

# Genome Sequencing & Assembly

## Michael Schatz

Feb 24, 2016  
Fundamentals of Genome Informatics



# Genome Sequencing & Assembly

Michael Schatz

Feb 24, 2016  
Fundamentals of Genome Informatics



# Schatzlab Overview



## Human Genetics

Role of mutations in disease

Narzisi *et al.* (2015)  
Iossifov *et al.* (2014)



## Plant Biology

Genomes & Transcriptomes

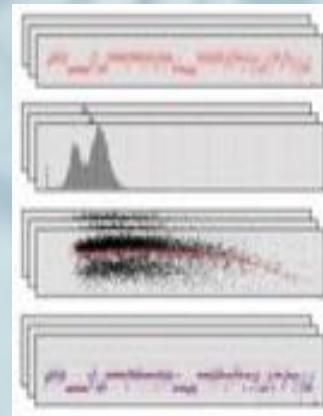
Ming *et al.* (2015)  
Schatz *et al.* (2014)



## Algorithmics & Systems Research

Ultra-large scale biocomputing

Stevens *et al.* (2015)  
Marcus *et al.* (2014)

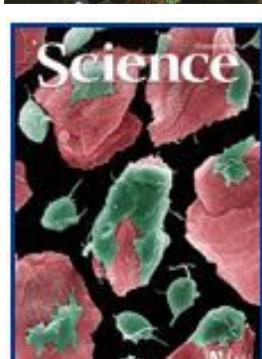
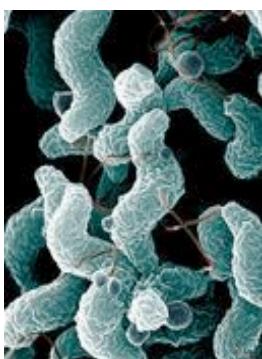
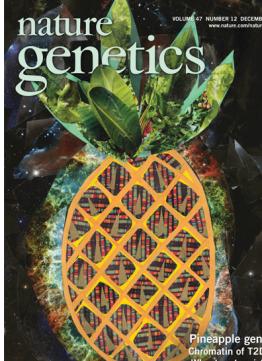
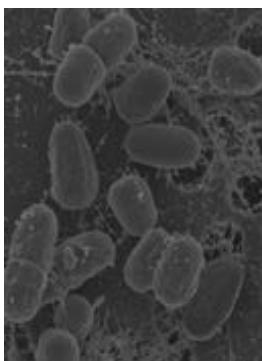


## Single Cell & Single Molecule

CNVs, SVs, & Cell Phylogenetics

Garvin *et al.* (2015)  
Goodwin *et al.* (2015)

# Genomics Across the Tree of Life





# 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



# Outline

## **I. 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

# 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?
  - $5 \text{ copies} \times 138,656 \text{ words} / 5 \text{ words per fragment} = 138k \text{ fragments}$
  - 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

- $D_k = (V, E)$ 
  - $V = \text{All length-}k \text{ subfragments } (k < l)$
  - $E = \text{Directed edges between consecutive subfragments}$ 
    - Nodes overlap by  $k-1$  words

Original Fragment

It was the best of

Directed Edge

It was the best → was the best of

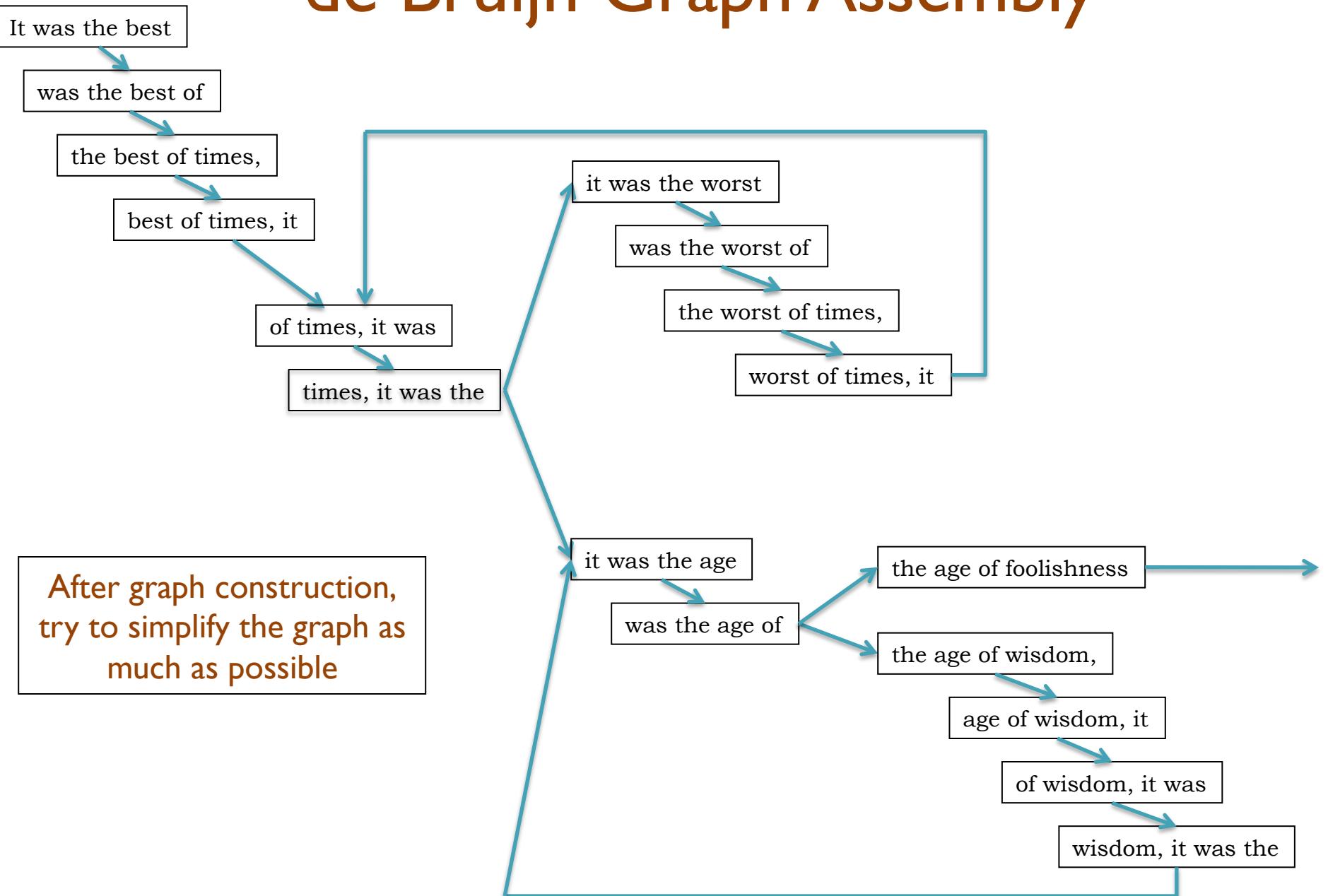
- Locally constructed graph reveals the global sequence structure
  - Overlaps between sequences implicitly computed

de Bruijn, 1946

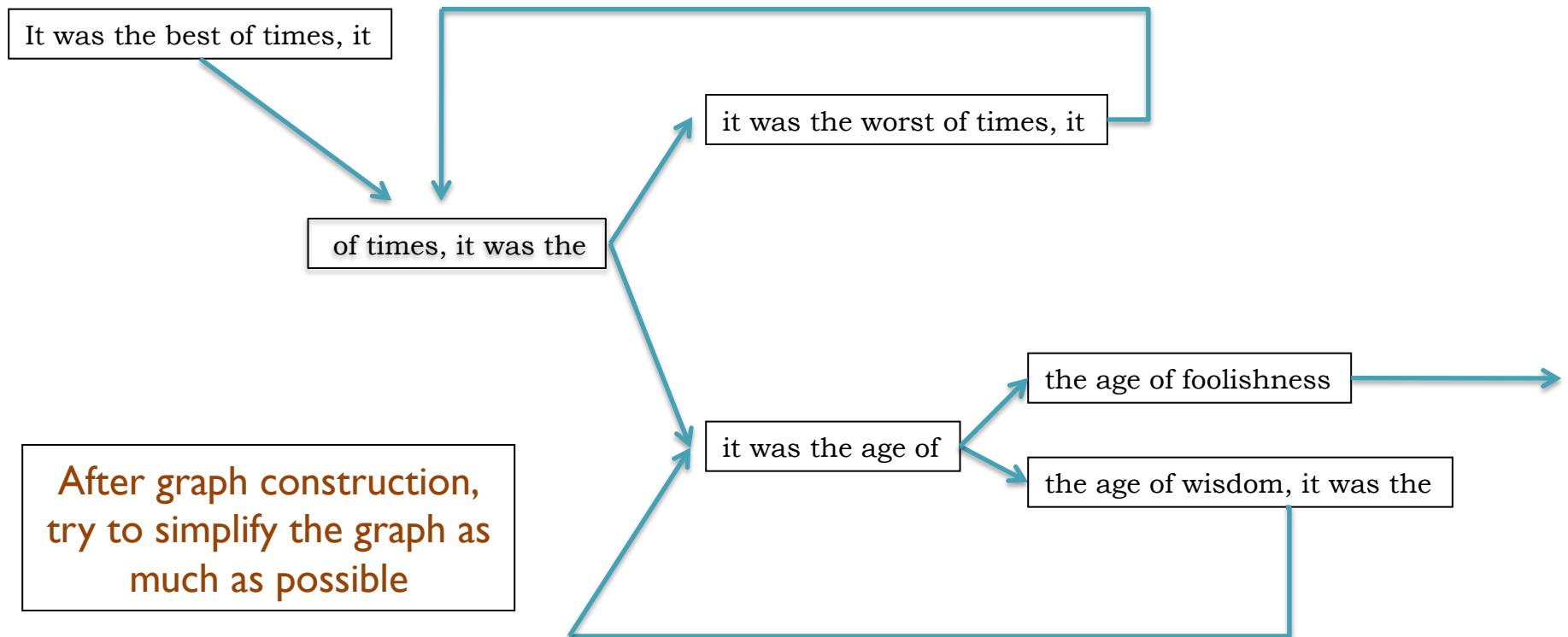
Idury and Waterman, 1995

Pevzner, Tang, Waterman, 2001

# de Bruijn Graph Assembly

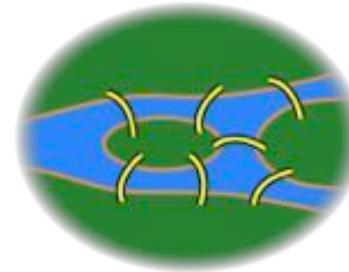
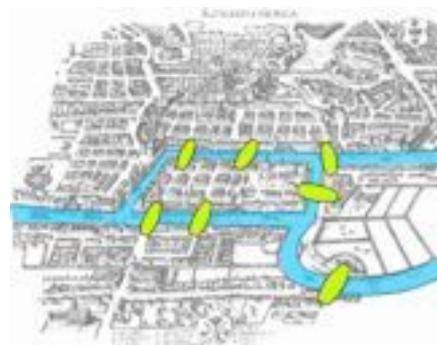


# de Bruijn Graph Assembly

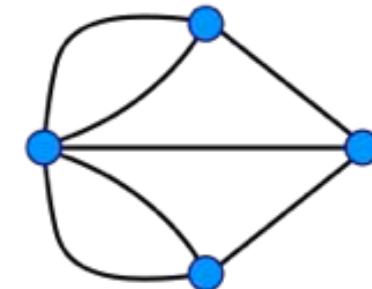


# Eulerian Cycle Problem

- **Seven Bridges of Königsberg**
  - Find a cycle that visits every **edge** exactly once



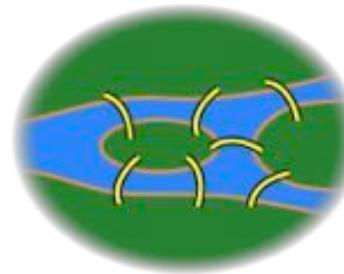
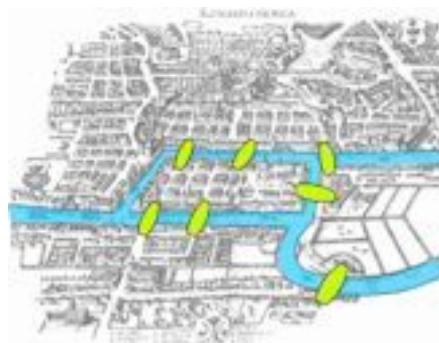
[Can you find the cycle?]



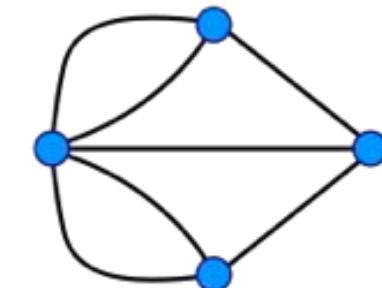
[bioalgorithms.info](http://bioalgorithms.info)

# Eulerian Cycle Problem

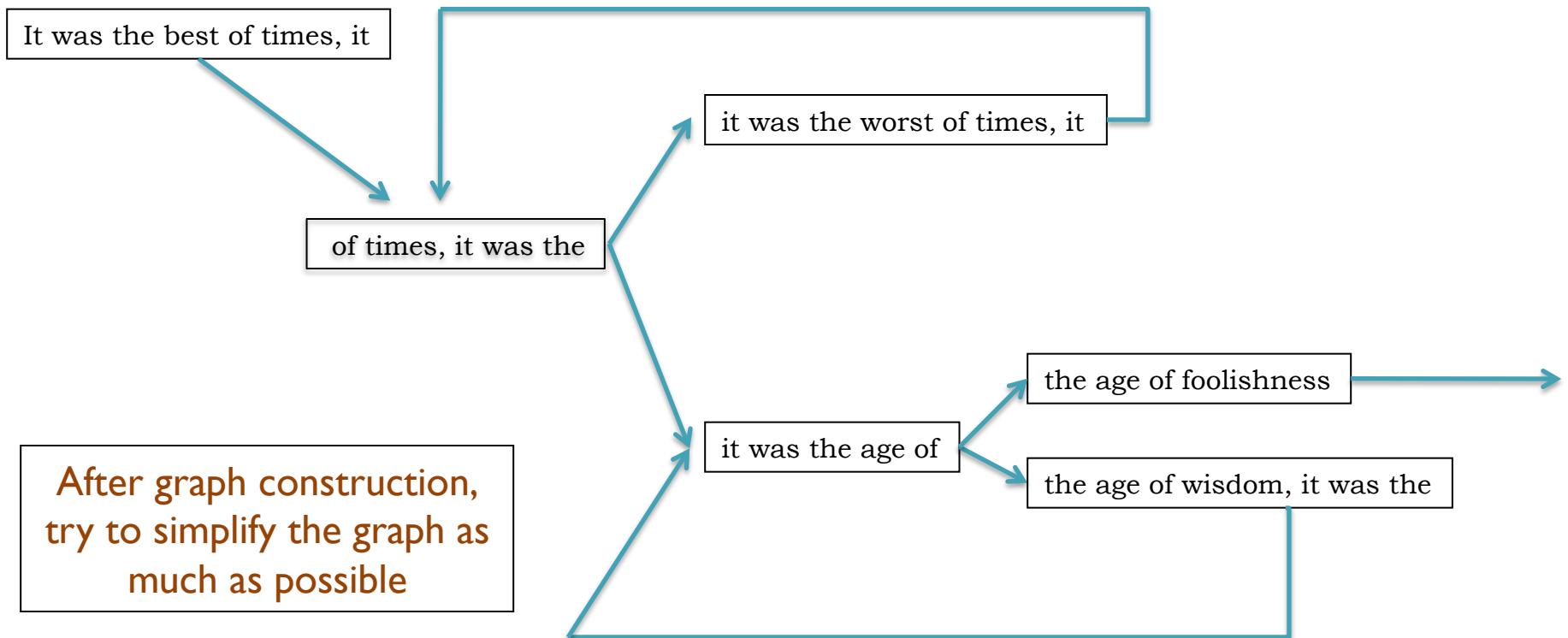
- **Seven Bridges of Königsberg**
  - Find a cycle that visits every **edge** exactly once
  - **Easy test:** every node has same number of incoming and outgoing edges (directed graph) or even degree (undirected graph)
  - **Easy search:** arbitrarily walk to find first cycle, then walk again from an unmarked edge



[Can you find the cycle?]



