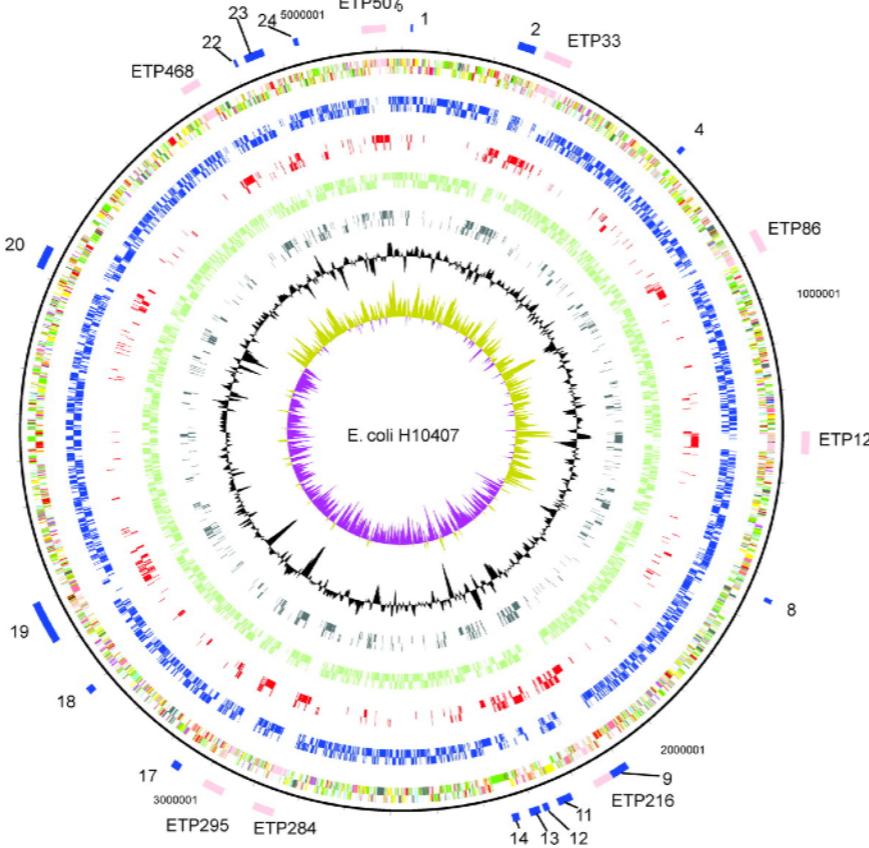
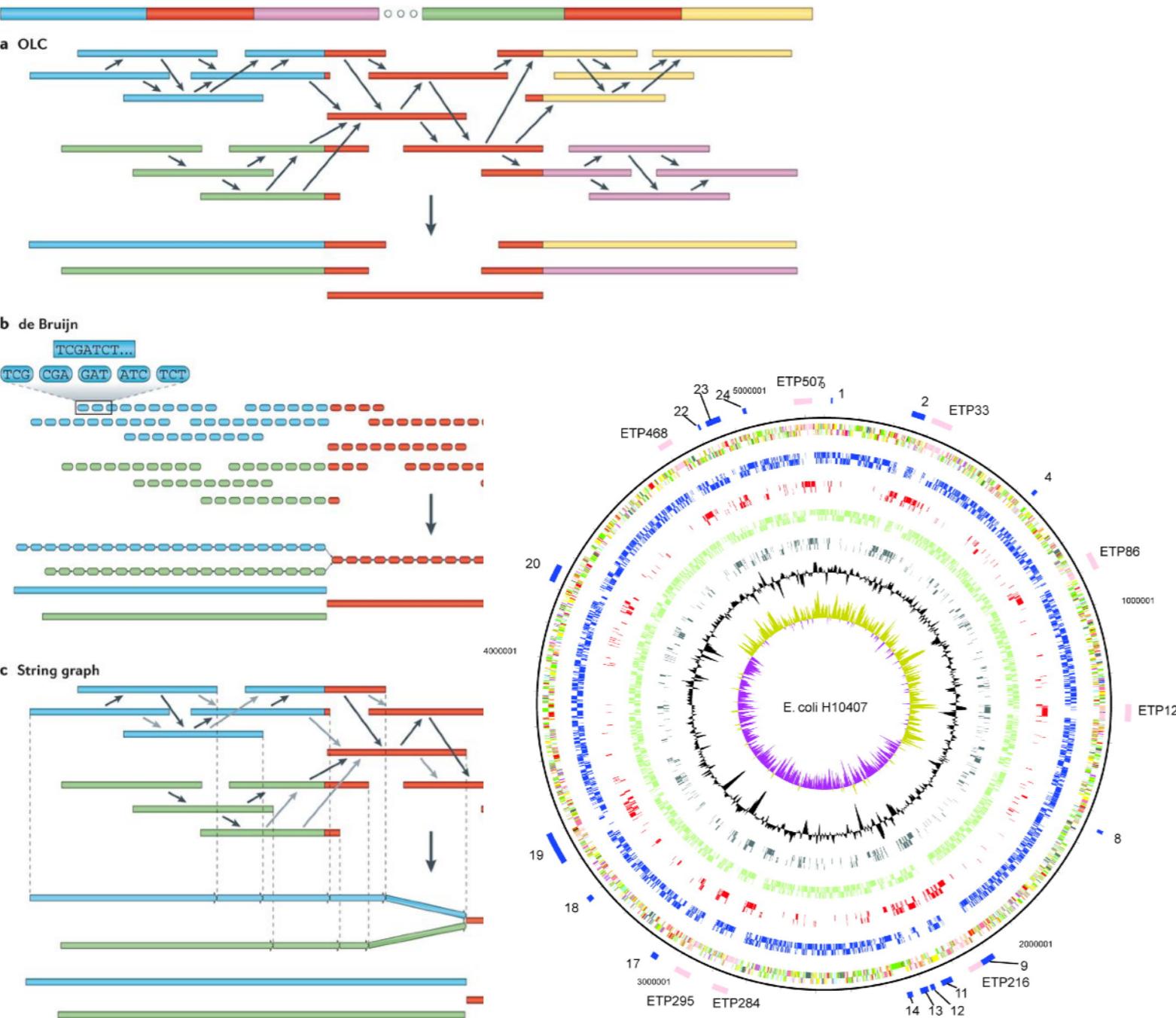


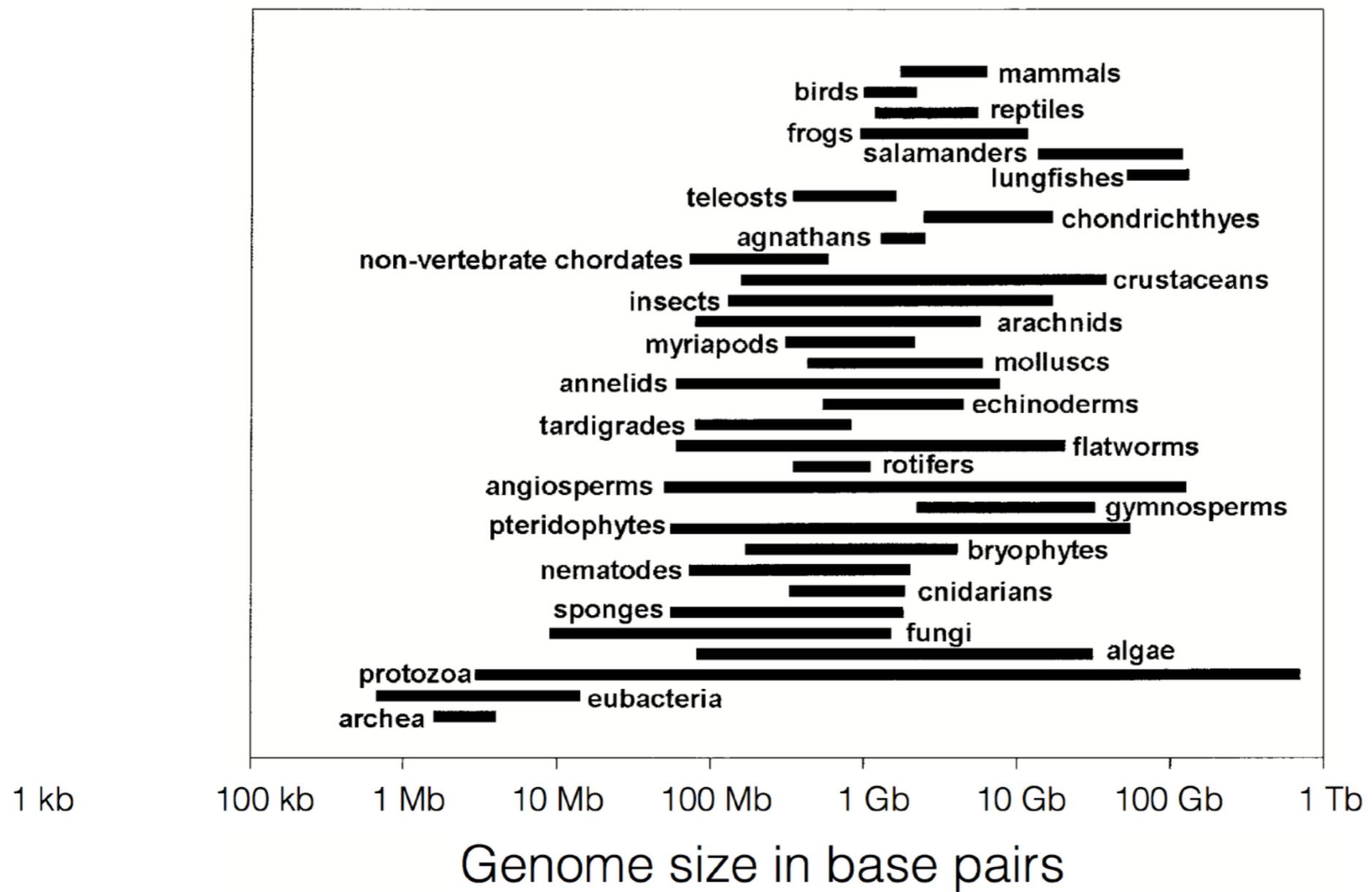
Assembly Algorithms



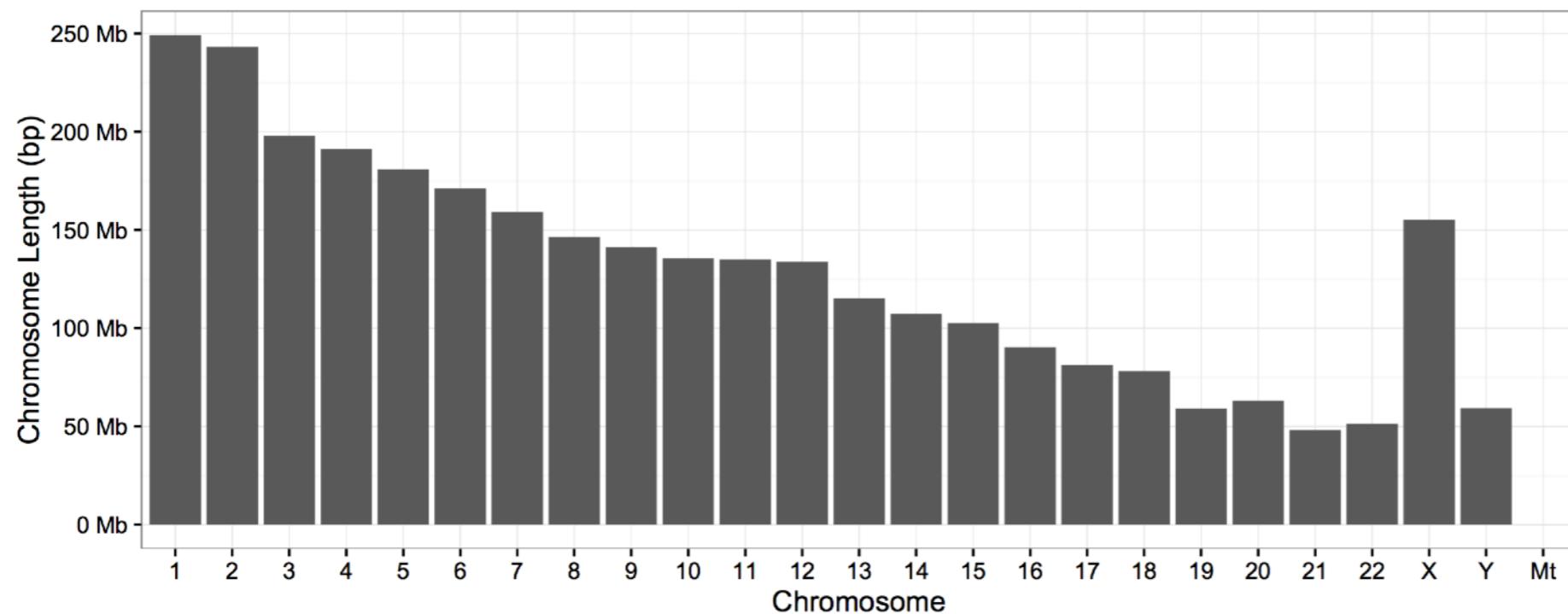
**HTS Workshop Genomics
& Transcriptomics
KAUST 2019**

Robert Lehmann
Octavio Salazar

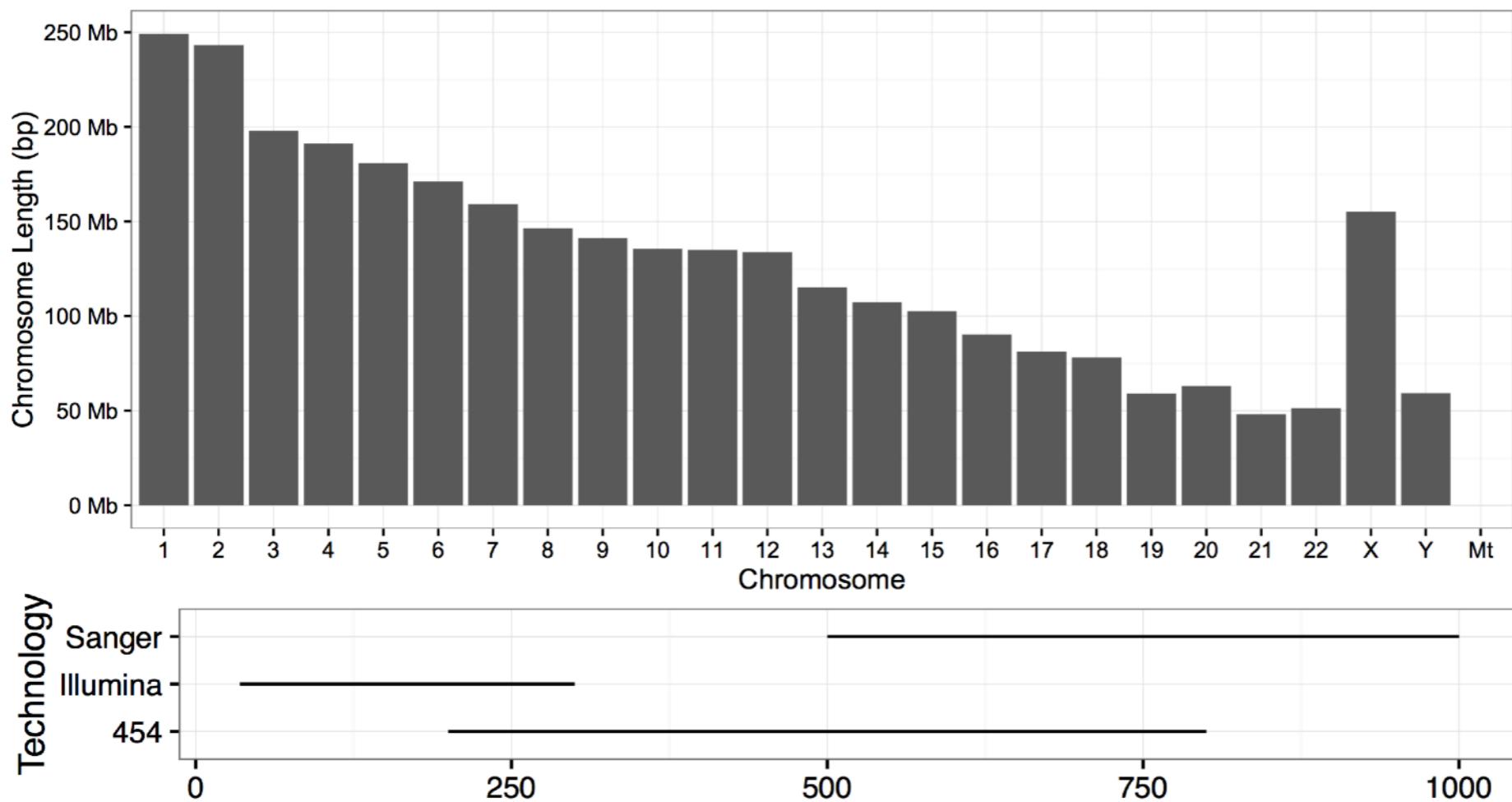
Complex organism = large genome?



