

Principles and problems of de novo genome assembly

Karin Lagesen
Norwegian Veterinary Institute



Material adapted from slides
provided by Lex Nederbragt



1

What is this thing called 'genome assembly'?

2

What is a genome assembly?

A hierarchical data structure
that maps the sequence data
to a putative reconstruction of the target

Miller et al 2010, Genomics 95 (6): 315-327

3

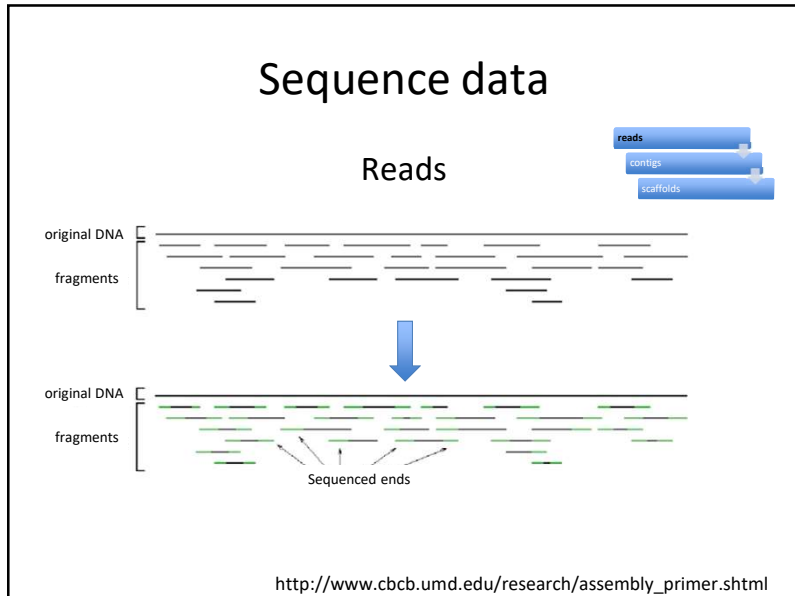
Hierarchical structure

reads

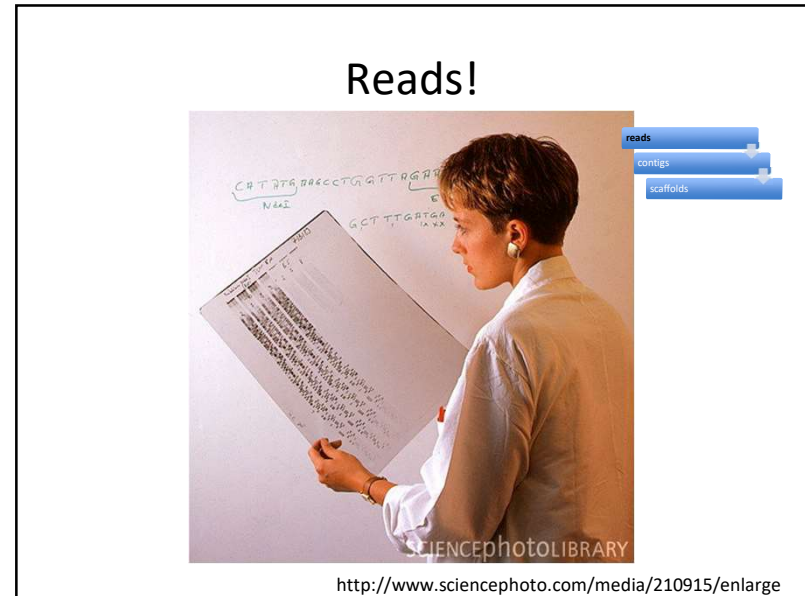
contigs

scaffolds

