

Genome Assembly

Michael Schatz

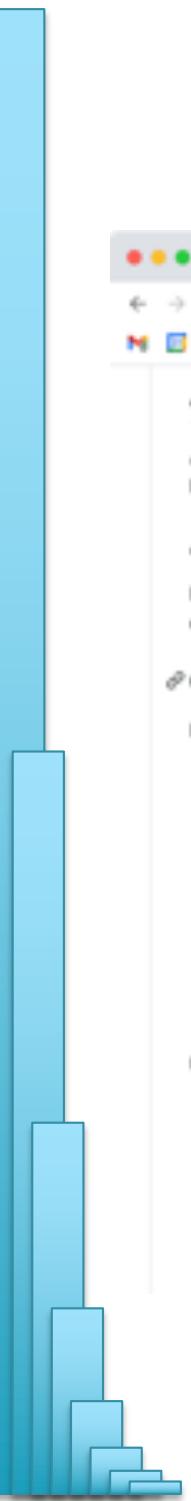
Feb 1, 2021

Lecture 3: Applied Comparative Genomics



Assignment 1: Chromosome Structures

Due Feb 3 @ 11:59pm



A screenshot of a web browser window showing the assignment details. The title bar says "appliedgenomics2021/assignment1". The address bar shows "github.com/schatzlab/appliedgenomics2021/tree/master/assignments/assignment1". The page content includes the assignment title, date, overview, and a list of species with links to their genome size files.

Assignment 1: Chromosome Structures

Assignment Date: Wednesday, Jan 27, 2021
Due Date: Wednesday, Feb. 3, 2021 @ 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

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

1. [Arabidopsis thaliana \(TAIR10\)](#) - An important plant model species [\[info\]](#)
2. [Tomato \(Solanum lycopersicum v4.00\)](#) - One of the most important food crops [\[info\]](#)
3. [E. coli \(Escherichia coli K12\)](#) - One of the most commonly studied bacteria [\[info\]](#)
4. [Fruit Fly \(Drosophila melanogaster, dm3\)](#) - One of the most important model species for genetics [\[info\]](#)
5. [Human \(hg38\) - us::\)](#) [\[info\]](#)
6. [Wheat \(Triticum aestivum, IWGSC\)](#) - The food crop which takes up the largest land area [\[info\]](#)
7. [Worm \(Caenorhabditis elegans, ce10\)](#) - One of the most important animal model species [\[info\]](#)
8. [Yeast \(Saccharomyces cerevisiae, sacCer3\)](#) - an important eukaryotic model species, also good for bread and beer [\[info\]](#)

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

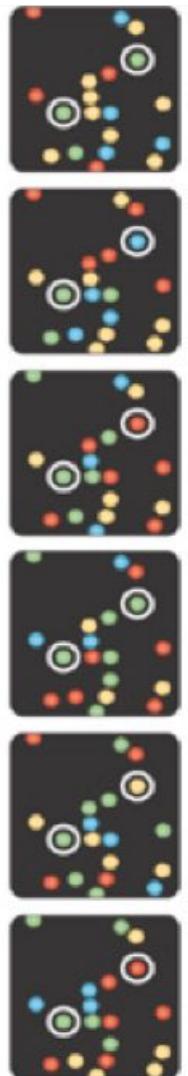
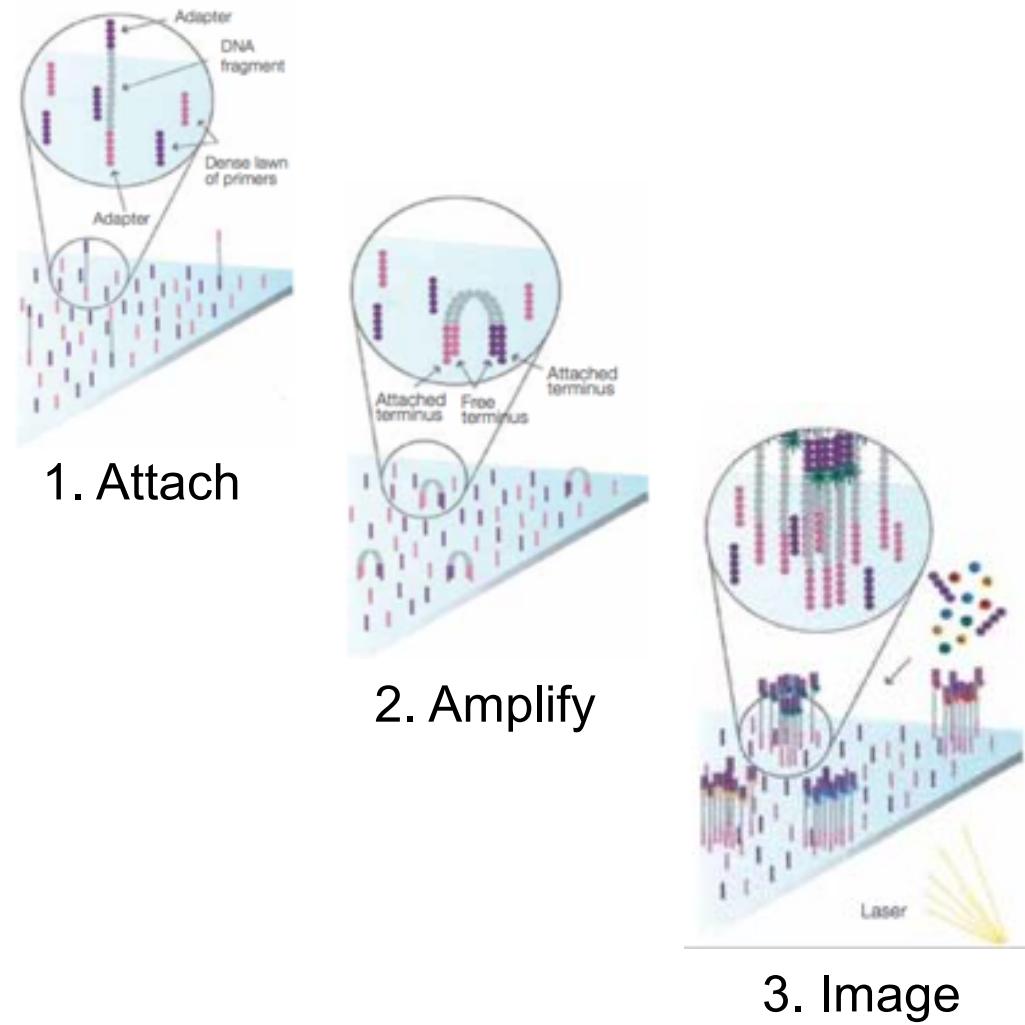
Part I: Recap

Second Generation Sequencing



Illumina HiSeq 2000
Sequencing by Synthesis

>60Gbp / day

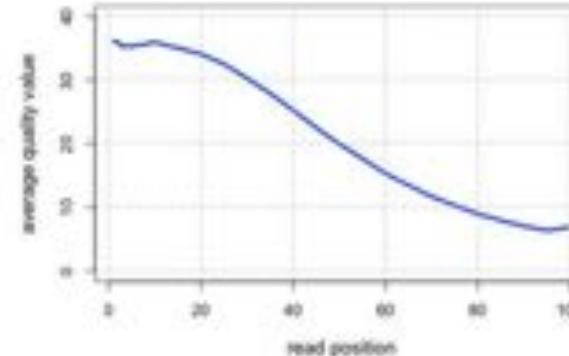


Metzker (2010) Nature Reviews Genetics 11:31-46
<https://www.youtube.com/watch?v=fCd6B5HRaZ8>

Illumina Quality

QV	P _{error}
40	1/10000
30	1/1000
20	1/100
10	1/10

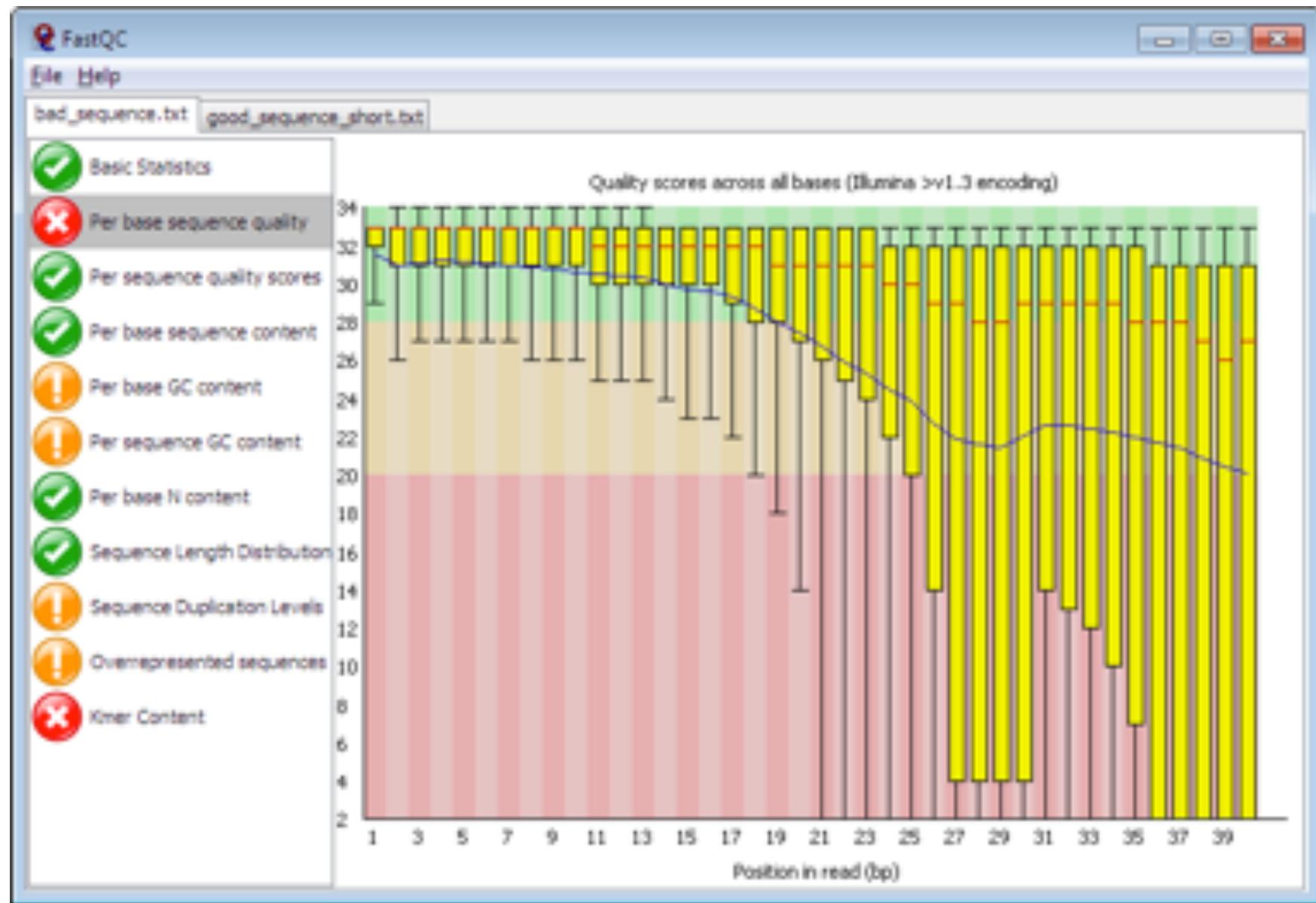
$$Q_{\text{sanger}} = -10 \log_{10} p$$



.....SS.....
.....XXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXX.....
.....IIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIII.....
.....JJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJ.....
.....LLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLL.....
! "#\$%&'()*+,./0123456789!;<=>?@ABCDEFGHIJKLMNOPQRSTUVWXYZ{`}^_`abcdefgijklmnopqrstuvwxyz{|}`^_`
| | | | |
33 59 64 73 104 126

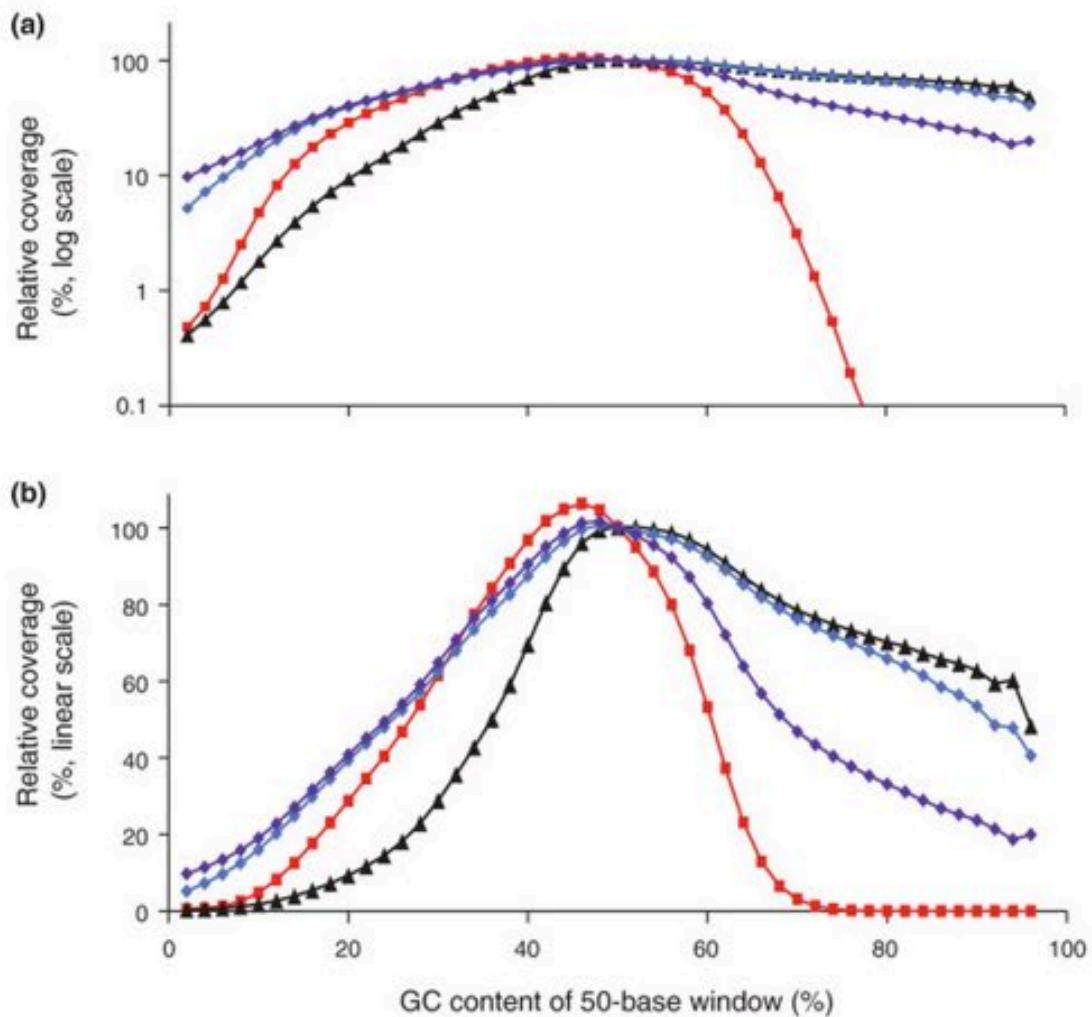
S - Sanger Phred+33, raw reads typically (0, 40)
X - Solexa Solexa+64, raw reads typically (-5, 40)
I - Illumina 1.3+ Phred+64, raw reads typically (0, 40)
J - Illumina 1.5+ Phred+64, raw reads typically (3, 40)
with 0=unused, 1=unused, 2=Read Segment Quality Control Indicator (bold)
(Note: See discussion above).
L - Illumina 1.8+ Phred+33, raw reads typically (0, 41)

FASTQC: Is my data any good?



<http://www.bioinformatics.babraham.ac.uk/projects/fastqc/>

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.

Question?

We would love to generate
longer and longer reads with this technology

What can we do?

