

Practical Assembly pt2

Michael Schatz

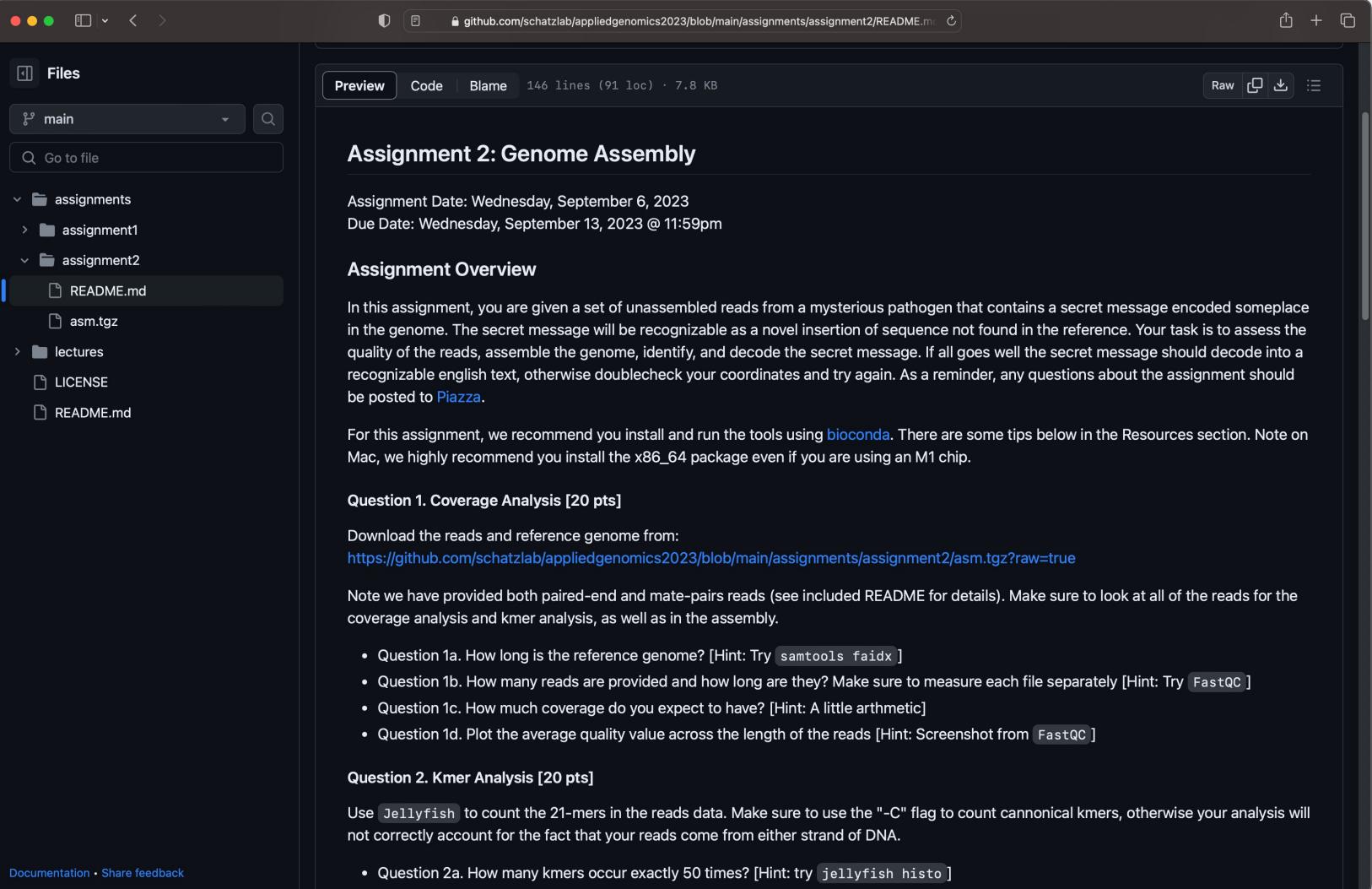
Sept 13, 2023

Lecture 5: Applied Comparative Genomics



Assignment 2: Genome Assembly

Due Wednesday Sept 13 by 11:59pm



The screenshot shows a GitHub repository page for 'assignment2'. The left sidebar lists files and folders: main, assignments (assignment1, assignment2), lectures, LICENSE, and README.md. The README.md file is selected. The main content area shows the file's content:

Assignment 2: Genome Assembly

Assignment Date: Wednesday, September 6, 2023
Due Date: Wednesday, September 13, 2023 @ 11:59pm

Assignment Overview

In this assignment, you are given a set of unassembled reads from a mysterious pathogen that contains a secret message encoded someplace in the genome. The secret message will be recognizable as a novel insertion of sequence not found in the reference. Your task is to assess the quality of the reads, assemble the genome, identify, and decode the secret message. If all goes well the secret message should decode into a recognizable english text, otherwise doublecheck your coordinates and try again. As a reminder, any questions about the assignment should be posted to [Piazza](#).

For this assignment, we recommend you install and run the tools using [bioconda](#). There are some tips below in the Resources section. Note on Mac, we highly recommend you install the x86_64 package even if you are using an M1 chip.

Question 1. Coverage Analysis [20 pts]

Download the reads and reference genome from:
<https://github.com/schatzlab/appliedgenomics2023/blob/main/assignments/assignment2/asm.tgz?raw=true>

Note we have provided both paired-end and mate-pairs reads (see included README for details). Make sure to look at all of the reads for the coverage analysis and kmer analysis, as well as in the assembly.

- Question 1a. How long is the reference genome? [Hint: Try `samtools faidx`]
- Question 1b. How many reads are provided and how long are they? Make sure to measure each file separately [Hint: Try `FastQC`]
- Question 1c. How much coverage do you expect to have? [Hint: A little arithmetic]
- Question 1d. Plot the average quality value across the length of the reads [Hint: Screenshot from `FastQC`]

Question 2. Kmer Analysis [20 pts]

Use `Jellyfish` to count the 21-mers in the reads data. Make sure to use the "-C" flag to count canonical kmers, otherwise your analysis will not correctly account for the fact that your reads come from either strand of DNA.

- Question 2a. How many kmers occur exactly 50 times? [Hint: try `jellyfish histo`]

<https://github.com/schatzlab/appliedgenomics2023/tree/main/assignments/assignment2>

Check Piazza for questions!

Pop Quiz!

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

I need $10\text{Mbp} \times 24\text{x} = 240\text{Mbp}$ 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?

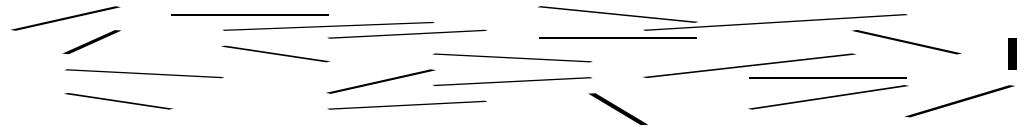
Find X such that $X - 2\sqrt{X} = 24$

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

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

Assembling a Genome

I. Shear & Sequence DNA



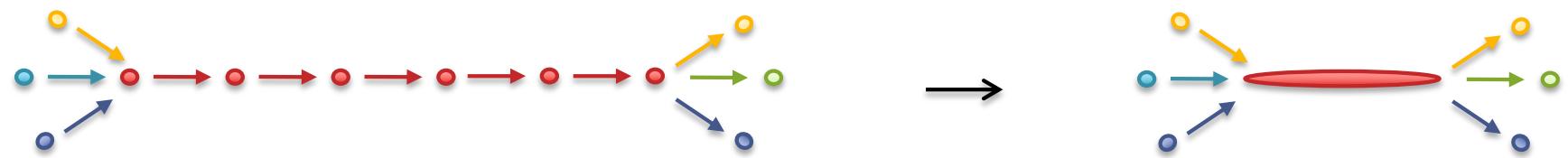
2. Construct assembly graph from reads (de Bruijn / overlap graph)

...AGCCTAG**GGATGCGCGACACGT**

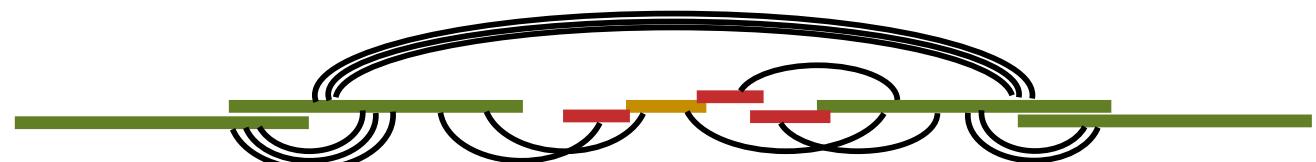
GGATGCGCGACACGTCGCATATCCGGTTTGGT**CAACCTCGGACGGAC**

CAACCTCGGACGGACCTCAGCGAA...

