

Human Genome & Genome Assembly

Michael Schatz

Sept 25, 2017 – Lecture 7

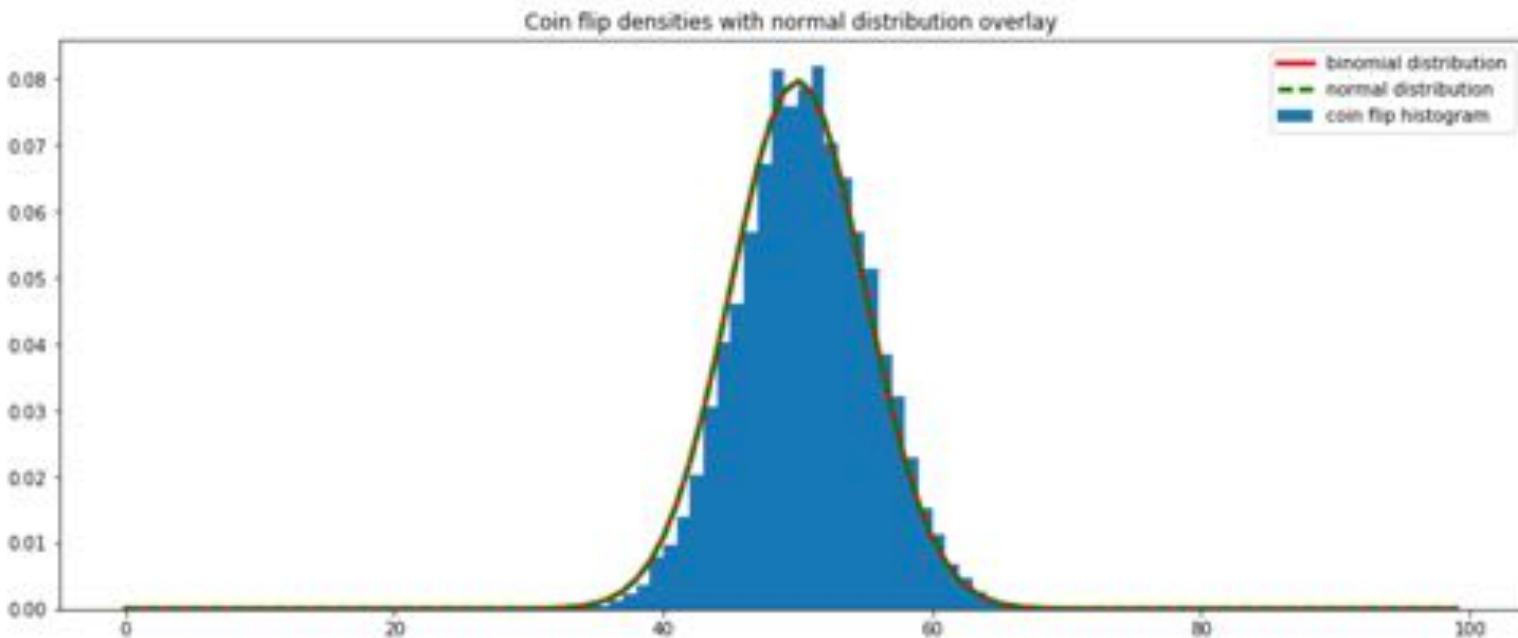
EN.601.452 Computational Biomedical Research

AS.020.415 Advanced Biomedical Research

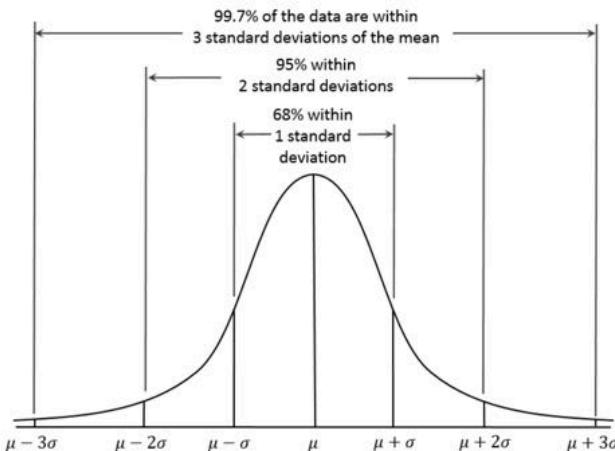




Exercise I: Due tonight @ 11:59pm

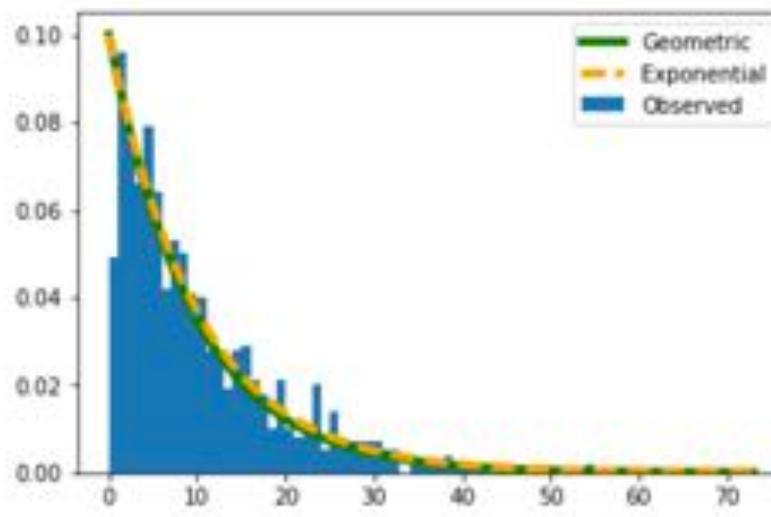


The binomial distribution is closely related to the normal distribution (aka Gaussian distribution)



The probability density for the Gaussian distribution is

$$p(x) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\mu)^2}{2\sigma^2}},$$





Project Proposal: Due in 1 week

Project Proposal

Assignment Date: September 25, 2017

Due Date: Monday, October 2, 2017 @ 11:59pm

Please email a PDF of your project proposal (1/2 to 1 page) to "jhbiomedicalresearch at gmail dot com" by 11:59pm on Monday October 2, 2017.

The proposal should have the following components:

- Short title for your proposal
- Your name
- Email addresses
- Description of what you hope to do and how you will do it:
 - What is the key question you hope to address?
 - What data will you use to study it? Are all the data you need available now? When will they be available?
 - What techniques (experimental or computational) will you use to generate and analyze the data?
 - What are the desired results?
- References to relevant papers
- References/URLs to datasets that you will be studying (Note you can also use simulated data)

After submitting your proposal, we will schedule a time to discuss your proposal, especially to ensure you have access to the data that you need and you have a clear path forward for analyzing it. The sooner that you submit your proposal, the sooner we can schedule the meeting. No late days can be used for this project.

Later, you will present your project in class at the end of the semester. You will also submit a written report (5-7 pages) of your project, formatting as scholarly article with separate sections for the Abstract, Introduction, Methods, Results, Discussion, and References. More details will be provided later in the semester.

***Focus on the problems where you have a unique advantage
and will have impact on the world***



The human genome



“Without a doubt, this is the most important, most wondrous map ever produced by humankind.”

*Bill Clinton
June 26, 2000*

The human genome





Outline

1. ***Assembly theory***

- Assembly by analogy

2. ***Practical Issues***

- Coverage, read length, errors, and repeats

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

- PacBio/ONT projects

Genome Assembly by Analogy

A TALE OF TWO CITIES in Three Books

BOOK THE FIRST. RECALLED TO LIFE

CHAPTER I THE PERIOD

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; we had everything before us, we had nothing before us, we were all going direct to Heaven, we were all going direct the other way—in short, the period was so far present period, that some of its noisiest authorities

had received, for good or for evil, in the super-

a large jaw and a queen smile.

there were a king and a

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

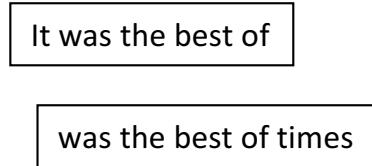
- It was the best of times, it was the [worst/age]

Model the assembly problem as a graph problem

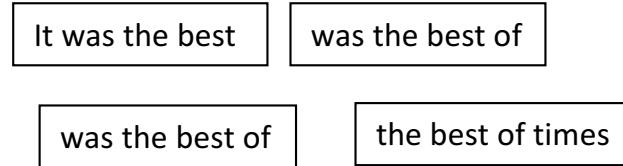
de Bruijn Graph Construction

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

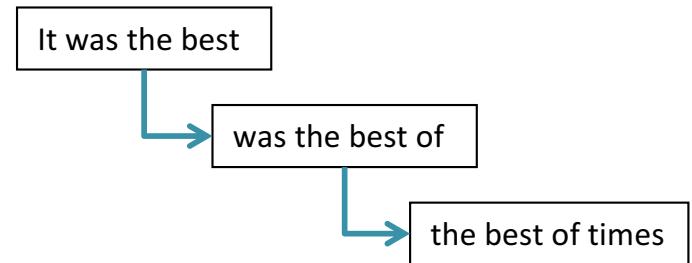
Fragments $|f|=5$



Sub-fragment $k=4$



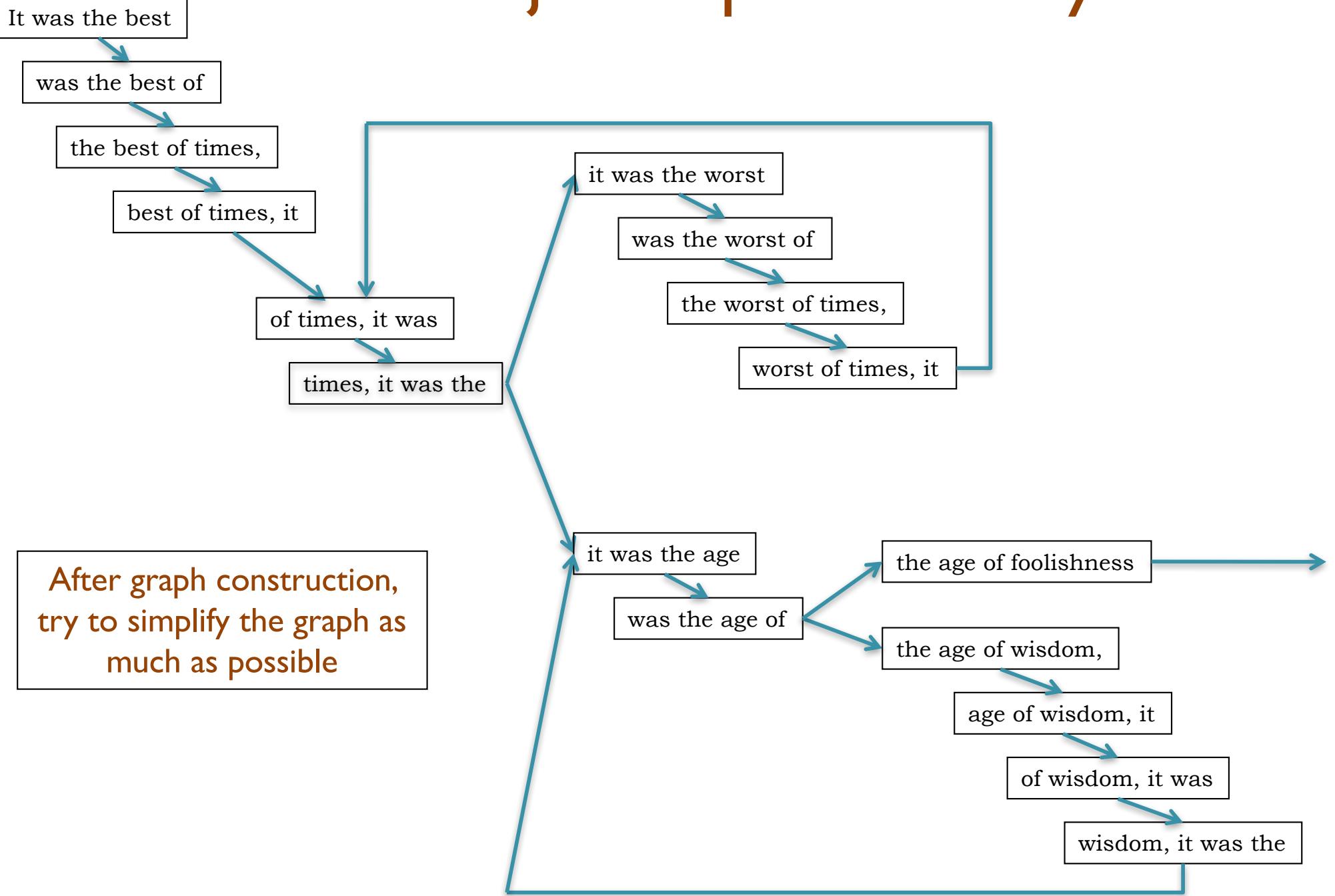
Directed edges (overlap by $k-1$)



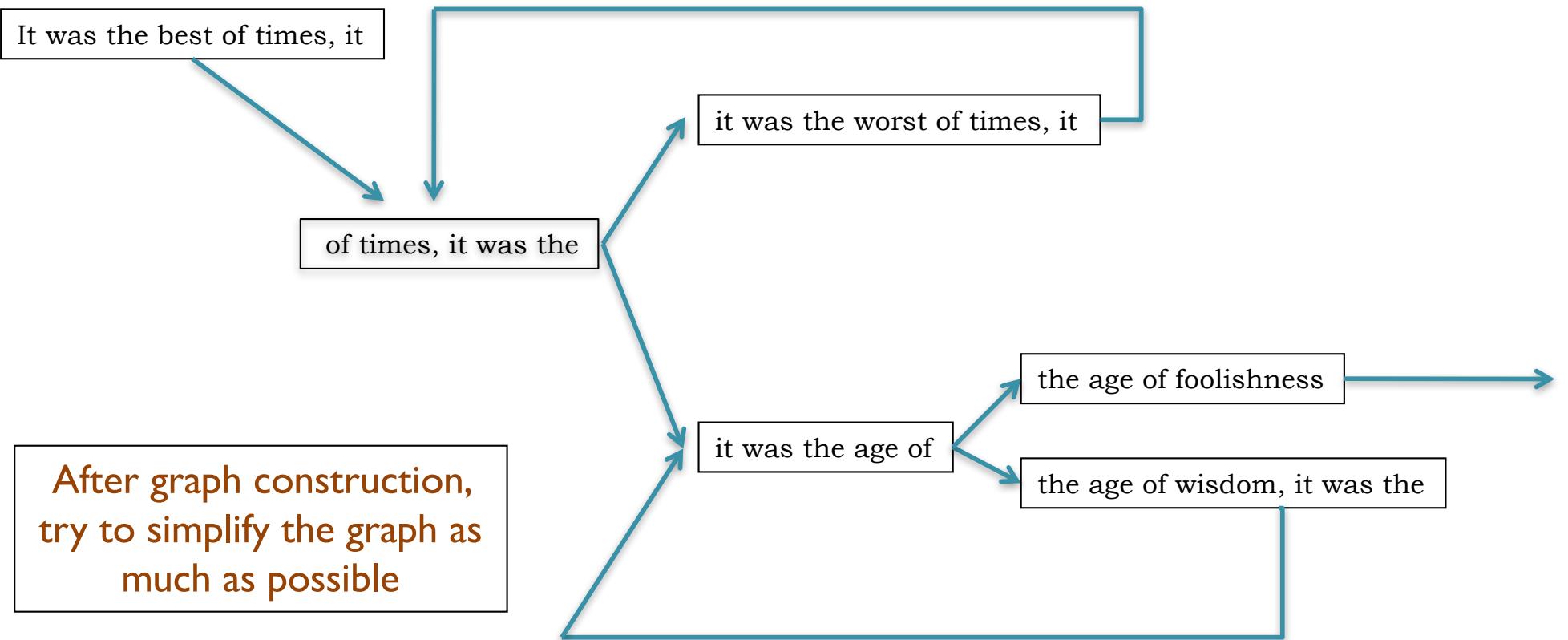
- Overlaps between fragments are implicitly computed