# de Bruijn Graph Assembly



# The full tale

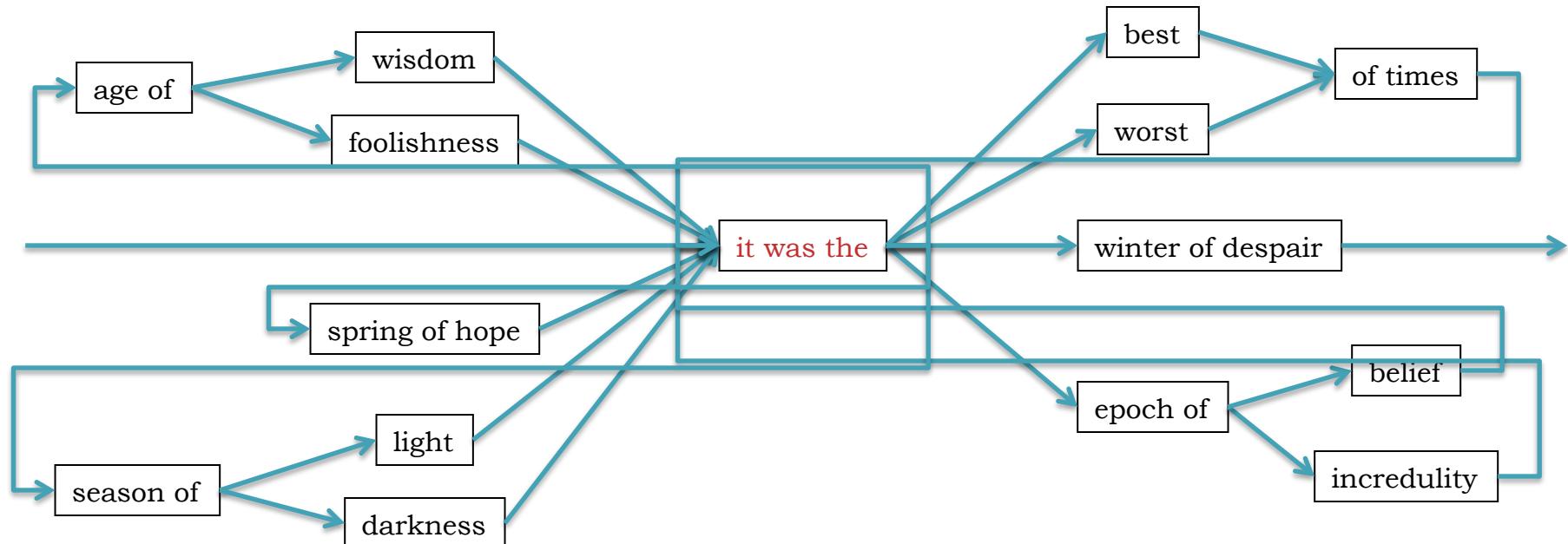
... it was the best of times it was the worst of times ...

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

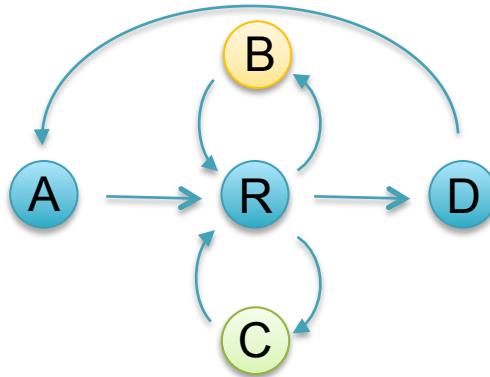
... it was the epoch of belief it was the epoch of incredulity ...

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

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



# Counting Eulerian Cycles



ARBRCRD  
or  
ARCRBRD

Generally an exponential number of compatible sequences

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

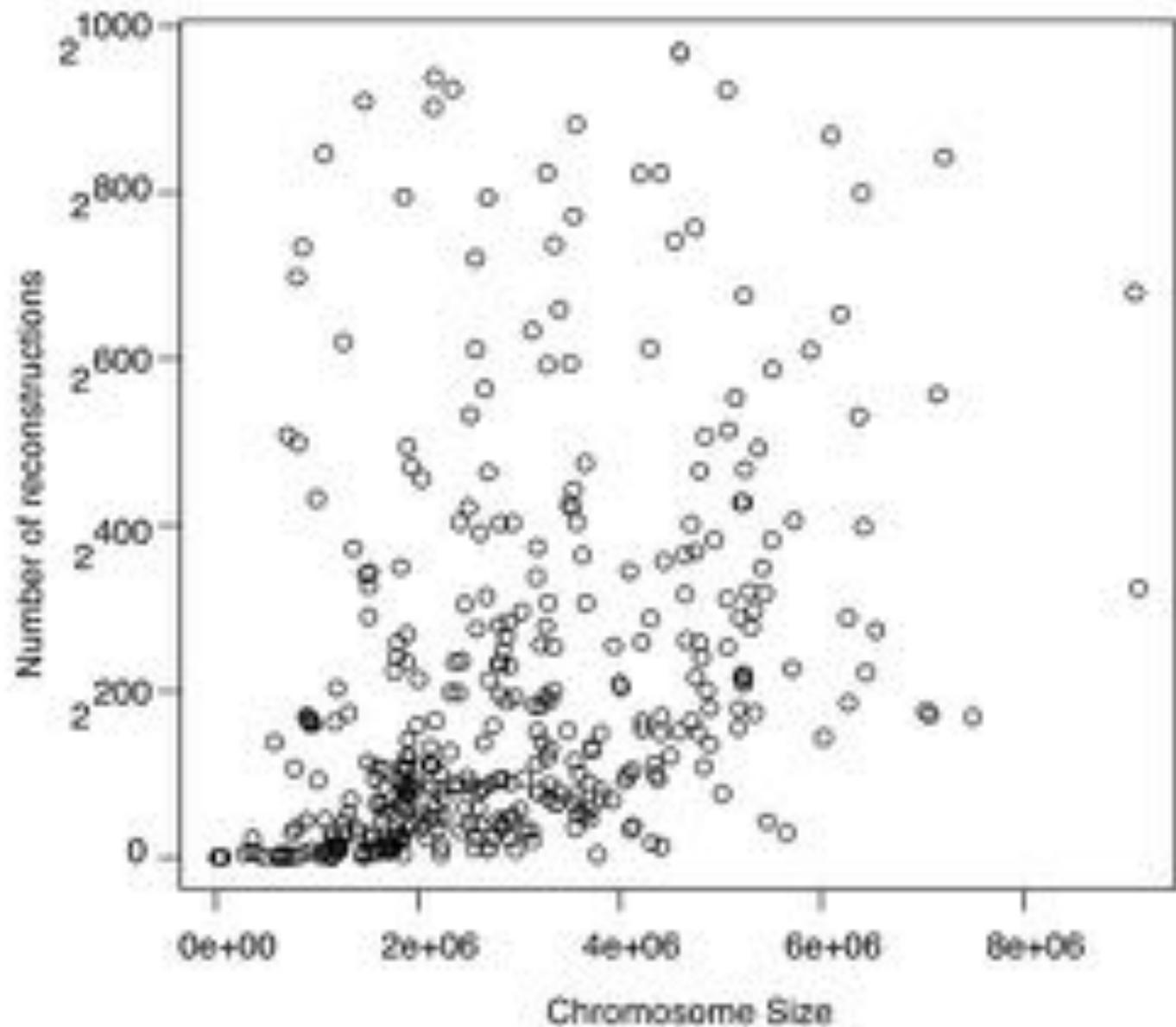
$$\mathcal{W}(G, t) = (\det L) \left\{ \prod_{u \in V} (r_u - 1)! \right\} \left\{ \prod_{(u,v) \in E} a_{uv}! \right\}^{-1}$$

$L$  =  $n \times n$  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

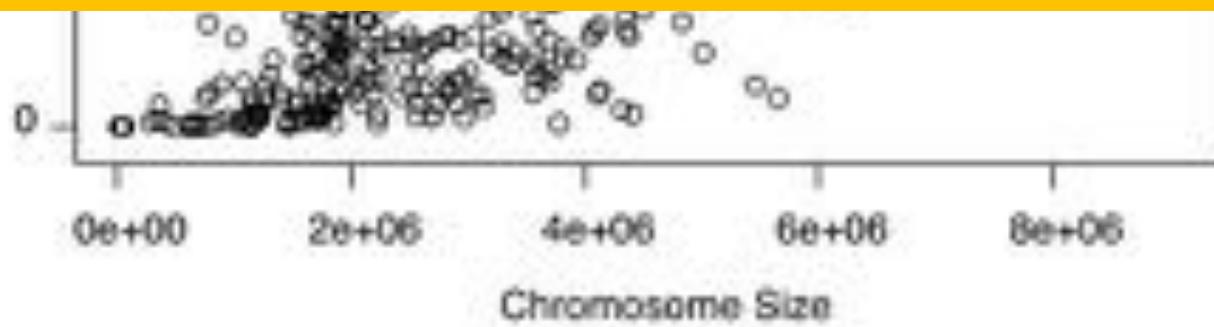
$a_{uv}$  = multiplicity of edge from  $u$  to  $v$

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



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

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

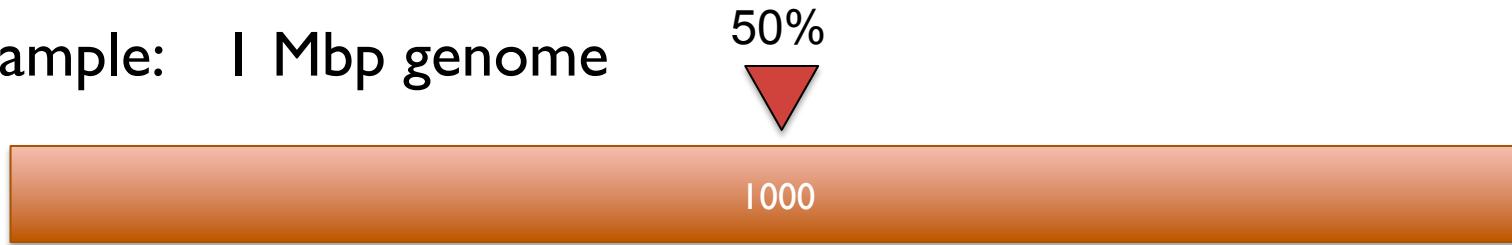


**Assembly Complexity of Prokaryotic Genomes using Short Reads.**  
Kingsford C, Schatz MC, Pop M (2010) *BMC Bioinformatics*.

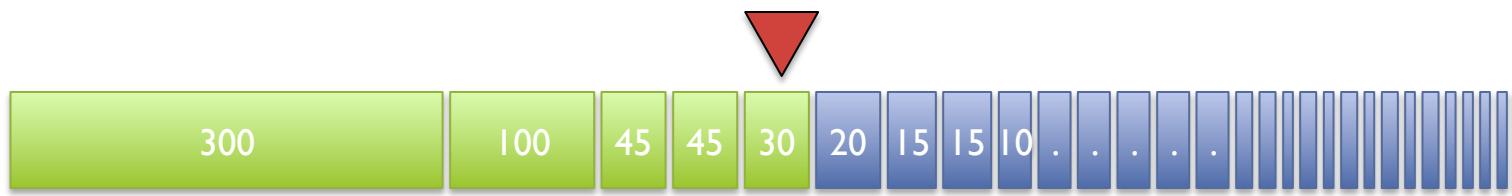
# Contig N50

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

## Example: 1 Mbp genome

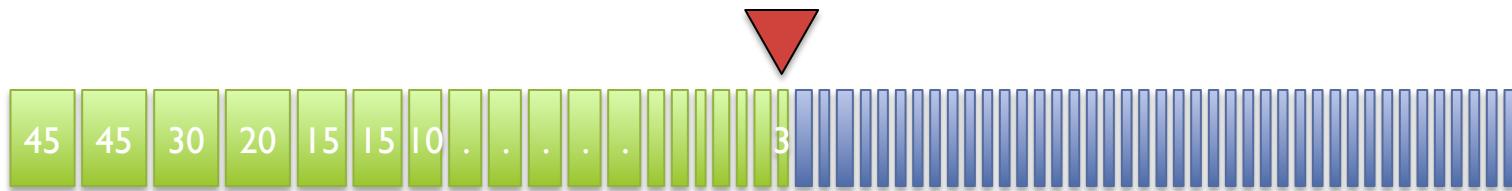


A



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

## ***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*

- Canu: recommended for PacBio/ONT project

# Milestones in Genome Assembly

Science Vol. 207, February 20, 1980  
articles

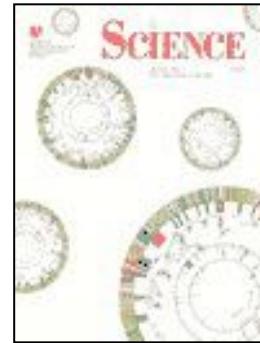
Nucleotide sequence of bacteriophage  $\phi$ X174 DNA

E.Sanger, G.M.Air, R.G.Bonell, S.E.Brown, A.R.Cochise, J.C.Fiddes,  
C.A.Hinchliffe, B.H.P.M.Shorey\* & M.Smith\*

\*MRC Laboratory of Molecular Biology, 40 Cambridge Heath Road, London E1 2AP, U.K.

In 1977, Sanger et al. published the first genome sequence of a bacteriophage,  $\phi$ X174. This was done by sequencing overlapping DNA fragments from cloned phage DNA. The authors used the rapid and simple "socalled" method. The authors describe some of the features responsible for the production of the same or similar sequences for the synthesis of proteins and enzymes. They also describe the properties and functions of  $\phi$ X174 using different cloning strategies.

The authors of the paper state: "The use of these results as a reference for genome technology?". Is it still useful today?

