

Principles and problems of de novo genome assembly

Karin Lagesen

Norwegian Veterinary Institute



Material adapted from slides
provided by Lex Nederbragt



What is this thing called ‘genome assembly’?

What is a genome assembly?

A hierarchical data structure

that maps the sequence data

to a putative reconstruction of the target

Hierarchical structure

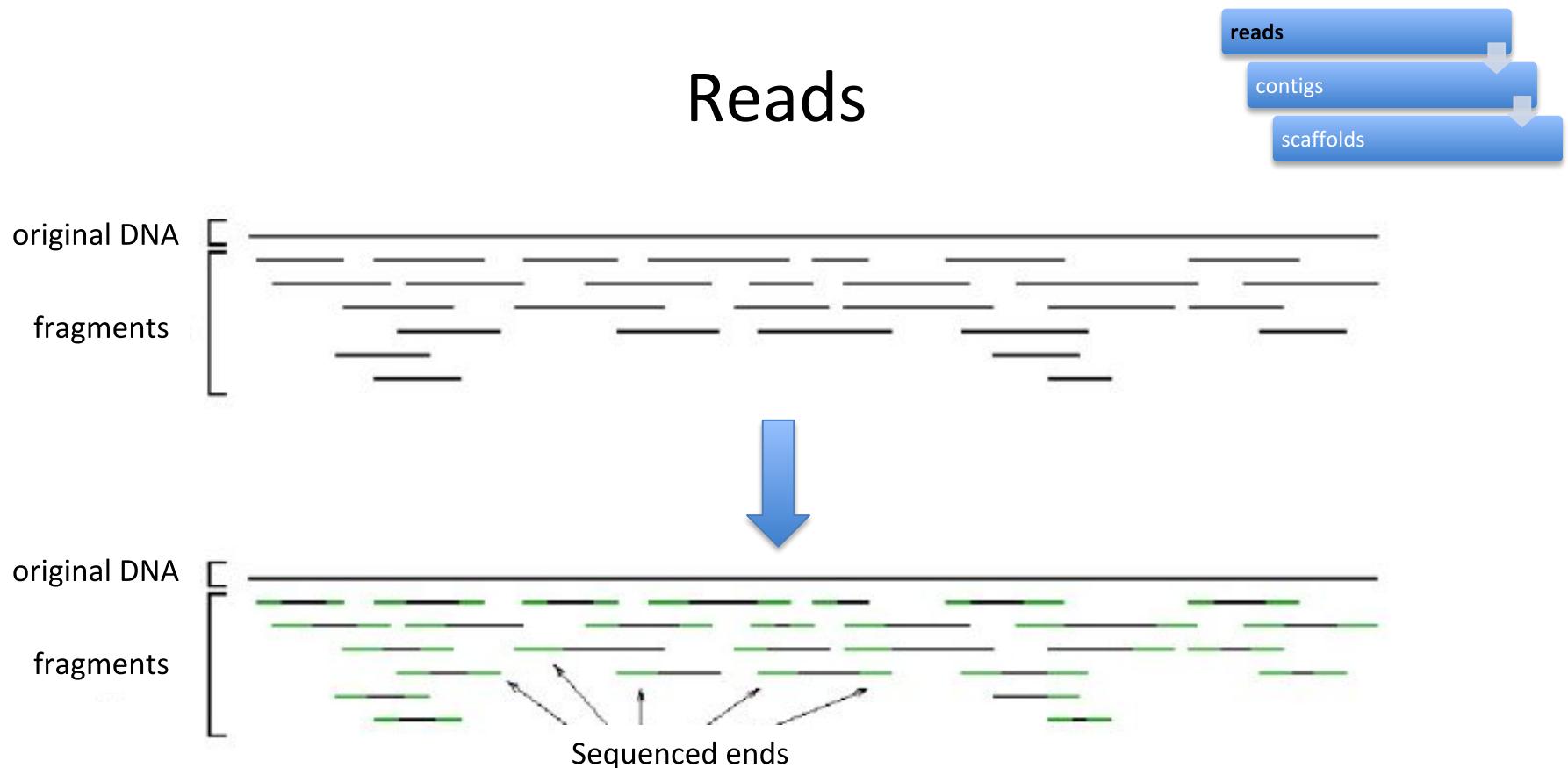
reads

contigs

scaffolds



Sequence data



Contigs

Building contigs



Aligned reads

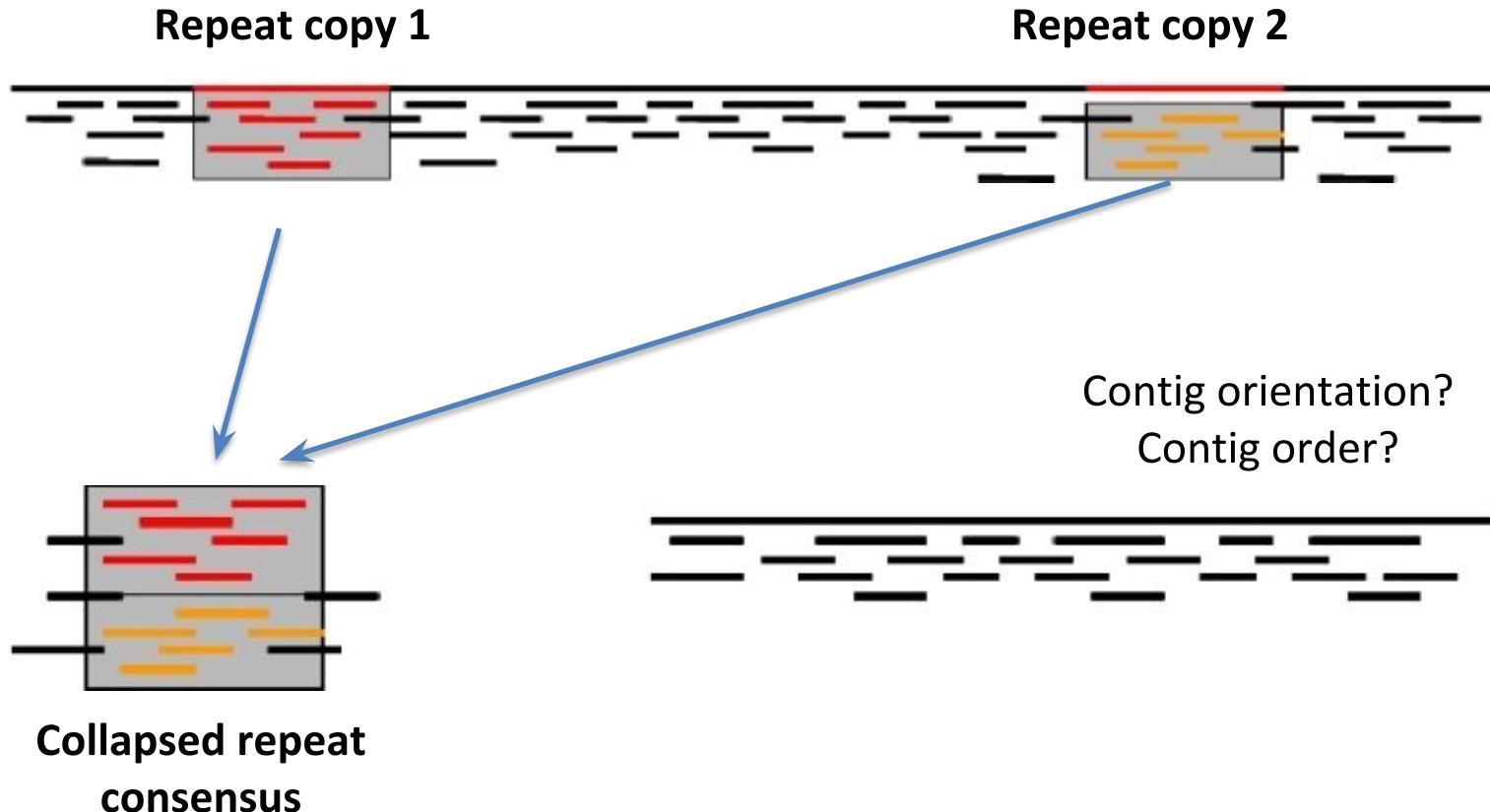
```
ACGCGATTCAAGGTACCACG  
GCGATTCAAGGTACCACGCG  
GATTCAAGGTACCACGCGTA  
TTCAGGTACCAACGCGTAGC  
CAGGTACCAACGCGTAGCGC  
GGTTACCAACGCGTAGCGCAT  
TTACCAACGCGTAGCGCATTA  
ACCACGCGTAGCGCATTACA  
CACGCGTAGCGCATTACAC  
CGCGTAGCGCATTACACAGA  
CGTAGCGCATTACACAGATT  
TAGCGCATTACACAGATTAG
```

Consensus contig

```
ACGCGATTCAAGGTACCACGCGTAGCGCATTACACAGATTAG
```

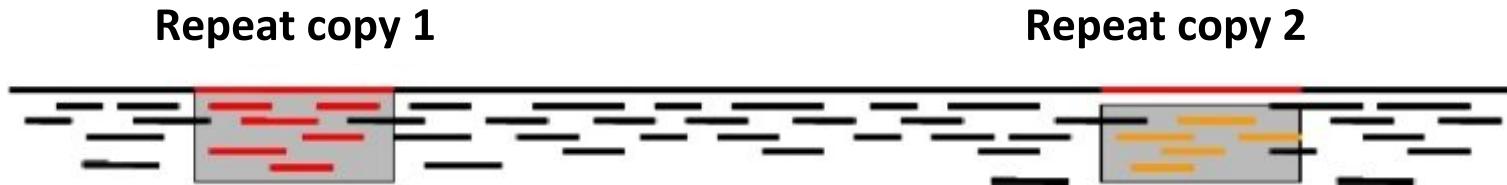
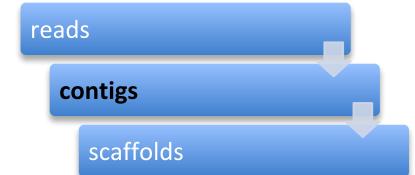
Contigs

