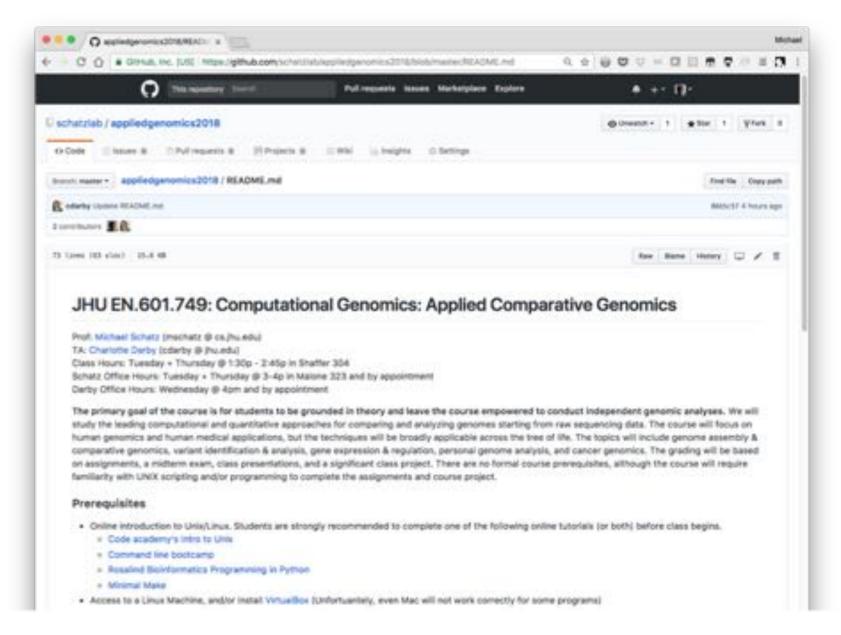# Genome Assembly

## Michael Schatz

Feb 6, 2018
Lecture 3: Applied Comparative Genomics
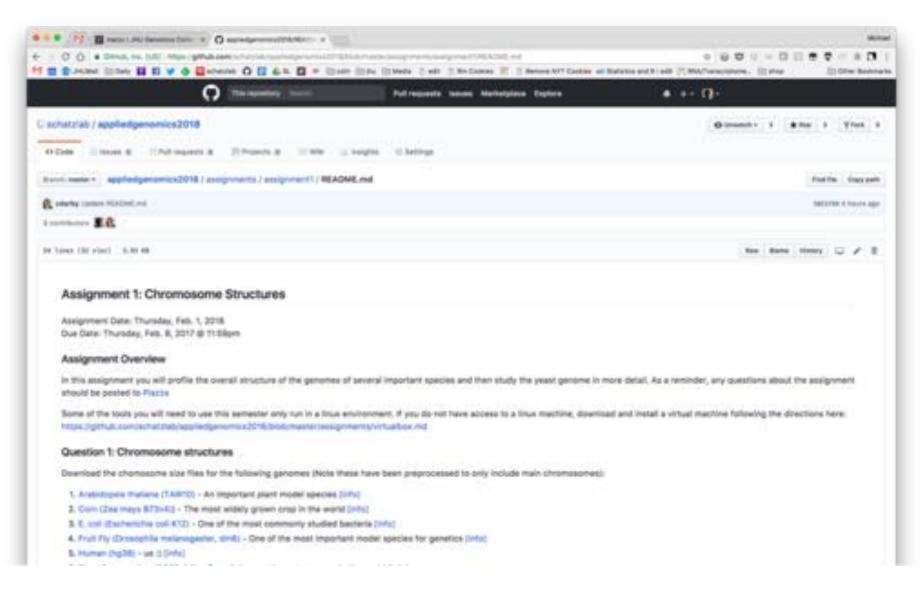
# Course Webpage



