

Lecture 4: Genome assembly & annotation

- Genome assembly
 - de Bruijn graphs
- Genome annotation
 - Hidden Markov Models

Genome sequencing

Why sequence the genome?

- Determine the "complete" sequence of a haploid genome.
 - Previously "snippets" of the genome were available.
- Identify the sequence and location of every protein coding gene.
- Use as a "map" with which to track the location and frequency of genetic variation.
- Unravel the genetic architecture of inherited and somatic traits/diseases.
- To understand genome and species evolution.

"All the News
That's Fit to Print"

The New York Times

Late Edition

New York: Today, afternoon thunderstorms, high 88. Tonight, showers end, low 67. Tomorrow, partly cloudy with showers late, high 81. Yesterday, high 88, low 74. Weather map, Page D8.

VOL. CXLIX . . No. 51,422

Copyright © 2000 The New York Times

NEW YORK, TUESDAY, JUNE 27, 2000

\$1 beyond the greater New York metropolitan area.

75 CENTS

Genetic Code of Human Life Is Cracked by Scientists

JUSTICES REAFFIRM MIRANDA RULE, 7-2; A PART OF 'CULTURE'

By LINDA GREENHOUSE

WASHINGTON, June 26 — The Supreme Court reaffirmed the Miranda decision today by a 7-to-2 vote that erased a shadow over one of the most famous rulings of modern times and acknowledged that the Miranda warnings "have become part of our national culture."

The court said in an opinion by Chief Justice William H. Rehnquist that because the 1966 Miranda decision "announced a constitutional rule," a statute by which Congress had sought to overrule the decision was itself unconstitutional.

Miranda had appeared to be in jeopardy, both because of that long-ignored but recently rediscovered law, by which Congress had tried to overrule Miranda 32 years ago, and because of the court's perceived hostility to the original decision.

The chief justice said, though, that the 1968 law, which replaced the Miranda warnings with a case-by-case test of whether a confession was voluntary, could be upheld only if the Supreme Court decided to overturn Miranda. But with Miranda having

Justices Antonin Scalia and Clarence Thomas cast the dissenting votes.

The decision overturned a ruling last year by the federal appeals court in Richmond, Va., which held that Congress was entitled to the last word because Miranda's presumption that a confession was not voluntary unless preceded by the warnings was not required by the Constitution.

The decision today — only 14 pages long, in Chief Justice Rehnquist's typically spare style — brought an abrupt end to one of the odder episodes in the court's recent history, an intense and strangely delayed re-fighting of a previous generation's battle over the rights of criminal suspects. Miranda v. Arizona was a landmark of the Warren Court, and Chief Justice Rehnquist, despite his record as an early and tenacious critic of the decision, evidently did not want its repudiation to be an imprint of his own tenure.

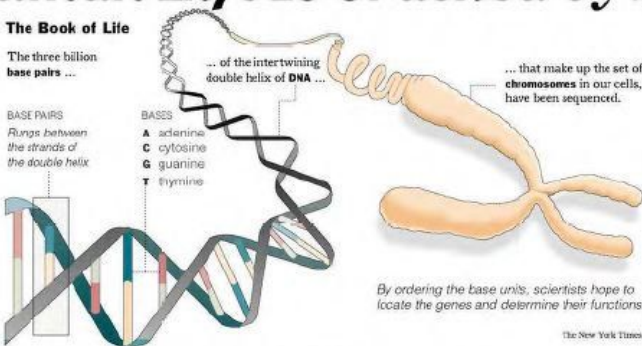
There was considerable drama in the courtroom today as the chief justice announced that he would die.

The Book of Life

The three billion base pairs ...

BASE PAIRS
Rungs between the strands of the double helix

BASES
A adenine
C cytosine
G guanine
T thymine



By ordering the base units, scientists hope to locate the genes and determine their functions.

The New York Times

Science Times A special issue

- Putting the genome to work.
- Some information has already paid research dividends.
- Two research methods, two results.
- From Mendel to helix to genome.
- More articles, charts and photos of the genome effort.

Section F

Francis S. Collins, head of the Human Genome Project, left, with J. Craig Venter, head of Celera Genomics, after the announcement yesterday that they had finished the first survey of the human genome.



Paul Heston/The New York Times

A SHARED SUCCESS

2 Rivals' Announcement Marks New Medical Era, Risks and All

By NICHOLAS WADE

WASHINGTON, June 26 — In an achievement that represents a pinnacle of human self-knowledge, two rival groups of scientists said today that they had deciphered the hereditary script, the set of instructions that defines the human organism.

