

Lecture 14: DNA Assembly

Study Chapter 8.9
Midterm on Wednesday 3/4
Open book, open notes, no computer
Study Session on 3/5 in SN014 from 5pm-7pm

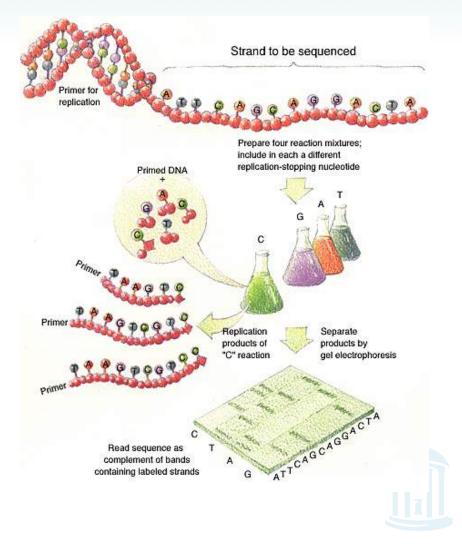
3/2/15

Comp 555

Spring 2015

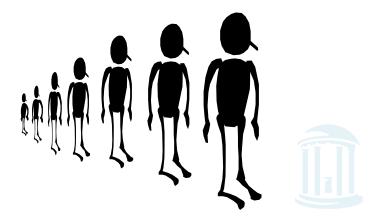
DNA Sequencing

- Shear DNA into millions of small fragments
- Read 500 700
 nucleotides at a time
 from the small
 fragments
 (Sanger method)



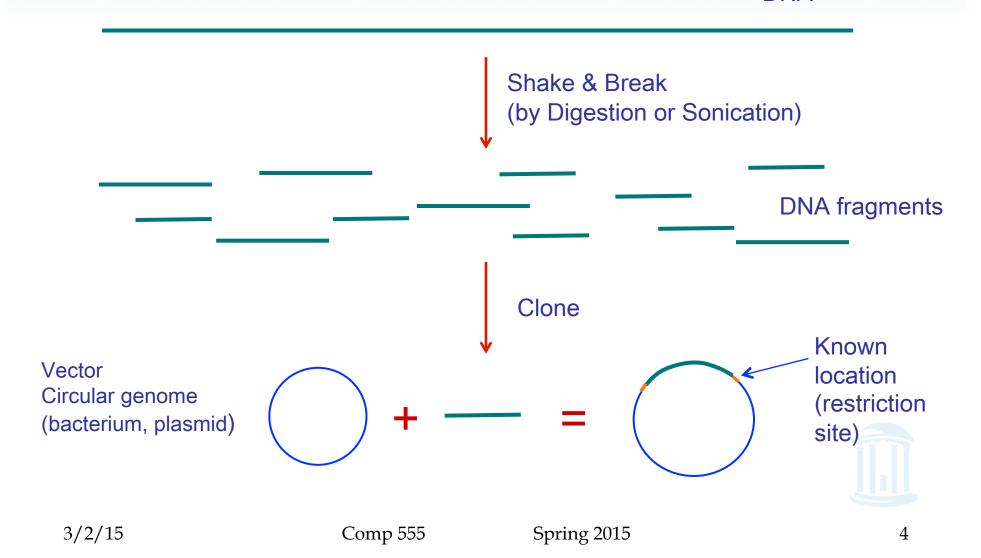
Fragment Assembly

- Assembles the individual overlapping short fragments (reads) into a genomic sequence
- Shortest Superstring problem from last time is an overly simplified abstraction
- Problems:
 - DNA read error rate of 1% to 3%
 - Can't separate strands
 - DNA is **full** of repeats
- Let's take a closer look



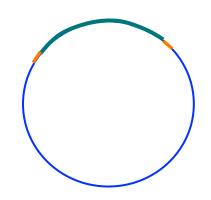
Traditional DNA Sequencing

PNA DNA



Different Types of Vectors

<u>VECTOR</u>	Size of insert (bp)	
Plasmid	2,000 - 10,000	
Cosmid	40,000	
BAC (Bacterial Artificial Chromosome)	70,000 - 300,000	
YAC (Yeast Artificial Chromosome)	> 300,000 Not used much recently	





Dideoxy (Sanger) Sequencing

Template strand - g t a a g a c t g t
Coding strand - c a t t c t g a c a

ddT Reaction - c a t t
c a t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t
c a t t