Building contigs



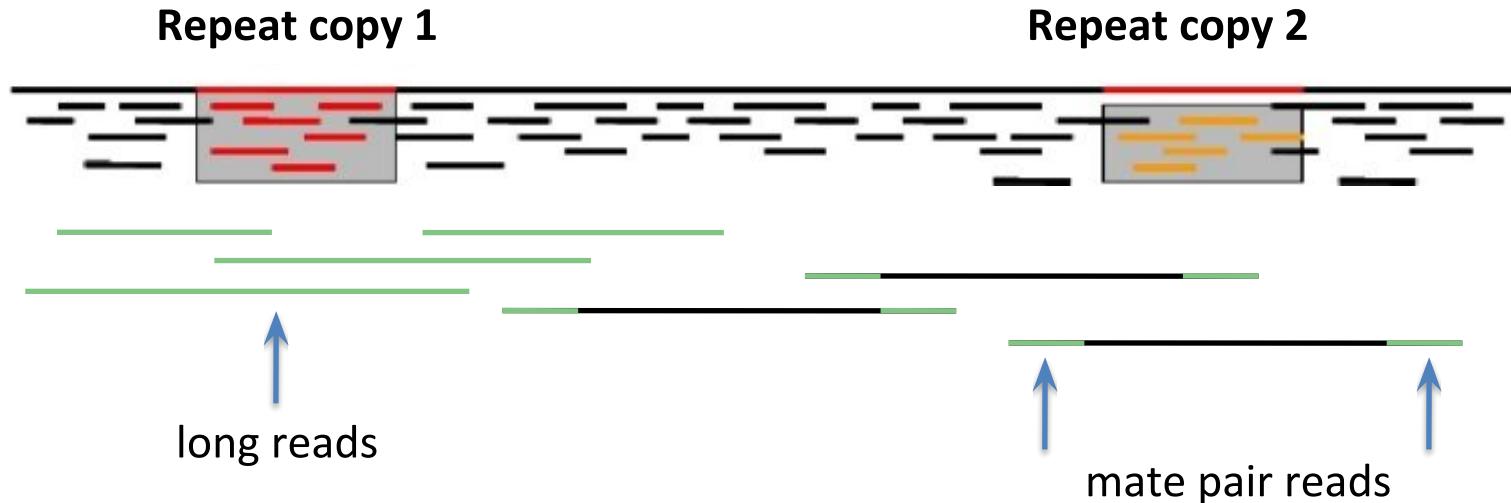
Contigs

Repeats: major problem



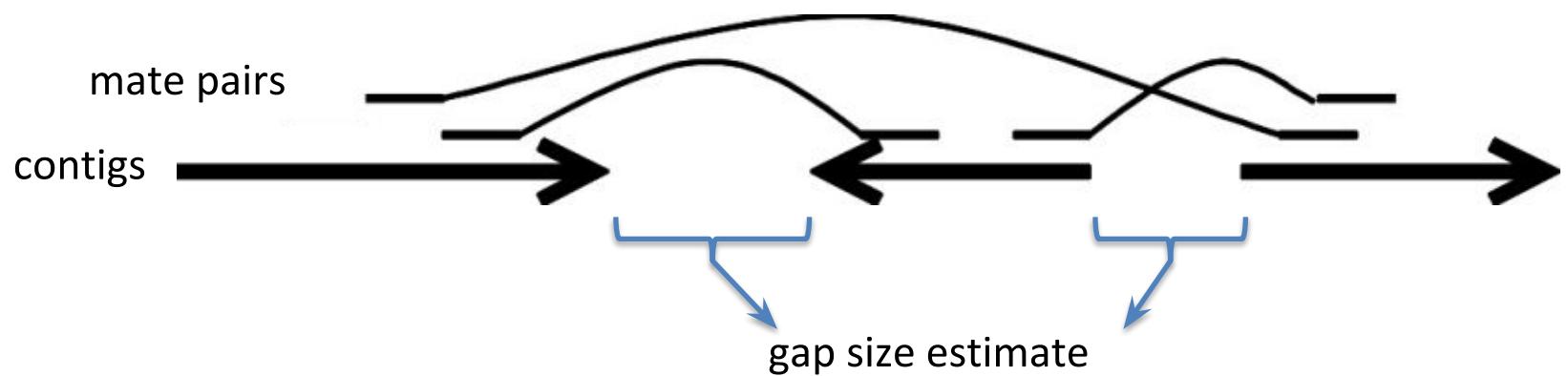
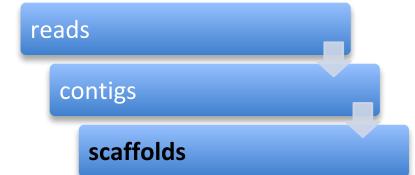
Mate pairs/long reads

Other read type



Scaffolds

Ordered, oriented contigs



Hierarchical structure



Assembly

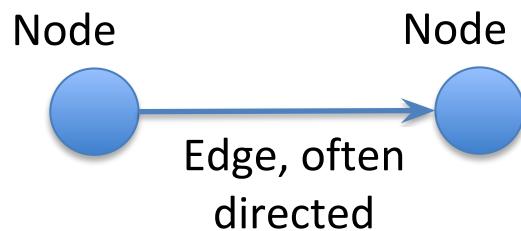
How to do this?



Algorithms

Algorithms

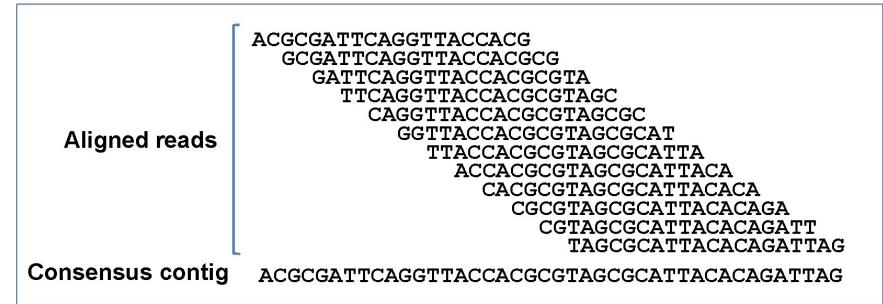
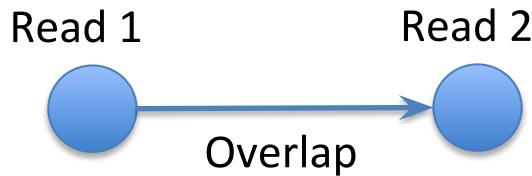
All are graph-based



Graph-theory!

Algorithms

All are graph-based

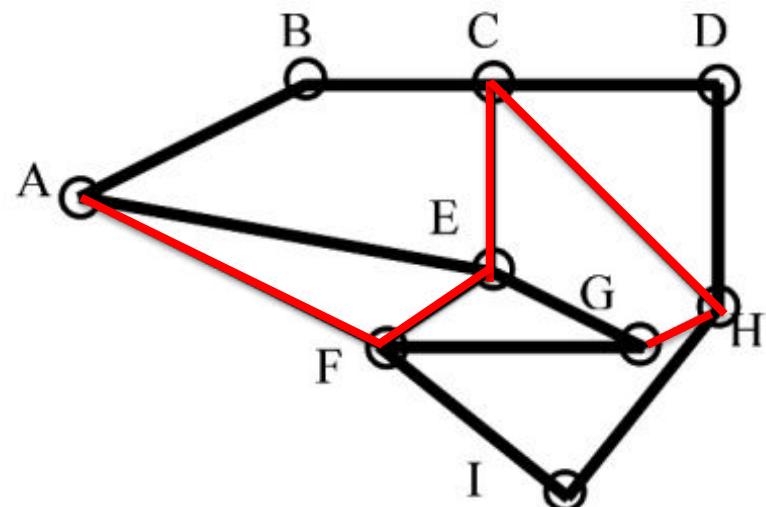
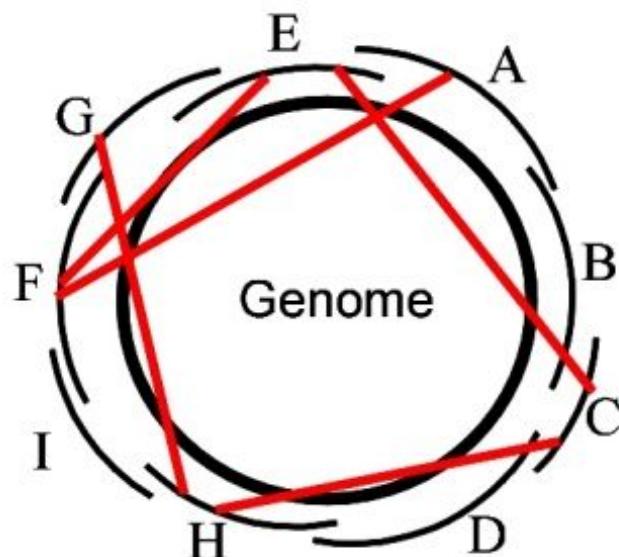


Graph-theory!