https://github.com/schatzlab/appliedgenomics2018

# Assignment 1: Chromosome Structures
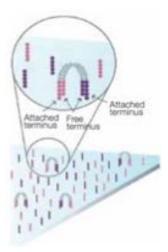## Due Feb 8 @ 11:59pm



https://github.com/schatzlab/appliedgenomics2018

# Part 1: Recap

# Second Generation Sequencing
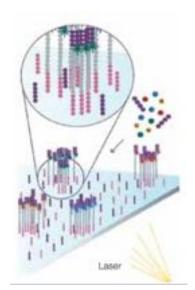


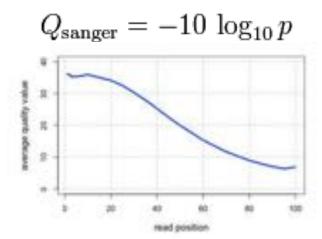1. Attach

2. Amplify

3. Image

**Illumina HiSeq 2000**
*Sequencing by Synthesis*

>60Gbp / day
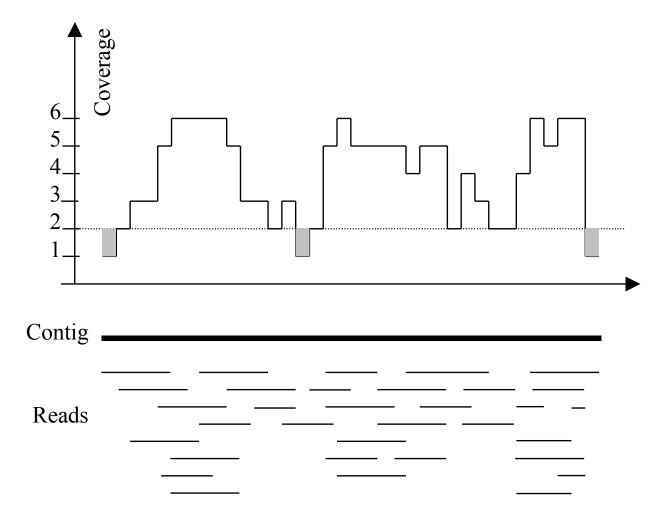
Metzker (2010) Nature Reviews Genetics 11:31-46
https://www.youtube.com/watch?v=fCd6B5HRaZ8

# Illumina Quality

| QV | $p_{error}$ |
|----|-------------|
| 40 | 1/10000 |
| 30 | 1/1000 |
| 20 | 1/100 |
| 10 | 1/10 |

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



```
SSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSS......................................
.........................XXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXX.....
.........................IIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIII....
................................JJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJ....
LLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLL......................................
!"#$%&'()*+,-./0123456789:;<=>?@ABCDEFGHIJKLMNOPQRSTUVWXYZ[\]^_`abcdefghijklmnopqrstuvwxyz{|}~
|                        |    |        |                                    |       |
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)
```

http://en.wikipedia.org/wiki/FASTQ_format

# 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?

# Illumina Sequencing Summary

## Advantages:

- Best throughput, accuracy and read length for any 2nd gen. sequencer

- Fast & robust library preparation

## Disadvantages:

- Inherent limits to read length (practically, 150bp)

- Some runs are error prone

- Requires amplification, sequences a population of molecules

**Illumina HiSeq**
~3 billion paired 100bp reads
~600Gb, $10K, 8 days
(or "rapid run" ~90Gb in 1-2 days)

**Illumina X Ten**
~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**

**Ira Hall**

# Part 2: De novo genome assembly

# Outline

1. **Assembly theory**
   - Assembly by analogy

2. **Practical Issues**
   - Coverage, read length, errors, and repeats
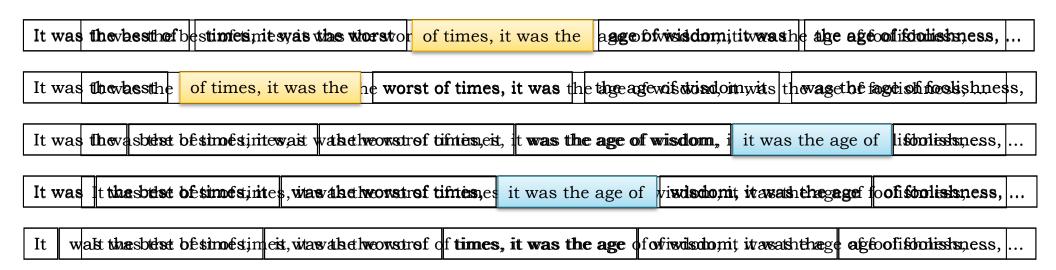
3. **Next-next-gen Assembly**
   - Canu: recommended for PacBio/ONT project

