

CAP 5510: Introduction to Bioinformatics
CGS 5166: Bioinformatics Tools

Giri Narasimhan

ECS 254; Phone: x3748

giri@cis.fiu.edu

www.cis.fiu.edu/~giri/teach/BioinfS13.html

Syllabus

- Fundamentals of Biology, Statistics, the Internet, and Bioinformatics
- Databases and Data Integration; Programming with BioPerl and BioPython;
- Sequence Alignment, Multiple Sequence Alignment
- Sequencing; Next Generation Sequencing & Applications
- Discovery, Learning, Prediction & Inference: Nucleotide and Protein Sequences
- Machine Learning: NN, HMM, SOM, SVM, etc.
- Gene Regulation; Regulatory Elements; microRNA; Regulatory networks
- Transcriptomics: Analysis of Gene Expression Data
- Gene Ontology and Pathways; Protein-protein interactions
- Genomics, Proteomics, Comparative Genomics
- Phylogenetic Analysis
- Molecular Structural Analysis: RNA and Proteins
- Genetics and Genome-Wide Association Schemes
- Single Nucleotide Polymorphisms
- Miscellanea: Omics; Alternative Splicing; Epigenetics; Translational bioinformatics; visualization;

Evaluation

<input type="checkbox"/> Semester Project	(45 %)
<input type="checkbox"/> Homework Assignments	(20 %)
<input type="checkbox"/> Exam	(15 %)
<input type="checkbox"/> Quizzes	(10 %)
<input type="checkbox"/> Summary Reports of Interest	(5 %)
<input type="checkbox"/> Class Participation	(5 %)

Course Homepage

<http://www.cis.fiu.edu/~giri/teach/BioinfS13.html>

- Lecture notes, required reading material, homework, announcements, etc.*

History

- ❑ What major world event took place on **26 June, 2000**?
- ❑ What major discovery was made in **1953**?
- ❑ 1975: Sanger Sequencing
- ❑ 1977: first bacteriophage sequenced
- ❑ 1990: HGP initiated

Introduction

1. What is Bioinformatics?

- Analysis of biological data with computing & statistical tools.

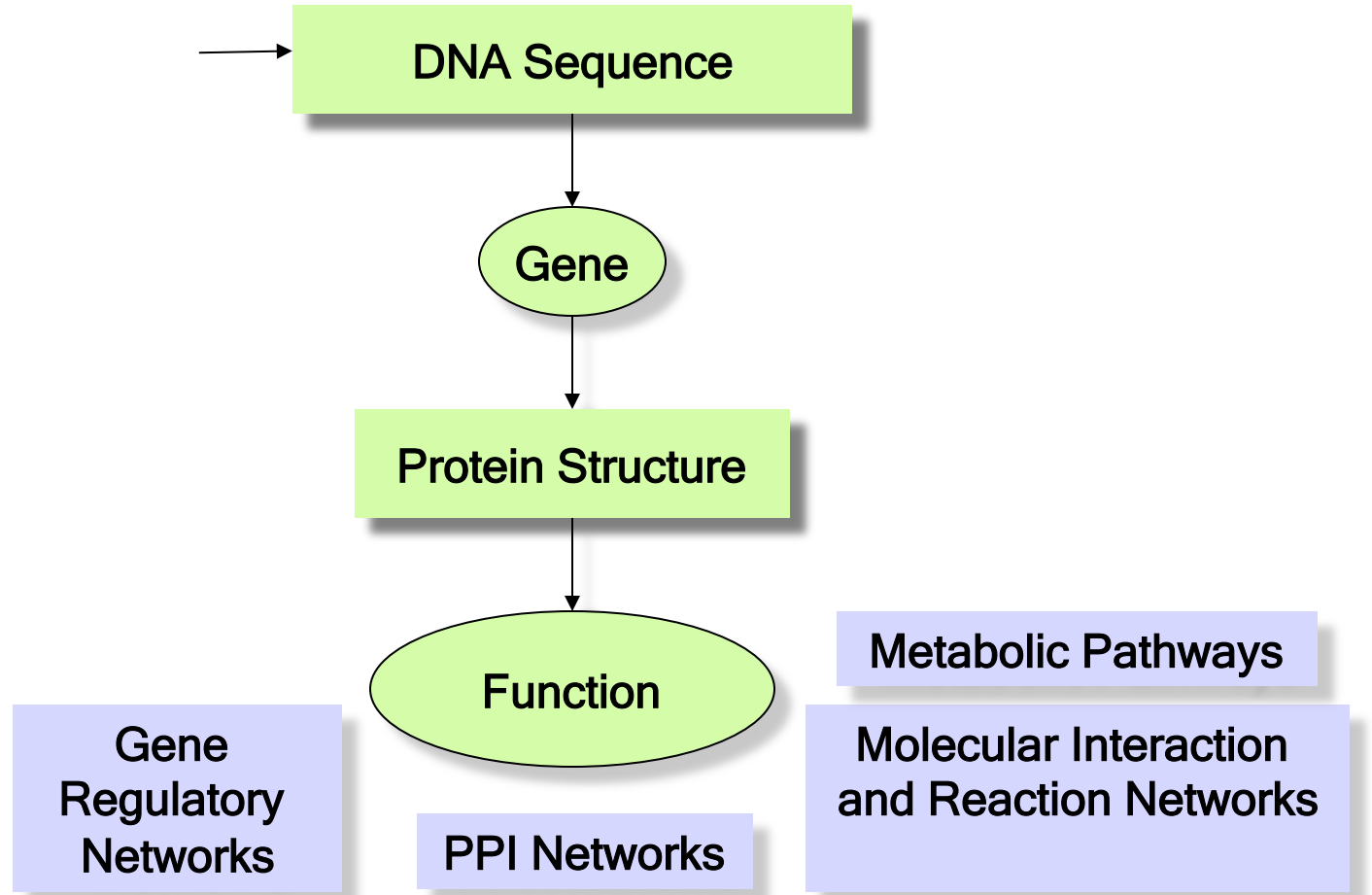
2. The different aspects of Informatics?