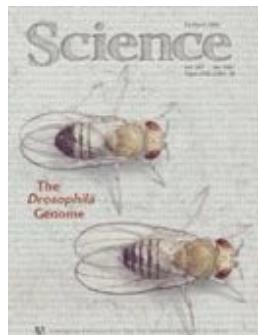


1977. Sanger et al.  
1<sup>st</sup> Complete Organism  
5375 bp

1995. Fleischmann et al.  
1<sup>st</sup> Free Living Organism  
TIGR Assembler. 1.8Mbp



1998. C.elegans SC  
1<sup>st</sup> Multicellular Organism  
BAC-by-BAC Phrap. 97Mbp



2000. Myers et al.  
1<sup>st</sup> Large WGS Assembly.  
Celera Assembler. 116 Mbp



2001. Venter et al., IHGSC  
Human Genome  
Celera Assembler/GigaAssembler. 2.9 Gbp



2010. Li et al.  
1<sup>st</sup> Large SGS Assembly.  
SOAPdenovo 2.2 Gbp

Like Dickens, we must computationally reconstruct a genome from short fragments

# Assembly Applications

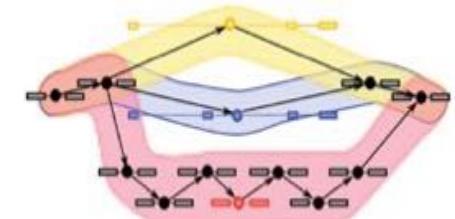
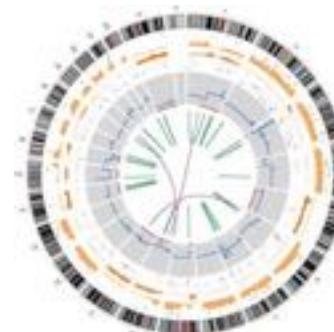
- Novel genomes



- Metagenomes

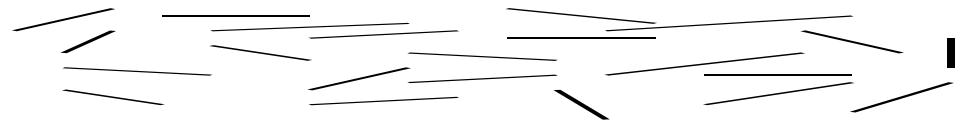


- Sequencing assays
  - Structural variations
  - Transcript assembly
  - ...



# Assembling a Genome

## 1. Shear & Sequence DNA



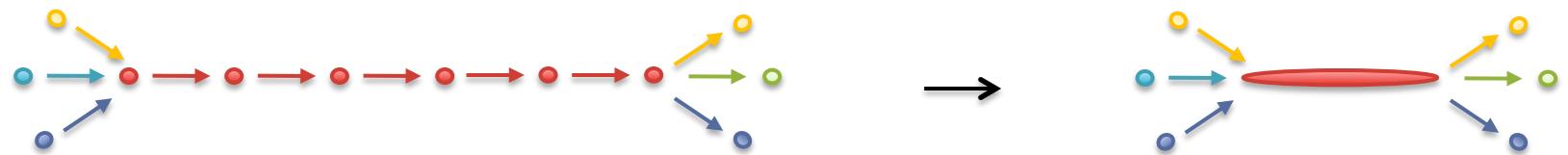
## 2. Construct assembly graph from reads (de Bruijn / overlap graph)

...AGCCTAGGGATGCGCGACACGT

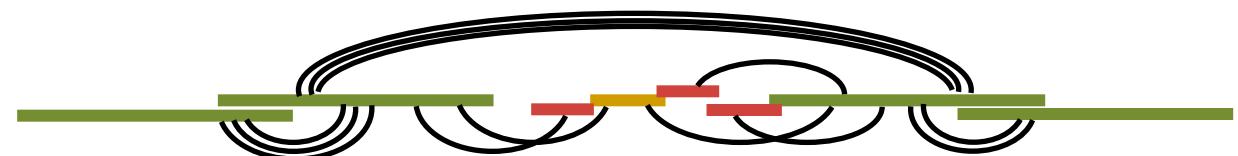
GGATGCGCGACACGT CGCATATCCGGTTGGT CAACCTCGGACGGAC

CAACCTCGGACGGAC CTCAGCGAA...

## 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:**

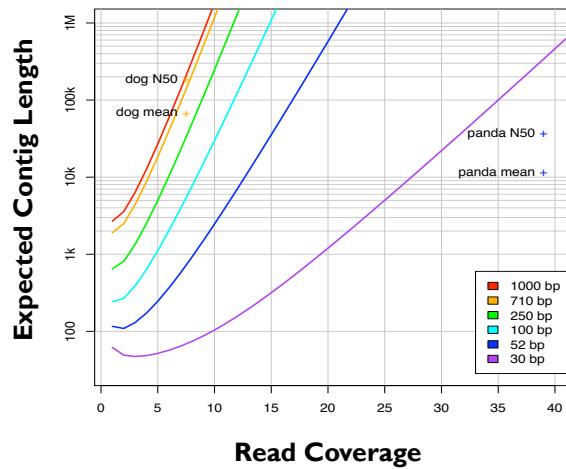
- (Very) Large genomes, complex structure

## 4. **Accuracy:**

- (Very) Hard to assess correctness

# Ingredients for a good assembly

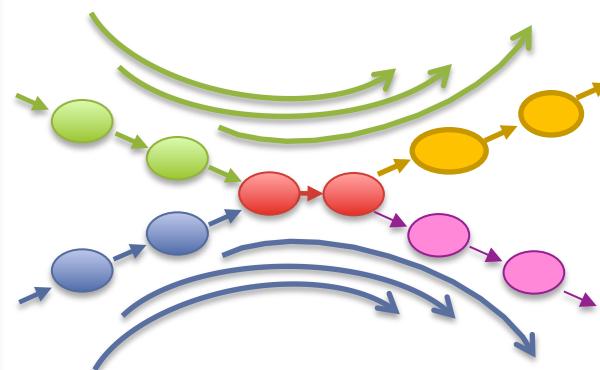
## Coverage



### High coverage is required

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

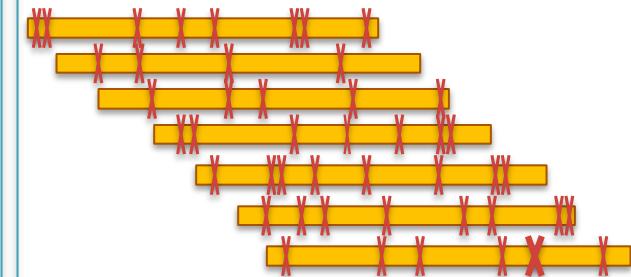
## Read Length



### 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

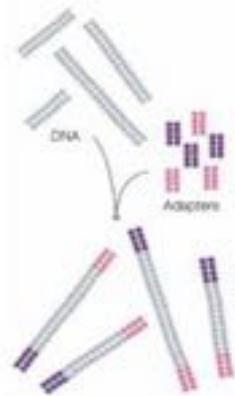


### Errors obscure overlaps

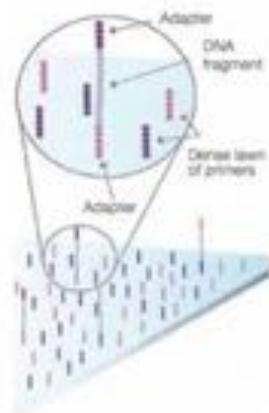
- Reads are assembled by finding kmers shared in pair of reads
- High error rate requires very short seeds, increasing complexity and forming assembly hairballs

**Current challenges in *de novo* plant genome sequencing and assembly**  
Schatz MC, Witkowski, McCombie, WR (2012) *Genome Biology*. 12:243

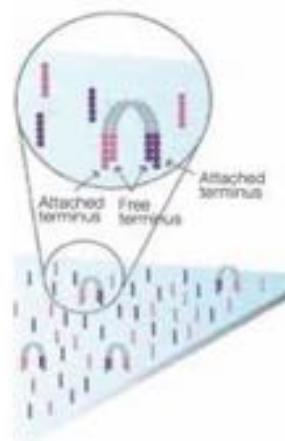
# Illumina Sequencing by Synthesis



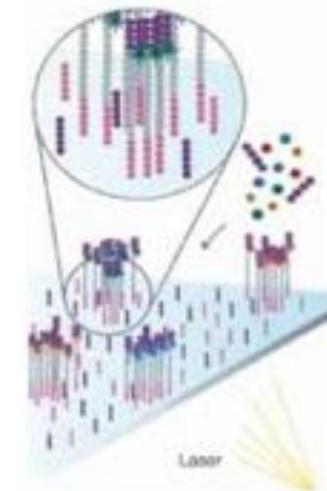
1. Prepare



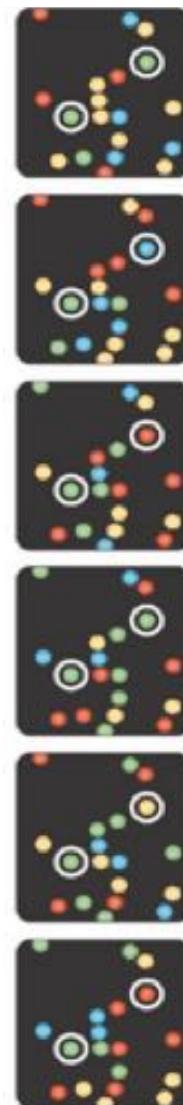
2. Attach



3. Amplify



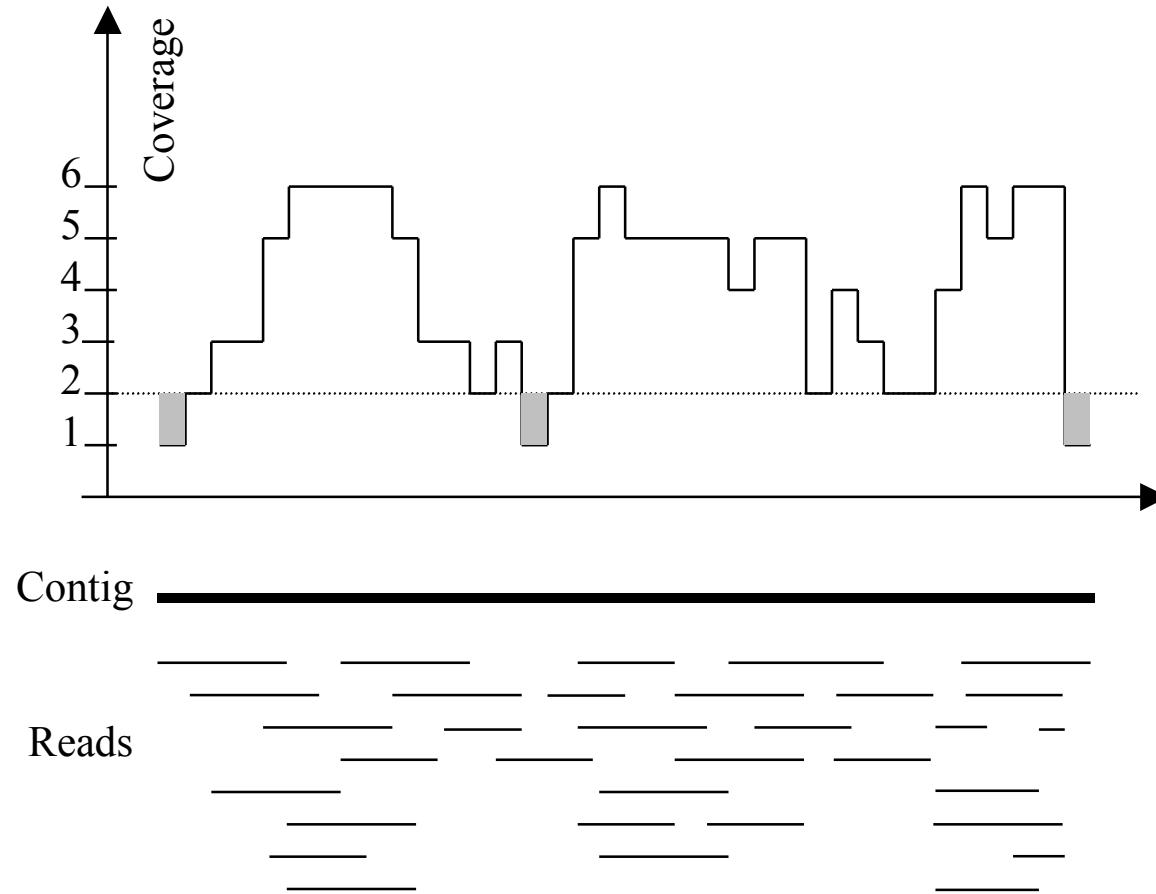
4. Image



5. Basecall

Metzker (2010) Nature Reviews Genetics 11:31-46  
<http://www.youtube.com/watch?v=l99aKKHcxC4>