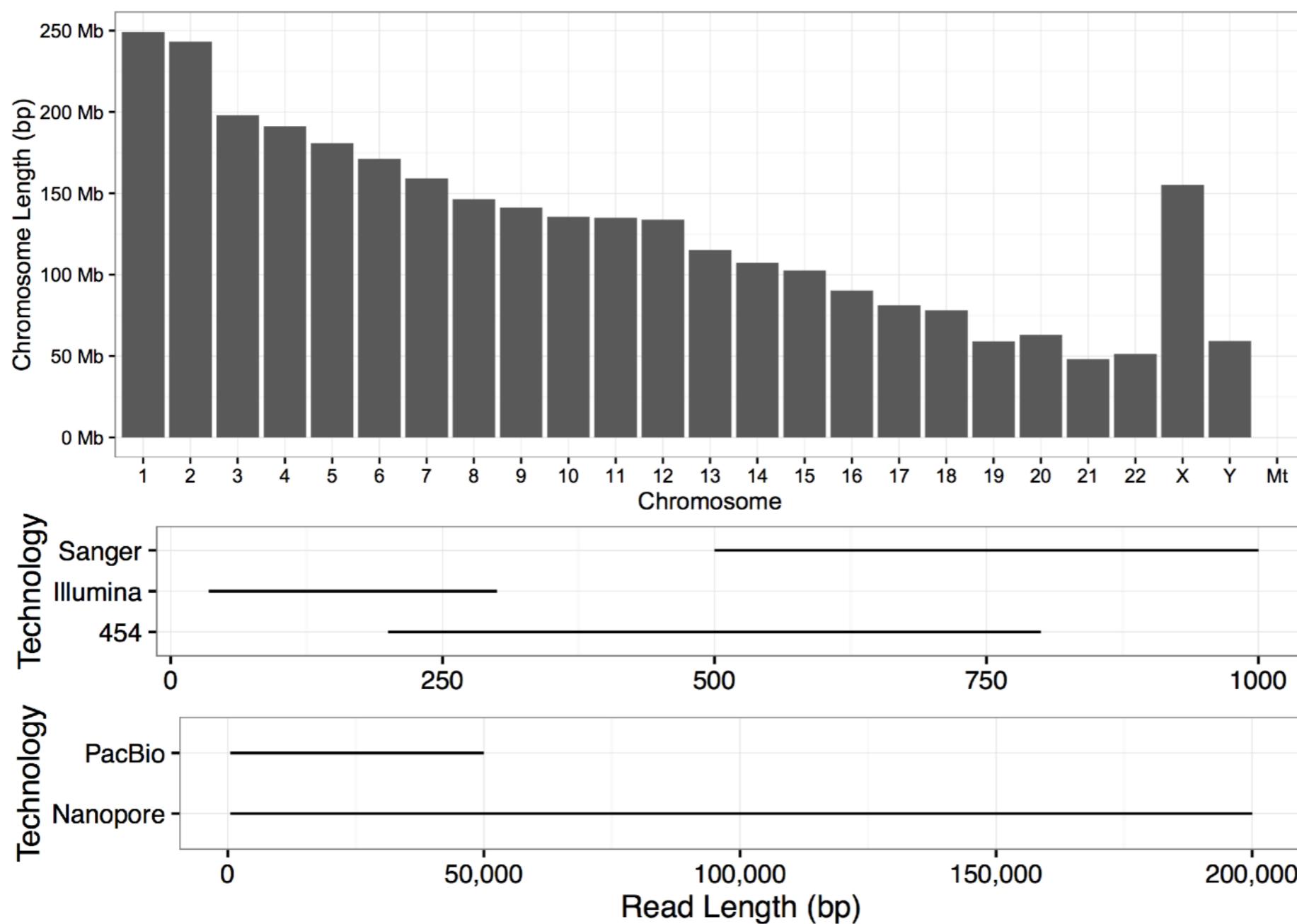
Read length and chromosome length



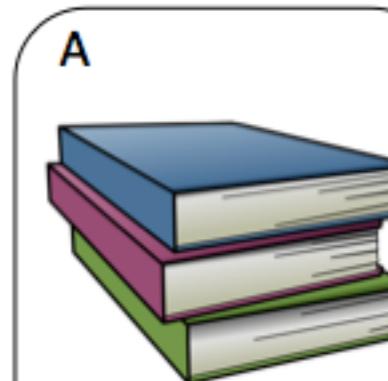
Read length and chromosome length



Read length vs. chromosome length



Assembling the pieces



B

This is an analogy for genome assembly. This page will be torn into horizontal strips.

C

This i
is a
s an

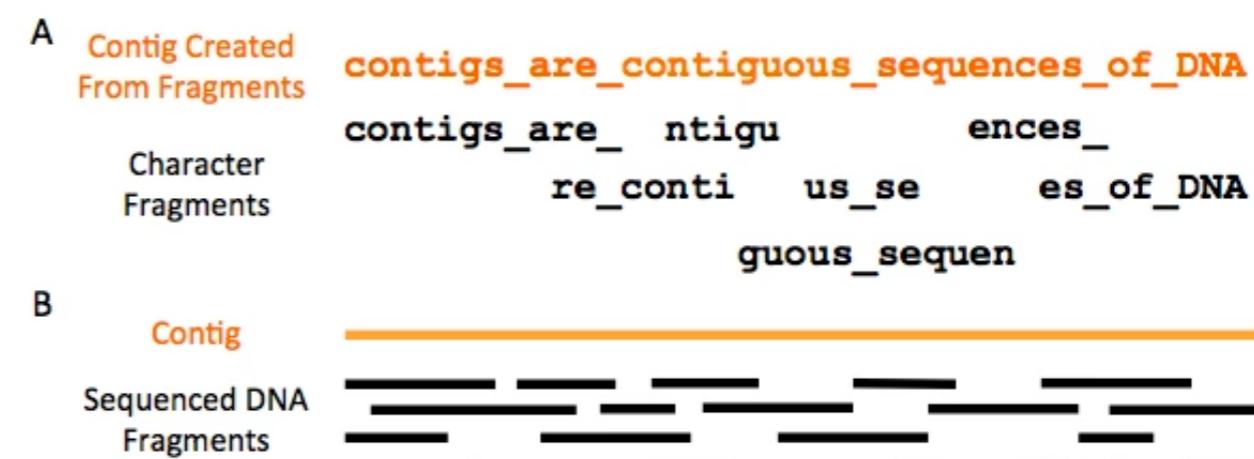
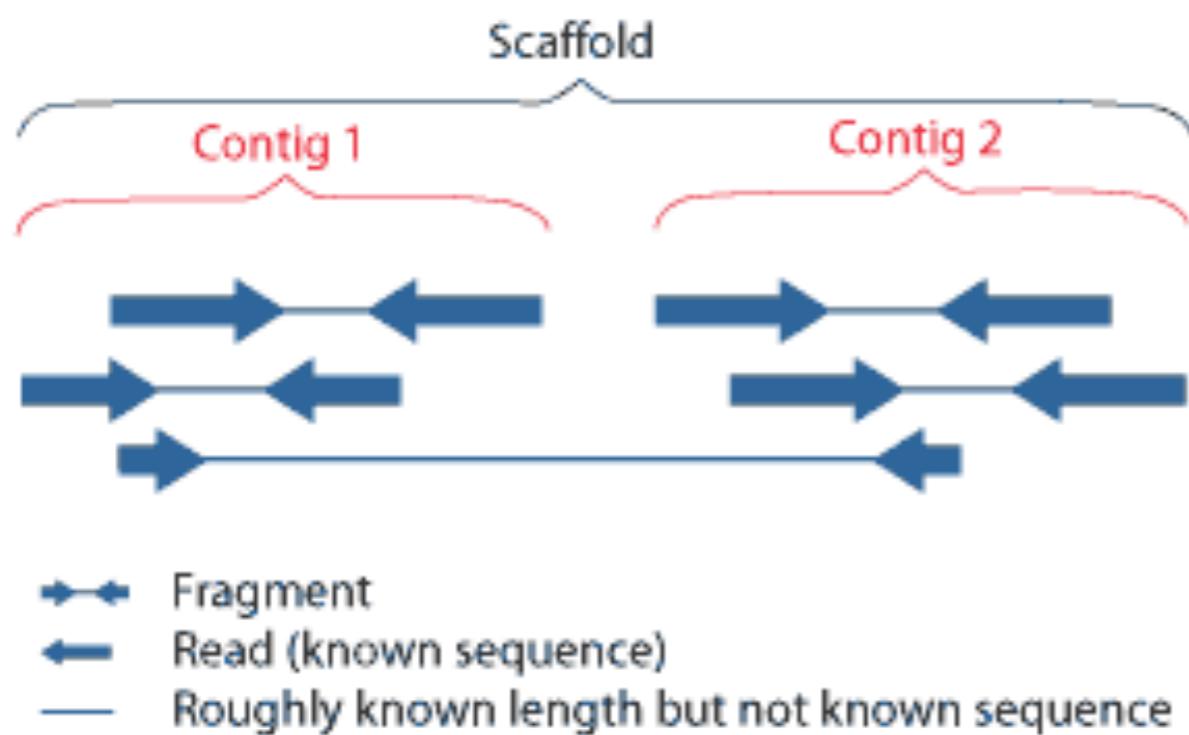
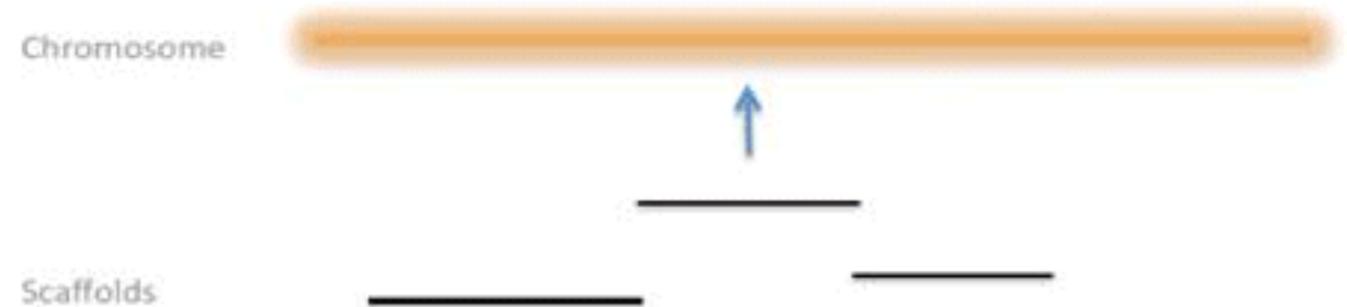


E

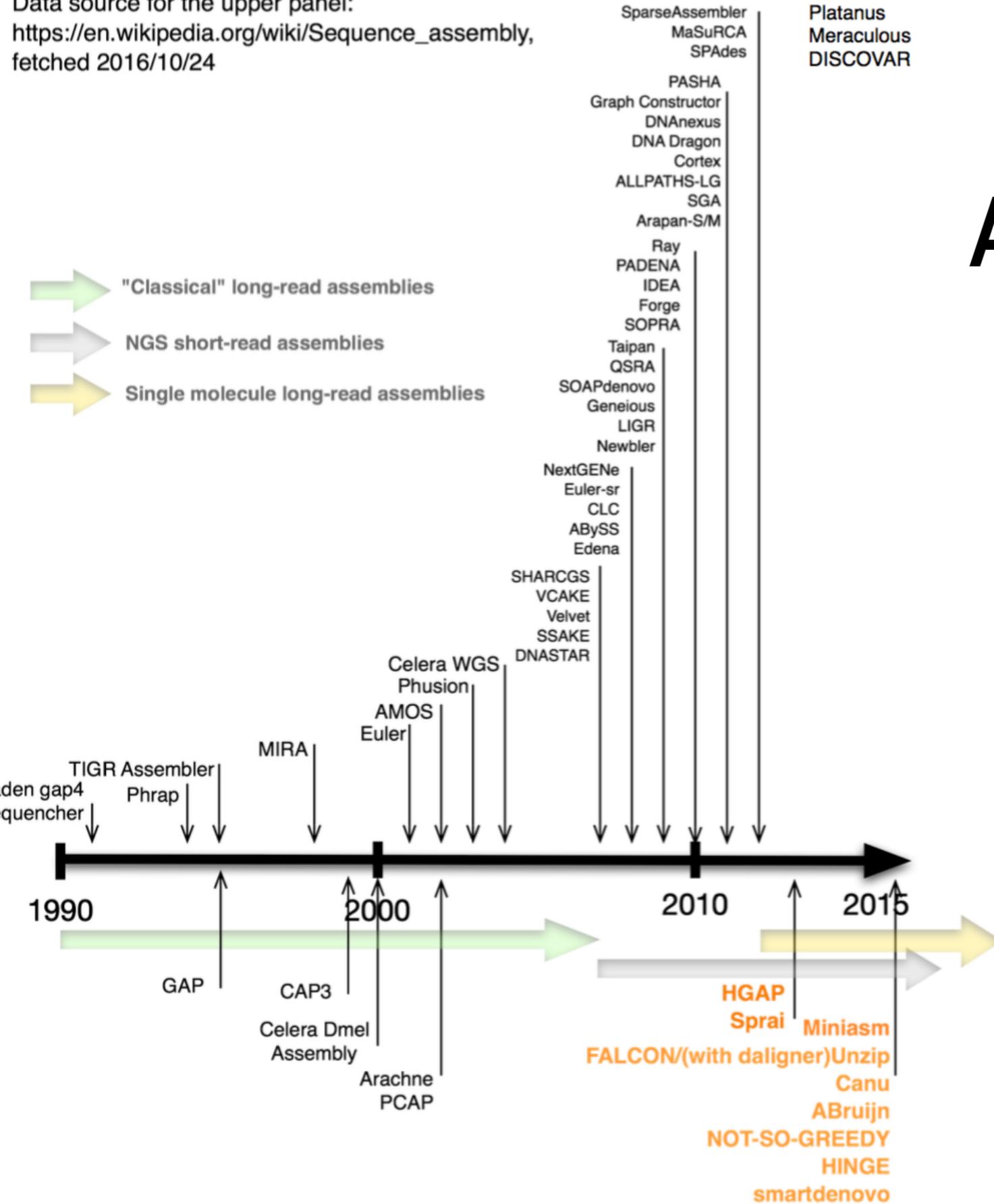
This is an

This i
is a
s an

Assembling the pieces



Data source for the upper panel:
https://en.wikipedia.org/wiki/Sequence_assembly,
fetched 2016/10/24