- Data Management (Database Technology, Internet Programming)
- Data Analysis (Data Mining, Modeling, Statistics)
- Development of Efficient Algorithms
- Visualization and Interface Design (HCI, Graphics)

3. How to assist biological research?

- Build databases for data
- Build efficient tools for search, retrieval, analysis, & visualization
- Propose models and efficient tools to verify the model using known data
- use predicted information to narrow down search
- propose new experiments based on model or analysis
- Build smart, hyperlinked, integrated mining environments

Overall Goals



Perspective of Bioinformatics

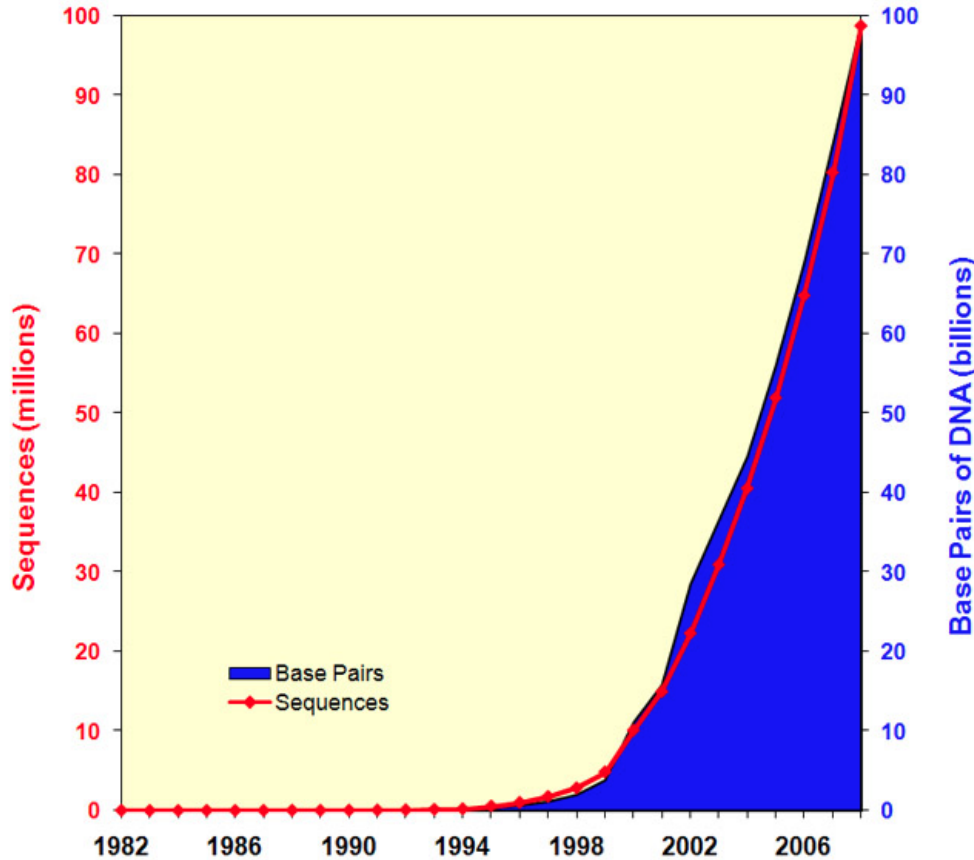
- ❑ Study of the cell: DNA, genes, proteins
- ❑ Study of the organism: genome, changes over time, over body regions, or over physiological or pathological states
- ❑ Study of all life: Tree of Life, Phylogeny, Variations, comparative genomics

