

# **CENG 465**

## **Introduction to Bioinformatics**

Spring 2008-2009

Tolga Can (Office: B-109)  
e-mail: [tcan@ceng.metu.edu.tr](mailto:tcan@ceng.metu.edu.tr)

Course Web Page:  
[http://www.ceng.metu.edu.tr/~tcan/ceng465\\_s0809/](http://www.ceng.metu.edu.tr/~tcan/ceng465_s0809/)

# Teaching Assistant

- Dr. Ahmet Saçan
- e-mail: [ahmet@ceng.metu.edu.tr](mailto:ahmet@ceng.metu.edu.tr)
- Office: A-206

# Goals of the course

- Working at the interface of computer science and biology
  - New motivation
  - New data and new demands
  - Real impact
- Introduction to main issues in computational biology
- Opportunity to interact with algorithms, tools, data in current practice

# High level overview of the course

- A general introduction
  - what problems are people working on?
  - how people solve these problems?
  - what key computational techniques are needed?
  - how much help computing has provided to biological research?
- A way of thinking -- tackling “biological problems” computationally
  - how to look at a “biological problem” from a computational point of view?
  - how to formulate a computational problem to address a biological issue?
  - how to collect statistics from biological data?
  - how to build a “computational” model?
  - how to solve a computational modeling problem?
  - how to test and evaluate a computational algorithm?

# Course outline

- Motivation and introduction to biology (1 week)
- Sequence analysis (4 weeks)
  - Analyze DNA and protein sequences for clues regarding function
  - Identification of homologues
    - Pairwise sequence alignment
  - Statistical significance of sequence alignments
  - Sequence Motifs
  - Suffix trees
  - Multiple sequence alignment
- Phylogenetic trees, clustering methods (1 week)

# Course outline

- Protein structures (4 weeks)
  - Analyze protein structures for clues regarding function
    - Structure alignment
  - Structure prediction (secondary, tertiary)
  - Structural motifs, active sites, docking
  - Multiple structural alignment, geometric hashing
- Microarray data analysis (2 weeks)
  - Correlations, clustering
  - Inference of function
- Gene/Protein networks, pathways (2 weeks)
  - Protein-protein, protein/DNA interactions
  - Construction and analysis of large scale networks

# Grading

- Midterm exam - 25%
- Final exam - 35%
- Written assignments - 20%
- Programming assignments - 20%

# Miscellaneous

- Course webpage
  - [http://www.ceng.metu.edu.tr/~tcan/ceng465\\_s0809/](http://www.ceng.metu.edu.tr/~tcan/ceng465_s0809/)
  - Lecture slides and reading materials
  - Assignments
  - Other relevant information
- Newsgroup
  - metu.ceng.course.465
  - You should follow the newsgroup for course related announcements
  - Students from other departments should get a CENG account for this semester (Room: A-210) in order to access the newsgroup



# What is Bioinformatics?

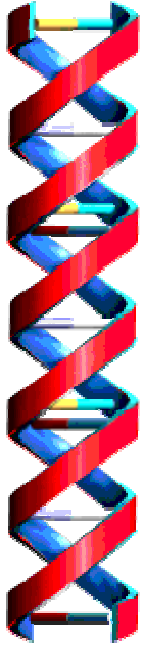
- (*Molecular*) **Bio - informatics**

- One idea for a definition?

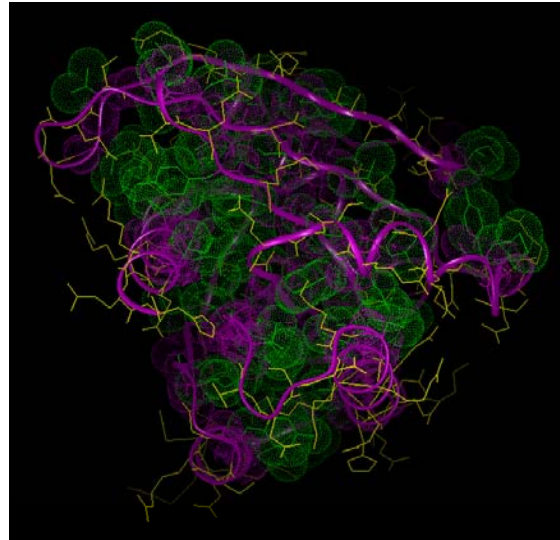
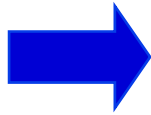
Bioinformatics is conceptualizing **biology in terms of molecules** (in the sense of physical-chemistry) and then applying **“informatics” techniques** (derived from disciplines such as applied math, CS, and statistics) to understand and **organize the information associated** with these molecules, **on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

- Bioinformatics is a practical discipline with many **applications.**

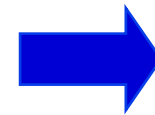
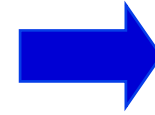
# Introductory Biology



DNA  
(Genotype)

