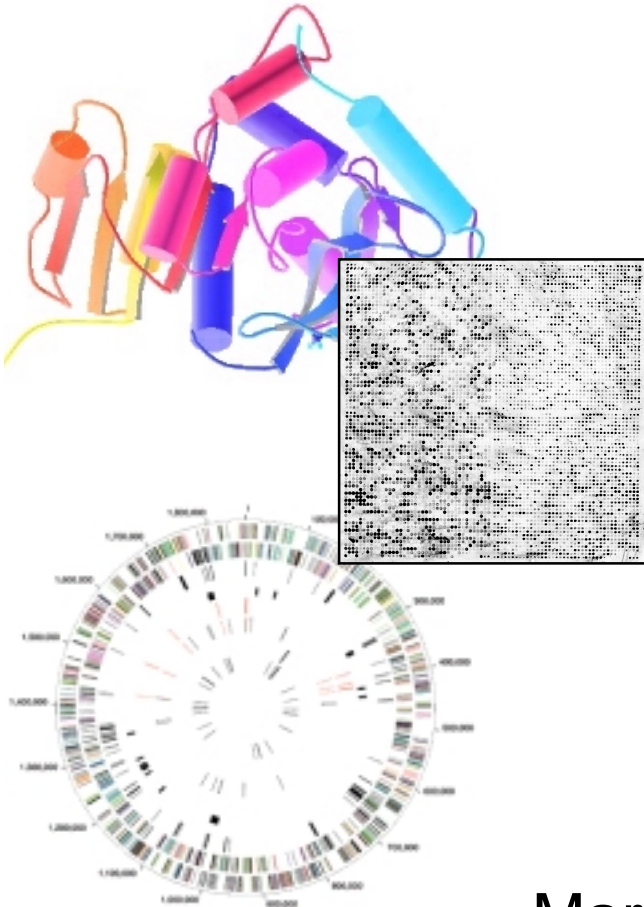


BIOINFORMATICS

Introduction



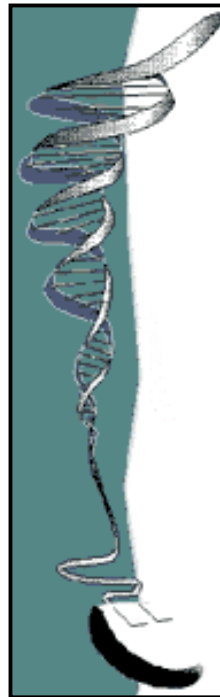
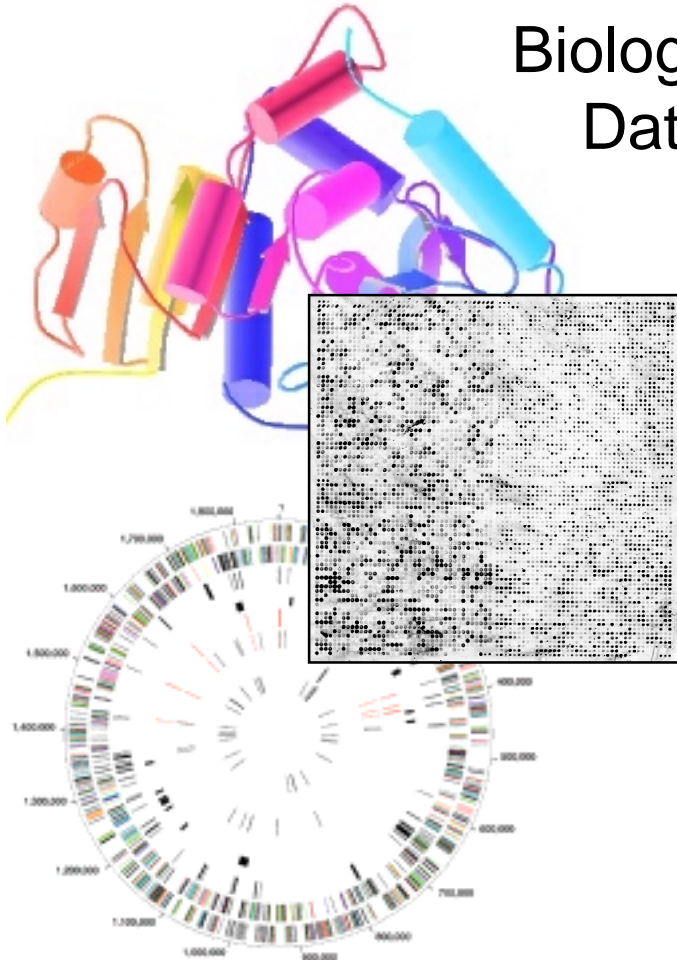
Mark Gerstein, Yale University
bioinfo.mbb.yale.edu/mbb452a

Bioinformatics

Biological
Data

+

Computer
Calculations



What is Bioinformatics?

- (*Molecular*) **Bio - informatics**
- One idea for a definition?
Bioinformatics is conceptualizing **biology in terms of molecules** (in the sense of physical-chemistry) and then applying **“informatics” techniques** (derived from disciplines such as applied math, CS, and statistics) to understand and **organize the information associated** with these molecules, **on a large-scale.**
- Bioinformatics is “MIS” for Molecular Biology Information

Molecular Biology: an Information Science

- Central Dogma of Molecular Biology

DNA

-> RNA

-> Protein

-> Phenotype

-> DNA

- Molecules

◇ Sequence, Structure, Function

- Processes

◇ Mechanism, Specificity, Regulation

- Central Paradigm for Bioinformatics

Genomic Sequence Information

-> mRNA (level)

-> Protein Sequence

-> Protein Structure

-> Protein Function

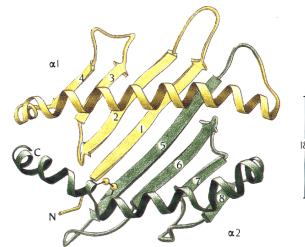
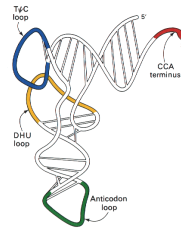
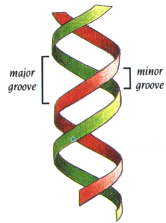
-> Phenotype

- Large Amounts of Information

◇ Standardized

◇ Statistical

(idea from D Brutlag, Stanford, graphics from S Strobel)



- Most cellular functions are performed or facilitated by proteins.

- Primary biocatalyst

- Cofactor transport/storage

- Mechanical motion/support

- Immune protection

- Control of growth/differentiation

- Genetic material

- Information transfer (mRNA)

- Protein synthesis (tRNA/mRNA)

- Some catalytic activity

Molecular Biology Information - DNA

- Raw DNA Sequence
 - ◇ Coding or Not?
 - ◇ Parse into genes?
 - ◇ 4 bases: AGCT
 - ◇ ~1 K in a gene,
~2 M in genome

```
atggcaattaaaattgggtatcaatgggttttggctcgtatcggccgcatcgtattccggtgca
gcacaacaccgctgatgacattgaagttgtaggtattaacgacttaatcgacggttgaatac
atggcttatatggtgaaatatgattcaactcacggctcgtttcgcacggcactggttgaagtg
aaagatggtaacttagtgggttaatggtaaaactatccggtgtaactgcagaacgtgatcca
gcaaacttaaactggggtgcaatcgggtgtgatatcgctggtgaagcactgggttattc
ttaactgatgaaactgctcgtaaacatatcactgcaggcgcaaaaaaagtgtattaact
ggcccatctaaagatgcaaccctatggttcggtcgtgggtgtaaaacttcaacgcatacgca
ggtcaagatatcgtttctaacgcactcttgtacaacaaactgttagctccttagcactg
gttggtcatgaaactttcgggtatcaaagatgggttaatgaccactggtcacgcaacgact
gcaactcaaaaaactgtggatggtccatcagctaaagactggcgcggcggcccggtgca
tcacaaaacatcattccatcttcaacaggtgcagcgaagcagtaggtaaagtattacct
gcattaaacggtaaatctggtatggctttccggtggtccaacgcgcaaacgtatctggt
gttgatttaacagttaatcttgaaaaaccagcttcttatgatgcaatcaacaagcaatc
aaagatgcagcgggaaggtaaaacgttcaatggcgaattaaaaggcgtattaggttacact
gaagatgctgtgtgttctactgacttcaacgggtgtgctttaaactctgtatttgatgca
gacgctgggtatcgcatctaactgattctttcgttaaattgggtac . . .
```

```
. . . caaaaatagggttaatatgaatctcgcattctccatcttctggttcattcgtattcaa
caacaagccaaaactcgtacaaatatgaccgcacttcgctataaagaacacggcttgtgg
cgagatatctcttggaaaaactttcaagagcaactcaatcaactttctcgagcattgctt
gctcacaatatgacgtacaagataaaaatcgccatttttgccataatatggaacggttgg
gttggtcatgaaactttcgggtatcaaagatgggttaatgaccactggtcacgcaacgact
acaatcgttgacattgacaccttacaattcgagcaatcacagtgccattttacgcaacc
aatcacgcccagcaagcagaatttatcctaaatcacgccgatgtaaaaattctcttcgctc
ggcgatcaagagcaatcagatcaaacattggaattgctcatattgtccaaaattacaa
aaaattgtagcaatgaaatccaccattcaattacaacaagatcctctttctgcacttgg
```

