

Computational Genomics

10-810/02-710, Spring 2009

DNA sequencing and genome assembly

Eric Xing

Lecture 3, January 21, 2009

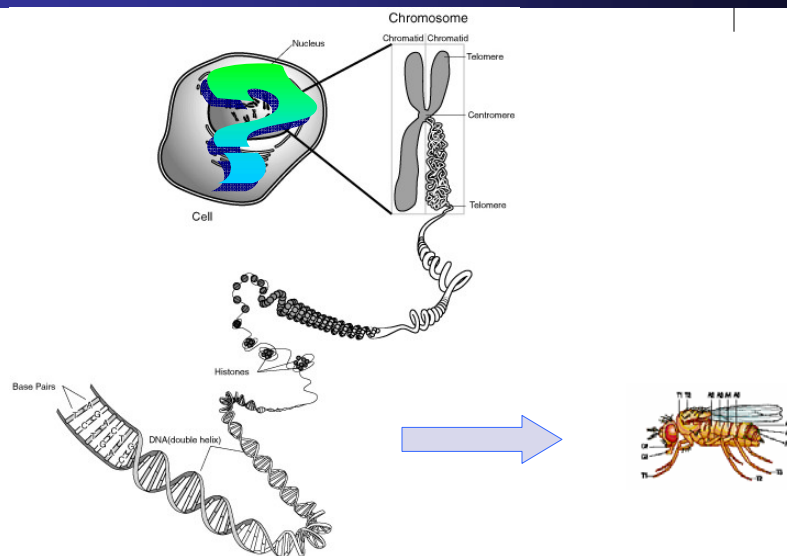
Reading: class assignment



© Eric Xing @ CMU, 2005-2009

1

DECODING the Genome



© Eric Xing @ CMU, 2005-2009

2

What & Why?

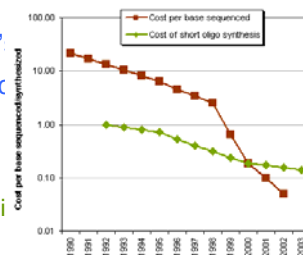
- “**Sequencing**” means finding the order of nucleotides on a piece of DNA
- Nucleotide order determines
 - amino acid composition, and by extension, protein structure and function (proteomics)
 - Transcription factor binding sites and their organizations, and thereby, gene expression regulation
 - ...
- An alteration in a DNA sequence can lead to
 - altered regulatory program,
 - altered protein structure/function,
 - and therefore, phenotypes or harmful effects in a plant or animal
- Understanding a particular DNA sequence can
 - shed light on a genetic condition and offer hope for the eventual development of treatment of diseases
 - environmental, agricultural applications
 - forensic applications
 - ...

© Eric Xing @ CMU, 2005-2009

3

DNA Sequencing – Overview

- Chain-termination methods (a.k.a. Dideox, termination method)
 - Predominant method based on Sanger’ (Nobel price in Chemistry, 1980, his sec Do you know how he got his first NP?)
 - Key ideas:
 - Terminate elongation via dideoxynucleoti
 - Size separation via Gel electrophoresis (now Capillary electrophoresis)
- Whole genome strategies
 - Physical mapping
 - Walking



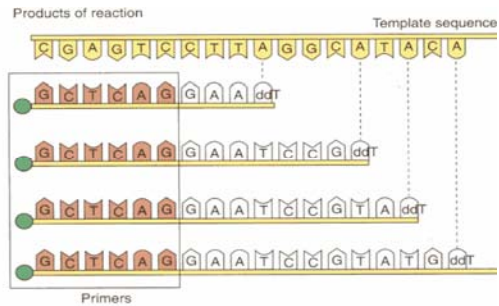
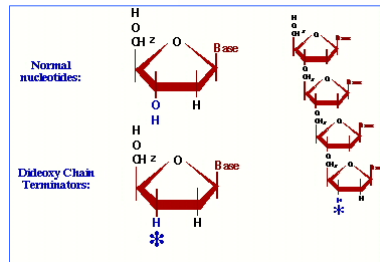
1975

2015

© Eric Xing @ CMU, 2005-2009

4

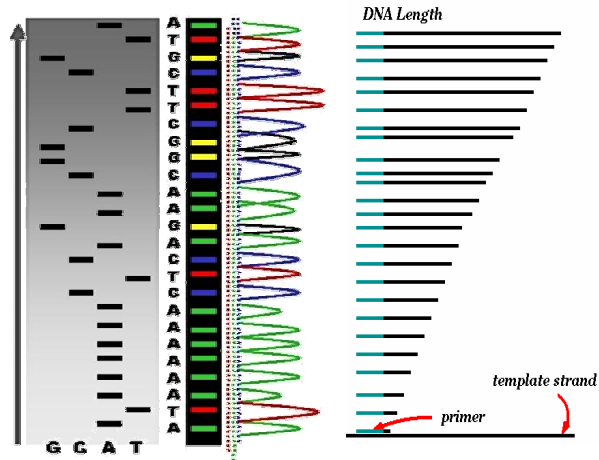
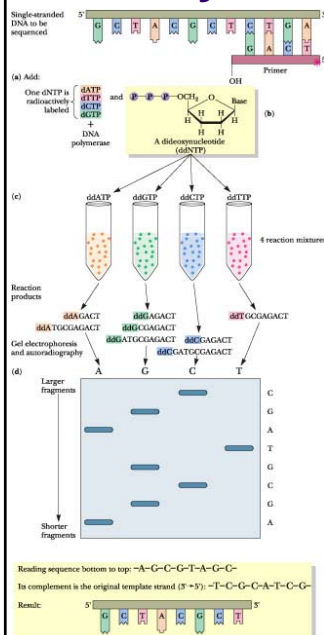
Dideoxy termination method



© Eric Xing @ CMU, 2005-2009

5

Dideoxy termination method



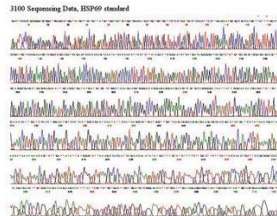
© Eric Xing @ CMU, 2005-2009

6

DNA Sequencing



- **Goal:**
Find the complete sequence of A, C, G, T's in DNA, at a Genome scale
- **Challenge:**
There is no machine that takes long DNA as an input, and gives the complete sequence as output
- Can only sequence ~500 letters at a time

