

Course summary

- Goals
- Organization
- Labs
- Project
- Reading

Today

- DNA sequencing.
- Genome Projects.
- New DNA sequencing technologies.

Obtaining molecular data

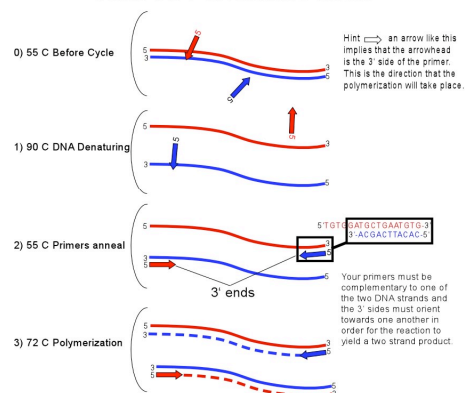
- PCR
 - Typically used in empirical molecular evolution studies
- Sequencing
 - Sequence chromatograms and base calling
- Genome sequencing (data used by Mol. Evol.)
 - Shotgun versus Hierarchical
- Cloning vectors
 - Plasmids and BACs
- The human genome

PCR

Polymerase chain reaction

- Template DNA
- Primers (must know sequence)
- dNTPs (dCTP, dGTP, dATP, dTTP)
- Polymerase
- Buffer
- $MgCl_2$

Inside the PCR reaction tube...



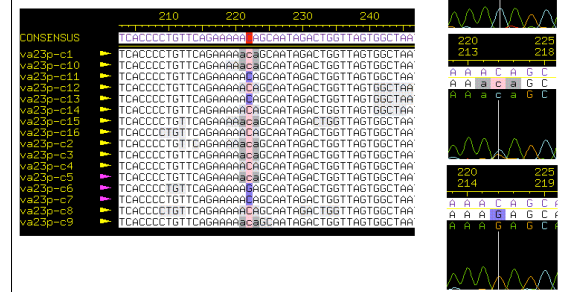
Sequencing

- Template DNA
- Primers
- dNTPs (dCTP, dGTP, dATP, dTTP)
- **Labeled ddNTPs**
- Polymerase
- Buffer
- $MgCl_2$

What about polymorphism in sequence data?

- Directly sequenced PCR products can be from diploid individuals with potential for heterozygotes.

SNP detection by PolyPhred. View of a Consed window with a tag (red=highest ranking SNP tag) marking the consensus position of the SNP in the traces and genotype tags marking each of the samples below (purple=homozygote, pink=heterozygote). On the right trace windows for alternate homozygotes (C/C (top) and G/G (bottom)) and a heterozygote (C/G (middle)).



Genome Projects

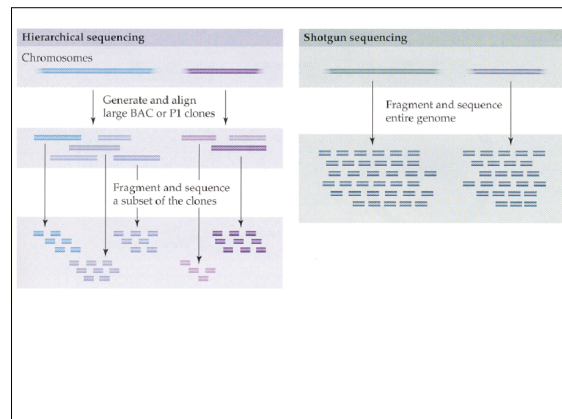
- Drosophila
- C. elegans
- Human, Chimpanzee
- Mouse, Rat
- Dog
- Plants
- Cow, Rhesus,

Goals of Human Genome Project

- Technology development
 - High throughput, inexpensive DNA sequencing
- Complete Sequence
 - Provide the genetic sequence
- Gene Identification
 - Identify the “parts” list. Also non-coding functional regions
- Polymorphism Data
 - Association studies, evolutionary studies
- Functional characterization
 - ENCODE