Molecular Biology Information: Protein Sequence

- 20 letter alphabet
 - ◇ ACDEFGHIKLMNPQRSTVWY but not BJOUXZ
- Strings of ~300 aa in an average protein (in bacteria),
~200 aa in a domain
- ~200 K known protein sequences

```
d1dhfa_ LNCIVAVSQNMGIGKNGDLPWPPLRNEFRYFQRM TTTSSVEGKQ-NLVIMGKKTWFSI
d8dfr__ LNSIVAVCQNMGIGKDG NLPWPPLRNEYKYFQRM TSTSHVEGKQ-NAVIMGKKTWFSI
d4dfra_ ISLIAALAVDRVIGMENAMPWN-LPADLAWFKRNTL-----NKPVIMGRHTWESI
d3dfr__ TAFLWAQDRDGLIGKDGHL PWH-LPDDLHYFRAQTV-----GKIMVVGRRTYESF
```

```
d1dhfa_ LNCIVAVSQNMGIGKNGDLPWPPLRNEFRYFQRM TTTSSVEGKQ-NLVIMGKKTWFSI
d8dfr__ LNSIVAVCQNMGIGKDG NLPWPPLRNEYKYFQRM TSTSHVEGKQ-NAVIMGKKTWFSI
d4dfra_ ISLIAALAVDRVIGMENAMPW-NLPADLAWFKRNTLD-----KPVIMGRHTWESI
d3dfr__ TAFLWAQDRNGLIGKDGHL PW-HLPDDLHYFRAQTVG-----KIMVVGRRTYESF
```

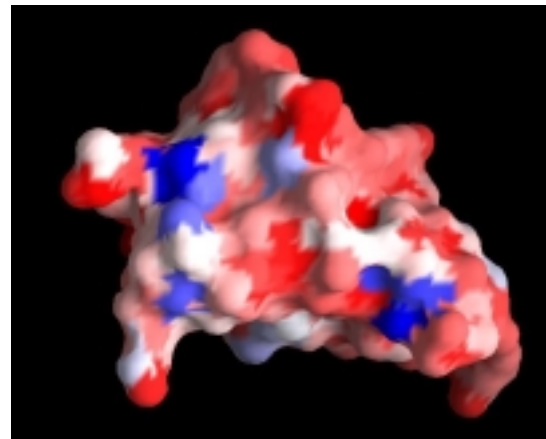
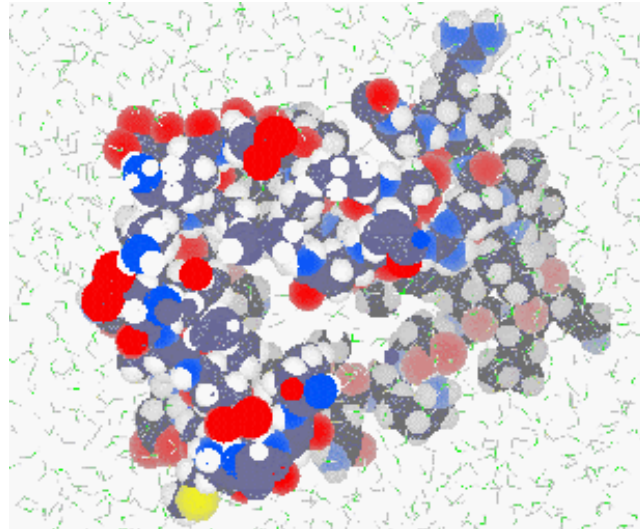
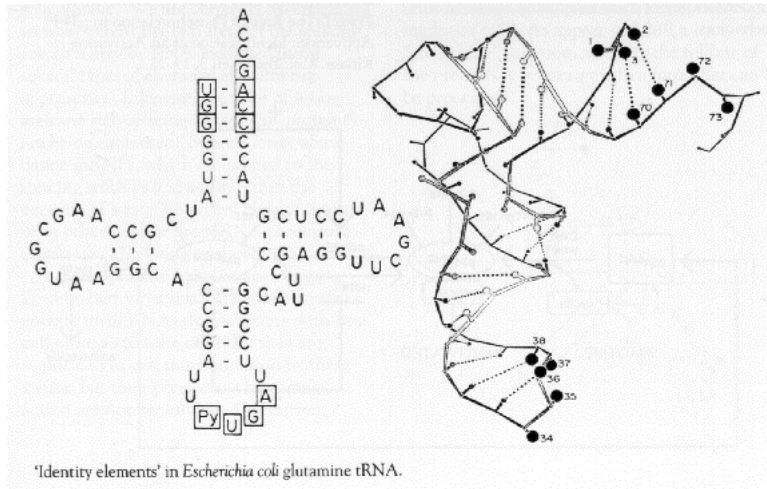
```
d1dhfa_ VPEKNRPLKGRINLVLSRELKEPPQGAHF LSRSLDDALKLTEQPELANKVDMVWIVGGSSVYKEAMNHP
d8dfr__ VPEKNRPLKDRINIVLSRELKEAPKGAHYLSKSLDDALALLDSPELKS KVD MVWIVGGTAVYKAAMEKP
d4dfra_ ---G-RPLPGRKNII LS-SQPGTDDR V-TWVKSVD EAIACGDV P-----EIMVIGGGRVYEQFLPKA
d3dfr__ ---PKRPLPERTNVVLTHQEDYQAQGA-VVVDVA AVFAYAKQHLDQ----ELVIAGGAQIFTAFKDDV
```

```
d1dhfa_ -PEKNRPLKGRINLVLSRELKEPPQGAHF LSRSLDDALKLTEQPELANKVDMVWIVGGSSVYKEAMNHP
d8dfr__ -PEKNRPLKDRINIVLSRELKEAPKGAHYLSKSLDDALALLDSPELKS KVD MVWIVGGTAVYKAAMEKP
d4dfra_ -G---RPLPGRKNII LSSSQPGTDDR V-TWVKSVD EAIACGDV PE-----IMVIGGGRVYEQFLPKA
d3dfr__ -P---KRPLPERTNVVLTHQEDYQAQGA-VVVDVA AVFAYAKQHLD----QELVIAGGAQIFTAFKDDV
```

Molecular Biology Information: Macromolecular Structure

- DNA/RNA/Protein
 - ◊ Almost all protein

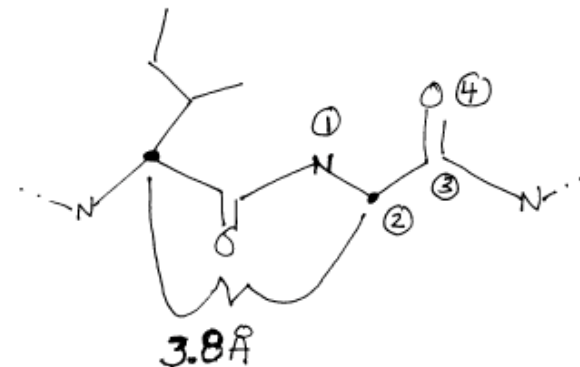
(RNA Adapted From D Soll Web Page,
Right Hand Top Protein from M Levitt web page)



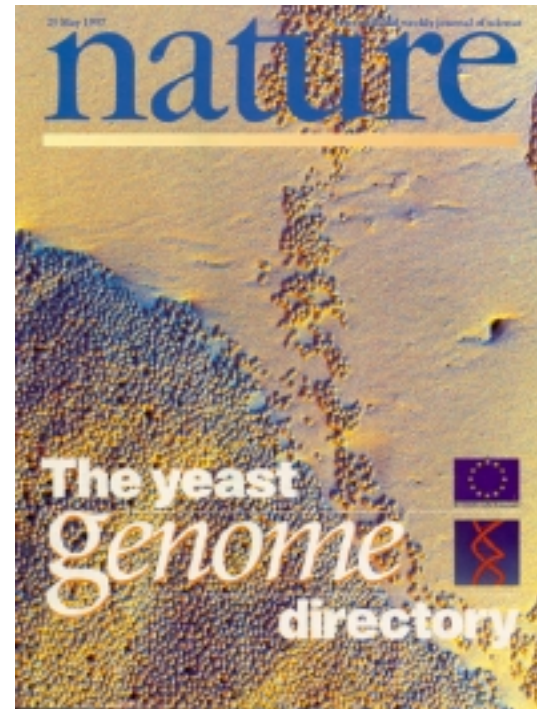
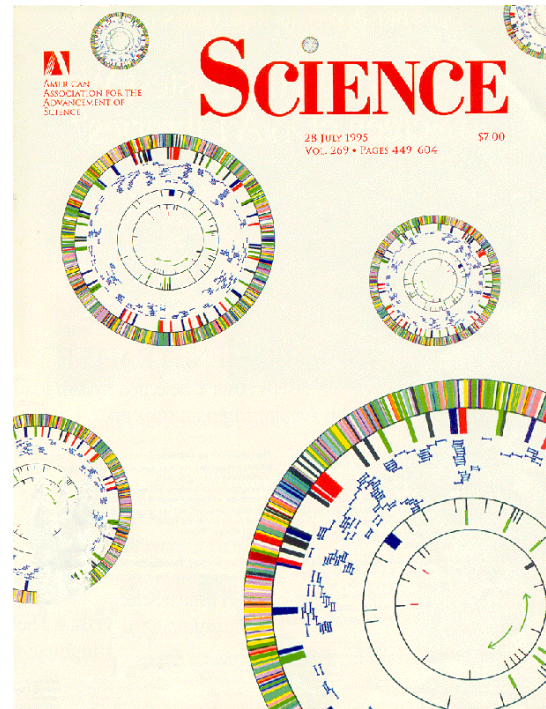
Molecular Biology Information: Protein Structure Details

