

C G T A C G T A
A C G T A C G T

Sequence Assembly Intro

Adam M. Phillippy

April 30, 2019

@aphillippy 



National Human Genome
Research Institute

—
The **Forefront**
of **Genomics**[®]
—

Slides courtesy of: Michael Schatz

Feb 4, 2019

Lecture 3: Applied Comparative Genomics

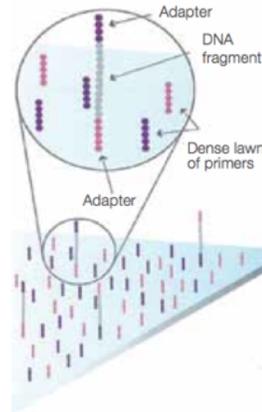


Second Generation Sequencing

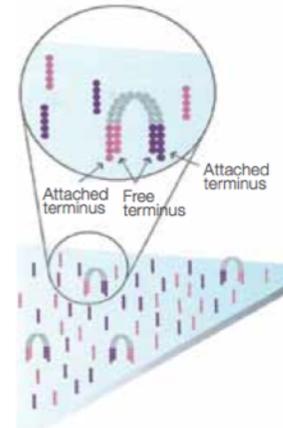


Illumina HiSeq 2000
Sequencing by Synthesis

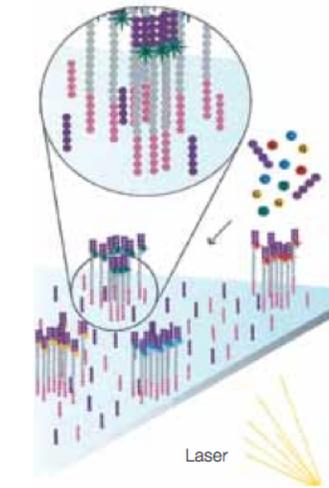
>60Gbp / day



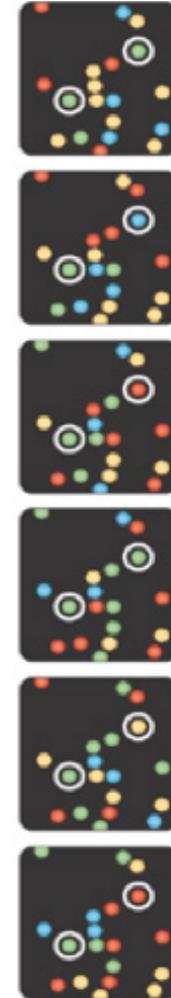
1. Attach



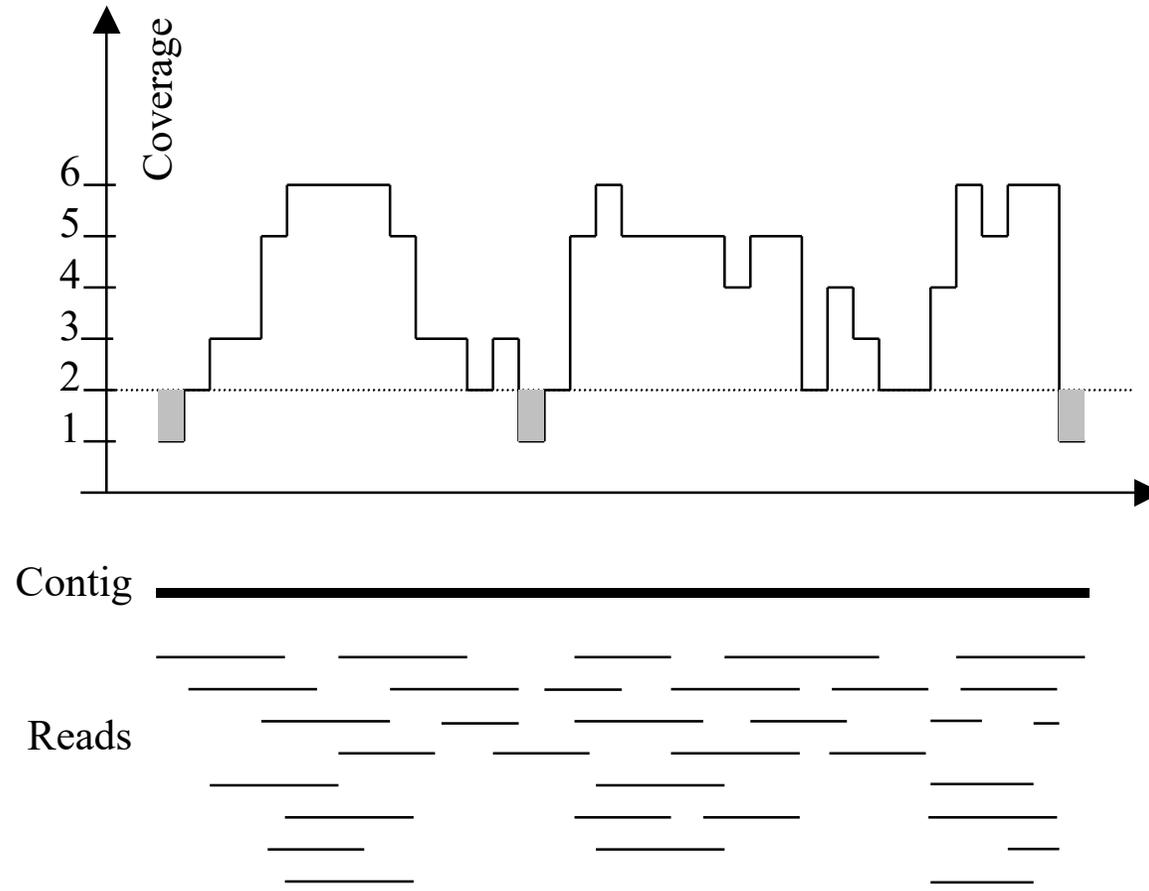
2. Amplify



3. Image



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?

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.

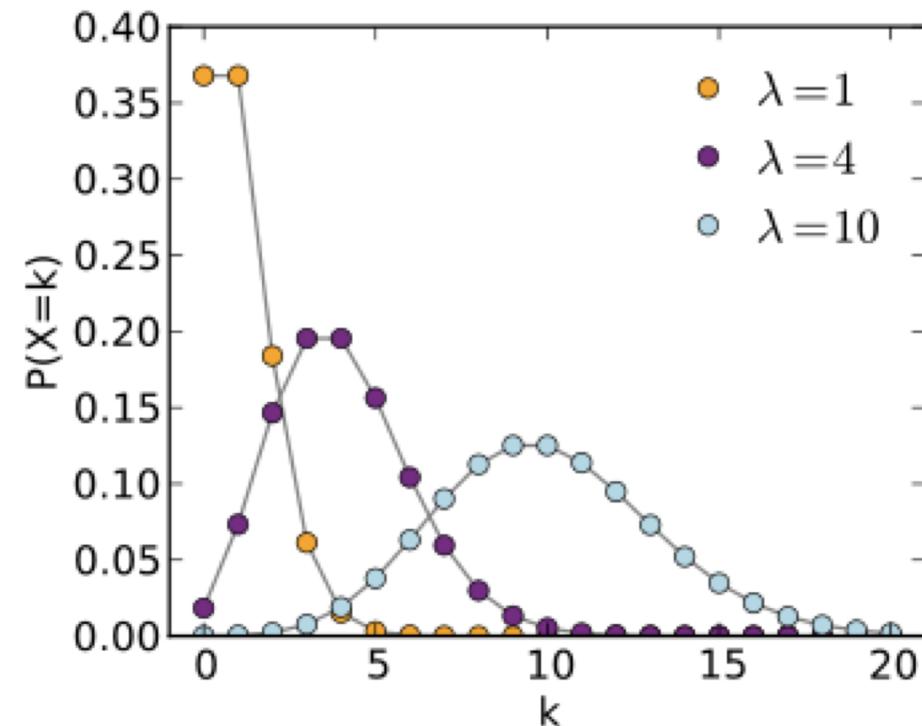
Formulation comes from the limit of the binomial equation

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

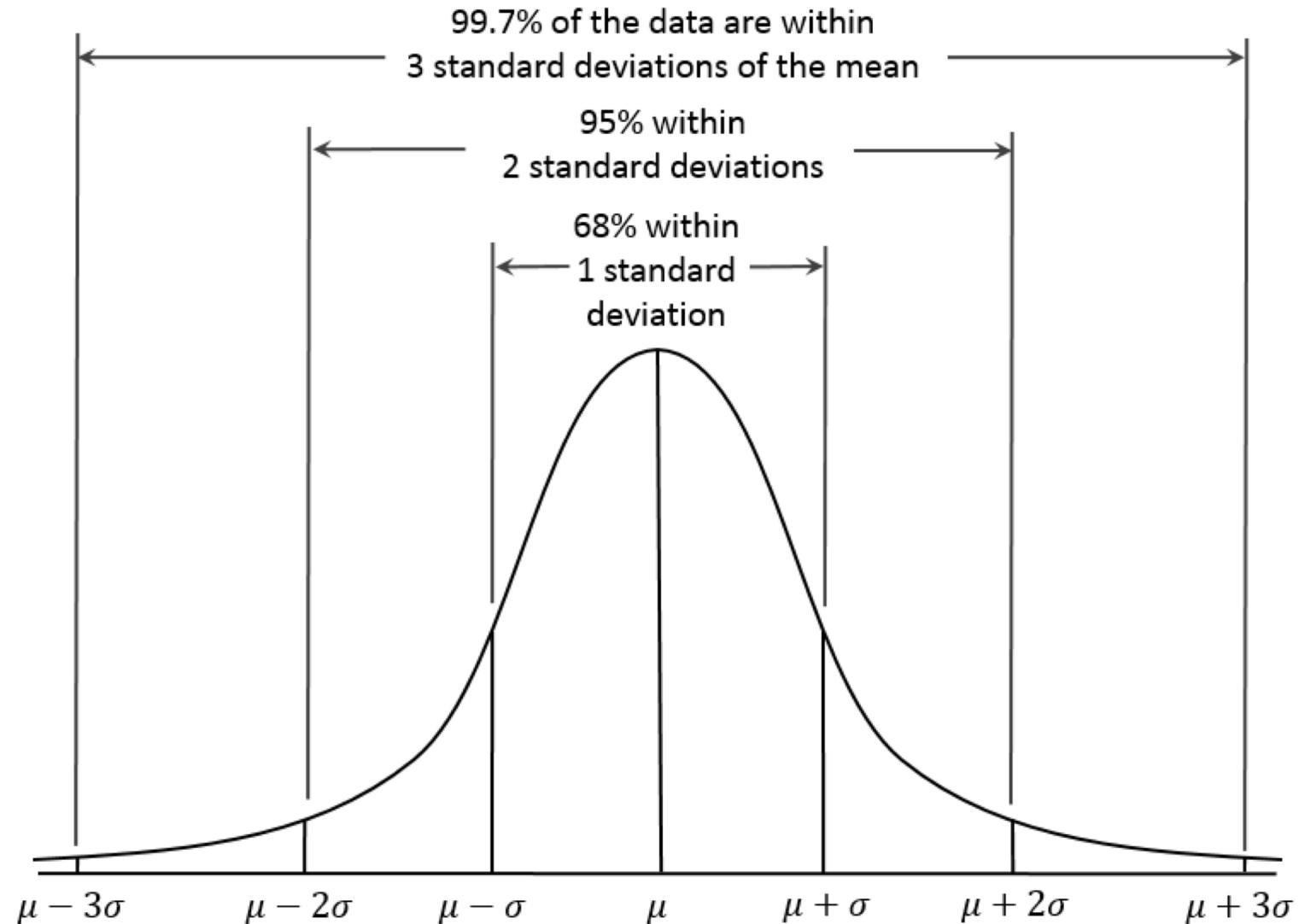
Key properties:

- ***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 need $10\text{Mbp} \times 24x = 240\text{Mbp}$ of data
 $240\text{Mbp} / 120\text{bp} / \text{read} = 2\text{M}$ reads

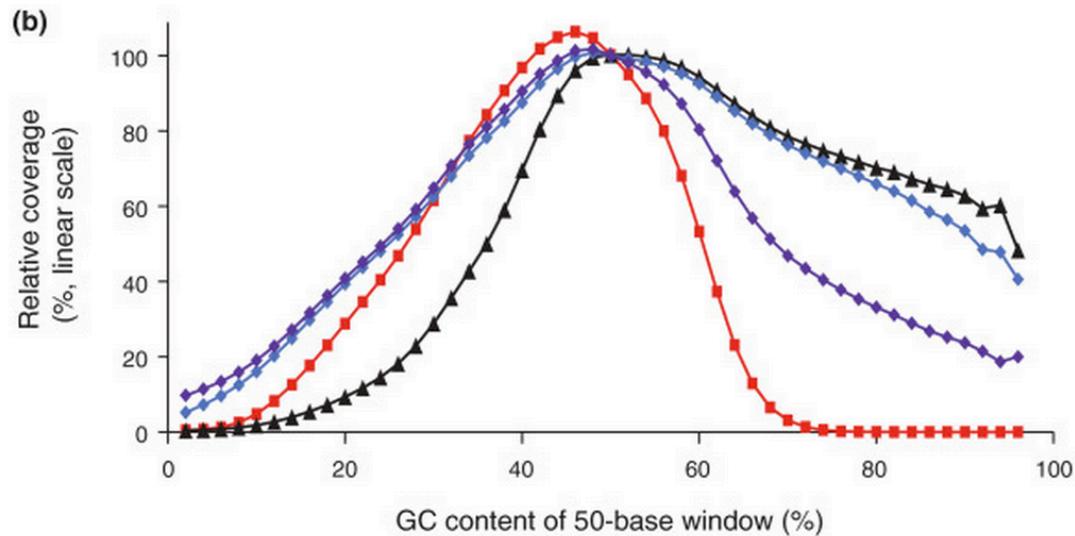
I want to sequence a 10Mbp genome so that
>97.5% of the genome has at least 24x coverage.
How many 120bp reads do I need?

