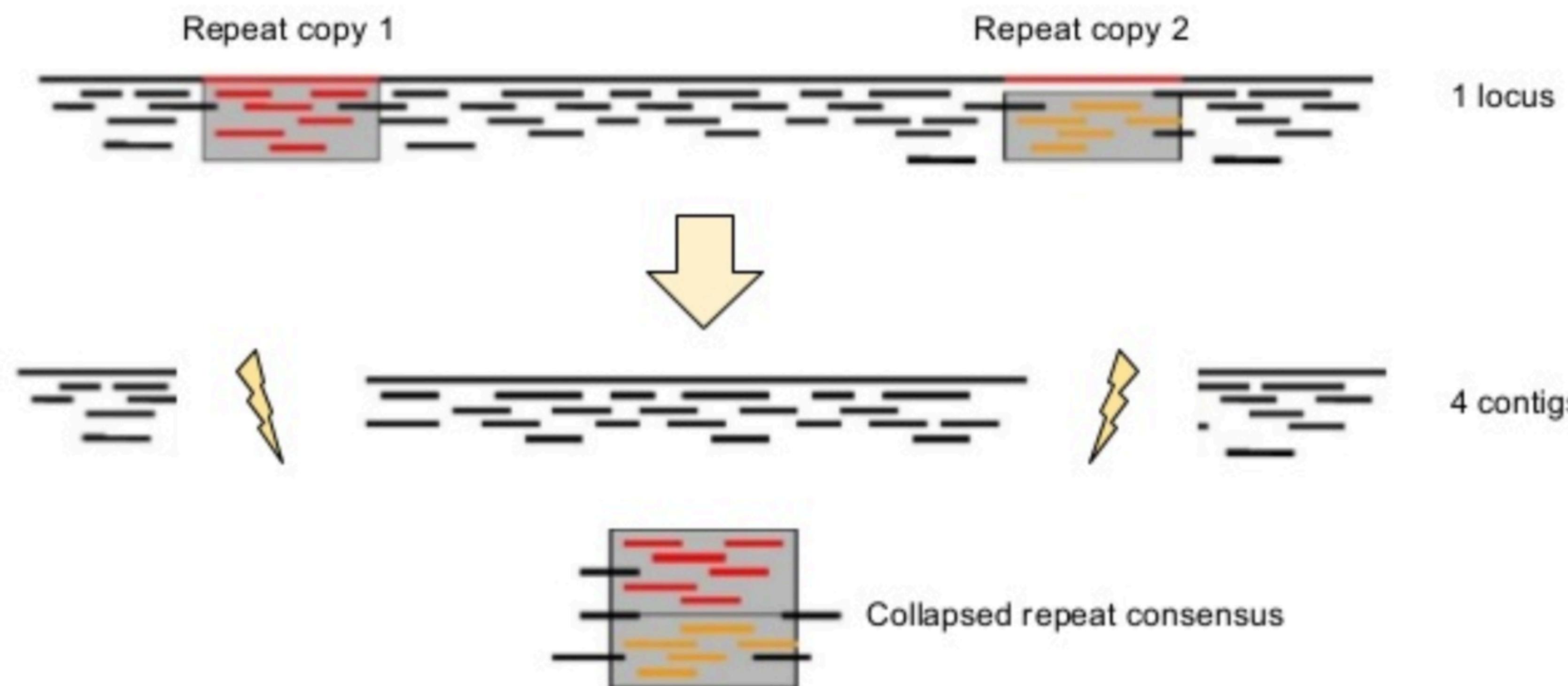
**Nodes** = kmers

**Links** = two kmers overlap by k-1