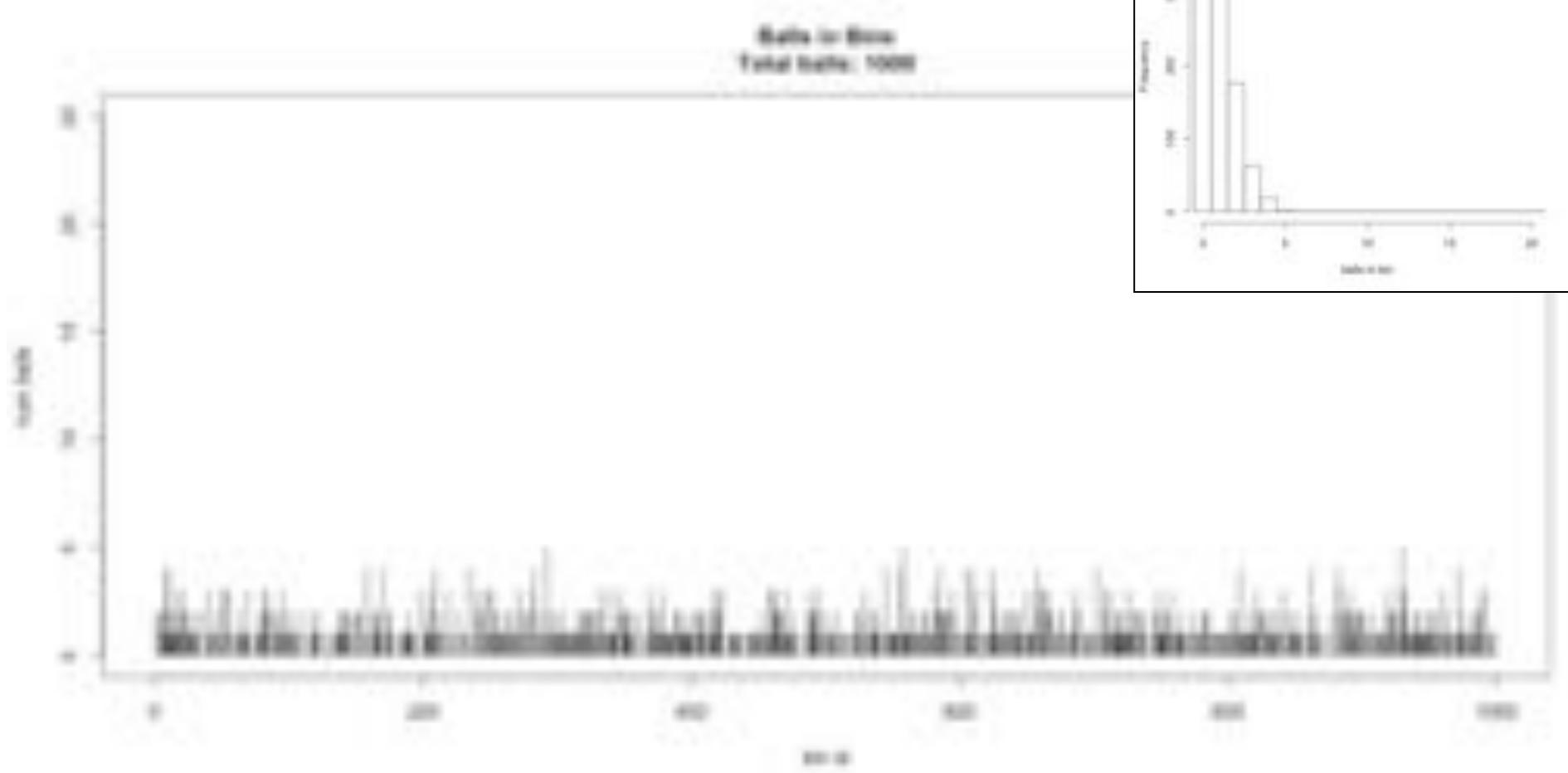
# Typical sequencing coverage



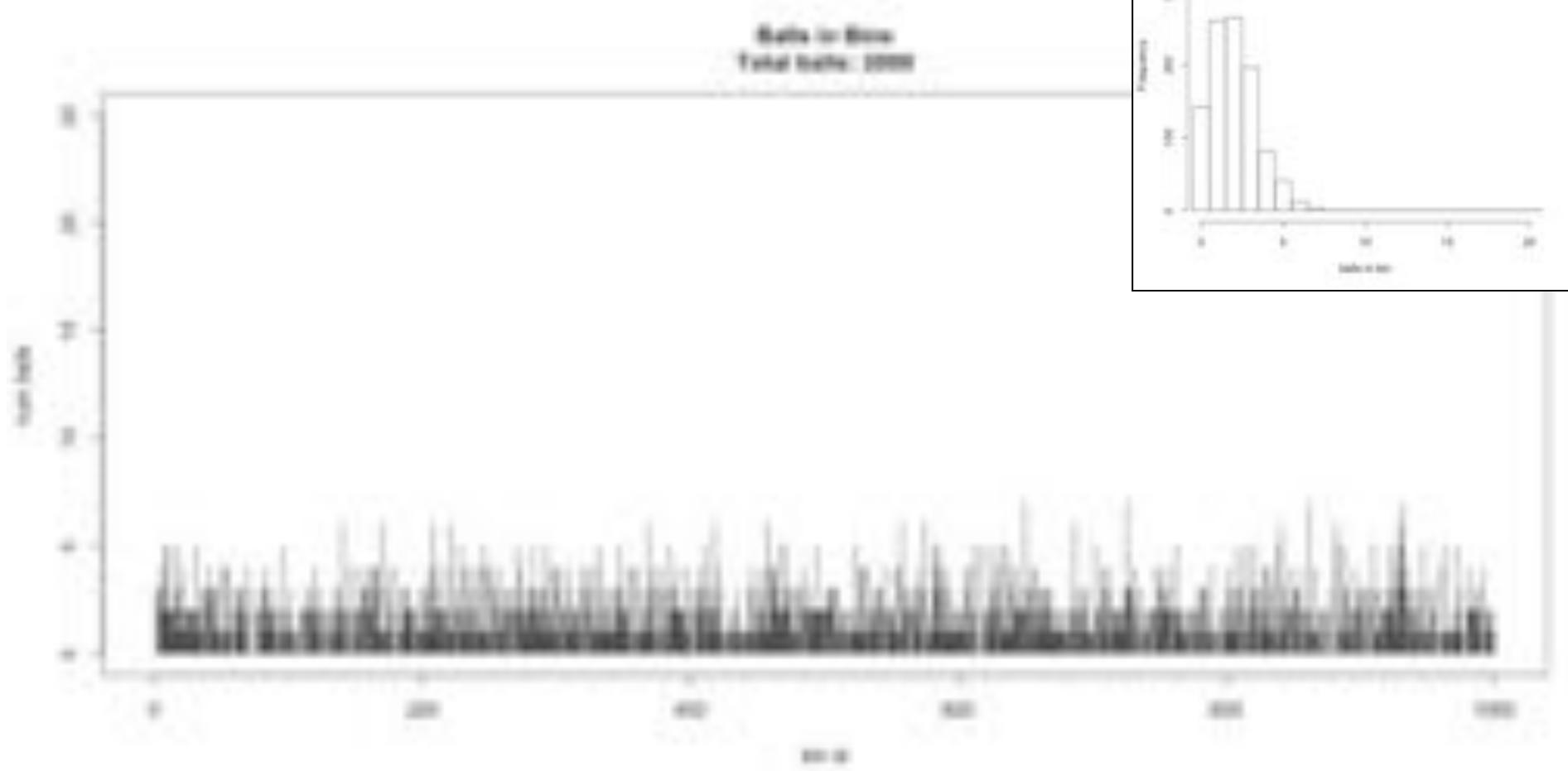
Imagine raindrops on a sidewalk

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

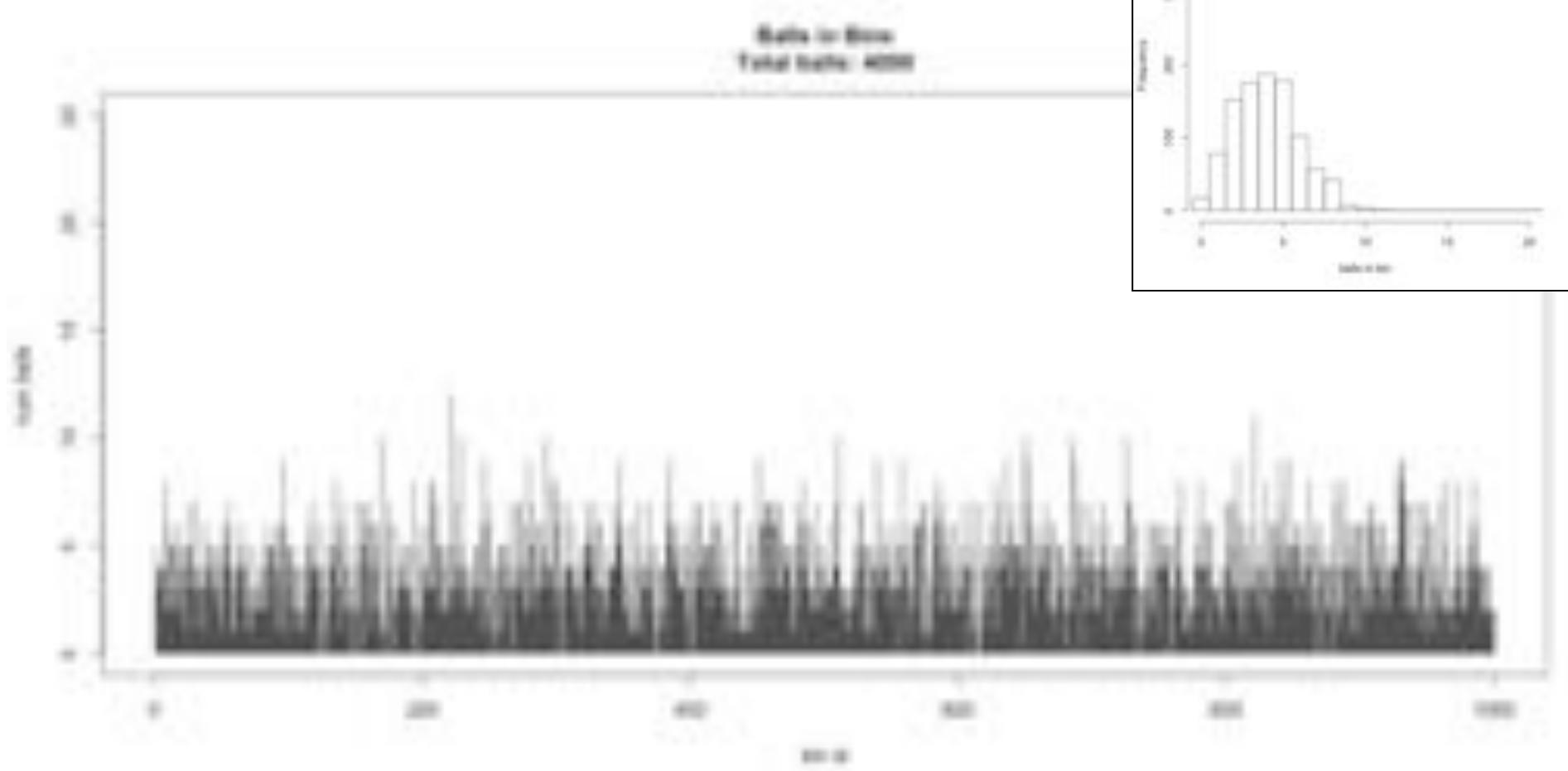
# Ix sequencing



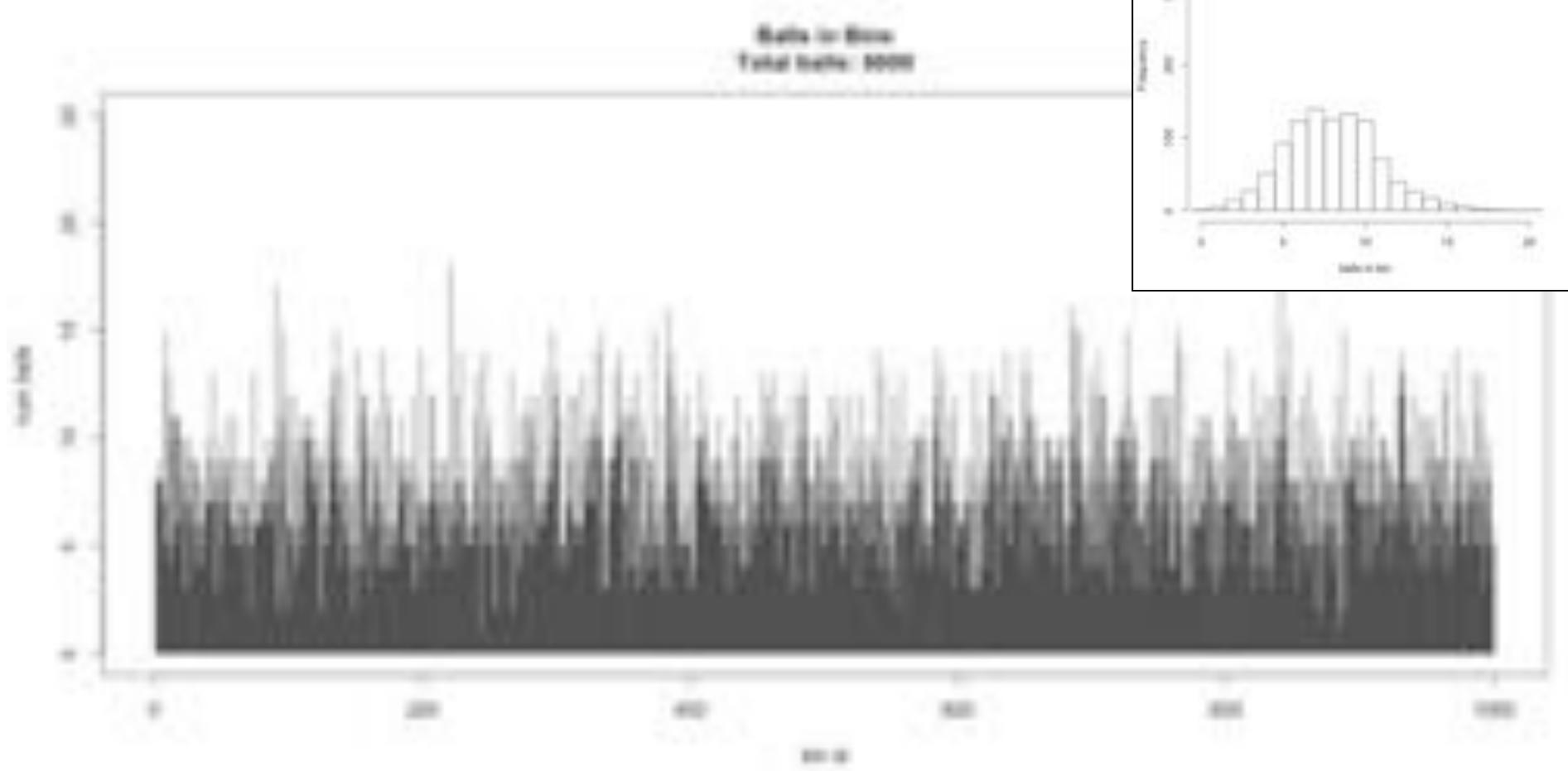
# 2x sequencing



# 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.

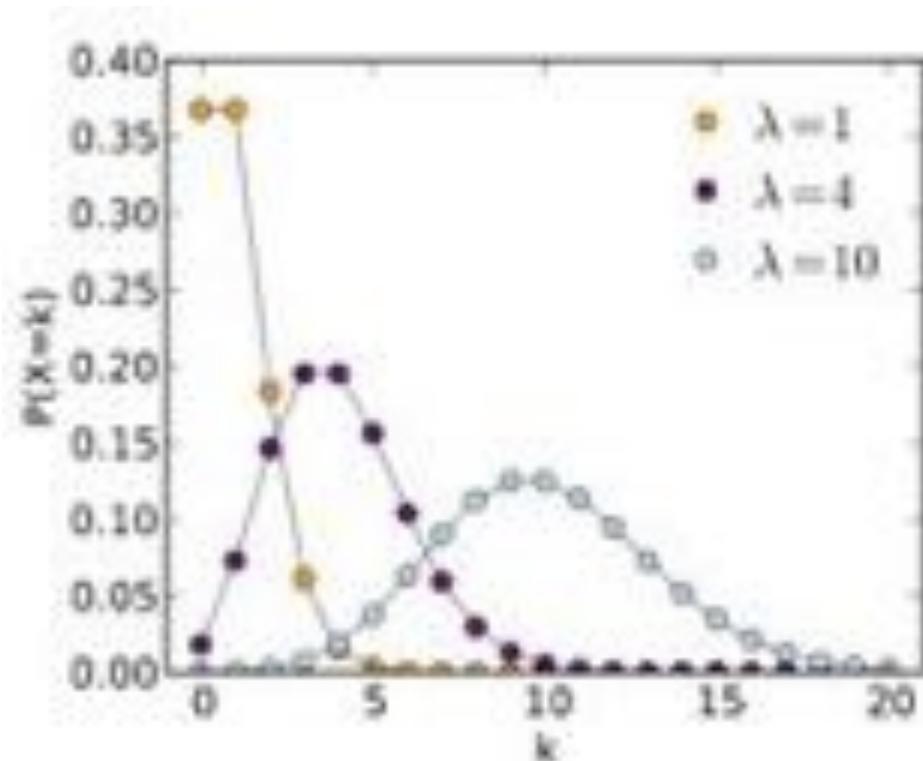
Formulation comes from the limit of the binomial equation

