



NORWEGIAN SEQUENCING CENTRE

Principles and problems of de novo genome assembly

Lex Nederbragt

Norwegian High-Throughput Sequencing Centre (NSC)

and

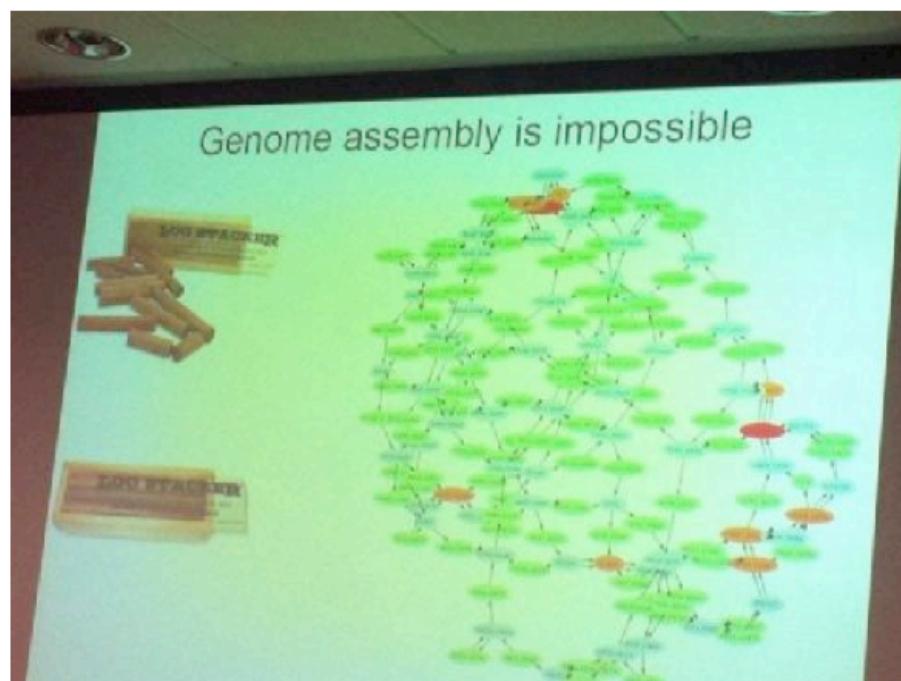
Centre for Ecological and Evolutionary Synthesis
(CEES)



@gilbertjacka

Jack A Gilbert
iobioimaging

Finally someone (Mihai Pop) admits it!!
yfrog.com/nxen3ohj



YFrog

Flag this media

5 hours ago via Twitter for iPhone  Favorite  Undo Retweet  Reply

Retweeted by [lexnederbragt](#)



What is this thing called ‘genome assembly’?

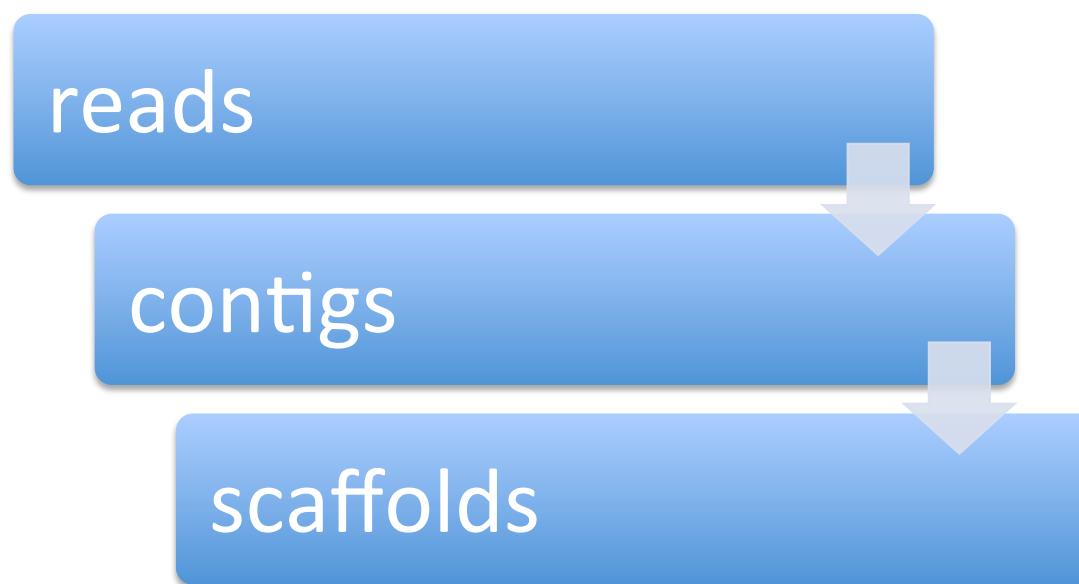
What is a genome assembly?