Algorithms

Hamiltonian path

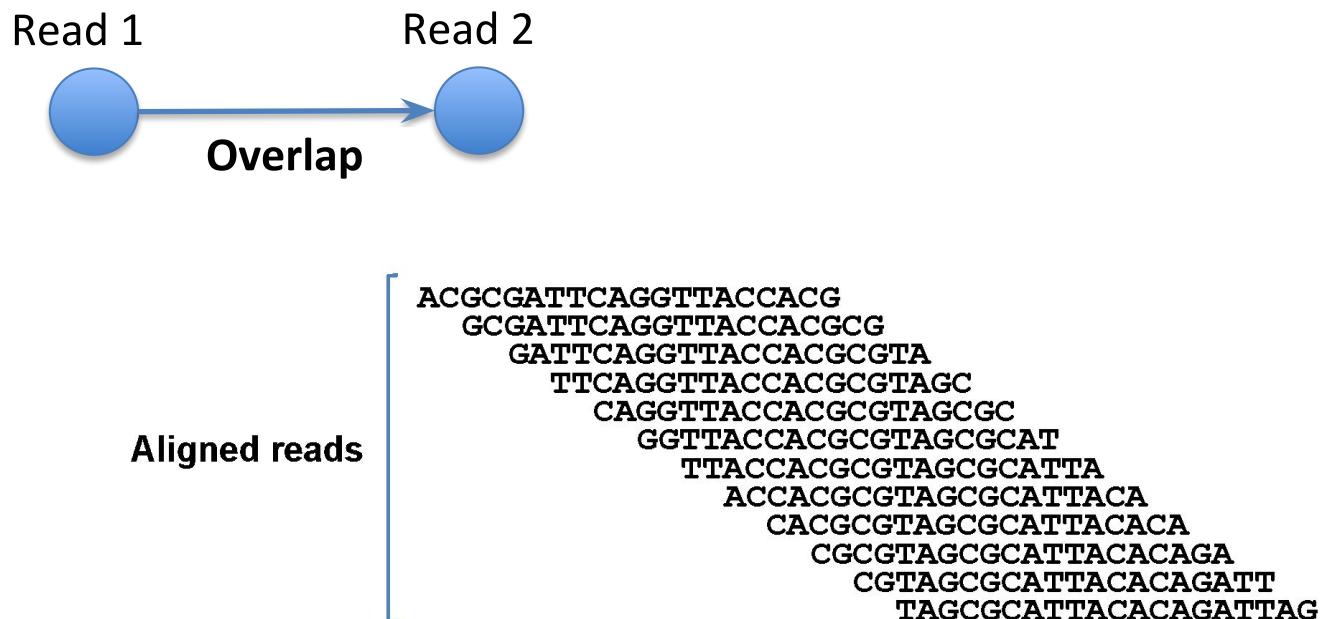
– a path that contains all the nodes



Algorithms

Overlap calculation (alignment)

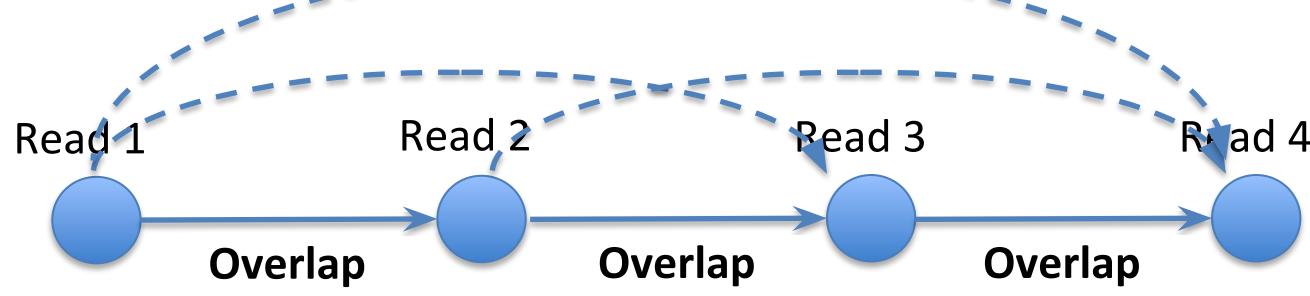
- computationally intensive



Algorithms

Path through the graph

□ contig



Aligned reads

ACGCGATTCAAGTTACCACG
GCGATTCAAGTTACCACGCG
GATTCAAGTTACCACCGCGTA
TTCAGGTACACCGCTAGC
CAGGTACACCGCTAGCGC
GGTTACACCGTAGCGCAT
TTACCACCGTAGCGCATTAA
ACCACCGTAGCGCATTACACA
CACCGTAGCGCATTACACAGA
CGCGTAGCGCATTACACAGATT
CGTAGCGCATTACACAGATT
TAGCGCATTACACAGATTAG

Consensus contig

ACGCGATTCAAGTTACCACCGTAGCGCATTACACAGATTAG

Algorithms

Many flavors



Abandoned

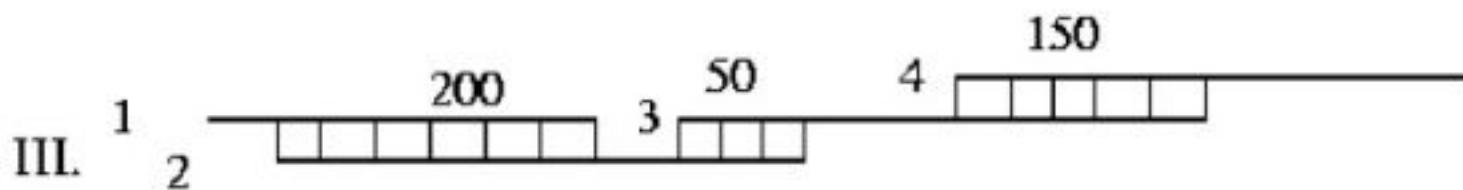
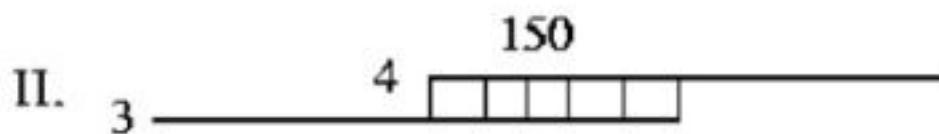
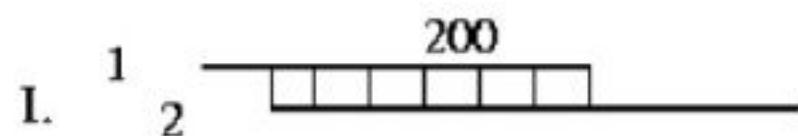
- Greedy extension

Two most used

- Overlap Layout Consensus
- de Bruijn graph

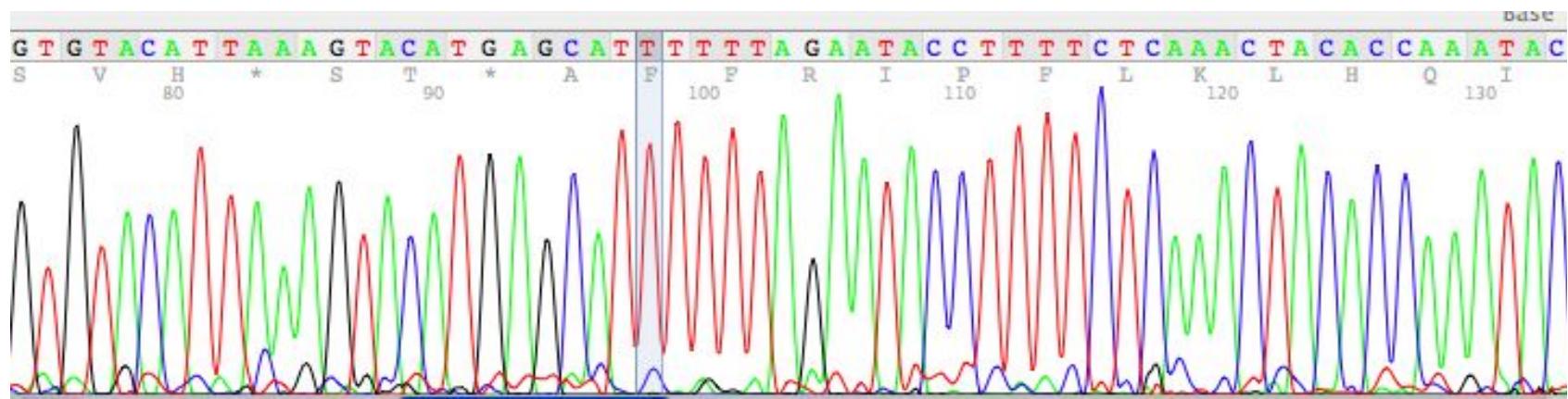
Greedy extension

Oldest



Overlap-Layout-Consensus

Developed for Sanger-type reads (longer reads)



Overlap-Layout-Consensus

Steps

- Overlap computation
- Layout: graph simplification
- Consensus: sequence

Overlap-Layout-Consensus

Overlap phase: find “similar enough” reads

Comparing all against all: expensive

Trick for finding “similar enough” reads:

- Split reads into k-mers

ACGCGATTCAAGGTTACCCACG

K-mer: substring of length k from a longer string

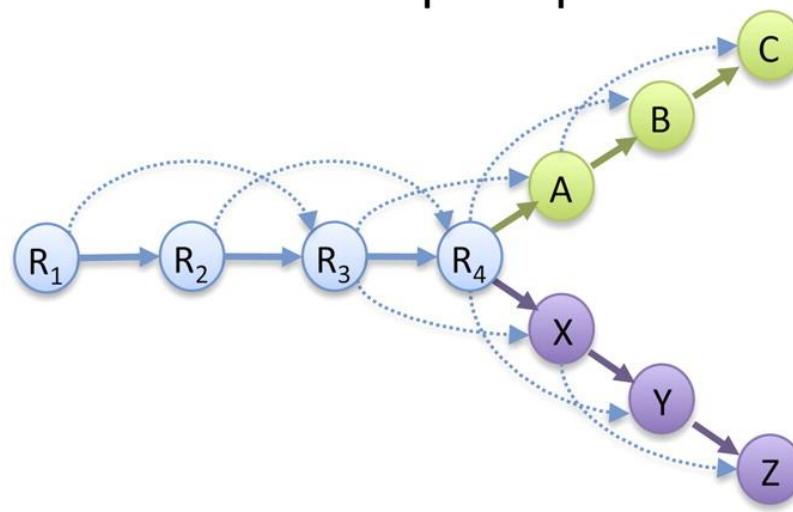
- Make list over which read has which k-mers
- If two reads share k-mers, test for similarity

Overlap-Layout-Consensus

A Read Layout

R ₁ :	GACCTACA
R ₂ :	ACCTACAA
R ₃ :	CCTACAAG
R ₄ :	CTACAAGT
A:	TACAAGTT
B:	ACAAGTTA
C:	CAAGTTAG
X:	TACAAGTC
Y:	ACAAGTCC
Z:	CAAGTCCG

