# Contig Assembly

ATCGATGCGTAGCAGACTACCGTTACGATGCCTT...
TAGCTACGCATCGTCTGATGGCAATGCTACGGAA..

TAGCTACGCATCGT
ATCGATGCGTAGC
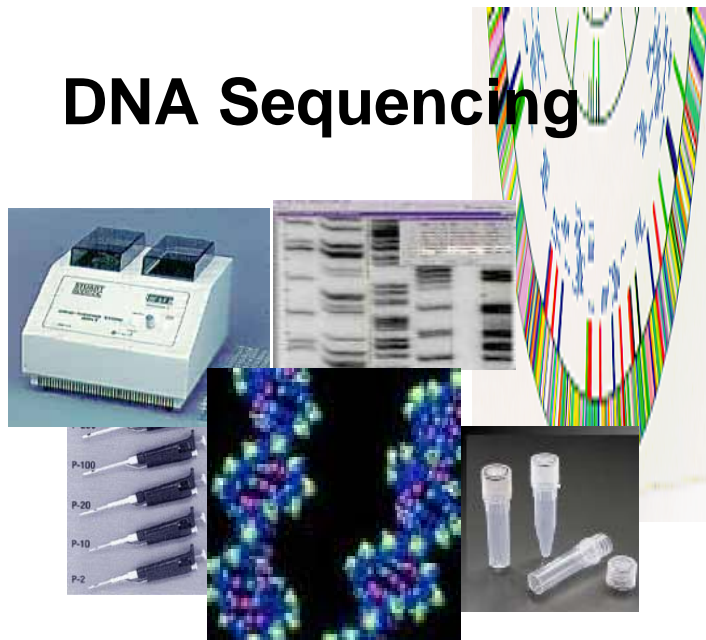TAGCAGACTACCGTT
GTTACGATGCCTT

**David Wishart, Ath 3-41**
**david.wishart@ualberta.ca**

# DNA Sequencing

# Principles of DNA Sequencing

**DNA fragment**

**Amp**

**PBR322**

**Tet**

**Ori**

**Denature with heat to produce ssDNA**
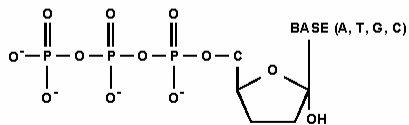
**Primer**

**Klenow + ddNTP + dNTP + primers**

---

# The Secret to Sanger Sequencing

• Structure of the di_deoxynucleotide

BASE (A, T, G, C)

• structure of a dNTP

BASE (A, T, G, C)