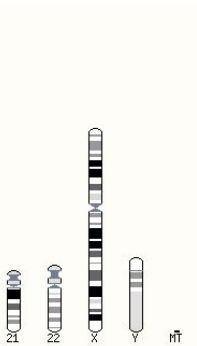
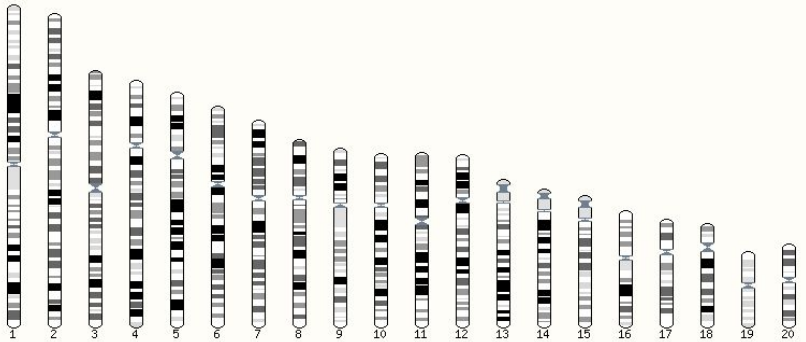
"Today we are learning the language in which God created life," President Clinton said at a White House ceremony attended by members of the two teams, Dr. James D. Watson, codiscoverer of the structure of DNA, and, via satellite, Prime Minister Tony Blair of Britain. [Excerpts, Page D8.]

The teams' leaders, Dr. J. Craig Venter, president of Celera Genomics, and Dr. Francis S. Collins, director of the National Human Genome Research Institute, praised each other's contributions and signaled a spirit of cooperation from now on, even though the two efforts will remain firmly independent.

The human genome, the ancient script that has now been deciphered, consists of two sets of 23 giant DNA

Genome sequencing

The Human Genome – Summary



http://useast.ensembl.org/Homo_sapiens/Location/Genome

Assembly	GRCh38.p12 (Genome Refere
Base Pairs	3,609,003,417
Golden Path Length	3,096,649,726
Annotation provider	Ensembl
Annotation method	Full genebuild
Genebuild started	Jan 2014
Genebuild released	Jul 2014
Genebuild last updated/patched	Jul 2018
Database version	95.38
Gencode version	GENCODE 29

Gene counts (Primary assembly)

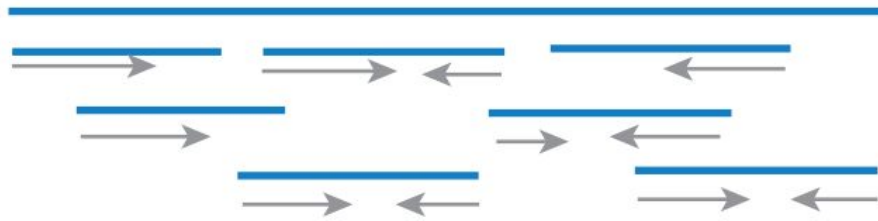
Coding genes	20,418 (incl 650 readthrough)
Non coding genes	22,107
Small non coding genes	4,871
Long non coding genes	15,014 (incl 284 readthrough)
Misc non coding genes	2,222
Pseudogenes	15,195 (incl 8 readthrough)
Gene transcripts	206,762

Timeline of the Human Genome Project (1865-2003):

- 1865:** Mendel discovers laws of genetics
- 1900:** Rediscovery of Mendel's work
- 1905:** Garrod formulates the concept of human inborn errors of metabolism
- 1913:** Shattworth makes the first linear map of a gene
- 1944:** Avery, McClelland, and McCarty demonstrate DNA is the hereditary material
- 1953:** Watson and Crick describe the double helical structure of DNA
- 1966:** Nirenberg, Khorana, and Holley determine the genetic code
- 1972:** Cohen and Boyer develop recombinant DNA technology
- 1974:** Issuing of Belmont Report on the use of human subjects in research
- 1977:** Sanger and Maxam & Gilbert develop DNA sequencing methods
- 1982:** GenBank database established
- 1983:** First human disease gene mapped (Huntington disease)
- 1984:** First public discussion of sequencing the human genome
- 1985:** PCR invented
- 1986:** First automated DNA sequencing instrument developed
- 1987:** First-generation human genetic map developed
- 1988:** Human Genome Organization (HGO) formed
- 1990:** First draft of the human genome published
- 2003:** Completion of the Human Genome Project