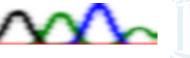
ddC Reaction -

ddG Reaction -

Good for up to
1000 base pairs

```
c a t t c c a t t c t g a c
```

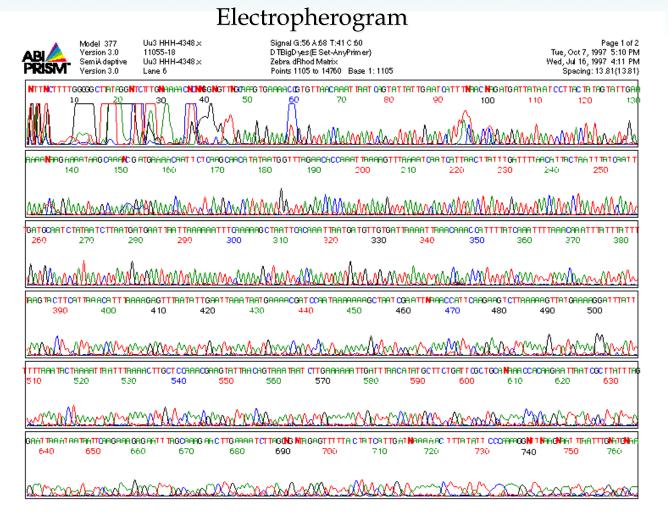
```
ca
cattctga
cattctgaca
```





Challenging to Read Answers







Reading an Electropherogram



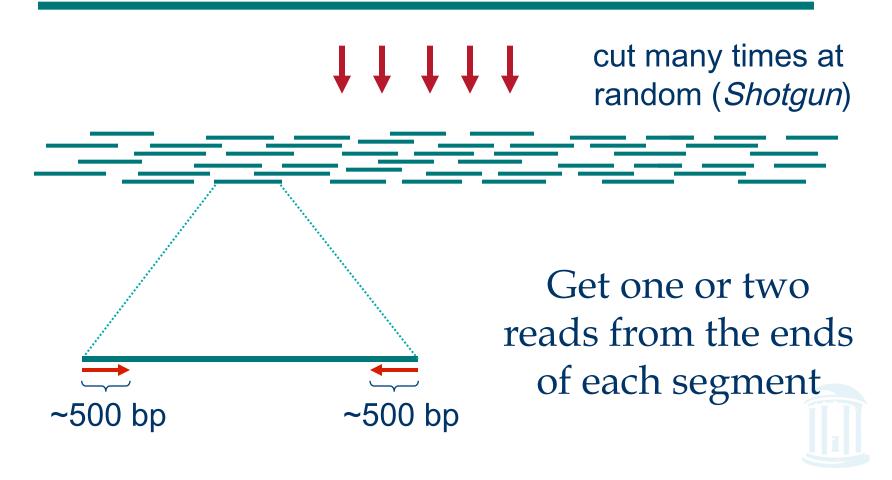
Issues

- Noisy start up due to anomalous migration of short fragments that carry bulky dyes
- Traces become less uniform as run proceeds
- Large dye responses can overwhelm succeeding lower amplitude responses
- Occasional mismatches of reaction with template
- Methods for calling the nucleotides: PHRED
 - Base calls
 - Maintains quality scores
 - Monitors peak positions



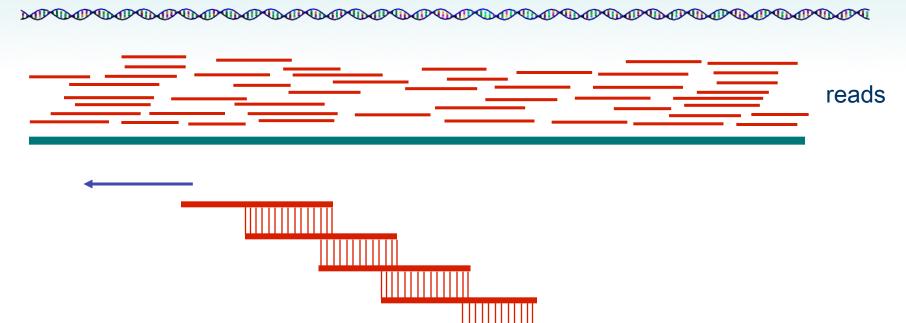
Shotgun Sequencing

genomic segment



3/2/15

Fragment Assembly



Cover region with ~7-fold redundancy

Overlap reads and extend to reconstruct the original genomic region

Read Coverage





Length of genomic segment: *L*

Number of reads: n Coverage C = n l/L

Length of each read:

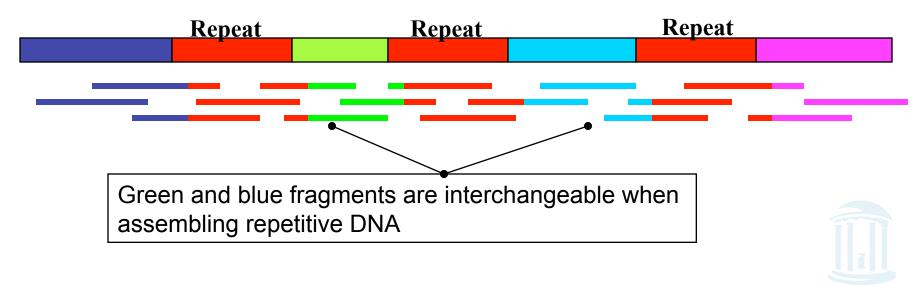
How much coverage is enough?