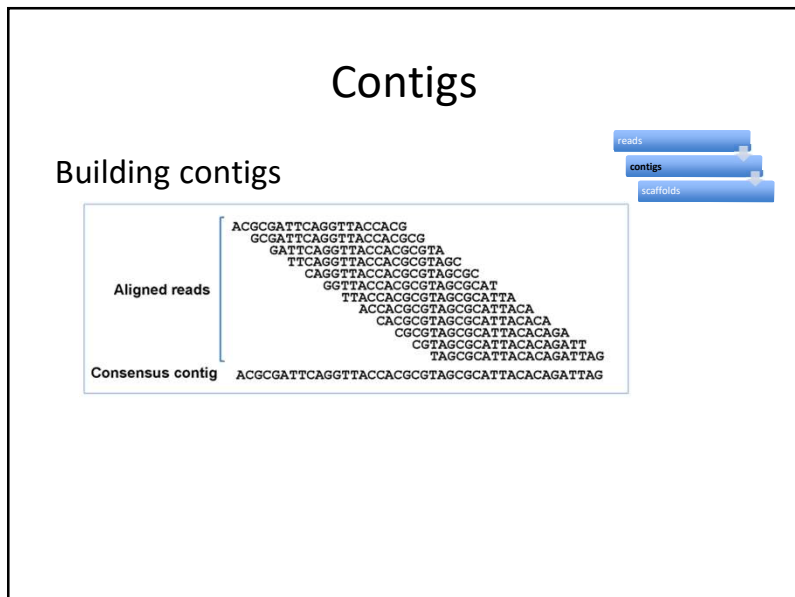
4



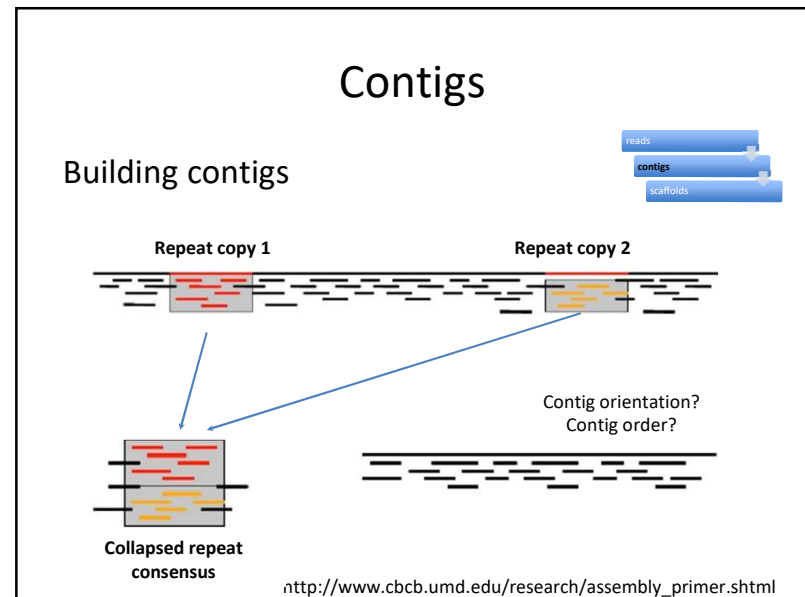
5



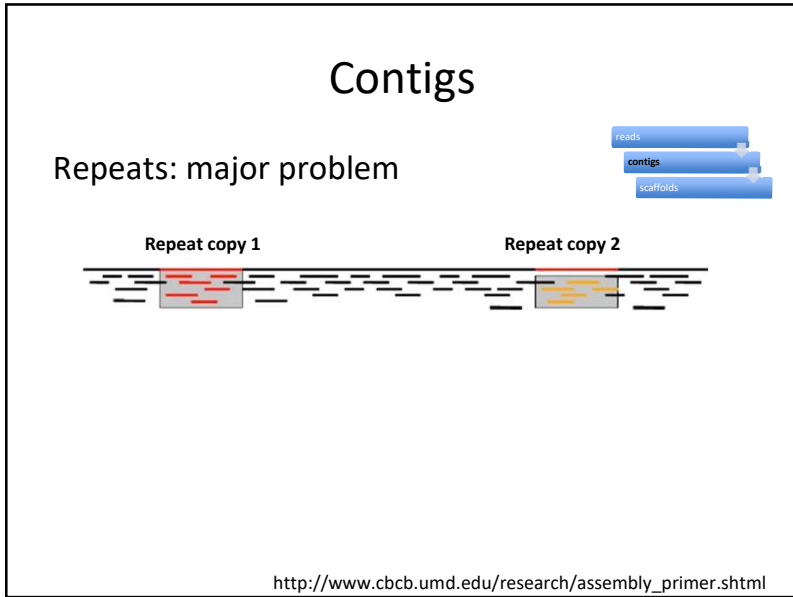
6



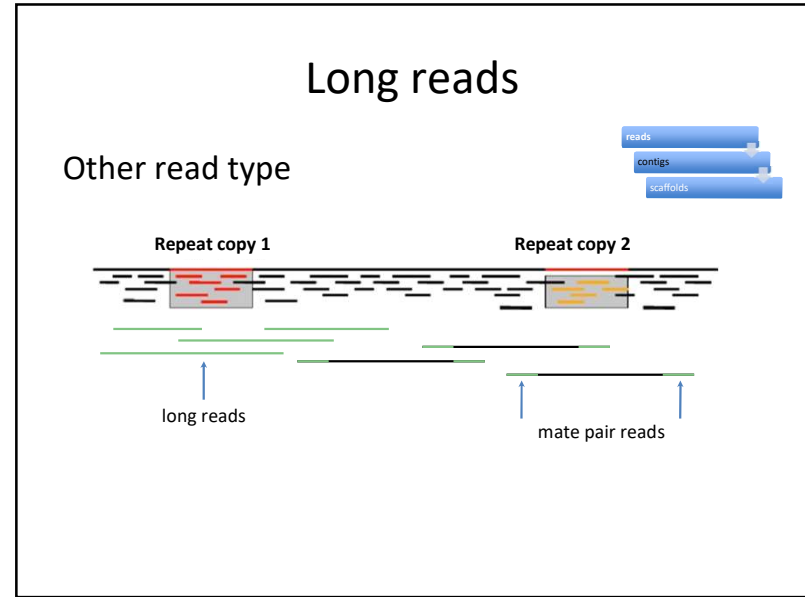
7



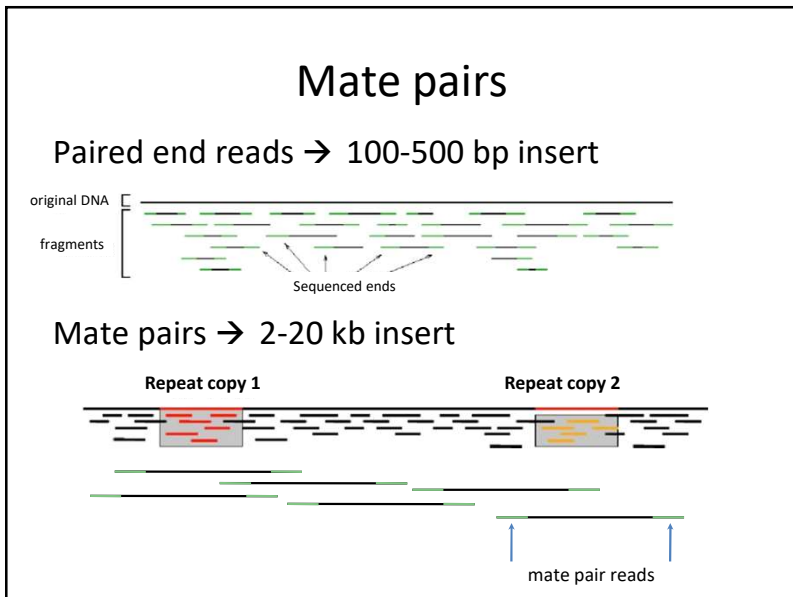
8



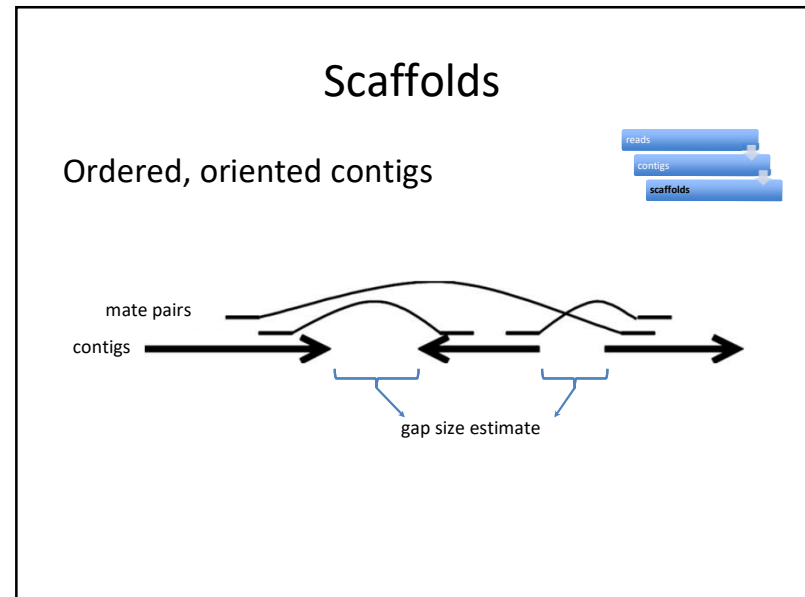
9



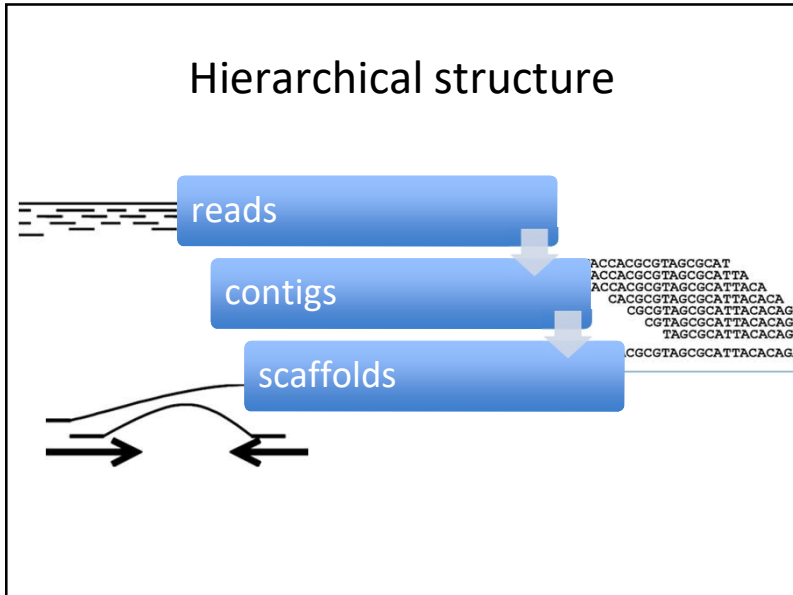
10



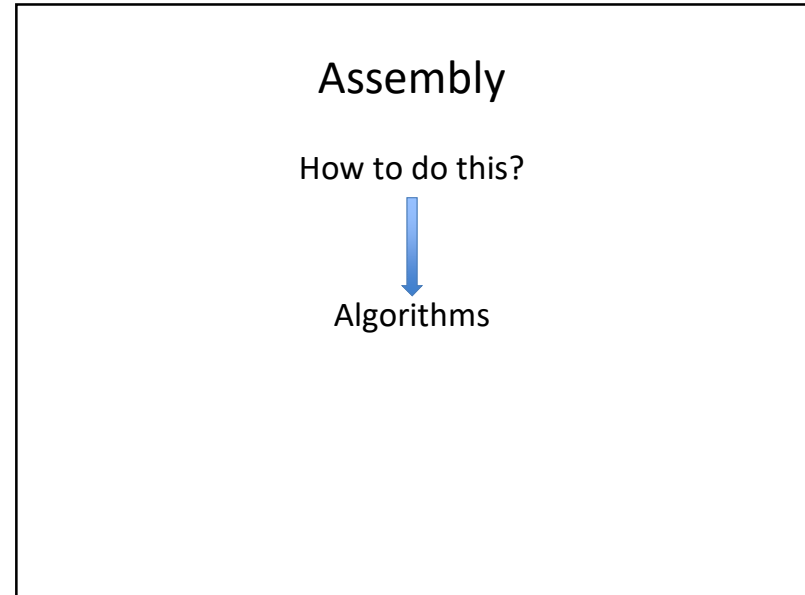
11



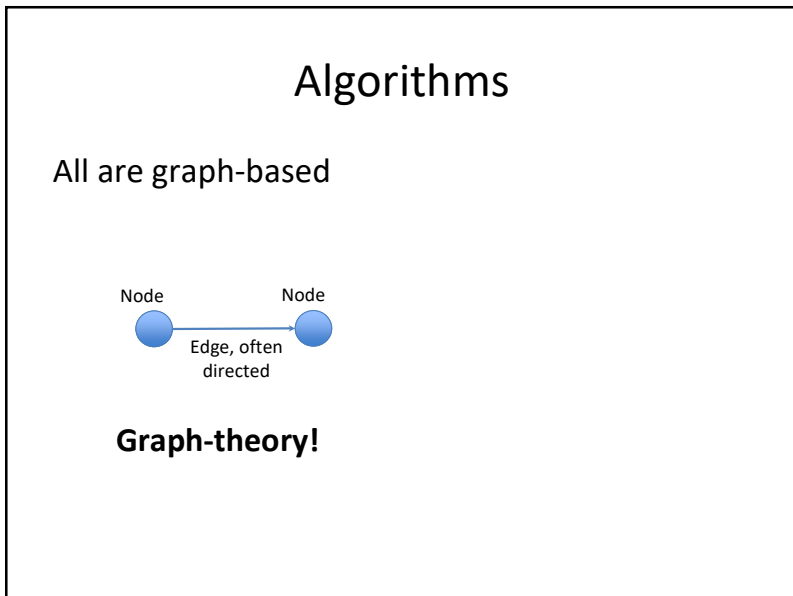
12



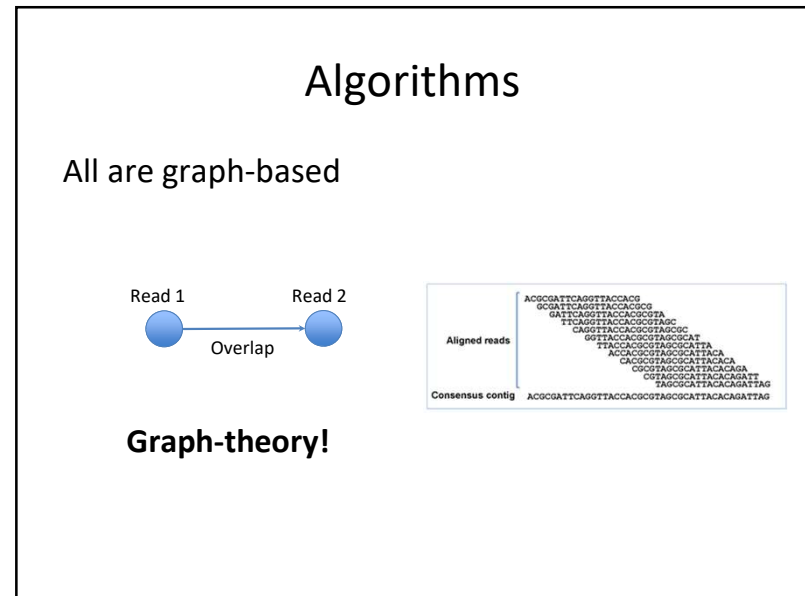
13



14



15

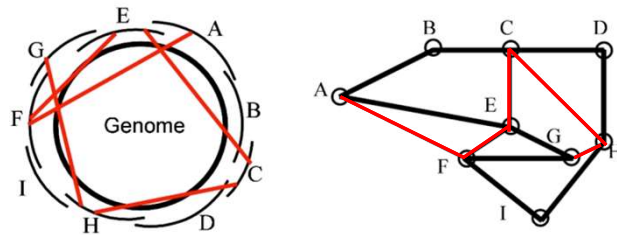


16

Algorithms

Hamiltonian path