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

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

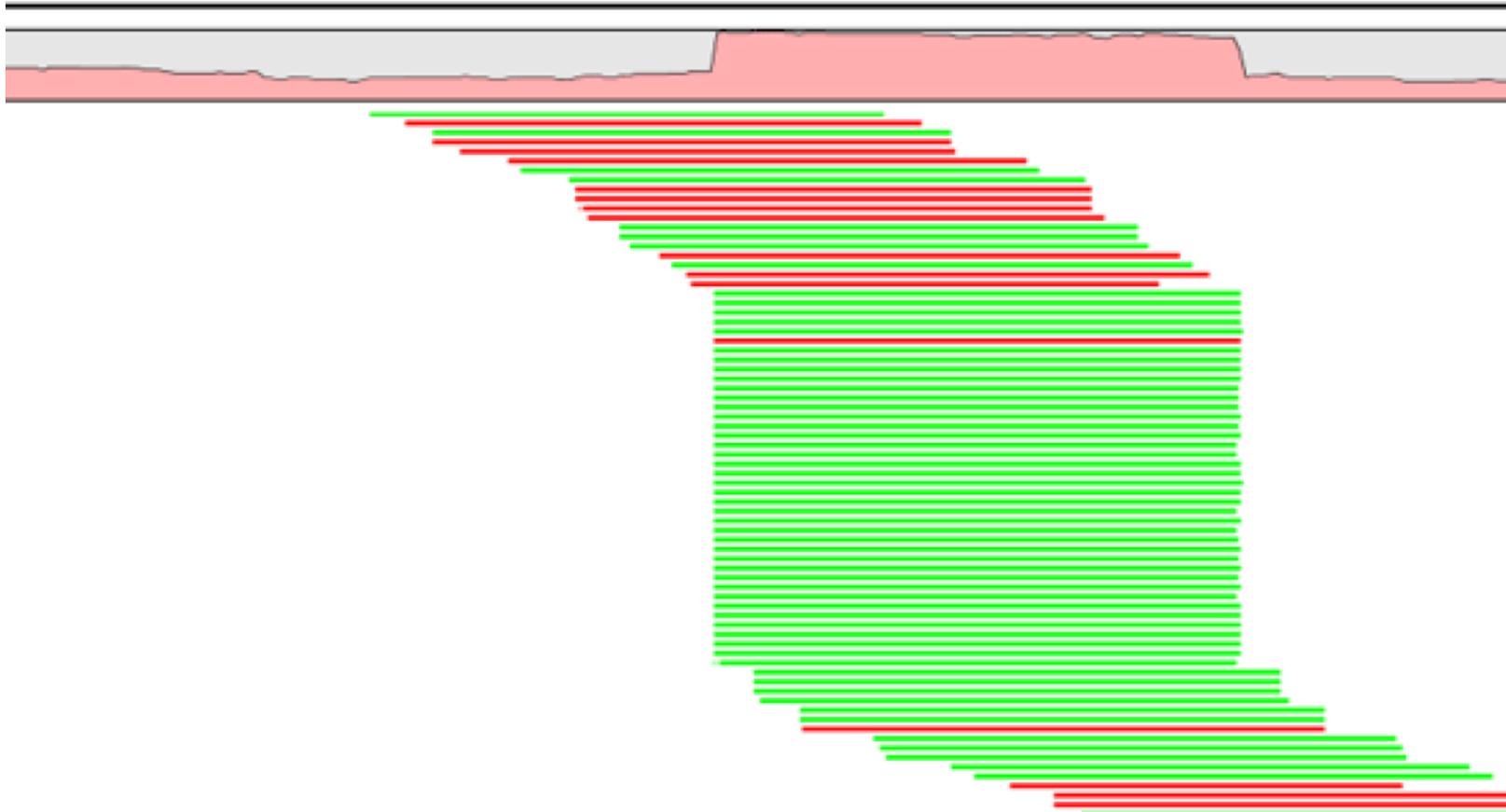
I need $10\text{Mbp} \times 36x = 360\text{Mbp}$ of data
 $360\text{Mbp} / 120\text{bp} / \text{read} = 3\text{M}$ reads

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

Beware of Duplicate Reads

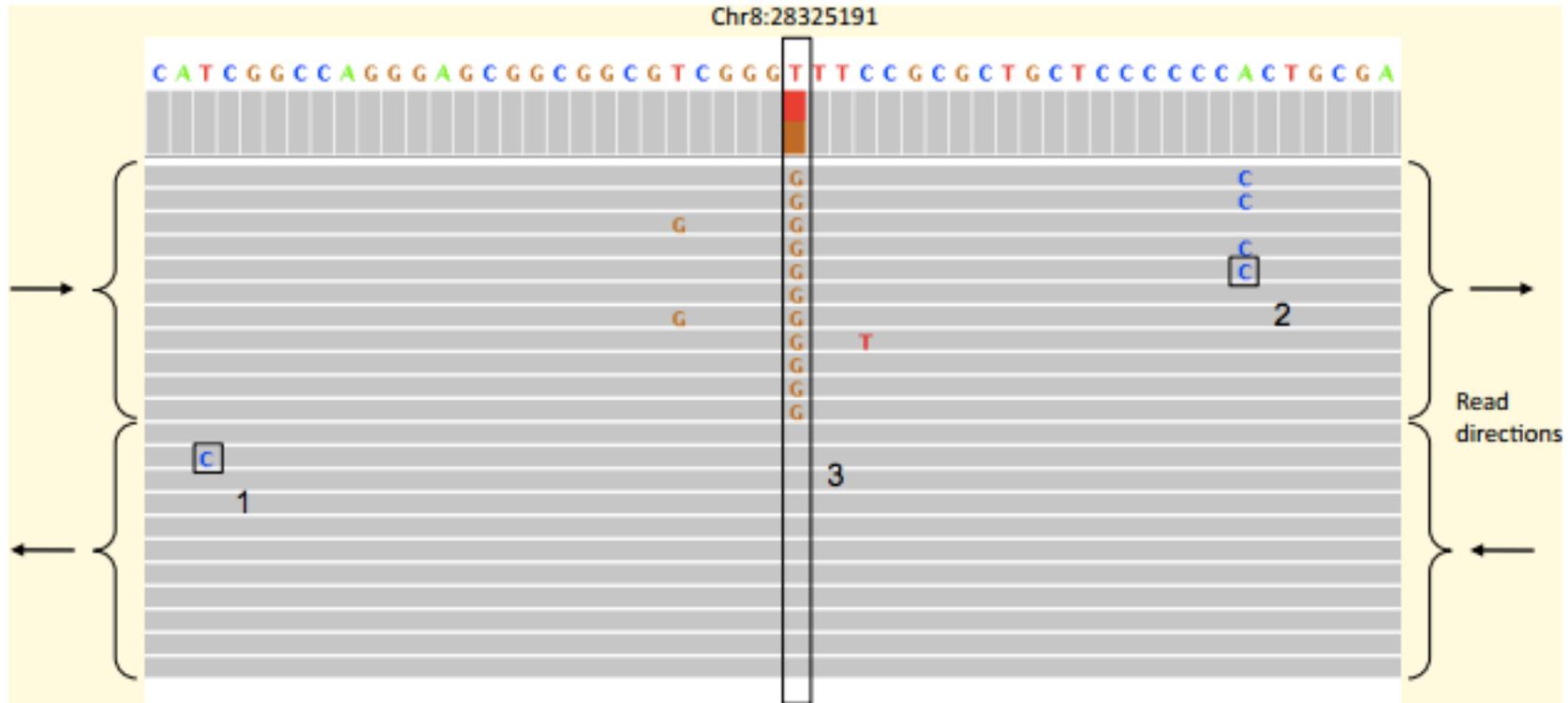


The Sequence alignment/map (SAM) format and SAMtools.

Li et al. (2009) *Bioinformatics*. 25:2078-9

Picard: <http://picard.sourceforge.net>

Beware of (Systematic) Errors



Identification and correction of systematic error in high-throughput sequence data

Meacham et al. (2011) *BMC Bioinformatics*. 12:451

A closer look at RNA editing.

Lior Pachter (2012) *Nature Biotechnology*. 30:246-247

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

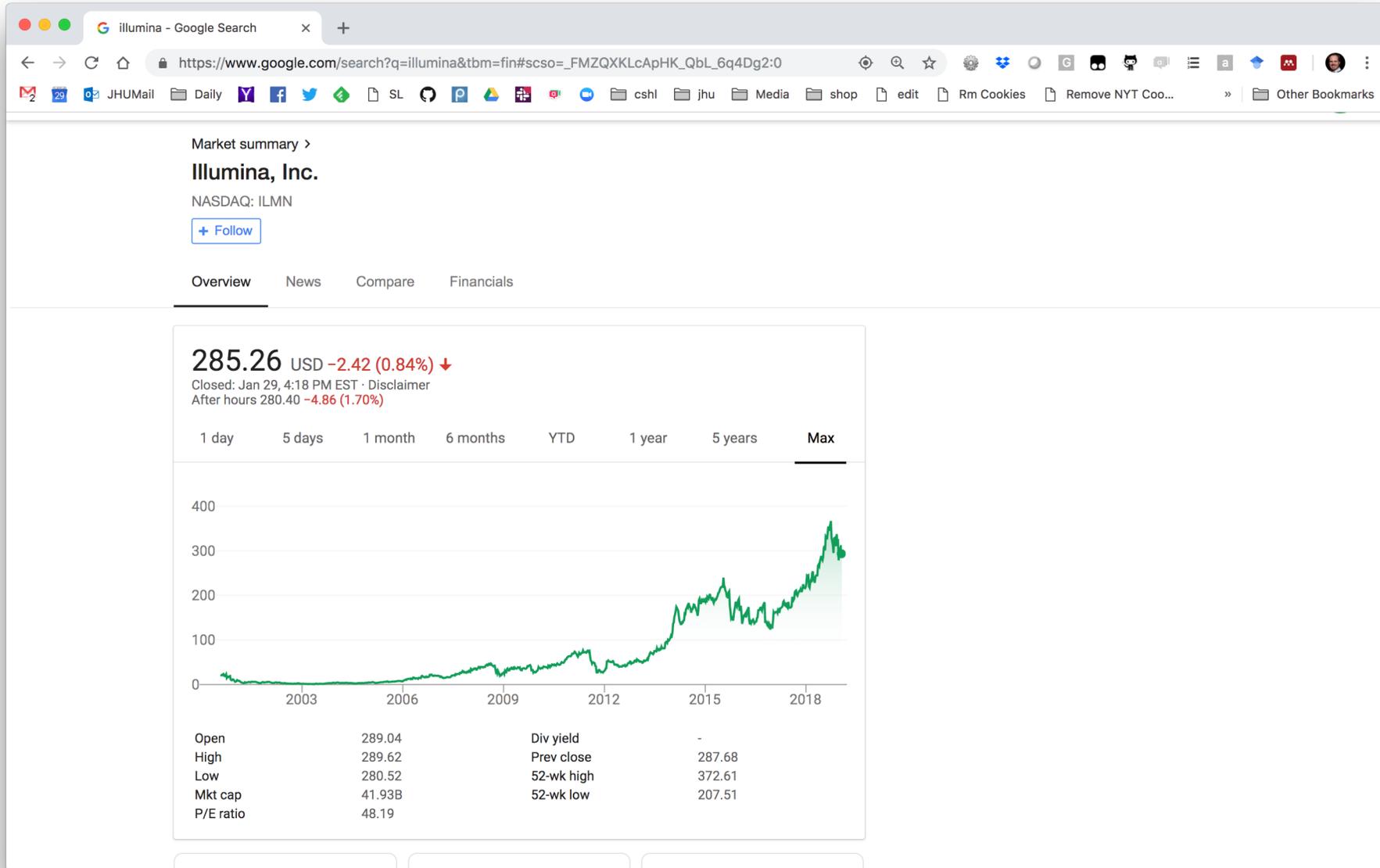
~3 billion paired 100bp reads
~600Gb, \$10K, 8 days
(or “rapid run” ~90Gb in 1-2 days)

Illumina X Ten

~6 billion paired 150bp reads
1.8Tb, <3 days, ~1000 / genome(\$\$)
(or “rapid run” ~90Gb in 1-2 days)

Illumina NovaSeq

Population-scale sequencing





De novo genome assembly

Outline

1. *Assembly theory*

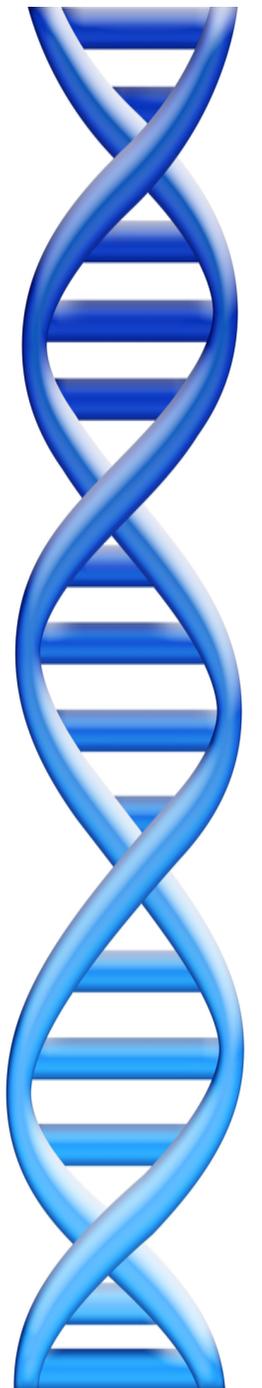
- Assembly by analogy

2. *Practical issues*

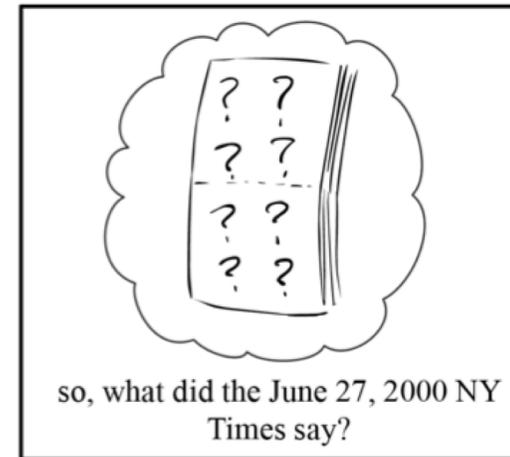
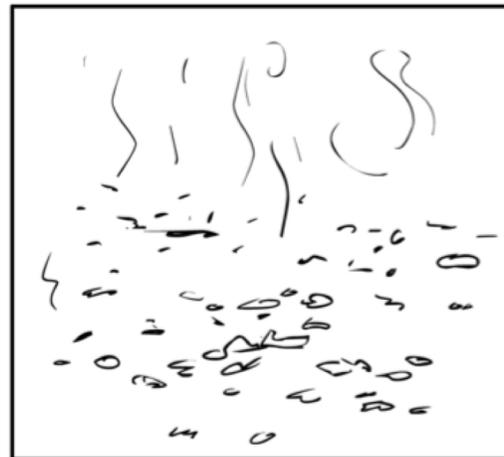
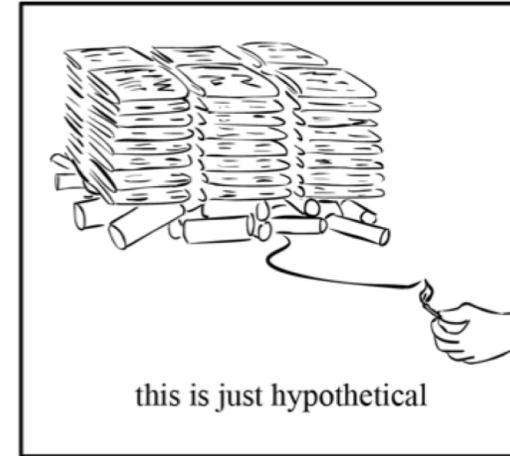
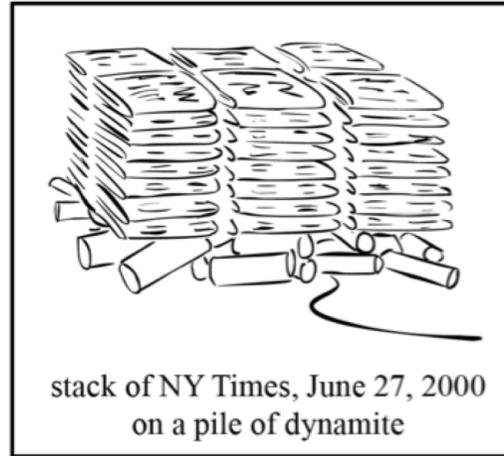
- Coverage, read length, errors, and repeats

3. *Recent advances in assembly*

- PacBio, Nanopore, and Canu
- Dr. Sergey Koren (Thursday, May 2)



The exploding newspapers problem



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

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]

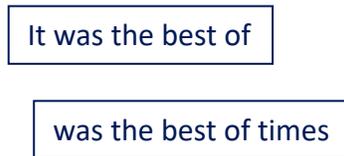
Model the assembly problem as a graph problem

How long will it take to compute the overlaps?