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

de Bruijn Graph Assembly



de Bruijn Graph Assembly



The full tale

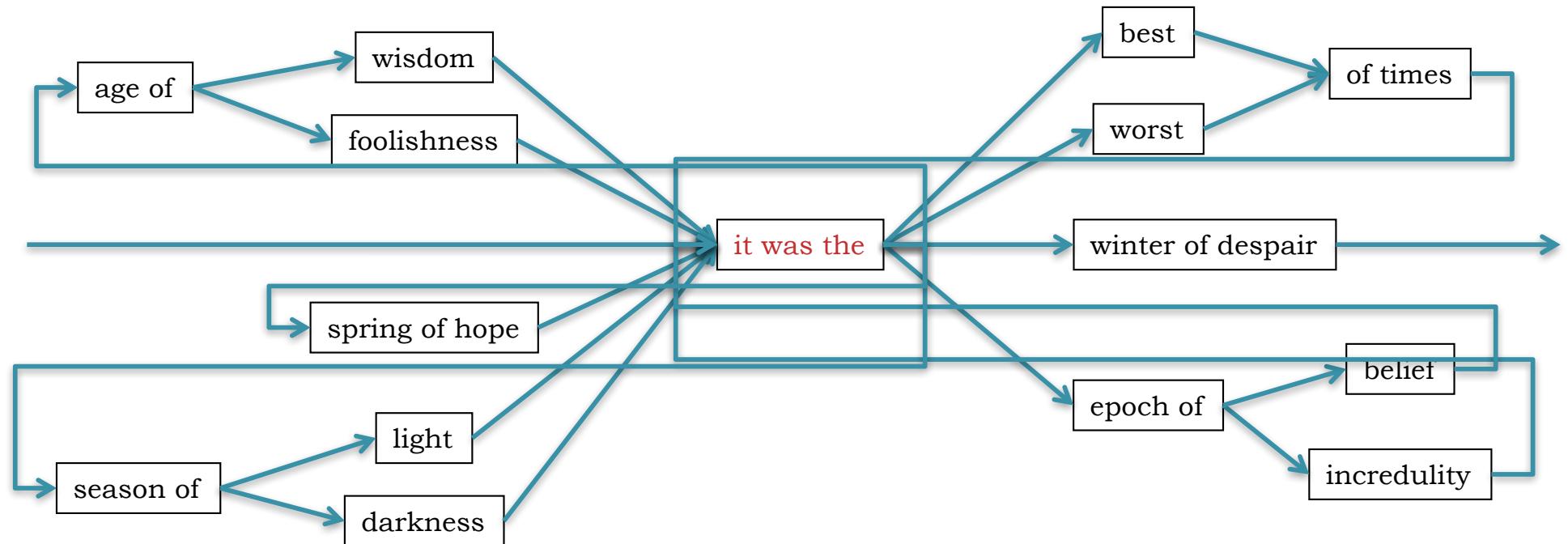
... it was the best of times it was the worst of times ...

... 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 ...

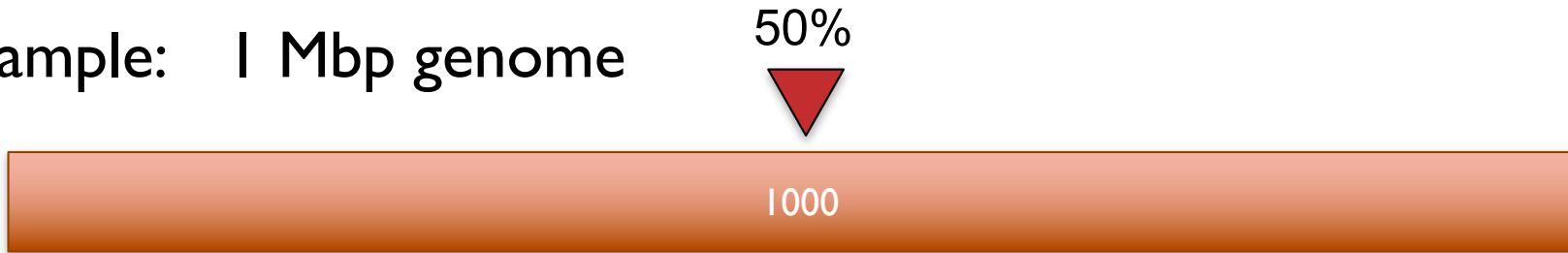
... it was the spring of hope it was the winter of despair ...



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



Outline

1. ***Assembly theory***

- Assembly by analogy

2. ***Practical Issues***

- Coverage, read length, errors, and repeats

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

- PacBio/ONT projects

Assembly Applications

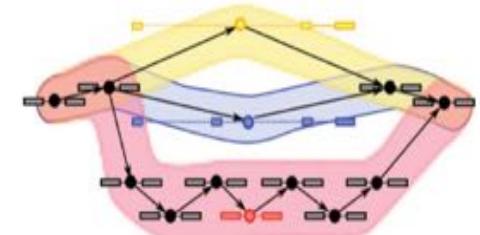
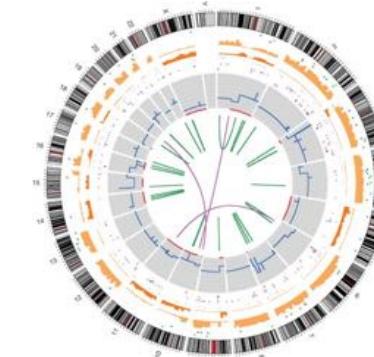
- Novel genomes



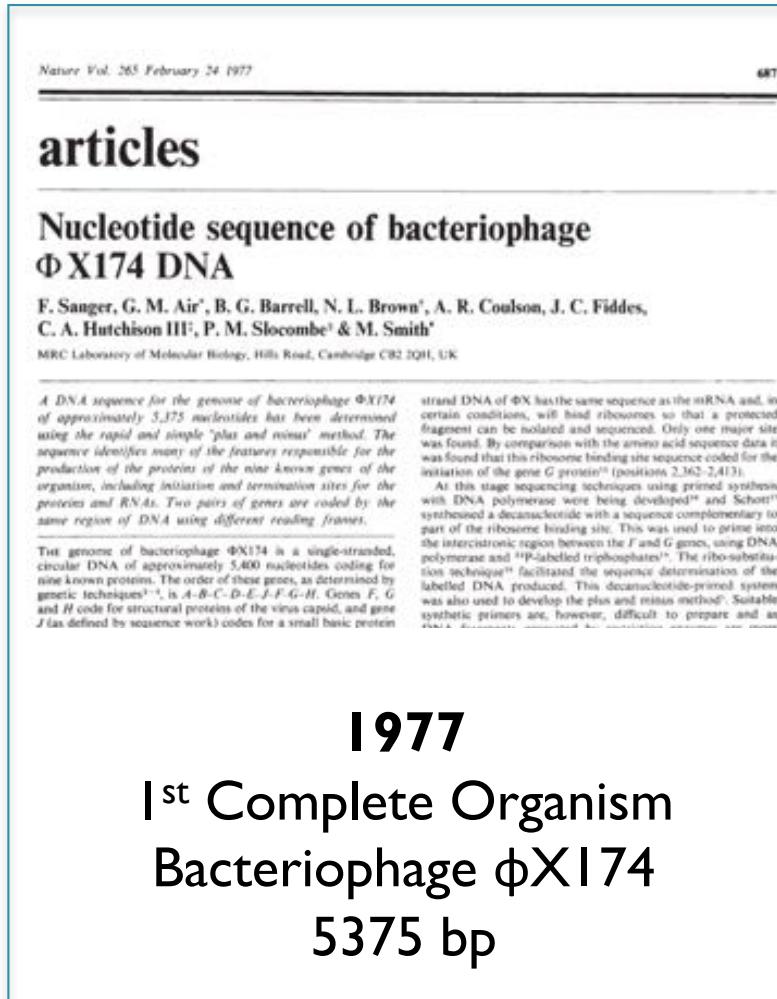
- Metagenomes



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



Milestones in Molecular Biology

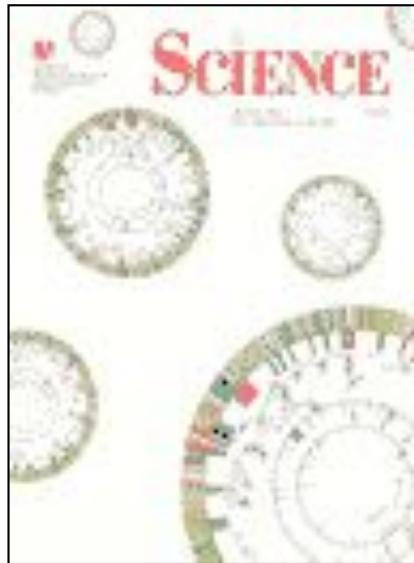


Radioactive Chain Termination
5000bp / week / person

<http://en.wikipedia.org/wiki/File:Sequencing.jpg>
<http://www.answers.com/topic/automated-sequencer>

Nucleotide sequence of bacteriophage φ X174 DNA
Sanger, F. et al. (1977) Nature. 265: 687 - 695

Milestones in Molecular Biology



1995

Fleischmann et al.
1st Free Living Organism
TIGR Assembler. 1.8Mbp



2000

Myers et al.
1st Large WGS Assembly.
Celera Assembler. 116 Mbp



2001

Venter et al. / IHGSC
Human Genome
Celera Assembler. 2.9 Gbp

ABI 3700: 500 bp reads x 768 samples / day = 384,000 bp / day.

"The machine was so revolutionary that it could decode in a single day the same amount of genetic material that most DNA labs could produce in a year." J. Craig Venter

Milestones in Molecular Biology



2004
454/Roche Titanium
Pyrosequencing
1M 400bp reads / run =
1Gbp / day

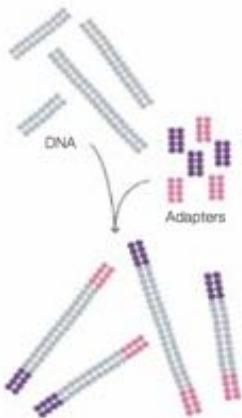


2007
Illumina HiSeq 2000
Sequencing by Synthesis
2.5B 2x100bp reads / run =
60Gbp / day

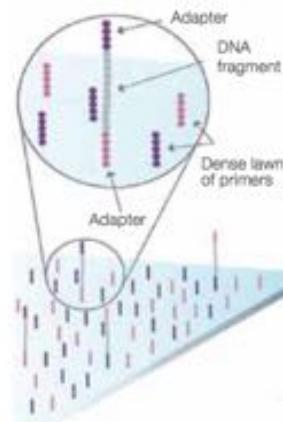


2017
Illumina NovaSeq
Sequencing by Synthesis
12B 2x150bp reads / run =
>1TB / day

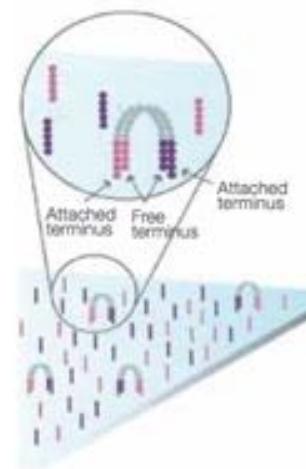
Illumina Sequencing by Synthesis



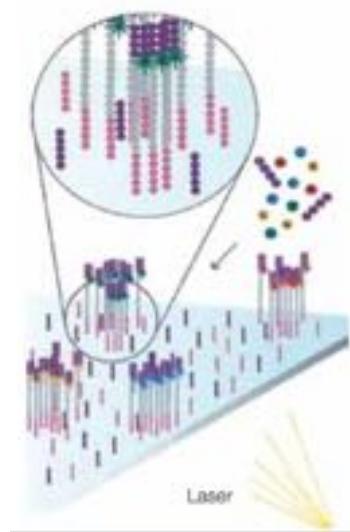
1. Prepare



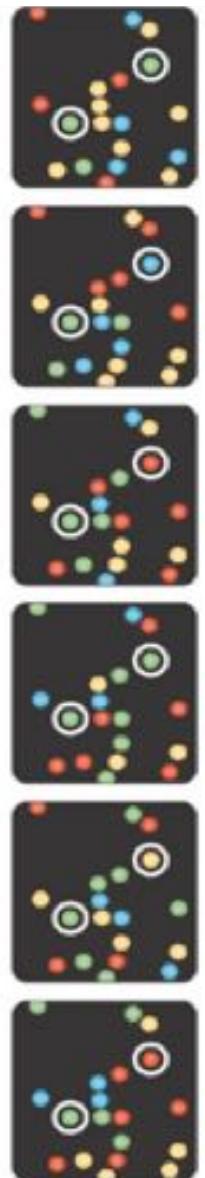
2. Attach



3. Amplify



4. Image

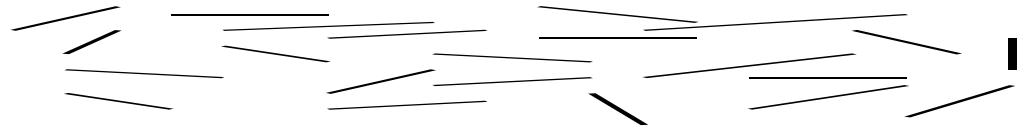


5. Basecall

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

Assembling a Genome

1. Shear & Sequence DNA

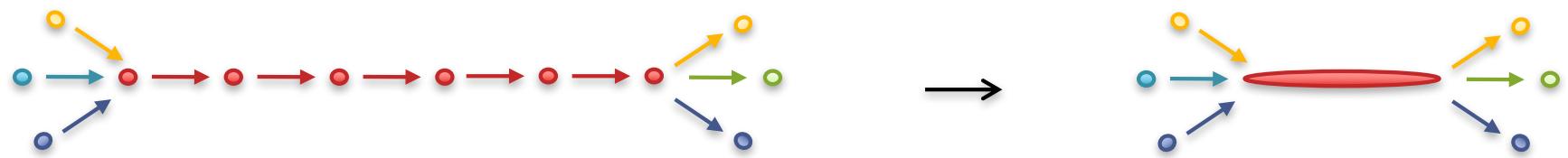


2. Construct assembly graph from reads (de Bruijn / overlap graph)

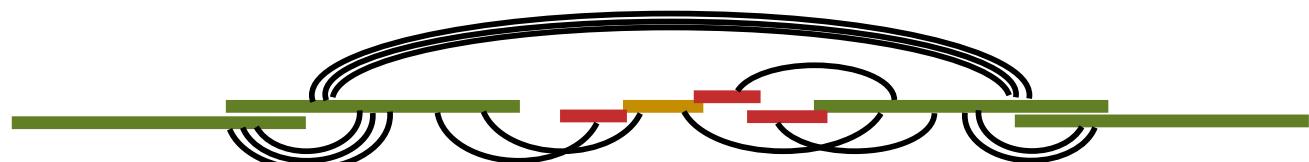
...AGCCTAG**GGATGCGCGACACGT**

GGATGCGCGACACGTCGCATATCCGGTTTGGT**CAACCTCGGACGGAC**
CAACCTCGGACGGACCTCAGCGAA...

