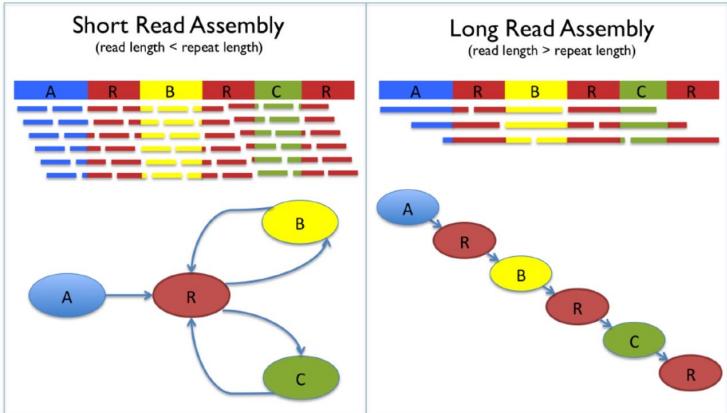


PACBIO®

Genome assembly

Jacob B. Landis

July 24th, 2022



What do all these things have in common?

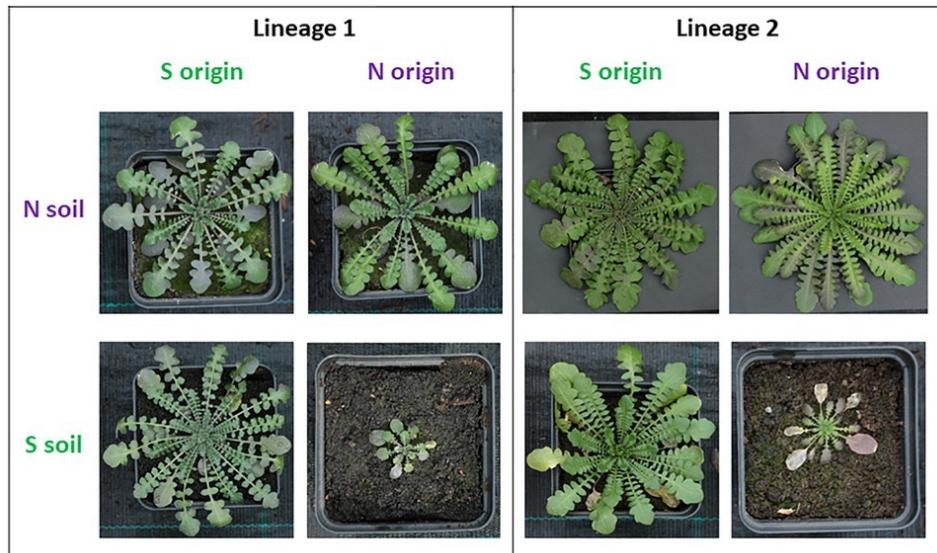
Shifts in pollinators



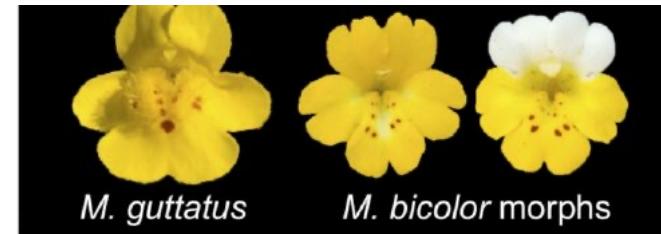
Evolution of carnivorous plants



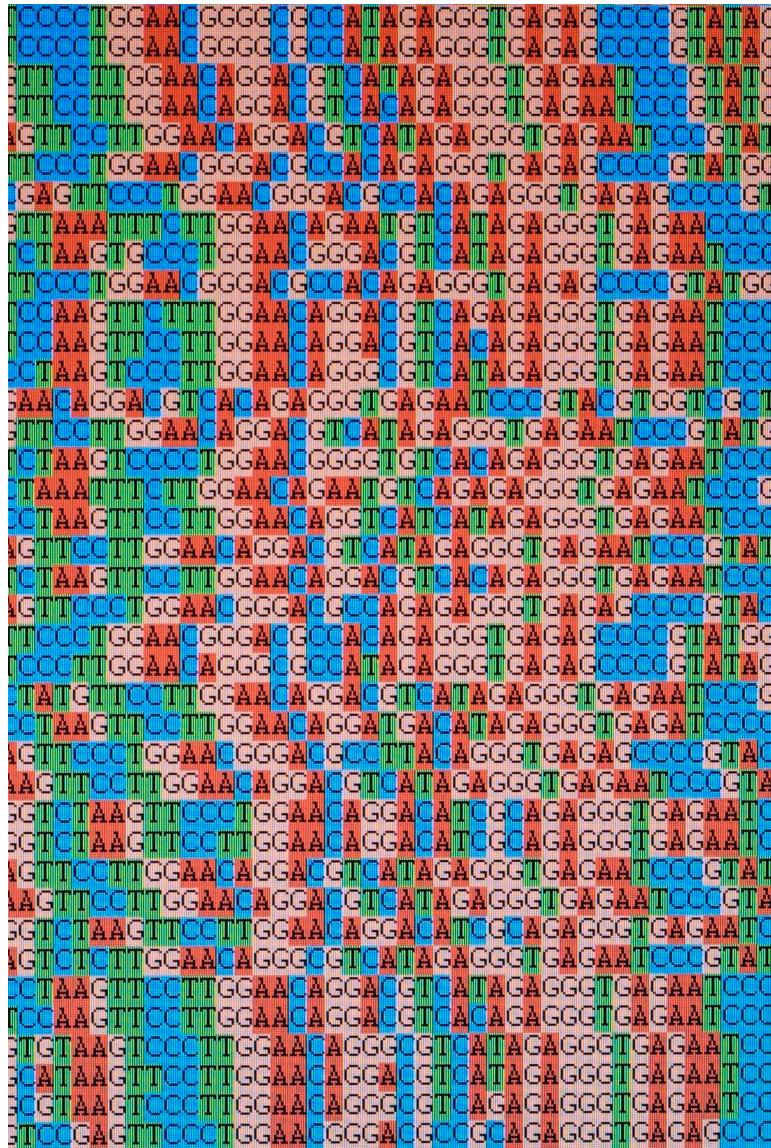
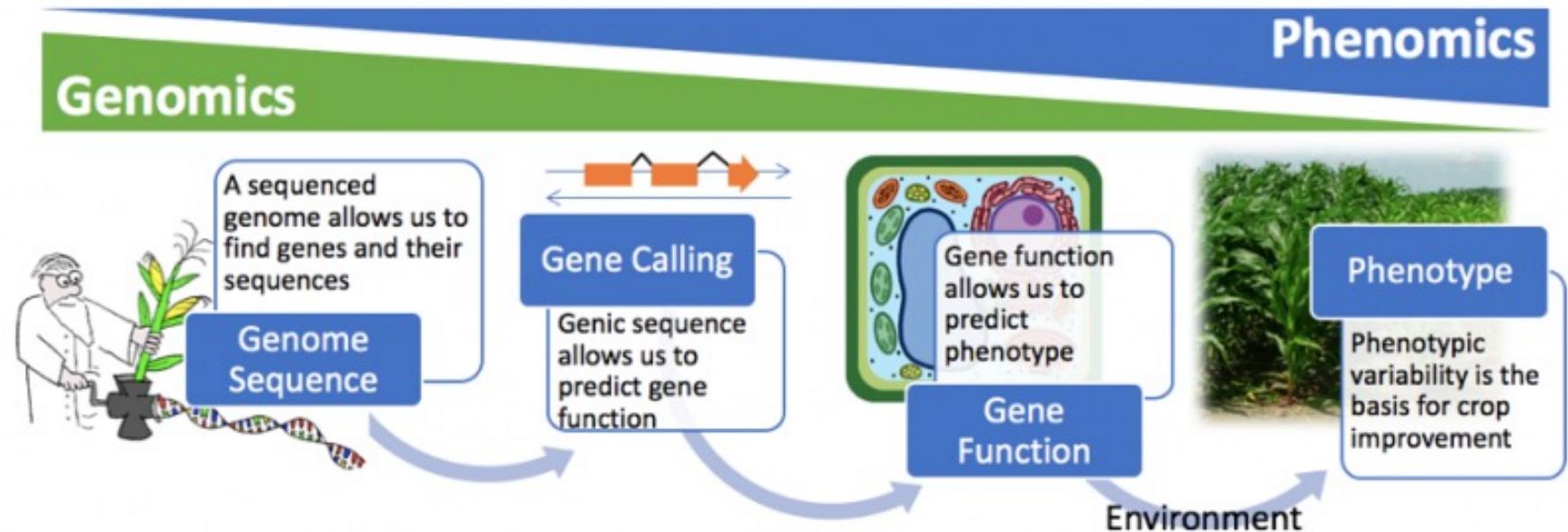
Response to nutrients



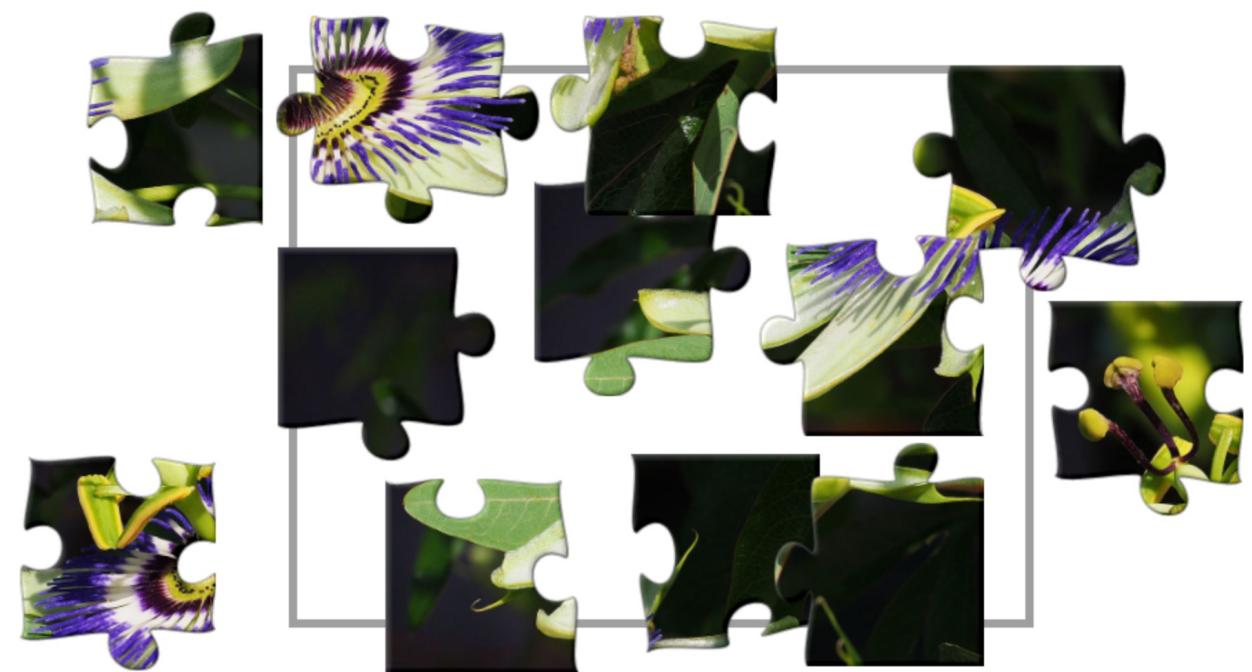
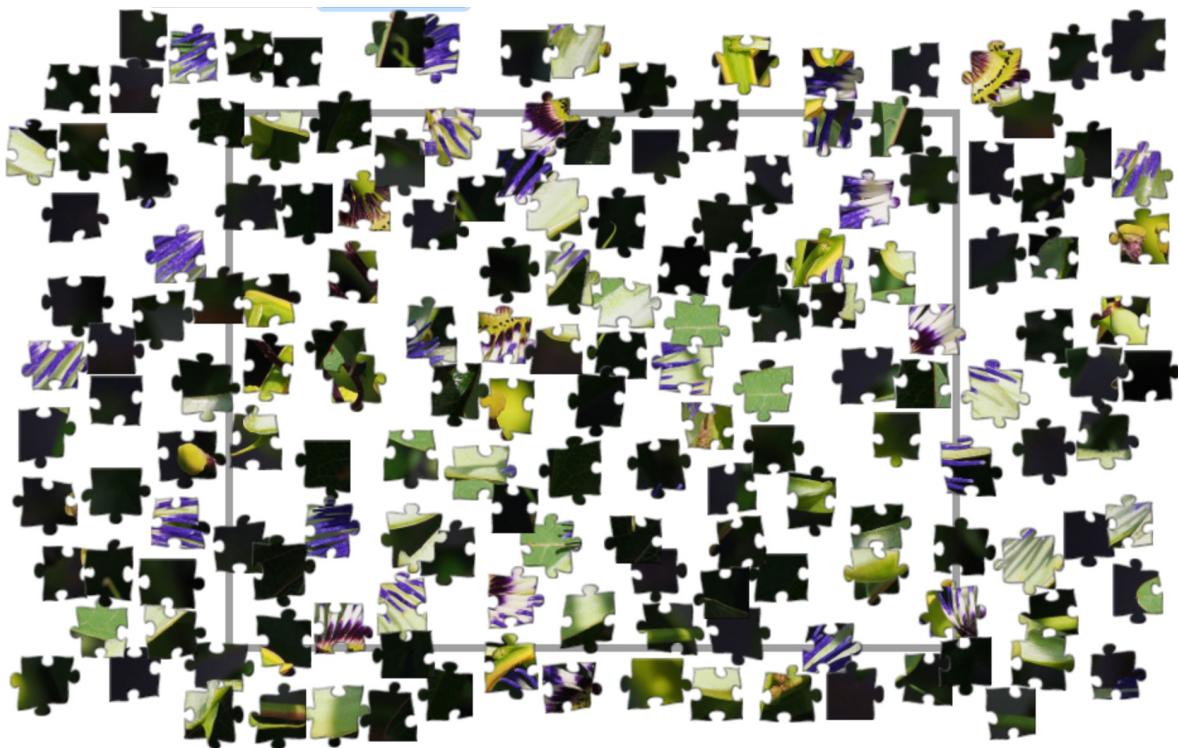
Flower color morphs