– a path that contains all the nodes



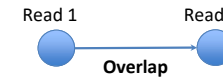
https://www.cbcb.umd.edu/research/assembly_primer

17

Algorithms

Overlap calculation (alignment)

– computationally intensive



Aligned reads

```

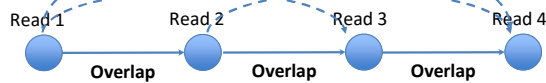
ACGCGATTTCAGGTTACCACG
GCATTTCAGGTTACCACGCG
GATTCAGGTTACCACGCGTA
TTCAGGTTACCACGCGTAGC
CAGGTTACCACGCGTAGCCG
GGTTACCACGCGTAGCCAT
TTACCACGCGTAGCCGATTA
ACCACGCGTAGCCGATTACA
CACGCGTAGCCGATTACACA
CCGCGTAGCCGATTACACAGA
CGTAGCCGATTACACAGATT
TAGCGCATTACACAGATTAG
    
```

18

Algorithms

Path through the graph

→ contig



Aligned reads

```

ACGCGATTTCAGGTTACCACG
GCATTTCAGGTTACCACGCG
GATTCAGGTTACCACGCGTA
TTCAGGTTACCACGCGTAGC
CAGGTTACCACGCGTAGCCG
GGTTACCACGCGTAGCCAT
TTACCACGCGTAGCCGATTA
ACCACGCGTAGCCGATTACA
CACGCGTAGCCGATTACACA
CGCGTAGCCGATTACACAGA
CGTAGCCGATTACACAGATT
TAGCGCATTACACAGATTAG
    
```

Consensus contig ACGCGATTTCAGGTTACCACGCGTAGCCGATTACACAGATTAG

19

Algorithms

Many flavors



Abandoned

- Greedy extension

Two most used

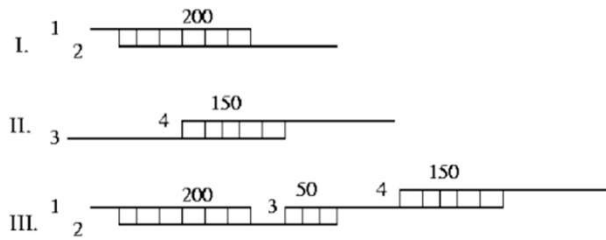
- Overlap Layout Consensus
- de Bruijn graph