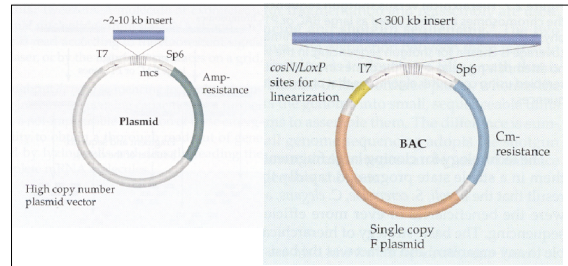
Two strategies for genome sequencing

- Whole Genome Shotgun
 - Random reads assembled by computer
- Hierarchical
 - Map bases, clone by clone



Genome sequencing Terminology

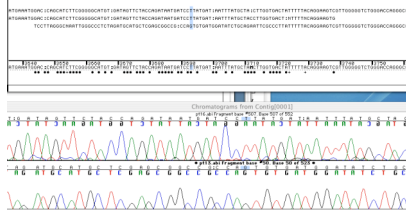
- Coverage
 - Estimated number of times a region of the genome is sequenced (typically 5 - 10 fold). Note on distribution.
- Consensus Sequence
 - Sequence based upon multiple reads
- Contig
 - Contiguous sequence aligned
- Scaffold
 - Ordered contigs based upon BAC-end sequencing or other knowledge



Plasmid multiple copy to allow easy purification large quantities
BACs single copy to prevent recombination

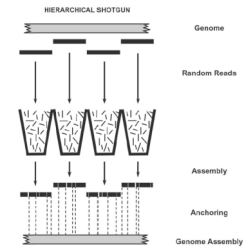
Sequencing clones

- Usually get better sequence than sequencing PCR products directly.
- Must remove vector sequence.
- Sometimes clones are chimeric



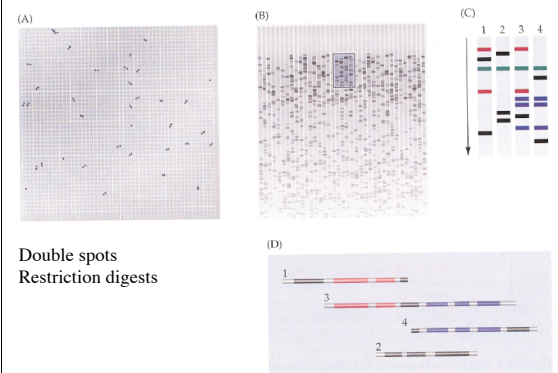
Hierarchical sequencing

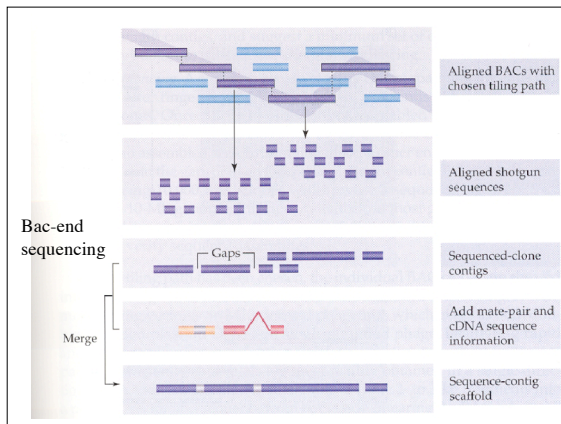
- Break genome into smaller aligned chunks and sequence.



Aligning BACs

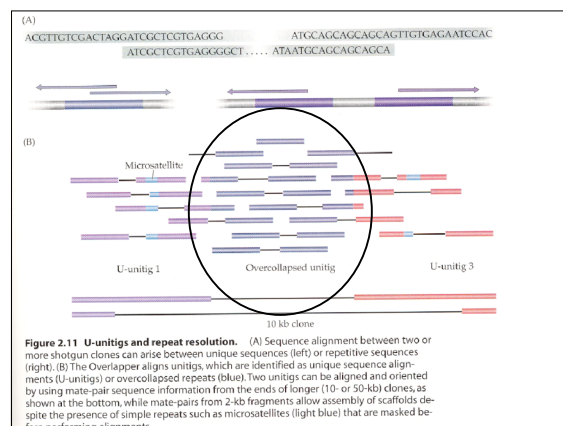
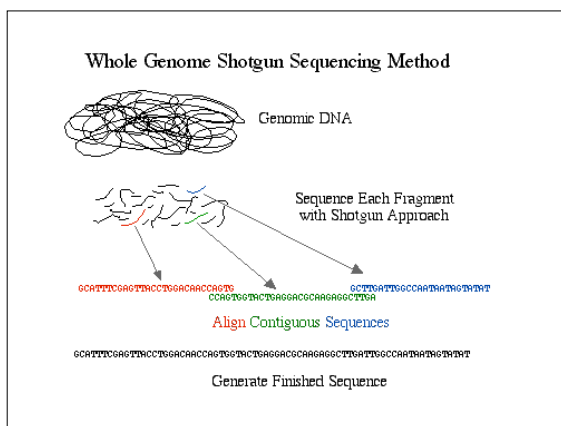
- Hybridize arrayed library with unique sequence.
- Isolate all hybridizing BACs.
- Perform restriction based mapping.
- Align BACs and make tilling path.