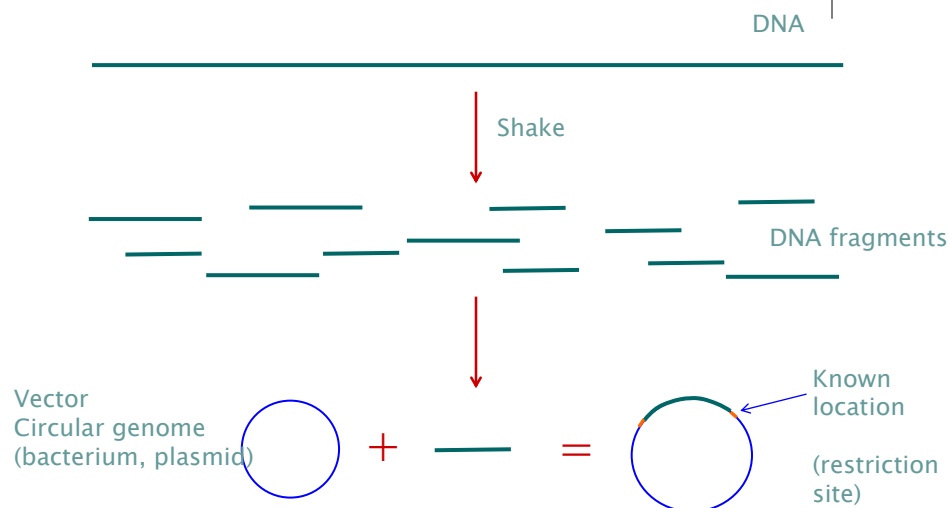


© Eric Xing @ CMU, 2005-2009



7

DNA Sequencing – vectors

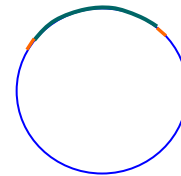


© Eric Xing @ CMU, 2005-2009

8

Different types of vectors

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



Different vector can be used to carry DNA pieces of different lengths

© Eric Xing @ CMU, 2005-2009

9

Method to sequence longer genome regions

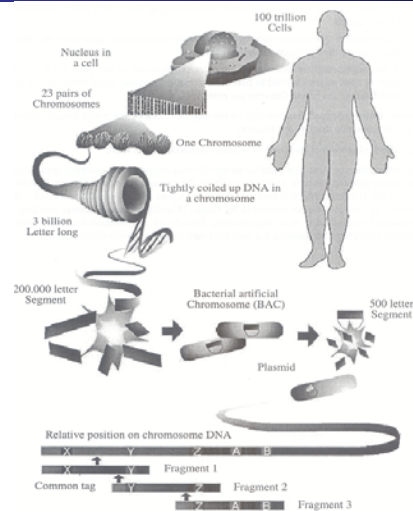


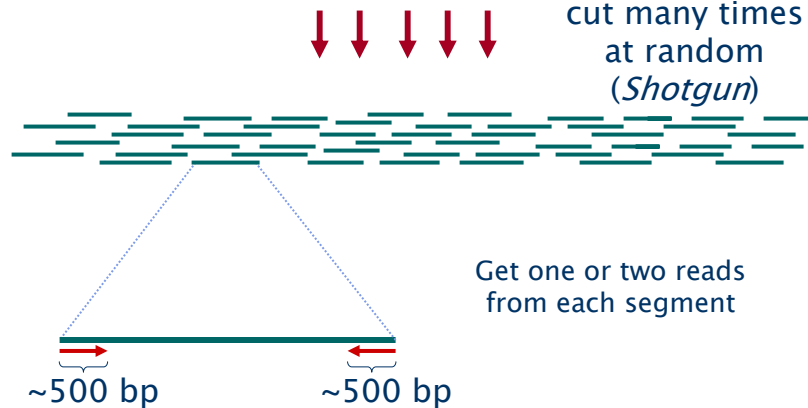
Fig (18) Methodology
© Eric Xing @ CMU, 2005-2009

10

Method to sequence longer genome regions



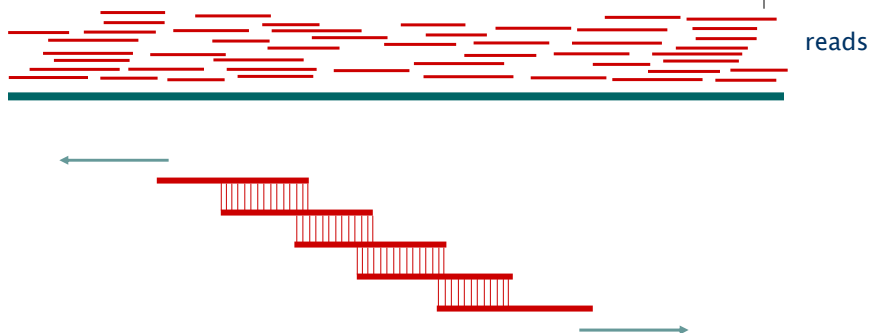
genomic segment



© Eric Xing @ CMU, 2005-2009

11

Reconstructing the Sequence (Fragment Assembly)



Cover region with ~7-fold redundancy (7X)

Overlap reads and extend to reconstruct the original genomic region

© Eric Xing @ CMU, 2005-2009

12

Definition of Coverage



Length of genomic segment: L

Number of reads: n

Length of each read: l

Definition: Coverage $C = n l / L$

How much coverage is enough?

Lander-Waterman model:

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

© Eric Xing @ CMU, 2005-2009

13

Are we done?



© Eric Xing @ CMU, 2005-2009

14

Difficulties

- How to fragment?
- How to efficiently find overlaps?
- What if there are **repeats** in the genome?