- Statistics on Number of XYZ triplets
 - ◇ 200 residues/domain → 200 CA atoms, separated by 3.8 Å
 - ◇ Avg. Residue is Leu: 4 backbone atoms + 4 sidechain atoms, 150 cubic Å
 - => ~1500 xyz triplets (=8x200) per protein domain
 - ◇ 10 K known domain, ~300 folds

ATOM	1	C	ACE	0	9.401	30.166	60.595	1.00	49.88	1GKY	67
ATOM	2	O	ACE	0	10.432	30.832	60.722	1.00	50.35	1GKY	68
ATOM	3	CH3	ACE	0	8.876	29.767	59.226	1.00	50.04	1GKY	69
ATOM	4	N	SER	1	8.753	29.755	61.685	1.00	49.13	1GKY	70
ATOM	5	CA	SER	1	9.242	30.200	62.974	1.00	46.62	1GKY	71
ATOM	6	C	SER	1	10.453	29.500	63.579	1.00	41.99	1GKY	72
ATOM	7	O	SER	1	10.593	29.607	64.814	1.00	43.24	1GKY	73
ATOM	8	CB	SER	1	8.052	30.189	63.974	1.00	53.00	1GKY	74
ATOM	9	OG	SER	1	7.294	31.409	63.930	1.00	57.79	1GKY	75
ATOM	10	N	ARG	2	11.360	28.819	62.827	1.00	36.48	1GKY	76
ATOM	11	CA	ARG	2	12.548	28.316	63.532	1.00	30.20	1GKY	77
ATOM	12	C	ARG	2	13.502	29.501	63.500	1.00	25.54	1GKY	78
...											
ATOM	1444	CB	LYS	186	13.836	22.263	57.567	1.00	55.06	1GKY1510	
ATOM	1445	CG	LYS	186	12.422	22.452	58.180	1.00	53.45	1GKY1511	
ATOM	1446	CD	LYS	186	11.531	21.198	58.185	1.00	49.88	1GKY1512	
ATOM	1447	CE	LYS	186	11.452	20.402	56.860	1.00	48.15	1GKY1513	
ATOM	1448	NZ	LYS	186	10.735	21.104	55.811	1.00	48.41	1GKY1514	
ATOM	1449	OXT	LYS	186	16.887	23.841	56.647	1.00	62.94	1GKY1515	
TER	1450		LYS	186						1GKY1516	



Genomes
highlight
the
Finiteness
of the
World of
Sequences

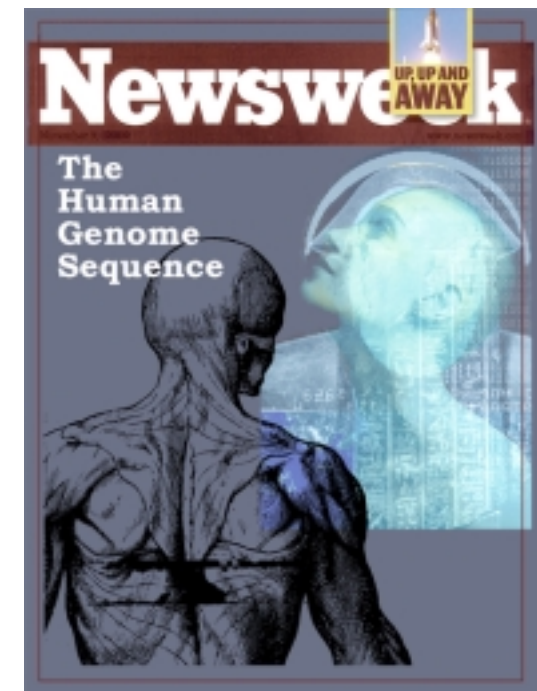


1995

Bacteria, 1.6 Mb, ~1600 genes [Science 269: 496]

1997

Eukaryote, 13 Mb, ~6K genes [Nature 387: 1]



1998

Animal, ~100 Mb, ~20K genes [Science 282: 1945]

2000?

Human, ~3 Gb, ~100K genes [???

Molecular Biology Information: Whole Genomes

- The Revolution Driving Everything

Fleischmann, R. D., Adams, M. D., White, O., Clayton, R. A., Kirkness, E. F., Kerlavage, A. R., Bult, C. J., Tomb, J. F., Dougherty, B. A., Merrick, J. M., McKenney, K., Sutton, G., Fitzhugh, W., Fields, C., Gocayne, J. D., Scott, J., Shirley, R., Liu, L. I., Glodek, A., Kelley, J. M., Weidman, J. F., Phillips, C. A., Spriggs, T., Hedblom, E., Cotton, M. D., Utterback, T. R., Hanna, M. C., Nguyen, D. T., Saudek, D. M., Brandon, R. C., Fine, L. D., Fritchman, J. L., Fuhrmann, J. L., Geoghagen, N. S. M., Gnehm, C. L., McDonald, L. A., Small,

K. V., Fraser, C. M., Smith, H. O. & Venter, J. C. (1995). "Whole-genome

random sequencing and assembly of *Haemophilus influenzae* rd."

Science 269: 496-512.

(Picture adapted from TIGR website,
<http://www.tigr.org>)

- Integrative Data

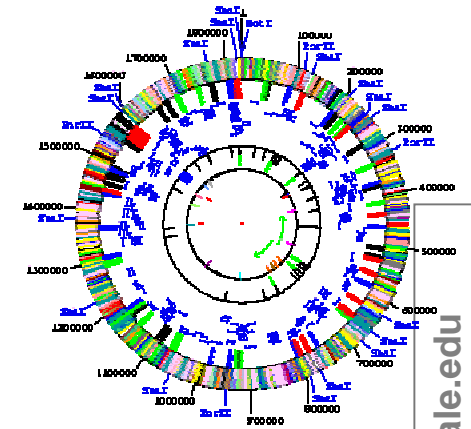
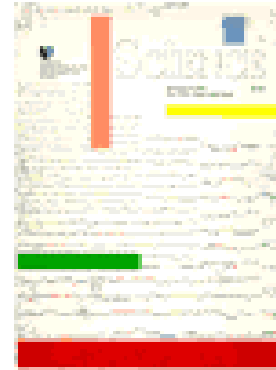
1995, HI (bacteria): 1.6 Mb & 1600 genes done

1997, yeast: 13 Mb & ~6000 genes for yeast

1998, worm: ~100Mb with 19 K genes

1999: >30 completed genomes!

2003, human: 3 Gb & 100 K genes...



Genome sequence now accumulate so quickly that, in less than a week, a single laboratory can produce more bits of data than Shakespeare managed in a lifetime, although the latter make better reading.

-- G A Pekso, *Nature* **401**: 115-116 (1999)

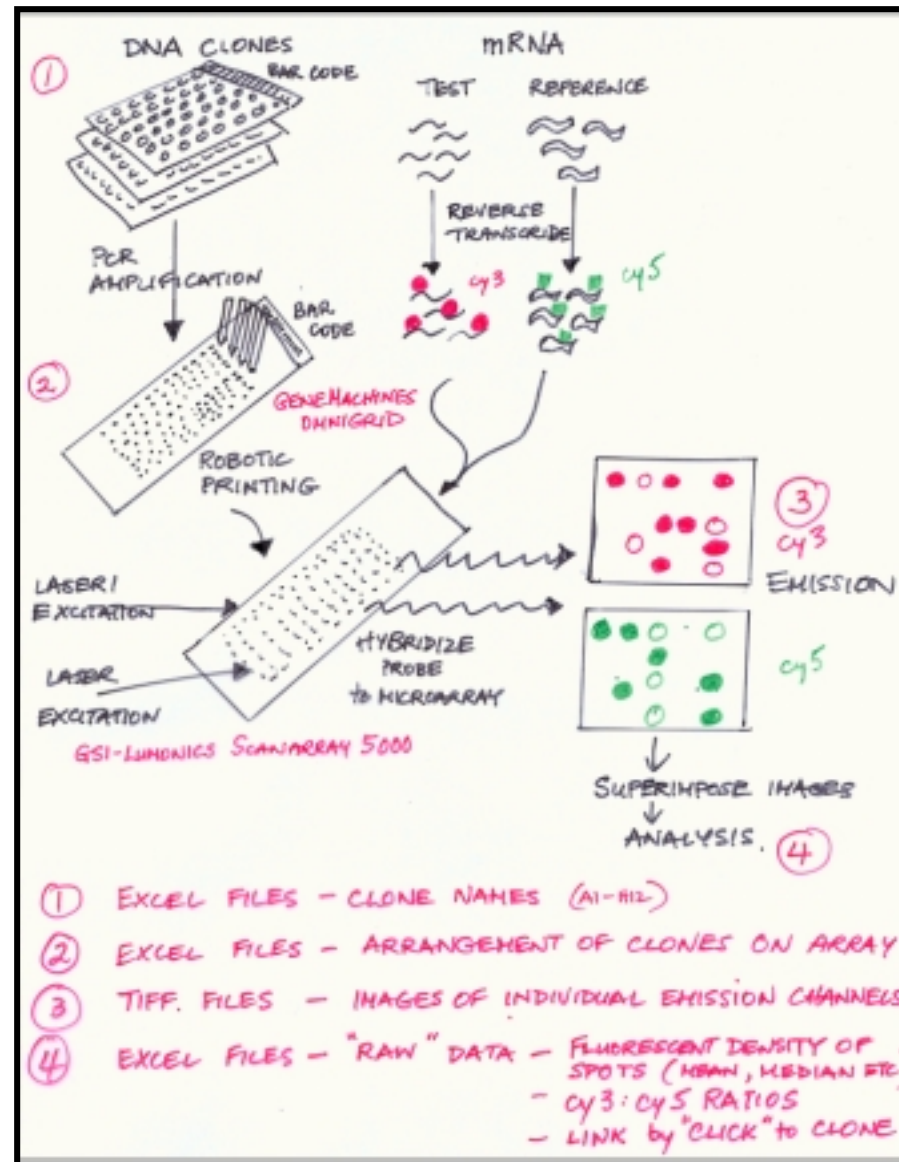
Array Data

Yeast Expression Data in Academia:
levels for all 6000 genes!