Lander-Waterman model:

Assuming uniform distribution of reads, *C*=10 results in 1 gapped region per 1,000,000 nucleotides

Challenges in Fragment Assembly

- Repeats: A major problem for fragment assembly
- > 50% of human genome is repeats:
 - over 1 million *Alu* repeats (about 300 bp)
 - about 200,000 LINE repeats (1000 bp and longer)



Types of Genome Assemblies

- De Novo –
 An assembly based entirely on self-consitency or self-similarity of short reads (contigs).
- Comparative –
 Refers an assembly of a genome using the sequence of a close relative as a scaffold or reference. Sometimes called a "template assembly" or "a resequencing"
- Confounding problem for both types: Repeats

Repeat Types

- Low-Complexity DNA
- (e.g. ATATATATACATA...)
- Microsatellite repeats $(a_1...a_k)^N$ where $k \sim 3-6$ (e.g. CAGCAGTAGCAGCACCAG)
- Transposons/retrotransposons
 - SINE

Short Interspersed Nuclear Elements (e.g., *Alu*: ~300 bp long, >10⁶ in human)

- LINE

- Long Interspersed Nuclear Elements ~500 5,000 bp long, > 200,000 in human
- LTR retroposons
- Long Terminal Repeats (~700 bp) at each end

Gene Families

- genes duplicate & then diverge
- Segmental duplications
- ~very long, very similar copies



Overlap-Layout-Consensus

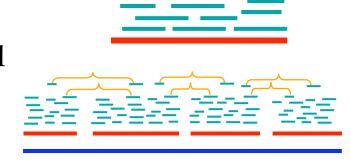
Assembler programs ARACHNE, PHRAP, CAP, TIGR, CELERA

Common Approach:

Overlap: find potentially overlapping reads



Layout: merge reads into contigs and then combine contigs into supercontigs



Consensus: requires many overlapping reads to derive the DNA sequence and to correct for read errors

..ACGATTACAATAGGTT..

Overlap

- Find the best match between the suffix of one read and the prefix of another (shortest superstring)
- Due to sequencing errors, most algorithms use dynamic programming to find the optimal overlap alignment
- Filter out fragment pairs that do not share a significantly long common substring



Overlapping Reads

- Make an index of all k-mers of all reads $(k \sim 21-27)$
- Find read-pairs sharing a k-mer
- Extend alignment –
 throw away if not >95% similar



Histogram Example

v = tagattacacagattattga

Histogram of 3-mers (18 total)

	A_2	C_2	G_2	T_2
	A ₃ :C ₃ :G ₃ :T ₃	A ₃ :C ₃ :G ₃ :T ₃	A ₃ :C ₃ :G ₃ :T ₃	A ₃ :C ₃ :G ₃ :T ₃
A_1	0:0:0:0	2:0:0:0	2:0:0:0	0:0:0:3
C_1	0:1:1:0	0:0:0:0	0:0:0:0	0:0:0:0
G_1	0:0:0:2	0:0:0:0	0:0:0:0	0:0:0:0
T_1	0:1:1:1	0:0:0:0	1:0:0:0	2:0:1:0



Overlapping Reads and Repeats

- Does this really speed up the process?
- A k-mer that appears N times, initiates N^2 comparisons (you consider all pairs of reads that share the k-mer substring)
- For an Alu that appears 10^6 times $\rightarrow 10^{12}$ comparisons too much
- How to avoid repeats:

Discard all k-mers that appear more than $t \times \text{Coverage}$, $(t \sim 10)$



Finding Overlapping Reads

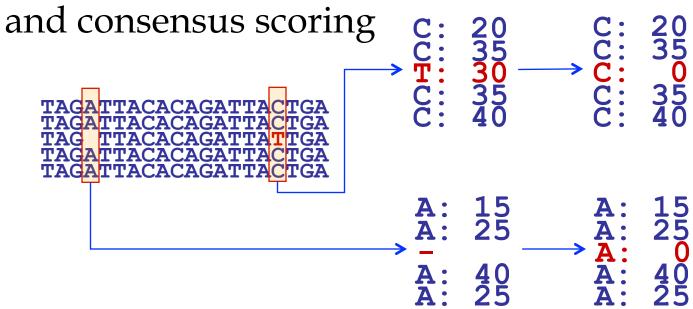
k-mer table makes it easy to create local multiple alignments from the overlapping reads

TAGATTACACAGATTACTGA
TAGATTACACAGATTACTGA
TAG TTACACAGATTACTGA
TAGATTACACAGATTACTGA
TAGATTACACAGATTACTGA
TAGATTACACAGATTACTGA
TAG TTACACAGATTACTGA
TAGATTACACAGATTACTGA
TAGATTACACAGATTACTGA



Finding Overlapping Reads (cont'd)

• Correct errors using multiple alignment



- Score alignments
- Accept alignments with good scores

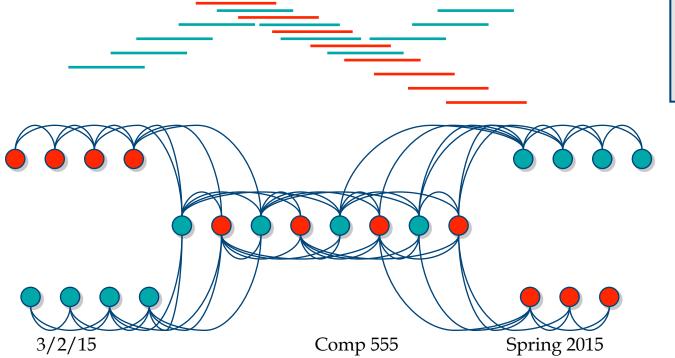


Layout

- Repeats are still a major challenge
- Do two aligned fragments really overlap, or are they from two copies of a repeat?
- Solution: repeat masking hide the repeats!!!
- Masking results in high rate of misassembly (up to 20%)
- Misassembly means alot more work at the finishing step