General Information

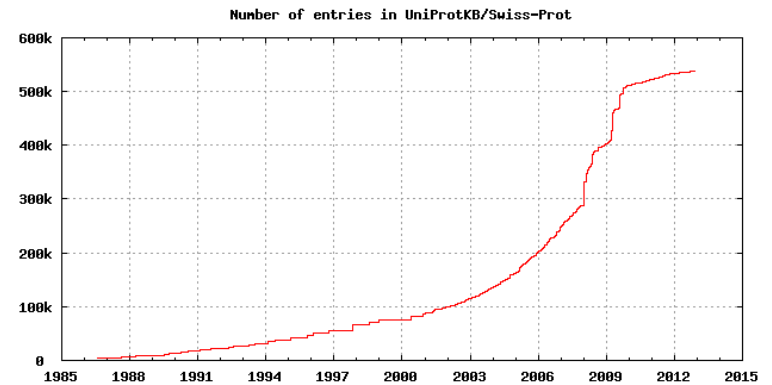
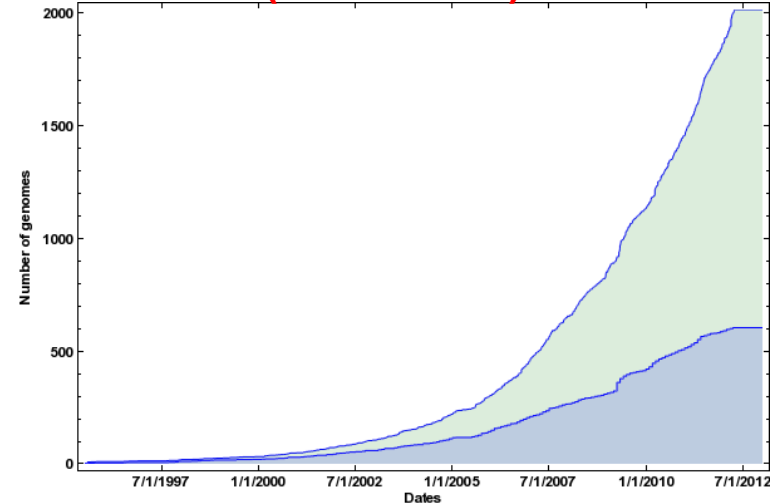
- ❑ **GenBank** Release 157/175/193 (Dec 2006/09/12) contains over 64/112/161 million sequence entries totaling over 69/110/126 Gb from over 2,500/?/9000 organisms [<http://www.ncbi.nlm.nih.gov>] (Storage: ~600 GBytes uncompressed)
- ❑ **Human Genome** has ~3 billion bp with 32,000+ genes.
- ❑ 435/624/3880 complete **microbial** genomes sequenced (684/914/12847 more in progress)
- ❑ 2540/3250 **Viral** genomes (300bp - 300Kb) (1st 1978: Simian virus; 5Kb).
- ❑ 22/180 complete **eukaryotic** genomes sequenced (175/613 more in progress):
 - Caenorhabditis elegans, Arabidopsis thaliana, Saccharomyces cerevisiae, Mus musculus, Homo sapiens, Oryza sativa, Plasmodium falciparum, Drosophila melanogaster, Anopheles gambiae, Macaca mulatta, Bos taurus, Felis catus, Gallus gallus*
- ❑ **Swiss-Prot** Release 51.3/54.7/2012_11 (Dec'06/Jan'08/Nov 2012): 250K/333K/550K entries; 91/120/191 million amino acids.

Growth of Sequence Databases

Growth of GenBank (1982 - 2008)



Microbial Genome Growth (1995-2012)



Genome Sizes

Organism	Size	Date	Est. # genes
<i>HIV type 1</i>	9.2 Kb	1997	9
<i>H. influenzae</i>	1.8 Mb	1995	1,740
<i>M. genitalium</i>	0.58 Mb	1998	525
<i>E. coli</i>	4.7 Mb	1997	4,000
<i>S. cerevisiae</i>	12.1 Mb	1996	6,034
<i>C. elegans</i>	97 Mb	1998	19,099
<i>A. thaliana</i>	100 Mb	2000	25,000
<i>D. melanogaster</i>	180 Mb	2000	13,061
<i>M. musculus</i>	3 Gb	2002	~30,000
<i>H. sapiens</i>	3 Gb	2001	32,000+

Short Homework: Curious Facts

□ What is interesting about the following organisms?

