

ERNEST ORLANDO LAWRENCE BERKELEY NATIONAL LABORATORY

Computational Biology and High Performance Computing 2000

Horst D. Simon, Manfred D. Zorn, Sylvia J. Spengler, Brian K. Shoichet, Craig Stewart, Inna L. Dubchak, and Adam P. Arkin

National Energy Research Scientific Computing Division

October 2000

To be presented at *Supercomputing 2000*, Dallas, TX, November 6–10, 2000, and to be published in the Proceedings

DISCLAIMER

This document was prepared as an account of work sponsored by the United States Government. While this document is believed to contain correct information, neither the United States Government nor any agency thereof, nor The Regents of the University of California, nor any of their employees, makes any warranty, express or implied, or assumes any legal responsibility for the accuracy, completeness, or usefulness of any information, apparatus, product, or process disclosed, or represents that its use would not infringe privately owned rights. Reference herein to any specific commercial product, process, or service by its trade name, trademark, manufacturer, or otherwise, does not necessarily constitute or imply its endorsement, recommendation, or favoring by the United States Government or any agency thereof, or The Regents of the University of California. The views and opinions of authors expressed herein do not necessarily state or reflect those of the United States Government or any agency thereof, or The Regents of the University of California.

Ernest Orlando Lawrence Berkeley National Laboratory is an equal opportunity employer.

Computational Biology and High Performance Computing 2000

Horst D. Simon, Manfred D. Zorn, Sylvia J. Spengler, Brian K. Shoichet, Craig Stewart, Inna L. Dubchak, and Adam P. Arkin

National Energy Research Scientific Computing Division Ernest Orlando Lawrence Berkeley National Laboratory University of California Berkeley, California 94720

October 2000



Computational Biology and High Performance Computing 2000

Tutorial M4 a.m. November 6, 2000 SC'2000, Dallas, Texas



Abstract



The pace of extraordinary advances in molecular biology has accelerated in the past decade due in large part to discoveries coming from genome projects on human and model organisms. The advances in the genome project so far, happening well ahead of schedule and under budget, have exceeded any dreams by its protagonists, let alone formal expectations. Biologists expect the next phase of the genome project to be even more startling in terms of dramatic breakthroughs in our understanding of human biology, the biology of health and of disease. Only today can biologists begin to envision the necessary experimental, computational and theoretical steps necessary to exploit genome sequence information for its medical impact, its contribution to biotechnology and economic competitiveness, and its ultimate contribution to environmental quality. High performance computing has become one of the critical enabling technologies, which will help to translate this vision of future advances in biology into reality. Biologists are increasingly becoming aware of the potential of high performance computing. The goal of this tutorial is to introduce the exciting new developments in computational biology and genomics to the high performance computing community.



Introduction

Horst Simon HDSimon@lbl.gov NERSC



Computational Biology and High Performance Computing



† Presenters:

- † Horst D. Simon
 - * Director, NERSC
- * Manfred Zorn
 - † Co-Head, Center of Bioinformatics and Computational Genomics, NERSC
- † Sylvia J. Spengler
 - † Co-Head, Center of Bioinformatics and Computational Genomics, NERSC and Program Director, NSF
- † Craig Stewart
 - † Director, Research & Academic Computing, Indiana University
- † Inna Dubchak
 - † Staff Scientist, NERSC
- * Organizer:
 - † Manfred D. Zorn
- t November 6, 2000



Tutorial Outline



- * 8:30 a.m. 12:00 p.m.
 - * Introduction to Biology
 - † Overview Computational Biology
 - * DNA sequences
- † 1:30 p.m. 5:00 p.m.
 - * Protein Sequences
 - † Phylogeny
 - **† Specialized Databases**

Computational Biology @ SC 2000

ERSC

Tutorial Outline: Morning



* 8:30 a.m. - 8:45 a.m.

Introduction

* 8:45 a.m. - 10:00 a.m.

Biology

† 10:00 a.m. - 10:30 a.m.

BREAK

† 10:30 a.m. - 12:00 p.m.

Working with DNA



Tutorial Outline



- * Introduction
- * Brief Introduction into Biology
- † DNA
 - * What is DNA and how does it work?
 - † What can you do with it?
- * Proteins
 - * What are proteins?
 - * What do we need to know?
- * Phylogeny
- * Specialized Databases

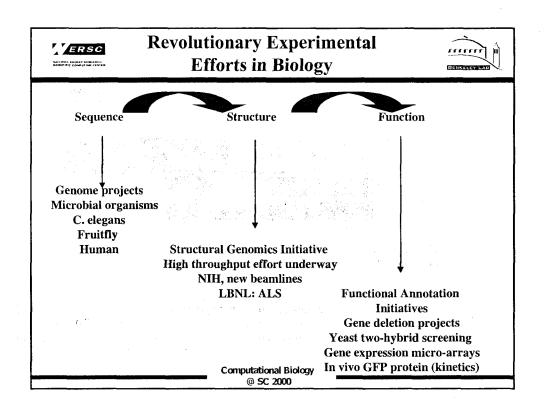
Computational Biology @ SC 2000



Slide Credits



- * Adam Arkin, LBNL
- * Brian Shoichet, NorthWestern Univ.
- * Teresa Head-Gordon, LBNL
- * Sylvia J. Spengler, LBNL
- * Manfred Zorn, LBNL
- * Dodson-Hoagland: "The Way Life Works"
- * National Museum of Health http://www.accessexcellence.org/
- * B. Alberts et al.: "Essential Cell Biology" http://www.essentialcellbiology.com/
- * L. Stryer: Biochemistry
- * Genome Annotation Consortium
- * Bob Robbins, FHCRC





Computational Biology White Paper



http://cbcg.lbl.gov/ssi-csb

A technical document to define areas of biology exhibiting computational problems of scale

Organization:

Introduction to biological complexity and needs for advanced computing (1) Scientific areas (2-6)

Computing hardware, software, CSET issues (7) Appendices

For each scientific chapter:

illustrate with state of the art application (current generation hpc platform) define algorithmic kernals

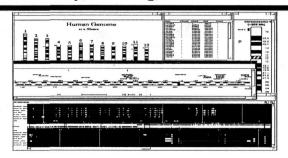
deficiencies of methodologies

define what can be accomplished with 100 teraflop computing



High-Throughput Genome Sequence Assembly, Modeling, and Annotation





The Genome Channel Browser to access and visualize current data flow, analysis and modeling. (Manfred Zorn, NERSC)

Genome sequencing and annotation

Bioinformatics

100,000 human genes; genes from other organism

Structure/functional annotation at the sequence level

Computation to determine regions of a genome that might yield new folds

Experimental Structural Genomics Initiative

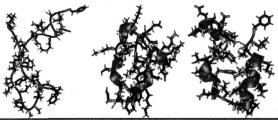
Functional annotation at the structure level by experiment

Computational Biology @ SC 2000



Low Resolution Fold Topologies to High Resolution Structure





One microsecond simulation of a fragment of the protein, Villin. Duan & Kollman, Science 1998

Low Resolution Structures from Predicted Fold Topology

Fold class gives some idea of biological function, but....

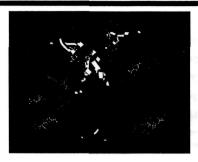


Higher Resolution Structures with Biochemical Relevance Drug design, bioremediation, diseases of new pathogen



Simulating Molecular Recognition/Docking





Changes in the structure of DNA that can be induced by proteins. Through such mechanisms proteins regulate genes, repair DNA, and carry out other cellular functions.

Improvements in Methodology and Algorithms of Higher Resolution Structure Breaking down size, time, lengthscale bottlenecks (IT², algorithms, teraflop computing)

Protein, DNA recognition, binding affinity, mechanism with which drugs bind to proteins

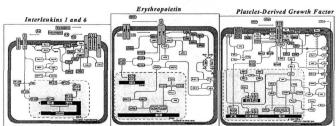
Simulating two-hybrid yeast experiments Protein-protein and Protein-nucleic acid docking

> Computational Biology @ SC 2000



Modeling the Cellular Program





Three mammalian signal transduction pathway that share common molecular elements (i.e. they cross-talk). From the Signaling PAthway Database (SPAD) (http://www.grt.kyushu-u.ac.jp/spad/)

Integrating Computational/Experimental Data at all levels

Sequence, structural functional annotation (Virtually all biological initiatives) Simulating biochemical/genetic networks to mode cellular decisions

Modeling of network connectivity (sets of reactions: proteins, small molecules,

Functional analysis of that network (kinetics of the interactions)



The Need for Advanced Computing for Computational Biology



Computational Complexity arises from inherent factors:

100,000 gene products just from human; genes from many other organisms

Experimental data is accumulating rapidly

N², N³, N⁴, etc. interactions between gene products

Combinatorial libraries of potential drugs/ligands

New materials that elaborate on native gene products from many organisms

Algorithmic Issues to make it tractable

Objective Functions

Optimization

Treatment of Long-ranged Interactions

Overcoming Size and Time scale bottlenecks

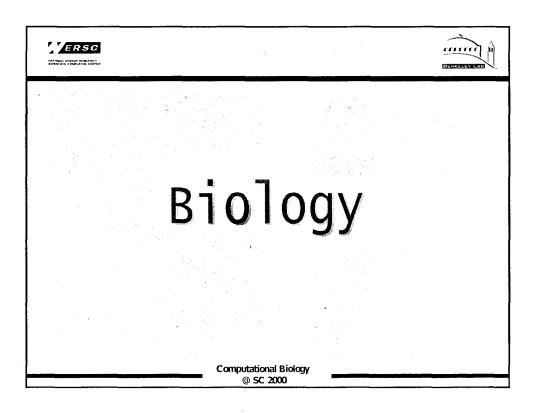
Statistics

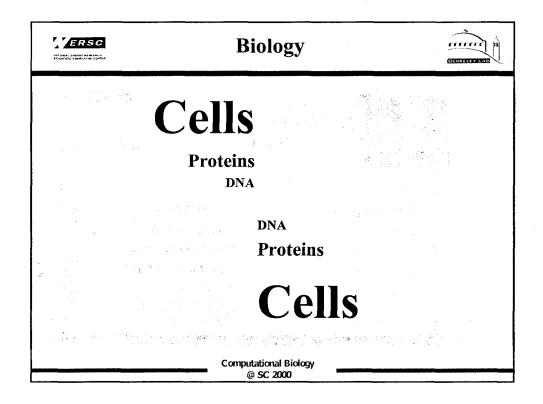
Computational Biology @ SC 2000

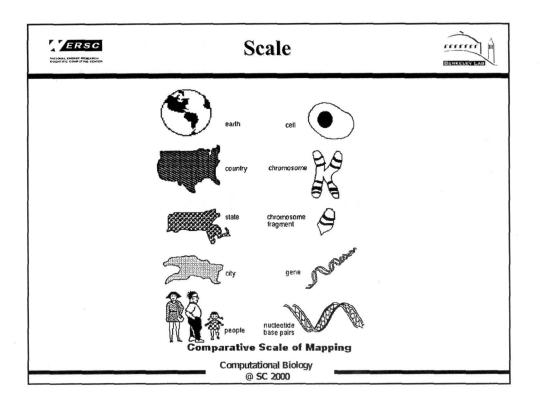


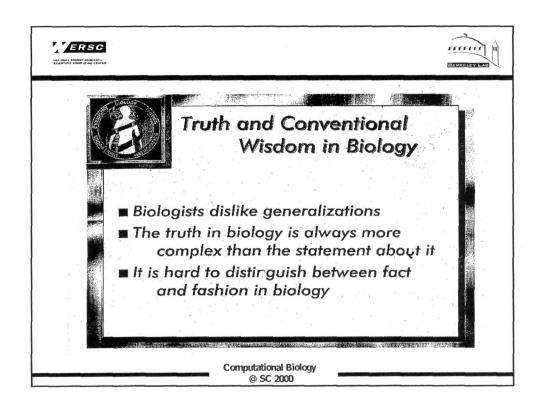
Introduction to Biology

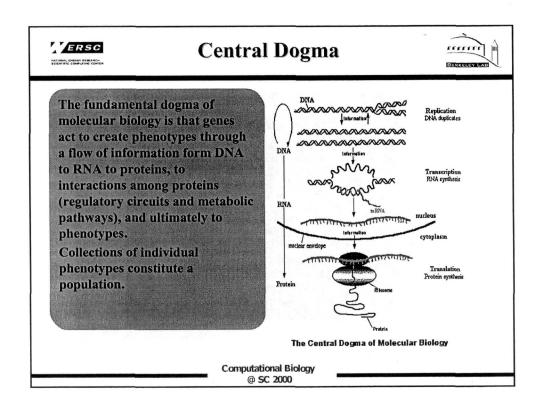
Sylvia Spengler SJSpengler@lbl.gov NERSC













Biology is Special

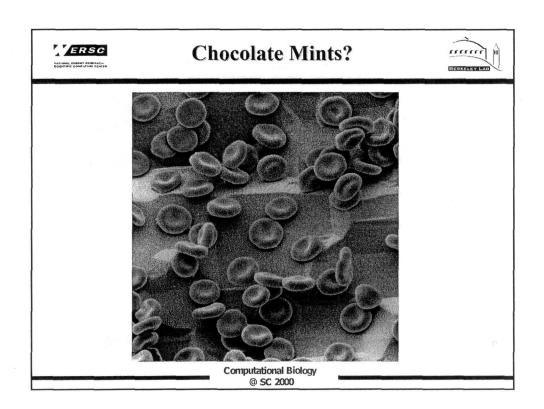


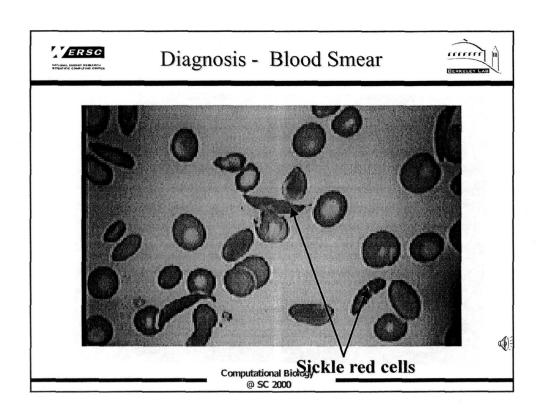
Life is characterized by

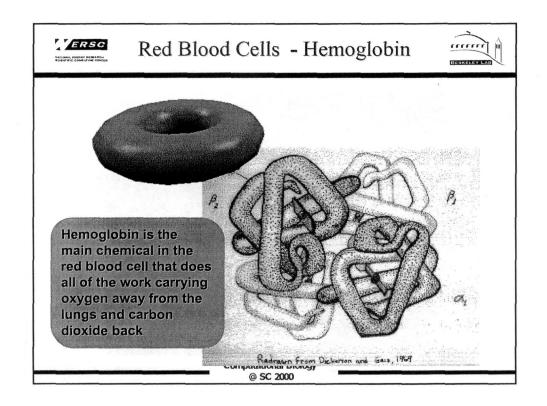
- * Individuality
- * Historicity
- * Contingency
- * high (digital) information content

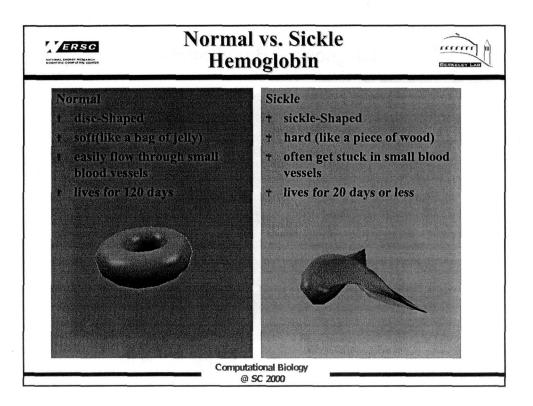
No law of large numbers, since every living thing is genuinely unique.