© Eric Xing @ CMU, 2005-2009

15

Repeats

- **Low-Complexity DNA** (e.g. ATATATATACATA...)
- **Microsatellite repeats** $(a_1 \dots a_k)^N$ where $k \sim 3-6$
(e.g. CAGCAGTAGCAGCACCAG)
- **Transposons**
 - **SINE** (Short Interspersed Nuclear Elements), e.g., ALU: ~300-long, 10^6 copies
 - **LINE** (Long Interspersed Nuclear Elements), ~4000-long, 200,000 copies
 - **LTR retroposons** (Long Terminal Repeats (~700 bp) at each end), cousins of HIV
- **Gene Families** genes duplicate & then diverge (paralogs)
- **Recent duplications** ~100,000-long, very similar copies

Bacterial genomes: 5%

Mammals: 50%

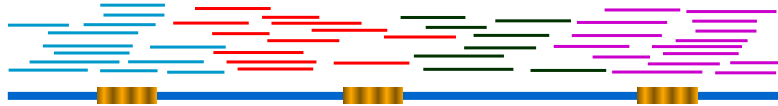
© Eric Xing @ CMU, 2005-2009

16

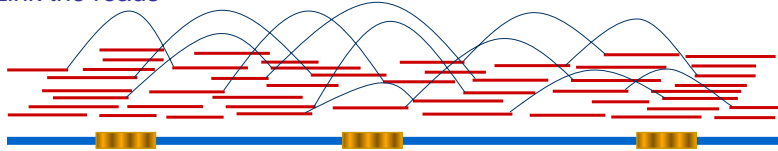
What can we do about repeats?

Two main approaches:

- Cluster the reads



- Link the reads



© Eric Xing @ CMU, 2005-2009

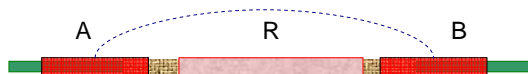
17

Sequencing and Fragment Assembly

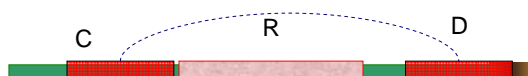


```
AGTAGCACAGA
CTACGACGAGA
CGATCGTGGGA
GCGACGGCGTA
GGTGGCTGTAC
TGTCTGTGTG
TGTACTCTCCT
```

3×10^9 nucleotides



ARB, CRD



or
~~ARD, CRB ?~~

© Eric Xing @ CMU, 2005-2009

18

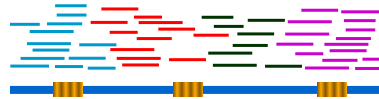
Strategies for whole-genome sequencing



1. Hierarchical – Clone-by-clone

- Break genome into many long pieces
- Map each long piece onto the genome
- Sequence each piece with shotgun

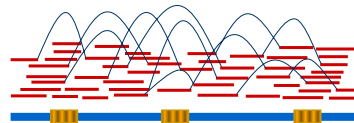
Example: Yeast, Worm, Human, Rat



2. Online version of (1) – Walking

- Break genome into many long pieces
- Start sequencing each piece with shotgun
- Construct map as you go

Example: Rice genome



3. Whole genome shotgun

One large shotgun pass on the whole genome

Example: Drosophila, Human (Celera),
Neurospora, Mouse, Rat, Dog

© Eric Xing @ CMU, 2005-2009

19

Whole Genome Shotgun



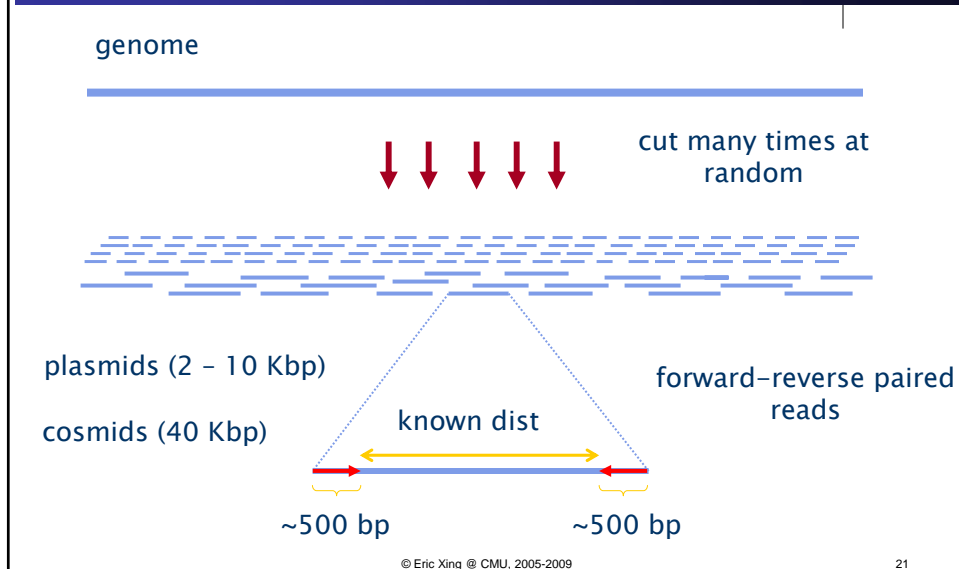
- Break up the entire genome into pieces
- Sequence ends, and assemble using a computer
- LW statistics & Repeats argue against the success of such an approach



© Eric Xing @ CMU, 2005-2009

20

Whole Genome Shotgun Sequencing



Steps to Assemble a Genome



Some Terminology

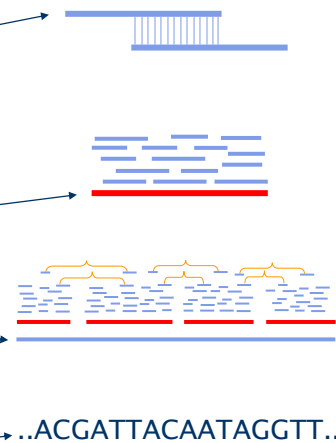
read a 500-900 long word that comes out of sequencer

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

supercontig (scaffold) an ordered and oriented set of contigs, usually by mate pairs

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

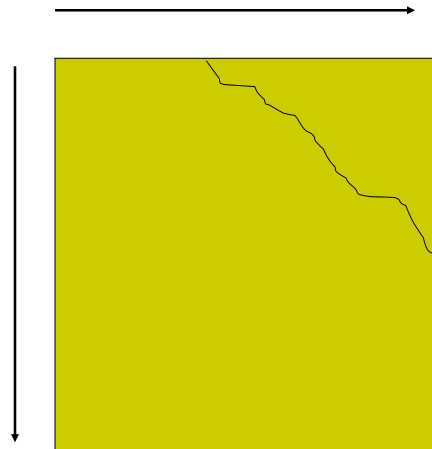


© Eric Xing @ CMU, 2005-2009

22

1. Find Overlapping Reads

- Given a pair of fragments s_1 and s_2 , do they belong together?
- Yes, if a prefix of s_2 matches a suffix of s_1
- How would you compute such a match?