3. Simplify assembly graph



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



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

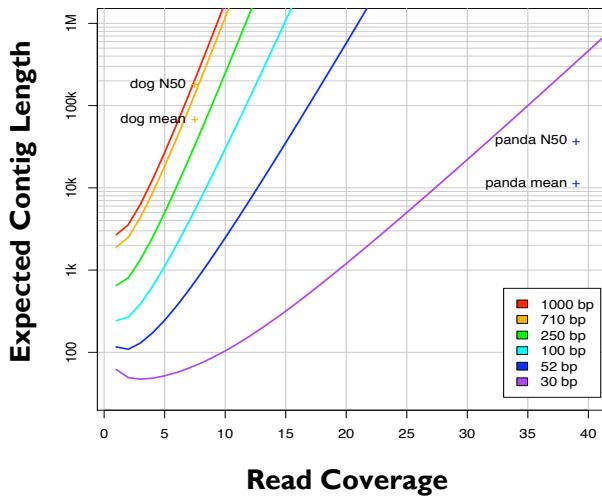
4. *Accuracy:*

- (Very) Hard to assess correctness



Ingredients for a good assembly

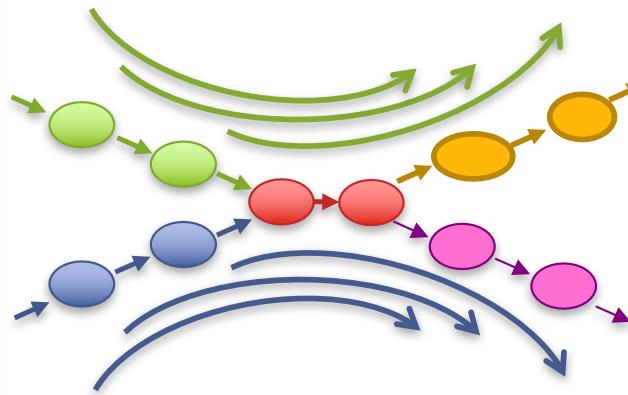
Coverage



High coverage is required

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

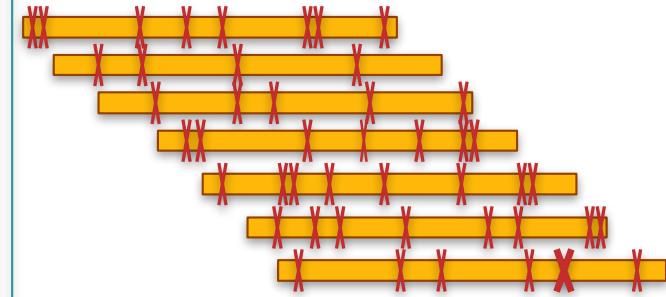
Read Length



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



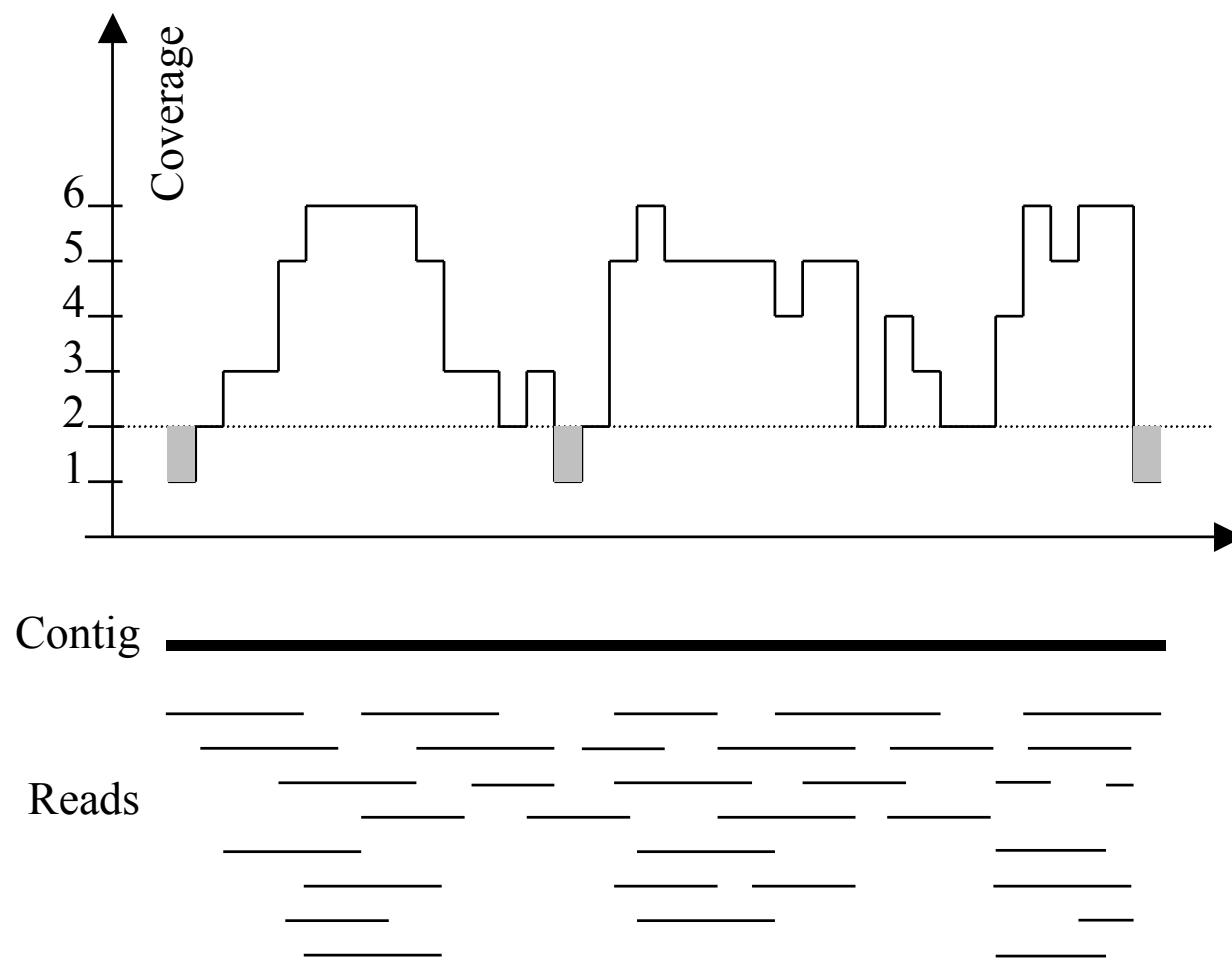
Errors obscure overlaps

- 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

Typical sequencing coverage

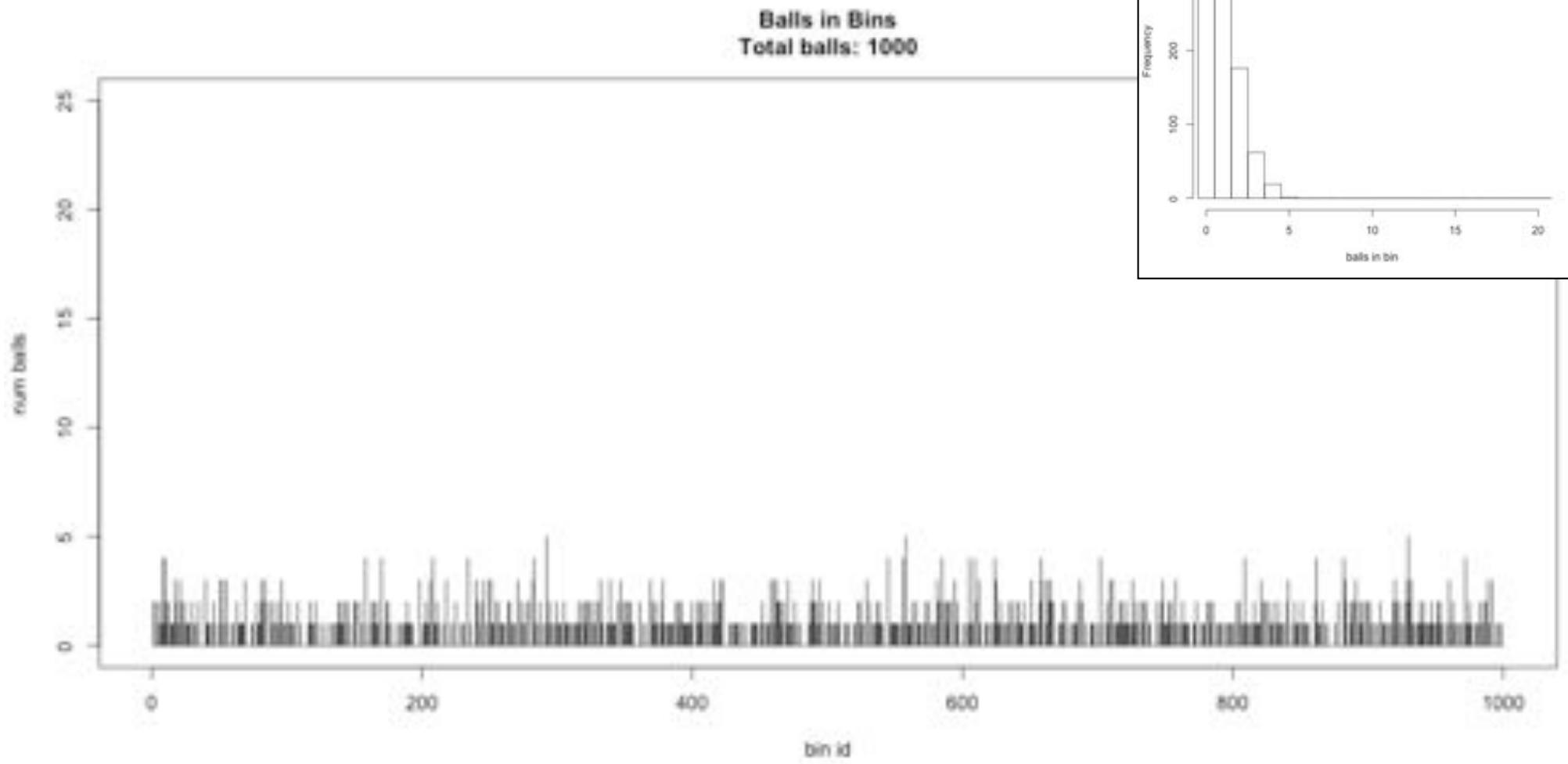


Imagine raindrops on a sidewalk

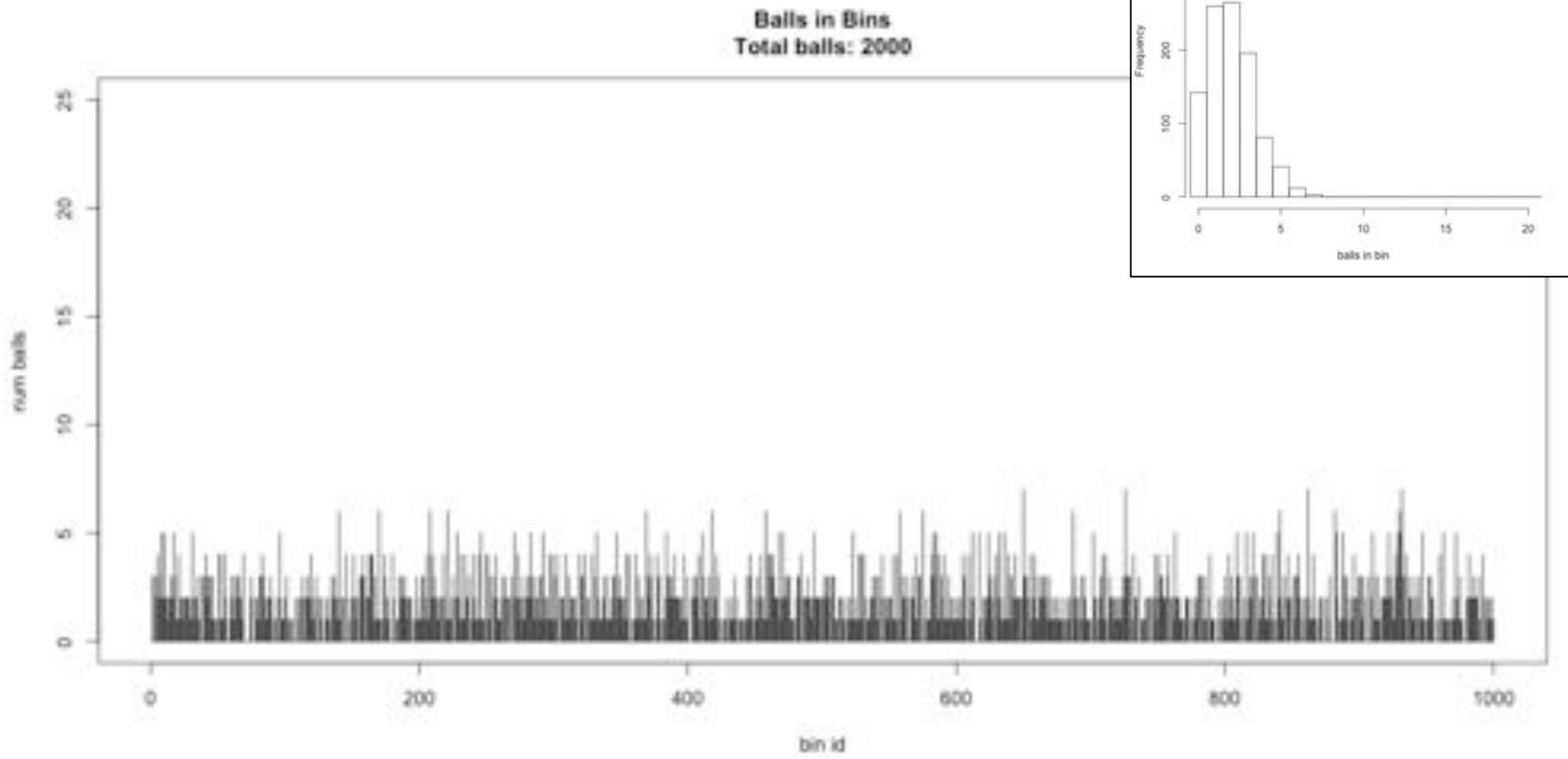
We want to cover the entire sidewalk but each drop costs \$1

If the genome is 100 Mbp, should we sequence 1M 100bp reads?