Shotgun sequencing

- Screen sequence for known repeats
 - Note there will still be duplications
- Finding overlapping sequence
- Link contigs into scaffolds
- Fill in gaps by BAC sequencing

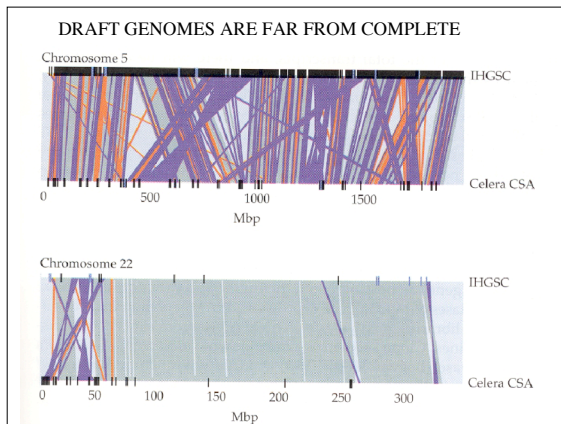


Sequence verification

- Completeness
 - Repetitive regions, heterochromatin
- Accuracy
 - Based upon probability scores for base calls (Phred scores)
- Validity of assembly
 - Difficult to assess. Predict restriction profiles and compare to observed. Correct spacing of pair-end sequences (correct for repeats), consistent read depth.

Is there any difference between two draft human genomes?

- International public consortium
 - Used hierarchical sequencing
 - Published in Nature 2001
- Private company Celera
 - Used Whole Genome Shotgun
 - Published in Science 2001



How effective was the whole genome shotgun approach?

Celera spiked their data with “shredded” reads from the human genome project

Celera used “shredded” reads from public data.
This was not done randomly as would be obtained from WGS

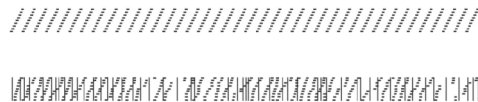


Fig. 2. Random vs. perfect spacing in 100 kb. (Upper) Two-fold coverage in perfectly spaced reads. (Lower) Two-fold coverage in randomly selected reads. There are frequent regions in which the adjacent reads either fail to overlap or the overlap is too small to allow reliable detection (< 40 bp). These breaks in continuity are indicated by vertical lines.

From Waterston et al. 2002 *PNAS* **99**, 3712 - 3716

Table 1. Reported statistics for genome assemblies in the HGP and Celera papers

Category	WGS	Celera		HGP
		Faux WGS	Faux CSA	
Sequence coverage	5.1 × Celera	5.1 × Celera + 7.5 × HGP	5.1 × Celera + 7.5 × HGP	7.5 × HGP
Length (in Gb) of draft genome assembly, counting only bases with known sequence*	NR	12.6 × total 2,587	12.6 × total 2,654	2,693
Length (in Gb) of draft genome assembly, including unknown nucleotides in gaps†	NR	2,848	2,906	2,916
Proportion of sequence in euchromatic genome present in draft genome assembly, %‡	NR	89	91	92
Number of contigs§	NR	221,036	170,833	149,821
Number of scaffolds¶	NR	118,968	53,591	87,757
Number of components, to be anchored in genome¶¶	NR	118,968	3,845	942