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

**Key property:**

- ***The standard deviation is the square root of the mean.***

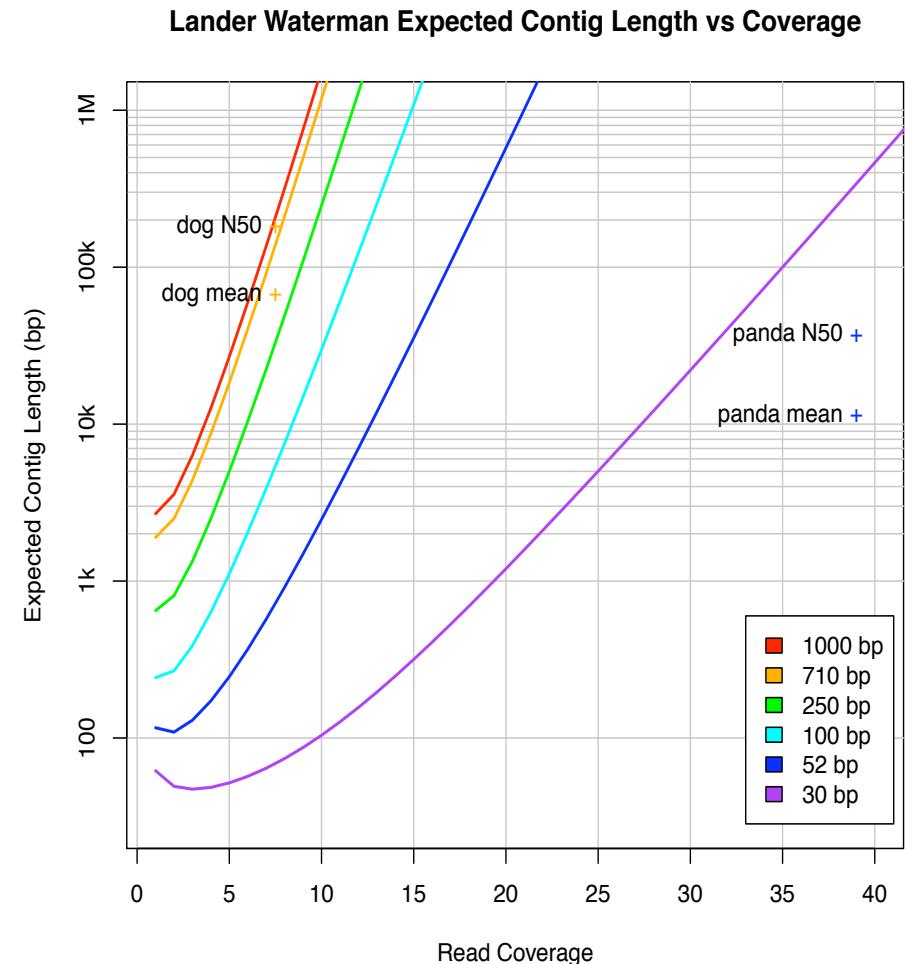
$$P(k) = \frac{\lambda^k}{k!} e^{-\lambda}$$



# Coverage and Read Length

## Idealized Lander-Waterman model

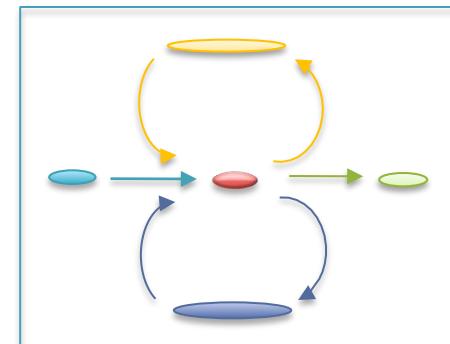
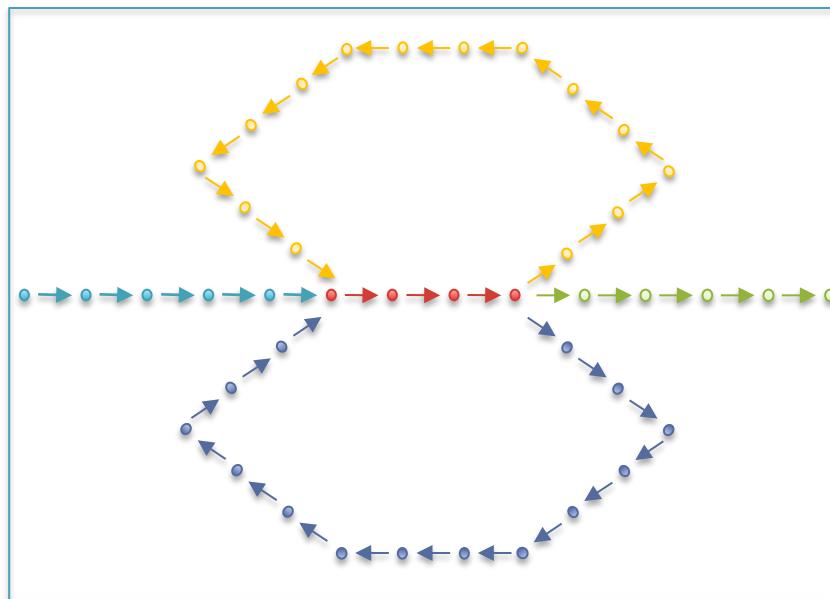
- Reads start at perfectly random positions
- Contig length is a function of coverage and read length
  - Short reads require much higher coverage to reach same expected contig length
- Need even high coverage for higher ploidy, sequencing errors, sequencing biases
  - Recommend 100x coverage



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

# 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



# Errors in the graph



(Chaisson, 2009)

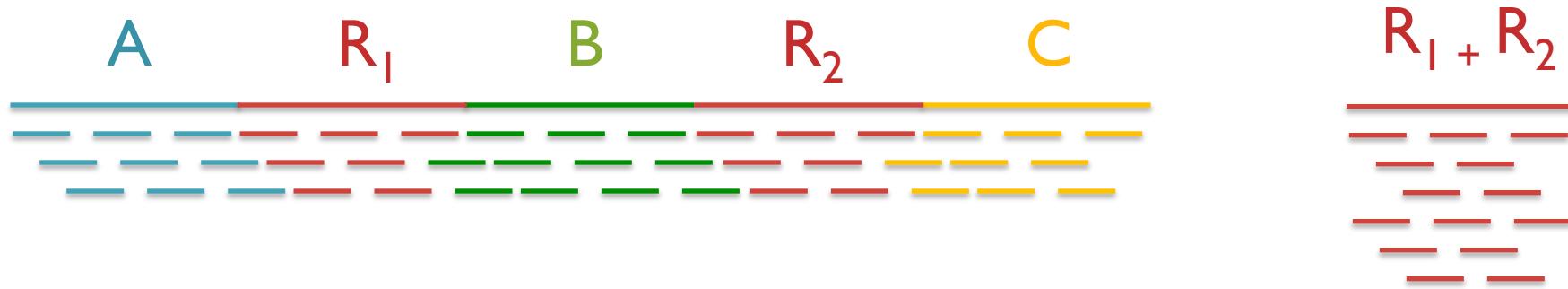
Clip Tips	Pop Bubbles
<p>was the worst of times,</p> <p>was the worst of <b>tymes</b>,</p> <p>the worst of times, it</p>	<p>was the worst of times,</p> <p>was the worst of <b>tymes</b>,</p> <p>times, it was the age</p> <p><b>tymes</b>, it was the age</p>
<p>the worst of <b>tymes</b>,</p> <p>was the worst of</p> <p>the worst of times,</p> <p>worst of times, it</p>	<p><b>tymes</b>,</p> <p>was the worst of</p> <p>it was the age</p> <p>times,</p>

# 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  $\Delta$ .
  - If we see many more reads than  $k$  (if the arrival rate is  $> A$ ) , it is likely to be a collapsed repeat

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

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

**The fragment assembly string graph**  
Myers, EW (2005) Bioinformatics. 21(suppl 2): ii79-85.

# Paired-end and Mate-pairs

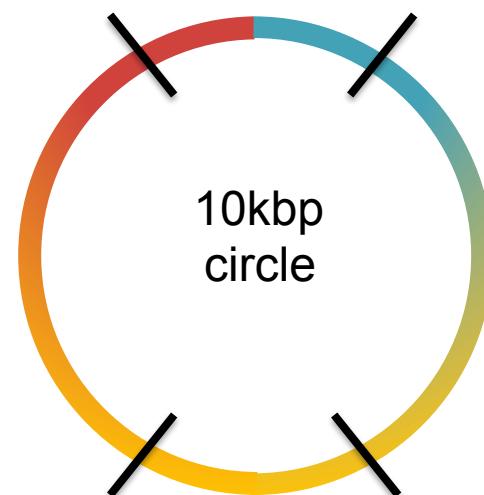
## Paired-end sequencing

- 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)

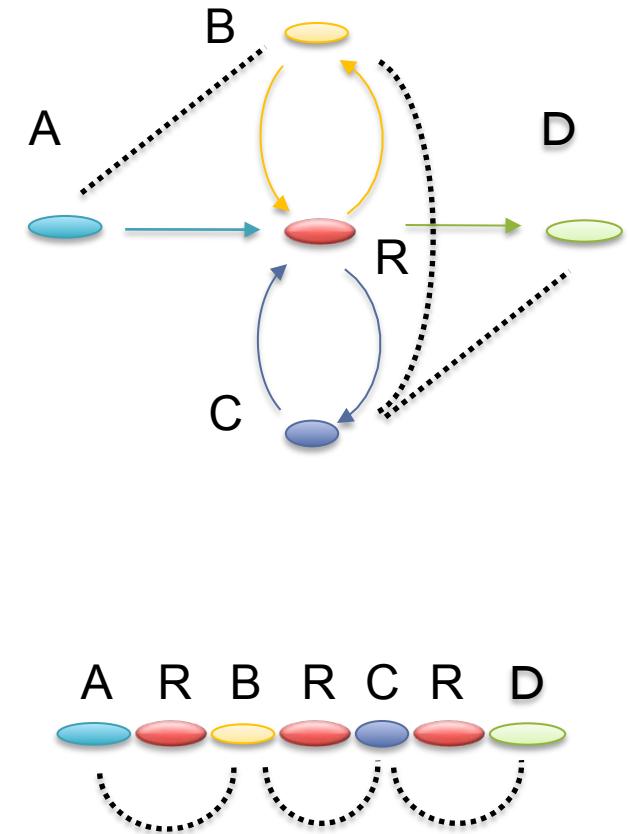


2x100 @ 300bp (innies)



# 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





# Outline

## *1. Assembly theory*

- Assembly by analogy

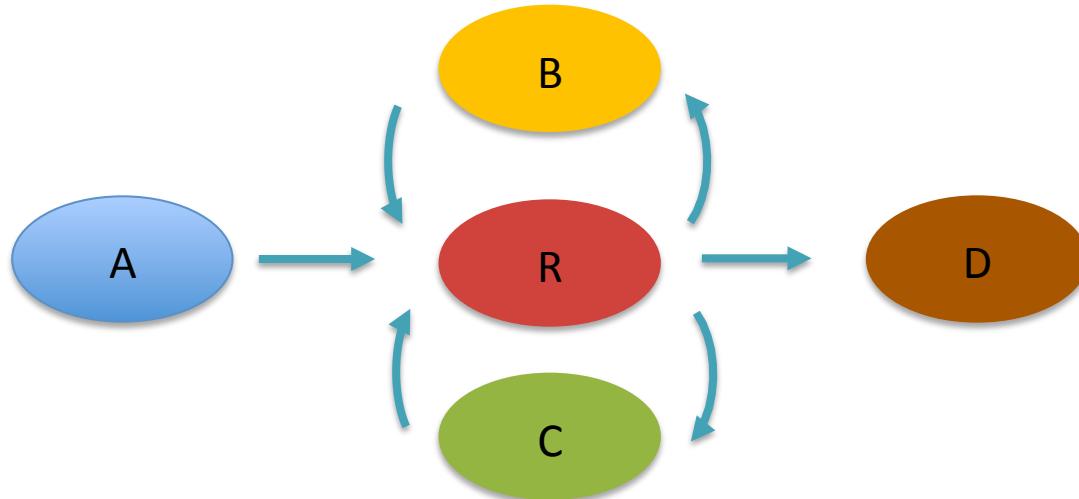
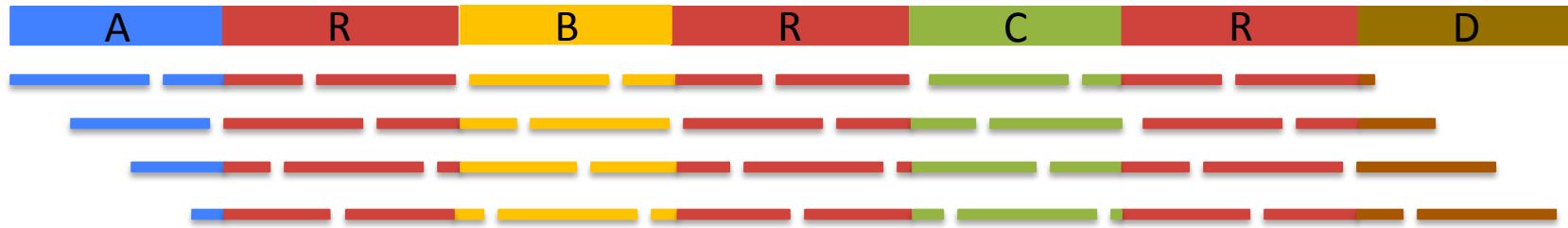
## *2. Practical Issues*

- Coverage, read length, errors, and repeats

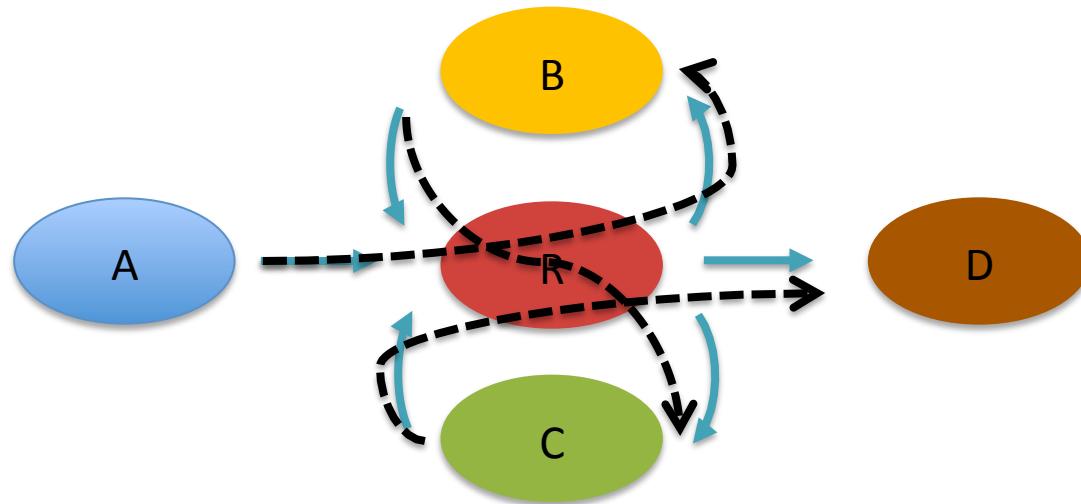
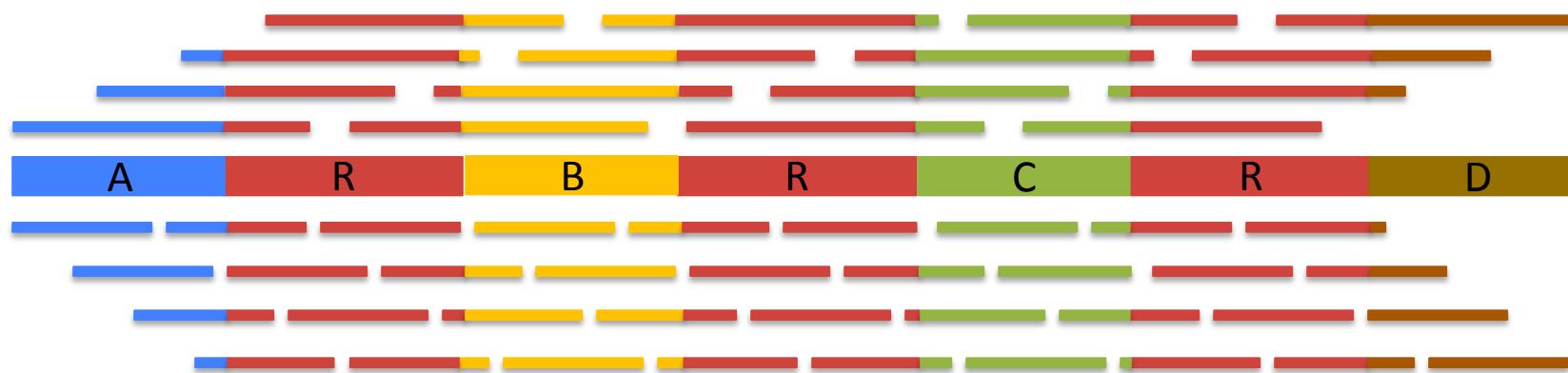
## **3. Next-next-gen Assembly**