Illumina Hacking

BIOINFORMATICS ORIGINAL PAPER

Genome analysis

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

Inanc Birol^{1,2,3*}, Anthony Raymond¹, Shaun D. Jackman¹, Stephen Pleasance¹, Robin Coop¹, Greg A. Taylor¹, Macaire Man Saint Yuen⁴, Christopher I. Keeling⁴, Dana Brand¹, Benjamin P. Vandervalk¹, Heather Kirk¹, Pawan Pandoh¹, Richard A. Moore¹, Yongjun Zhao¹, Andrew J. Mungall⁵, Barry Jaquish⁵, Alvin Yanchuk⁶, Carol Ritland⁶, Brian Boyle⁷, Jean Bousquet^{7,8}, Kermit Ritland⁶, John MacKay^{7,8}, Jörg B. Steven J.M. Jones^{1,2,9}

¹Genome Sciences Centre, British Columbia Cancer Agency, Vancouver, BC V5Z 4S6, Canada; ²University of British Columbia, Vancouver, BC V6H 3N1, Canada; ³School of Computer Science, Fraser University, Burnaby, BC V5A 1S6, Canada; ⁴Michael Smith Laboratories, University of British Columbia, Vancouver, BC V6T 1Z4, Canada; ⁵British Columbia Ministry of Forests, Lands and Natural Resource Operations, Victoria, BC V8W 9C2, Canada; ⁶Department of Forest Sciences, University of British Columbia, V1Z 1Z4, Canada; ⁷Institute for Systems and Integrative Biology, Université Laval, Québec, QC G1V 0A6, Canada; ⁸Department of Wood and Forest Sciences, Université Laval, Québec, QC G1V 0A6, Canada; ⁹Molecular Biology and Biochemistry, Simon Fraser University, Burnaby, BC V5A 1S6, Canada

Associate Editor: Michael Brudno

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/8 software. We report a 20.8 giga base pairs draft genome in 4.9 million scaffolds, with a scaffold N50 of 20354 bp. 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 note that scalable bioinformatics tools are instrumental in providing rapid draft assemblies.

Availability: The *Picea glauca* genome sequencing and assembly data is available through NCBI (Accession#: ALWZ01000000000 PID: PTUNAB3439, <http://www.ncbi.nlm.nih.gov/bioproject/83435>).

Correspondence: ibiro@bcrc.ca

Supplementary Information: Supplementary data are available at [Bioinformatics](http://bioinformatics.oxfordjournals.org) online.

Received on March 20, 2013; revised on April 10, 2013; accepted on April 11, 2013

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 properly cited.

Vol. 29 no. 12 2013, pages 1490–1497
doi:10.1093/bioinformatics/btt178

Advance Access publication May 22, 2013

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 large 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 Range, Vernon, British Columbia, Canada. Genomic DNA was extracted from 60 gm tissue by BioS&T (<http://www.biost.com/>), Montreal, QC, Canada using an organic extraction method yielding 300 µg of high quality pure nuclear DNA.

2.2 Library preparation and sequencing

DNA quality was assessed by spectrophotometry and gel electrophoresis before library construction. DNA was sheared for 5' using an E210 sonicator (Cavitar) and then analyzed 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 500 µl of elution buffer [5:1 (vol/vol) LTote buffer (Takara) + Tris-HCl (pH 7.5), 0.2 mM EDTA] / 7.5 M ammonium acetate. The 500 bp size fraction was eluted into a DNase-free tube and ethanol precipitation. Genomic libraries were prepared using a modified paired-end tag (PET) protocol supplied by Illumina Inc. This involved DNA end repair and formation of 5' adenine overhang using the Klenow fragment of DNA polymerase I (5'-exonuclease minus) and ligation to Illumina PE adaptors (with 5' overhang). 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 purified by a Maxwell® 24 automated liquid handling system using a Quant-iT 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-132-1011). The genomic DNA sample was simultaneously fragmented and tagged with a biotin containing mate pair junction adaptor, which left a short sequence gap in the fragmented DNA. The gap was filled by a strand displacement reaction using a polymerase to ensure that all fragments were flush and ready for circulation. After an AMPure bead cleanup, size selection was done on a 0.6% agarose gel to excise 6–9 kb and 9–14 kb fractions, which were purified using a Zymo Clean & Filter DNA Recovery Kit. The fragments were size-selected by ligation followed by shearing and capture using linear molecules and left circulatory DNA for shearing. The sheared DNA fragments that contain the biotynlated junction adaptor (mate pair fragments) were purified by means of binding to streptavidin magnetic beads, and the unwanted unbiotynylated molecules were washed away. The DNA fragments were then end-repaired and A-tailed following the

protocol and ligated to indexed TruSeq adaptors. The final library was emulsion sequenced for 20 cycles at speed code 12 using a Hybrid-seq Marlowe (Marlborough, MA) equipped with a large assembly mode. Illumina DNA was loaded on a 1% agarose gel, and fragments 18 kb were extracted. Biotinylation circulation adaptors (TruSeq Paired-end Adapters set (454 Life Sciences, Roche CT)) were added to each of the generated Mate-Pair and MPET libraries. Ligation ends were converted with Cea-mouse recombinant Biotin, Ipswich, MA, and linear molecules remaining were removed with Plasmid Safe (Epicenter, Madison, WI). DNA molecules were fragmented using GS Rapid Library Nebulizer (Science Roche, Bradford, CT), and fragment end-repair tailing was performed with the GS Rapid Library preparation Sciences Roche, Bradford, CT). Truseq Adapters (Illumina CA) were ligated to the repaired/A-tailed ends. Biotinylation was enriched using Streptavidin-coupled Dynabeads (Life Sciences Grand Island, NY) and amplified by PCR using Illumiprimer primers.

Random bacterial artificial chromosome (BAC) was sequenced using DNA from the same genome on a Titanium with 4-kb paired-end libraries at the PacBio Genome Sequencer of the Institute for Systems and Integrative Biology, Université Laval, Québec City, QC). A single paired-end preparation was used for 15 BACs (representing ~13.8 Mb of genome). The DNA was sequenced on a Titanium with the following modifications: 15 kb was fragmented using a Hydroaser with a standard assembly mode at speed code 18. 6–10-kb fragments were extracted from GS-FLX library adaptors were ligated to the repaired/A-tailed ends. GS-FLX sequencing was performed using the titanium chemistry according to manufacturer's instructions (454 Life Sciences Grand Island, CT). Sanger sequencing method was used to obtain BAC sequencing data as previously described (Hamberger Keeling *et al.*, 2013).

2.3 MiSeq modification

In sequencing the spruce genome, we generated longer runs modifying the MiSeq platform. The MiSeq uses a clamshell style cartridge (Supplementary Fig. S1A) to hold reagent tubes in an array covered by the MiSeq's sippers. Most of the reagents are at length independent steps such as denaturation and cluster formation, the Scan, Cleavage and Incorporation mixed at each cycle. Although the MiSeq allows any read specified in the control software, the reagent cartridge carries during the run without stopping it. Increasing the read length requires increasing the quantity of the length-dependent cartridge. This led to the solution of combining the long reagents of two kits into one.

A tool was designed that opens the snap-hook latches carrying together. We show in Supplementary Figs S1B and S2, give the reagent tubes, yet allowing the cartridge to be pulled up without damage to its components (Supplementary Fig. S4) and, the stock length-dependent reagent container allows runs of ~450 cycles in total. To maximize the potential of kit approach, a new reagent tray with 70 ml wells was placed in a modified clamshell base.

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

Inanc Birol^{1,2,3*}, Anthony Raymond¹, Shaun D Jackman¹, Stephen Pleasance¹, Robin Coop¹, Greg A Taylor¹, Macaire Man Saint Yuen⁴, Christopher I. Keeling⁴, Dana Brand¹, Benjamin P. Vandervalk¹, Heather Kirk¹, Pawan Pandoh¹, Richard A. Moore¹, Yongjun Zhao¹, Andrew J. Mungall⁵, Barry Jaquish⁵, Alvin Yanchuk⁶, Carol Ritland⁶, Brian Boyle⁷, Jean Bousquet^{7,8}, Kermit Ritland⁶, John MacKay^{7,8}, Jörg Bohlmann^{4,6}, Steven J.M. Jones¹

¹British Columbia Cancer Agency, Genome Sciences Centre, Vancouver, BC V5Z 4S6

²University of British Columbia, Department of Medical Genetics, Vancouver, BC V6H 3N1

³Simon Fraser University, School of Computing Science, Burnaby, BC V5A 1S6

⁴University of British Columbia, Michael Smith Laboratories, Vancouver, BC V6T 1Z4

⁵British Columbia Ministry of Forests, Lands and Natural Resource Operations, Victoria, BC V8W 9C2

⁶University of British Columbia, Department of Forest Sciences, Vancouver, BC V6T 1Z4

⁷Université Laval, Institute for Systems and Integrative Biology, Québec, QC G1V 0A6

⁸Université Laval, Department of Wood and Forest Sciences, Québec, QC G1V 0A6

⁹Simon Fraser University, Department of Molecular Biology and Biochemistry, Burnaby, BC V5A 1S6

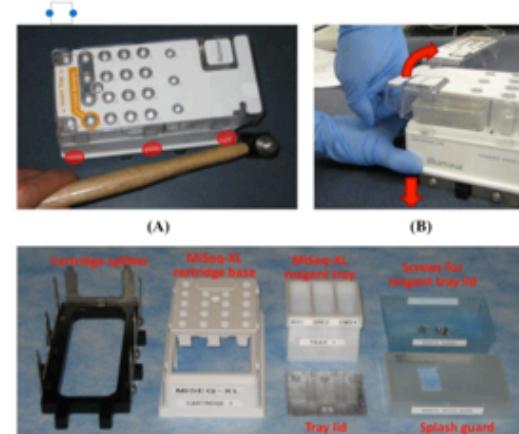


Figure S1. Modification of the MiSeq cartridge. MiSeq reagent cartridge was modified to allow for longer read lengths. (A, B) Opening of the clamshell style cartridge. (C) Contents of the modified cartridge. This was initially used to combine two PE150 kits for PE300 runs. When Illumina introduced the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

Paired-end and Mate-pairs

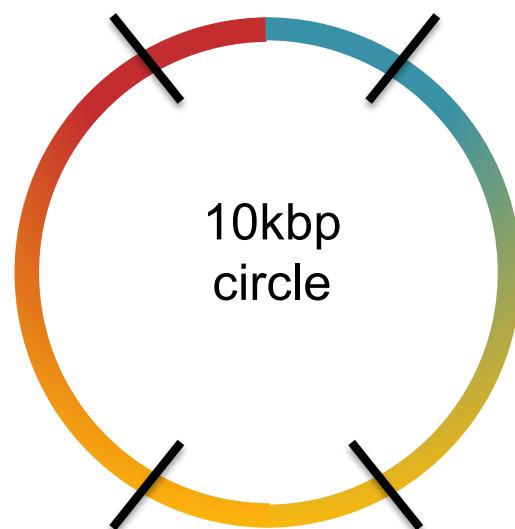
Paired-end sequencing

- 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)



FASTQ Files



```
@SEQ_ID  
GATTGGGGTTCAAAGCAGTATCGATCAAATAGTAAATCCATTGTTCAACTCACAGTTT  
+  
! ' * ( ( ( ****+ ) ) % % % ++ ) ( % % % ) . 1 * * * - + * ' ' ) ) * * 55CCF>>>>>cccccccc65
```

@Identifier
Sequence
+Separator
Quality Values
...

Illumina Sequencing Summary

Advantages:

- Best throughput, accuracy and read length for any 2nd gen. sequencer
- Fast & robust library preparation



Illumina HiSeq

~3 billion paired 100bp reads

~600Gb, \$10K, 8 days

(or “rapid run” ~90Gb in 1-2 days)

Illumina X Ten / NovaSeq

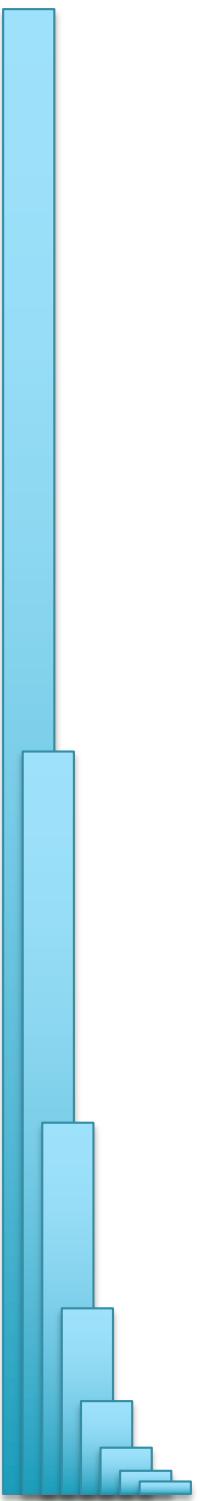
~6 billion paired 150bp reads

1.8Tb, <3 days, ~1000 / genome(\$\$)

(or “rapid run” ~90Gb in 1-2 days)

