

Practical Assembly

Michael Schatz

Sept 8, 2025

Lecture 4: Applied Comparative Genomics



Assignment I

The screenshot shows a GitHub repository interface for 'assignment1'. The left sidebar displays a file tree with various genome size files (e.g., TAIR10.chrom.sizes, ce10.chrom.sizes) and other files like README.md and LICENSE. The main content area is titled 'Assignment 1: Chromosome Structures' and includes assignment details: 'Assignment Date: Wednesday, August 27, 2025' and 'Due Date: Wednesday, Sept. 3, 2025 @ 11:59pm'. It also contains an 'Assignment Overview' section describing the task of profiling genome structures and studying human chromosome 22. A 'Question 1: Chromosome structures [10 pts]' section lists eight species with links to their genome size files, followed by instructions to create a table with specific information per species. A 'Question 2: Coverage simulator [20 pts]' section is partially visible at the bottom.

Assignment Date: Wednesday, August 27, 2025
Due Date: Wednesday, Sept. 3, 2025 @ 11:59pm

Assignment Overview

In this assignment you will profile the overall structure of the genomes of several important species and then study human chromosome 22 in more detail. As a reminder, any questions about the assignment should be posted to [Piazza](#).

Question 1: Chromosome structures [10 pts]

Download the chromosome size files for the following genomes (Note these have been preprocessed to only include main chromosomes):

1. [E. coli \(Escherichia coli K12\)](#) - One of the most commonly studied bacteria [\[info\]](#)
2. [Yeast \(Saccharomyces cerevisiae, sacCer3\)](#) - an important eukaryotic model species, also good for bread and beer [\[info\]](#)
3. [Worm \(Caenorhabditis elegans, ce10\)](#) - One of the most important animal model species [\[info\]](#)
4. [Fruit Fly \(Drosophila melanogaster, dm6\)](#) - One of the most important model species for genetics [\[info\]](#)
5. [Arabidopsis thaliana \(TAIR10\)](#) - An important plant model species [\[info\]](#)
6. [Tomato \(Solanum lycopersicum v4.00\)](#) - One of the most important food crops [\[info\]](#)
7. [Human \(hg38\)](#) - us :) [\[info\]](#)
8. [Wheat \(Triticum aestivum, IWGSC\)](#) - The food crop which takes up the largest land area [\[info\]](#)

Using these files, make a table with the following information per species:

- Question 1.1. Total genome size
- Question 1.2. Number of chromosomes
- Question 1.3. Largest chromosome size and name
- Question 1.4. Smallest chromosome size and name
- Question 1.5. Mean chromosome length

Question 2: Coverage simulator [20 pts]

<https://github.com/schatzlab/appliedgenomics2025/tree/main/assignments/assignment1>
Due end of day on Wednesday Sept 3 (right before midnight)

Part 0: Recap

K-mers and K-mer counting

GATTACATACACATTGGATG

GAT : 2 CAT : 2 ATG : 1 TGG : 1

ACA : 3 CAC : 1 TTA : 1 TAC : 2

ATT : 2 TTG : 1 ATA : 1 GGA : 1

1 : 7 (ATG, TGG, ...)

2 : 4 (GAT, CAT, ATT, TAC)

3 : 1 (ACA)

See HW1

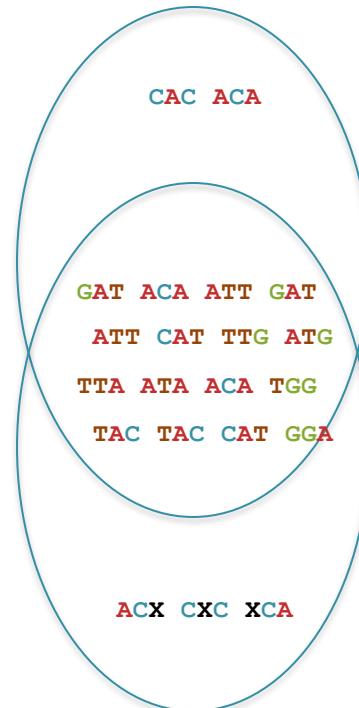
GATTACATAACACATTGGATG
 GAT ACA ACA ATT GAT
 ATT CAT CAC TTG ATG
 TTA ATA ACA TGG
 TAC TAC CAT GGA

GATTACATAACXCATGGATG
 GAT ACA ACX ATT GAT
 ATT CAT CXC TTG ATG
 TTA ATA XCA TGG
 TAC TAC CAT GGA

$$J(A, B) = \frac{|A \cap B|}{|A \cup B|} = \frac{16}{21} = 76\%$$

$$\text{ANI} \approx 1 + \frac{1}{k} \ln \left(\frac{2J}{1+J} \right)$$

ANI = 95.1%



ERRORS!
Thank you to
Piazza!

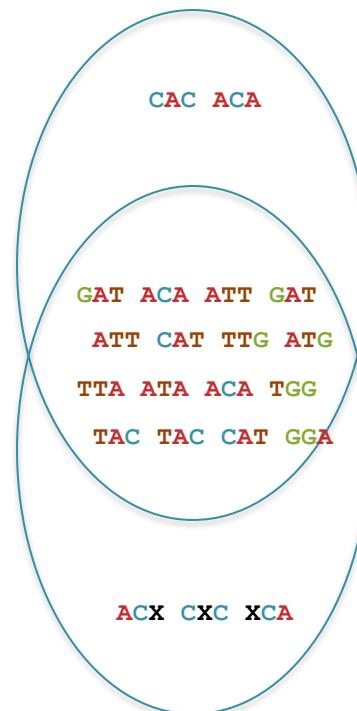
GATTACATAACACATTGGATG
GAT ACA ACA ATT GAT
ATT CAT CAC TTG ATG
TTA ATA ACA TGG
TAC TAC CAT GGA

GATTACATAACXCATGGATG
GAT ACA ACX ATT GAT
ATT CAT CXC TTG ATG
TTA ATA XCA TGG
TAC TAC CAT GGA

$$J(A, B) = \frac{|A \cap B|}{|A \cup B|} = \frac{16}{21} = 76\%$$

$$\text{ANI} \approx 1 + \frac{1}{k} \ln \left(\frac{2J}{1+J} \right)$$

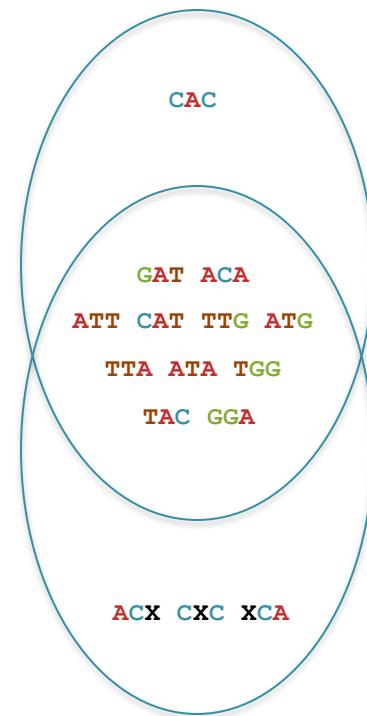
ANI = 95.1%



GATTACATAACACATTGGATG
 GAT ACA ACA ATT GAT
 ATT CAT CAC TTG ATG
 TTA ATA ACA TGG
 TAC TAC CAT GGA

GATTACATAACXCATTTGGATG
 GAT ACA ACX ATT GAT
 ATT CAT CXC TTG ATG
 TTA ATA XCA TGG
 TAC TAC CAT GGA

Corrected analysis



$$J(A, B) = \frac{|A \cap B|}{|A \cup B|} = \frac{11}{15} = 73.3\%$$

$$\text{ANI} \approx 1 + \frac{1}{k} \ln \left(\frac{2J}{1+J} \right)$$

$$\text{ANI} = 94.4\%$$

de Bruijn Graph Construction

- $G_k = (V, E)$
 - V = Length- k sub-fragments
 - E = Directed edges between consecutive sub-fragments
 - Sub-fragments overlap by $k-1$ words

Fragments $|f|=5$

It was the best of

was the best of times

Sub-fragment $k=4$

It was the best

was the best of

was the best of

the best of times

Directed edges (overlap by $k-1$)

It was the best

was the best of

the best of times

- Overlaps between fragments are implicitly computed

How to pronounce:

https://forvo.com/word/de_bruijn/

de Bruijn, 1946
Idury et al., 1995
Pevzner et al., 2001

de Bruijn Graph Assembly

It was the best

was the best of

the best of times,

best of times, it

of times, it was

times, it was the

it was the worst

was the worst of

the worst of times,

worst of times, it

After graph construction,
try to simplify the graph as
much as possible

it was the age

was the age of

the age of foolishness

the age of wisdom,

age of wisdom, it

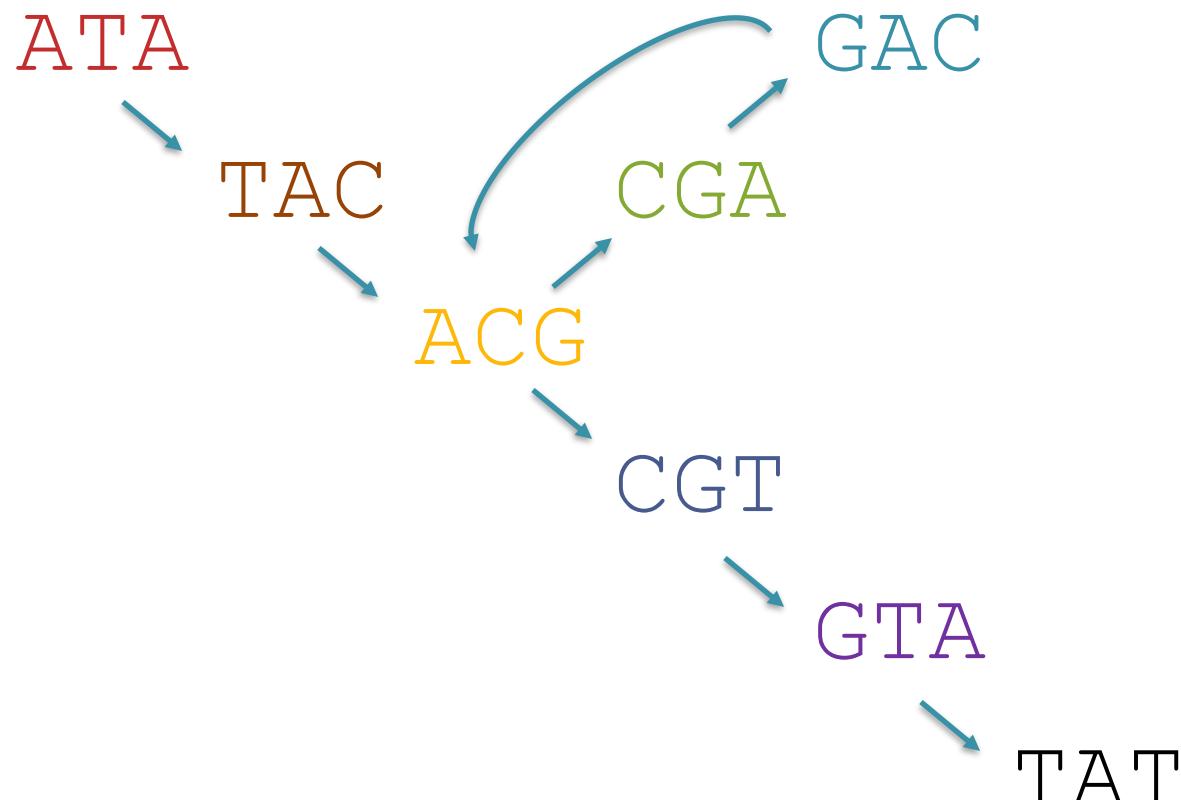
of wisdom, it was

wisdom, it was the

Pop Quiz 2

Assemble these reads using a de Bruijn graph approach (k=3):

~~ACGA~~
~~ACGT~~
~~ATAC~~
~~CGAC~~
~~CGTA~~
~~GACG~~
~~GTAT~~
~~TACG~~



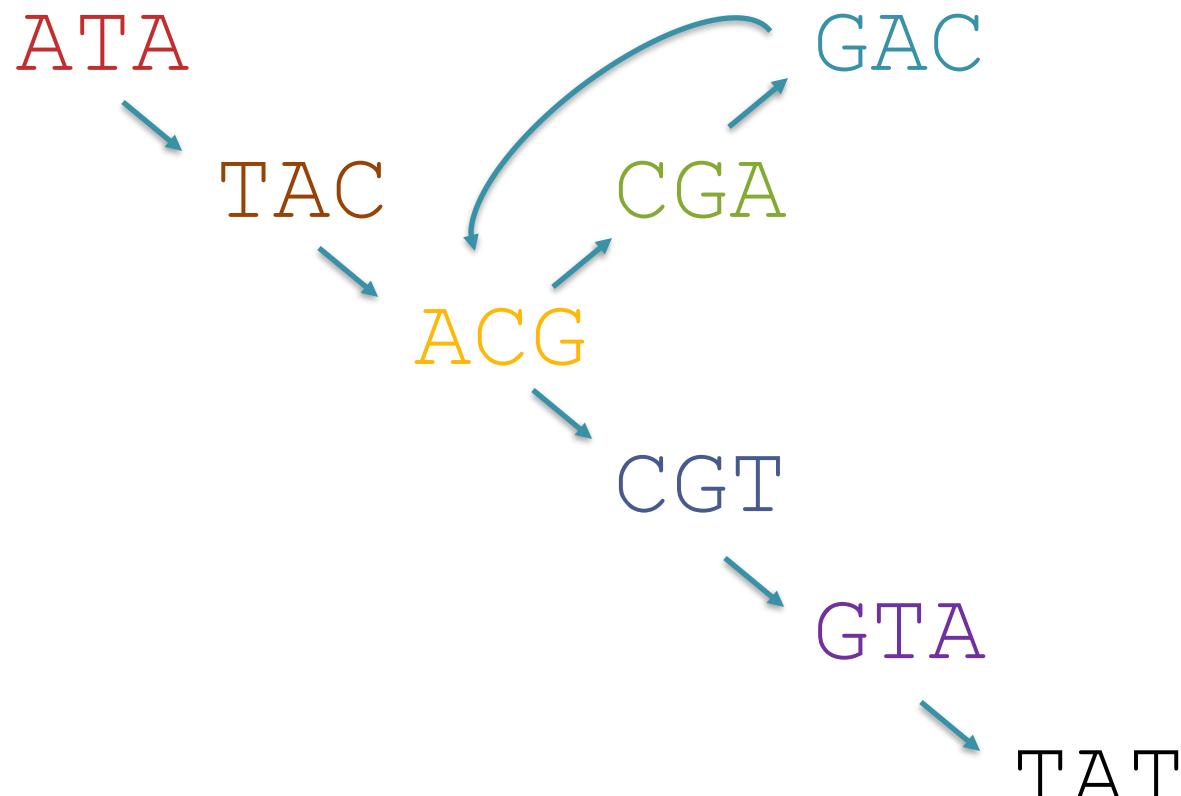
Whats another possible genome?

ATACGACGTAT

Pop Quiz 2

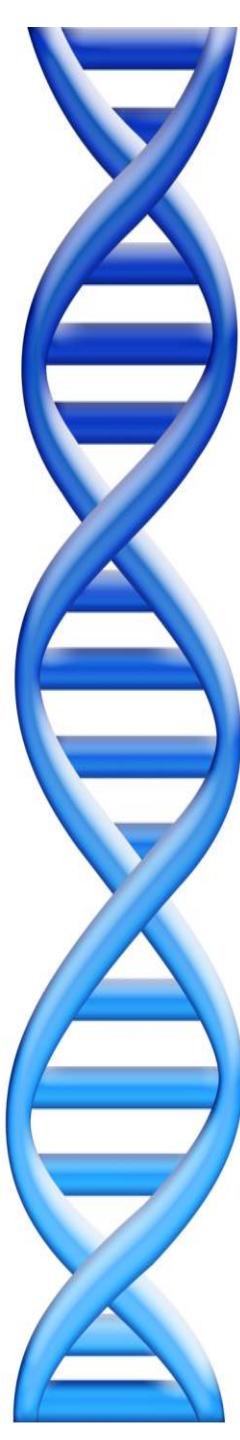
Assemble these reads using a de Bruijn graph approach ($k=3$):

~~ACGA~~
~~ACGT~~
~~ATAC~~
~~CGAC~~
~~CGTA~~
~~GACG~~
~~GTAT~~
~~TACG~~



Should we add the edge $TAT \rightarrow ATA$?

ATACGACGTAT



Outline

1. Assembly theory

- Assembly by analogy

2. Practical Issues

- Coverage, read length, errors, and repeats

3. Whole Genome Alignment

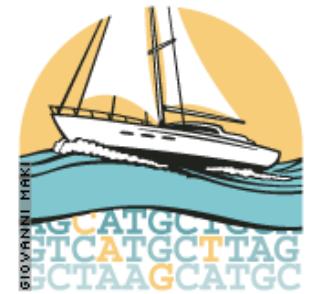
- MUMmer recommended

Assembly Applications

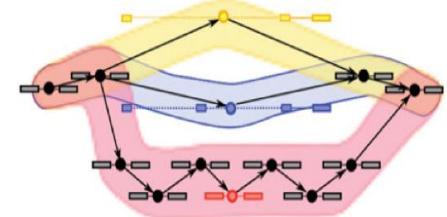
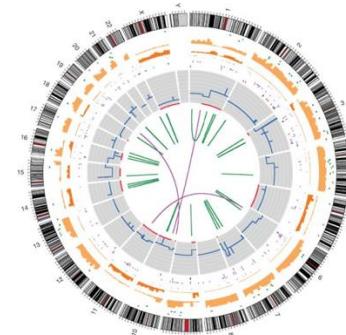
- Novel genomes



- Metagenomes



- Sequencing assays
 - Structural variations
 - Transcript assembly
 - ...



Why are genomes hard to assemble?