Lots of Assemblers developed

Algorithm Feature	OLC Assemblers	DBG Assemblers
<i>Modeled features of reads</i>		
Base substitutions	CABOG	Euler, AllPaths, SOAP
Homopolymer miscount		Euler
Concentrated error in 3' end	Newbler	
Flow space	Shorty	Velvet
Color space		
<i>Removal of erroneous reads</i>		
Based on K-mer frequencies		Euler, Velvet, AllPaths
Based on K-mer freq and QV		AllPaths
For multiple values of K	CABOG	AllPaths
By alignment to other reads		
By alignment and QV		
<i>Correction of erroneous base calls</i>		
Based on K-mer frequencies	CABOG	Euler, SOAP
Based on Kmer freq and QV		AllPaths
Based on alignments		
<i>Approaches to graph construction</i>		
Implicit		
Reads as graph nodes	CABOG, Newbler, Edena	Euler, Velvet, ABySS, SOAP
K-mers as graph nodes		AllPaths
Simple paths as graph nodes		Euler
Multiple values of K		
Multiple overlap stringencies		
<i>Approaches to graph reduction</i>		
Filter overlaps	CABOG	Euler, Velvet, SOAP
Greedy contig extension		Euler, Velvet, AllPaths, SOAP
Collapse simple paths	CABOG, Newbler	
Erosion of spurs	CABOG, Edena	
Transitive overlap reduction	Edena	
Bubble smoothing	Edena	Euler, Velvet, SOAP
Bubble detection		AllPaths
Reads separate tangled paths		Euler, SOAP
Break at low coverage		Velvet, SOAP
Break at high coverage	CABOG	Euler
High coverage indicates repeat	CABOG	Velvet
Special use of long reads	Shorty	Velvet
<i>Graph partitions</i>		
Partition by K-mers		ABySS
Partition by scaffolds		AllPaths
<i>Uses for mate pairs</i>		
Constrain path searches		Euler, Velvet, AllPaths
Guide path selection		Euler, Allpaths
Detect misassembled contigs	CABOG, Shorty	Velvet, ABySS, SOAP
Merge contigs or fill gaps	CABOG, Shorty	SOAP
Transitive link reduction	CABOG	Velvet, SOAP
Detect, avoid repeat contigs	CABOG	Euler, Velvet, AllPaths, SOAP
Create scaffolds	CABOG, Shorty	

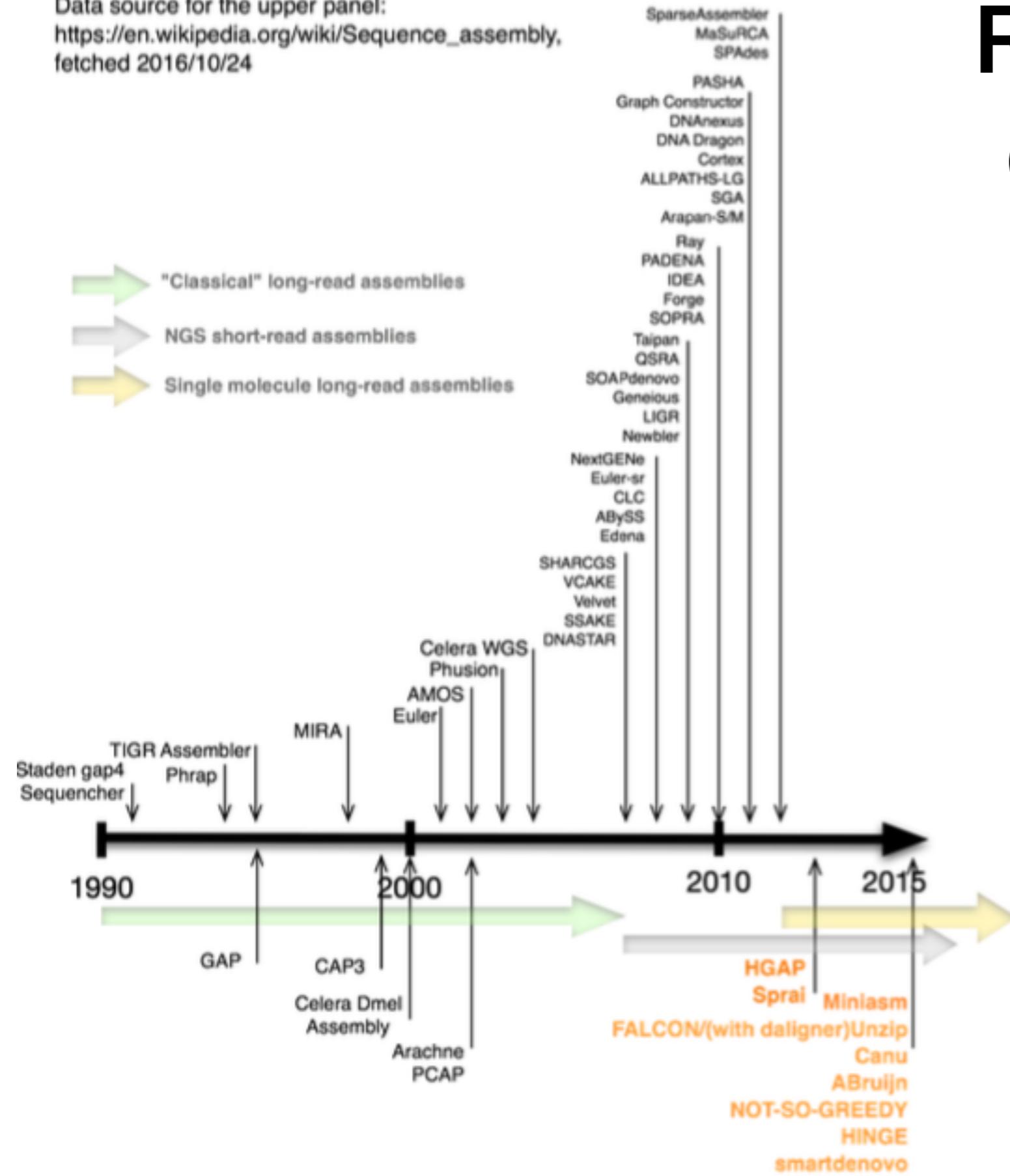
Different approaches to key problems

Miller, J. R., Koren, S., & Sutton, G. (2010). Assembly algorithms for next-generation sequencing data, 95(6), 315–327

Why so many different strategies?

- We don't know the true sequence - which algorithm is right?
- Simulated data often doesn't match reality
- Assemblers depend highly on input data
- Speed optimization vs. correctness
- ...

Data source for the upper panel:
https://en.wikipedia.org/wiki/Sequence_assembly,
fetched 2016/10/24



Read length determines strategy

- **Classical long read:**
 - Sanger ~600 bp
 - Overlap - Layout - Consensus
- **NGS:**
 - Illumina ~ 150 bp
 - De-Brujin graph
- **Single molecule long read:**
 - Pacbio/Nanopore/Minion ~ 10kb
 - Overlap - Layout - Consensus

Overlap - Layout - Consensus

Assembler paradigms

CTAGGCCCTCAATTTTT
GGCGTCTATATCT
CTCTAGGCCCTCAATTTTT
TCTATATCTCGGCTCTAGG
GGCTCTAGGCCCTCATTTTTT
CTCGGCTCTAGCCCCTCATTTT
TATCTCGACTCTAGGCCCTCA
GGCGTCGATATCT
TATCTCGACTCTAGGCC
GGCGTCTATATCTCG

Overlap - Layout - Consensus

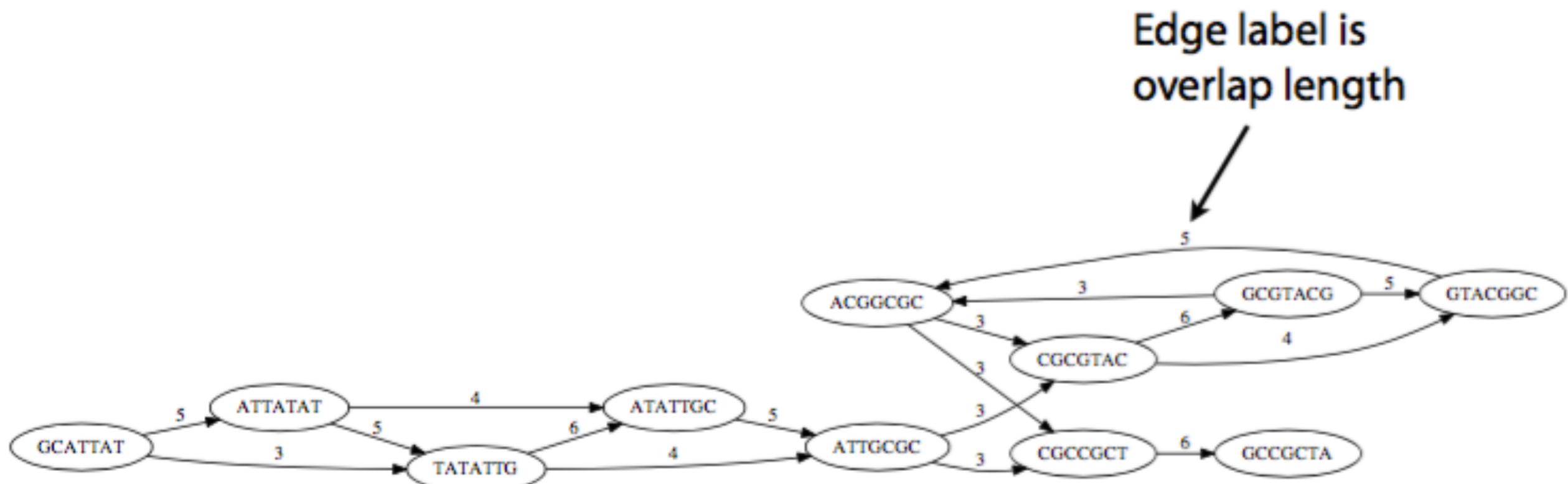
Assembler paradigms

CTCGGCTCTAGCCCCTCATT
||| ||| | | | | | | | | | |
GGCTCTAGGCCCTCATTTTT

- CTAGGCCCTCAATTTTT
- GGCGTCTATATCT
- CTCTAGGCCCTCAATTTTT
- TCTATATCTCGGCTCTAGG
- GGCTCTAGGCCCTCATTTTTT
- CTCGGCTCTAGCCCCTCATT
- TATCTCGACTCTAGGCCCTCA
- GGCGTCGATATCT
- TATCTCGACTCTAGGCC
- GGCGTCTATATCTCG

Overlap - Layout - Consensus

Assembler paradigms



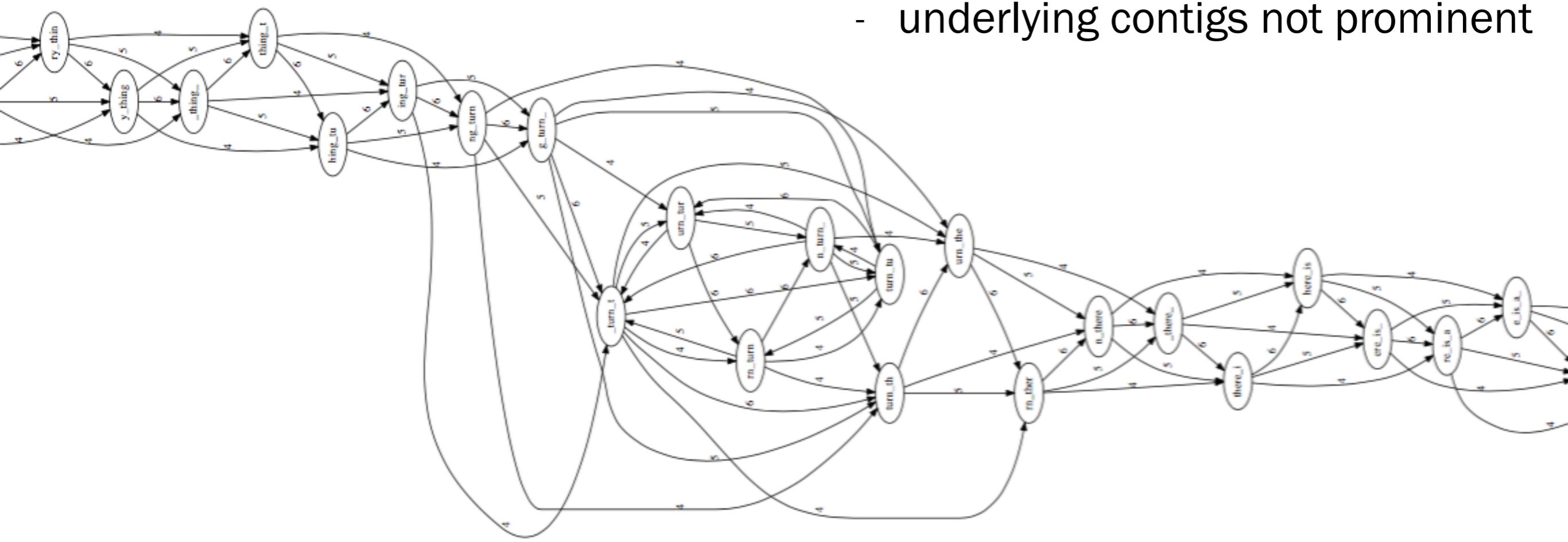
Original string: GCATTATATATTGCGCGTACGGCGCCGCTACA

Overlap - Layout - Consensus

Assembler paradigms

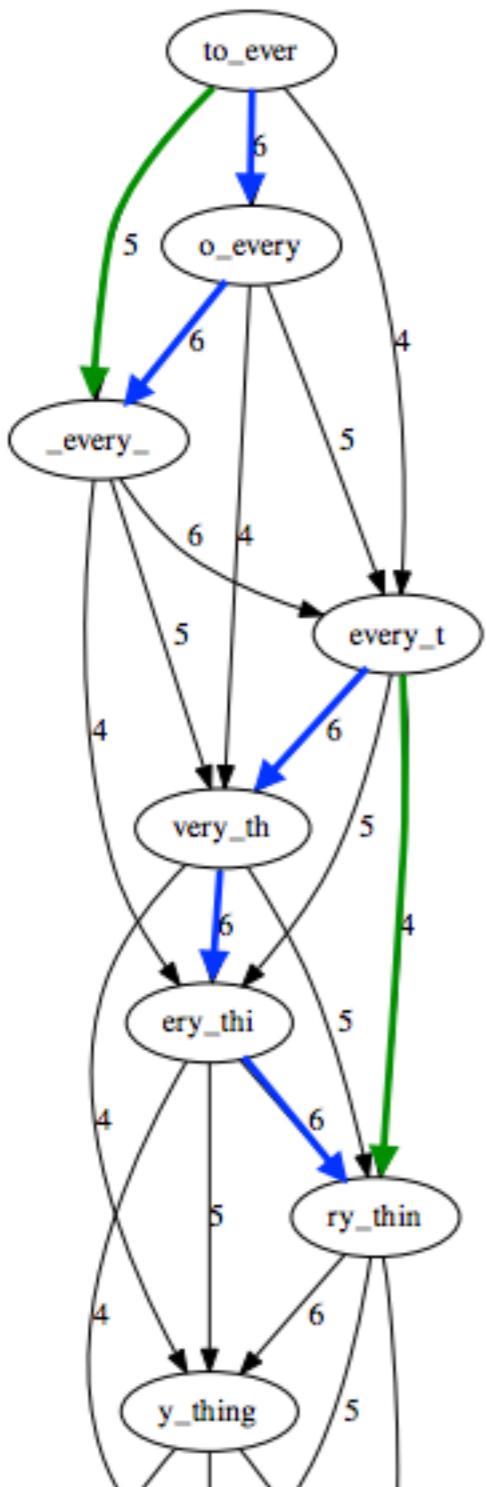
[to_every_thing_turn_turn_turn_there_is_a_season](#)