3. Simplify assembly graph

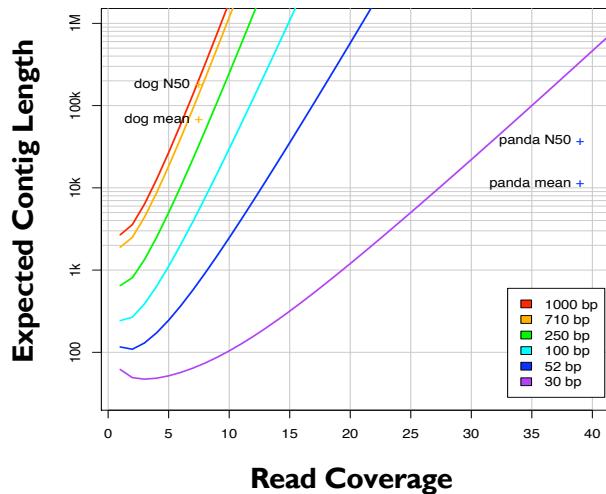


4. Detangle graph with long reads, mates, and other links



Ingredients for a good assembly

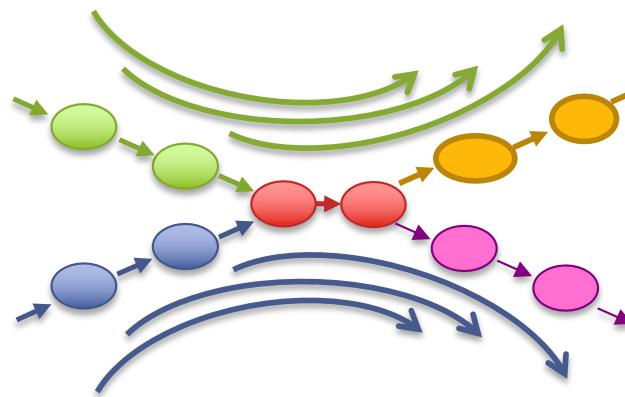
Coverage



High coverage is required

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

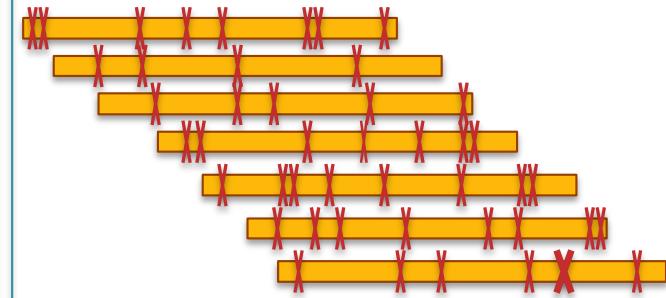
Read Length



Reads & mates must be longer than the repeats

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

Quality



Errors obscure overlaps

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

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

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

Coverage Statistics

$$\text{sequencing_coverage} = \frac{\text{total_bases_sequenced}}{\text{genome_size}}$$

$$\text{genome_size} = \frac{\text{total_bases_sequenced}}{\text{sequencing_coverage}}$$

$$\text{genome_size} = \frac{100\text{Gb}}{50x} = 2\text{Gb}$$

But how can you figure out
the coverage without a genome?

K-mer counting in heterozygous genomes

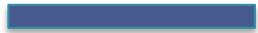
Sequencing read
from homologous
chromosome 1A



Sequencing read
from homologous
chromosome 1B



K-mer counting in heterozygous genomes



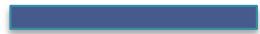
Sequencing read
from homologous
chromosome 1A



Sequencing read
from homologous
chromosome 1B



K-mer counting in heterozygous genomes



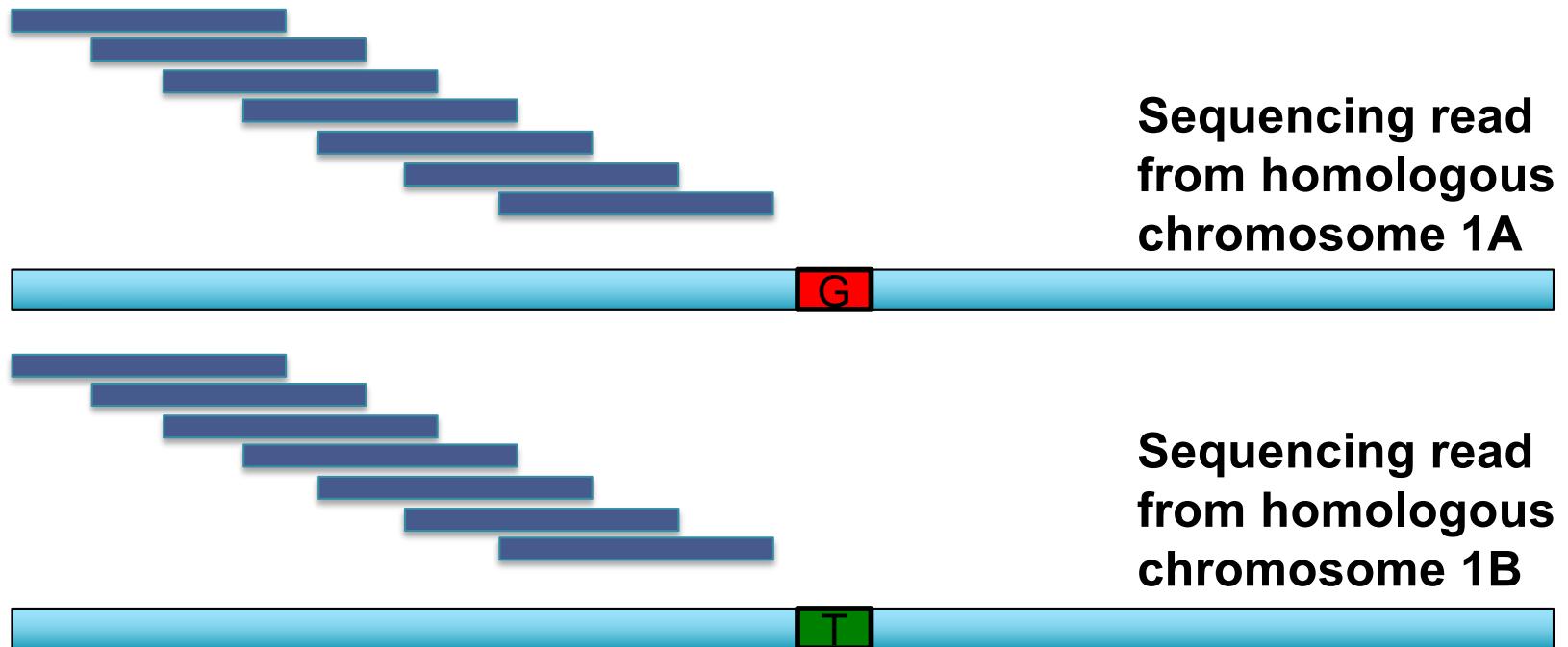
Sequencing read
from homologous
chromosome 1A



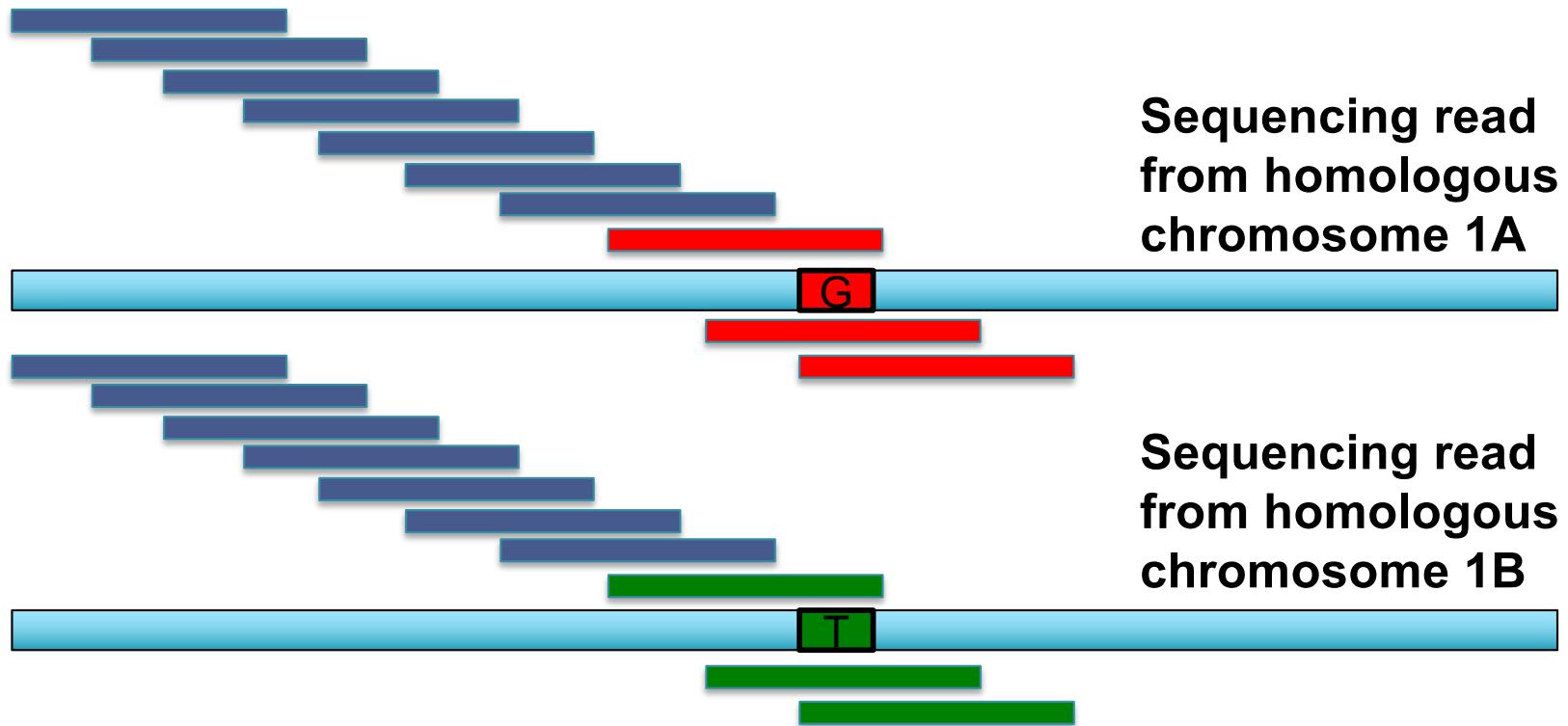
Sequencing read
from homologous
chromosome 1B