<http://www.waiiauasodaworks.com/images/flavors2009.jpg>

20

Greedy extension

Oldest

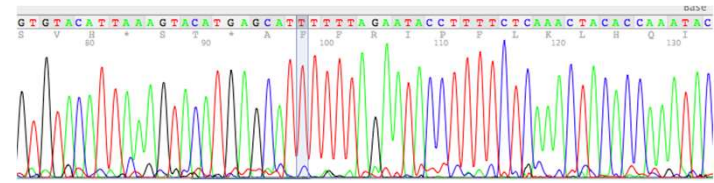


https://www.cbcb.umd.edu/research/assembly_primer

21

Overlap-Layout-Consensus

Developed for Sanger-type reads (longer reads)



22

Overlap-Layout-Consensus

Steps

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

23

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 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

ACGC GATT CAGG TAC CACG

24

Overlap-Layout-Consensus

A Read Layout

R₁: GACCTACA
 R₂: ACCTACAA
 R₃: CCTACAAG
 R₄: CTACAAGT

A: TACAAGTT
 B: ACAAGTTA
 C: CAAGTTAG
 X: TACAAGTC
 Y: ACAAGTCC
 Z: CAAGTCCG

B Overlap Graph

Schatz M C et al. Genome Res. 2010;20:1165-1173

25

de Bruijn graphs

Developed outside of DNA-related work
 – Best solution for short(er) reads

Read GACCTACA
 GAC
 ACC
 CCT
 CTA
 TAC
 ACA

K-mers (K=3)