**Number of nodes** = genome size + errors

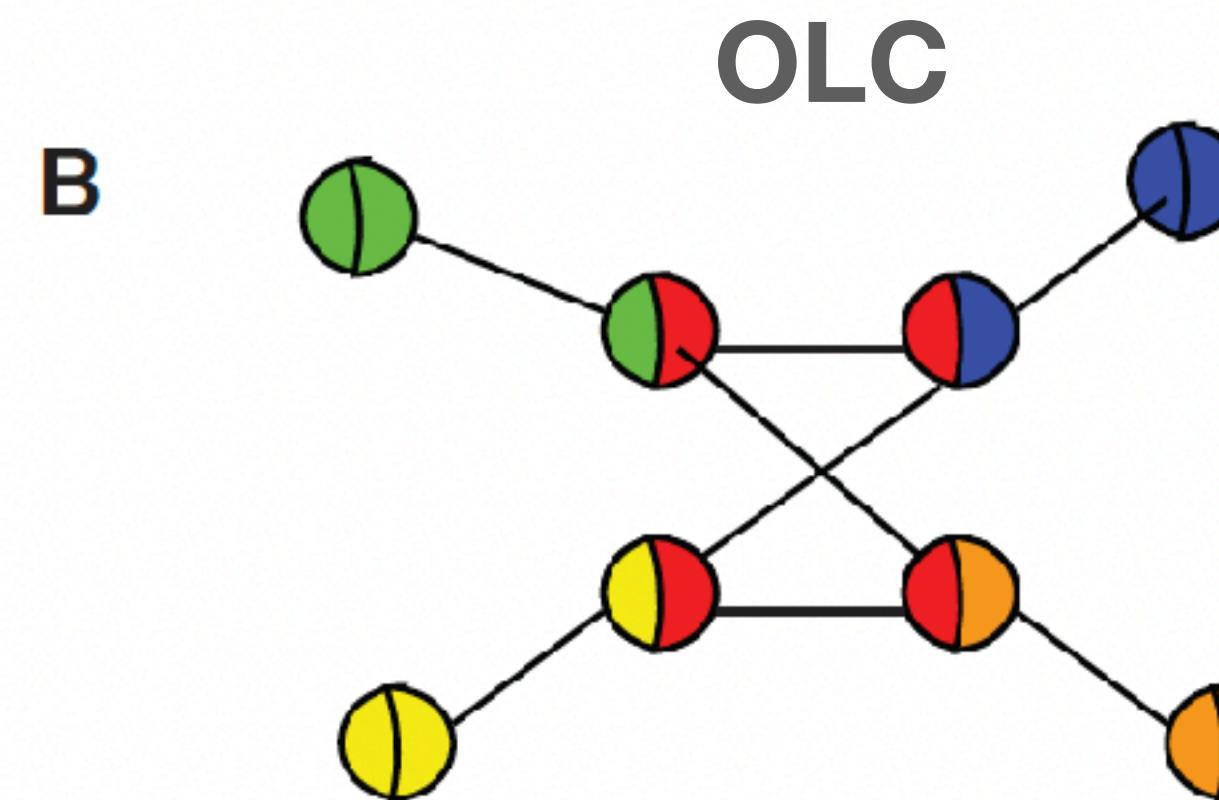
