



Lecture 14: DNA Sequencing

Study Chapter 8.9

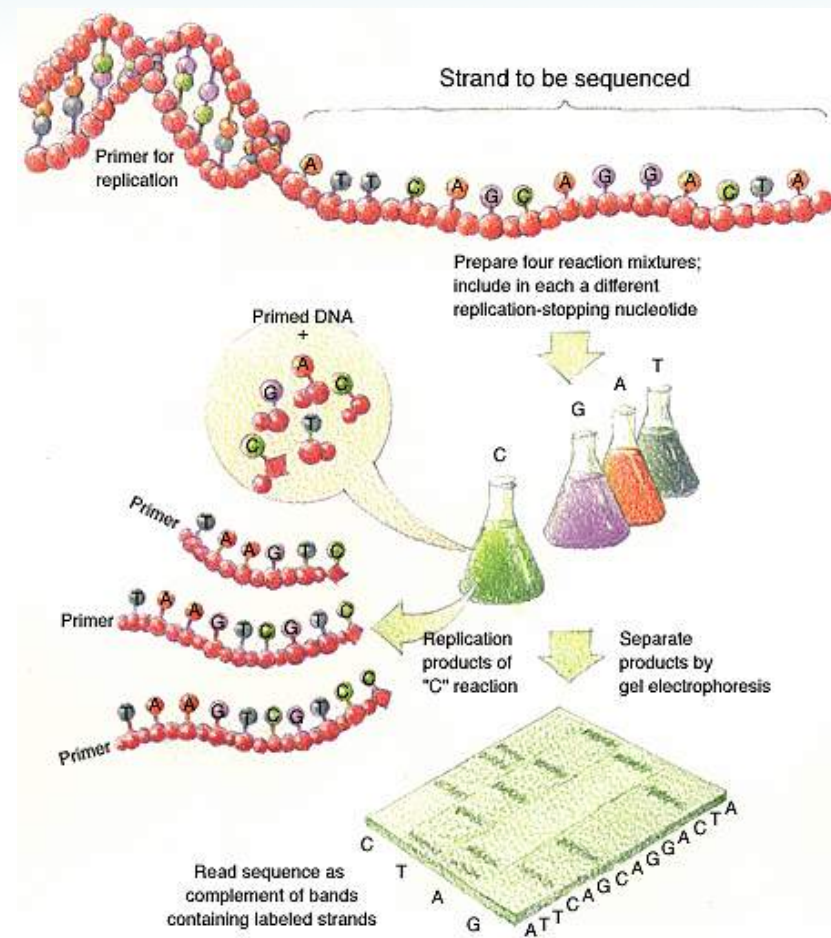
Midterm on Tuesday 10/15

Open book, open notes, no computer

Study Session on 10/14 in FB008 from 5pm-7pm

DNA Sequencing

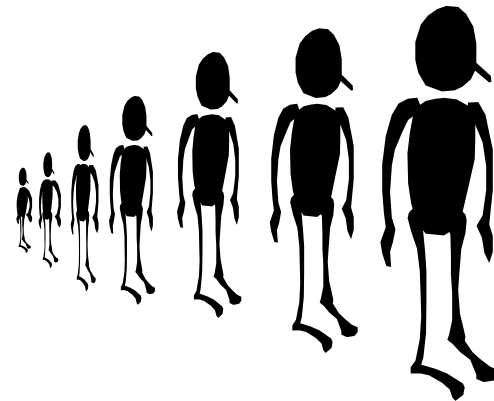
- Shear DNA into millions of small fragments
- Read 500 – 700 nucleotides at a time from the small fragments (Sanger method)



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



DNA

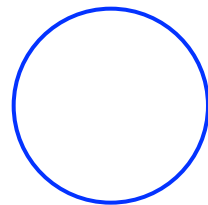
Shake & Break
(by Digestion or Sonication)



DNA fragments

Clone

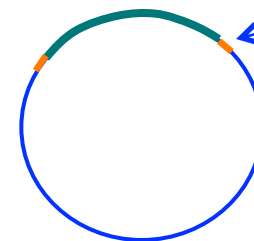
Vector
Circular genome
(bacterium, plasmid)



+



=



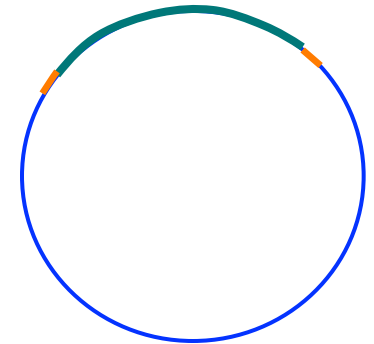
Known
location
(restriction
site)



Different Types of Vectors



<u>VECTOR</u>	<u>Size of insert (bp)</u>
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							
c	a	t	t						
c	a	t	t	c	t				

ddC Reaction -

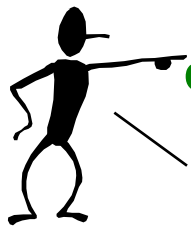
c									
c	a	t	t	c					
c	a	t	t	c	t	g	a	c	