Can only sequence genome once but can do an infinite variety of these array experiments

at 10 time points,
 $6000 \times 10 = 60K$ floats

telling signal from background



(courtesy of J Hager)

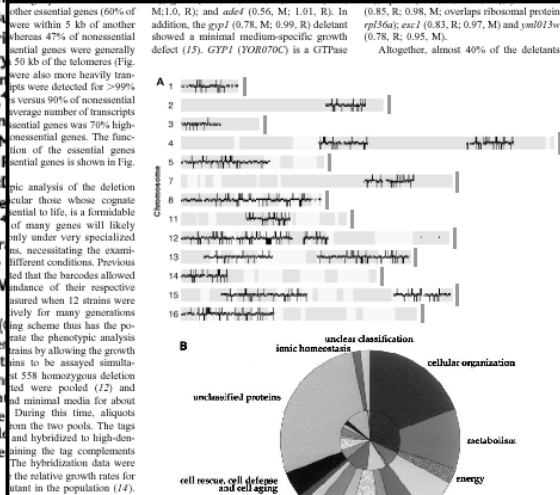
Functional Characterization of the *S. cerevisiae* Genome by Gene Deletion and Parallel Analysis

Elizabeth A. Winzeler,^{1*} Daniel D. Shoemaker,^{2*} Anna Astromoff,^{1*} Hong Liang,^{1*} Keith Anderson,¹ Bruno Andre,³ Rhonda Bangham,⁴ Rocio Benito,⁵ Jef D. Boeke,⁶ Howard B. Carla Connelly,⁶ Karen Davis,⁷ Fred Dietrich,⁸ Mohamed El Bakkoury,⁹ Françoise Foury,¹⁰ Erik Gentalen,¹¹ Guri Giaever,¹ Johan Ted Jones,¹ Michael Laub,¹ Hong Liao,¹ David J. Lockhart,¹¹ Anca Lucau-Danilasi,¹² Nasih M'Rabet,³ Patrice Menard,⁷ Chai Pai,¹ Corinne Rebischung,⁸ Jose L. Christopher J. Roberts,² Petra Ross-Macdonald,¹³ Michael Snyder,⁴ Sharon Sookhai-Mahadeo,¹⁴ Steve Véronneau,⁷ Marleen Voet,¹⁵ Teresa R. Ward,² Robert Wysocki,¹⁰ G. Katja Zimmermann,¹² Peter Mark Johnston,¹³ Ronald

The functions of many open reading frames in sequencing projects are unknown. New, whole-genome approaches to systematically determine their function. A set of 558 homozygous deletion strains were constructed, by a high-throughput parallel deletion of one of 2026 ORFs (more than 1% of the genome). Of the deleted ORFs, 17 percent were essential. The phenotypes of more than 500 deletion strains were analyzed in parallel. Of the deletion strains, 40 percent showed either rich or minimal medium growth.

that serve as strain identifiers (6, 7). We show that these barcodes allow large numbers of deletion strains to be pooled and analyzed in parallel in competitive growth assays. This direct, simultaneous, competitive assay of fitness increases the sensitivity, accuracy and speed with which growth defects can be detected relative to conventional methods.

To take full advantage of this approach and to accelerate the pace of progress, an international consortium was organized to



of 1620 nonessential (short black bars) and 356 essential genes (all genes are generated in consecutive groups on multiple chromosomes. A lighter grey bar indicates the location of chromosomal duplication blocks (23). For 15 of the 356 essential genes, a null mutant had been previously reported in the literature.

Other Whole-Genome Experiments

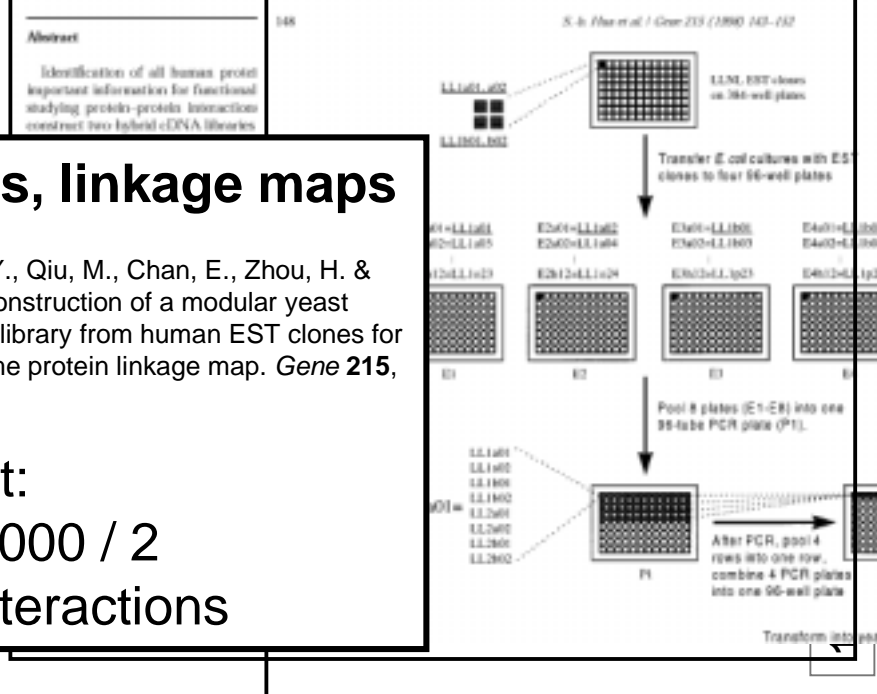
GENE
AN INTERNATIONAL JOURNAL OF GENES AND GENOMES

Elsevier
Gene 215 (1998) 143-152

Construction of a modular yeast two-hybrid cDNA library from human EST clones for the human genome protein linkage map

Shao-bing Hua L.^{*}, Ying Luo L.J., Mengsheng Qiu L.J., Eva Chan L., Helen Zhou L., Li Zhu L.
GeneNet Group, CLONTECH Laboratories Inc., 1020 East Meadow Circle, Palo Alto, CA 94303, USA

Received 1 February 1998; received in revised form 28 April 1998; accepted 28 April 1998; Received by E.Y. Chen



Systematic Knockouts

Winzeler, E. A., Shoemaker, D. D., Astromoff, A., Liang, H., Anderson, K., Andre, B., Bangham, R., Benito, R., Boeke, J. D., Bussey, H., Chu, A. M., Connelly, C., Davis, K., Dietrich, F., Dow, S. W., El Bakkoury, M., Foury, F., Friend, S. H., Gentalen, E., Giaever, G., Hegemann, J. H., Jones, T., Laub, M., Liao, H., Davis, R. W. & et al. (1999). Functional characterization of the *S. cerevisiae* genome by gene deletion and parallel analysis. *Science* **285**, 901-6

2 hybrids, linkage maps

Hua, S. B., Luo, Y., Qiu, M., Chan, E., Zhou, H. & Zhu, L. (1998). Construction of a modular yeast two-hybrid cDNA library from human EST clones for the human genome protein linkage map. *Gene* **215**, 143-52