1. ***Biological:***

- (Very) High ploidy, heterozygosity, repeat content

2. ***Sequencing:***

- (Very) large genomes, imperfect sequencing

3. ***Computational:***

- (Very) Large genomes, complex structure

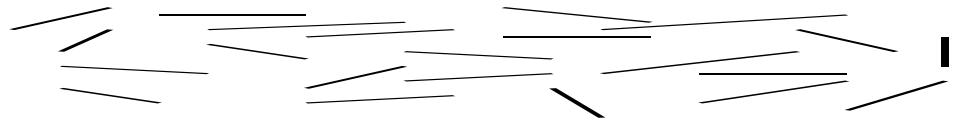
4. ***Accuracy:***

- (Very) Hard to assess correctness



Assembling a Genome

1. Shear & Sequence DNA



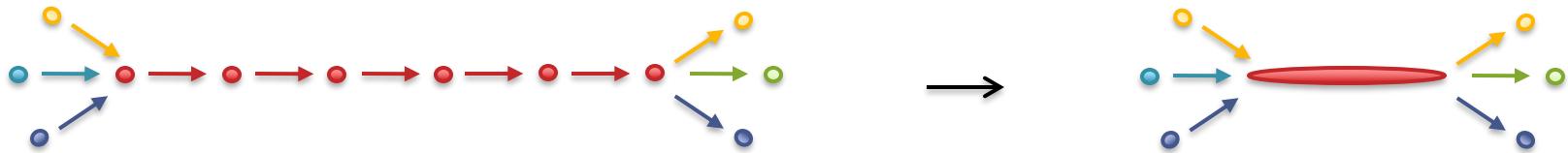
2. Construct assembly graph from reads (de Bruijn / overlap graph)

...AGCCTAG**GGATGCGCGACACGT**

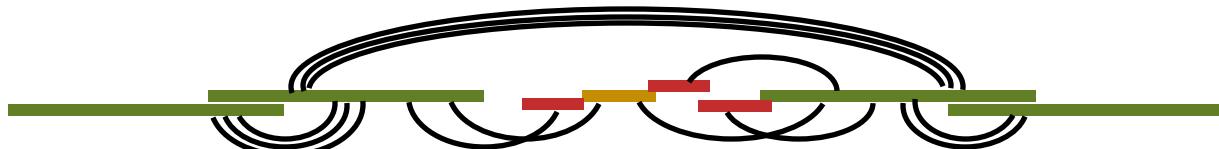
GGATGCGCGACACGTCGCATATCCGGTTGGT**CAACCTCGGACGGAC**

CAACCTCGGACGGACCTCAGCGAA...

3. Simplify assembly graph

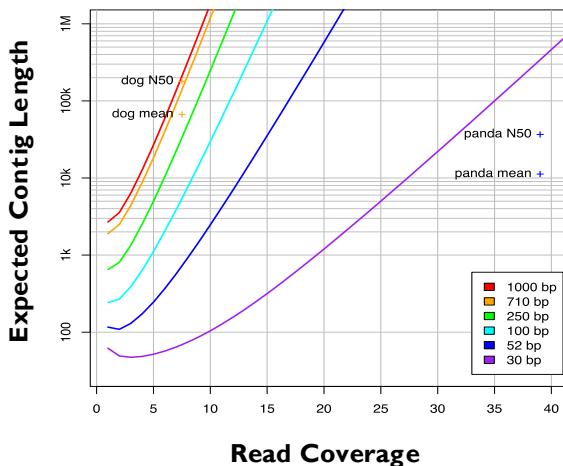


4. Detangle graph with long reads, mates, and other links



Ingredients for a good assembly

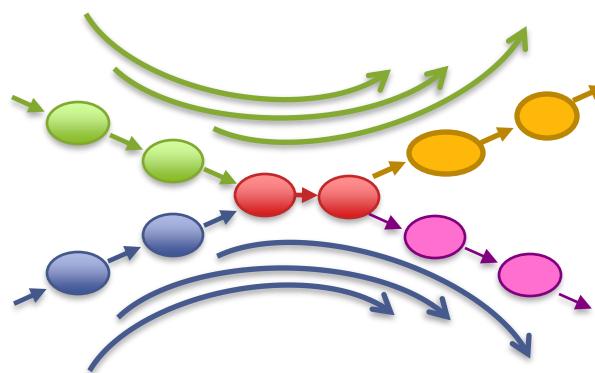
Coverage



High coverage is required

- Oversample the genome to ensure every base is sequenced with long overlaps between reads
- Biased coverage will also fragment assembly

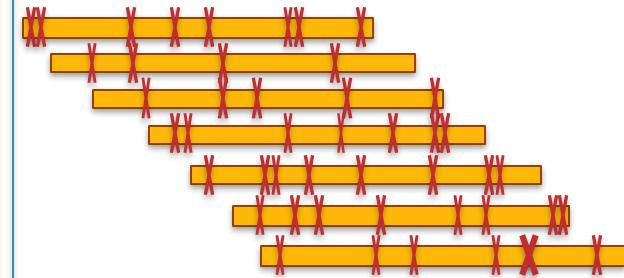
Read Length



Reads & mates must be longer than the repeats

- Short reads will have **false overlaps** forming hairball assembly graphs
- With long enough reads, assemble entire chromosomes into contigs

Quality



Errors obscure overlaps

- Reads are assembled by finding kmers shared in pair of reads
- High error rate requires very short seeds, increasing complexity and forming assembly hairballs

Current challenges in *de novo* plant genome sequencing and assembly

Schatz MC, Witkowski, McCombie, WR (2012) *Genome Biology*. 12:243

Coverage Statistics

$$\text{sequencing_coverage} = \frac{\text{total_bases_sequenced}}{\text{genome_size}}$$

$$\text{genome_size} = \frac{\text{total_bases_sequenced}}{\text{sequencing_coverage}}$$

$$\text{genome_size} = \frac{100\text{Gb}}{50x} = 2\text{Gb}$$

But how can you figure out
the coverage without a genome?

K-mer counting

Kmer-ize

Read 1: GATTACA => GAT, ATT, TTA, TAC, ACA
Read 2: TACAGAG => TAC, ACA, CAG, AGA, GAG
Read 3: TTACAGA => TTA, TAC, ACA, CAG, AGA

list

GAT	ACA	ACA : 3
ATT	ACA	
TTA	ACA	
TAC	AGA	AGA : 2
ACA	AGA	
TAC	ATT	ATT : 1
ACA	CAG	CAG : 2
CAG	CAG	
AGA	GAG	GAG : 1
GAG	GAT	GAT : 1
TTA	TAC	TAC : 3
TAC	TAC	
ACA	TAC	
CAG	TTA	TTA : 2
AGA	TTA	

3 kmers occur 1x
3 kmers occur 2x
2 kmers occur 3x

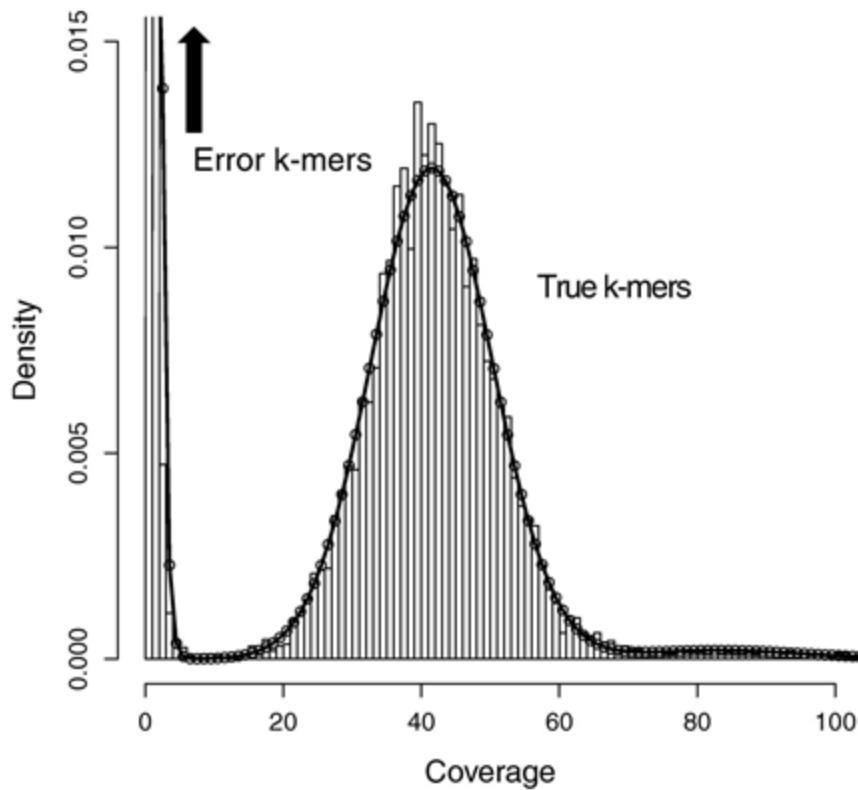
tally

sort count

From read k-mers alone, can learn something about how frequently different sequences occur (aka coverage)

Fast to compute even over huge datasets

K-mer counting in real genomes



- The tally of k-mer counts in real genomes reveals the coverage distribution.
- Here we sequenced 120Gb of reads from a female human (haploid human genome size is 3Gb), and indeed we see a clear peak centered at 40x coverage
- There are also many kmers that only occur <5 times. These are from errors in the reads
- There are also kmers that occur many times (>>70 times). These are repeats in the genome