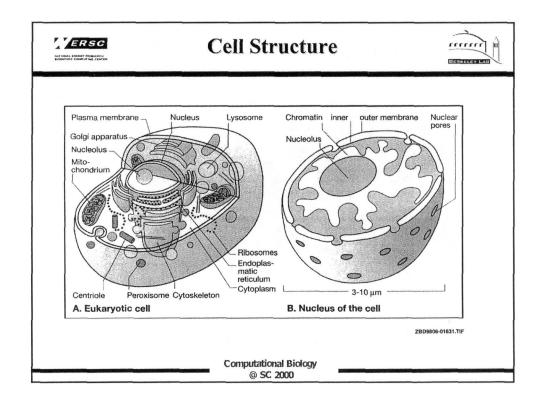


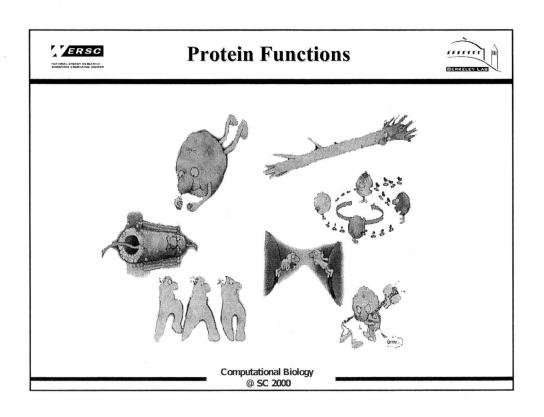


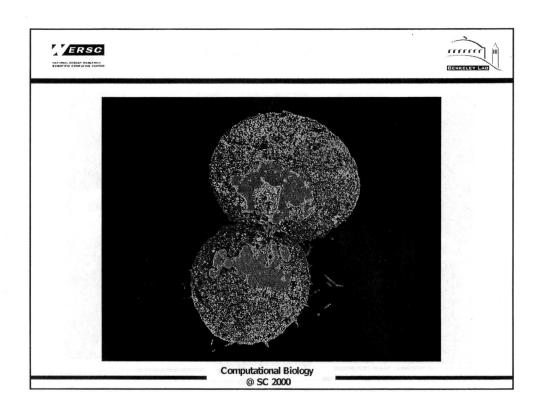


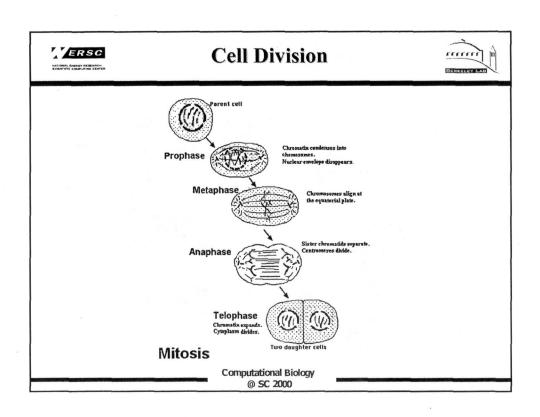


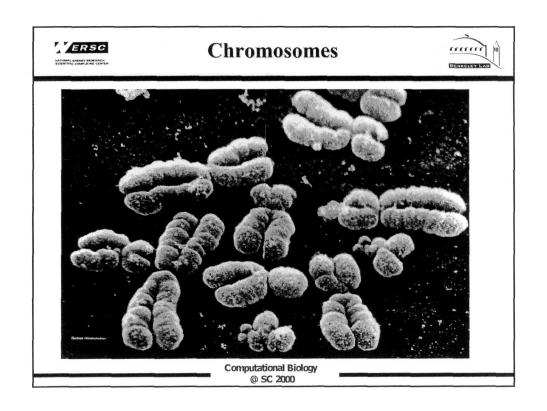


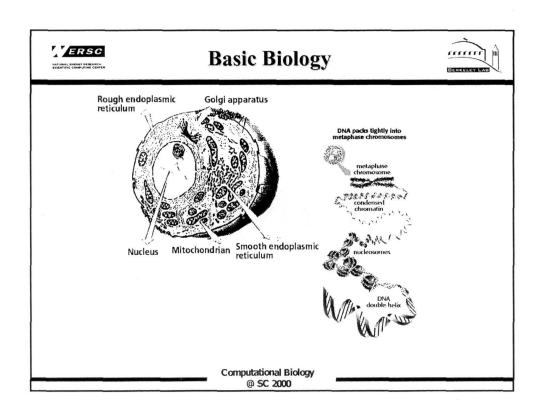


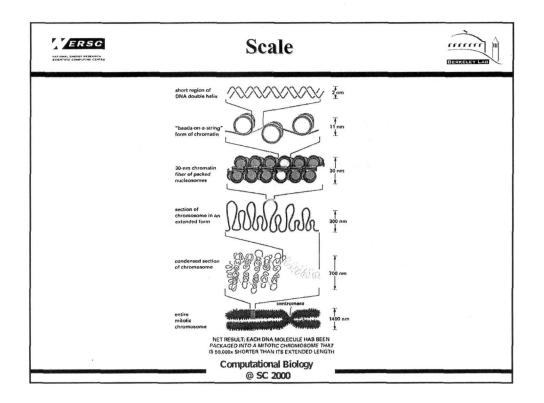








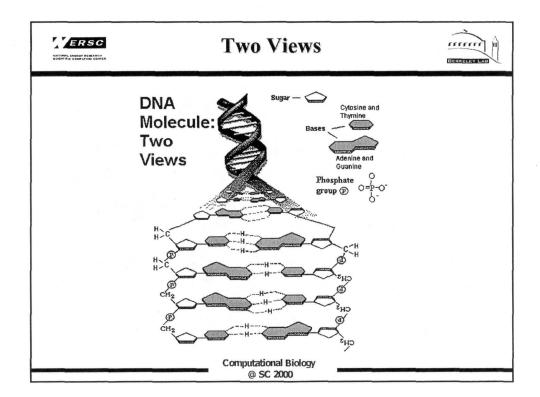


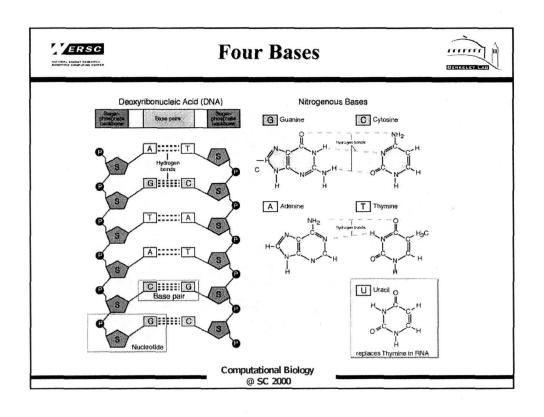


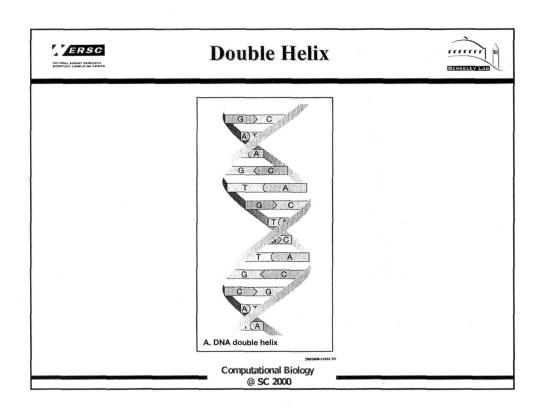


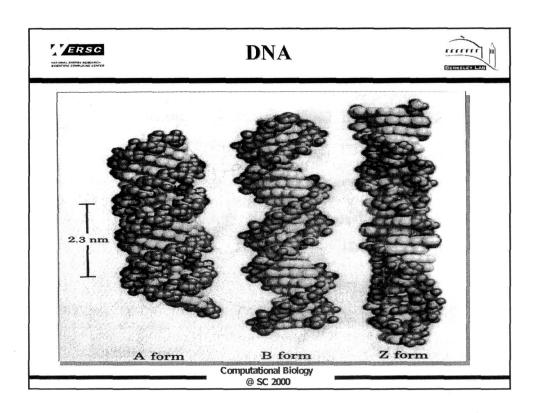


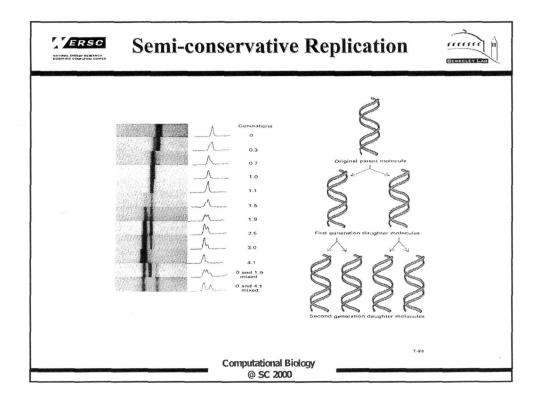
DNA

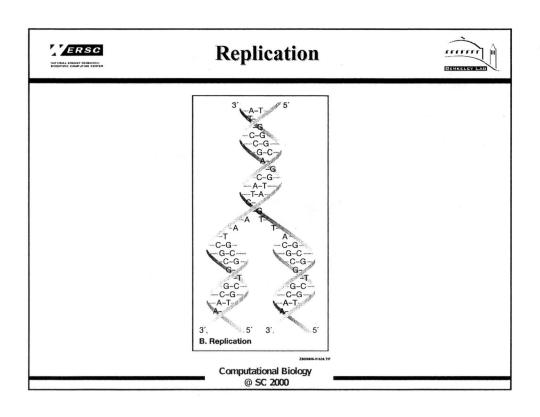


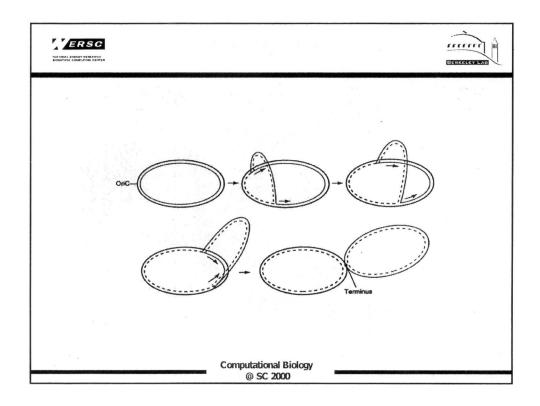


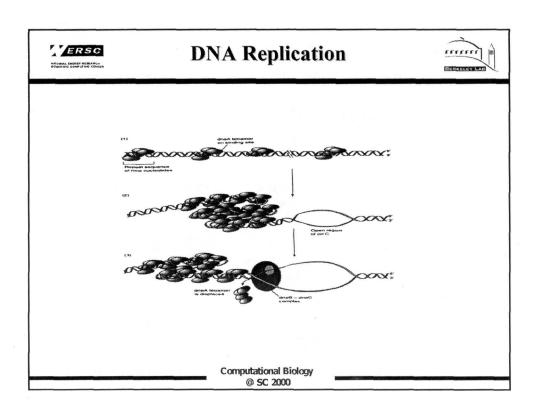


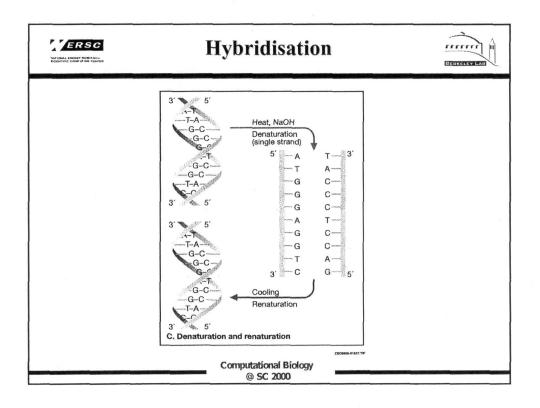


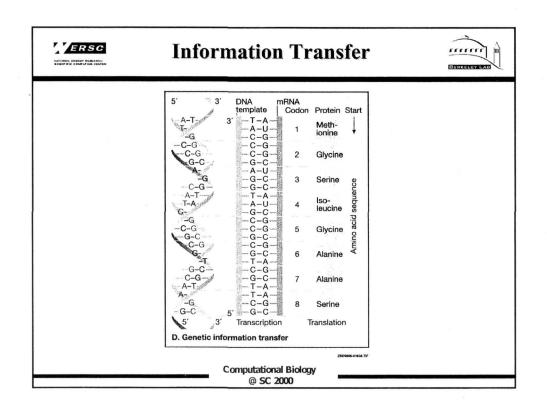


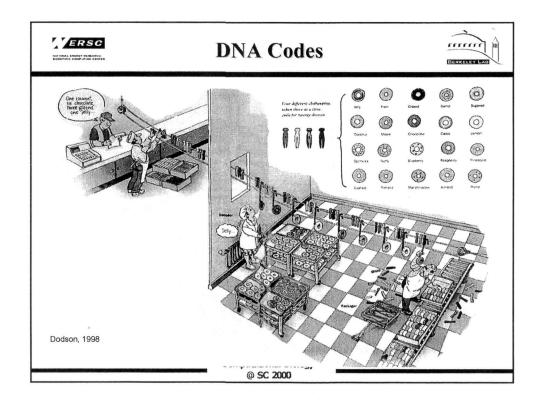


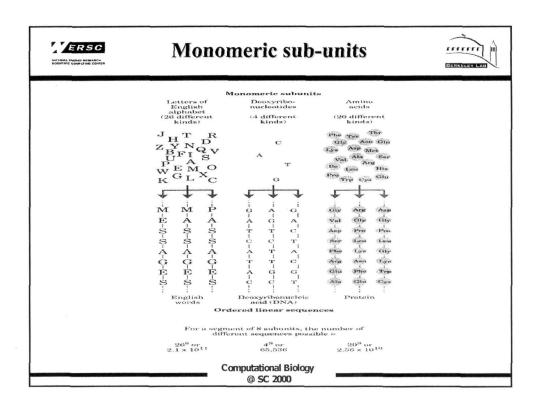


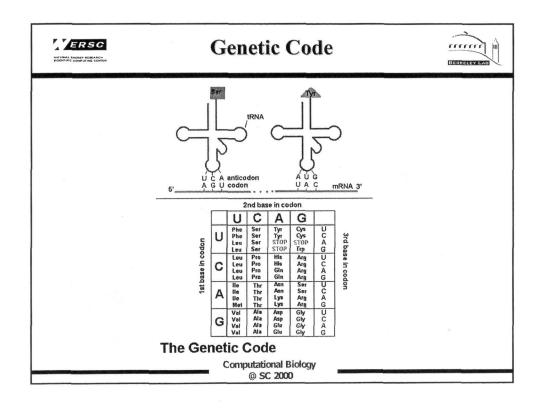


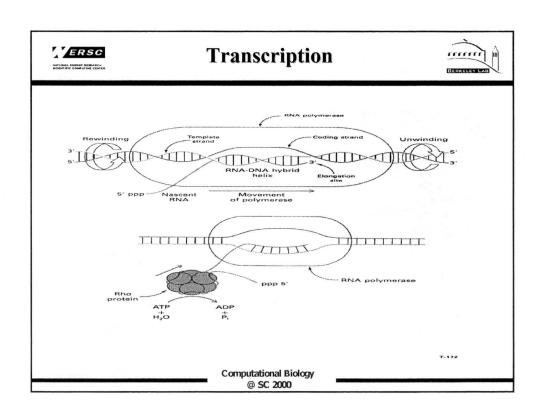


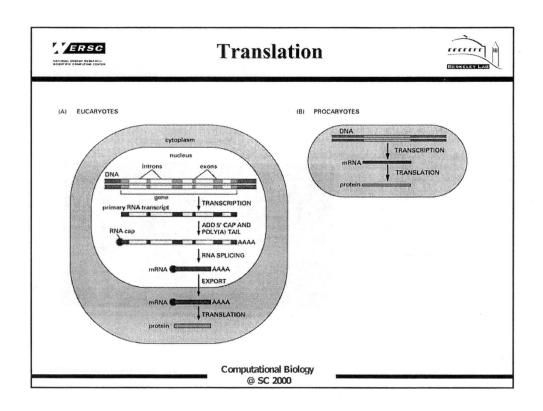


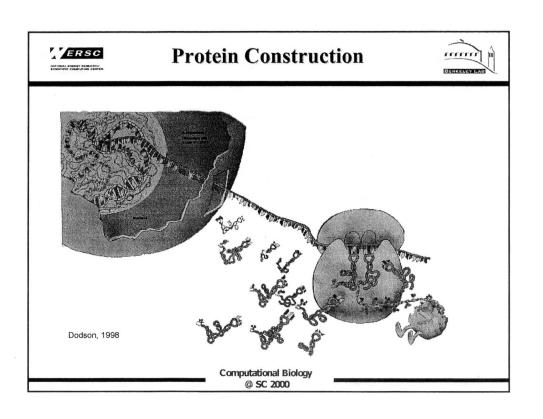


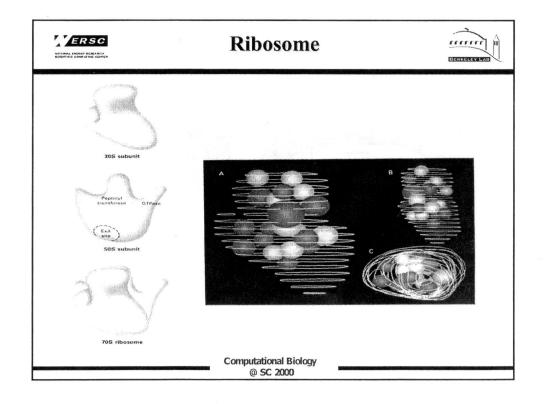


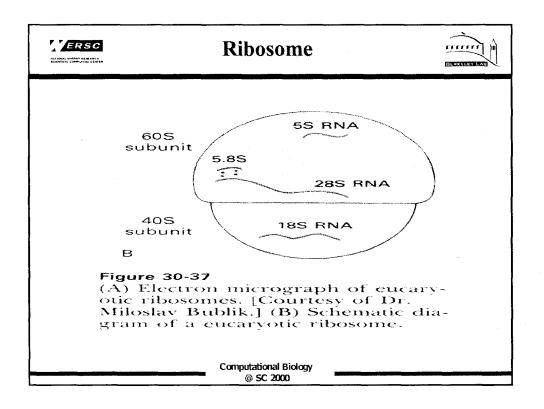


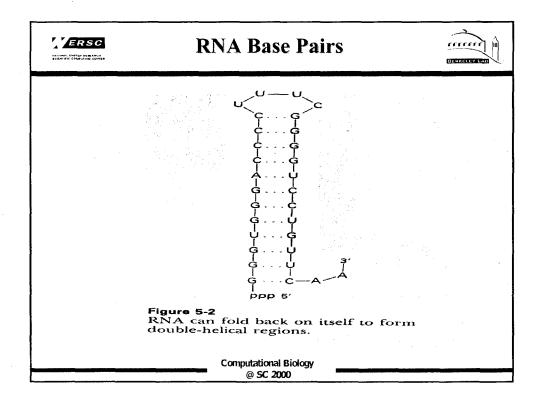


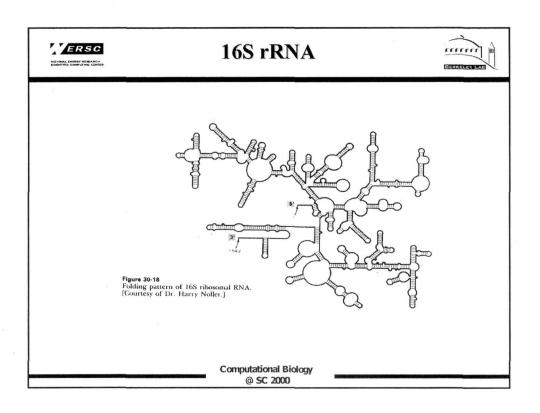


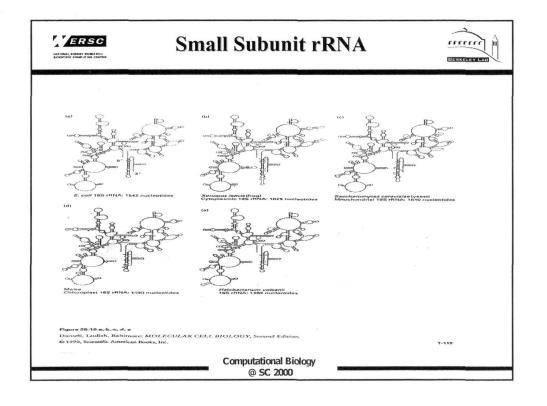


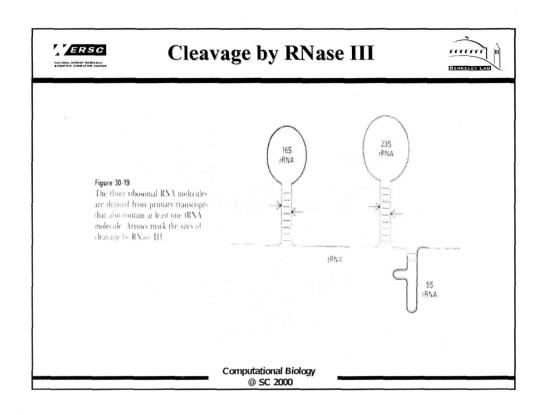


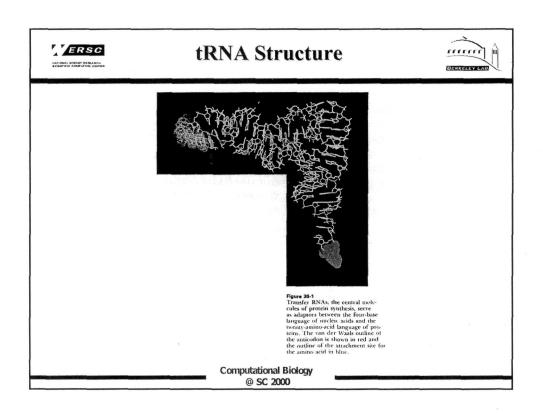


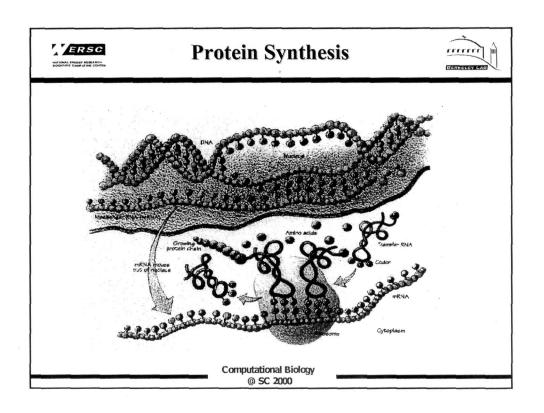


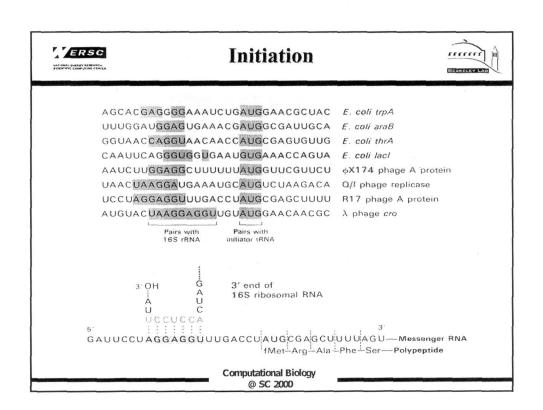


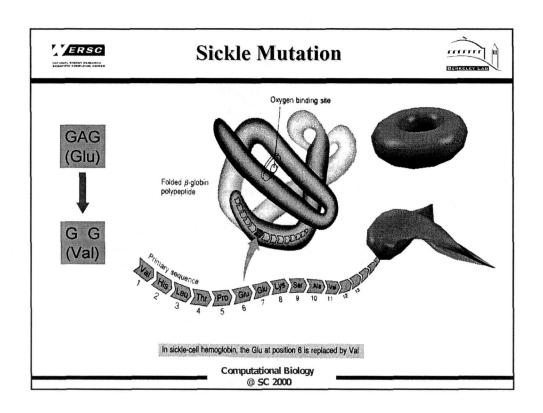


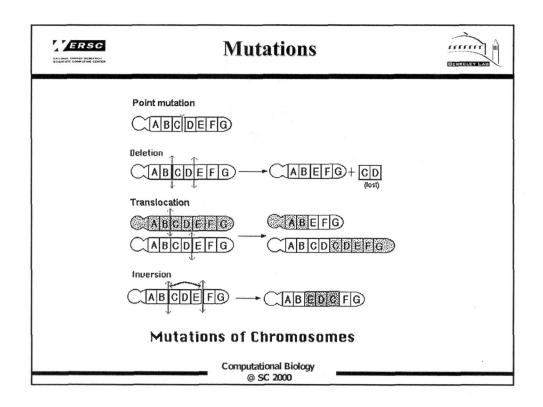


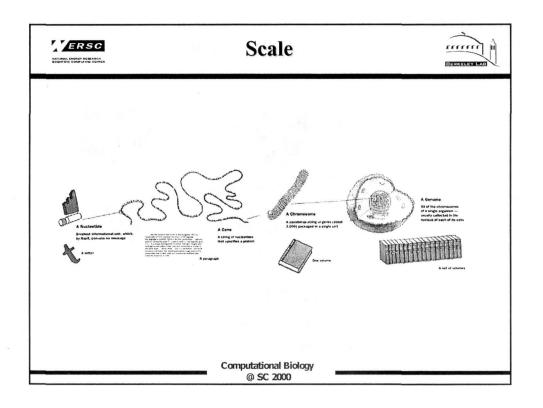


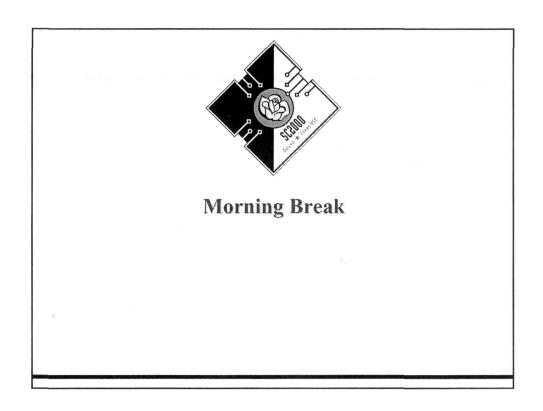














Nucleomics

Manfred Zorn MDZorn@lbl.gov NERSC



Genome Project Timeline



- † 1984
 - * Department of Energy and Intl. Commission on Protection Against Environmental Mutagens and Carcinogens in Alta, Utah.
- † 1986
 - * DOE announces Human Genome Initiative
- † 1987
 - * NIH Director establishes Office of Genome Research
- † 1988
 - * NRC Mapping and Sequencing the Human Genome
 - * Berkeley Lab launches Human Genome Center
- † 1990 Human Genome I