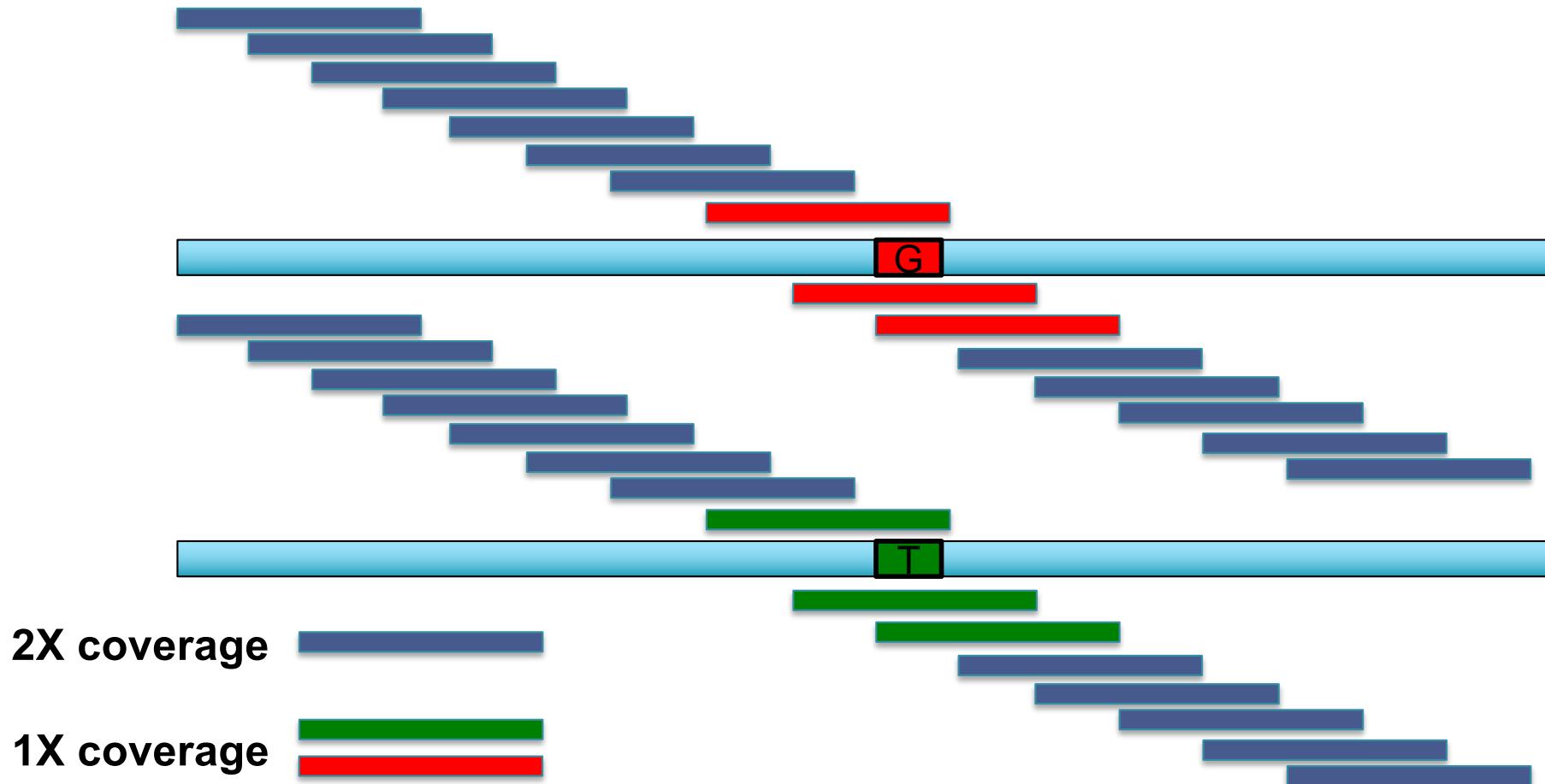
K-mer counting in heterozygous genomes



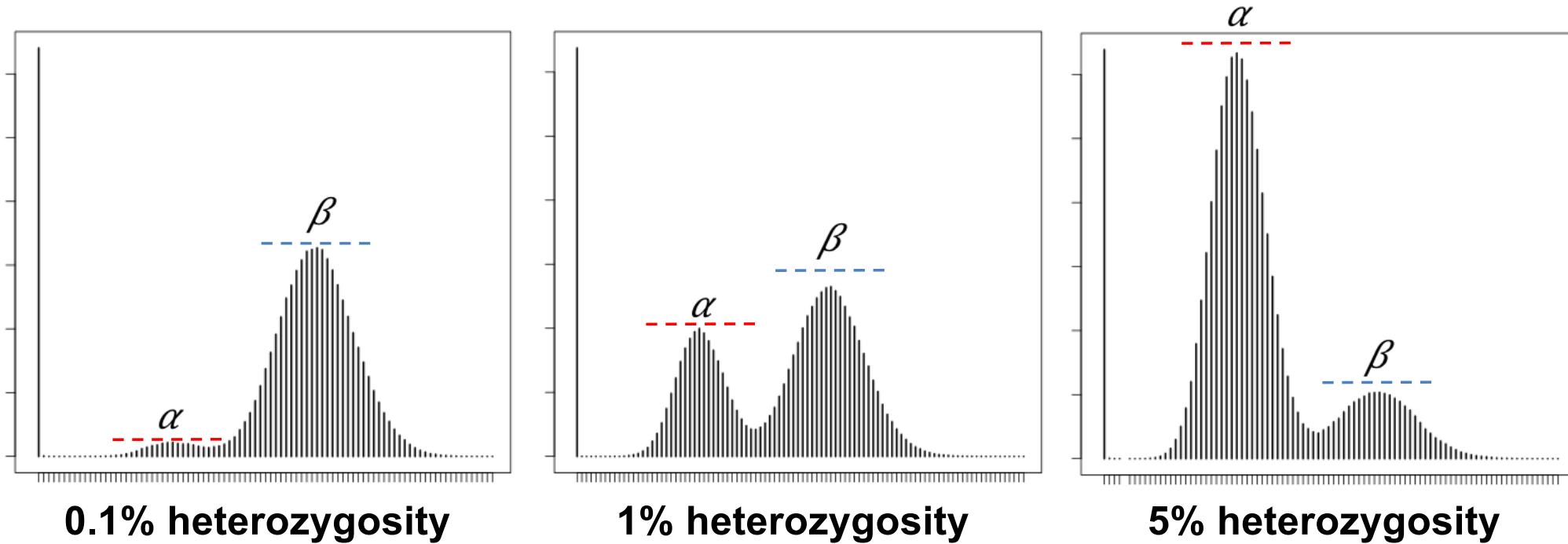
K-mer counting in heterozygous genomes



K-mer counting in heterozygous genomes



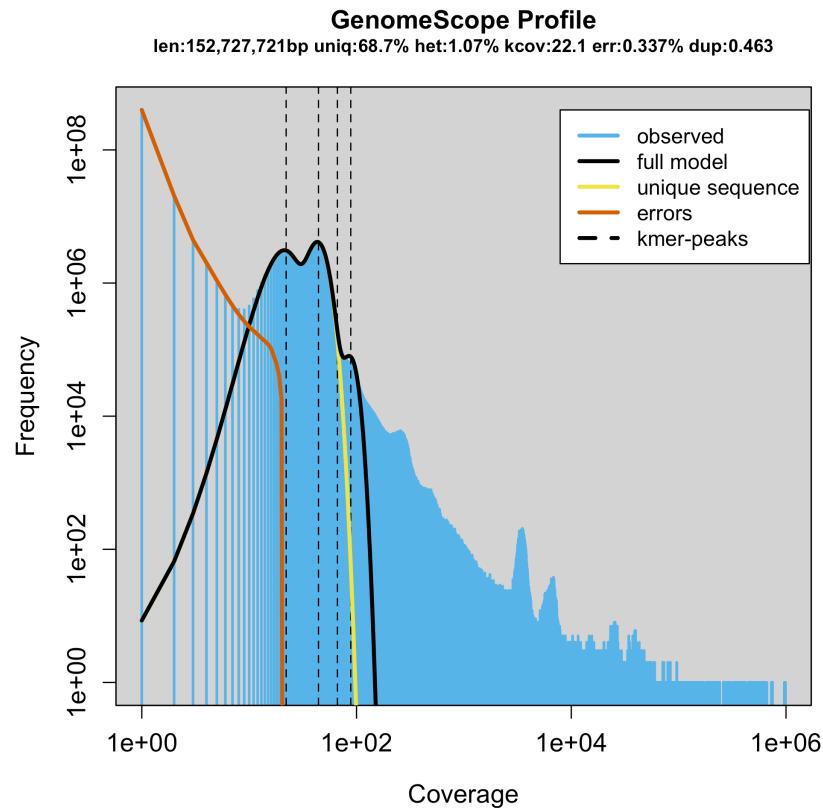
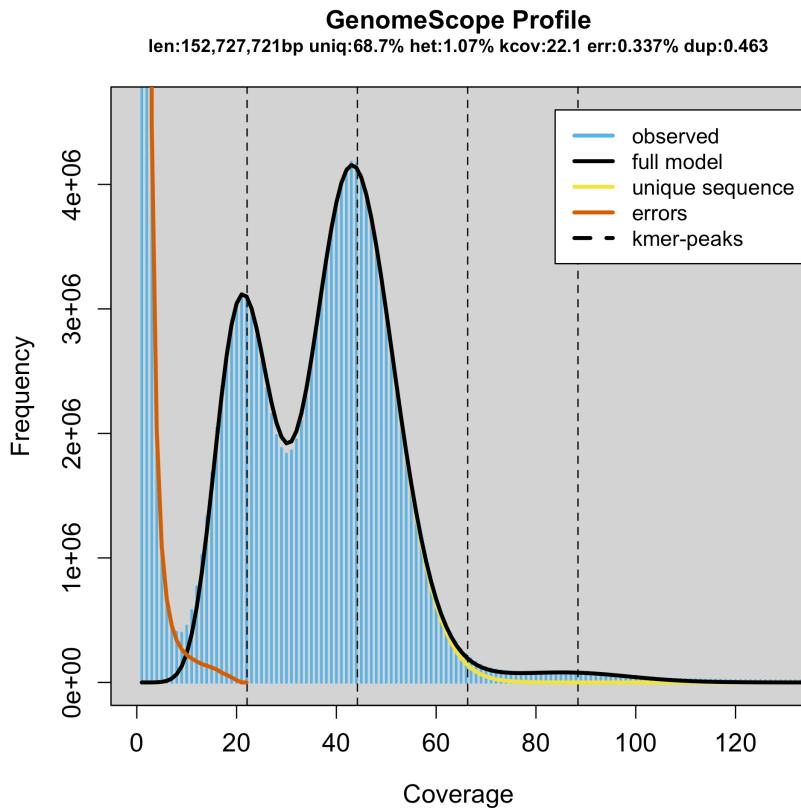
Heterozygous Kmer Profiles



- ***Heterozygosity creates a characteristic “double-peak” in the Kmer profile***
 - Second peak at twice k-mer coverage as the first: heterozygous kmers average 50x coverage, homozygous kmers average 100x coverage
- ***Relative heights of the peaks is directly proportional to the heterozygosity rate***
 - The peaks are balanced at around 1.25% because each heterozygous SNP creates 2^k heterozygous kmers (typically $k = 21$)