Genome as a puzzle



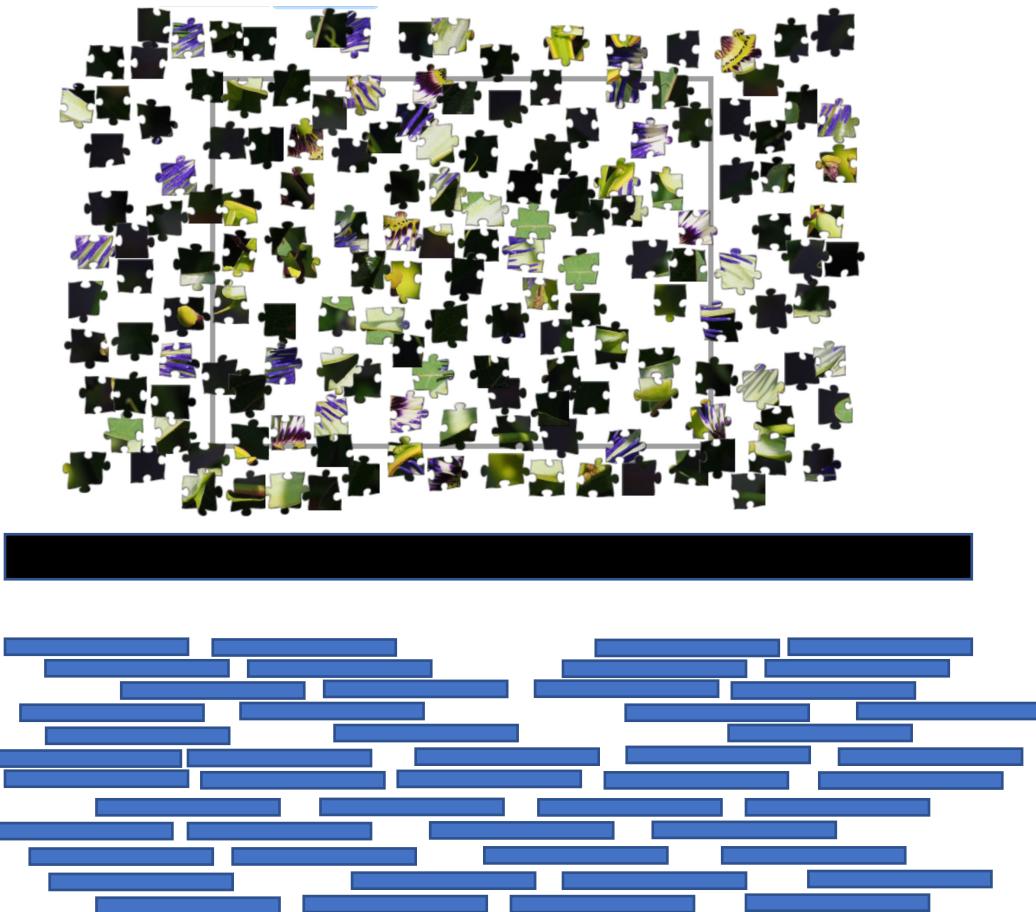
Which puzzle is easier to put together?



Which puzzle is easier to put together?

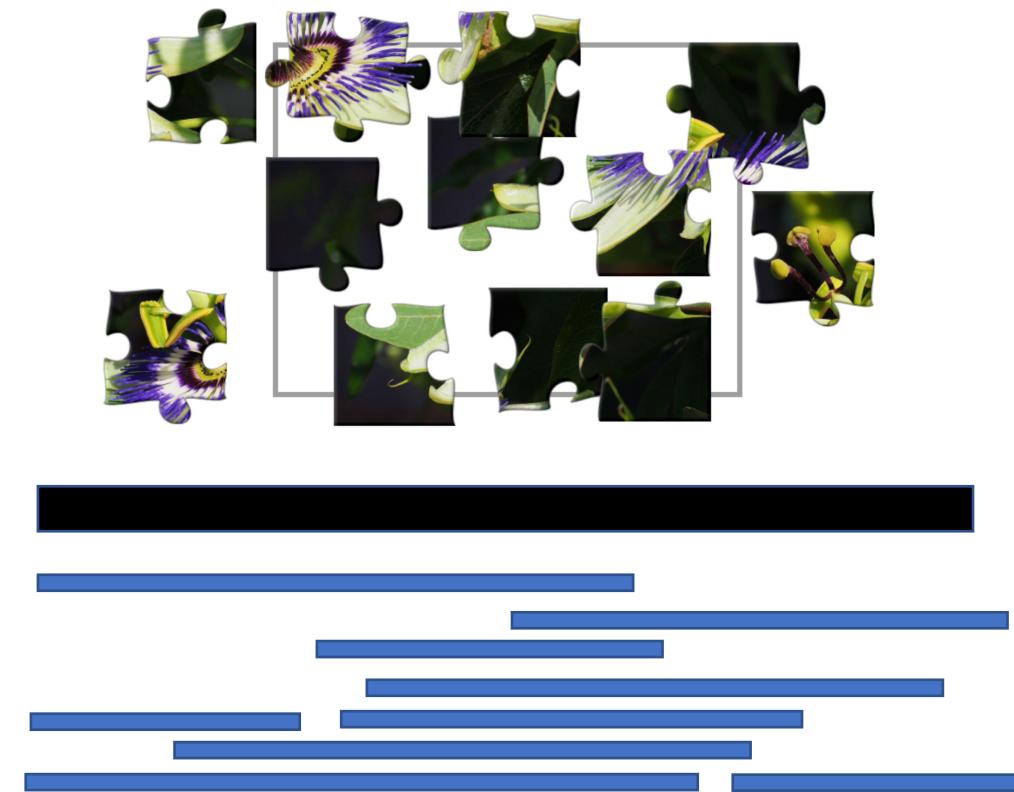


In the world of genomes



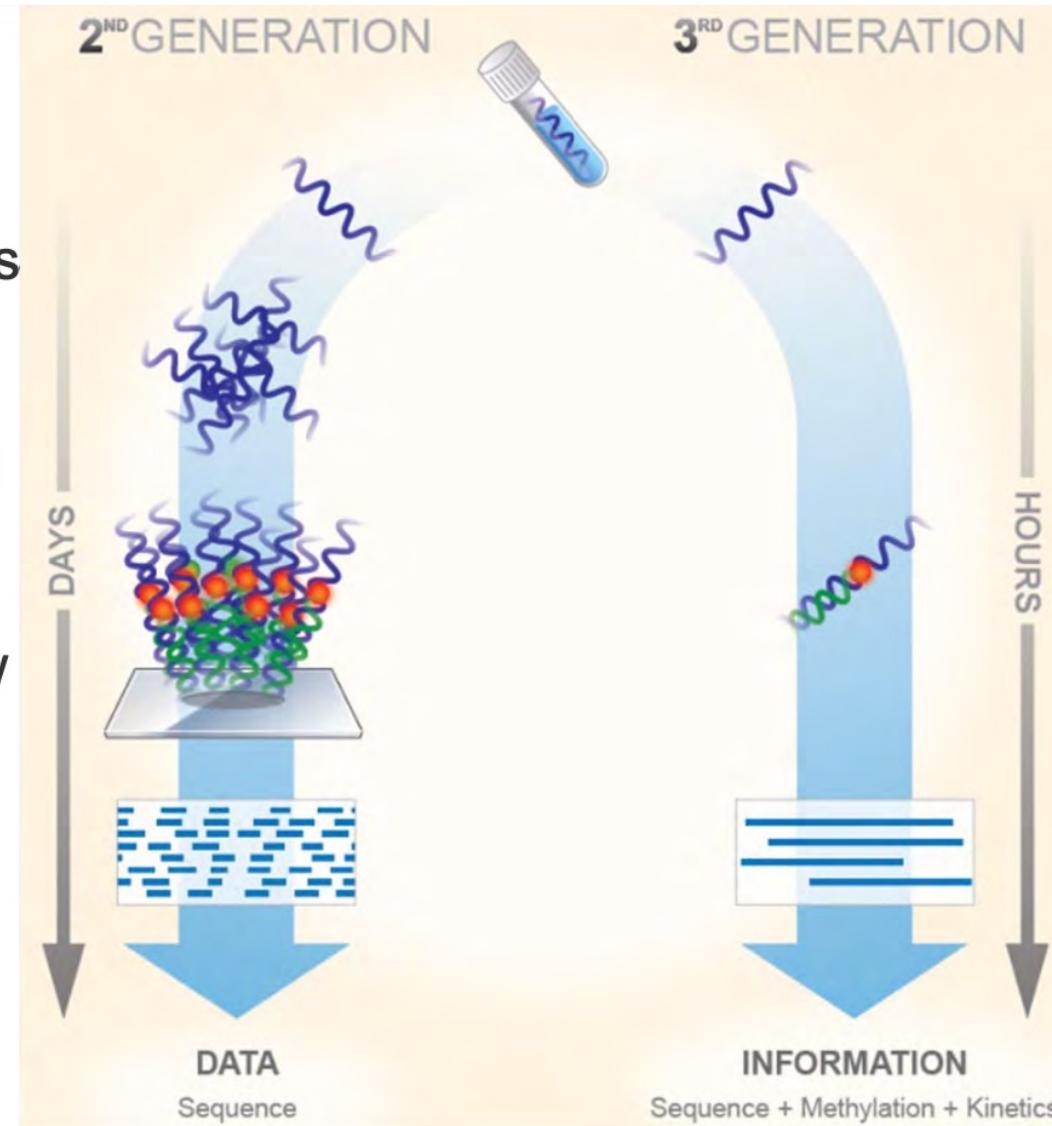
Reference
Genome

Sequencing
Reads



Short vs Long-reads

- Short reads
- Amplification errors and bias
- Several enzymatic steps
- Multi-molecule raw accuracy
- Errors tend to be systematic
- More coverage required