Genome Timeline cont'd

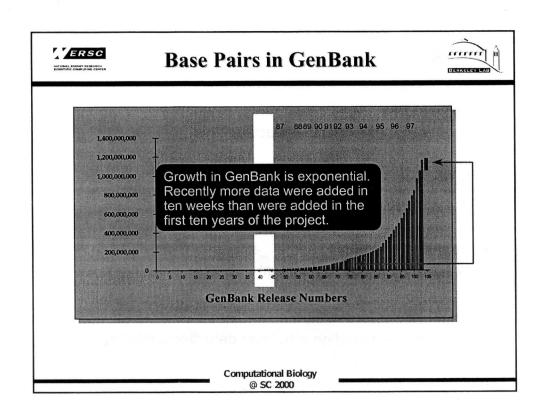


- * September 1994
 - * First complete map of all human chromosomes one year ahead of schedule.
- * May 1995
 - * First genome sequenced: H. inf.
- † May 1998
 - † Celera announces commercial project
 - † Public effort regroups to five major centers
- † June 2000



Computational Biology @ SC 2000

Genome Projects	GENERAL
1995 H. influenzae	2 Mb
1996 S. cerevisiae	12 Mb
1997 E. coli	5 Mb
1998 C. elegans	100 Mb
1999 Human Chromosome 22	34 Mb
2000 D. melanogaster	140 Mb
2000 H. sapiens	3,000 Mb



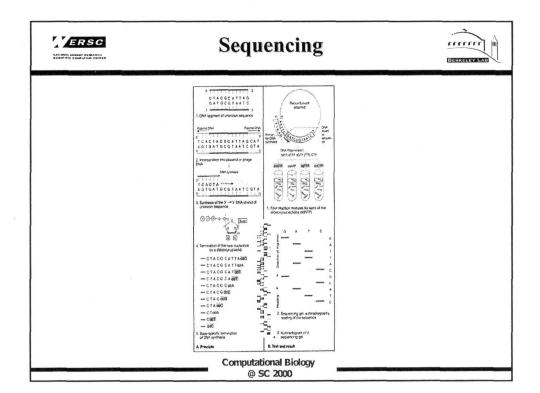
Read base code from storage medium! † Read length: About 600 bases at once † Reader capacity † 100 lanes in parallel in about 2-5 hours † 1000 lanes in parallel in about 2 hours * Computational Biology © SC 2000

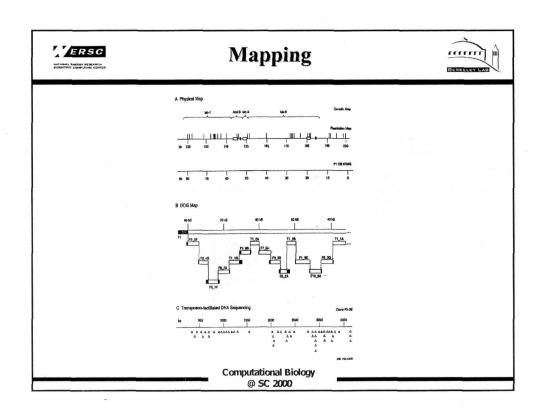


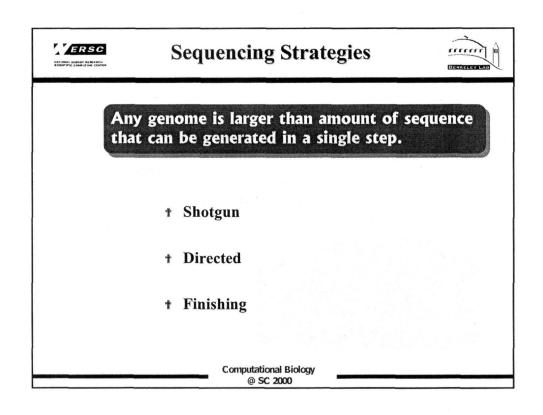
Sequencing: "bird's eye view"



- † Prepare DNA
 - * about a trillion DNA molecules
- * Do the sequencing reactions
 - * synthesize a new strand with terminators
- * Separate fragments
 - * by time, length = constant
- * Sequence determination
 - * automatic reading with laser detection systems









Shotgun



- * Break DNA into manageable pieces
- † Sequence each piece
- † Use sequence to reassemble original DNA

Uniform process
Easily automatable

Computational Biology @ SC 2000



Coverage



Coverage = $\frac{\text{Number x Size of clone}}{\text{Genome size}}$

Expected gaps ~ Number e-coverage

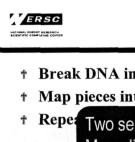
Mapping project (Olson et al. 1986):

N=4,946 L=15,000

G=20,000,000

1,422 contigs vs. 1,457 predicted

Lander-Waterman 1988



Directed



- * Break DNA into manageable pieces
- * Map pieces into tiling path
- Two separate processes: mapping and sequencing More difficult to automate Hard to integrate map information into assembly
- † Transposon mediated sequencing









* Use maps to assemble original DNA

Computational Biology @ SC 2000



Finishing



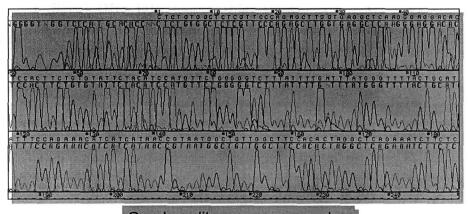
- * Special cases that drop out of the pipeline
- * Gap closing
- Difficult stretches

- † Primer walking
- * Different strains, vectors, chemistry
 - † Creative solutions,



Sequence Traces





Good quality sequence needs about 10X Coverage

Computational biolog @ SC 2000



Base Calling



- * Machine records intensities in each channel
- † Vendor software translates values into smooth signal for each base
- * Base calling software "calls" the sequence
 - † Modern base callers use peak shape, size, and spacing as well as heuristics to improve quality of calls, i.e., fewer N's and better confidence.
 - † Quality values carry base quality to the assembly step.



Phred - Base-caller



- † Developed by Phil Green and Brent Ewing
- † Better base calling accuracy
 - † 40-50% lower error rates than ABI software on large test data sets
- * Error probabilities for each base call
 - † More accurate consensus sequences
 - * Automatic identification of areas that require "finishing" efforts
 - † Identification of repeat sequences in during assembly

Computational Biology

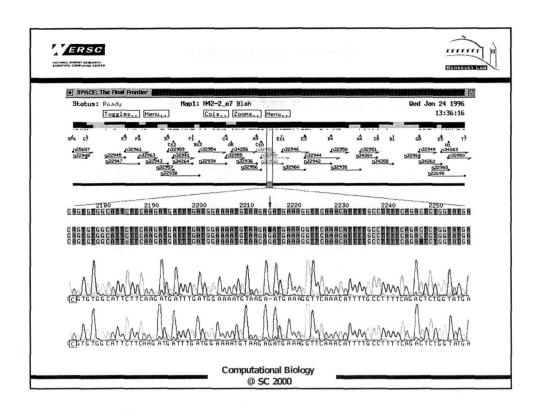
ERSC

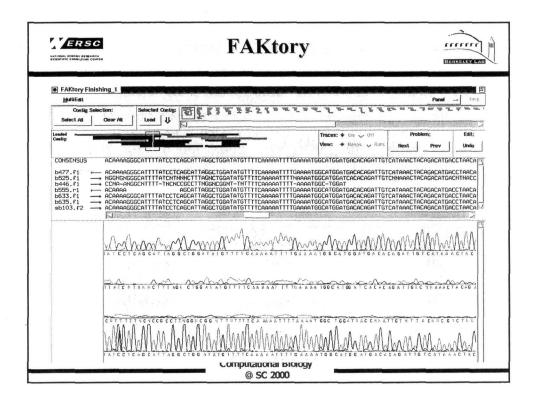
Phred's quality scores



After calling bases, Phred examines the peaks around each base call to assign a quality score to each base call. Quality scores range from 4 to about 60, with higher values corresponding to higher quality. The quality scores are logarithmically linked to error probabilities.

Quality score	Probability of wrong call	Accuracy
10	1 in 10	90%
20	1 in 100	99%
30	1 in 1,000	99.9%
40	1 in 10,000	99.99%
50	1 in 100,000	99.999%







Assembly



Putting humpty-dumpty together again!

- † Overlap
 - † Find overlapping fragments
- * Layout
 - * Order and orientation of fragments
- † Consensus
 - * Determining the consensus sequence
- † Use of constraints

Computational Biology @ SC 2000



Assembly Features



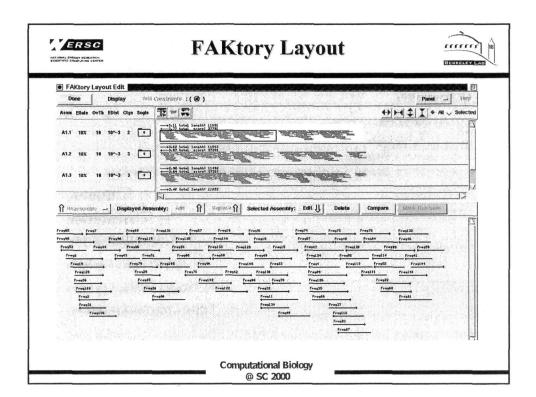
- † Repeats,
 - * repeats,
 - * repeats,
 - † Repeats
 - * 200 bp Alu repeat every ~4,000 bp with 5% -15% error
 - † Clipping
 - † Orientation
 - **†** Contamination
 - * Rearrangements
 - **†** Sequencing errors
 - * True Polymorphisms



Phrap - Assembler



- † Fast assemblies
 - † Projects with several hundred to two thousand reads typically take only minutes
- * Accurate consensus sequences from mosaic
 - * Examines all individual sequences at a given position, and generally uses the highest quality sequence to build the consensus.
- † Consensus quality estimates
 - **†** Quality information of individual sequences yields the quality of the consensus sequence
 - † Other available information about sequencing chemistry (dye terminator or dye primer) and confirmation by "other strand" reads used in estimating the consensus quality.





More assembly



- † Finishing: closing gaps
- † Building chromosomes from large contigs that are consistent with map information

Computational Biology @ SC 2000



What is a Gene?



† Definition: An inheritable trait associated with a region of DNA that codes for a polypeptide chain or specifies an RNA molecule which in turn have an influence on some characteristic phenotype of the organism.

Abstract concept that describes a complex phenomenon



What is Annotation?



t Definition: Extraction, definition, and interpretation of features on the genome sequence derived by integrating computational tools and biological knowledge.

Identifiable features in the sequence

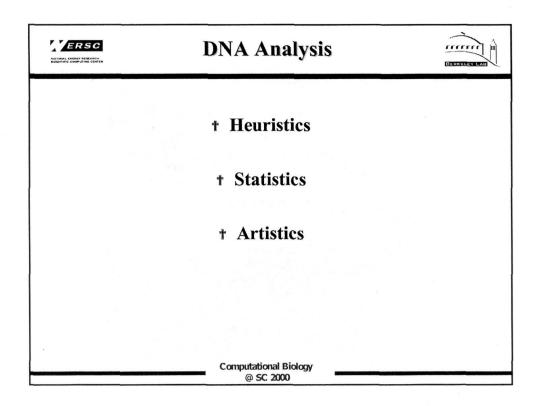
Computational Biology @ SC 2000



How does an annotation differ from a gene?



- * Many annotations describe features that constitute a gene.
- † Other annotations may not always directly correspond in this way, e.g., an STS, or sequence overlap



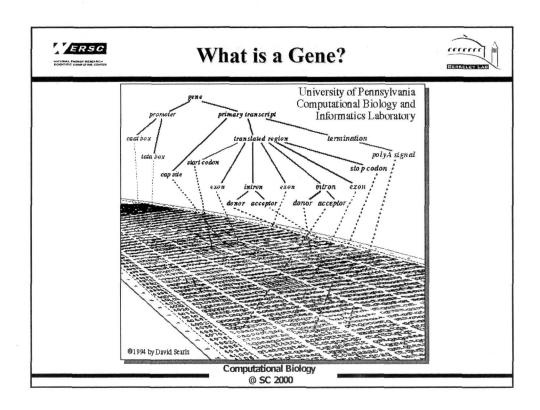


DNA Analysis



Disassemble the base code!

- † Find the genes
 - * Heuristic signals
 - * Inherent features
 - * Intelligent methods
- † Characterize each gene
 - * Compare with other genes
 - * Find functional components
 - * Predict features





Heuristic Signals



DNA contains various recognition sites for internal machinery

- * Promoter signals
- * Transcription start signals
- † Start Codon
- † Exon, Intron boundaries
- * Transcription termination signals



Heuristic Signals



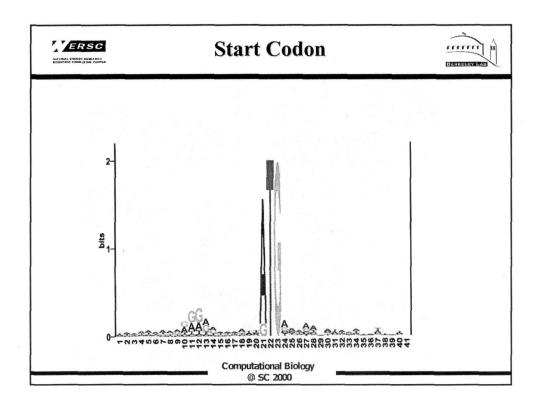
@ SC 2000



Heuristic Signals



@ SC 2000



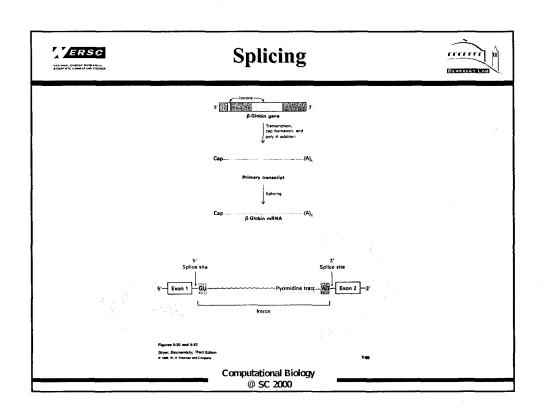


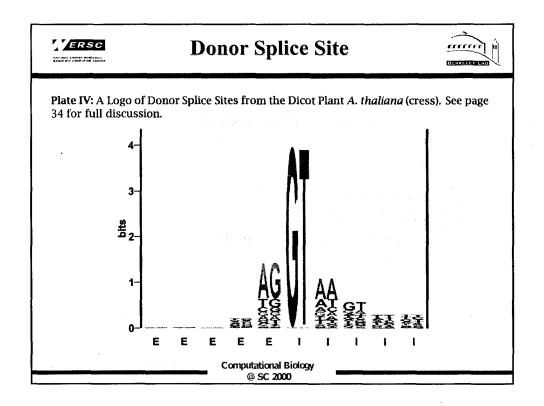
Inherent Features

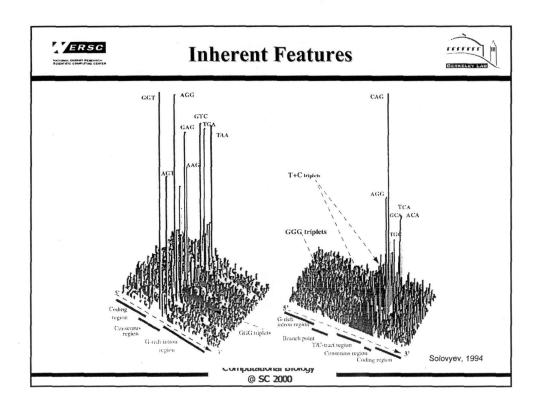


DNA exhibits certain biases that can be exploited to locate coding regions

- † Uneven distribution of bases
- * Codon bias
- † CpG islands
- † In-phase words
- † Encoded amino acid sequence
- * Imperfect periodicity
- * Other global patterns







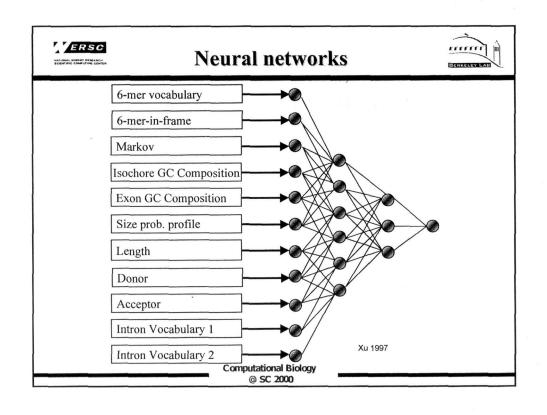


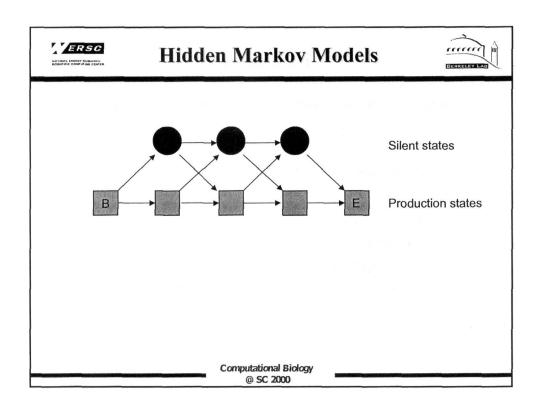
Intelligent Methods



Pattern recognition methods weigh inputs and predict gene location

- * Neural Networks
- † Hidden Markov Models
- † Stochastic Context-Free Grammer







Characterize a Gene



Collect clues for potential function

- † Comparison with other known genes, proteins
- † Predict secondary structure
- **†** Fold classification
- **†** Gene Expression
- † Gene Regulatory Networks
- † Phylogenetic comparisons
- † Metabolic pathways

Computational Biology @ SC 2000



Comparison with other sequences



- † Dynamic programming
 - † Needleman Wunsch
 - * Smith Waterman
 - * Evolution
- * Speed vs. sensitivity
 - * Hashing
 - * Statistical considerations
 - * Suffix trees



Terminology



- * Homology
 - * Common ancestry
 - * Sequence (and usually structure) conservation
 - Homology is not a measurable quantity, but can be inferred, under suitable conditions
- † Identity
 - * Objective and well defined
 - * Can be quantified by several methods:
 - † Percent
 - † The number of identical matches divided by the length of the aligned region
- * Similarity
 - * Most common method used
 - * Not so well defined
 - † Depends on the parameters used (alphabet, scoring matrix, etc.)

Computational Biology @ SC 2000



Alignment



- † An alignment is an arrangement of two sequences opposite one another
- † It shows where they are different and where they are similar

We want to find the optimal alignment - the most similarity and the least differences



Alignment



- † Alignments have two aspects:
 - † Quantity: To what degree are the sequences similar (percentage, other scoring method)
 - † Quality: Regions of similarity in a given sequence

Computational Biology



How is an alignment done?



- t When we compare sequences, we take two strings of letters (nucleotides or amino acids) and align them.
- t Where the characters are identical, we give them a positive score, and where they differ, a negative value.
- t We count the identical and nonidentical characters, and give the alignment a score (usually called the quality)



Dynamic Programming



- * Sequence A
- * Sequence B
- * Substitution
- † Deletion
- * Insertion
- † Matrix Element

$$A = (A_1, ... A_m)$$

$$B = (B_1, ... B_n)$$

$$\omega(A_i, B_j)$$

$$\omega(A_i, \Delta)$$

$$\omega(\Delta, B_j)$$

$$H_{i,j} = \max \begin{cases} H_{i-1,j-1} + \omega_{A_i, B_j} \\ H_{i,j-1} + \omega_{A_i, \Delta} \\ H_{i-1,j} + \omega_{\Delta, B_i} \end{cases}$$

Computational Biology @ SC 2000





Differences in the sequence can be caused by deletions or insertions in the DNA, or by point mutations. These changes can be seen at the protein level as well (changes in the translation of the protein

This scheme works fine as long as you assume that all possible mutations occur at the same frequency. However, nature doesn't work this way. It has been found that in DNA, transitions occur more often than transversions.



Scoring Matrices



- † Identity scoring
- † Genetic code scoring
- † Physical chemical similarities
- t Observed substitutions
 - t Dayhoff matrix (PAM)
 - * BLOSUM

Computational Biology @ SC 2000



The Gap Penalty



Consider the two following alignments:

VITKLGTCVGS VITKLGTCVGS VIT...TCVGS V.TK.GTCV.S

According to the algorithm these 2 cases will get the same gap penalty. However nature is different. In most cases insertions/deletions are longer than a single residue, even for very homologous sequences.





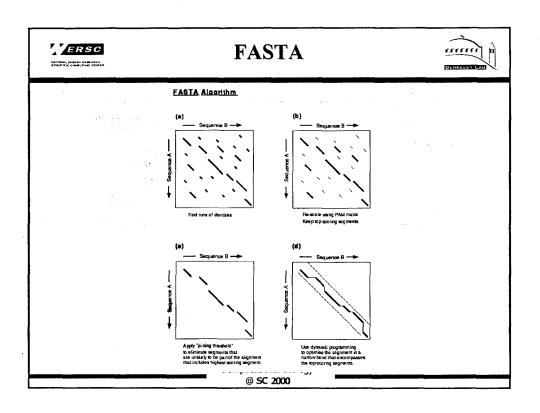
- † To compensate for this, and to differentiate between cases like the one above, the gap penalty is made up of two factors:
 - t The gap creation penalty subtracted from the alignment quality whenever a gap is opened.
 - † The gap extension penalty subtracted from the alignment quality according to the length of the gap.