A hierarchical data structure

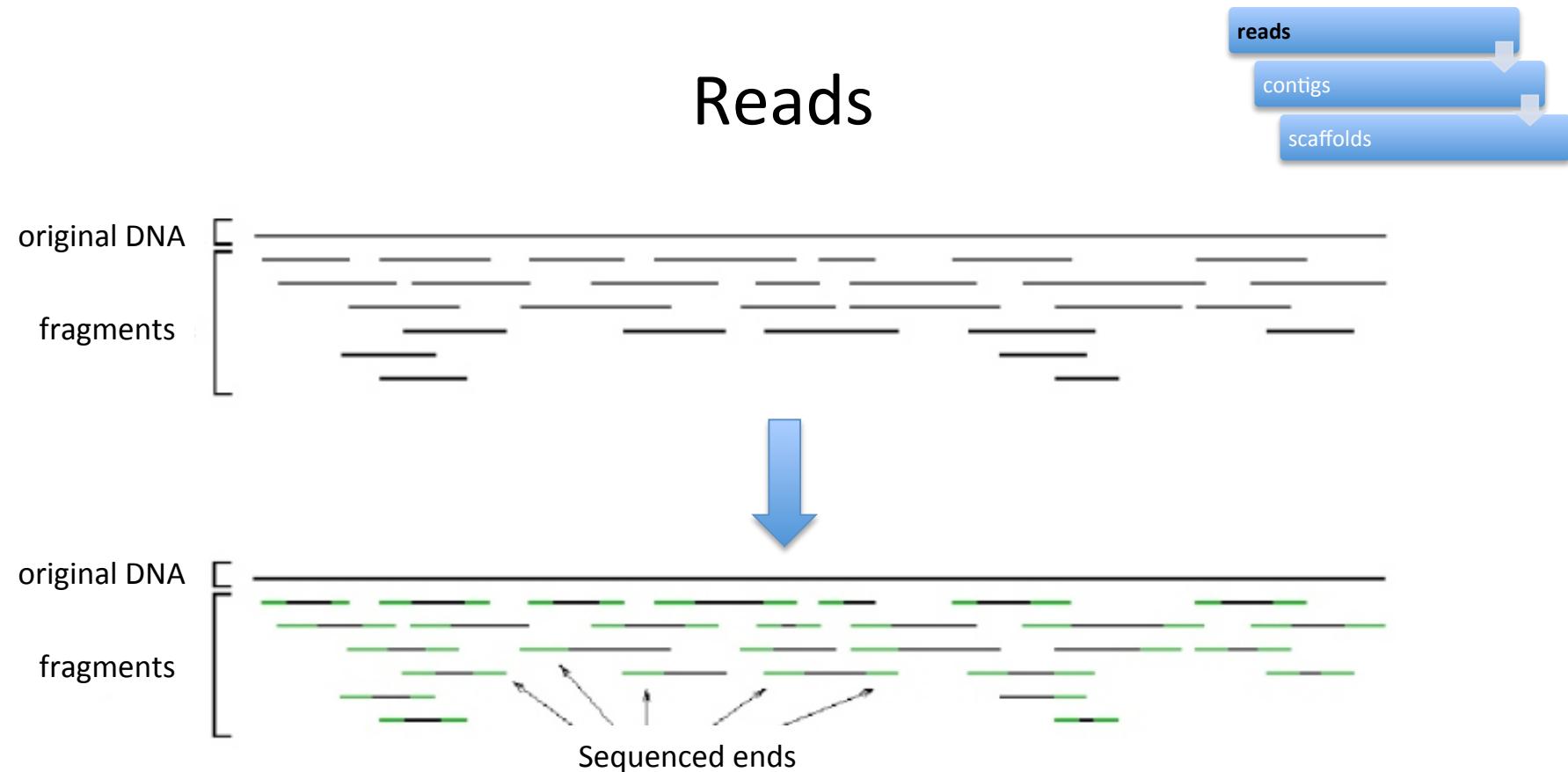
that maps the sequence data

to a putative reconstruction of the target

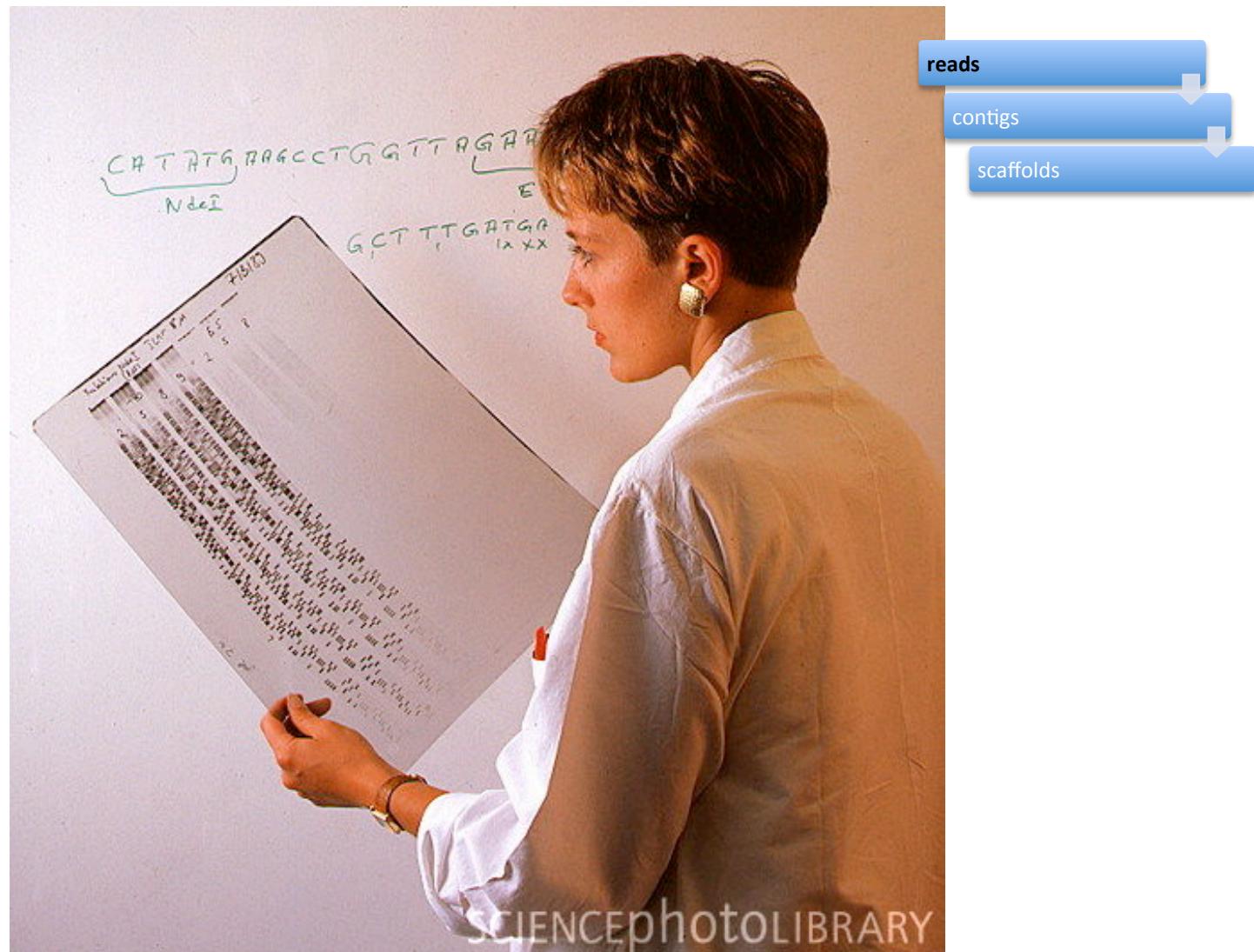
Hierarchical structure



Sequence data



Reads!