• structure of a ddNTP

# Principles of DNA Sequencing

5'  G  C  A  T  G  C                          3' Template

5' Primer

dATP
dCTP
dGTP
dTTP
ddCTP

dATP
dCTP
dGTP
dTTP
ddATP

dATP
dCTP
dGTP
dTTP
ddTTP

dATP
dCTP
dGTP
dTTP
ddGTP

GddC

GCddA

GCAddT

ddG

GCATGddC

GCATddG

# Principles of DNA Sequencing

G          T

C          A

_          _ short

G
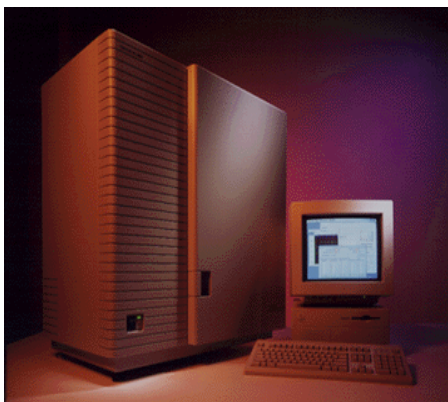C
A
T
G
C

+          +  long

# Capillary Electrophoresis
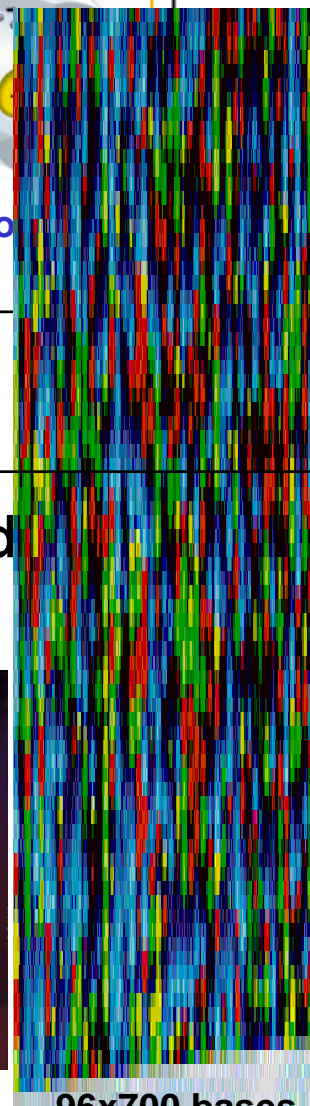


**Separation by Electro**

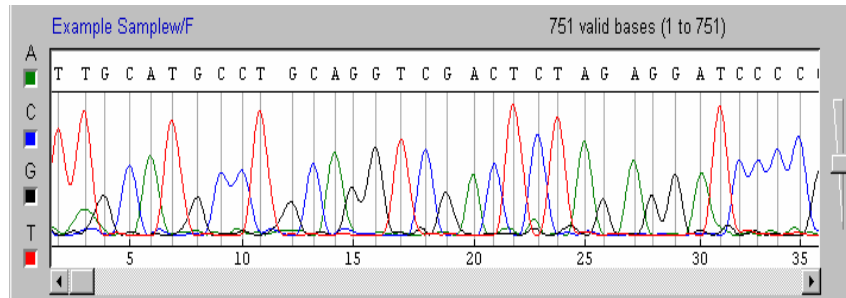# Multiplexed Fluorescent



**ABI 3700**          **96x700 bases**
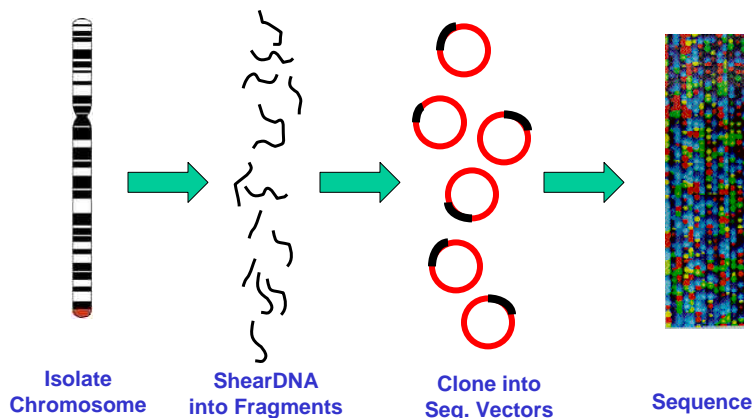
# High Throughput DNA Sequencing



---

# Large Scale Sequencing

- **Goal is to determine the nucleic acid sequence of molecules ranging in size from a few hundred bp to >$10^9$ bp**

- **The methodology requires an extensive computational analysis of raw data to yield the final sequence result**
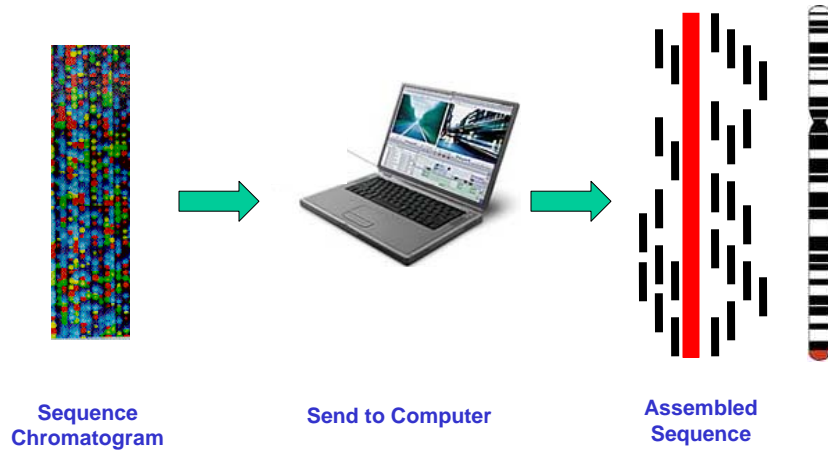
# Shotgun Sequencing

- **High throughput sequencing method that employs automated sequencing of random DNA fragments**
- **Automated DNA sequencing yields sequences of 500 to 1000 bp in length**
- **To determine longer sequences you obtain fragmentary sequences and then join them together by overlapping**
- **Overlapping is an alignment problem, but different from those we have discussed up to now**

# Shotgun Sequencing

| Isolate Chromosome | ShearDNA into Fragments | Clone into Seq. Vectors | Sequence |
|---|---|---|---|

# Shotgun Sequencing

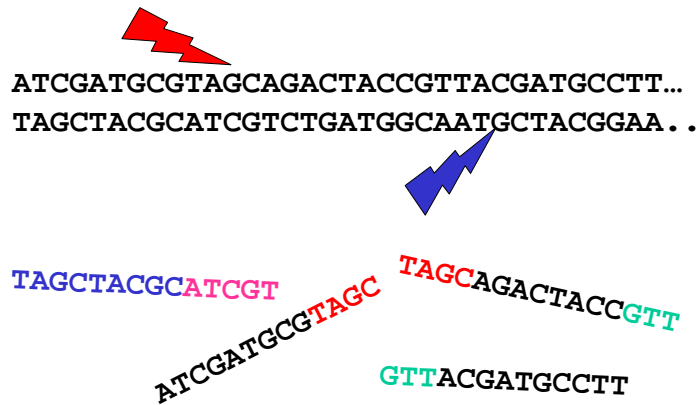**Sequence Chromatogram**  **Send to Computer**  **Assembled Sequence**

# Analogy

- **You have 10 copies of a movie**
- **The film has been cut into short pieces with about 240 frames per piece (10 seconds of film), at random**
- **Reconstruct the film**

# Multi-alignment & Contig Assembly

ATCGATGCGTAGCAGACTACCGTTACGATGCCTT...
TAGCTACGCATCGTCTGATGGCAATGCTACGGAA..