4. **Whole Genome Alignment**
   - MUMmer recommended

# Shredded Book Reconstruction

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



- 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

# Greedy Reconstruction

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

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)

| It was the best of |

| was the best of times |

| It was the best |     | was the best of |

| was the best of |     | the best of times |

| It was the best |

| was the best of |

| the best of times |

– Overlaps between fragments are implicitly computed

*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
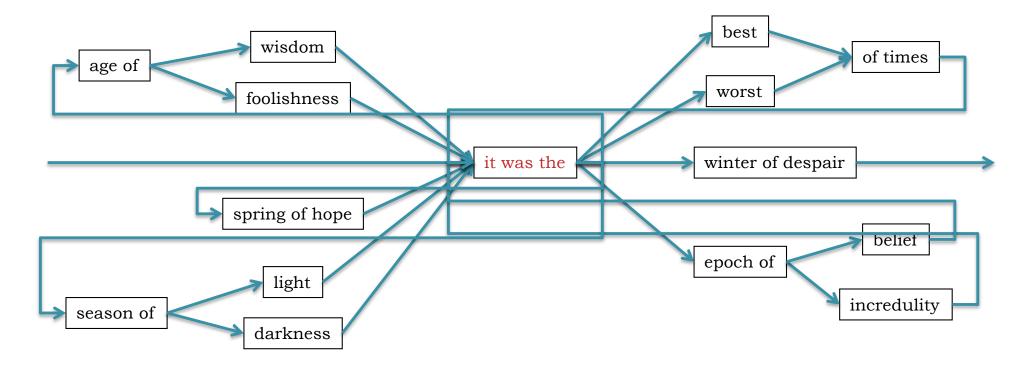
the age of wisdom,

age of wisdom, it

of wisdom, it was

wisdom, it was the

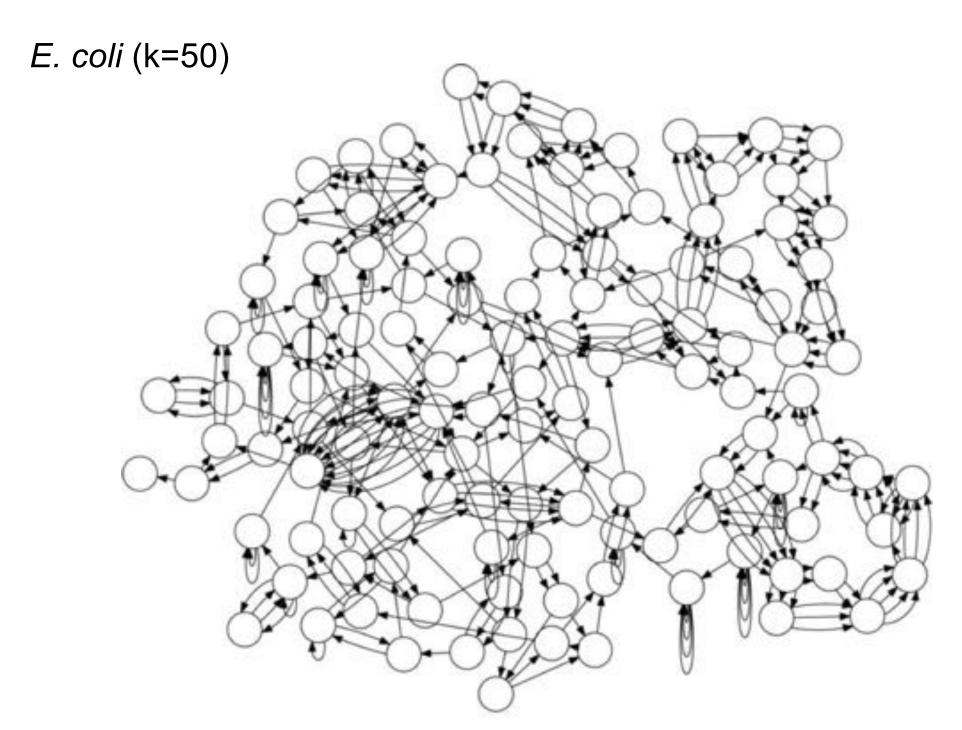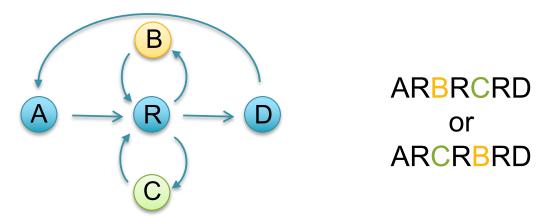After graph construction, try to simplify the graph as much as possible