B Overlap Graph

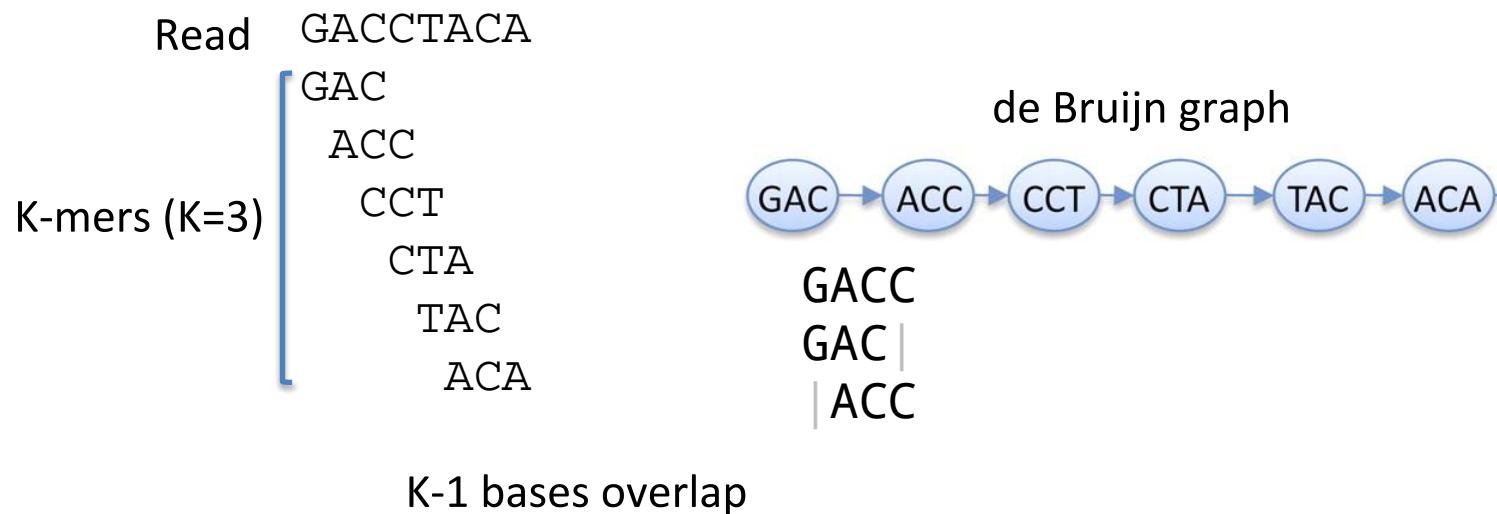




de Bruijn graphs

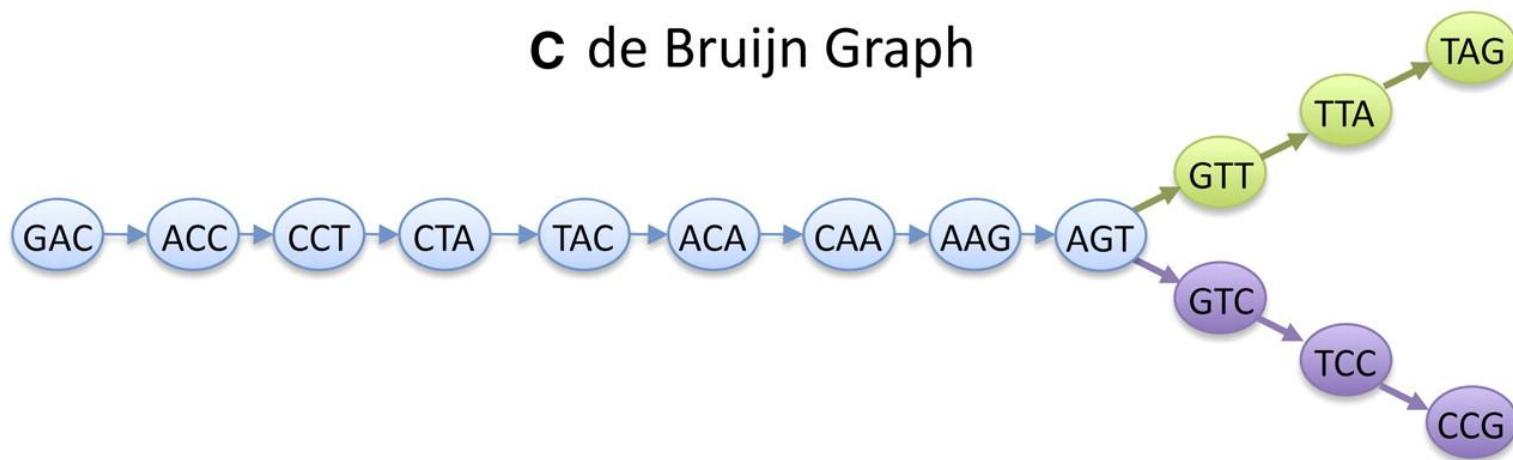
Developed outside of DNA-related work

- Best solution for short(er) reads



Graphs

C de Bruijn Graph

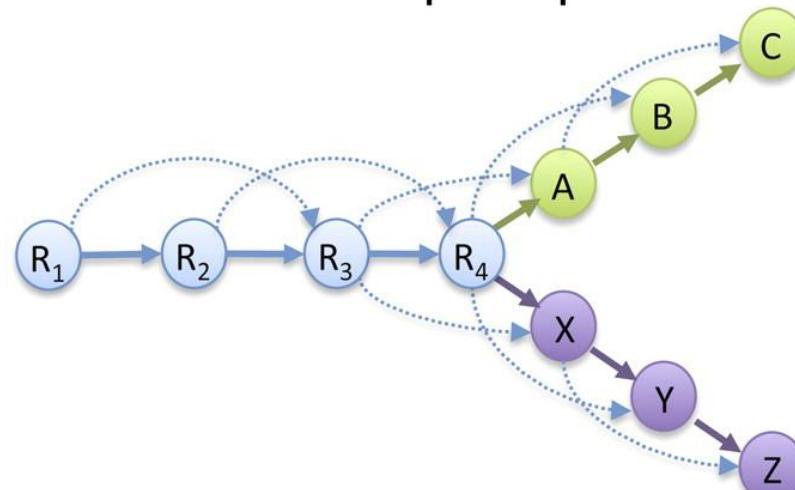


Graphs

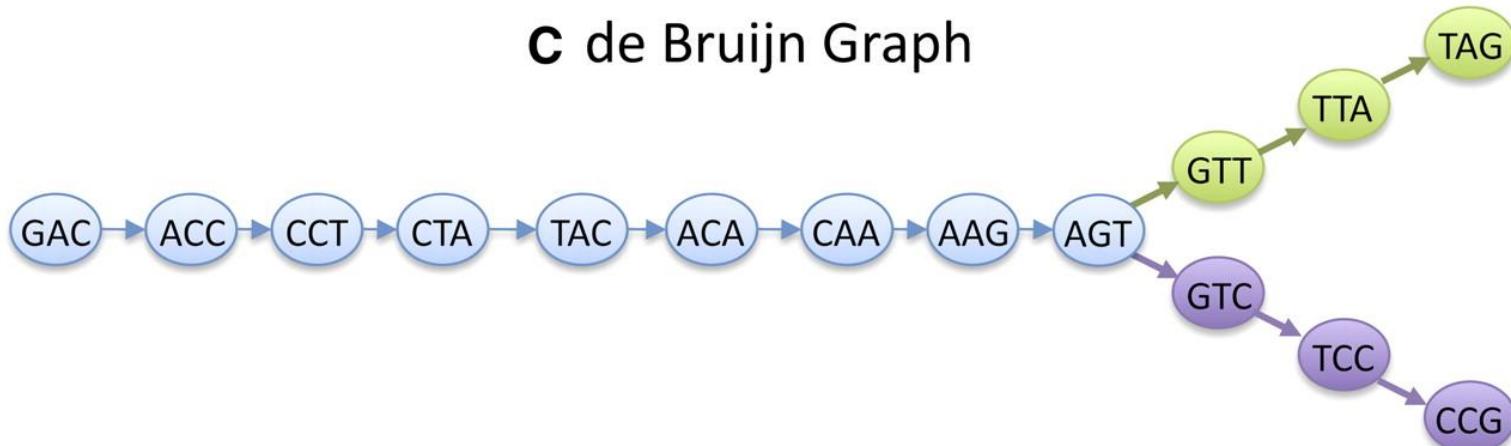
A Read Layout

R_1 :	GACCTACA
R_2 :	ACCTACAA
R_3 :	CCTACAAG
R_4 :	CTACAAGT
A:	TACAAGTT
B:	ACAAGTTA
C:	CAAGTTAG
X:	TACAAGTC
Y:	ACAAGTCC
Z:	CAAGTCCG