TAGCTACGCATCGT

ATCGATGCGTAGC

TAGCAGACTACCGTT

GTTACGATGCCTT

---

# Multiple Sequence Alignment

| Consensus: | CSNLSTCVLGKLSQDLHKLQTFPRT--GAG-P |
|---|---|
| 1: sockeye | CSNLSTCVLGKLSQDLHKLQTFPRTNTGAGVP |
| 2: chum | CSNLSTCVLGKLSQDLHKLQTFPRTNTGAGVP |
| 3: pink | CSNLSTCVLGKLSQDLHKLQTFPRTNTGAGVP |
| 4: coho | CSNLSTCMLGKLSQDLHKLQTFPRTNTGAGVP |
| 5: pig | CSNLSTCVLSAYWRNLNNFHRFSGMGFGPETP |
| 6: bovine | CSNLSTCVLSAYWKDLNNYHRFSGMGFGPETP |
| 7: eel | CSNLSTCVLGKLSQELHKLQTYPRTDVGAGTP |

**Multiple alignment of Calcitonins**

8

# Multiple Sequence Alignment

- **A general method to align and compare more than 2 sequences**
- **Typically done as a hierarchical clustering/alignment process where you match the two most similar sequences and then use the combined consensus sequence to identify the next closest sequence with which to align**

# Multiple Alignment Algorithm

- *Take all "n" sequences and perform all possible pairwise (n/2(n-1)) alignments*
- *Identify highest scoring pair, perform an alignment & create a consensus sequence*
- *Select next most similar sequence and align it to the initial consensus, regenerate a second consensus*
- *Repeat step 3 until finished*

# Multiple Sequence Alignment

- **Developed and refined by many (Doolittle, Barton, Corpet) through the 1980's**
- **Used extensively for extracting hidden phylogenetic relationships and identifying sequence families**
- **Powerful tool for extracting new sequence motifs and signature sequences**
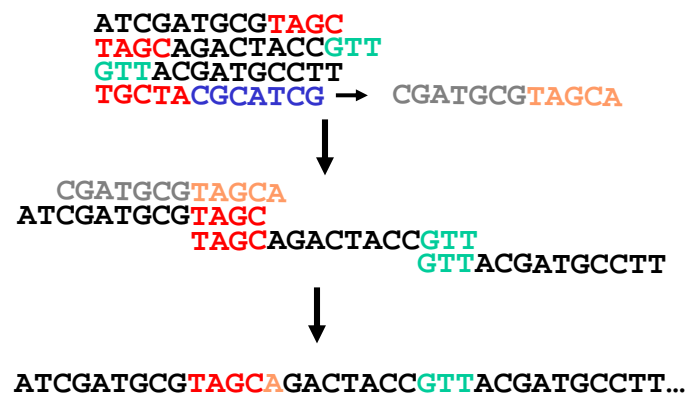- **Also applicable to DNA contig assembly**

# Contig Assembly ⇌ Multiple Alignment

1. **Only accept a very high sequence identity**
2. **Accept unlimited number of "end" gaps**
3. **Very high cost for opening "internal" gaps**
4. **A short match with high score/residue is preferred over a long match with low score/residue**
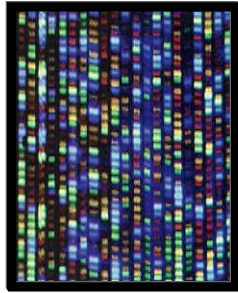
# Contig Assembly Algorithm

- *Read, edit & trim DNA chromatograms*
- *Remove overlaps & ambiguous calls*
- *Read in all sequence files (10-10,000)*
- *Reverse complement all sequences (doubles # of sequences to align)*
- *Remove vector sequences (vector trim)*
- *Remove regions of low complexity*
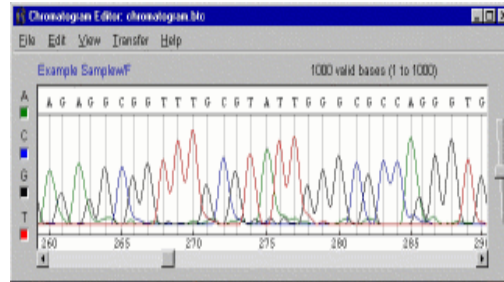- *Perform multiple sequence alignment*