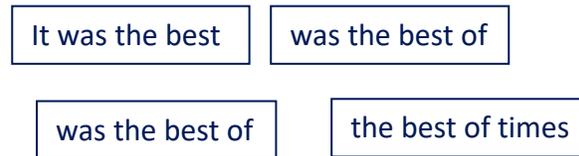
de Bruijn Graph Construction

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

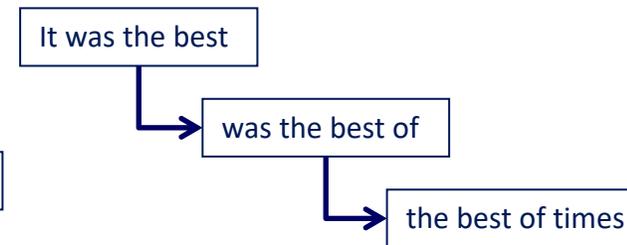
Fragments $|f|=5$



Sub-fragment $k=4$



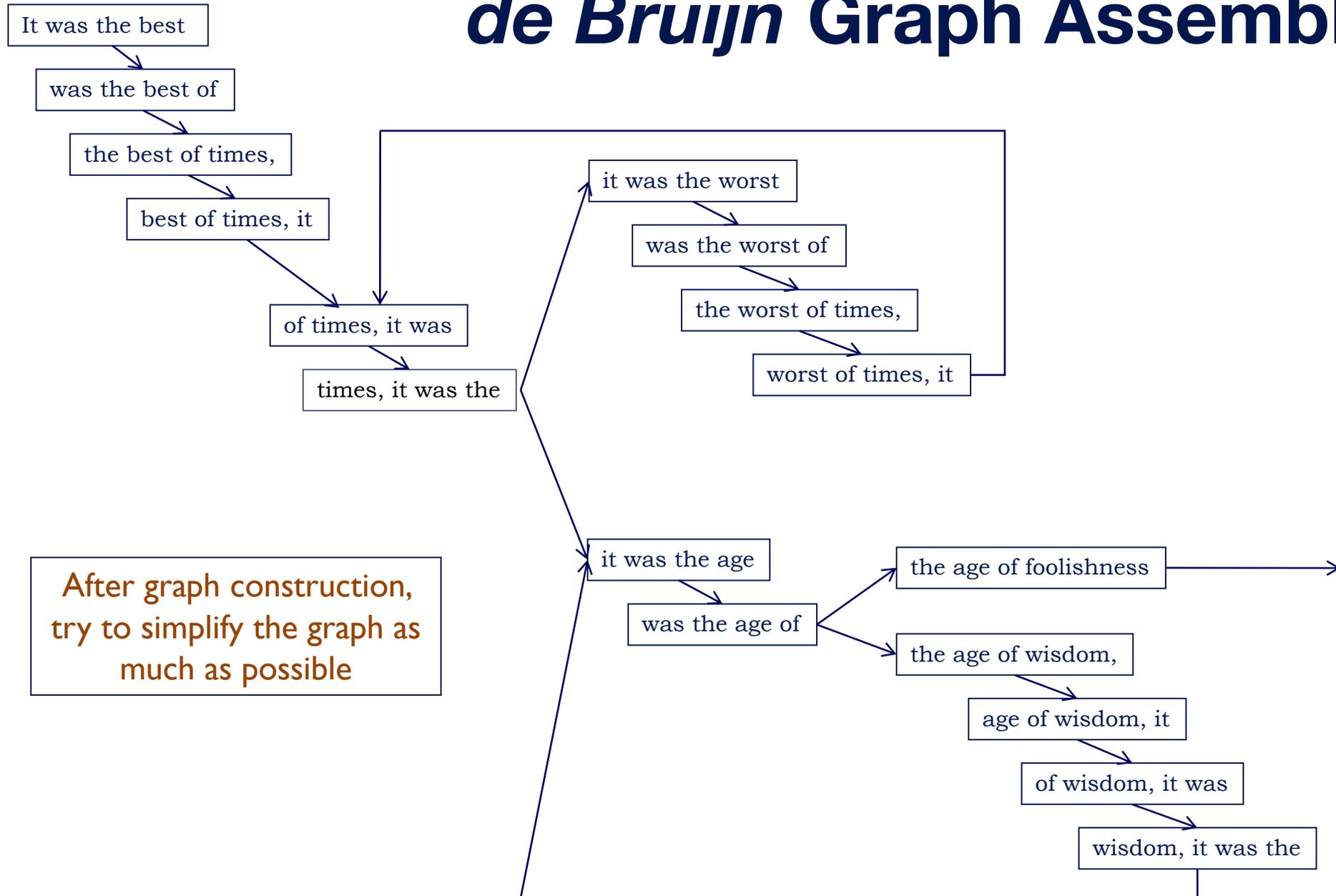
Directed edges (overlap by $k-1$)



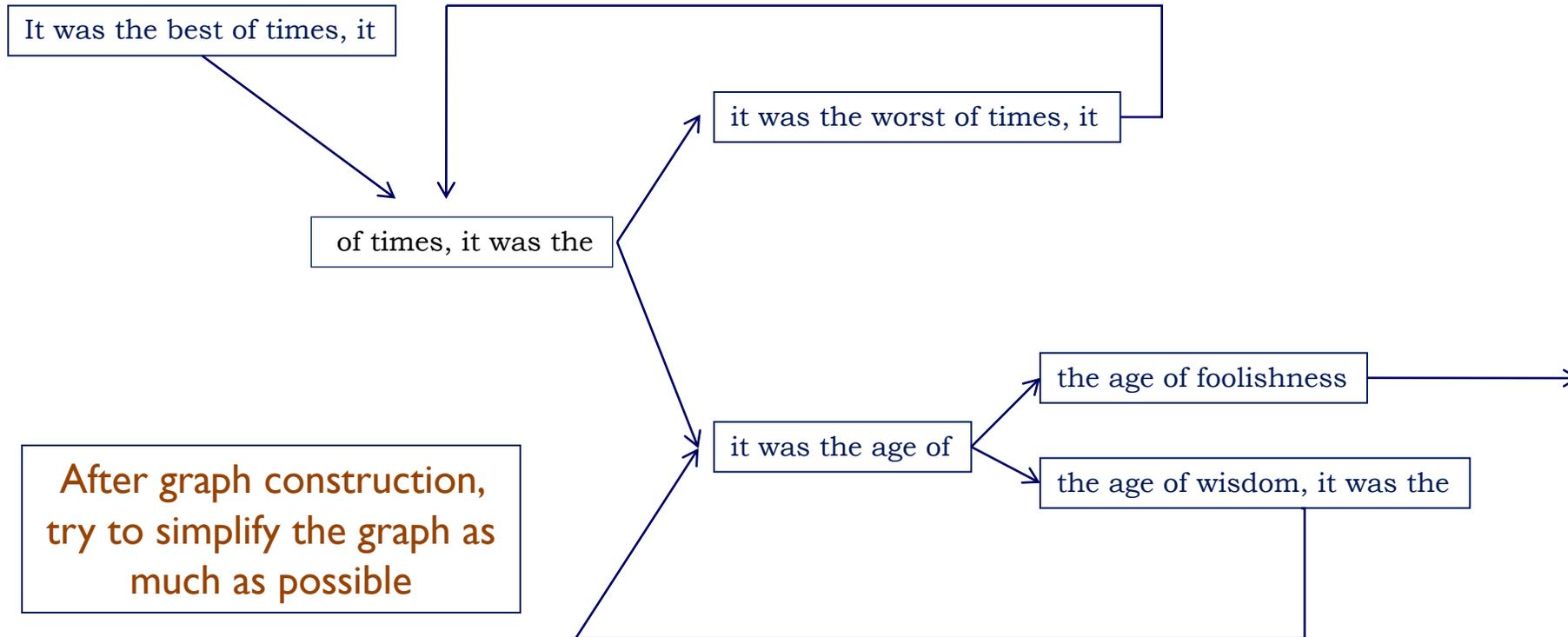
- Overlaps between fragments are implicitly computed

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

de Bruijn Graph Assembly



de Bruijn Graph Assembly



The full tale

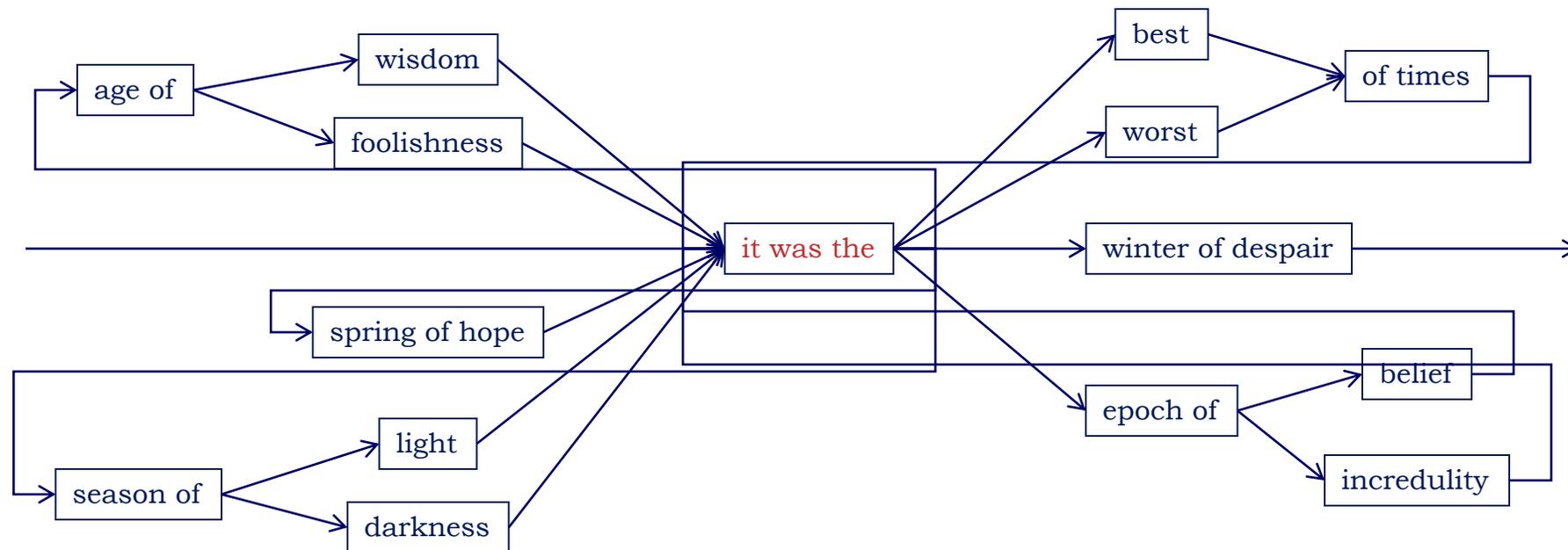
... it was the best of times it was the worst of times ...

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

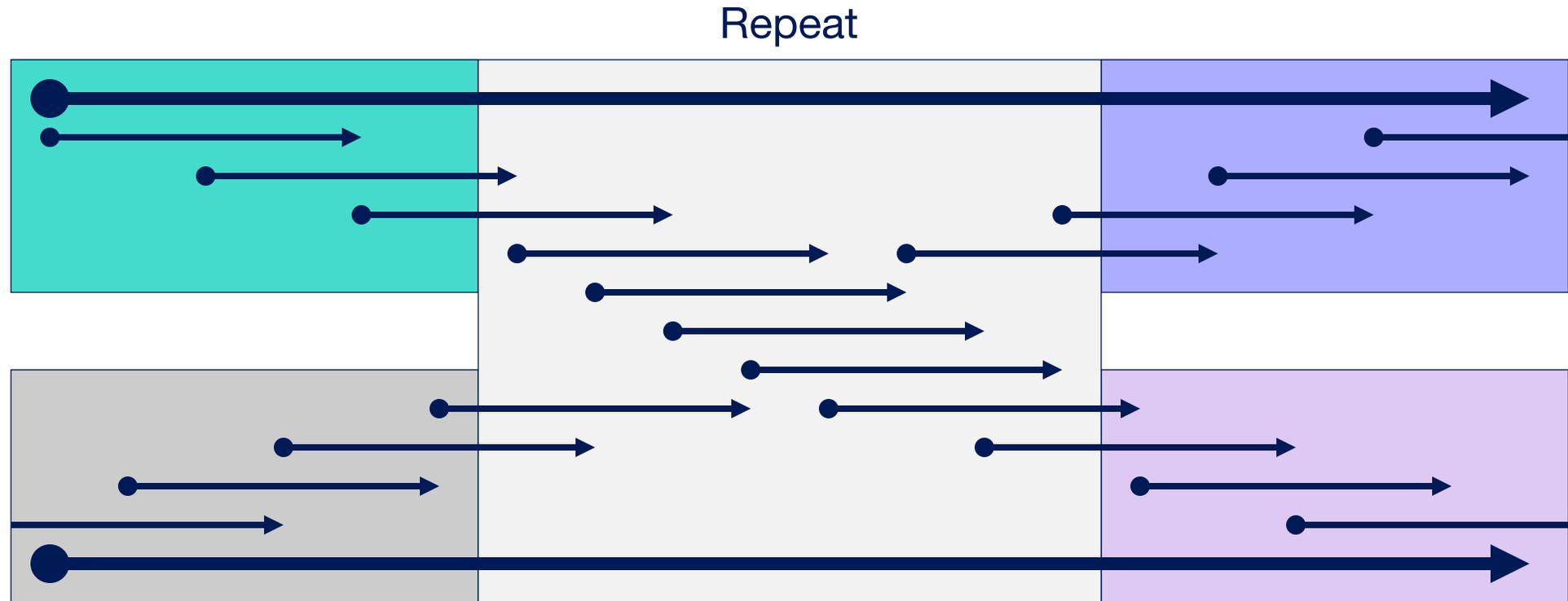
... it was the epoch of belief it was the epoch of incredulity ...