- Complex overlap graph
- underlying contigs not prominent



Overlap - Layout - Consensus

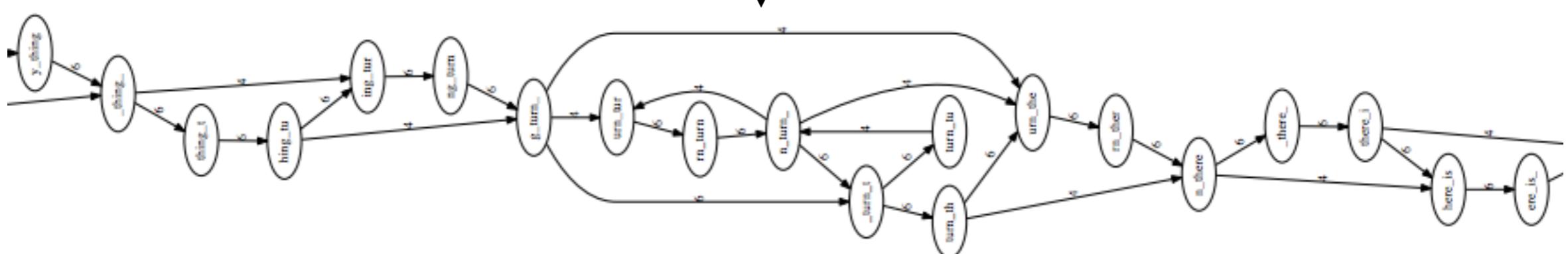
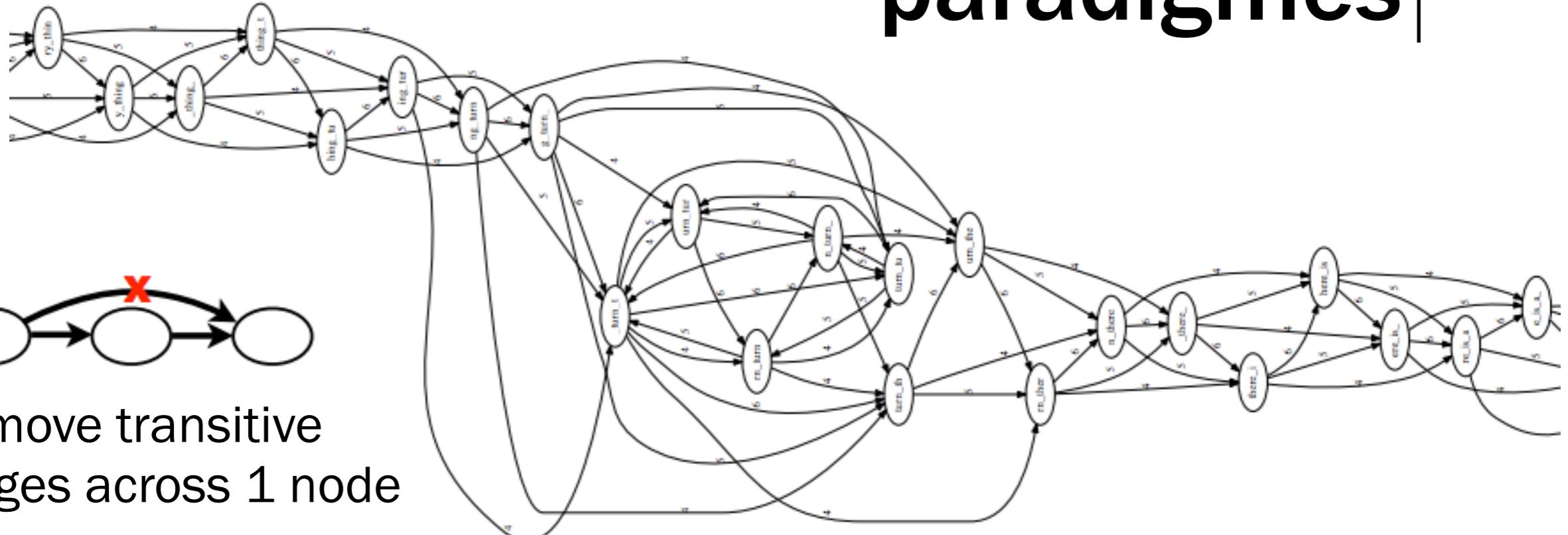
Assembler paradigms



- Some edges are implied by others
- **green** edge can be inferred by **blue**

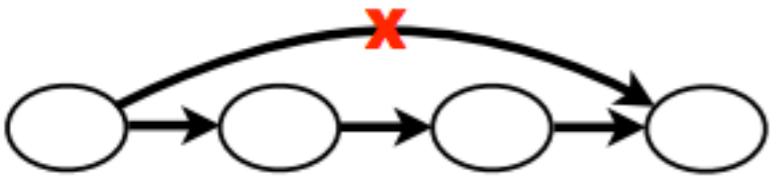
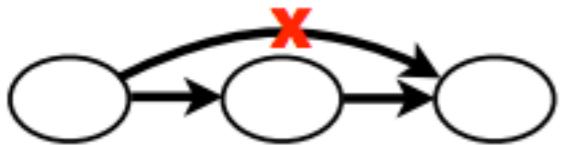
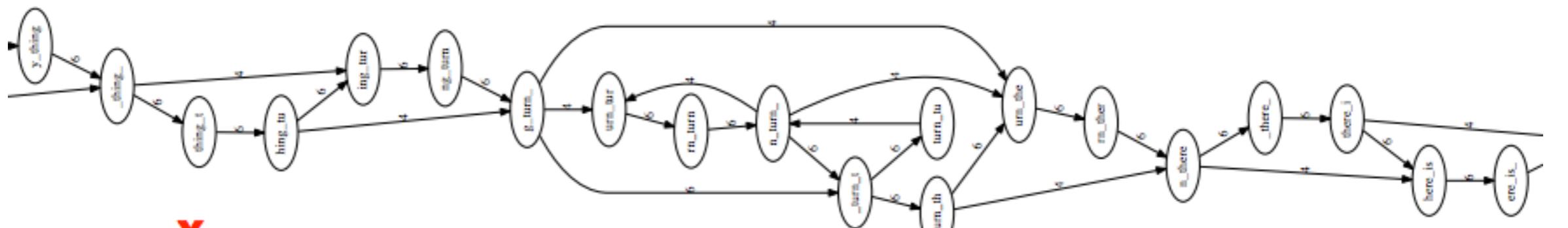
Overlap - Layout - Consensus

Assembler paradigms



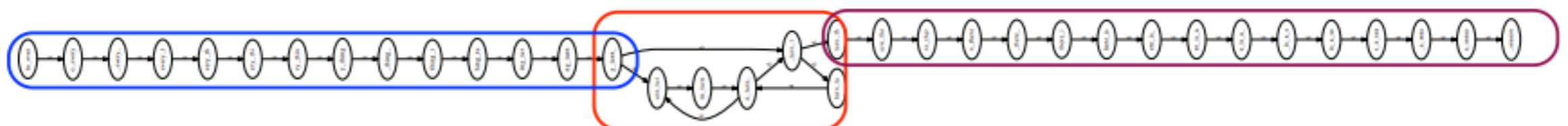
Overlap - Layout - Consensus

Assembler paradigms



Overlap - Layout - Consensus

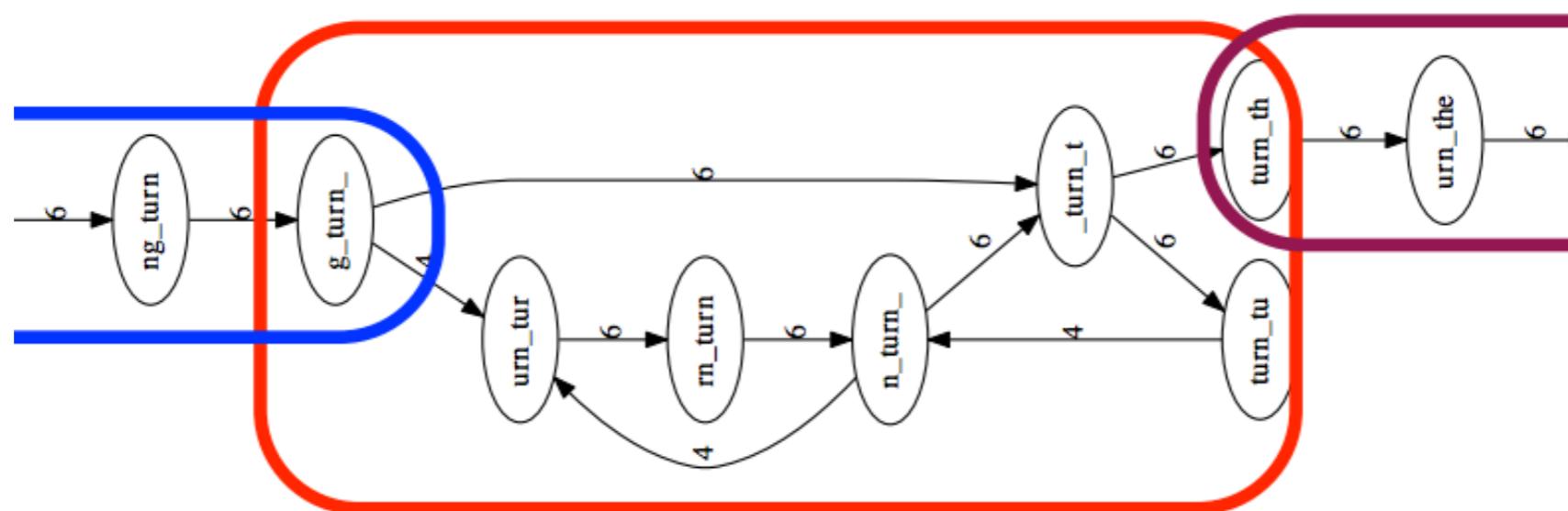
Assembler paradigms



Contig 1
to_every_thing_turn_

Contig 2
turn_there_is_a_season

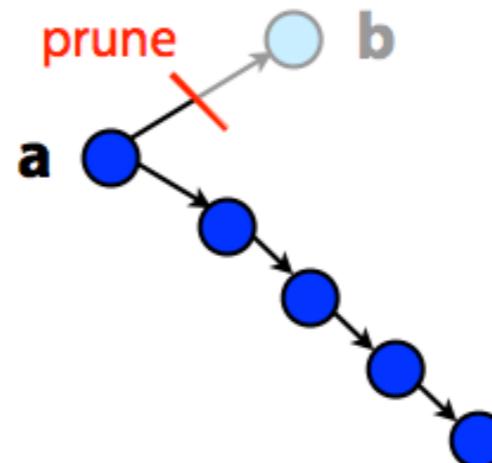
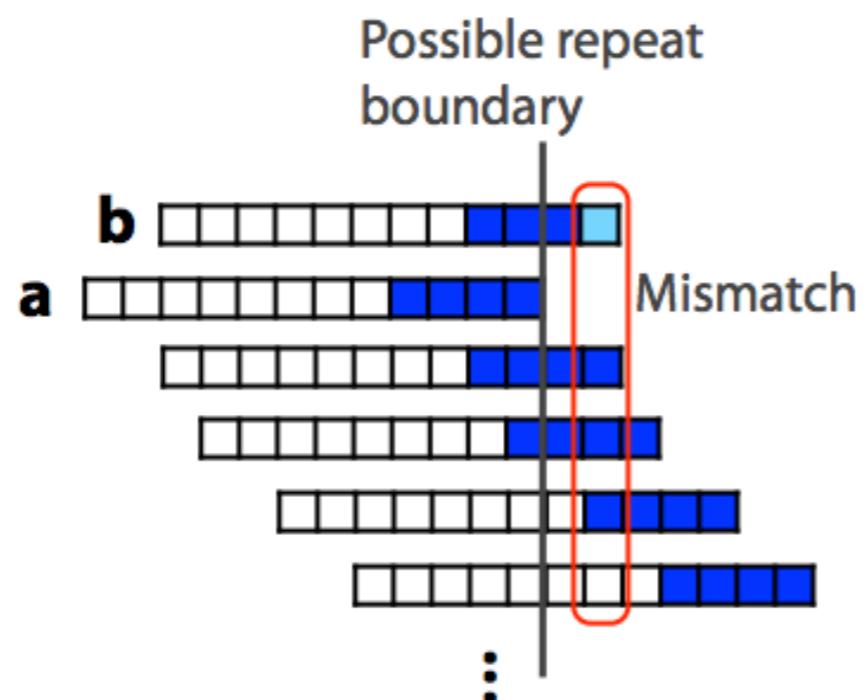
Unresolvable repeat



Overlap - Layout - Consensus

Assembler paradigms

Must handle subgraphs that are spurious, e.g. because of sequencing error



Mismatch could be due to sequencing error or repeat. Since the path through **b** ends abruptly we might conclude it's an error and prune **b**.

Overlap - Layout - Consensus

Assembler paradigms

TAGATTACACAGATTACTGA TTGATGGCGTAA CTA
TAGATTACACAGATTACTGACTTGATGGCGTAAACTA
TAG TTACACAGATTATTGACTT~~C~~ATGGCGTAA CTA
TAGATTACACAGATTACTGACTTGATGGCGTAA CTA
TAGATTACACAGATTACTGACTTGATGGCGTAA CTA

↓ ↓ ↓ ↓

TAGATTACACAGATTACTGACTTGATGGCGTAA CTA

Take reads that make up a contig and line them up

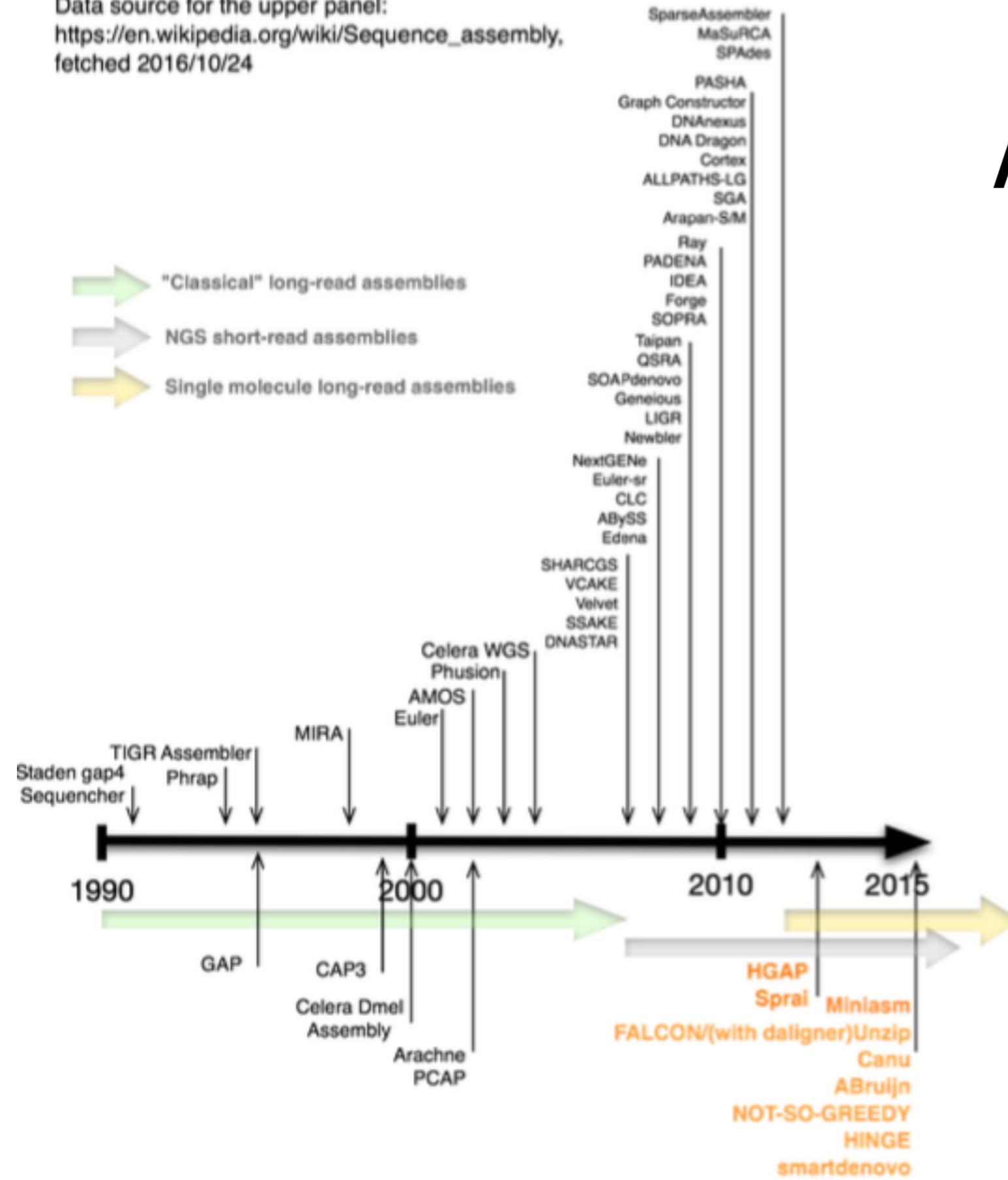
Take *consensus*, i.e. majority vote

Complications: (a) sequencing error, (b) ploidy

Lots of

Assemblers developed

Data source for the upper panel:
https://en.wikipedia.org/wiki/Sequence_assembly,
fetched 2016/10/24

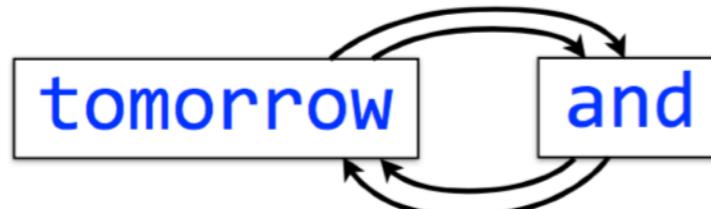


- **Classical long read:**
 - Sanger ~600 bp
 - Overlap - Layout - Consensus
- **NGS:**
 - Illumina ~ 150 bp
 - De-Brujin graph
- **Single molecule long read:**
 - Pacbio/Nanopore/Minion ~ 10kb
 - Overlap - Layout - Consensus

Assembler paradigmes

How to represent
complex string as
ordered
fragments?