GenomeScope: Fast genome analysis from short reads

<http://genomescope.org>

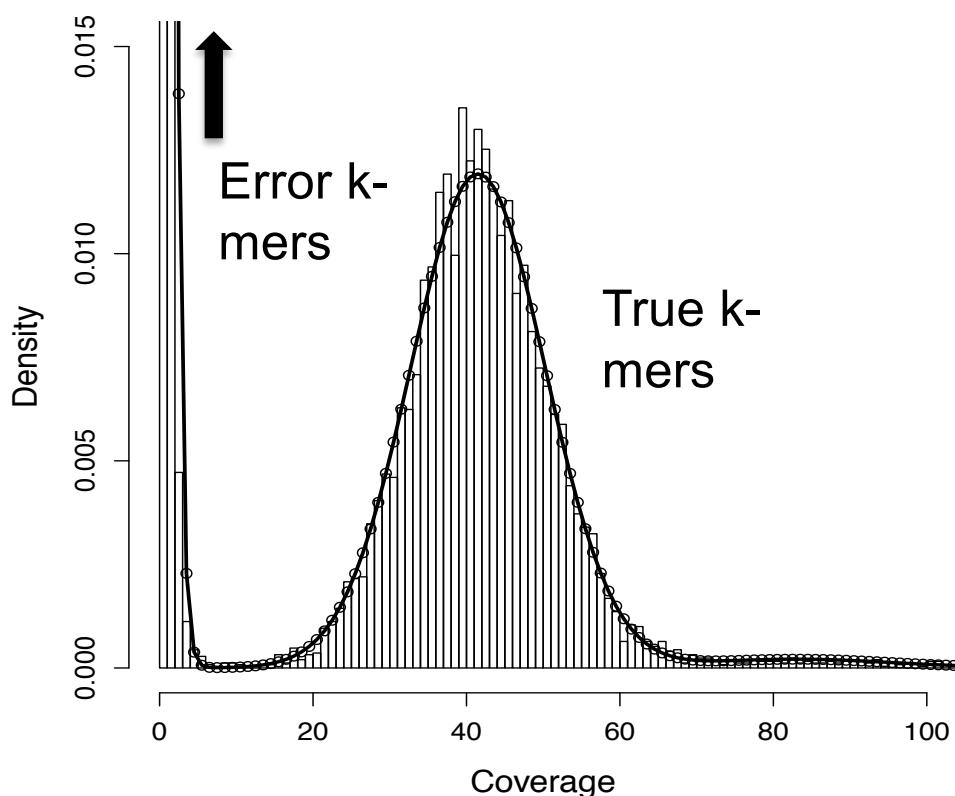


- Theoretical model agrees well with published results:
 - Rate of heterozygosity is higher than reported by other approaches but likely correct.
 - Genome size of plants inflated by organelle sequences (exclude very high freq. kmers)

Error Correction with Quake

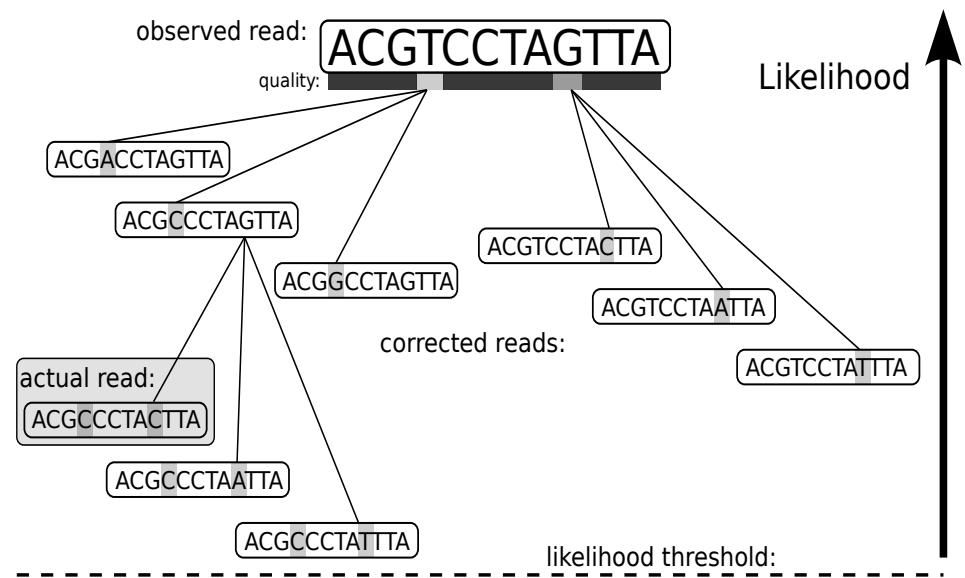
1. Count all “Q-mers” in reads

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



2. Correction Algorithm

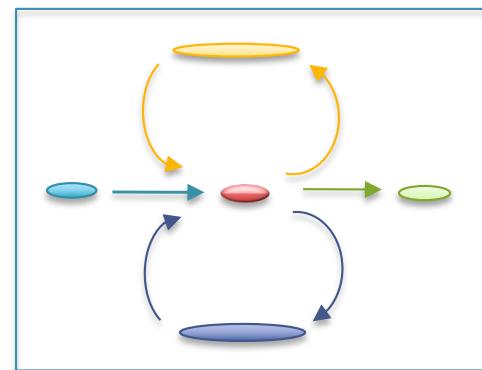
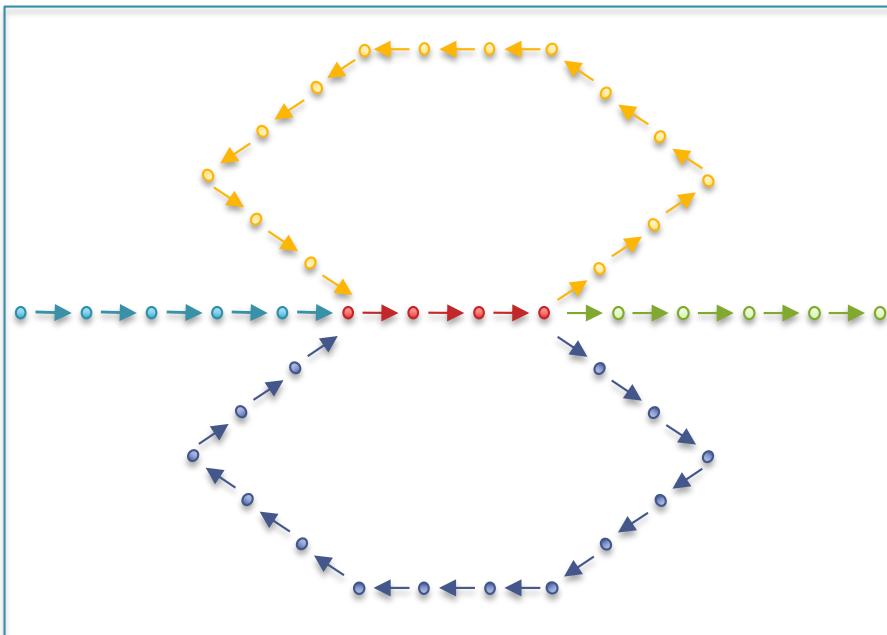
- Considers editing erroneous kmers into trusted kmers in decreasing likelihood
- Includes quality values, nucleotide/nucleotide substitution rate



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

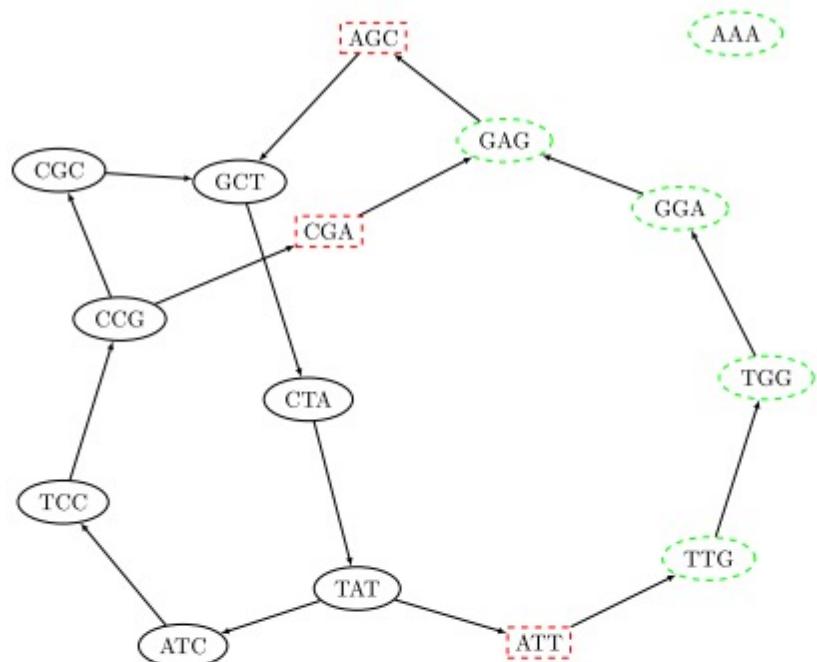
Unitigging / Unipathing

- 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



(a)

$a_1 \dots a_k$	$\sum_{i=1}^k a_i^i \bmod 10$	Bloom filter
ATC	0	0
CCG	0	0
TCC	5	1
CGC	6	1
...	...	0

(b)

(c)

Nodes self-information:
 $\lceil \log_2 \binom{4^3}{7} \rceil = 30 \text{ bits}$

Structure size:
 $\underbrace{10}_{\text{Bloom}} + \underbrace{3 \cdot 6}_{\text{False positives}} = 28 \text{ bits}$

(d)

Space-efficient and exact de Bruijn graph representation based on a Bloom filter
Chikhi and Rizk (2013) *Algorithms for Molecular Biology*. 8:22

Table 2 de novo human genome (NA18507) assemblies

Method	Minia	C. & B.	ABySS	SOAPdenovo
Value of k chosen	27	27	27	25
Number of contigs (M)	3.49	7.69	4.35	-
Longest contig (kbp)	18.6	22.0	15.9	-
Contig N50 (bp)	1156	250	870	886
Sum (Gbp)	2.09	1.72	2.10	2.08
Nb of nodes/cores	1/1	1/8	21/168	1/16
Time (wall-clock, h)	23	50	15	33
Memory (sum of nodes, GB)	5.7	32	336	140