K-mer counting in heterozygous genomes

Sequencing read
from homologous
chromosome 1A



Sequencing read
from homologous
chromosome 1B



K-mer counting in heterozygous genomes



Sequencing read
from homologous
chromosome 1A



Sequencing read
from homologous
chromosome 1B



K-mer counting in heterozygous genomes



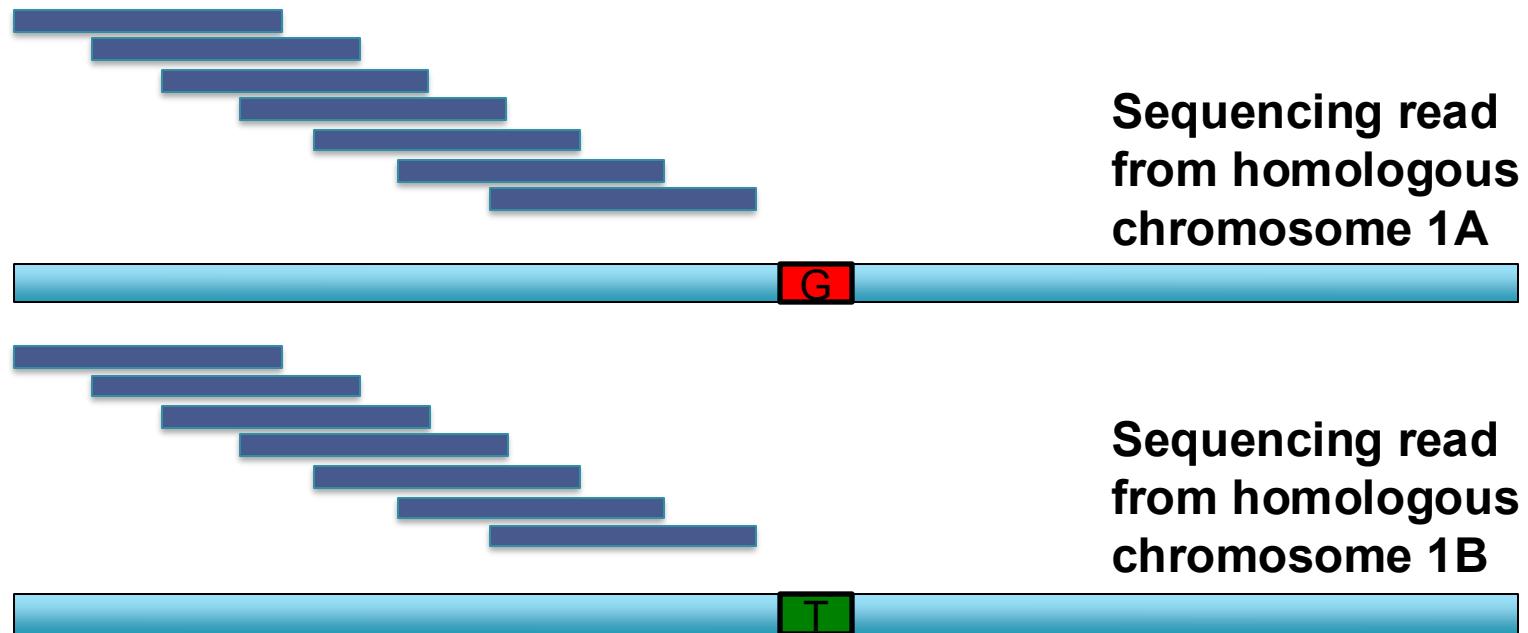
Sequencing read
from homologous
chromosome 1A



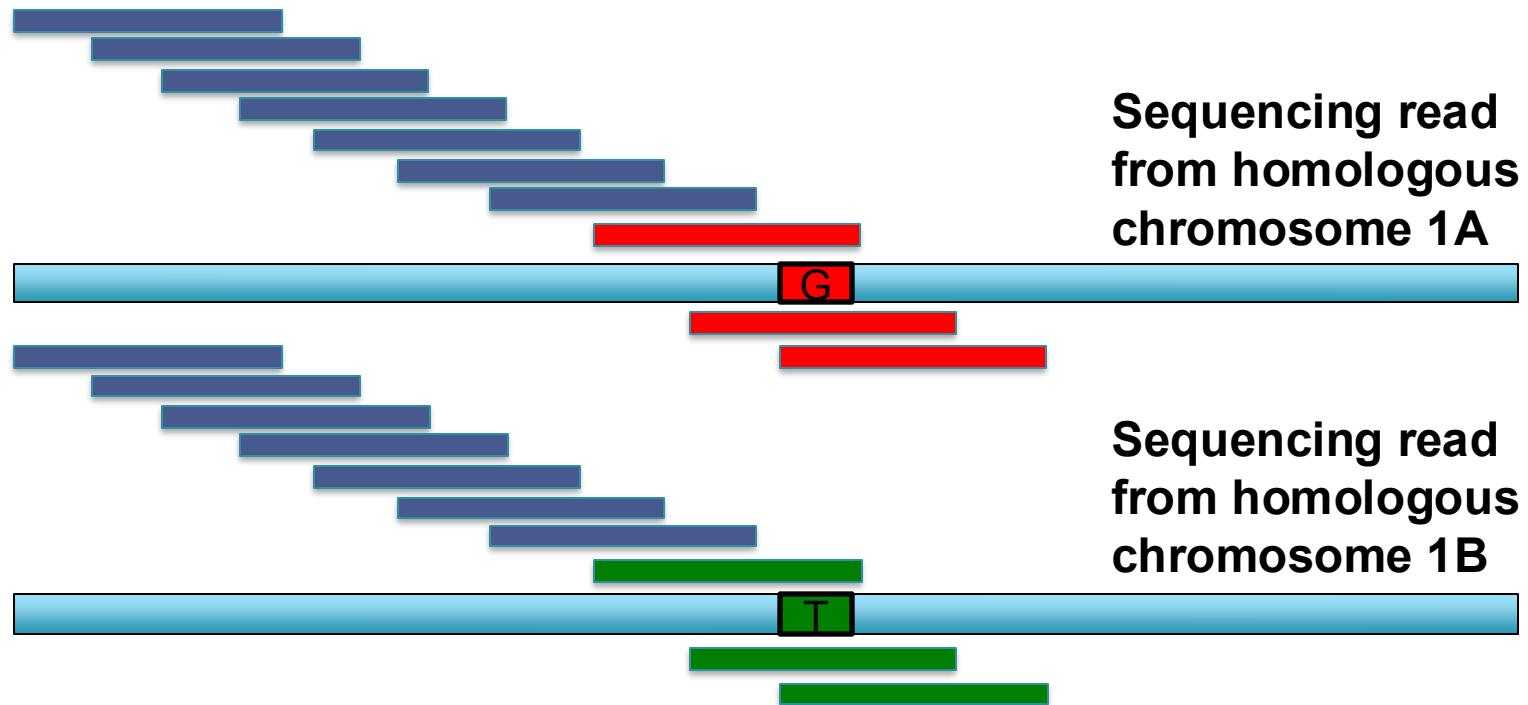
Sequencing read
from homologous
chromosome 1B