B Overlap Graph

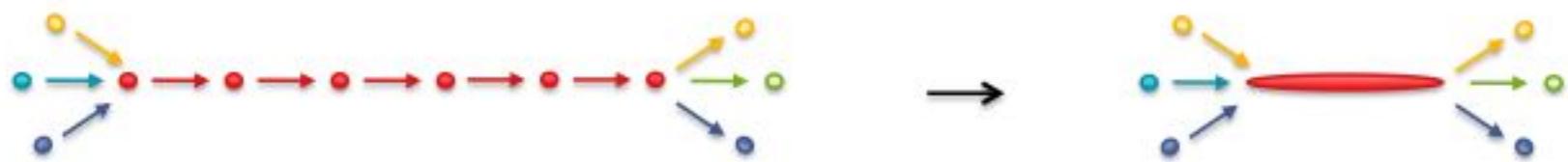


C de Bruijn Graph

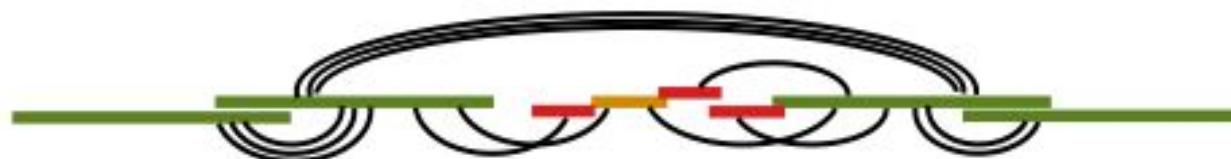


Graphs

Simplify the graph

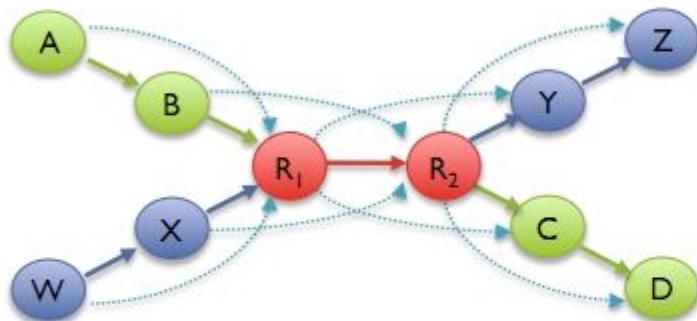


Add scaffolding information



de Bruijn Graphs

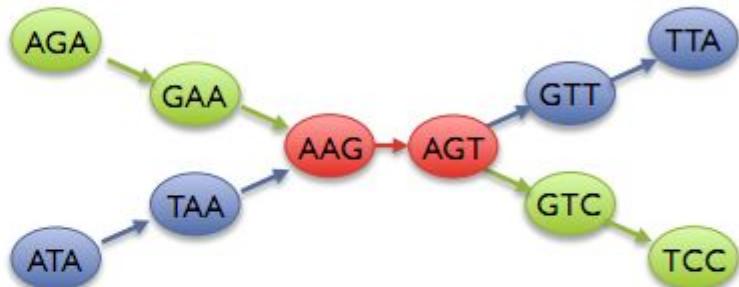
Overlap Graph



Long read assemblers

- Repeats depends on read length
- Read coherency, placements kept
- Tangled by high coverage

de Bruijn Graph

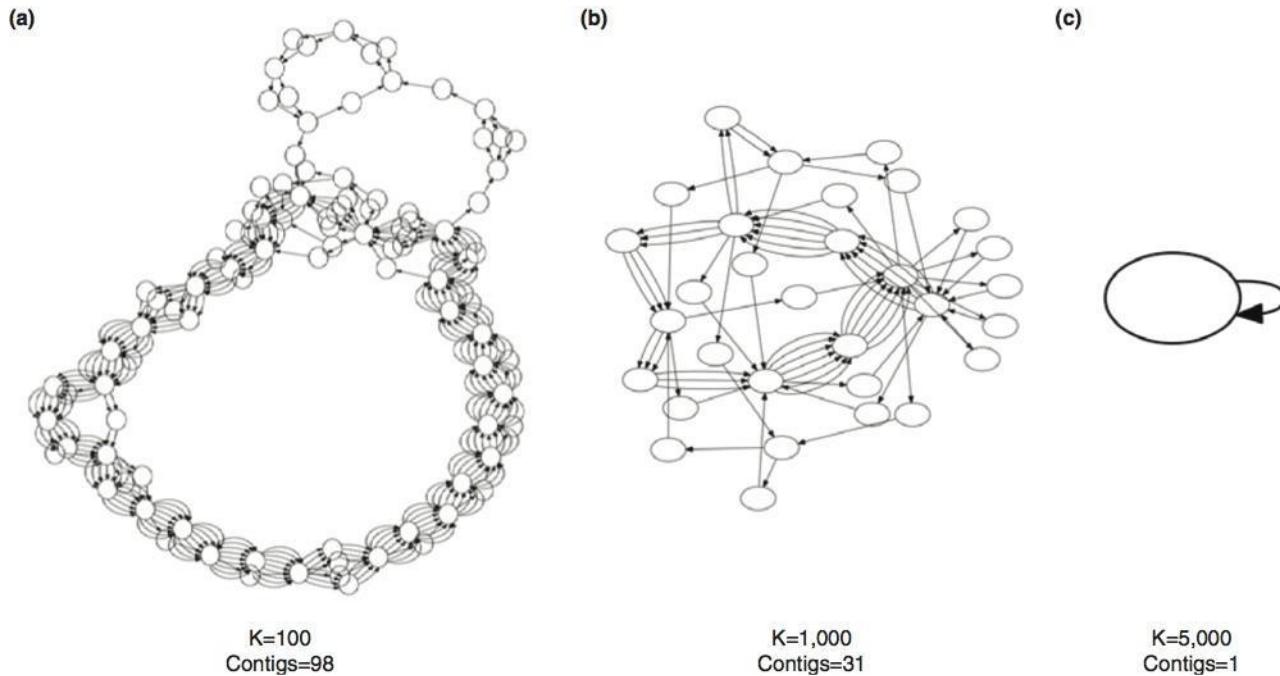


Short read assemblers

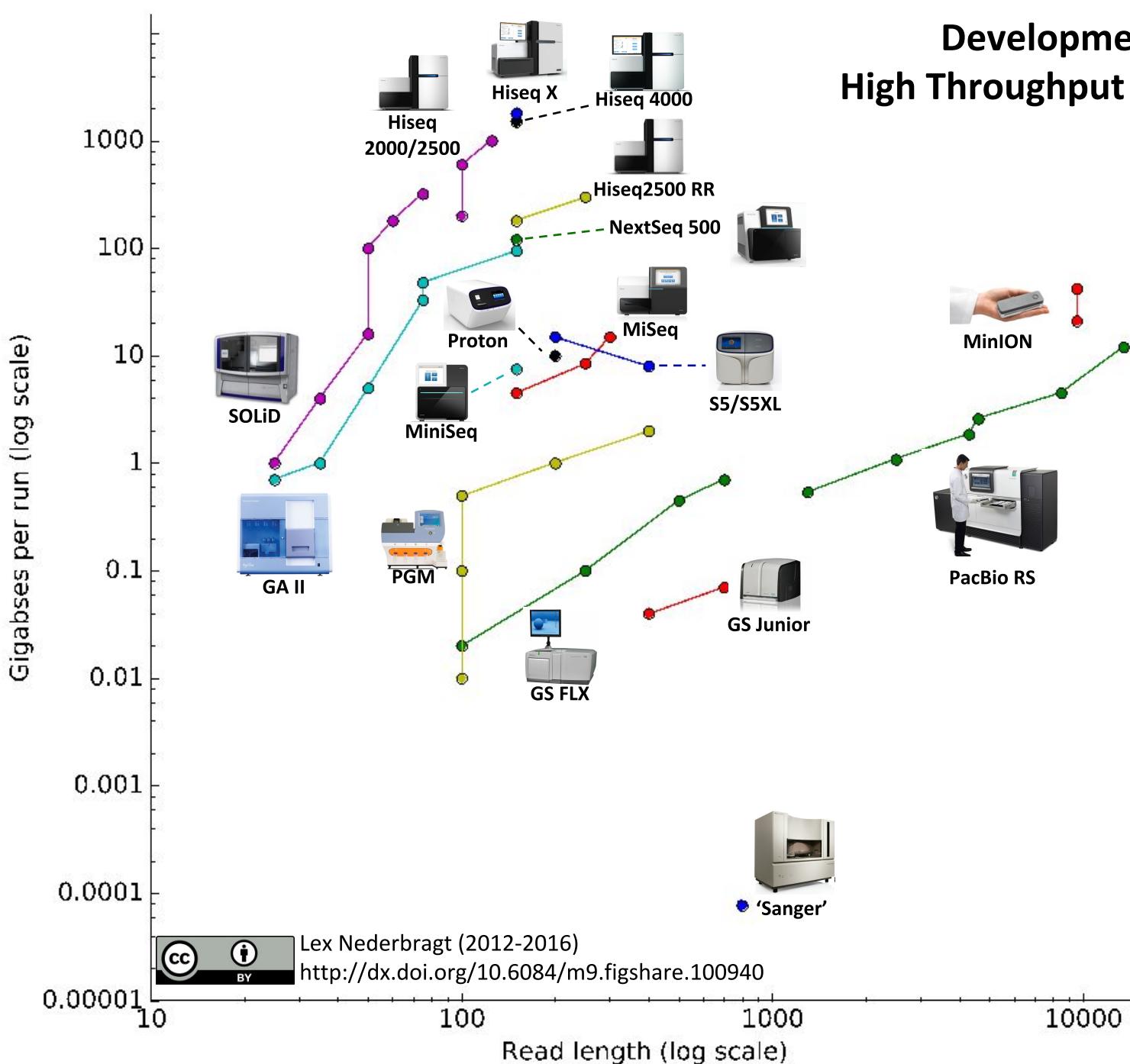
- Repeats depends on word length
- Read coherency, placements lost
- Robust to high coverage