SCIENCEphotOLIBRARY

<http://www.sciencephoto.com/media/210915/enlarge>

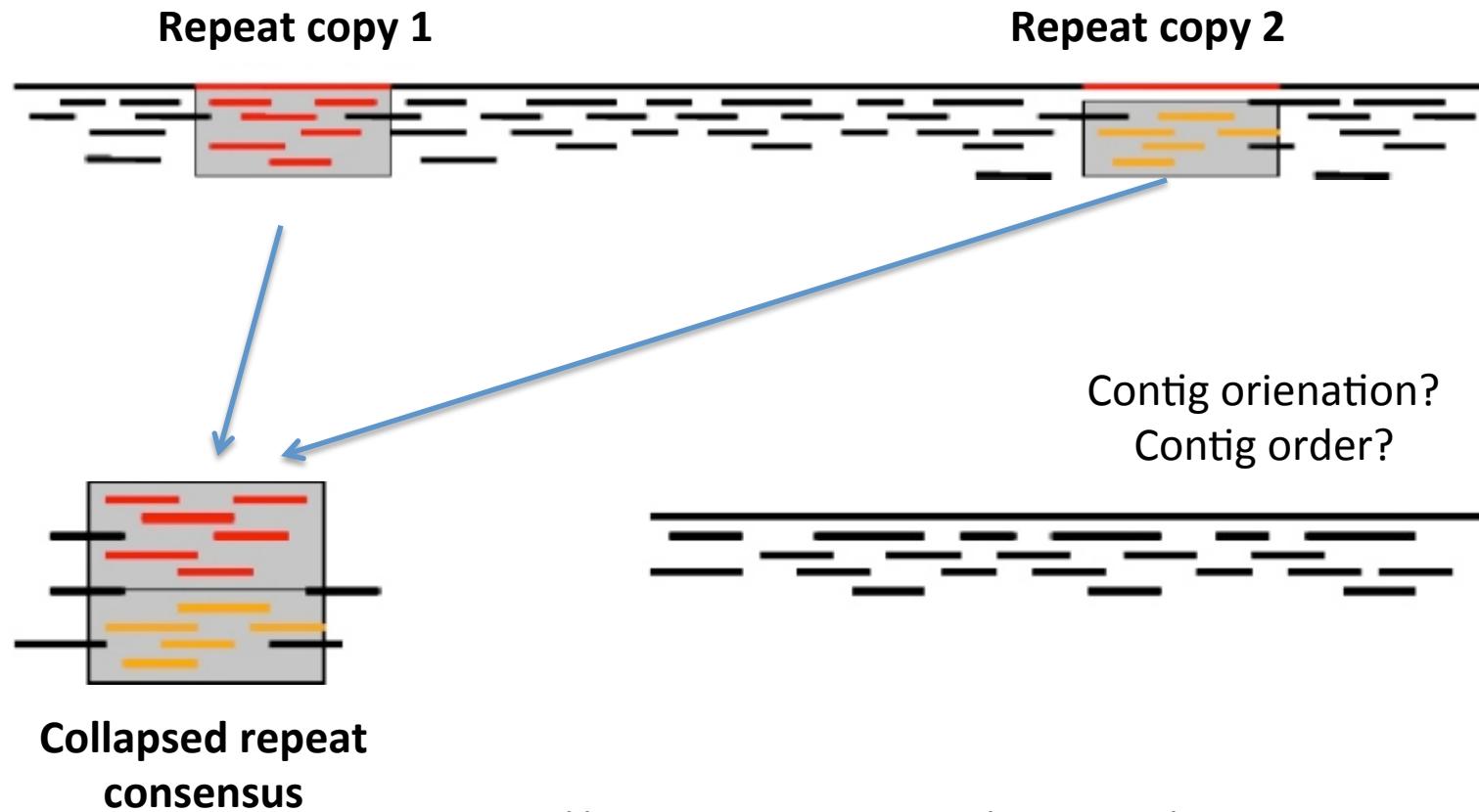
Contigs

Building contigs



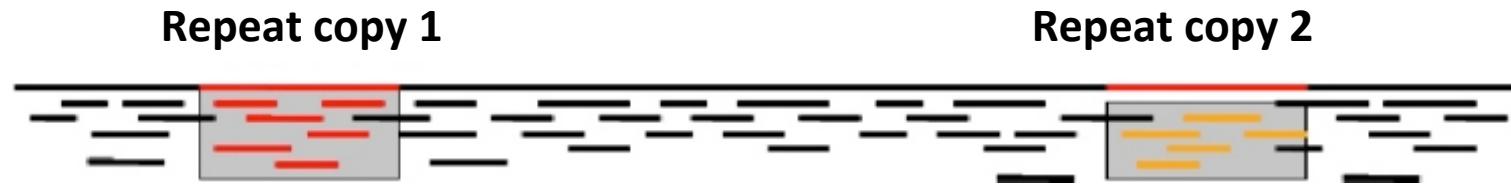
Contigs

Building contigs



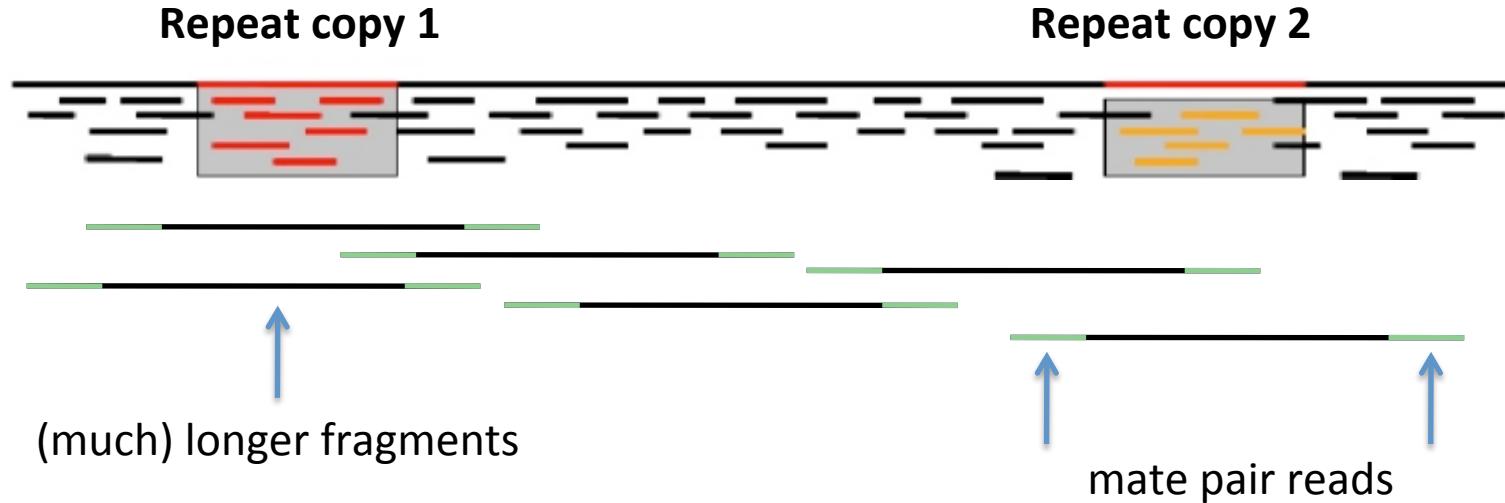
Contigs

Repeats: major problem



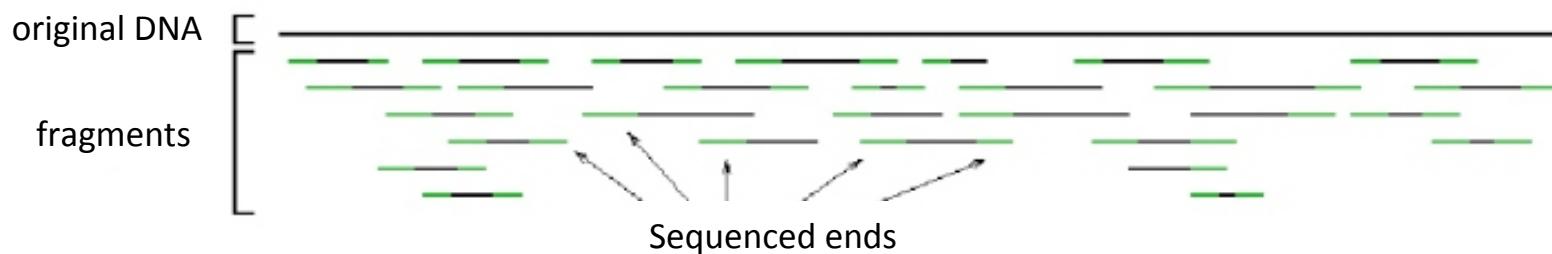
Mate pairs

Other read type

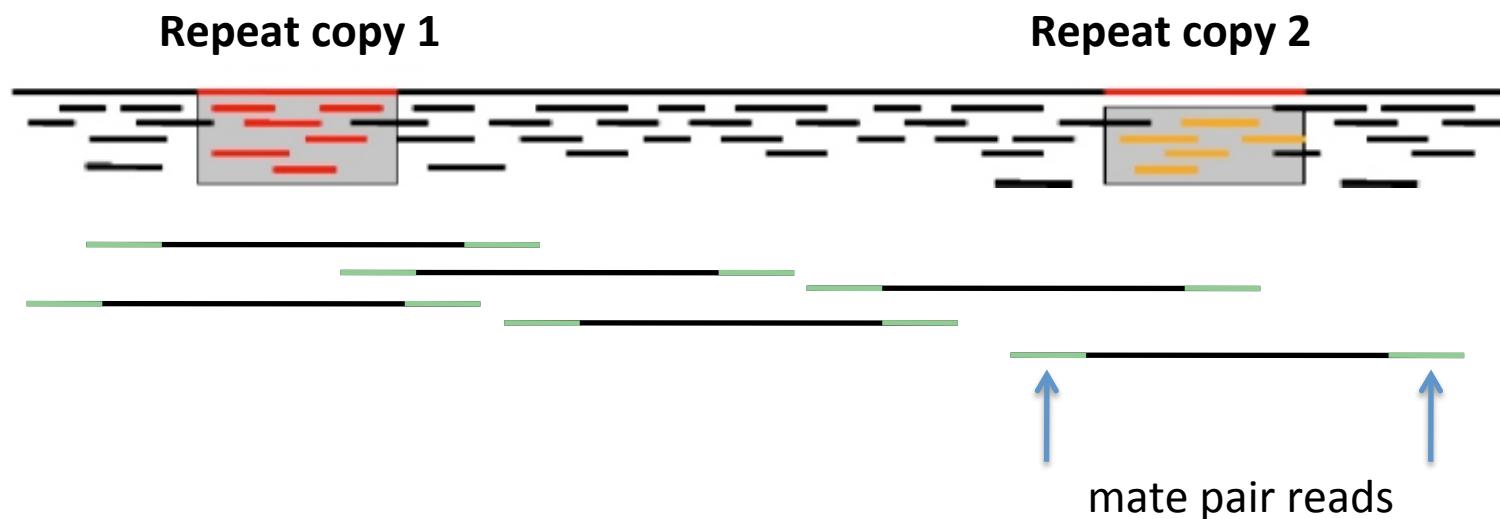


Mate pairs

Paired end reads → 100-500 bp insert

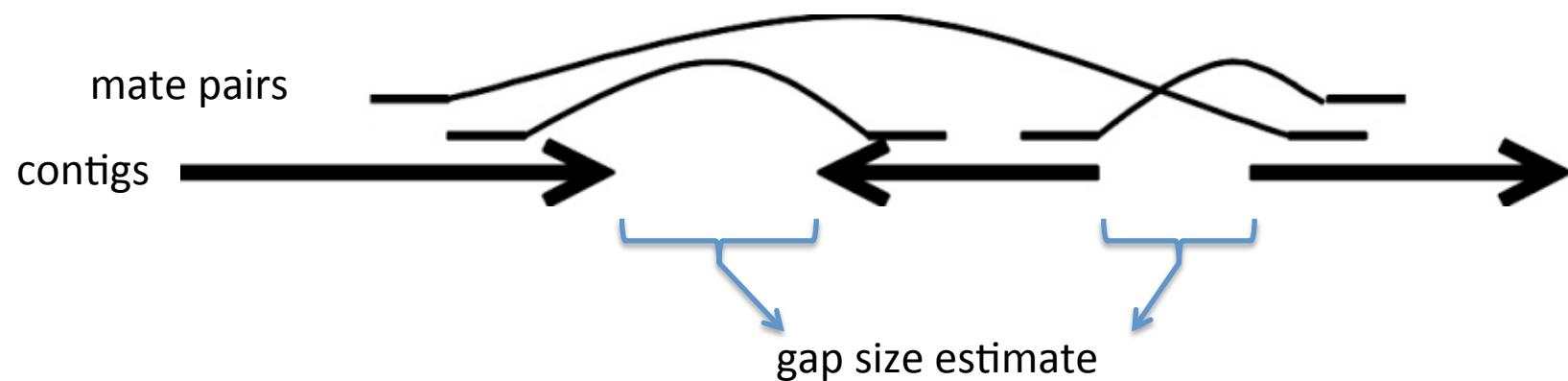


Mate pairs → 2-20 kb insert



Scaffolds

Ordered, oriented contigs

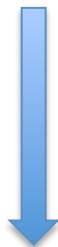


Hierarchical structure



Assembly

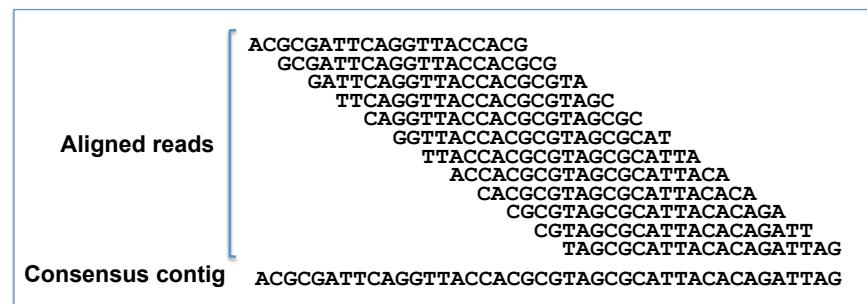
How to do this?