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

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



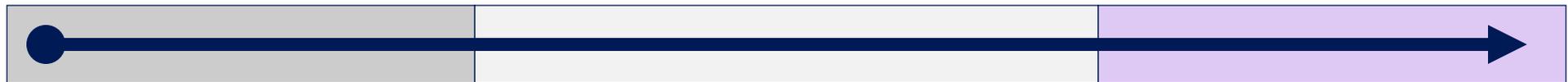
Repeats, repeats, repeats...



Repeats, repeats, repeats...

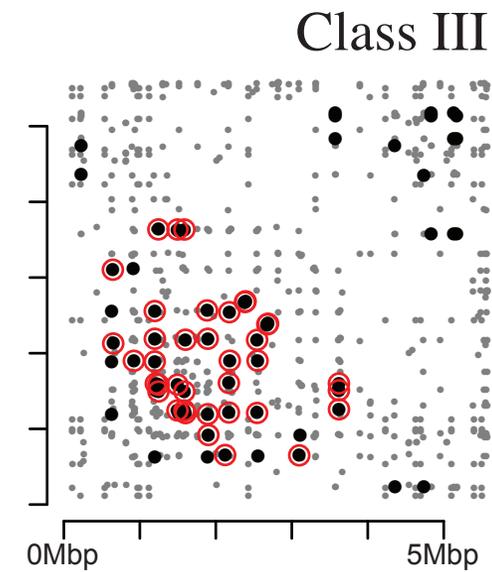
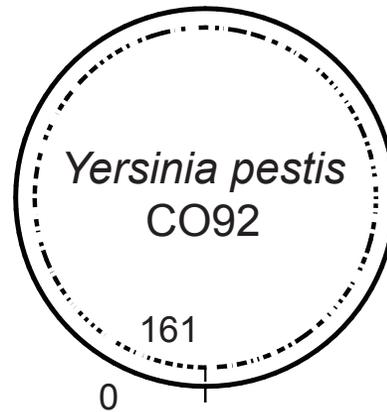
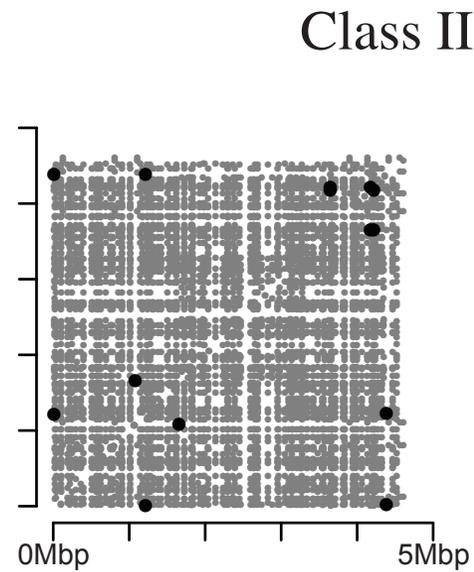
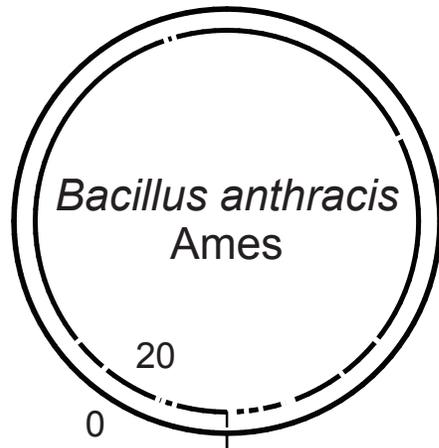
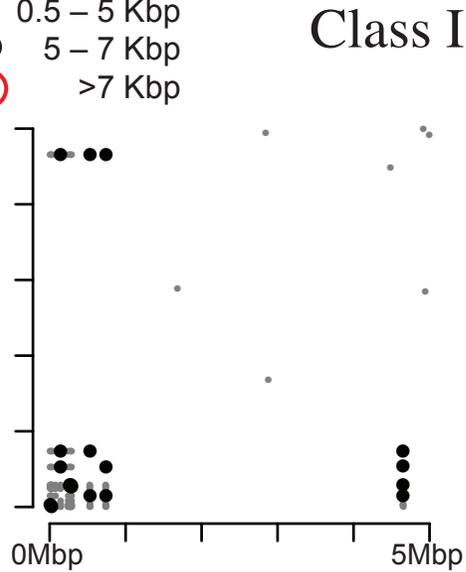


Repeats only matter if longer than the k-mer length



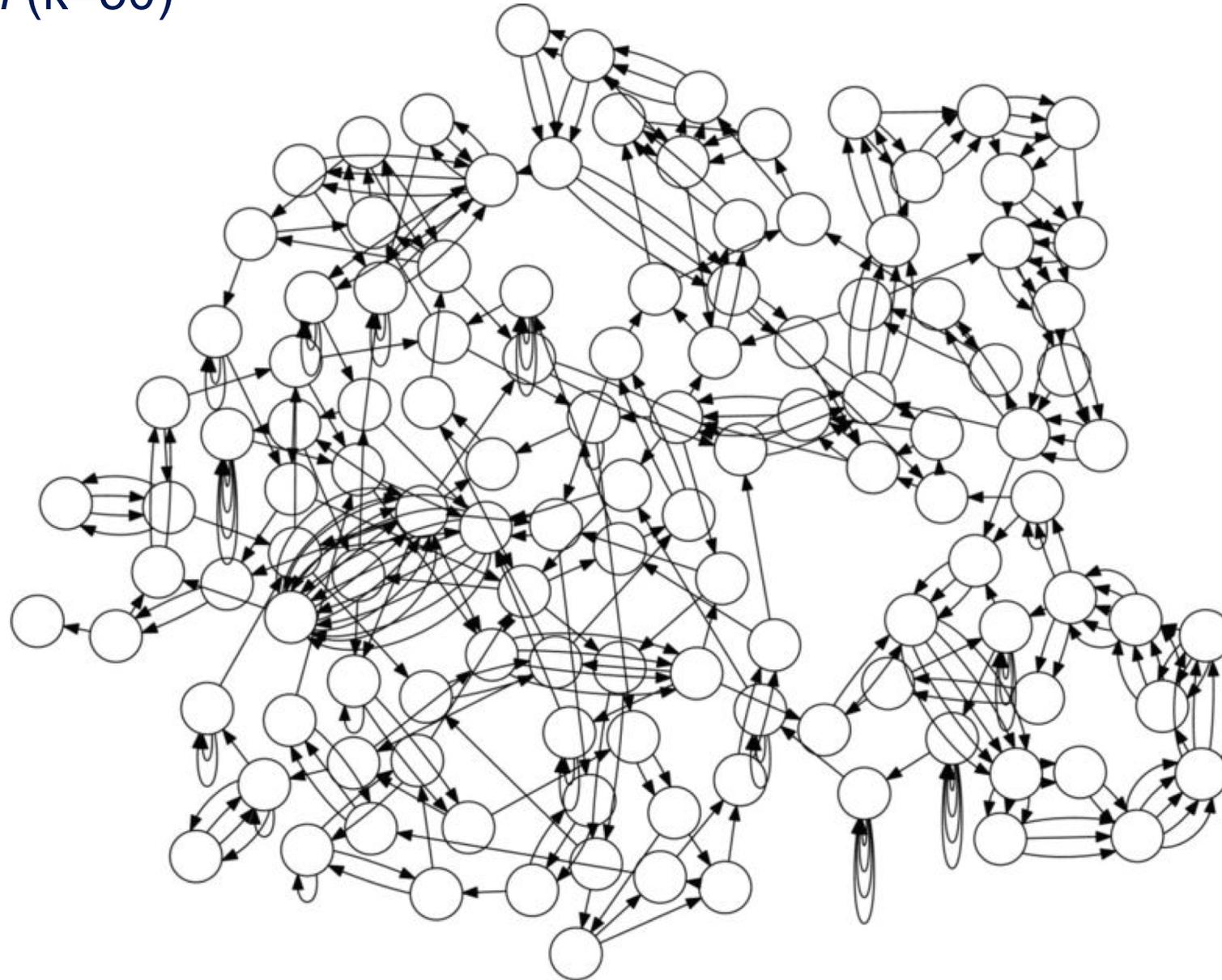
Three classes of complexity

- 0.5 – 5 Kbp
- 5 – 7 Kbp
- >7 Kbp



Reducing assembly complexity of microbial genomes with single-molecule sequencing. Koren *et al.* (2013) *Genome Biology*.

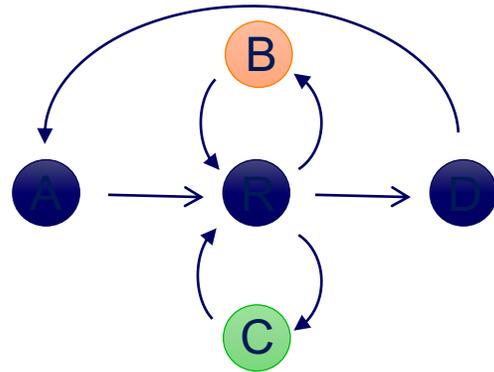
E. coli (k=50)



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>

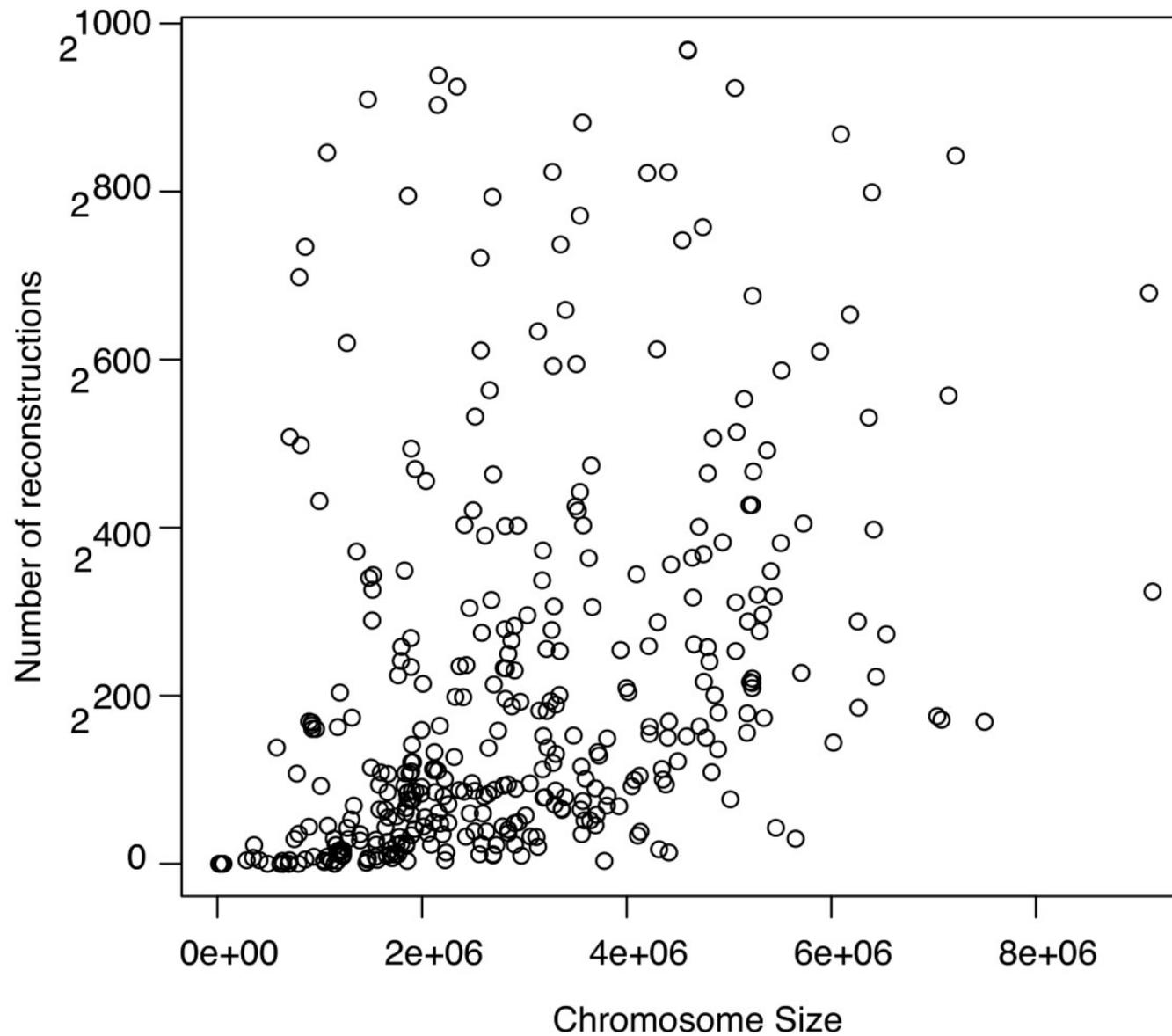
Eulerian Tours



ARBRCD
or
ARCRBD

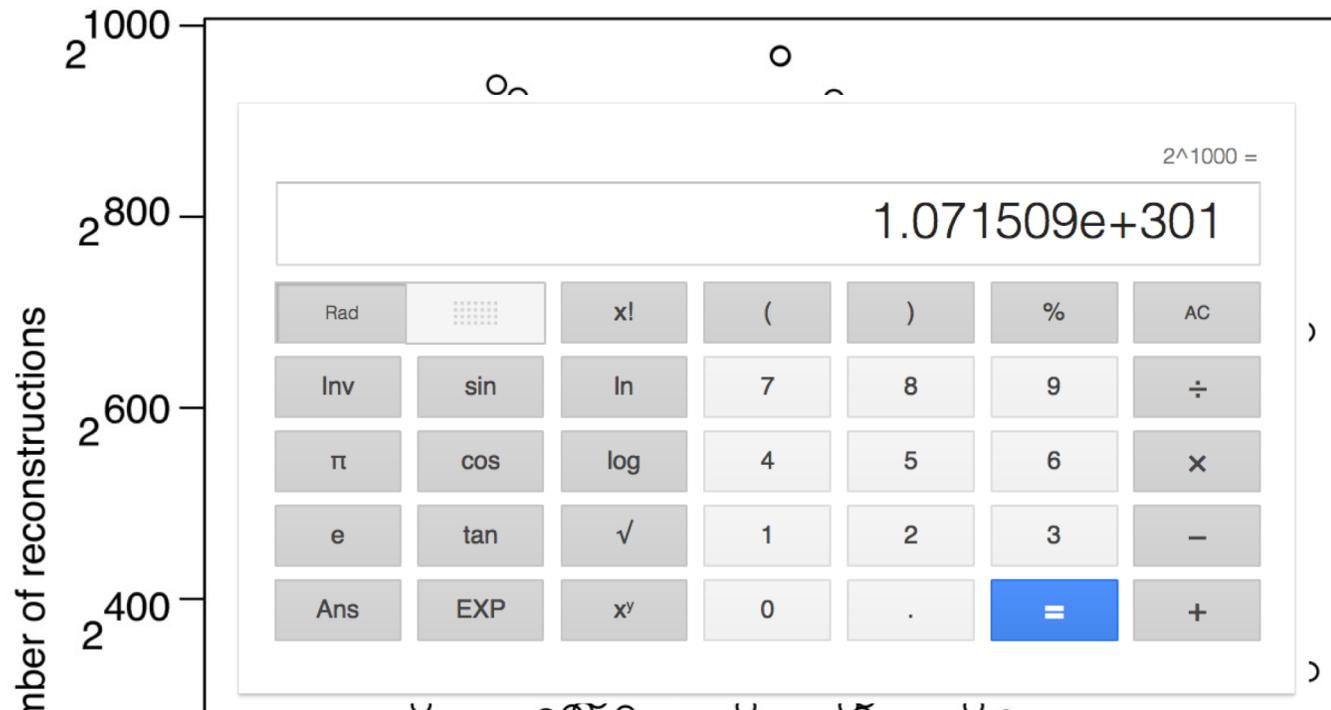
Generally an exponential number of compatible sequences

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



Assembly Complexity of Prokaryotic Genomes using Short Reads.