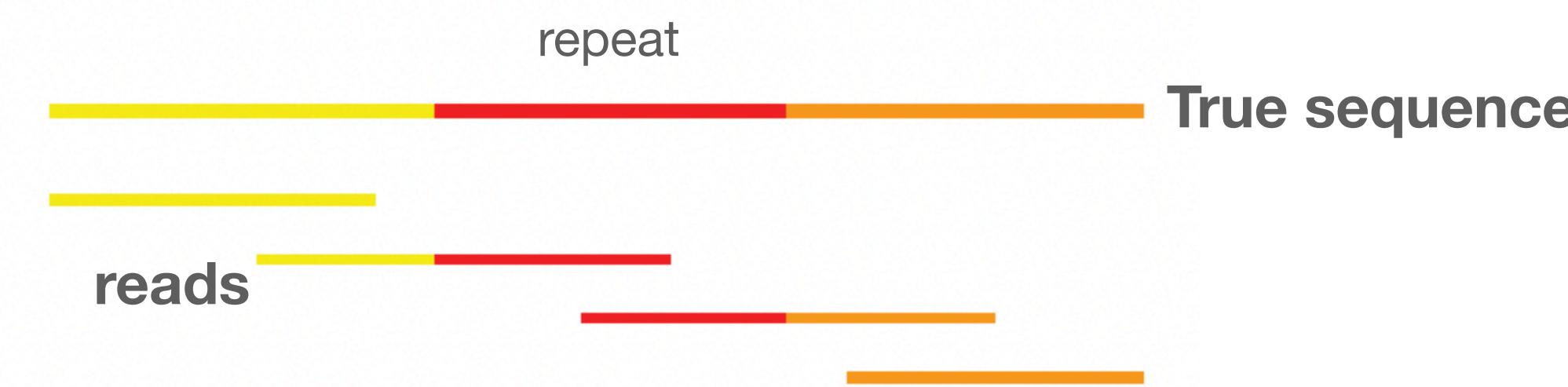
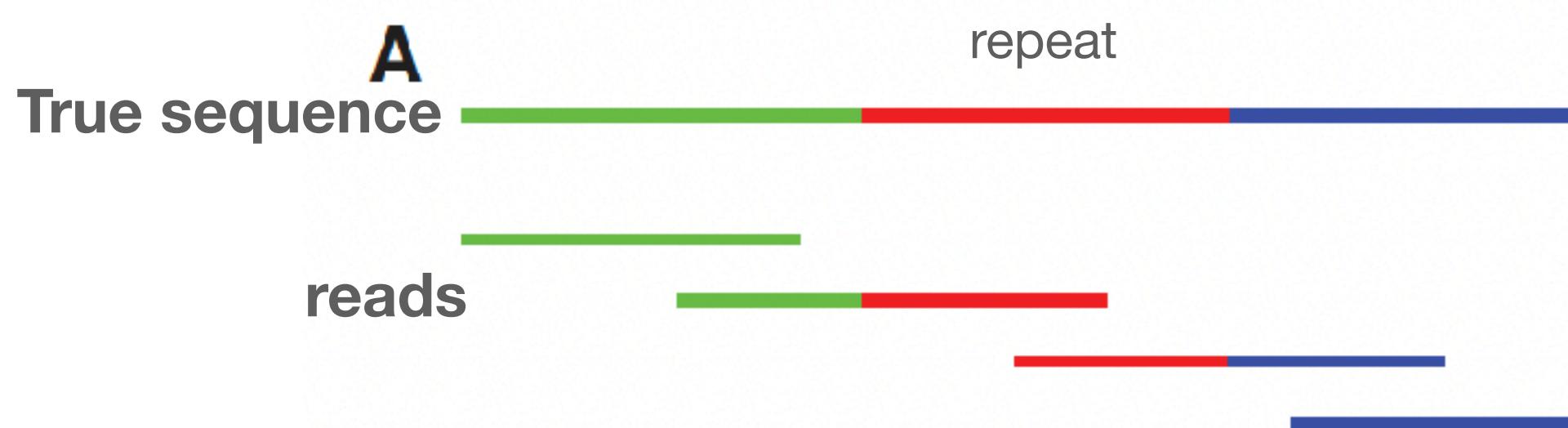
# Assembling repeats