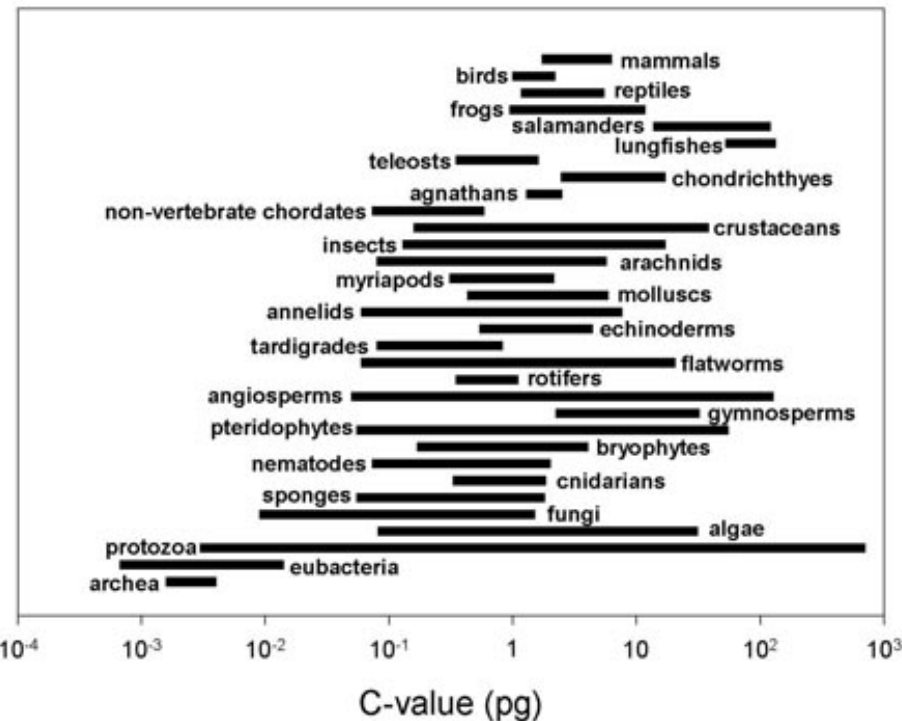
● *Paris japonica*, *Protopterus aethiopicus*, *Encephalozoon intestinalis*, *Ameoba dubia*, *Arabidopsis lyrata*

□ Ferrari of the virus world?

□ 1001 Genomes Project

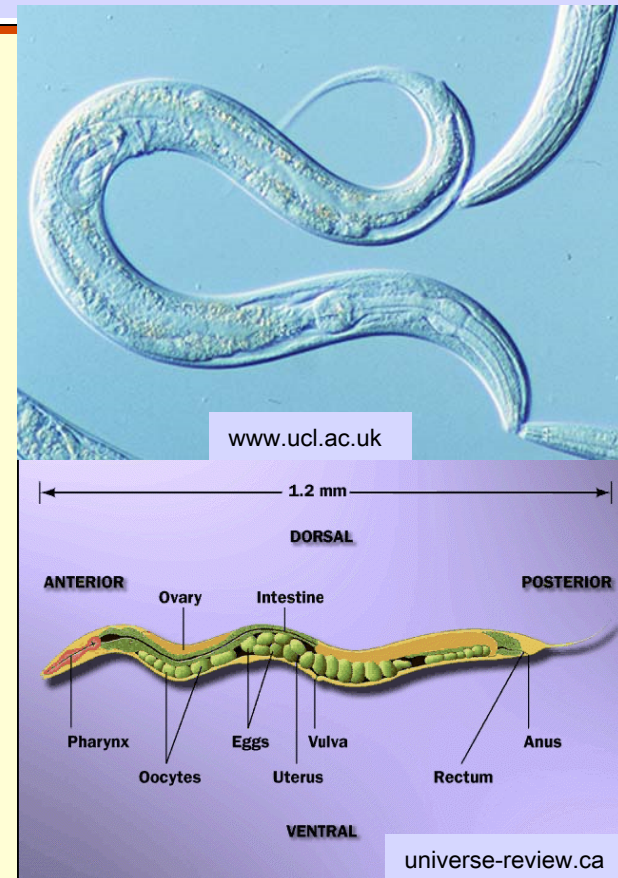
□ App: LeafSnap

Gregory, T.R. (2004a).
Macroevolution, hierarchy
theory, and the C-value enigma.
Paleobiology 30: 179-202.