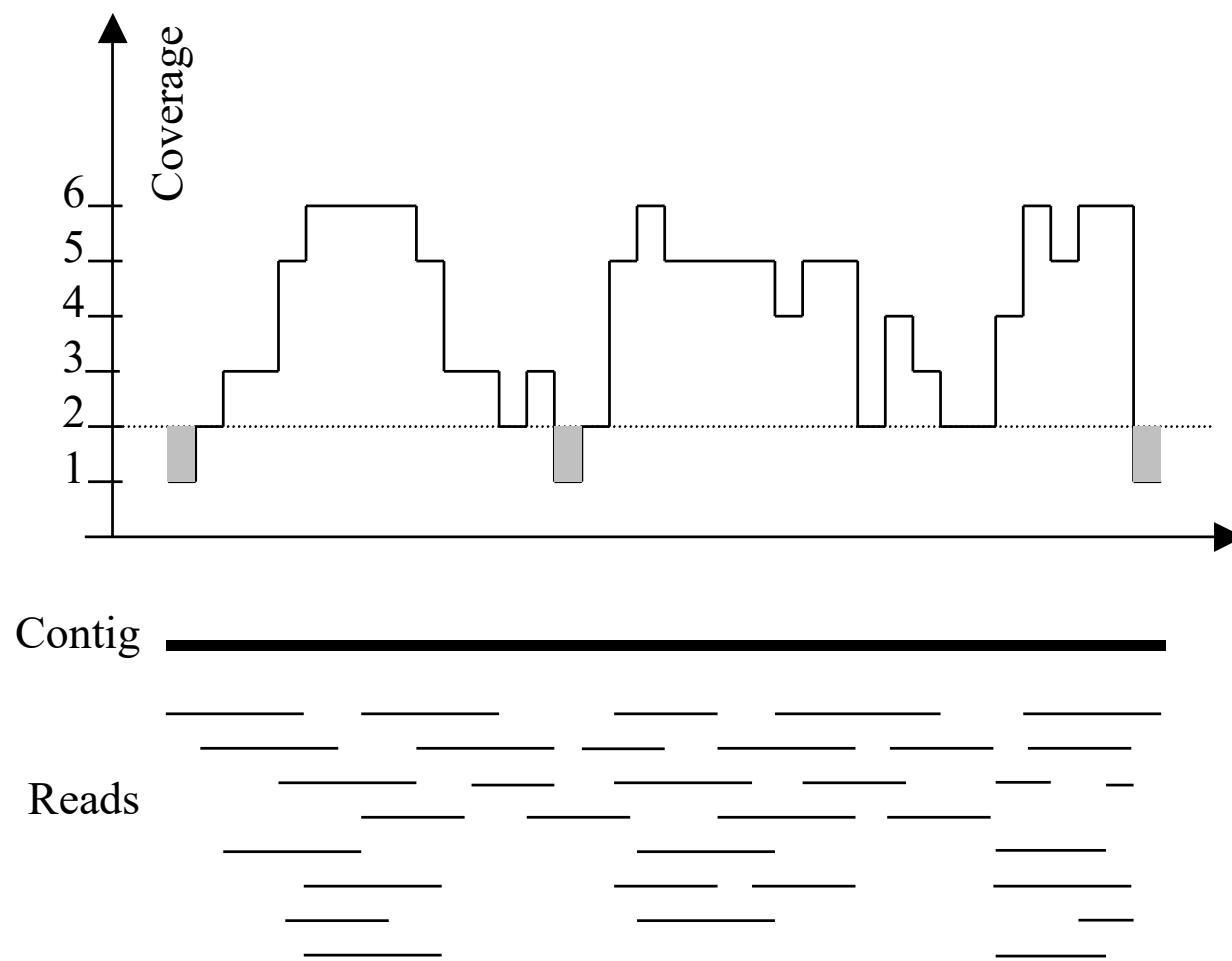
Illumina NextSeq

One human genome in <30 hours



Part 2: Coverage

Typical sequencing coverage

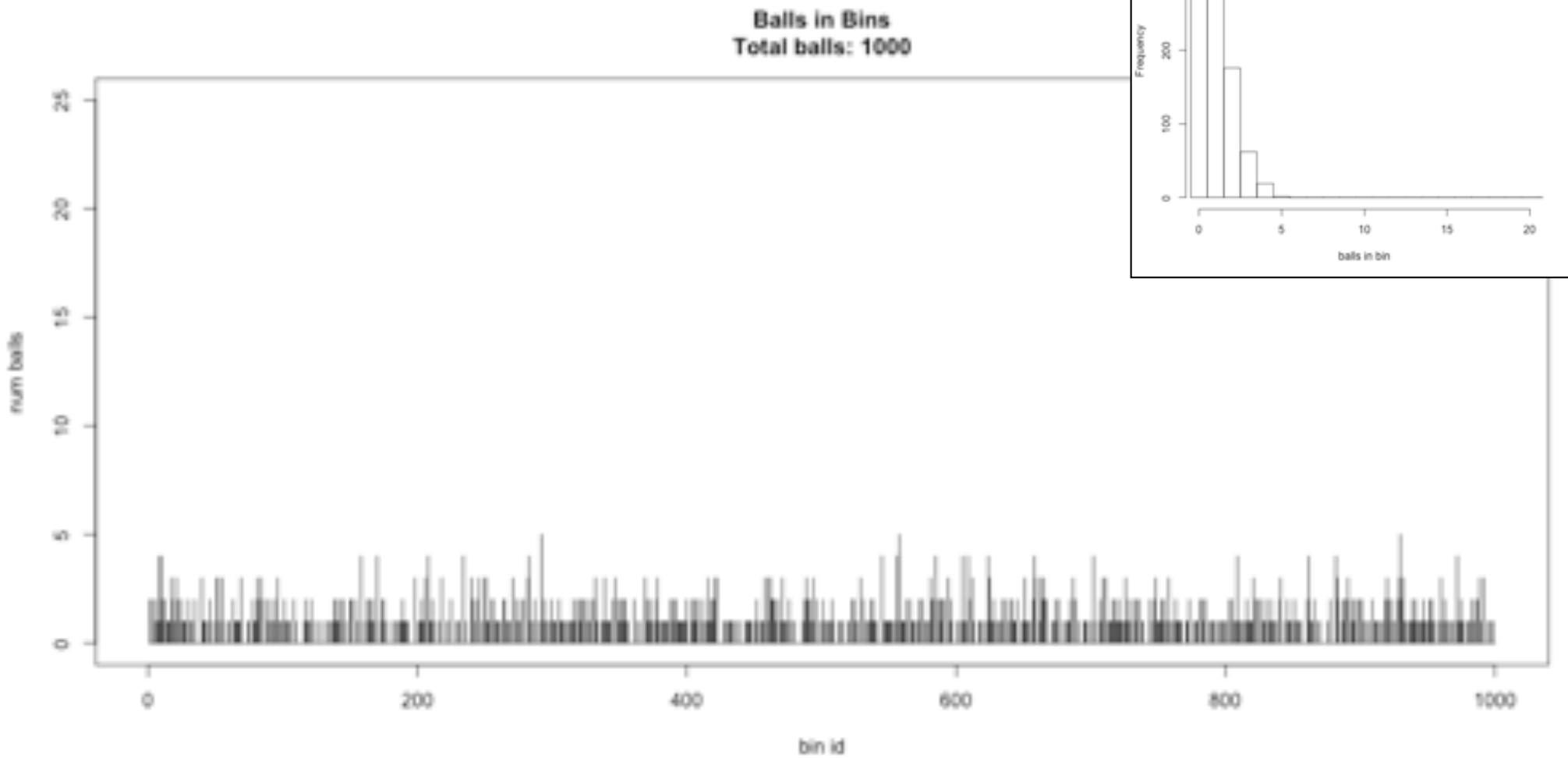


Imagine raindrops on a sidewalk

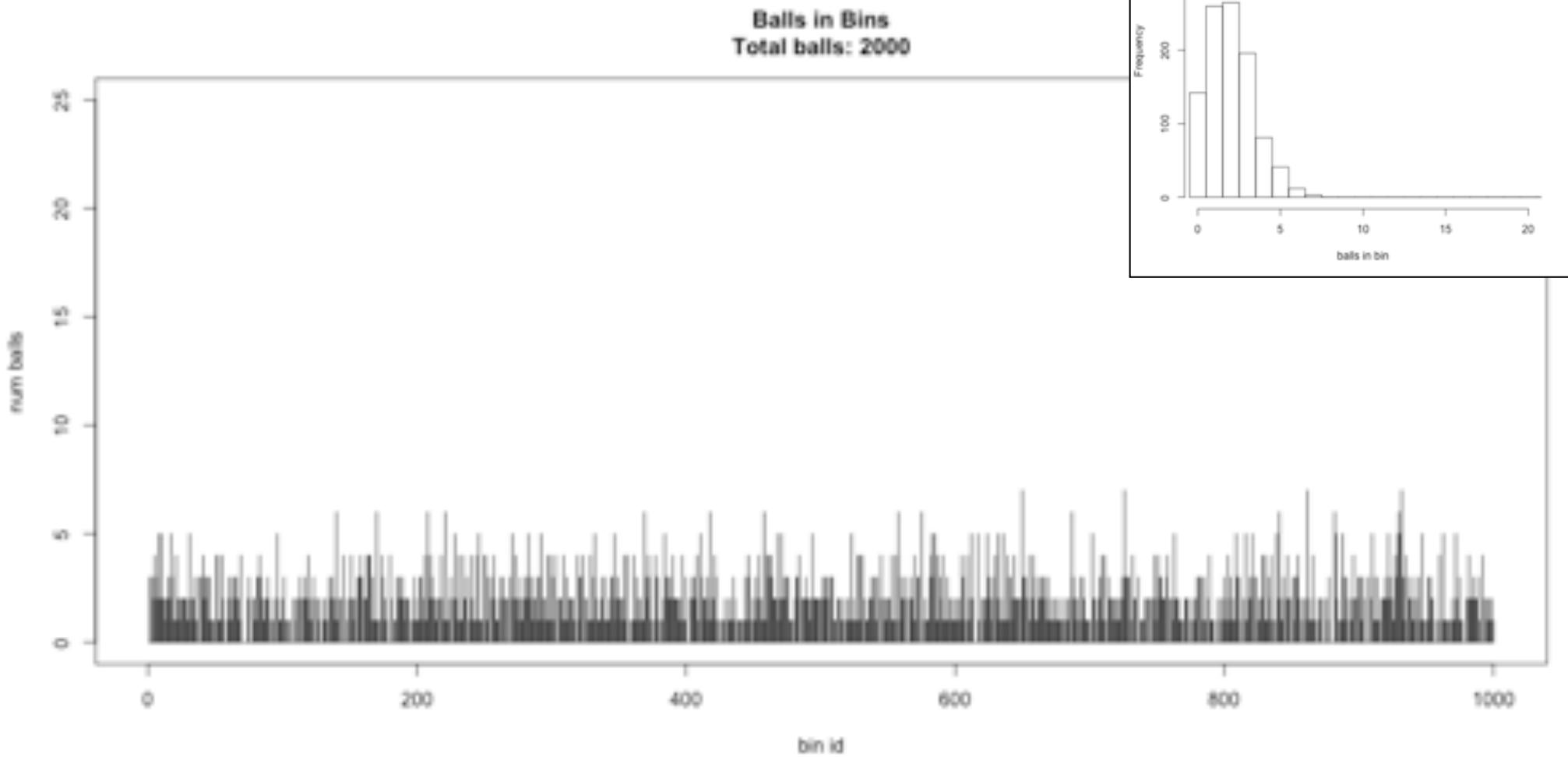
We want to cover the entire sidewalk but each drop costs \$1

If the genome is 10 Mbp, should we sequence 100k 100bp reads?

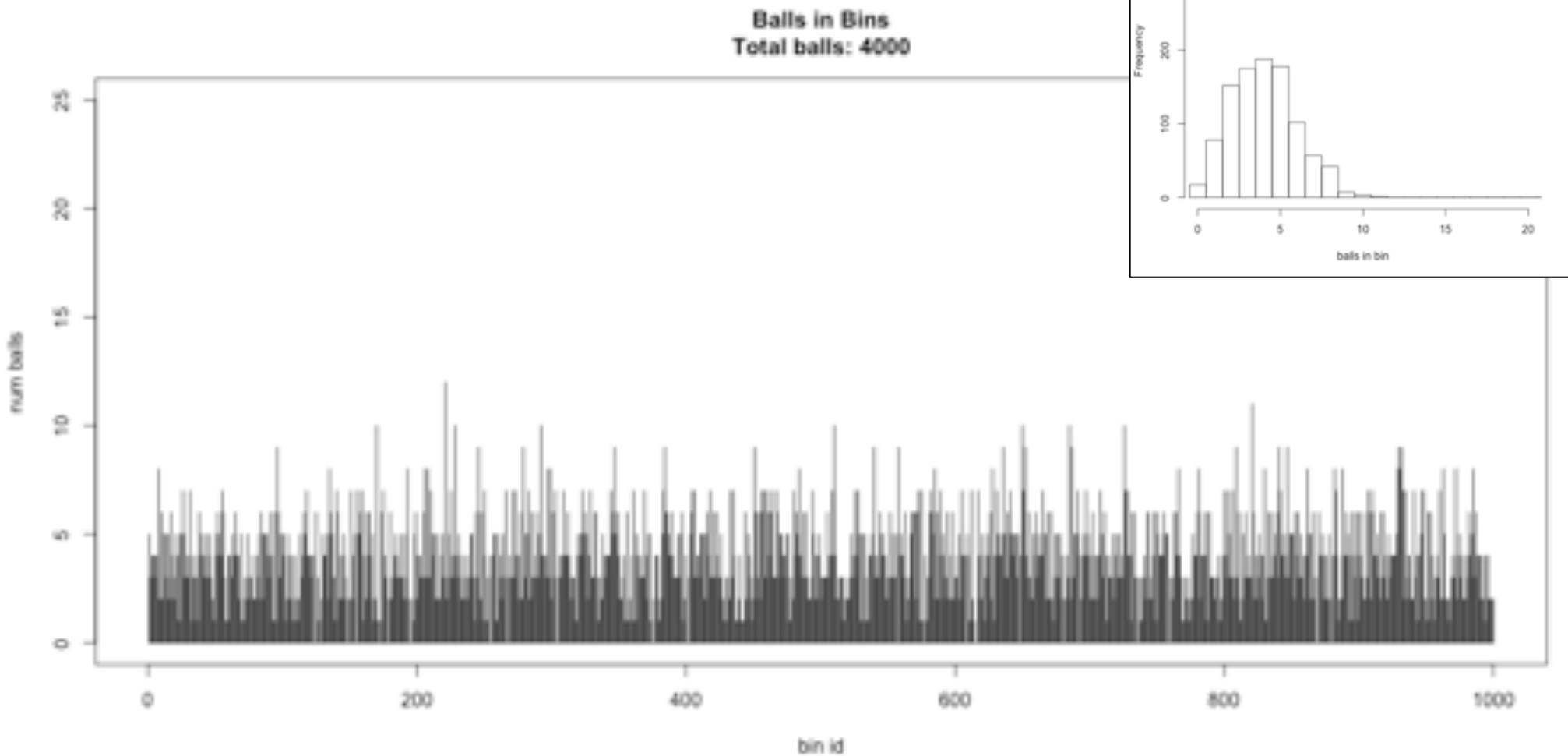
Ix sequencing



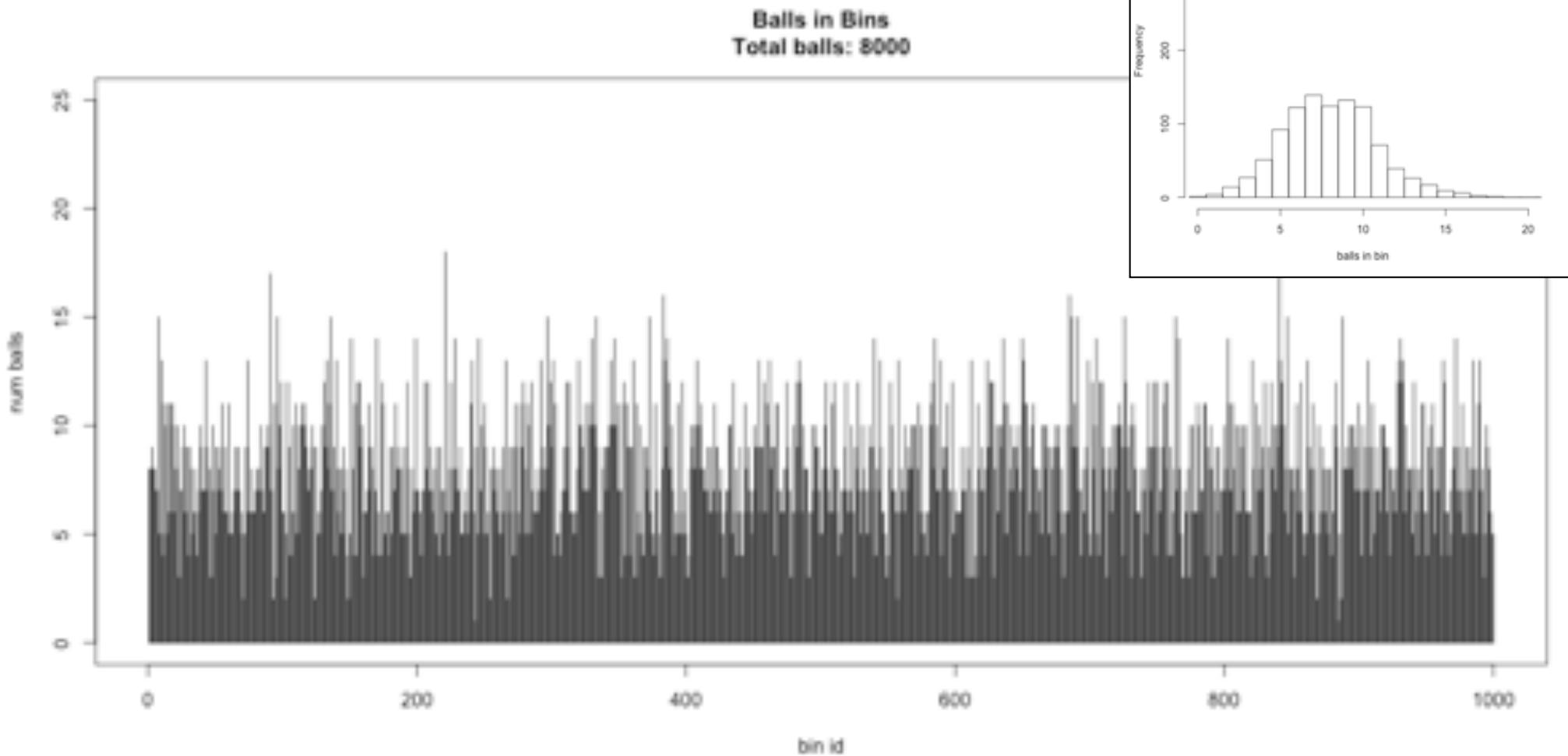
2x sequencing



4x sequencing



8x sequencing



Poisson Distribution

The probability of a given number of events occurring in a fixed interval of time and/or space if these events occur with a known average rate and independently of the time since the last event.