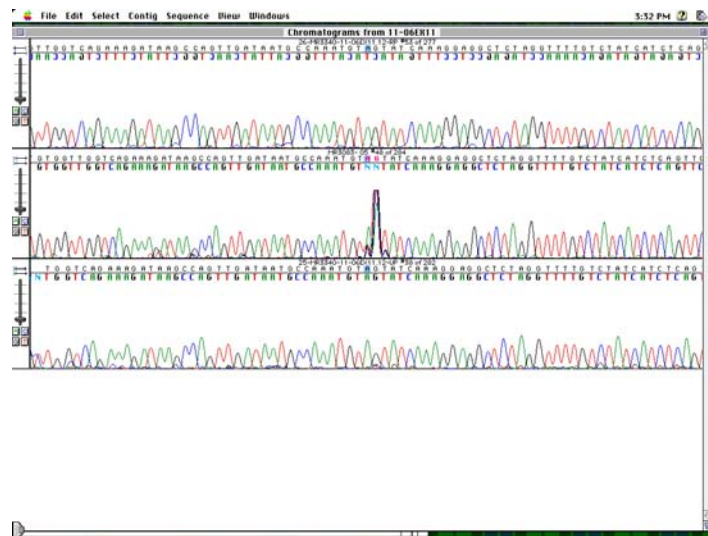# Contig Alignment - Process

# Reading DNA Chromatograms



**Gel**                    **ABI Chromatogram**
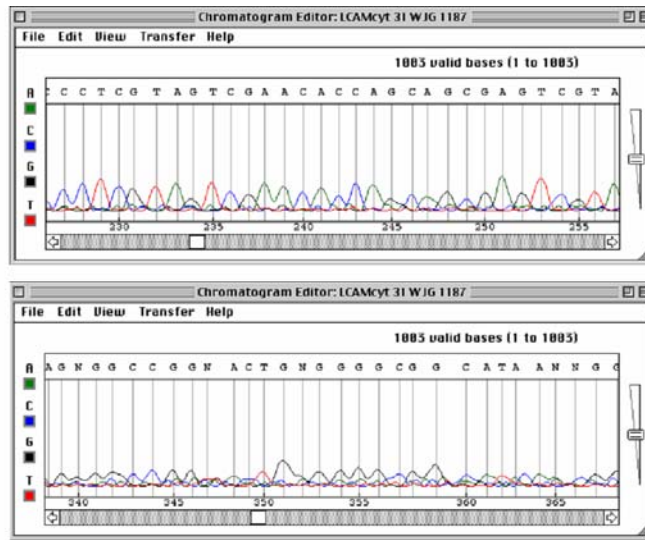
# Typical Raw Data

# Chromatograms (Problems)

- **Degradation of gel resolution (Pile-up or Band Broadening)**
- **Diminishment or excess of fluorescence intensity (too little or too much DNA tmplte)**
- **Differential overlap (large peak followed by a small one , ie. "G" dropouts (small G following a big A peak)**
- **Homopolymeric stretches of A's and T's**
- **Inappropriate spacing (contaminant  DNA or poor/noisy primers causing random priming)**
- **High GC content or GC rich regions**
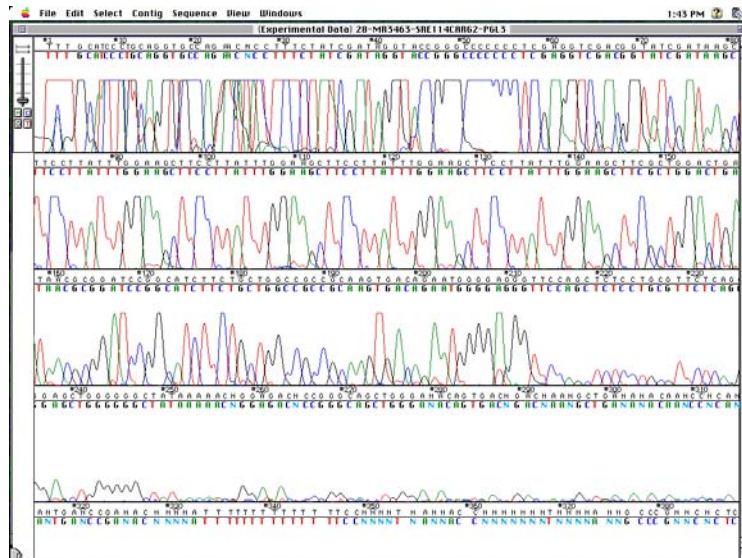- **Secondary structure or inverted repeats of the DNA**
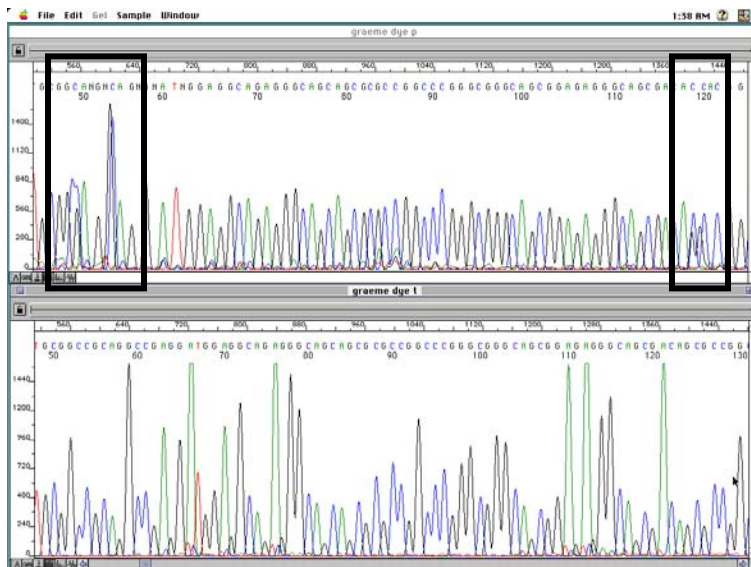
# Band Broadening