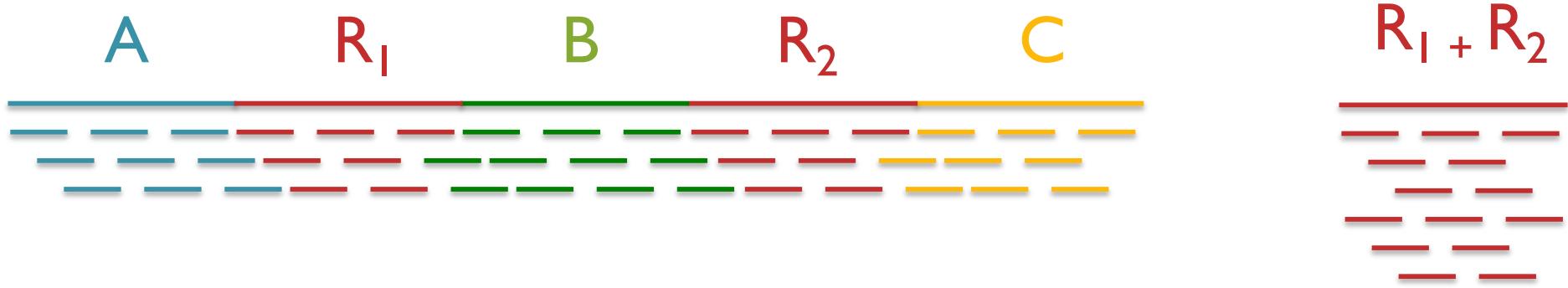
de novo human genome (NA18507) assemblies reported by our assembler (Minia), Conway and Bromage assembler [9], ABySS [8], and SOAPdenovo [7]. Contigs shorter than 100 bp were discarded. Assemblies were made without any pairing information.

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)	Alu 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 $> A$) , it is likely to be a collapsed repeat

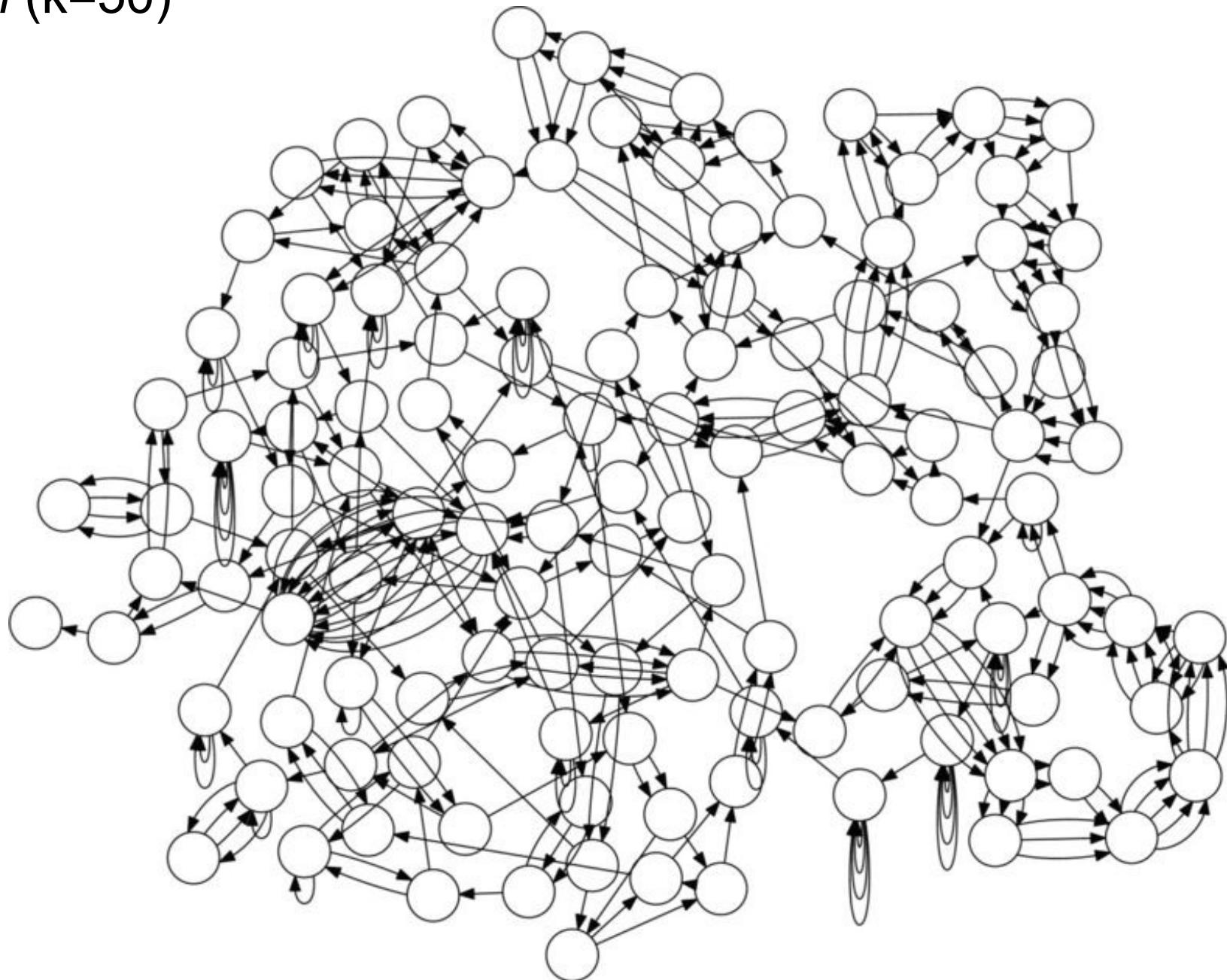
$$\Pr(X - copy) = \binom{n}{k} \left(\frac{X\Delta}{G} \right)^k \left(\frac{G - X\Delta}{G} \right)^{n-k}$$

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

The fragment assembly string graph

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

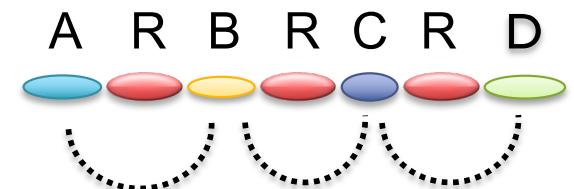
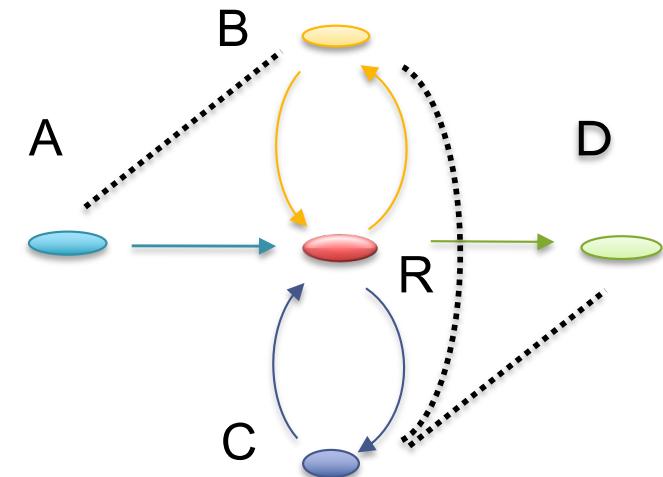
E. coli ($k=50$)



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>

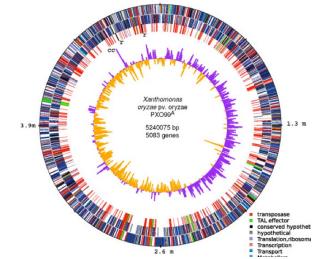
Scaffolding

- 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?