# de Bruijn Graph Assembly

It was the best of times, it

of times, it was the

it was the worst of times, it

it was the age of

the age of foolishness

the age of wisdom, it was the

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

# 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 winder of despair …

# E. coli (k=50)



**Reducing assembly complexity of microbial genomes with single-molecule sequencing**

# 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 = $n$ x $n$ 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

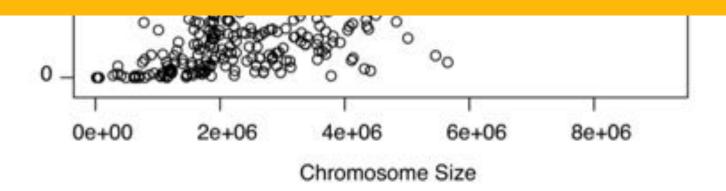$a_{uv}$ = multiplicity of edge from $u$ to $v$

**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.**
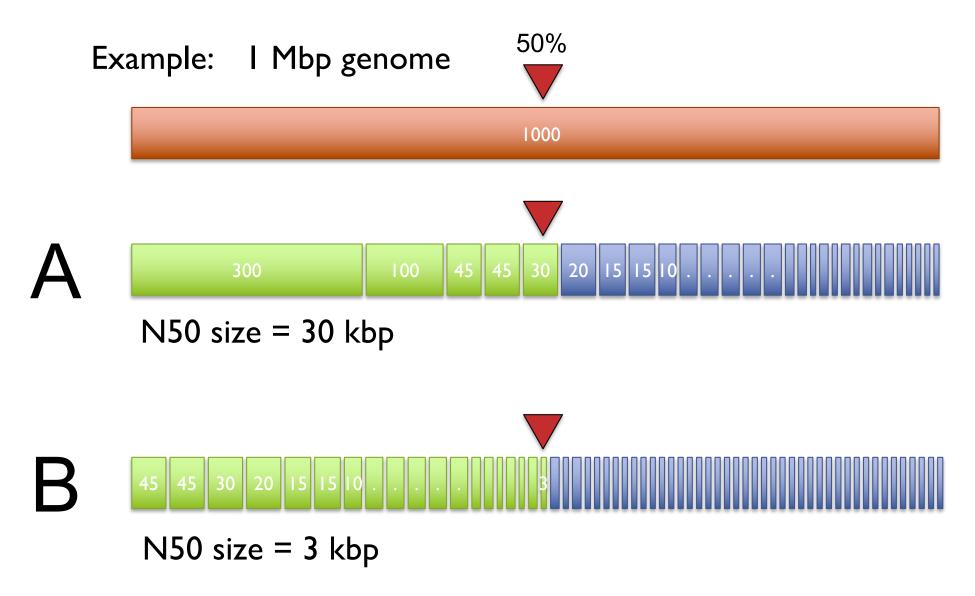Kingsford C, Schatz MC, Pop M (2010) *BMC Bioinformatics*.