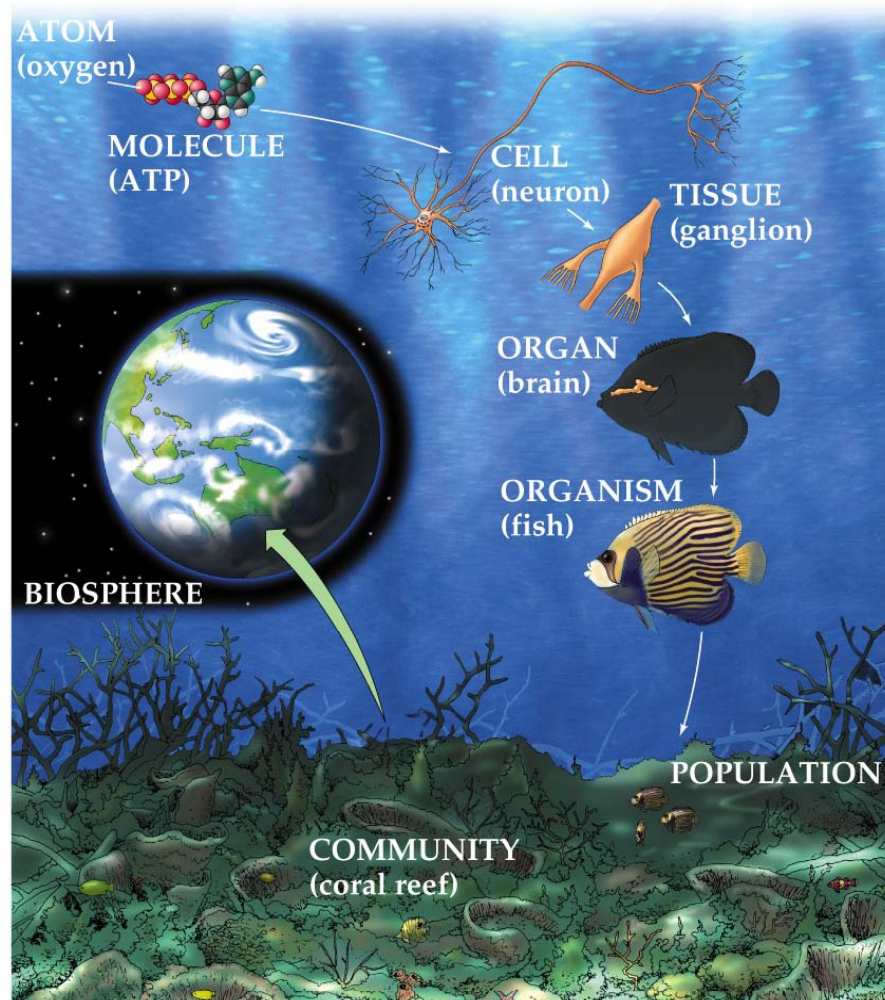


Protein



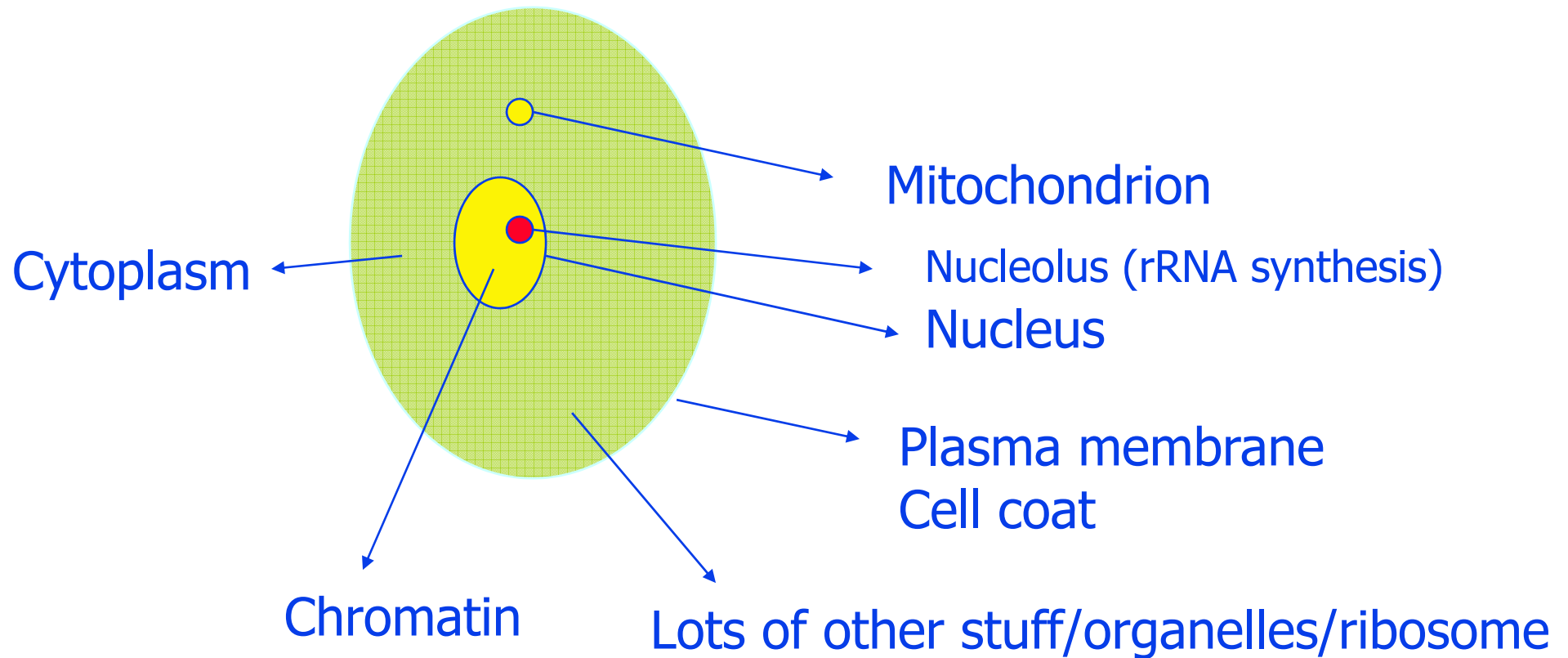
Phenotype

# Scales of life

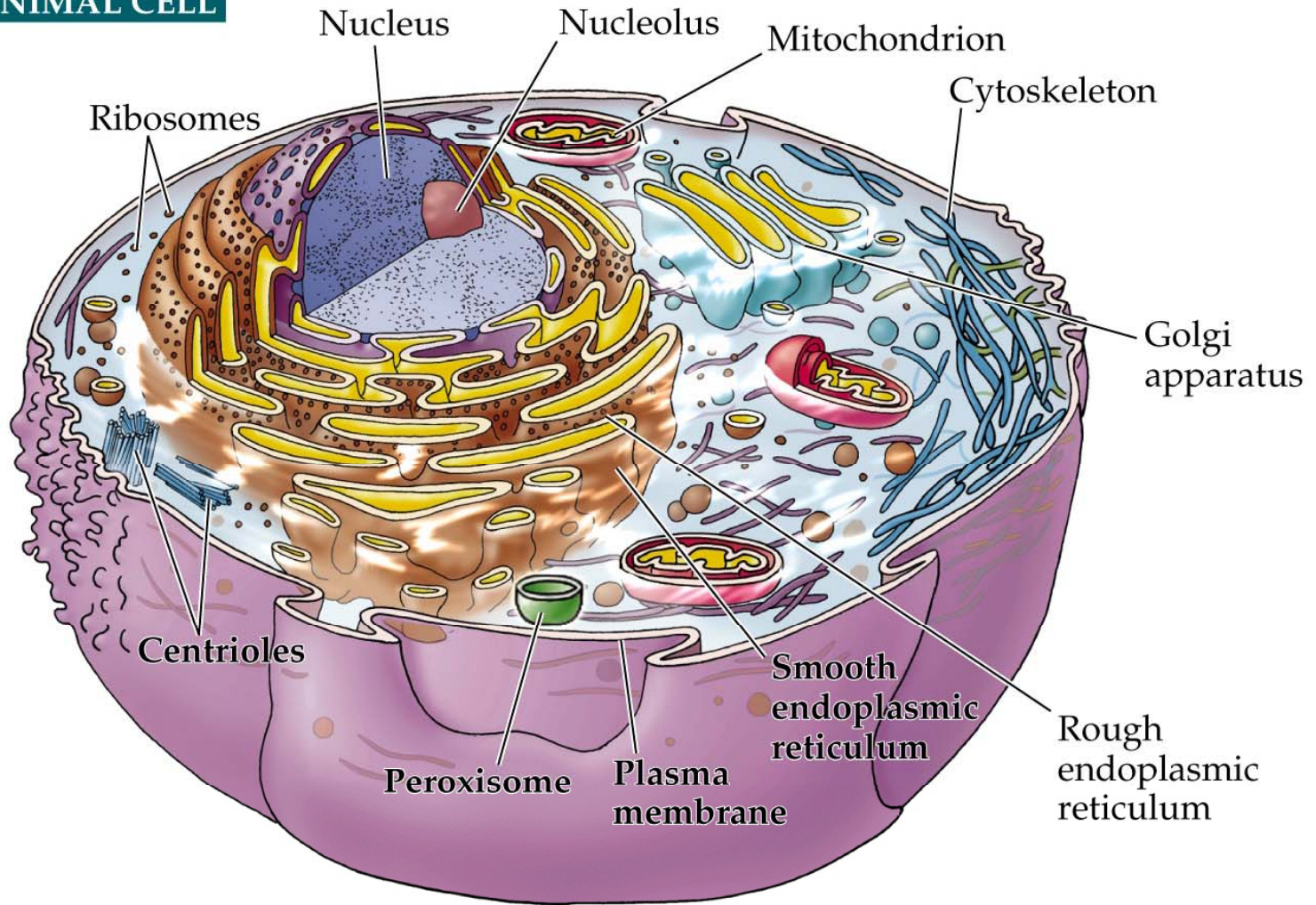


© 2001 Sinauer Associates, Inc.

# Animal Cell



## AN ANIMAL CELL



© 2001 Sinauer Associates, Inc.

# Two kinds of Cells

- Prokaryotes – no nucleus (bacteria)
  - Their genomes are circular
- Eukaryotes – have nucleus (animal, plants)
  - Linear genomes with multiple chromosomes in pairs. When pairing up, they look like



Middle: centromere

Top: p-arm

Bottom: q-arm

# Molecular Biology Information - DNA

- Raw DNA Sequence

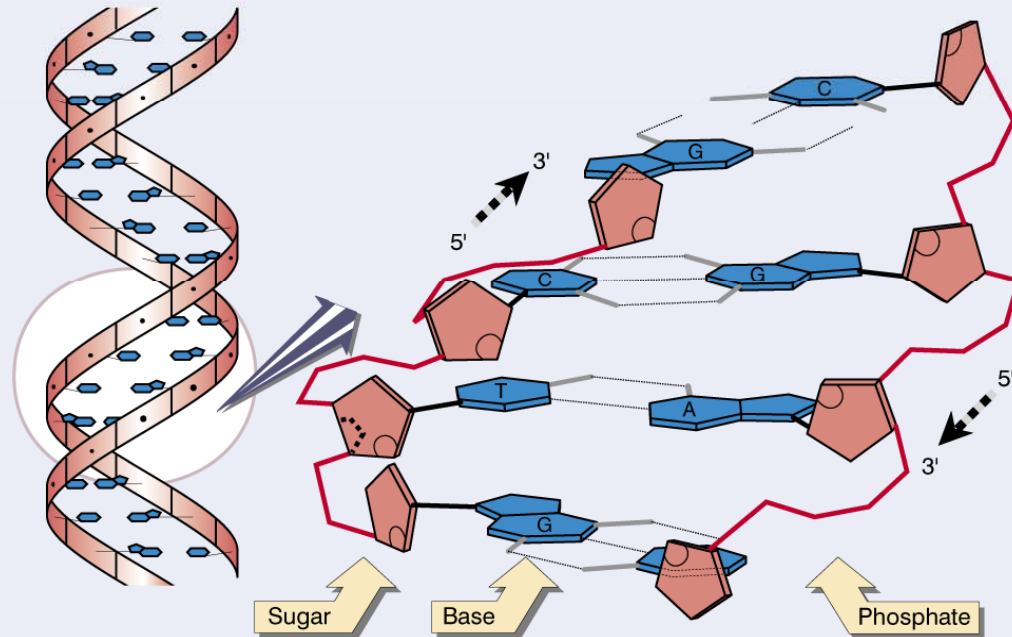
- Coding or Not?
- Parse into genes?
- 4 bases: AGCT
- ~1 Kb in a gene, ~2 Mb in genome
- ~3 Gb Human

```
atggcaattaaaattgggatcaatgggttttggtcgatcggccgatcgtattccgtgca
gcacaacaccgtgatgacattgaagttgtaggtattaacgacttaatcgacggttgaatac
atggcttataatggtgaaatatgattcaactcacggctcgttttcgacggcactggtgaagtg
aaagatggtaacttagtgggttaatggtaaaactatccgtgtaactgcagaacgtgatcca
gcaactttaaactgggggtgcaatcgggtgttgatcgcctggtgaagcgaactggtttattc
ttaactgatgaaactgctcgtaaacatatcactgcaggcgcaaaaaaagtgtattaact