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

Contig N50

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

Example: 1 Mbp genome

50%



A



N50 size = 30 kbp

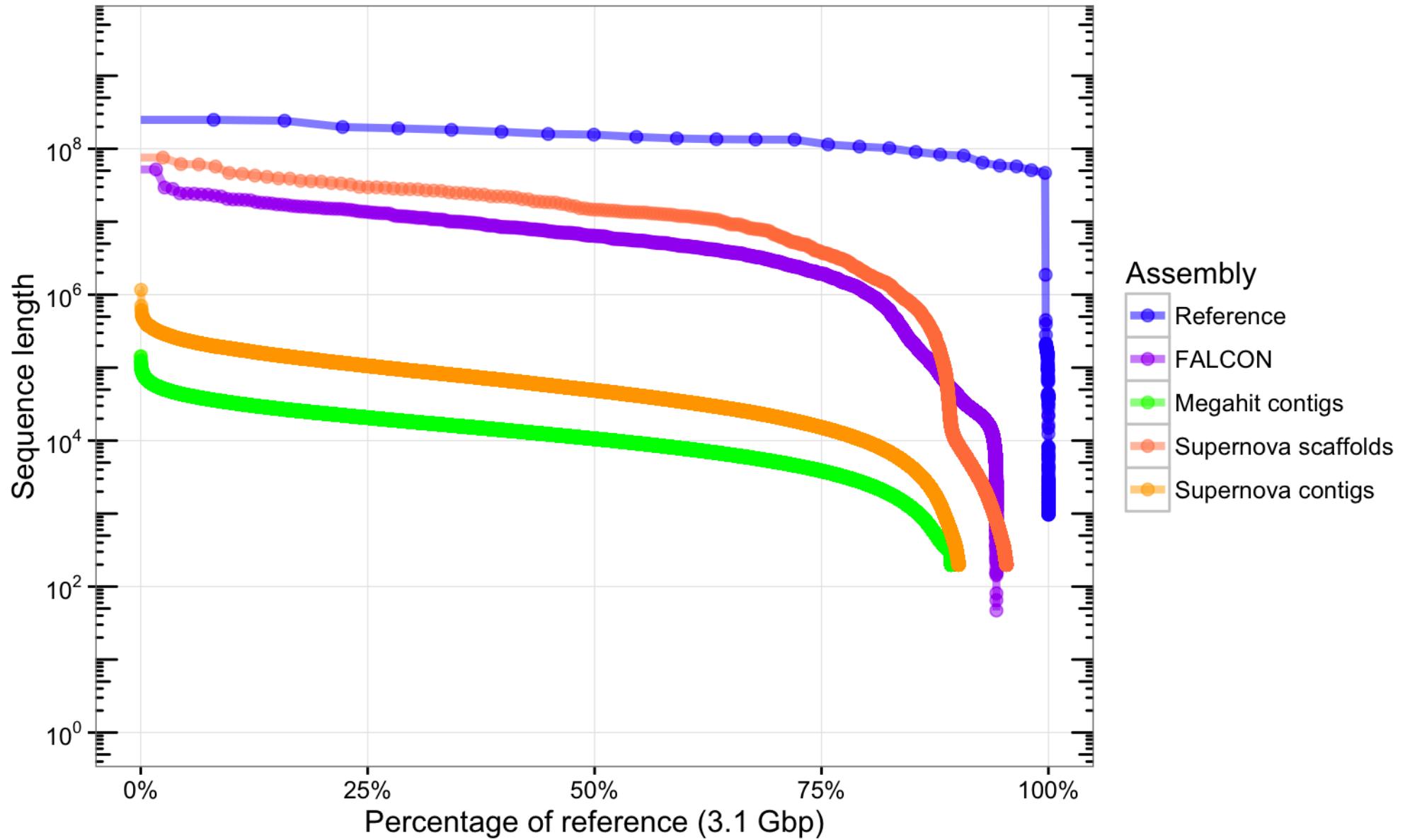
B



N50 size = 3 kbp

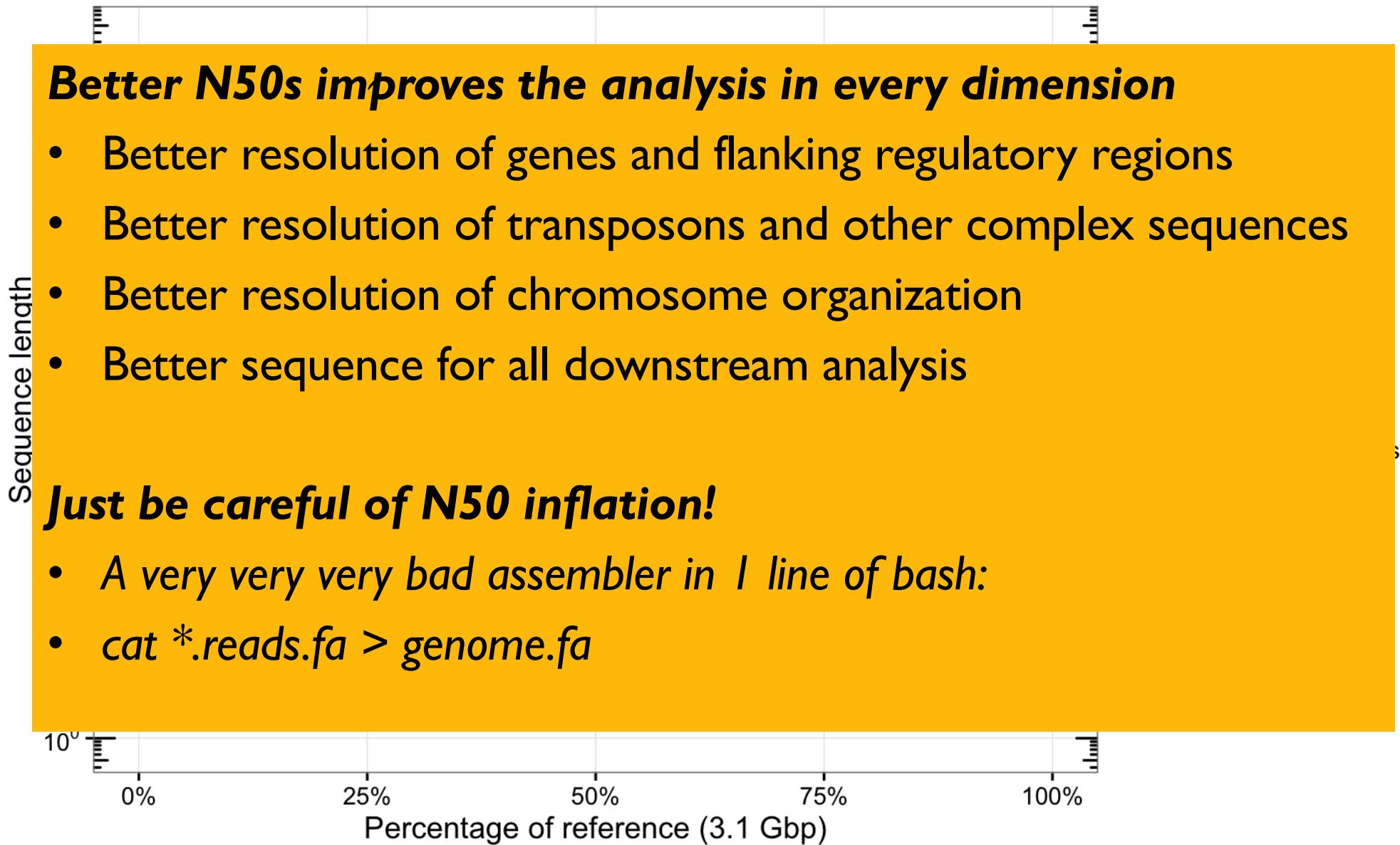
Contig Nchart

Cumulative sequence length



Contig Nchart

Cumulative sequence length



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

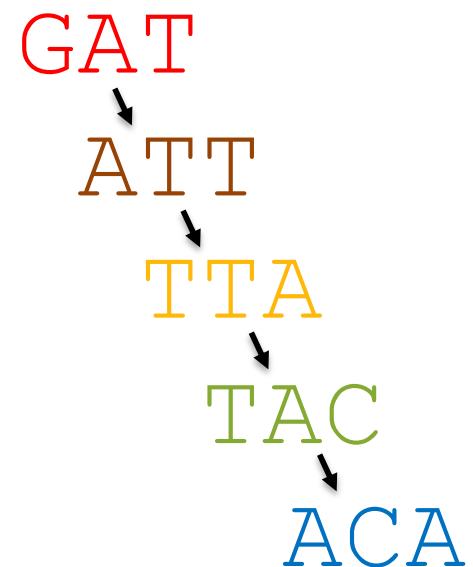
Assemble these reads using a de Bruijn graph approach (k=3):