ddG Reaction -

c	a	t	t	c	t	g			
---	---	---	---	---	---	---	--	--	--

ddA Reaction -

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



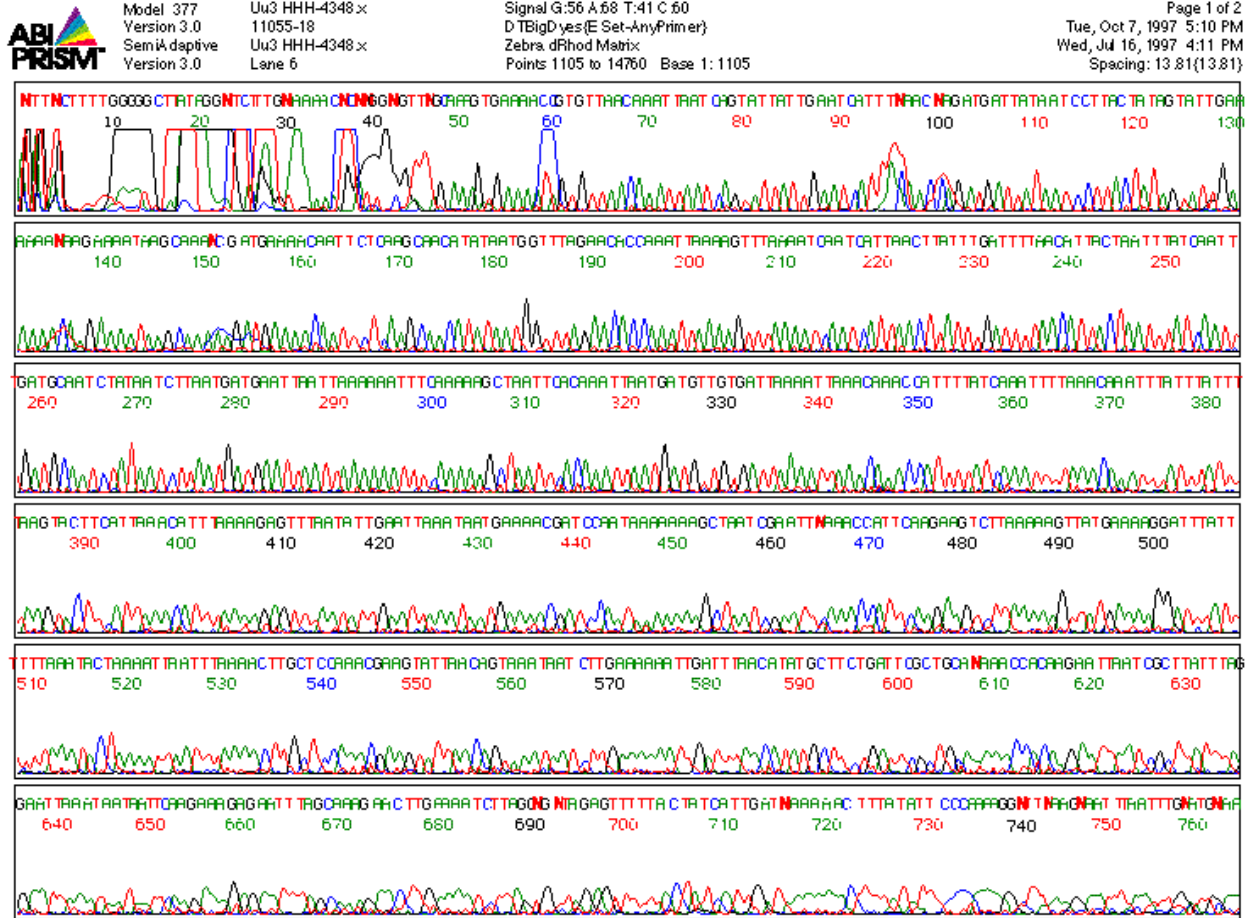
Good for up to
1000 base pairs



Challenging to Read Answers



Electropherogram



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



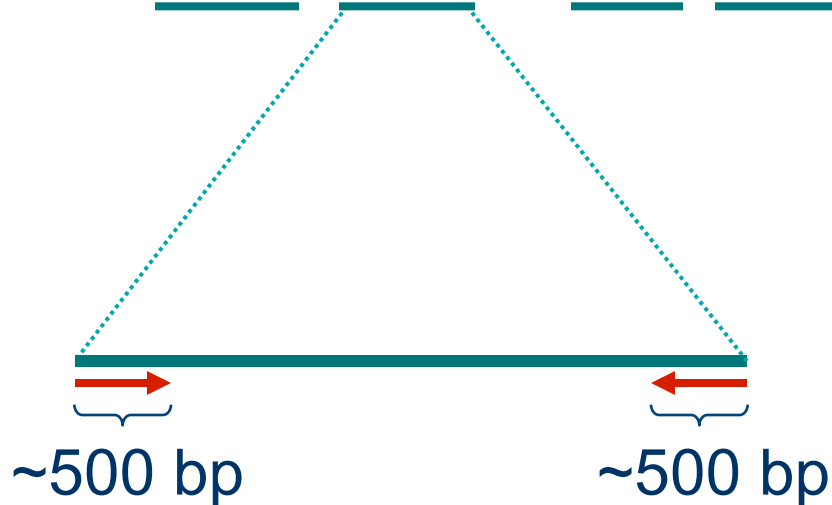
Shotgun Sequencing



genomic segment



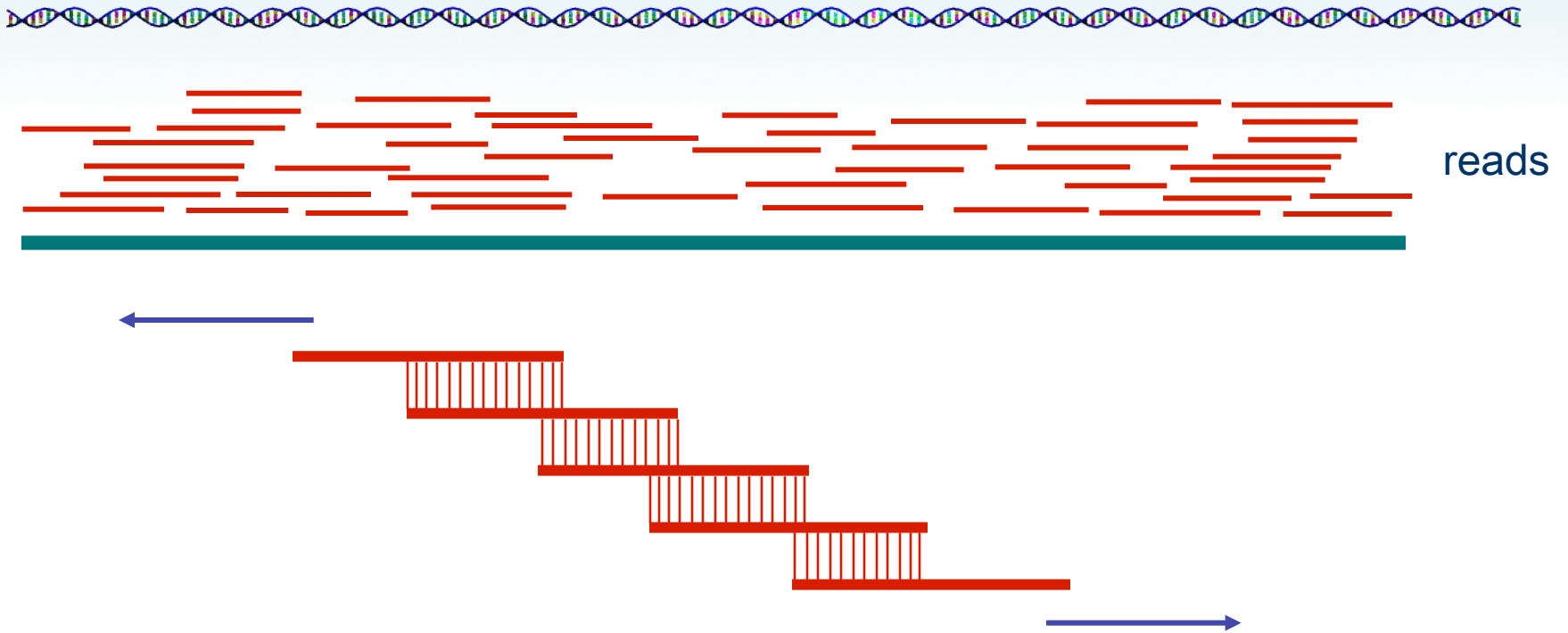
cut many times at
random (*Shotgun*)



Get one or two
reads from the ends
of each segment



Fragment Assembly

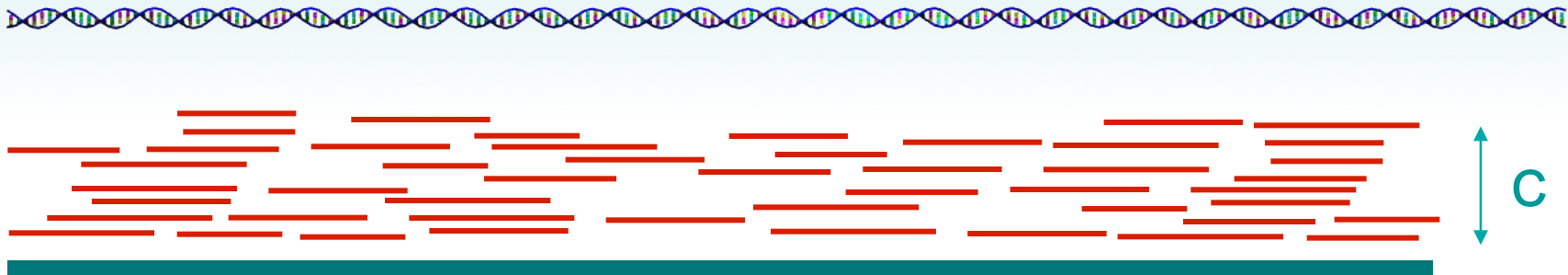


Cover region with ~ 7 -fold redundancy

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: l

How much coverage is enough?

Lander-Waterman model:

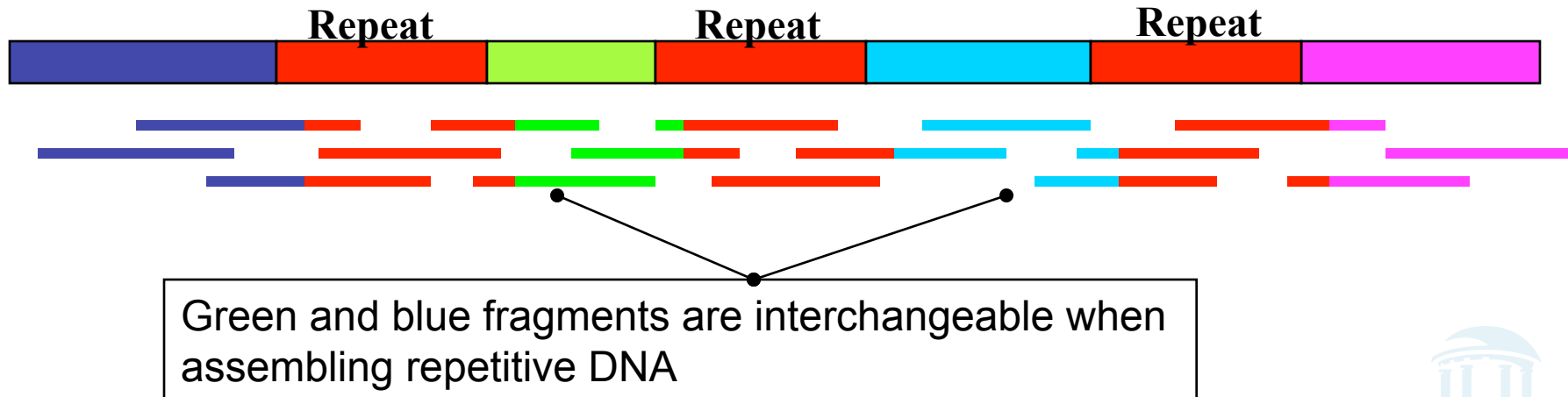
Assuming uniform distribution of reads, $C=10$ results in 1 gapped region per 1,000,000 nucleotides



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-consistency 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** $(a_1 \dots a_k)^N$ where $k \sim 3-6$
(e.g. CAGCAGTAGCAGCACCAG)
- **Transposons/retrotransposons**
 - **SINE** Short Interspersed Nuclear Elements
(e.g., *Alu*: ~300 bp long, $>10^6$ 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



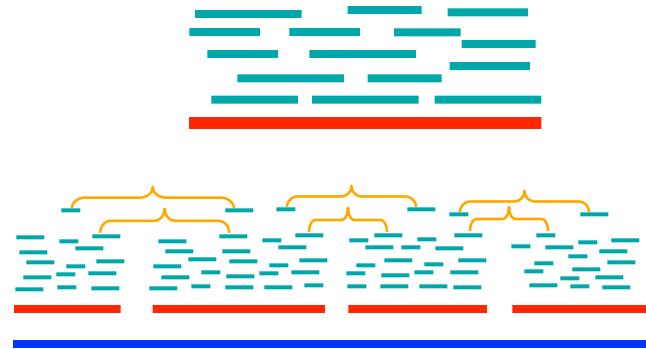
Assembler programs ARACHNE, PHRAP, CAP, TIGR, CELERA

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
($k \sim 20-24$)
- Find read-pairs sharing a k -mer
- Extend alignment –
throw away if not $>95\%$ similar