Algorithms

Algorithms

All are graph-based

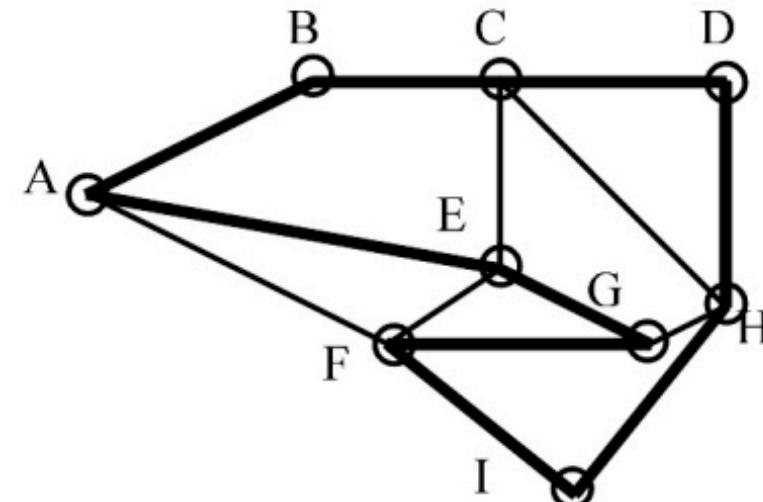
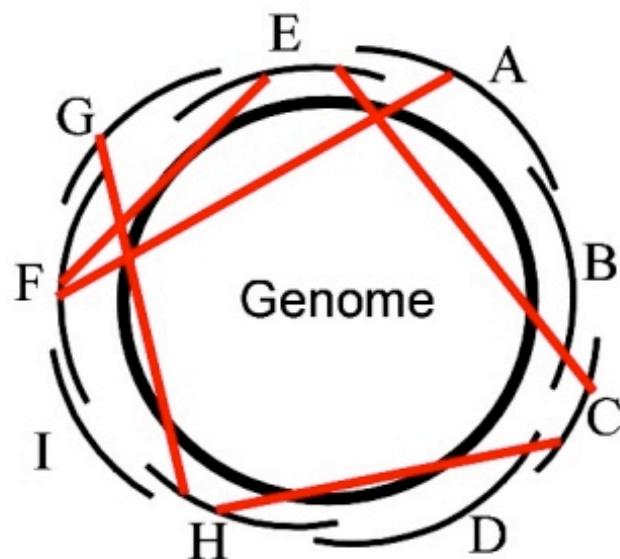


Graph-theory!

Algorithms

Hamiltonian path

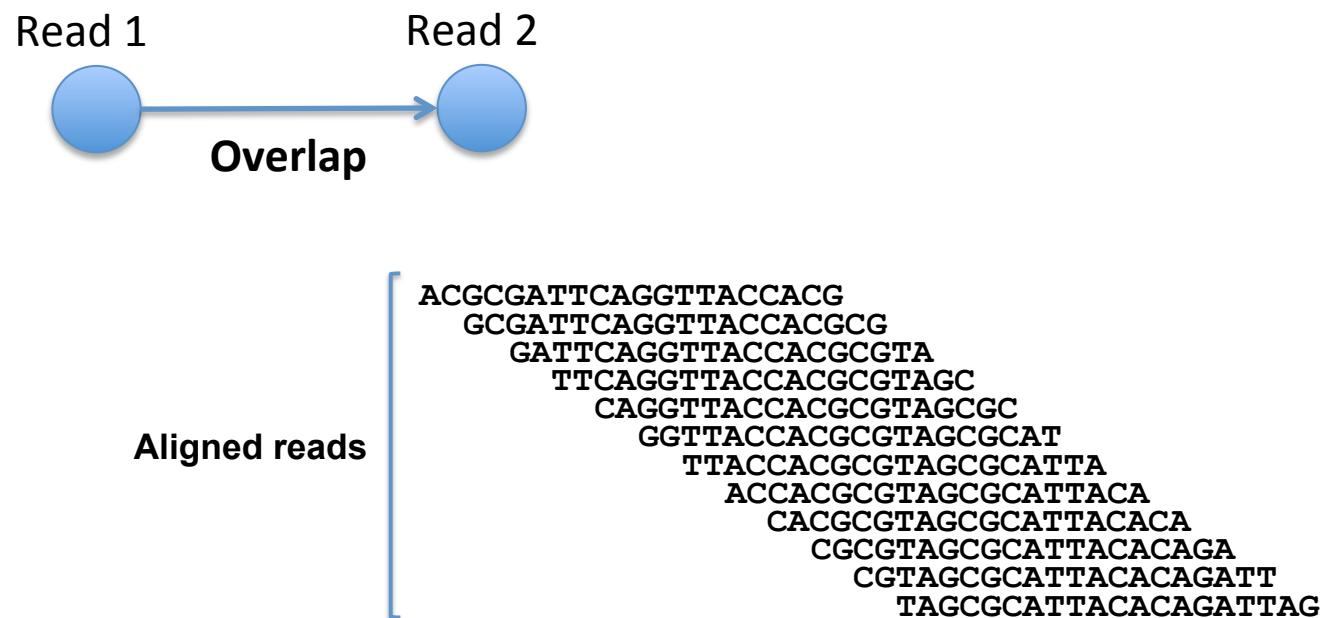
– a path that contains all the nodes



Algorithms

Overlap calculation (alignment)