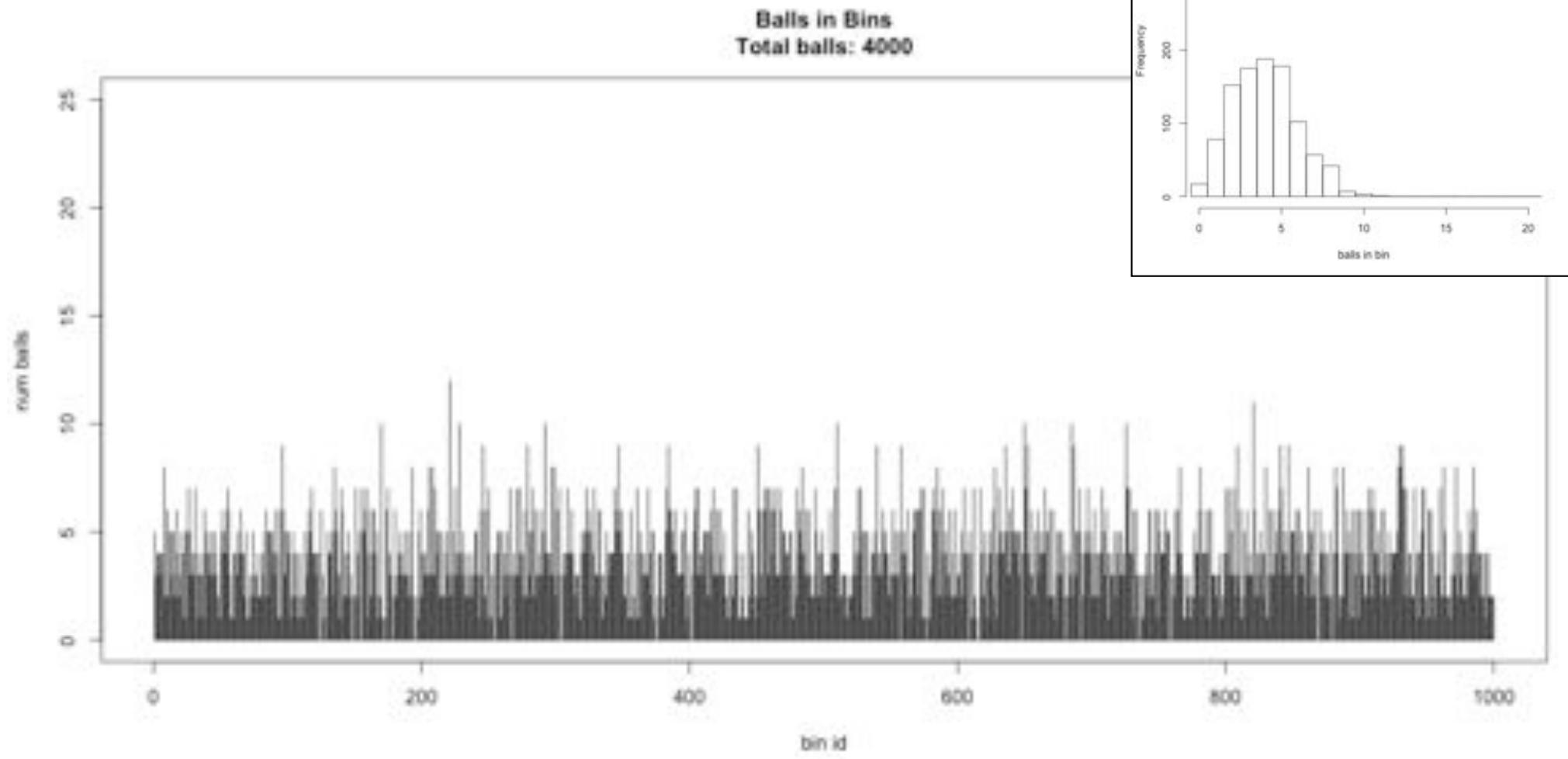
Ix sequencing



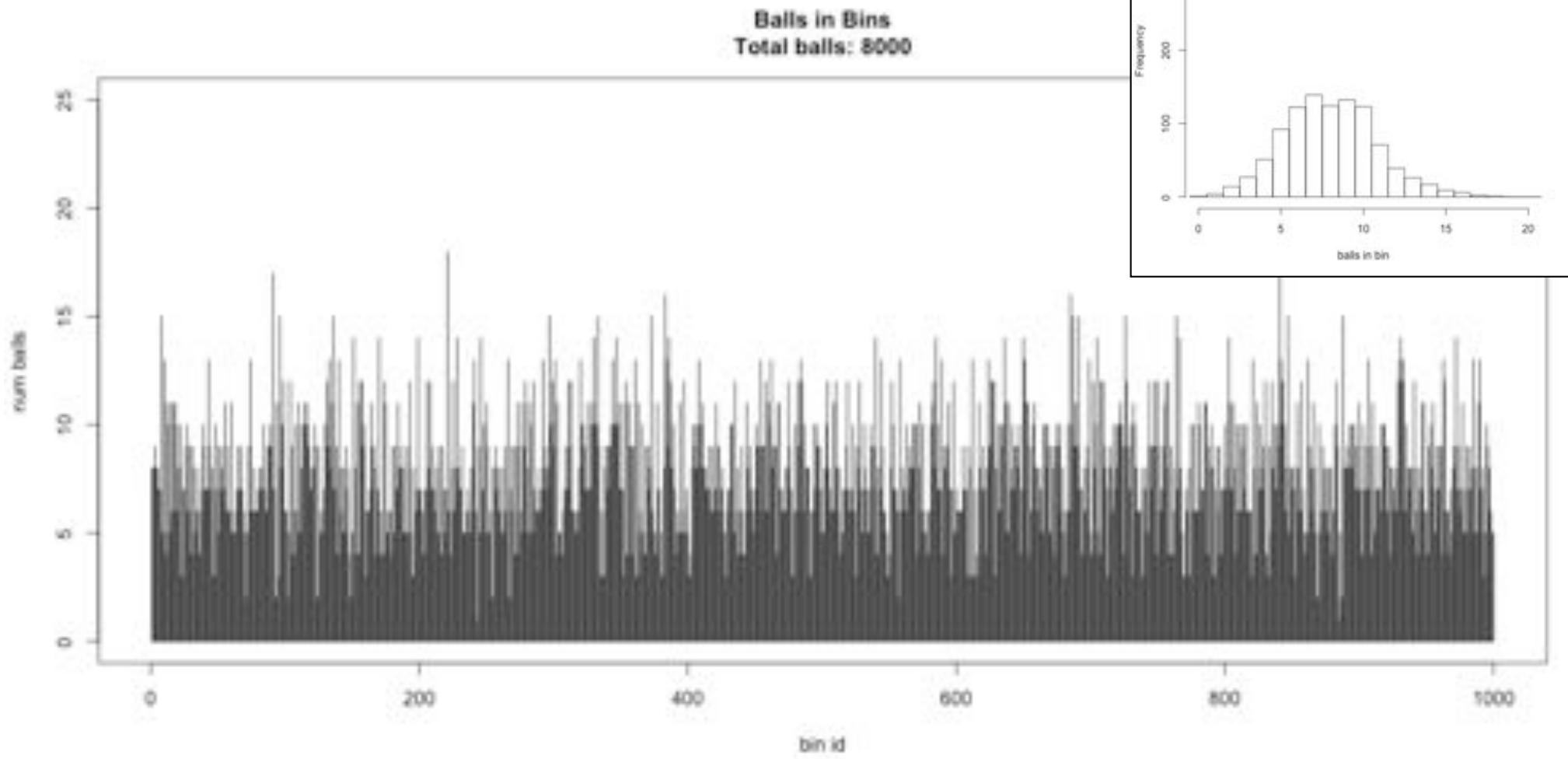
2x sequencing



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.

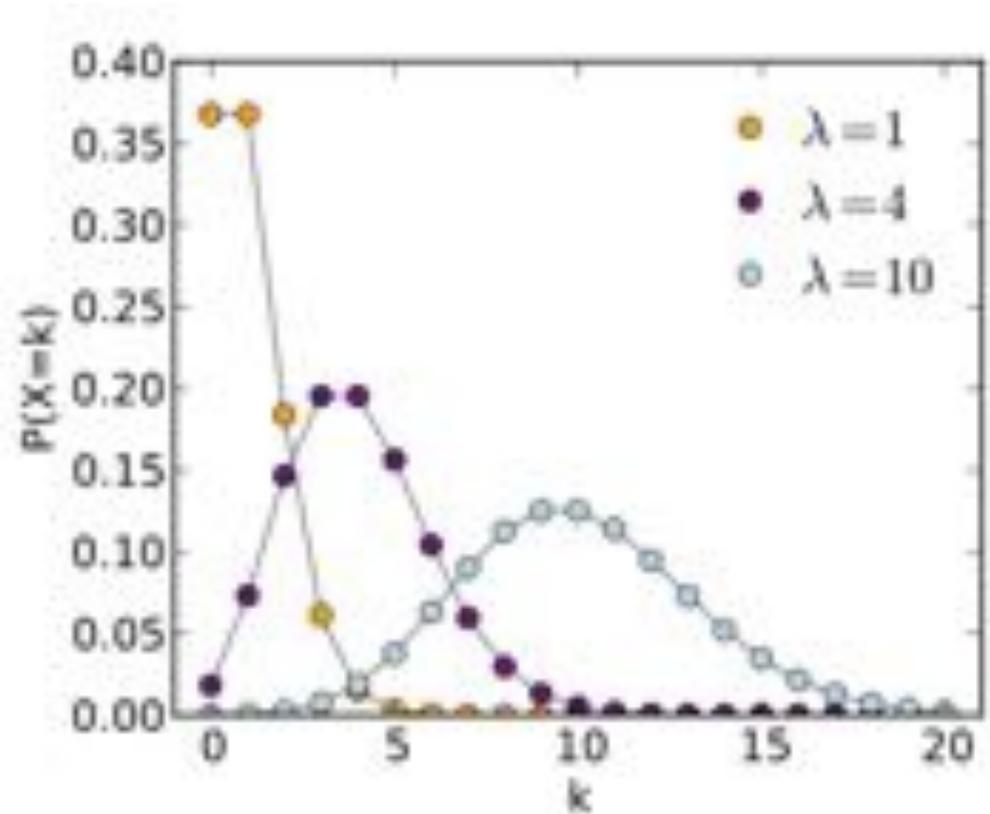
Formulation comes from the limit of the binomial equation

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

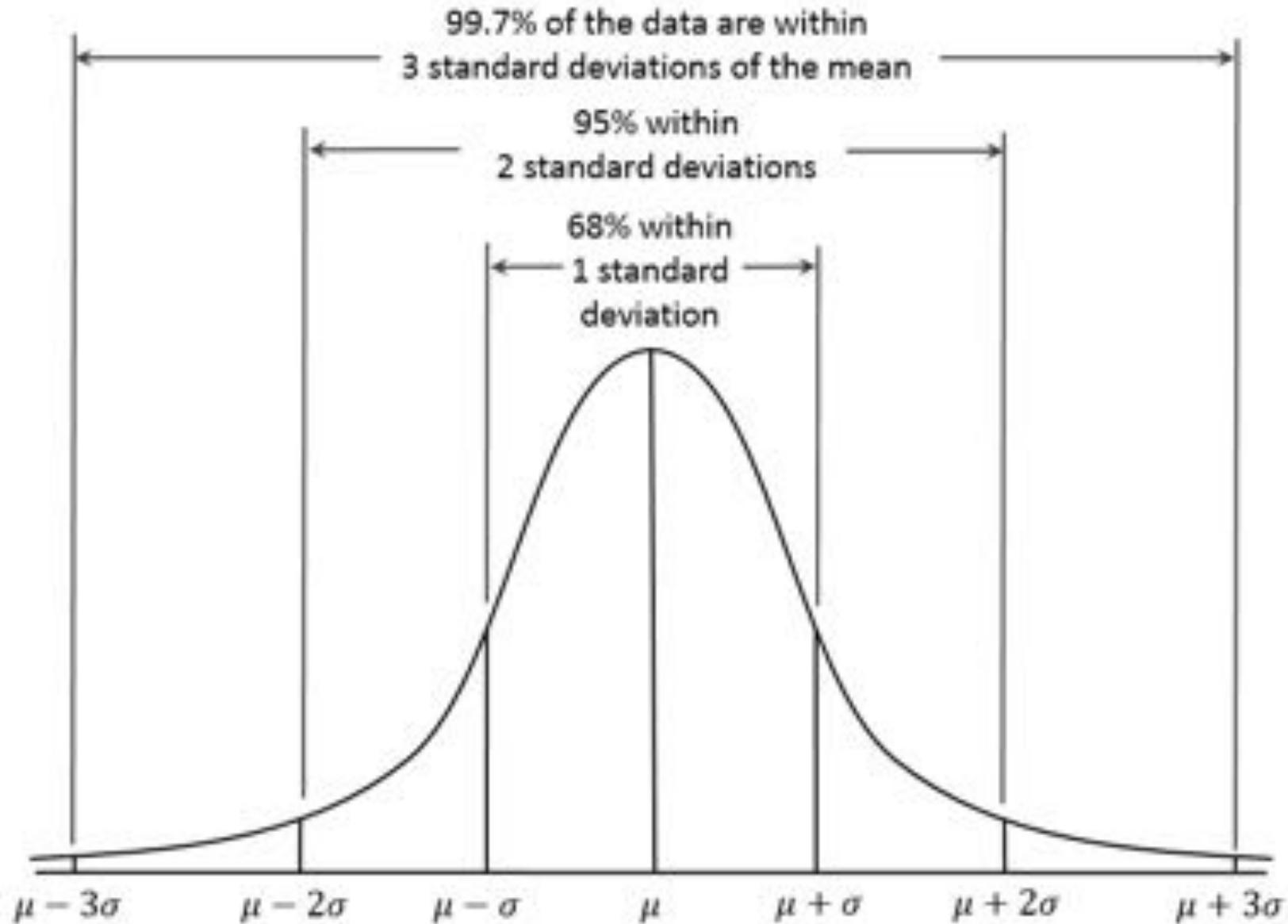
Key properties:

- ***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 150bp reads do I need?

I need $10\text{Mbp} \times 24\text{x} = 240\text{Mbp}$ of data
 $240\text{Mbp} / 150\text{bp} / \text{read} = 1.6\text{M reads}$

I want to sequence a 10Mbp genome so that
>97.5% of the genome has at least 24x coverage.
How many 150bp reads do I need?