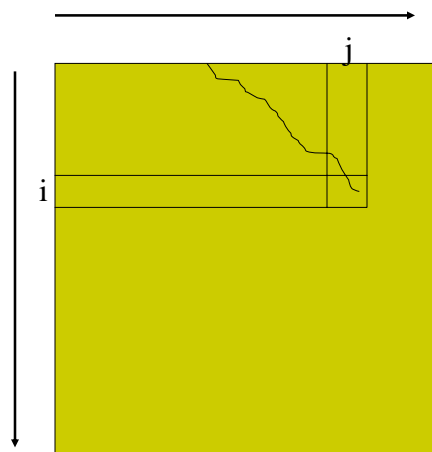


© Eric Xing @ CMU, 2005-2009

23

1. Find Overlapping Reads (cont'd)

- $S[i,j]$ = optimum score of an alignment of $s_1[1..i]$ against a suffix of $s_2[1..j]$
- The best prefix-suffix alignment is given by:
- $\text{Max}_i \{S[i,n]\}$



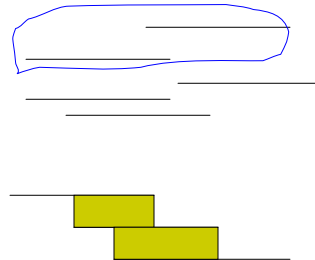
© Eric Xing @ CMU, 2005-2009

24

1. Find Overlapping Reads (cont'd)



- Compute the best prefix-suffix alignments between each pair of fragments.
- Keep the “high-scoring” ones as evidence of true overlap.
- What is the problem?



© Eric Xing @ CMU, 2005-2009

25

1. Find Overlapping Reads (cont'd)



- Consider the number of fragments. The LW statistics say that we need good coverage ($c=8, 10$) to get most of the base-pairs.
 - $G = 3000\text{Mb}$, $L=500$
 - Coverage $LN/G = 10$
 - $N = 10^3 \cdot 10^9 / 500 = 6 \cdot 10^7$
 - Number of comparisons needed = $3.6 \cdot 10^{15}$
 - Not good! (Only a small fraction are true overlaps)

© Eric Xing @ CMU, 2005-2009

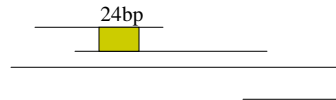
26

1. Find Overlapping Reads (cont'd)



k-mer based overlap

- Consider a k -bp sequence ($k \sim 24$)
 - Expected number of occurrences in the genome
 - $3 \cdot 10^9 \cdot 4^{-24} = 8 \cdot 10^{-6}$
- A 24-bp sequence appears is unique to the genome!
- Two overlapping sequences should share a 24-mer
- Two non-overlapping



© Eric Xing @ CMU, 2005-2009

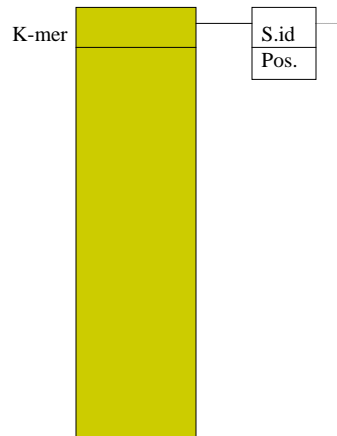
27

1. Find Overlapping Reads (cont'd)



Sorting k-mers

- Build a list of k-mers that appear in the sequences and their reverse complements
- Create a record with 4 entries:
 - K-mer
 - Sequence number
 - Position in the sequence
 - Reverse complementation flag
- Sort a vector of these according to k-mer
- How many records per k-mer are expected?
- If number of records exceeds threshold, discard (why?)



© Eric Xing @ CMU, 2005-2009

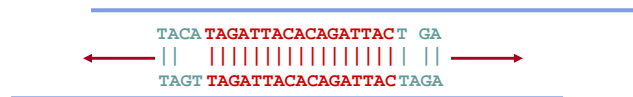
28

1. Find Overlapping Reads (cont'd)



Alignment

- Find pairs of reads sharing a k-mer, $k \sim 24$
- Extend to full alignment – throw away if not >98% similar



- Coalesce k-mer hits into longer, gap-free partial alignments.
- These extended k-mer hits are saved.
- For each pair of sequences, form a directed graph.
- For each maximal path in the graph, construct an alignment.
- Refine alignment via banded DP

© Eric Xing @ CMU, 2005-2009

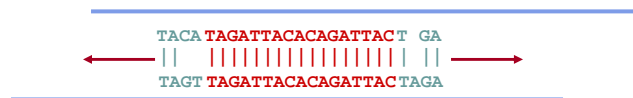
29

1. Find Overlapping Reads (cont'd)



Alignment

- Find pairs of reads sharing a k-mer, $k \sim 24$
- Extend to full alignment – throw away if not >98% similar



- **Caveat: repeats**
 - A k-mer that occurs N times, causes $O(N^2)$ read/read comparisons
 - ALU k-mers could cause up to 1,000,000 comparisons
- **Solution:**
 - Discard all k-mers that occur "too often"
 - Set cutoff to balance sensitivity/speed tradeoff, according to genome at hand and computing resources available