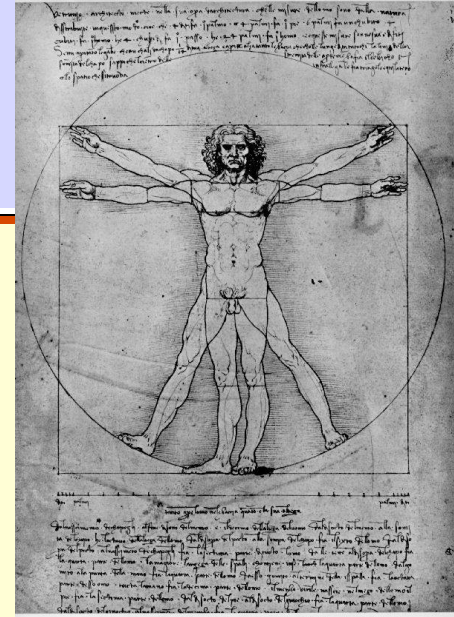


Caenorhabditis Elegans

- ❑ Entire genome - 1998; 8 year effort
- ❑ 1st animal; 2nd eukaryote (after yeast)
- ❑ Nematode (phylum)
- ❑ Easy to experiment with; Easily observable
- ❑ 97 million bases; 20,000 genes;
- ❑ 12,000 with known function; 6 Chromosomes;
- ❑ GC content 36%
- ❑ 959 cells; 302-cell nervous system
- ❑ 36% of proteins common with human
- ❑ 15 Kb mitochondrial genome
- ❑ Results in **ACeDB**
- ❑ 25% of genes in operons
- ❑ Important for HGP: technology, software, scale/efficiency
- ❑ 182 genes with alternative splice variants



Homo sapiens



- ❑ Sequenced - 2001; 15 year effort
- ❑ 3 billion bases, 500 gaps
- ❑ Variable density of **Genes, SNPs, CpG islands**
- ❑ ~ 1.1% of genome codes for proteins; **99%?**
- ❑ ~ 40-48% of the genome consists of repeat sequences
- ❑ ~ 10 % of the genome consists of repeats called ALUs
- ❑ ~ 5 % of the genome consists of long repeats (>1 Kb)
- ❑ 223 genes common with bacteria that are missing from worm, fly or yeast.

Sequence Alignment – Why?

```
>gi|12643549|sp|O18381|PAX6_DROME Paired box protein Pax-6 (Eyeless protein)
MRNLPCLGTAGGSGLGGIAGKPSPTMEAVEASTASHRHSTSSYFATTYYHLTDDECHSGVNQLGGVVFVGG
RPLPDSTRQKIVELAHSGARPCDISRILQVSNQCVSKILGRYYETGSIRPRAIGGSKPRVATAEVSISKIS
QYKRECPSIFAWAIRDRLLQENVCTNDNIPSVSSINRVLRLNLAQKEQQSTGSGSSSTSAGNSISAKVSV
SIGGNVSNVASGSRGTLSSSTDLMQTATPLNSSSESGGASNSGEGSEQEAIYEKLRLLNTQHAAGPGPLEP
ARAAPLVGQSPNHLGTRSSHPQLVHGNHQALQQHQQQSWPPRHYSWSWYPTSLSEIPISSAPNIASVTAY
ASGPSLAHSLSPNDIESLASIGHQRNCPVATEDIHLKKELDGHQSDETGSGEGENSNGGASNIGNTEDD
QARLILKRKLQRNRTSFTNDQIDSLEKEFERETHYPDVFARERLAGKIGLPEARIQVWFSNRRAKWRREEK
LRNQRRTPNSTGASATSSSTSATASLTDSPNSLSACSSLLSGSAGGPSVSTINGLSSPSTLSTNVNAPT
GAGIDSSSEPTPIPHIRPSCTSDNDNGRQSEDCRRVCSPCPLGVGGHQNTHHIQSNGHAQGHALVPAISP
RLNFNSGSGFGAMYSNMHHTALSMSDSYGAVTPIPSFNHSAVGPLAPPSPIPQQGDLTPSSLYPCHMTLRP
PPMAPAHHHIVPGDGGRPAGVGLGSGQSANLGASCSSGSGYEVLSAYALPPPPMASSAADSSFSAASSAS
ANVTPHHTIAQESCSPCSSASHFGVAHSSGFSSDPISPAVSSYAHMSYNYASSANTMTTPSSASGTSAHV
APGKQQFFASCFYSPWV
```

```
>gi|6174889|PAX6_HUMAN Paired box protein (Oculorhombin) (Aniridia, type II protein)
MQNSHSGVNQLGGVVFVNGRPLPDSTRQKIVELAHSGARPCDISRILQVSNQCVSKILGRYYETGSIRPRA
IGGSKPRVATPEVVSIAQYKRECPSIFAWAIRDRLLSEGVCTNDNIPSVSSINRVLRLNLASEKQQMGAD
GMYDKLRMLNGQTGSWGTRPGWYPGTSVPGQPTQDGCQQQEGGGENTNSISSNGEDSDEAQMRLQLKRKL
QRNRTSFTQEQIEALEKEFERETHYPDVFARERLAAKIDLPEARIQVWFSNRRAKWRREEKLRNQRRQASN
TPSHIPISSSFSTSVYQPIPQPTTPVSSFTSGSMLGRTDTALTNTYSALPPMPSFTMANNLPMQPPVPSQ
TSSYSCMLPTSPSVNGRSYDITYTPPHMQTHMNSQPMGTSGETTSTGLISPGVSVPVQVPGSEPDMSQYWPR
LQ
```