Histogram Example



$v = \text{tagattacacagattattga}$

- 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



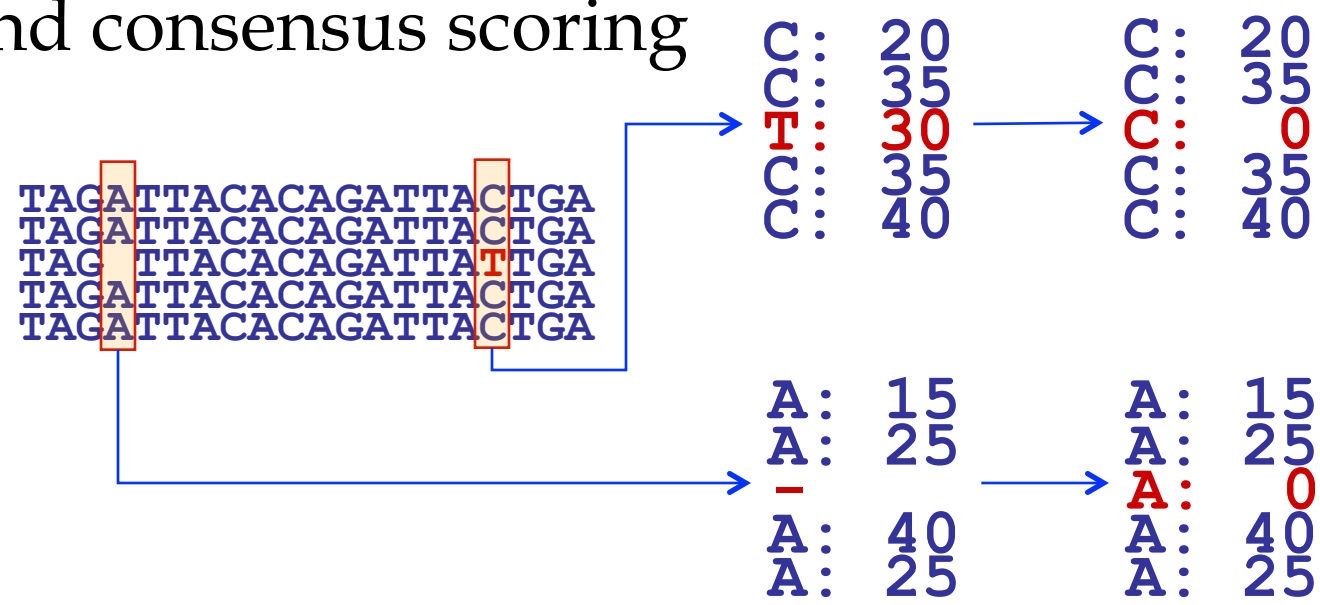
k-mer table makes it easy to create local multiple alignments from the overlapping reads



Finding Overlapping Reads (cont'd)



- Correct errors using multiple alignment and consensus scoring



- Score alignments
- Accept alignments with good scores



Layout



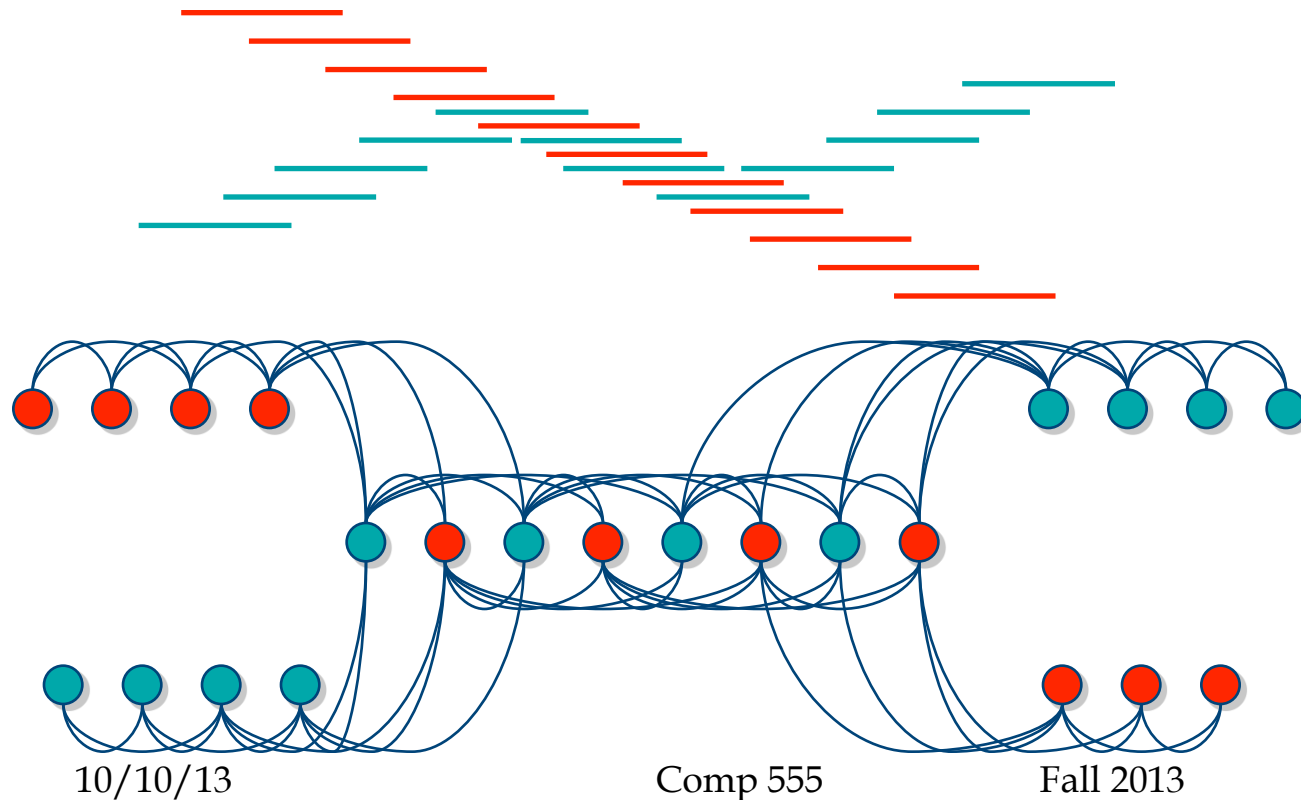
- Repeats are still a major challenge
- Do two aligned fragments really overlap, or are they from two copies of a repeat?
- Solution: repeat masking – hide the repeats!!!
- Masking results in high rate of misassembly (up to 20%)
- Misassembly means alot more work at the finishing step



2. Merge Reads into Contigs



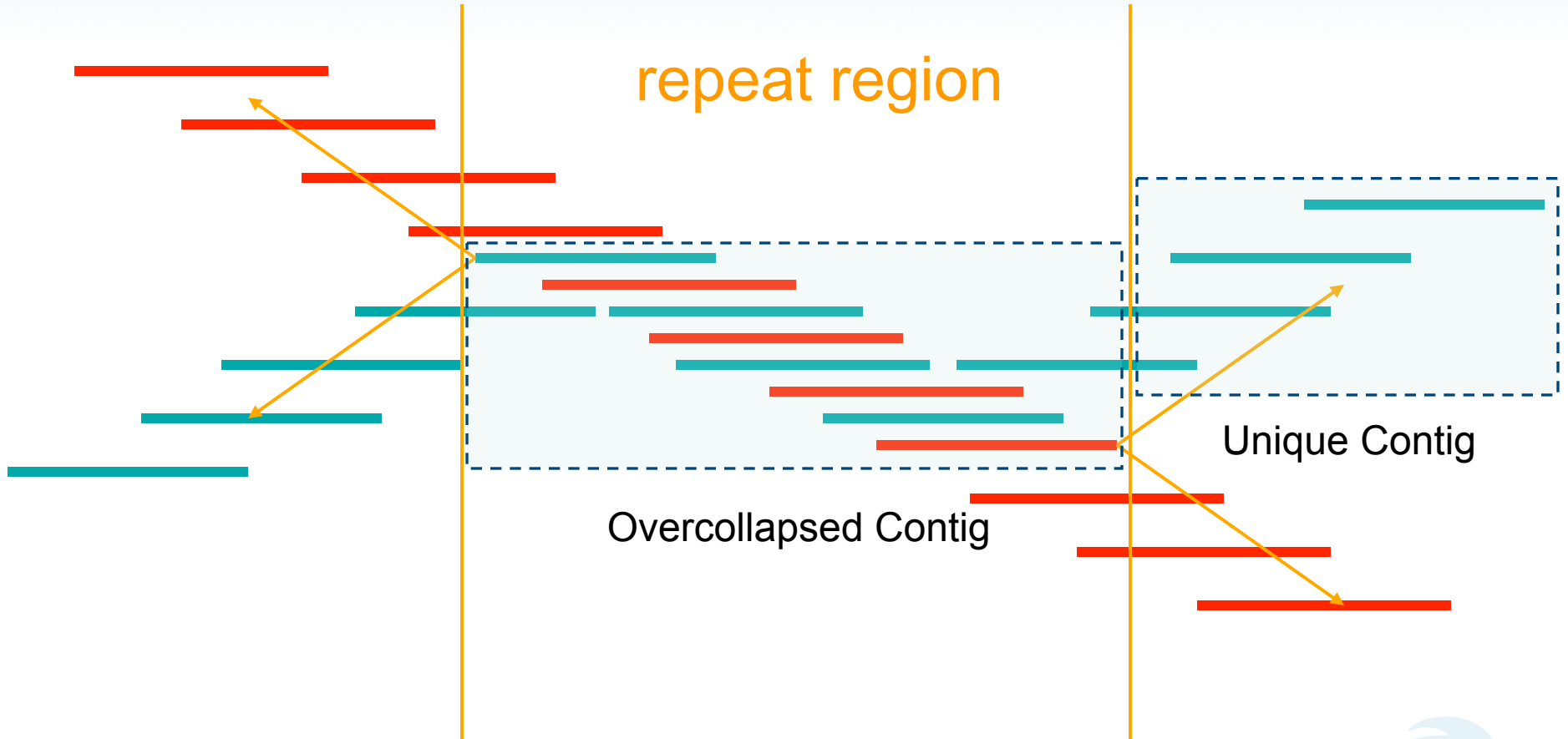
- Overlap graph:
 - Nodes: reads $r_1 \dots r_n$
 - Edges: overlaps (r_i, r_j , shift, orientation, score)