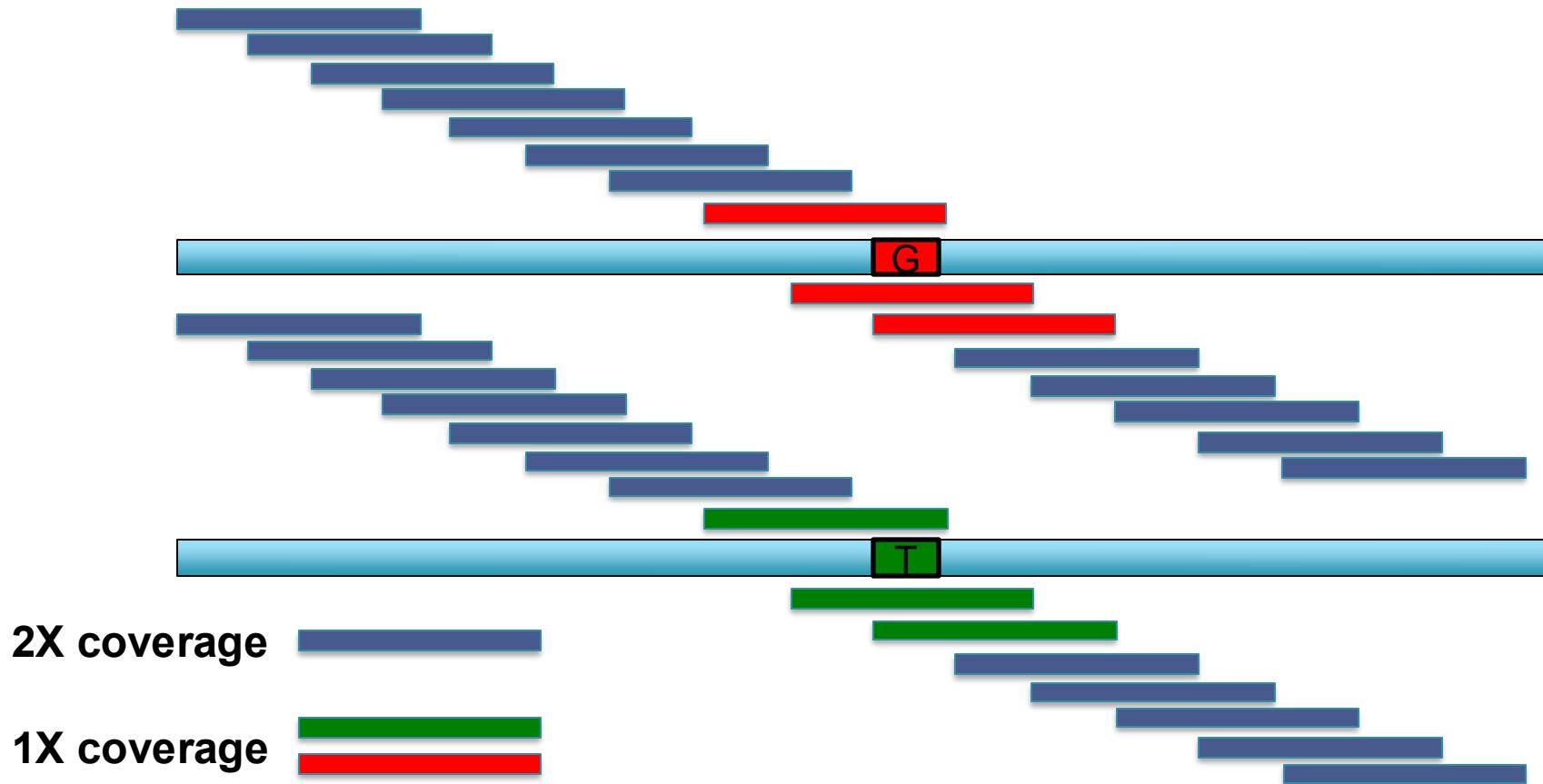
K-mer counting in heterozygous genomes



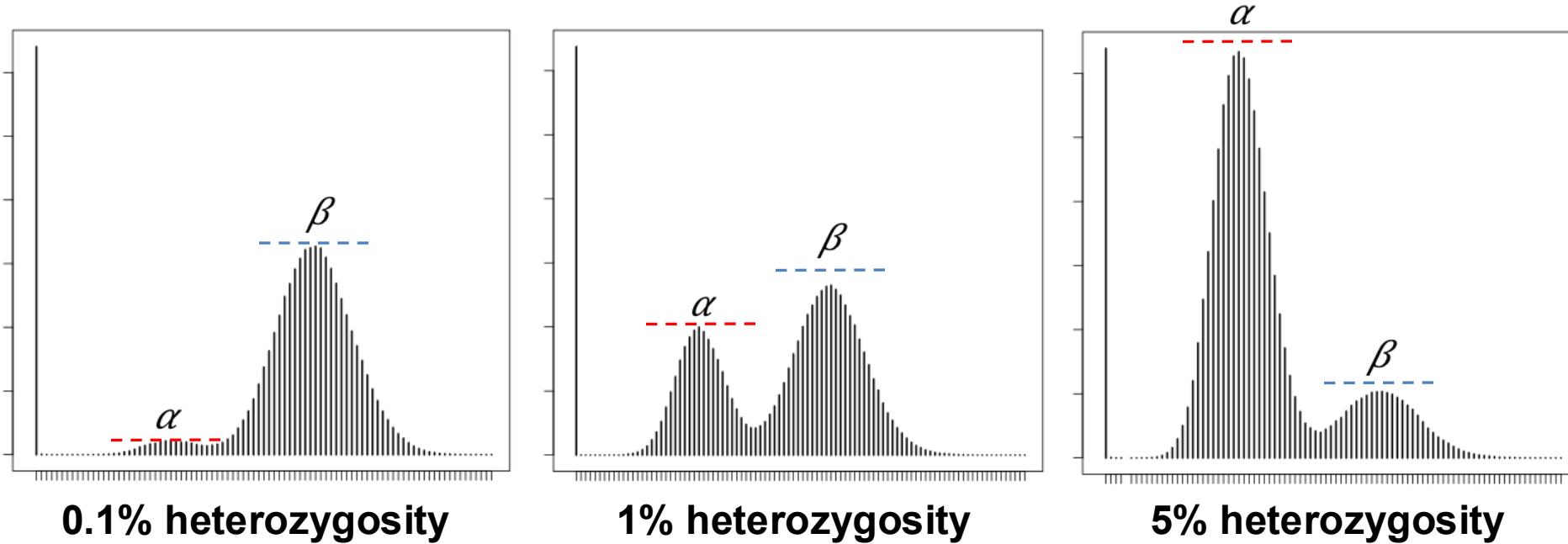
K-mer counting in heterozygous genomes



K-mer counting in heterozygous genomes

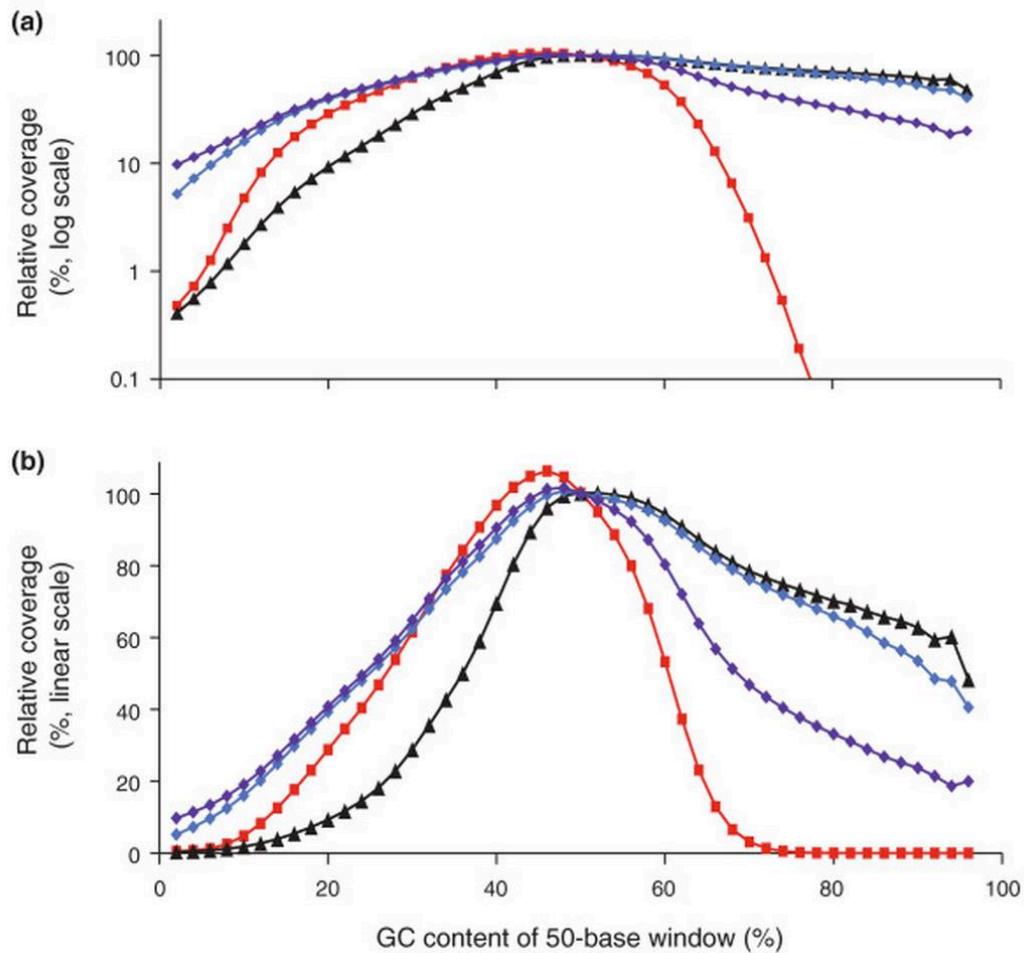


Heterozygous Kmer Profiles



- **Heterozygosity creates a characteristic “double-peak” in the Kmer profile**
 - Second peak at twice k-mer coverage as the first: heterozygous kmers average 50x coverage, homozygous kmers average 100x coverage
- **Relative heights of the peaks is directly proportional to the heterozygosity rate**
 - The peaks are balanced at around 1.25% because each heterozygous SNP creates 2^k heterozygous kmers (typically $k = 21$)

Beware of GC Biases



Illumina sequencing does not produce uniform coverage over the genome

- Coverage of extremely high or extremely low GC content will have reduced coverage in Illumina sequencing
- Biases primarily introduced during PCR; lower temperatures, slower heating, and fewer rounds minimize biases
- This makes it very difficult to identify variants (SNPs, CNVs, etc) in certain regions of the genome

Analyzing and minimizing PCR amplification bias in Illumina sequencing libraries.

Aird et al. (2011) *Genome Biology*. 12:R18.

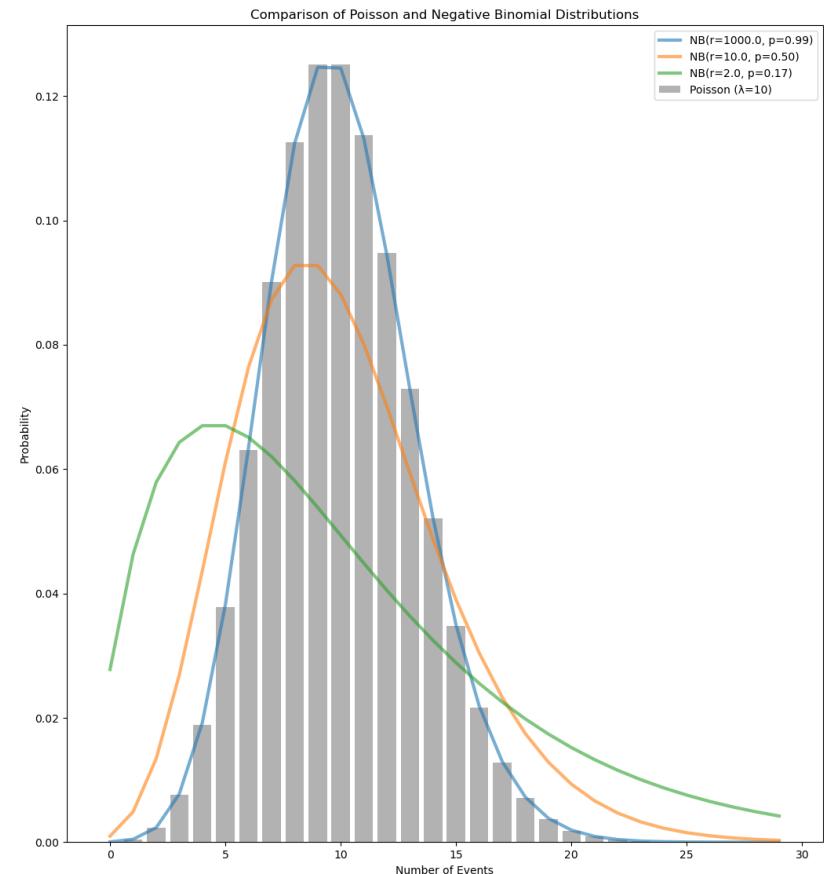
Negative Binomial Distribution

Models the number of failures in a sequence of independent and identically distributed Bernoulli trials before a specified number of successes (r) occurs

- Commonly used to model over-dispersed count data

Also arises as a continuous mixture of Poisson distributions where the mixing distribution of the Poisson rate is a gamma distribution.