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

ACG
 ↑
CGA

Pop Quiz 2

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

~~-ACGA~~

~~-ACGT~~

ATAC

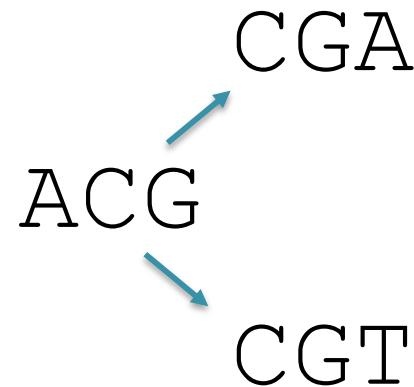
CGAC

CGTA

GACG

GTAT

TACG



Pop Quiz 2

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

~~ACGA~~

~~ACGT~~

ATAC

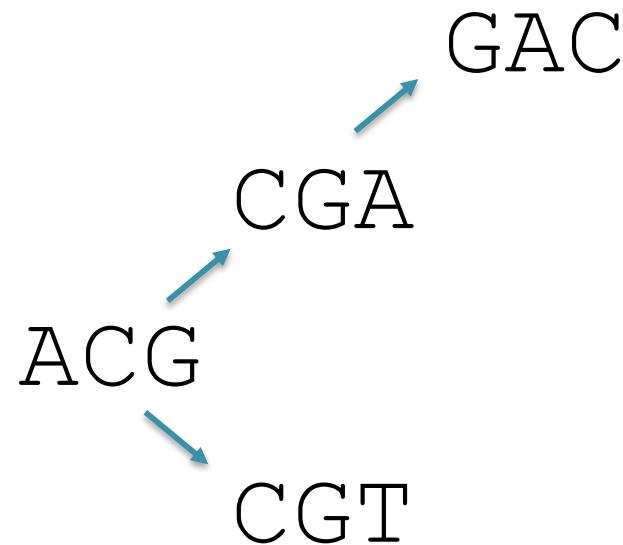
~~CGAC~~

CGTA

GACG

GTAT

TACG



Pop Quiz 2

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

~~ACGA~~

~~ACGT~~

ATAC

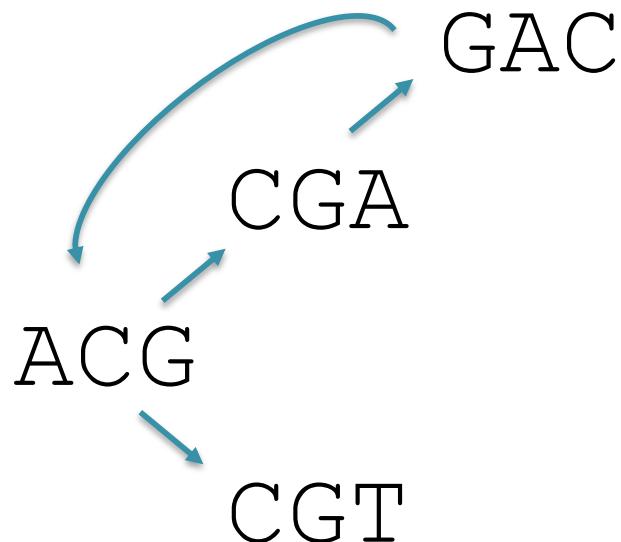
~~CGAC~~

CGTA

~~GACG~~

GTAT

TACG



Pop Quiz 2

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

~~-ACGA~~

~~-ACGT~~

ATAC

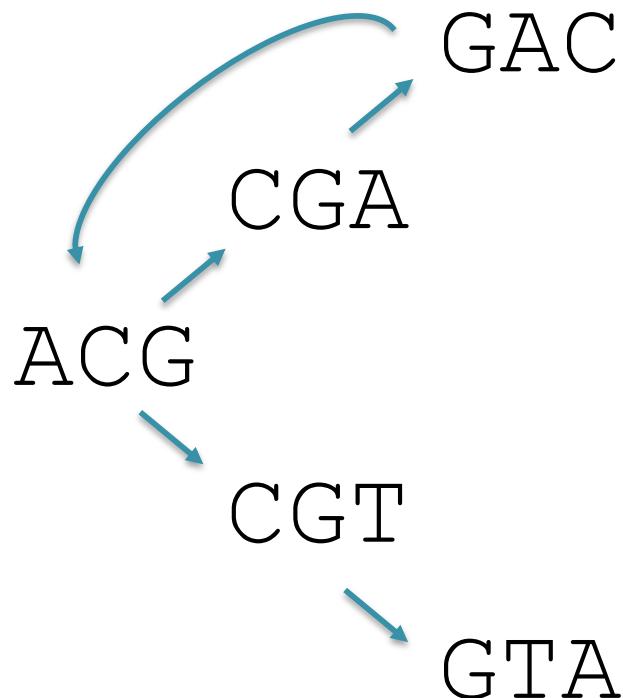
~~-CGAC~~

~~-CGTA~~

~~-GACG~~

GTAT

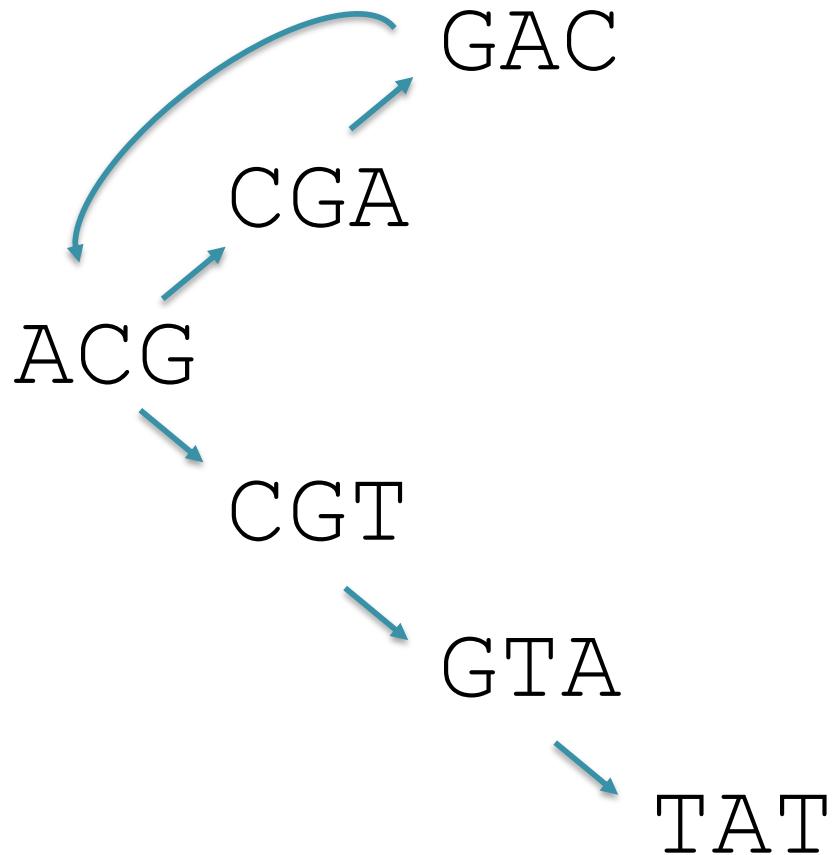
TACG



Pop Quiz 2

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

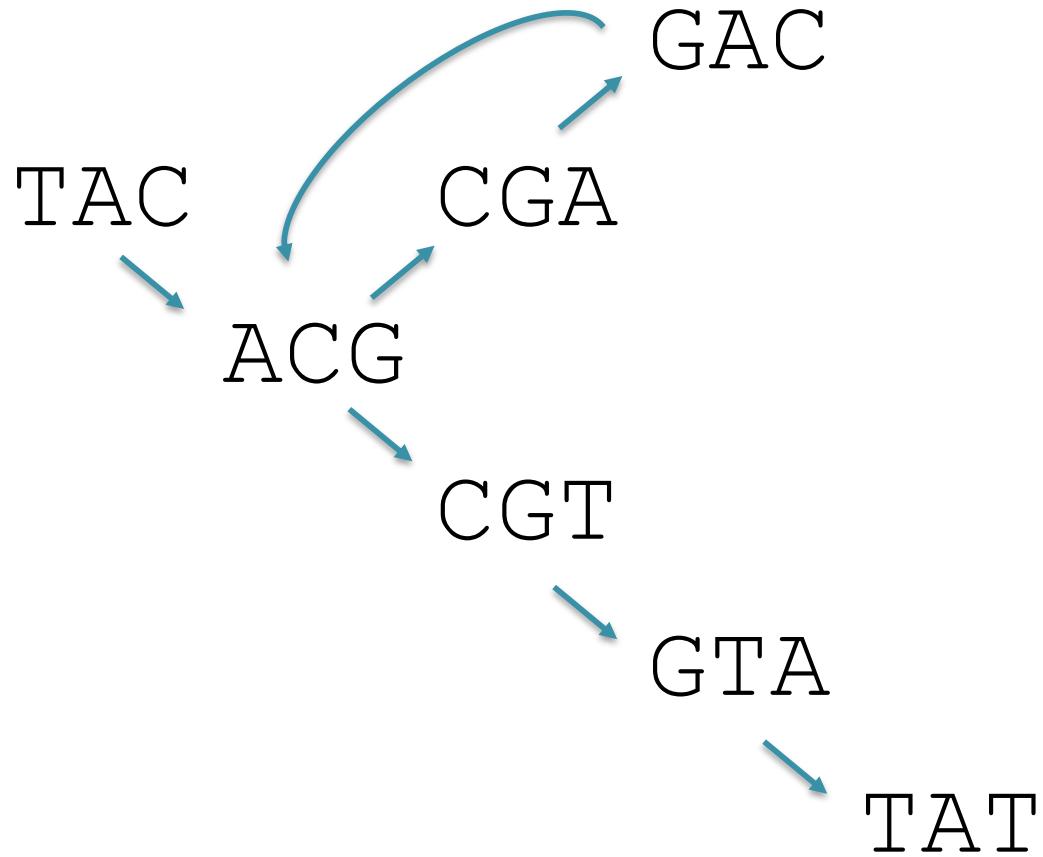
~~-ACGA~~
~~-ACGT~~
ATAC
~~-CGAC~~
~~-CGTA~~
~~-GACG~~
~~-GTAT~~
TACG



Pop Quiz 2

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

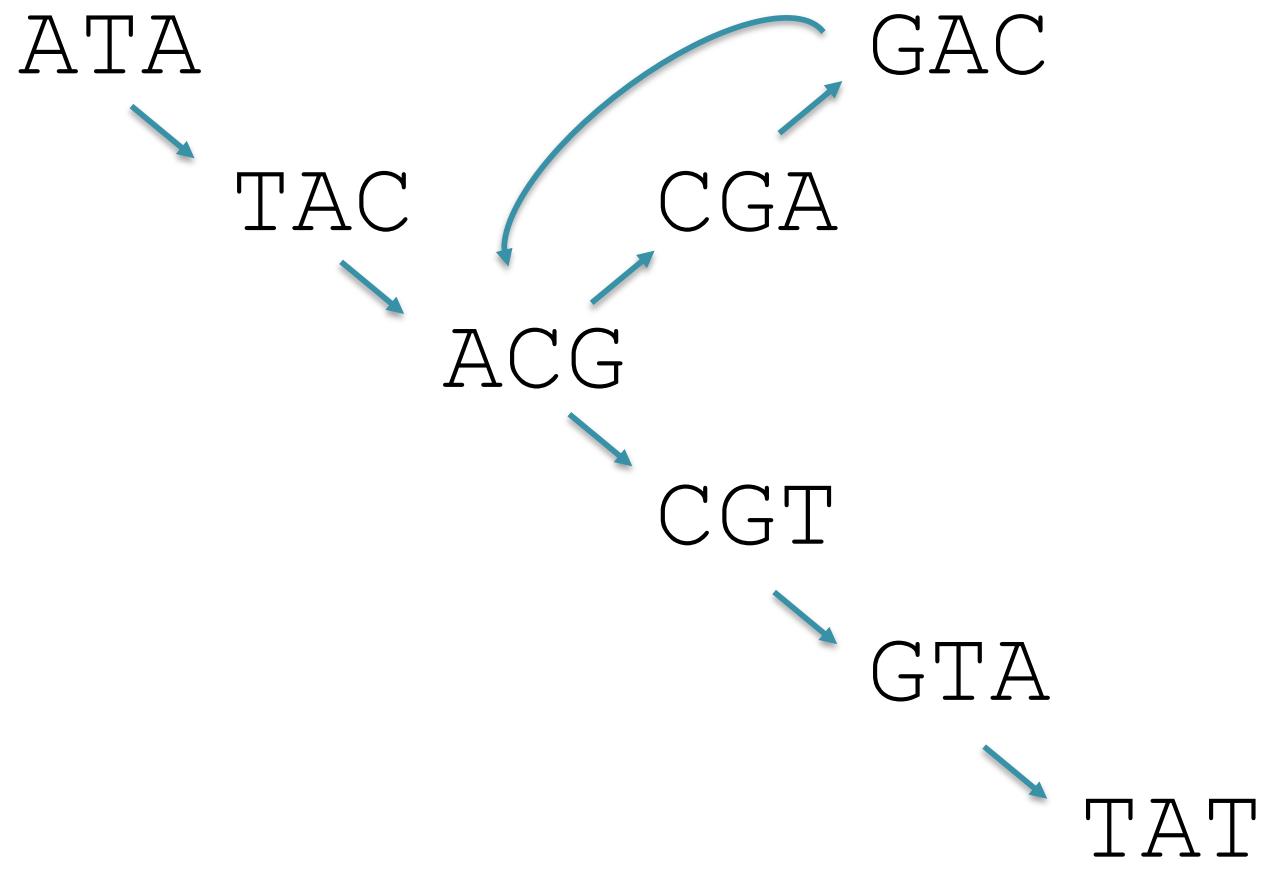
~~-ACGA~~
~~-ACGT~~
ATAC
~~-CGAC~~
~~-CGTA~~
~~-GACG~~
~~-GTAT~~
~~-TACG~~