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

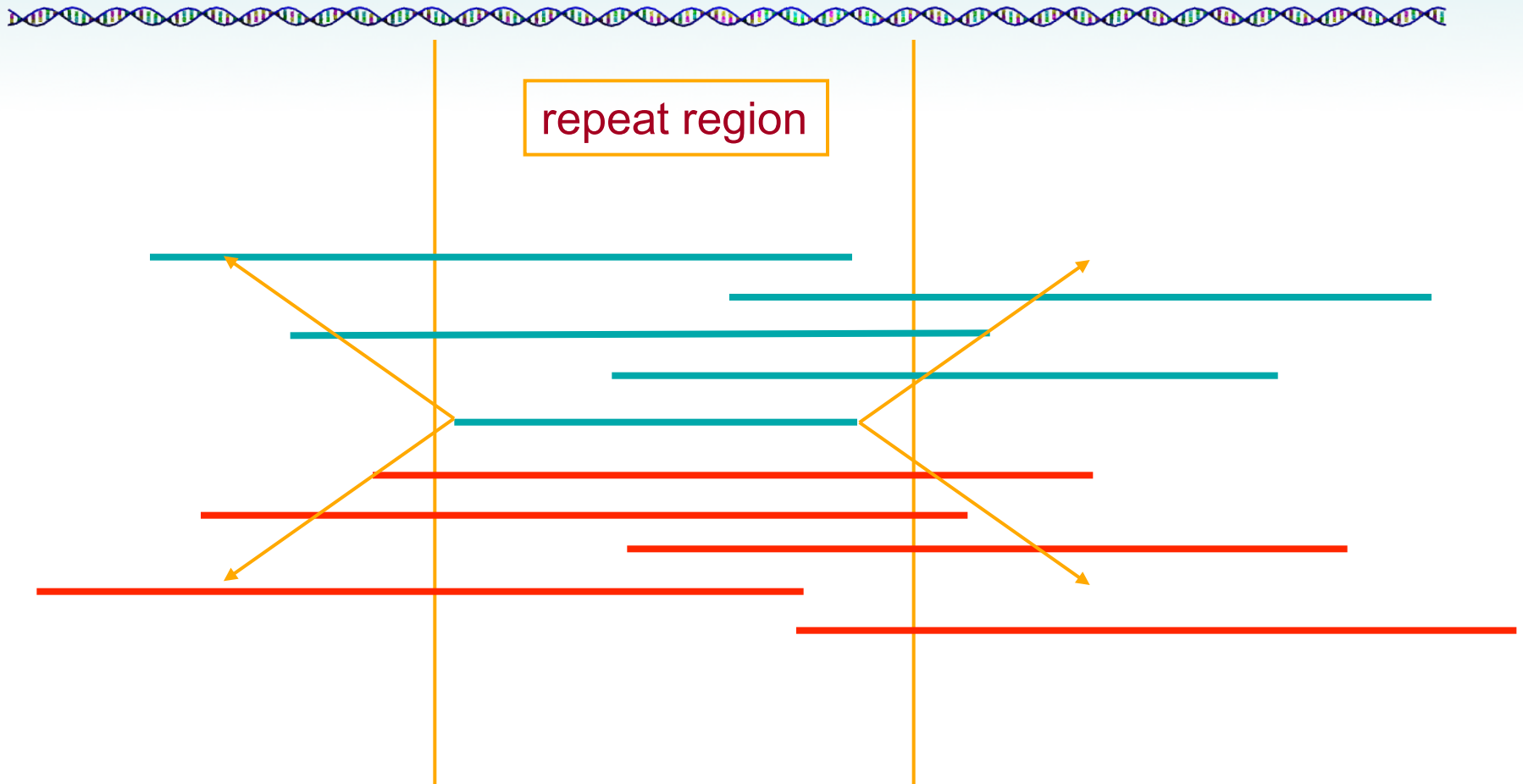
2. Merge Reads into Contigs



We want to merge reads up to potential repeat boundaries



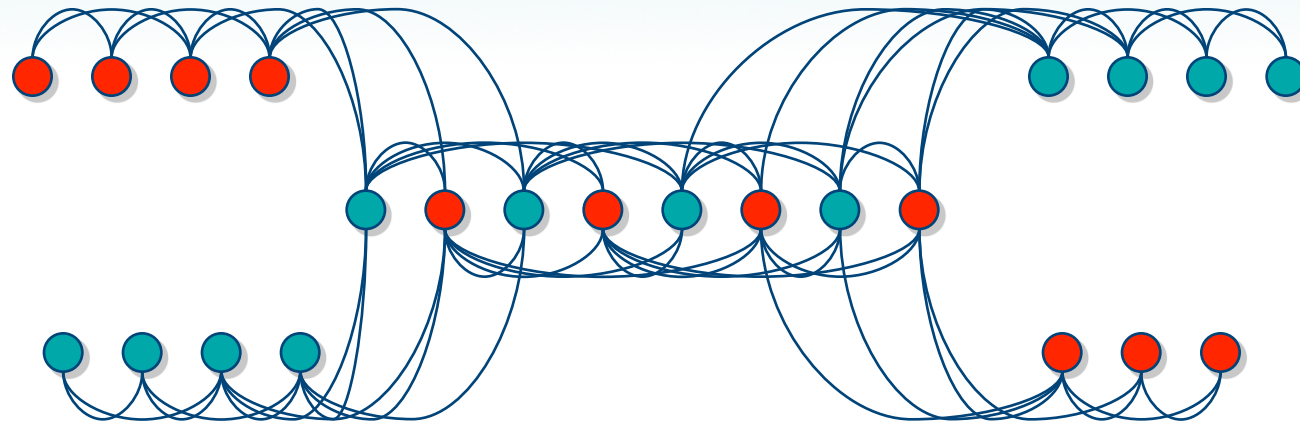
2. Merge Reads into Contigs



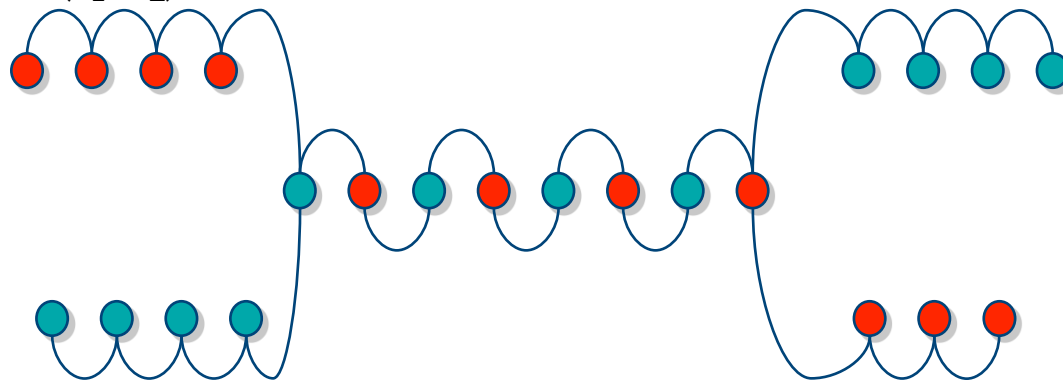
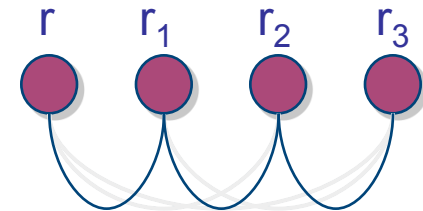
- Ignore non-maximal reads
- Merge only maximal reads into contigs