Genome assembly & annotation – Overview



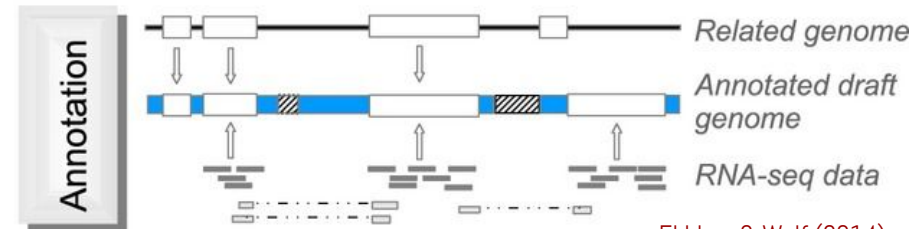
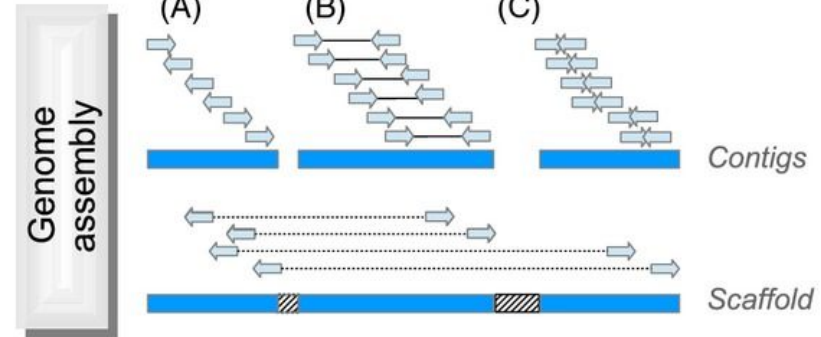
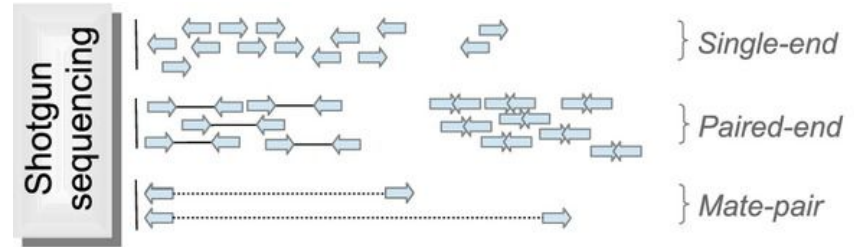
read: a short/long word that comes out of sequencer

mate pair: a pair of reads from two ends of the same insert fragment

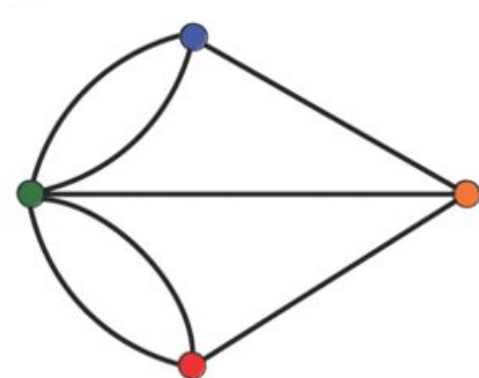
contig: a contiguous sequence formed by several overlapping reads with no gaps

scaffold: an ordered and oriented set of contigs, usually by mate pairs

consensus sequence: derived from the sequence multiple alignment of reads in a contig



Introduction to graph theory



Bridges of Königsberg

Can every part of the city be visited by walking across each of the seven bridges exactly once and returning to one's starting location?



Graph: (Nodes, Edges, Weights)

(Eulerian) Path exists if the graph contains zero or two vertices that have an odd degree.

de Bruijn graphs: the 'superstring problem'

Find a shortest circular 'superstring' that contains all possible 'substrings' of length k (k -mers) over a given alphabet.

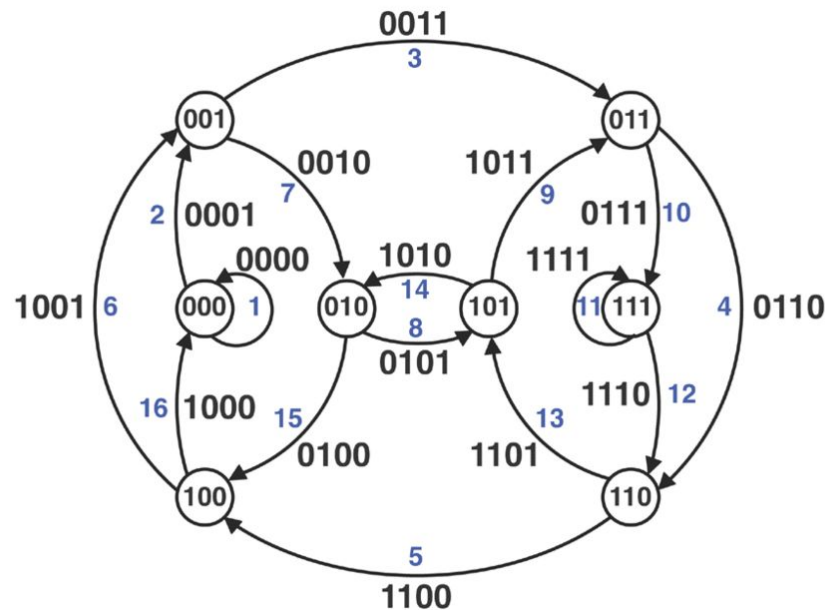
There exist n^k k -mers in an alphabet containing n symbols:

- alphabet: 0 & 1
- all 3-mers: 000, 001, 010, 011, 100, 101, 110, 111.
- The circular superstring **0001110100** contains all 3-mers & each 3-mer exactly once.



How can we construct a superstring for all k -mers in the case of an arbitrary value of k and an arbitrary alphabet?

Construct a **directed graph**: all prefixes & suffixes as nodes; all k -mers as edges.