For yeast:
6000 x 6000 / 2
~ 18M interactions

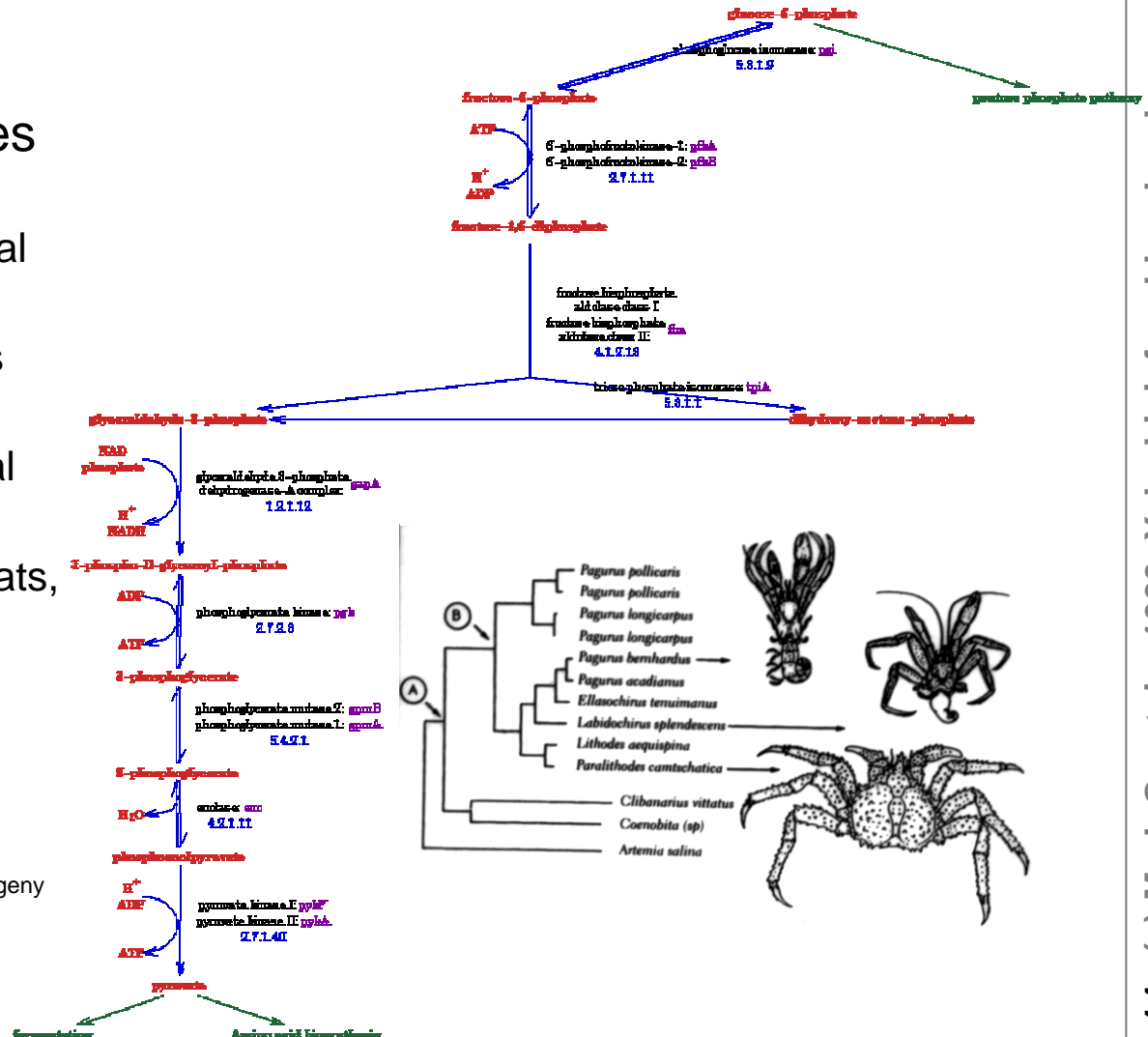
Molecular Biology Information: Other Integrative Data

- Information to understand genomes

- ◇ Metabolic Pathways (glycolysis), traditional biochemistry
- ◇ Regulatory Networks
- ◇ Whole Organisms Phylogeny, traditional zoology
- ◇ Environments, Habitats, ecology
- ◇ The Literature (MEDLINE)

- The Future....

(Pathway drawing from P Karp's EcoCyc, Phylogeny from S J Gould, Dinosaur in a Haystack)



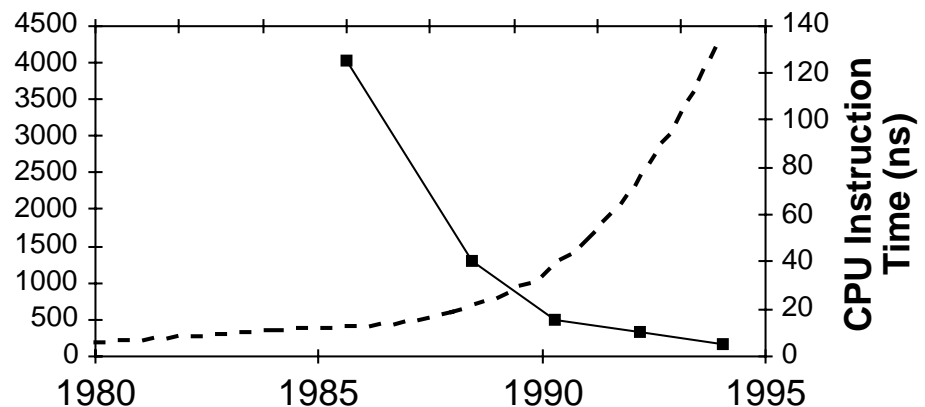
Exponential Growth of Data Matched by Development of Computer Technology

- CPU vs Disk & Net
 - ◇ As important as the increase in computer speed has been, the ability to store large amounts of information on computers is even more crucial

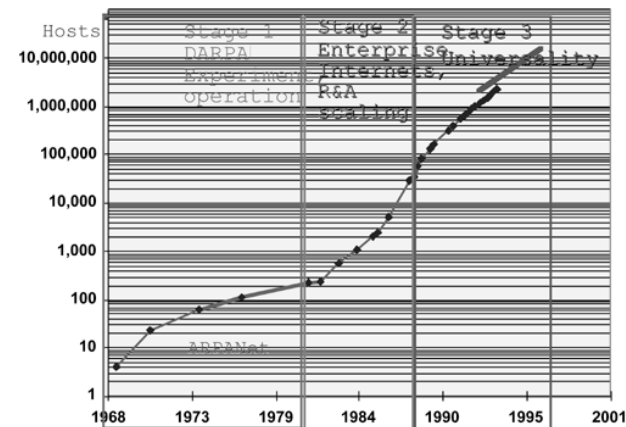
- Driving Force in Bioinformatics

(Internet picture adapted from D Brutlag, Stanford)

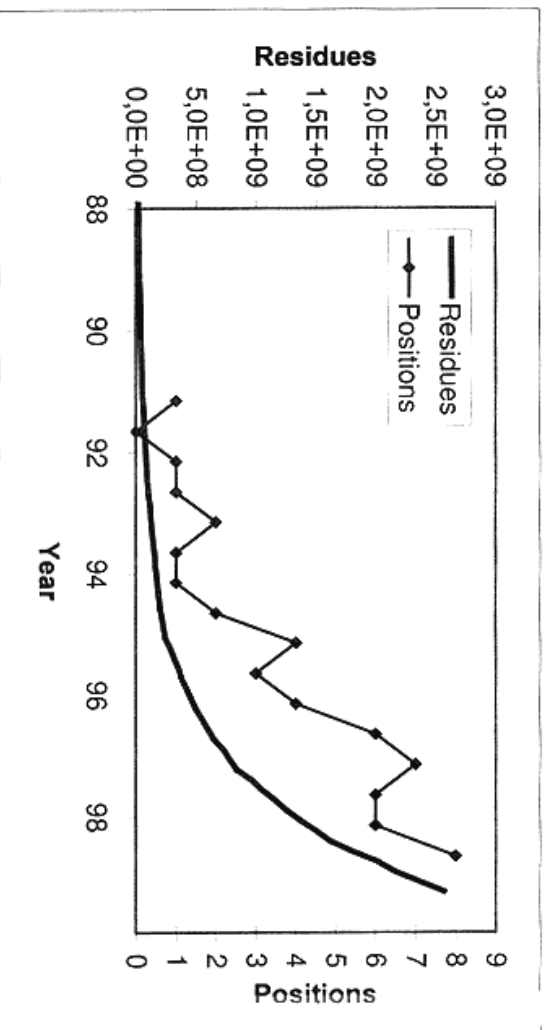
Num.
Protein
Domain
Structures



Internet
Hosts



Bioinformatics is born!



Growth in number of residues in Genbank, a central database for sequence data, compared to the request for people with competence in bioinformatics. The request for scientists is estimated from the number of relevant positions advertised in the first number of Nature in March and September of each year.



B. Watterson, "There's treasure everywhere", Andrews and McMeel, 1996.

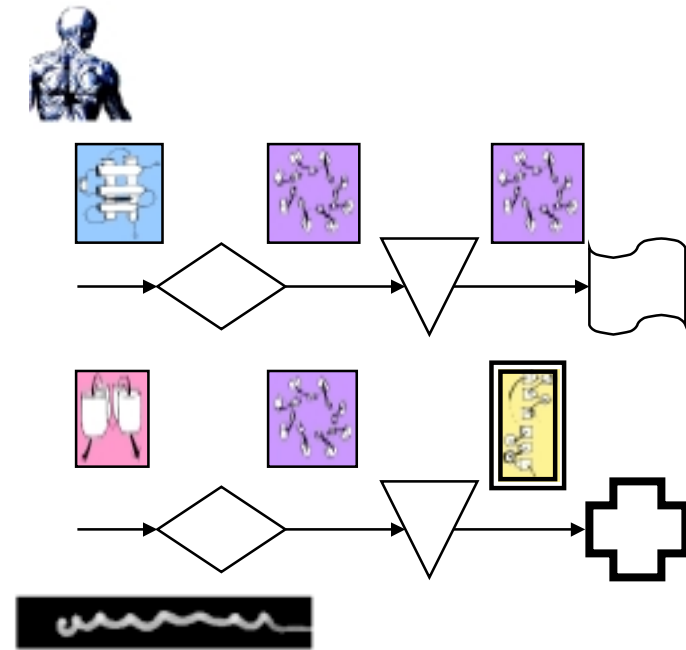
(courtesy of Finn Drablos)

Weber
Cartoon



The Character of Molecular Biology Information: Redundancy and Multiplicity

- Different Sequences Have the Same Structure
- Organism has many similar genes
- Single Gene May Have Multiple Functions
- Genes are grouped into Pathways
- Genomic Sequence Redundancy due to the Genetic Code
- **How do we find the similarities?.....**



Integrative Genomics -
genes ↔ structures ↔
functions ↔ **pathways** ↔
expression levels ↔
regulatory systems ↔