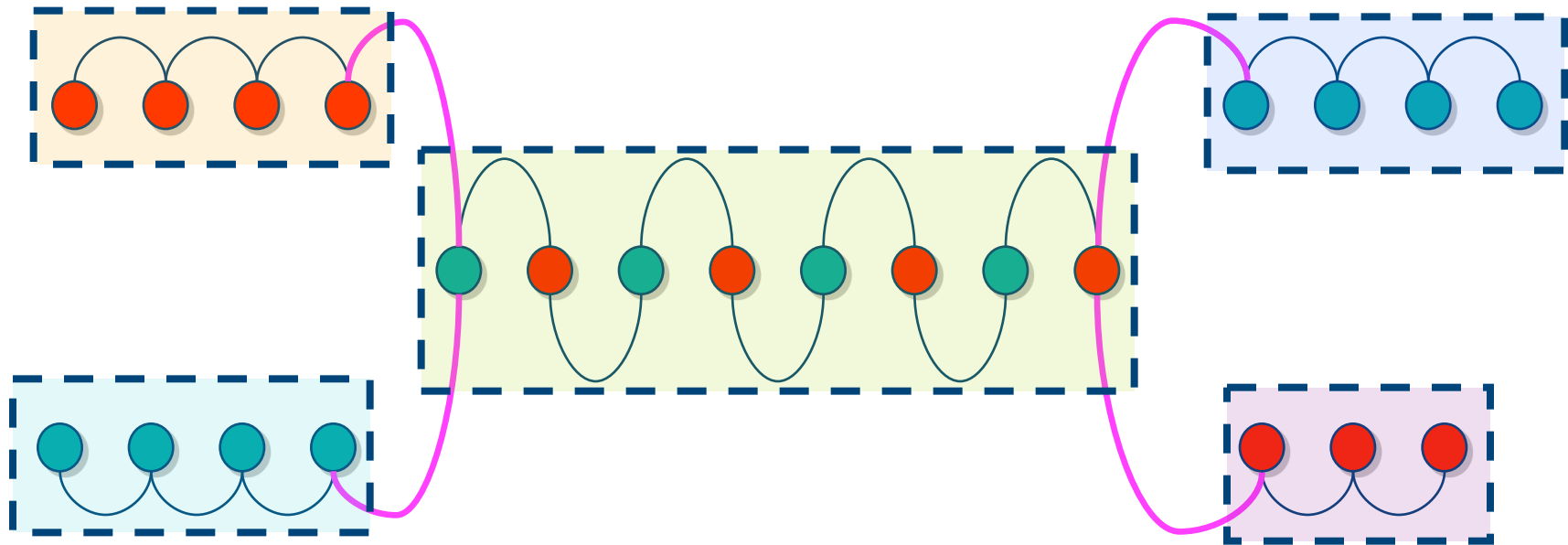
2. Merge Reads into Contigs



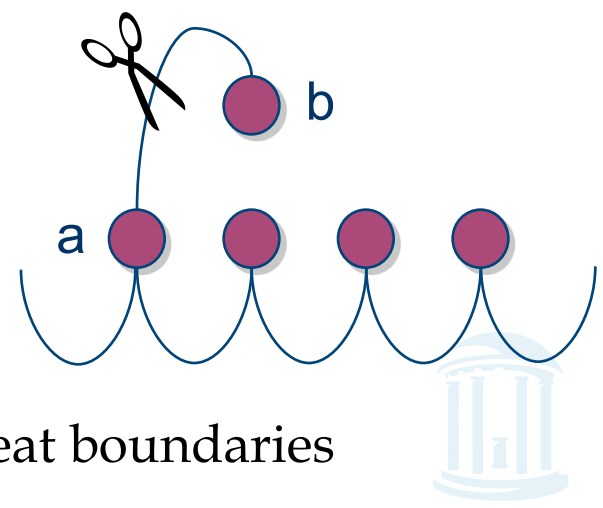
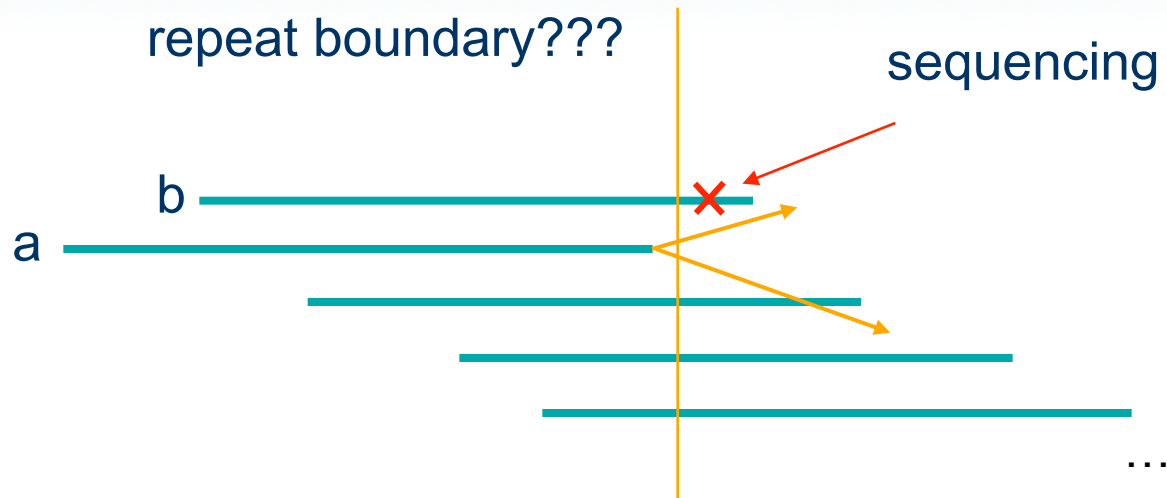
- Remove transitively inferable overlaps
 - If read r overlaps to the right reads r_1 , r_2 , and r_1 overlaps r_2 , then (r, r_2) can be inferred by (r, r_1) and (r_1, r_2)



2. Merge Reads into Contigs



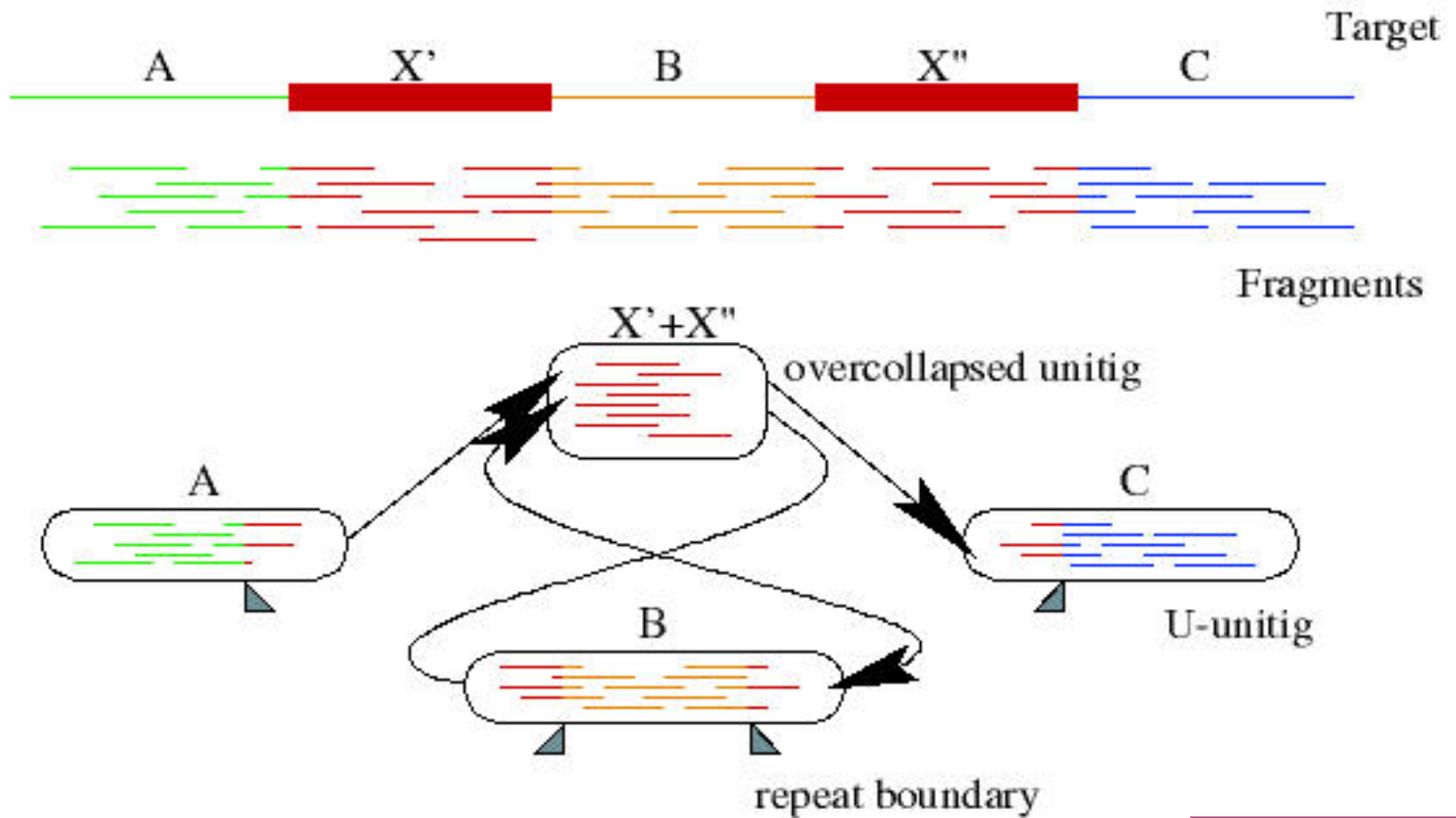
2. Merge Reads into Contigs



- Ignore “hanging” reads, when detecting repeat boundaries



Overlap graph after forming contigs



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:

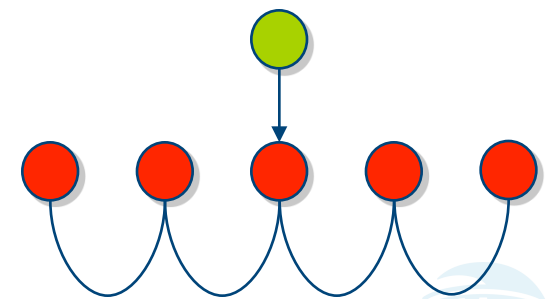
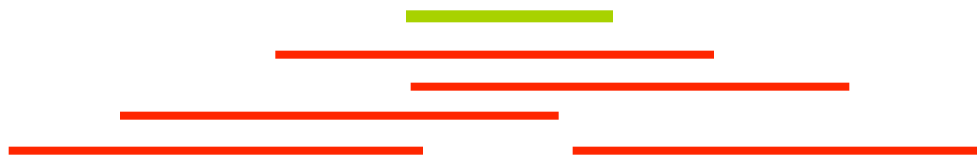
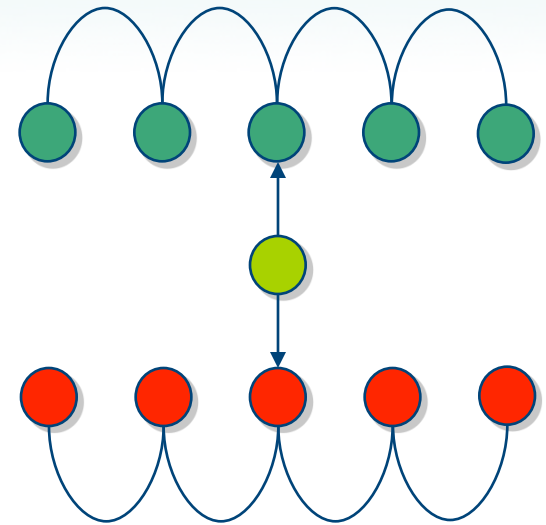
Discards up to 98% of single-letter sequencing errors
decreases error rate
⇒ decreases effective repeat content
⇒ increases contig length



2. Merge Reads into Contigs



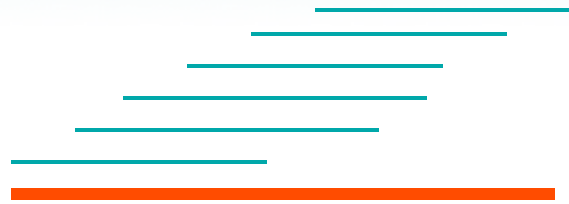
?????