# Diminishing Intensity



# Too Much DNA Template

# High G-C Content

- **>60% GC content may be difficult to sequence (leads to pile-up)**
- **Dye terminator performs better than dye primer**
- **Easiest modification is to add 5% DMSO final concentration to the reaction mix**
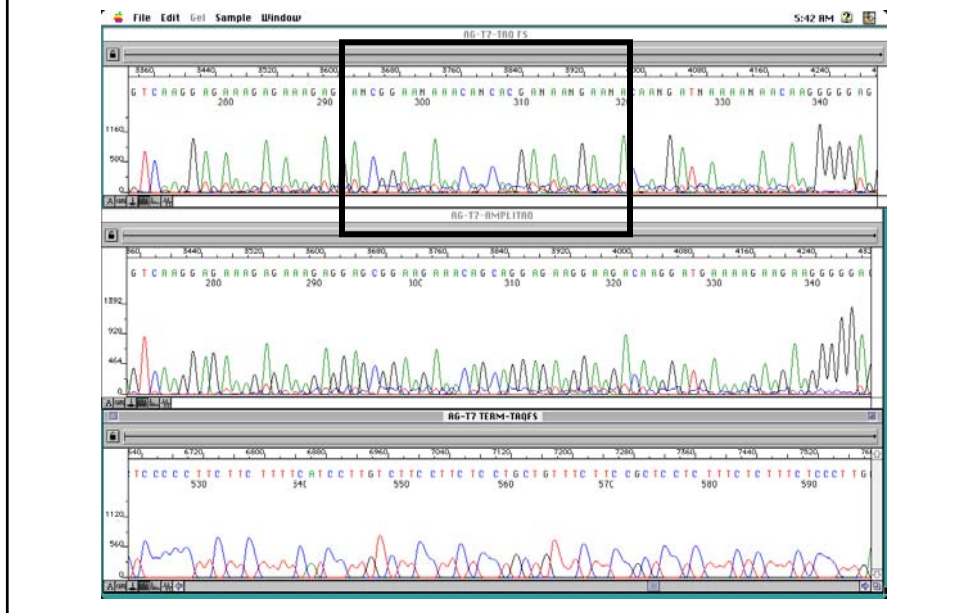- **Sequence the opposite strand to help resolve ambiguities**

# GC Pile Up

# Inverted/Extended Repeats

- **An abrupt loss of signal usually signifies a DNA sequence structure problem, due to the inability of the enzyme to proceed through the problem area**
- **5% DMSO sometimes helps**
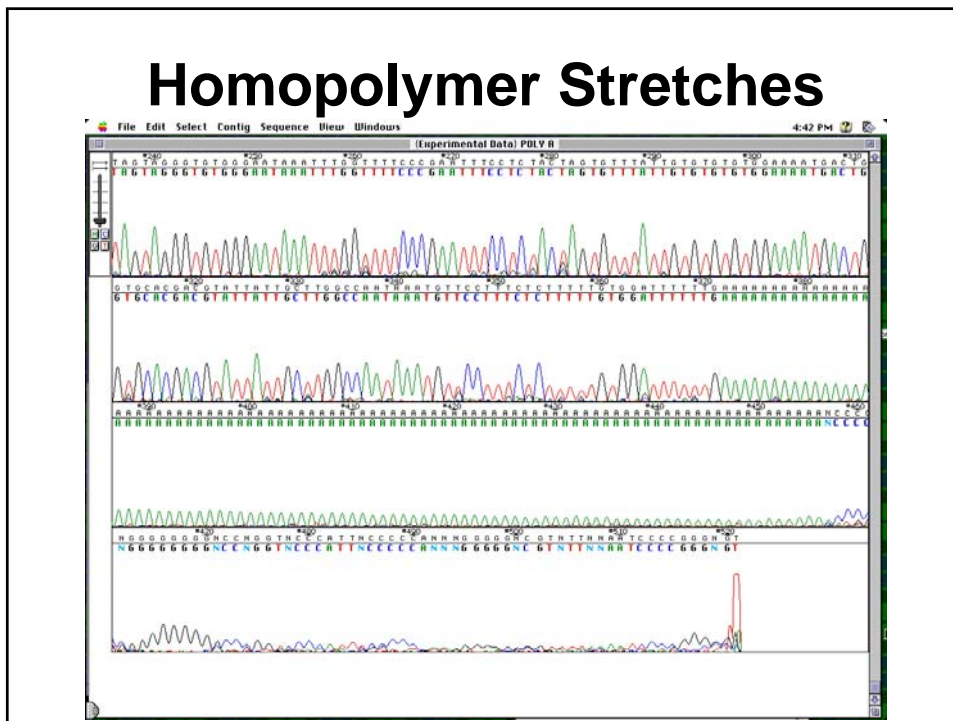- **Treat these the same way as high GC content regions**

# Repeats

- **Longer repeat sequences such as variable tandem repeats of 30 or more bases repeated many times are usually difficult to deal with**
- **AG repeat sequences can be problematic because Taq FS produces a weak G signal after A in terminator data**
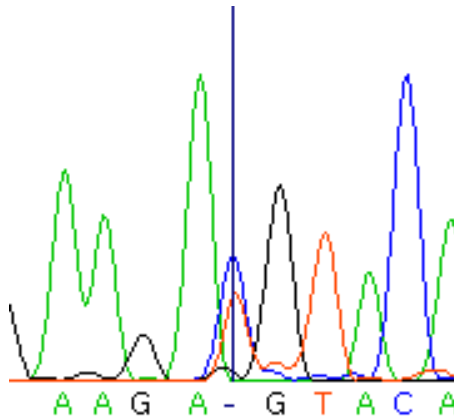- *More examples at*