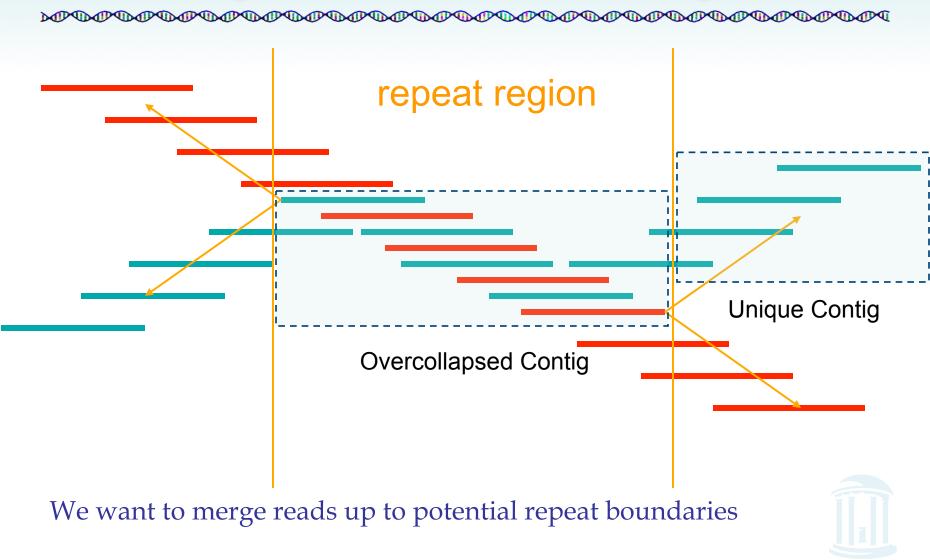
- Overlap graph:
 - Nodes: reads r_1 r_n
 - Edges: overlaps (r_i, r_j, shift, orientation, score)

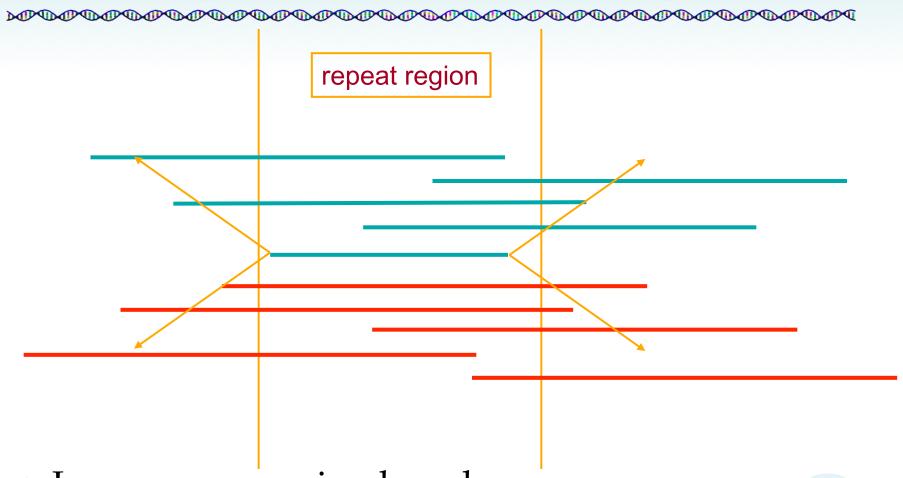


Reads that come from two regions of the genome (blue and red) that contain the same repeat

Note:

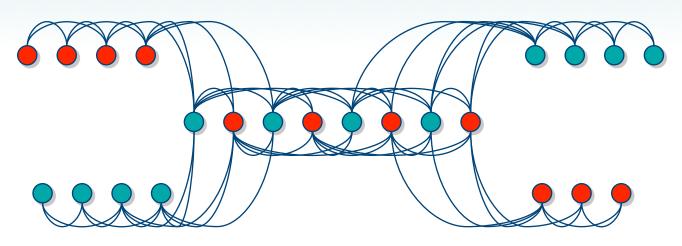
of course, we don't know the "color" of these nodes



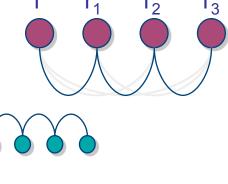


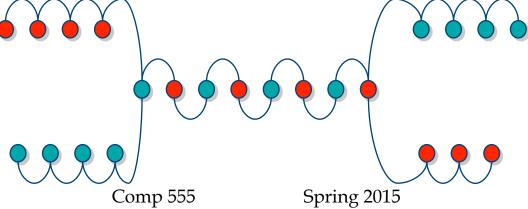
- Ignore non-maximal reads
- Merge only maximal reads into contigs