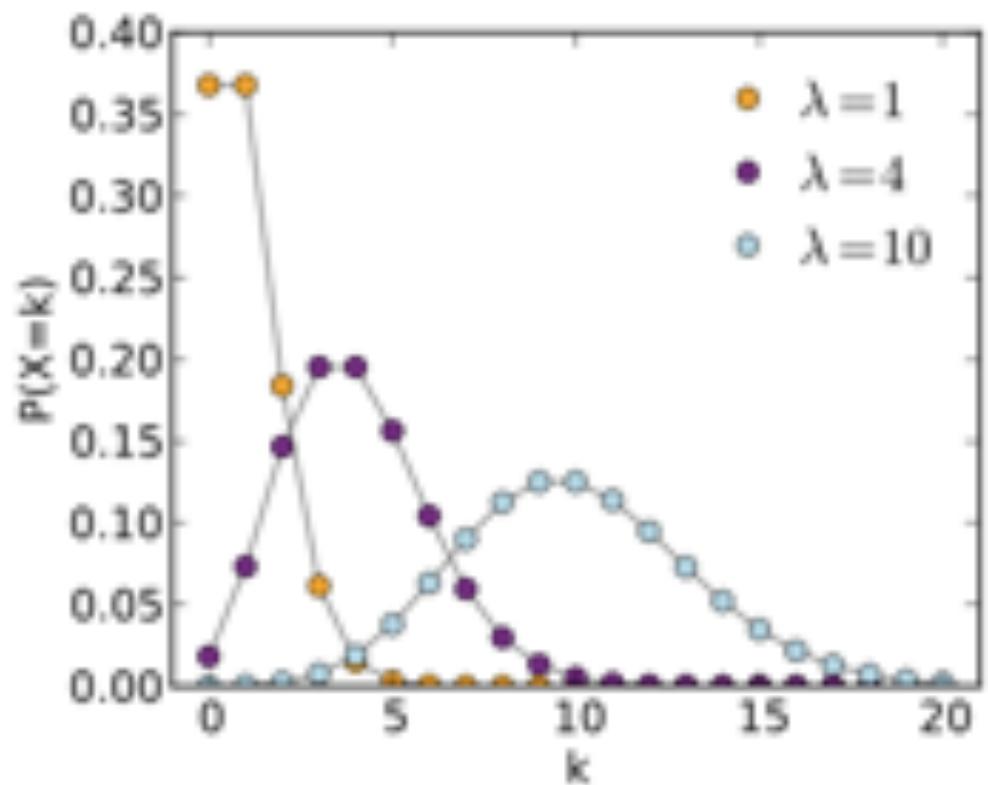
Formulation comes from the limit of the binomial equation

Resembles a normal distribution, but over the positive values, and with only a single parameter.

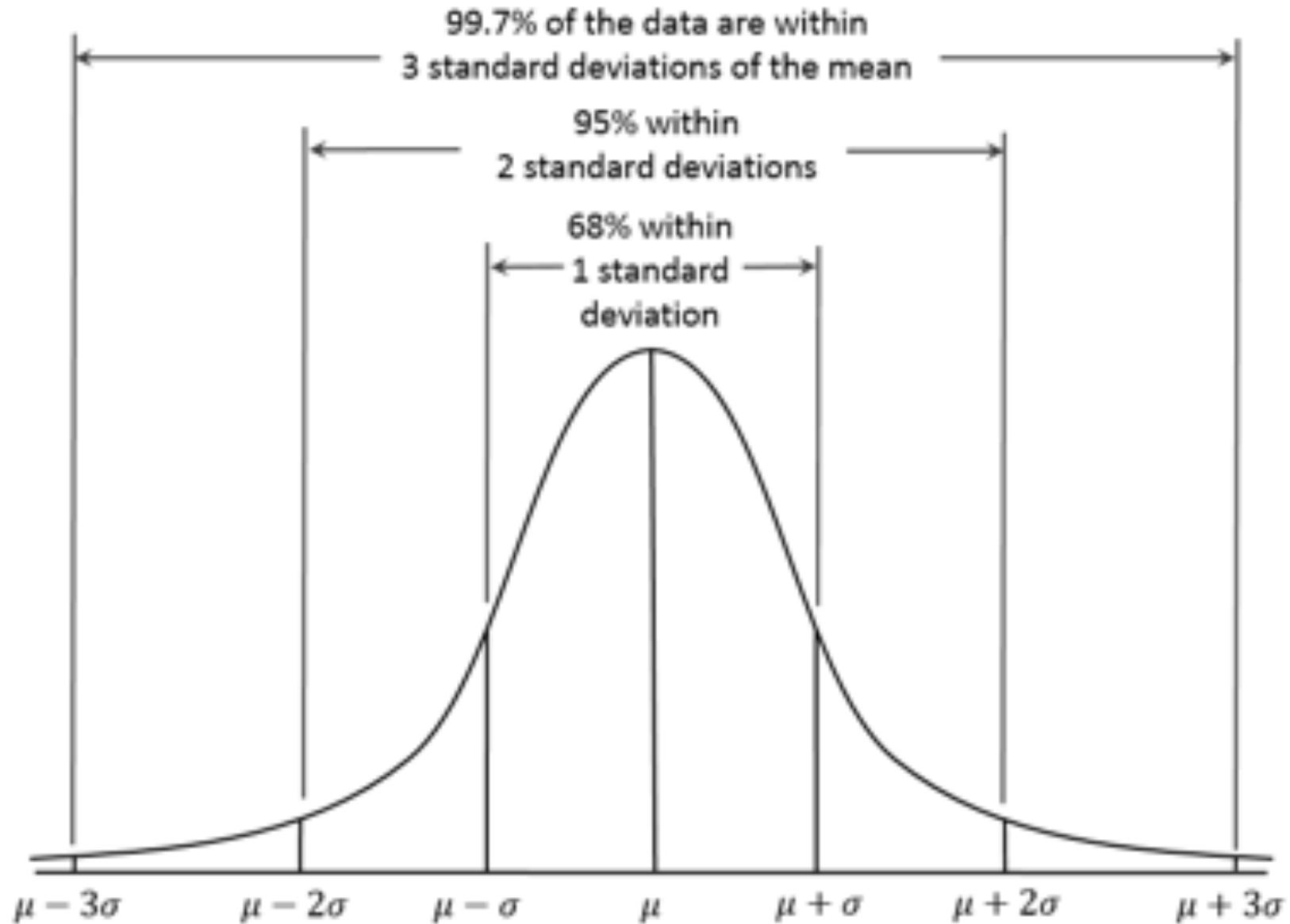
Key properties:

- ***The standard deviation is the square root of the mean.***
- ***For mean > 5, well approximated by a normal distribution***

$$P(k) = \frac{\lambda^k}{k!} e^{-\lambda}$$



Normal Approximation

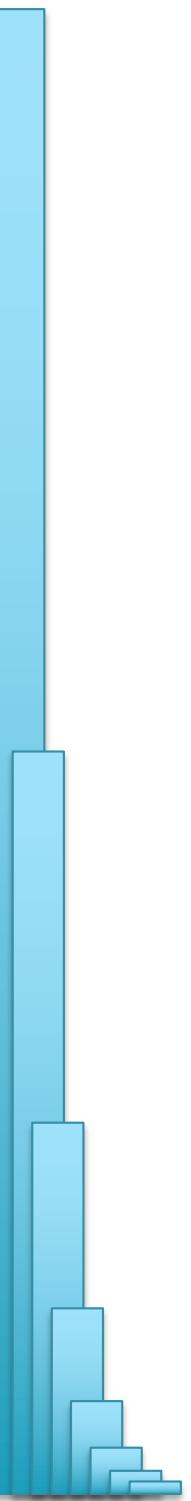


Can estimate Poisson distribution as a normal distribution when $\lambda > 10$

Pop Quiz!

I want to sequence a 10Mbp genome to 24x coverage.
How many 120bp reads do I need?

I want to sequence a 10Mbp genome so that
>97.5% of the genome has at least 24x coverage.
How many 120bp reads do I need?



Part 3: De novo genome assembly



Outline

1. ***Assembly theory***

- Assembly by analogy

2. ***Practical Issues***

- Coverage, read length, errors, and repeats

3. ***Next-next-gen Assembly***

- Canu: recommended for PacBio/ONT project

4. ***Whole Genome Alignment***

- MUMmer recommended

Shredded Book Reconstruction

- Dickens accidentally shreds the first printing of A Tale of Two Cities
 - Text printed on 5 long spools

It was the best of times, it was the worst of times, it was the age of wisdom; it was the age of foolishness, ...

It was the best of times, it was the worst of times, it was the age of wisdom; it was the age of foolishness, ...

It was the best of times, it was the worst of times, it was the age of wisdom; it was the age of foolishness, ...

It was the best of times, it was the worst of times, it was the age of wisdom; it was the age of foolishness, ...

It was the best of times, it was the worst of times, it was the age of wisdom; it was the age of foolishness, ...

- How can he reconstruct the text?

- 5 copies x 138,656 words / 5 words per fragment = 138k fragments
- The short fragments from every copy are mixed together
- Some fragments are identical

It was the best of
age of wisdom, it was
best of times, it was
it was the age of
it was the age of
it was the worst of
of times, it was the
of times, it was the
of wisdom, it was the
the age of wisdom, it
the best of times, it
the worst of times, it
times, it was the age
times, it was the worst
was the age of wisdom,
was the age of foolishness,
was the best of times,
was the worst of times,
wisdom, it was the age
worst of times, it was

Greedy Reconstruction

It was the best of
was the best of times,
the best of times, it
best of times, it was
of times, it was the
of times, it was the
times, it was the worst
times, it was the age

The repeated sequence make the correct reconstruction ambiguous

- It was the best of times, it was the [worst/age]

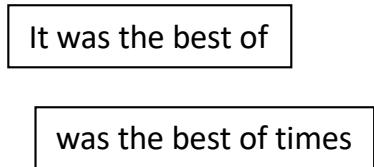
Model the assembly problem as a graph problem

How long will it take to compute the overlaps?

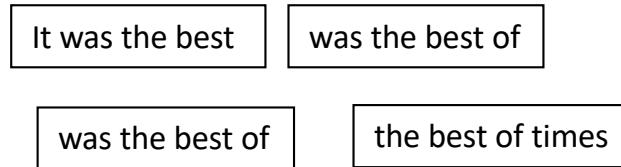
de Bruijn Graph Construction

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

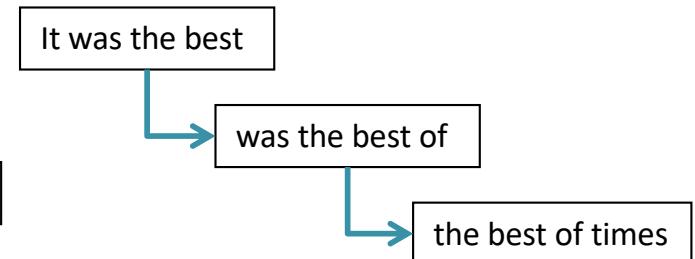
Fragments $|f|=5$



Sub-fragment $k=4$



Directed edges (overlap by $k-1$)



– Overlaps between fragments are implicitly computed

How to pronounce:

https://forvo.com/word/de_briuin/

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

it was the age

was the age of

the age of foolishness

the age of wisdom,

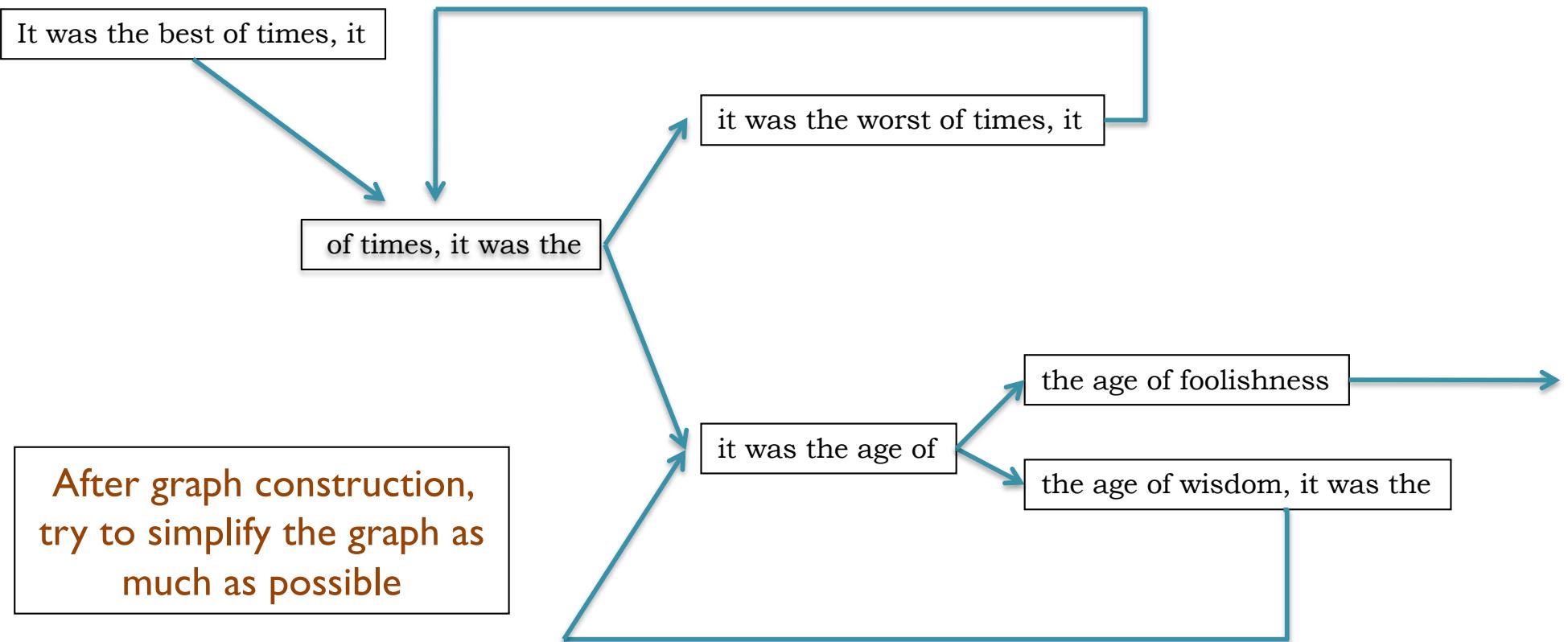
age of wisdom, it

of wisdom, it was

wisdom, it was the

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

de Bruijn Graph Assembly



The full tale

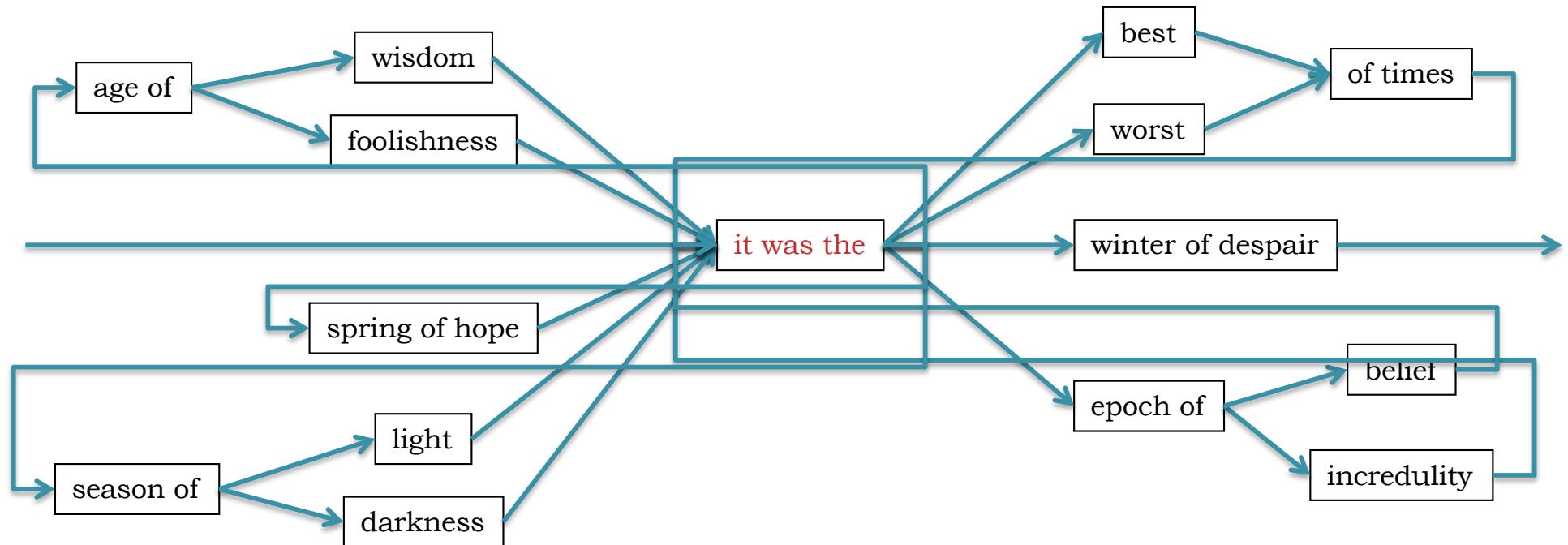
... it was the best of times it was the worst of times ...