  *http://www.abrf.org/Other/ABRFmeetings/ABRF96/tutorial4/*

# Weak G after A



# Homopolymer Stretches
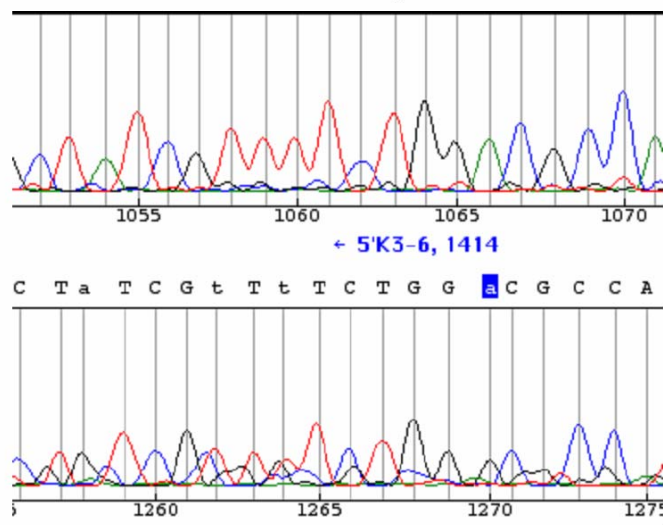
# Base Calling



A A G A - G T A C A

# Imperfect Raw Data

- **The data from sequencers varies in quality along the length of a single scan**
- **The base calls can be ambiguous, but there is still some information**
- **Need a quantitative analysis, not qualitative, to maximize information**

# Quality Factors

- **Simplest approach is human inspection, but not automatable**
- **Although computationally more difficult, quantitative factors provide a significant improvement in the assembly process**
- **Particularly important in high-throughput sequencing projects**

# Human Inspection

# Automated Base Calling with Phred

- **The Phred software reads DNA sequencing trace files, calls bases, and assigns a quality value to each called base**
- **The quality value is a log-transformed error probability, specifically**

$$Q = -10 \log_{10}( P_e )$$

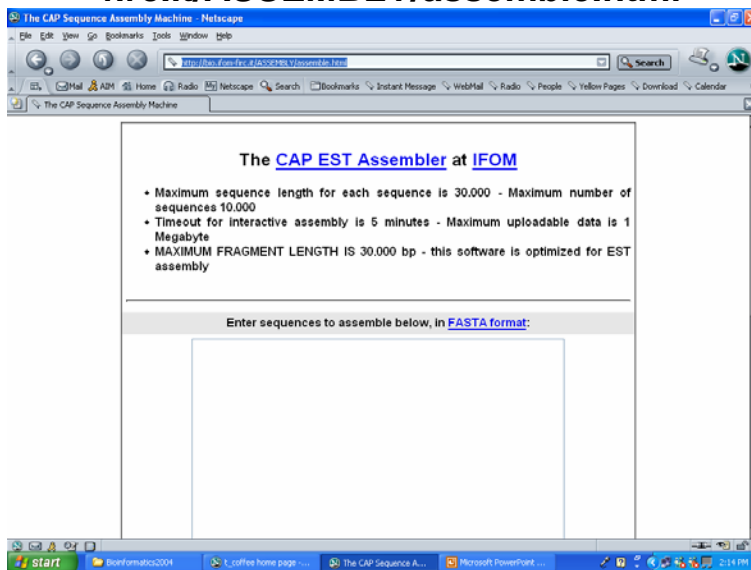- **where Q and Pe are respectively the quality value and error probability of a particular base call**

# Phred

- **The Phred quality values have been thoroughly tested for both accuracy and power to discriminate between correct and incorrect base-calls**
- **Phred can use the quality values to perform sequence trimming**

Ewing B, Green P: Basecalling of automated sequencer traces using phred. II. Error probabilities. Genome Research 8:186-194 (1998)

# Sequence Assembly Programs

- **Phred - base calling program that does detailed statistical analysis (UNIX)**
  http://www.phrap.org/
- **Phrap - sequence assembly program (UNIX)**
  http://www.phrap.org/
- **TIGR Assembler - microbial genomes (UNIX)**
  http://www.tigr.org/softlab/assembler/
- **The Staden Package (UNIX)**
  http://www.mrc-lmb.cam.ac.uk/pubseq/
- **GeneTool/ChromaTool/Sequencer (PC/Mac)**