Read length matters

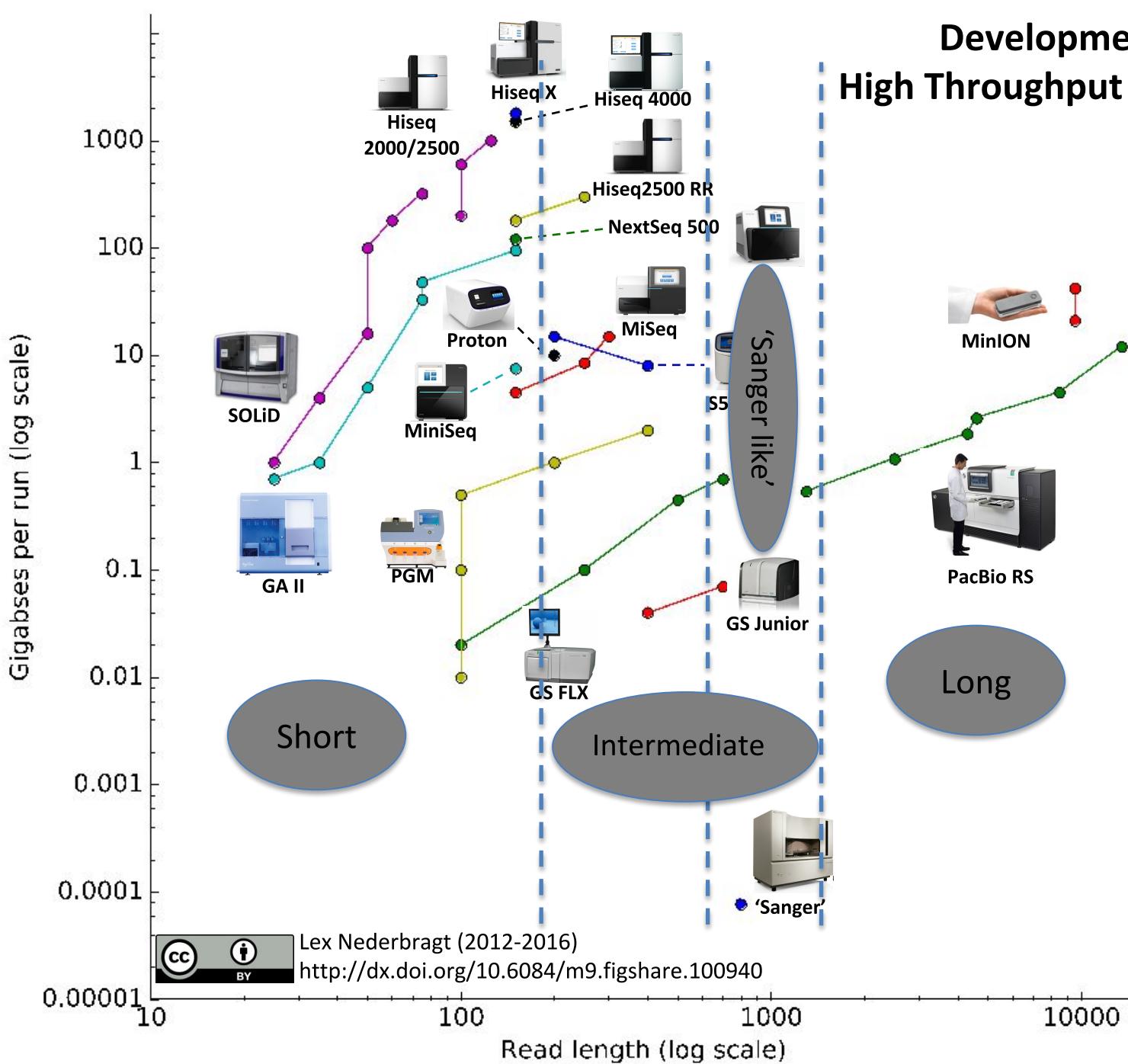
5.2 Mb circular genome, infinite error-free reads



Developments in High Throughput Sequencing

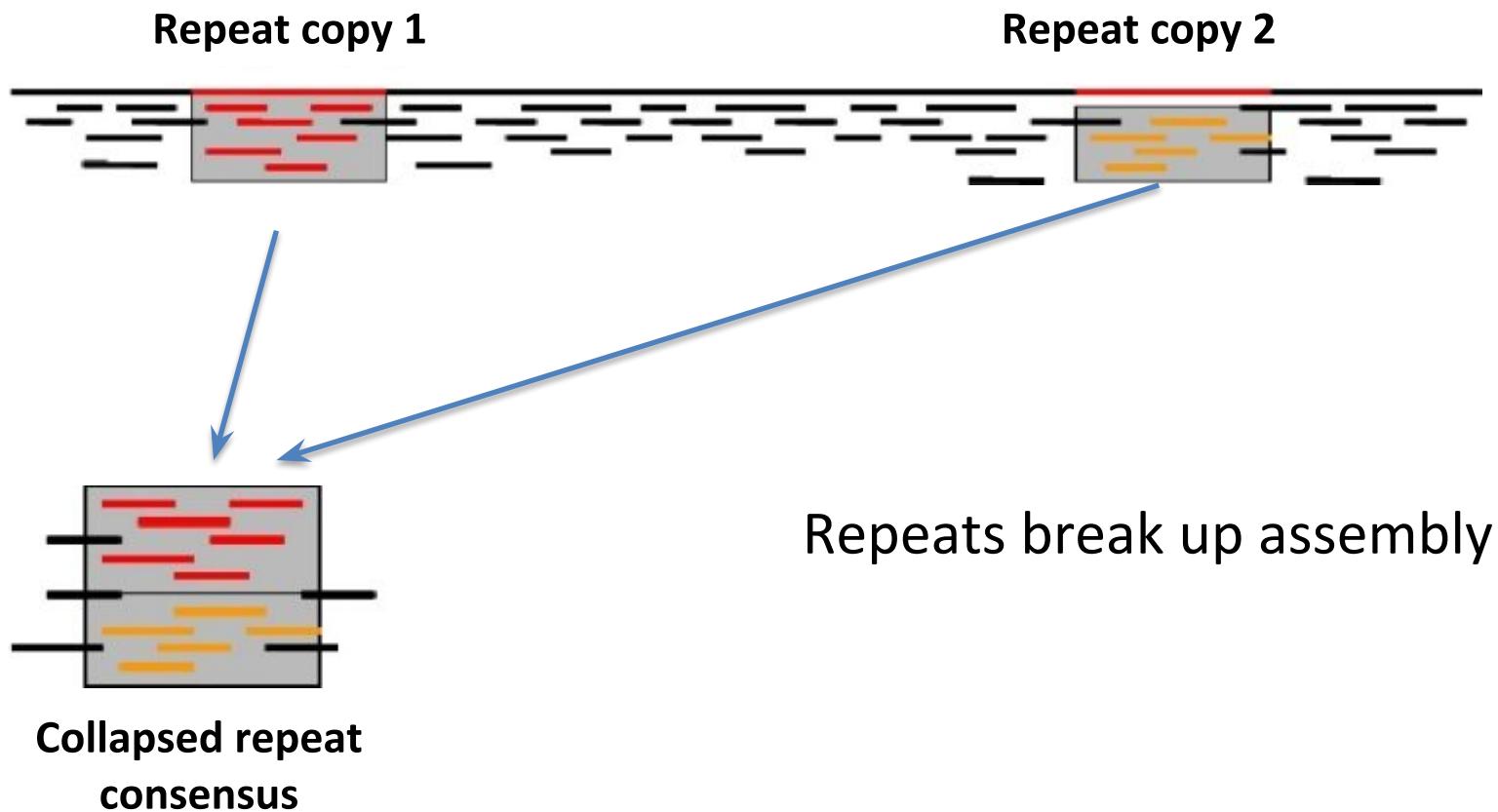


Developments in High Throughput Sequencing

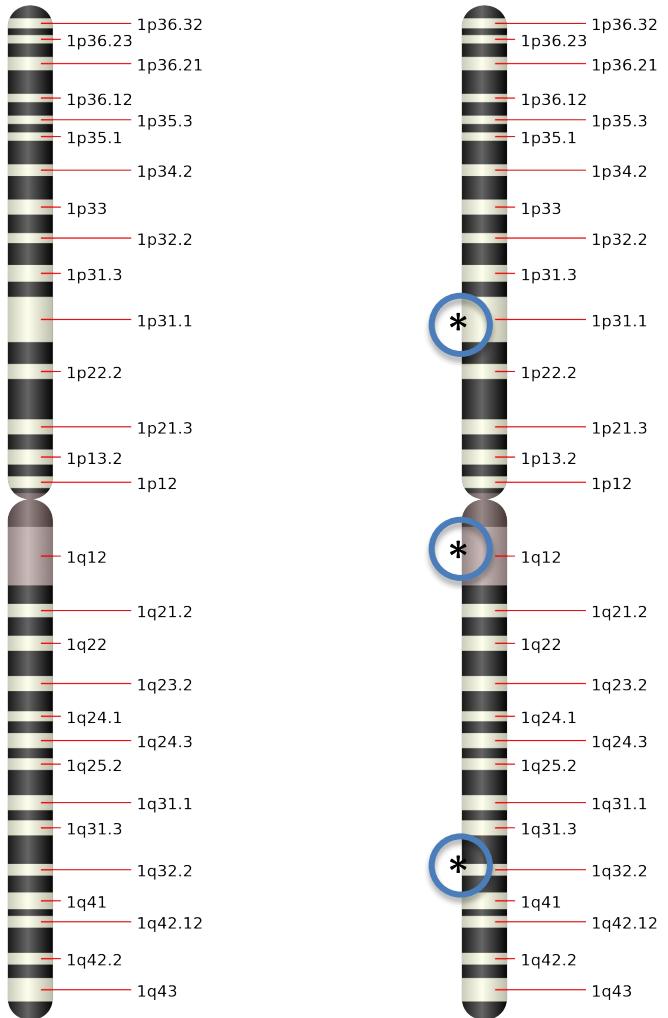


Why is genome assembly such
a difficult problem?