ggcccatctaaagatgcaaccctatgttcggttcggtggtgtaaaacttcaacgcatacgca
ggtcaagatatcgttttctaacgcactctgtacaacaaactgttttagctccttagcacgt
gttgttcatgaaactttcgggtatcaaagatgggtttaaatacactggttcacgcaacgact
gcaactcaaaaaactgtggatgggtccatcagctaaagactggcgcgccggcgcggtgca
tcacaaaacatcattccatcttcaacaggtgcagcgaagcagtaggtaaagtattacct
gcattaaacggtaaatctaactggtatggctttccggtgttccaacgccaaacgtatctgtt
gttgatttaacagttaatcttgaaaaaccagcttcttatgatgcaatcaacaagcaatc
aaagatgcagcgggaaggtaaaacgttcaatggcgaattaaaaggcgtattaggttacct
gaagatgctgttgttttctactgacttcaacgggtgtgctttaaacttctgtatttgatgca
gacgctgggtatcgcattaactgattcttttcgttaaattggtatc . . .
```

```
. . . caaaaataggggttaatatgaatctcgcattctccattttgttcatcgtattcaa
caacaagccaaaactcgtacaaatatgaccgcacttcgctataaagaacacggcttgtgg
cgagatatctcttggaaaaactttcaagagcaactcaatcaactttctcgagcattgctt
gctcacaatattgacgtacaagataaaaatcgccatttttgcccataatatggaacgttgg
gttgttcatgaaactttcgggtatcaaagatgggtttaaatacactggttcacgcaacgact
acaatcgttgacattgacacttacaatttcgagcaatcacagtgccattttacgcaacc
aatacagcccagcaagcagaatttatcctaaatcacgccgatgtaaaaattctcttcgctc
ggcgatcaagagcaatacgcatacaacattggaaattgctcatcattgtccaaaattacaa
aaaattgtagcaatgaaatccaccattcaattacaacaagatcctctttcttgcacttgg
```