---

# http://bio.ifom-firc.it/ASSEMBLY/assemble.html

# Contig Assembly Algorithm

- *Read, edit & trim DNA chromatograms*
- *Remove overlaps & ambiguous calls*
- *Read in all sequence files (10-10,000)*
- *Reverse complement all sequences (doubles # of sequences to align)*
- *Remove vector sequences (vector trim)*
- *Remove regions of low complexity*
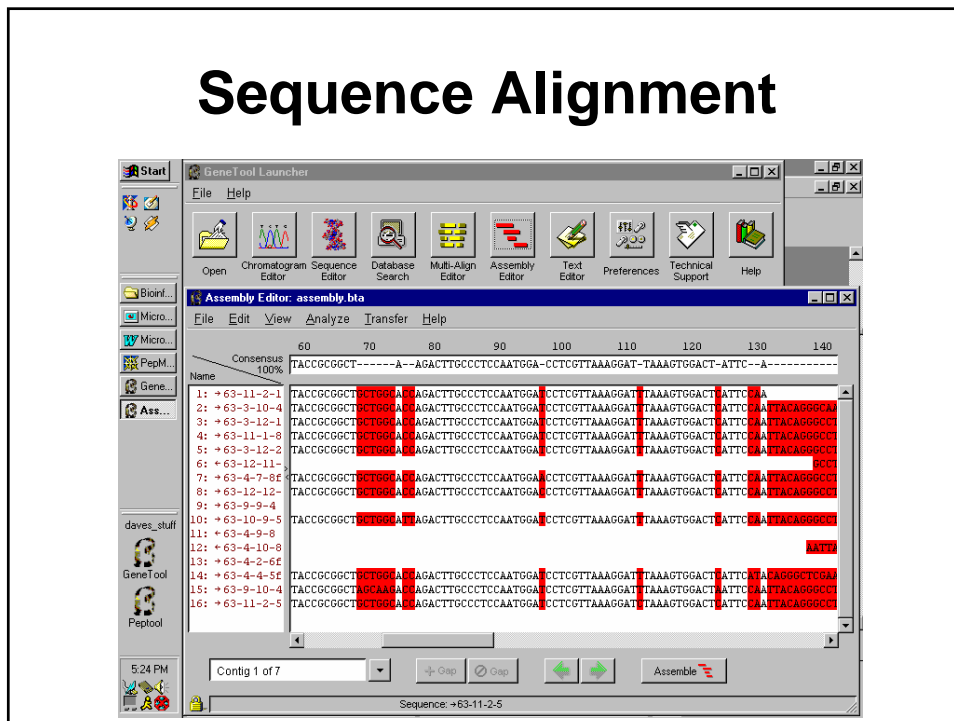- *Perform multiple sequence alignment*

# Chromatogram Editing

# Sequence Loading



# Sequence Alignment

# Assembly Parameters

- **User-selected parameters**
  1. **minimum length of overlap**
  2. **percent identity within overlap**
- **Non-adjustable parameters**
  1. **sequence "quality" factors**

# Phrap

- **Phrap is a program for assembling shotgun DNA sequence data**
- **Uses a combination of user-supplied and internally computed data quality information to improve assembly accuracy in the presence of repeats**
- **Constructs the contig sequence as a mosaic of the highest quality read segments rather than a consensus**
- **Handles large datasets**

# Problems for Assembly

- **Repeat regions**
  - Capture sequences from non-contiguous regions
- **Polymorphisms**
  - Cause failure to join correct regions
- **Large data volume**
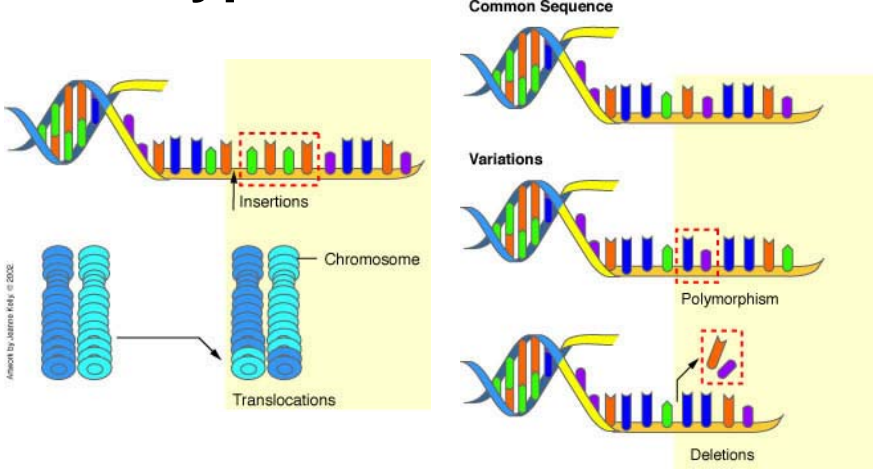  - Requires large numbers of pair-wise comparisons
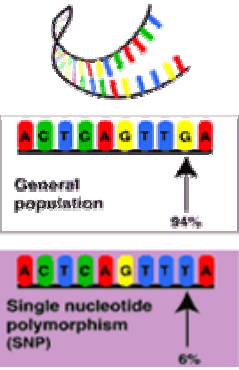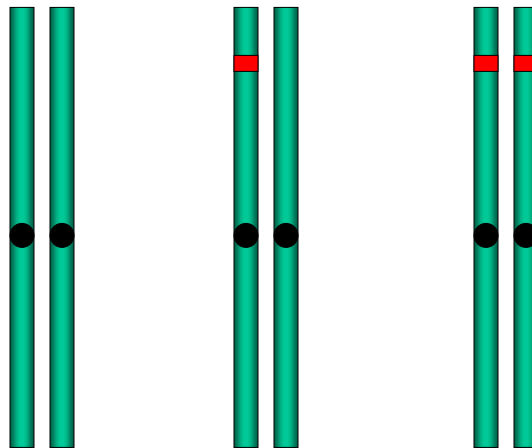
# Mutation Detection

# Types of Mutations



# SNPs & Polymorphisms

# SNPs (Single Nucleotide Polymorphisms)

- **Single nucleotide polymorphisms or SNPs are DNA sequence variations that occur when a single nucleotide (A,T,C or G) in the genome sequence is altered**