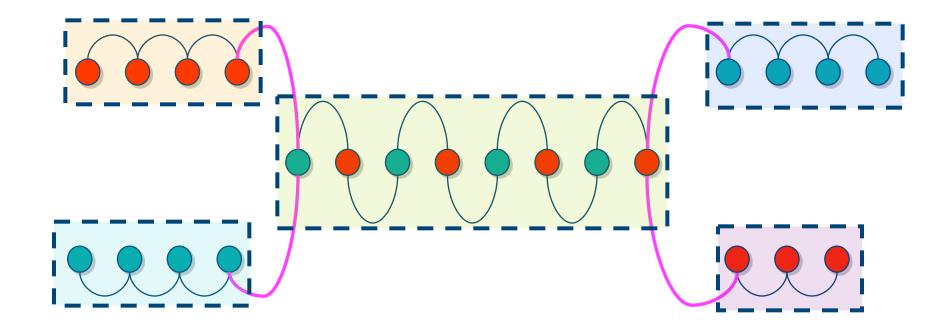


- Remove transitively inferable overlaps
 - If read r overlaps to the right reads r_1 , r_2 , and r_1 overlaps r_2 , then (r, r_2) can be inferred by (r, r_1) and (r_1, r_2)

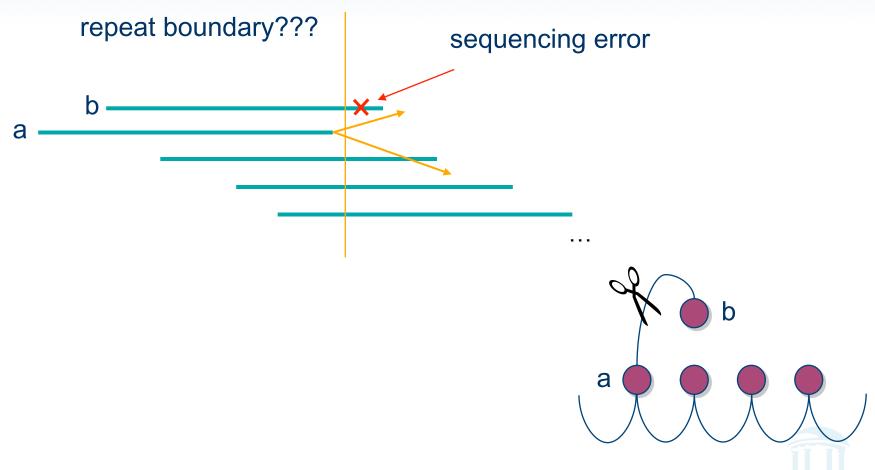








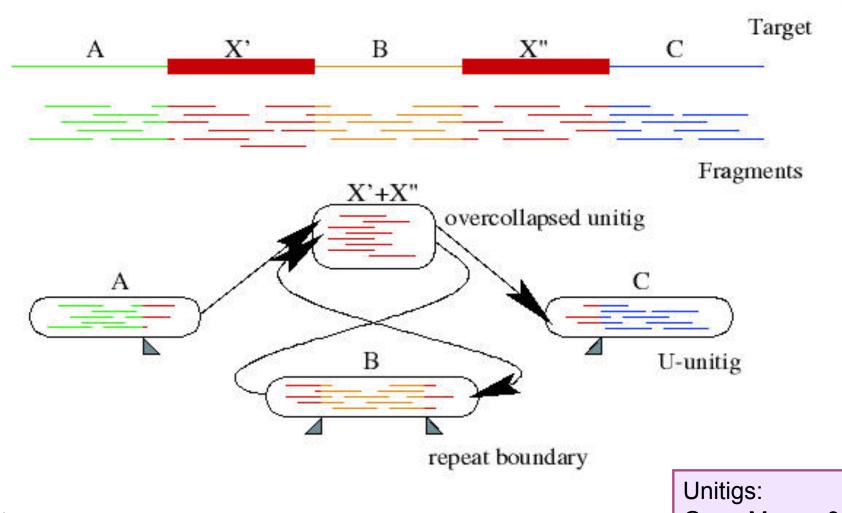




• Ignore "hanging" reads, when detecting repeat boundaries

Overlap graph after forming

proposed contigs proposed proposed contigs proposed proposed contigs proposed propos



3/2/15

Comp 555

Spring 2015

Gene Myers, 95

Repeats, errors, and contig lengths

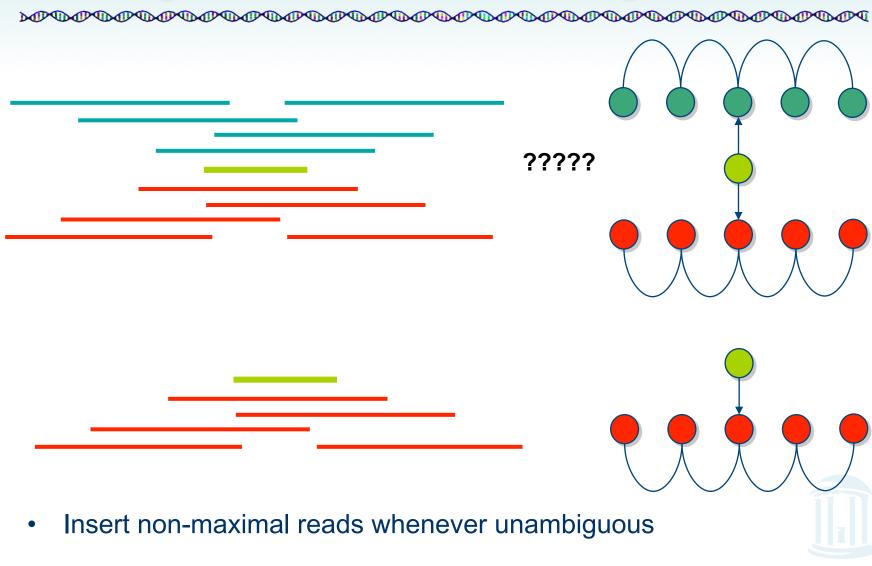
- Repeats shorter than read length are easily resolved
 - Read that spans across a repeat disambiguates order of flanking regions
- Repeats with more base pair diffs than sequencing error rate are OK
 - We throw overlaps between two reads in different copies of the repeat
- To make the genome appear less repetitive, try to:
 - Increase read length
 - Decrease sequencing error rate