New Paradigm for Scientific Computing

- Because of increase in data and improvement in computers, new calculations become possible
- But Bioinformatics has a new style of calculation...
 - ◇ Two Paradigms
- Physics
 - ◇ Prediction based on physical principles
 - ◇ Exact Determination of Rocket Trajectory
 - ◇ Supercomputer, CPU
- Biology
 - ◇ Classifying information and discovering unexpected relationships
 - ◇ globin ~ colicin~ plastocyanin~ repressor
 - ◇ networks, “federated” database

General Types of “Informatics” in Bioinformatics

- Databases
 - ◇ Building, Querying
 - ◇ Object DB
- Text String Comparison
 - ◇ Text Search
 - ◇ 1D Alignment
 - ◇ Significance Statistics
 - ◇ Alta Vista, grep
- Finding Patterns
 - ◇ AI / Machine Learning
 - ◇ Clustering
 - ◇ Datamining
- Geometry
 - ◇ Robotics
 - ◇ Graphics (Surfaces, Volumes)
 - ◇ Comparison and 3D Matching (Vision, recognition)
- Physical Simulation
 - ◇ Newtonian Mechanics
 - ◇ Electrostatics
 - ◇ Numerical Algorithms
 - ◇ Simulation

Bioinformatics Topics -- Genome Sequence

- Finding Genes in Genomic DNA
 - ◇ introns
 - ◇ exons
 - ◇ promoters
- Characterizing Repeats in Genomic DNA
 - ◇ Statistics
 - ◇ Patterns
- Duplications in the Genome

- Sequence Alignment
 - ◇ non-exact string matching, gaps
 - ◇ How to align two strings optimally via Dynamic Programming
 - ◇ Local vs Global Alignment
 - ◇ Suboptimal Alignment
 - ◇ Hashing to increase speed (BLAST, FASTA)
 - ◇ Amino acid substitution scoring matrices
- Multiple Alignment and Consensus Patterns
 - ◇ How to align more than one sequence and then fuse the result in a consensus representation
 - ◇ Transitive Comparisons
 - ◇ HMMs, Profiles
 - ◇ Motifs

Bioinformatics

Topics --

Protein Sequence

- Scoring schemes and Matching statistics
 - ◇ How to tell if a given alignment or match is statistically significant
 - ◇ A P-value (or an e-value)?
 - ◇ Score Distributions (extreme val. dist.)
 - ◇ Low Complexity Sequences

Bioinformatics

Topics -- Sequence / Structure

- Secondary Structure
“Prediction”

- ◇ via Propensities
- ◇ Neural Networks, Genetic Alg.
- ◇ Simple Statistics
- ◇ TM-helix finding
- ◇ Assessing Secondary Structure Prediction

- Tertiary Structure Prediction

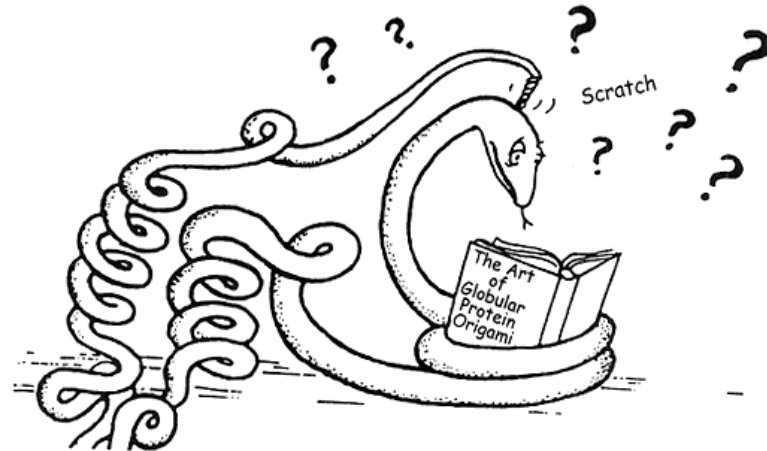
- ◇ Fold Recognition
- ◇ Threading
- ◇ Ab initio

- Function Prediction

- ◇ Active site identification

- Relation of Sequence Similarity to Structural Similarity

“Now collapse down hydrophobic core, and fold over helix 'A' to dotted line, bringing charged residues of 'A' into close proximity to ionic groups on outer surface of helix 'B' ...”



Reproduced in U. Tollemar, “Protein Engineering i USA”, Sveriges Tekniska Attach er, 1988

Topics -- Structures

- Basic Protein Geometry and Least-Squares Fitting
 - ◇ Distances, Angles, Axes, Rotations
 - Calculating a helix axis in 3D via fitting a line
 - ◇ LSQ fit of 2 structures
 - ◇ Molecular Graphics
- Calculation of Volume and Surface
 - ◇ How to represent a plane
 - ◇ How to represent a solid
 - ◇ How to calculate an area
 - ◇ Docking and Drug Design as Surface Matching
 - ◇ Packing Measurement
- Structural Alignment
 - ◇ Aligning sequences on the basis of 3D structure.
 - ◇ DP does not converge, unlike sequences, what to do?
 - ◇ Other Approaches: Distance Matrices, Hashing
 - ◇ Fold Library

Topics -- Databases

- Relational Database Concepts

- ◇ Keys, Foreign Keys
- ◇ SQL, OODBMS, views, forms, transactions, reports, indexes
- ◇ Joining Tables, Normalization
 - Natural Join as "where" selection on cross product
 - Array Referencing (perl/dbm)
- ◇ Forms and Reports
- ◇ Cross-tabulation

- Protein Units?

- ◇ What are the units of biological information?
 - sequence, structure
 - motifs, modules, domains
- ◇ How classified: folds, motions, pathways, functions?

- Clustering and Trees

- ◇ Basic clustering
 - UPGMA
 - single-linkage
 - multiple linkage
- ◇ Other Methods
 - Parsimony, Maximum likelihood
- ◇ Evolutionary implications

- The Bias Problem

- ◇ sequence weighting
- ◇ sampling

Topics -- Genomics

- Expression Analysis
 - ◇ Time Courses clustering
 - ◇ Measuring differences
 - ◇ Identifying Regulatory Regions
- Large scale cross referencing of information
- Function Classification and Orthologs
- The Genomic vs. Single-molecule Perspective
- Genome Comparisons
 - ◇ Ortholog Families, pathways
 - ◇ Large-scale censuses
 - ◇ Frequent Words Analysis
 - ◇ Genome Annotation
 - ◇ Trees from Genomes
 - ◇ Identification of interacting proteins
- Structural Genomics
 - ◇ Folds in Genomes, shared & common folds
 - ◇ Bulk Structure Prediction
- Genome Trees
-

Topics -- Simulation

- Molecular Simulation
 - ◇ Geometry -> Energy -> Forces
 - ◇ Basic interactions, potential energy functions
 - ◇ Electrostatics
 - ◇ VDW Forces
 - ◇ Bonds as Springs
 - ◇ How structure changes over time?
 - How to measure the change in a vector (gradient)
 - ◇ Molecular Dynamics & MC
 - ◇ Energy Minimization
- Parameter Sets
- Number Density
- Poisson-Boltzman Equation
- Lattice Models and Simplification

Bioinformatics Schematic

		<u>Breadth</u> : Homologues, Large-scale Surveys, Informatics →		
		pairwise comparison, sequence & structure alignment	multiple alignment, patterns, templates, trees	databases, scoring schemes, censuses
<u>Depth</u> : Rational Drug Design, physics ↓	gene finding			
		Protein Sequence		
	structure prediction	↓		
		Protein Structure		
	geometry calculation	↓		
		Protein Surface		
	molecular simulation	↓		
		Force Field		
	structure docking	↓		
		Ligand Complex		

Background

	Math	Biology
Need to Know Today	Calculation of Standard Deviation, a Bell-shaped Distribution (of test scores), a 3D vector	DNA, RNA, alpha-helix, the cell nucleus, ATP
What You'll Learn	Force is the Derivative (grad) of Energy, Rotation Matrices (3D), a P-value of .01 and an Extreme Value Distribution	Proteins are tightly packed, sequence homology twilight zone, protein families
Not really necessary....	Poisson-Boltzman Equation, Design a Hashing Function, Write a Recursive Descent Parser	What GroEL does, a worm is a metazoa, E. coli is gram negative, what chemokines are

Are They or Aren't They Bioinformatics? (#1)

- Digital Libraries
 - ◇ Automated Bibliographic Search and Textual Comparison
 - ◇ Knowledge bases for biological literature
- Motif Discovery Using Gibb's Sampling
- Methods for Structure Determination
 - ◇ Computational Crystallography
 - Refinement
 - ◇ NMR Structure Determination
 - Distance Geometry
- Metabolic Pathway Simulation
- The DNA Computer

Are They or Aren't They Bioinformatics? (#1, Answers)

- **(YES?)** Digital Libraries
 - ◇ Automated Bibliographic Search and Textual Comparison
 - ◇ Knowledge bases for biological literature
- **(YES)** Motif Discovery Using Gibb's Sampling
- **(NO?)** Methods for Structure Determination
 - ◇ Computational Crystallography
 - Refinement
 - ◇ NMR Structure Determination
 - **(YES)** Distance Geometry
- **(YES)** Metabolic Pathway Simulation
- **(NO)** The DNA Computer

Are They or Aren't They Bioinformatics? (#2)

- Gene identification by sequence inspection
 - ◇ Prediction of splice sites
- DNA methods in forensics
- Modeling of Populations of Organisms
 - ◇ Ecological Modeling
- Genomic Sequencing Methods
 - ◇ Assembling Contigs
 - ◇ Physical and genetic mapping
- Linkage Analysis
 - ◇ Linking specific genes to various traits