- **For a variation to be considered a SNP, it must occur in at least 1% of the population**

- **If the frequency is less than 1% (although this is somewhat arbitrary) then this variation is called a mutation**

- **SNPs are classified in three different ways…**

# Zygosity and SNPs

**Homozygous WT**      **Heterozygous**      **Homozygous Var.**

# SNPs

- **SNPs account for about 90% of all human genetic variation and are believed to occur every 100 to 300 bases along the 3-billion-base human genome**
- **Approximately 5 million of the ~10 million human SNPs have been catalogued**
- **SNPs may occur in exons, introns (non coding regions between exons) and intergenic regions (regions between genes)**
- **SNPs may lead to coding or amino acid sequence changes (non-synonymous) or they may leave the sequence unchanged (synonymous)**

# Synonymous vs. Non-Synonymous SNPs



**Hardy Weinberg Equilibrium**

# Hardy Weinberg Equilibrium

- **True SNPs should follow Hardy Weinberg Equilibrium in that**
- **The choice of a mate is not influenced by his/her genotype at the locus/gene (random mating or panmixia)**
- **The locus/gene/SNP does not affect the chance of mating at all, either by altering fertility or decreasing survival to reproductive age**
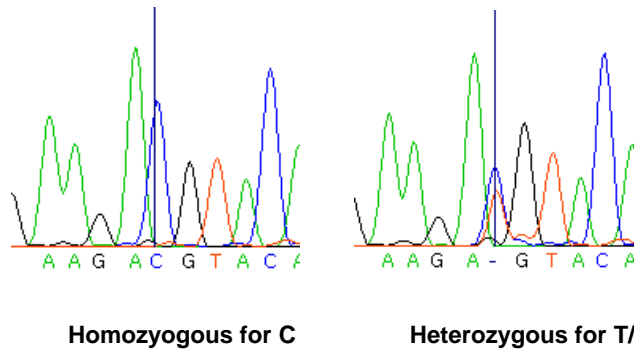
# Deviations from HWE

- **Marital assortment: "like marrying like"**
- **Inbreeding**
- **Population stratification: multiple subgroups are present within the population, each of which mates only within its own group (homogamy)**
- **Decreased viability of a particular genotype (hemophilia)**

# Measuring SNPs

- **Classical sequencing (homozygotes)**
- **Chromatogram analysis (heterozygotes)**
- **Denaturing HPLC**
- **Rolling Circle Amplification**
- **Antibody-based detection**
- **Enzyme- or cleavage-based detection**
- **Mass spectrometry**
- **SNP chips or microarrays**

# Polymorphism in Connexin26 (CX26) – Common Cause of Deafness -- ID by Sequencing



**Homozyogous for C**          **Heterozygous for T/C**

# The Finished Product

```
GATTACAGATTACAGATTACAGATTACAGATTACAG
ATTACAGATTACAGATTACAGATTACAGATTACAGA
TTACAGATTACAGATTACAGATTACAGATTACAGAT
TACAGATTAGAGATTACAGATTACAGATTACAGATT
ACAGATTACAGATTACAGATTACAGATTACAGATTA
CAGATTACAGATTACAGATTACAGATTACAGATTAC
AGATTACAGATTACAGATTACAGATTACAGATTACA
GATTACAGATTACAGATTACAGATTACAGATTACAG
ATTACAGATTACAGATTACAGATTACAGATTACAGA
TTACAGATTACAGATTACAGATTACAGATTACAGAT
```

# Shotgun Sequencing Summary

- **Very efficient process for small-scale (~10 kb) sequencing (preferred method)**
- **First applied to whole genome sequencing in 1995 (*H. influenzae*)**
- **Now standard for all prokaryotic genome sequencing projects**
- **Successfully applied to *D. melanogaster***
- **Moderately successful for *H. sapiens***

# NCBI Mapping & Assembly

- **Shotgun assembly doesn't always work (as was the case for the human genome)**
- **http://www.ncbi.nlm.nih.gov/genome/guide/build.html**
- **Describes the process used in the NCBI genome assembly and annotation process**

# Sequencing Successes



**T7 bacteriophage
completed in 1983
39,937 bp, 59 coded proteins**

**Escherichia coli
completed in 1998
4,639,221 bp, 4293 ORFs**

**Sacchoromyces cerevisae
completed in 1996
12,069,252 bp, 5800 genes**

# Sequencing Successses



**Caenorhabditis elegans
completed in 1998
95,078,296 bp, 19,099 genes**

**Drosophila melanogaster
completed in 2000
116,117,226 bp, 13,601 genes**

**Homo sapiens
Final draft completed in 2003
3,201,762,515 bp, 31,780 genes**

# Genomes to Date

- **8 vertebrates (**human, mouse, rat, fugu, zebrafish**)**
- **2 plants (arabadopsis, rice)**
- **2 insects (fruit fly, mosquito)**
- **2 nematodes (C. elegans, C. briggsae)**
- **1 sea squirt**
- **4 parasites (plasmodium, guillardia)**
- **4 fungi (S. cerevisae, S. pombe)**
- **200 bacteria and archaebacteria**
- **1900+ viruses**

# Sequenced Genomes



**http://www.genomenewsnetwork.org/**