Drosophila Eyeless vs. Human Aniridia

Query: 57 HSGVNQLGGV FVGG RPLPDSTRQKIVELAHSGARPCDISRILQVSN GCVSKILGRYYETG 116

HSGVNQLGGV FV GRPLPDSTRQKIVELAHSGARPCDISRILQVSN GCVSKILGRYYETG

Sbjct: 5 HSGVNQLGGV FVNGRPLPDSTRQKIVELAHSGARPCDISRILQVSN GCVSKILGRYYETG 64

Query: 117 SIRPRAIGGSKPRVATAE VVSKISQYKRECPSIFAW EIRDRL LQENVCTNDNIPSVSSIN 176

SIRPRAIGGSKPRVAT EVVSKI+QYKRECPSIFAW EIRDRL L E VCTNDNIPSVSSIN

Sbjct: 65 SIRPRAIGGSKPRVATPE VVSKIAQYKRECPSIFAW EIRDRL LSEG VCTNDNIPSVSSIN 124

Query: 177 RVLRLNLA AQKEQ 188

RVLRLNLA++K+Q

Sbjct: 125 RVLRLNLA SEKQQ 136

Query: 417 TEDDQARLILKRKLQRNRTSFTNDQIDSLEKEFER THYPDVFARERLAGKIGLPEARIQV 476

+++ Q RL LKRKLQRNRTSFT +QI++LEKEFER THYPDVFARERLA KI LPEARIQV

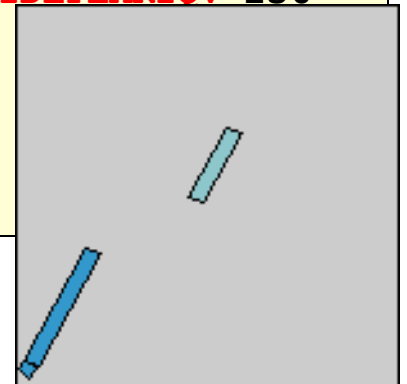
Sbjct: 197 SDEAQMRLQLKRKLQRNRTSFTQE QIEALEKEFER THYPDVFARERLAAKIDLPEARIQV 256