- Insert non-maximal reads whenever unambiguous



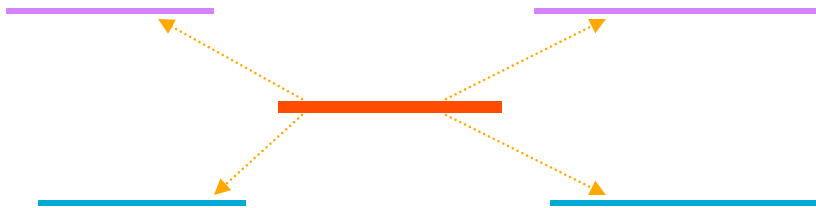
Link Contigs into Supercontigs



Normal density



Too dense:
Overcollapsed?



Inconsistent links:
Overcollapsed?



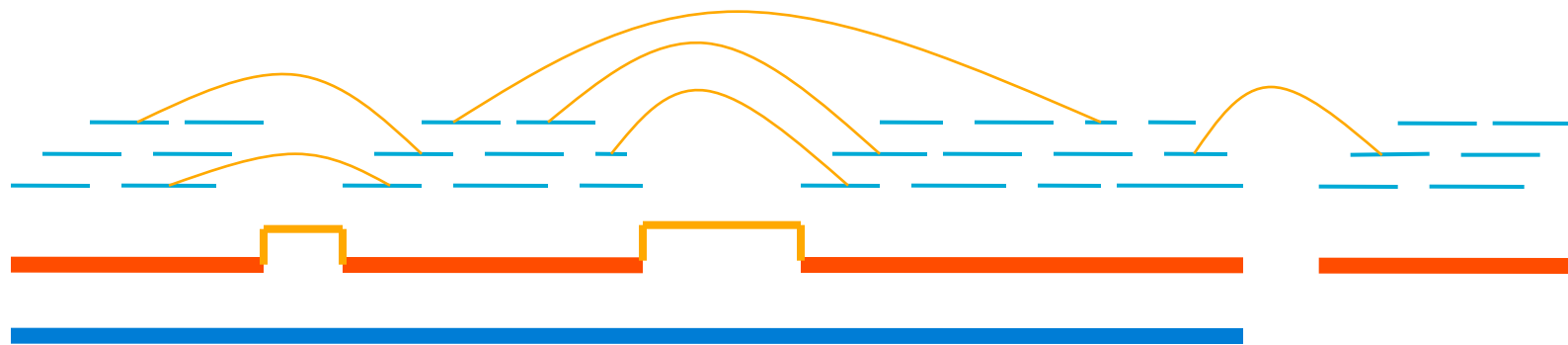
Link Contigs into Supercontigs

(cont'd)



Find all links between unique contigs

Connect contigs incrementally, if ≥ 2 links

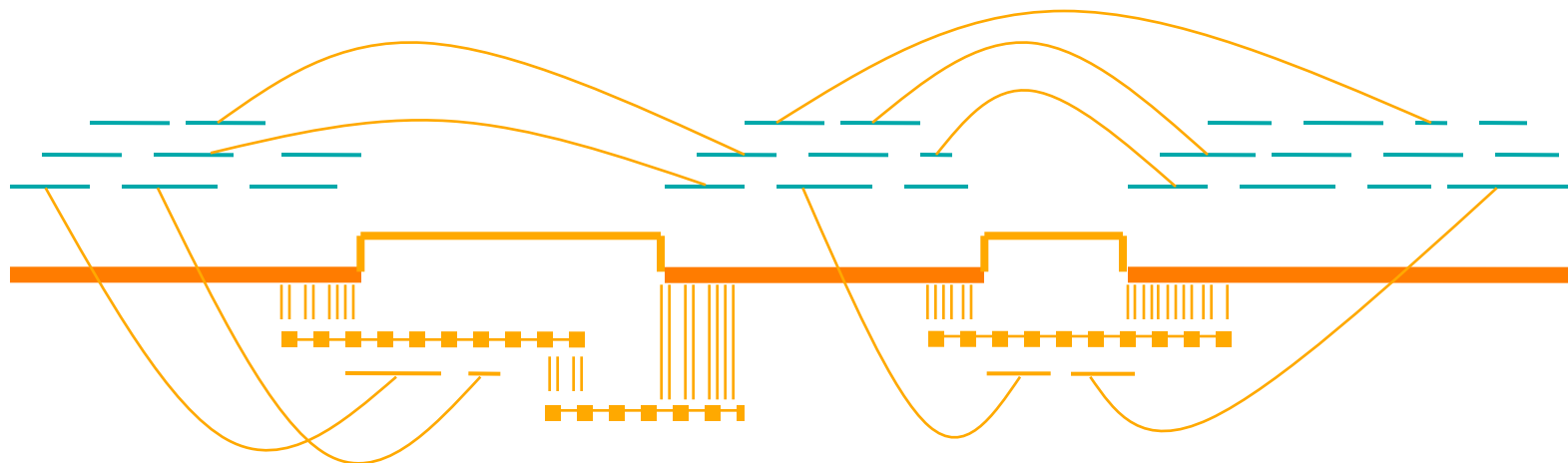


Link Contigs into Supercontigs

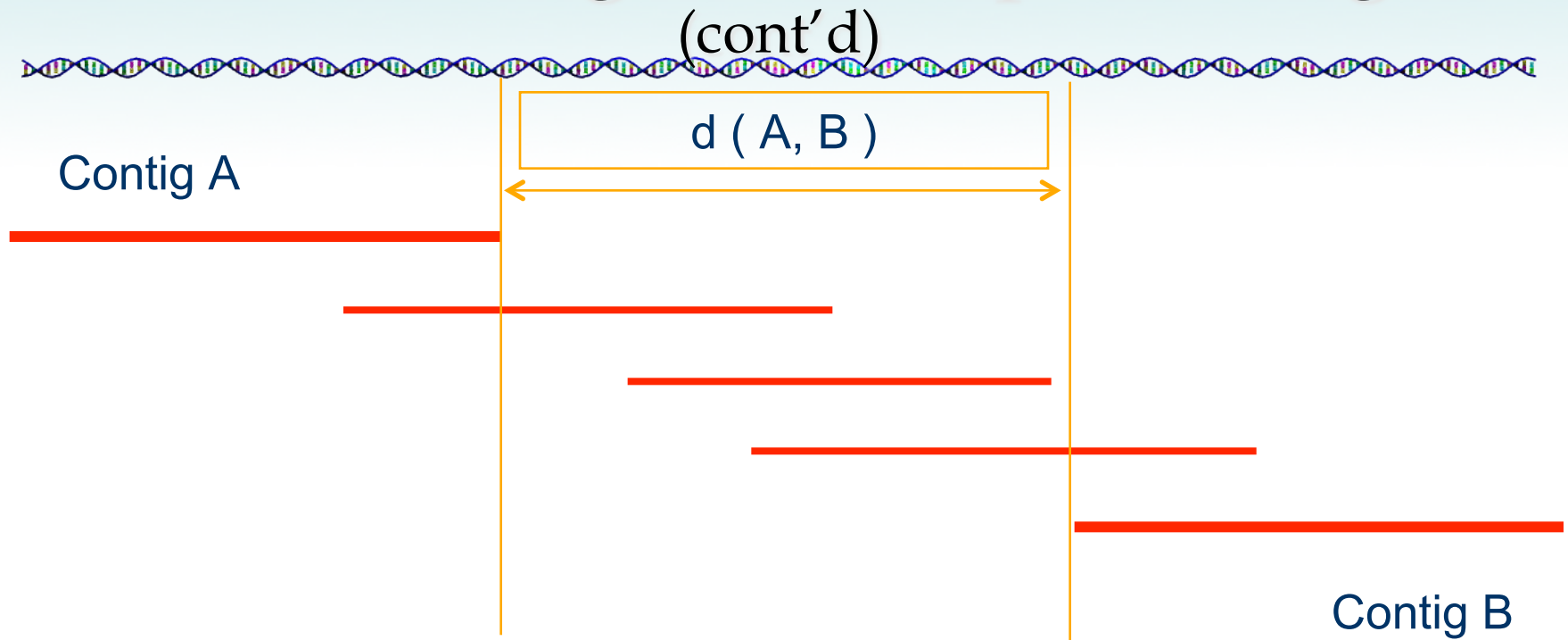
(cont'd)