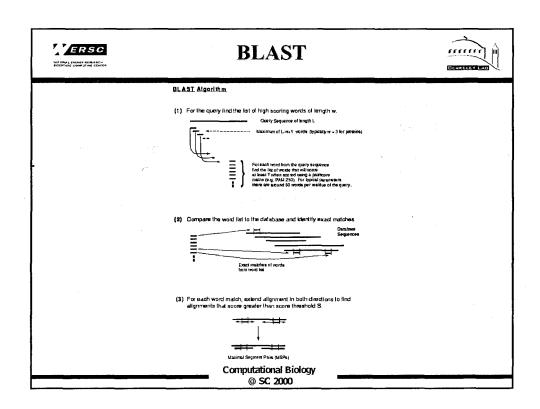
Computational Biology @ SC 2000

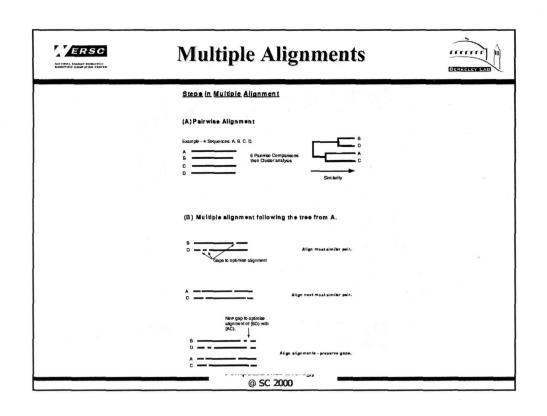


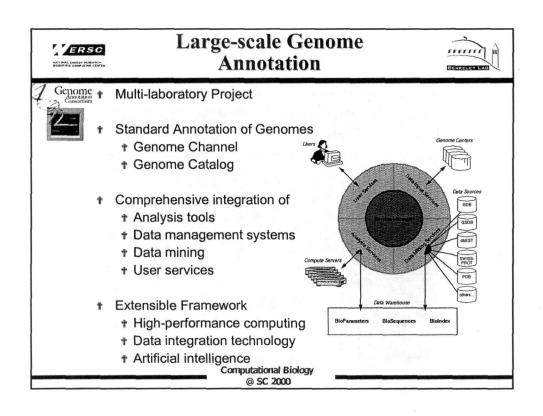


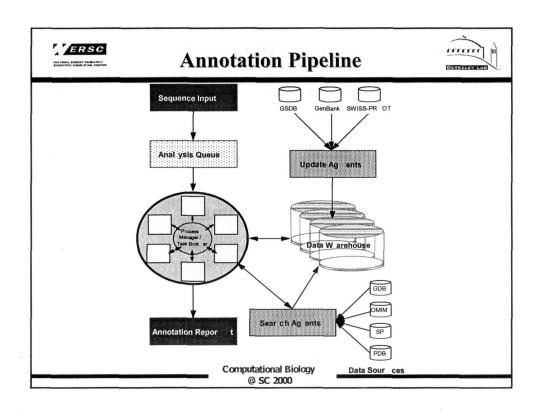
- † Thus we have:
 - † Quality = matches (mismatches + gap penalty)
 - * Gap penalty = gap creation penalty + (gap extension penalty X gap length)

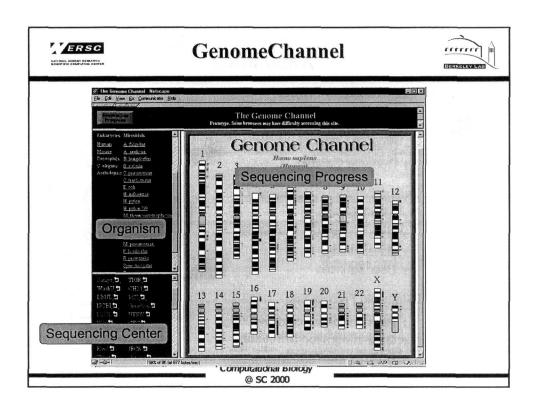


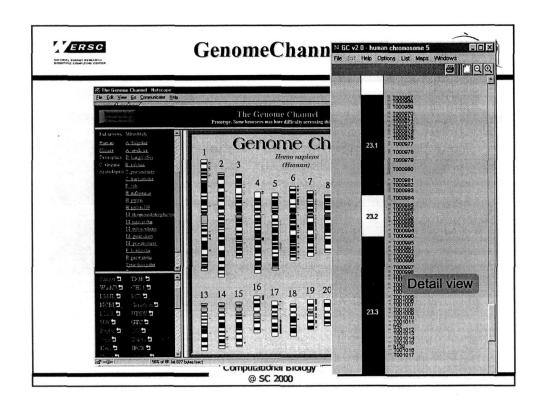


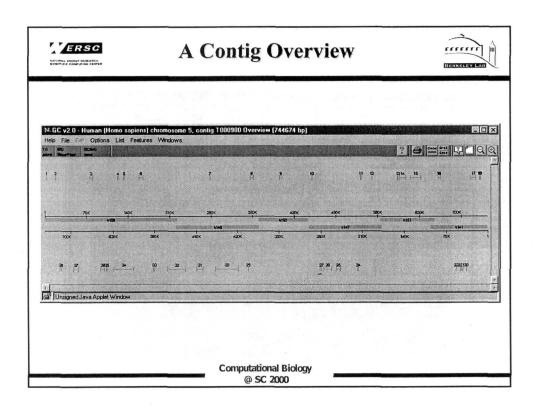


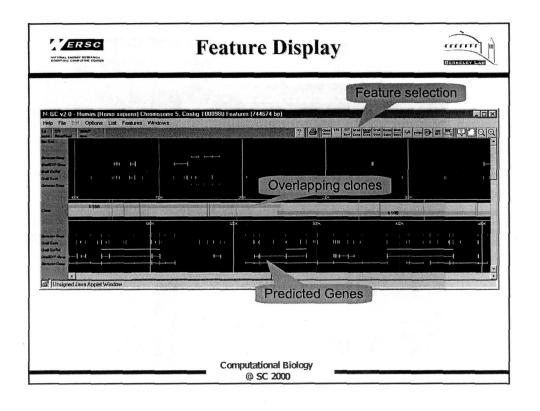


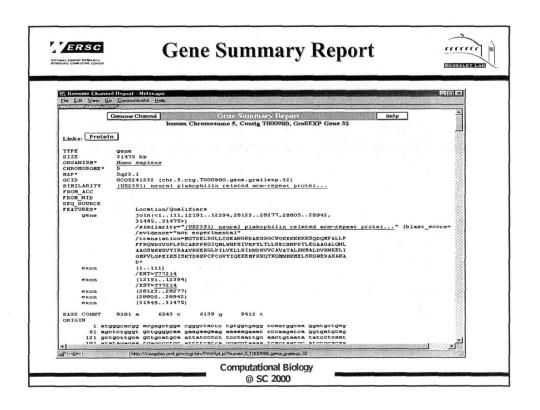






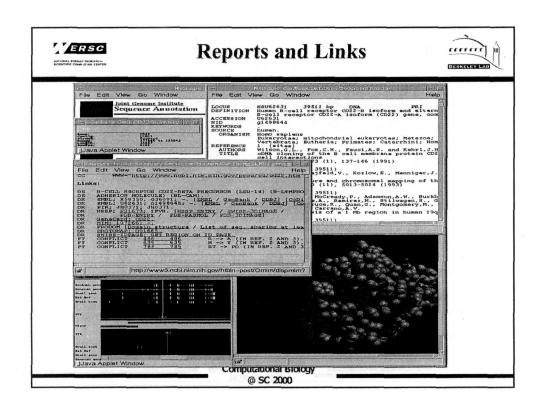


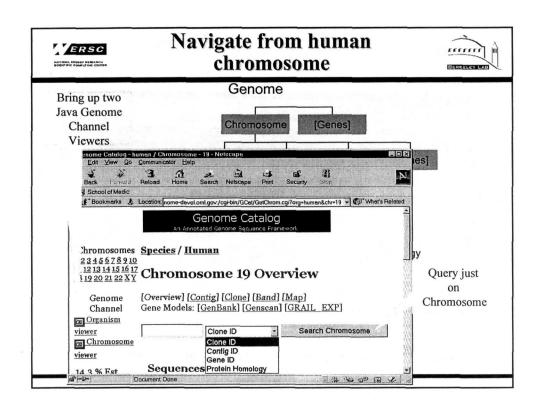


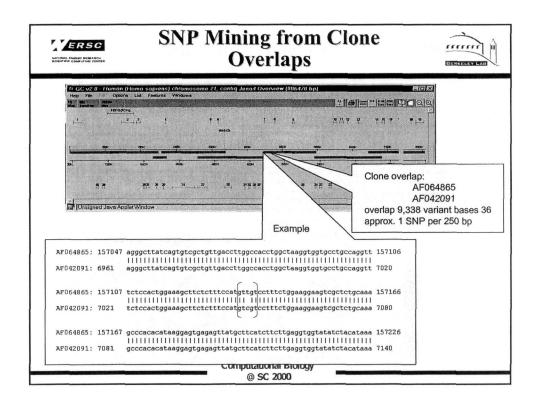


ERSC ALL ENTRY DESERBOON WHITE COMMUTING CENTER	BEAUTY - Gene Search Results	DERKELEY L.
Ele Edi View Go Communicator		
Edition of	Distribution of 29 Blast Hits on the Query Sequence	-
2822195 (US2	351) neural plakophilin related arm-repeat proteinS= 253 E-	6e-67
	Color Key for Alignment Scores	
QUERY #		
QUEKT P	100 200	

A 10.000 (10.0		
Sequences producing sign	Score E	
gi[2822195 (U52351) neura	l plakophilin related arm-repeat protei 253 6e-67	- S
g1[3712673(U96136) delta- g1[2580537(U90331) neura	-catenin [Homo sapiens] 249 9e-66 1 plakophilin related arm-repeat protei 236 9e-62	
g1 1702924 gnl PID e2592	79 (X81889) p0071 protein [Homo sapiens] 163 3e-40	
gi[1932727(US1269) armad gi[2253569(US2626) delta-	illo repeat protein [Homo sapiens] 109 1e-23 -catenin [Homo sapiens] 106 1e-22	
g1 3152867 (AFO62344) p120	catenin isoform 4B [Homo sapiens] 82 4e-15	
	Catenin isoform 2ABC [Homo sapiens] 82 4e-15 Catenin isoform 4A [Homo sapiens] >di 82 4e-15	31
41		· · · · · · · · · · · · · · · · · · ·
http://grad	lsd.ornl.gov/GC/human/chromosome/5/conlig/T000990/gens/gralexp/seerch/beauty/32.html#28221	5
	Computational Biology	
	@ SC 2000	



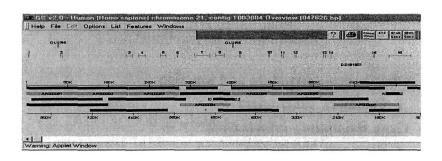






SNP Mining from Clone Overlaps





Coverage includes clones from different sources 1 SNP per 250 bases 160,000 SNPs in 408 Mb dataset

> Computational Biology @ SC 2000



What's supercomputing got to do with it?



- † Complexity of the information
- * Amount of data
- * Most applications are trivially parallel

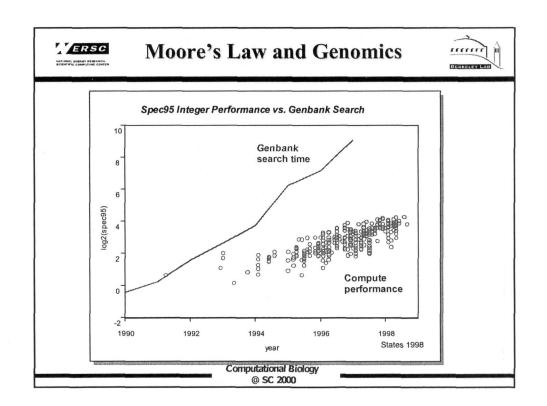


Layers of Information



The same base sequence contains many layered instructions!

- † Chromosome structure and function
 - * Telomers, centromers
- † Gene Regulatory information
 - * Enancers, promoters
- † Instructions for gene structure
- † Instructions for protein
- † Instructions for protein post-processing and localization

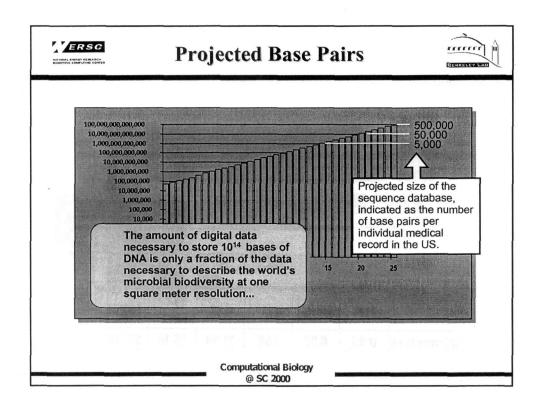




CPU Requirements



- * Current annotation
 - * 250 Mbases DNA yield ~125 Gbytes of data
 - * It takes ~ 7.5 days on 20 workstations ~3,600nhr
- **†** Celera Sequencing
 - * Assembly of 1.7 Million reads in 25 hrs
 - † Annotation 8-10 Mbases per months with 6 FTE
 - † Assembly of Human Genome: expected ~ 3 months

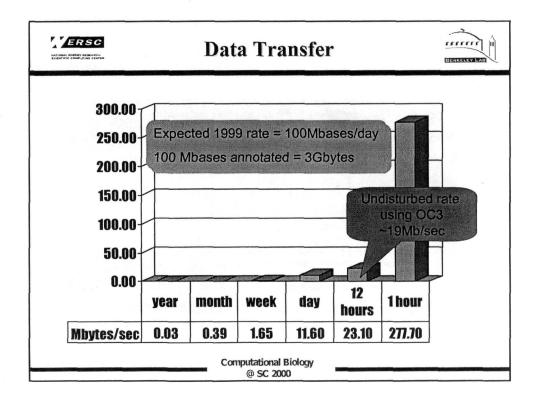




Sequence Assembly



- † Complexity
 - * Adding a day's read of 100 Mb to a billion base pairs of contig would require 100 Pops operations
 - * A 1 Tops machine would take about one day to process 100 Mbases

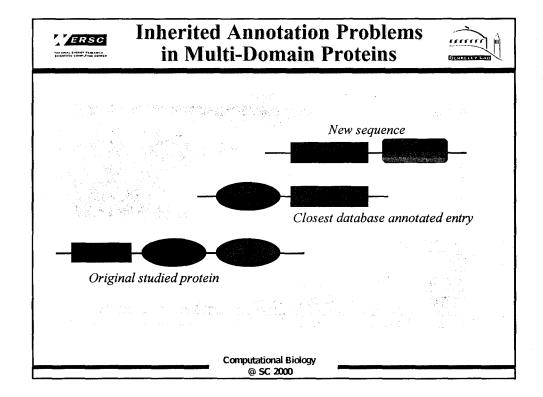


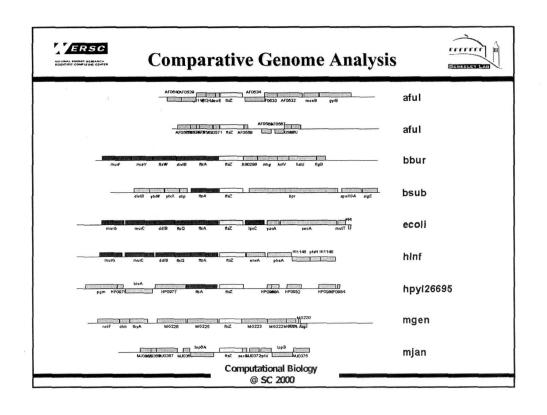


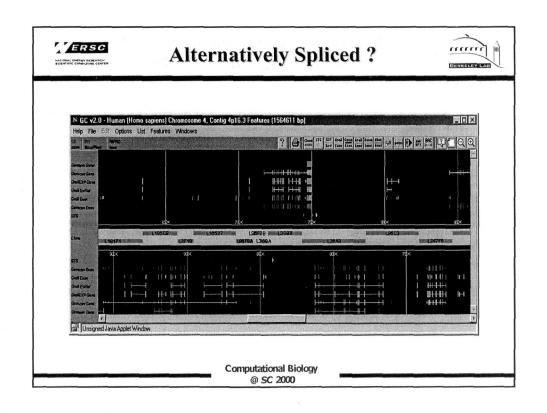
Challenges

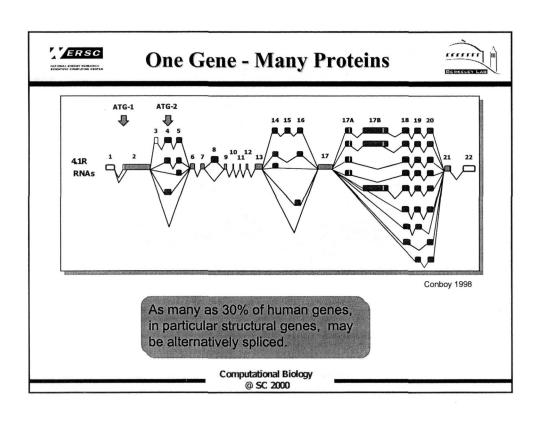


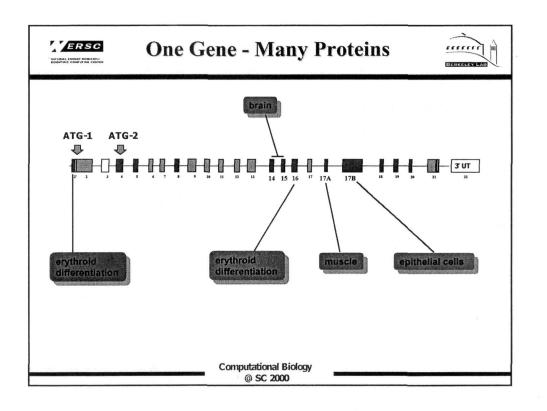
- † Discovering new biology
- † Lack of software integration
- † Beginning to build high-performance applications
- † Shortage of personnel

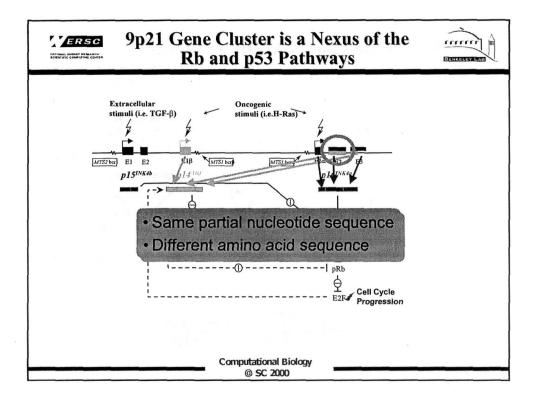












* NERSC/LBNL	† ORNL
† John Conboy	† Ed Uberbacher
† Donn Davy	† Richard Mural
† Inna Dubchak	† Phil LoCascio
† Sylvia Spengler † Denise Wolf	† Sergey Petrov † Manesh Shah
† Eric P. Xing	* Morey Parang
* Manfred Zorn	



Computational Biology and High Performance Computing 2000

Tutorial M4 p.m. November 6, 2000 SC'2000, Dallas, Texas



Tutorial Outline



- * 8:30 a.m. 12:00 p.m.
 - + Introduction to Biology
 - **†** Overview Computational Biology
 - † DNA sequences
- † 1:30 p.m. 5:00 p.m.
 - * Protein Sequences
 - * Phylogeny
 - * Specialized Databases



Tutorial Outline: Afternoon



† 1:30 p.m. - 2:00 p.m.

Working with Proteins

† 2:00 p.m. - 3:00 p.m.

Phylogeny

† 3:00 p.m. - 3:30 p.m.

BREAK

Called Market Carlotte

† 3:30 p.m. - 4:30 p.m.

Specialized Databases

† 4:30 p.m. - 5:00 p.m.