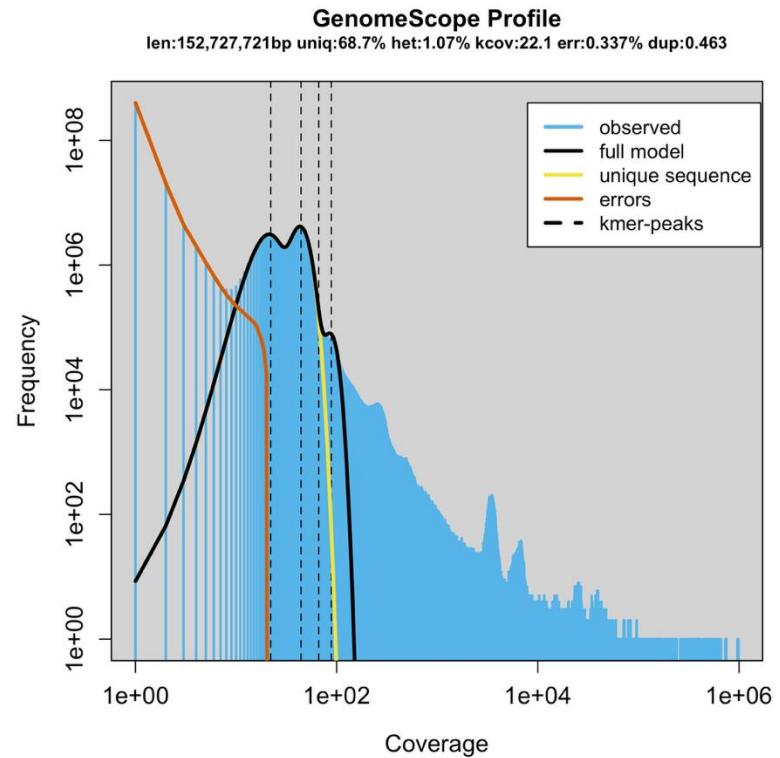
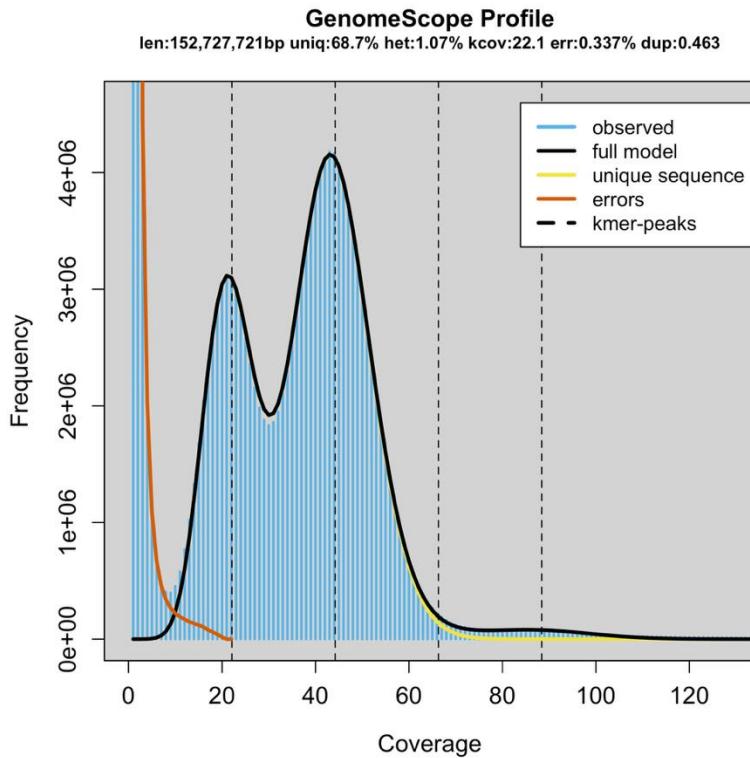
- Sequence coverage distribution where there is a non-uniform probability of a read starting at each position



$$f(k; r, p) \equiv \Pr(X = k) = \binom{k + r - 1}{k} (1 - p)^k p^r$$

GenomeScope: Fast genome analysis from short reads

<http://genomescope.org>



$$f(k) = \alpha \cdot \text{NB}(k; r_1, p_1) + \beta \cdot \text{NB}(k; r_2, p_2)$$

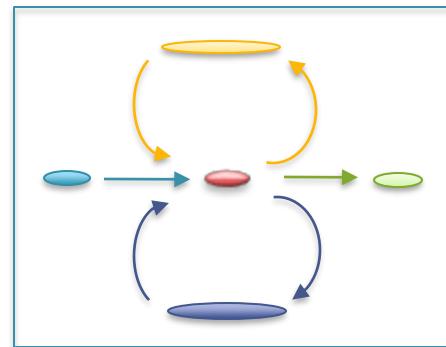
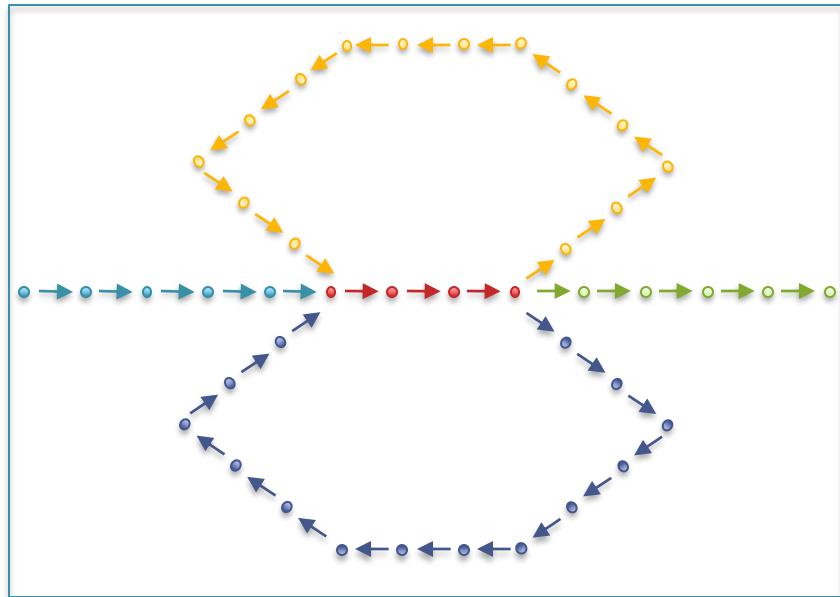
- Theoretical model agrees well with published results:
 - Quickly estimate genome size, rate of heterozygosity, and other genome properties
 - Generalized to higher ploidies by introducing additional terms
 - “Reference-free analysis” does not require the use of an assembled genome

Vulture, GW*, Sedlazeck FJ*, et al. (2017) *Bioinformatics*

Ranallo-Benavidez, TR. et al. (2020) *Nature Communication*

Unitigging / Unipathing

- After simplification and correction, compress graph down to its non-branching initial contigs
 - Aka “unitigs”, “unipaths”



Why do contigs end?

- (1) End of chromosome! ☺, (2) lack of coverage, (3) errors,
(4) heterozygosity and (5) repeats

Errors in the graph



(Chaisson, 2009)

Clip Tips

was the worst of times,

was the worst of **tymes**,

the worst of times, it

Pop Bubbles

was the worst of times,

was the worst of **tymes**,

times, it was the age

tymes, it was the age

the worst of **tymes**,

was the worst of

the worst of times,

worst of times, it

tymes,

was the worst of

it was the age

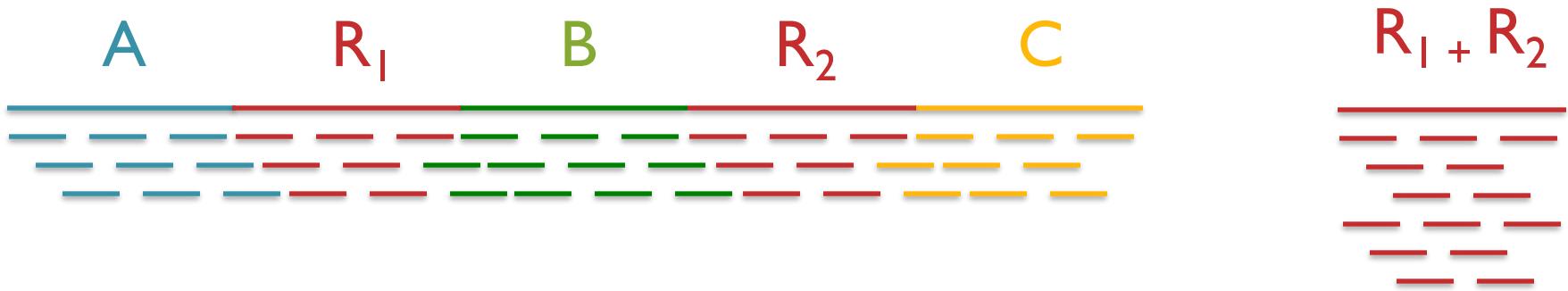
times,

Repetitive regions

Repeat Type	Definition / Example	Prevalence
Low-complexity DNA / Microsatellites	$(b_1 b_2 \dots b_k)^N$ where $1 \leq k \leq 6$ CACACACACACACACACA	2%
SINEs (Short Interspersed Nuclear Elements)	<i>Alu</i> sequence (~280 bp) Mariner elements (~80 bp)	13%
LINEs (Long Interspersed Nuclear Elements)	~500 – 5,000 bp	21%
LTR (long terminal repeat) retrotransposons	Ty1-copia, Ty3-gypsy, Pao-BEL (~100 – 5,000 bp)	8%
Other DNA transposons		3%
Gene families & segmental duplications		4%

- Over 50% of mammalian genomes are repetitive
 - Large plant genomes tend to be even worse
 - Wheat: 16 Gbp; Pine: 24 Gbp

Repeats and Coverage Statistics



- If n reads are a uniform random sample of the genome of length G , we expect $k = n \Delta/G$ reads to start in a region of length Δ .
 - If we see many more reads than k (if the arrival rate is $> A$) , it is likely to be a collapsed repeat

$$\Pr(X - \text{copy}) = \binom{n}{k} \left(\frac{X\Delta}{G} \right)^k \left(\frac{G - X\Delta}{G} \right)^{n-k}$$

$$A(\Delta, k) = \ln \left(\frac{\Pr(1 - \text{copy})}{\Pr(2 - \text{copy})} \right) = \ln \left(\frac{\frac{(\Delta n / G)^k e^{-\Delta n}}{k!}}{\frac{(2\Delta n / G)^k e^{-2\Delta n}}{k!}} \right) = \frac{n\Delta}{G} - k \ln 2$$

The fragment assembly string graph

Myers, EW (2005) Bioinformatics. 21(suppl 2): ii79-85.