# DNA structure

**Figure 1.7** Flat base pairs lie perpendicular to the sugar-phosphate backbone.





# Molecular Biology Information: Protein Sequence

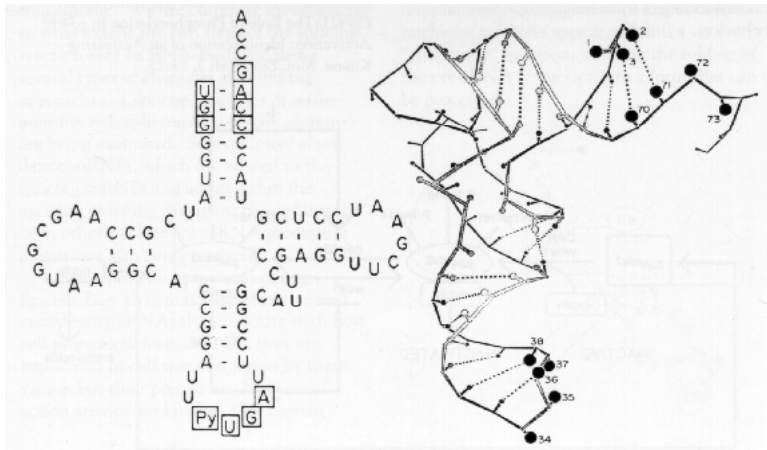
- 20 letter alphabet
  - ACDEFGHIKLMNPQRSTVWY but not BJOUXZ
- Strings of ~300 aa in an average protein (in bacteria),
  - ~200 aa in a domain
- ~1M known protein sequences

```
d1dhfa_ LNCIVAVSQNMGIGKNGDLPWPPLRNEFRYFQRMTTTSSVEGKQ-NLVIMGKKTWFSI
d8dfr_  LNSIVAVCQNMGIGKDG NLPWPPLRNEYKYFQRM TSTSHVEGKQ-NAVIMGKKTWFSI
d4dfra_ ISLIAALAVDRVIGMENAMPWN-LPADLAWFKRNTL-----NKPVIMGRHTWESI
d3dfr_  TAFLWAQDRDGLIGKDGHL PW-HLPDDLHYFRAQTV-----GKIMVGRRTYESF
```

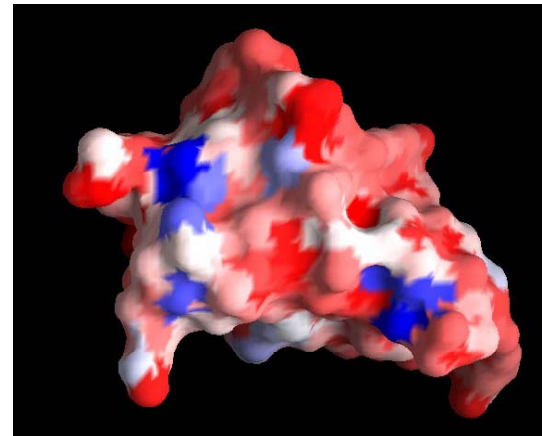
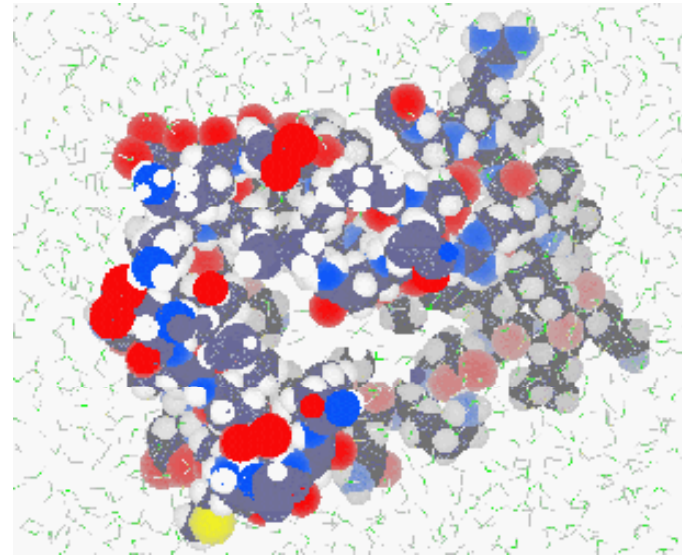
```
d1dhfa_ LNCIVAVSQNMGIGKNGDLPWPPLRNEFRYFQRMTTTSSVEGKQ-NLVIMGKKTWFSI
d8dfr_  LNSIVAVCQNMGIGKDG NLPWPPLRNEYKYFQRM TSTSHVEGKQ-NAVIMGKKTWFSI
d4dfra_ ISLIAALAVDRVIGMENAMPW-NLPADLAWFKRNTLD-----KPVIMGRHTWESI
d3dfr_  TAFLWAQDRNGLIGKDGHL PW-HLPDDLHYFRAQTVG-----KIMVGRRTYESF
```

# Molecular Biology Information: Macromolecular Structure

- DNA/RNA/Protein
  - Almost all protein

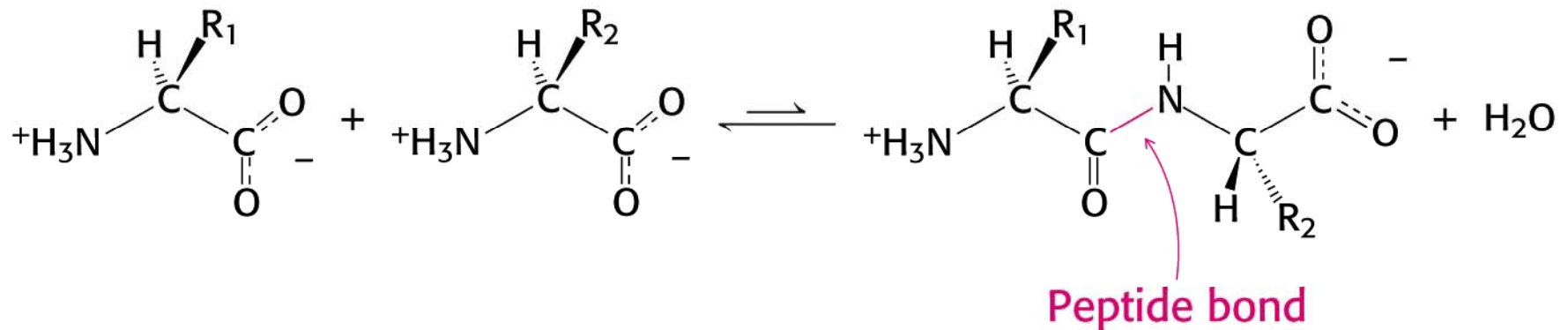


'Identity elements' in *Escherichia coli* glutamine tRNA.