Role of error correction:

Discards up to 98% of single-letter sequencing errors decreases error rate

- ⇒ decreases effective repeat content
- ⇒ increases contig length

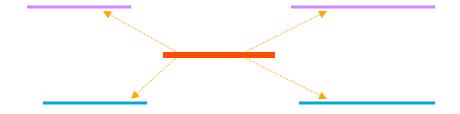






Normal density

Too dense: Overcollapsed?



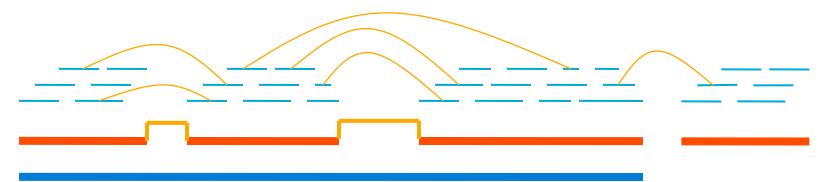
Inconsistent links: Overcollapsed?



(cont'd)

Find all links between unique contigs

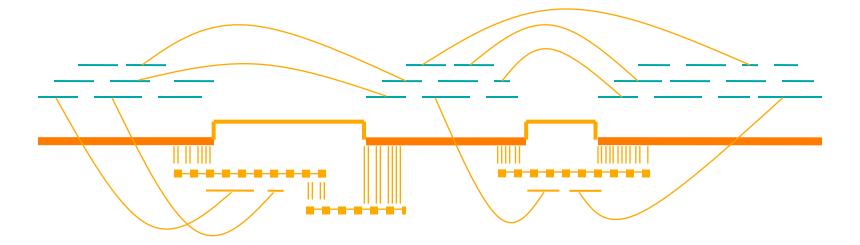
Connect contigs incrementally, if ≥ 2 links



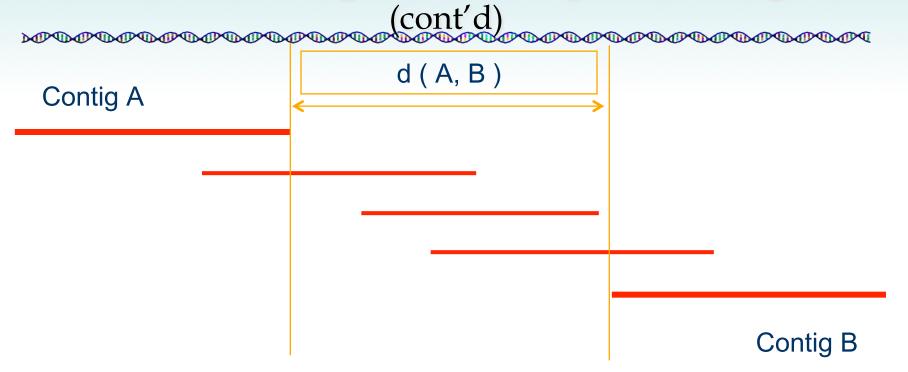


(cont'd)

Fill gaps in supercontigs with paths of overcollapsed contigs







Define G = (V, E)

V := contigs

E := (A, B) such that d(A, B) < C

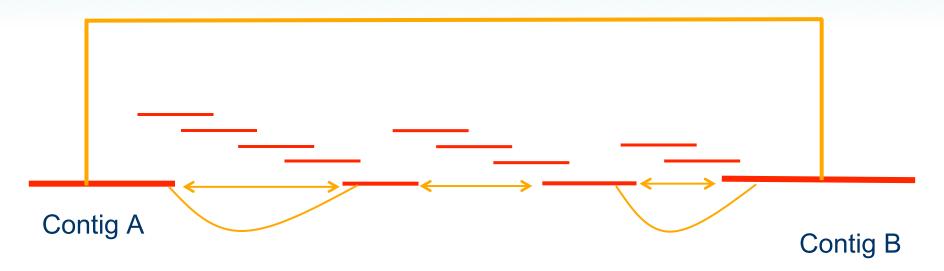
Reason to do so: Efficiency; full shortest paths cannot be computed

3/2/15

Comp 555

Spring 2015





Define T: contigs linked to either A or B

Fill gap between A and B if there is a path in G passing only from contigs in T



Consensus

 A consensus sequence is derived from a profile of the assembled fragments

• A sufficient number of reads is required to ensure a statistically significant consensus