Illumina Hacking

BIOINFORMATICS ORIGINAL PAPER

Vol. 29 no. 12 2013, pages 1492–1497
doi:10.1093/bioinformatics/btt178

Genome analysis

Advance Access publication May 22, 2013

Assembling the 20Gb white spruce (*Picea glauca*) genome from whole-genome shotgun sequencing data

Inanc Bird^{1,2,3*}, Anthony Raymond¹, Shaun D. Jackman¹, Stephen Pleasance¹, Robin Coope¹, Greg A. Taylor¹, Macaire Mai Saint Yuen¹, Christopher I. Keeling⁴, Dana Brand¹, Benjamin P. Vandervalk¹, Heather Kirk¹, Pawan Pandori¹, Richard A. Moran¹, Yongjun Zhao¹, Andrew J. Mungall⁵, Barry Jaquish⁵, Alvin Yanchuk⁶, Cara Brian Boyle¹, Jean Bousquet^{7,8}, Kermit Ritland⁶, John Mackay^{7,8}, Jörg B. Steven J.M. Jones^{1,2,9}

Associate Editor: Michael Brown

ABSTRACT

White spruce (*Picea glauca*) is a dominant conifer of the boreal forests of North America, and providing genomic resources for this commercially valuable tree will help improve forest management and conservation efforts. Sequencing and assembling the large and highly repetitive spruce genome though pushes the boundaries of the current technology. Here, we describe a whole-genome shotgun sequencing strategy using two Illumina sequencing platforms and an assembly approach using the ABySS software. We report a 20.8 giga base pair draft genome in 4.9 million scaffolds, with a scaffold N50 of 2035bp. We demonstrate how recent improvements in the sequencing technology, especially increasing read lengths and paired end reads from longer fragments have a major impact on the assembly contiguity. We also show that standard bioinformatics tools are instrumental in producing rapid draft assemblies.

Received on May 11, 2012; revised on July 11, 2012; accepted on August 1, 2012.
Published online in *PRJNA* 83435).

1 INTRODUCTION

The assembly of short reads to develop genomic resources for non-model species remains an active area of development (Schatz *et al.*, 2012). The feasibility of the approach and its scalability to

*To whom correspondence should be addressed.

© The Author 2013. Published by Oxford University Press.
This is an Open Access article distributed under the terms of the Creative Commons Attribution Non-Commercial License (<http://creativecommons.org/licenses/by-nc/3.0/>), which permits non-commercial re-use, distribution, and reproduction in any medium, provided the original work is [re-used](http://creativecommons.org/licenses/by-nc/3.0/), please contact journals.permissions@oxfordjournals.org. <http://dx.doi.org/10.1093/ehess/ehs005>

assemble the spruce genome, we used the Abyss algorithm (Simpson *et al.*, 2009), which captures a representation of read-to-read overlaps by a distributed de Bruijn graph and uses parallel computations to build the target genome. The modular nature of the tool allowed us to execute a large number of tests to tune the message passing interface for a successful execution, train the assembly parameters for an optimal assembly and quantify the utility of long reads for genome assemblies. To the best of our knowledge, the Abyss algorithm is unique in its ability to enable genome assemblies of this scale using whole-genome shotgun sequencing data.

2 METHODS

2.1 Sample collection

Apical shoot tissues were collected in April 2006 from a single white spruce (*Picea glauca*, genotype PG29) tree at the Kalamalka Research Station of the British Columbia Ministry of Forests and Ranges, Vernon, British Columbia, Canada. Genomic DNA was extracted from 60 gm tissue by Bio-S&T (<http://www.biost.com/>), Montreal, QC, Canada, using an organelle exclusion method yielding 300 µg of high quality pure nuclear DNA.

2.2 Library preparation and sequencing

2.1. Library preparation and sequencing
DNA quality was assessed by spectrophotometry and gel electrophoresis before library construction. DNA was sheared for 45 s using an E210 sonicator (Covaris) and then analysed on 8% PAGE gels. The 200–300 bp (for libraries with 250 bp insert size) or 450–550 bp (for libraries with 500 bp insert size) DNA size fractions were excised and eluted from the gel slices overnight at 4°C in 300 µl of elution buffer [5:1 vol/vol

LoTE buffer [3 mM Tris-HCl (pH 7.5), 0.2 mM EDTA, 7.5 mM ammonium acetate] and was purified with a Spin-X Filter Tube (Fisher Scientific) and ethanol precipitation. Genome libraries were prepared using a modified paired-end tag (PET) protocol supplied by Illumina Inc. This involved DNA end repair and formation of 3' adenosine overhangs using the Klenow fragment of DNA polymerase I (3'-5' exonuclease minus) and ligation to Illumina PE adaptors (with 5' overhangs). Adaptor-ligated products were purified on QIAquick spin columns.

Adapter-ligated products were purified on QIAquick spin columns (Qiagen) and amplified using Phusion DNA polymerase (NEB) and 10 PCR cycles with the PE primer 1.0 and 2.0 (Illumina). PCR products of the desired size range were purified from adapter ligation artifacts using 8% PAGE gels. DNA quality was assessed and quantified using an Agilent DNA 1000 series II assay (Agilent) and Nanodrop 7500 spectrophotometer (Nanodrop). DNA was subsequently diluted to 8 nM. The final concentration was confirmed using a Quant-iT T dsDNA HS assay kit and Qubit fluorometer (Invitrogen).

The mate pair (MPET, a.k.a. jumping) libraries were constructed using 4 µg of genomic DNA with the Illumina Nextera Mate-Pair library construction protocol and reagent (FC-132-1001). The genomic DNA sample was simultaneously fragmented and tagged with a barcode consisting of a unique sequence identifier (UID) and a sequencing index tag. The tag sequence is a 12 bp sequence unique to each library. The DNA sample was digested with a restriction enzyme to generate fragments of approximately 6 kb. The gels were run by a standard displacement reaction using a polyacrylamide gel to ensure that all fragments were flush and ready for circularization. After an AmPure Bead cleanup, size selection was done on a 0.6% agarose gel to exclude 6-kb and 9-13-kb fractions, which were purified using a Zymoclean Large Fragment DNA Recovery Kit. The fragments were circularized by ligation, followed by a digestion to remove any linear molecules and left circularized DNA for shearing. The sheared DNA fragments (approximately 300 bp) were then bisulfite converted (Chemicon International) prior to the step of biotinylation (using biotinylated mag3 beads) and the unwanted unbiotinylated molecules were washed away. The DNA fragments were then end-mirrored and A-tailed, following the

protocol and ligated to indexed TruSeq adapters. The final library was enriched by a 10-cycle PCR and purified by AMPure bead clean-up. Library quality and size were assessed by Agilent DNA 1000 series II assay and KAPA Library Quantification protocol. The two fractions were pooled for sequencing paired end 100 bp using Illumina HiSeq2000.

The construction of the 12 kb mate pair libraries was achieved by a hybrid 454/Illumina procedure. Briefly, 50 µg of genomic DNA was fragmented for 20 cycles at speed code 12 using a Herculase II MSK cutter (Agilent). The fragmented DNA was then ligated to modified DNA ends on a 1% agarose gel, and fragments >15 kb were extracted. Biotinylated circularization adaptors (T7E1 Tailed Adapter set (454 Life Sciences/Roche CT)) were added to ends of the gel-extracted fragments. Recombination of the ends was performed with Cre recombinase (Stratagene) and the resulting linearized plasmid fragments were removed with Plasmid Safe (Epientre, Madison, WI). molecules were fragmented using GS Rapid Library Nebulizer (Sciences/Roche, Branford, CT), and fragment end-repair for tailing was performed with the GS Rapid Library preparation kit (Sciences/Roche, Branford, CT). T7E1 Adapters (Illumina) were ligated to the repaired, A-tailed ends. Biotinylation was enriched using Streptavidin-coated Dynabeads (Life Tech Grand Island, NY) and amplified by PCR using Illumina primers.