© Eric Xing @ CMU, 2005-2009

30

1. Find Overlapping Reads (cont'd)



Create local multiple alignments from the overlapping reads



© Eric Xing @ CMU, 2005-2009

31

1. Find Overlapping Reads (cont'd)



- Correct errors using multiple alignment

```
TAGATTACACAGATTACTGA
TAGATTACACAGATTACTGA
TAGATTACACAGATTACTGA
TAGATTACACAGATTACTGA
TAGATTACACAGATTACTGA
```

insert A

replace T with C

```
TAGATTACACAGATTACTGA
TAGATTACACAGATTACTGA
TAGATTACACAGATTACTGA
TAGATTACACAGATTACTGA
TAGATTACACAGATTACTGA
```

correlated errors—
probably caused by repeats
⇒ disentangle overlaps

```
TAGATTACACAGATTACTGA
TAGATTACACAGATTACTGA
TAGATTACACAGATTACTGA
```

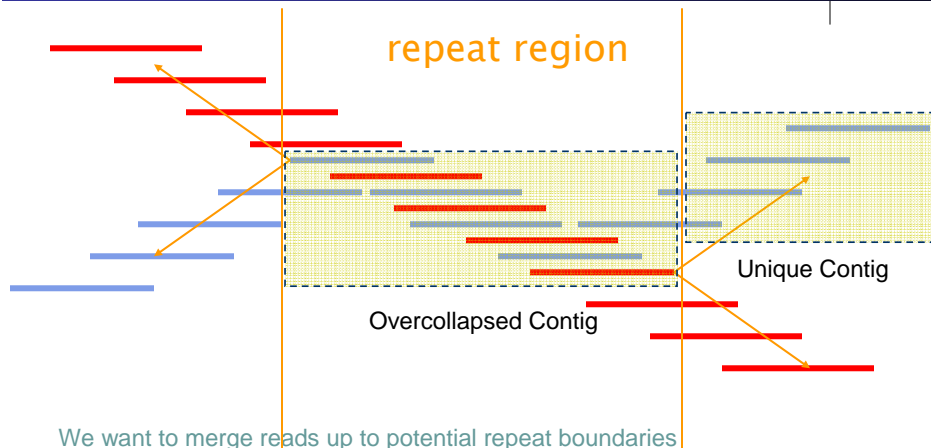
In practice, error correction removes
up to 98% of the errors

```
TAG-TTACACAGATTATTGA
TAG-TTACACAGATTATTGA
```

© Eric Xing @ CMU, 2005-2009

32

2. Merge Reads into Contigs



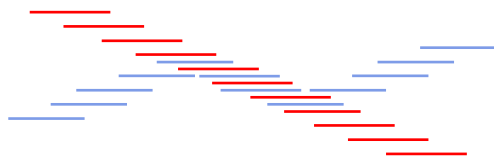
© Eric Xing @ CMU, 2005-2009

33

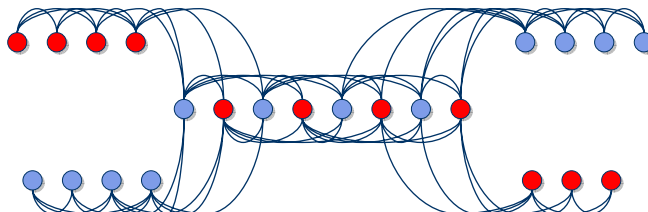
2. Merge Reads into Contigs



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



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

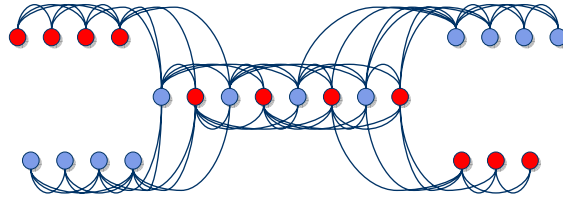


Note:
of course, we don't know the "color" of these nodes

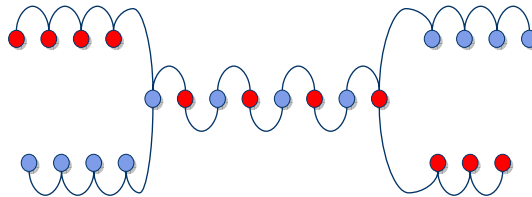
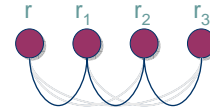
© Eric Xing @ CMU, 2005-2009

34

2. Merge Reads into Contigs



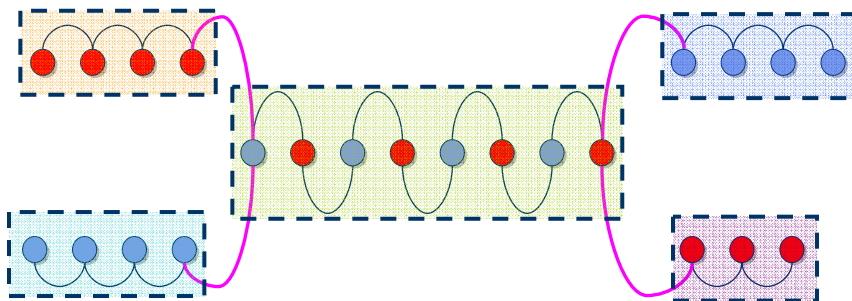
- Remove transitively inferable overlaps
 - If read r overlaps to the right reads r_1 , r_2 , and r_1 overlaps r_2 , then (r, r_2) can be inferred by (r, r_1) and (r_1, r_2)



© Eric Xing @ CMU, 2005-2009

35

2. Merge Reads into Contigs



© Eric Xing @ CMU, 2005-2009

36

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

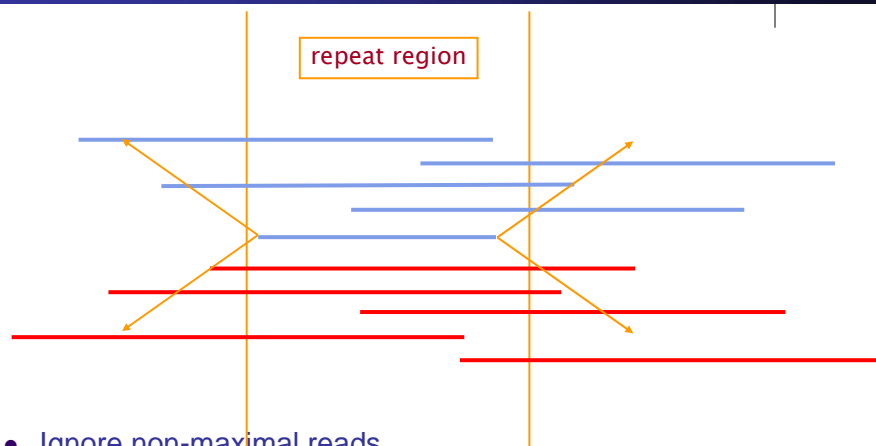
Role of error correction:

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

© Eric Xing @ CMU, 2005-2009

37

2. Merge Reads into Contigs

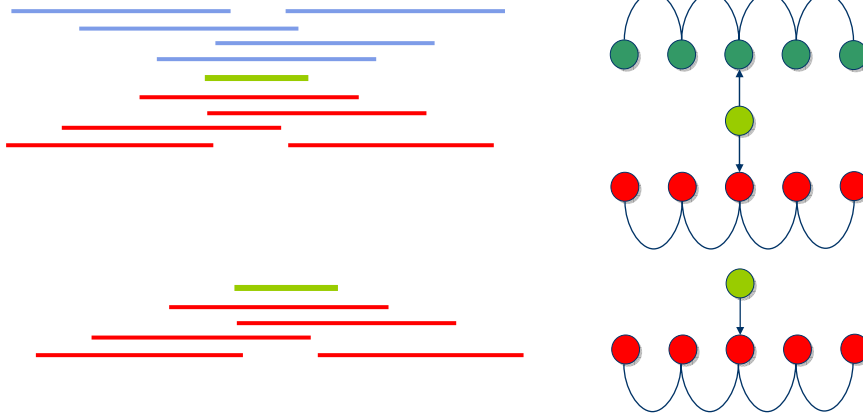


- Ignore non-maximal reads
- Merge only maximal reads into contigs

© Eric Xing @ CMU, 2005-2009

38

2. Merge Reads into Contigs

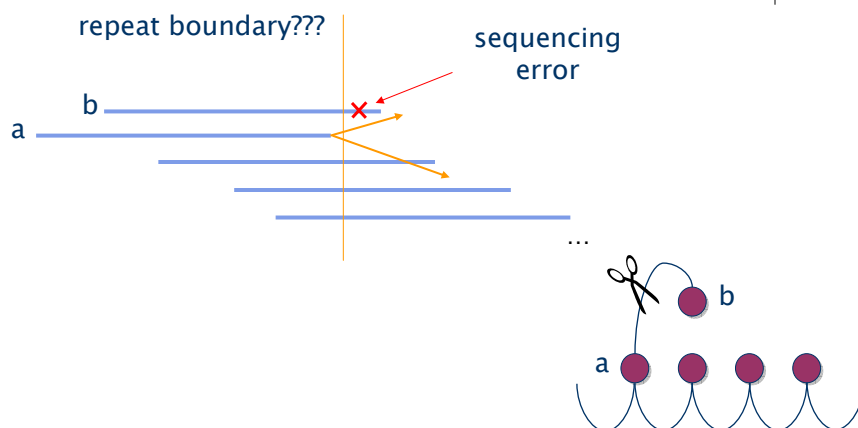


- Insert non-maximal reads whenever unambiguous

© Eric Xing @ CMU, 2005-2009

39

2. Merge Reads into Contigs

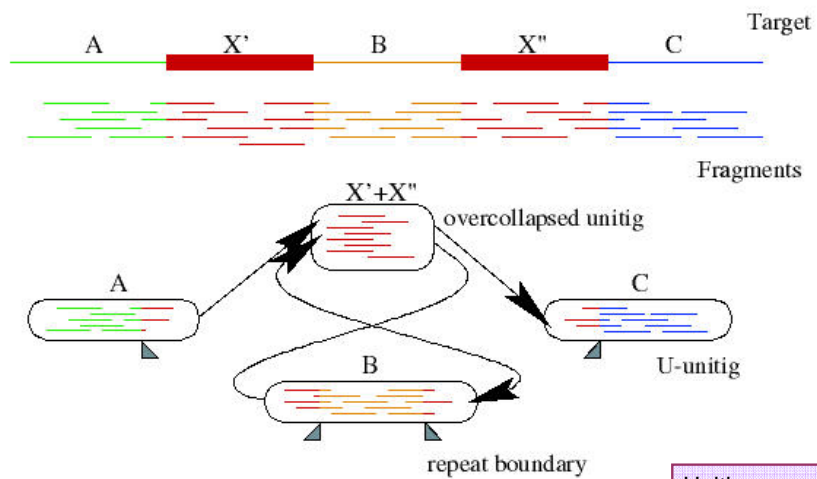


- Ignore “hanging” reads, when detecting repeat boundaries

© Eric Xing @ CMU, 2005-2009

40

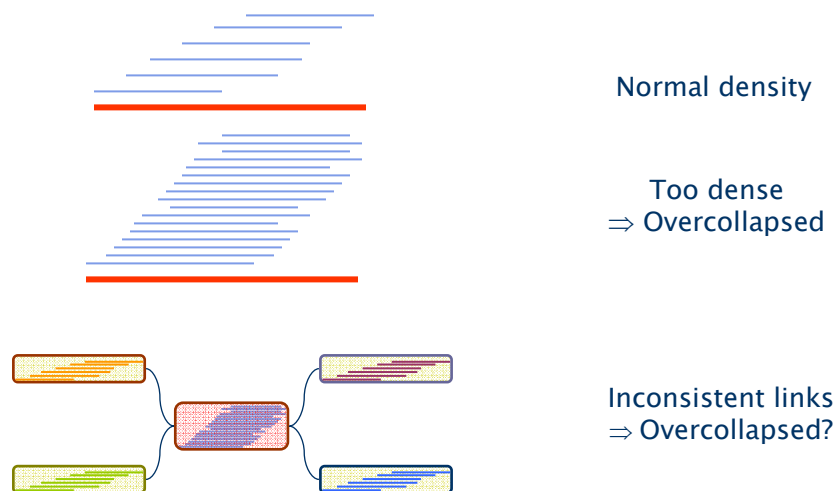
Overlap graph after forming contigs



Unitigs:
Gene Myers, 95

© Eric Xing @ CMU, 2005-2009

3. Link Contigs into Supercontigs



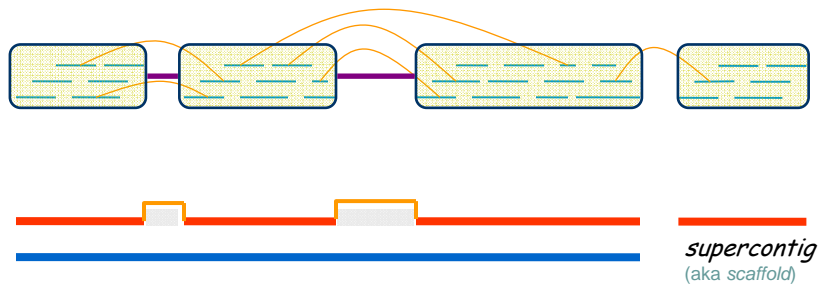
© Eric Xing @ CMU, 2005-2009

42

3. Link Contigs into Supercontigs



- Find all links between unique contigs
- Connect contigs incrementally, if ≥ 2 links



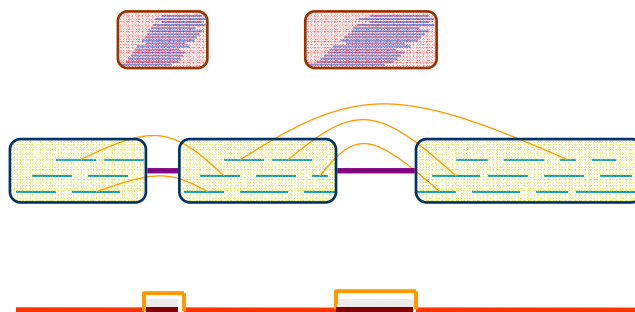
© Eric Xing @ CMU, 2005-2009

43

3. Link Contigs into Supercontigs