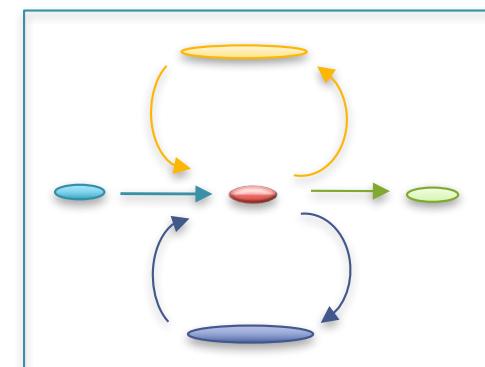
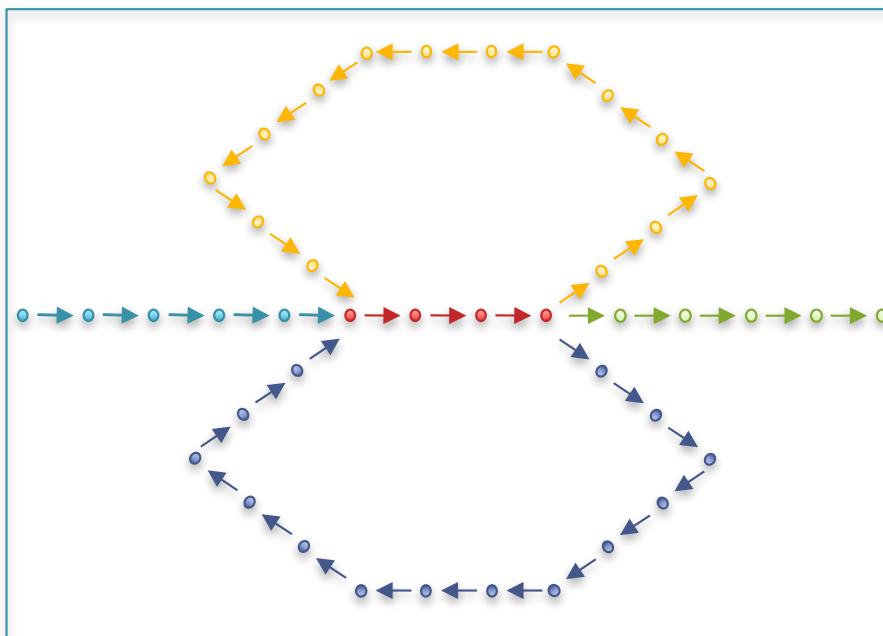
Find X such that $X - 2\sqrt{X} = 24$

$$36 - 2\sqrt{36} = 24$$

I need $10\text{Mbp} \times 36\text{x} = 360\text{Mbp}$ of data
 $360\text{Mbp} / 150\text{bp} / \text{read} = 2.4\text{M reads}$

Unitigging / Unipathing

- After simplification and correction, compress graph down to its non-branching initial contigs
 - Aka “unitigs”, “unipaths”
 - Unitigs end because of (1) lack of coverage, (2) errors, (3) heterozygosity and (4) repeats

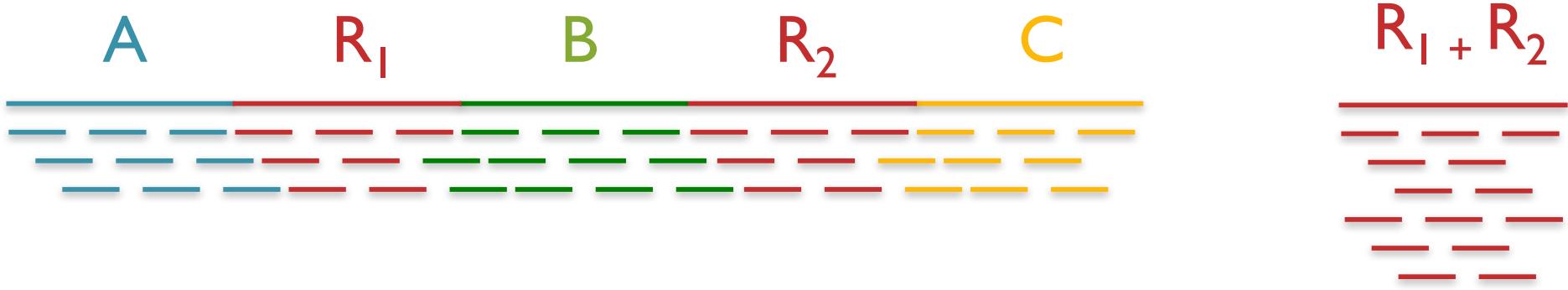


Repetitive regions

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

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

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

The fragment assembly string graph

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

Paired-end and Mate-pairs

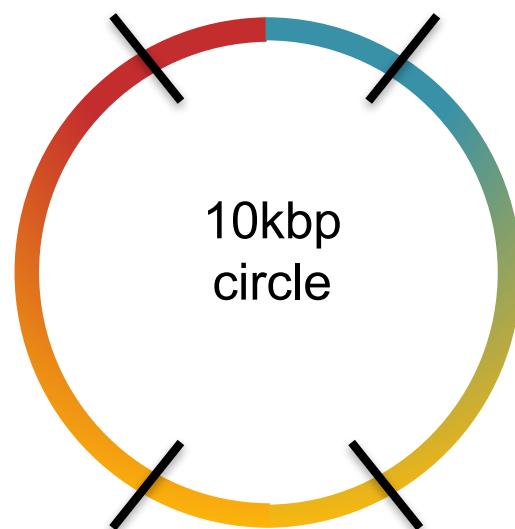
Paired-end sequencing

- 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)

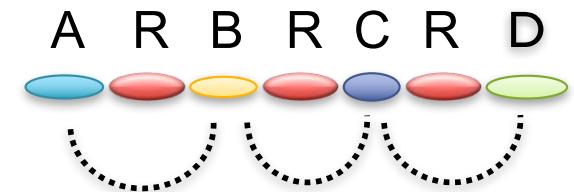
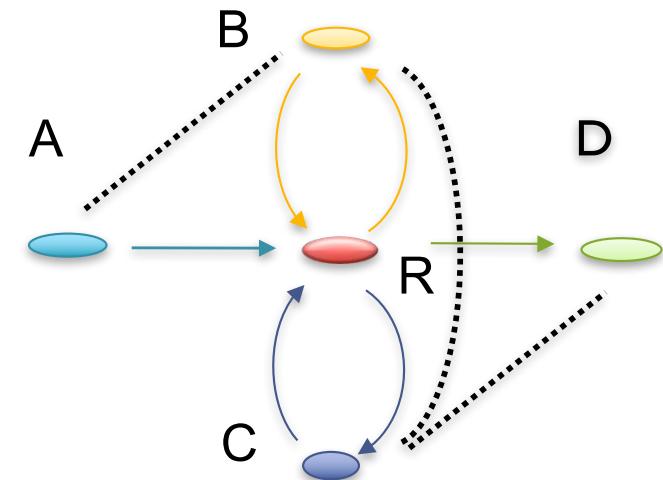


2x100 @ 300bp (innies)

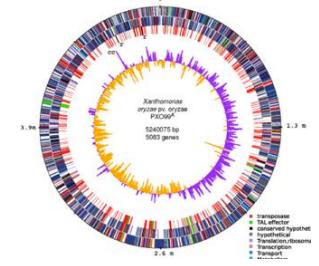


Scaffolding

- 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



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



Outline

1. ***Assembly theory***

- Assembly by analogy

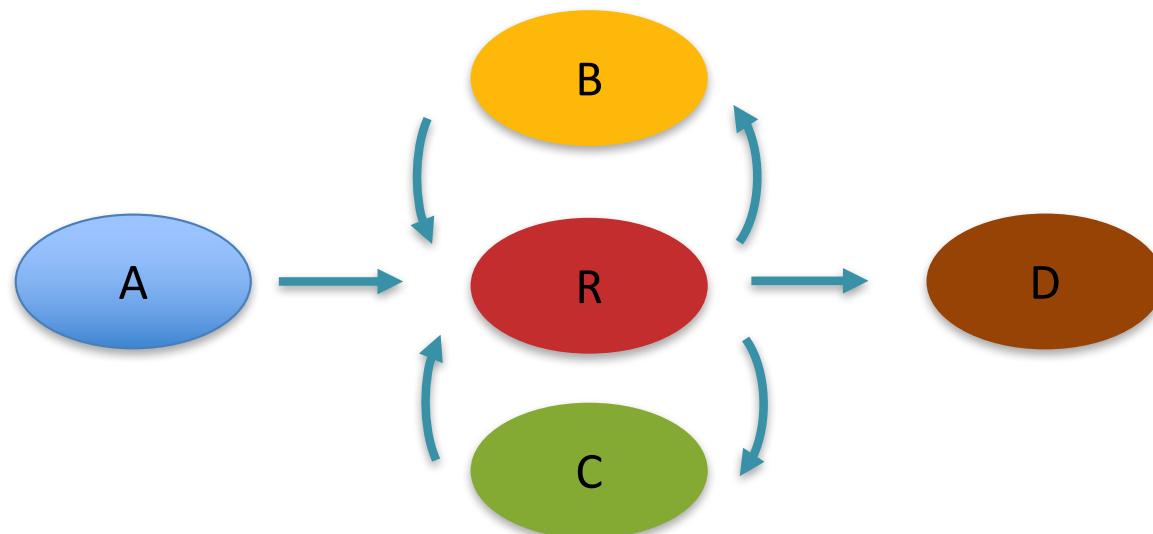
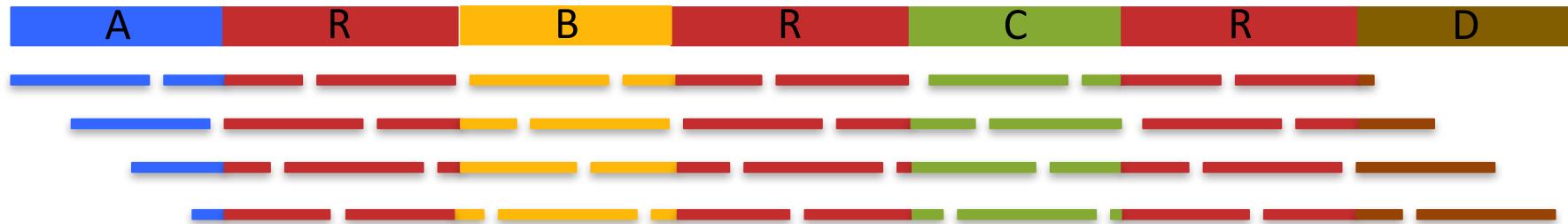
2. ***Practical Issues***