# More on Macromolecular Structure

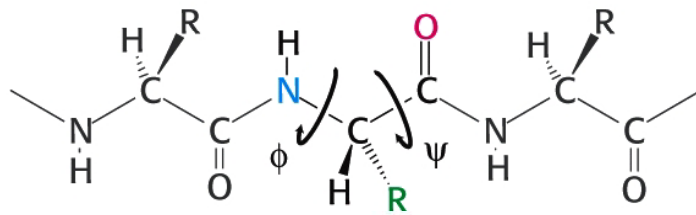
- Primary structure of proteins
  - Linear polymers linked by peptide bonds
  - Sense of direction



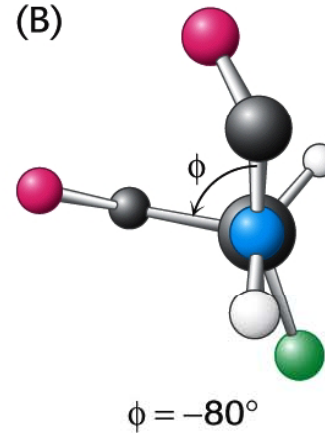
# Secondary Structure

- Polypeptide chains fold into regular local structures
  - alpha helix, beta sheet, turn, loop
  - based on energy considerations
  - Ramachandran plots

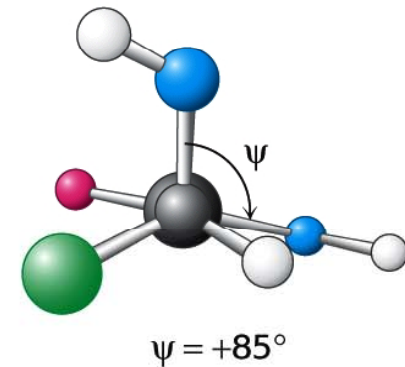
(A)



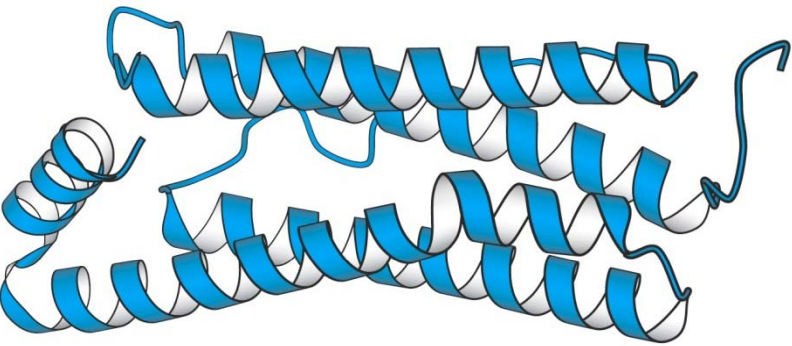
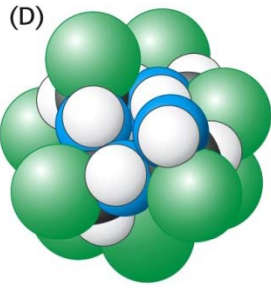
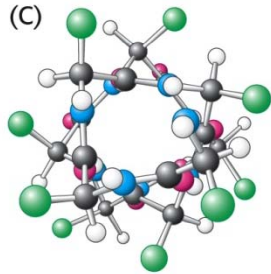
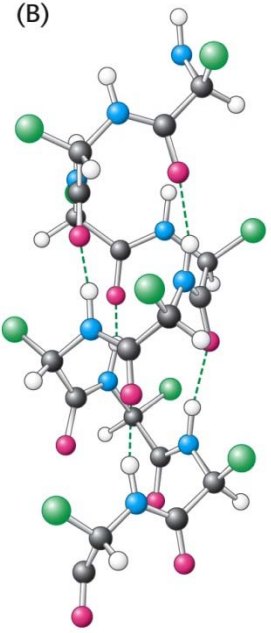
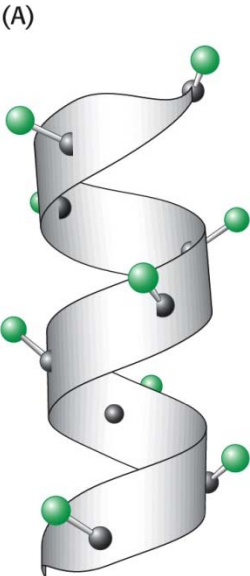
(B)



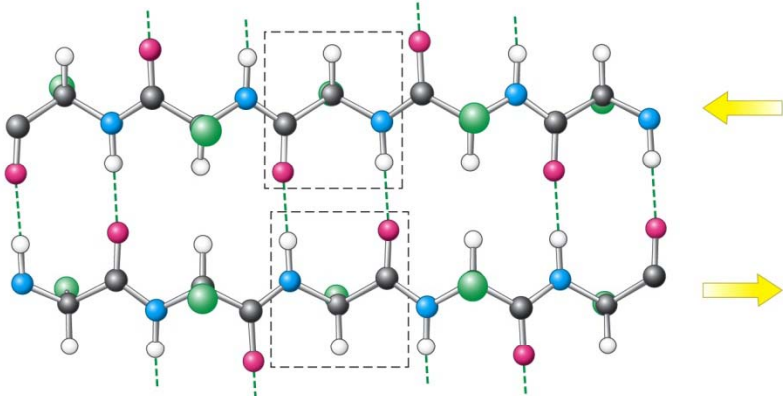
(C)



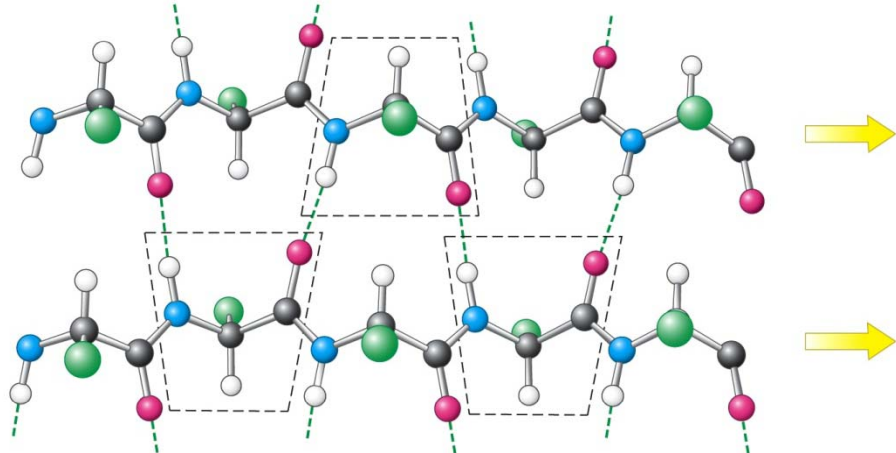
# Alpha helix



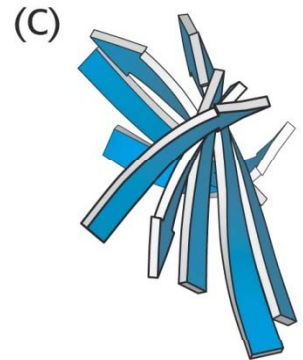
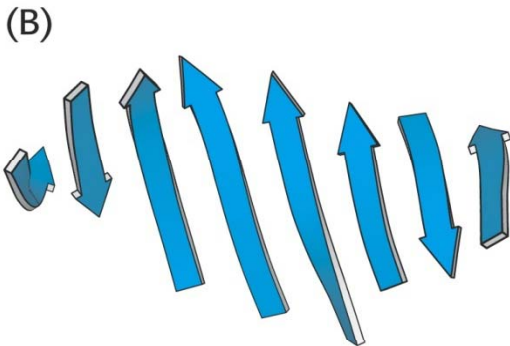
# Beta sheet