- Canu: recommended for PacBio/ONT project

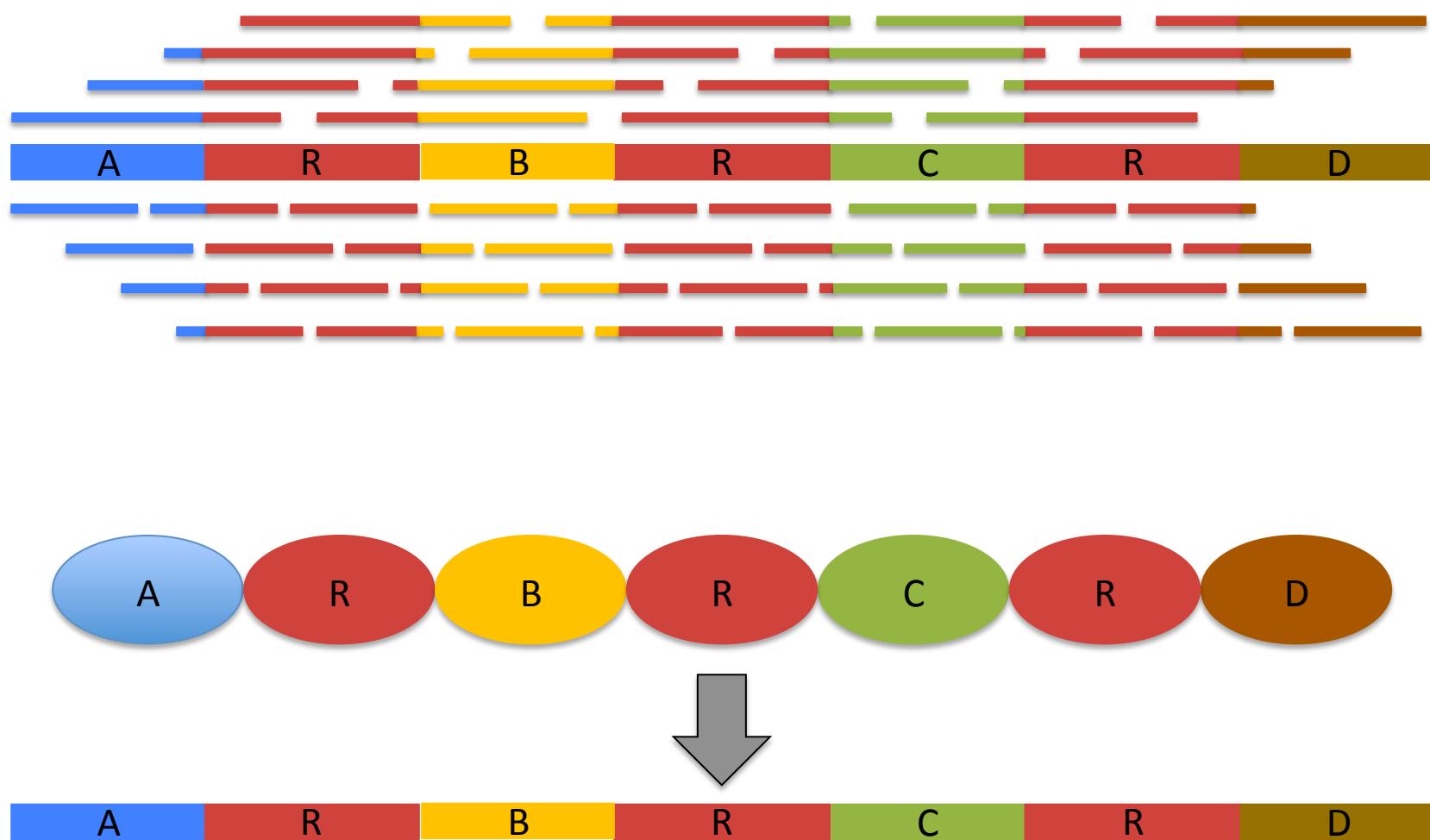
# Assembly Complexity



# 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 2016

## Long Read Sequencing: De novo assembly, SV analysis, phasing

**Illumina/Moleculo**



(Kuleshov et al. 2014)

**Pacific Biosciences**



(Berlin et al, 2014)

**Oxford Nanopore**



(Quick et al, 2014)

## Long Span Sequencing: Chromosome Scaffolding, SV analysis, phasing

**Molecular Barcoding**



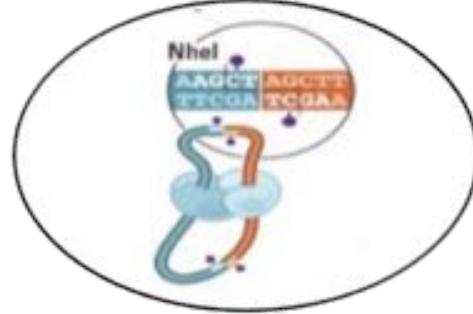
(10Xgenomics.com)

**Optical Mapping**



(Cao et al, 2014)

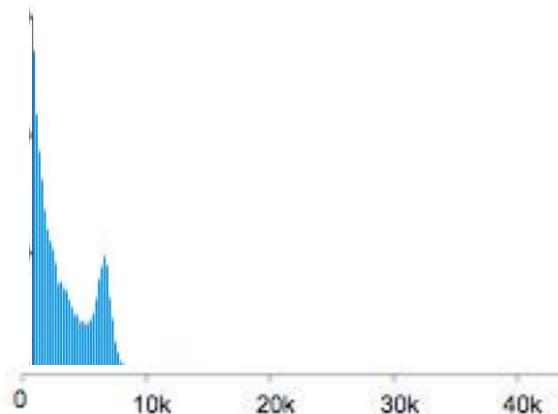
**Chromatin Assays**



(Putnam et al, 2015)

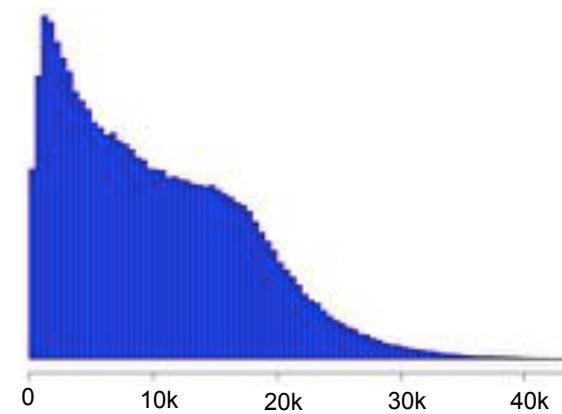
# Long Read Sequencing Technology

Moleculo



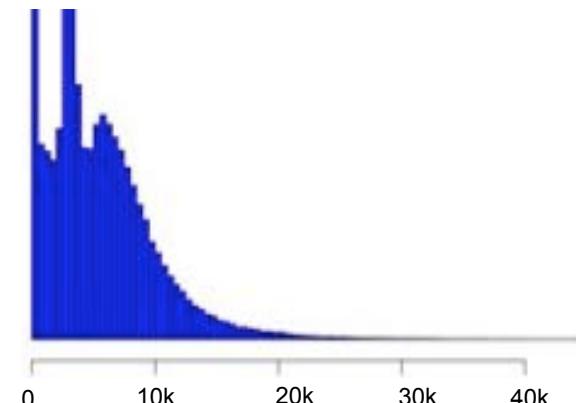
(Voskoboynik et al. 2013)