- Fill gaps in supercontigs with paths of repeat contigs



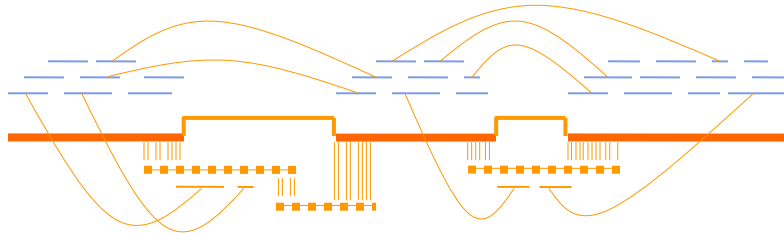
© Eric Xing @ CMU, 2005-2009

44

3. Link Contigs into Supercontigs



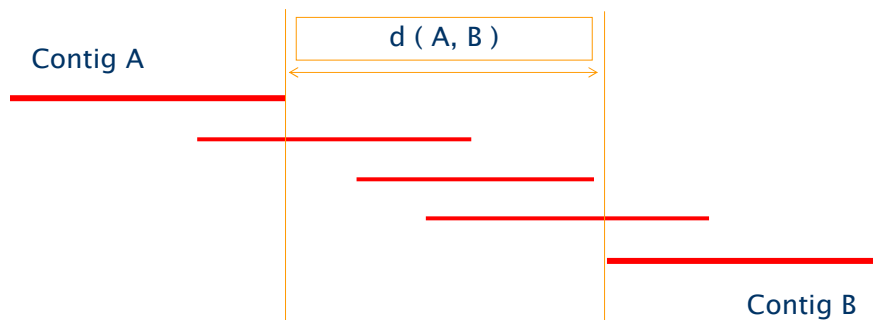
- Fill gaps in supercontigs with paths of repeat contigs



© Eric Xing @ CMU, 2005-2009

45

3. Link Contigs into Supercontigs



Define $G = (V, E)$

$V :=$ contigs

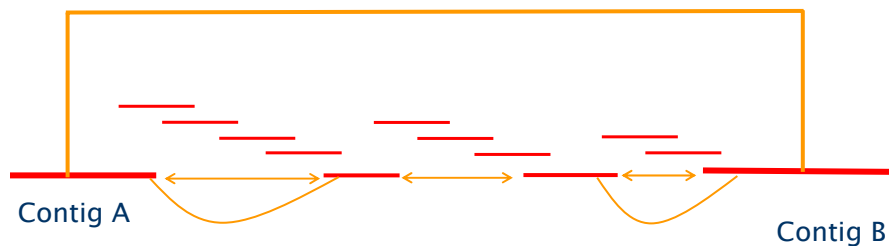
$E := (A, B)$ such that $d(A, B) < C$

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

© Eric Xing @ CMU, 2005-2009

46

3. Link Contigs into Supercontigs



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

© Eric Xing @ CMU, 2005-2009

47

4. Derive Consensus Sequence



```
TAGATTACACAGATTACTGA TTGATGGCGTAA CTA
TAGATTACACAGATTACTGACTTGATGGCGTAACTA
TAG TTACACAGATTATGACTTCATGGCGTAA CTA
TAGATTACACAGATTACTGACTTGATGGCGTAA CTA
TAGATTACACAGATTACTGACTTGATGGGGTAA CTA
```

↓ ↓ ↓ ↓ ↓

```
TAGATTACACAGATTACTGACTTGATGGCGTAA CTA
```

- Derive **multiple alignment** from pairwise read alignments
- Derive each consensus base by weighted voting
- (Alternative: take maximum-quality letter)