Fill gaps in supercontigs with paths of overcollapsed contigs



Link Contigs into Supercontigs



Define $G = (V, E)$

$V :=$ contigs

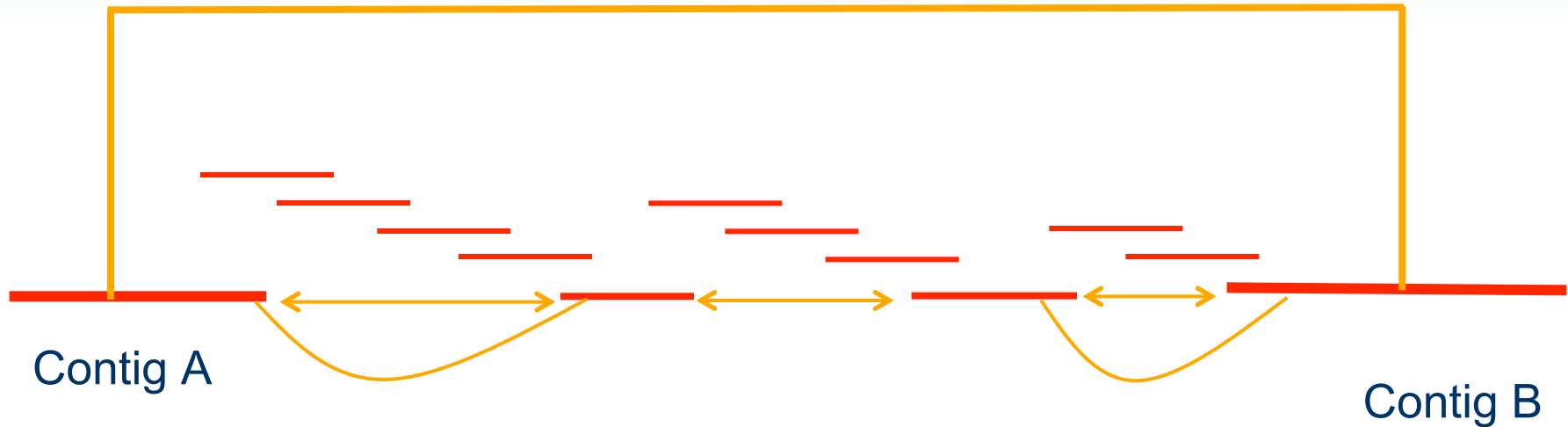
$E := (A, B)$ such that $d(A, B) < C$

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



Link Contigs into Supercontigs

(cont'd)



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



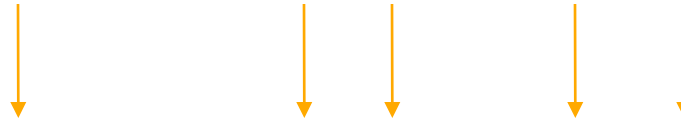
- A consensus sequence is derived from a profile of the assembled fragments
- A sufficient number of reads is required to ensure a statistically significant consensus
- Reading errors are corrected



Derive Consensus Sequence



```
TAGATTACACAGATTACTGA TTGATGGCGTAA CTA
TAGATTACACAGATTACTGACTTGGATGGCGTAACTA
TAG TTACACAGATTATTGACTTCATGGCGTAA CTA
TAGATTACACAGATTACTGACTTGGATGGCGTAA CTA
TAGATTACACAGATTACTGACTTGGGGTAA CTA
```



```
TAGATTACACAGATTACTGACTTGGCGTAA CTA
```

Derive **multiple alignment** from pairwise read alignments

Derive each consensus base by weighted voting



Some Assemblers



- PHRAP
 - Early assembler, widely used, good model of read errors
 - Overlap $O(n^2)$ → layout (no mate pairs) → consensus
- Celera
 - First assembler to handle large genomes (fly, human, mouse)
 - Overlap → layout → consensus
- Arachne
 - Public assembler (mouse, several fungi)
 - Overlap → layout → consensus
- Phusion
 - Overlap → clustering → PHRAP → assemblage → consensus
- Euler
 - Indexing → Euler graph → layout by picking paths → consensus



EULER Fragment Assembly



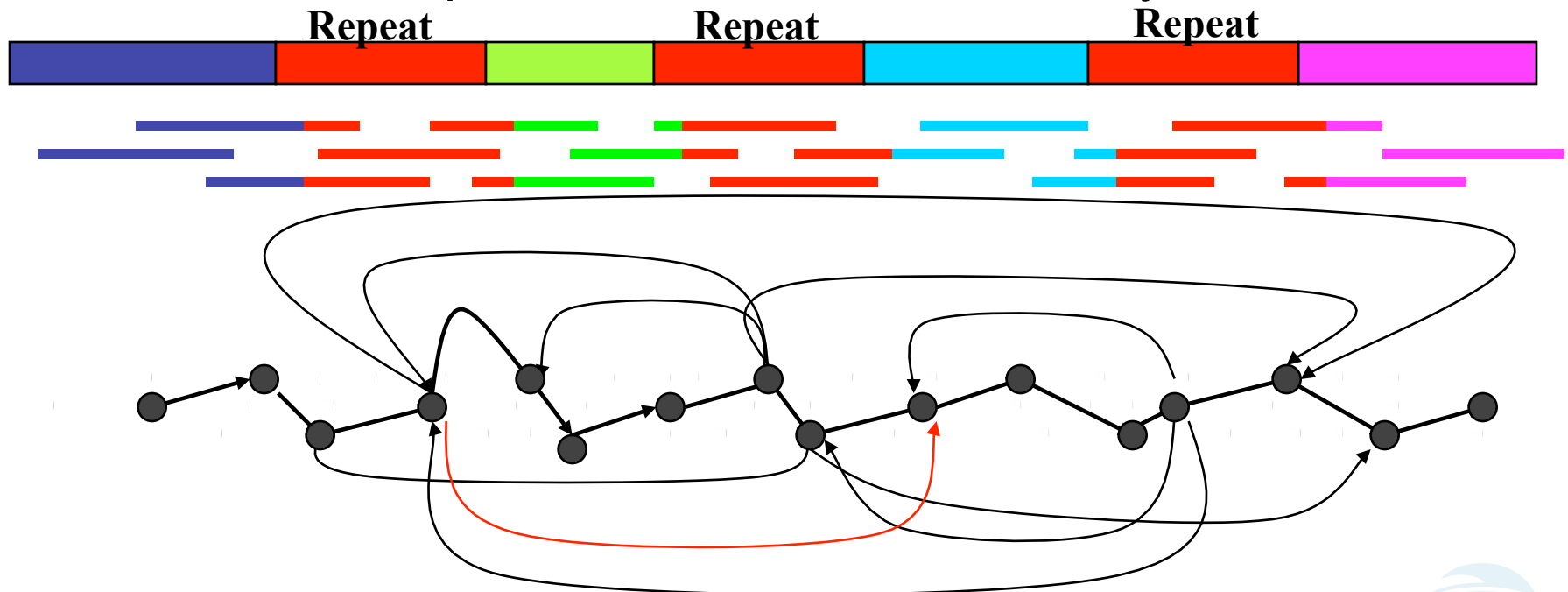
- Traditional “overlap-layout-consensus” technique has a high rate of mis-assembly
- EULER uses the Eulerian Path approach borrowed from the SBH problem
- Fragment assembly without repeat masking can be done in linear time with greater accuracy



Overlap Graph: Hamiltonian Approach



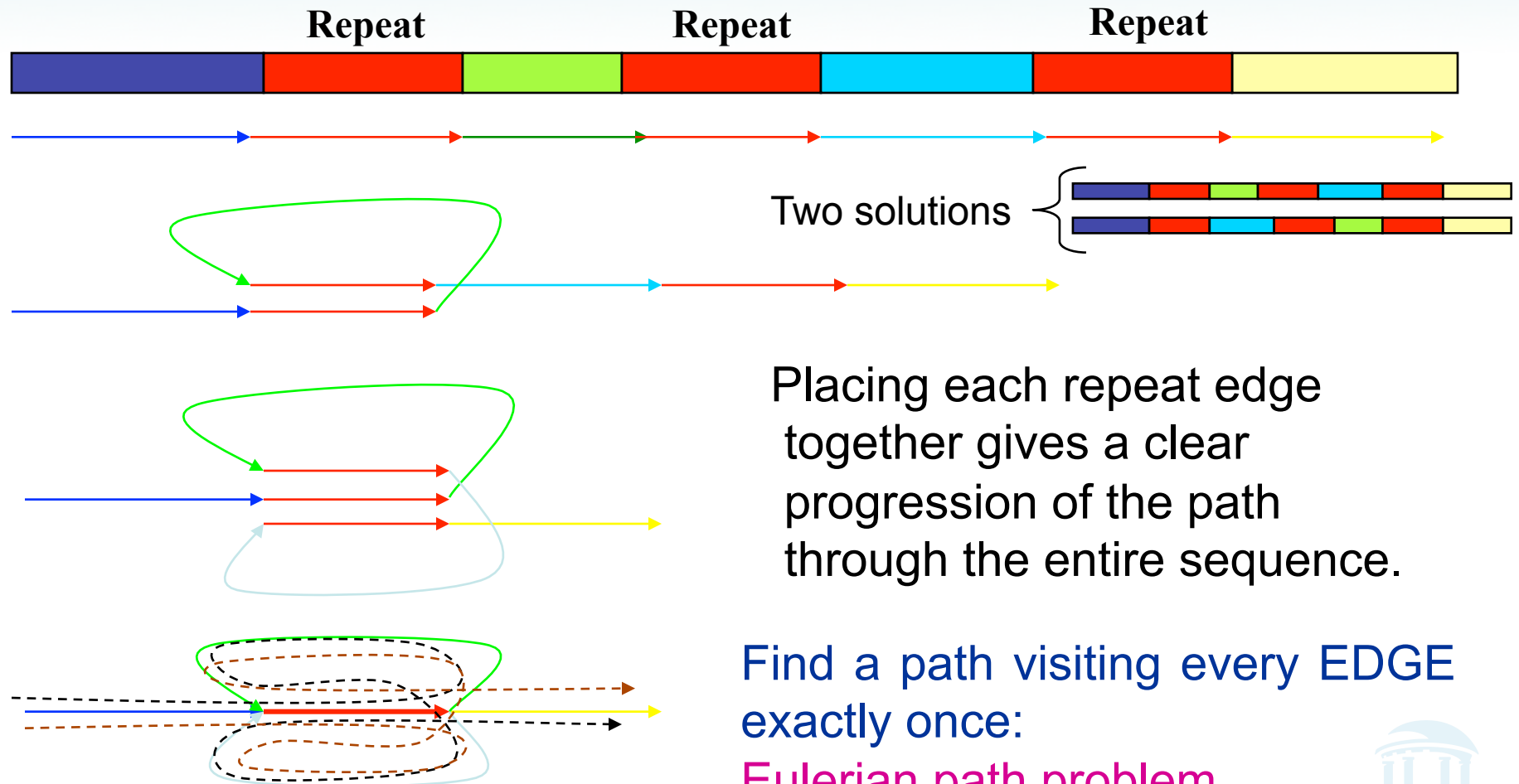
Each vertex represents a read from the original sequence.
Vertices from repeats are connected to many others.