anti-parallel



parallel

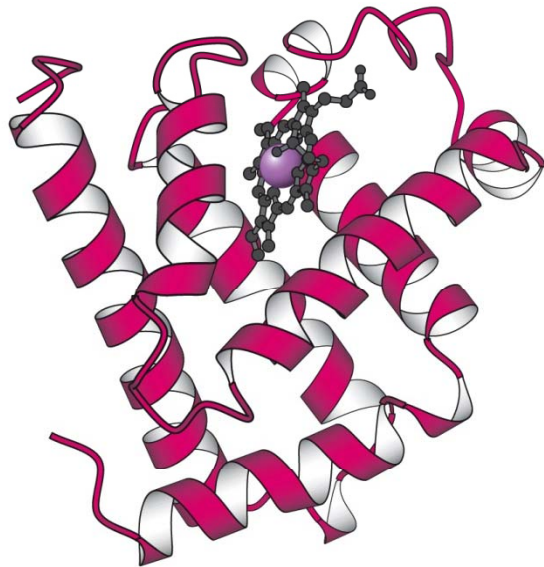


schematic

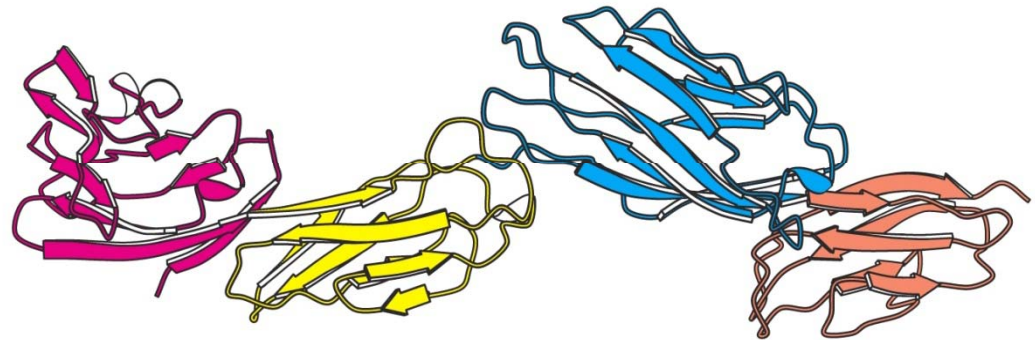
# Tertiary Structure

- 3-d structure of a polypeptide sequence
  - interactions between non-local and foreign atoms
  - often separated into domains

(B)



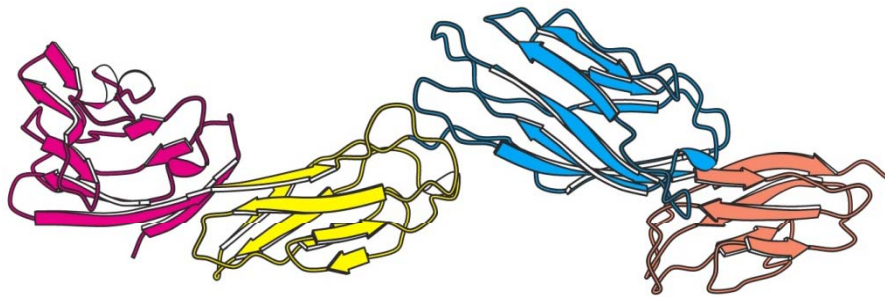
tertiary structure of  
myoglobin



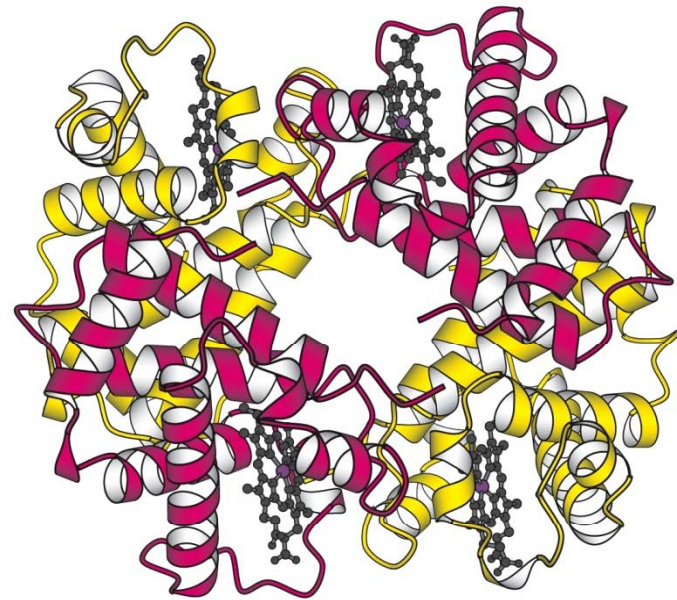
domains of CD4

# Quaternary Structure

- Arrangement of protein subunits
  - dimers, tetramers



quaternary structure  
of Cro



human hemoglobin  
tetramer



# Structure summary

- 3-d structure determined by protein sequence
- Cooperative and progressive stabilization
- Prediction remains a challenge
  - ab-initio (energy minimization)
  - knowledge-based
    - Chou-Fasman and GOR methods for SSE prediction
    - Comparative modeling and protein threading for tertiary structure prediction
- Diseases caused by misfolded proteins
  - Mad cow disease
- Classification of protein structures

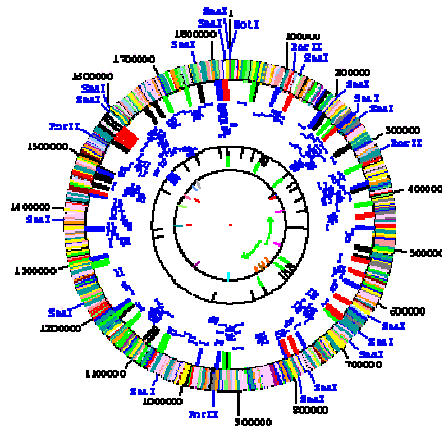
# Genes and Proteins

- One gene encodes one\* protein.
- Like a program, it starts with start codon (e.g. ATG), then each three code one amino acid. Then a stop codon (e.g. TGA) signifies end of the gene.
- Sometimes, in the middle of a (eukaryotic) gene, there are introns that are spliced out (as junk) during transcription. Good parts are called exons. This is the task of gene finding.