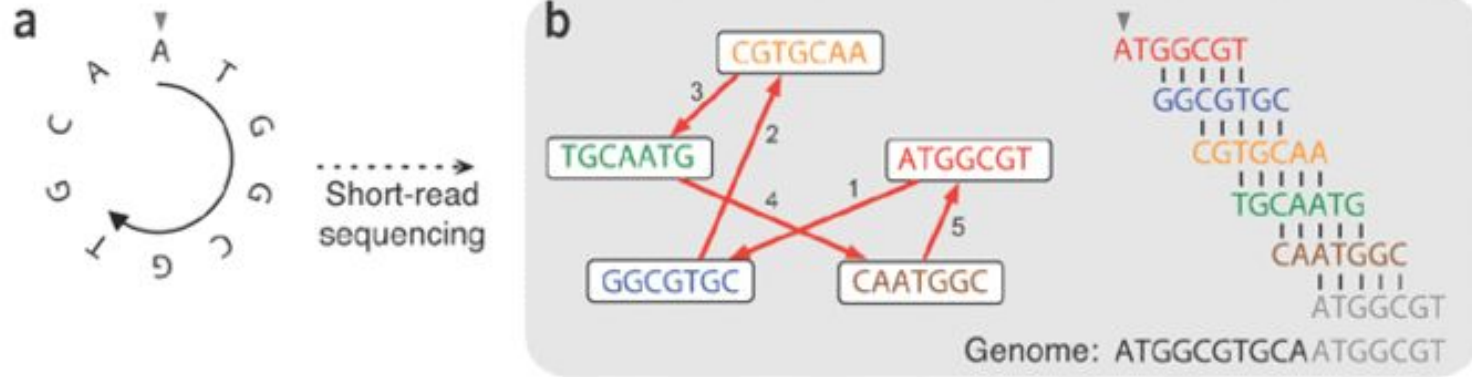


$k = 4$ | Two-character alphabet: digits 0 & 1

Does this graph have an Eulerian cycle? [Balanced?]

Following the blue numbered edges in order from 1 to 16 traces the cyclic superstring 0000110010111101.

de Bruijn graphs for sequence assembly



Nodes: Reads

Edges: Alignments between reads (≥ 5 bases)

Genome: Walk along a Hamiltonian cycle (visit each vertex exactly once) to combine alignments between successive reads and reconstruct the full circular genome.

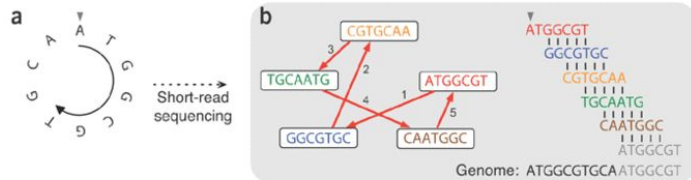
de Bruijn graphs for sequence assembly

Four hidden assumptions that **do not hold** for real sequencing:

1. We can generate all k-mers present in the genome.
2. All k-mers are error free.
3. Each k-mer appears at most once in the genome.
4. The genome consists of a single circular chromosome.

E.g., a technology that generates 100-nucleotide long reads:

- may miss some 100-mers present in the genome (even if the read coverage is high)
- the 100-mers that it does generate typically have errors.



de Bruijn graphs for sequence assembly

Break reads into shorter k-mers!

Resulting k-mers often represent nearly all k-mers from the genome for sufficiently small k.

Nodes: k-mers

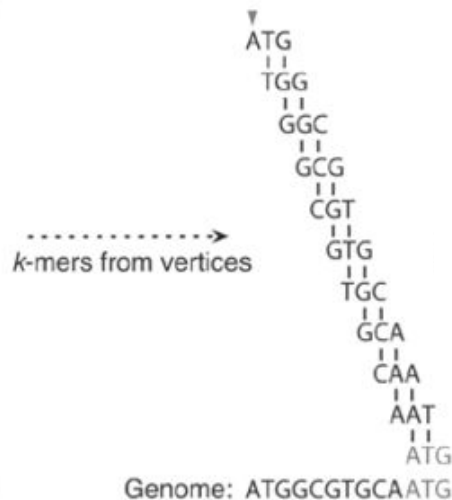
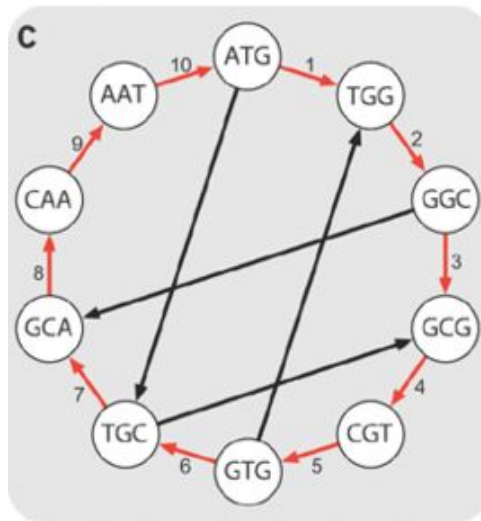
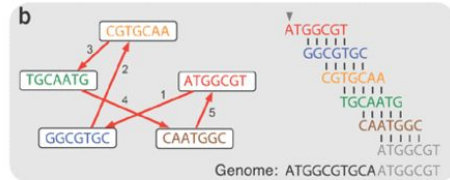
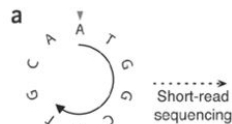
Edges: Overlap between k-mers

Genome: Walk the Hamiltonian cycle

Large genomes result in too many reads:

- 10^6 reads $\rightarrow 10^{12}$ pairwise alignments.
- 10^9 reads $\rightarrow 10^{18}$ alignments.