- Coverage, read length, errors, and repeats

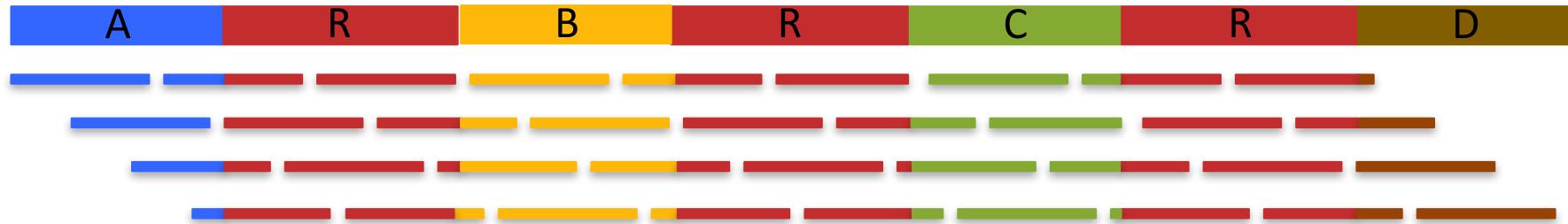
3. ***Next-next-gen Assembly***

- PacBio/ONT projects

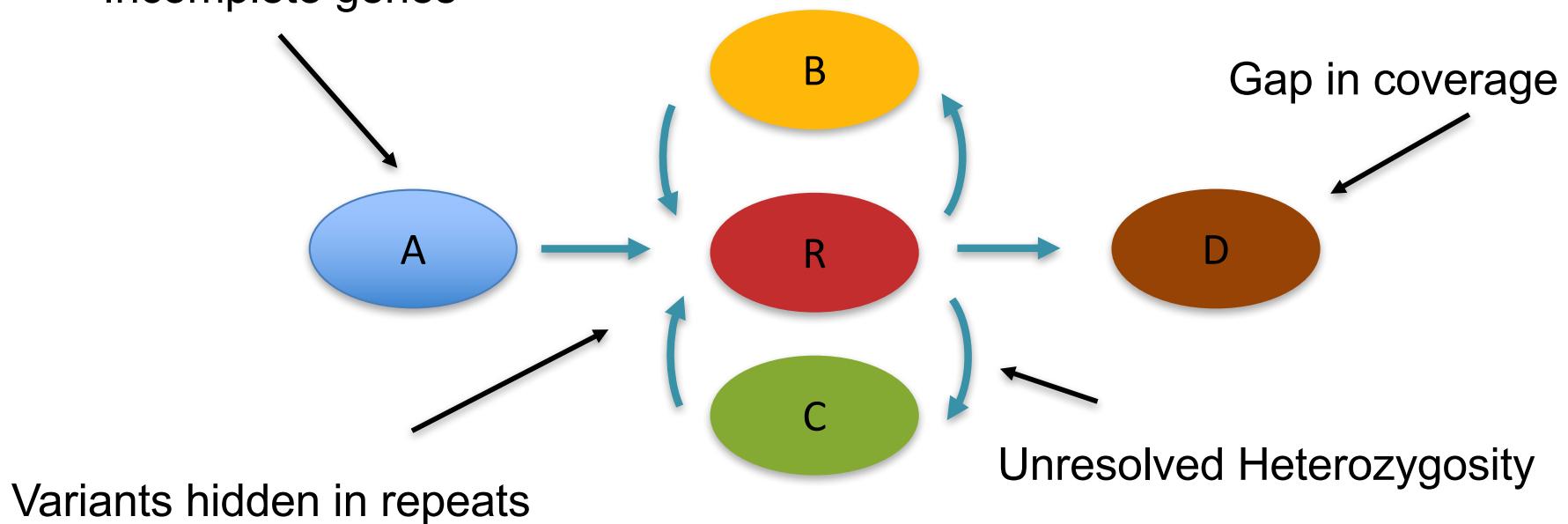
Assembly Complexity



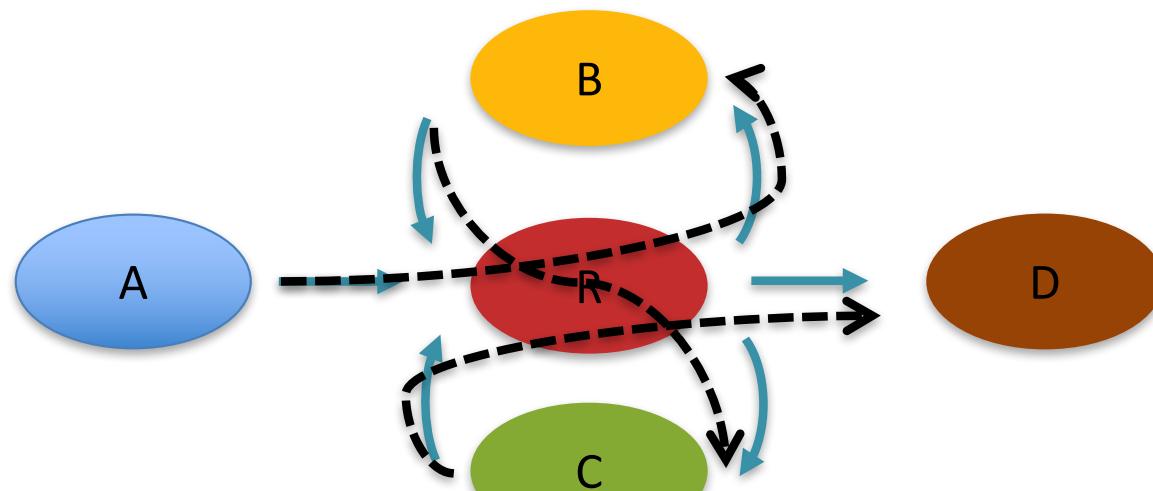
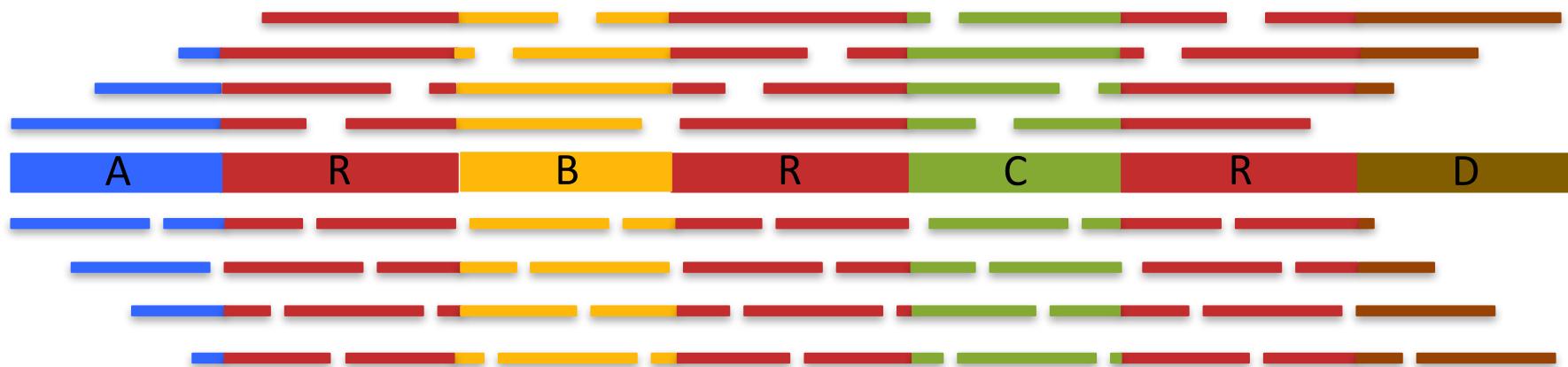
Assembly Complexity



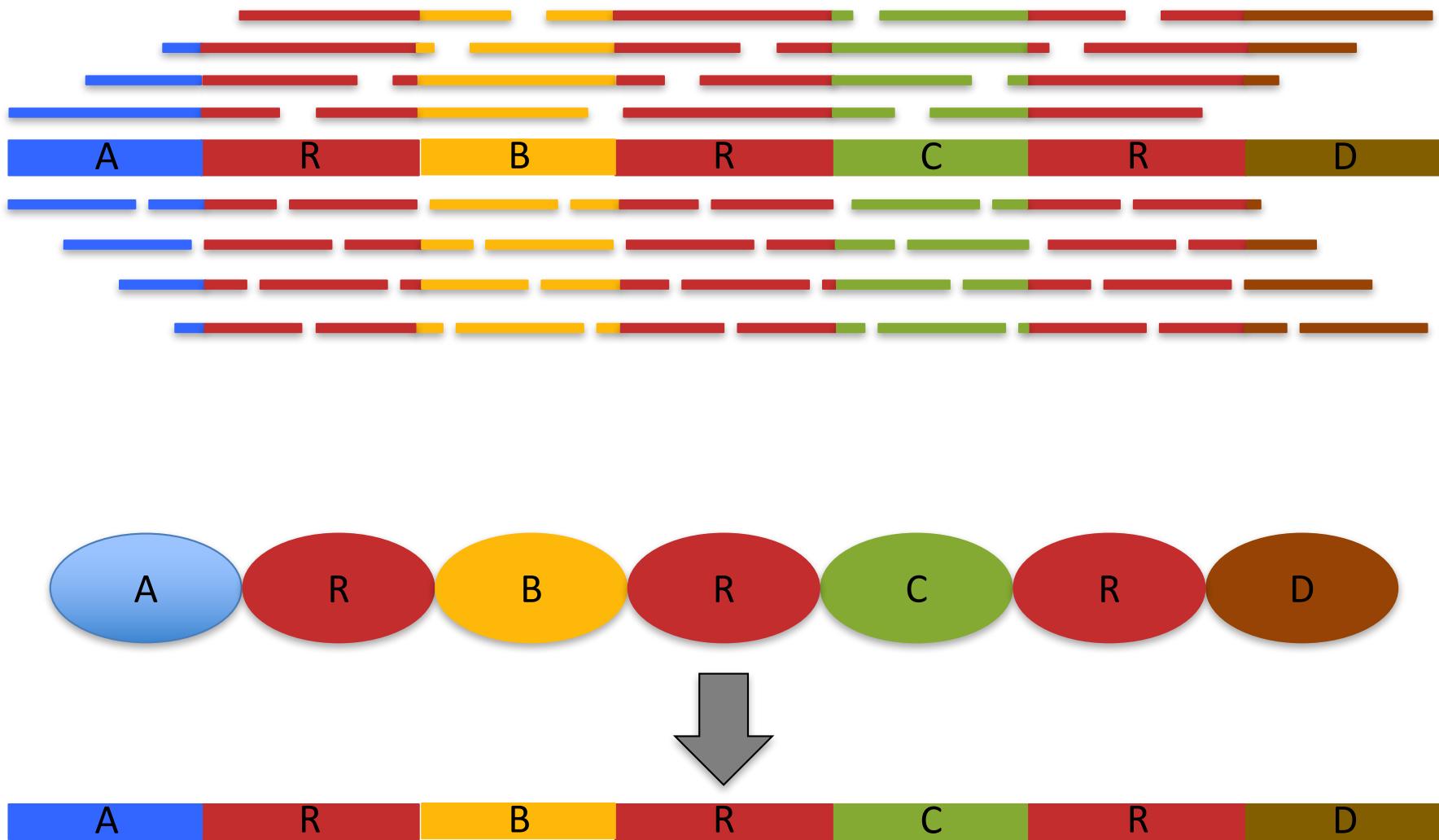
Short contigs &
Incomplete genes



Assembly Complexity



Assembly Complexity



The advantages of SMRT sequencing

Roberts, RJ, Carneiro, MO, Schatz, MC (2013) *Genome Biology*. 14:405

Genomics Arsenal in the year 2017

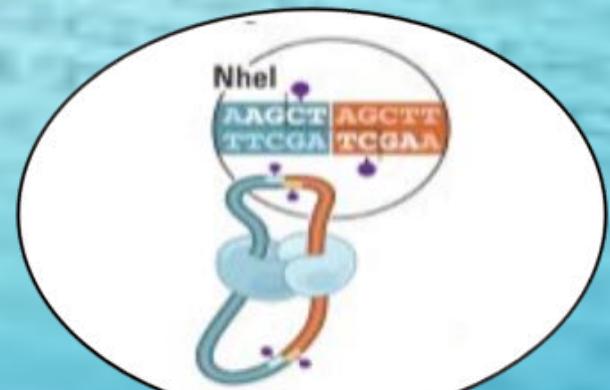
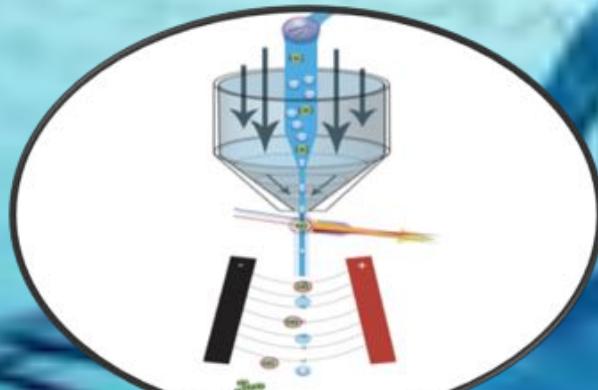
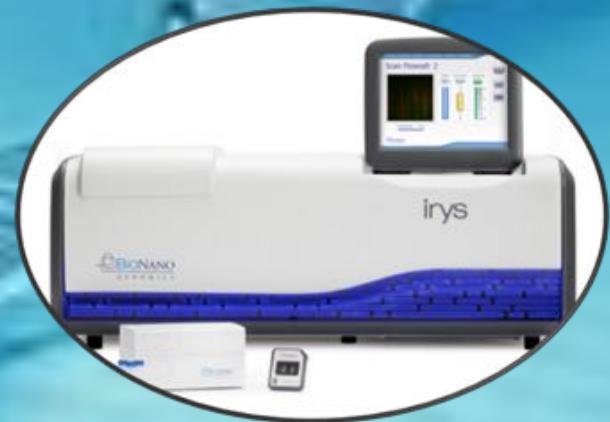
Sample Preparation



Sequencing

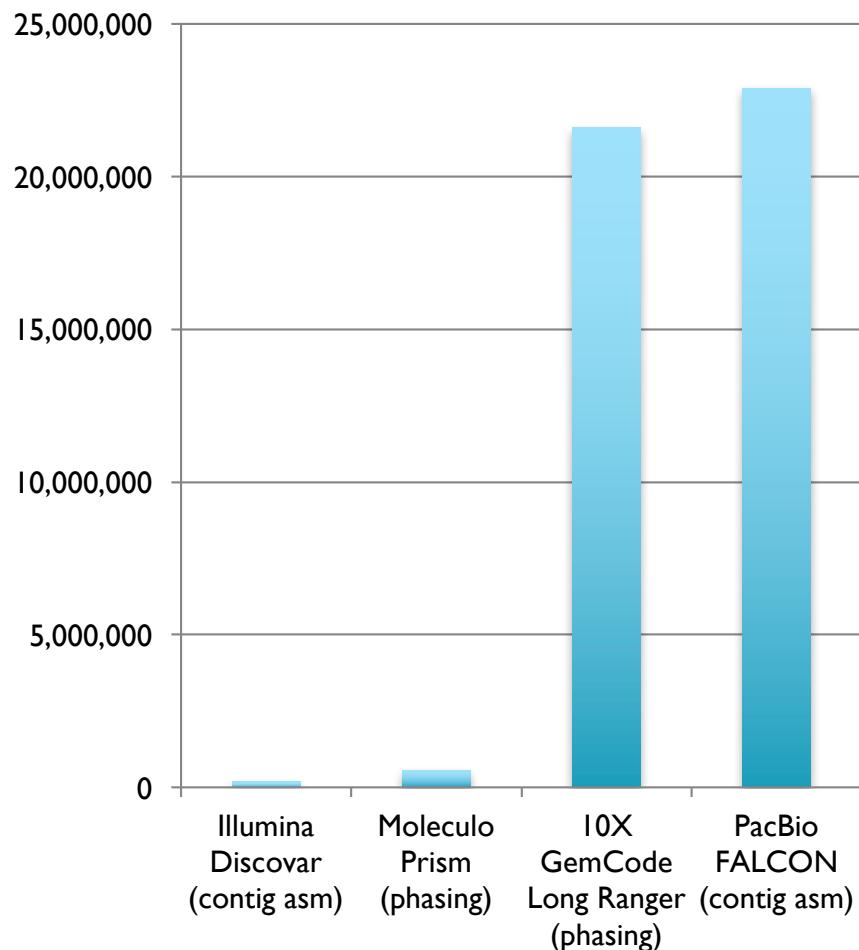


Chromosome Mapping



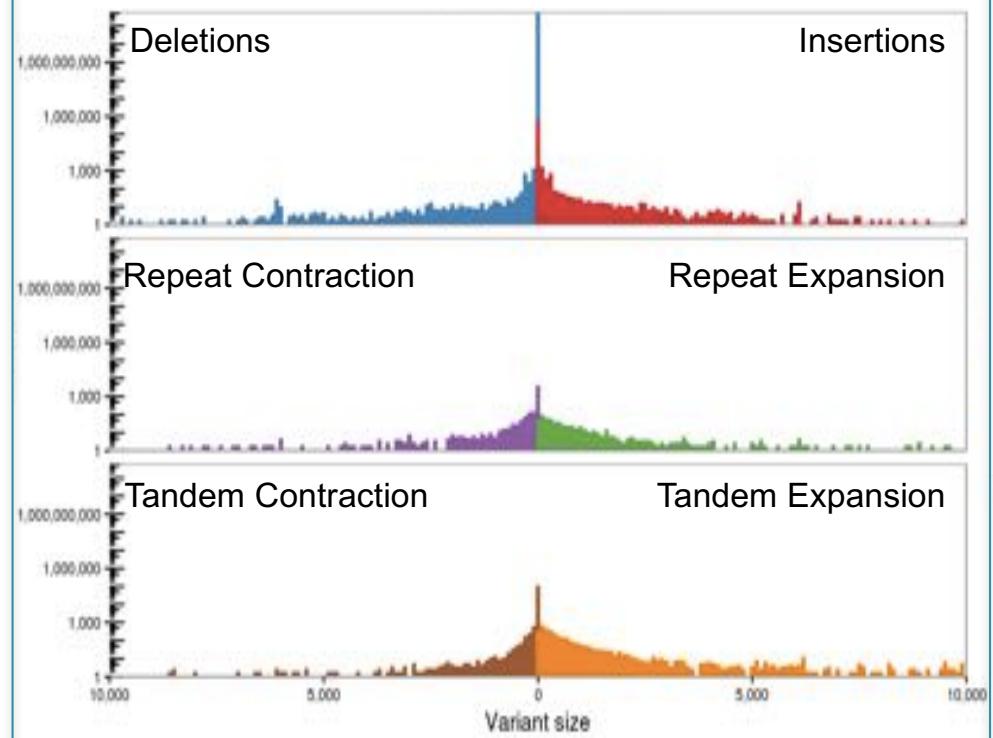
Recent Long Read Assemblies

Human Analysis N50 Sizes



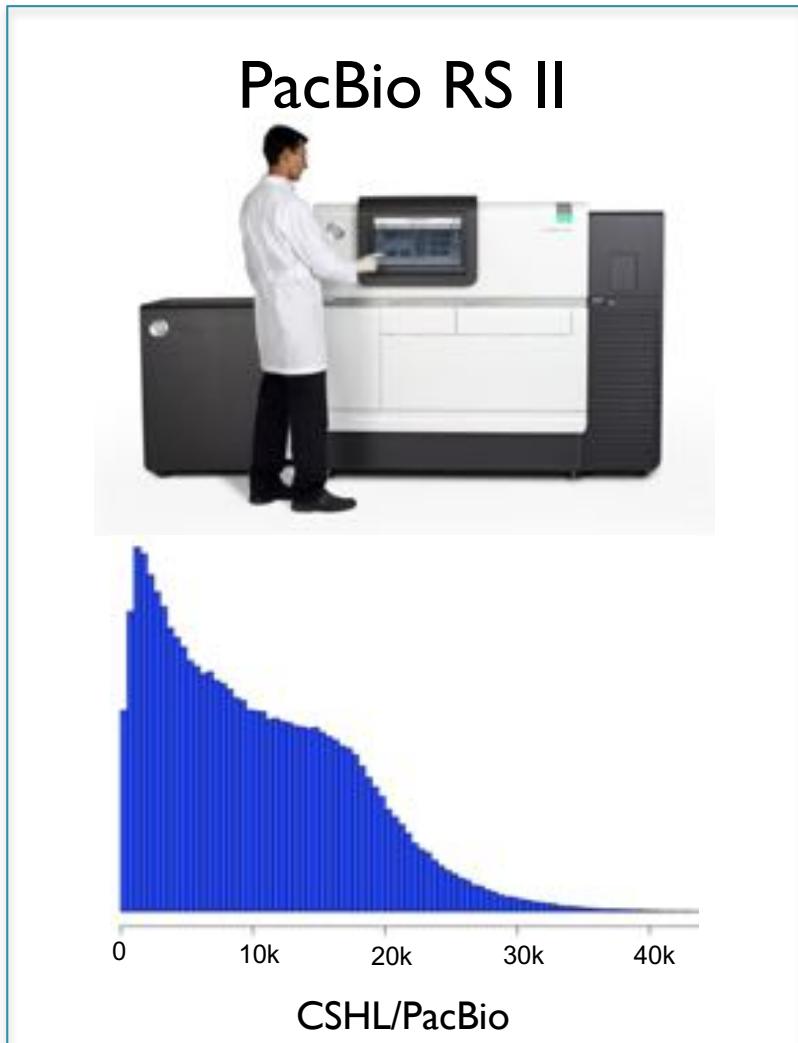
Third-generation sequencing and the future of genomics
Lee et al (2016) *bioRxiv*
doi: <http://dx.doi.org/10.1101/048603>