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



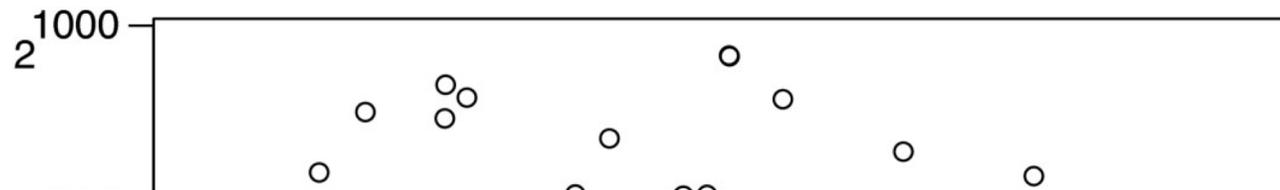
It is believed 74% of the mass of the Milky Way, for example, is in the form of hydrogen atoms. The Sun contains approximately 10^{57} atoms of hydrogen. If you multiple the number of atoms per star (10^{57}) times the estimated number of stars in the universe (10^{23}), you get a value of 10^{80} atoms in the known universe. Nov 5, 2017



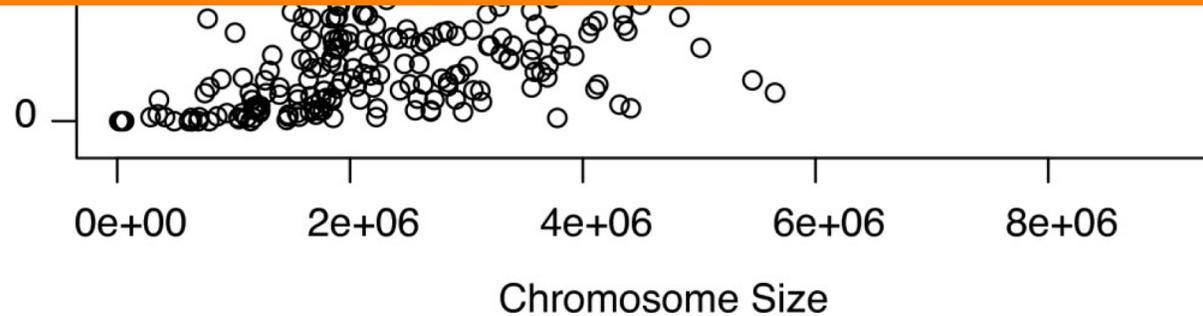
[How Many Atoms Are There in the Universe? - ThoughtCo](https://www.thoughtco.com/number-of-atoms-in-the-universe-603795)
<https://www.thoughtco.com/number-of-atoms-in-the-universe-603795>

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 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

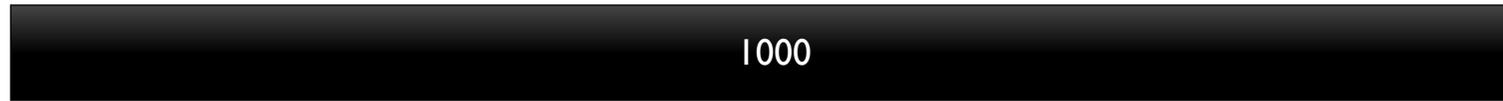
Def: 50% of the genome is in contigs as large as the N50 value

Example: 1 Mbp genome

50%



1000



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 I

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

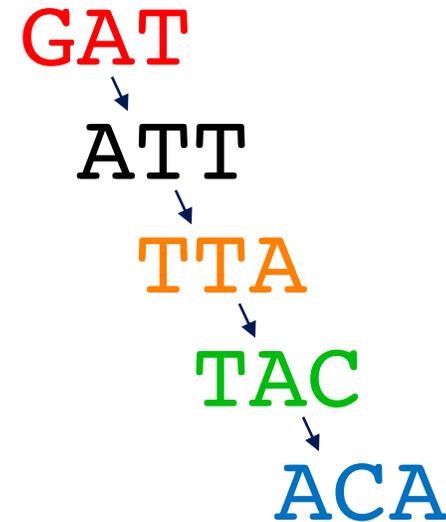
Assemble these reads using a de Bruijn graph approach (k=3):