PacBio RS II



CSHL/PacBio

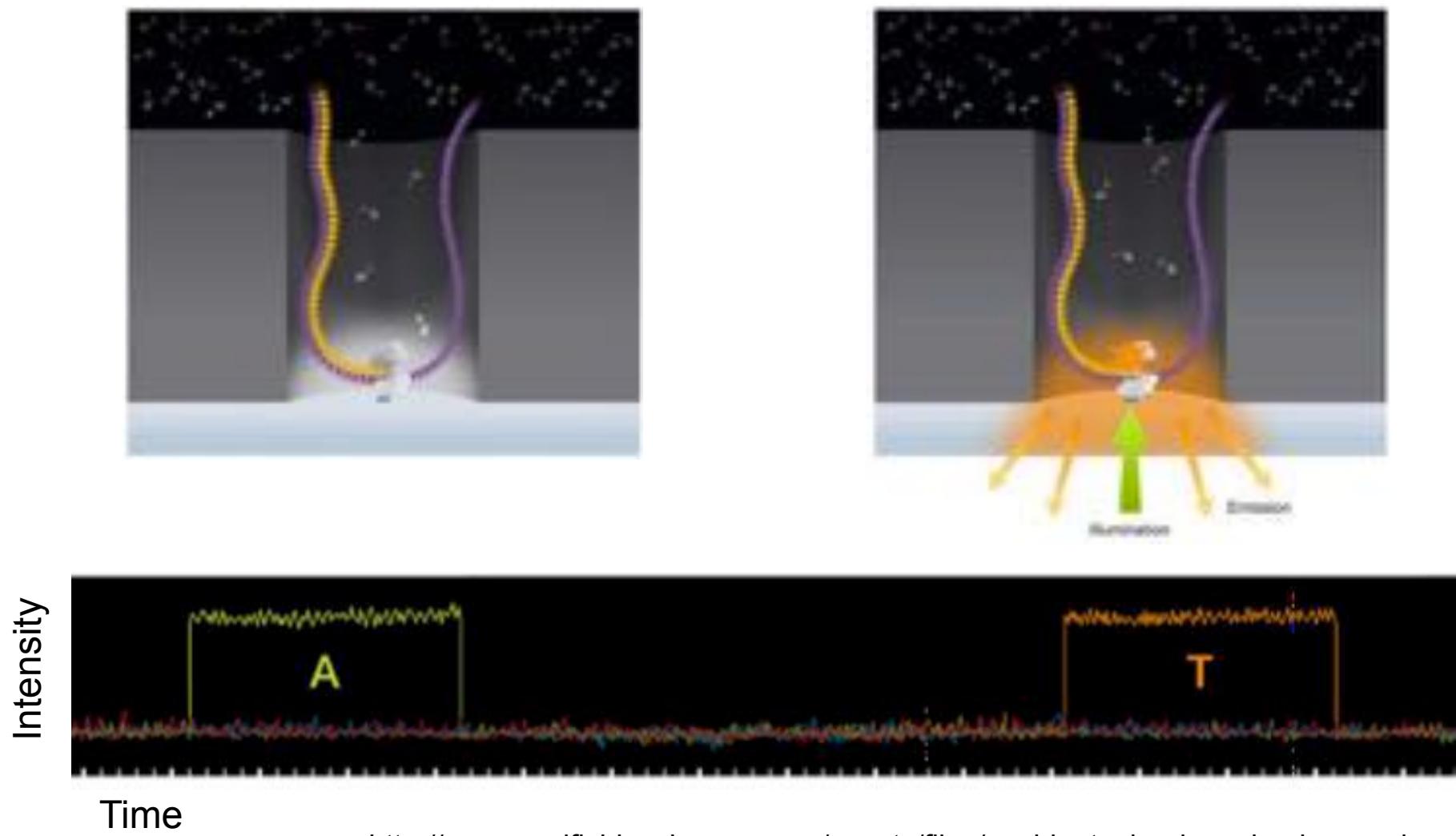
Oxford Nanopore



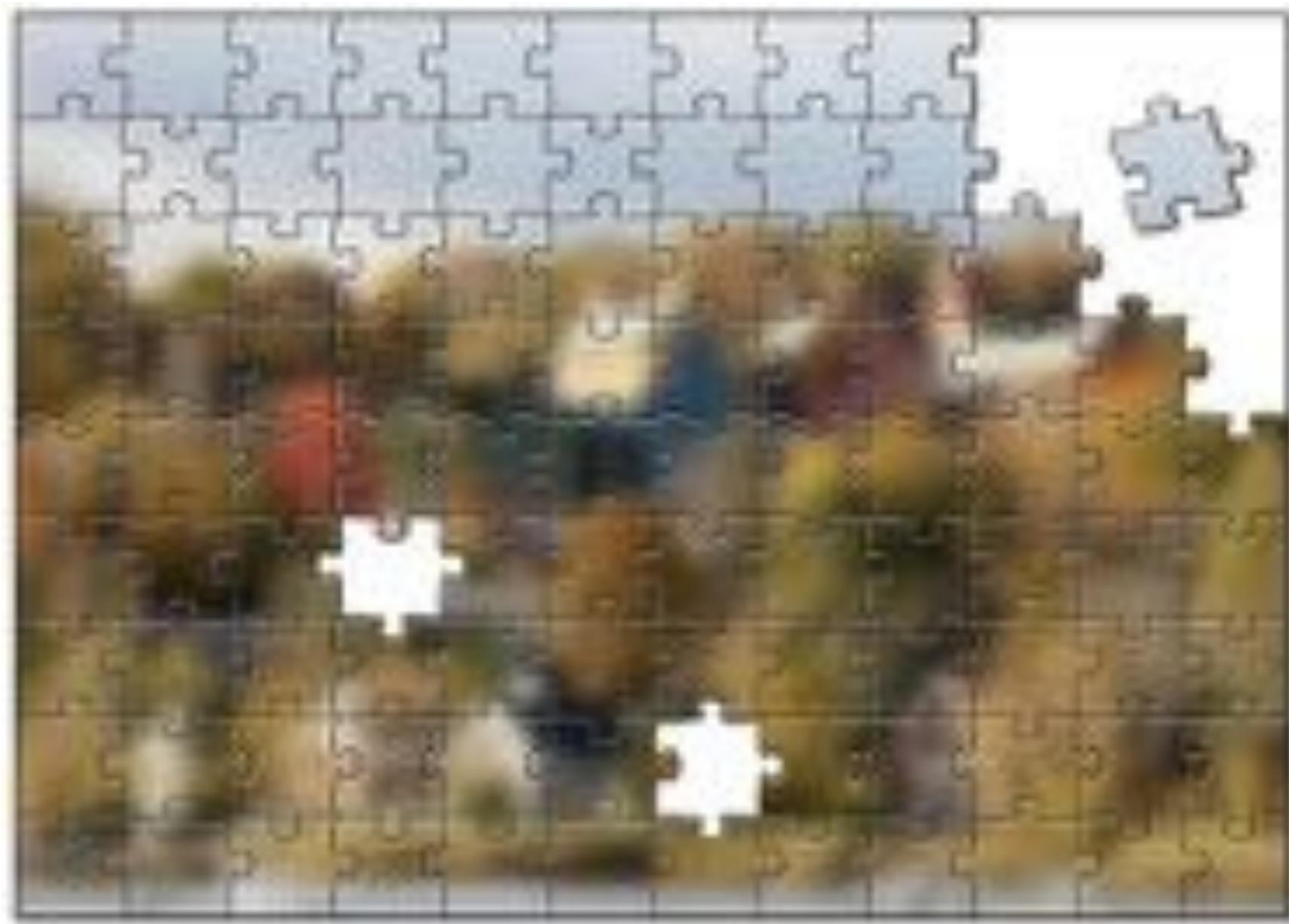
CSHL/ONT

# PacBio SMRT Sequencing

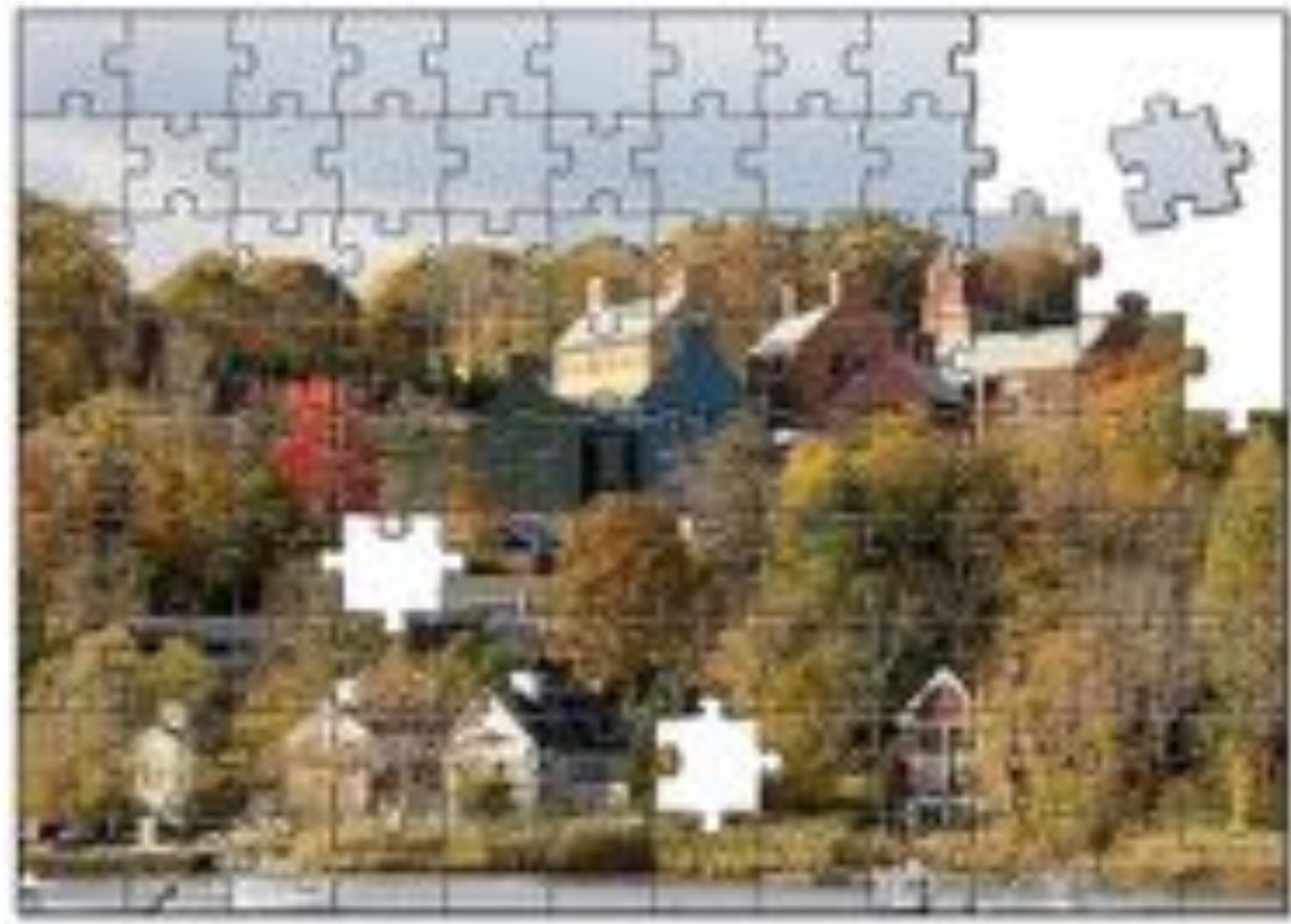
Imaging of fluorescently phospholinked labeled nucleotides as they are incorporated by a polymerase anchored to a Zero-Mode Waveguide (ZMW).



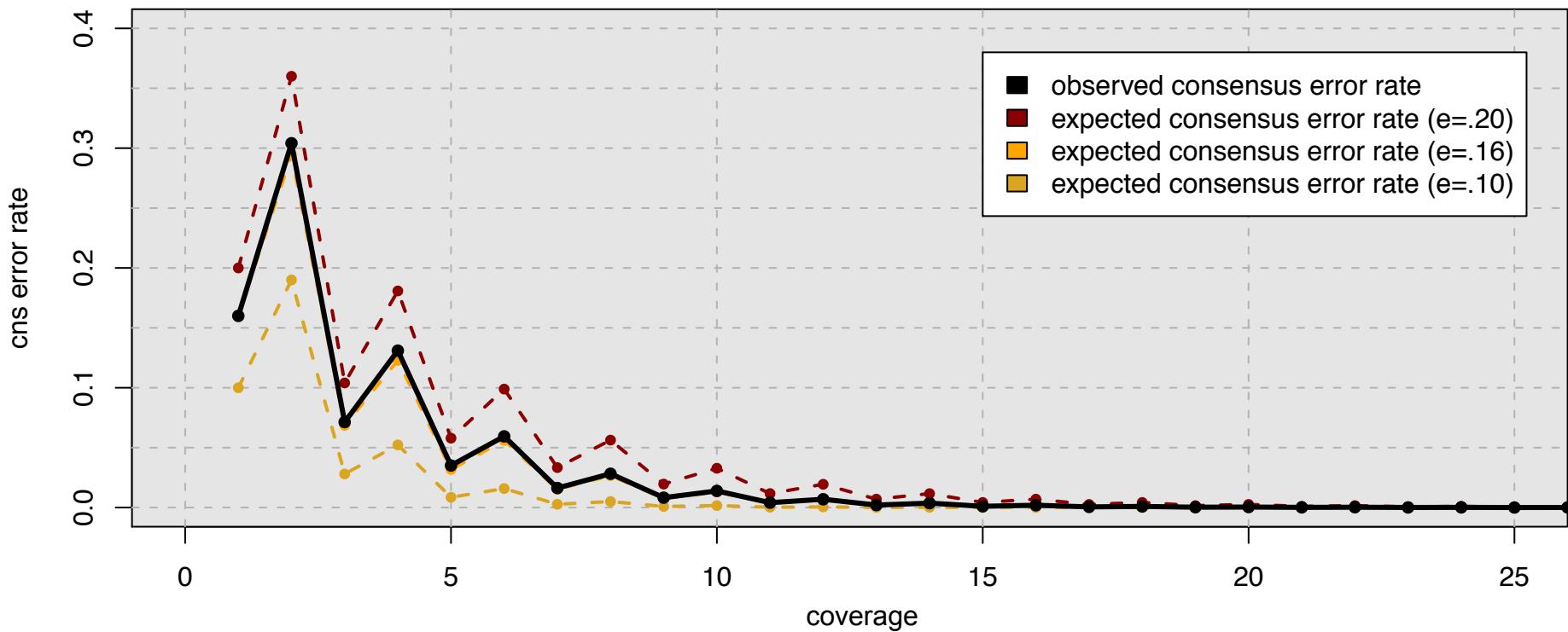
# Single Molecule Sequences



# “Corrective Lens” for Sequencing



# Consensus Accuracy and Coverage



Coverage can overcome random errors

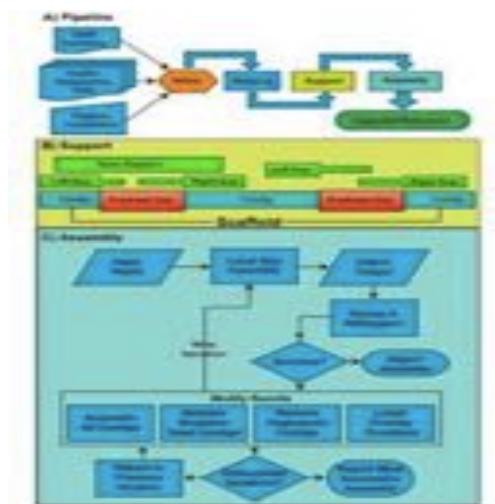
- Dashed: error model from binomial sampling
- Solid: observed accuracy

Koren, Schatz, et al (2012)  
Nature Biotechnology. 30:693–700

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

# PacBio Assembly Algorithms

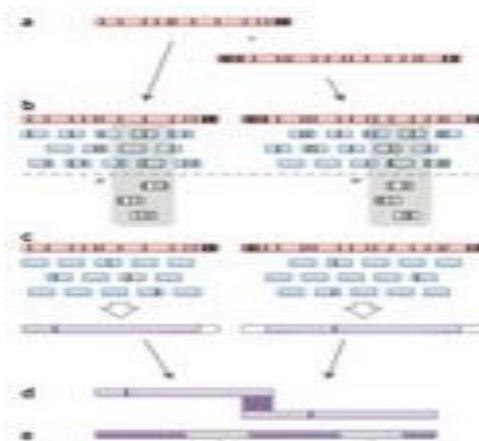
## PBJelly



### Gap Filling and Assembly Upgrade

English et al (2012)  
PLOS One. 7(11): e47768

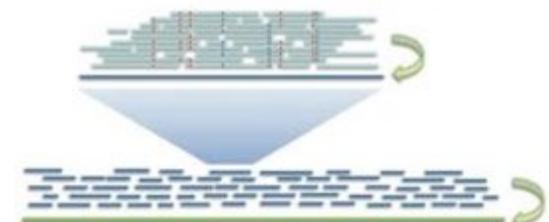
## PacBioToCA & ECTools



### Hybrid/PB-only Error Correction

Koren, Schatz, et al (2012)  
Nature Biotechnology. 30:693–700

## HGAP/MHAP/Falcon & Quiver



$$\Pr(\mathbf{R} \mid T)$$
$$\Pr(\mathbf{R} \mid T) = \prod_k \Pr(R_k \mid T)$$

Quiver Performance Results Comparison to Reference Genome ( <i>M. ruber</i> ; 3.1 MB; SMRT® Cells)		
	Initial Assembly	Quiver Consensus
QV	43.4	54.5
Accuracy	99.99540%	99.99964%
Differences	141	11