Query: 477 WFSNRRAKWRREEKLRNQR 496

WFSNRRAKWRREEKLRNQR

Sbjct: 257 WFSNRRAKWRREEKLRNQR 276

E-Value = $2e^{-31}$

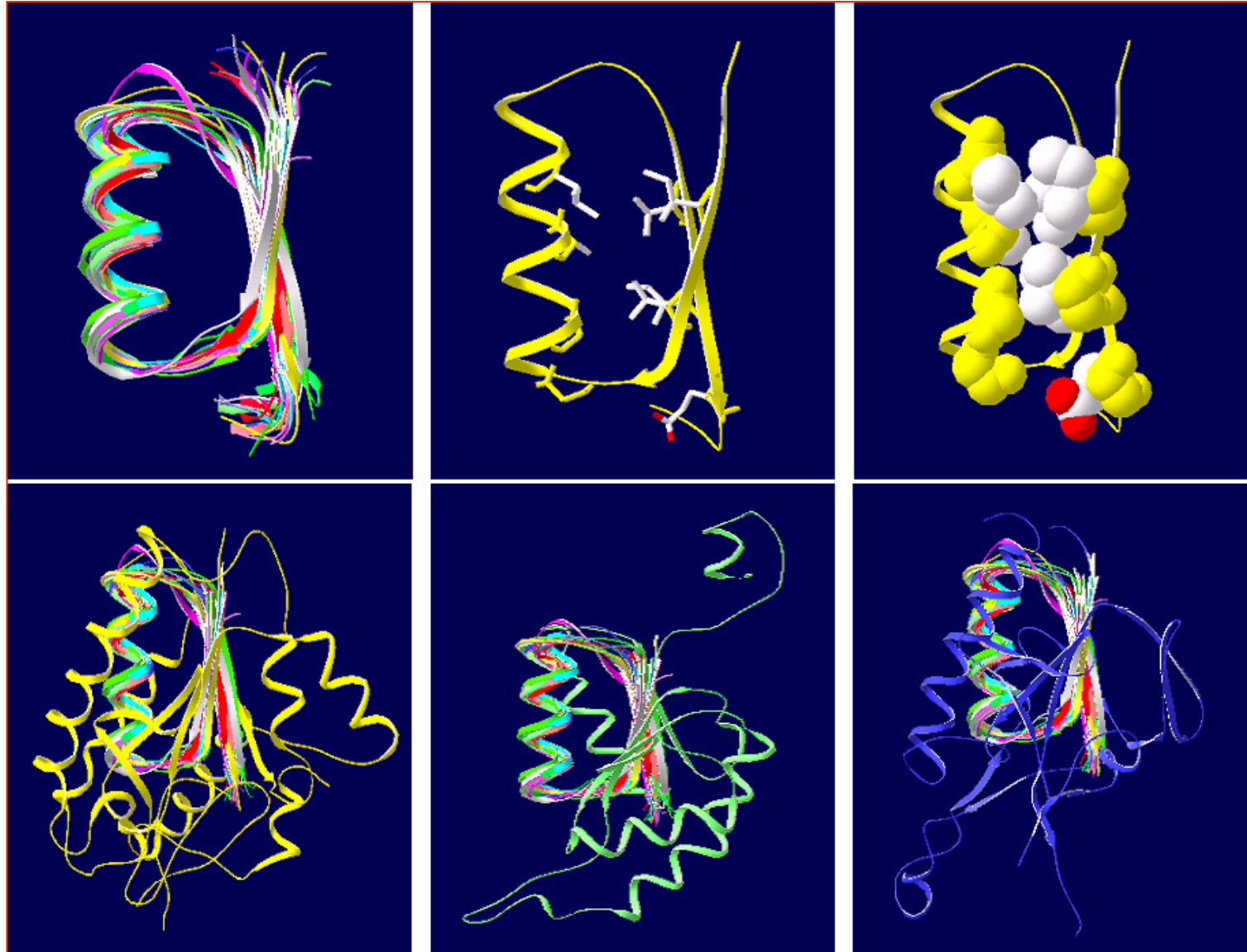


Motif Detection in Protein Sequences

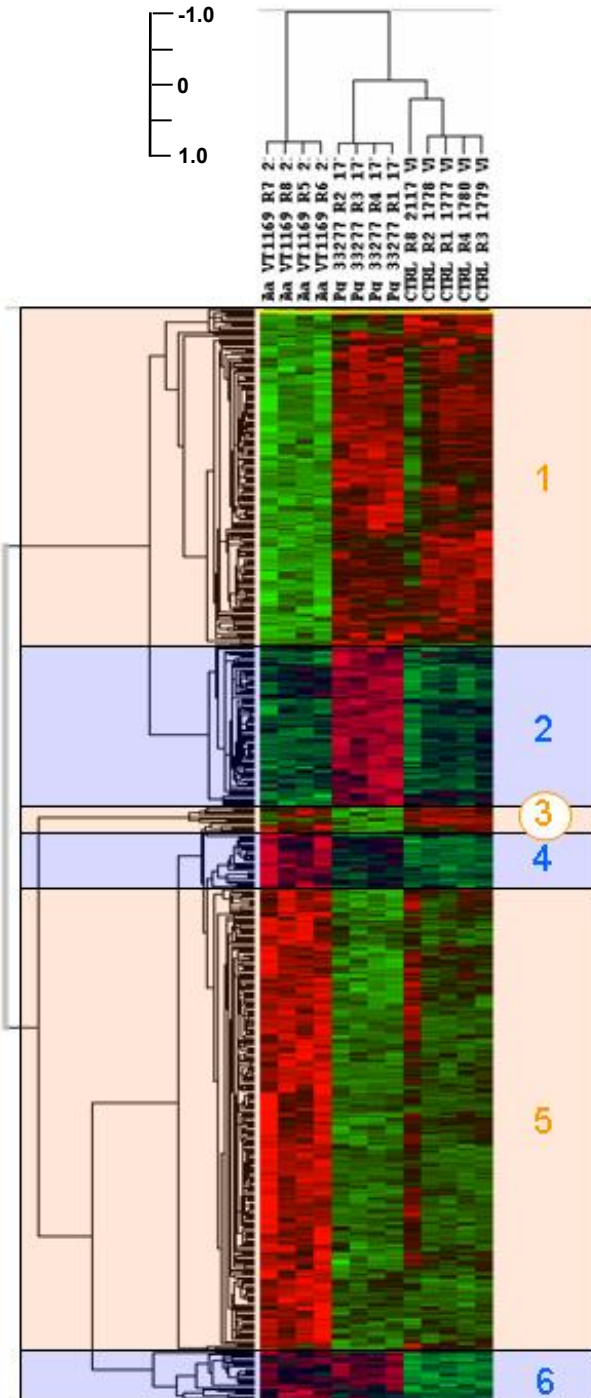
□ MTDKMQSLALAPVGNLDSYIRAANAWPMLSAD EERALAEK LHYHGDLEAA
KTLILSHLRFVVHIARNYAGYGLPQADLIQEGNIGLMKAVRRFNPEVGVR
LVSFVHWIKAEIHEYVLRNWRIVKVATTKAQRKLFNLRKTKQRLGWFN
QDEVEMVARELGVT SKDVREME SRMAAQDMTFDLS SDDSDS QPMAPVLY
LQDKSSNFADGIEDDNWEEQAANRLTDAMQGLDERSQDIIRARWLDEDNK
STLQELADRYGVSAERVRQLEKNAMKKLRAAIEA