## Repeats

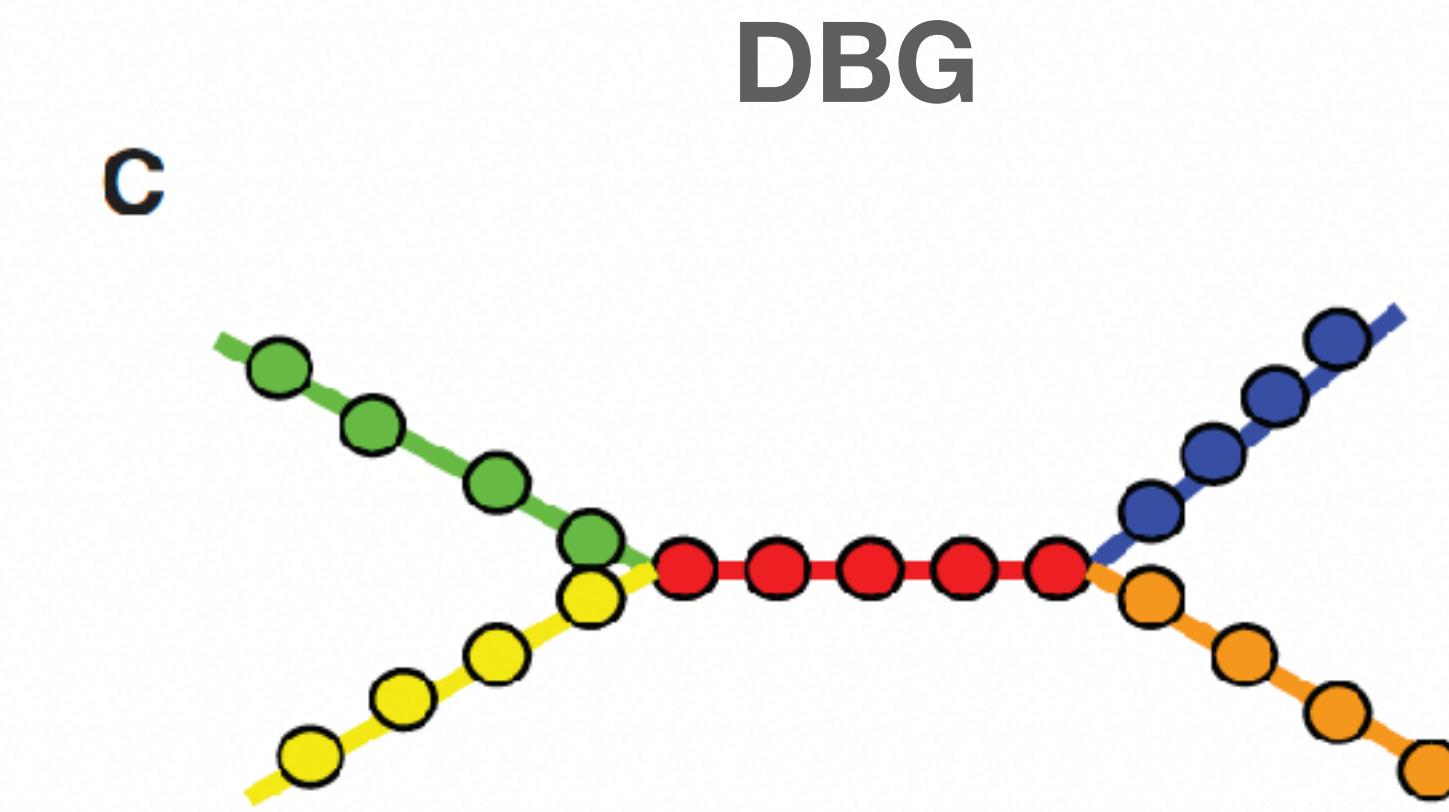


# OLC vs de Bruijn graph assembly

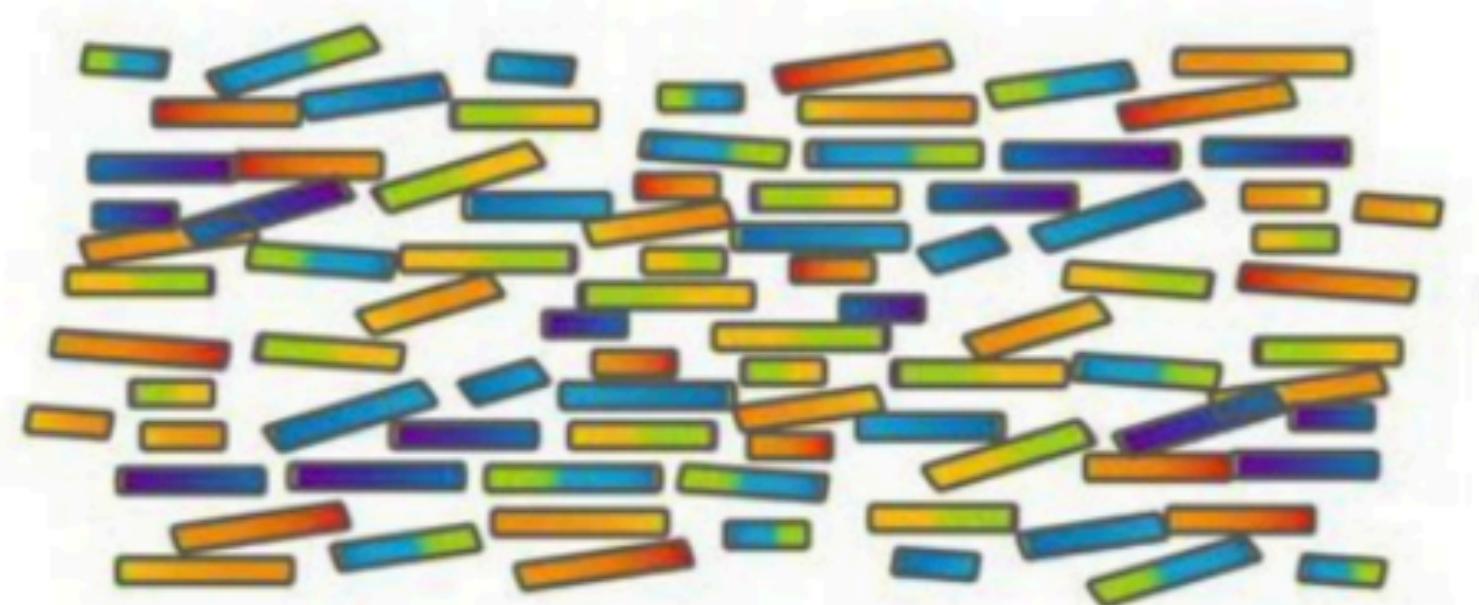