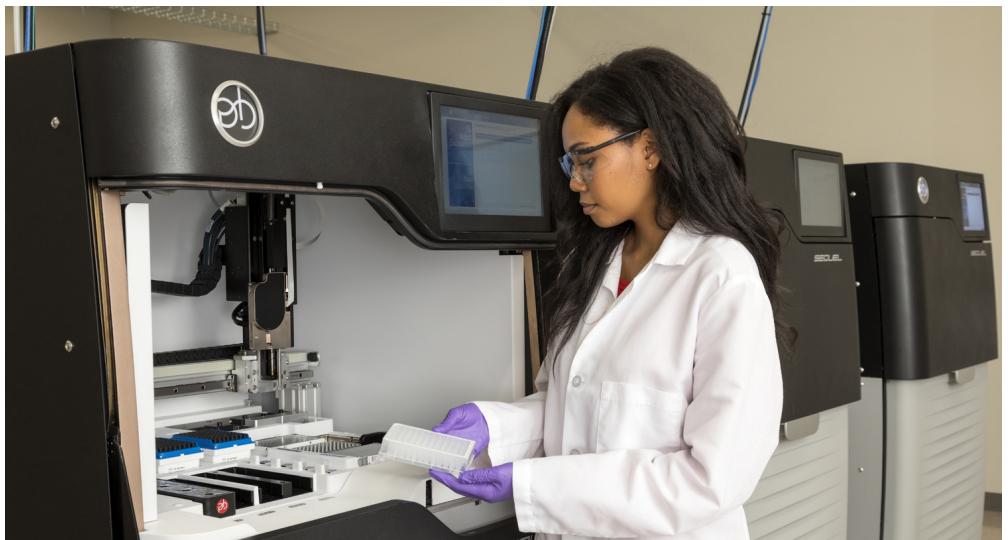


- Long reads
- No required amplification
- Simple sample prep
- Single molecule raw accuracy
- Errors tend to be random (vs. systematic)
- Less coverage required

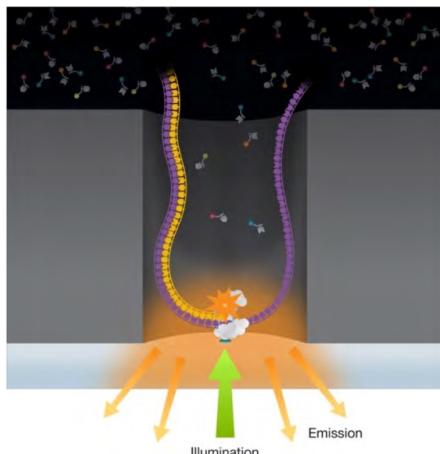
Long-read options

	PacBio ¹		Oxford Nanopore ²	
Instrument Specifications	RS II (P6-C4)	Sequel	MinION	PromethION
Average read length	10 – 15 kb	10 – 15 kb	Variable (up to 900 kb) ^{3,4}	*
Error rate	10 – 15 %	10 – 15 %	5 – 15 % ^{4,5}	*
Output	500 Mb – 1 Gb	5 Gb – 10 Gb	~5 Gb ⁴	*
# of reads	~50k	~500k	Variable (up to 1M) ^{6,7}	*
Instrument price/Access fee ^a	\$700k	\$350k	\$1000 ⁸	\$135k bundle ⁹
Run price	~\$400	~\$850	\$500-\$900 ⁷	*

PacBio



SMRT Cell



Science, Vol 299, Jan 31 2003, pp682-686
J. Appl. Phys. 103, 034301 (2008)

1. generate amplicon

2. ligate adaptors

3. sequence

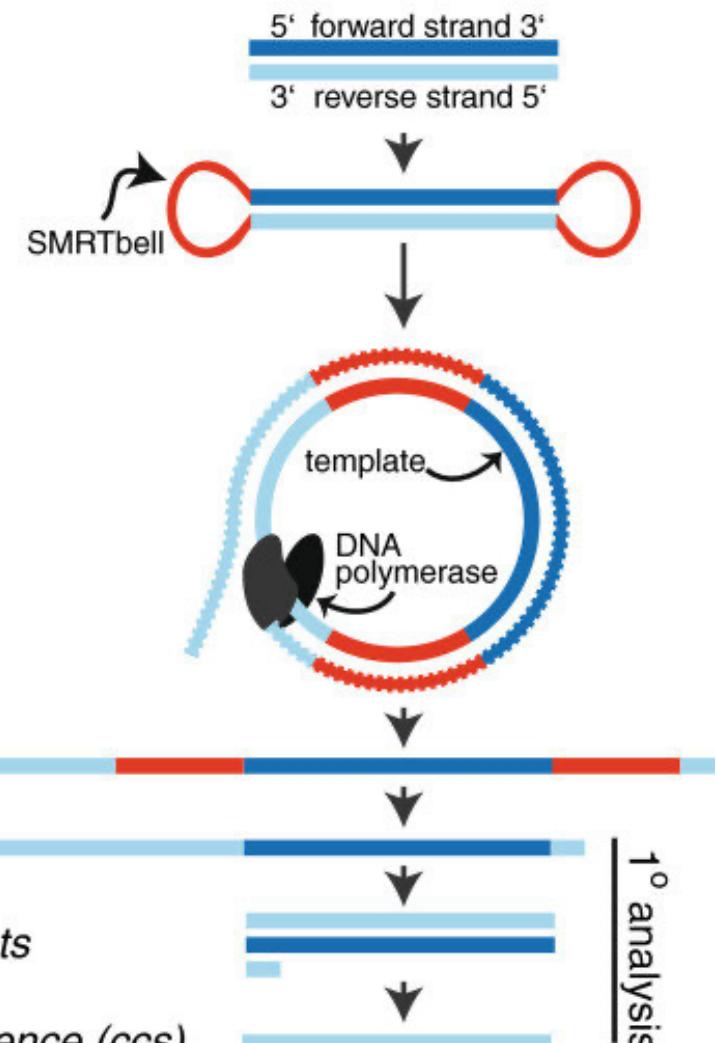
4. data analysis

raw long read

processed long read

single-molecule fragments

circular consensus sequence (ccs)