“tomorrow and tomorrow and tomorrow”

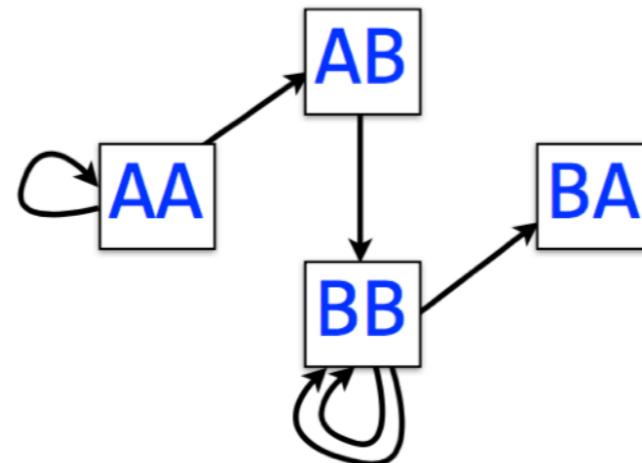
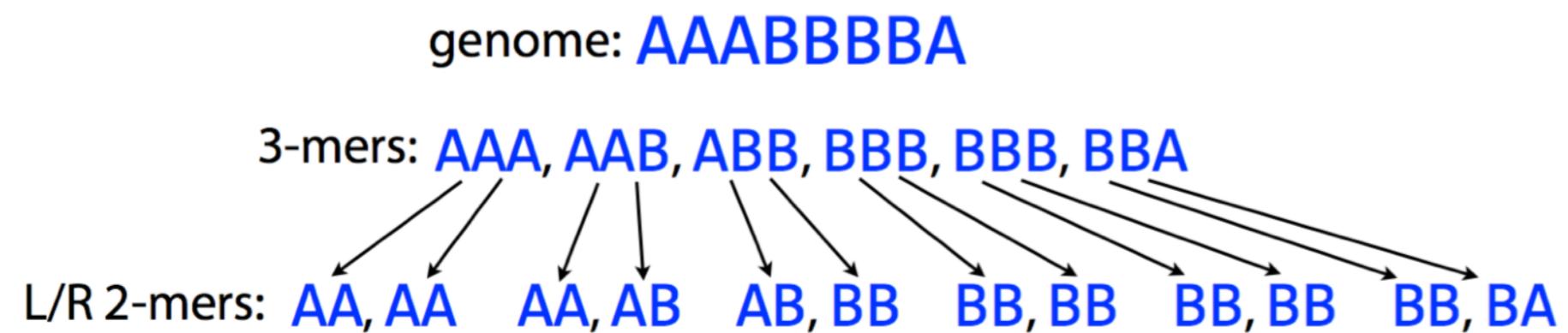


- graph of text fragments, e.g. words
- edge represents order of words in text
- Multigraph: multiple edges between nodes

De-Bruijn Graph

- Method to represent string as graph
- edge = k-mer
- node = k-1 mer

Assembler paradigms



Assembler paradigms

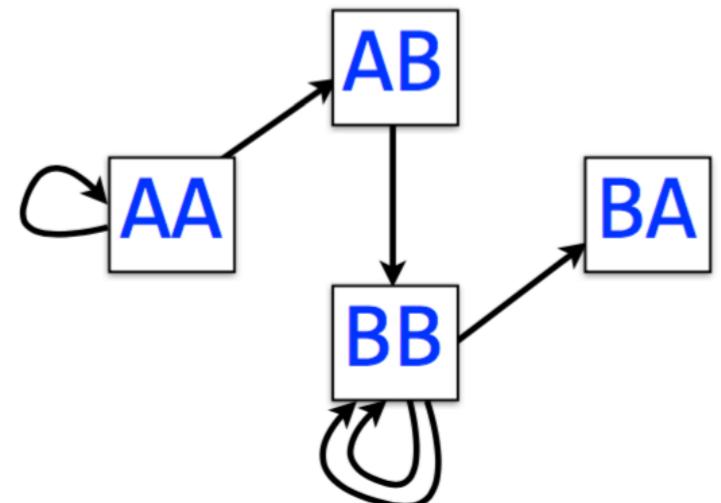
- Node is balanced if indegree equals outdegree
- Node is semi-balanced if indegree differs from outdegree by 1

Graph is connected if each node can be reached by some other node

- Eulerian walk visits each edge exactly once

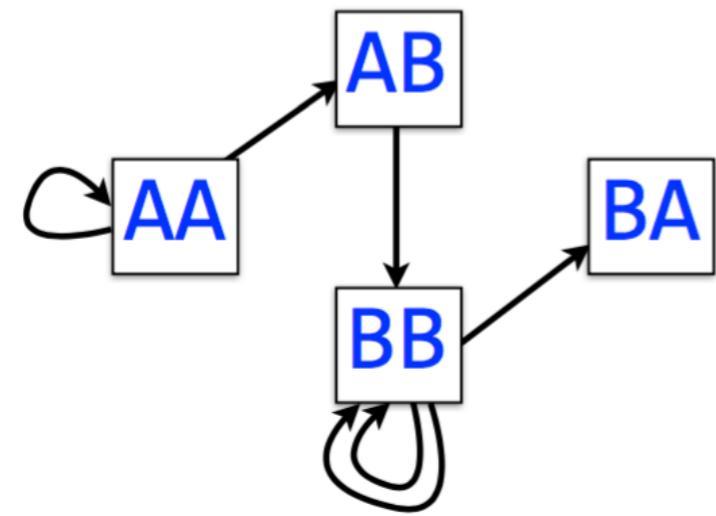
Not all graphs have Eulerian walks. Graphs that do are Eulerian.

- A directed, connected graph is Eulerian if and only if it has at most 2 semi-balanced nodes and all other nodes are balanced



De-Bruijn Graph

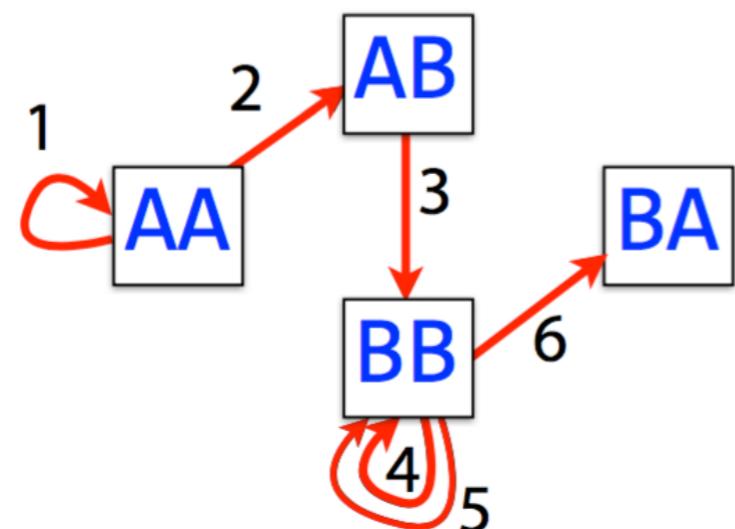
Assembler paradigmes



- De-Bruijn graphs have Eulerian path
- walking along each edge once (in right direction) yields the original sequence

De-Bruijn Graph

Assembler paradigms

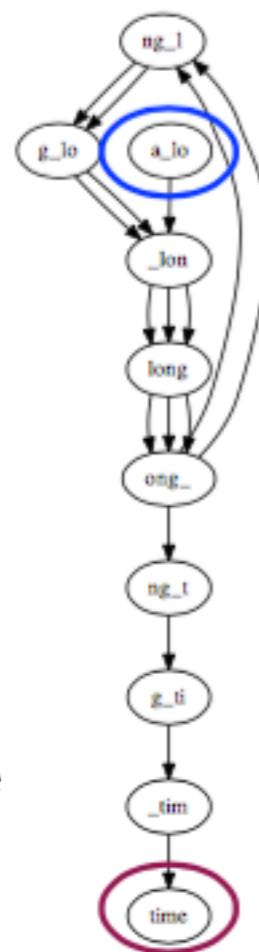
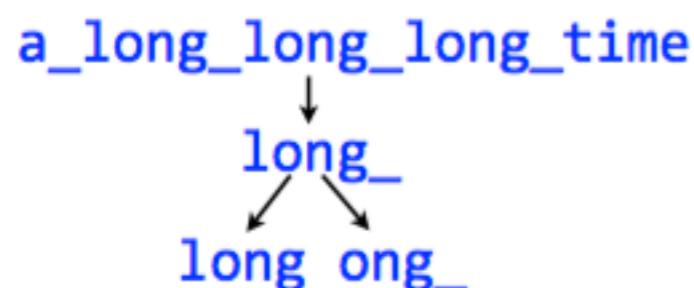


AAABBBBA

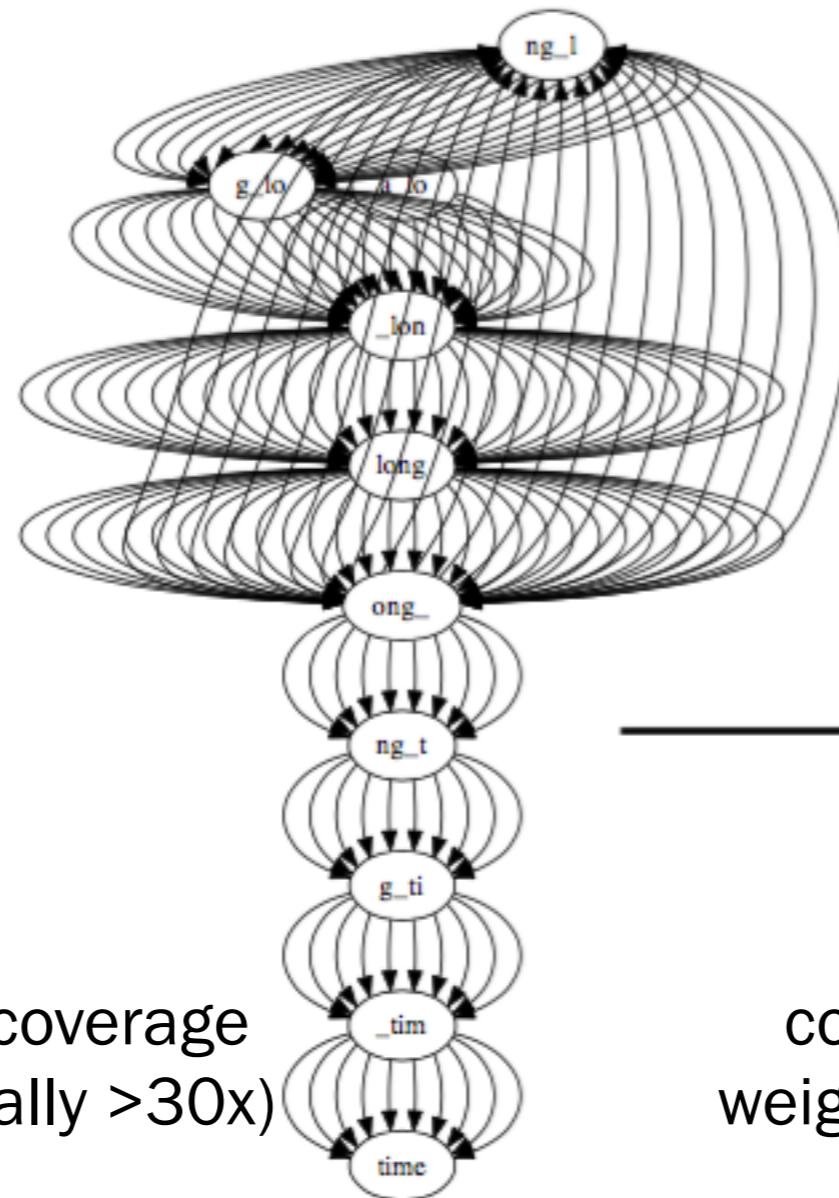
- De-Bruijn graphs have Eulerian path
- walking along each edge once (in right direction) yields the original sequence
- not possible with all graphs