Reading errors are corrected



Derive Consensus Sequence

POPODO DO POPODO POPODO DO POPODO POPO



TAGATTACACAGATTACTGACTTGATGGCGTAA CTA

Derive multiple alignment from pairwise read alignments

Derive each consensus base by weighted voting



Some Assemblers

DO CONTRACTO DE CO

PHRAP

- Early assembler, widely used, good model of read errors
- Overlap $O(n^2) \rightarrow layout$ (no mate pairs) \rightarrow consensus

Celera

- First assembler to handle large genomes (fly, human, mouse)
- Overlap → layout → consensus

Arachne

- Public assembler (mouse, several fungi)
- Overlap \rightarrow layout \rightarrow consensus

Phusion

- Overlap \rightarrow clustering \rightarrow PHRAP \rightarrow assemblage \rightarrow consensus
- Euler
 - Indexing → Euler graph → layout by picking paths → consensus

EULER Fragment Assembly

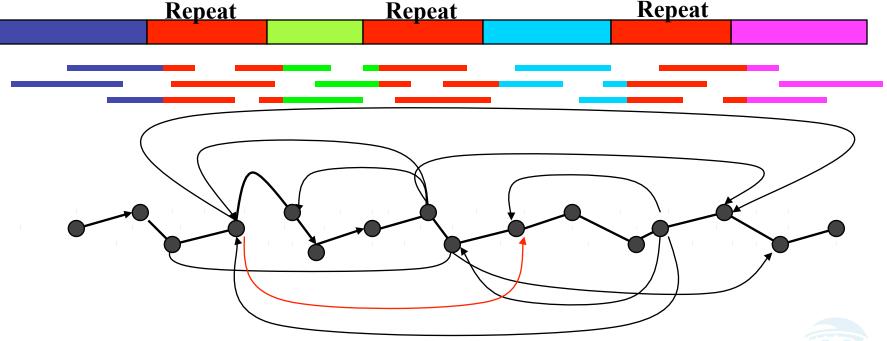


- Traditional "overlap-layout-consensus" technique has a high rate of mis-assembly
- EULER uses the Eulerian Path approach borrowed from the SBH problem
- Fragment assembly without repeat masking can be done in linear time with greater accuracy



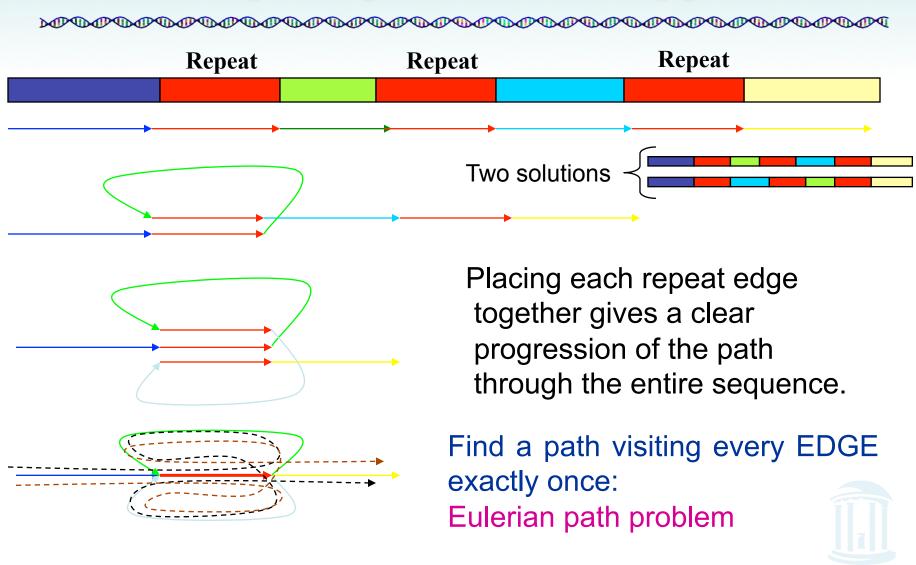
Overlap Graph: Hamiltonian Approach

Each vertex represents a read from the original sequence. Vertices from repeats are connected to many others.