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



Overlap Graph: Eulerian Approach



Placing each repeat edge together gives a clear progression of the path through the entire sequence.

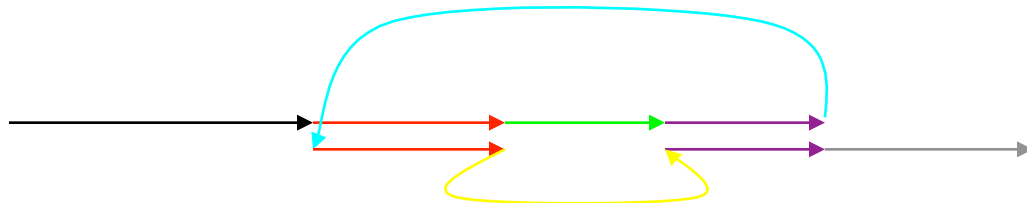
Find a path visiting every **EDGE** exactly once:
Eulerian path problem



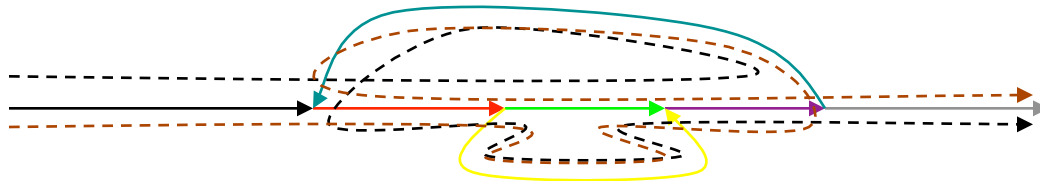
Multiple Repeats



Two solutions



Can be easily constructed with any number of repeats



Construction of Repeat Graph



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

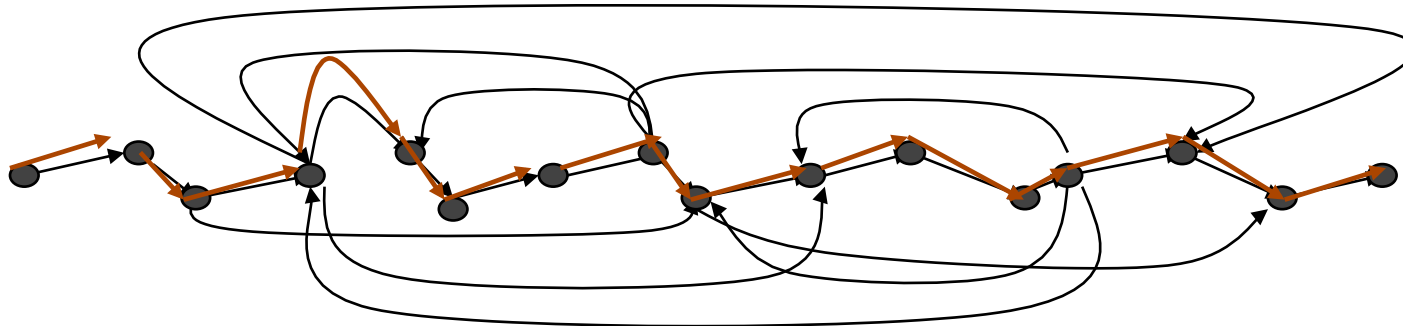


Approaches to Fragment Assembly



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

Hamiltonian path problem



NP-complete: algorithms unknown



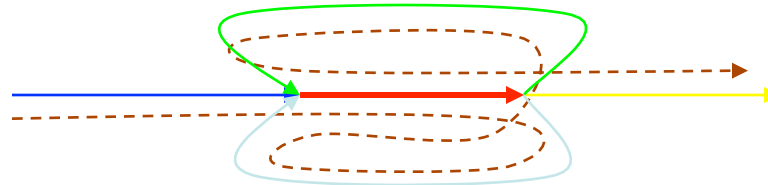
Approaches to Fragment Assembly

(cont'd)



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

Eulerian path problem



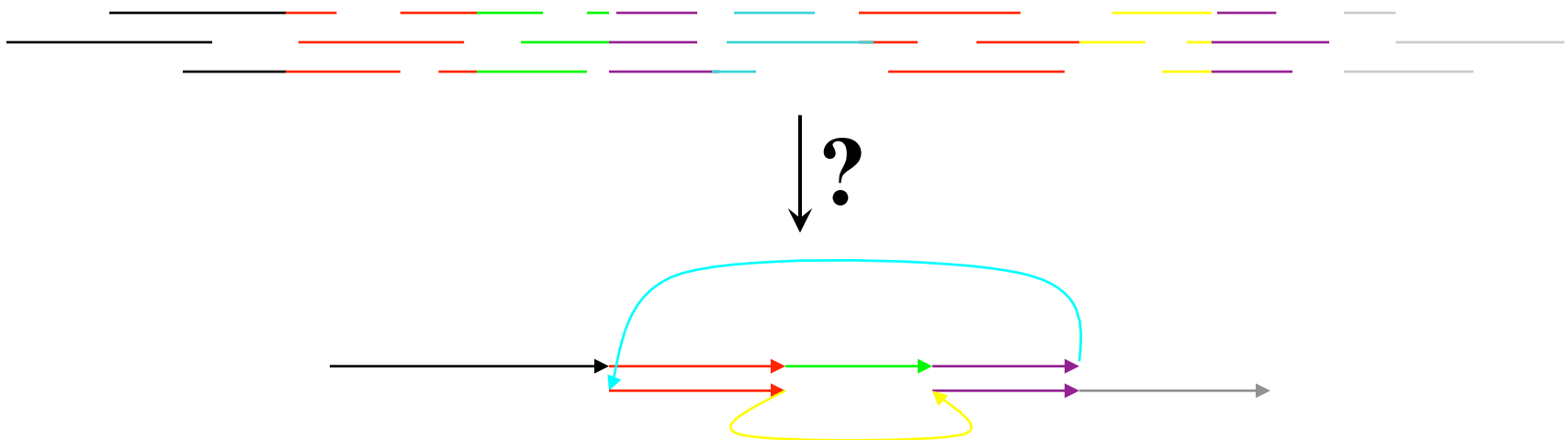
Linear time algorithms are known



Making Repeat Graph Without DNA



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



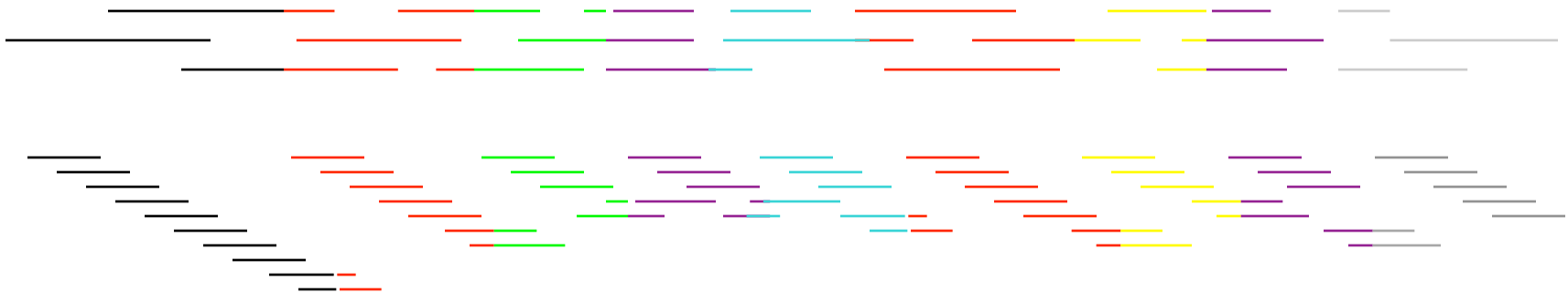
- Solution: Break the reads into smaller pieces.



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/