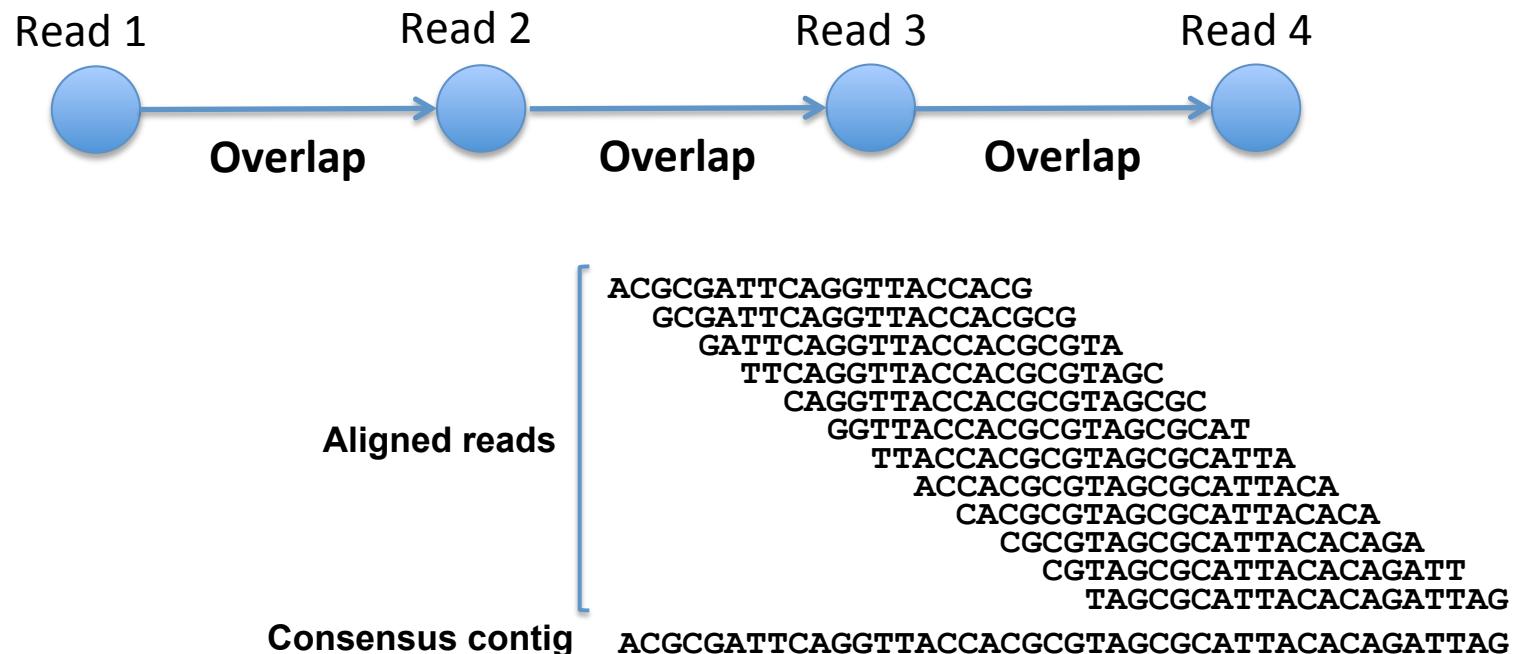
– computationally intensive



Algorithms

Path through the graph

→contig



Algorithms

Many flavors

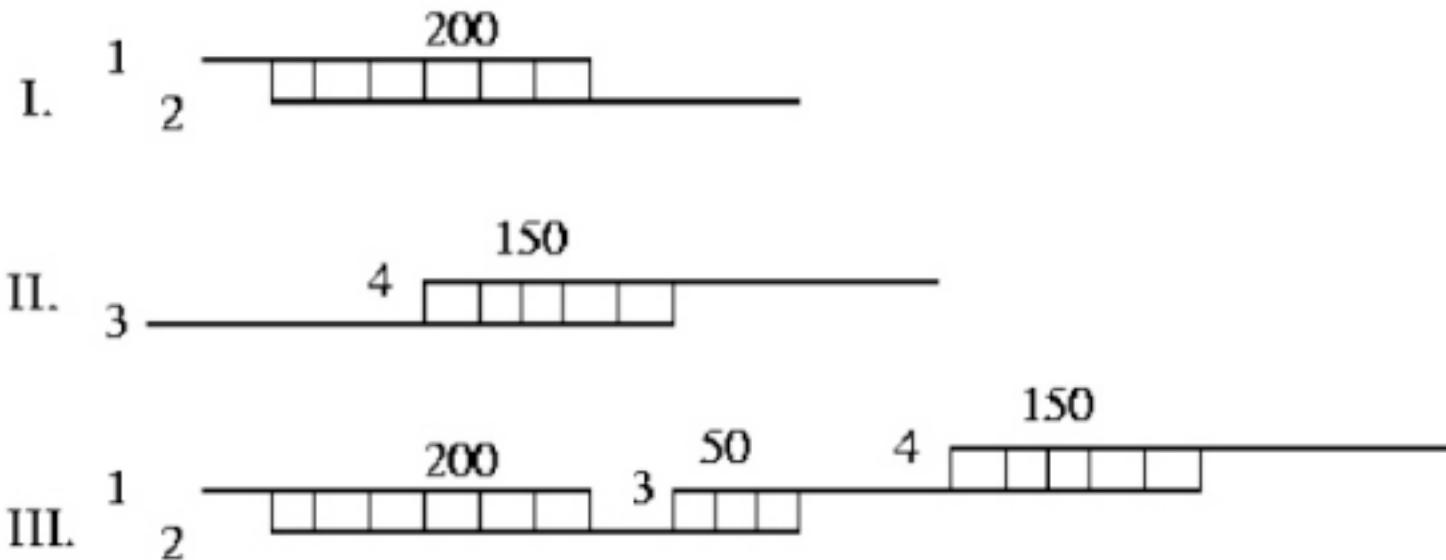


Abandoned
→ Greedy extension

Two most used
→ Overlap Layout Consensus
→ de Bruijn graph

Greedy extension

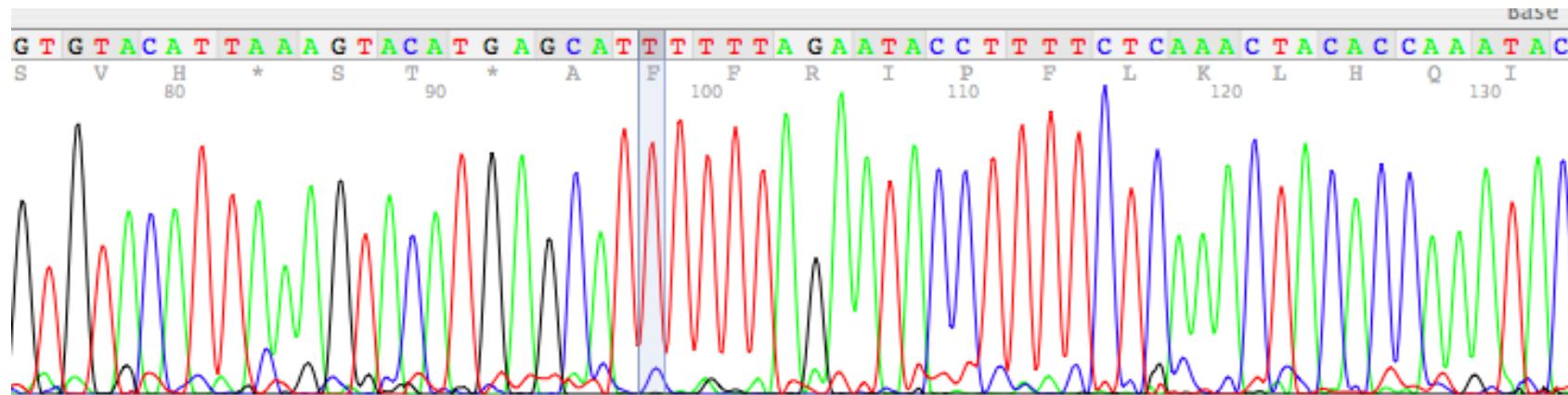
Oldest



Overlap-Layout-Consensus

Typical for Sanger-type reads

– also used by newbler from 454 Life Sciences



Overlap-Layout-Consensus

Steps

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

Overlap-Layout-Consensus

Overlap phase:

- K-mer seeds initiate overlap

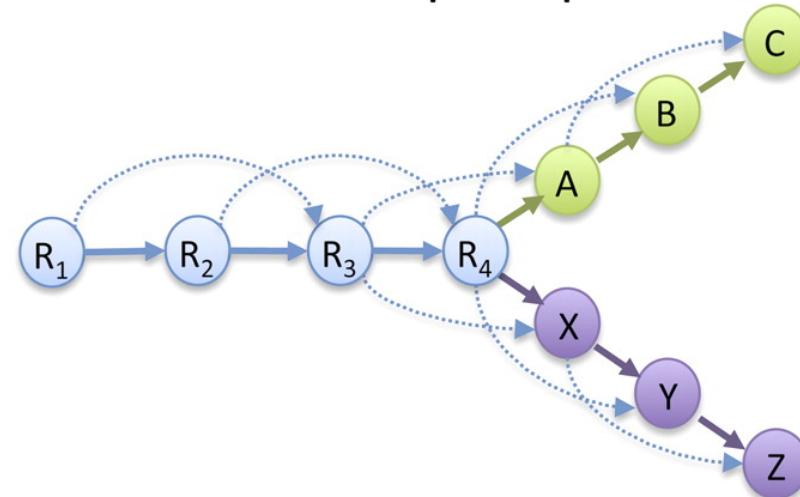


Overlap-Layout-Consensus

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:	CAAGTCG

B Overlap Graph

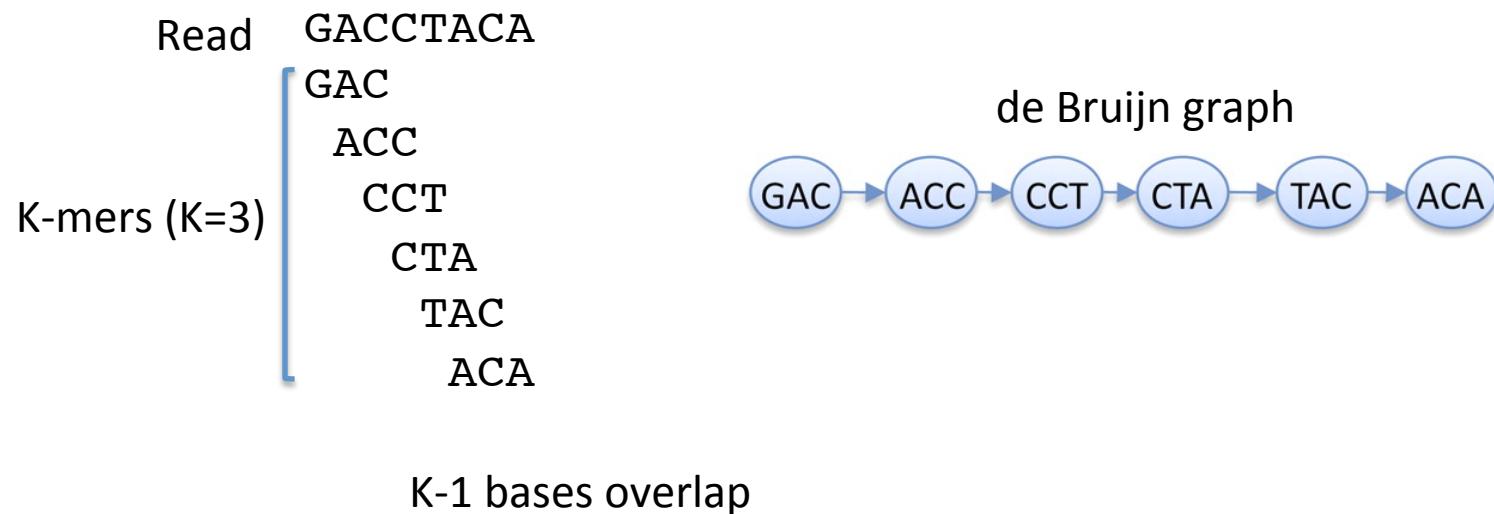




de Bruijn graphs

Developed outside of DNA-related work

- Best solution for very short reads ≤ 100 nt

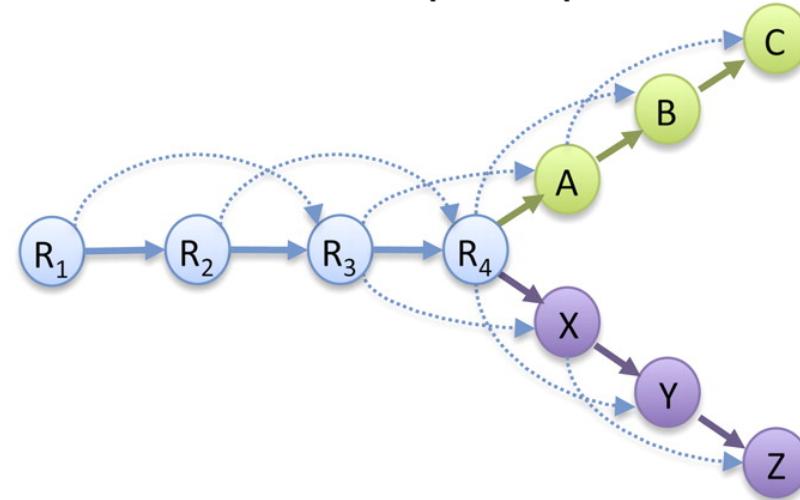


Graphs

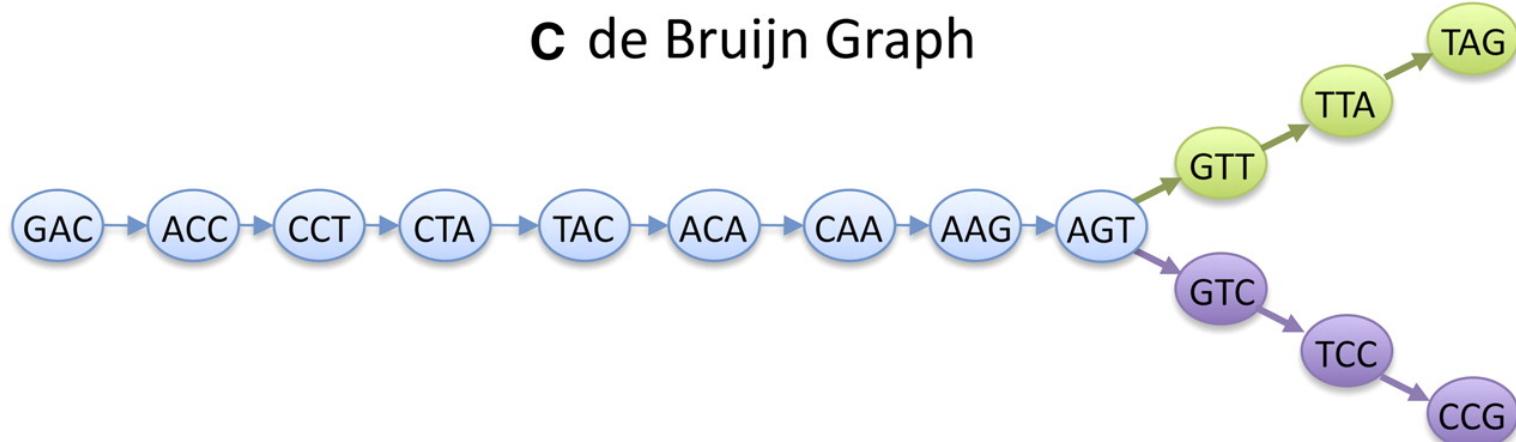
A Read Layout

$R_1:$	GACCTACA
$R_2:$	ACCTACAA
$R_3:$	CCTACAAG
$R_4:$	CTACAAGT
A:	TACAAGT T
B:	ACAAGT TA
C:	CAAGT TAG
X:	TACAAGTC
Y:	ACAAGTCC
Z:	CAAGTCCG

B Overlap Graph



C de Bruijn Graph



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

Graphs

Simplify the graph



Add scaffolding information



Read length matters

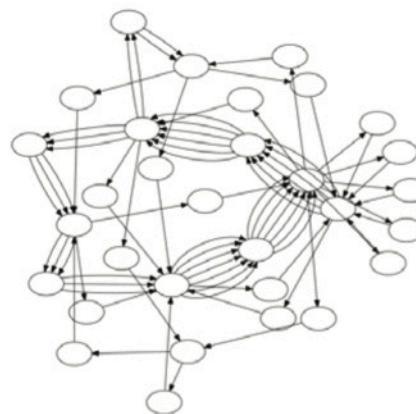
5.2 Mb circular genome, infinite error-free reads

(a)



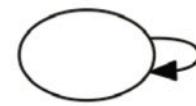
K=100
Contigs=98

(b)



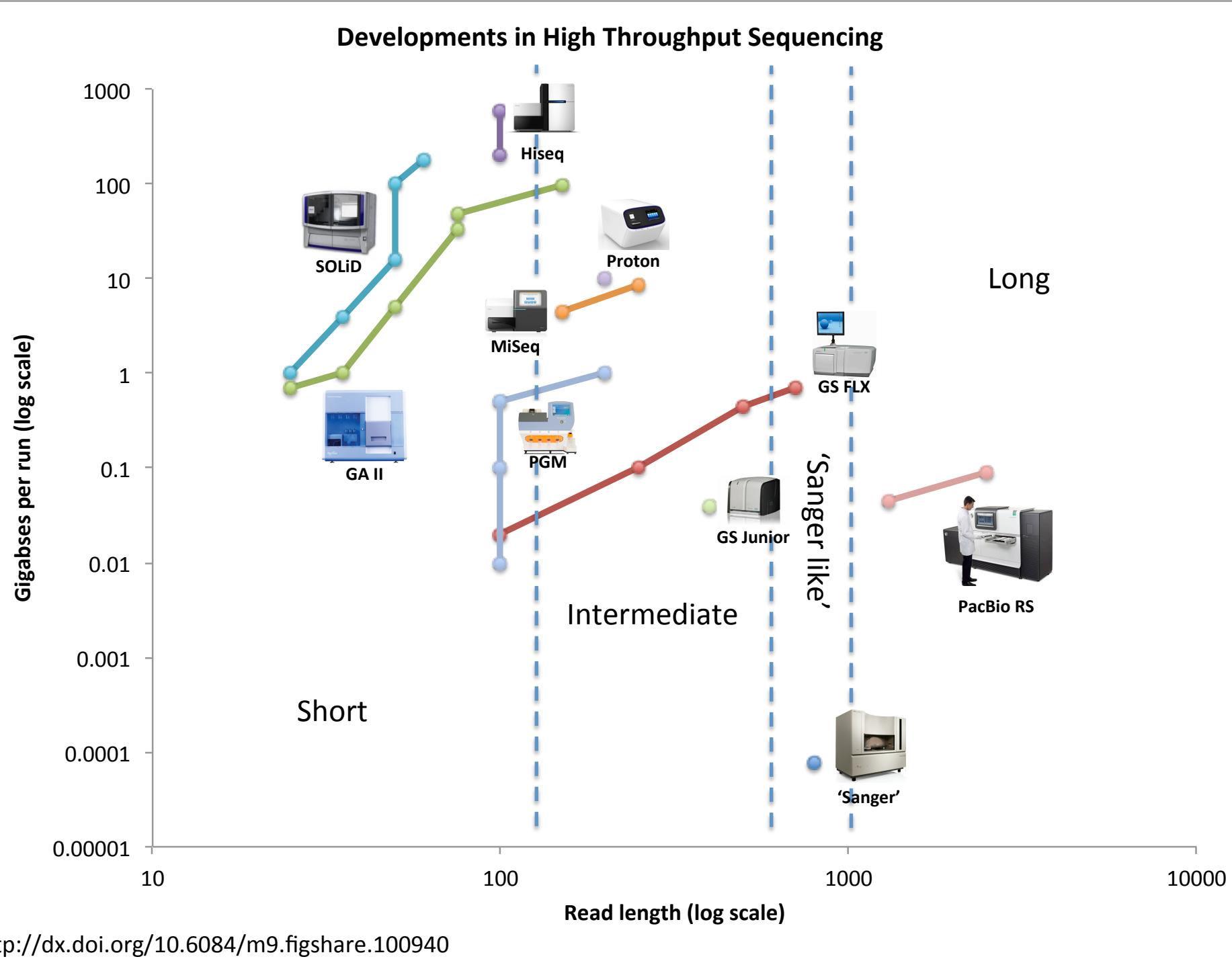
K=1,000
Contigs=31

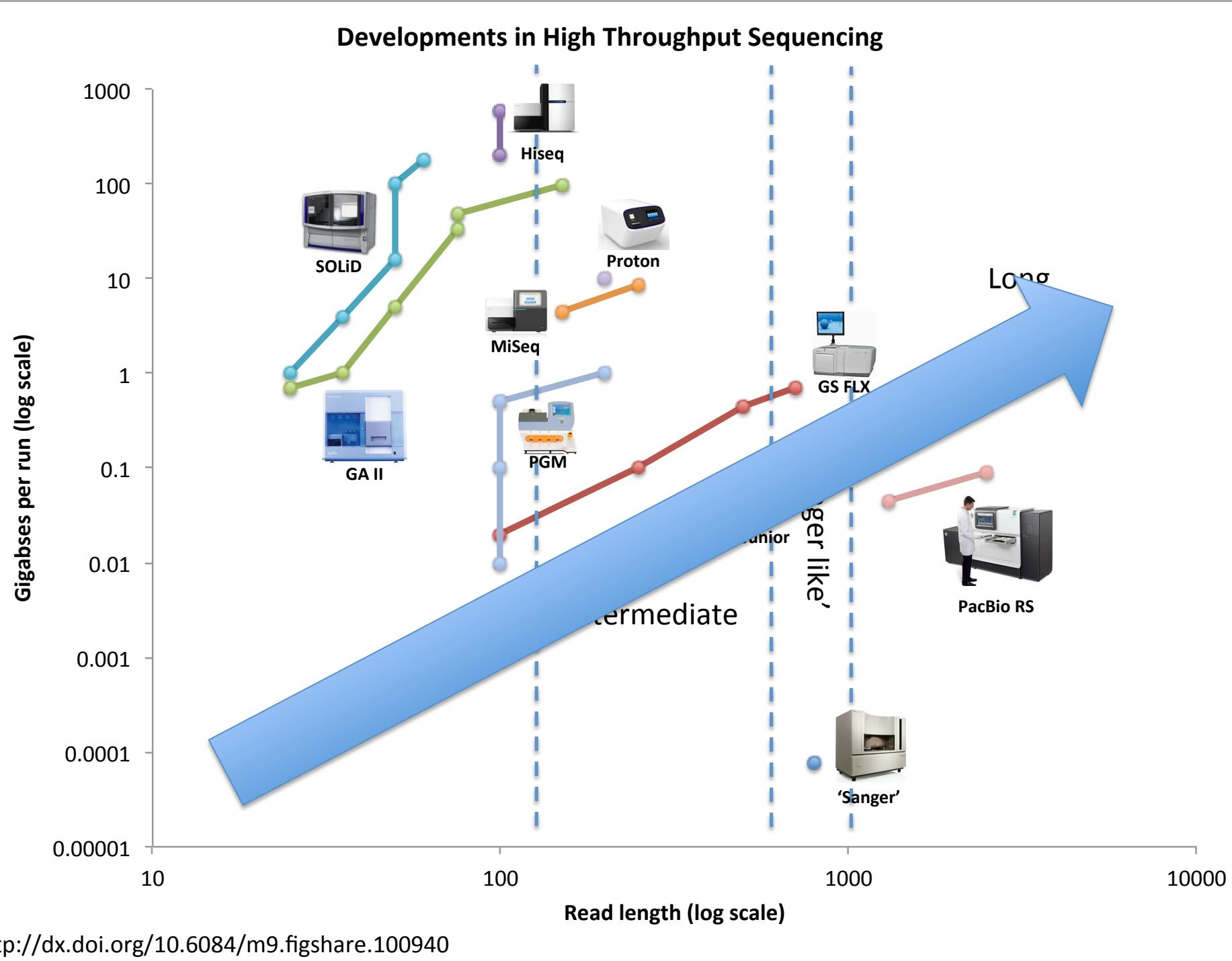
(c)