© Eric Xing @ CMU, 2005-2009

48

Simulated Whole Genome Shotgun

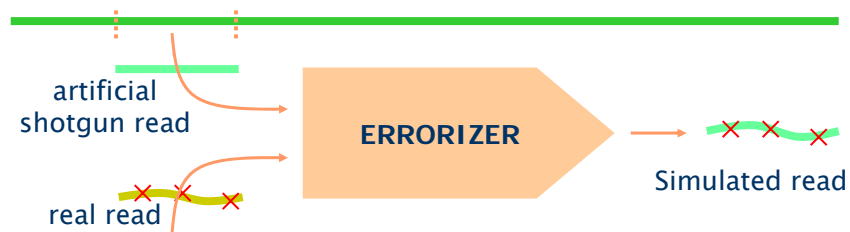


- Known genomes
Flu, yeast, fly, Human chromosomes 21, 22
- Make “realistic” shotgun reads
- Run assembly program
- Align output with genome and compare

© Eric Xing @ CMU, 2005-2009

49

Making a Simulated Read

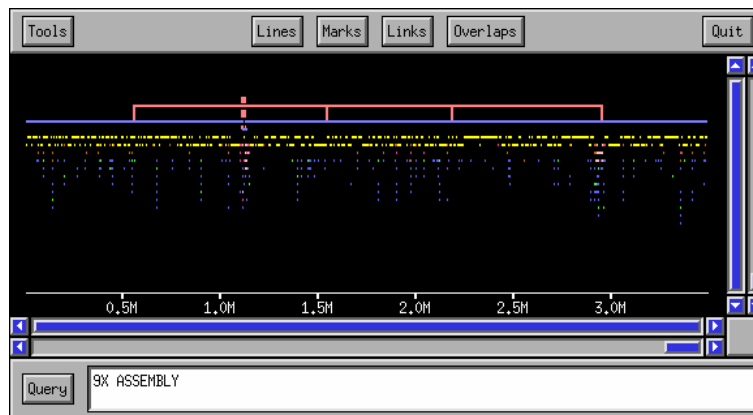
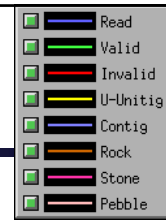


Simulated reads have error patterns taken from random real reads

© Eric Xing @ CMU, 2005-2009

50

Assembly Progression (Macro View)



© Eric Xing @ CMU, 2005-2009

51

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

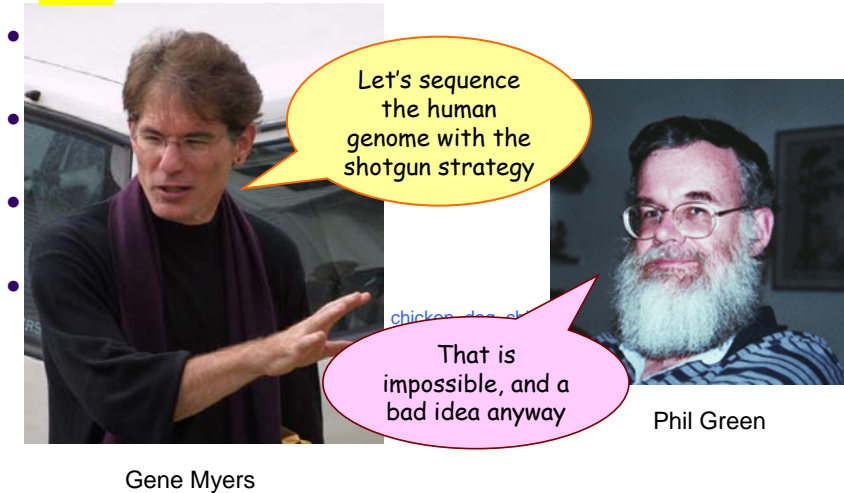
© Eric Xing @ CMU, 2005-2009

52

History of WGA



1997



Gene Myers

Phil Green

© Eric Xing @ CMU, 2005-2009

53

New Sequencing Methods



- Sequencing by MALDI-TOF Mass Spectrometry
- Sequencing by Hybridization
- Pyrosequencing
- Atomic-Force Microscopy
- Single-Molecule Fluorescence Microscopy
- Nanopore Sequencing

© Eric Xing @ CMU, 2005-2009

54

Some future directions for sequencing



1. Personalized genome sequencing

- Find your ~1,000,000 single nucleotide polymorphisms (SNPs)
- Find your rearrangements
- Goals:
 - Link genome with phenotype
 - Provide personalized diet and medicine
 - (???) designer babies, big-brother insurance companies
- Timeline:
 - Inexpensive sequencing: 2010-2015
 - Genotype–phenotype association: 2010-???
 - Personalized drugs: 2015-???

© Eric Xing @ CMU, 2005-2009

55

Some future directions for sequencing



2. Environmental sequencing

- Find your flora: organisms living in your body
 - External organs: skin, mucous membranes
 - Gut, mouth, etc.
- Normal flora: >200 species, >trillions of individuals
- Flora–disease, flora–non-optimal health associations
- Timeline:
 - Inexpensive research sequencing: today
 - Research & associations: within next 10 years
 - Personalized sequencing: 2015+
- Find diversity of organisms living in different environments
 - Hard to isolate
 - Assembly of all organisms at once

© Eric Xing @ CMU, 2005-2009

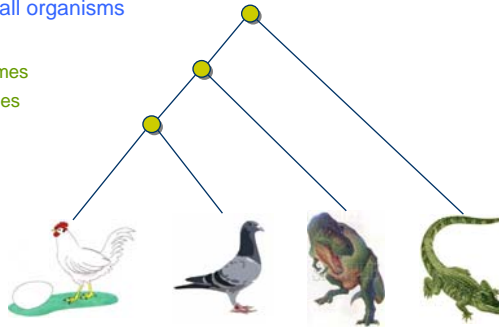
56

Some future directions for sequencing



3. Organism sequencing

- Sequence a large fraction of all organisms
- Deduce ancestors
 - Reconstruct ancestral genomes
 - *Synthesize* ancestral genomes
 - Clone—Jurassic park!
- Study evolution of function
 - Find functional elements within a genome
 - How those evolved in different organisms
 - Find how modules/machines composed of many genes evolved



© Eric Xing @ CMU, 2005-2009

57

Acknowledgment



- This set of slides is based on materials from
 - Serafim Batzoglou @ stanford
 - Sorin Istrail @ Brown
 - Vineet Bafna @UCSD

© Eric Xing @ CMU, 2005-2009

58