... it was the age of wisdom it was the age of foolishness ...

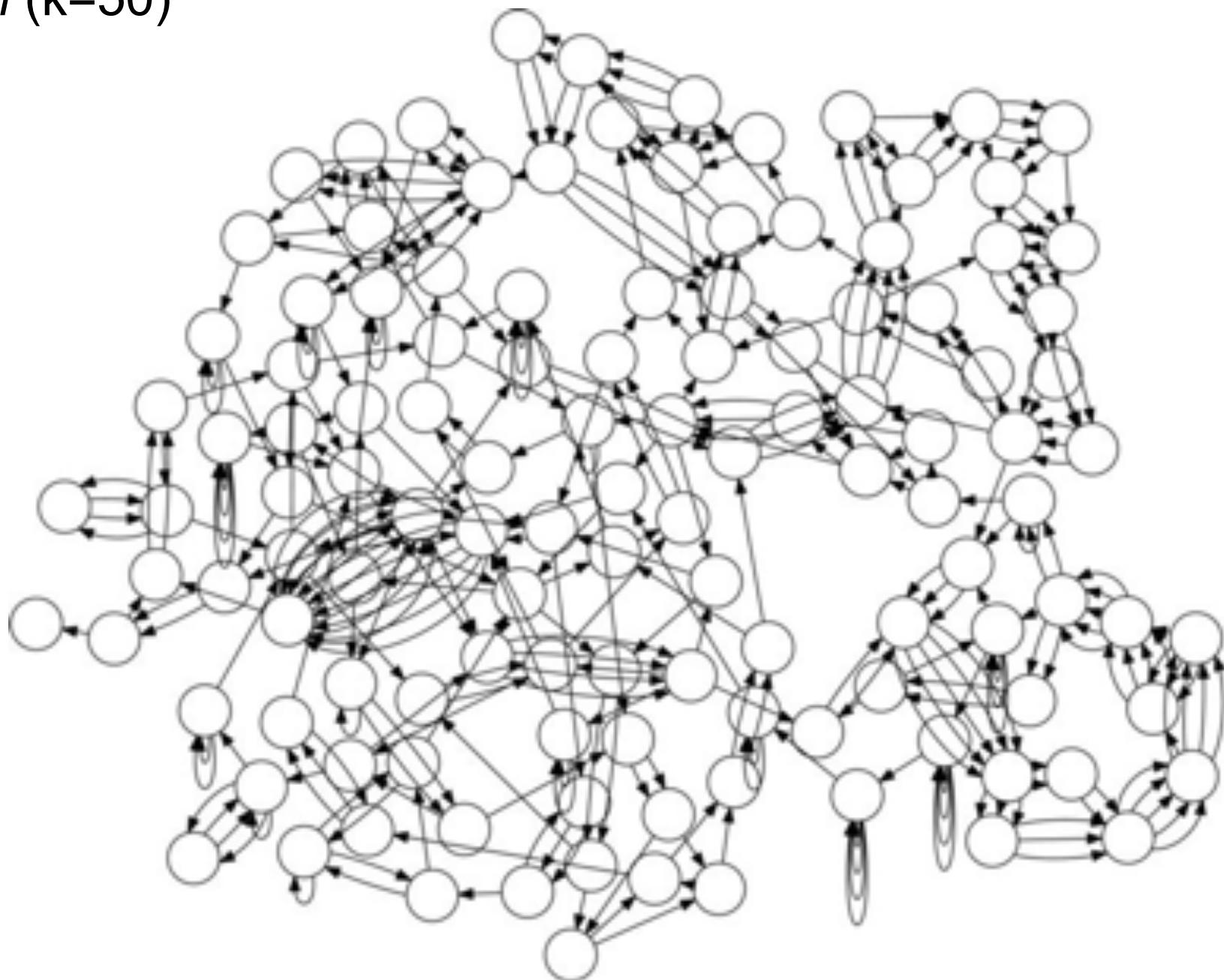
... it was the epoch of belief it was the epoch of incredulity ...

... it was the season of light it was the season of darkness ...

... it was the spring of hope it was the winter of despair ...

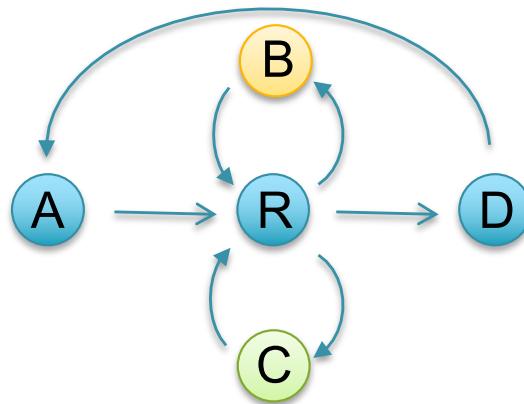


E. coli ($k=50$)



Reducing assembly complexity of microbial genomes with single-molecule sequencing
Koren et al (2013) Genome Biology. 14:R101 <https://doi.org/10.1186/gb-2013-14-9-r101>

Counting Eulerian Cycles



ARBRCRD
or
ARCRBRD

Generally an exponential number of compatible sequences

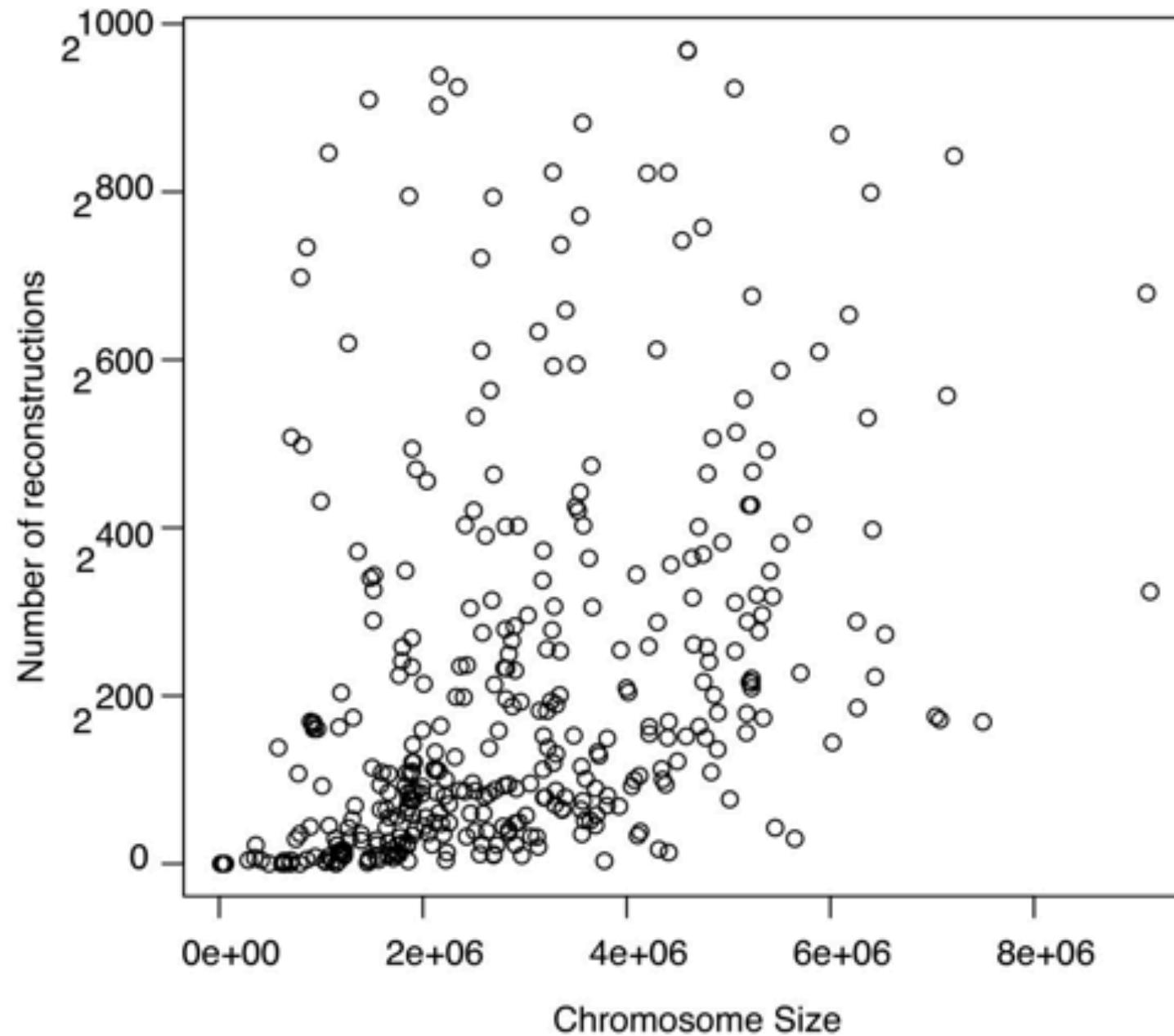
- Value computed by application of the BEST theorem (Hutchinson, 1975)

$$\mathcal{W}(G, t) = (\det L) \left\{ \prod_{u \in V} (r_u - 1)! \right\} \left\{ \prod_{(u,v) \in E} a_{uv}! \right\}^{-1}$$

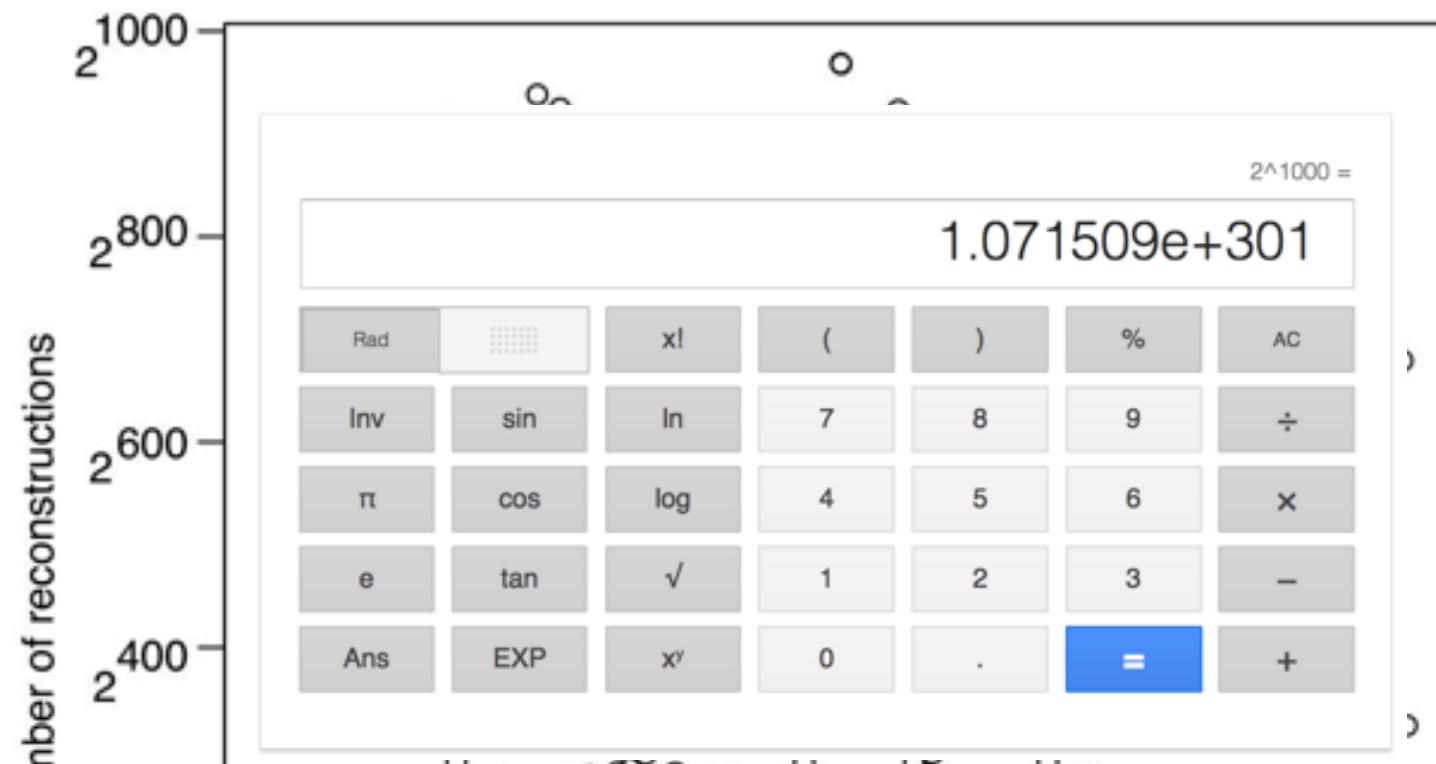
L = $n \times n$ matrix with $r_u - a_{uu}$ along the diagonal and $-a_{uv}$ in entry uv

$r_u = d^+(u) + 1$ if $u=t$, or $d^+(u)$ otherwise

a_{uv} = multiplicity of edge from u to v



Assembly Complexity of Prokaryotic Genomes using Short Reads.
Kingsford C, Schatz MC, Pop M (2010) *BMC Bioinformatics*.

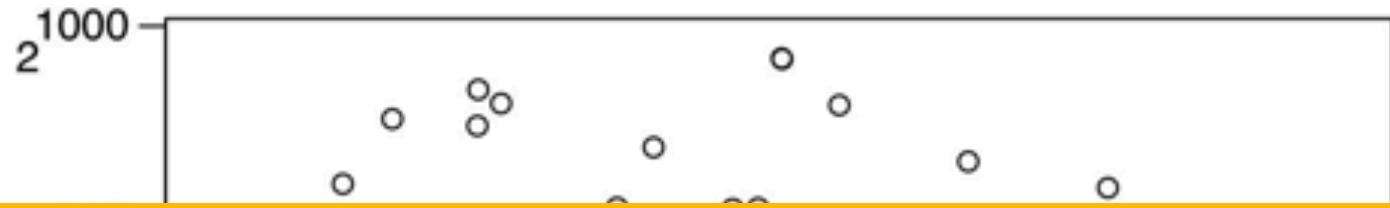


It is believed 74% of the mass of the Milky Way, for example, is in the form of hydrogen atoms. The Sun contains approximately **10⁵⁷ atoms** of hydrogen. If you multiple the number of atoms per star (10⁵⁷) times the estimated number of stars in the universe (10²³), you get a value of **10⁸⁰ atoms** in the known universe. Nov 5, 2017

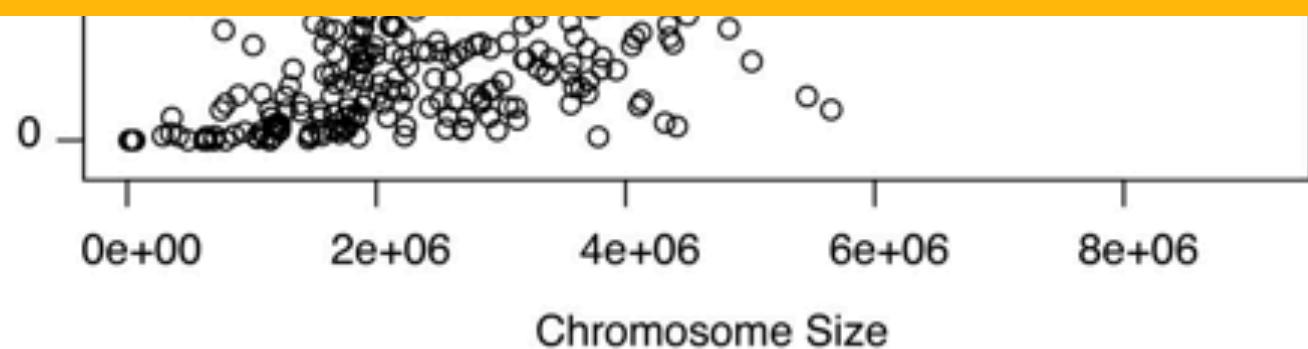


[How Many Atoms Are There in the Universe? - ThoughtCo](https://www.thoughtco.com/number-of-atoms-in-the-universe-603795)
<https://www.thoughtco.com/number-of-atoms-in-the-universe-603795>

Assembly Complexity of Prokaryotic Genomes using Short Reads.
Kingsford C, Schatz MC, Pop M (2010) *BMC Bioinformatics*.