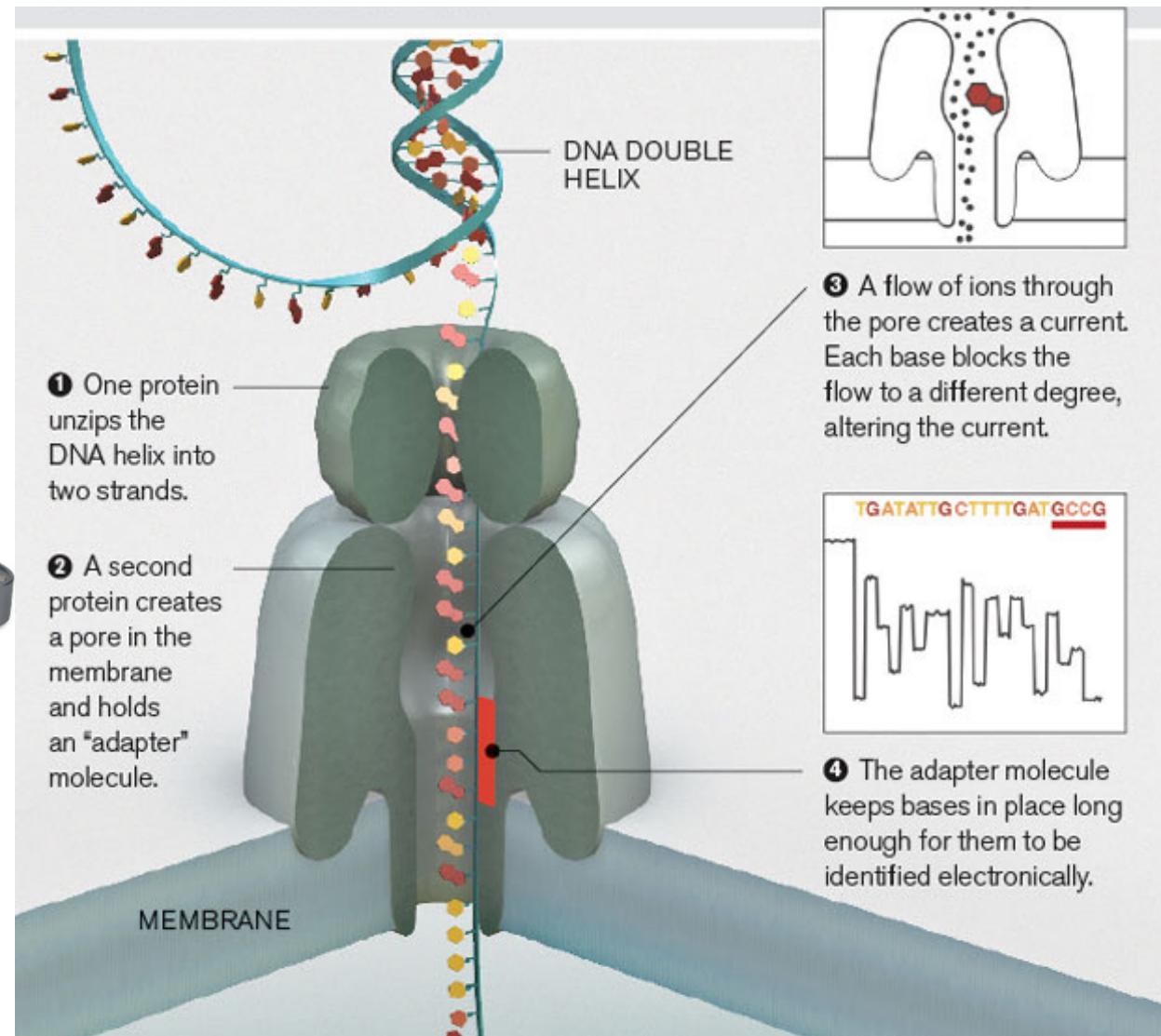
Fichot and Norman 2013; *Microbiome*

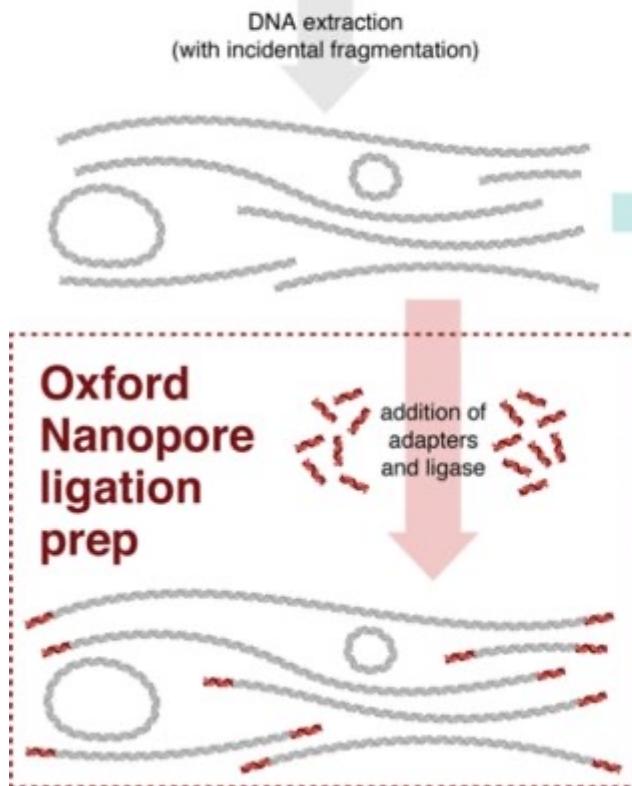
Nanopore



Flowcell

MinION device



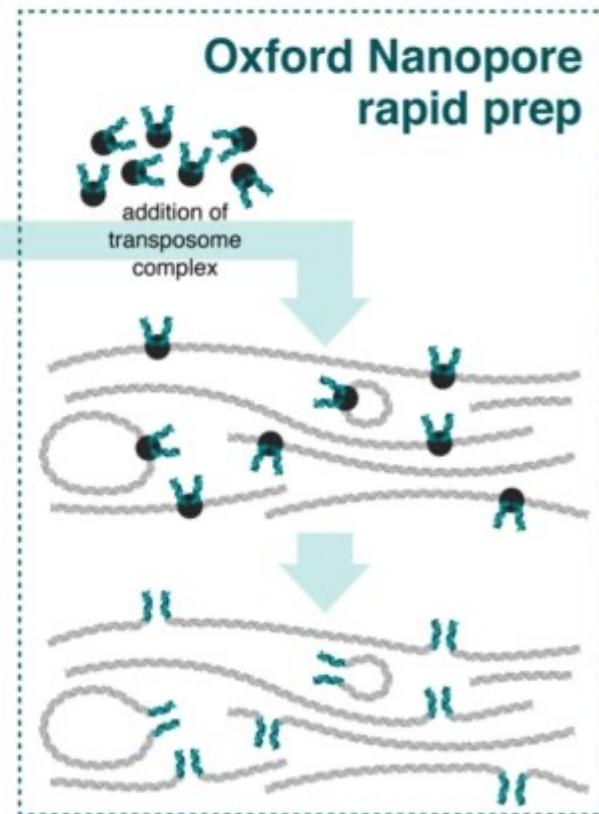


Library Prep ~2.5 hours

Output 10-20 GB in 96 hours

Nanopore approaches

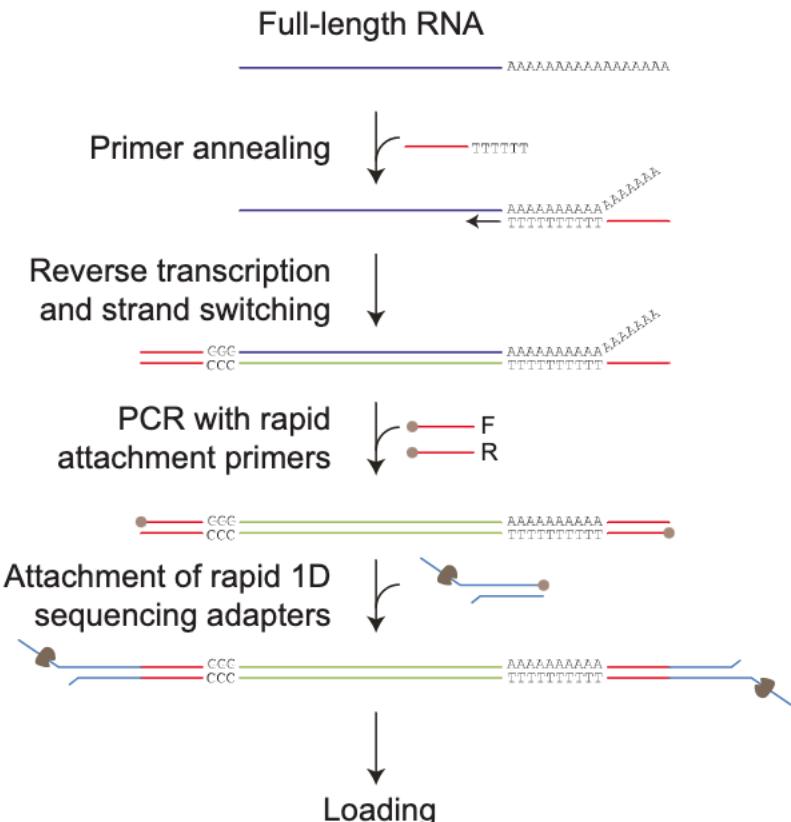
Genomic DNA



10 minutes

8-10 GB in 96 hours

RNA



165 minutes

5 – 7 Million transcripts in 48 hours

de novo assembly

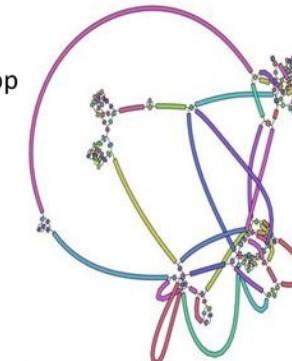
- Illumina only
 - High quality reads with fewer errors
- Hybrid option
 - Nanopore or PacBio + Illumina
 - Either raw or error corrected long-reads
- Long-read only
 - Raw typically works better
 - Need to polish after with Illumina data to fix errors
- Recommendation is 50x coverage short-reads and 50x long-reads
- So how much data do I need?

Short-read

Assemblies

- Fragmented
- Small N50:
10s–100s of kbp
- Very accurate

Uses



Long-read

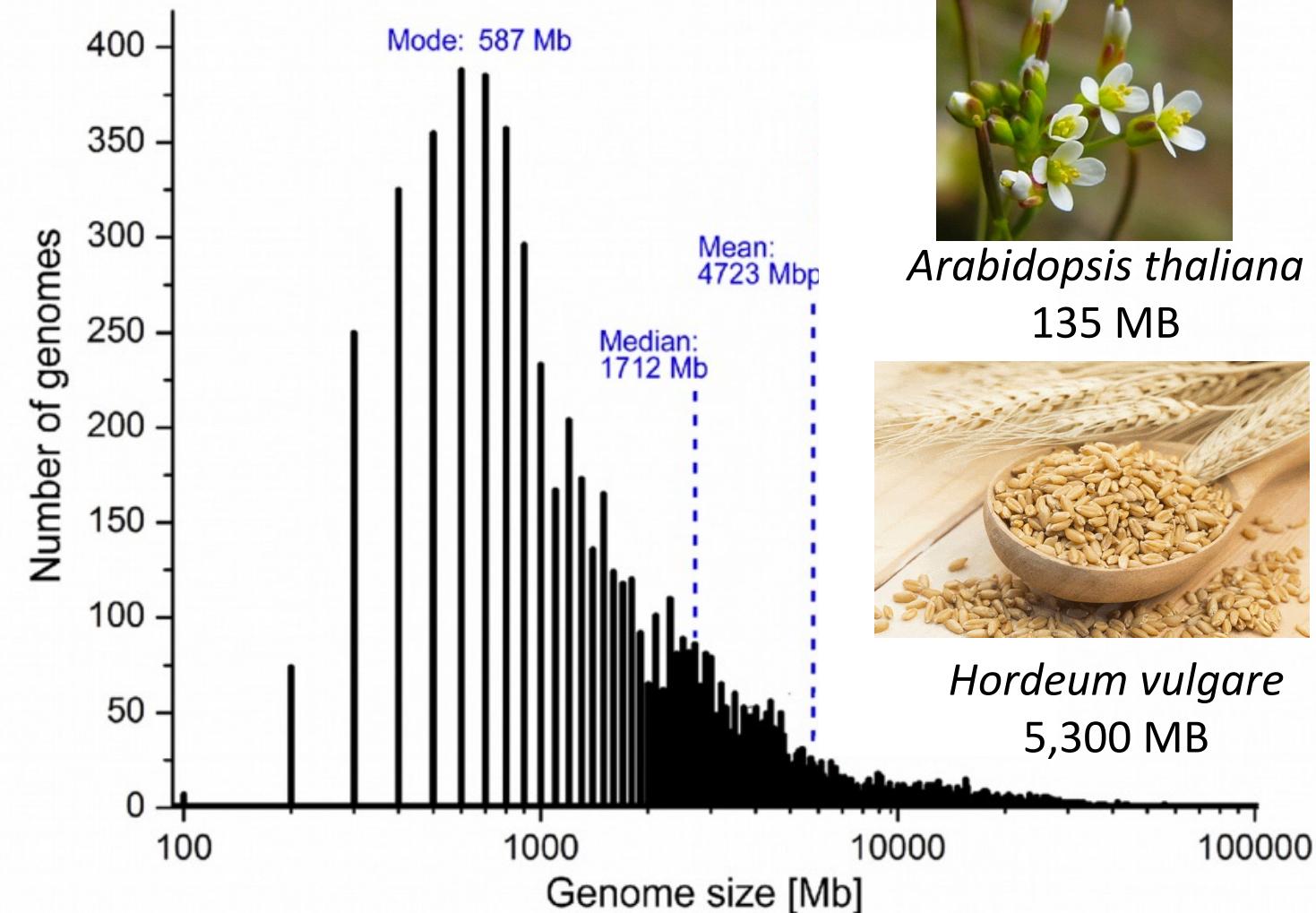
Assemblies

- Complete
- 98+% accuracy

Uses



Genome size of angiosperms



Arabidopsis thaliana
135 MB



Hordeum vulgare
5,300 MB



Oryza sativa
430 MB



Allium cepa
16,000 MB



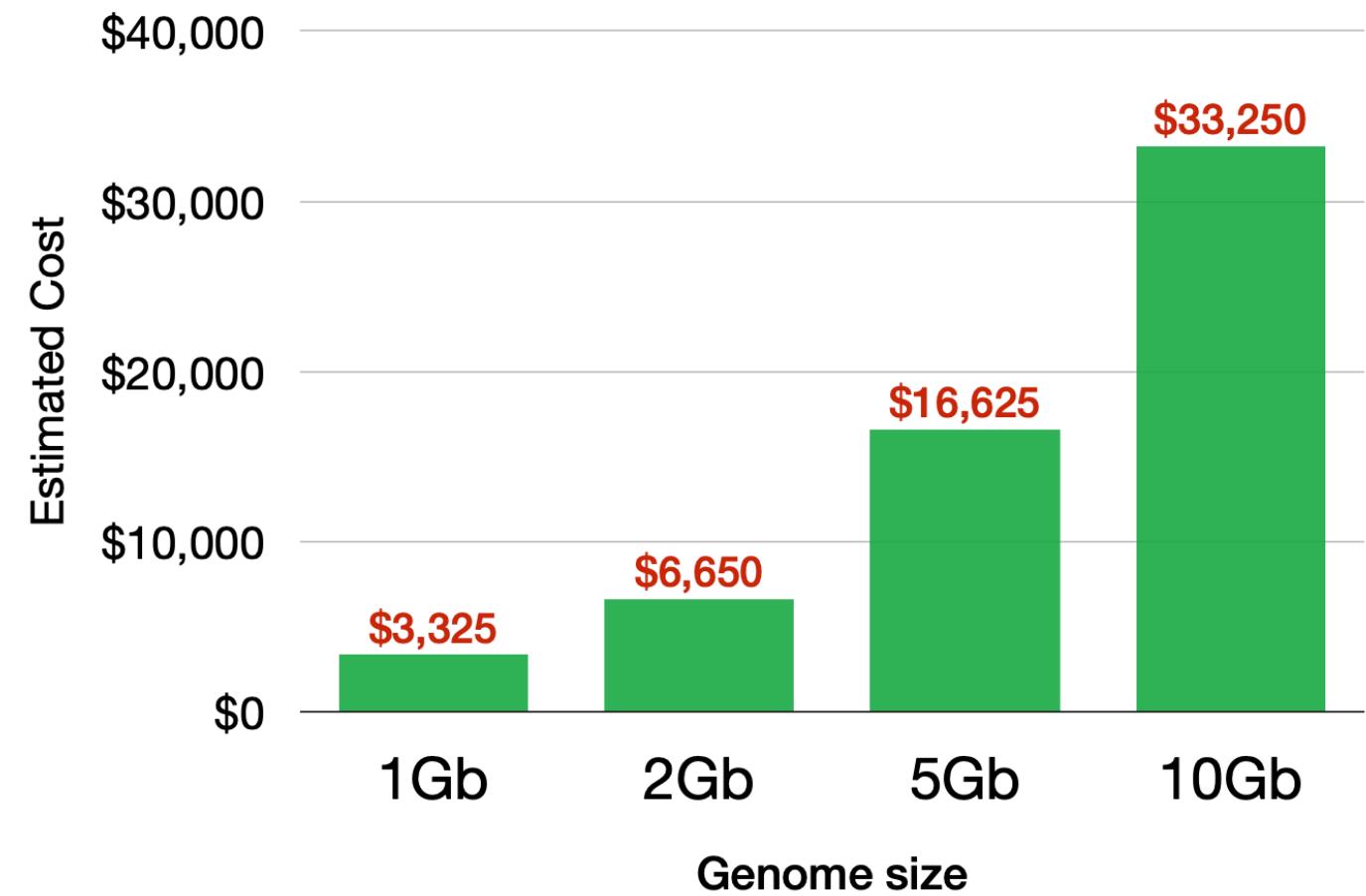
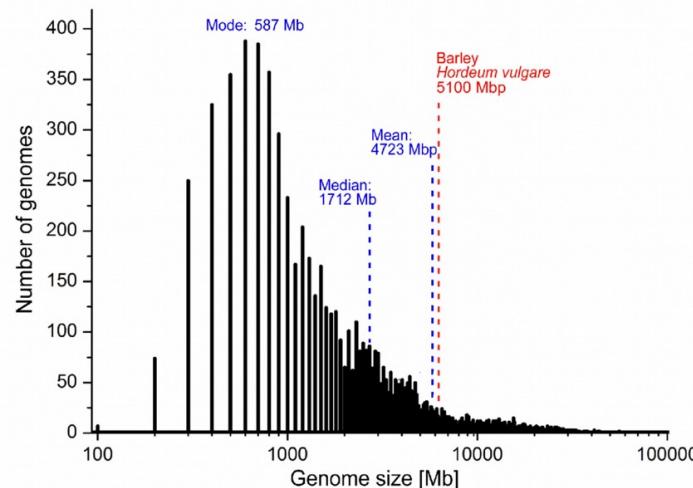
Zingiber officinale
1,582 MB