De-Bruijn Graph

Assembler paradigms

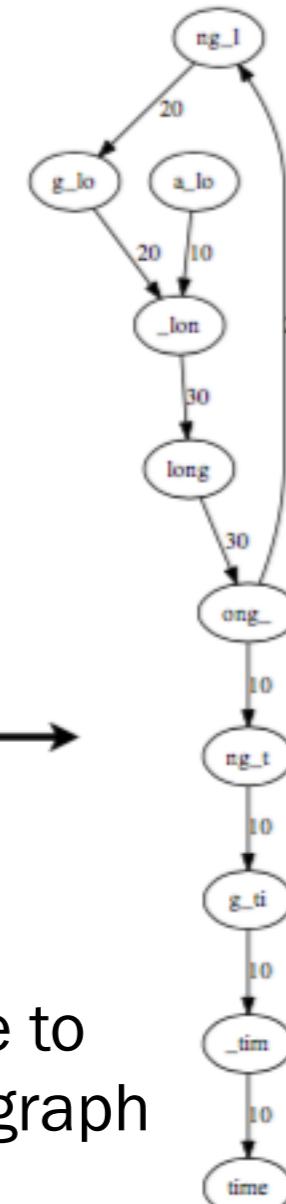


1x coverage

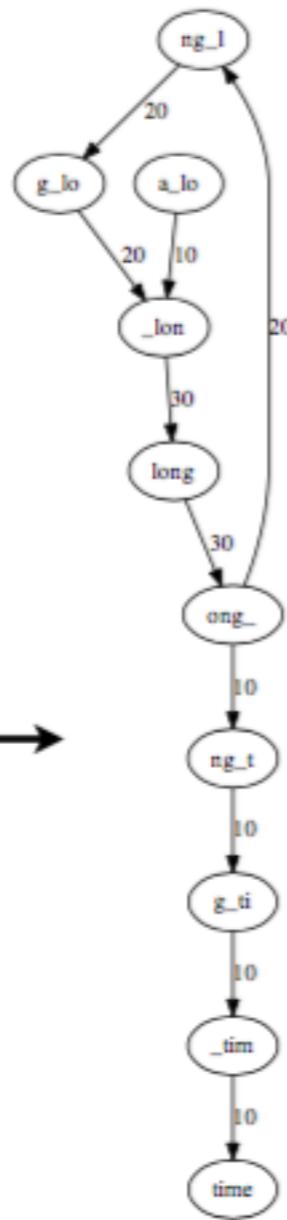
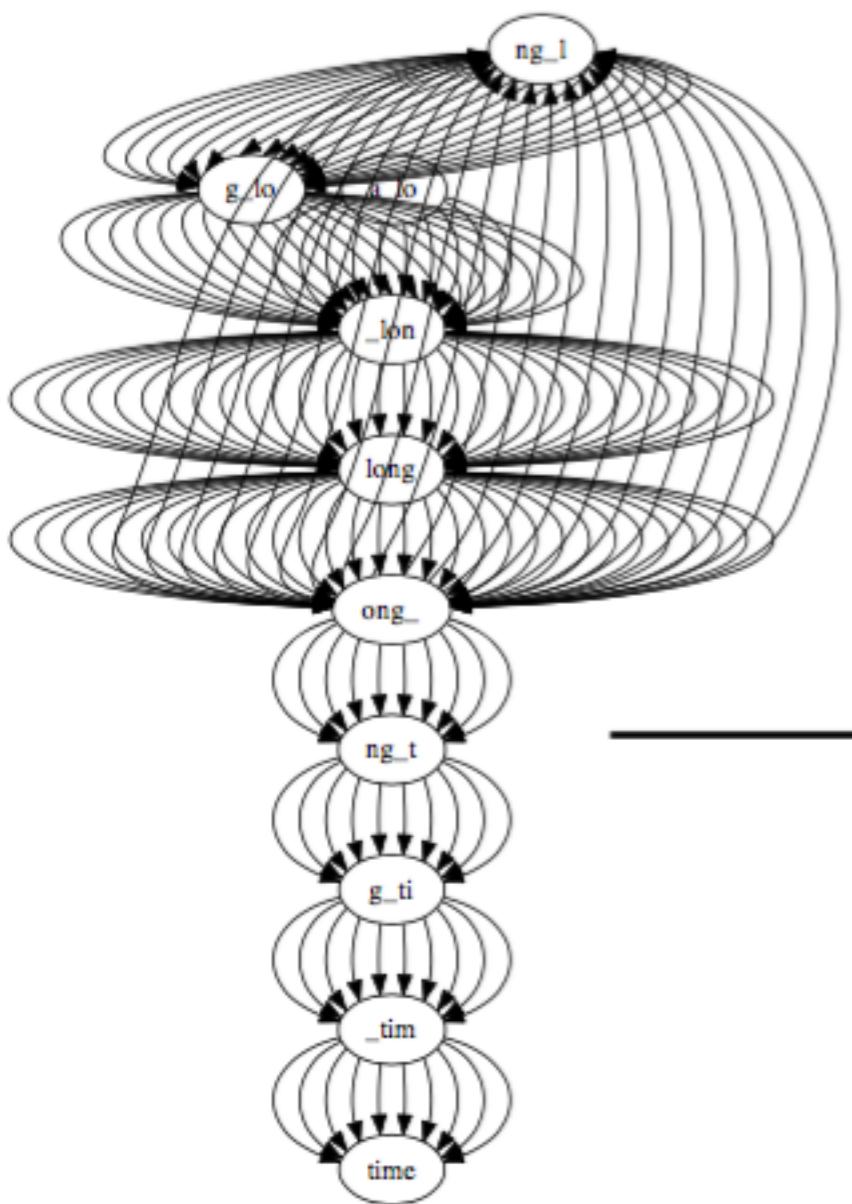


10x coverage
(typically >30x)

collapse to
weighted graph

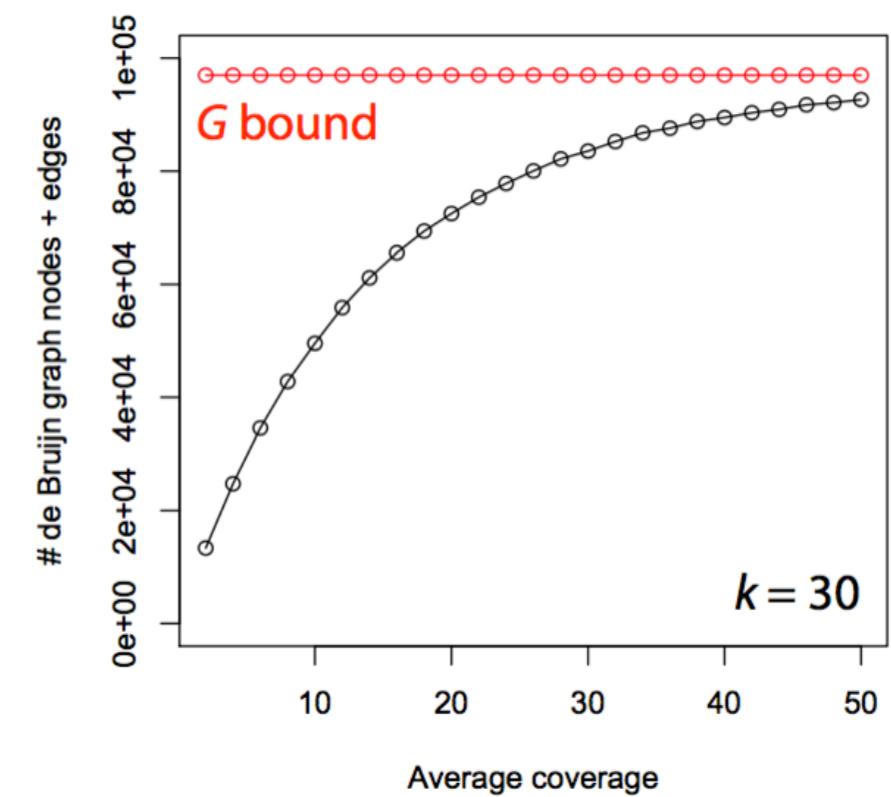


De-Brujin Graph



Assembler paradigms

Number of k-mer's depends on genome, not number of reads -> good for high coverage datasets (advantage over OLC)



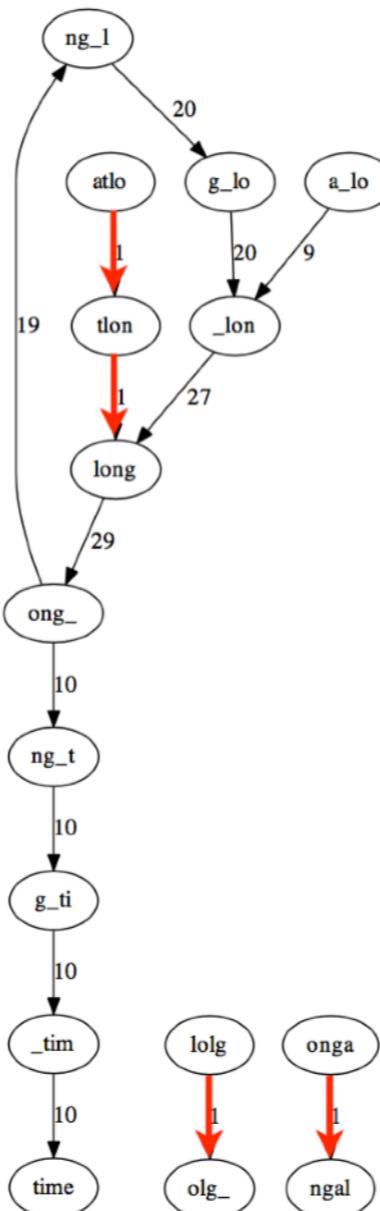
De-Brujin Graph

What about Errors?

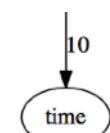
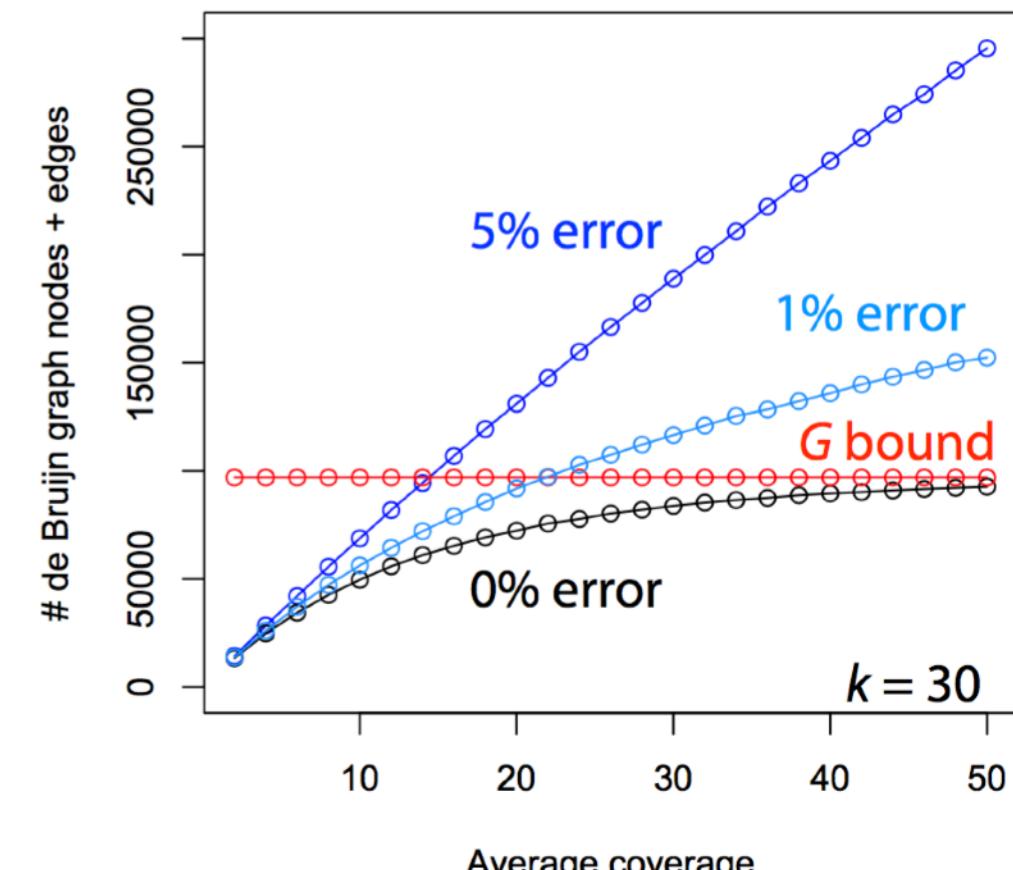
a_long_long_long_time
atlon
tlong
lolg_
ongal

long_

long ong_



Assembler paradigms



De-Brujin Graph

How many sequencing errors?

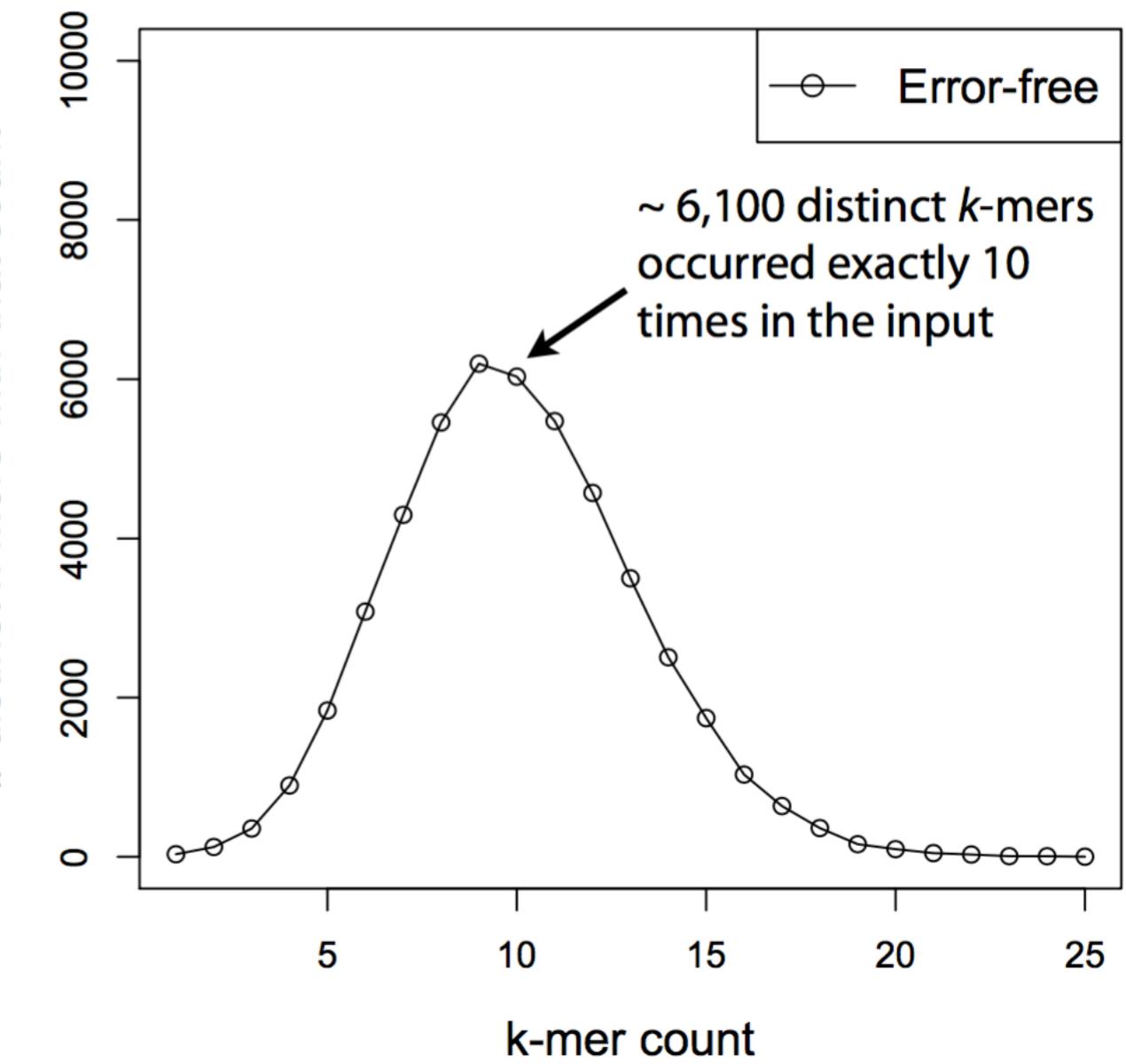
Read:

GCGTATTACGCGTCTGGCCT
GCGTATTA: 8
CGTATTAC: 8
GTATTACG: 9
TATTACGC: 9
ATTACGCG: 9
TTACGCGT: 12
TACGCGTC: 9
ACGCGTCT: 8
CGCGTCTG: 10
GCGTCTGG: 10
CGTCTGGC: 11
GTCTGGCC: 9
TCTGGCCT: 8