- **Finding possible assemblies is easy!**
- **However, there is an astronomical genomic number of possible paths!**
- **Hopeless to figure out the whole genome/chromosome, figure out the parts that you can**

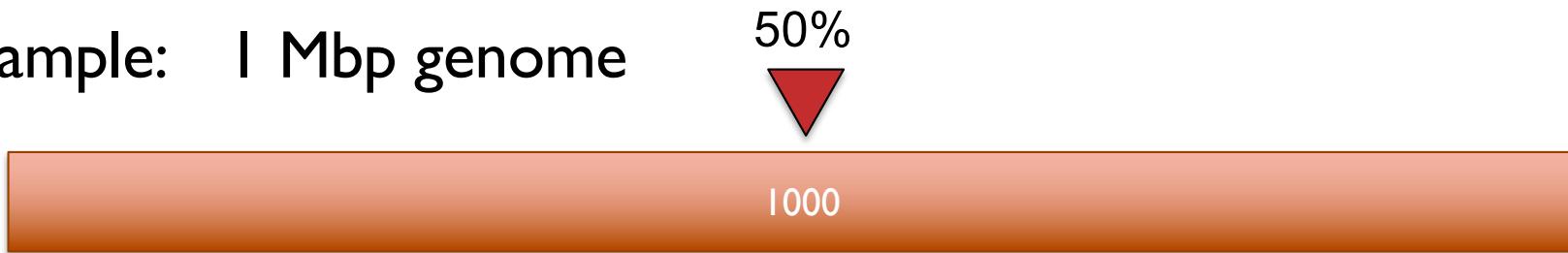


Assembly Complexity of Prokaryotic Genomes using Short Reads.
Kingsford C, Schatz MC, Pop M (2010) *BMC Bioinformatics*.

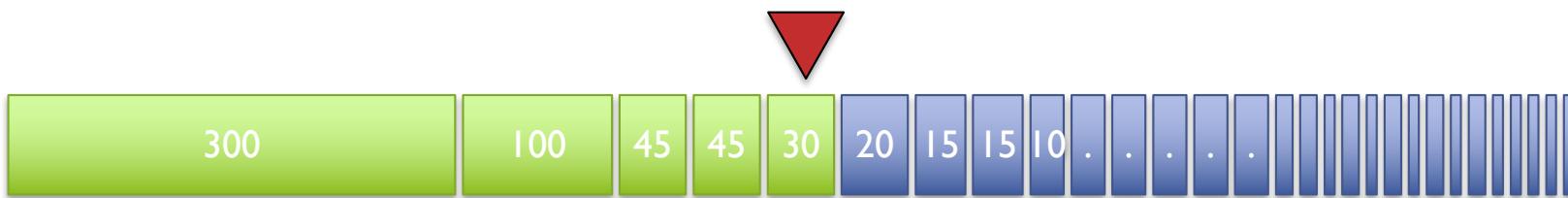
Contig N50

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 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

Pop Quiz I

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

ATTA

GATT

TACA

TTAC

Pop Quiz 2

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

ACGA

ACGT

ATAC

CGAC

CGTA

GACG

GTAT

TACG



Titus Brown

@ctitusbrown

Following



Wow, this could double as life philosophy, too!

Michael Schatz @mike_schatz

Replying to @ZaminIqbal @nomad421 and 4 others

Yep, very easy to find *a* path, very hard to find *the* path

11:40 AM - 22 Jan 2018

4 Retweets 17 Likes



2

4

17





Outline

1. *Assembly theory*

- Assembly by analogy

2. **Practical Issues**

- Coverage, read length, errors, and repeats

3. *Next-next-gen Assembly*

- Canu: recommended for PacBio/ONT project

4. Whole Genome Alignment

- MUMmer recommended

Assembly Applications

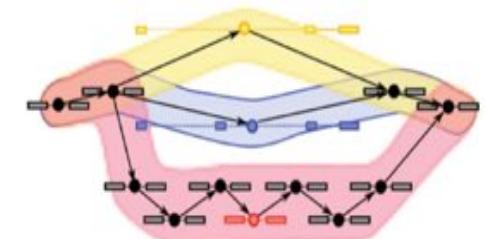
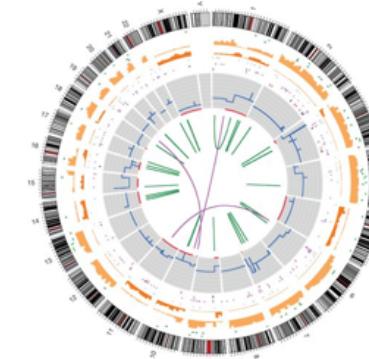
- Novel genomes



- 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

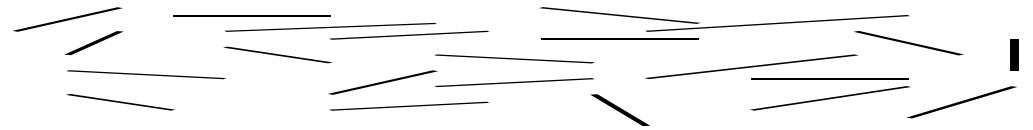
4. *Accuracy:*

- (Very) Hard to assess correctness



Assembling a Genome

I. Shear & Sequence DNA

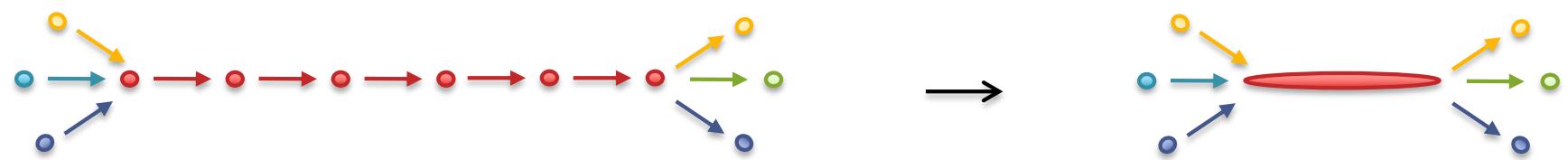


2. Construct assembly graph from reads (de Bruijn / overlap graph)

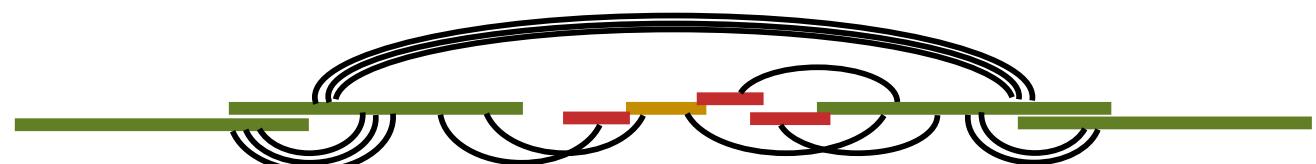
...AGCCTAG**GGATGCGCGACACGT**

GGATGCGCGACACGTCGCATATCCGGTTTGGT**CAACCTCGGACGGAC**
CAACCTCGGACGGACCTCAGCGAA...

3. Simplify assembly graph

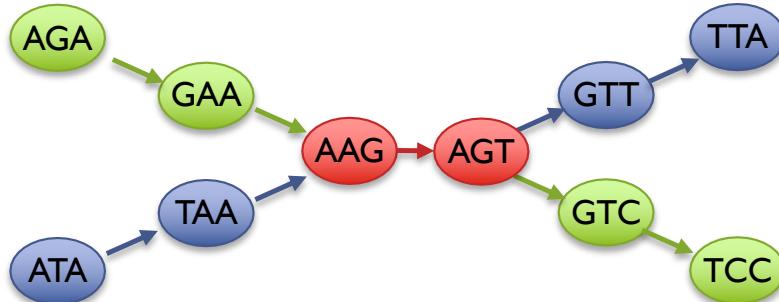


4. Detangle graph with long reads, mates, and other links



Two Paradigms for Assembly

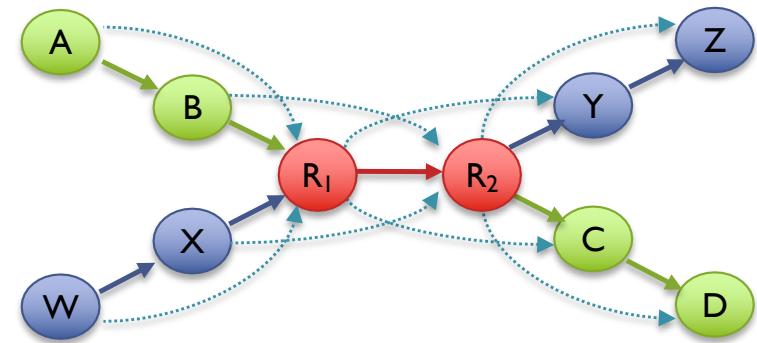
de Bruijn Graph



Short read assemblers

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

Overlap Graph



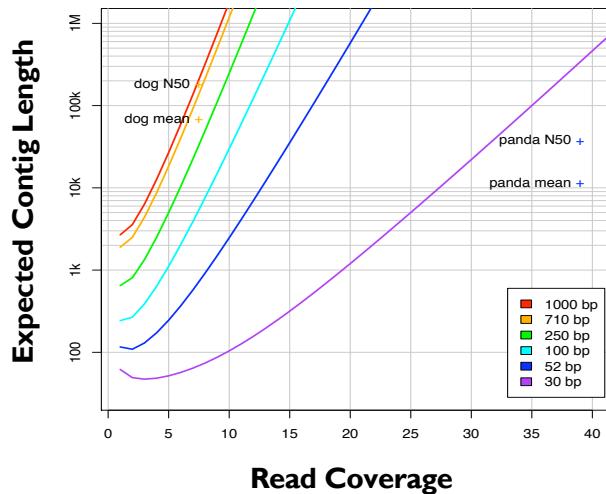
Long read assemblers

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

Assembly of Large Genomes using Second Generation Sequencing
Schatz MC, Delcher AL, Salzberg SL (2010) *Genome Research*. 20:1165-1173.

Ingredients for a good assembly

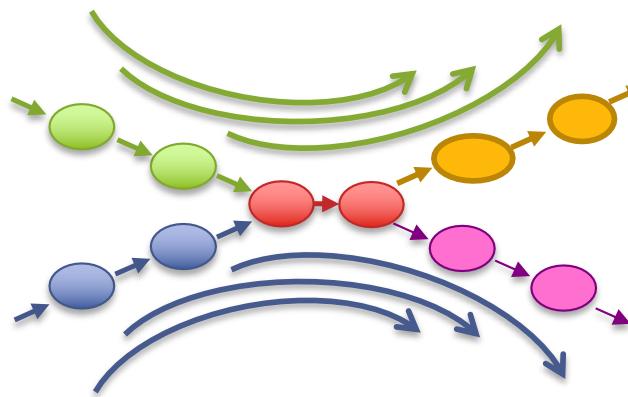
Coverage



High coverage is required

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

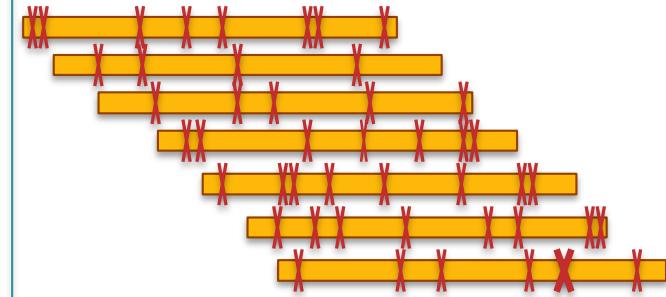
Read Length



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



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	

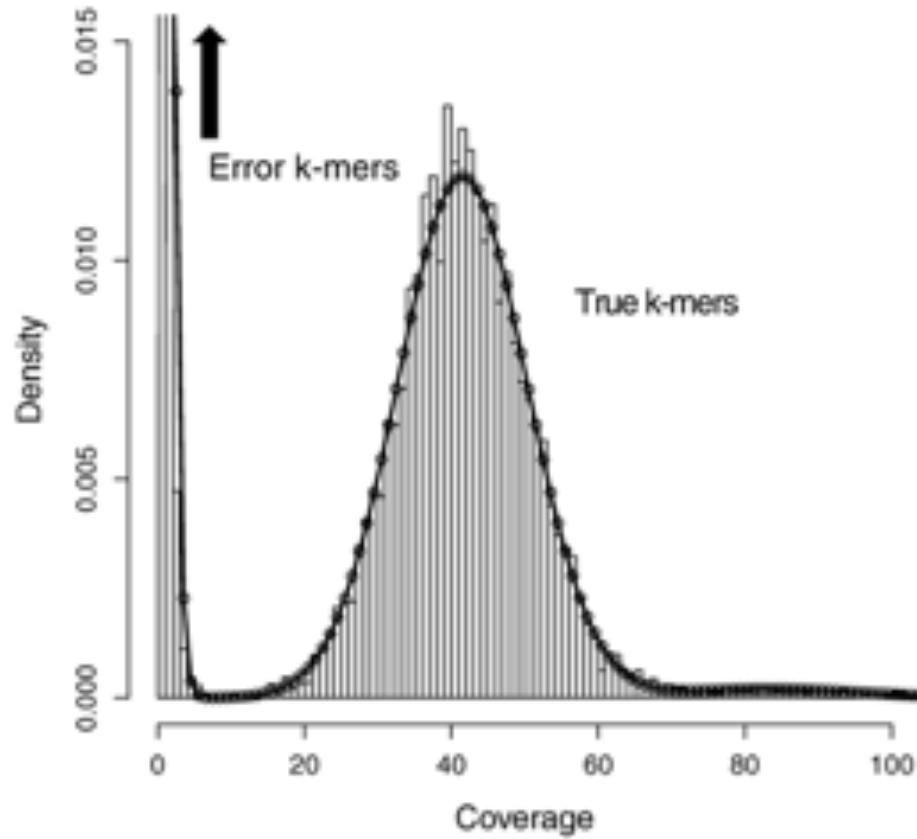
3 kmers occur 1x
3 kmers occur 2x
2 kmers occur 3x

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

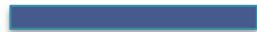
Sequencing read
from homologous
chromosome 1A



Sequencing read
from homologous
chromosome 1B



K-mer counting in heterozygous genomes



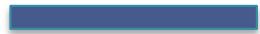
Sequencing read
from homologous
chromosome 1A



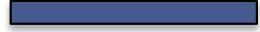
Sequencing read
from homologous
chromosome 1B



K-mer counting in heterozygous genomes



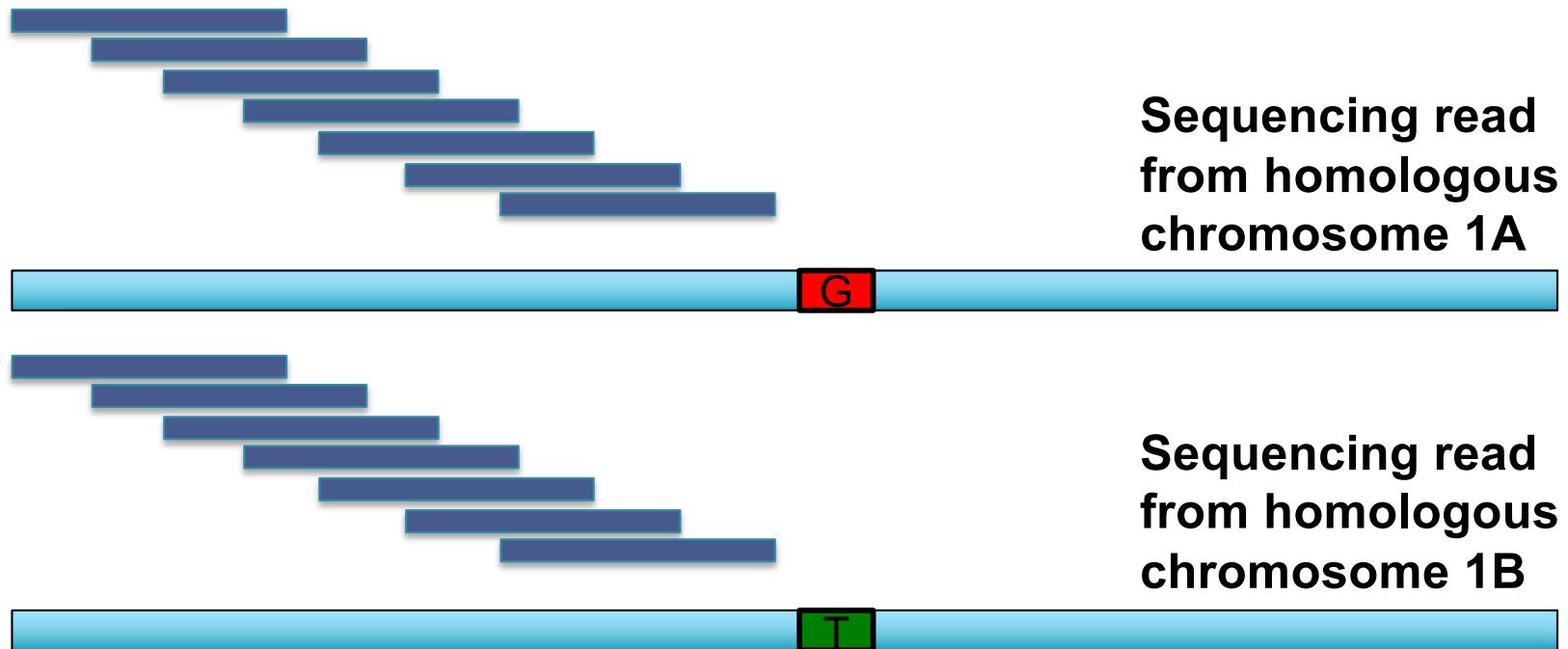
Sequencing read
from homologous
chromosome 1A