Tulipa sylvestris
59,241 MB

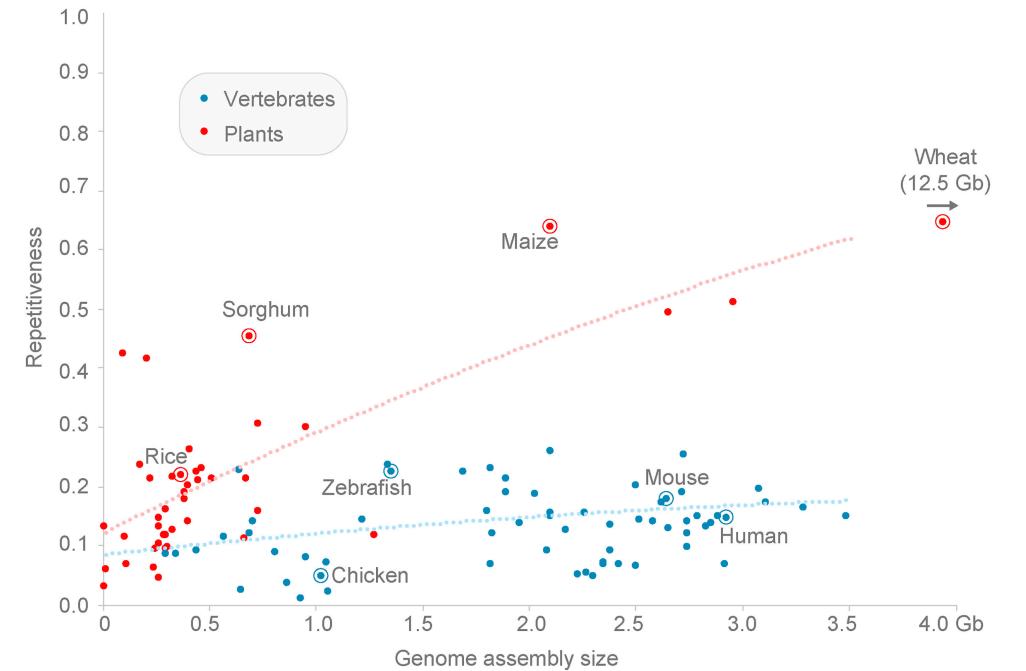
Cost of sequencing genomes

- 50X Illumina:
 - $50\text{Gb} \times \$26.5/\text{Gb} = \$1,325$
 - 50X nanopore:
 - $50\text{Gb} \times \$40/\text{Gb} = \$2,000$
-
- \$3,325**

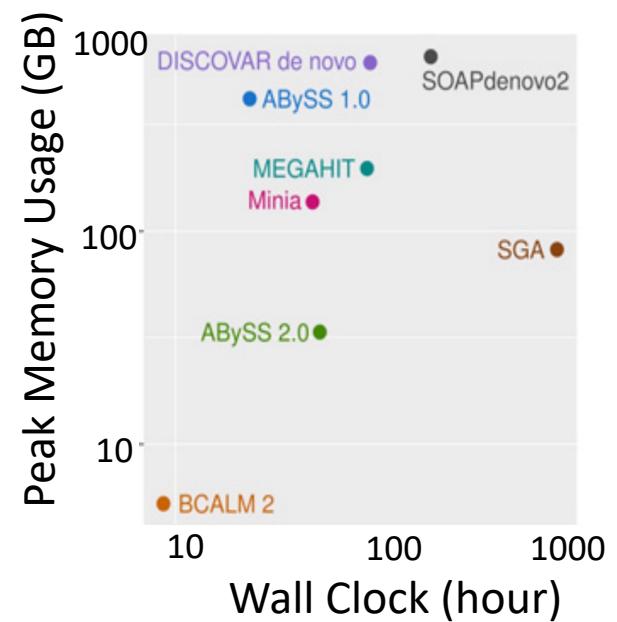


Assembling large genomes

- Many assemblers are designed for genomes equal or smaller to the human genome (3 GB)
- Larger genomes are more repetitive but generally have similar number of genes
- Computation resources intense for deep coverage need
 - memory and wall-time