de Bruijn graph

GAC
ACC

K-1 bases overlap

26

Graphs

C de Bruijn Graph

Schatz M C et al. Genome Res. 2010;20:1165-1173

27

Graphs

A Read Layout

R₁: GACCTACA
 R₂: ACCTACAA
 R₃: CCTACAAG
 R₄: CTACAAGT

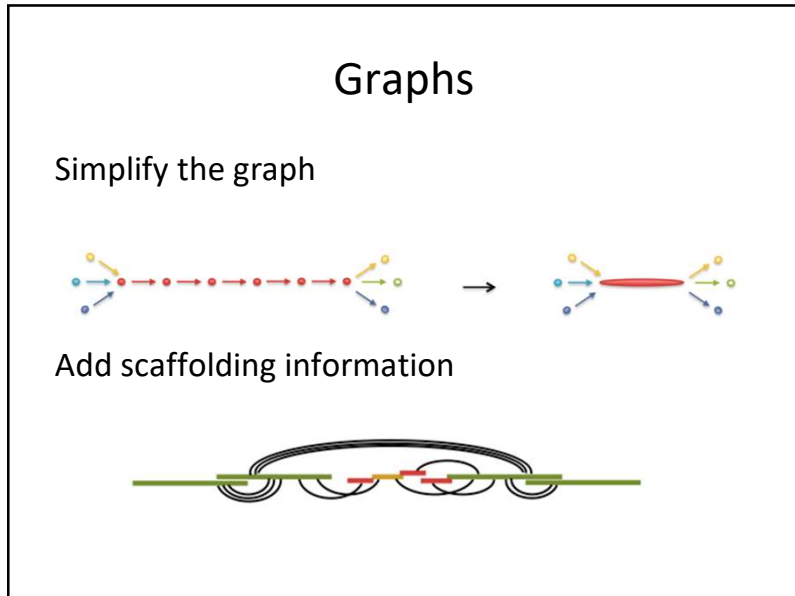
A: TACAAGTT
 B: ACAAGTTA
 C: CAAGTTAG
 X: TACAAGTC
 Y: ACAAGTCC
 Z: CAAGTCCG

B Overlap Graph

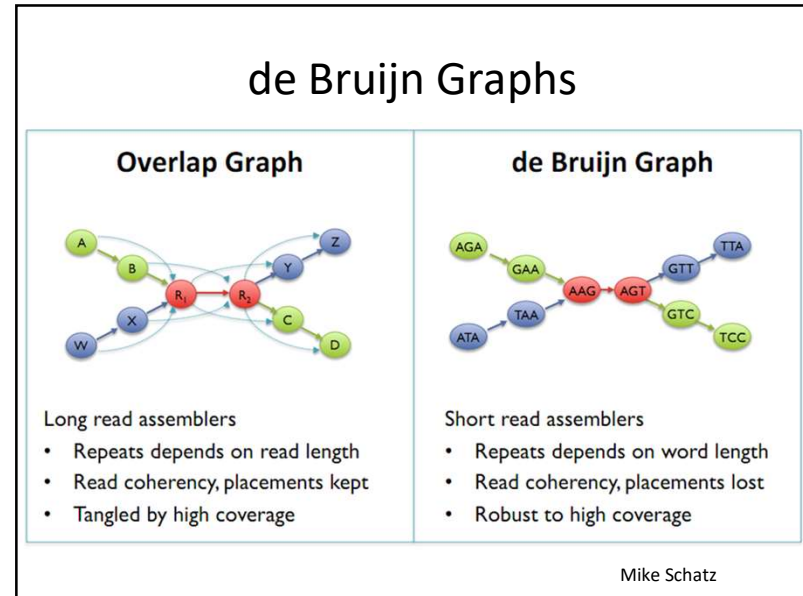
C de Bruijn Graph

Schatz M C et al. Genome Res. 2010;20:1165-1173