There is no known efficient algorithm for finding a Hamiltonian cycle in a large graph with millions (let alone billions) of nodes.



de Bruijn graphs for sequence assembly

Finding a path that visits all *edges* of a graph exactly once is much easier!

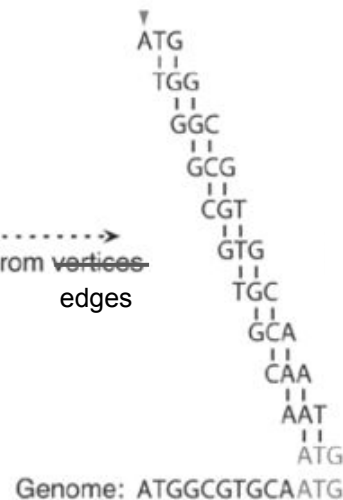
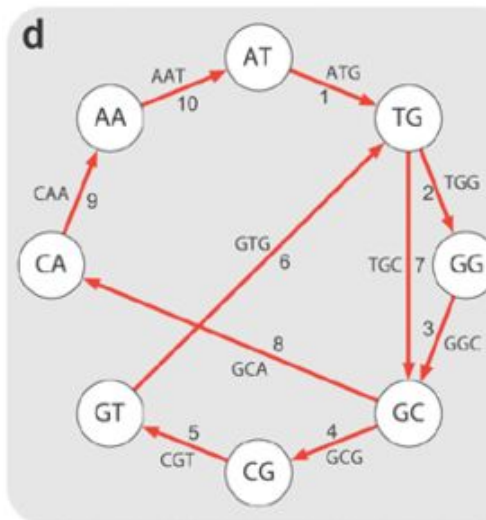
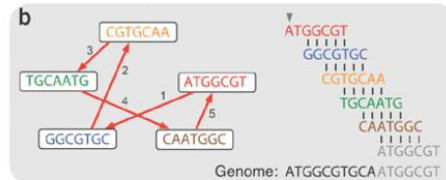
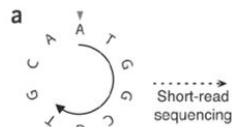
Nodes: (k-1)-mers [in- & out-degrees?]

Edges: k-mers

Genome: Walking the Eulerian cycle

Send out an ant, find a cycle; Note edges traversed, send out another ant, ... until all of the graph's edges have been explored. Combine all cycles to form an Eulerian cycle!

Computationally tractable; time roughly proportional to the number of edges.

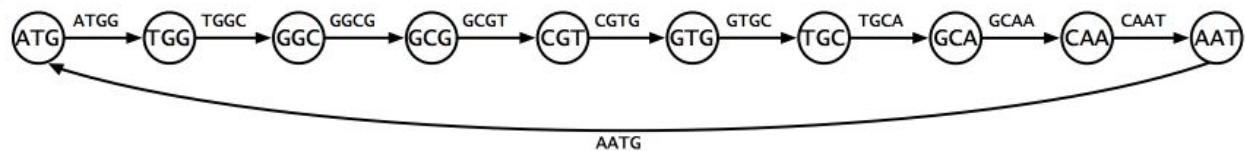


de Bruijn graphs for sequence assembly – not so easy w/ real data

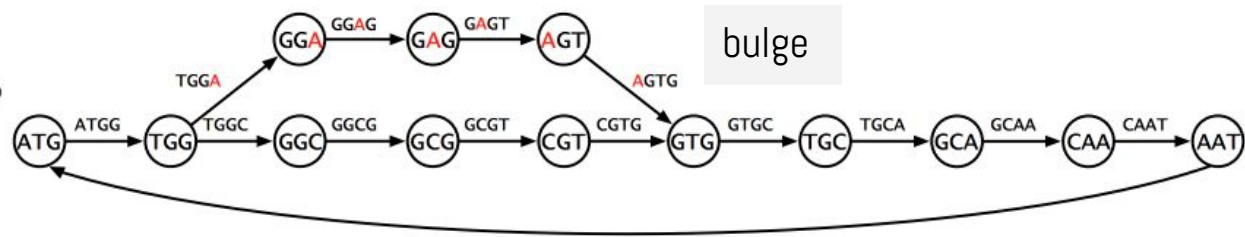
- Not all k-mers in the genome
 - Read breaking
- Errors in reads
 - Error correcting reads
 - Removing bulges in de Bruijn graphs
- DNA repeats
 - Incorporating k-mer multiplicity
- Multiple and linear chromosomes
 - Cycles to paths
- Unsequenced regions
 - Scaffolds