□ MTDKMQSLALAPVGNLDSYIRAANAWPMLSAD EERALAEK LHYHGDLEAA
KTLILSHLRFVVHIARNYAGYGLPQADLIQEGNIGLMKAVRRFNPEVGVR
LVSFVHWIKAEIHEYVLRNWRIVKVATTKAQRKLFNLRKTKQRLGWFN
Q DEVEMVARELGVT SKDVREME SRMAAQDMTFDLS SDDSDS QPMAPVLY
LQDKSSNFADGIEDDNWEEQAANRLTDAMQGLDERSQDIIRARWLDEDNK
STLQELADRYGVSAERVRQLEK NAMKKLRAAIEA

Patterns in Protein Structures



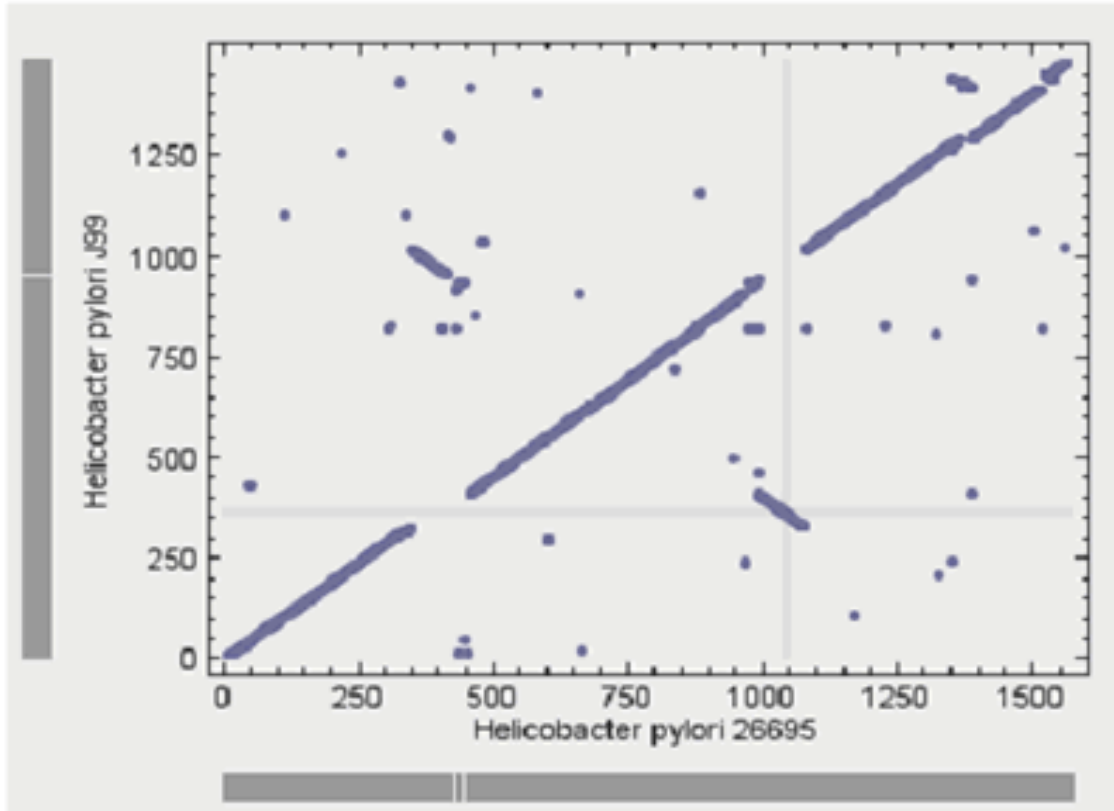
Microarray Analysis



Different patterns of gene expression of oral epithelial IHGK cells upon co-culture with *A. actinomycetemcomitans* or *P. gingivalis*.

Tools: GenePlot

1491 proteins total



Comparison of proteins from two strains of *Helicobacter Pylori*, 26695 and J99. Each point represents a pair of proteins from the two organisms showing a symmetrical best BLAST score; the coordinates of each point correspond to the position of the protein genes in the 2 genomes. Note the juxtaposition and inversion of two segments of the genome between the two strains.

SIDS



- ❑ 18000 Amish people in Pennsylvania
- ❑ Mostly intermarried due to religious doctrine
- ❑ rare recessive diseases occurred with high frequencies.
- ❑ SIDS: 3000 deaths/year (US); 21 deaths (Amish community)
- ❑ Many research centers failed to identify cause
- ❑ Collaboration between Affymetrix, TGEN & Clinic for special children solved the problem in 2 months
- ❑ Studied 10000 SNPs using microarray technology
- ❑ Their experiments showed that all the sick infants had two mutant copies of a specific gene, and their parents were carriers of the mutant gene.
- ❑ Conclusion: **Disease caused by 2 abnormal copies of TSPYL gene**
- ❑ Identified genes expressed in key organs (brainstem, testes)
- ❑ http://www.affymetrix.com/community/wayahead/modern_miracle.affx