WGS: Whole genome shotgun

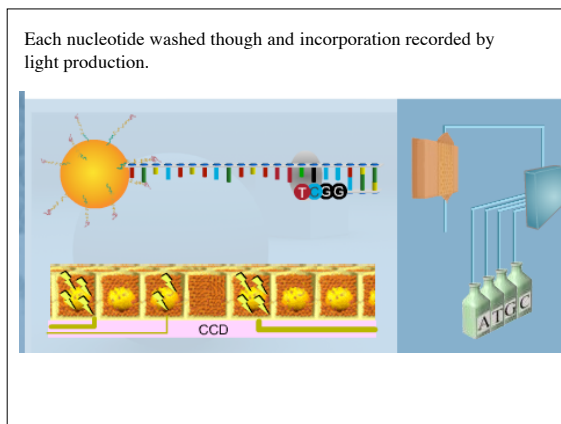
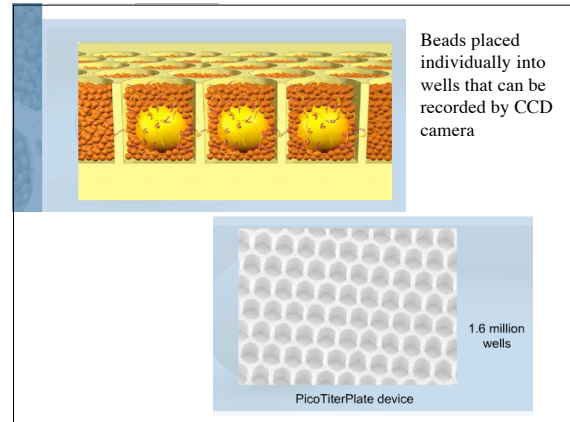
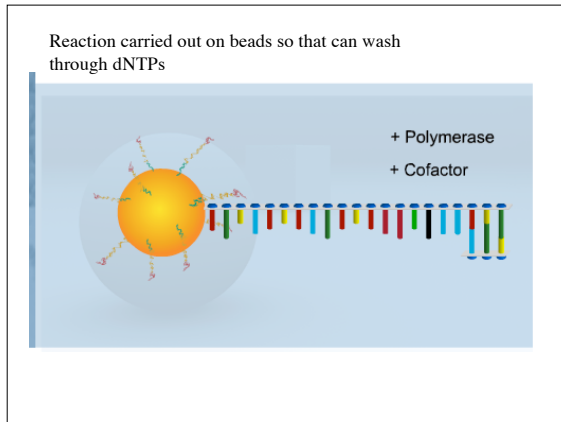
CSA: Compartmentalized sequence assembly

HGP: Human genome Project (hierarchical sequencing)

New sequencing technology Pyrosequencing (454)

- 4 enzymes
 - DNA polymerase
 - DNA synthesis
 - ATP sulfurylase
 - Converts PPi to ATP
 - Luciferase
 - Uses ATP to generate light
 - Apyrase
 - Degrades dNTPs
- Cheaper, higher throughput
 - 20 million bases per run, versus 67K
 - Moving towards the \$1,000 genome

Show chemistry



Some examples of pyrosequencing

- Marine Microbes
 - Large amount of sequencing to get “unculturable” microbes
- Extinct cave bears
 - Test case for Neanderthals
- Neanderthals
 - Small amount of highly degraded DNA, lots of microbe contamination, need to sequence a lot.

Disadvantages

- Short reads
 - Only 100 bps compared to 800+
- Homo-polymer runs difficult
 - Beyond 2 or 3, can not quantify
- Needs for reference genome
 - Reads potentially too short for shotgun assembly. No possibility of clone based sequencing
- Initial investment high
 - Instrument cost.

1. Describe DNA sequencing.
2. What is the difference between a BAC and Plasmid?
3. What are two methods for sequencing genomes?
4. What is hierarchical sequencing?
5. What is whole genome shotgun sequencing?
6. Describing a tilling path.
7. What is a phred score?
8. What is the difference between a scaffold and contig sequence?
9. What is the difference between dNTP and ddNTP?
10. Describe pyrosequencing? What is advantage/disadvantage?