de Bruijn graphs from reads with sequencing errors

a



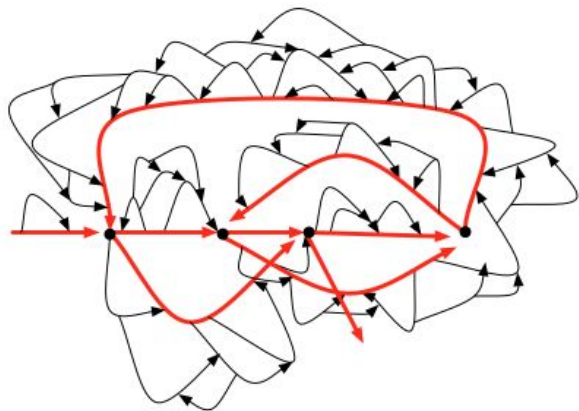
b



TGGAGTG (incorrect)

TGGCGTG (correct)

c

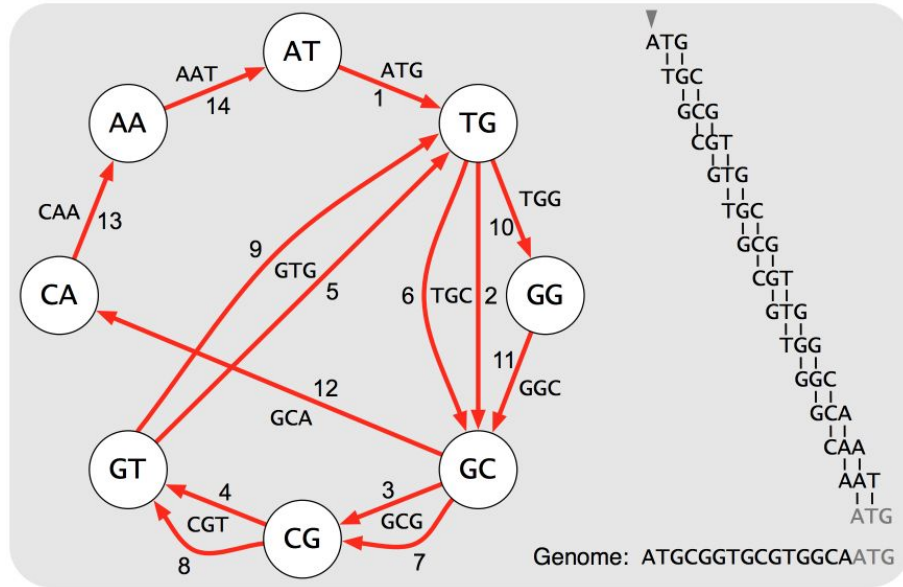


The process of bulge removal should leave only the red edges remaining, yielding an Eulerian path in the resulting graph.

de Bruijn graphs for sequence assembly – not so easy w/ real data

- Not all k-mers in the genome
 - Read breaking
- Errors in reads
 - Error correcting reads
 - Removing bulges in de Bruijn graphs
- DNA repeats [E.g., **ATGCATGC** → four 3-mers: **ATG**, **TGC**, **GCA** & **CAT**]
 - Incorporating k-mer multiplicity
- Multiple and linear chromosomes
 - Cycles to paths
- Unsequenced regions
 - Scaffolds

de Bruijn graphs for dealing with repeats



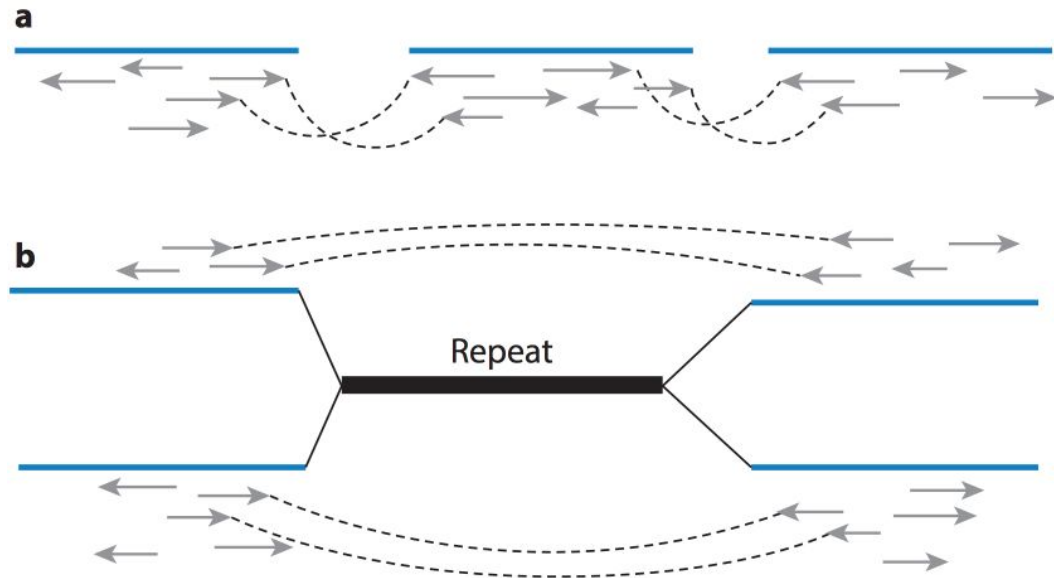
Genome: **ATGCGTGCGTGGCA**

Incorporate k-mer multiplicity:

- Four 3-mers **TGC**, **GCG**, **CGT**, and **GTG**: multiplicity 2
- Six 3-mers **ATG**, **TGG**, **GGC**, **GCA**, **CAA**, and **AAT**: multiplicity 1

de Bruijn graphs for sequence assembly – not so easy w/ real data

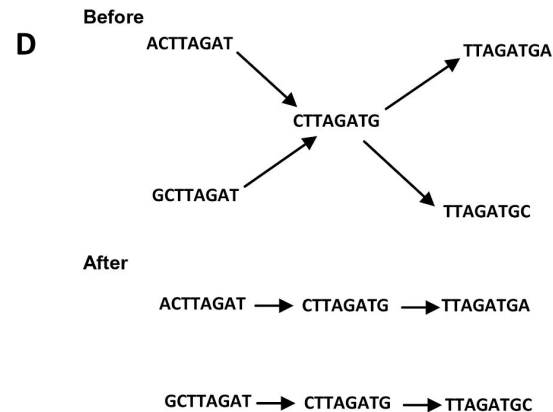
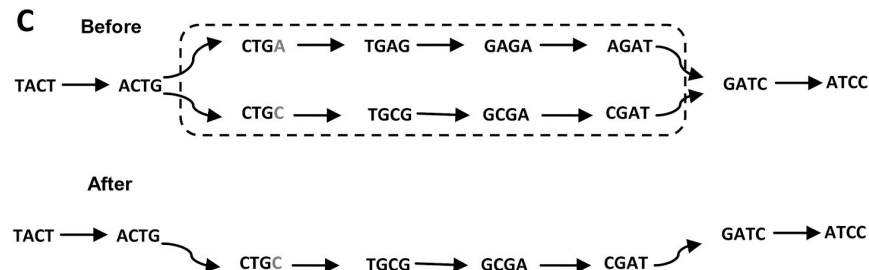
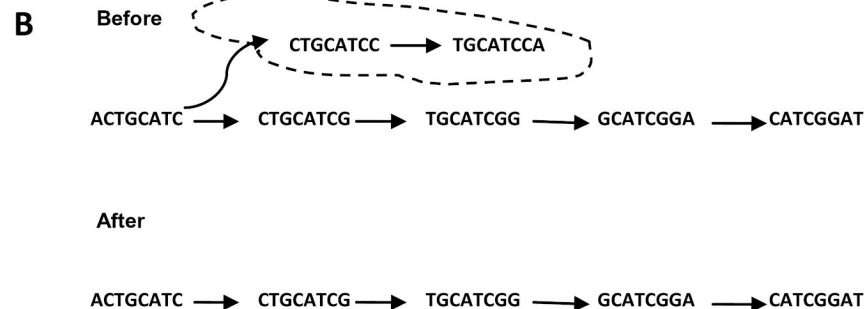
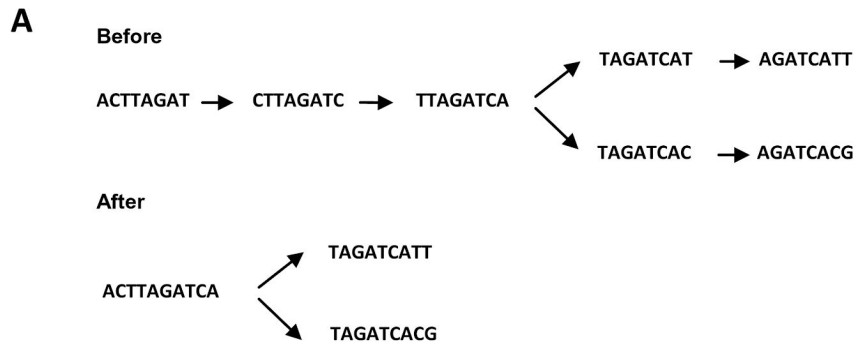
- Not all k-mers in the genome
 - Read breaking
- Errors in reads
 - Error correcting reads
 - Removing bulges in de Bruijn graphs
- DNA repeats
 - Incorporating k-mer multiplicity
- Multiple and linear chromosomes
 - Cycles to paths
- Unsequenced regions
 - Scaffolds