Genetic Networks

Computational Biology



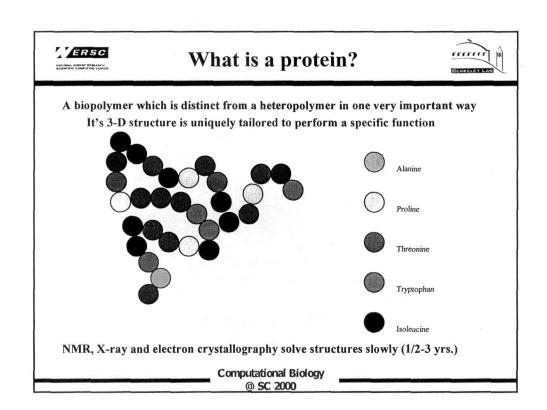
Proteins

Manfred Zorn MDZorn@lbl.gov NERSC

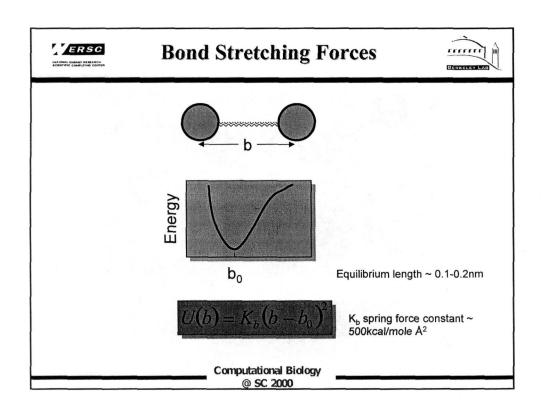


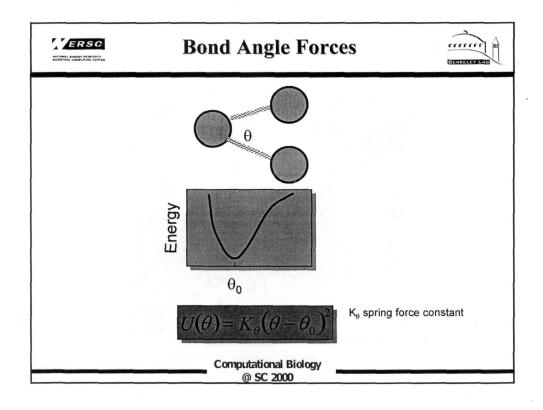


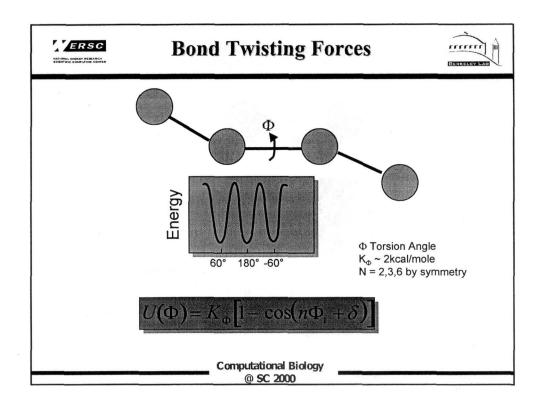
Proteins

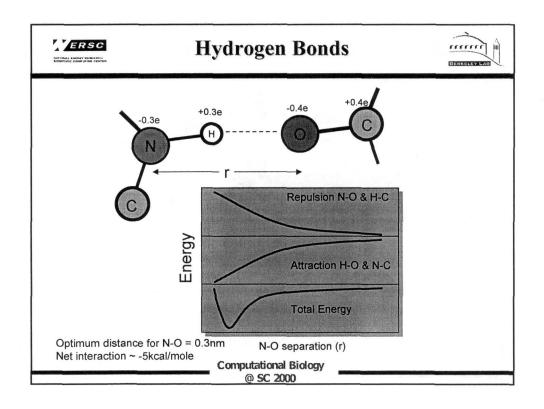


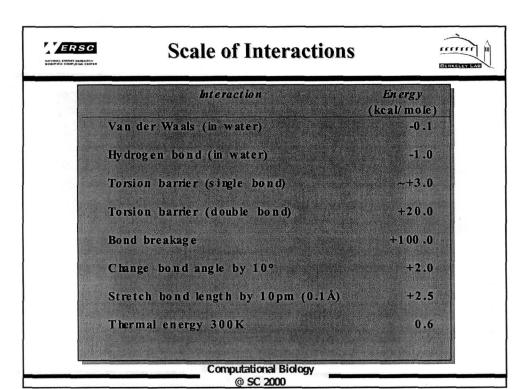
Forces Between Atoms Force between pair of atoms is the same in isolation or when part of a big molecule * Basic assumptions: * Energy contributions are strictly additive * Energy is independent of neighbors; transferability * Quantum mechanics is insignificant as long as no bonds are broken Computational Biology © SC 2000



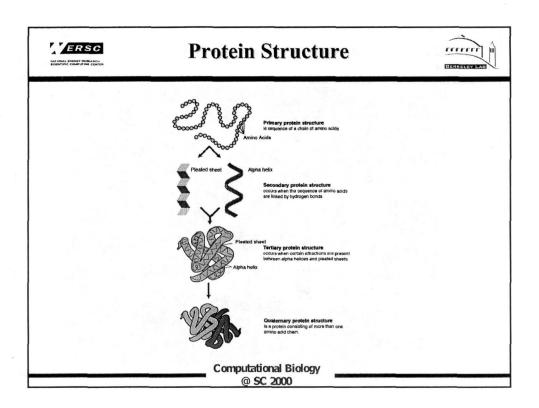


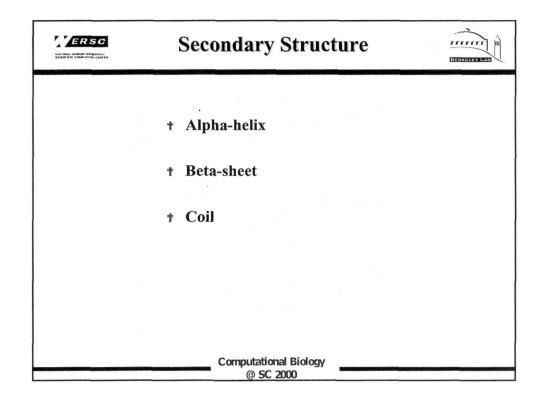


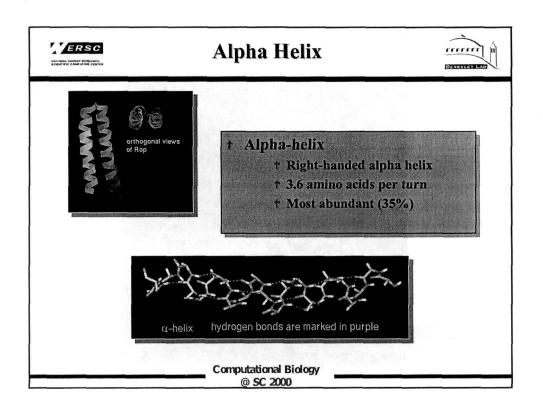


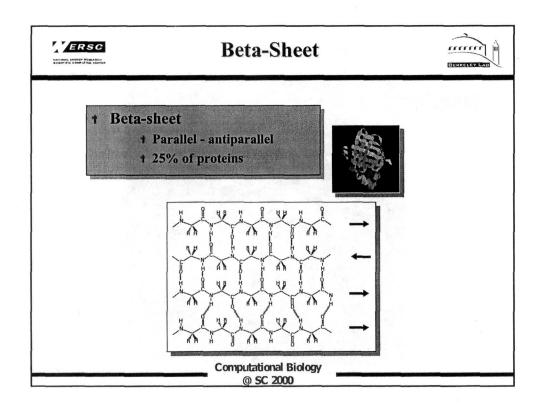


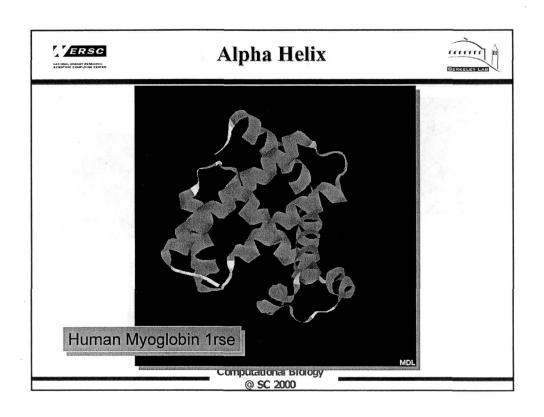
VERSC	Aro	mat	ic Ar	nino A	cids ——	BERKELEY LAG
	Amin o Aci d	pK,'s ²	Pro Stetue	Che mical Str ut uel	3-DStretue ¹	
	PhenylalminePhe,F Nocharge absorbs Uy hyskophoik (25) Molec. Wt. 447 Mole % 3-5	N=9.13 C=1.83 pI=5.48	a =1.16 B =1.33 t =0.59	H,N+ CO;		
	Tyr caine, Tyr, Y weak ahge a horb LIV by dogen b dongl no bydropille (0.08) Molec Wt. ~ 163 Mole % = 3.5	N=9.11 C=2.20 R=10.07 pI=5.66	a =0.74 B =1.45 t =0.76	H.W1		
	Typtop han, Tr ₄ AV largat main oacid rec stamin oacid nocharge a borb LIV hydogen bómgl hydophoùb (15)	N=9.39 C=2.38 p1=5.89	a =1.02 B =1.35 t =0.65	H-9/1° -1 -1 -1 -1 -1 -1 -1 -1 -1 -1 -1 -1 -1		
	Mole c Wt. = 18 6 Mole % = 1.1			• >		
			Charles S. Gasser 19 nputationa @ SC 20	al Biology		

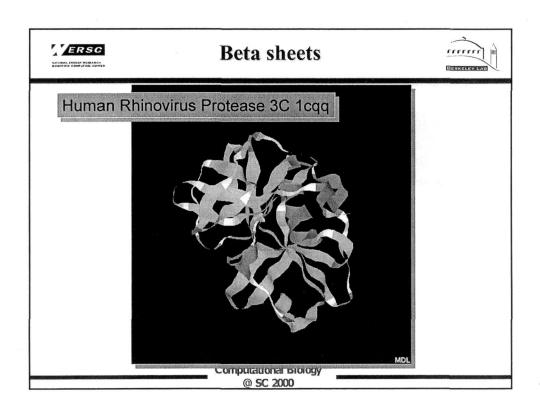














SCOP: Structural Classification of Proteins



- 1. All alpha proteins (a)
- 2. All beta proteins (b)
- 3. Alpha and beta proteins (a/b)
 - † Mainly parallel beta sheets (beta-alpha-beta units)
- 4. Alpha and beta proteins (a+b)
 - † Mainly antiparallel beta sheets (segregated alpha and beta regions)
- 5. Multi-domain proteins (alpha and beta)
 - t Folds consisting of two or more domains belonging to different classes
- 6. Membrane and cell surface proteins and peptides
 - t Does not include proteins in the immune system
- 7. Small proteins
 - t Usually dominated by metal ligand, heme, and/or disulfide bridges
- 8. Coiled coil proteins
- 9. Low resolution protein structures
- 10. Peptides
- 11. Designed proteins

Computational Biology @ SC 2000



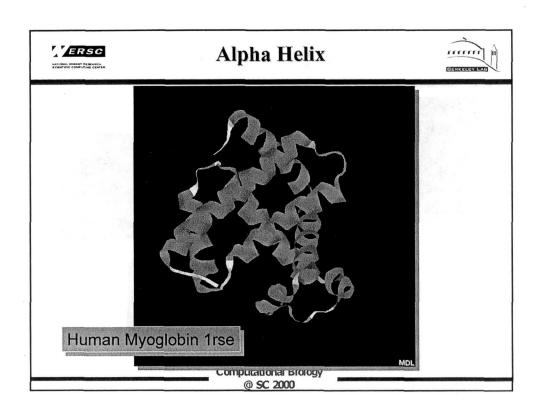
SCOP Classifications

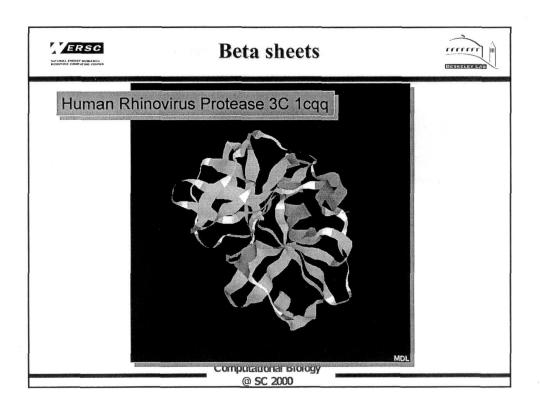


Class	Number of fold	N with er of sup erfam lies	N mmb er of fami keis	
All lapha rpteins	128	197	296	
All bla pro tins	87	158	251	
Alphand eta pro eins (a/)	93	153	323	
Alp haan deta proteins (a+b)	168	237	345	
Mu ti-d o mai n pro tei ns	25	25	32	
Membareandedl surface poteins	11	17	19	
Sm all pote in	52	72	102	
T dal	564	859	1368	

SCOP: Structural Classification of Proteins, 1.53 release 11410 PDB Entries (1 Jul 2000). 26219 Domains.

Copyright © 1994-2000 The scop authors / scop@mrc-lmb.cam.ac.uk September 2000







SCOP: Structural Classification of Proteins



- 1. All alpha proteins (a)
- 2. All beta proteins (b)
- 3. Alpha and beta proteins (a/b)
 - † Mainly parallel beta sheets (beta-alpha-beta units)
- 4. Alpha and beta proteins (a+b)
 - † Mainly antiparallel beta sheets (segregated alpha and beta regions)
- 5. Multi-domain proteins (alpha and beta)
 - † Folds consisting of two or more domains belonging to different classes
- 6. Membrane and cell surface proteins and peptides
 - t Does not include proteins in the immune system
- 7. Small proteins
 - t Usually dominated by metal ligand, heme, and/or disulfide bridges
- 8. Coiled coil proteins
- 9. Low resolution protein structures
- 10. Peptides
- 11. Designed proteins

Computational Biology @ SC 2000



SCOP Classifications



Cl as	Number of fold	N umber of superfamilies	N who er of fami kis 296 251	
All lp ha rptein s	128	197		
All bta pro ei rs	87	158		
Alphaand eta proteins (a/)	93	153	323	
Alp h aan d etta pro eins (a+b)	168	237	345 32	
Mu ti-d o mai n pro tei ns	25	25		
Membareandedl surface poteins	11	17	19	
Sm all pote in	52	72	102	
T dal	564	859	1368	

SCOP: Structural Classification of Proteins, 1.53 release 11410 PDB Entries (1 Jul 2000).

26219 Domains.
Copyright © 1994-2000 The scop authors / scop@mre-lmb.cam.ac.uk
September 2000



Protein Fold Recognition, Structure Prediction, and Folding



- † Drawing analogies with known protein structures
 - † Sequence homology, Structural Homology
 - † Inverse Folding, Threading
- † Ab initio folding: the ability to follow kinetics, mechanism
 - † robust objective function
 - † severe time-scale problem
 - * proper treatment of long-ranged interactions
- † Ab initio prediction: the ability to extrapolate to unknown folds
 - * multiple minima problem
 - t robust objective function
 - **† Stochastic Perturbation and Soft Constraints**
- † Simplified Models that Capture the Essence of Real Proteins
 - † Lattice and Off-Lattice Simulations
 - † Off-Lattice Model that Connect to Experiments: Whole Genomes?

Computational Biology



Protein Fold Predictions: Neural Network Structure Classifications



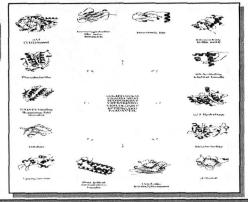
- **† Protein fold predictor based on global descriptors of amino acid sequence**
- * Empirical prediction using a database of known folds in machine learning
- * Databases
 - † 3D-ALI (83 folds)
 - * SCOP (used ~120 folds)
- † Representation of protein sequence in terms of physical, chemical, and structural properties of amino acids
- * Feed forward neural network for machine learning



Protein Fold Recognition: Threading



Sequence Assignments to Protein Fold Topology (David Eisenberg, UCLA)



Take a sequence with unknown structure and align onto structural template of a given fold Score how compatible that sequence is based on empirical knowledge of protein structure Right now 25-30% of new sequences can be assigned with high confidence to fold class

100,000's of sequences and 10,000's of structures (each of order 102-103 amino acids long)

Computational Biology @ SC 2000



Protein Fold Recognition: Threading



Computational Approach:

Dynamic programming: capable of finding optimal alignments if optimal alignments of subsequences can be extended to optimal alignments of whole objective functions that are one-dimensional $E=\Sigma\ V_i+\Sigma\ V_{gap}$

Complexity: all to all comparison of sequence to structure scales as ${\bf L}^2$ Whole human genome: ${\bf 10^{13}}$ flops

Improve Objective function:

Take into account structural environment

3D→1D: dynamic programming, L²

Build pairwise or multi-body objective function

NP-hard if: variable-length gaps and model nonlocal effects such as distance dependence

Recursive dynamic programming, Hidden markov models, stochastic grammers

Complexity: all to all comparison of sequence to structure scales as ${\rm L}^3$ Whole human genome: ${\sim}10^{16}$ flops

> Computational Biology @ SC 2000

> > 13



Computational Protein Folding



One microsecond simulation of a fragment of the protein. Villin. (Duan & Kollman, Science 1998)







- ✓ robust objective function
 all atom simulation with molecular water present: some structure present
- ✓ severe time-scale problem required 10° energy and force evaluations: parallelization (spatial decomposition)
- proper treatment of long-ranged interactions

 cut-off interactions at 8Å, poor by known simulation standards
- Statistics (1 trajectory is anecdotal)
- Many trajectories required to characterize kinetics and thermodynamics

Computational Biology



Computational Protein Folding



(1) Size-scaling bottlenecks: Depends on complexity of energy function, V

Empirical (less accurate): cN^2 ; ab initio (more accurate): CN^3 or worse; c << C empirical force field used

"long-ranged interactions" truncated so cM^2 scaling; $M \leq N$ spatial decomposition, linked lists

(2) Time-Scale of motions bottlenecks (Δt)

$$r_{i}(t + \Delta t) = 2r_{i}(t) - r_{i}(t - \Delta t) + \frac{f_{i}(t)(\Delta t)^{2}}{m_{i}} + O[(\Delta t)^{4}]v_{i}(t) = \frac{r_{i}(t + \Delta t) - r_{i}(t - \Delta t)}{2\Delta t} + O[(\Delta t)^{3}]$$

$$f_{i} = m_{i}a_{i} = -\nabla_{i}V(r_{1}, r_{2}, ..., r_{N})$$

Use timestep commensurate with fastest timescale in your system

bond vibrations: 0.01Å amplitude: 10-15 seconds (1fs)

Shake/Rattle bonds (2fs)

Multiple timescale algorithms (~5fs) (not used here)

@ SC 2000



Ab Initio Protein Structure Prediction



Primary Squence and an Energy function → Tertiary structure

Empirical energy functions:

(1) Detailed, Atomic description: leads to enormous difficulties!

$$V_{MM} = \sum_{i}^{\# Bonds} k_{b} (b_{i} - b_{o})^{2} + \sum_{i}^{\# Angles} k_{\theta} (\theta_{i} - \theta_{o})^{2} + \sum_{i}^{\# Impropers} k_{\tau} (\tau_{i} - \tau_{o})^{2} + \sum_{i}^{\# Impropers} k_{\tau}$$

(1) Multiple minima problem is fierce

Find a way to effectively overcome the multiple minima problem

(2) Objective Functions: Replaceable algorithmic component?

Global energy minimum should be native structure, misfolds higher in energy

Computational Biology

@ SC 2000



The Objective (Energy) Function



Empirical Protein Force Fields: AMBER, CHARMM, ECEPP "gas phase"





CATH protein classification: http://pdb.pdb.bnl.gov/bsm/cath

 α -helical sequence/ β -sheet structure

β-sheet sequence/a-helical structure

Energies the same! Makes energy minimization difficult!

Add penalty for exposing hydrophobic surface: favors more compact structures

 $E_{native \ folds} < E_{misfolds}$ for a few test cases

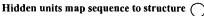
Solvent accessible surface area functions: Numerically difficult to use in optimization

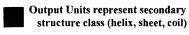


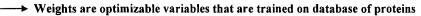
Neural Networks for 2° **Structure Prediction**



-) Input units represent amino acid sequence







Poorly designed networks result in overfitting, inadequate generalization to test set

Neural network design

input and output representation number of hidden neurons

weight connection patterns that detect structural features

Computational Biology @ SC 2000



Neural Network Results



No sequence homology through multiple alignments

Train

Test

Total predicted correctly = 66%

Total predicted correctly = 62.5%

Helix: 51% C_a=0.42

Helix: 48% C_a=0.38

Sheet: 38% $C_b = 0.39$

Sheet: 28% $C_b = 0.31$

Coil: 82% $C_c = 0.36$

Coil: 84% $C_c = 0.35$

Network with Design: Yu and Head-Gordon, Phys. Rev. E 1995

Train

Test

Total predicted correctly = 67%

Total predicted correctly = 66.5%

Helix: 66% C_a=0.52

Helix: 64% C_a=0.48

Sheet: 63% $C_b = 0.46$

Sheet: 53% $C_b = 0.43$

Coil: 69% C_c=0.43

Coil: 73% C_c =0.44

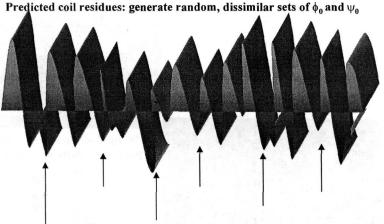
Combine networks of Yu and Head-Gordon with multiple alignments



Neural Networks Used To Guide Global Optimization Methods



Generate expanded tree of configurations



Explore tree configuration in depth:

Global Optimization in sub-space of coil residues: walk through barriers, move downhill

Computational Biology @ SC 2000

ERSC

Hierarchical Parallel Implementation of Global Optimization Algorithm



Static vs. Dynamic Load Balancing of Tasks

Central Processor

Central Processor: Assigns starting coordinates to GOPT's

Task time is highly variable

GOPT's: Divide up sub-space into N regions for global search

Task time is variable

Workers: Generate sample points; find best minimizer in region (Number of workers depends on sub-space)

Dynamical load balancing of tasks: reassigning GOPT/workers to GOPT/workers

Gain in efficiency of a factor of 5-10 Computational Biology

@ SC 2000



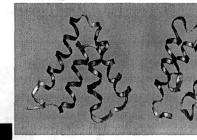
Global Optimization Predictions of α -Helical Proteins



Crystal (left), Prediction (right) R.M.S. 7.0Å



1pou: 72 aa DNA binding protein



Prediction (left) and crystal (right) R.M.S. 6.3Å

2utg_A: 70aa α -chain of uteroglobin:



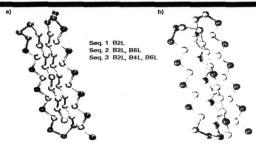
Still have not reached crystal energy yet!