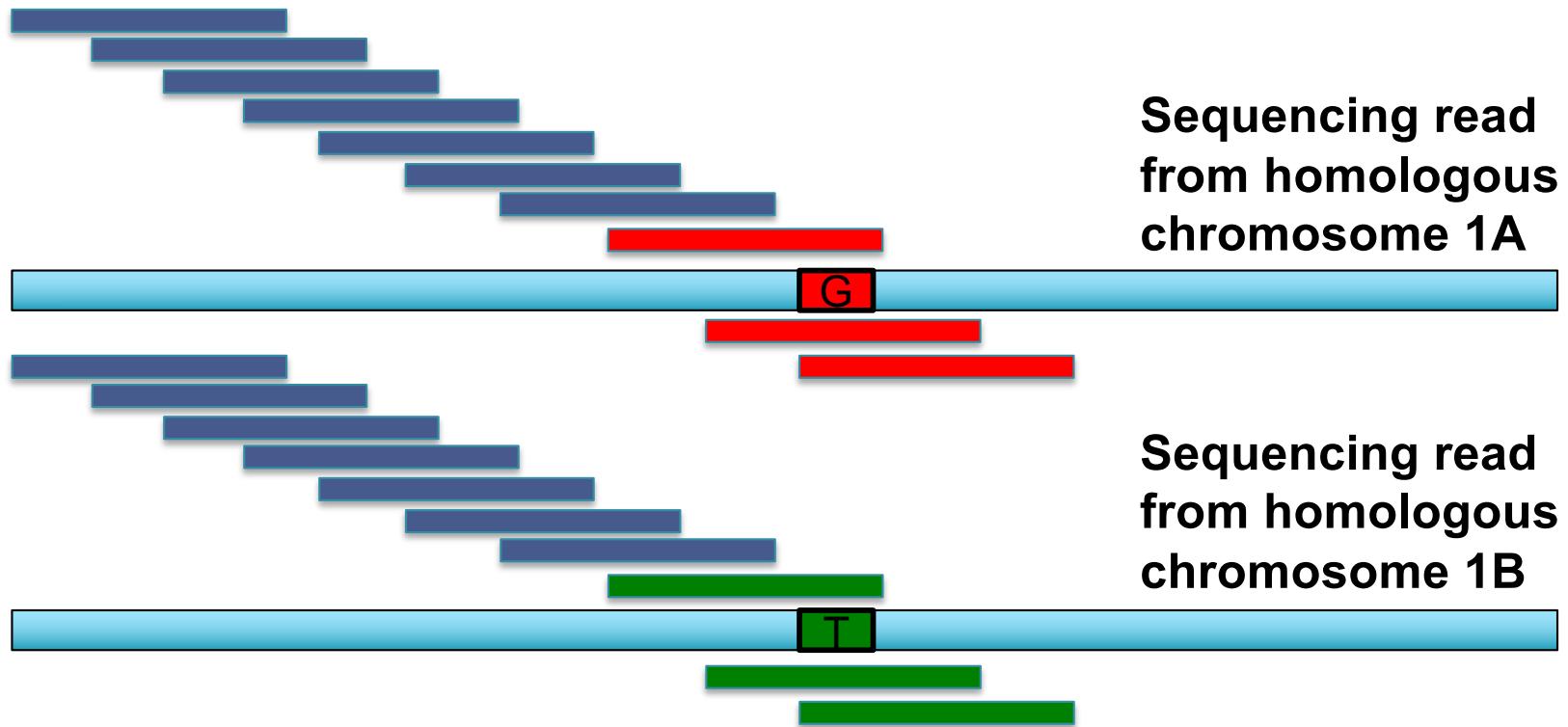
Sequencing read
from homologous
chromosome 1B



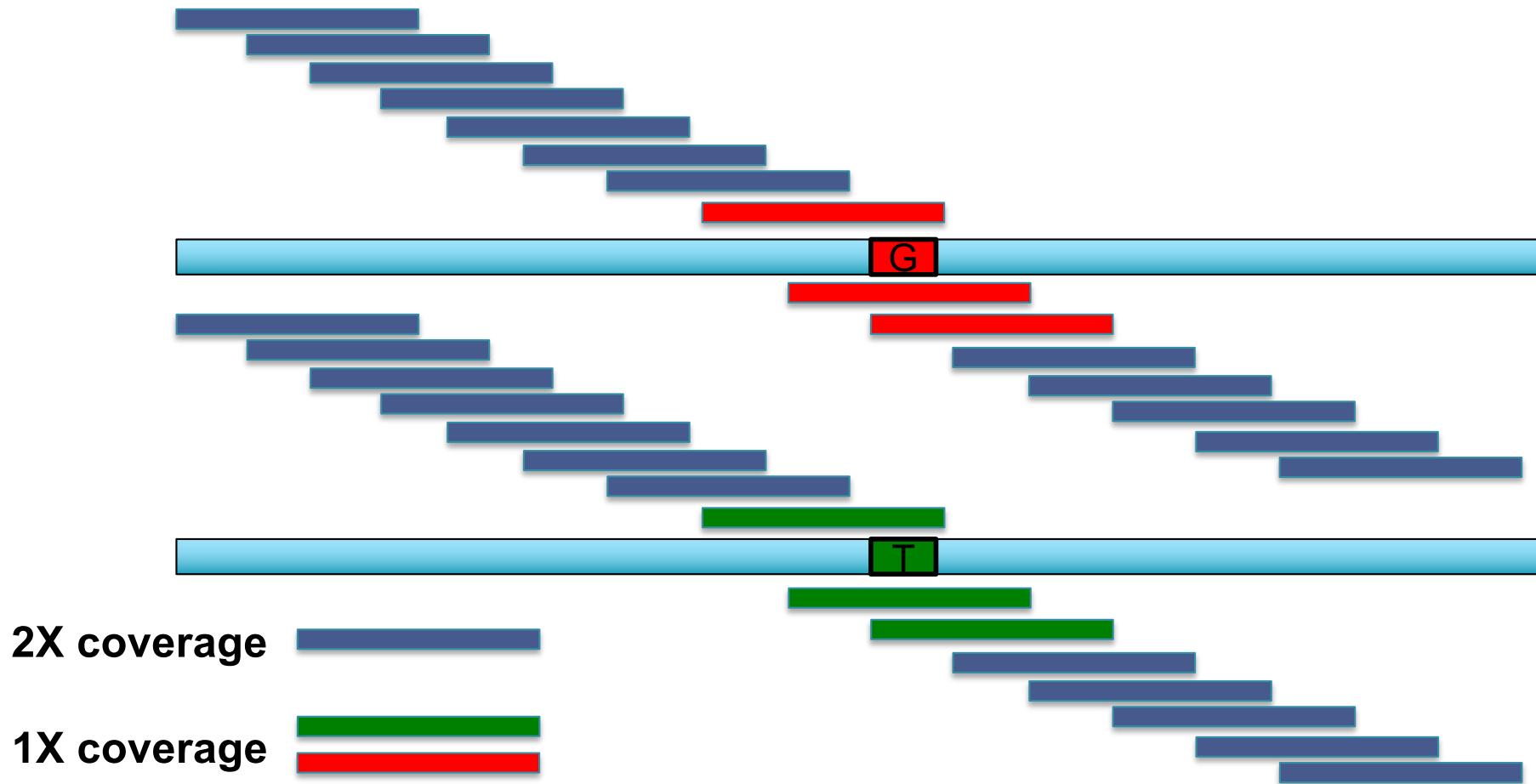
K-mer counting in heterozygous genomes



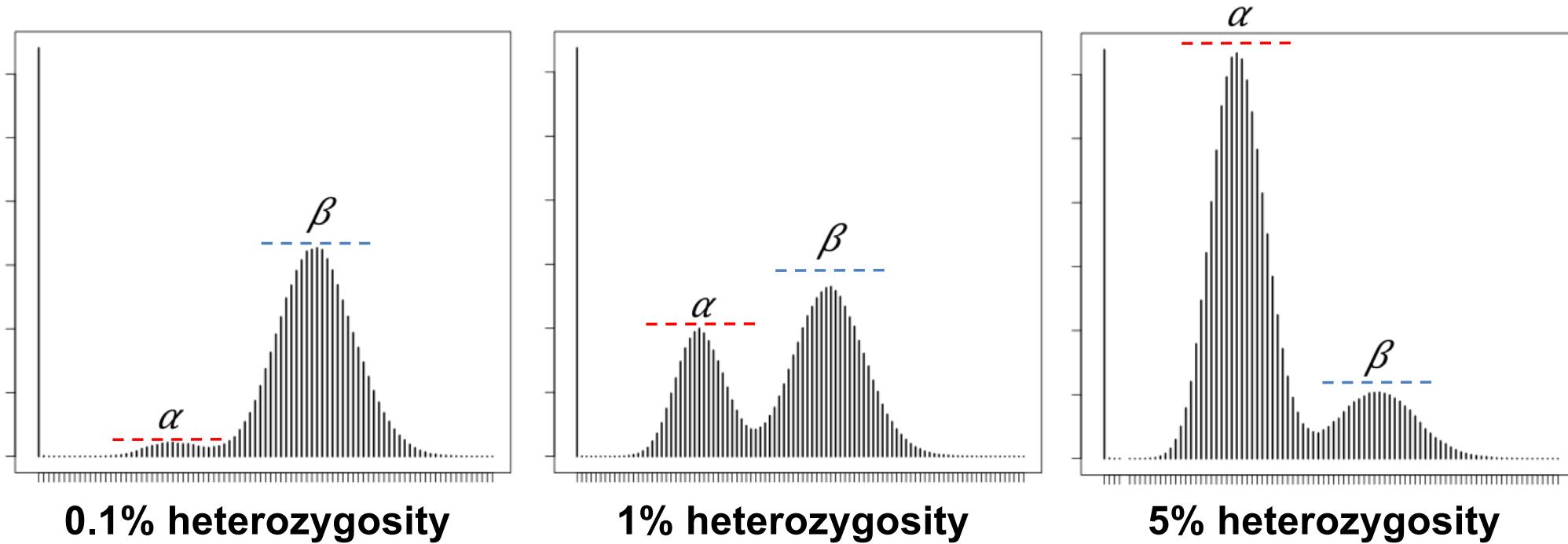
K-mer counting in heterozygous genomes



K-mer counting in heterozygous genomes



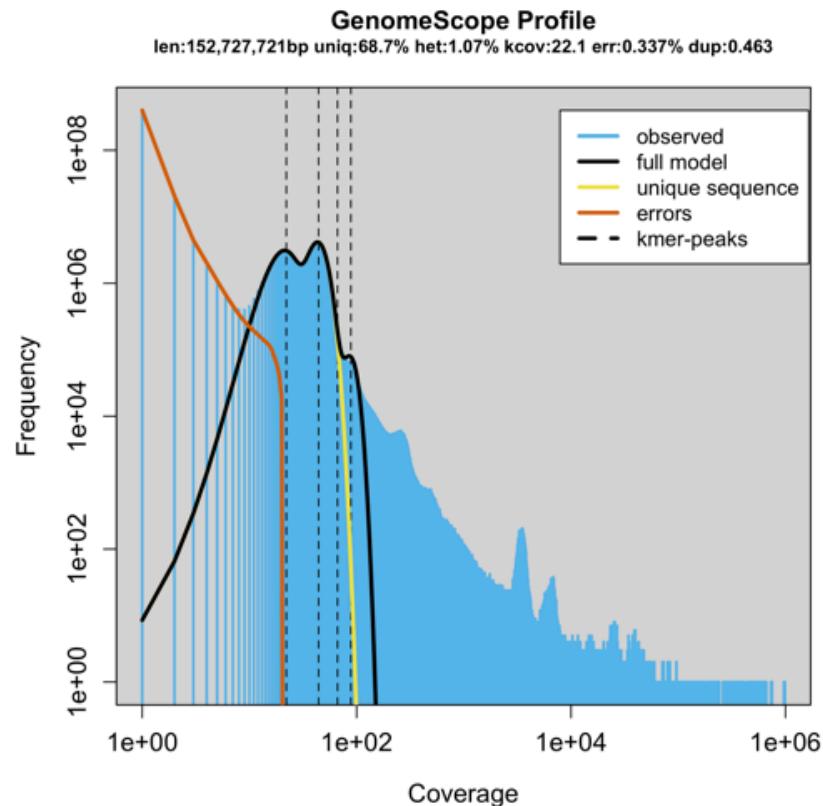
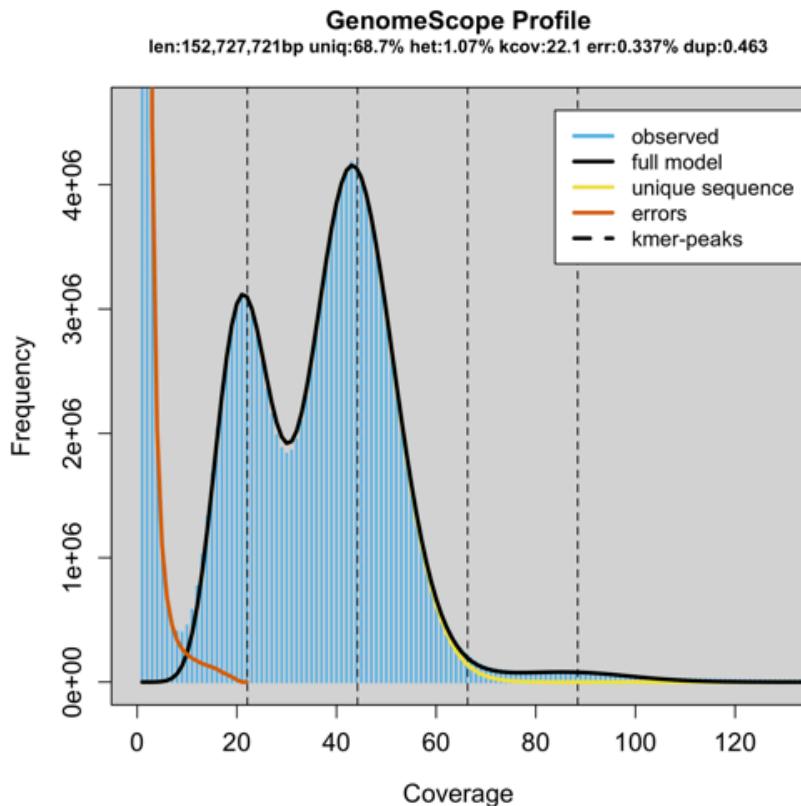
Heterozygous Kmer Profiles



- ***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 $k = 21$)