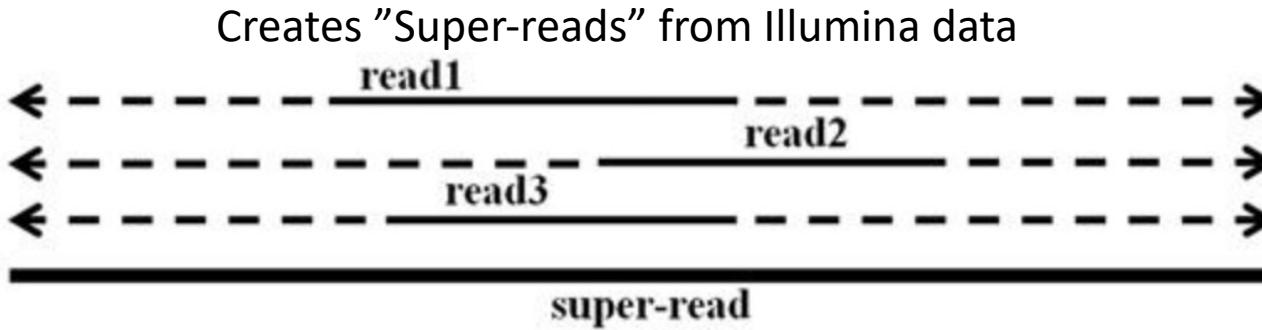
Jiao and Schneeberger 2017



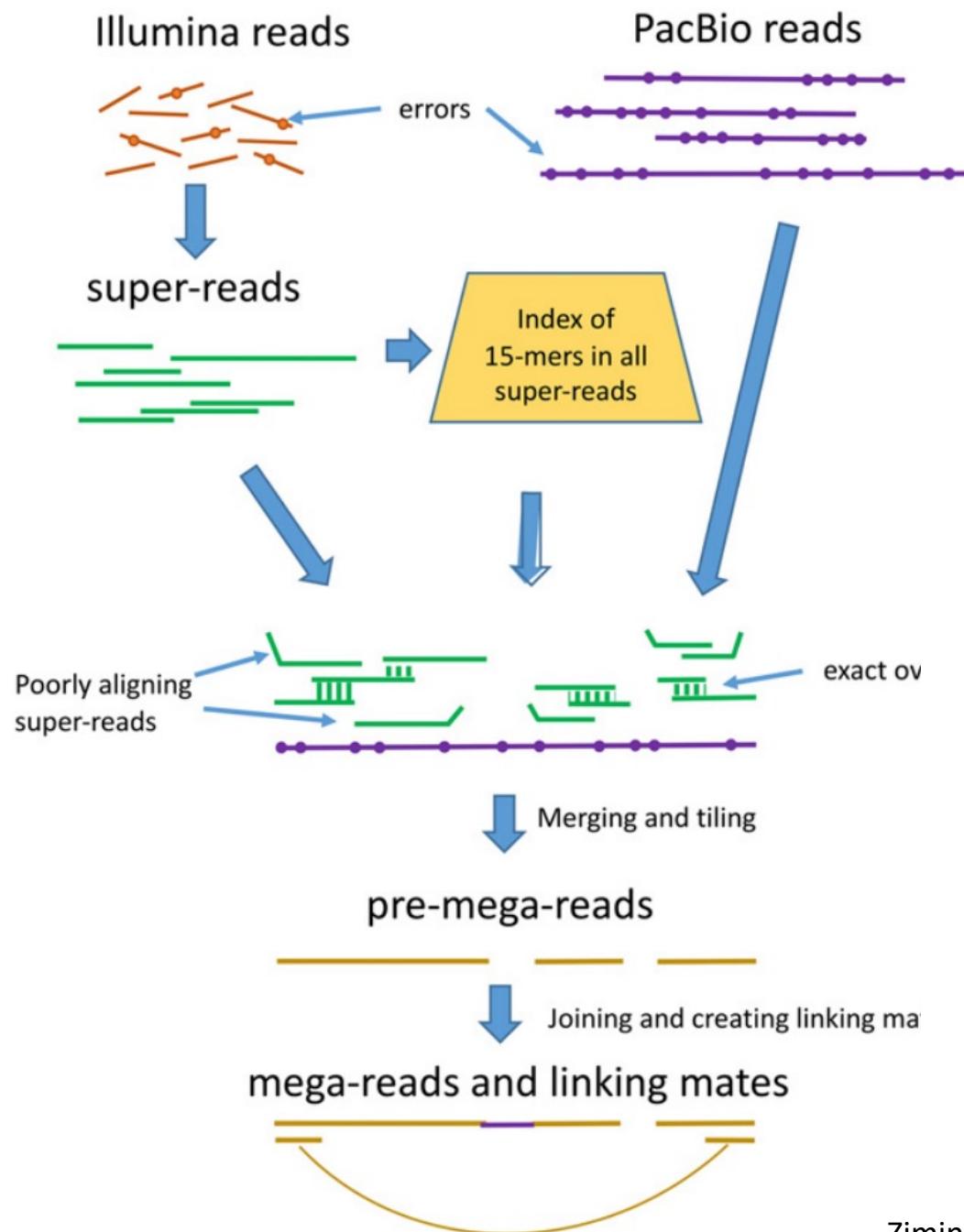
Adapted from Jackman et al. 2017

MaSuRCA – short-read and hybrid

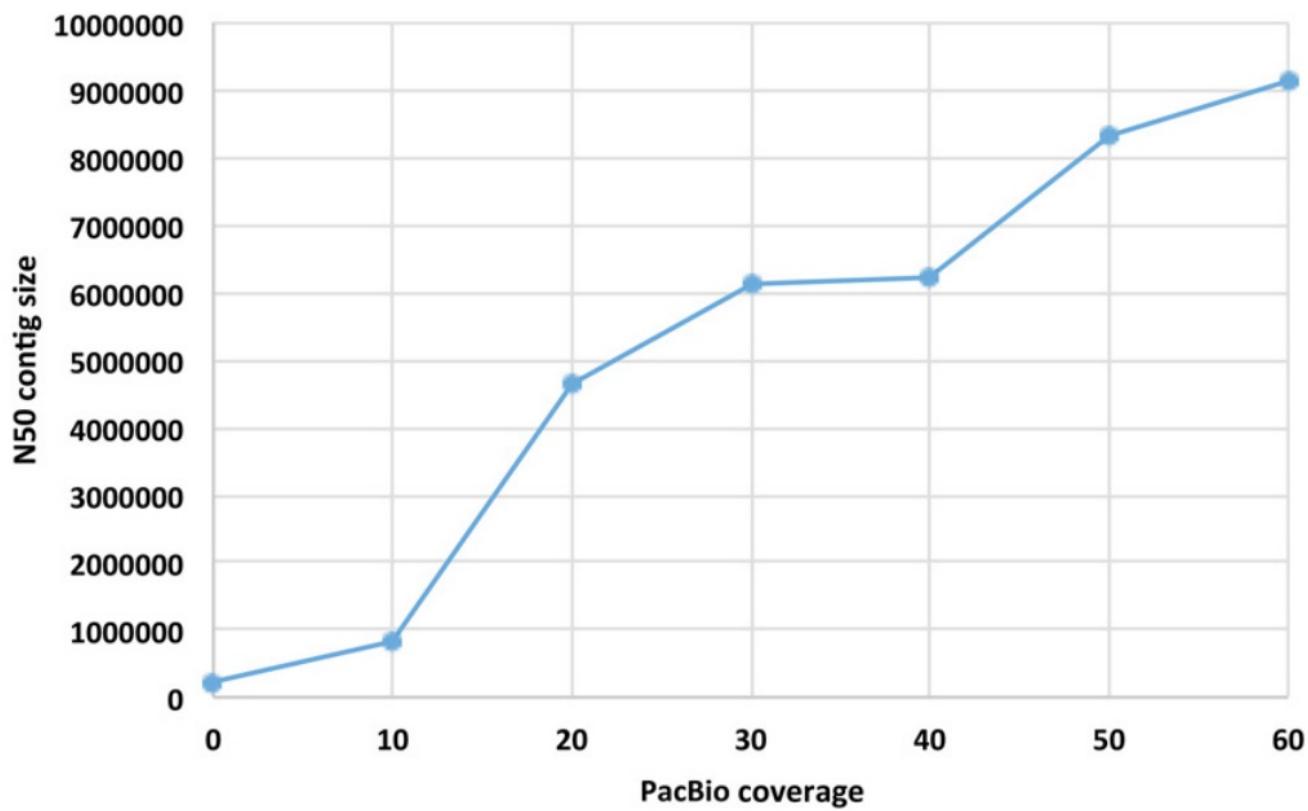
- A fast and accurate option, produces longer N50s and more BUSCO genes than other assemblers given the same data
- Combines de Bruijn graph and Overlap-Layout-Consensus



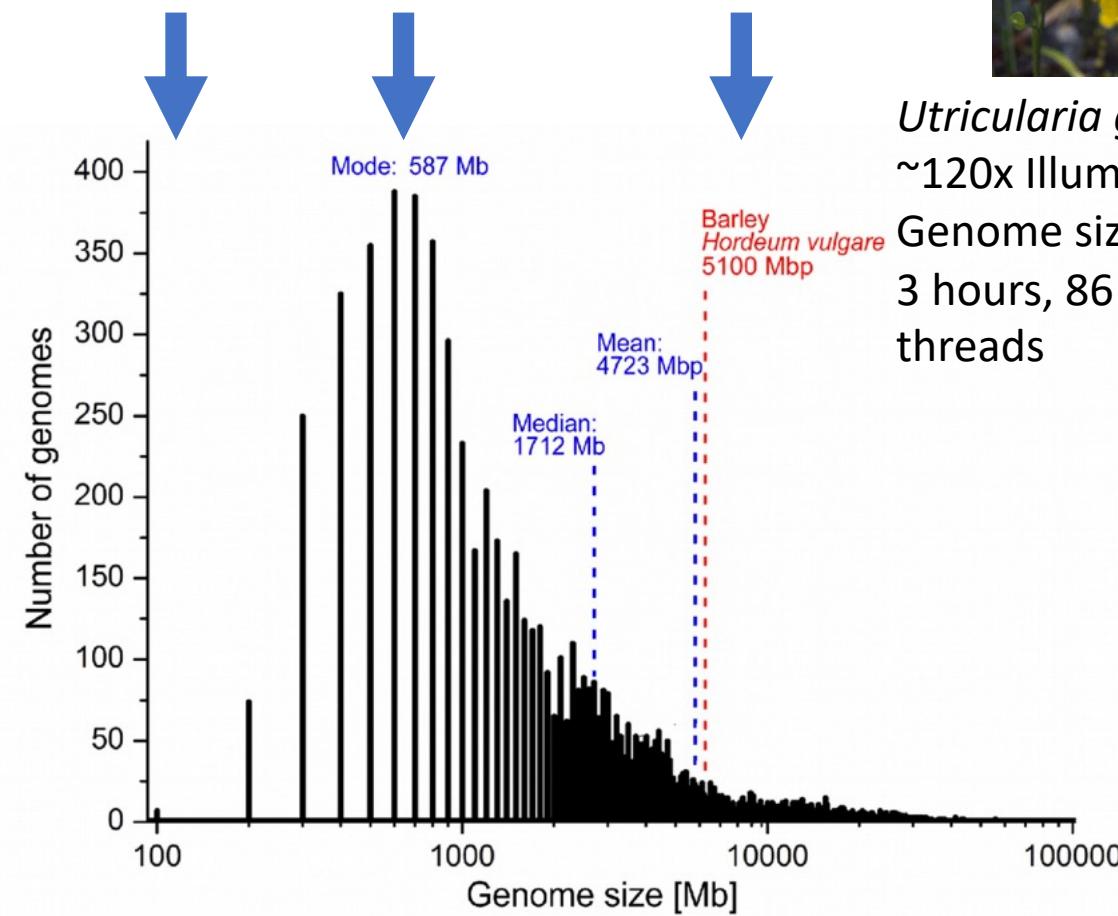
Assembler	Quast contig NGA50	Quast contig misassemblies	NGA50 scaffold (Kb)	Scaffold misassemblies/ MB
Allpaths-LG	28	175	261	0.03
SOAPdenovo2	8	369	1828	0.17
MaSuRCA	56	283	3445	0.19
Assemblies including some Long Read (LR) data				
MaSuRCA + 1× LR	70	256	4472	0.04
MaSuRCA + 2× LR	82	248	3704	0.21
MaSuRCA + 4× LR	102	246	4511	0.21



MaSuRCA –hybrid



Examples for requirements



Utricularia gibba
~120x Illumina, 100x PacBio
Genome size: 78 MB
3 hours, 86 GB storage, 16 threads



Costus spiralis
80x Illumina, 20x Nanopore
Genome size: 1 GB
3 days, 200 GB storage, 28 threads

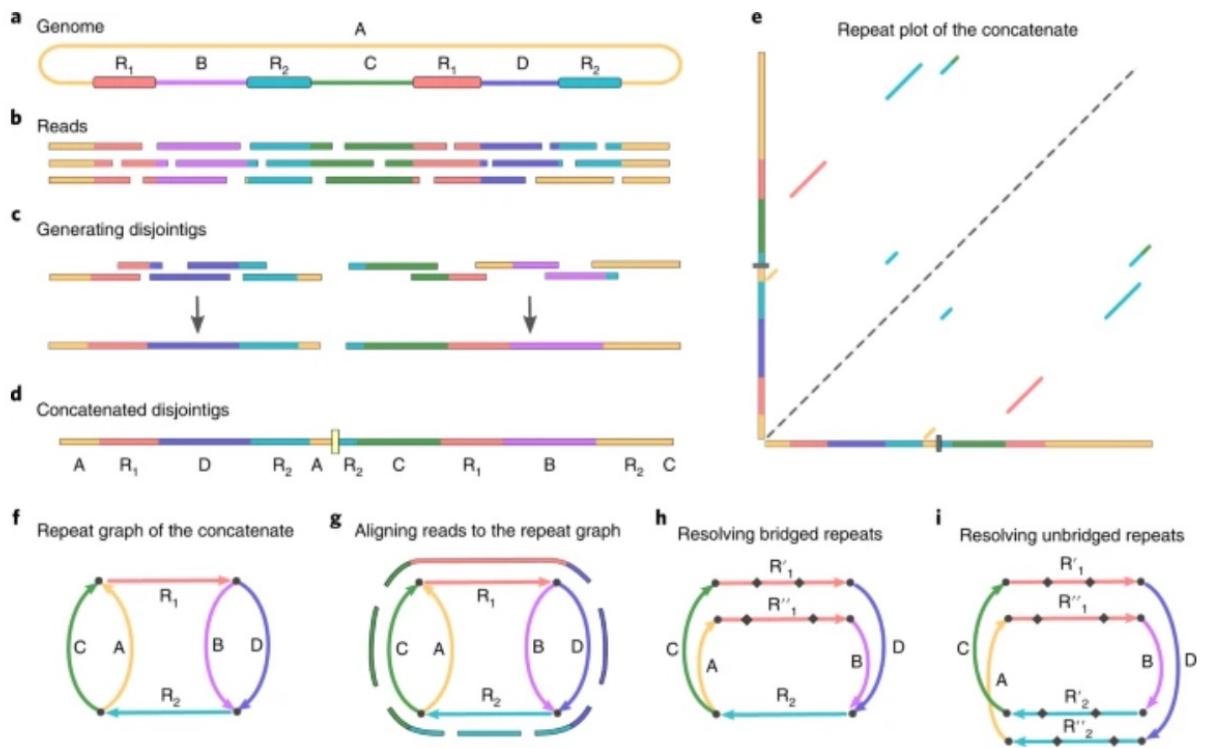


Calochortus venustus
40x Illumina, 4x Nanopore
Genome size: 5.5 GB
22 days, 3.7 TB storage, 28 threads

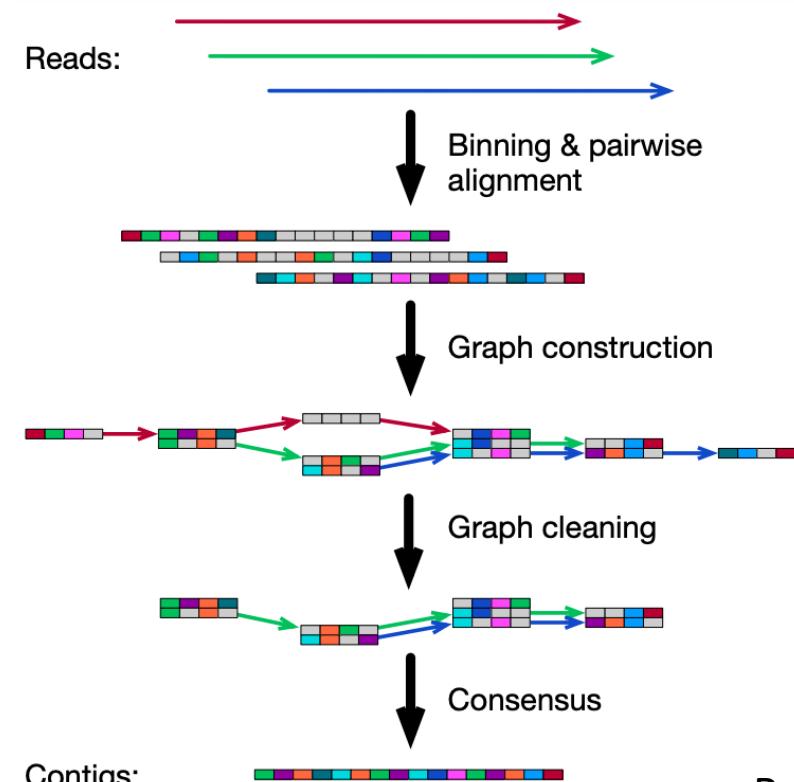
Long-read assemblers

- Many options out there; size of the genome can “make or break”