Two regions of the genome, same repeat



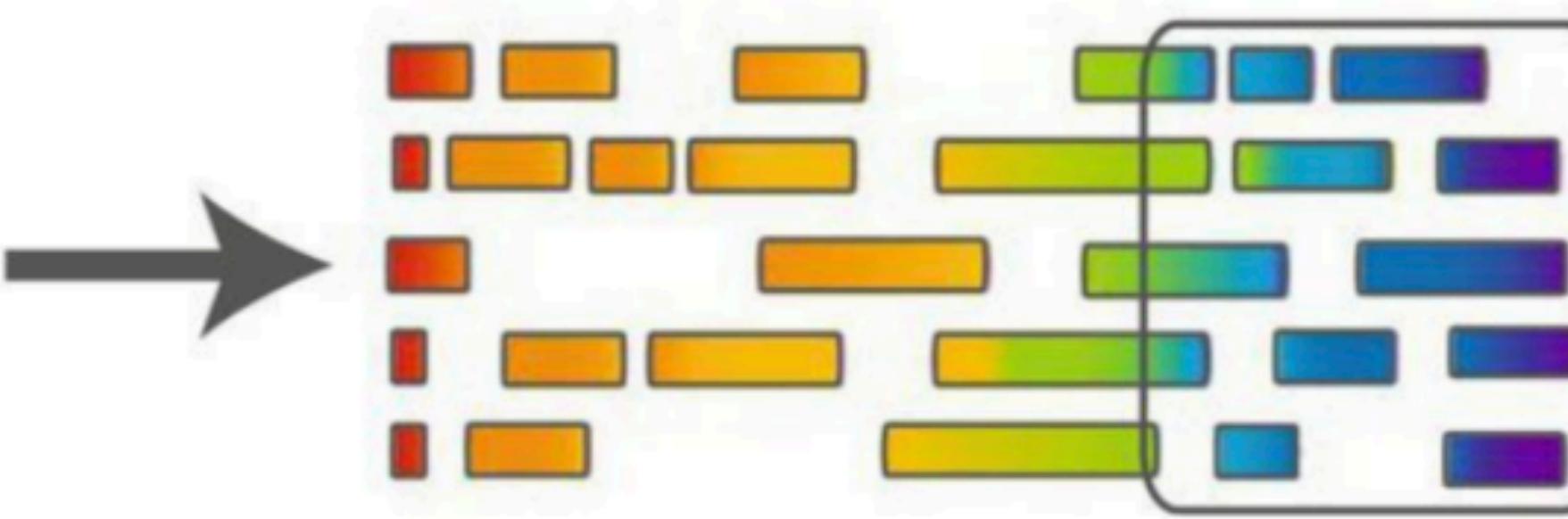
Keep increasing read length!