### PB-only Correction & Polishing

Chin et al (2013)  
Nature Methods. 10:563–569

< 5x

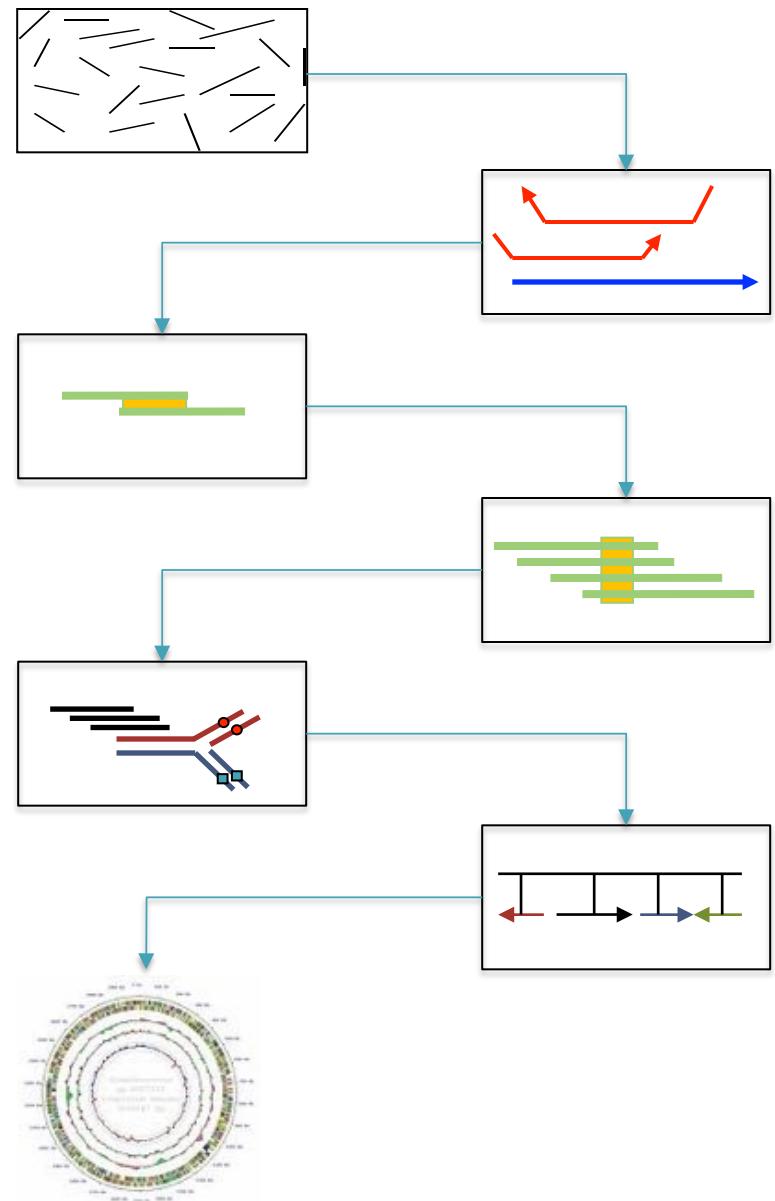
PacBio Coverage

> 50x

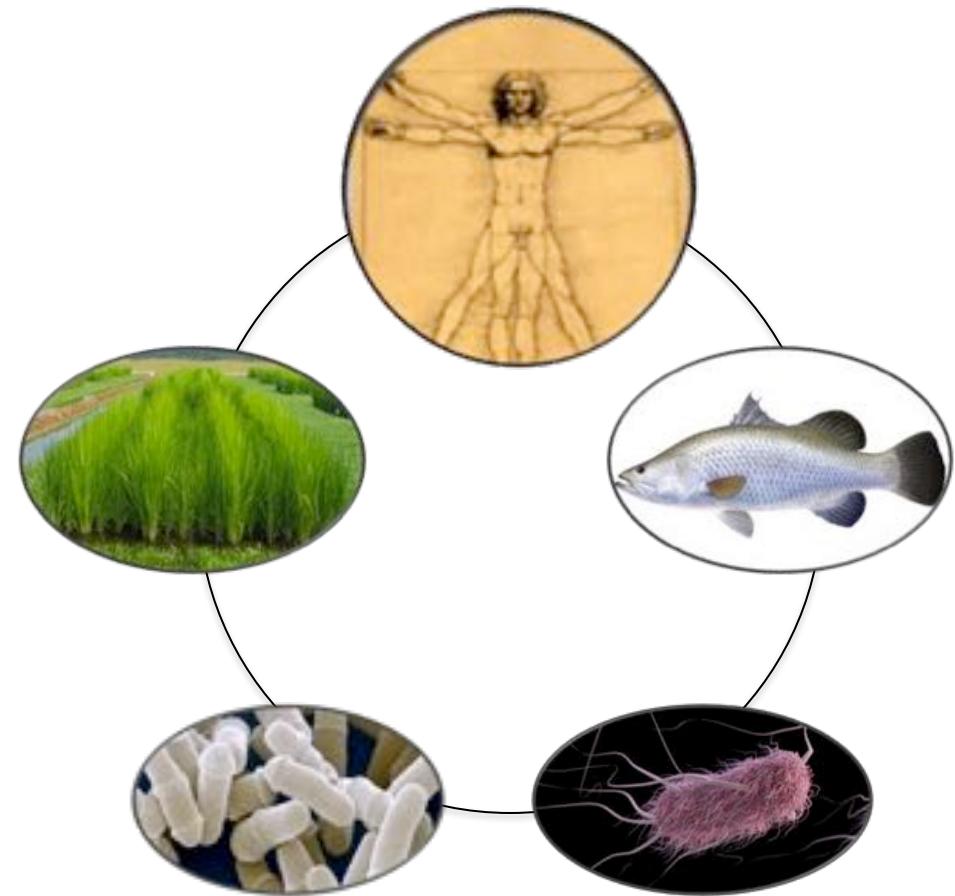
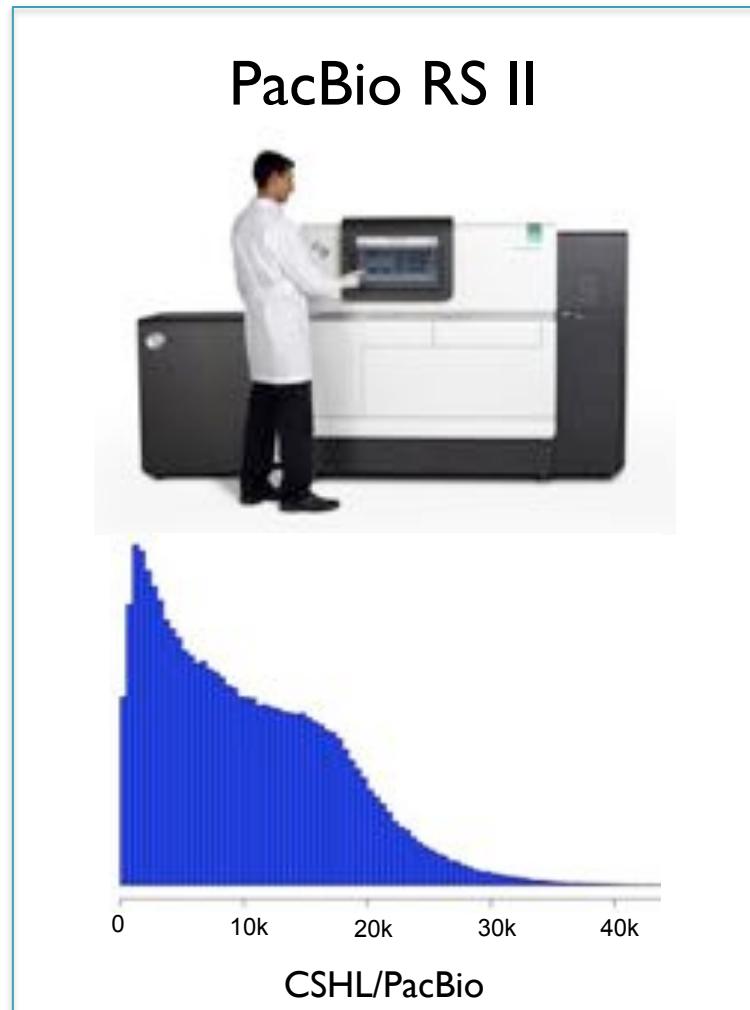
# Celera Assembler aka CANU

<https://github.com/marbl/canu>

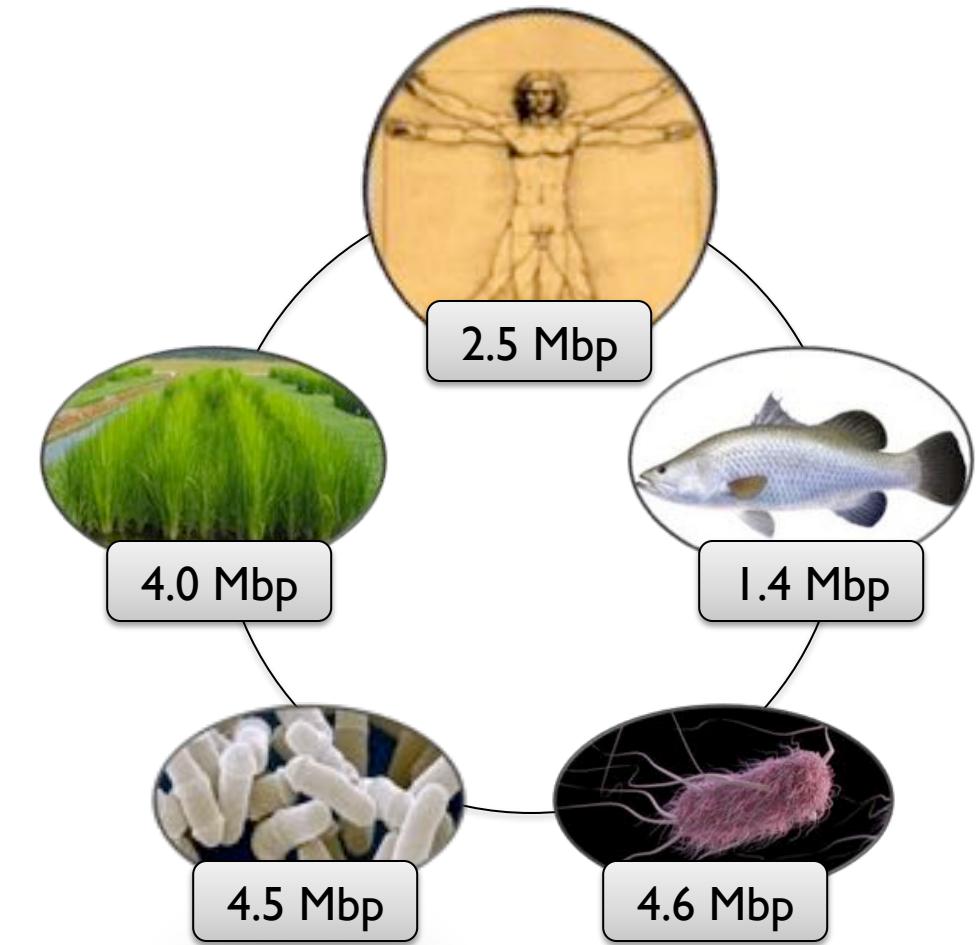
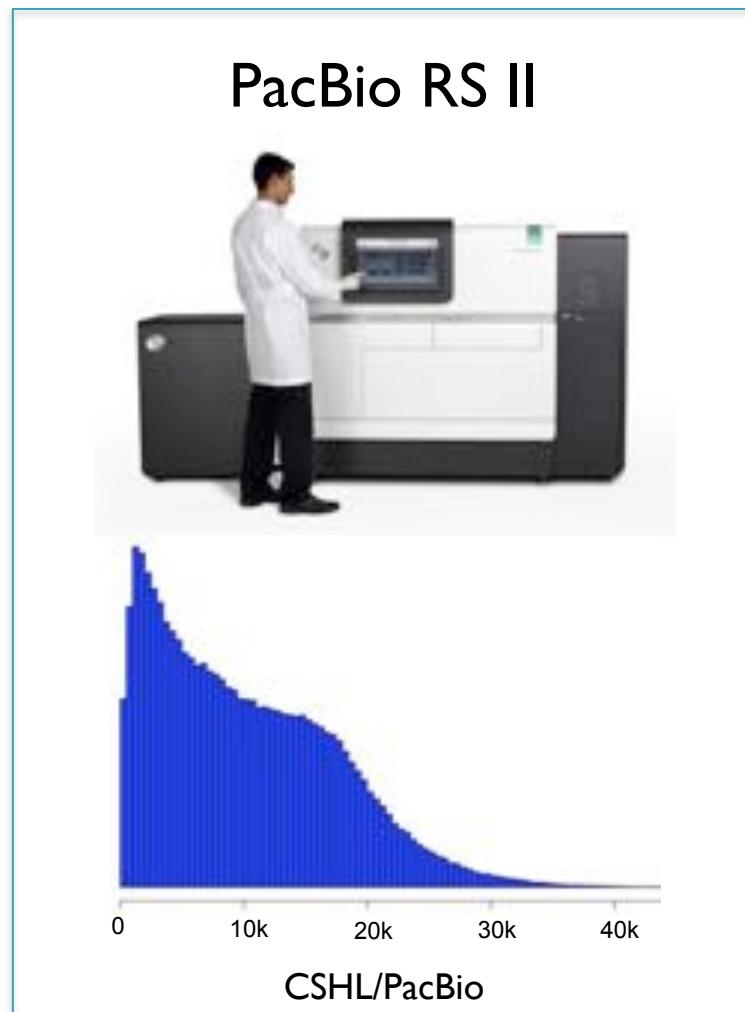
1. Pre-overlap
  - Consistency checks
2. Trimming
  - Quality trimming & partial overlaps
3. Compute Overlaps
  - Find high quality overlaps
4. Error Correction
  - Evaluate difference in context of overlapping reads
5. Unitigging
  - Merge consistent reads
6. Scaffolding
  - Bundle mates, Order & Orient
7. Finalize Data
  - Build final consensus sequences



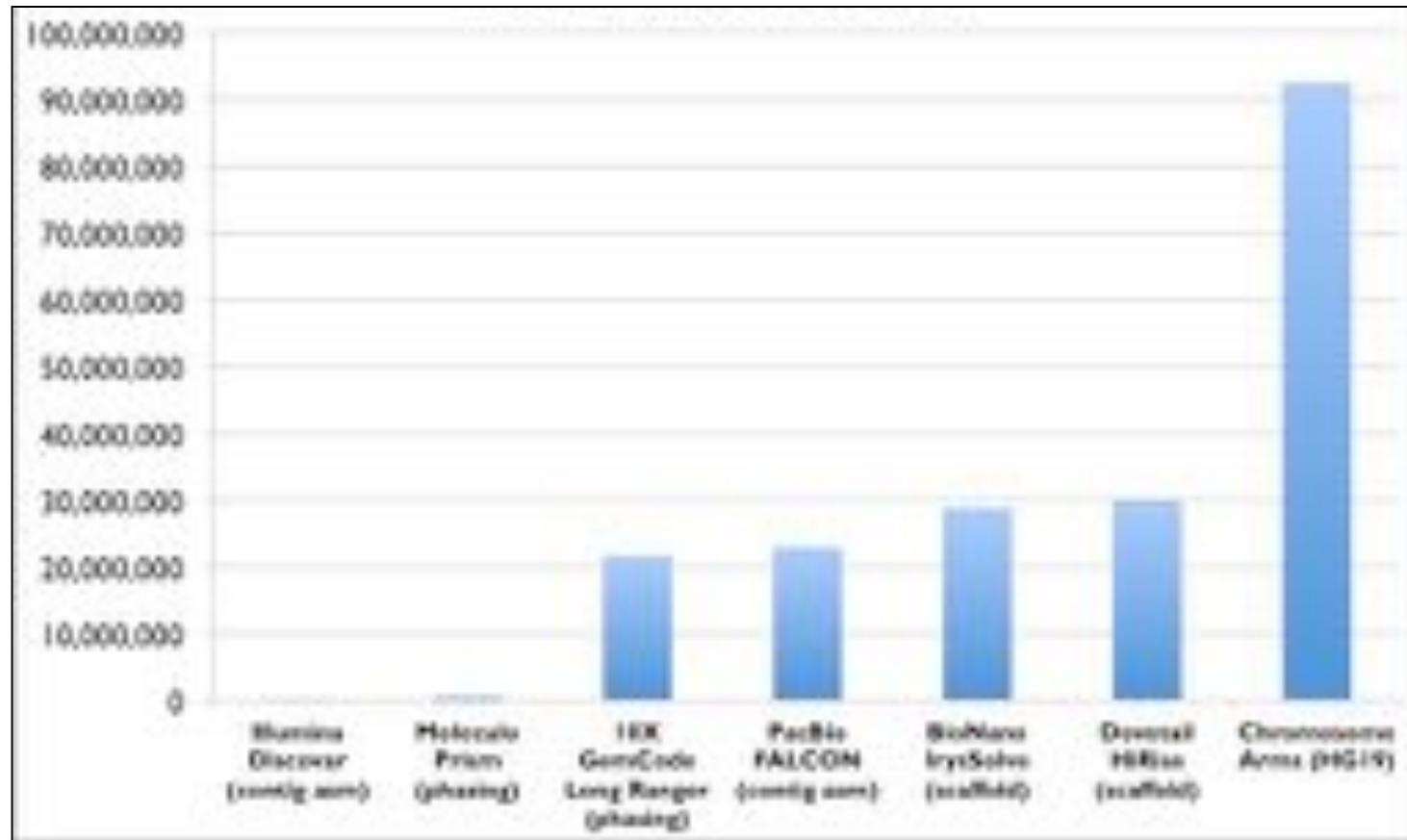
# 3<sup>rd</sup> Gen Long Read Sequencing



# 3<sup>rd</sup> Gen Long Read Sequencing



# Human Analysis N50s\*



Technology	Application	N50	Sample	Citation
Illumina Discover	contig asm	178,000	NA12877	Putnam <i>et al.</i> (2015) arXiv:1502.05331
Moleculo Prism	phasing	563,801	NA12878	Kuleshov <i>et al.</i> (2014) Nature BioTech. doi:10.1038/nbt.2833
10X GemCode Long Ranger	phasing	21,600,000	GIAB	Zook <i>et al.</i> (2015) bioRxiv. doi: <a href="http://dx.doi.org/10.1101/026468">http://dx.doi.org/10.1101/026468</a>
PacBio FALCON	contig asm	22,900,000	JCV-1	Jason Chin, PAG2016
BioNano IrysSolve	scaffold	28,800,000	NA12878	Pendleton <i>et al.</i> (2015) Nature Methods. doi:10.1038/nmeth.3454
Dovetail HiRise	scaffold	29,900,000	NA12878	Putnam <i>et al.</i> (2015) arXiv:1502.05331

\*Cross analysis of different applications

# PacBio Roadmap



## ***PacBio RS II***

\$750k instrument cost  
1895 lbs

~\$75k / human @ 50x



## ***SMRTcell***

150k Zero Mode Waveguides  
~10kb average read length  
~1 GB / SMRTcell  
~\$500 / SMRTcell

# PacBio Roadmap



## ***PacBio Sequel***

\$350k instrument cost  
841 lbs

~\$15k / human @ 50x



## ***SMRTcell v2***

1M Zero Mode Waveguides  
~15kb average read length  
~10 GB / SMRTcell  
~\$1000 / SMRTcell

# Oxford Nanopore



## MinION

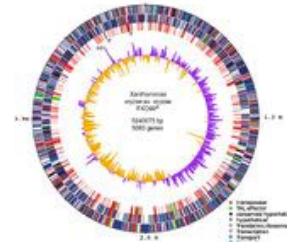
\$2k / instrument  
1 GB / day  
~\$300k / human @ 50x

## 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

# 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

# Acknowledgements

## **Schatz Lab**

Rahul Amin  
Eric Biggers  
Han Fang  
Tyler Gavin  
James Gurtowski  
Ke Jiang  
Hayan Lee  
Zak Lemmon  
Shoshana Marcus  
Giuseppe Narzisi  
Maria Nattestad  
Aspyn Palatnick  
Srividya  
Ramakrishnan  
Fritz Sedlazeck  
Rachel Sherman  
Greg Verture  
Alejandro Wences

## **CSHL**

Hannon Lab  
Gingeras Lab  
Jackson Lab  
Hicks Lab  
Iossifov Lab  
Levy Lab  
Lippman Lab  
Lyon Lab  
Martienssen Lab  
McCombie Lab  
Tuveson Lab  
Ware Lab  
Wigler Lab

## **OICR**

Karen Ng  
Timothy Beck  
Yogi Sundaravadanam  
John McPherson

## **NBACC**

Adam Phillippy  
Serge Koren



# ***Biological Data Science***

Bonnie Berger, Jeff Leek, Michael Schatz

Oct 26 - 29, 2016



# Thank you

<http://www.cs.jhu.edu/~mschatz>

@mike\_schatz