Structural Variants in CHM1



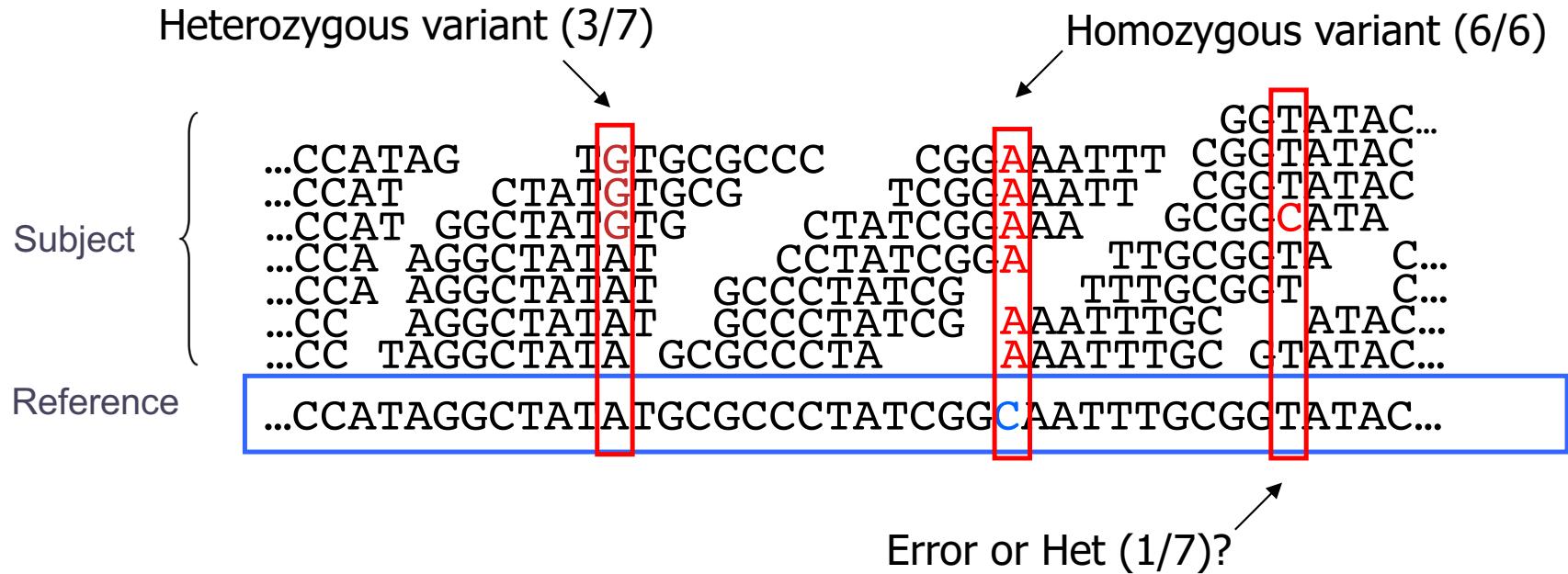
Assemblytics: a web analytics tool for the detection of variants from an assembly
Nattestad & Schatz (2016) *Bioinformatics*.
doi: [10.1093/bioinformatics/btw369](https://doi.org/10.1093/bioinformatics/btw369)

Single Molecule Sequencing



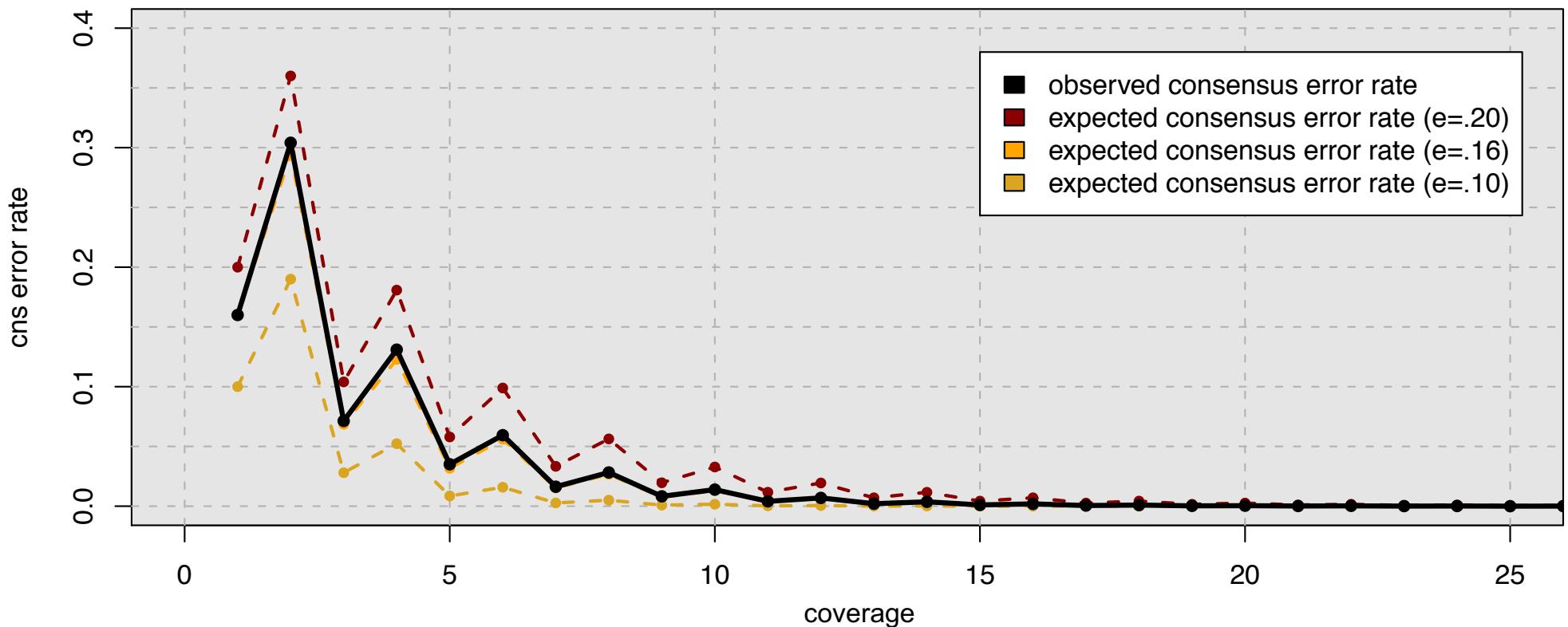
Sample of 100k reads aligned with BLASR requiring >100bp alignment
Average overall accuracy 83.7%: 11.5% insertions, 3.4% deletions, 1.4% mismatch

Genotyping Theory



- If there were no sequencing errors, identifying SNPs would be trivial:
 - Any time a read disagrees with the reference, it must be a variant!
- A single read of many differing from the reference is probably just an error, but it becomes more likely to be real as we see it multiple times
 - Use binomial test to evaluate prob. of heterozygosity vs. prob of error
 - Coverage (oversampling) is our main tool to improve accuracy

Consensus Accuracy and Coverage

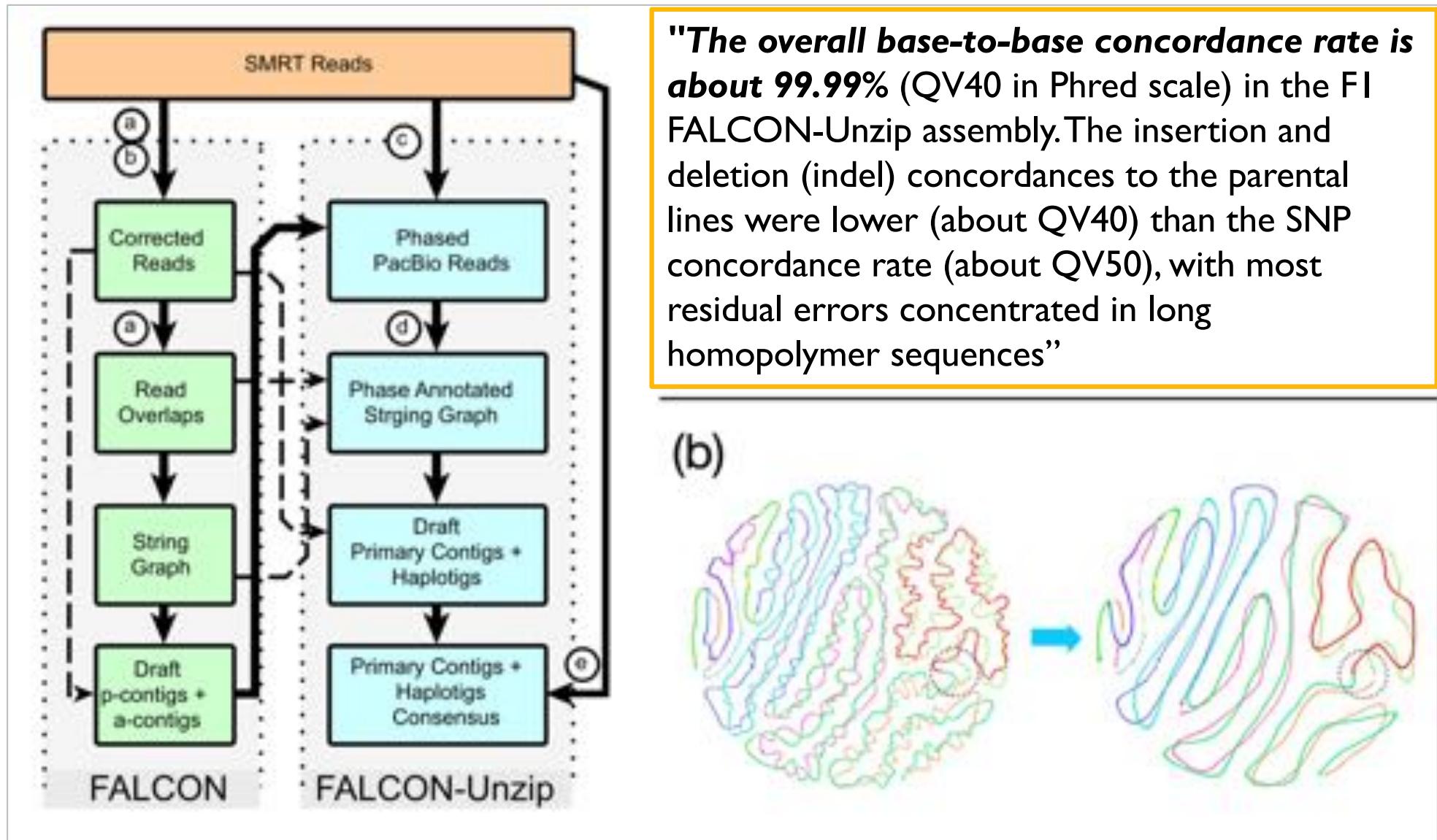


Coverage can overcome random errors

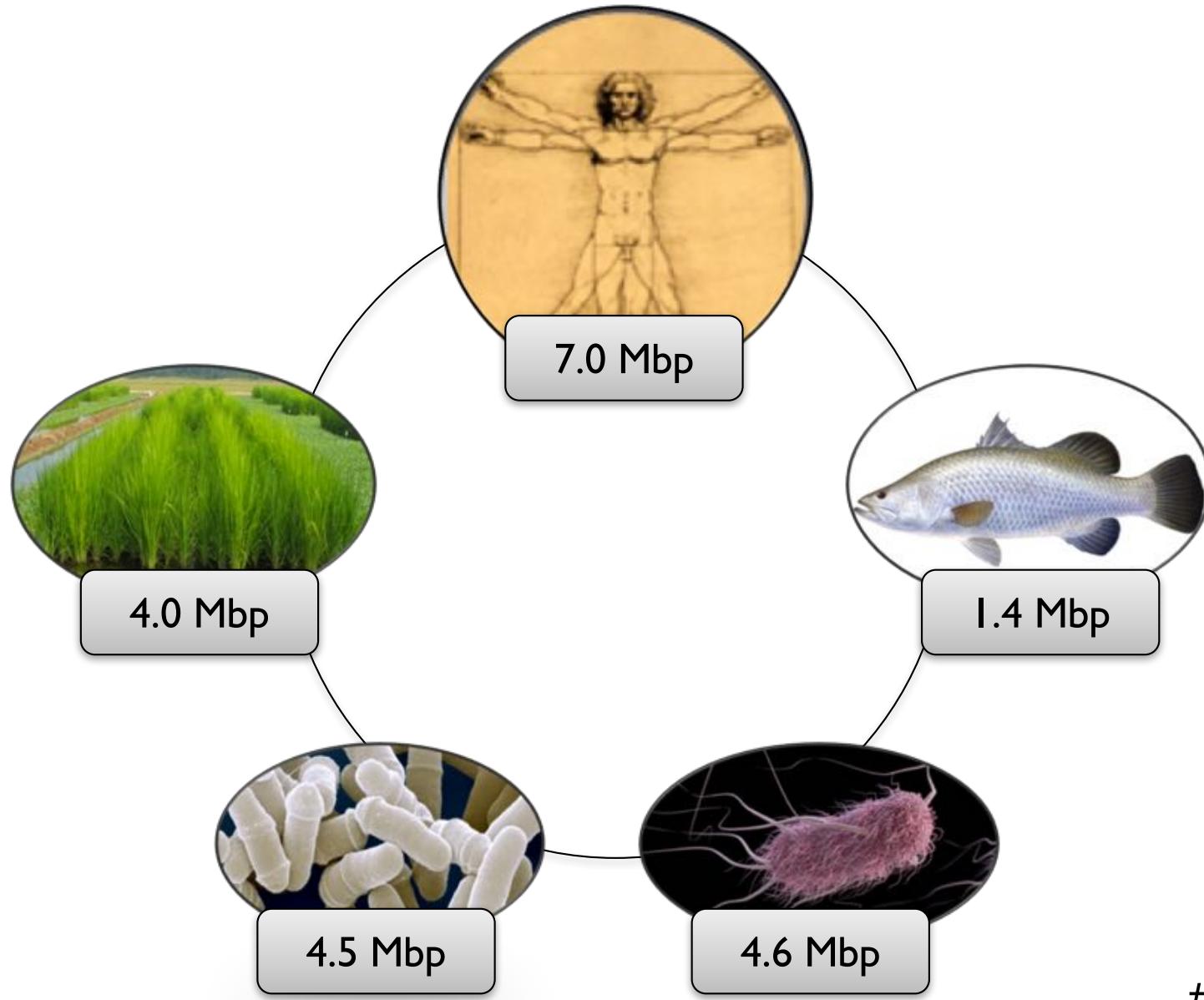
- Dashed: error model from binomial sampling
- Solid: observed accuracy