Assembling the 20 Gb white spruce genome

Assembling the 20 Gb white spruce (*Picea glauca*) genome from whole-genome shotgun sequencing data

Nanci Birol^{1,2,3,*}, Anthony Raymond¹, Shaun Jackman¹, Stephen Pleasance¹, Robin Cooper¹, Greg A Taylor¹, Macaire Man Saint Yuen⁴, Christopher Kieling⁴, Dena Brandl¹, Benjamin Vandervalk¹, Heather Kirk¹, Pawan Pandoh¹, Richard A Moore¹, Yongjun Zhao¹, Andrew J Mungall¹, Barry Jaquish¹, Alvin Yanuchuk¹, Carol Ritland^{4,6}, Brian Boyle¹, Jean Bousquet^{1,8}, Christen Britland^{1,8}, John MacKay^{7,8}, Ira Bohmlein^{4,6}, Steven IM Jones^{2,9}

- British Columbia Cancer Agency, Genome Sciences Centre, Vancouver, BC V5Z 4S6
University of British Columbia, Department of Medical Genetics, Vancouver, BC V6H 3N1
Simon Fraser University, School of Computing Science, Burnaby, BC V5A 1S6
University of British Columbia, Michael Smith Laboratories, Vancouver, BC V6T 1Z4
British Columbia Ministry of Forests, Lands and Natural Resource Operations, Victoria, BC V8W 9C2
University of British Columbia, Department of Forest Sciences, Vancouver, BC V6T 1Z4
Université Laval, Institute for Systems and Integrative Biology, Québec, QC G1V 0A6
Université Laval, Department of Wood and Forest Sciences, Québec, QC G1V 0A6
Simon Fraser University, Department of Molecular Biology and Biochemistry, Burnaby, BC V5A 1S6

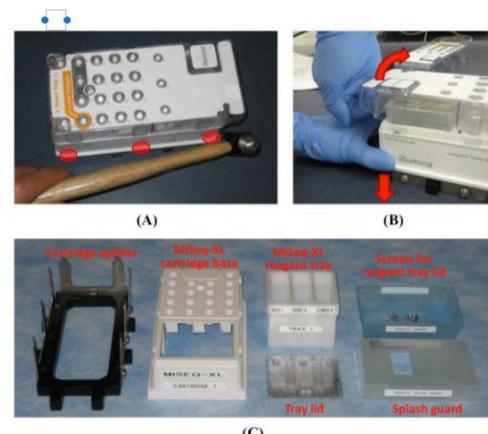


Figure S1. Modification of the MiSeq cartridge. MiSeq reagent cartridge was modified to allow for longer read lengths. (A, B) Opening of the clamshell style cartridge. (C) Contents of the modified cartridge. This was initially used to combine two PE150 kits for PE300 runs. When Illumina introduced the P250 kit, the same apparatus was used to enable PE500 runs.

Paired-end and Mate-pairs

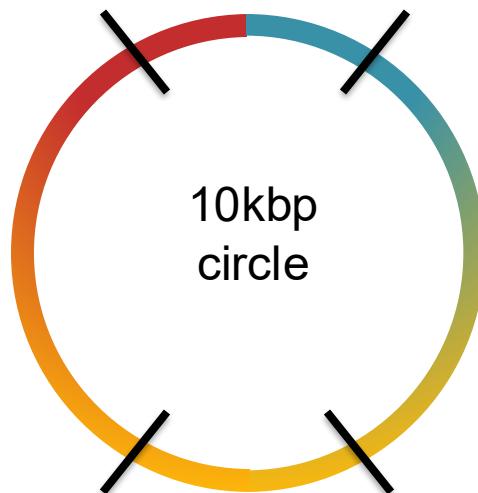
Paired-end sequencing

- Read one end of the molecule, flip, and read the other end
- Generate pair of reads separated by up to 500bp with inward orientation



Mate-pair sequencing

- Circularize long molecules (1-10kbp), shear into fragments, & sequence
- Mate failures create short paired-end reads



2x100 @ ~10kbp (outies)

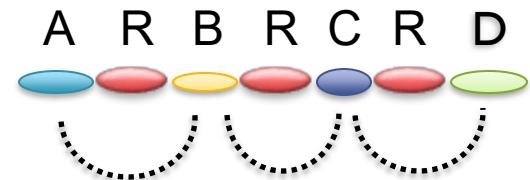
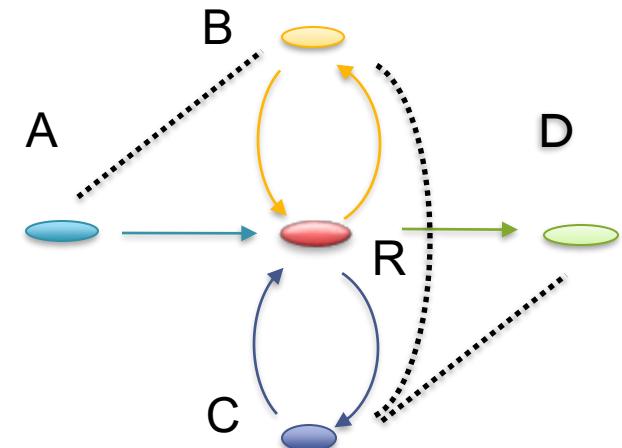


2x100 @ 300bp (innies)



Scaffolding

- Initial contigs (aka unipaths, unitigs) terminate at
 - Coverage gaps: especially extreme GC
 - Conflicts: errors, repeat boundaries
- Use mate-pairs to resolve correct order through assembly graph
 - Place sequence to satisfy the mate constraints
 - Mates through repeat nodes are tangled
- Final scaffold may have internal gaps called sequencing gaps
 - We know the order, orientation, and spacing, but just not the bases. Fill with Ns instead

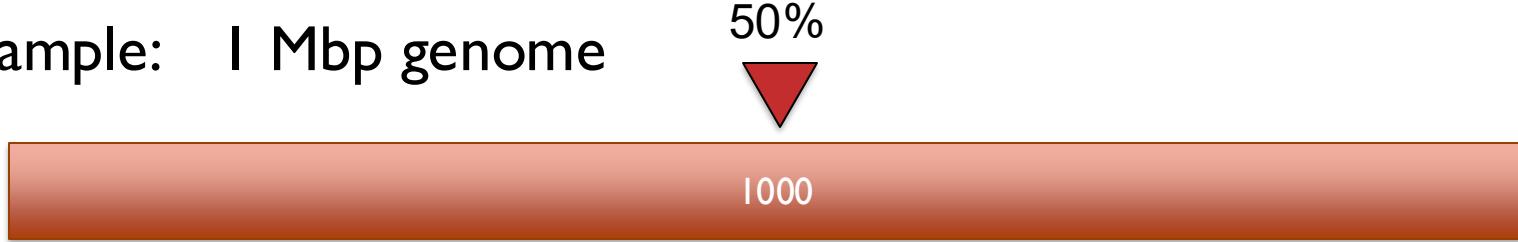


Why do scaffolds end?

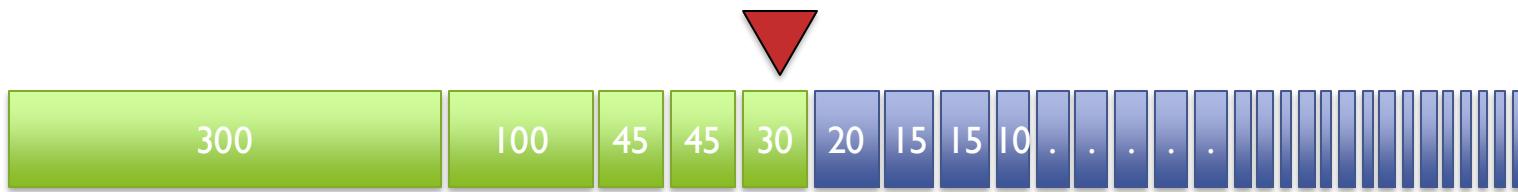
Contig/Scaffold N50

Def: 50% of the genome is in contigs as large as the N50 value

Example: 1 Mbp genome

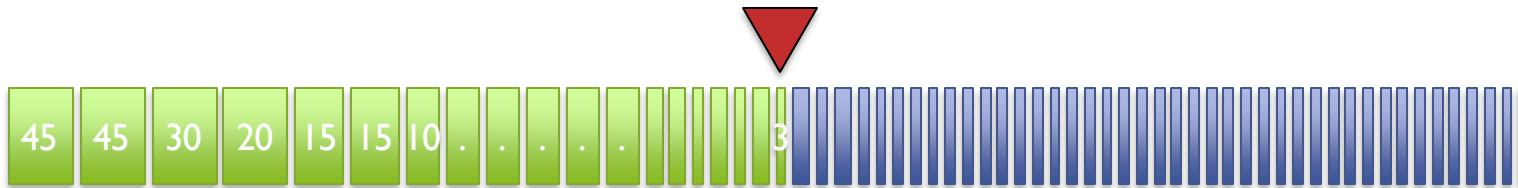


A



N50 size = 30 kbp

B



N50 size = 3 kbp

Contig/Scaffold N50

Def: 50% of the genome is in contigs as large as the N50 value

50%

Better N50s improves the analysis in every dimension

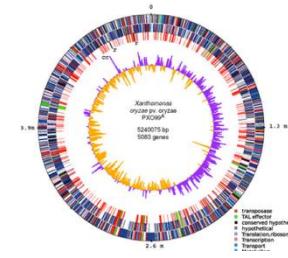
- Better resolution of genes and flanking regulatory regions
- Better resolution of transposons and other complex sequences
- Better resolution of chromosome organization
- Better sequence for all downstream analysis

Just be careful of N50 inflation!

- A very very very bad assembler in 1 line of bash:
- `cat *.reads.fa > genome.fa`

N50 size = 3 kbp

Assembly Summary



Assembly quality depends on

- Assembly is a hierarchical, starting from individual reads, build high confidence contigs/unitigs, incorporate the mates to build scaffolds
 - Extensive error correction is the key to getting the best assembly possible from a given data set
 - Recommend spades for short read assembly
 - Integrates error correction and scaffolding