Flye – repeat graphs

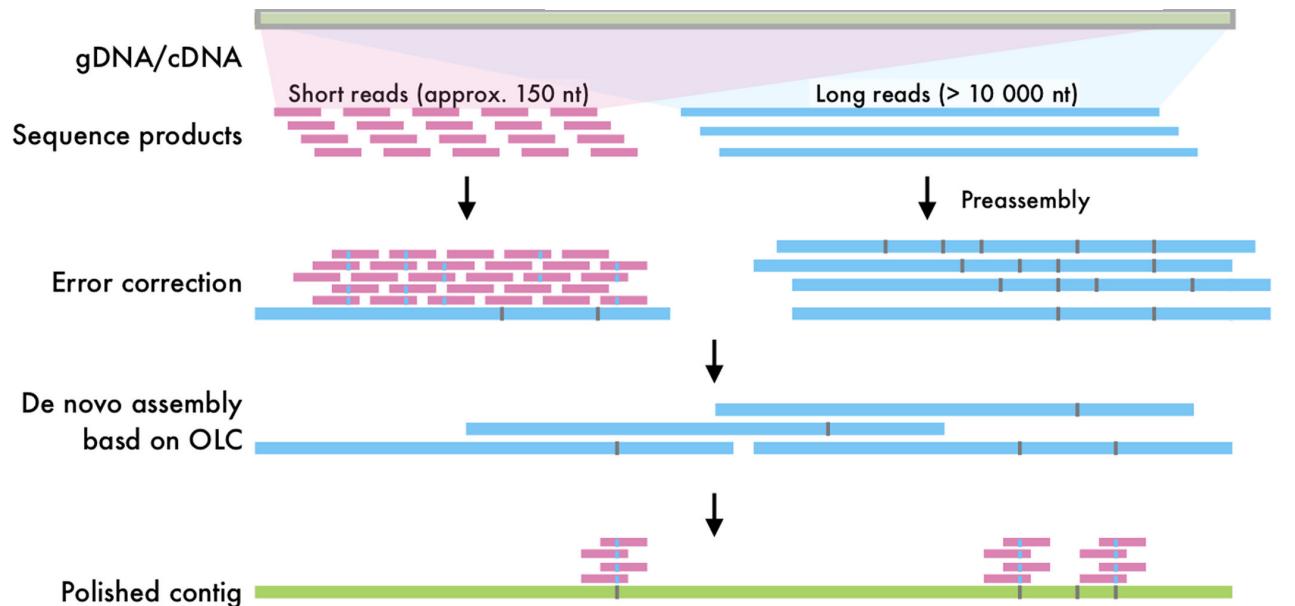
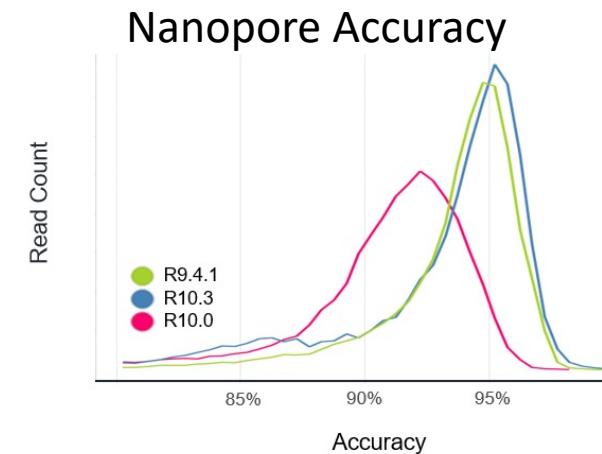


wtdbg2 – fuzzy-Brujin assembly graph



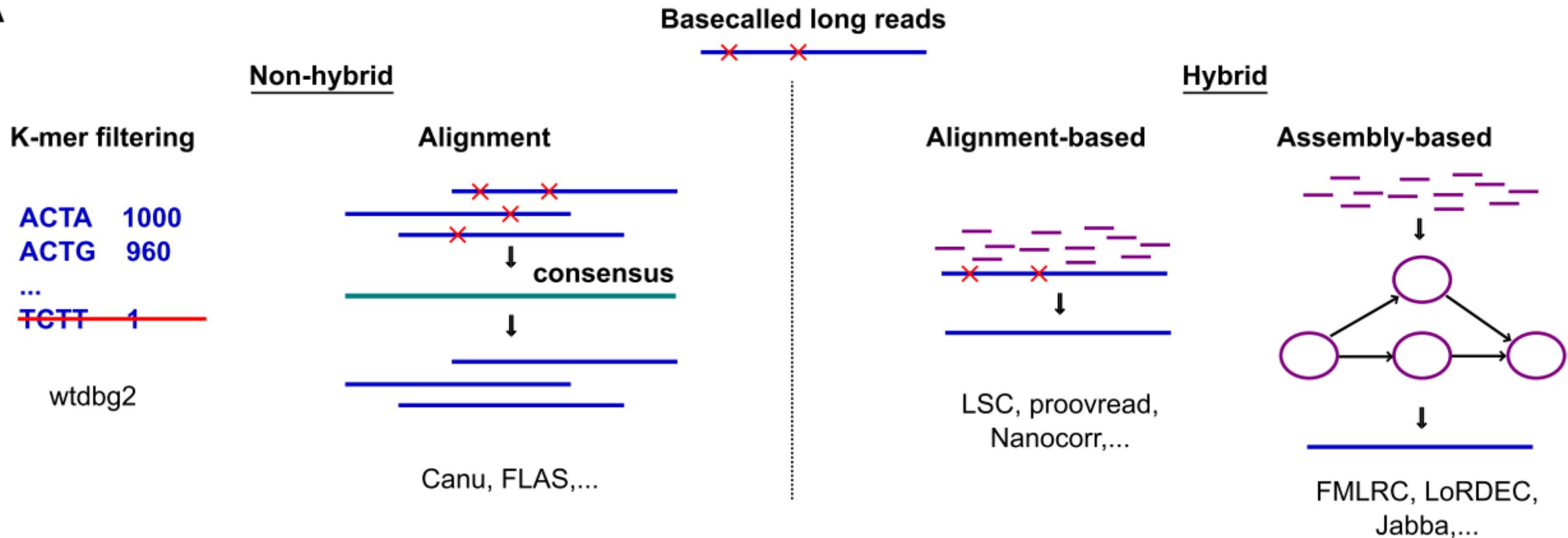
Is the error rate a problem?

- Definitely can be, but both Nanopore and PacBio have made recent improvements
- Recent PacBio HiFi reads (consensus sequences) are nearly as accurate as Illumina data
- How can we fix it?
 - Error correct the long reads prior to assembly
 - Polish the assembly afterwards
- Illumina data very useful but not mandatory



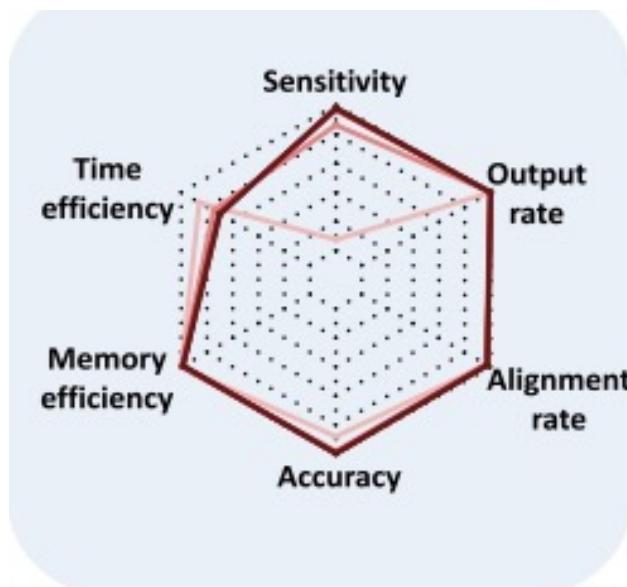
Error Correcting long-reads

A

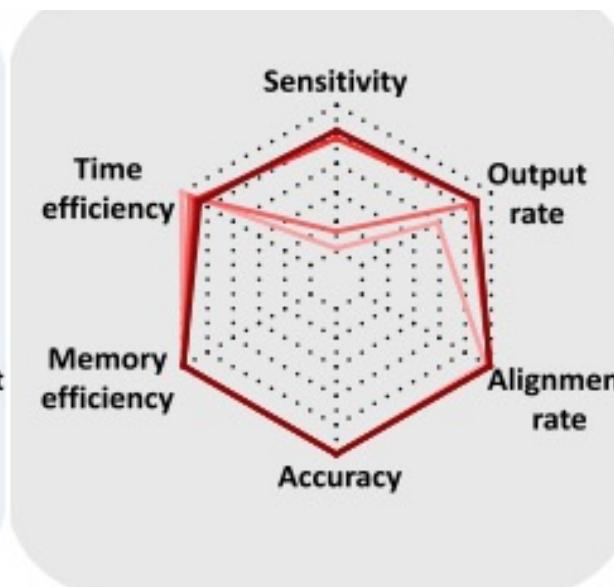


Error correction

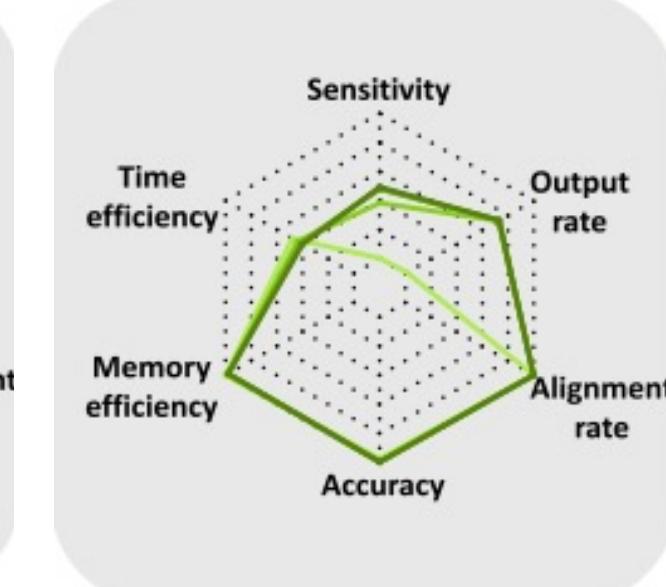
- Hybrid (using Illumina data) works better than just using depth of long-read
- Many factors to consider when comparing methods