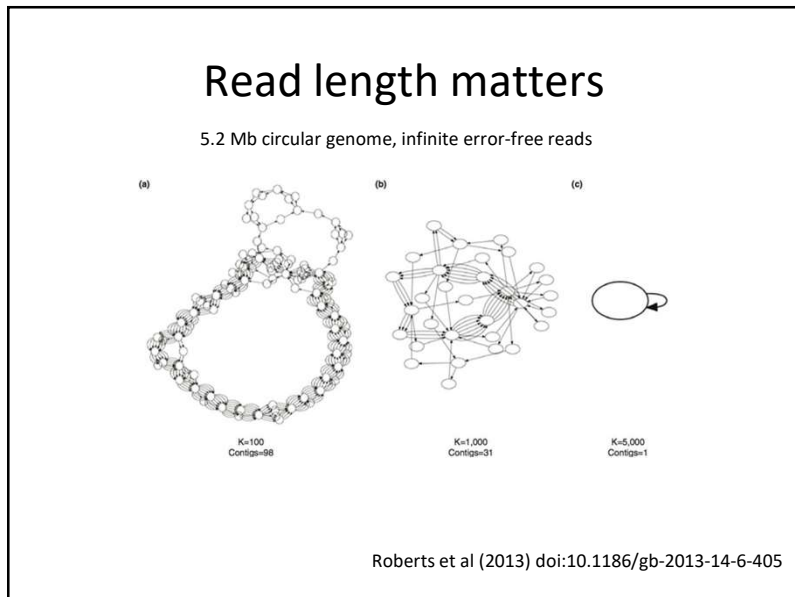
28



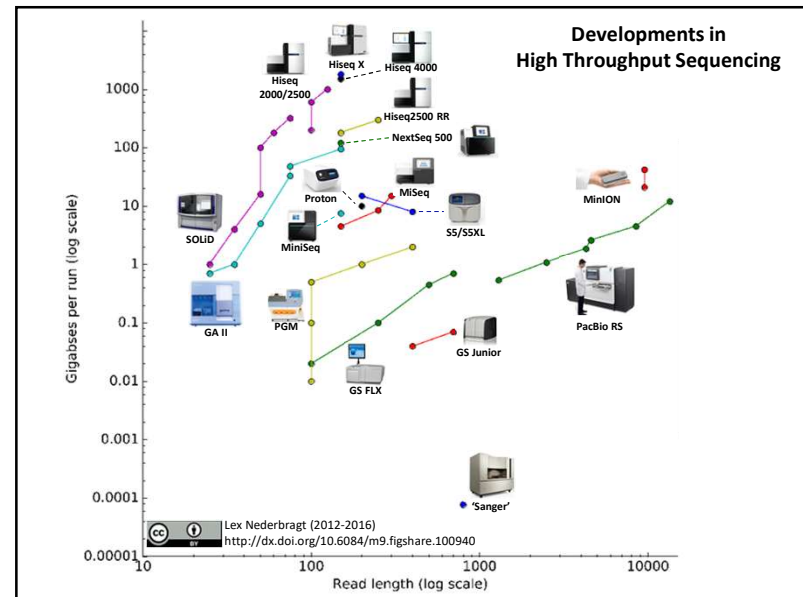
29



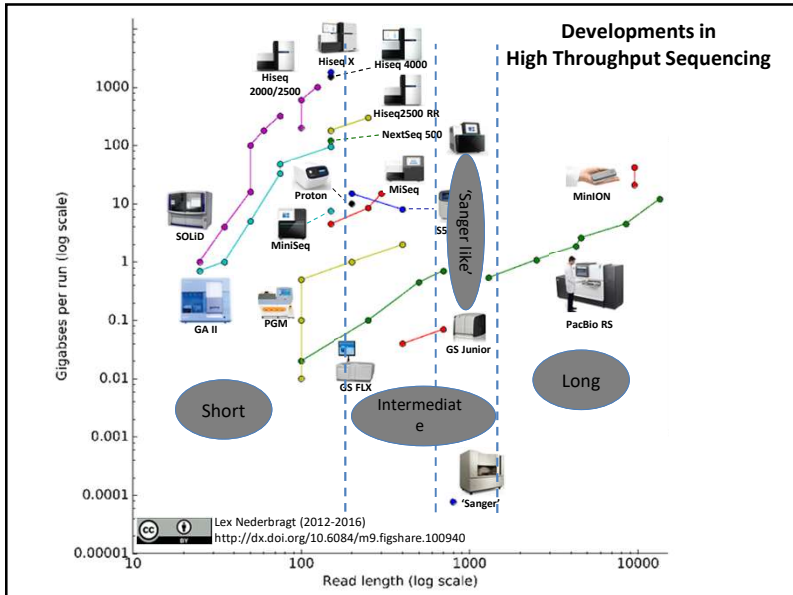
30



31



32



33

Quality matters

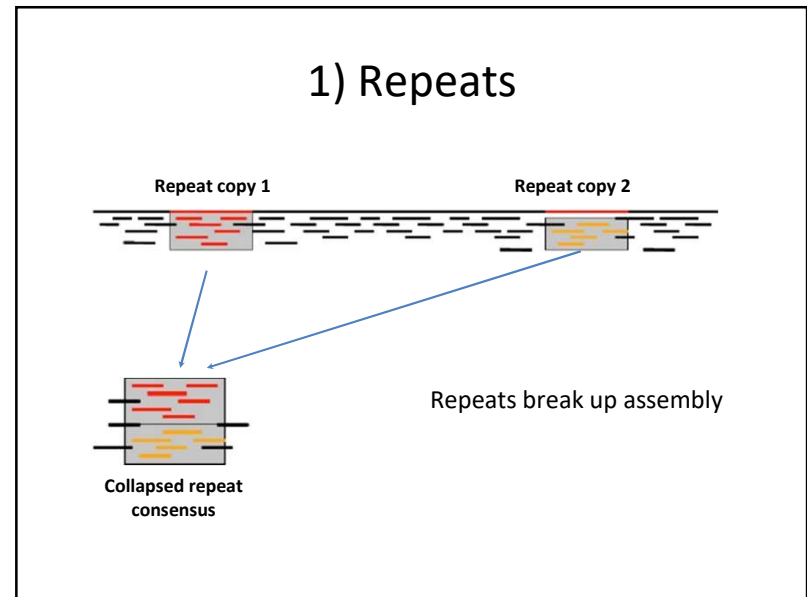
Too many errors → hard to find overlaps

<http://schatzlab.cshl.edu/presentations/2012-01-17.PAG.SMRTassembly.pdf>

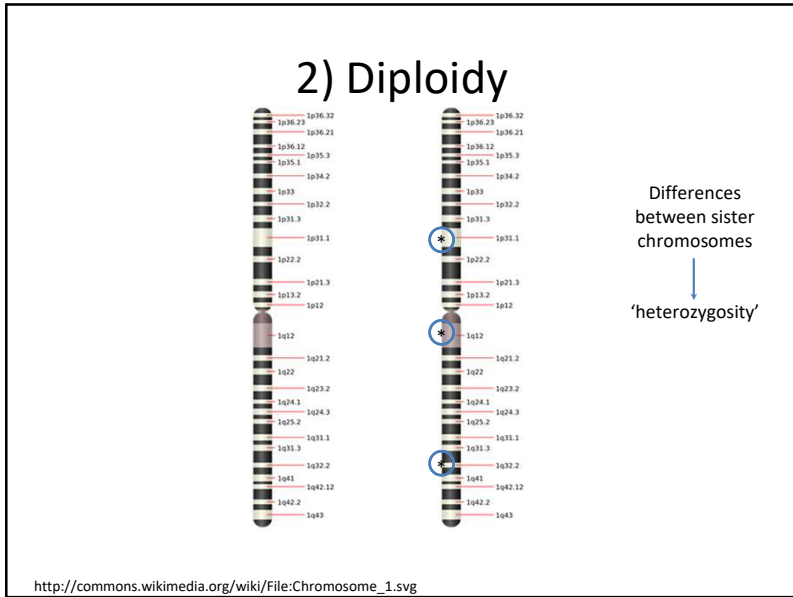
34

Why is genome assembly such a difficult problem?