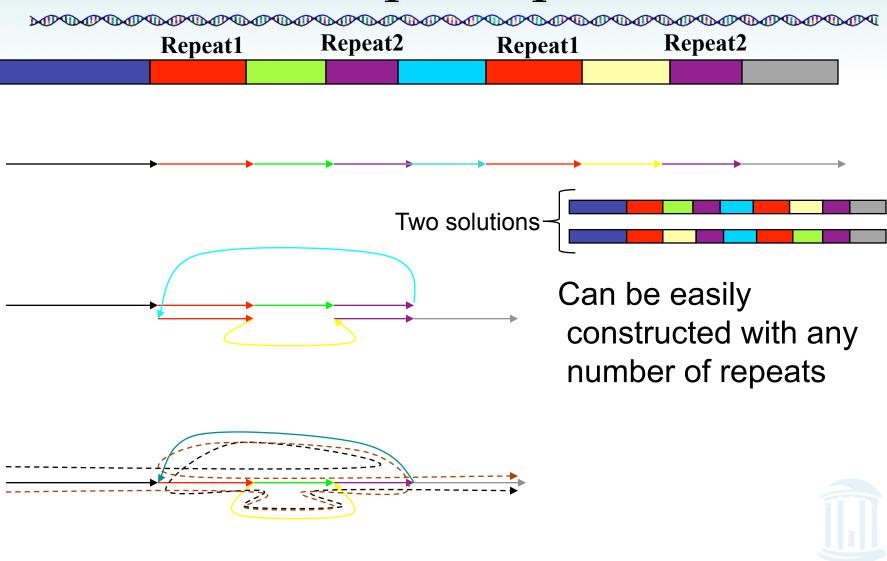


Find a path visiting every VERTEX exactly once: Hamiltonian path problem

Overlap Graph: Eulerian Approach



Multiple Repeats



Construction of Repeat Graph

 Construction of repeat graph from k – mers: emulates an SBH experiment with a huge (virtual) DNA chip.

• <u>Breaking reads into *k* – mers</u>: Transform sequencing data into virtual DNA chip data.



Construction of Repeat Graph (cont'd)

 Error correction in reads: "consensus first" approach to fragment assembly. Makes reads (almost) error-free BEFORE the assembly even starts.

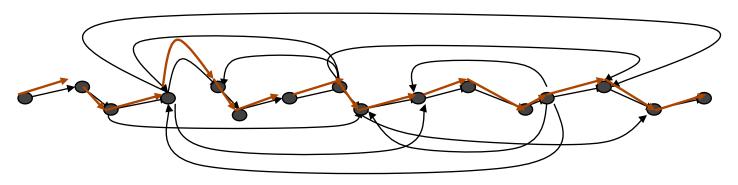
• Using reads and mate-pairs to simplify the repeat graph (Eulerian Superpath Problem).



Approaches to Fragment Assembly

Find a path visiting every VERTEX exactly once in the OVERLAP graph:

Hamiltonian path problem



NP-complete: algorithms unknown

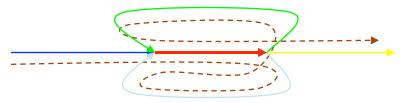


Approaches to Fragment Assembly

(cont'd)

Find a path visiting every EDGE exactly once in the REPEAT graph:

Eulerian path problem

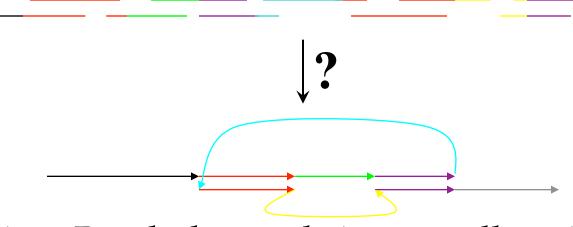


Linear time algorithms are known



Making Repeat Graph Without DNA

• Problem: Construct the repeat graph from a collection of reads.

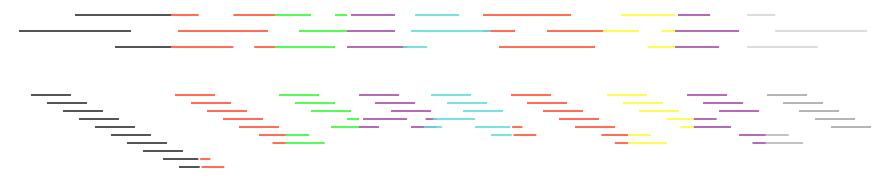


Solution: Break the reads into smaller pieces.



Repeat Sequences: Emulating a DNA Chip

 Virtual DNA chip allows the biological problem to be solved within the technological constraints.





Repeat Sequences: Emulating a DNA Chip (cont'd)

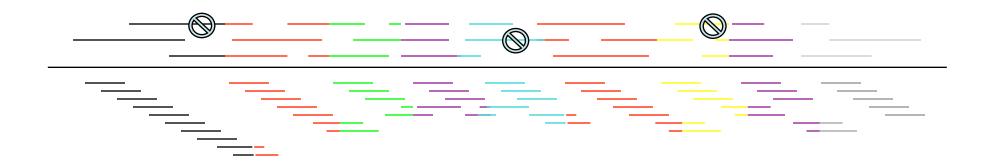
 Reads are constructed from an original sequence in lengths that allow biologists a high level of certainty.

 They are then broken again to allow the technology to sequence each within a reasonable array.



Minimizing Errors

• If an error exists in one of the 20-mer reads, the error will be perpetuated among all of the smaller pieces broken from that read.





Minimizing Errors (cont'd)

• However, that error will not be present in the other instances of the 20-mer read.

 So it is possible to eliminate most point mutation errors before reconstructing the original sequence.



Conclusions

Graph theory is a vital tool for solving biological problems

 Wide range of applications, including sequencing, motif finding, protein networks, and many more



References



Simons, Robert W. Advanced Molecular Genetics Course, UCLA (2002).
 http://www.mimg.ucla.edu/bobs/C159/Presentations/Benzer.pdf

• Batzoglou, S. *Computational Genomics Course*, Stanford University (2006). http://ai.stanford.edu/~serafim/CS262_2006/