K=5,000
Contigs=1

Roberts et al (2013) doi:10.1186/gb-2013-14-6-405





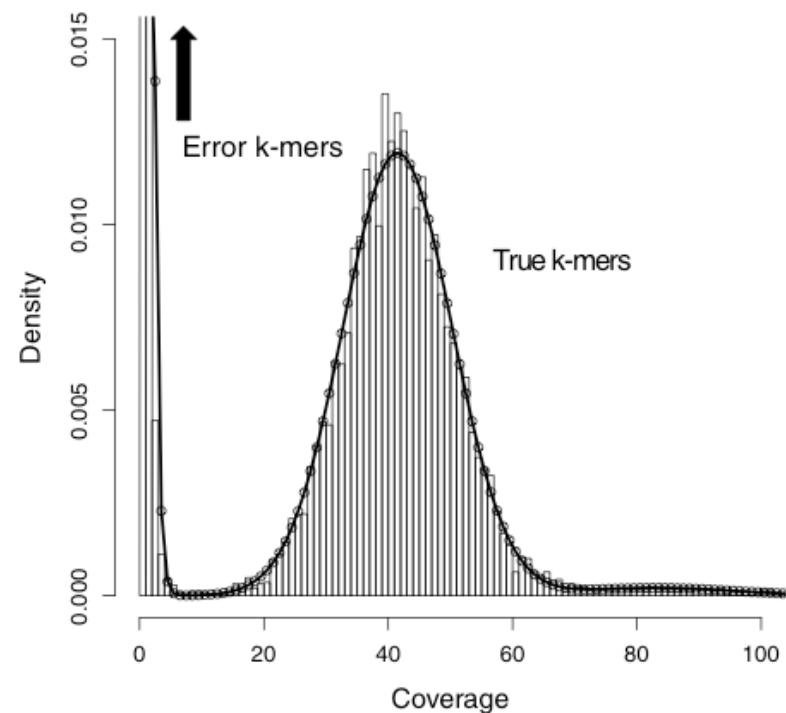
Quality matters

Sequencing errors

- add complexity to graph
- create new k-mers

Correction of errors

- k-mer frequency

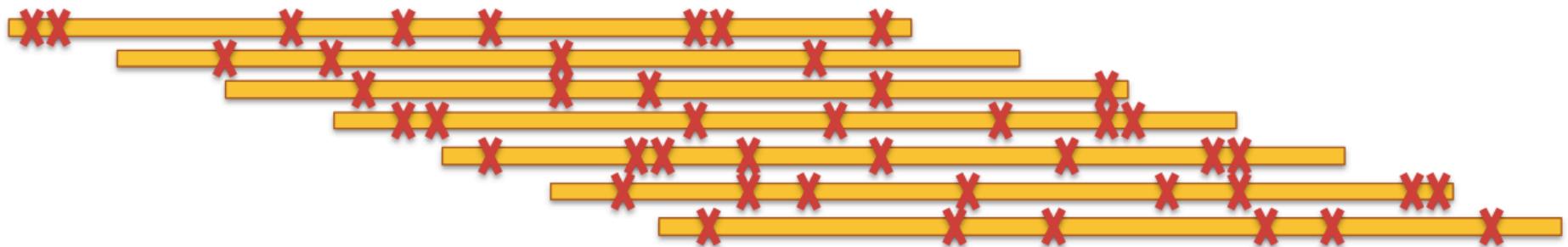


Kelley *et al.* *Genome Biology* 2010 **11**:R116

Quality matters

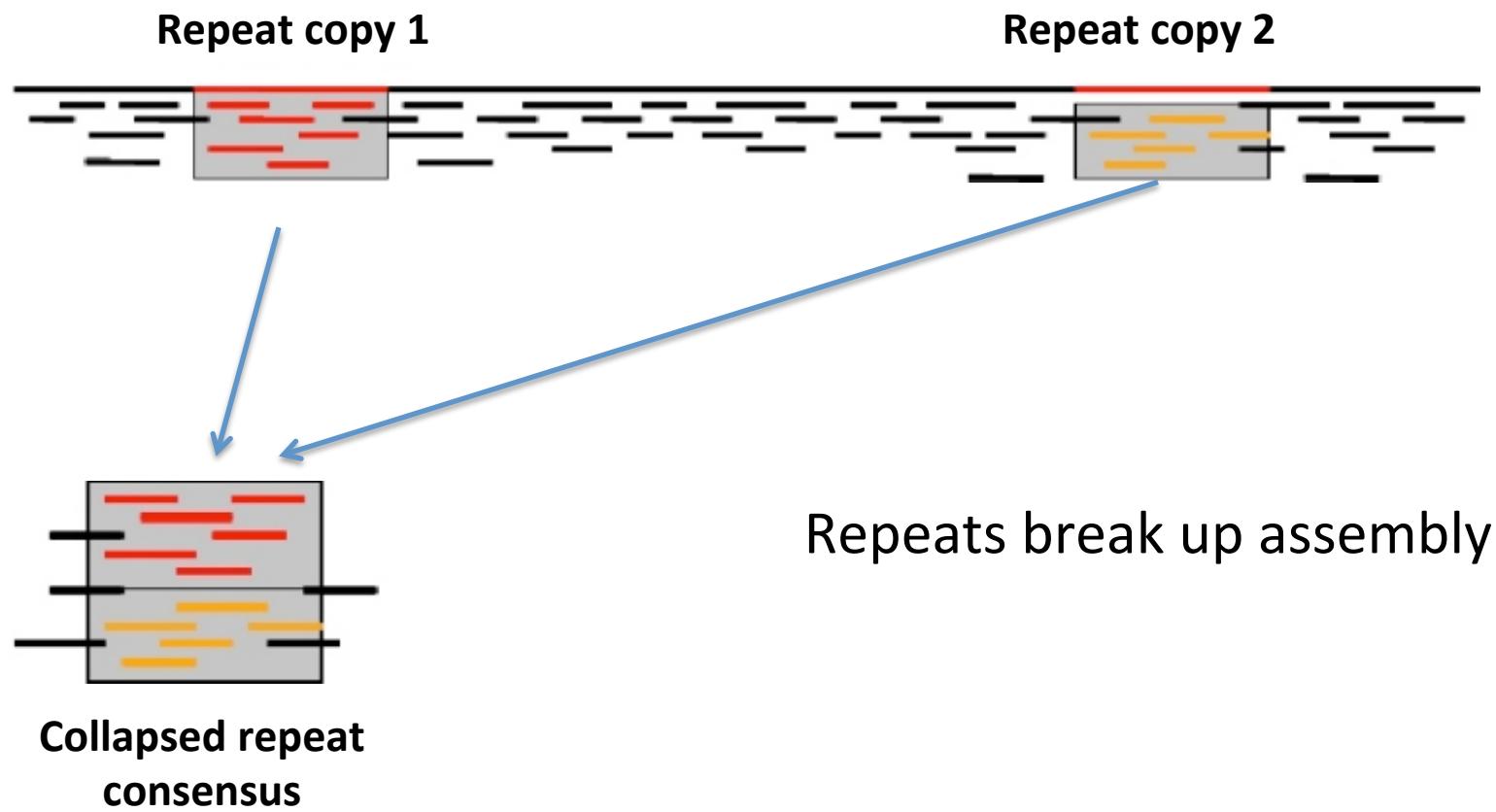


Too many errors → hard to find overlaps

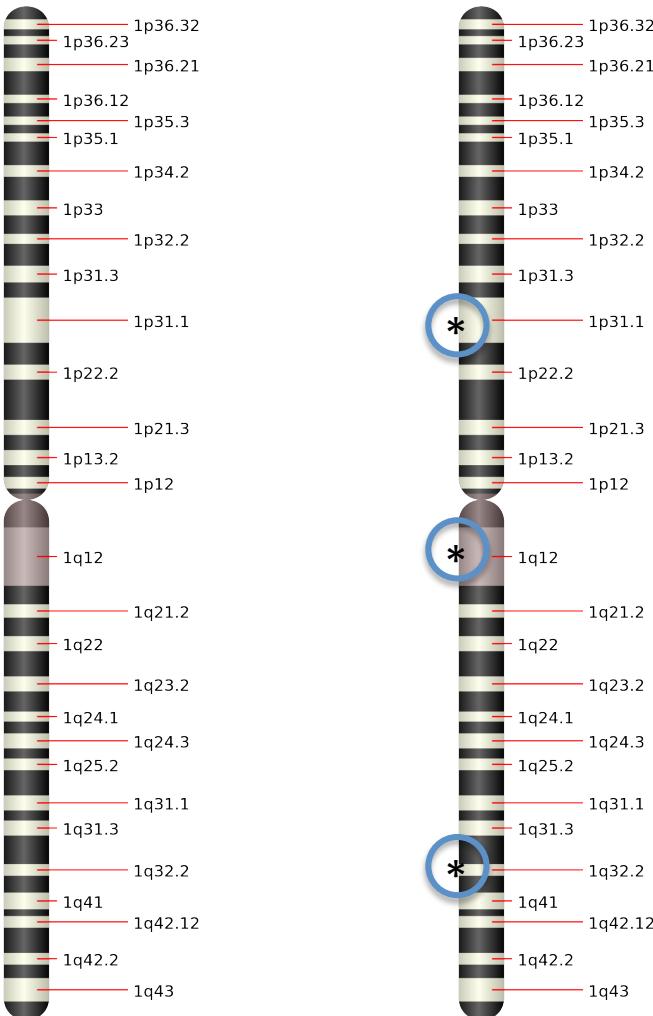


Why is genome assembly such
a difficult problem?

1) Repeats



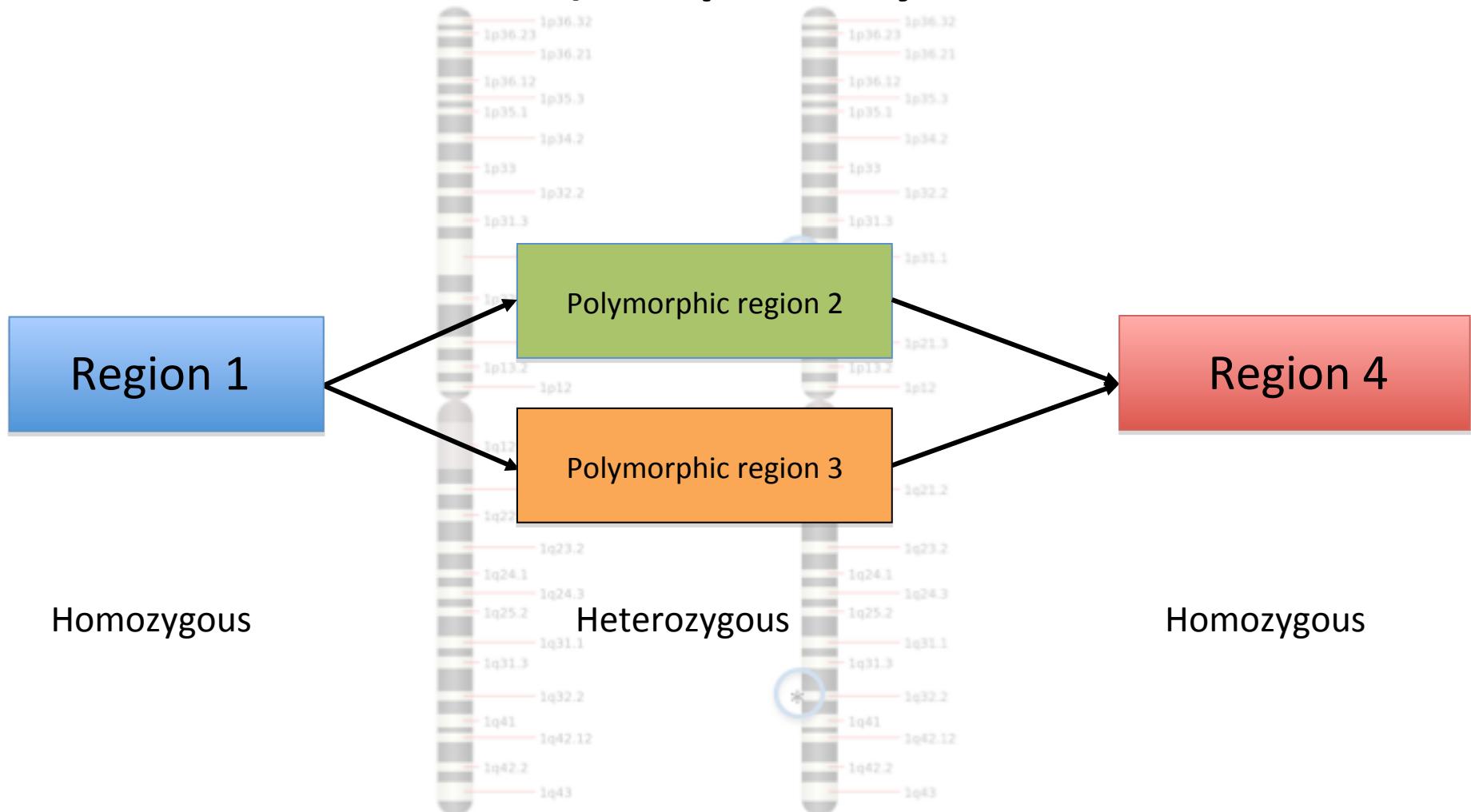
2) Diploidy



Differences
between sister
chromosomes

↓
'heterozygosity'

2) Diploidy



2) Diploidy

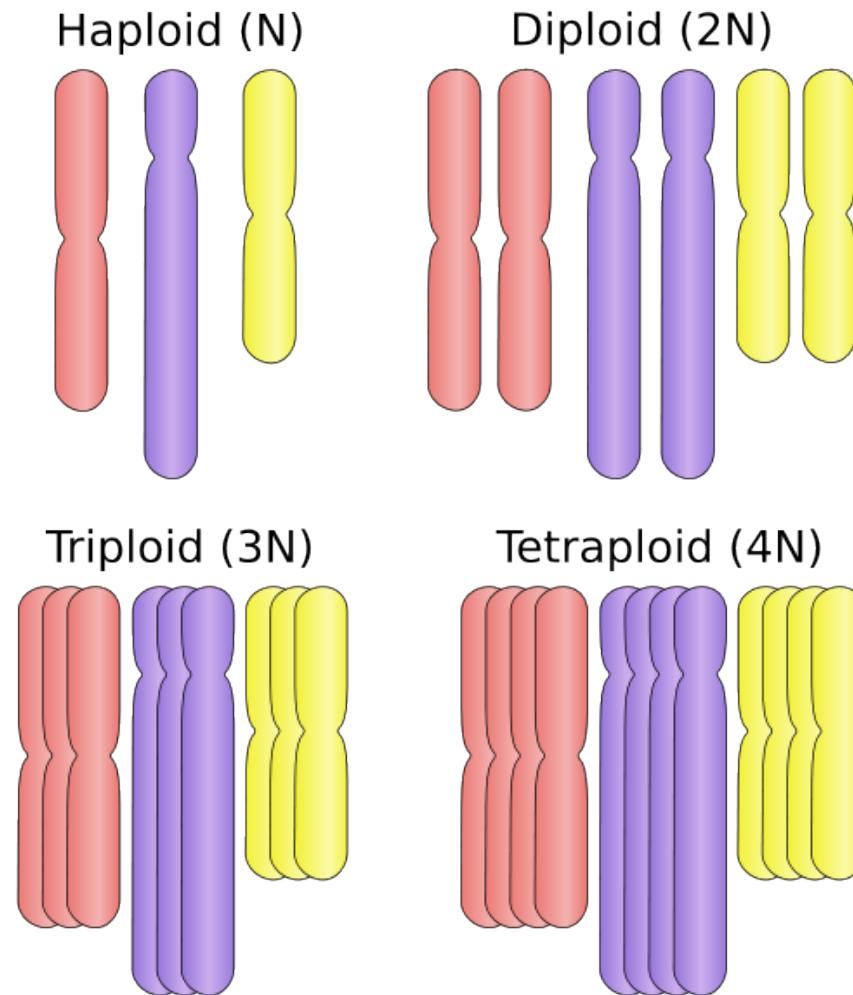


HETEROZYGOATS

Just allele uneven.

<http://www.astraeans.com/borderwars/wp-content/uploads/2012/04/heterozygoats.jpg>
and many other sites

3) Polyploidy



4) Many programs to choose from



Shameless self-promotion



A screenshot of a Twitter profile card. At the top is a small circular portrait of a man with blonde hair and a beard, smiling. Below the portrait, the name "Lex Nederbragt" is displayed in a large, bold, white font. Underneath the name is the handle "@lexnederbragt". A bio follows, starting with "Husband, father of two, biologist, bioinformatician, researcher, Dutchman." It continues with "Views expressed here are my own. An RT does not necessarily mean endorsement." At the bottom of the card, it says "Oslo · flavors.me/flxlex".

7,638
TWEETS

220
FOLLOWING

1,052
FOLLOWERS



Edit profile



A screenshot of a blog header. The URL "flxlexblog.wordpress.com" is visible in the address bar. The main title "In between lines of code" is prominently displayed in large, bold, dark letters. Below the main title, a subtitle in a smaller, italicized font reads "Biology, sequencing, bioinformatics and more".