ATTA: ATT -> TTA

GATT: GAT -> ATT

TACA: TAC -> ACA

TTAC: TTA -> TAC



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

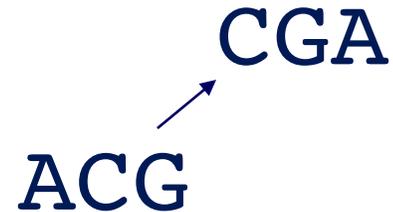
CGAC

CGTA

GACG

GTAT

TACG



Pop Quiz 2

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

~~ACGA~~

~~ACGT~~

ATAC

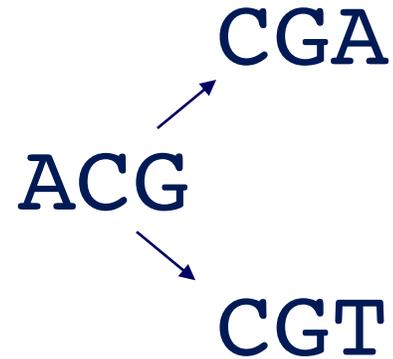
CGAC

CGTA

GACG

GTAT

TACG



Pop Quiz 2

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

~~ACGA~~

~~ACGT~~

ATAC

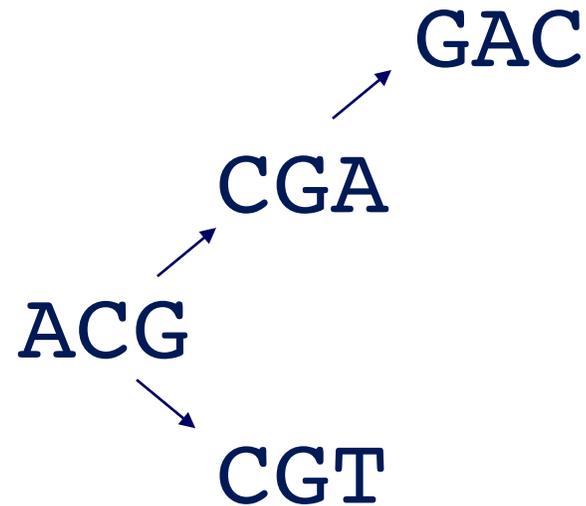
~~CGAC~~

CGTA

GACG

GTAT

TACG



Pop Quiz 2

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

~~ACGA~~

~~ACGT~~

ATAC

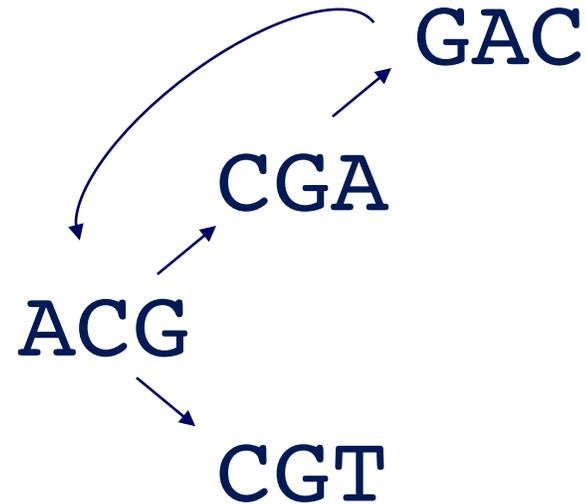
~~CGAC~~

CGTA

~~GACG~~

GTAT

TACG



Pop Quiz 2

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

~~ACGA~~

~~ACGT~~

ATAC

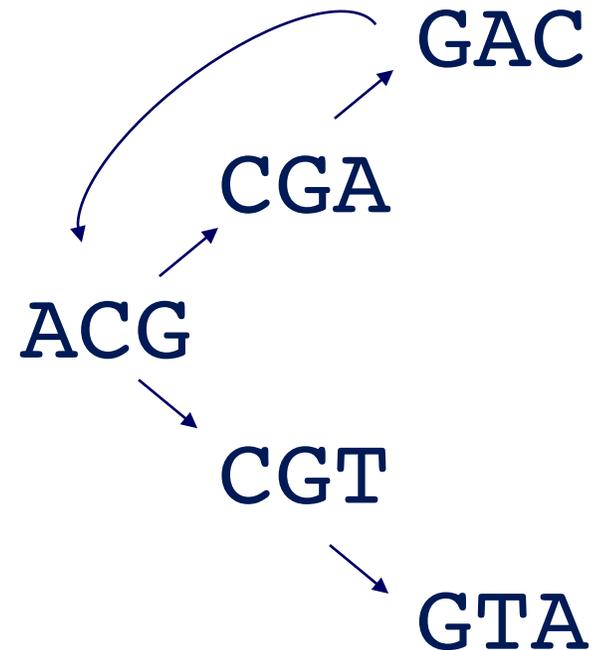
~~CGAC~~

~~CGTA~~

~~GACG~~

GTAT

TACG



Pop Quiz 2

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

~~ACGA~~

~~ACGT~~

ATAC

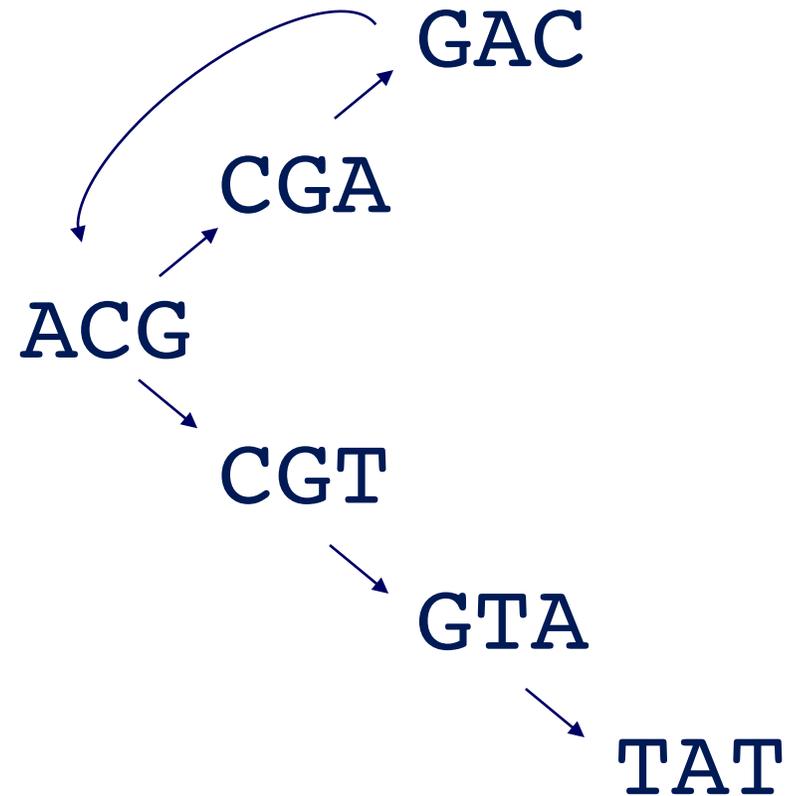
~~CGAC~~

~~CGTA~~

~~GACG~~

~~GTAT~~

TACG



Pop Quiz 2

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

~~ACGA~~

~~ACGT~~

ATAC

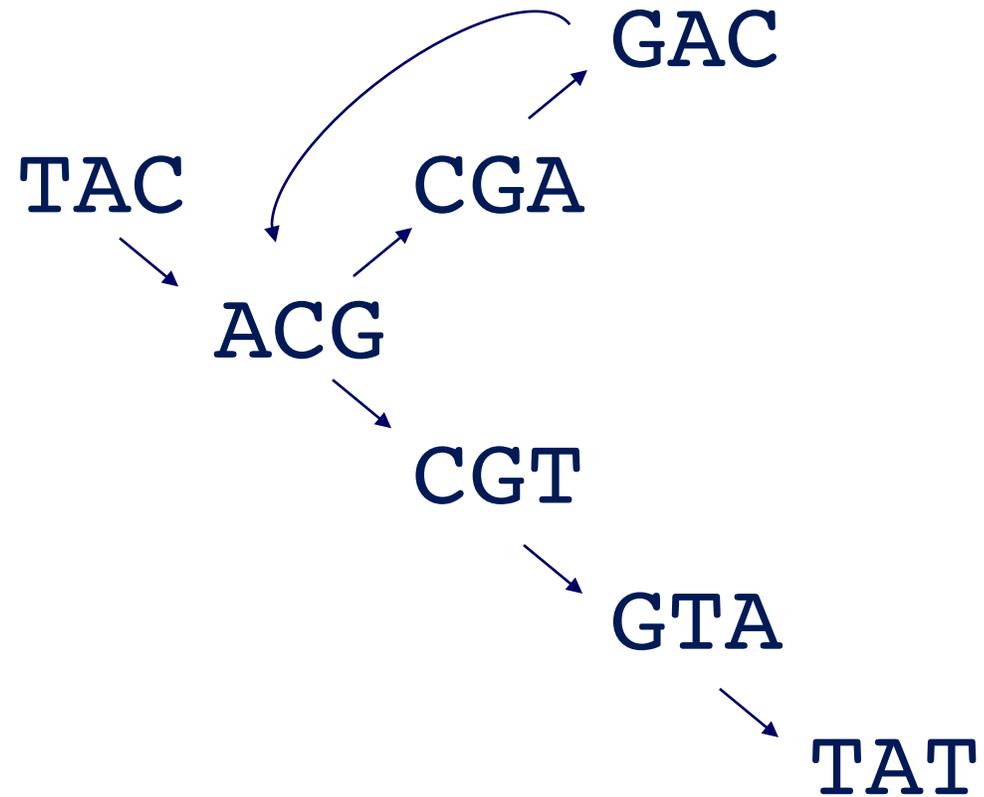
~~CGAC~~

~~CGTA~~

~~GACG~~

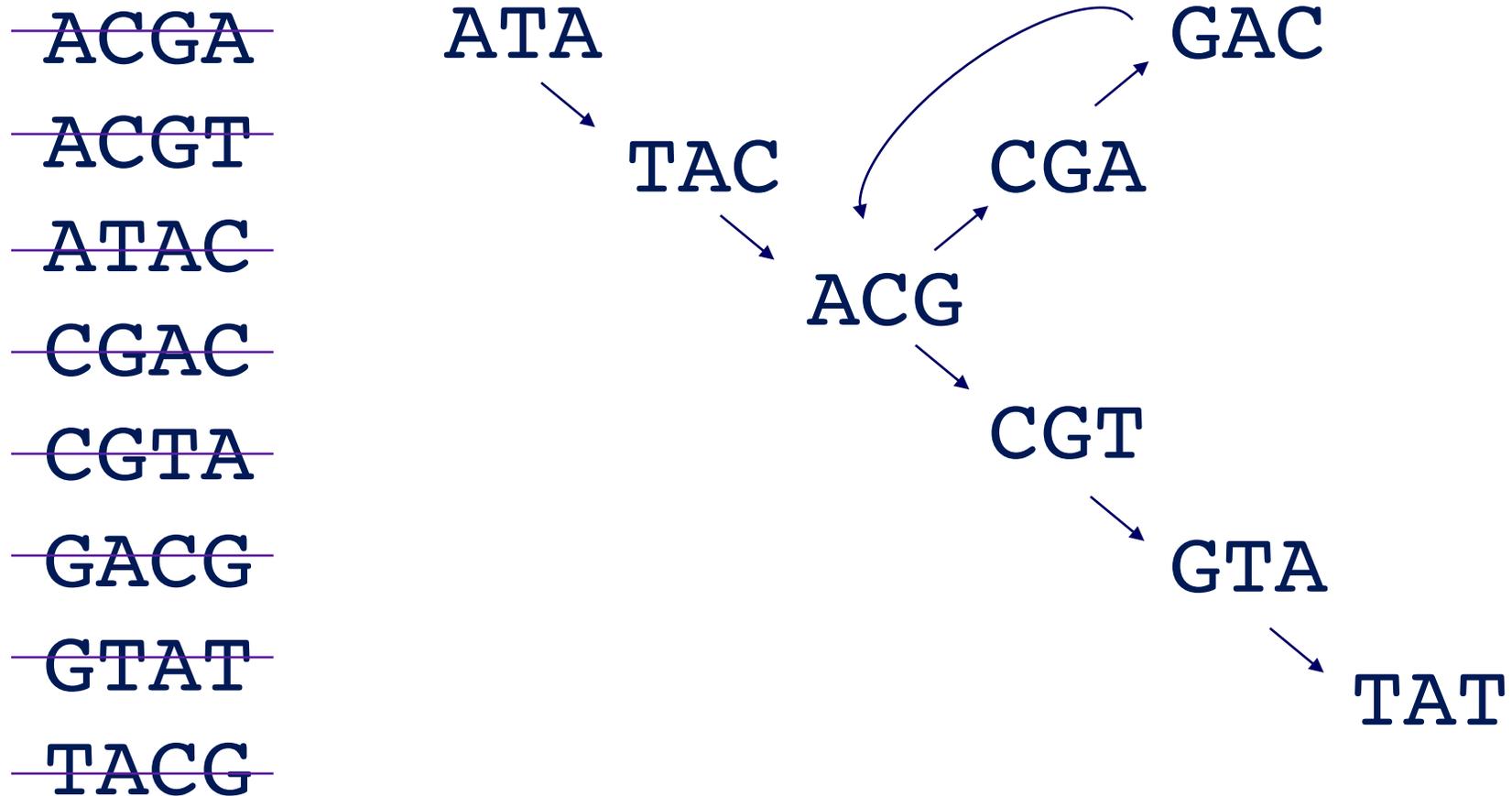
~~GTAT~~

~~TACG~~



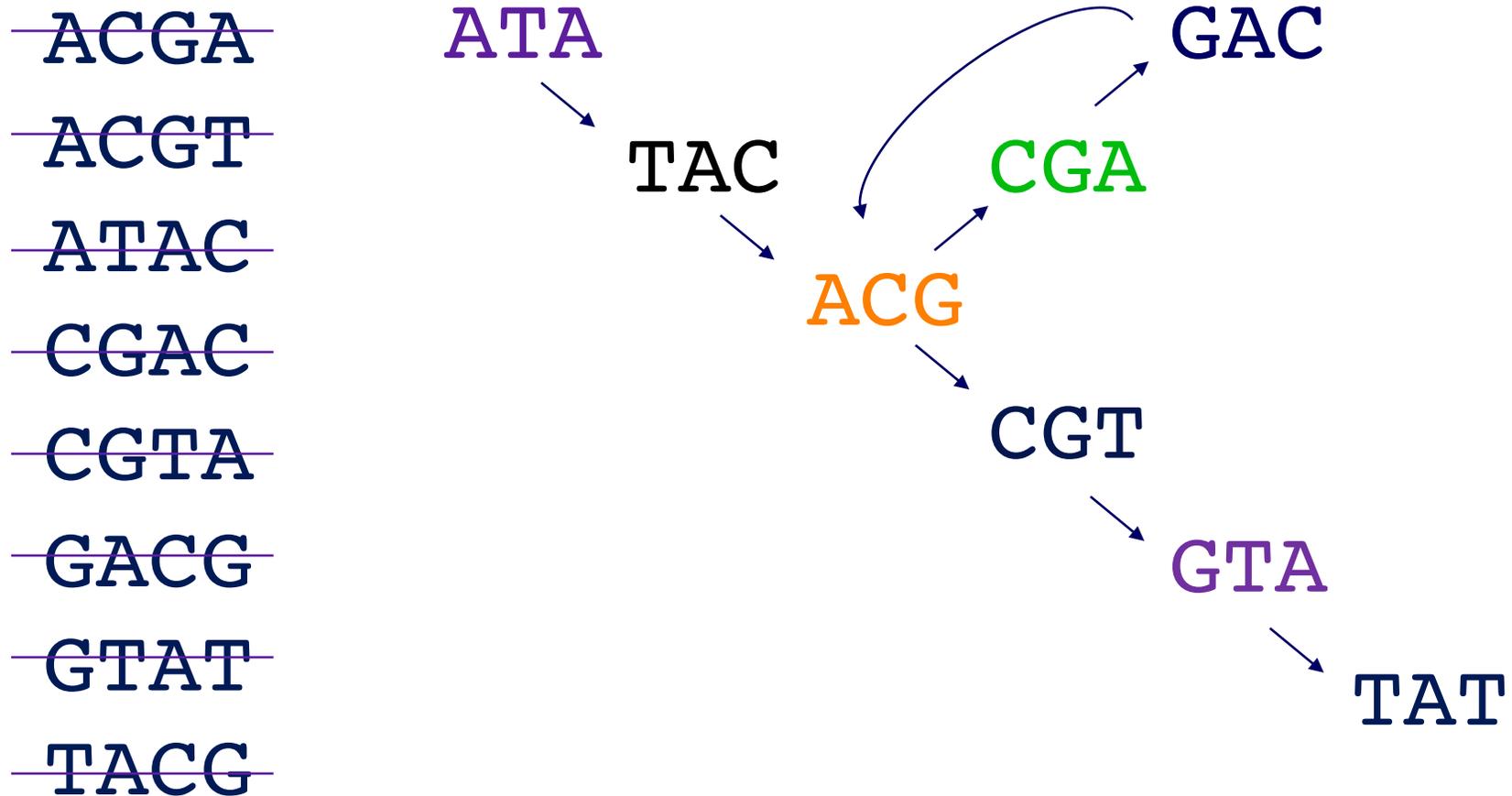
Pop Quiz 2

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



Pop Quiz 2

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

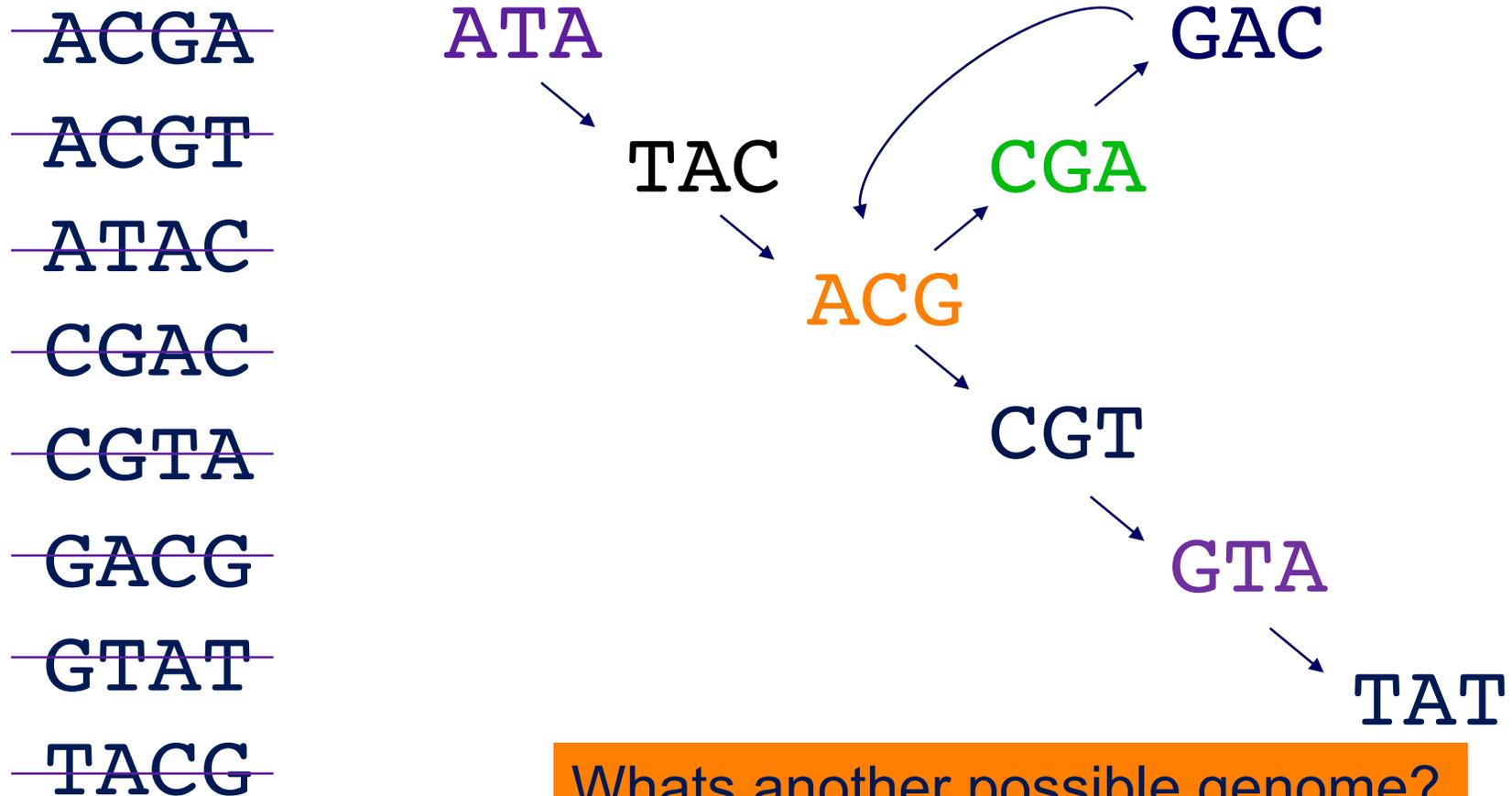


~~ACGA~~
~~ACGT~~
~~ATAC~~
~~CGAC~~
~~CGTA~~
~~GACG~~
~~GTAT~~
~~TACG~~

ATACGACGTAT

Pop Quiz 2

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



Whats another possible genome?

ATACGACGTAT

Assembly Applications

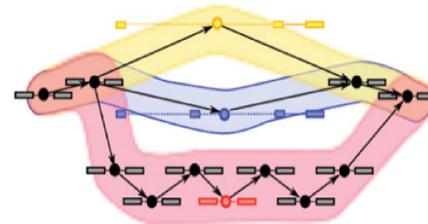
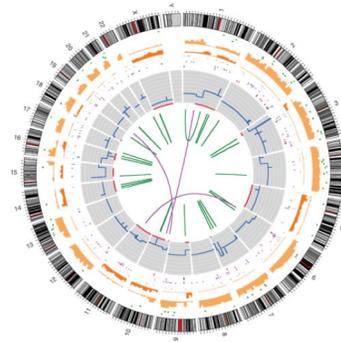
- Novel genomes



- 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

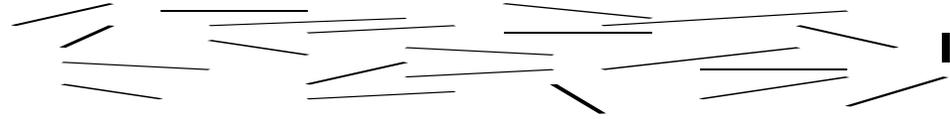
4. *Accuracy:*

- (Very) Hard to assess correctness



Assembling a Genome

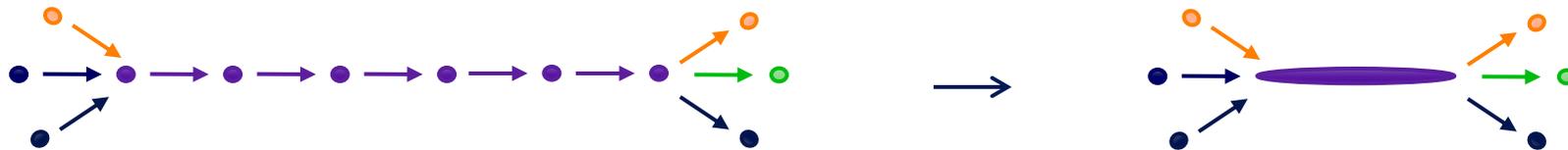
1. Shear & Sequence DNA



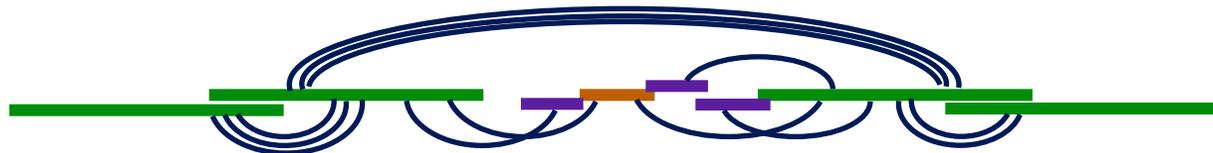
2. Construct assembly graph from reads (de Bruijn / overlap graph)

...AGCCTAGGGATGCGCGACACGT
GGATGCGCGACACGT CGCATATCCGGTTTGGTCAACCTCGGACGGAC
CAACCTCGGACGGACCTCAGCGAA...