**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~~ genomical 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

# Pop Quiz 1

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

ATTA

GATT

TACA

TTAC

# Pop Quiz I

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

```
ATTA: ATT -> TTA
GATT: GAT -> ATT
TACA: TAC -> ACA
TTAC: TTA -> TAC
```

# Pop Quiz I

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

ATTA: ATT -> TTA
GATT: GAT -> ATT
TACA: TAC -> ACA
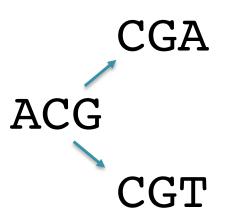TTAC: TTA -> TAC

GAT
ATT
TTA
TAC
ACA

GATTACA

# Pop Quiz 2

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

ACGA

ACGT

ATAC

CGAC

CGTA

GACG

GTAT

TACG

# Pop Quiz 2

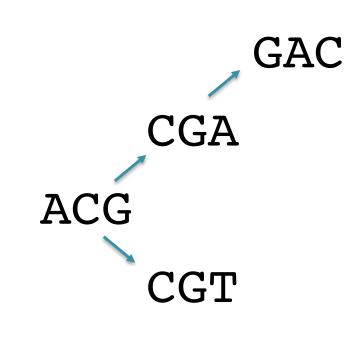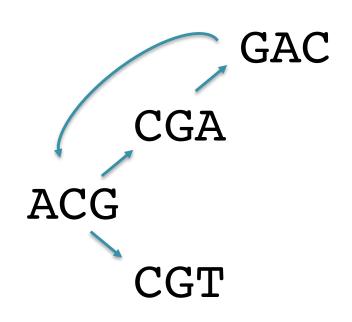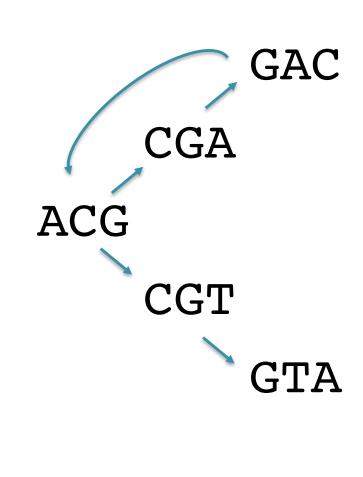Assemble these reads using a de Bruijn graph approach (k=3):

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

ACG → CGA

# Pop Quiz 2

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

~~ACGA~~

~~ACGT~~

ATAC

CGAC

CGTA

GACG

GTAT

TACG

```
         CGA
        ↗
  ACG
        ↘
         CGT
```

# Pop Quiz 2

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

~~ACGA~~

~~ACGT~~

ATAC

~~CGAC~~

CGTA

GACG

GTAT

TACG

GAC

CGA

ACG

CGT

# Pop Quiz 2

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

~~ACGA~~

~~ACGT~~

ATAC

~~CGAC~~

CGTA

~~GACG~~

GTAT

TACG

# Pop Quiz 2

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

~~ACGA~~

~~ACGT~~

ATAC

~~CGAC~~

~~CGTA~~

~~GACG~~

GTAT

TACG

# Pop Quiz 2

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

~~ACGA~~

~~ACGT~~

ATAC

~~CGAC~~

~~CGTA~~

~~GACG~~

~~GTAT~~

TACG

# Pop Quiz 2

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

~~ACGA~~

~~ACGT~~

ATAC

~~CGAC~~

~~CGTA~~

~~GACG~~

~~GTAT~~

~~TACG~~

# Pop Quiz 2

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

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

ATA → TAC → ACG → CGA → GAC

ACG → CGT → GTA → TAT

# Pop Quiz 2

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

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



ATACGACGTAT

# Pop Quiz 2

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

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

ATA → TAC → ACG → GAC
ACG → CGA
ACG → CGT → GTA → TAT

**Whats another possible genome?**

ATACGACGTAT

# Outline

# Assembly Applications
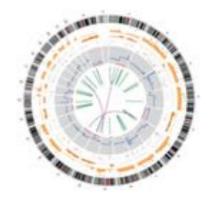
- Novel genomes
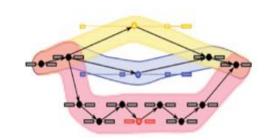
- Metagenomes