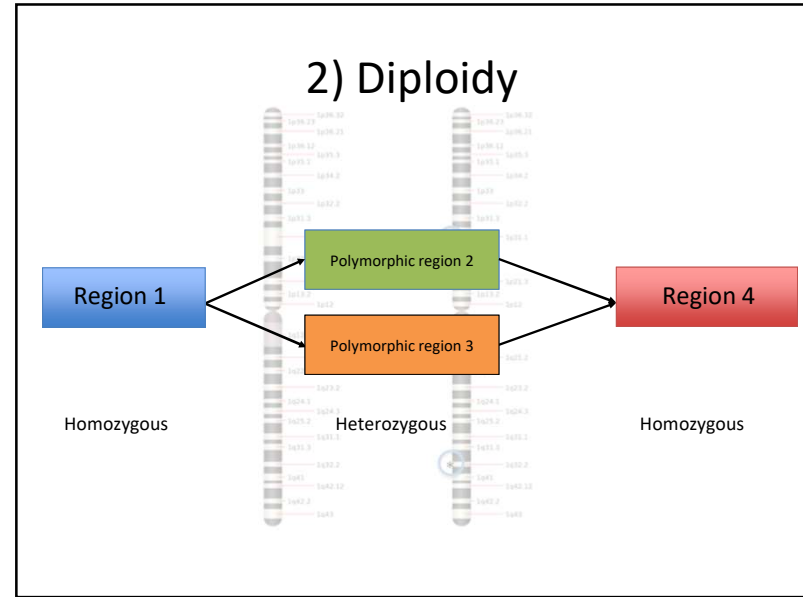
35



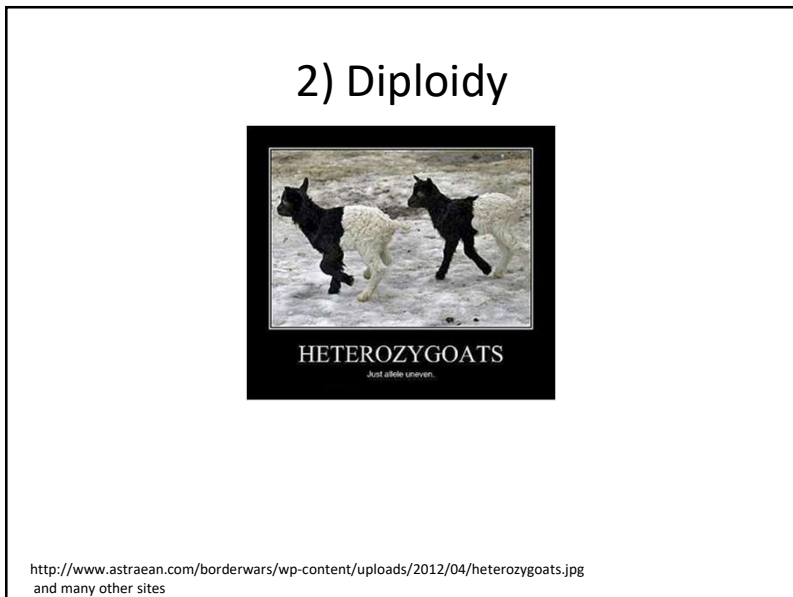
36



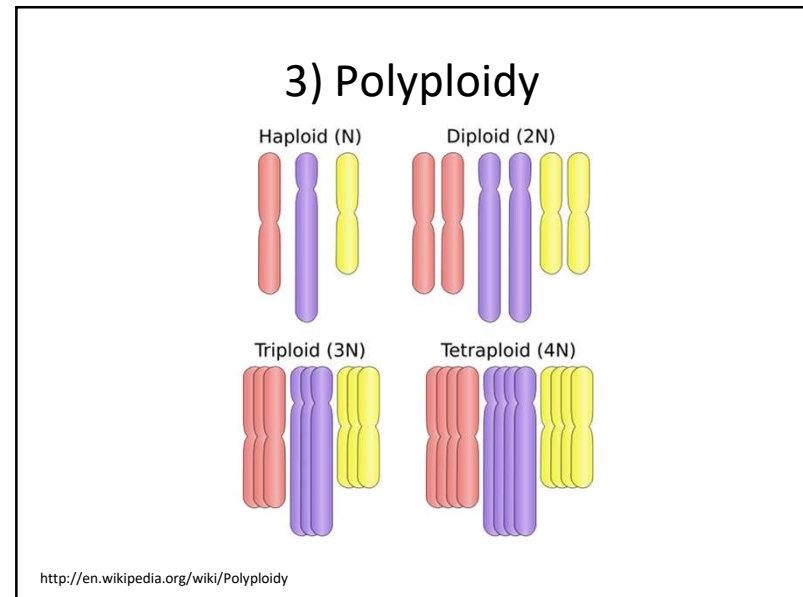
37



38

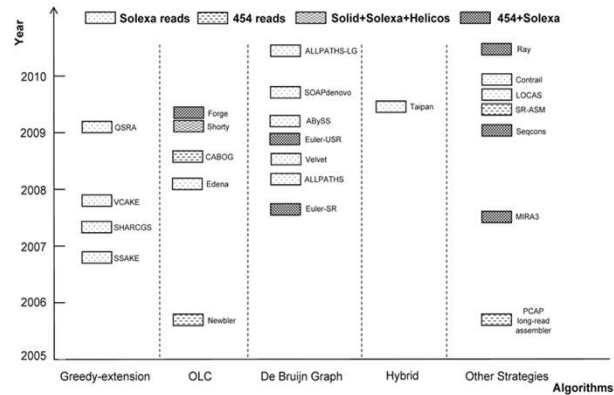


39



40

4) Many programs to choose from



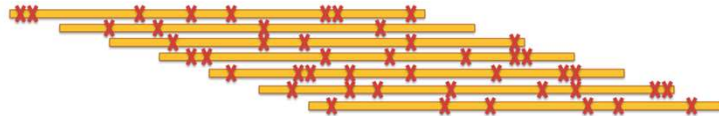
Zhang et al. PLoSOne 2011

41

Assembly with noisy single molecule sequencing data

42

Usage of long reads



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

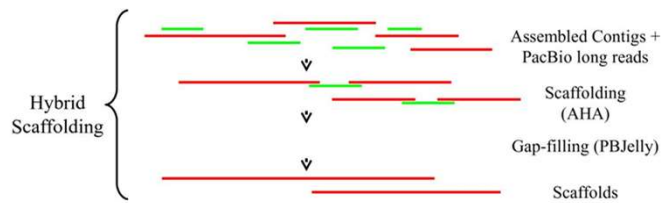
43

Long read assembly strategies

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

44

Scaffolding and gap closing (hybrid)



Powers *et al.*, BMC genomics 2013

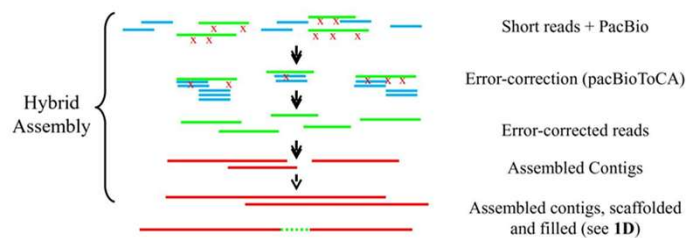
45

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

46

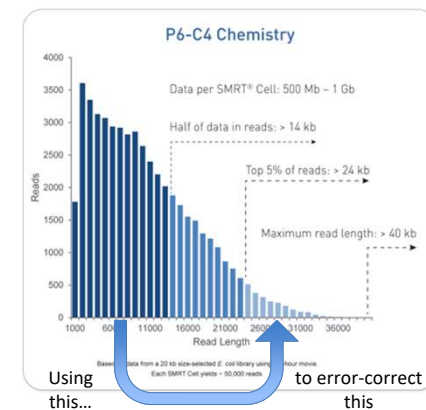
Mapping and error correcting (hybrid)



Powers *et al.*, BMC genomics 2013

47

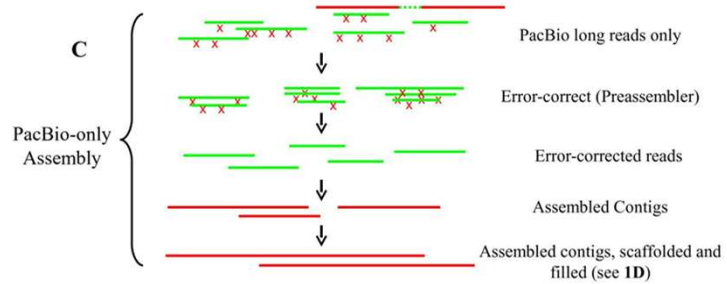
Hierarchical approach (self-correcting)



<https://genome.duke.edu/cores-and-services/sequencing-and-genomic-technologies/pacbio>

48

Short read error correction



Powers *et al.*, BMC genomics 2013

49

Assemble, correct

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

50

Questions?

51