3. Simplify assembly graph

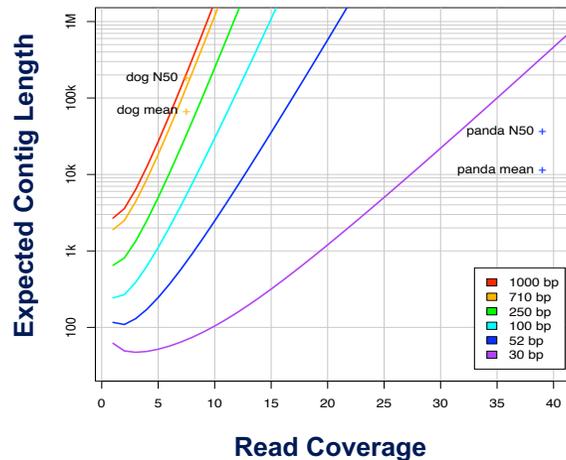


4. Detangle graph with long reads, mates, and other links



Ingredients for a good assembly

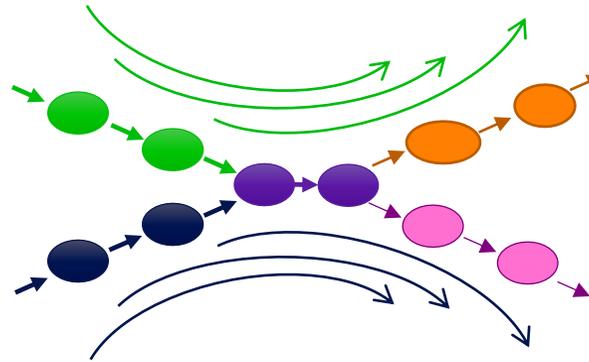
Coverage



High coverage is required

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

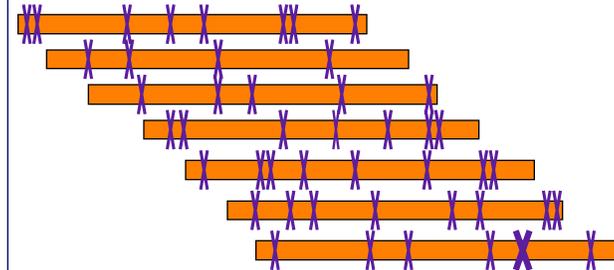
Read Length



Reads & mates must be longer than the repeats

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

Quality



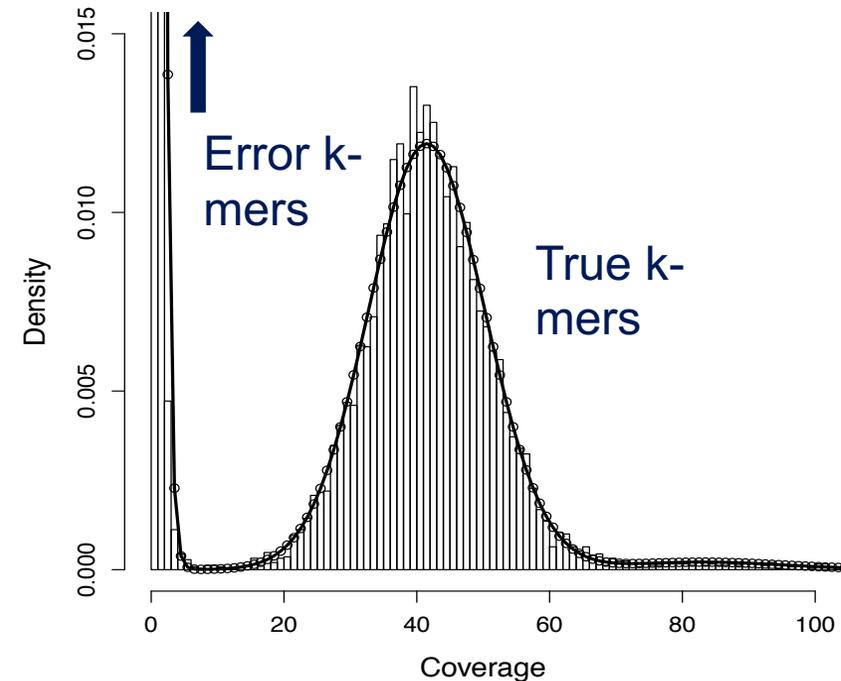
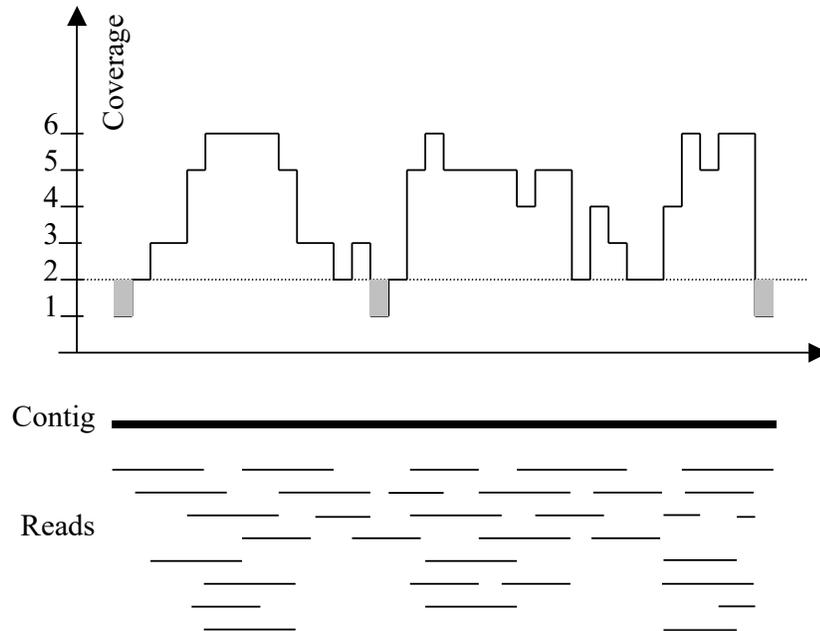
Errors obscure overlaps

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

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

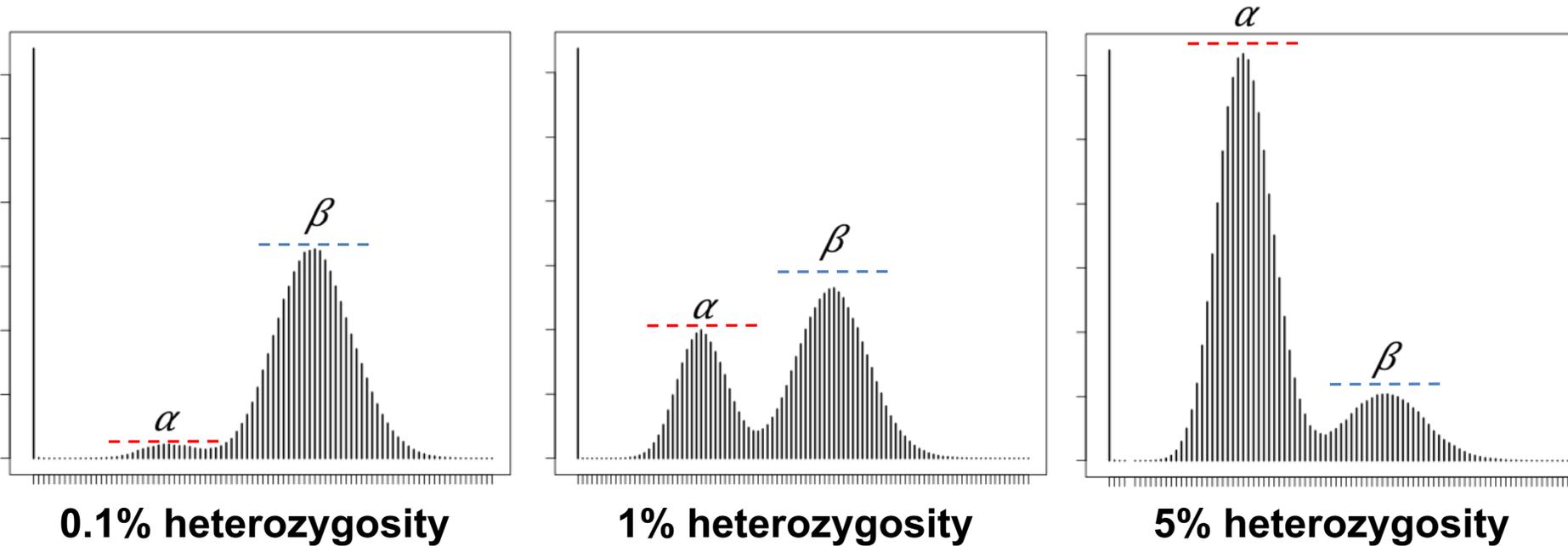
Schatz MC, Witkowski, McCombie, WVR (2012) *Genome Biology*. 12:243

Kmer-based Coverage Analysis



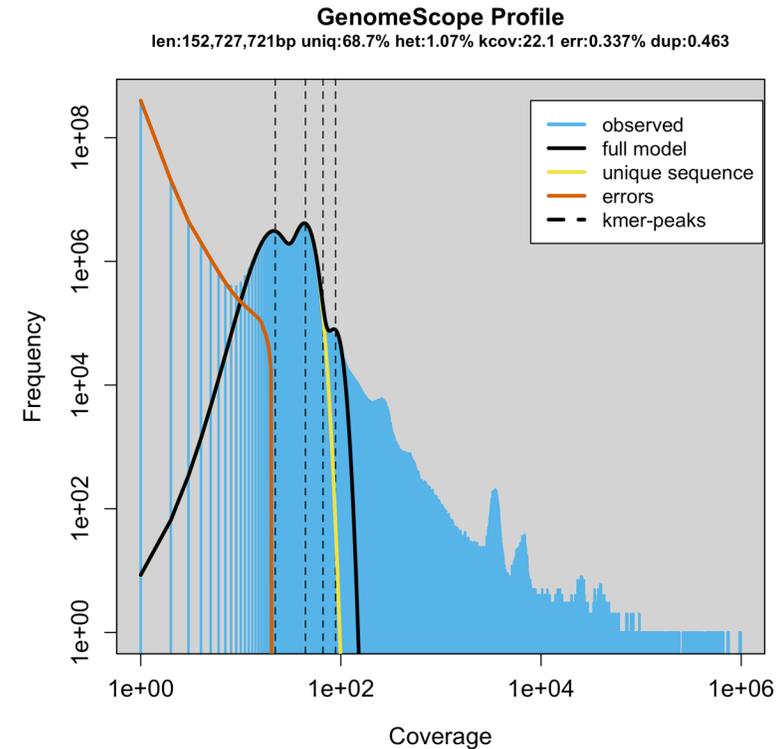
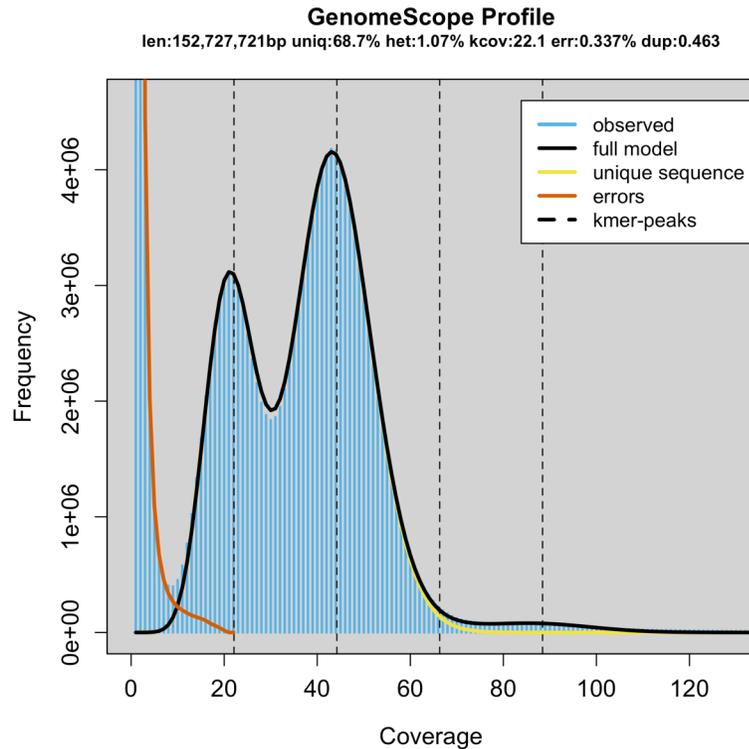
Even though the reads are not assembled or aligned (or reference available), Kmer counting is an effective technique to estimate coverage & other genome properties

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 $2 \cdot k$ heterozygous kmers (typically $k = 21$)

GenomeScope: Fast genome analysis from short reads <http://genomescope.org>

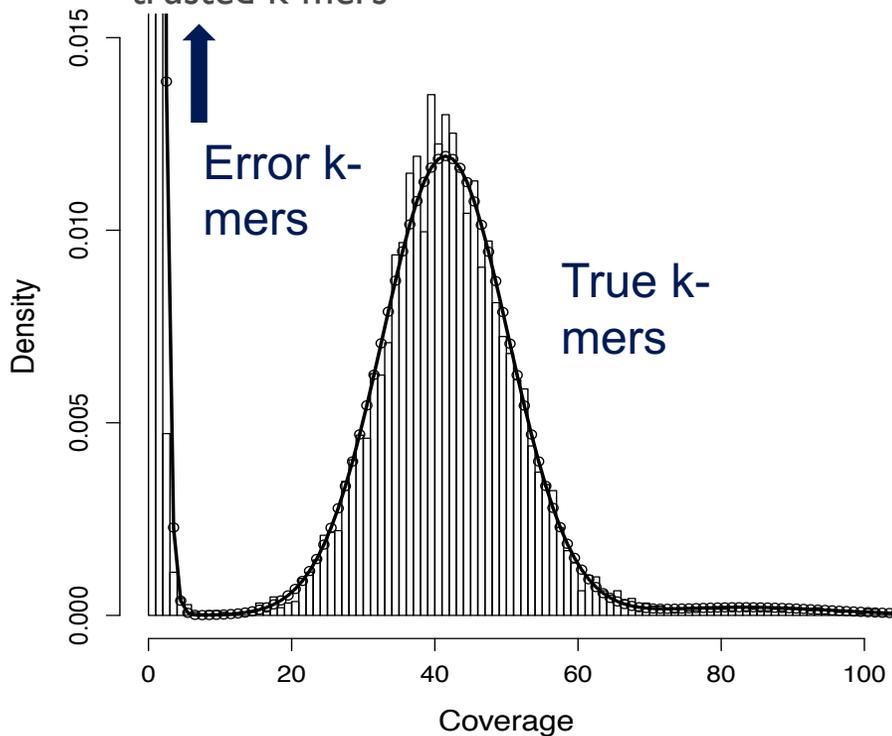


- Theoretical model agrees well with published results:
 - Rate of heterozygosity is higher than reported by other approaches but likely correct.
 - Genome size of plants inflated by organelle sequences (exclude very high freq. kmers)