- Sequencing assays
  - Structural variations
  - Transcript assembly
  - …

# Why are genomes hard to assemble?

1.  **Biological**:
    - (Very) High ploidy, heterozygosity, repeat content

2.  **Sequencing**:
    - (Very) large genomes, imperfect sequencing

3.  **Computational**:
    - (Very) Large genomes, complex structure

4.  **Accuracy**:
    - (Very) Hard to assess correctness
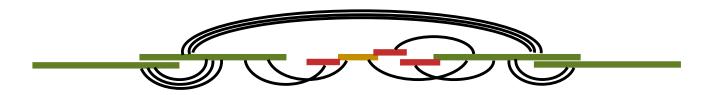
# Assembling a Genome

1. Shear & Sequence DNA

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

...AGCCTAGGGATGCGCGACACGT
        GGATGCGCGACACGTCGCATATCCGGTTTGGTCAACCTCGGACGGAC
                     CAACCTCGGACGGACCTCAGCGAA...
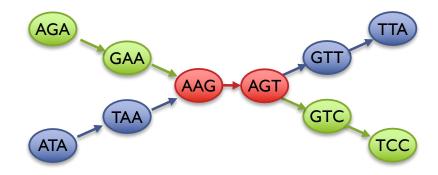
3. Simplify assembly graph

4. Detangle graph with long reads, mates, and other links
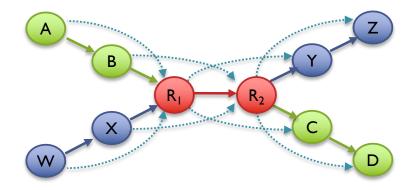
# Two Paradigms for Assembly

## de Bruijn Graph



Short read assemblers

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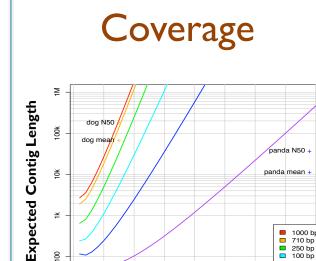- Robust to high coverage

## 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.
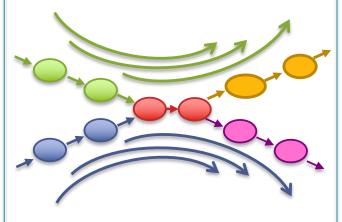
# Ingredients for a good assembly

## 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

# Kmer-based Coverage Analysis



Even though the reads are not assembled or aligned (or reference available),
Kmer counting is an effective technique to estimate coverage & other genome properties

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

# Heterozygous Kmer Profiles



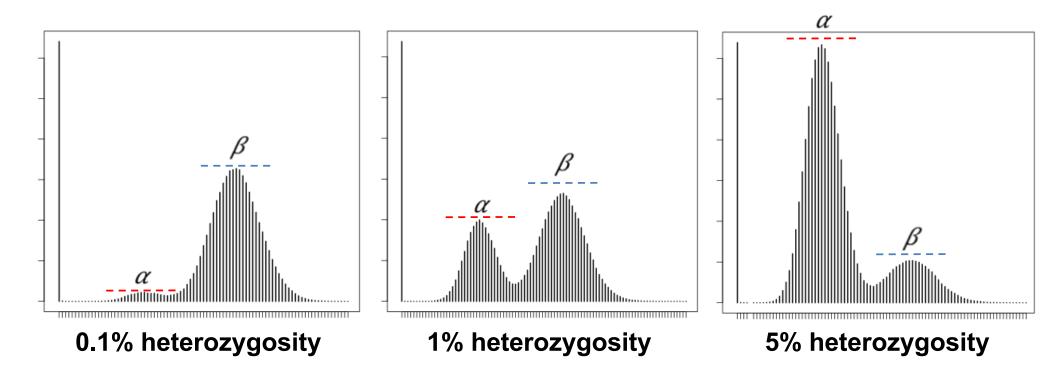**0.1% heterozygosity**     **1% heterozygosity**     **5% heterozygosity**

- *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 Model

$$f(x) = G\left\{\alpha NB(x, \lambda, \lambda/\rho) + \beta NB(x, 2\lambda, 2\lambda/\rho) + \gamma NB(x, 3\lambda, 3\lambda/\rho) + \delta NB(x, 4\lambda, 4\lambda/\rho)\right\}$$