Computational Biology @ SC 2000



Simplified Models for Simulating Protein Folding





Simplifies the "real" energy surface topology sufficiently that you can do

(1) Statistics ✓

Can do many trajectories to converge kinetics and thermodynamics

(2) severe time-scale problem√

characterize full folding pathway: mechanism, kinetics, thermodynamics

(3) proper treatment of long-ranged interactions ✓

all interactions are evaluated; no explicit electrostatics

(4) robust objective function?

good comparison to experiments

Computational Biology

@ SC 2000



Acknowledgements



Teresa Head-Gordon, Physical Biosciences Division, LBNL

Silvia Crivelli, Physical Biosciences and NERSC Divisions, LBNL

Betty Eskow, Richard Byrd, Bobby Schnabel, Dept. Computer Science, U. Colorado

Jon M. Sorenson, NSF Graduate Fellow, Dept. Chemistry UCB

Greg Hura, Graduate Group in Biophysics, UCB

Alan K. Soper, Rutherford Appleton Laboratory, UK

Alexander Pertsemlidis, Dept. of Biochemistry, U. Texas Southwestern Medical

Robert M. Glaeser, Mol. & Cell Biology, UCB and Life Sciences Division, LBNL

Funding Sources:

AFOSR, DOE (MICS), DOE/LDRD (LBNL), NIH, NERSC for cycles

Computational Biology @ SC 2000



Structure-Based Drug Discovery

Brian K. Shoichet, Ph.D Northwestern University, Dept of MPBC 303 E. Chicago Ave, Chicago, IL 60611-3008 Nov 15, 1999



Problems in Structure-Based Inhibitor Discovery & Design



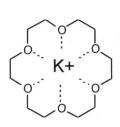
- † Balance of forces in binding
 - * Energies in condensed phases
 - † interaction energies
 - † desolvation
- † Problem scales badly with degrees of freedom
 - **†** Configuration
 - † configs α (prot-features)⁴ X (lig-features)⁴
 - * Conformation
 - * Ligand & Protein, confs a 3lbonds X 3pbonds
- † Sampling chemical space (scales very badly)
- † Defining binding sites

Computational Biology

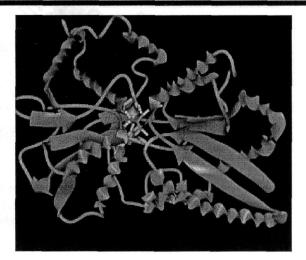


The Pros & Cons of Proteins

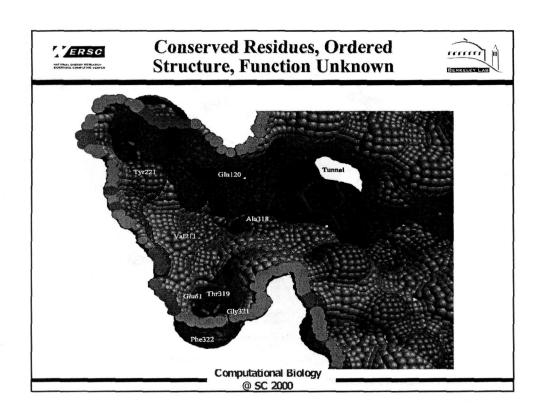




18 - Crown-6



sulfate binding protein





Inhibitor Discovery or Design?

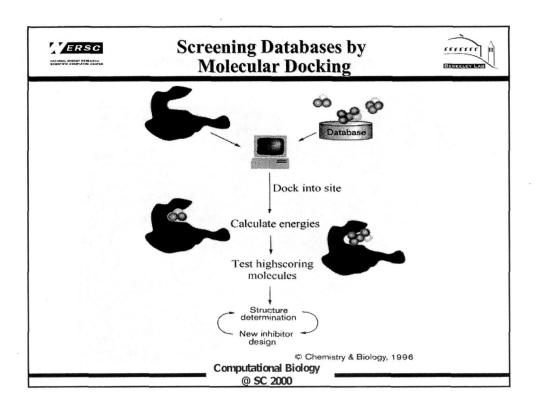


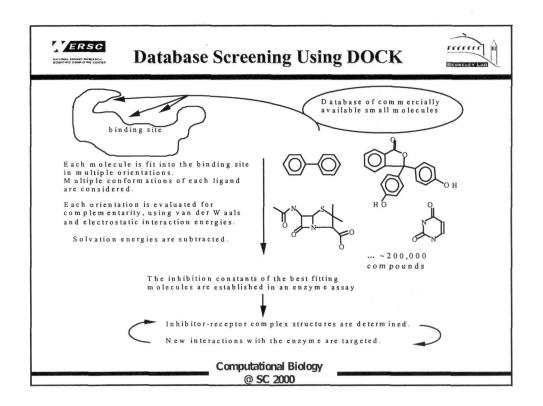
t Design ligands

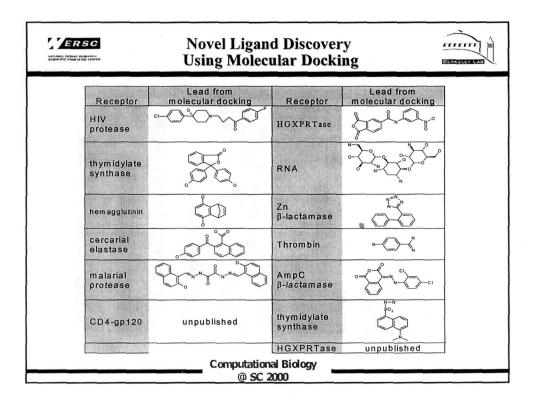
- † Ludi (Bohm)
- * Grow (Moon & Howe)
- **†** Builder (Roe & Kuntz)
- MCSS-Hook (Miranker & Karplus)
- * SMOG (DeWitte & Shaknovitch)
- + Others...

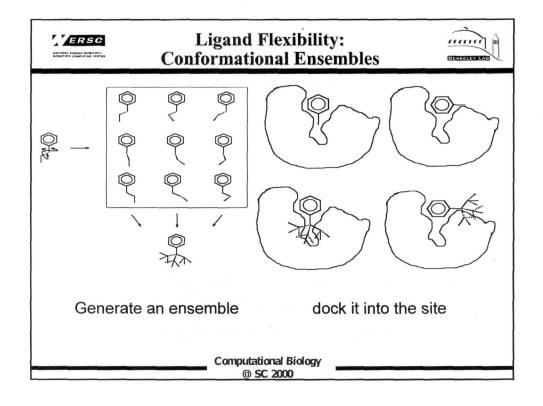
t Discover Ligands

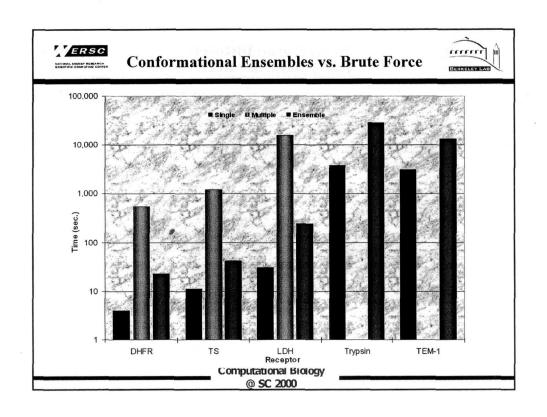
- + DOCK (Kuntz, et al., Shoichet)
- * CAVEAT (Bartlett)
- * Monte Carlo (Hart & Read)
- + AutoDock (Goodsell & Olson)
- * SPECITOPE (Kuhn et al)
- + Others...

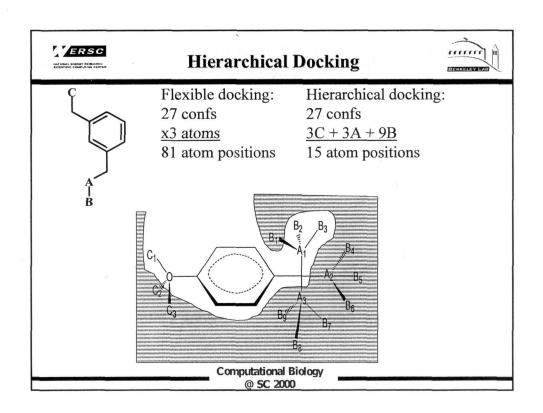














Computational Phylogenetics

Craig Stewart stewart@iu.edu Indiana University



Outline



- Evolution & Phylogenetics
- Alignment (brief)
- Why is phylogeny construction a HPC problem?
- Summary of methods and software for phylogenetics
- One example in detail: Maximum Likelihood analysis with fastDNAml
- Some interesting results and challenges for the future
- Caveat: this is an introduction, not an exhaustive review.



Why



- Curiosity: Anyone who as a child wandered through the dinosaur section of a natural history museum understands the inherent intellectual attraction of evolutionary biology
- Theoretical uses: testing hypotheses in evolutionary biology
- **■** Practical uses:
 - Medicine
 - Environmental management (biodiversity maintenance)

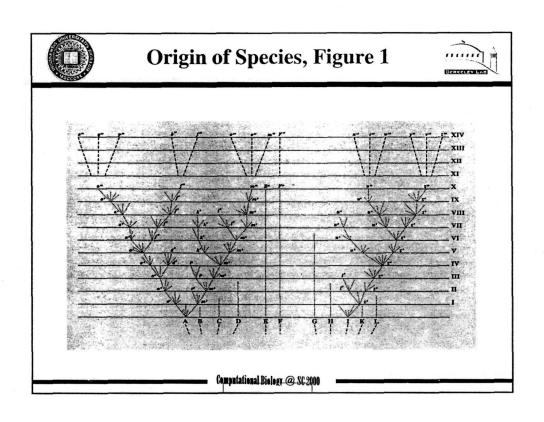
Computational Biology @ SC 2000

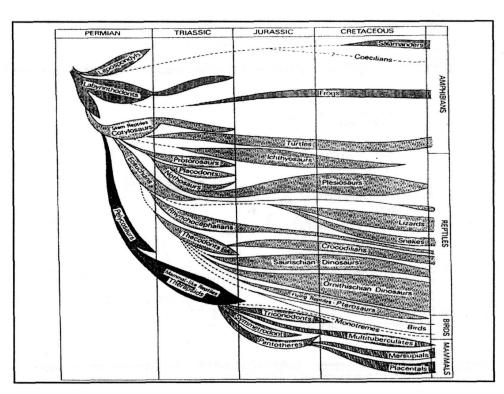


Phylogeny



- Evolution is an explicitly historical branch of biology, one in which the subjects are active players in the historical changes.
- A phylogeny, or phylogenetic tree, is a way of depicting evolutionary relationships among organisms, genes, or gene products.
- Modern evolutionary biology began with the publication of Darwin's *Origin of Species*, which included one figure a phylogenetic tree.



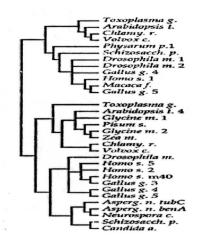




Building Phylogenetic Trees



- Goal: an objective means by which phylogenetic trees can be estimated in tolerable amounts of wallclock time, producing phylogenetic trees with measures of their uncertainty
- Closely related taxa (or genes) are grouped closely together. Lengths of tree branches correspond to amounts of genetic difference



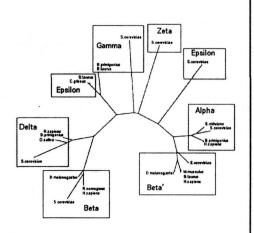
Computational Biology @ SC 2000



Basic Evolutionary Biology



- All evolutionary changes are described as bifurcating trees
 - evolutionary relationships among genes or gene products (trees of paralogues)
 - evolutionary relationships among organisms (trees of orthologues)
- **■** Basic rationale of phylogenetics
 - Groups of related organisms (or genes) share characters
 - Species (or genes) that share lots of characteristics are grouped closely together
 - Species (or genes) that share few characteristics are grouped far apart





Reconstructing history from DNA sequences



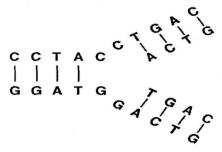
- DNA changes over time; much of this change is not expressed
- Changes in unexpressed DNA can be modeled as random process
- By comparing similar regions of DNA from different organisms (or different genes) one can infer the phylogenetic tree and evolutionary history that seems the best explanation of the current situation
- The process of creating phylogenies from DNA sequences is much like mapping relationships among different versions of the Bible one tracks transcription errors
- DNA transcription errors are sometimes corrected through a subsequent transcription error

Computational Biology @ SC 2006



DNA replication





Purines: Pyrimidines:

Adenine & Guanine Thymine & Cytosine



Changes in genetic information over time



■ Point mutations

DNA - sequences of the 4 nucleotides

CCTCTGAC

VS

TCTC**C**GAC

Protein - sequences of the 20 amino acids

GSAQVKGHGKK

VS

G**NPK**VK**A**HGKK

■ Insertions and deletions

DNA

CCTCT+GAC

VS

CCTCTTGAC

Computational Biology @ SC 2000



Alignment



- To build trees one compares and relates 'similar' segments of genetic data. Getting 'similar' right is absolutely critical!
- Methods:
 - dynamic programming
 - Hidden Markov Models
 - Pattern matching
- Some alignment packages:
 - BLAST http://www.ncbi.nlm.nih.gov/BLAST/
 - FASTA http://gcg.nhri.org.tw/fasta.html
 - MUSCA http://www.research.ibm.com/bioinformatics/home



Matching cost function



GCTAAATTC

++ x x

GC AAGTT

- Penalize for mismatches, for opening of gap, and for gap length
- This approach assumes independence of loci: good assumption for DNA, some problems with respect to amino acids, significant problems with RNA (RNA sequence alignment is a much more complicated matter)

Computational Biology @ SC 2000



Example of aligned sequences



Thermotoga	ATTTGCCCCA	${\tt GAAATTAAAG}$	CAAAAACCCC	${\tt AGTAAGTTGG}$	${\tt GGATGGCAAA}$
Tthermophi	${\bf ATTTGCCCCA}$	${\tt GGGGTTCCCG}$	CAAAAACCCC	${\bf AGTAAGTTGG}$	${\tt GGATGGCAGG}$
Taquaticus	ATTTGCCCCA	${\tt GGGGTTCCCG}$	CAAAAACCCC	AGTAAGTTGG	${\tt GGATGGCAGG}$
deinon	ATTTGCCCCA	GGGATTCCCG	CAAAAACCCC	AGTAAGTTGG	${\tt GGATGGCAGG}$
Chlamydi	ATTTTCCCCA	GAAATTCCCG	AAAAAACCCC	AATAAATTGG	${\tt GGATGGCAGG}$
flexistips	ATTTTCCCCA	CAAAAAAAAG	AAAAAACCCC	AGTAAGTTGG	GGATGGCAGG
borrelia-b	ATTTGCCCCA	GAAGTTAAAG	CAAAAACCCC	AATAAGTTGG	GGATGGCAGG
bacteroide	ATTTGCCCCA	GAAATTCCCG	CAAAAACCCC	AGTAAATTGG	GGATGGCAGG
Pseudom	ATTTGCCCCA	GGGATTCCCG	CAAAAACCCC	AGTAAGTTGG	GGATGGCAGG
ecoli	GTTTTCCCCA	GAAATTCCCG	CAAAAACCCC	AGTAAGTTGG	GGATGGCAGG
salmonella	++++++++	+++++++++		+++++++++	+++++++++
shewanella	GTTTGCCCCA	GCCATTCCCG	TAAAAACCCC	AGTAAGTTGG	GGATGGCAGG
bacillus	ATTTGCCCCA	GAAATTCCCG	CAAAAACCCC	AGCAAATTGG	GGATGGCAGG
myco-gent1	ATTTGCCCCG	GAAATTCCCG	CAAAAACCCC	AGTAAGTTGG	GGATGGCAAA



Sequences available



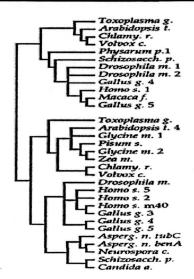
- DNA (sequences are series of the base molecules; aligned sequences will also contain +s for gaps)
- Amino acid sequences (series of letters indicating the 20 amino acids). Computational challenges more severe than with DNA sequences.
- RNA
- The availability of data at present exceeds the ability of researchers to analyze it!

Computational Biology @ SC 2000



Why is tree-building a HPC problem?





- The number of bifurcating unrooted trees for n taxa is (2n-5)!/{(n-3)! (2n-3)}
- for 50 taxa the number of possible trees is ~10⁷⁴; most scientists are interested in much larger problems
- The number of rooted trees is (2n-5)!



Phylogenetic methodologies



- Define a specific series of steps to produce the 'best' tree
 - Pair-group cluster analyses
 - Fast, but tend not to address underlying evolutionary mechanisms
- Define criteria for comparing different trees and judging which is better. Two steps:
 - Define the objective function (evolutionary biology)
 - Generate and compare trees (computation)
- All of the techniques described produce an unrooted tree.
- The trees produced likewise describe relationships among extant taxa, not the progress of evolution over time.
- Two computational approaches:
 - Distance-based methods
 - Character-based methods

Computational Biology @ SC 2000



Distance-based Tree-building methods



- Aligned sequences are compared, and analysis is based on the differences between sequences, rather than the original sequence data.
- Less computationally intensive than character-based methods
- Tend to be problematic when sequences are highly divergent



Distance-based Tree building methods, 2



- Cluster analysis Most common variant is Unweighted Pair Group Method with Arithmetic Mean (UPGMA) – join two closest neighbors, average pair, keep going. Problematic when highly diverged sequences are involved
- Additive tree methods Built on assumption that the lengths of branches can be summed to create some measure of overall evolution.
 - Fitch-Margoliash (FM) minimizes squared deviation between observed data and inferred tree.
 - Minimum evolution (ME) finds shortest tree consistent with data
- Of the distance methods, ME is the most widely implemented in computer programs

Computational Biology @ SC 2000



Character-based methods



- Use character data (actual sequences) rather than distance data
- Maximum parsimony. Creates shortest tree one with fewest changes. Inter-site rate heterogeneity creates difficulties for this approach.
- Maximum likelihood. Searches for the evolutionary model that has the highest likelihood value given the data. In simulation studies ML tends to outperform others, but is also computationally intensive.



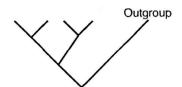
Rooting trees



■ If the assumption of a constant molecular clock holds, then the root is the midpoint of the longest span across the tree.



■ May be handled by including an 'outgroup' in the analysis



Computational Biology @ SC 2000



Evaluating trees



- Once a phylogenetic tree has been produced by some means, how do you test whether or not the tree represents evolutionary change, or just the results of a mathematical technique applied to a set of random data? These methods below can be used to perform a statistical significance test.
- Significance tests for MP trees:
 - Skewness tests. MP tree lengths produced from random data should be symmetric; tree lengths produced from data sets with real signal should be skewed.
- Significance tests for distance, MP, and ML trees:
 - Bootstrap. Recalculate trees using multiple samples from same data with resampling.
 - Jackknife. Recalculate trees using subsampling
- All of these methods are topics of active debate



Phylogenetic software



- Phylip. (J. Felsenstein). Collection of software packages that cover most types of analysis. One of the most popular software collections. Free.
- PAUP. (D. Swofford). Parsimony, distance, and ML methods. Also one of the most popular software collections. Not free, but not expensive.
- PAML. (Ziheng Yang). Maximum likelihood methods for DNA and proteins. Not as well suited for tree searching, but performs several analyses not generally available. Free.
- fastDNAml. (G. Olsen). Maximum likelihood method for DNA; becoming one of the more popular ML packages. MPI version available soon; well suited to tree searching in large data sets. Free.