Error Correction with Quake

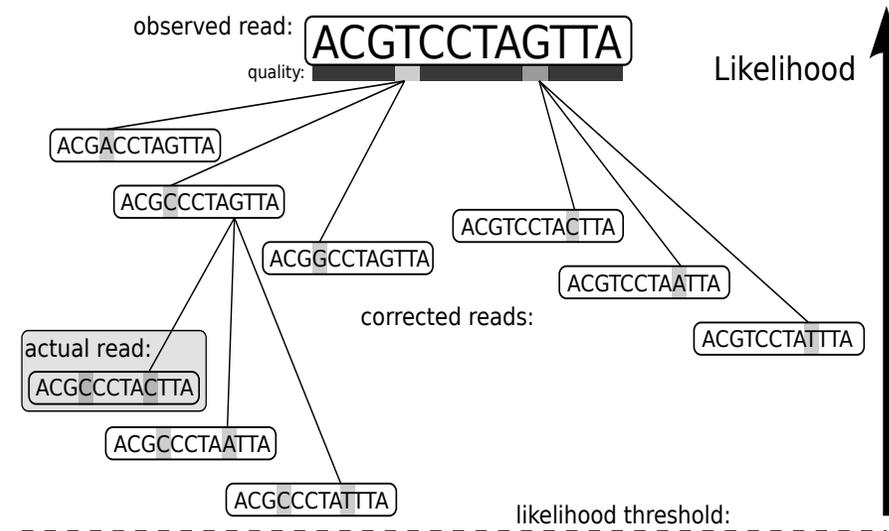
1. Count all “Q-mers” in reads

- 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

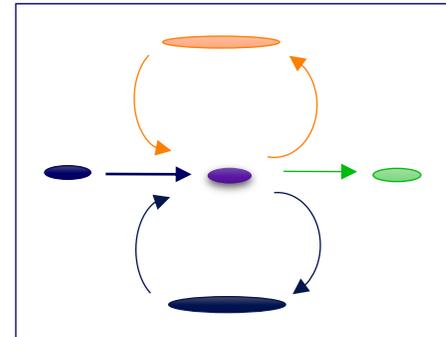
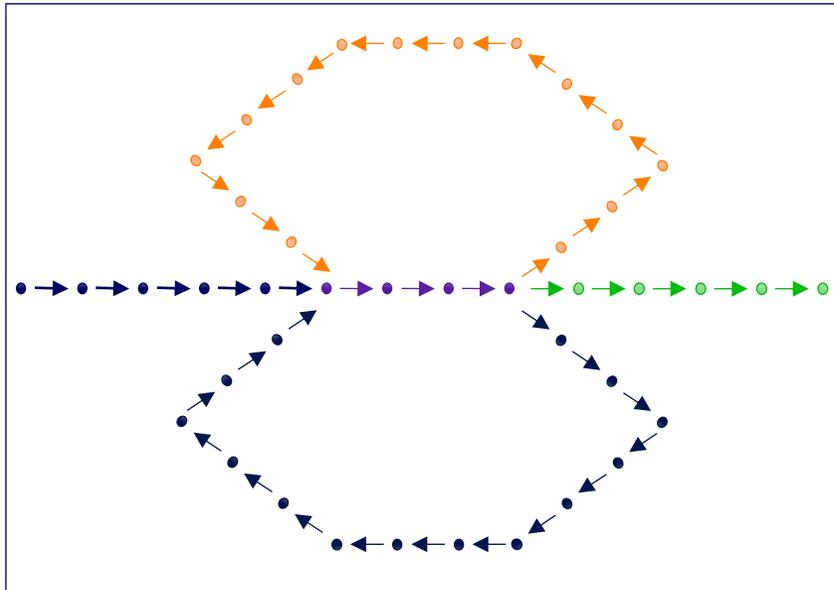


Quake: quality-aware detection and correction of sequencing reads.

Kelley, DR, Schatz, MC, Salzberg SL (2010) *Genome Biology*. 11:R116

Unitigging

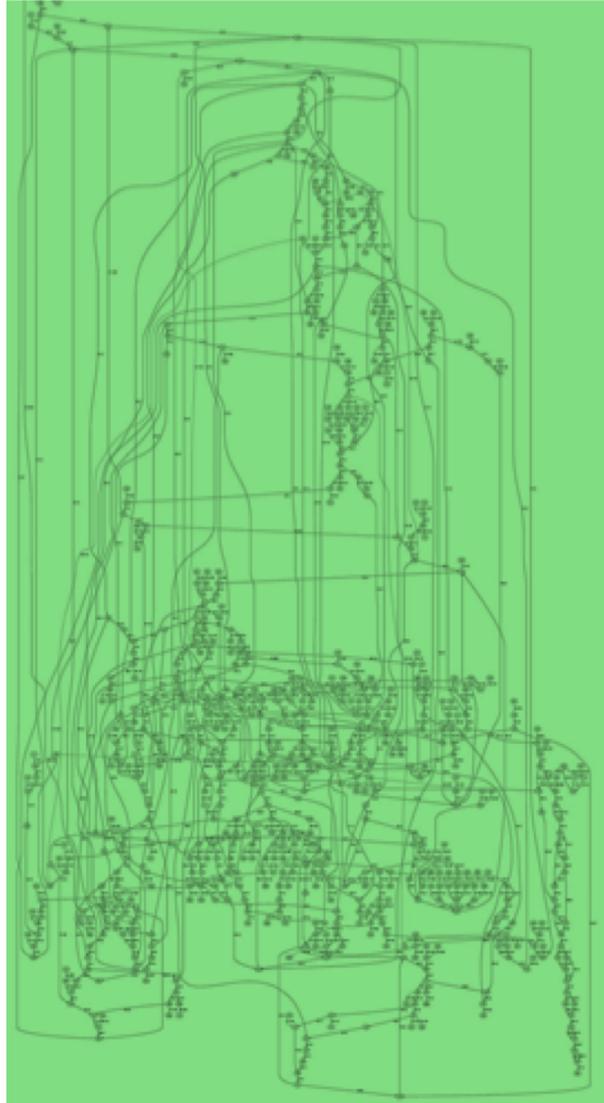
- After simplification and correction, compress graph down to its non-branching initial contigs
 - Aka “unitigs”, “unipaths”



Why do contigs end?

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

Errors in the graph



(Chaisson, 2009)

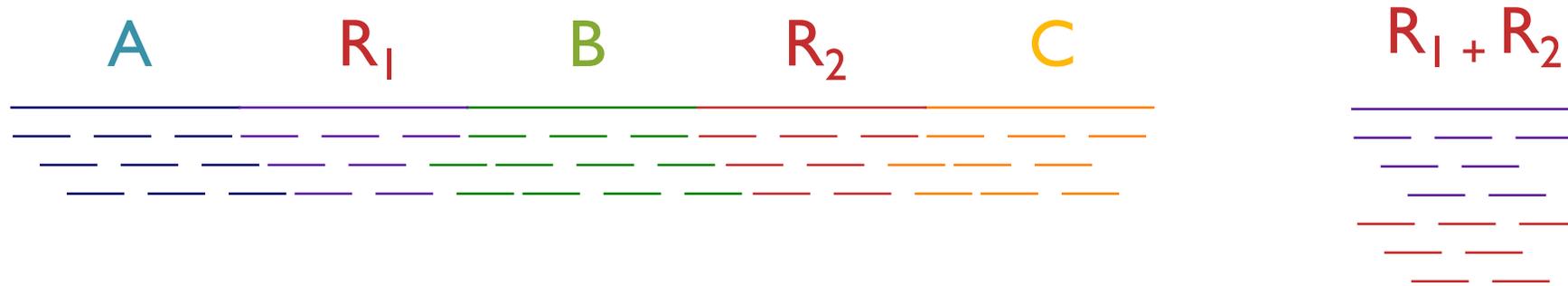
| Clip Tips | Pop Bubbles |
|---|--|
| <div data-bbox="1065 444 1472 501" style="border: 1px solid black; padding: 2px; margin-bottom: 10px;">was the worst of times,</div> <div data-bbox="1065 558 1472 615" style="border: 1px solid black; padding: 2px; margin-bottom: 10px;">was the worst of tymes,</div> <div data-bbox="1090 658 1447 715" style="border: 1px solid black; padding: 2px;">the worst of times, it</div> | <div data-bbox="1709 422 2117 479" style="border: 1px solid black; padding: 2px; margin-bottom: 10px;">was the worst of times,</div> <div data-bbox="1709 515 2117 572" style="border: 1px solid black; padding: 2px; margin-bottom: 10px;">was the worst of tymes,</div> <div data-bbox="1735 608 2091 665" style="border: 1px solid black; padding: 2px; margin-bottom: 10px;">times, it was the age</div> <div data-bbox="1735 701 2091 758" style="border: 1px solid black; padding: 2px;">tymes, it was the age</div> |
| <div data-bbox="1149 972 1493 1029" style="border: 1px solid black; padding: 2px; margin-bottom: 5px;">the worst of tymes,</div> <div data-bbox="1072 1065 1365 1122" style="border: 1px solid black; padding: 2px; margin-bottom: 5px;">was the worst of</div> <div data-bbox="1136 1158 1467 1215" style="border: 1px solid black; padding: 2px; margin-bottom: 5px;">the worst of times,</div> <div data-bbox="1238 1250 1544 1308" style="border: 1px solid black; padding: 2px;">worst of times, it</div> | <div data-bbox="1849 972 1989 1029" style="border: 1px solid black; padding: 2px; margin-bottom: 5px;">tymes,</div> <div data-bbox="1607 1072 1900 1129" style="border: 1px solid black; padding: 2px; margin-bottom: 5px;">was the worst of</div> <div data-bbox="1939 1072 2193 1129" style="border: 1px solid black; padding: 2px; margin-bottom: 5px;">it was the age</div> <div data-bbox="1837 1165 1977 1222" style="border: 1px solid black; padding: 2px;">times,</div> |

Repetitive regions

| Repeat Type | Definition / Example | Prevalence |
|---|---|------------|
| Low-complexity DNA / Microsatellites | $(b_1b_2\dots b_k)^N$ where $1 \leq k \leq 6$ CACACACACACACACACA | 2% |
| SINEs (Short Interspersed Nuclear Elements) | <i>Alu</i> sequence (~280 bp) Mariner elements (~80 bp) | 13% |
| LINEs (Long Interspersed Nuclear Elements) | ~500 – 5,000 bp | 21% |
| LTR (long terminal repeat) retrotransposons | Ty1-copia, Ty3-gypsy, Pao-BEL (~100 – 5,000 bp) | 8% |
| Other DNA transposons | | 3% |
| Gene families & segmental duplications | | 4% |

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

Repeats and Coverage Statistics



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

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

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

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



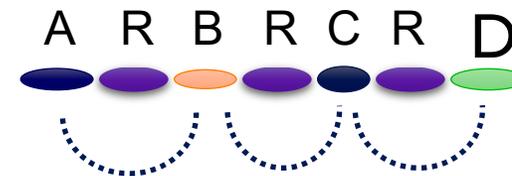
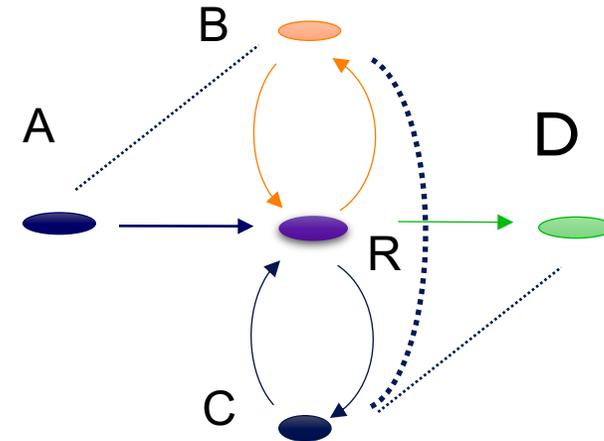
Mate-pair sequencing

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



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



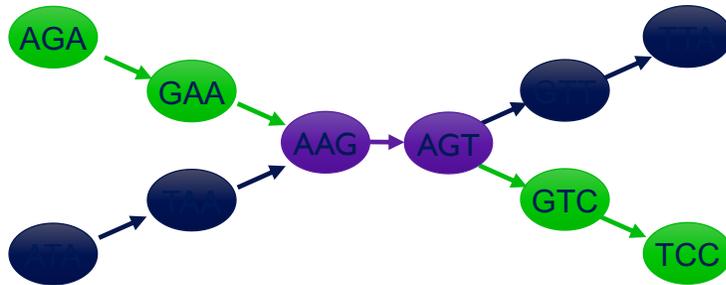
Assembly Summary

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

Two Paradigms for Assembly

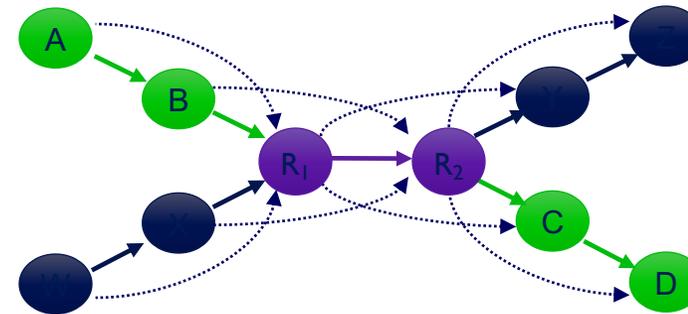
de Bruijn Graph



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.

De Bruijn graph vs. Overlap graph

- **De Bruijn**

- $O(N)$ complexity
- Depends on large k to overcome repeats
- Depends on small k to avoid errors

- **Overlap**

- $O(N^2)$ complexity with naive implementation
- Uses the full length of the reads
- More robust to errors