Analyze k-mer profiles using a mixture model of 4 negative binominal components
- Components centered at 1,2,3,4 * λ

- Four components capture heterozygous and homozygous unique (α,β) and 2 copy repeats (γ,δ). Higher order repeats do not contribute a significant number of kmers

- Negative binomial instead of Poisson to account for over dispersion observed in real data (especially PCR duplicates); variance modeled by ρ

$$\alpha = 2(1-d)(1-(1-r)^k) + 2d(1-(1-r)^k)^2 + 2d((1-r)^k)(1-(1-r)^k)$$

$$\beta = (1-d)((1-r)^k) + d(1-(1-r)^k)^2$$

$$\gamma = 2d((1-r)^k)(1-(1-r)^k)$$

$$\delta = d(1-r)^{2k}$$

$k$ is the *k-mer* length

$r$ is the rate of heterozygosity

$d$ represents the percentage of the genome that is two-copy repeat

***Fit model with nls, infer rate of heterozygosity, genome size, unique/repetitive content, sequencing error rate, rate of PCR duplicates***

# 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)

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

# 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 | Pop Bubbles |
|---|---|
| was the worst of times,<br><br>was the worst of t**y**mes,<br><br>the worst of times, it | was the worst of times,<br><br>was the worst of t**y**mes,<br><br>times, it was the age<br><br>t**y**mes, it was the age |

**Clip Tips (bottom):**

the worst of t**y**mes,

→ was the worst of

the worst of times,

worst of times, it

**Pop Bubbles (bottom):**

t**y**mes,

was the worst of → → it was the age

times,

# Repetitive regions

| Repeat Type | Definition / Example | Prevalence |
| --- | --- | --- |
| Low-complexity DNA / Microsatellites | $(b_1 b_2 \ldots b_k)^N$ where $1 \le k \le 6$ <br> CACACACACACACACACA | 2% |
| SINEs (Short Interspersed Nuclear Elements) | *Alu* sequence (~280 bp) <br> 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

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

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

**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
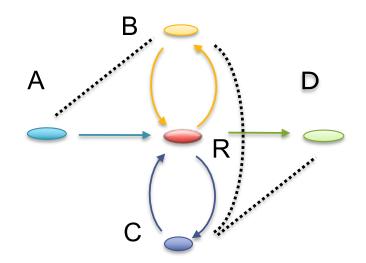
300bp

**Mate-pair sequencing**

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

10kbp

10kbp circle

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
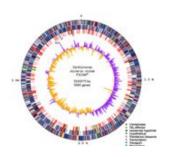
**Why do scaffolds end?**

# Assemblathon Results

| ID | Overall | CPNG50 | SPNG50 | Struct. | CC50 | Subs. | Copy. Num. | Cov. Tot. | Cov. CDS |
|---|---|---|---|---|---|---|---|---|---|
| BGI | 36 | ★ (yellow) | | | | | ☆ (gray) | ★ (yellow) | ☆ (gray) |
| Broad | 37 | ☆ (gray) | ★ (orange) | ★ (orange) | ★ (yellow) | | | | |
| WTSI-S | 46 | | ★ (yellow) | ☆ (gray) | ★ (orange) | ★ (yellow) | | | |
| CSHL | 52 | ★ (orange) | | | | | | | ☆ (gray) |
| BCCGSC | 53 | | | | | | | ☆ (gray) | ★ (yellow) |
| DOEJGI | 56 | | ☆ (gray) | ★ (yellow) | ☆ (gray) | ★ (orange) | | | |
| 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 1: A competitive assessment of de novo short read assembly methods*
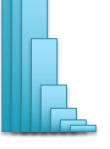Earl et al. (2011) Genome Research. 21: 2224-2241

# 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 1

    1.  Set up Linux, set up Virtual Machine
    2.  Set up Dropbox for yourself!
    3.  Get comfortable on the command line

**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    ✉️

*Welcome to Applied Comparative Genomics*
https://github.com/schatzlab/appliedgenomics2018

# Questions?