# A.A. Coding Table

Glycine (GLY)	GG*	Arginine (ARG)	CG*
Alanine(ALA)	GC*	Asparagine (ASN)	AAT, AAC
Valine (VAL)	GT*	Glutamine (GLN)	CAA, CAG
Leucine (LEU)	CT*	Cysteine (CYS)	TGT, TGC
Isoleucine (ILE)	AT(*-G)	Methionine (MET)	ATG
Serine (SER)	AGT, AGC	Phenylalanine (PHE)	TTT, TTC
Threonine (THR)	AC*	Tyrosine (TYR)	TAT, TAC
Aspartic Acid (ASP)	GAT, GAC	Tryptophan (TRP)	TGG
Glutamic Acid (GLU)	GAA, GAG	Histidine (HIS)	CAT, CAC
Lysine (LYS)	AAA, AAG	Proline (PRO)	CC*
Start: ATG, CTG, GTG		Stop	TGA, TAA, TAG

# Molecular Biology Information: Whole Genomes

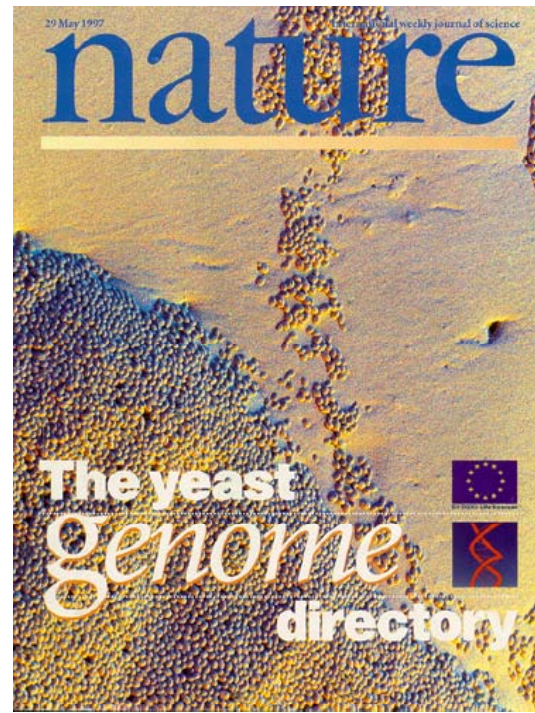


Genome sequences now accumulate so quickly that, in less than a week, a single laboratory can produce more bits of data than Shakespeare managed in a lifetime, although the latter make better reading.

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

**1995**

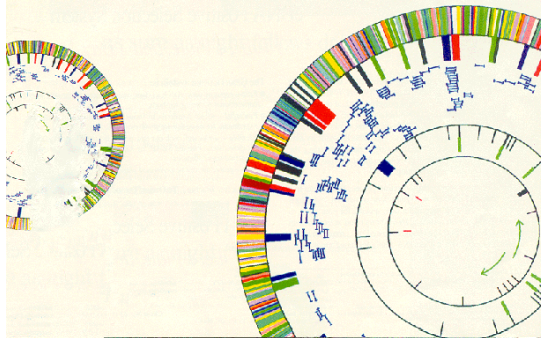
Bacteria,  
1.6 Mb,  
~1600 genes  
[*Science* 269: 496]



**Genomes highlight the Finiteness of the “Parts” in Biology**

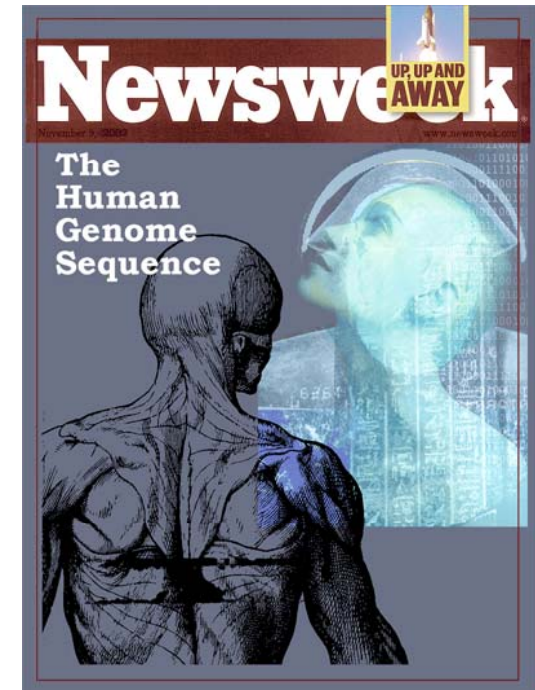
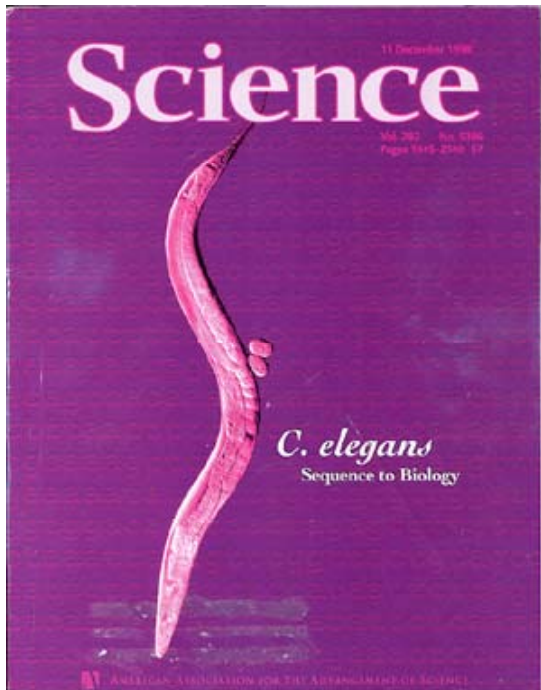
**1997**

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



**1998**

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



**2000?**

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



# Human Genome Project



**Impacting  
many  
disciplines**

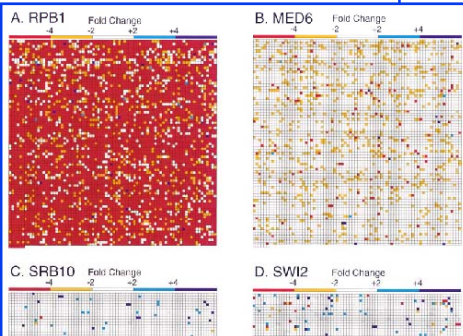
*Courtesy  
U.S. Department of Energy  
Human Genome Program*

***Global Carbon Cycles  
Industrial Resources • Bioremediation  
Evolutionary Biology • Biofuels • Agriculture • Forensics  
Molecular and Nuclear Medicine • Health Risks***

YGA 99-1133R

## Dissecting the Regulatory Circuitry of a Eukaryotic Genome

Frank C. P. Holstege,\* Ezra G. Jennings,\*<sup>1</sup> John J. Wyrick,\*<sup>1</sup> Tong Ihn Lee,\*<sup>1</sup> Christoph J. Hengartner,\*<sup>1</sup> Michael R. Green,\*<sup>1</sup> Todd R. Golub,\*<sup>5</sup> Eric S. Lander,\*<sup>1</sup> and Richard A. Young\*<sup>1||</sup>  
 \*Whitehead Institute for Biomedical Research  
 Cambridge, Massachusetts 02142  
<sup>1</sup>Department of Biology  
 Massachusetts Institute of Technology  
 Cambridge, Massachusetts 02139  
<sup>2</sup>Howard Hughes Medical Institute  
 Program in Molecular Medicine  
 University of Massachusetts Medical Center  
 Worcester, Massachusetts 01605  
<sup>3</sup>Dana-Farber Cancer Institute and  
 Harvard Medical School  
 Boston, Massachusetts 02115



## Young/Lander, Chips, Abs. Exp.

### Specific transcription factors, a novel mechanism

change in mRNA levels when a mutant is compared to its isogenic wild-type counterpart is presented in a grid format. In the left grid square represents the left-most gene on chromosome I, and the squares to its right represent adjacent genes, in fashion through chromosome I, then II, etc., until the last gene on the right arm of chromosome XVI is reached grid. The results are shown for (A) Rpb1, (B) Med6, (C) Srb10, and (D) Swi2.