$$CNS\ Error = \sum_{i=\lceil c/2 \rceil}^c \binom{c}{i} (e)^i (1-e)^{n-i}$$

FALCON Accuracy

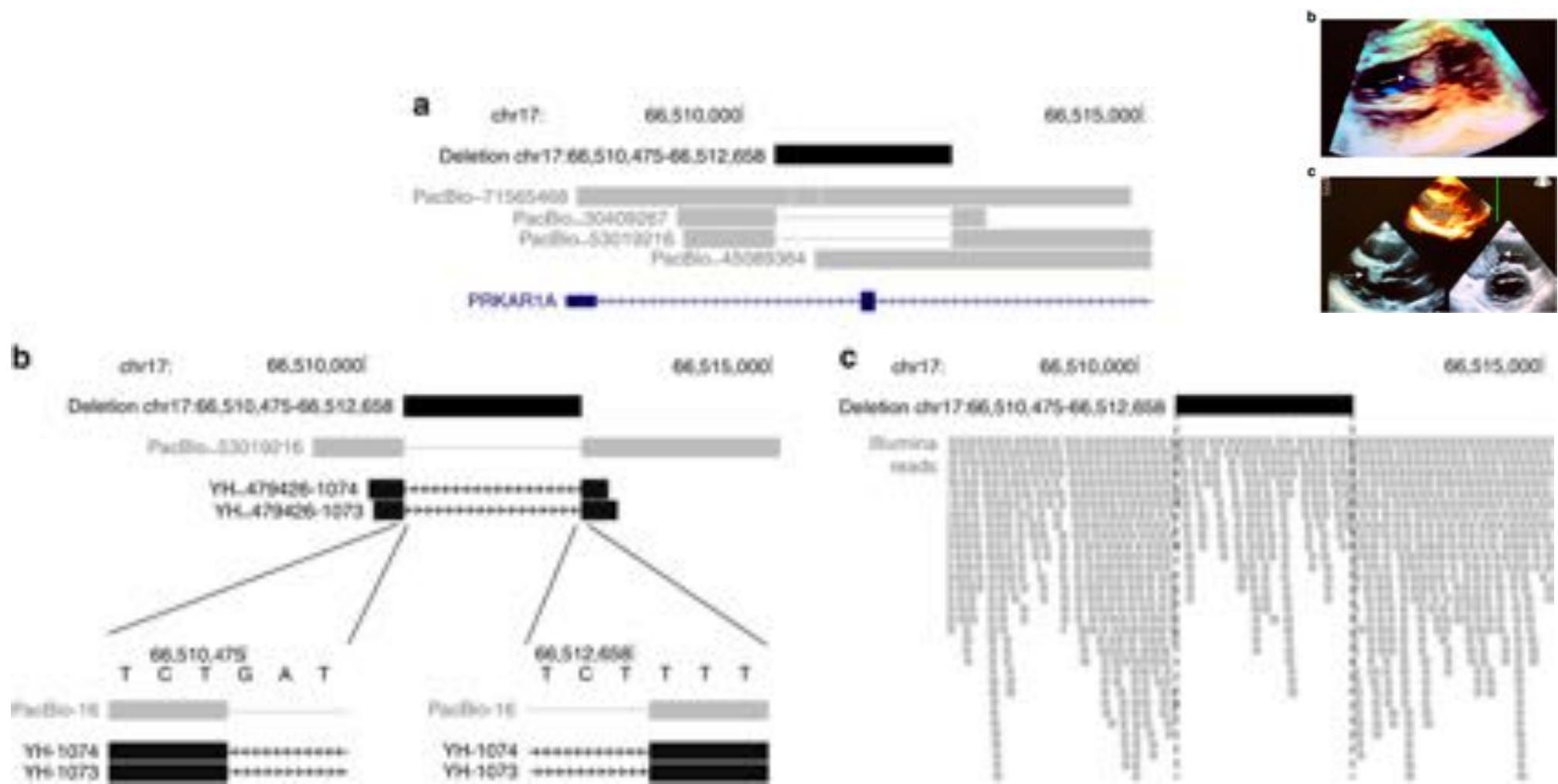


(A few) Recent PacBio Assemblies



#1mbctgclub

Structural Variations in Human Disease



Long-read genome sequencing identifies causal structural variation in a Mendelian disease
Merker et al (2017) *Genetics in Medicine*. doi:10.1038/gim.2017.86

NGMLR + Sniffles

BWA-MEM:



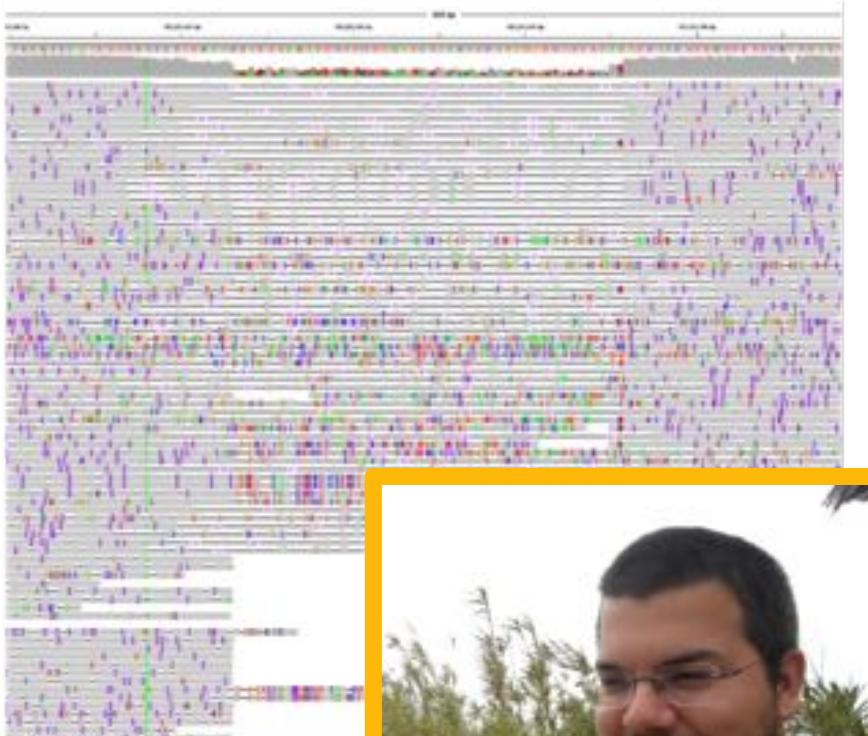
NGMLR:



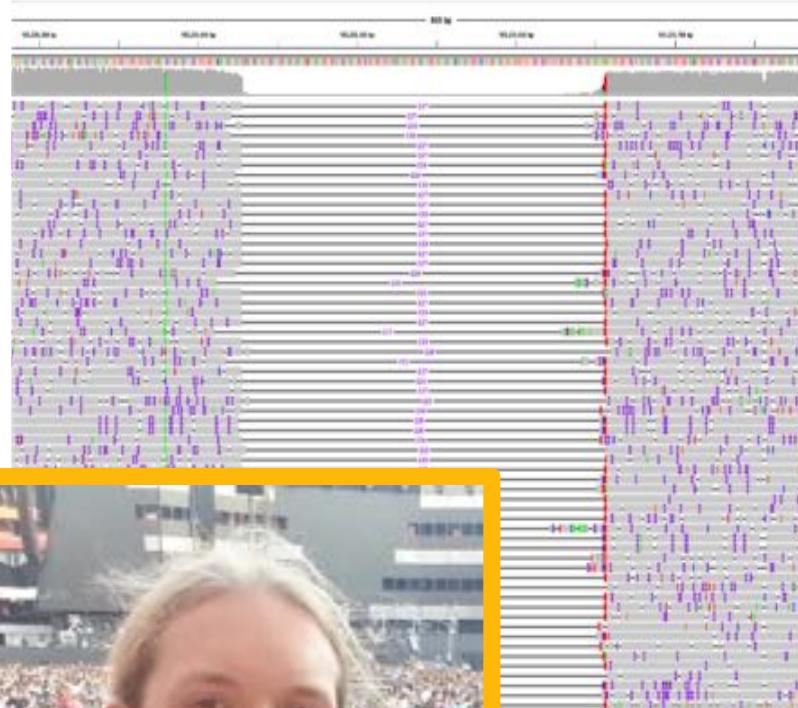
Accurate detection of complex structural variations using single molecule sequencing
Sedlazeck, Rescheneder et al (2017) *In preparation*

NGMLR + Sniffles

BWA-MEM:

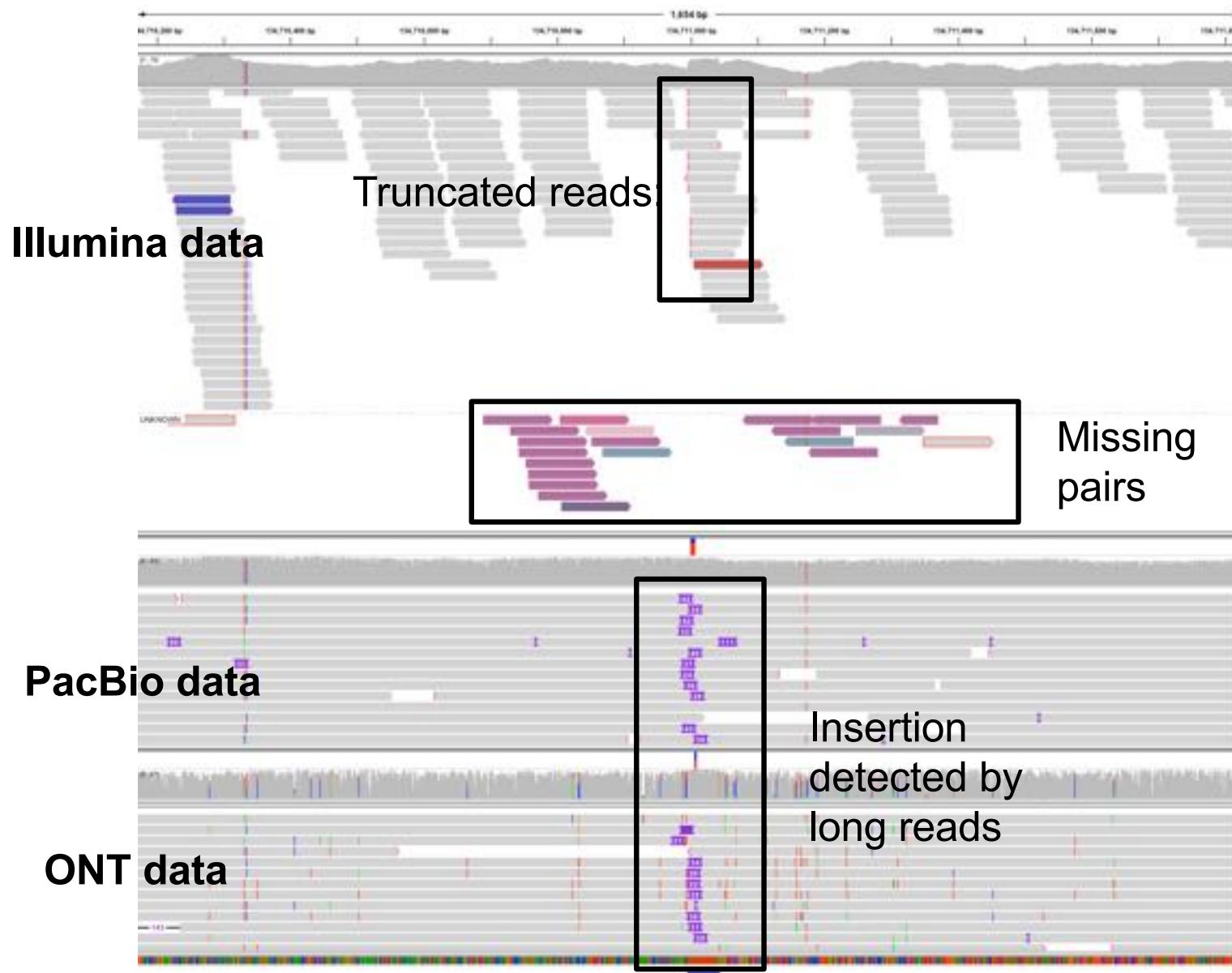


NGMLR:



Accurate detection of complex structural variations using single molecule sequencing
Sedlazeck, Rescheneder et al (2017) *In preparation*

No more false positives!



Accurate detection of complex structural variations using single molecule sequencing
Sedlazeck, Rescheneder et al (2017) *In preparation*

Illumina Roadmap



Illumina Novaseq

\$850k instrument cost
~\$1k / human @ 50x
Short reads, high throughput



10X Chromium

\$125k instrument costs
~\$2k / human
Linked reads, medium throughput

PacBio Roadmap



PacBio Sequel

\$350k instrument cost
~\$30k / human @ 50x
Long reads, Medium throughput



SMRTcell v2

1M Zero Mode Waveguides
~15kb average read length
~\$1000 / SMRTcell

Oxford Nanopore



MinION

\$1k / instrument
~\$30k / human @ 50x
Long reads, Low throughput

PromethION

\$75k / instrument
>>100GB / day
??? / human @ 50x

Oxford Nanopore sequencing, hybrid error correction, and de novo assembly of a eukaryotic genome
Goodwin, S, Gurtowski, J, Ethe-Sayers, S, Deshpande, P, Schatz MC* McCombie, WR* (2015) Genome Research doi: 10.1101/gr.191395.115