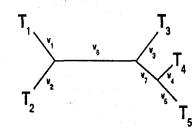
Computational Biology @ SC 2000



More on Maximum Likelihood methods



- Typical statistical inference: calculate probability of data given the hypothesis
- Tree, branch lengths, and associated likelihood values all calculated from the data.
- Likelihood values used to compare trees and determine which is best





Stochastic change of DNA



 Markov process, independent for each site: 4 x 4 matrix for DNA, 20 x 20 for amino acids

	A	C	\mathbf{G}	T
A	p(A->A)	p(A->C)	p(A->G)	p(A->T)
C	p(C->A)	p(C->C)	p(C->G)	p(C->T)
G	p(G->A)	p(G->C)	p(G->G)	p(G->T)
T	p(T->A)	p(T->C)	p(T->G)	p(T->T)

- **■** Transitions more probable than transversions.
- Must account for heterogeneity in substitution rates among sites (DNArates Olsen)

Computational Biology @ SC 2000



fastDNAml



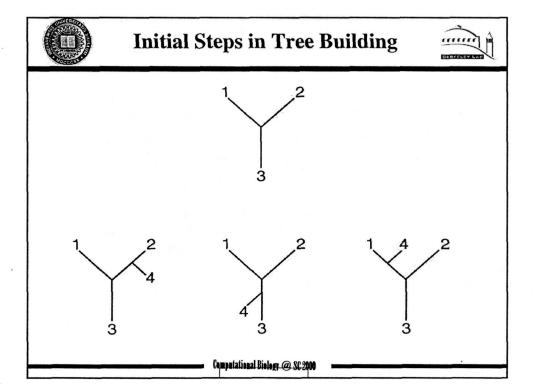
- **Developed by Gary Olsen**
- **■** Derived from Felsensteins's PHYLIP programs
- One of the more commonly used ML methods
- The first phylogenetic software implemented in a parallel program (at Argonne National Laboratory, using P4 libraries)
- Olsen, G.J., et al.1994. fastDNAml: a tool for construction of phylogenetic trees of DNA sequences using maximum likelihood. Computer Applications in Biosciences 10: 41-48
- MPI version produced in collaboration with Indiana University will be available soon

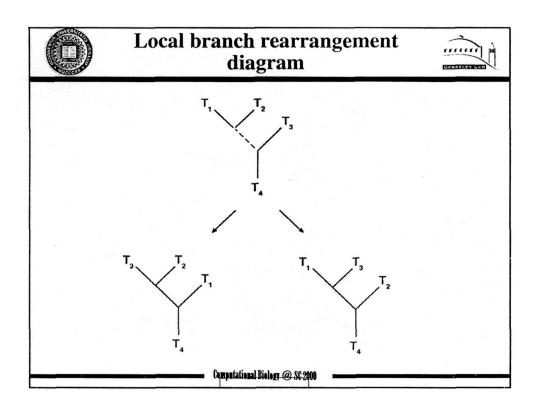


fastDNAml algorithm



- Compute the optimal tree for three taxa (chosen randomly) only one topology possible
- Randomly pick another taxon, and consider each of the 2i-5 trees possible by adding this taxon into the first, three-taxa tree.
- Keep the best (maximum likelihood tree)
- Local branch rearrangement: move any subtree to a neighboring branch (2i-6 possibilities)
- **■** Keep best resulting tree
- Repeat this step until local swapping no longer improves likelihood value







fastDNAml algorithm con't: Iterate

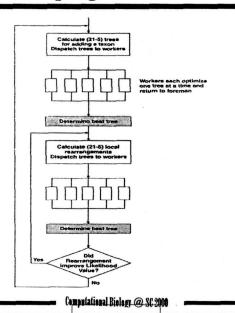


- Get sequence data for next taxon
- Add new taxa (2i-5)
- Keep best
- Local rearrangements (2i-6)
- Keep best
- Keep going....
- When all taxa have been added, perform a full tree check



Overview of parallel program flow



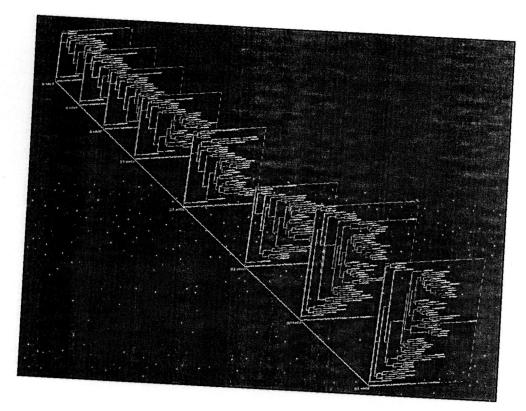


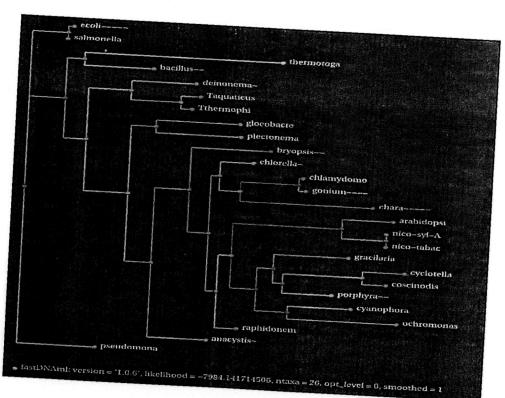


Because of local effects....



- Where you end up sometimes depends on where you start
- This process searches a huge space of possible trees, and is thus dependent upon the randomly selected initial taxa
- Can get stuck in local optimum, rather than global
- Must do multiple runs with different randomizations of taxon entry order, and compare the results
- Similar trees and likelihood values provide some confidence, but still the space of all possible trees has not been searched extensively







Grid computing with fastDNAml



- The high computation/communication ratio makes this program a good candidate for geographic distribution
- Time to completion is a constant forever and ever
- The key task is to combine geographically distributed resources so that large jobs can be completed in tolerable (for the biologist) amounts of wall clock time
- Handles timeouts, system crashes, etc.

	Foreman Task #		
	,		
		_	
		-	
_		_	
	3		
	2		
	1	-	

Sla	Slaves			
Machine #	Start Time			
n				
3				
2				
1				
Ready	queue			

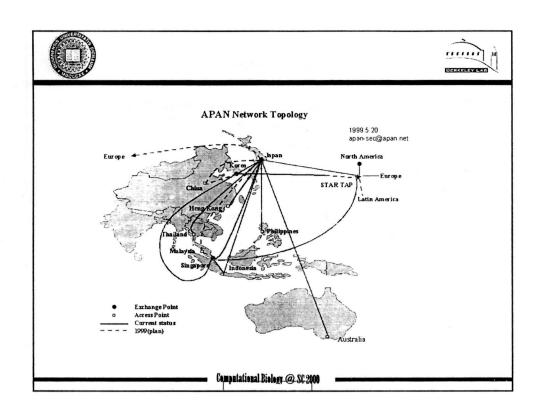
Computational Biology @ SC 2000

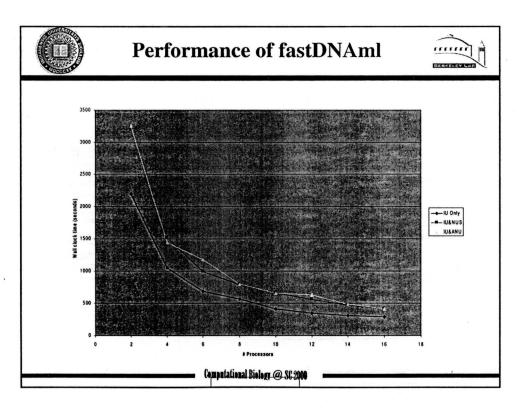


Demonstration at SC98



- Indiana University SP nodes
- NUS SP nodes
- ACSys DEC Workstations
- Immersadesk on the SC98 show floor as part of the IU/EVL iGRID demonstration



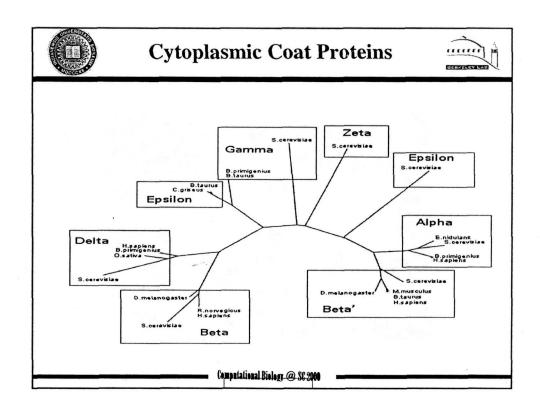




Applications & Interesting examples



- Better understanding of evolution (Ceolocanths, cyanobacterial origin of plastids)
- **■** Maintenance of biodiversity
- Medicine & molecular biology
 - our cousins, the fungi
 - Prediction of influenza vaccines
 - Cytoplasmic coat proteins
 - HIV





HIV



- Where did HIV come from, and how recent is it?
- Korber, et al. 2000. Timing the ancestor of the HIV-1 pandemic strains. Science 288:1789. (Online at www.sciencemag.org/cgi/content/full/288/5472/1789)
- Used completed HIV sequences from 159 individuals with known sampling dates (including one from 1959)
- Used a general-reversible (REV) base substitution model, accounting for different site-specific rates of evolution and base frequencies biased in favor of adenosine.
 Used modified version of fastDNAml.
- Used SIV as an outgroup
- Last common ancestor of main group of HIV-1 was 1931 (95% confidence interval: 1915-1941). Supports hypothesis that HIV has been around for some time and simply took a while to be common enough to be noticed.

Computational Biology @ SC 2000



Challenges for future



- HPC implementations of more phylogenetic techniques
- Better treatment of insertions and deletions (indels)
- Algorithms for more thorough searching of treespaces in incremental tree building processes (keep best n trees and keep looking)
- Techniques for not shaking the whole tree (that is, adding a taxa to a tree in a fashion that acknowledges damping of effect as you travel away from altered part of tree)
- Use of high-throughput techniques



Acknowledgements



- The phylogeny depicted in slide 6 is taken from E. Colbert. 1965. The age of reptiles. W.W. Norton, NY, NY.
- Some of the tree diagrams were adapted from Olsen et al. 1994.
- Les Teach [IU] created all other graphics for this talk
- Donald Berry is responsible for IU's programming efforts related to fastDNAml. David Hart and Richard Repasky are also involved in IU's efforts related to bioinformatics and fastDNAml.
- IU's work on parallel versions of fastDNAml has been facilitated by Shared University Research grants from IBM, Inc.
- IU's work with fastDNAml would be impossible without our collaboration with Gary Olsen, U. of Illinois, the creator of this program.

Computational Biology @ SC 2000



References



- Felsenstein, J. 1981. Evolutionary trees from DNA sequences: a maximum likelihood approach. Journal of Molecular Evolution 17:368-376
- Baxevanis, A.D., and B.F.F. Ouellette. 1998. Bioinformatics: a practical guide to the analysis of genes and proteins. Wiley-Interscience, NY.
- Swofford, D.L., and G.J. Olsen. Phylogeny reconstruction. pp. 411-501 IN D.M. Nillis & C. Mority (eds). Molecular systematics. Sinauer Associates, Sunderland, MA
- Durbin, R. et al. 1998. Biological sequence analysis. Cambridge University Press, Cambridge, UK.
- www.ucmp.berkely.edu/subway/phylogen
- evolution.genetics.washington.edu/phylip/software
- http://www.indiana.edu/uits/~rac



urls for phylogenetic software



- Phylip evolution.genetics.washington.edu/phylip/software.html
- PAUP www.lms.si.edu/PAUP/index.html
- PAML abacus.gene.ucl.ac.uk/software/paml.html
- fastDNAml geta.life.uiuc.edu/~gary/
 - MPI version available soon from www.indiana.edu/~uits/rac/bioinfo



Afternoon Break



Specialized biological databases and their role in building models of regulation

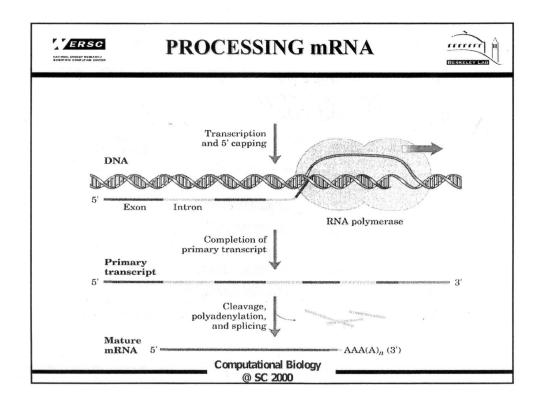
Inna Dubchak ILDubchak@lbl.gov NERSC



Overview of alternative splicing



- * What is alternative splicing?
- † What is possible to do computationally to better understand this complicated phenomenon?
 - * Frequency of alternative splicing
 - * Specialized databases
 - * Search for regulatory elements





The Nobel Prize in Physiology or Medicine 1993



The Nobel Assembly at the Karolinska Institute in Stockholm, Sweden, has awarded the Nobel Prize in Physiology or Medicine for 1993 jointly to Richard J. Roberts and Phillip A. Sharp for their discovery of split genes.





Computational Biology

@ SC 2000



a-Tropmyocin pre-mRNA



Alternative Splicing of a-tropomyocin pre-mRNA

Sense 5'	26-154 MG-188 M2-25	183.213 214.234	230-367 256-8	ALULAS BI 3 IST	100 k 15
Tissues	mRNA	Splicing			
Nonmuscle		P P BUMA			
Smooth muscle	St. Line	(enternant			
Striated muscle	200 D 315	松间间	148		
Striated muscle'	Mar El Six	of the following	14 15		
Hepatoma	S 24 2 5	667 836G	15		

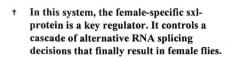
Computational Biology



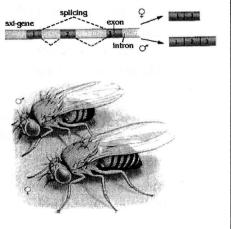
Gender in Drosophila



† A percursor-RNA may often be matured to mRNAs with alternative structures. An example where alternative splicing has a dramatic consequence is somatic sex determination in the fruit fly Drosophila melanogaster.



t Sex in Drosophila is largely determined by alternative splicing



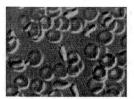
Computational Biology @ SC 2000



Splicing and diseases



- * Splicing errors cause thalassemia
- t Thalassemia, a form of anemia common in the Mediterranian countries, is caused by errors in the splicing process.



* Normal red blood cells contain correctly spliced beta-globin, an important component in hemoglobin that takes up oxygen in the lungs.





Information on alternative splicing in public databases:



- t Swiss-Prot (protein) database is well curated, but the information content is incomplete with reference to alternative splicing and does not allow for automatic retrieval of such entries.
- † Swiss-Prot entries just state the fact that a particular protein is one of the products of alternative splicing.
- † Some entries contain the information on the limited number of isoforms.

Computational Biology



Clustering procedure



Similarity analysis of two sequences

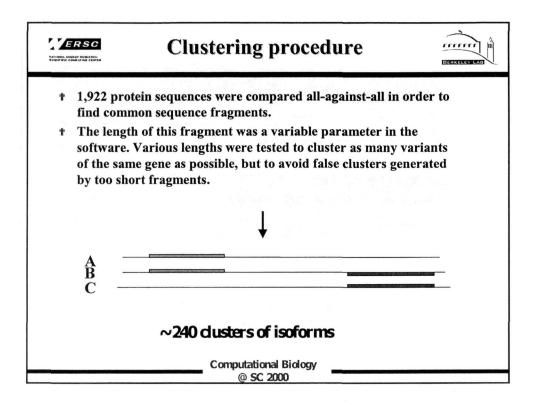
- Gene families
 multiple similar genes exist
 due too duplication and
 divergence of genes.
 - rgence of genes. than one way
- t Short similar fragments, a lot of mutations
- * Relatively long identical fragments

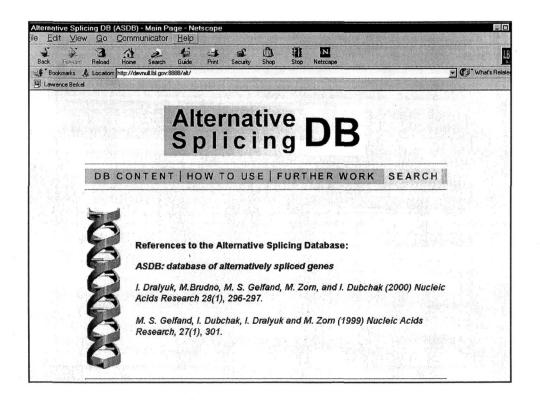
Alternative splicing

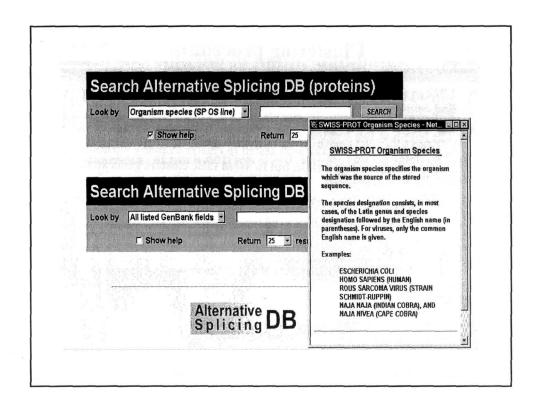
one gene but primary

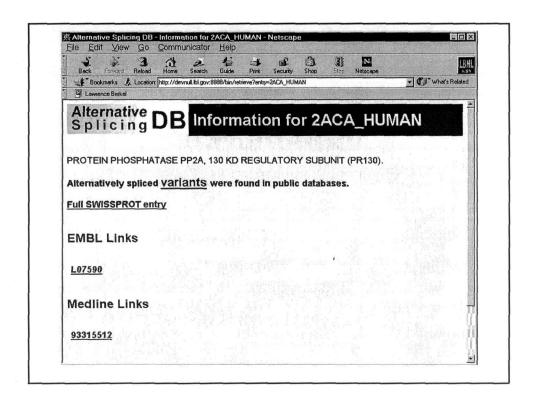
transcript spliced in more

Computational Biology

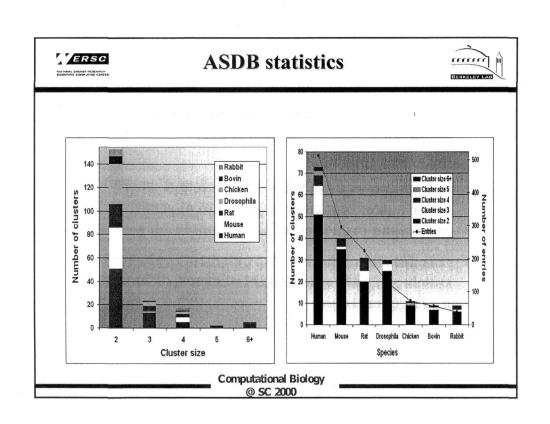








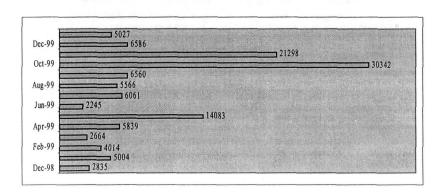
2 3	3 4 .	2 4		Shop Stop N	N letscape	LEK
many and the second second second	Location: http://devnull.li		min management of the second		· · · · · · · · · · · · · · · · · · ·	What's Related
Lawrence Berkel						
2ACA_HUMAN 2ACB HUMAN	TELOUDKENS	KKMDTVQSIP	NNSTNSLINL	EANDEKITKA	VQVQSQSLTM	- April 1 and 1 an
ZACB_HOHAN						
2ACA HUMAN	NPLENVSSDD	LMETLYIEEE	SDGKKALDKG	QKTENGPSHE	LLKVNEHRAE	
2ACB_HUMAN						
2ACA_HUMAN	FPEHATHLKK	CPTPMQNEIG	KIFEKSFVNL	PKEDCKSKVS	KFEEGDQRDF	
2ACB_HUMAN						
2ACA HUMAN	TNSSSORETD	KLLMDLESES	OKMETSLREP	LAKGKNSNFI	NSHSQLTGQT	
ZACB HUMAN						
2ACA_HUMAN	LVDLEPKSKV	SSPIEKVSPS	CLTRIIETNG	HKIEEEDRAL	LLRILESIED	
2ACB_HUMAN	*********				• • • • • • • • • • • • • • • • • • • •	
ZACA HUMAN	PAORTARCKS	SDCST.SOFKE	MMOTTORTT	TEROANT.SVC	RSPVGDKAKD	
2ACB HUMAN					FLARGCDEVL	
_						
2ACA_HUMAN	TTSAVLIQQT	PEVIKIONKP	EKKPGTPLPP	PATSPSSPRE	LSPVPHVNNV	
2ACB_HUMAN	PSRFKKRLKS	FQQTQIQNKP	EKKPGTPLPP	PATSPSSPRE	LSPVPHVNNV	
0101					BOULDTHONG	
2ACA_HUMAN 2ACB HUMAN					EQKADIYEMG EQKADIYEMG	3
ZACB_ROMAN	MARRITHIE	KETEPAGLED	TCSNREQIES	KIEIKENDIE	EQUADITERS	
2ACA HUMAN	KIAKVCGCPL	YWKAPMFRAA	GGEKTGFVTA	QSFIAMWRKI	LNNHHDDASK	
2ACB_HUMAN	KIAKVCGCPL	YWKAPMFRAA	GGEKTGFVTA	QSFIAMWRKI	LNNHHDDASK	
					FHSRYITTVI	





ASDB usage during 1999





Computational Biology

@ SC 2000



Study of Regulation



- **†** No systematic surveys to address the relative importance of such elements in the regulation of alternative splicing.
- t It is unknown as to whether regulatory words occur more frequently adjacent to alternative exons than in the rest of the genome.
- * It is not clear whether these elements enhance splicing of only a limited set of exons, or have a more general role.



Alternative Splicing Regulation



- **†** A number of genomic sequence regulatory elements have been identified outside of traditional splice sites.
- † The concept of splicing "enhancers" and "silencers" that promote or inhibit splicing at neighboring splice sites is well established.
- * Many alternative exons are probably regulated by a combination of silencers and enhancers.

Computational Biology



Data Collection



- * Automated processing of GenBank/Medline
- † Manual analysis of abstracts & articles
- † Collecting the sample

Computational Biology



BiSyCLES Search Options



- * BiSyCLES searches in the two databases, then establishes which of the retrieved entries are linked
 - * Medline: +"alternative splicing," tissue, muscle, brain, neuro*, heart, regul*, enhancer, silencer
 - * Genbank: +"alternative splicing" +"complete CDS"
- † Results:
 - * ~300 abstracts
 - * ~50 relevant papers

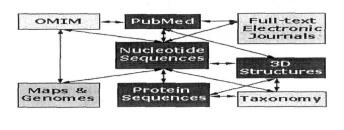
Computational Biology @ SC 2000



BiSyCLES: Biological System for Cross-Linked Entry Search



- * GenBank contains genomic data but little annotation
- * Medline (PubMed) contains abstracts from journals but no genomic data
- **†** NCBI's Entrez system keeps links between related entries in its databases





Word Counting



- * To calculate the confidence value of a particular word we select random subsets of a large dataset of constitutively spliced exons (1,504 exons; Burset & Guigo, 1996) equal in size to our alternative dataset.
- † We then calculate the fraction of these subsets in which the word is over-represented at a higher rate than in the alternative set.
- † (Over-representation is calculated as difference of frequencies)

Computational Biology @ SC 2000



Known Regulatory Elements



enhancers	<u>reference</u>			
UGCAUG	Huh & Hynes, 1994; Hedjran et al., 1997; Modafferi & Black, 1997; Kawamoto, 1996; Carlo et al., 1996			
CUG repeat	Ryan et al., 1996; Philips et al., 1998			
(A/U)GGG	Sirand-Pugnet et al., 1995a			
GGGGCUG	Carlo et al., 1996			
silencers				
UUCUCU	Chan & Black, 1995; Chan & Black, 1997; Ashiya & Grabowski, 199			

Computational Biology



Short summary



- † In the simple cases of splicing, introns are always introns and exons are always exons
- † During alternative splicing, within the same RNA, sequences can be recognized as either intron or exon under different conditions and the concept of exons and introns becomes rather empirical
- † RNAs are not spliced differently in the same cell at the same time but in different cells or in the same cell types at different times in development or under different conditions
- † A variety of patterns of alternate splicing have been observed.