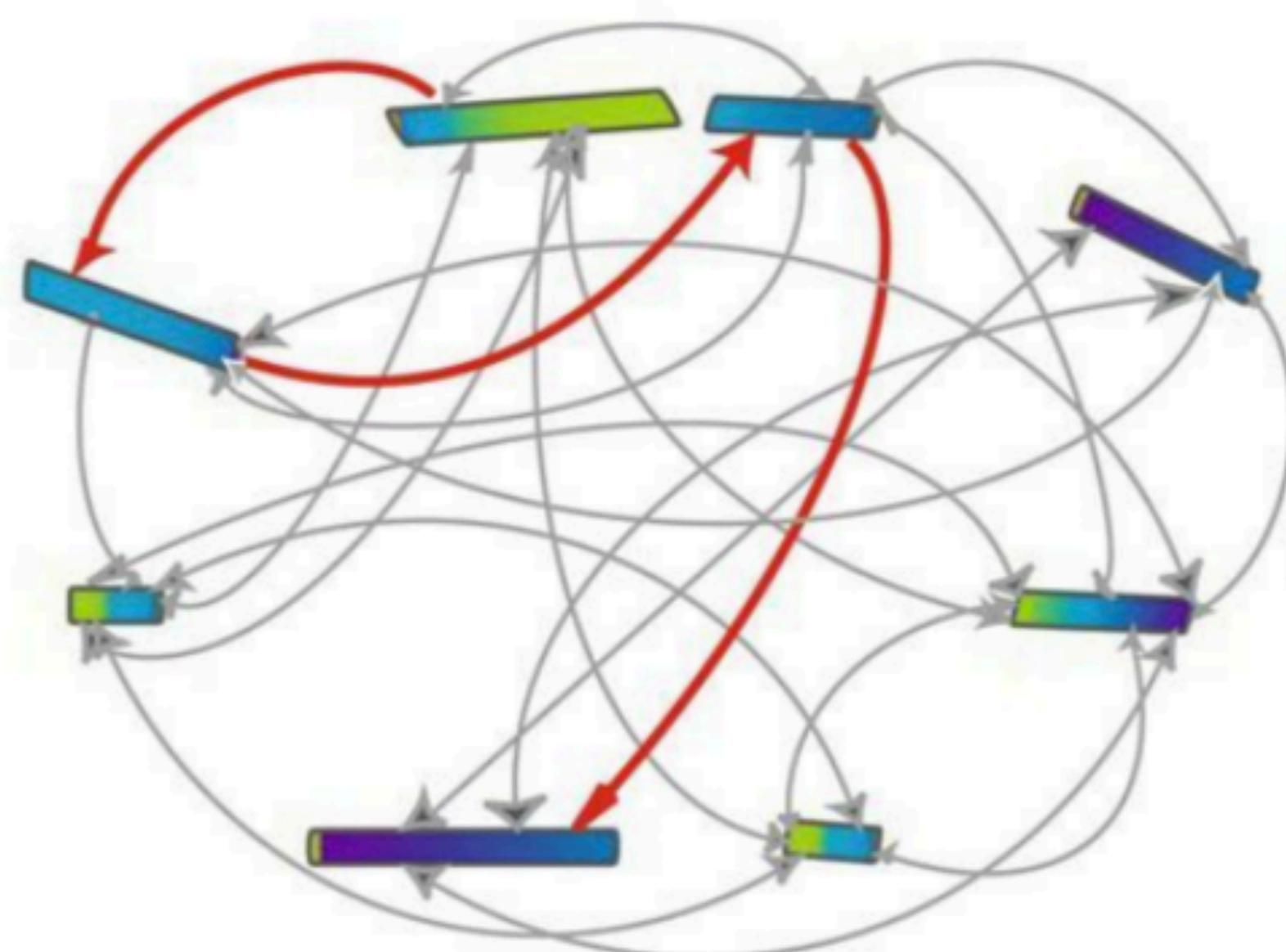
Can try to increase kmer size...