1) Repeats



2) Diploidy

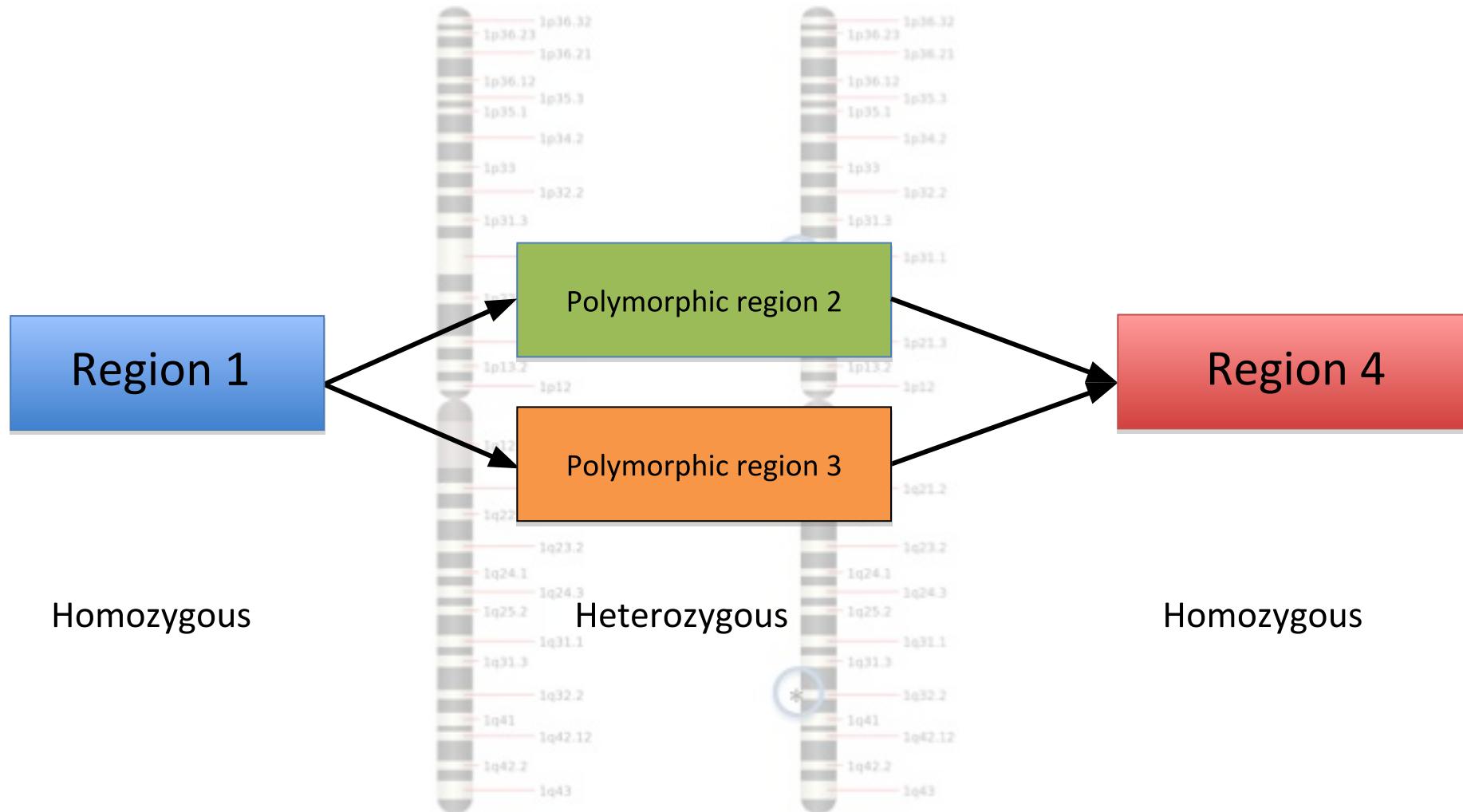


Differences
between sister
chromosomes

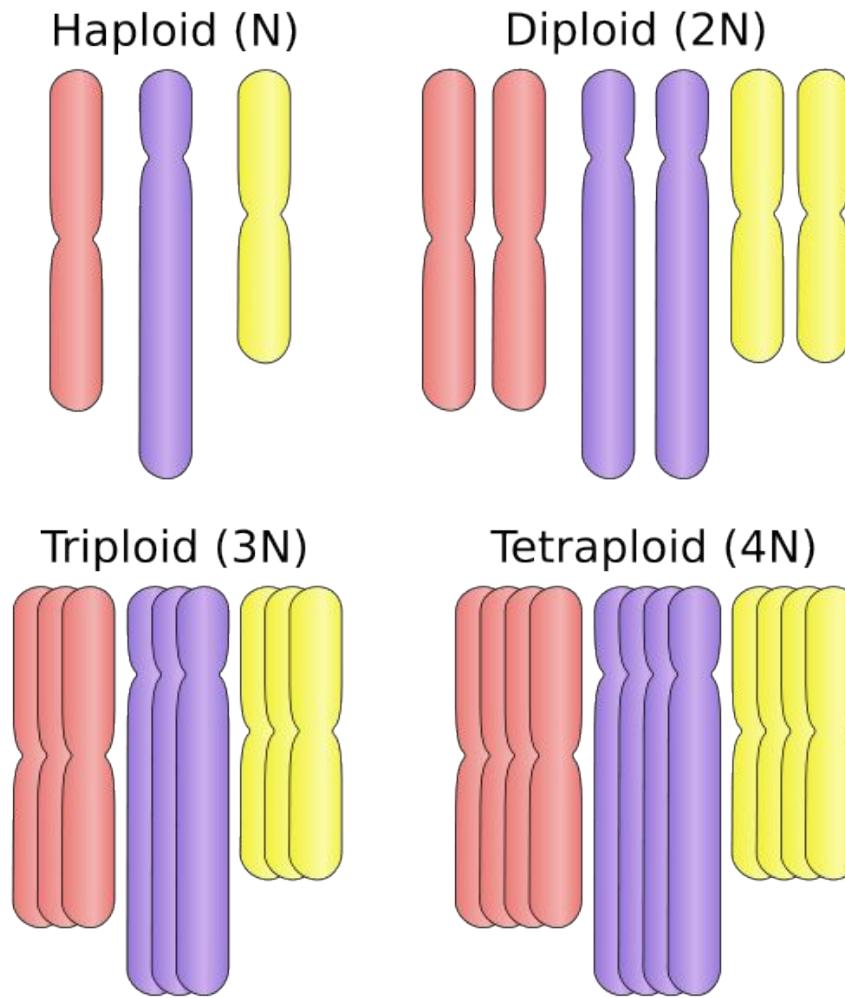


'heterozygosity'

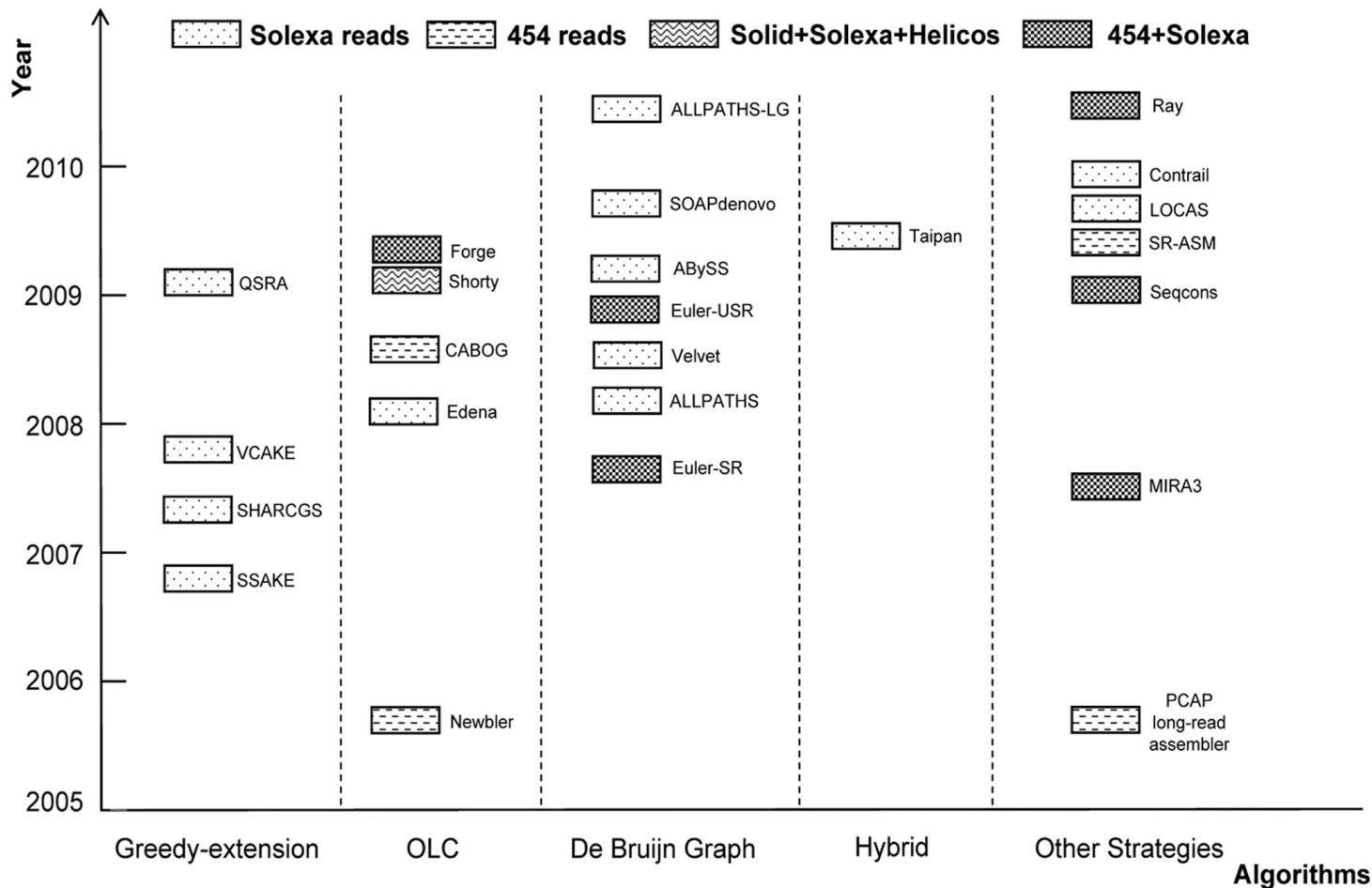
2) Diploidy



3) Polyploidy

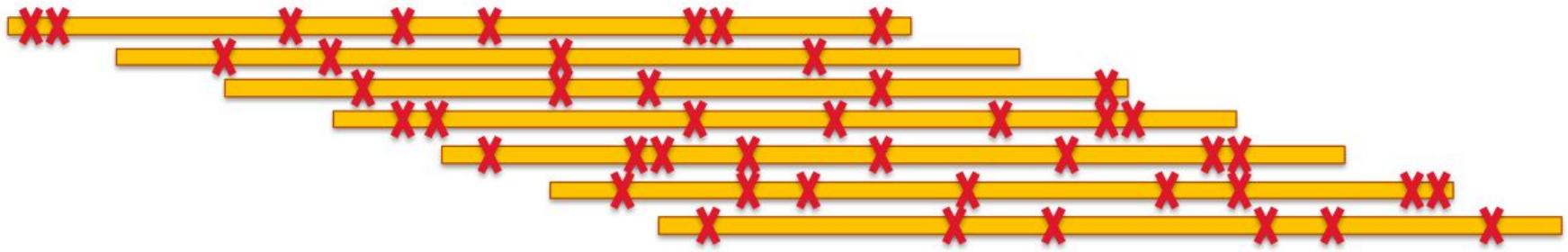


4) Many programs to choose from



Assembly with noisy single molecule sequencing data

Usage of long reads



- Problem: higher error rates
- Overlaps more difficult/expensive to find
- OLC more commonly used than for 2nd generation data