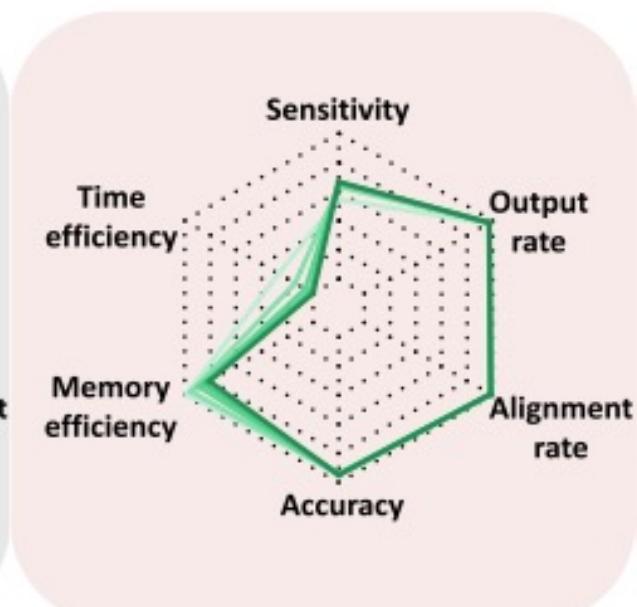
FMLRC



Jabba



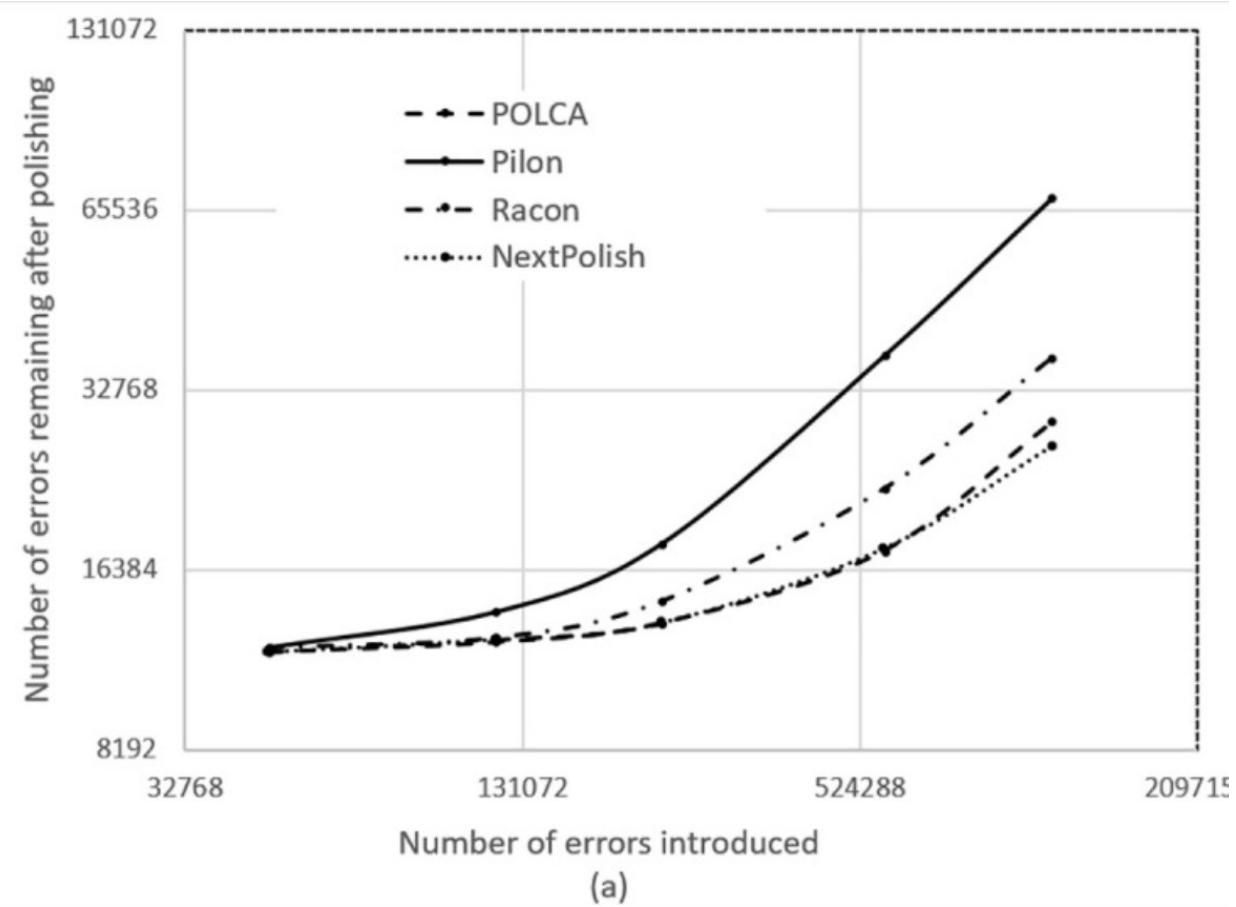
ECTools



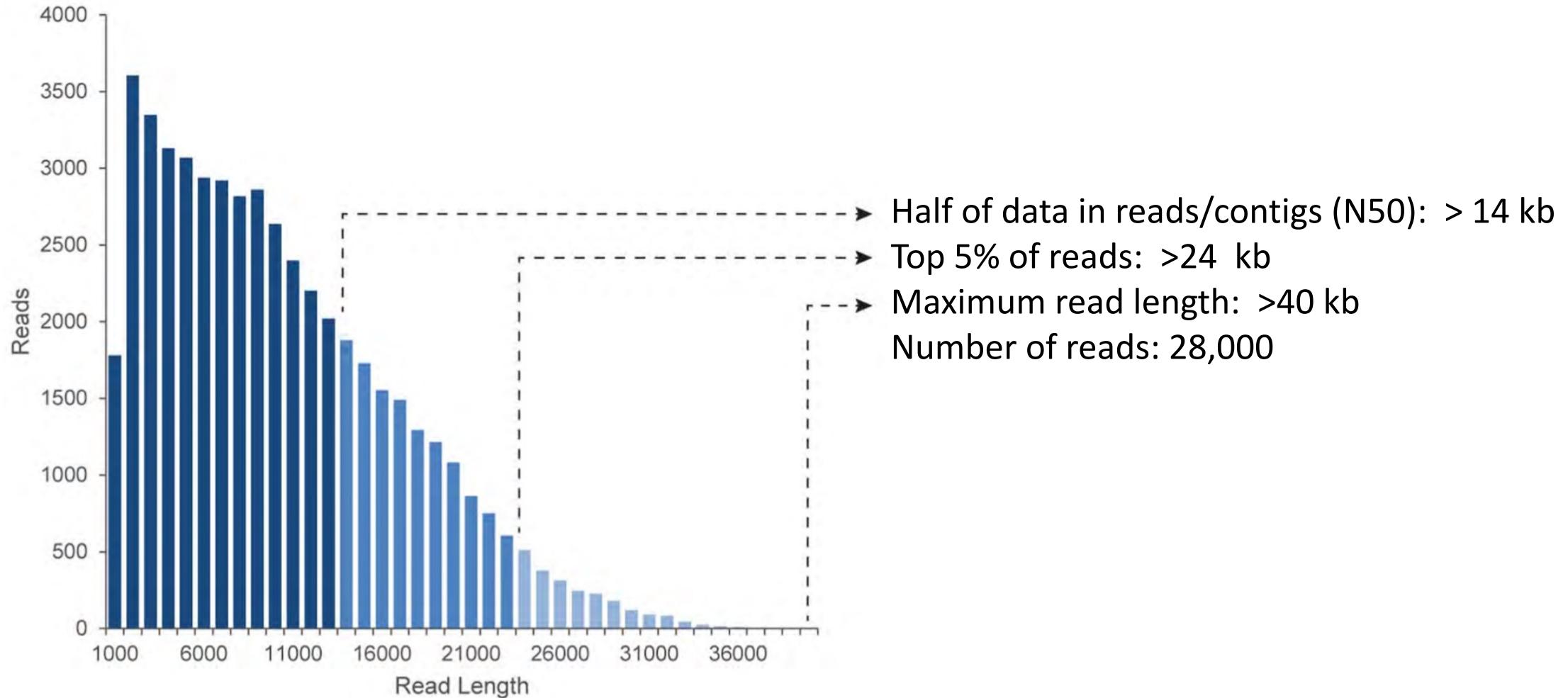
LSC

Polishing genome assemblies

- Many options for this as well; we'll use polca (part of the MaSuRCA package) to error correct the assemblies from the long-read assemblers

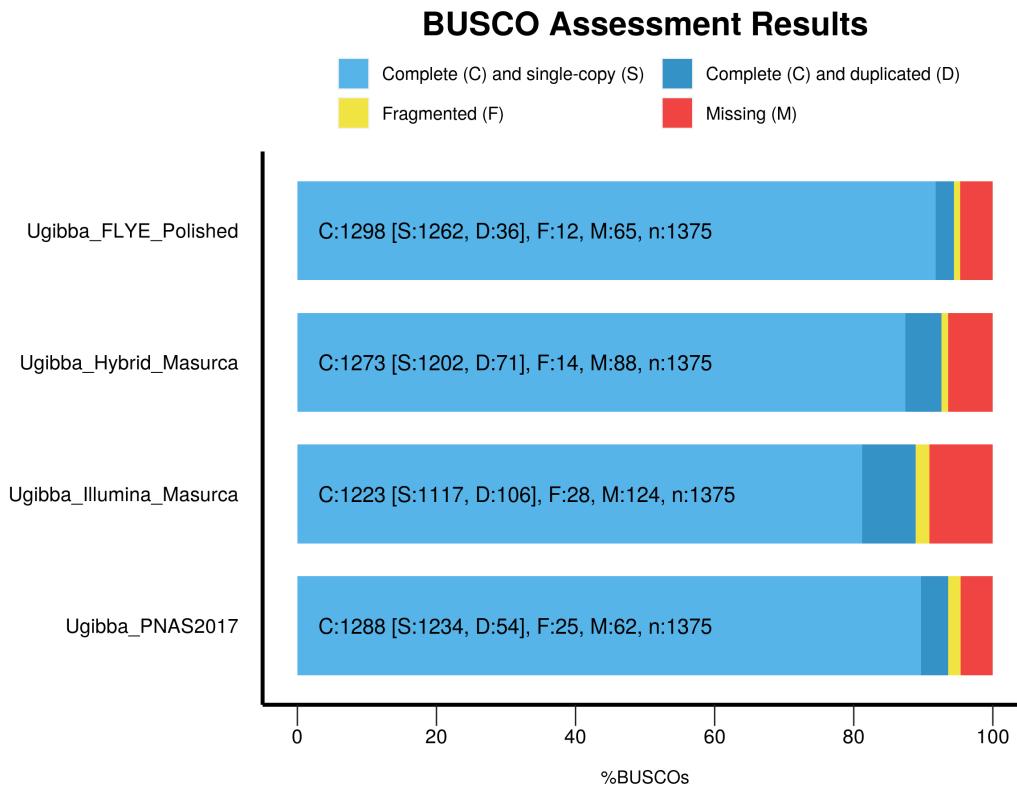


Quality Control Metrics



BUSCO

- Benchmarking Universal Single-Copy Orthologs
- A way to judge completeness of an assembly based on the number of single copy genes expected to find through BLAST
- Some lineage specific libraries, but most will use broader categories such as embryophyta or chlorophyta
 - Brassicales, Solanales, Poales, Fabales, or Eudicots are options



Long-read example data set

- For the walkthrough we will be using *Utricularia gibba*
 - Humped bladderwort; carnivorous aquatic species
 - Small genome size (76 MB) with only 3% of the genome non-coding
- Small data set that can be run on a local machine and all analyses should finish quickly
- Incorporates publicly available data using a high-quality genome assembly and RNA-Seq data for multiple organ types
 - Bladder, leaf, rhizoid, and stem



Utricularia gibba
Humped bladderwort



U. gibba traps

Tutorial

- Using many of the same programs I use for genome assembly and same steps, just a reduced data set

Assembly	Number of Contigs	N50 (bp)	Variants polished
MaSuRCA short-read only			
wtdbg2 long-read only			
Flye long-read only			
MaSuRCA hybrid (short-and long-read)			

Questions



@JLandisBotany



jbl256@cornell.edu