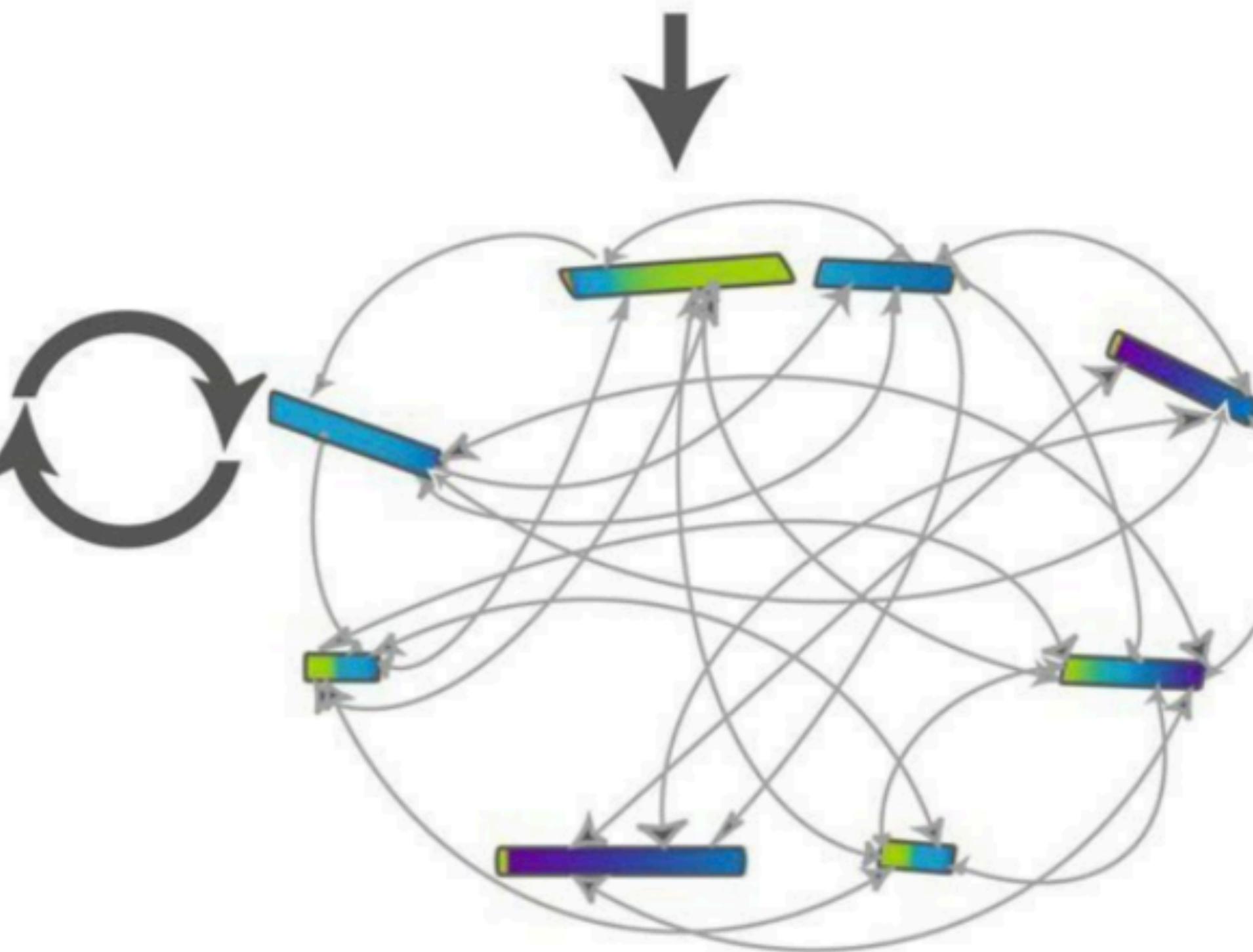
Reads provided to algorithm



Overlaps identified



Hamiltonian Path identified



Reads connected by overlaps



Consensus sequence

# **What complicates assembly**

**Biology can be annoying; data isn't perfect**

- Raw data error rates
  - Illumina <<< 1%
  - PacBio ~1% or less (for HiFi)
  - Nanopore 3-15% depending
- Repetitive elements
- How much coverage did you sequence?
- Heterozygosity
- Ploidy and flavor of ploidy
- Genome size

# What complicates assembly

## Repetitive Elements