GenomeScope: Fast genome analysis from short reads

<http://genomescope.org>



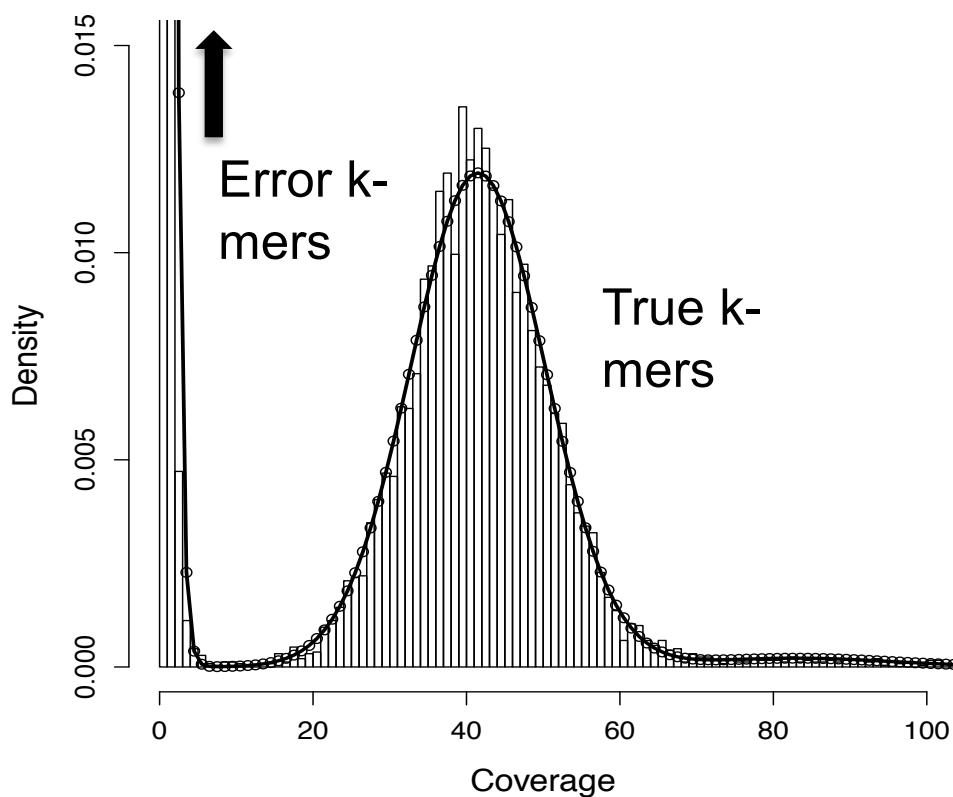
- Theoretical model agrees well with published results:
 - Rate of heterozygosity is higher than reported by other approaches but likely correct.
 - Genome size of plants inflated by organelle sequences (exclude very high freq. kmers)

Vulture, GW*, Sedlazeck FJ*, et al. (2017) *Bioinformatics*
Ranallo-Benavidez, TR. et al. (2020) *Nature Communication*

Error Correction with Quake

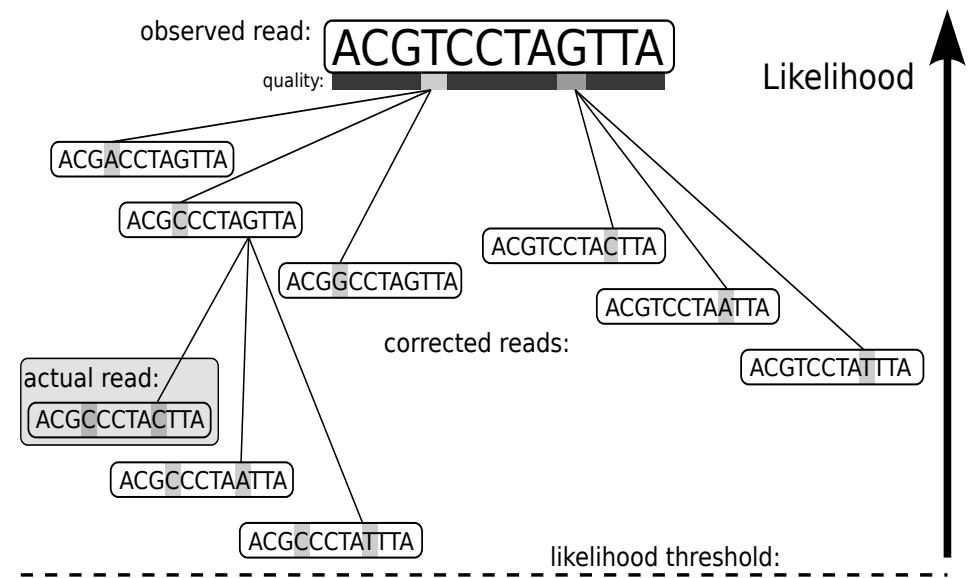
1. Count all “Q-mers” in reads

- Fit coverage distribution to mixture model of errors and regular coverage
- Automatically determines threshold for trusted k-mers



2. Correction Algorithm

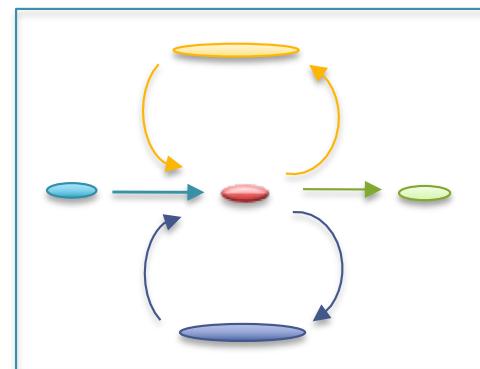
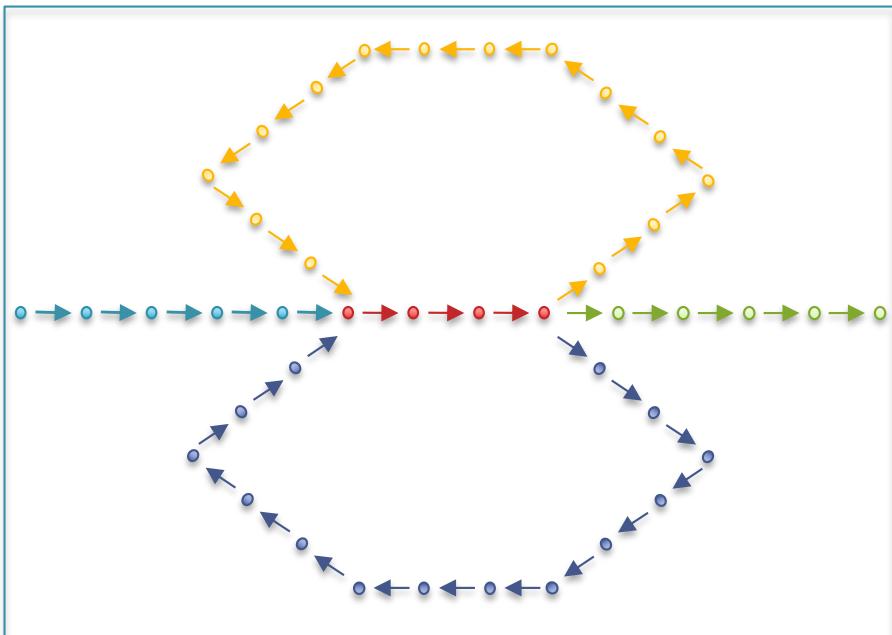
- Considers editing erroneous kmers into trusted kmers in decreasing likelihood
- Includes quality values, nucleotide/nucleotide substitution rate



Quake: quality-aware detection and correction of sequencing reads.
Kelley, DR, Schatz, MC, Salzberg SL (2010) *Genome Biology*. 11:R116

Unitigging / Unipathing

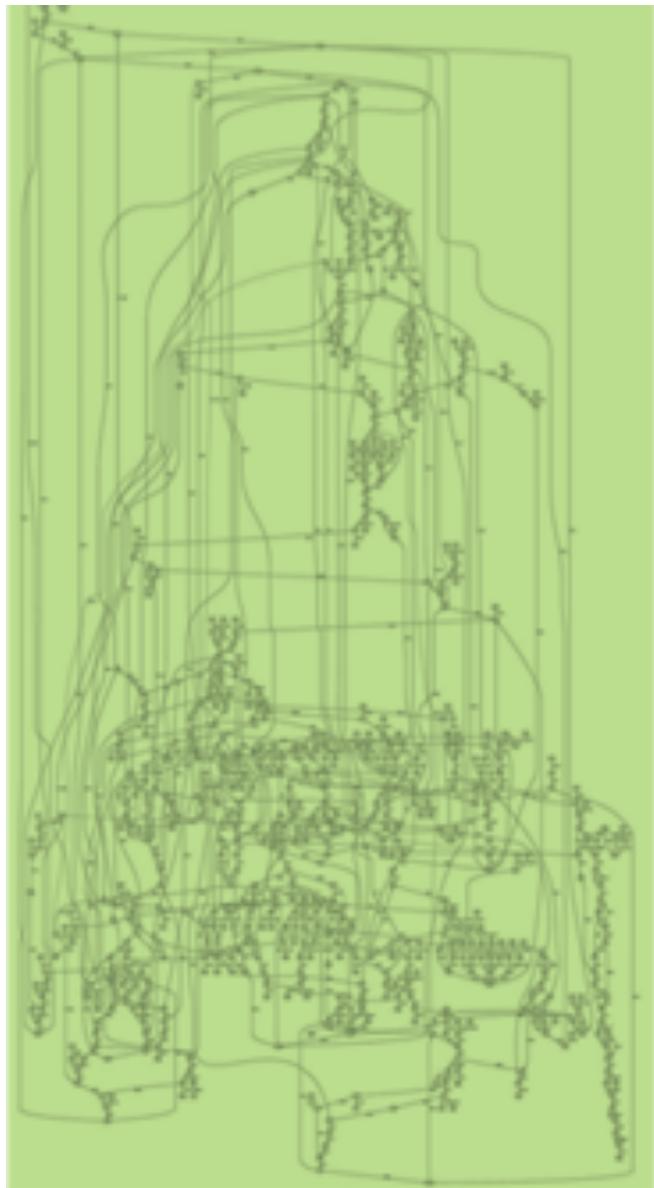
- After simplification and correction, compress graph down to its non-branching initial contigs
 - Aka “unitigs”, “unipaths”



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

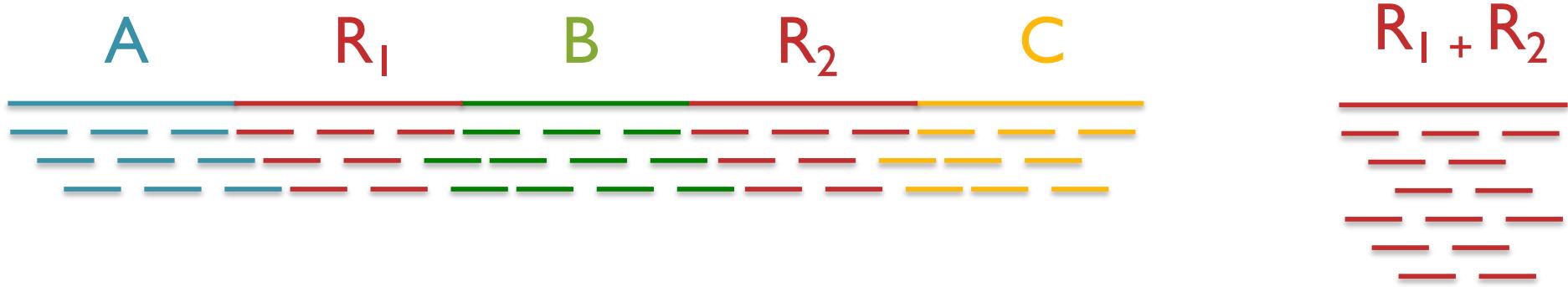
times,

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) <i>Mariner</i> 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

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

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

The fragment assembly string graph

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

Paired-end and Mate-pairs

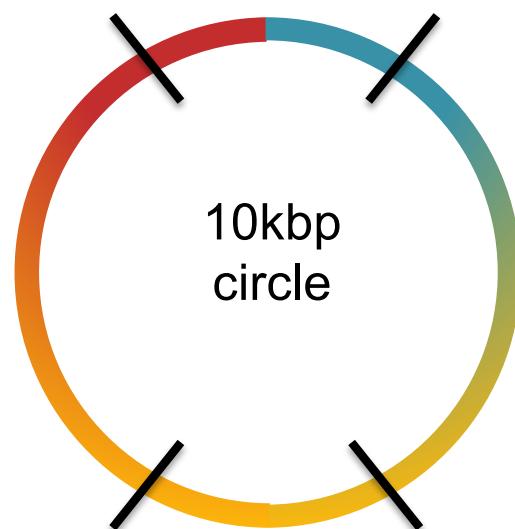
Paired-end sequencing

- 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)

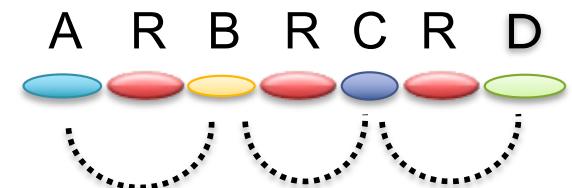
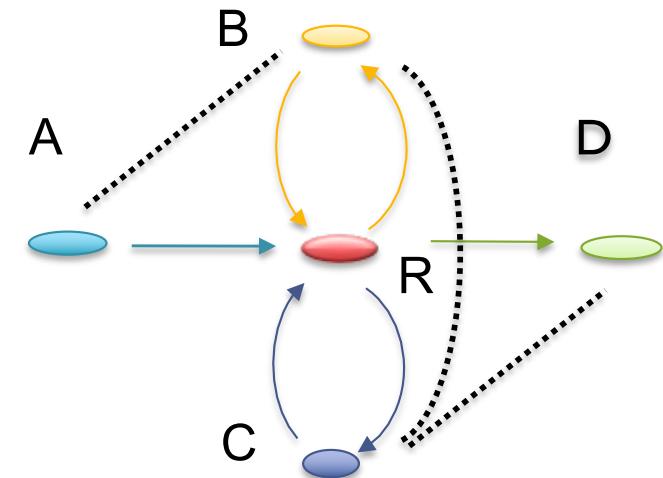


2x100 @ 300bp (innies)



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?

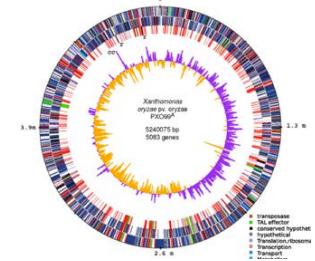
Assemblathon Results

ID	Overall	CPNG50	SPNG50	Struct.	CC50	Subs.	Copy. Num.	Cov. Tot.	Cov. CDS
BGI	36	★					★	★	★
Broad	37	★	★	★	★				
WTSI-S	46		★	★	★	★			
CSHL	52	★							★
BCCGSC	53						★	★	
DOEJGI	56		★	★	★	★			
RHUL	58								

- SOAPdenovo and ALLPATHS came out neck-and-neck followed closely behind by SGA, Celera Assembler, ABySS
- My recommendation for “typical” short read assembly is to use ALLPATHS or Spades

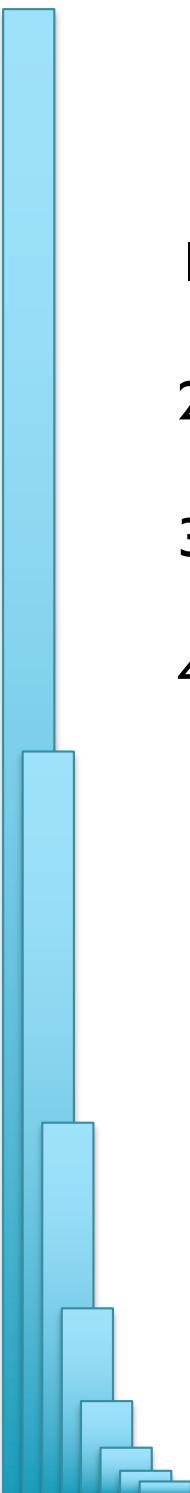
Assemblathon I: A competitive assessment of de novo short read assembly methods
Earl et al. (2011) Genome Research. 21: 2224-2241

Assembly Summary



Assembly quality depends on

1. **Coverage:** low coverage is mathematically hopeless
 2. **Repeat composition:** high repeat content is challenging
 3. **Read length:** longer reads help resolve repeats
 4. **Error rate:** errors reduce coverage, obscure true overlaps
-
- 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
 - Watch out for collapsed repeats & other misassemblies
 - Globally/Locally reassemble data from scratch with better parameters & stitch the 2 assemblies together



Next Steps

1. Reflect on the magic and power of DNA 😊
2. Check out the course webpage
3. Register on Piazza
4. Work on Assignment I
 1. Set up Linux, set up Virtual Machine
 2. Set up Dropbox for yourself!
 3. Get comfortable on the command line