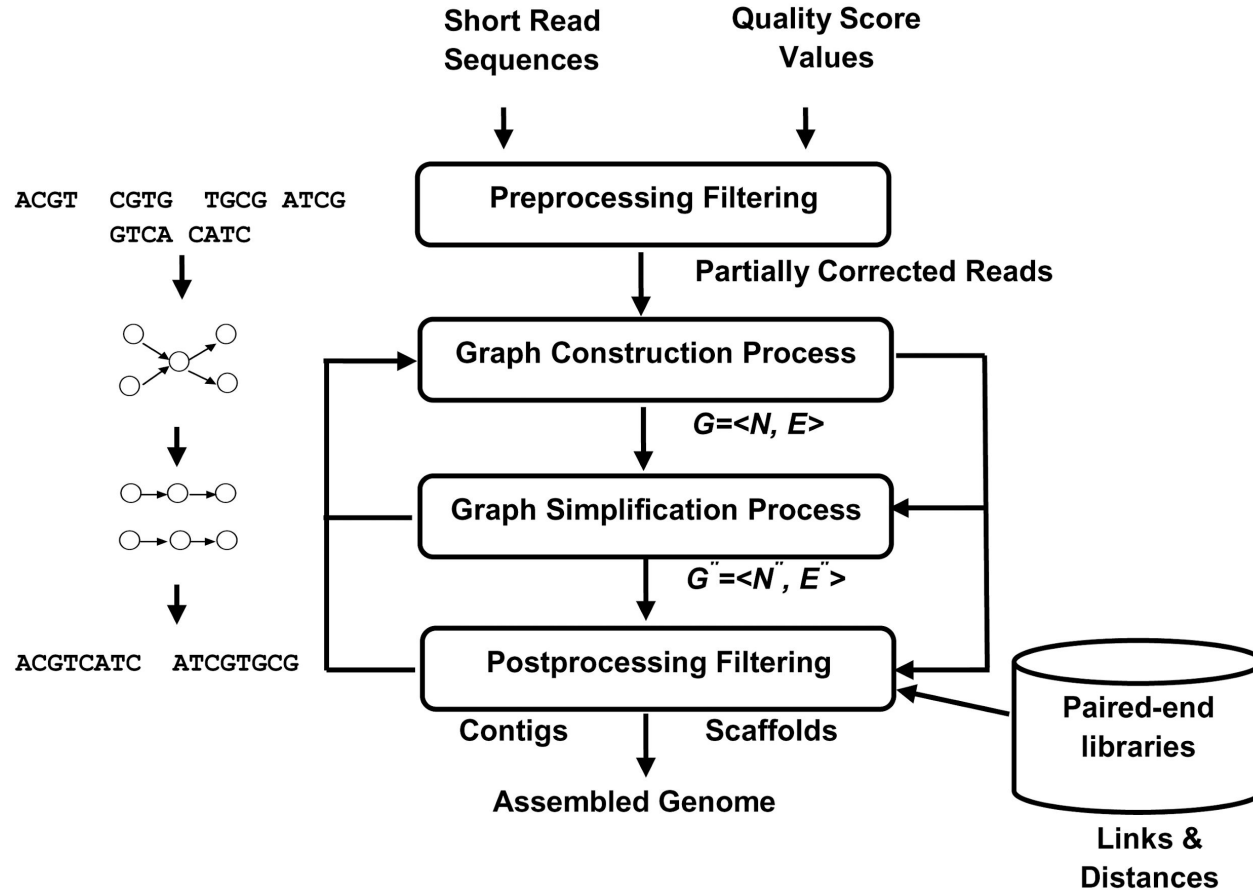
Algorithms for genome assembly

- Error detection and correction based on sequence composition of the reads.
- Graph construction to represent reads/k-mers and their shared sequence.
- Graph adjustments:
 - Reduction of simple non-intersecting paths to single nodes.
 - Removal of error-induced paths (recognized as spurs or bubbles).
 - Collapse of polymorphism-induced complexity (bubbles).
 - Simplification of tangles (using information outside the graph: individual, paired-end, or mate-pair reads to constraints on path distance & outcome).
- Conversion of reduced paths to contigs and scaffolds.
- Reduction of alignments to a consensus sequence.

Simplifying de Bruijn graphs



Algorithms for genome assembly



Algorithms for genome assembly

- New sequencing technologies → a different best computational strategy.
- Factors that influence the choice of algorithms:
 - Quantity of data (read length and coverage)
 - Quality of data (including error rates)
 - Genome structure (e.g., GC content and the number and size of repeated regions).
- de Bruijn graphs are best suited for current short-read sequencing technologies
 - Produce very large numbers of reads
 - Can represent genomes with repeats [Overlap methods need to mask repeats $>$ read length]
- Long-read technology growing at a rapid pace.

Some modern genome assemblers

Assemblers	Technology	Availability	Notes
<i>Genome assemblers</i>			
ALLPATHS-LG	Illumina, Pacific Biosciences	ftp://ftp.broadinstitute.org/pub/crd/ALLPATHS/Release-LG	Requires a specific sequencing recipe (BOX 3)
SOAPdenovo	Illumina	http://soap.genomics.org.cn/soapdenovo.html	Also used for transcriptome and metagenome assembly
Velvet	Illumina, SOLiD, 454, Sanger	http://www.ebi.ac.uk/~zerbino/velvet	May have substantial memory requirements for large genomes
ABYSS	Illumina, SOLiD, 454, Sanger	http://www.bcgsc.ca/platform/bioinfo/software/abyss	Also used for transcriptome assembly

Paper discussion

- Ask questions & offer answers/thoughts.
 - Let's collectively engage.
 - Good for you & the presenters.
- Also be ready with all the questions you had from your reading.