Are They or Aren't They Bioinformatics? (#2, Answers)

- **(YES)** Gene identification by sequence inspection
 - ◇ Prediction of splice sites
- **(YES)** DNA methods in forensics
- **(NO)** Modeling of Populations of Organisms
 - ◇ Ecological Modeling
- **(NO?)** Genomic Sequencing Methods
 - ◇ Assembling Contigs
 - ◇ Physical and genetic mapping
- **(YES)** Linkage Analysis
 - ◇ Linking specific genes to various traits

Are They or Aren't They Bioinformatics? (#3)

- RNA structure prediction
Identification in sequences
- Radiological Image Processing
 - ◇ Computational Representations for Human Anatomy (visible human)
- Artificial Life Simulations
 - ◇ Artificial Immunology / Computer Security
 - ◇ Genetic Algorithms in molecular biology
- Homology modeling
- Determination of Phylogenies Based on Non-molecular Organism Characteristics
- Computerized Diagnosis based on Genetic Analysis (Pedigrees)

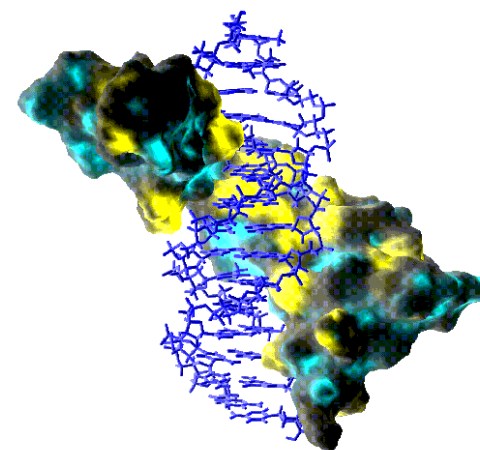
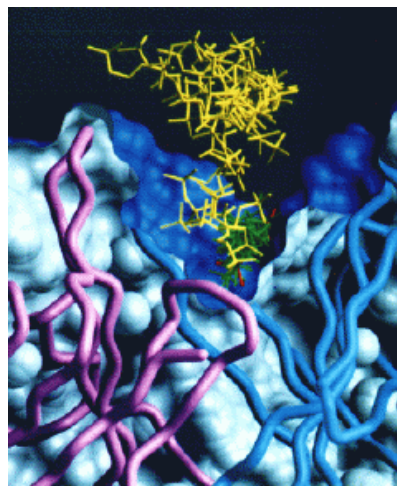
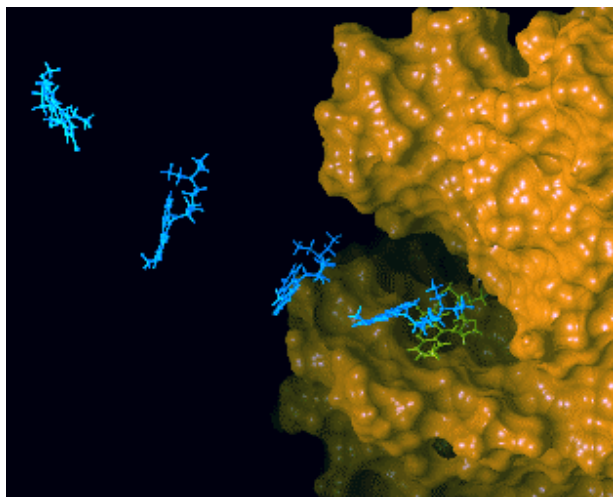
Are They or Aren't They Bioinformatics? (#3, Answers)

- **(YES)** RNA structure prediction
Identification in sequences
- **(NO)** Radiological Image Processing
 - ◇ Computational Representations for Human Anatomy (visible human)
- **(NO)** Artificial Life Simulations
 - ◇ Artificial Immunology / Computer Security
 - ◇ **(NO?)** Genetic Algorithms in molecular biology
- **(YES)** Homology modeling
- **(NO)** Determination of Phylogenies Based on Non-molecular Organism Characteristics
- **(NO)** Computerized Diagnosis based on Genetic Analysis (Pedigrees)

Major Application I: Designing Drugs

- Understanding How Structures Bind Other Molecules (Function)
- Designing Inhibitors
- Docking, Structure Modeling

(From left to right, figures adapted from Olsen Group Docking Page at Scripps, Dyson NMR Group Web page at Scripps, and from Computational Chemistry Page at Cornell Theory Center).



Major Application II: Finding Homologues

- Find Similar Ones in Different Organisms
- Human vs. Mouse vs. Yeast

◇ Easier to do Expts. on latter!

(Section from NCBI Disease Genes Database Reproduced Below.)

Best Sequence Similarity Matches to Date Between Positionally Cloned Human Genes and *S. cerevisiae* Proteins

Human Disease	MIM #	Human Gene	GenBank Acc# for Human cDNA	BLASTX P-value	Yeast Gene	GenBank Acc# for Yeast cDNA	Yeast Gene Description
Hereditary Non-polyposis Colon Cancer	120436	MSH2	U03911	9.2e-261	MSH2	M84170	DNA repair protein
Hereditary Non-polyposis Colon Cancer	120436	MH1	U07418	6.3e-196	MH1	U07187	DNA repair protein
Cystic Fibrosis	219700	CFTF	M28668	1.3e-167	YCF1	L35237	Metal resistance protein
Wilson Disease	277900	WND	U11700	5.9e-161	CCC2	L36317	Probable copper transporter
Glycerol Kinase Deficiency	307030	GK	L13943	1.8e-129	GUT1	X69049	Glycerol kinase
Bloom Syndrome	210900	BLM	U39817	2.6e-119	SGS1	U22341	Helicase
Adrenoleukodystrophy, X-linked	300100	ALD	Z21876	3.4e-107	PXA1	U17065	Peroxisomal ABC transporter
Ataxia Telangiectasia	208900	ATM	U26455	2.8e-90	TEL1	U31331	PI3 kinase
Amyotrophic Lateral Sclerosis	105400	SOD1	K00065	2.0e-58	SOD1	J03279	Superoxide dismutase
Myotonic Dystrophy	160900	DM	L19268	5.4e-53	YPK1	M21307	Serine/threonine protein kinase
Lowe Syndrome	309000	OCRL	M88162	1.2e-47	YIL002C	Z47047	Putative IPP-5-phosphatase
Neurofibromatosis, Type 1	162200	NF1	M89914	2.0e-46	IRA2	M33779	Inhibitory regulator protein
Choroideremia	303100	CHM	X78121	2.1e-42	GDI1	S69371	GDP dissociation inhibitor
Diastrophic Dysplasia	222600	DTD	U14528	7.2e-38	SUL1	X82013	Sulfate permease
Lissencephaly	247200	LIS1	L13385	1.7e-34	MET30	L26505	Methionine metabolism
Thomsen Disease	160800	CLC1	Z25884	7.9e-31	GEF1	Z23117	Voltage-gated chloride channel
Wilms Tumor	194070	WT1	X51630	1.1e-20	FZF1	X67787	Sulphite resistance protein
Achondroplasia	100800	FGFR3	M58051	2.0e-18	IPL1	U07163	Serine/threonine protein kinase
Menkes Syndrome	309400	MNK	X69208	2.1e-17	CCC2	L36317	Probable copper transporter

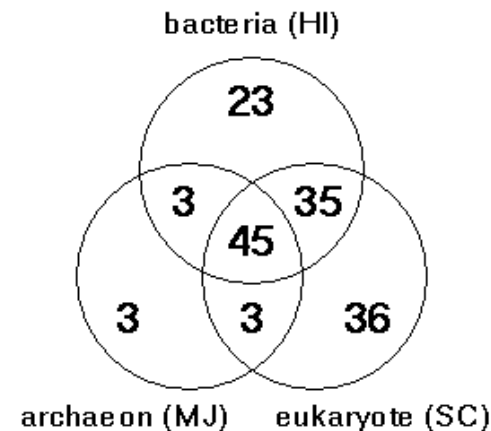
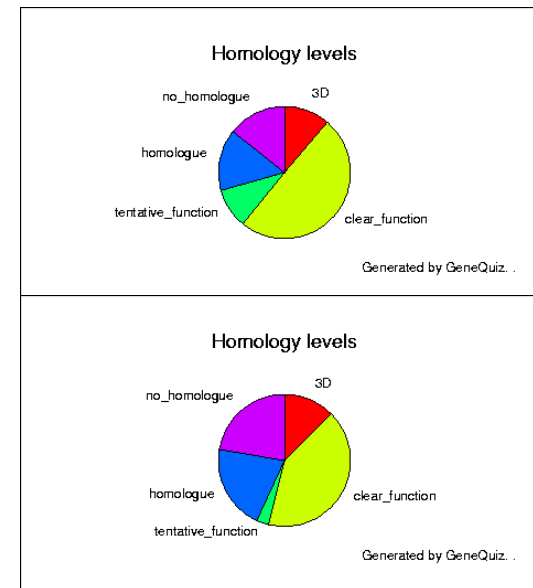
Major Application II: Finding Homologues (cont.)