Computational Biology

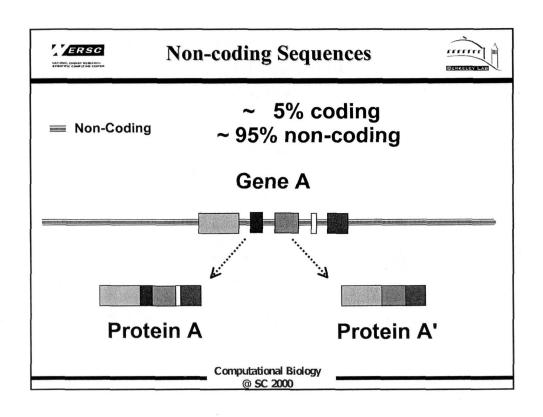


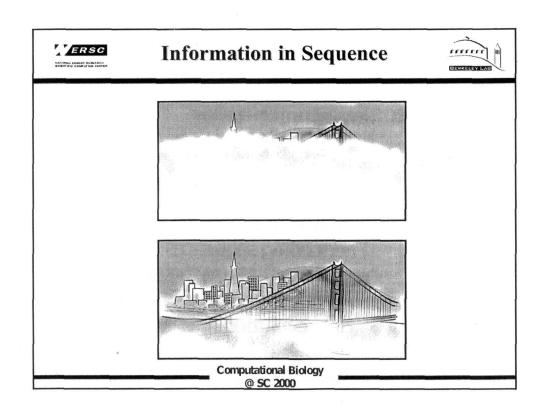
Evolutionarily conserved non-coding DNA sequences

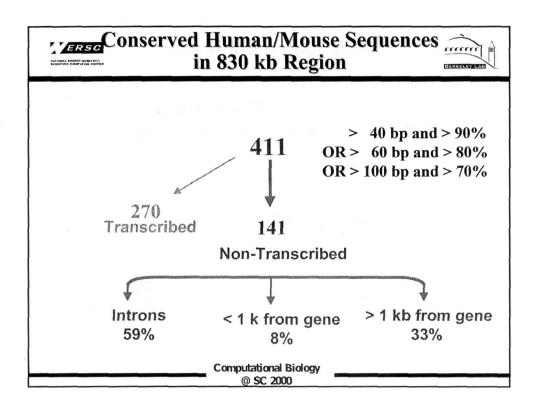


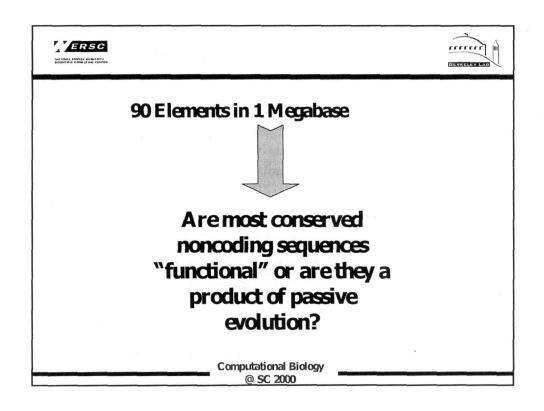
- * Discovering them in DNA sequence
- † Tools for their visualization
- * Biological importance

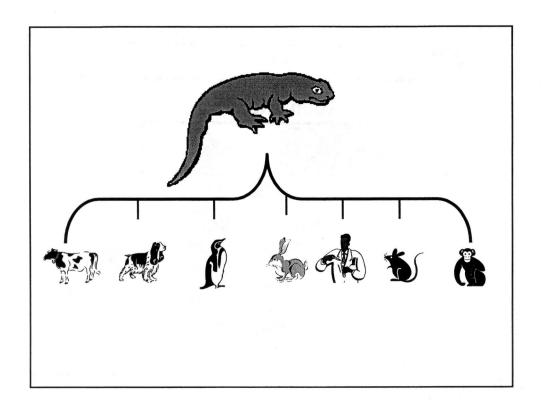
Computational Biology













Analysis of CNS-1



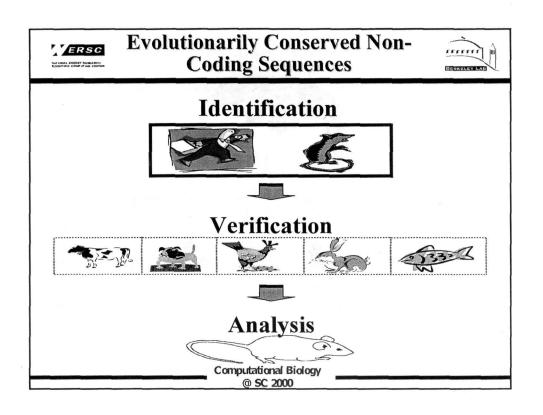
- * Present in other species:
 - † Cow (86%)
 - † Dog (81%)
 - * Rabbit (73%)
- * Genomic position conserved in human, mouse, dog and baboon

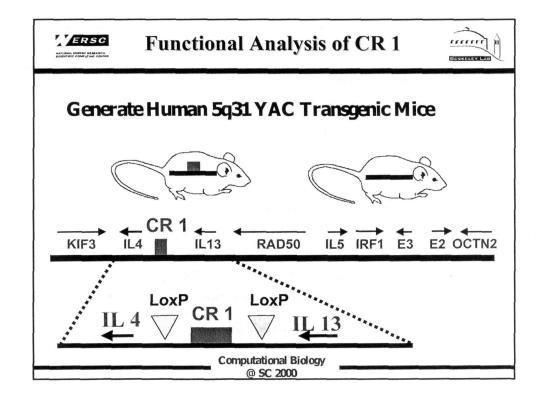
IL4

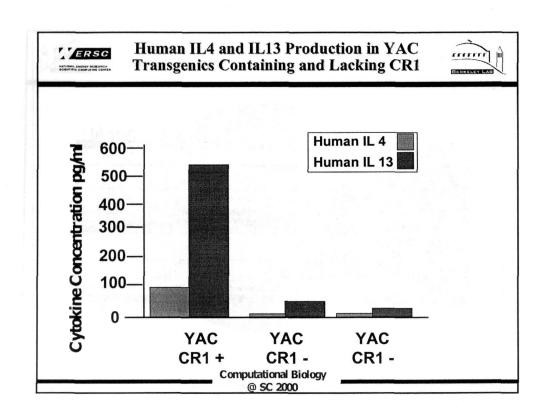
CNS-1

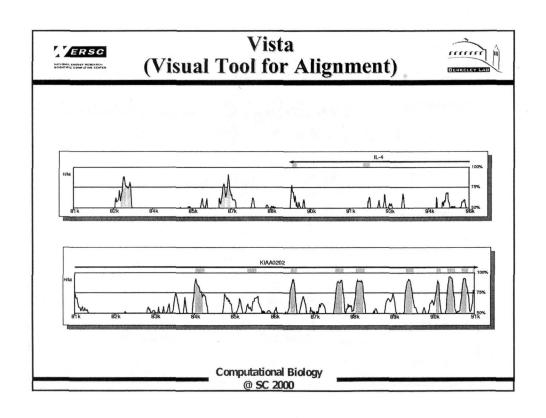
I<u>L</u> 13

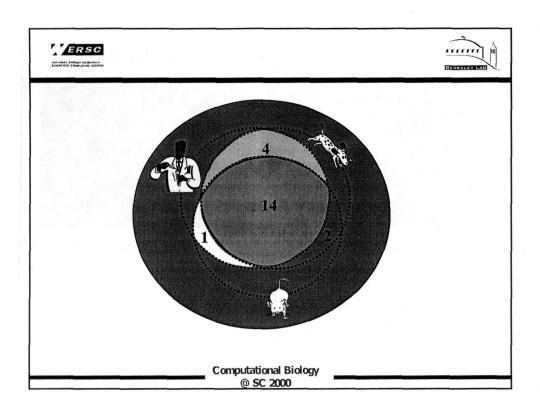
* Single copy in the human genome







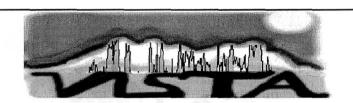






http://www-gsd.lbl.gov/vista/





Welcome to the VISTA, or VISualization Tool for Alignments home page

VISTA is an integrated system for global alignment and visualization, designed for comparative genomic analysis.

- 1. The visual output is clean and simple, allowing the user to easily identify conserved regions.
- Similarity scores are displayed for the entire sequence, thus allowing for the identification of shorter conserved regions, or regions with gaps.

Computational Biology

@ SC 2000



Gene Regulatory Networks and Cellular Processes

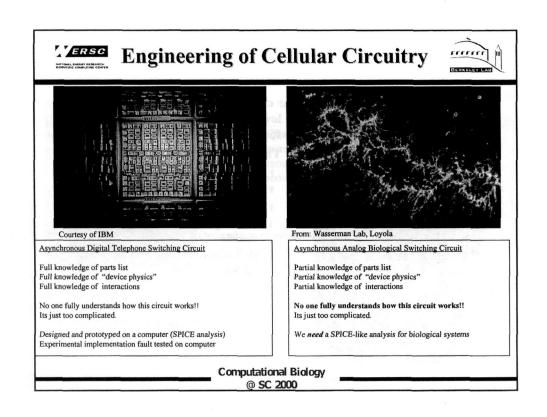
Adam Arkin APArkin@lbl.gov LBNL

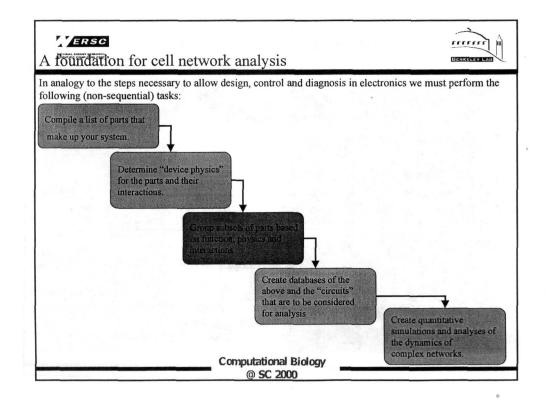




cells

Computational Biology @ SC 2000

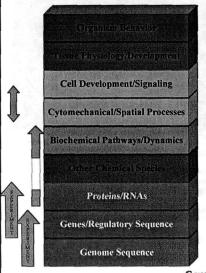






Analysis of Cell Function

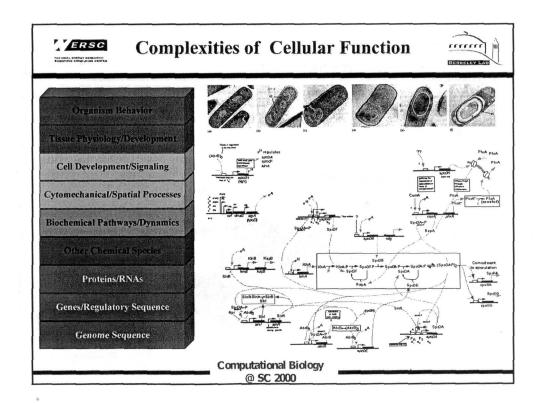


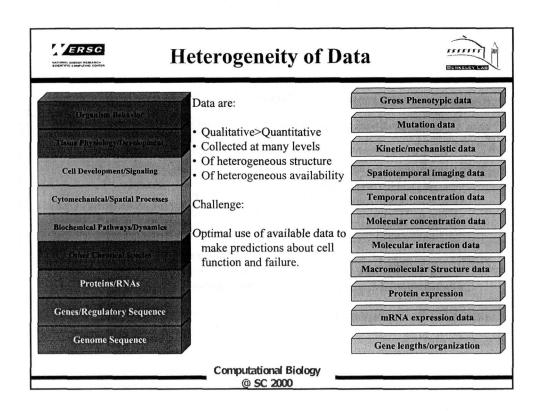


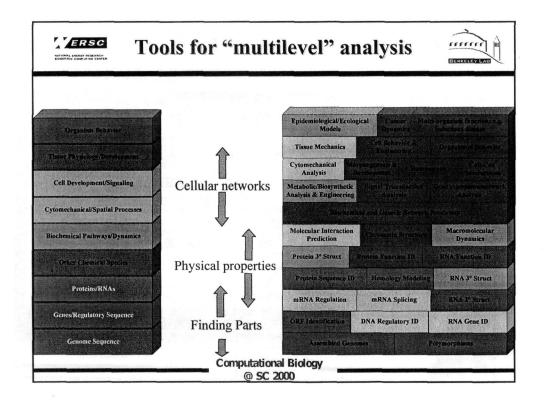
The challenge is to integrate data from all levels to produce a description of cellular function.

- † There are challenges in:
 - + Systematization and structuring of data
 - t Serving and query this data
 - t Representing the data
 - + Building multiscale, multi-resolution models
 - t Dynamic and static analysis of these models
- * Pay-off in
 - † Industrial bioengineering
 - † Rational pharmaceutical design
 - † Basic biological understanding

Computational Biology @ SC 2000







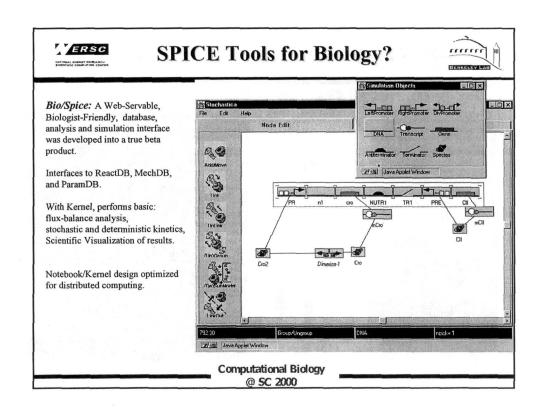


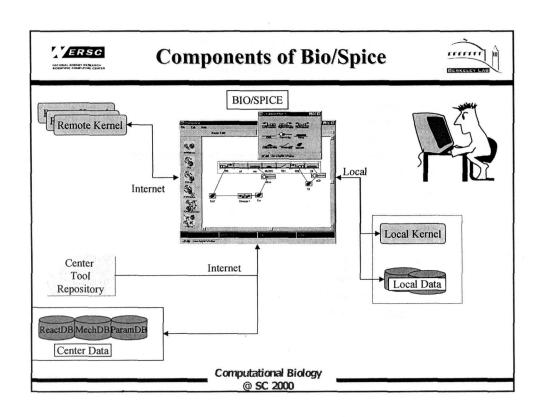
Why now?

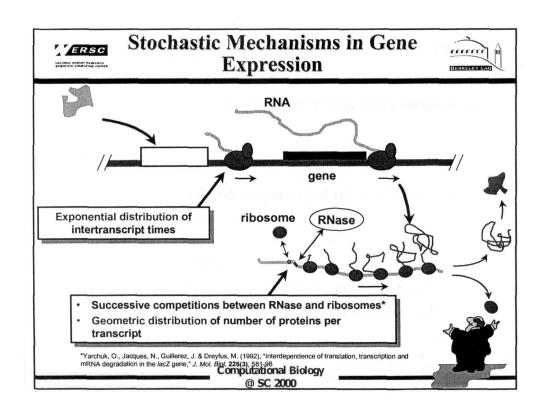


- •Genome projects are providing a large (but partial) list of parts
- •New measurement technologies are helping to identify further components, their interactions, and timings
 - · Gene microarrays
 - · Two-Hybrid library screens
 - High-throughput capillary electrophoresis arrays for DNA, proteins and metabolites
 - Fluorescent confocal imaging of live biological specimens
 - High-throughput protein structure determination
- •Data is being compiled, systematized, and served at an unprecedented rate
 - Growth of GenBank and PDB > polynomial
 - Proliferation of databases of everything from sequence to confocal images to literature
- •The tools for analyzing these various sorts of data are also multiplying at an astounding rate

@ SC 2000









Some Stochastic Cellular Phenomena



- t Lineage commitment in human hemopoiesis
- † Random, bimodal eukaryotic gene transcription in
 - † Activated T cells
 - * Steroid hormone activation of mouse mammary tumor virus
 - † HIV-1 virus
- * Clonal variation in:
 - † Bacterial chemotactic responses
 - † Cell cycle timing
- † E. coli type-1 pili expression
 - t Enhances virulence
- t Changing cell surface protein expression
 - t For immune response avoidance
- † Bacteriophage I lysis/lysogeny decision

Computational Biology

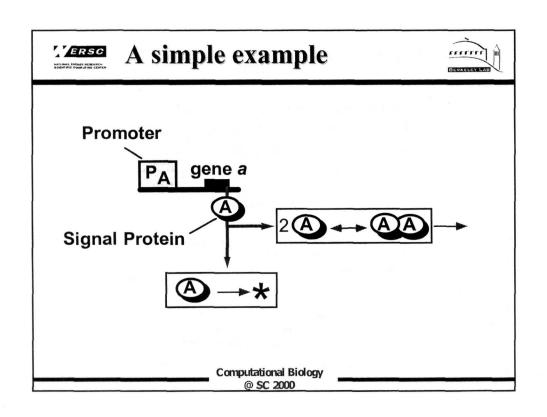


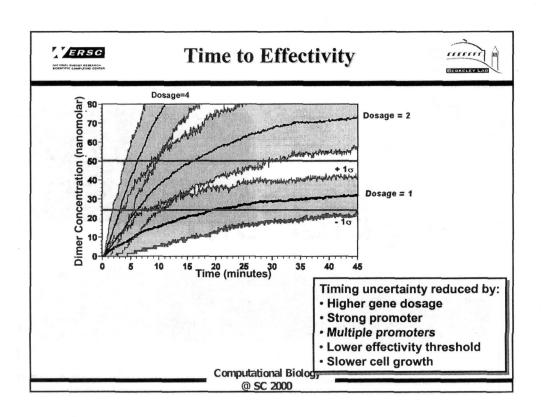
Where Noise Comes From

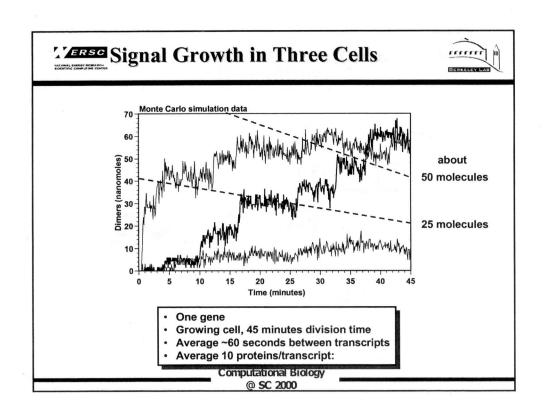


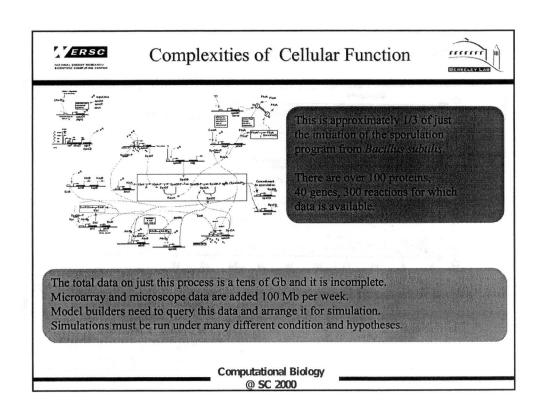
- † Random environmental influences
- * Mutations
- * Asymmetric partitioning at cell division
- † Stochastic mechanisms in gene expression
 - * Stochastic timing of gene expression
 - † Random variation in time for signal propagation
 - * Random variation total protein production

Computational Biology











The Need for Advanced Computing



* Data Handling:

The total data necessary for network analysis is huge. By nature it will be distributed and heterogeneous

We need:

- t Database standard and new query types
- t Means of secure, fast transmission of information
- † Means of quality control on data input

* Tool integration:

- † Centralization of computational biology tools and standards
- t Ability to use tools together to generate good network hypotheses
- t Good quality ratings on Tool outputs

* Advanced Simulation Tools:

- t Fast, distributed algorithms for dynamical simulation
- t Mixed mode systems (differential, Markov, algebraic, logical)
- t Spatially distributed systems

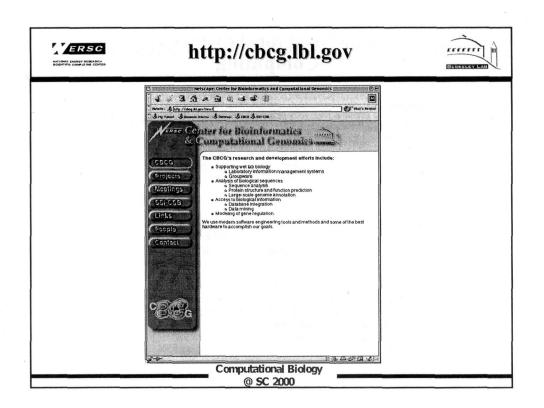
Computational Biology @ SC 2000





The End

Computational Biology @ SC 2000



ERNEST ORLANDO LAWRENCE BERKELEY NATIONAL LABORATORY ONE CYCLOTRON ROAD | BERKELEY, CALIFORNIA 94720

Prepared for the U.S. Department of Energy under Contract No. DE-AC03-76SF00098