8-mers:

distinct k-mers with that count

Assembler paradigms

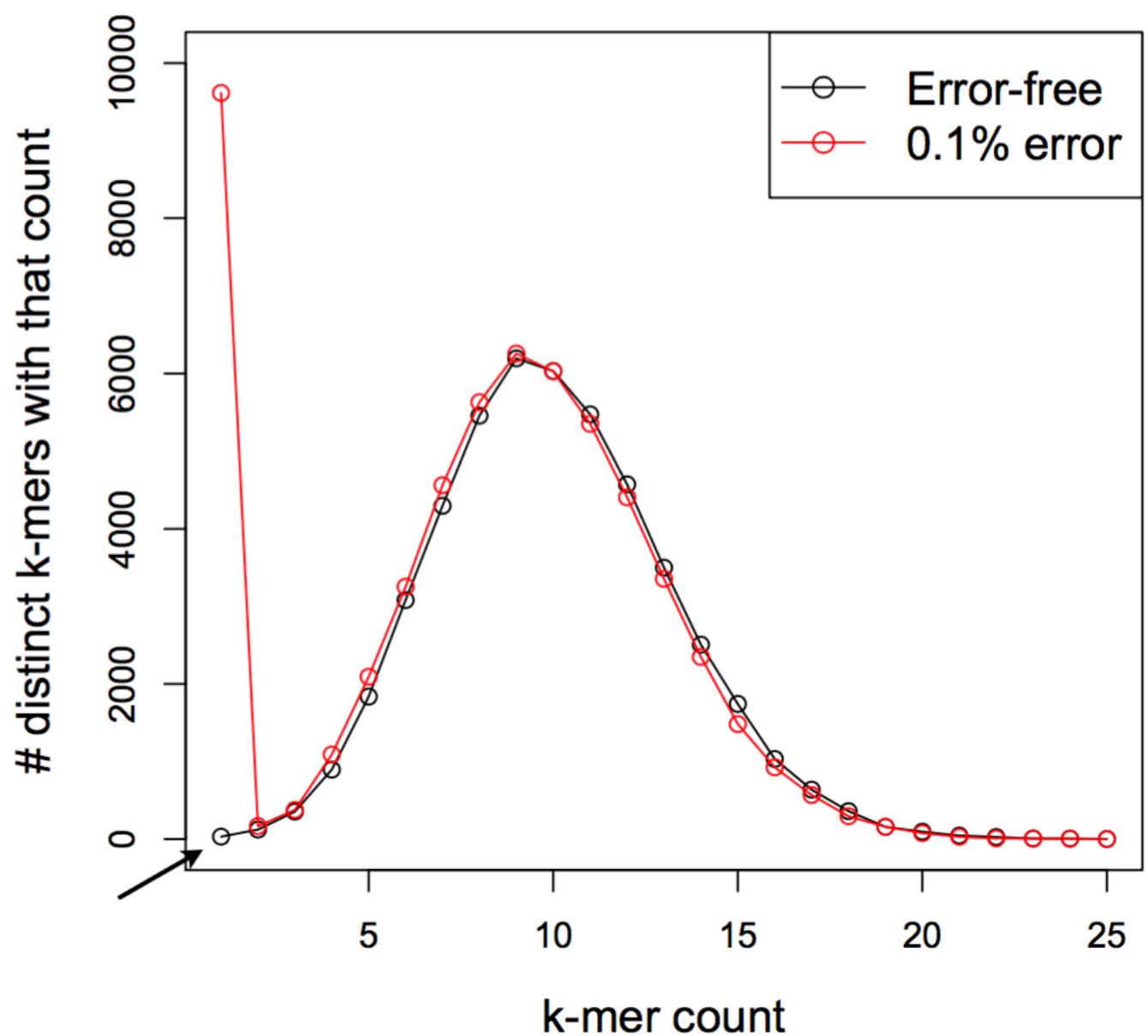


De-Brujin Graph

How many sequencing errors?

GCGTATTAC~~A~~CGTCTGGCCT
GCGTATTA: 8
CGTATTAC: 8
GTATTACA: 1
TATTACAC: 1
ATTACACG: 1
TTACACGT: 1
TAC~~A~~CGTC: 1
AC~~A~~CGTCT: 2
~~C~~ACGTCTG: 1
GCGTCTGG: 10
CGTCTGGC: 11
GTCTGGCC: 9
TCTGGCCT: 8

Assembler paradigms

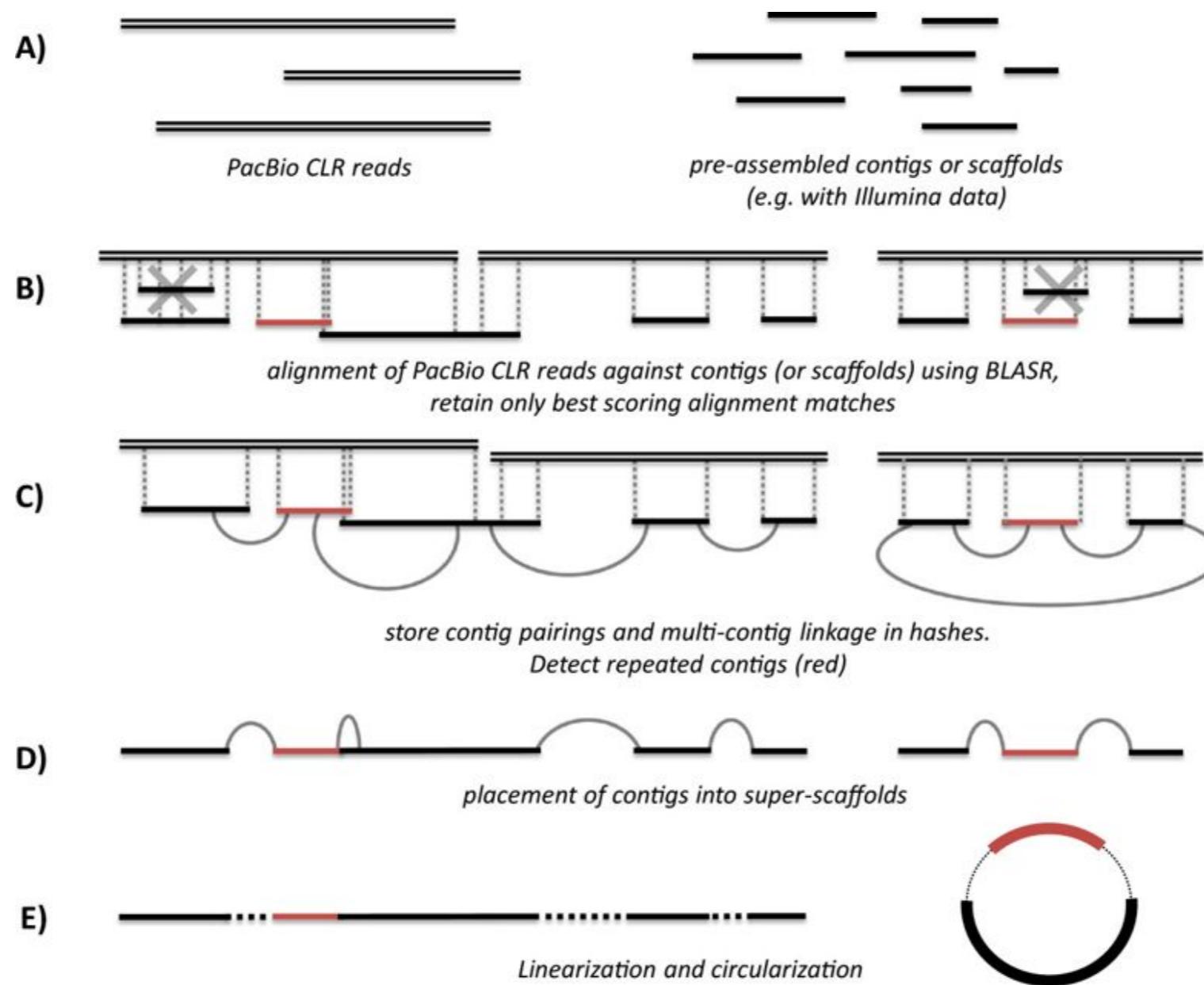


An assembly pipeline - SPAdes

De-Bruijn graph based

1. Read error correction based on k-mer frequencies with BayesHammer
2. De Bruijn graph assembly at multiple k-mer sizes
3. Merging different k-mer assemblies (good for varying coverage)
4. combine Contigs into scaffolds using paired end/mate pair reads
5. Repeat resolution from paired end/mate pair data using rectangle graphs
6. Contig error correction based on aligning the original reads with BWA back to contigs

Scaffolding



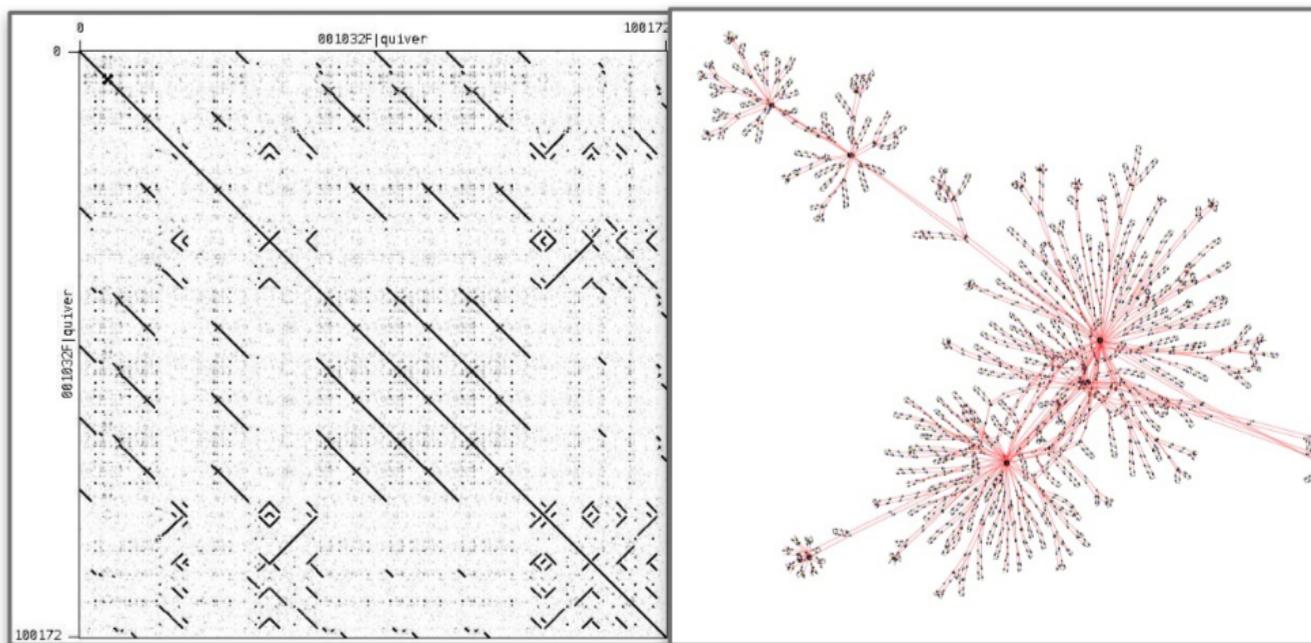
Open Problems: Repeats

- not enough known about true repeat structures across genomes

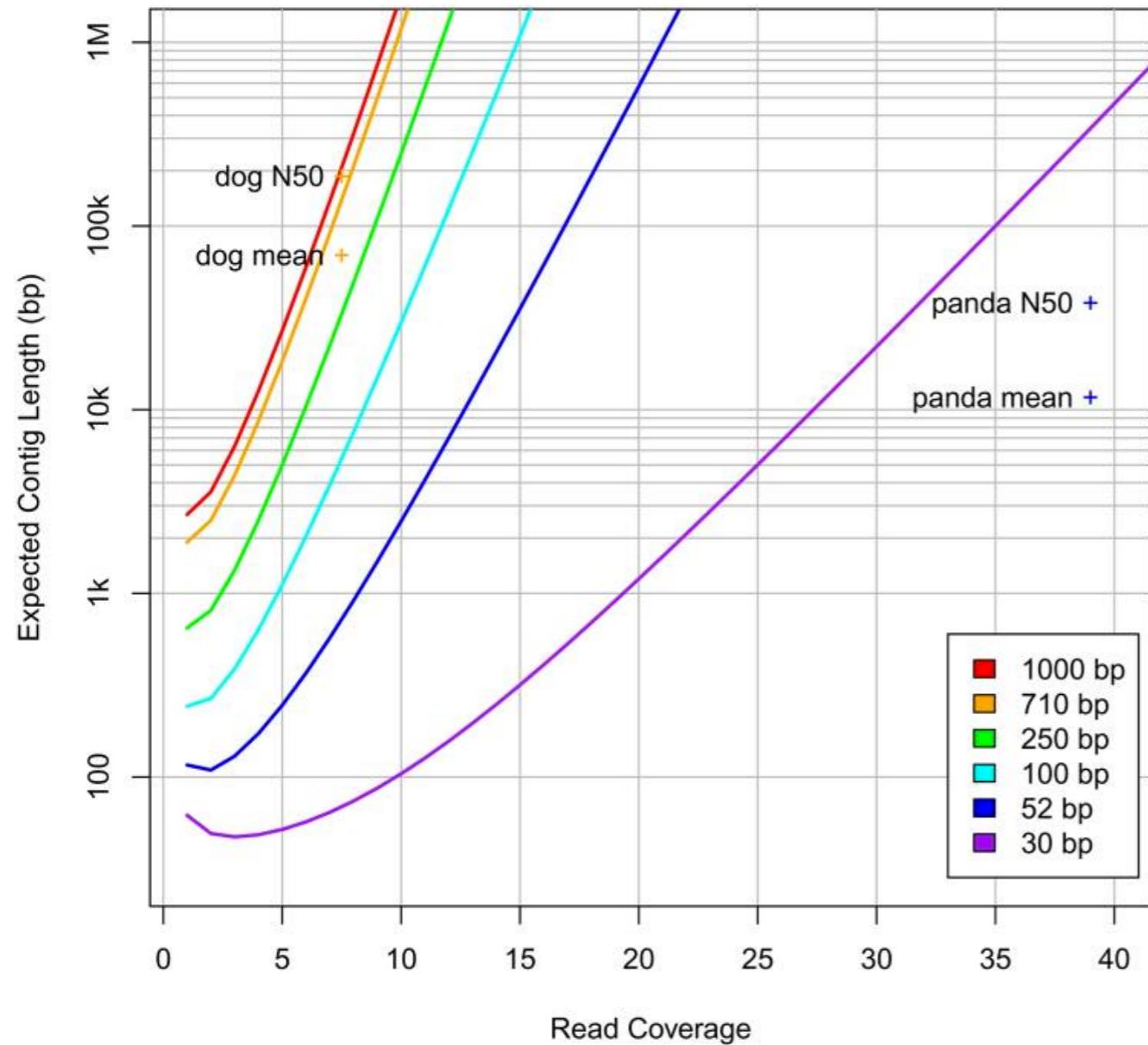
2019 - long read assembly - same problem

1984 Peltola, Söderlund, and Ukkonen.
Nucleic Acids Res.

Repetitive structures within the DNA sequence can confuse any automatic assembly process which uses only the gel readings. Without computer assistance we have earlier determined the sequence of a cDNA (20) which turns out to be an interesting benchmark for automatically assembling. This cDNA contains three copies of two separate perfectly repeated segments both more than hundred bases in length



Solution: longer Reads!