Long read assembly strategies

- Alt 0: Scaffolding short read asms
- Alt 1: Correct reads, then assemble
- Alt 2: Assemble reads, then correct

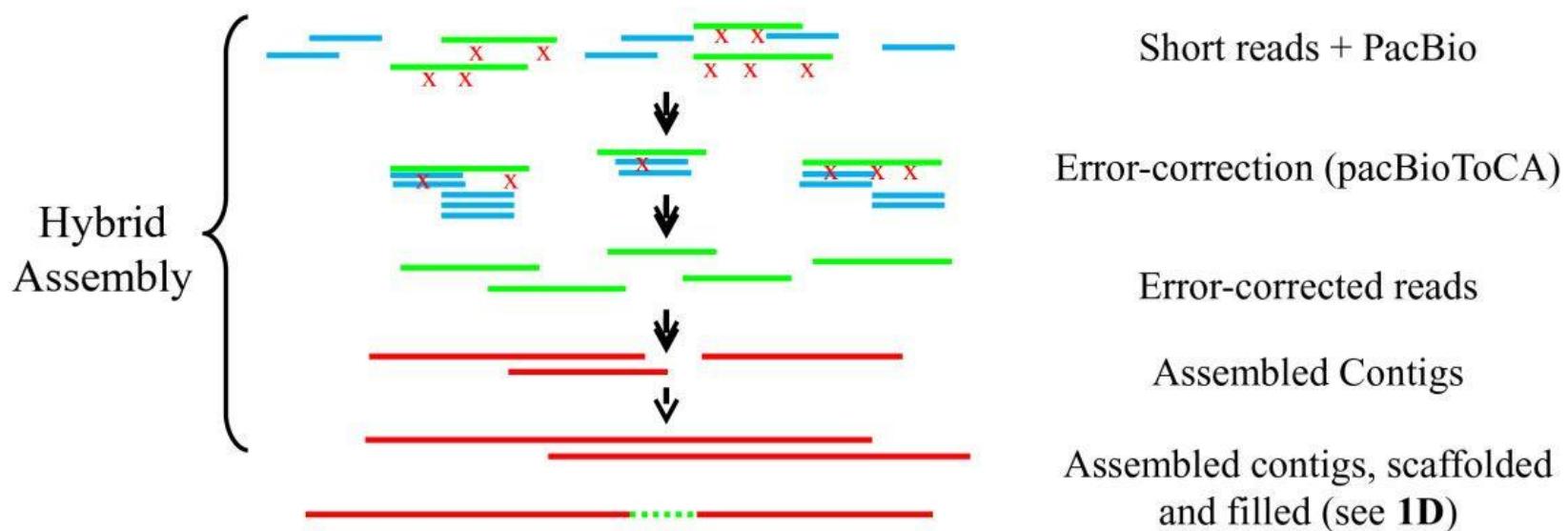
Scaffolding and gap closing (hybrid)



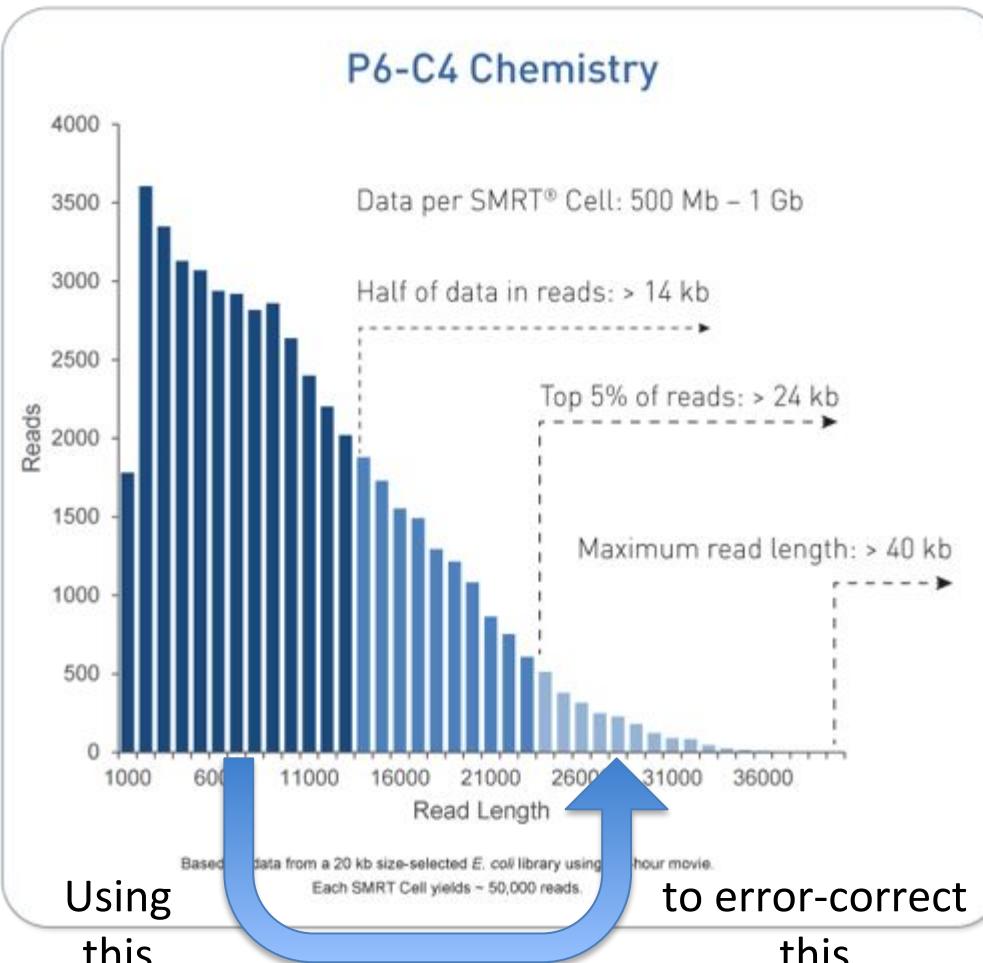
Correct, assemble

- Do pairwise comparison, find shorter reads that support the longer
- Align supporting reads, correct longer reads
- Overlap-Layout-Consensus on corrected reads
- Polish assembly

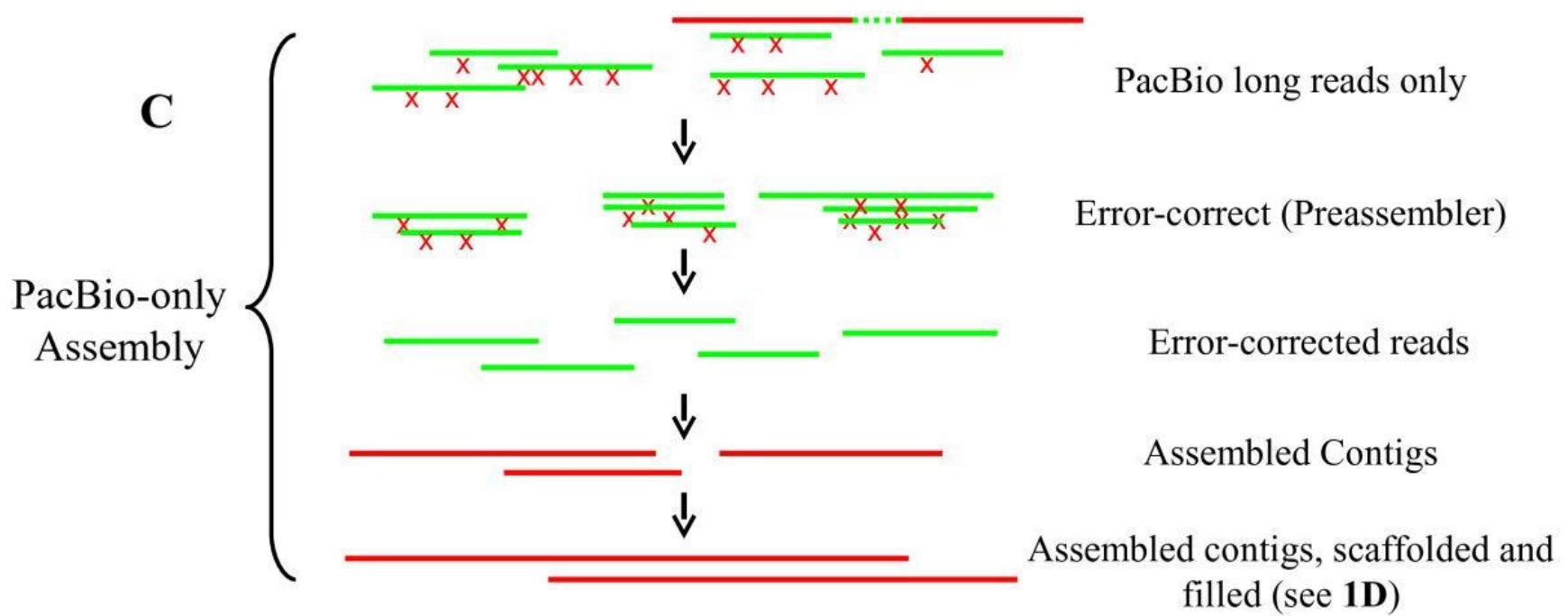
Mapping and error correcting (hybrid)



Hierarchical approach (self-correcting)



Short read error correction



Assemble, correct

- Compare reads, find overlaps
- Assemble reads, knowing things will be wrong
- Align reads to assembly
- Correct assembly

Questions?