Pop Quiz 2

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

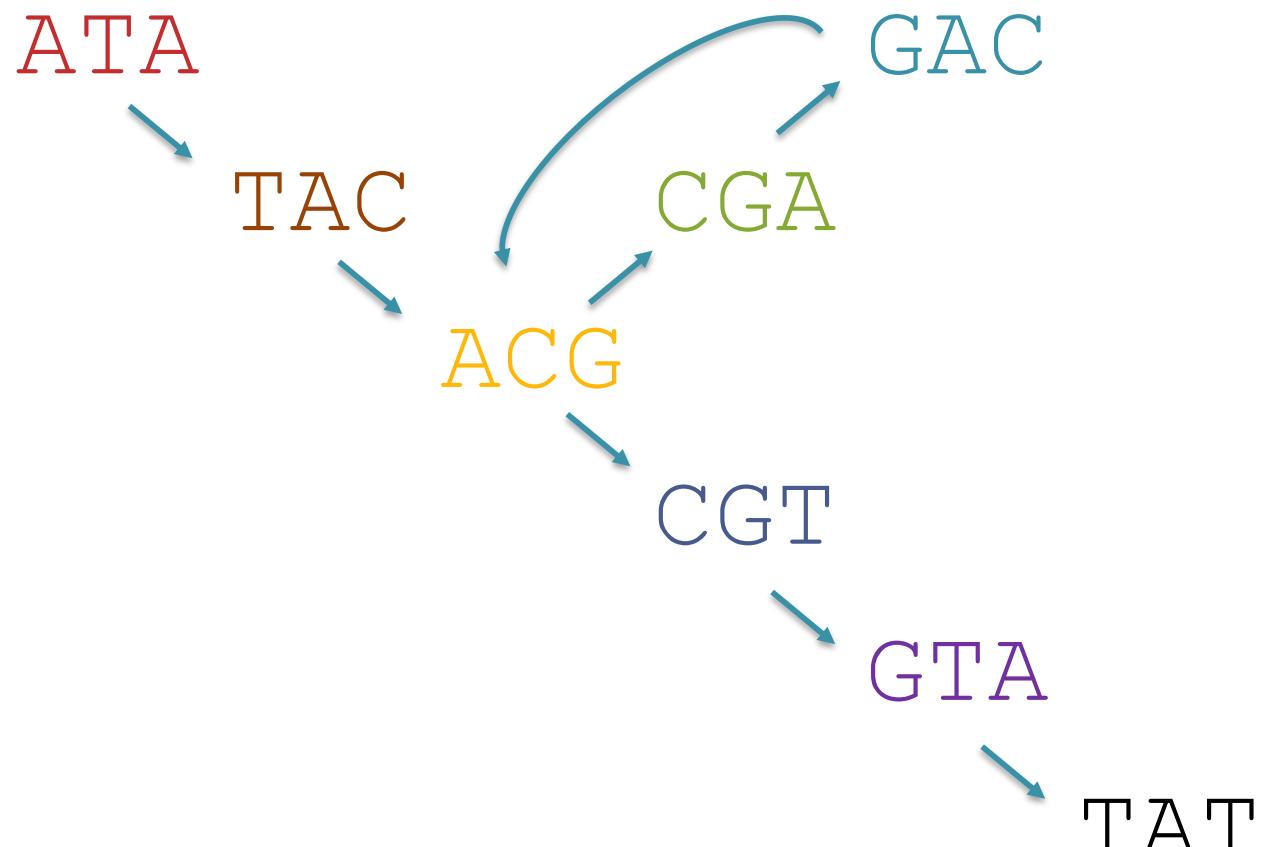
~~-ACGA~~
~~-ACGT~~
~~-ATAC~~
~~-CGAC~~
~~-CGTA~~
~~-GACG~~
~~-GTAT~~
~~-TACG~~



Pop Quiz 2

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

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

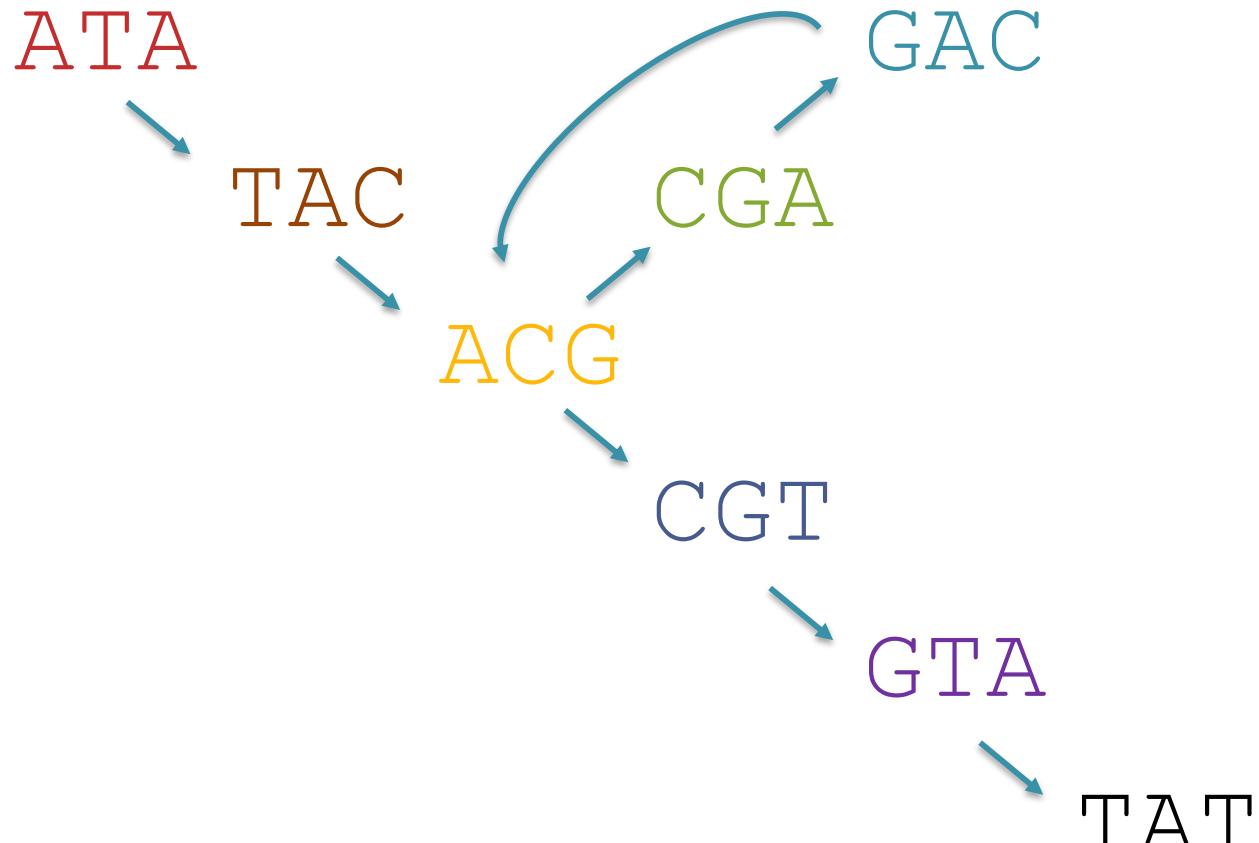


ATACGACGTAT

Pop Quiz 2

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

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



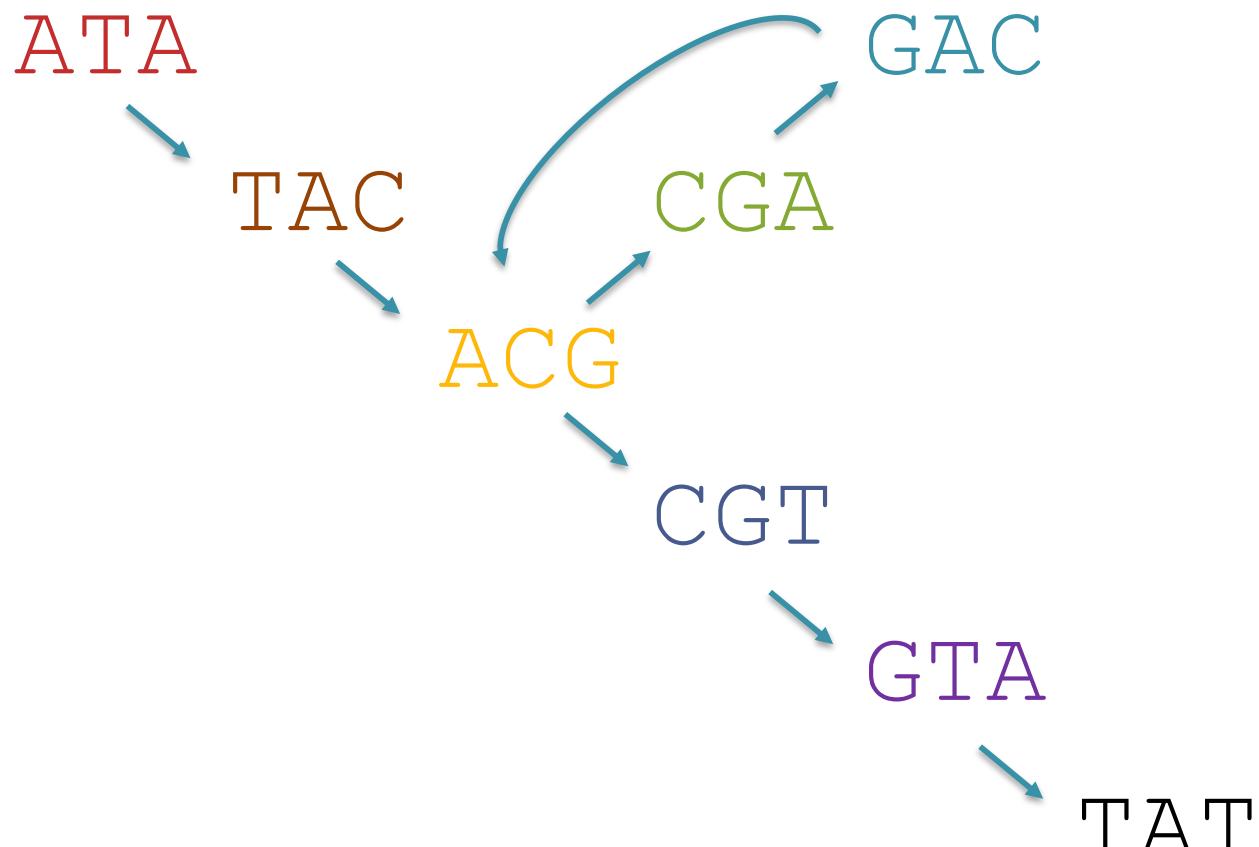
What's another possible genome?

ATACGACGTAT

Pop Quiz 2

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

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



Should we add the edge TAT -> ATA?

ATACGACGTAT