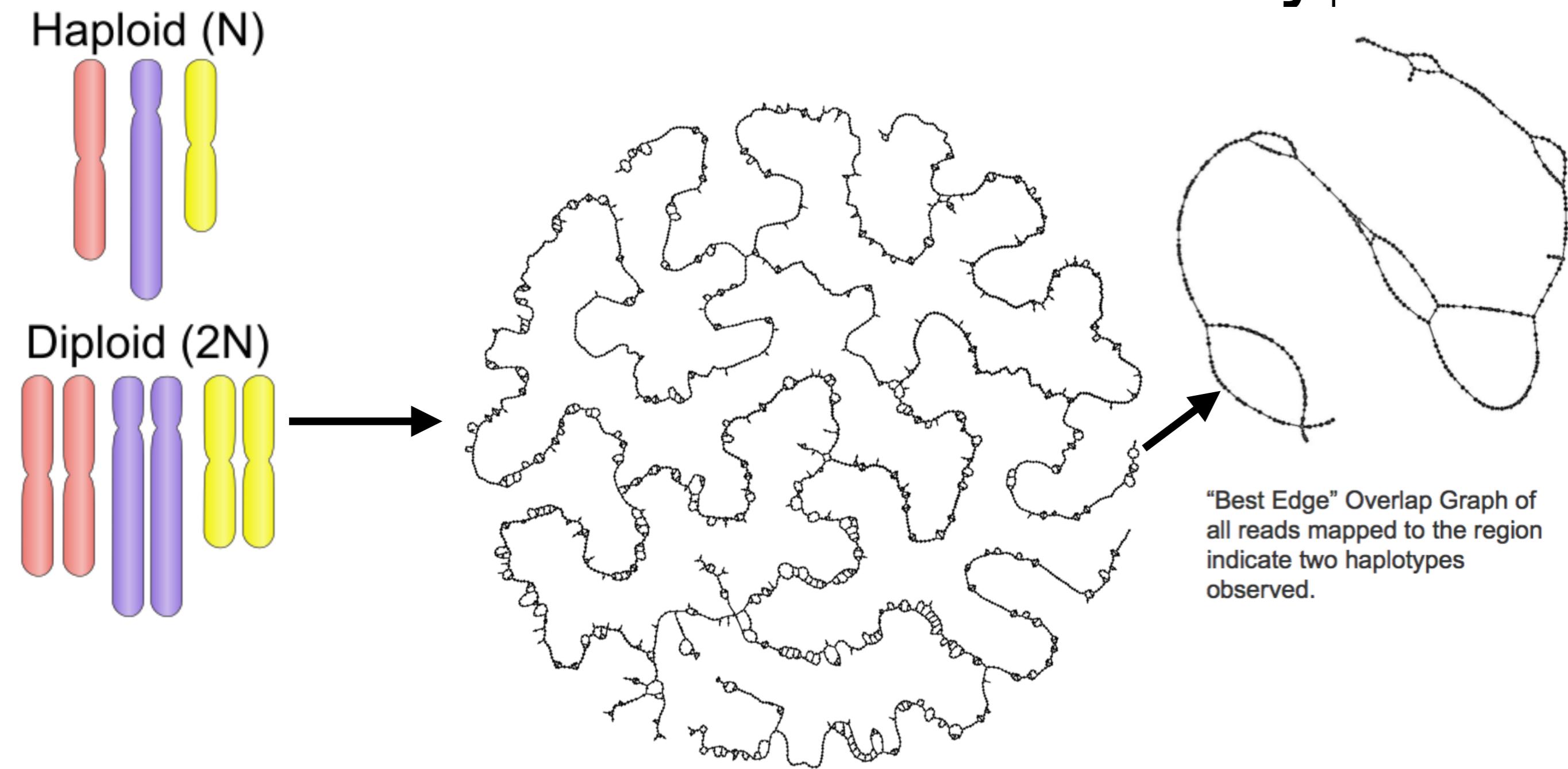


Open Problems: Repeats

Fig. 3: Expected average of contig length for a range of different read lengths and coverage values. Also shown are the average contig lengths and N50 lengths for the dog genome, assembled with 710-bp reads, and the panda genome, assembled with reads averaging 52 bp in length.

- not enough known structural difference between haplotypes for species with diploid or polyploid genomes

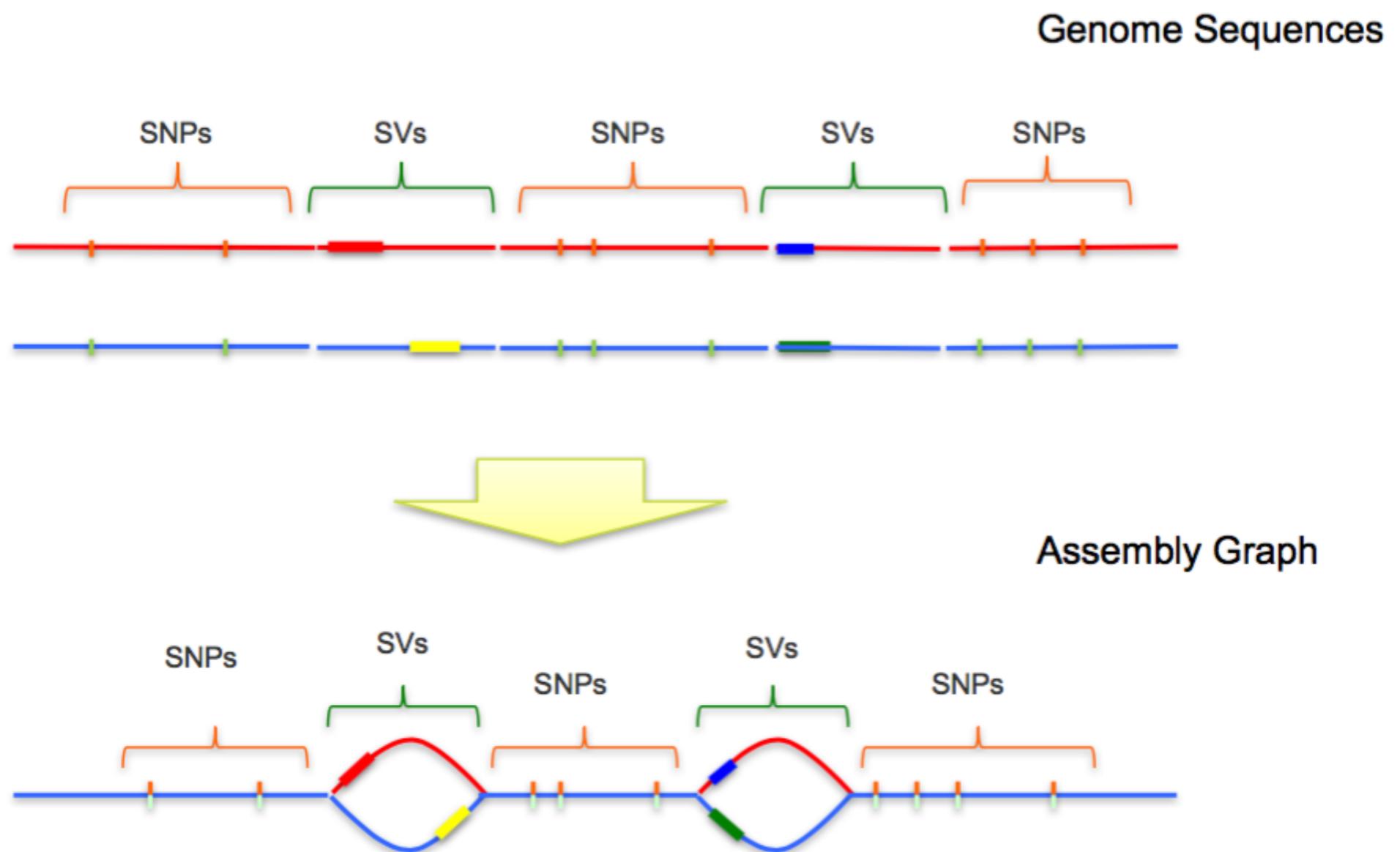
Open Problems: Ploidy

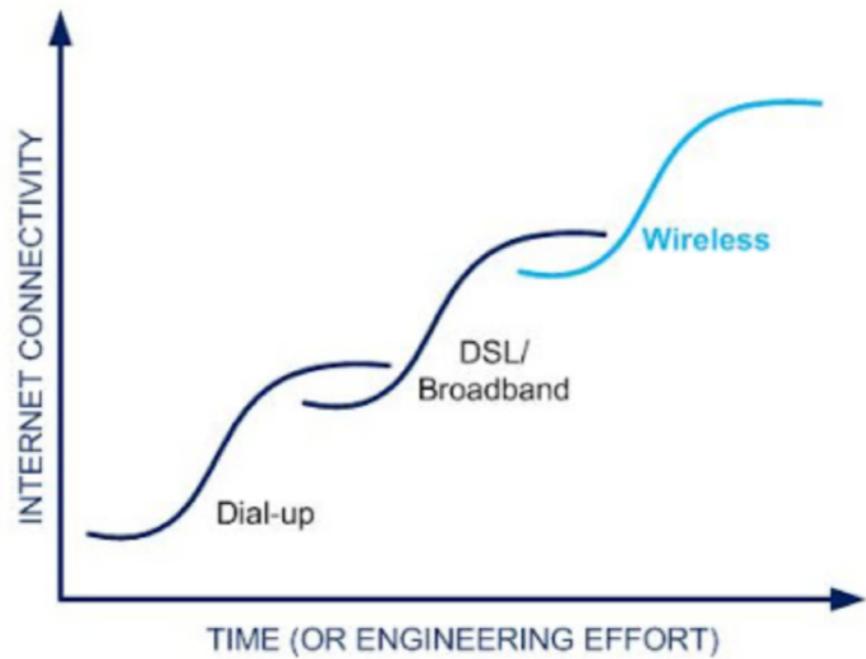


Open Problems: Ploidy

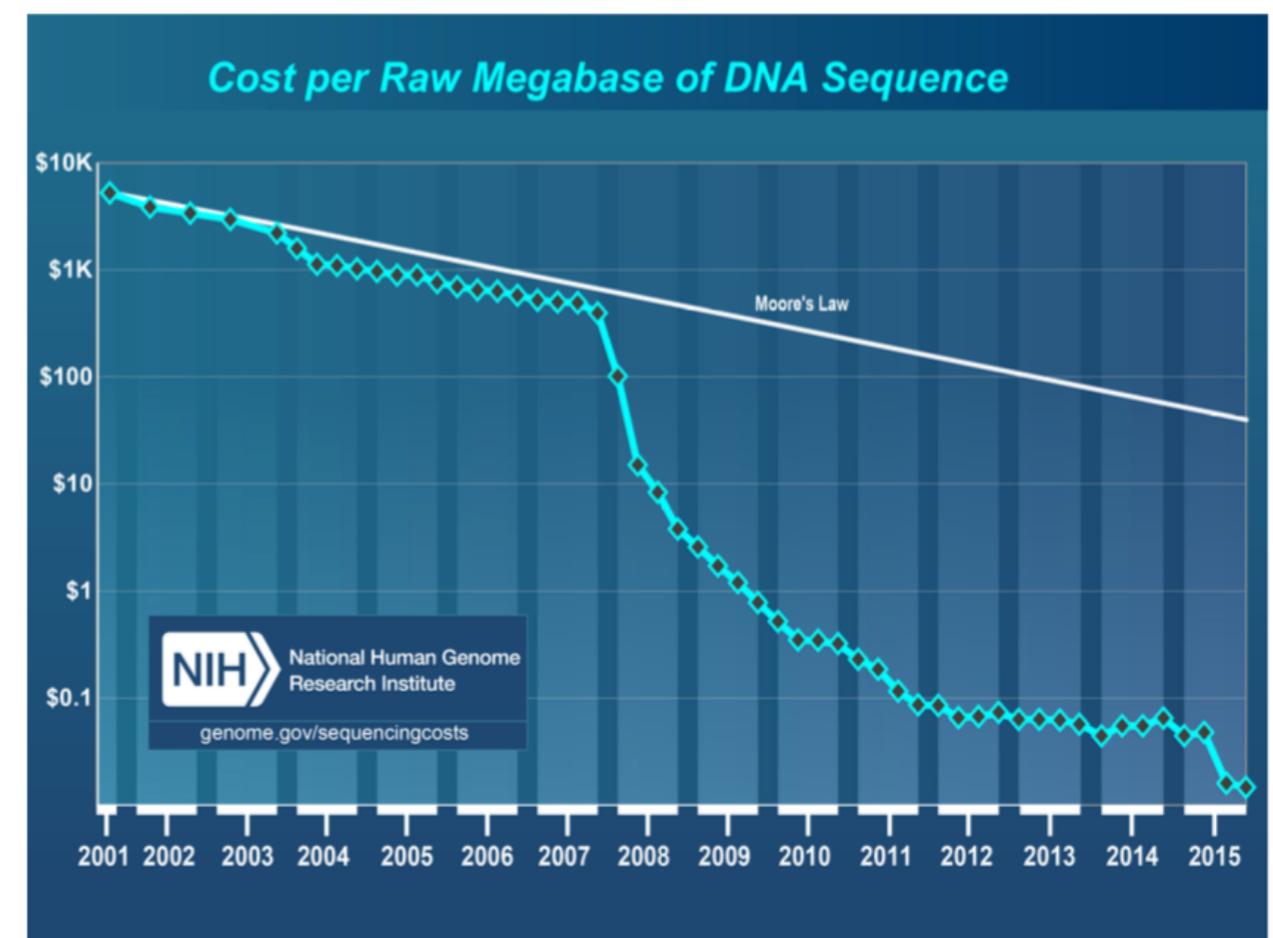
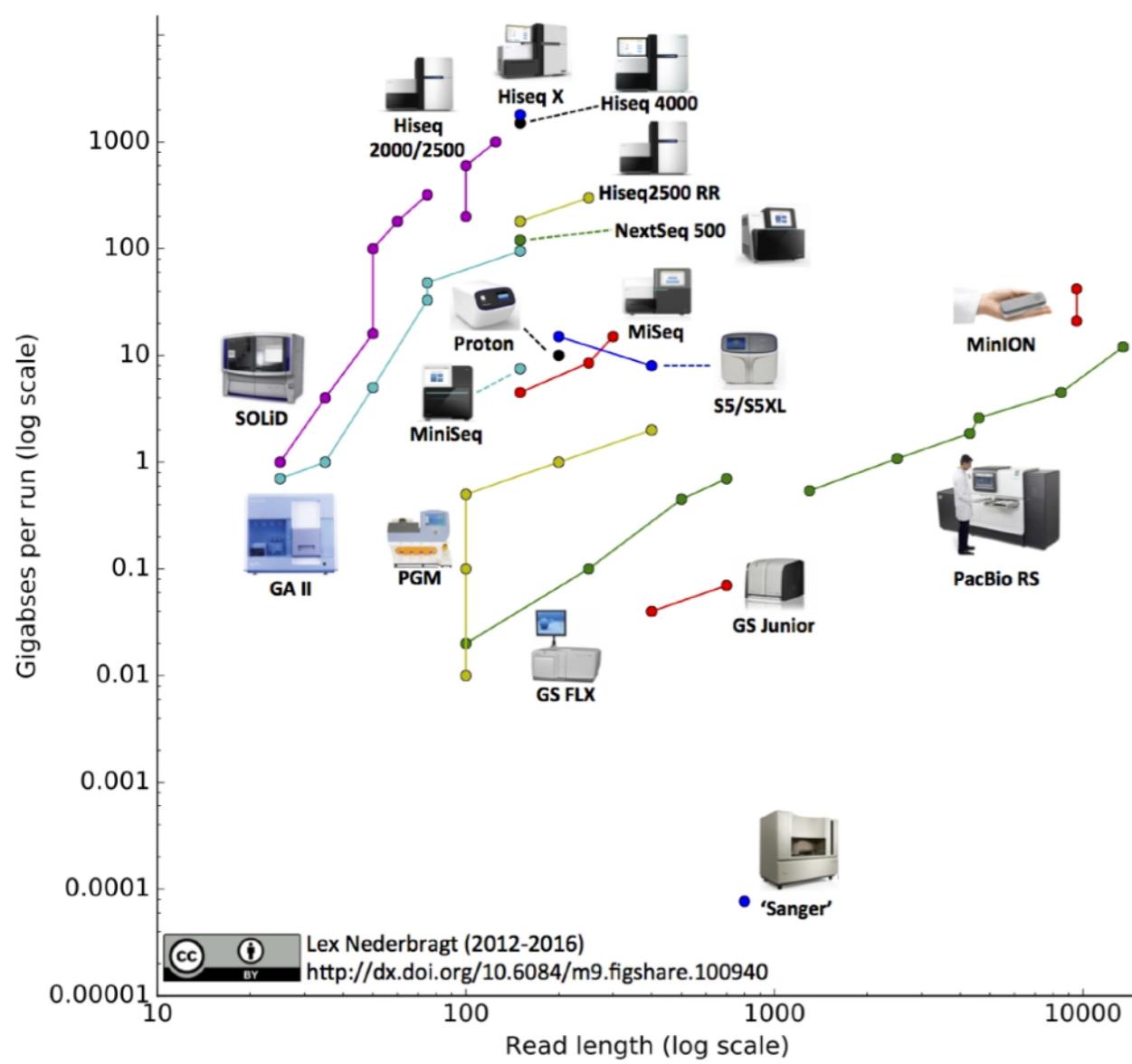
diploid-aware assembly algorithms, e.g.
Falcon_unzip

(Phased diploid genome assembly with single-molecule real-time sequencing. Nat Methods. 2016 Dec;13(12):1050-1054)





Reads are getting more and longer



Questions?