- Cross-Referencing, one thing to another thing
- Sequence Comparison and Scoring
- Analogous Problems for Structure Comparison
- Comparison has two parts:
 - (1) Optimally **Aligning** 2 entities to get a Comparison **Score**
 - (2) Assessing **Significance** of this score in a given **Context**
- **Integrated Presentation**
 - ◇ Align Sequences
 - ◇ Align Structures
 - ◇ Score in a Uniform Framework

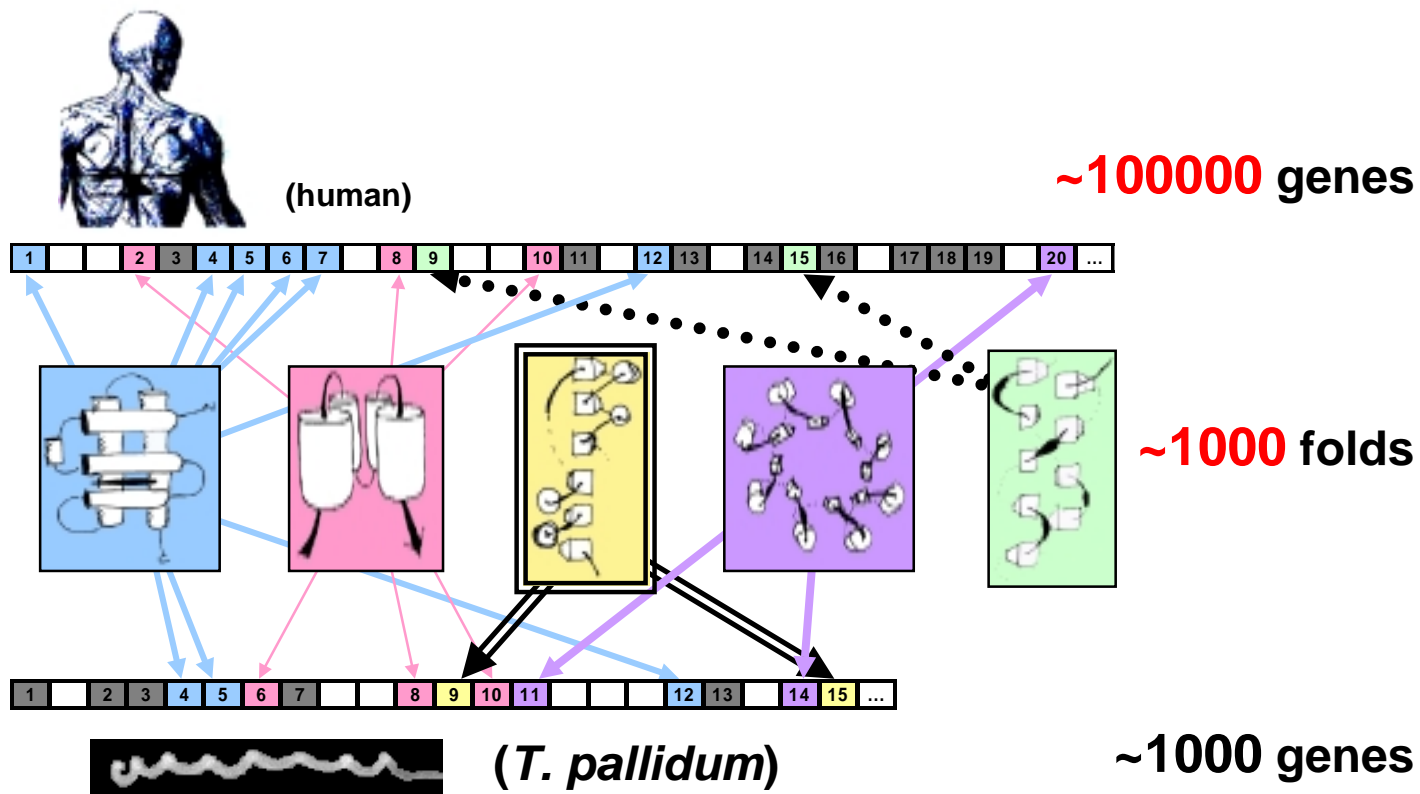
Major Application III: Overall Genome Characterization

- Overall Occurrence of a Certain Feature in the Genome
 - ◇ e.g. how many kinases in Yeast
- Compare Organisms and Tissues
 - ◇ Expression levels in Cancerous vs Normal Tissues
- Databases, Statistics

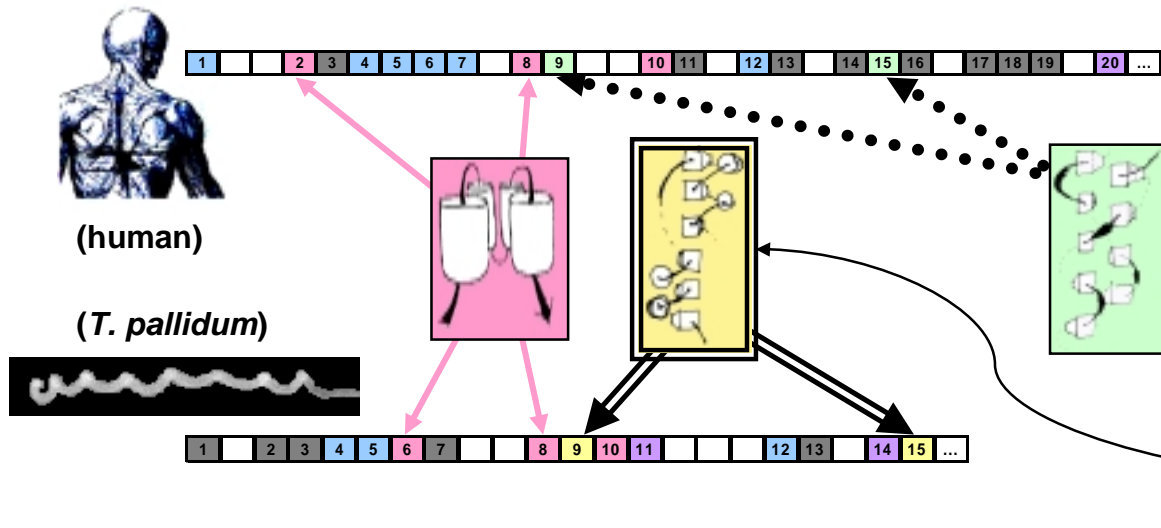
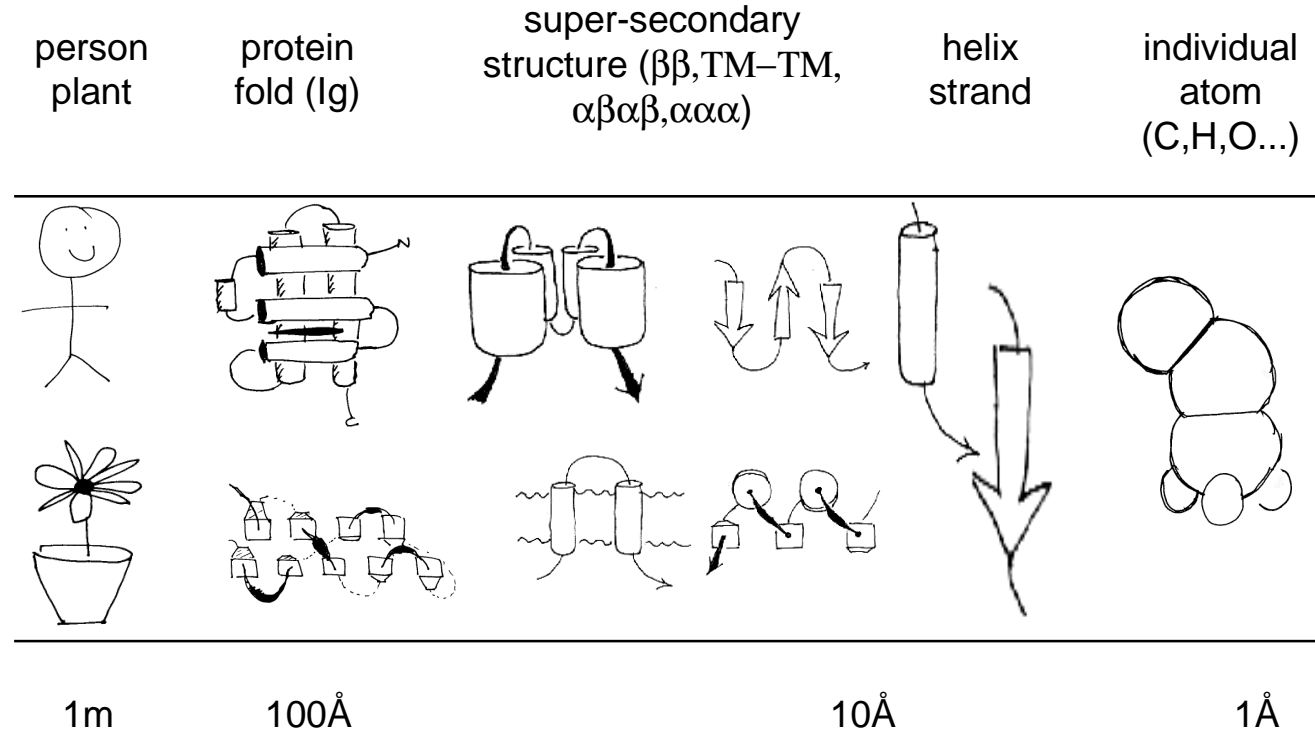
(Clock figures, yeast v. Synechocystis, adapted from GeneQuiz Web Page, Sander Group, EBI)



Simplifying Genomes with Folds, Pathways, &c



At What Structural Resolution Are Organisms Different?



Practical Relevance

(Pathogen only folds as possible targets)