Figure 2. Genome-Wide Expression Data for Selected Components of the RNA Polymerase II Holoenzyme

**The Brown Lab**  
 Stanford University Department of Biochemistry

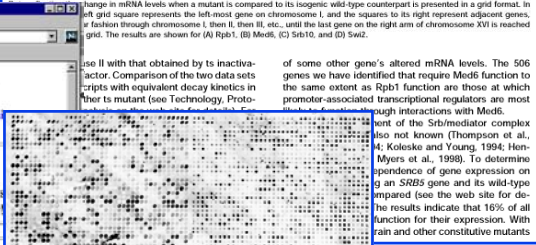
**The MGuide**  
 The Complete Guide to MicroArray  
 Build your own arrayer and scanner

**The transcriptional program in the responsive human fibroblasts to serum**  
 The web supplement to Iyer V.R. et al., (1997) Science 283: 82-87

**The Transcriptional Program of Sporulation in Budding Yeast**  
 The Web Companion to the Science Magazine Research Article

**Exploring the Gene Expression Database**

**See the entire Timecourse**



## Brown, μarray, Rel. Exp. over Timecourse

# Gene Expression Datasets: the Transcriptome

Proc. Natl. Acad. Sci. USA  
 Vol. 94, pp. 190-195, January 1997  
 Genetics

## A multipurpose transposon system for analyzing protein production, localization, and function in *Saccharomyces cerevisiae*

PETRA ROSS-MACDONALD, AMY SHEEHAN, G. SHIRLEEN ROEDER, AND MICHAEL SNYDER\*

Department of Biology, Yale University, P.O. Box 208103, New Haven, CT 06520-8103

Communicated by Gerald R. Fink, Whitehead Institute, Cambridge, MA, October 30, 1996 (received for review July 15, 1996)

**ABSTRACT** Analysis of the function of a particular gene product typically involves determining the expression profile of the gene, the subcellular location of the protein, and the phenotype of a null strain lacking the protein. Conditional alleles of the gene are often created as an additional tool. We have developed a multifunctional, transposon-based system that simultaneously generates constructs for all the above analyses and is suitable for mutagenesis of any given *Saccharomyces cerevisiae* gene. Depending on the transposon used, the yeast gene is fused to a coding region for β-galactosidase or green fluorescent protein. Gene expression can therefore be monitored by chemical or fluorescence assays. The transposons create insertion mutations in the target gene, allowing phenotypic analysis. The transposon can be reduced by Cre site-specific recombination to a smaller element that leaves a epitope tag inserted in the encoded protein. In addition to its utility for a variety of immunodetection purposes, the epitope tag element also has the potential to create conditional alleles of the target gene. We demonstrate these features of the transposons by mutagenesis of the *SP24*, *ARP100*, *SER1*, and *BDF1* genes.

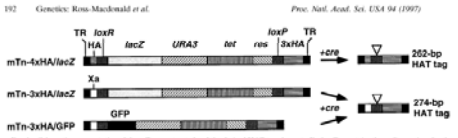


Fig. 1. Schematic representation of the mTn-constructed and derived HAT tag elements. Each mTn contains the coding region for the retroviral origin (ori) and Ura3 protein, and the ori element from Trf1-mTn-4xHA/lacZ and mTn-3xHA/lacZ contain a transposed lacZ gene. In mTn-3xHA/GFP, the coding region for GFP instead of lacZ. In each case, these are flanked by loxP sites and Trf1 terminal repeats (TR). Between loxP and the right TR, each transposon contains a sequence encoding three tandem copies of the HA epitope. Between the left TR and loxP is a sequence encoding either an additional copy of the HA epitope (mTn-4xHA/lacZ) or the factor Xa protease cleavage site (mTn-3xHA/lacZ, mTn-3xHA/GFP). Disruption of these transposons to Cre recombinase enables the formation of a similar element encoding HAT tag elements (see text). The size of the Cre recombinase product is indicated by a triangle. (Not drawn to scale.)

The yeast *Saccharomyces cerevisiae* has proved of great importance in characterizing basic biological processes. This utility can only become more marked now that the sequence of the entire yeast genome has been obtained, and additional homologs of yeast genes are identified in other organisms (1). Determination

## Also: SAGE; Samson and Church, Chips; Aebersold, Protein Expression

## Snyder, Transposons, Protein Exp.

transposition, and contain the Trf1 site for resolution of transposon integration. Trf1-mediated enzymes catalyzing transposition and resolution are provided in more. All three transposons contain the *URA3* and/or genes for selection in *S. cerevisiae* and *E. coli*, respectively. Transposons mTn-3xHA/lacZ and mTn-3xHA/GFP contain a *lacZ* gene that lacks an initiator methionine, while transposon mTn-3xHA/GFP contains the entire coding region for a mutant derivative of GFP that shows enhanced fluorescence (10, 11). These elements allow identification of in-frame fusions between a transposon and a yeast coding region by size of assays for either β-gal or fluorescence activity. Levels of both activities can be measured quantitatively and have been shown to provide indices of gene expression (refs. 4, 23, and 26).

A *loxP* site is located at one end of the transposon and a *loxP* element lies at the other end. These target sites for the Cre recombinase are divergent from one another and undergo low levels of spontaneous recombination. The *loxP* sites are internal to sequences encoding multiple copies of an epitope from the influenza virus hemagglutinin protein (HA epitope; ref. 25). The mTn-3xHA transposon also contains a factor Xa protease cleavage site (17) in the region external to the *loxP* site. Expression of the Cre recombinase induces recombination between the *loxP* sites resulting in excision of the central region of the transposon. The final product contains a 5-bp duplication caused by transposon insertion in addition to a 274-bp (mTn-3xHA) or 262-bp (mTn-3xHA) element. This element consists of a single, full-size and sequence encoding three or four copies of the HA epitope. Flanked by the Trf1 terminal repeats (Fig. 1). The mTn-3xHA-derived element also contains a sequence encoding the factor Xa cleavage site. When the transposon has inserted into a gene to generate an in-frame fusion of *lacZ* or GFP coding sequences, the excision event results in insertion of 93 amino acids (mTn-3xHA) or 89 amino acids (mTn-3xHA) into the protein. We designate these insertions HAT tags.

**Mutagenesis of Yeast Genes.** Transposons mTn-3xHA/lacZ and mTn-3xHA/GFP were tested by mutagenesis of the yeast *SP24* gene. *SP24* encodes a monomeric protein that localizes to sites of polarized growth and mutants exhibit defects in

was investigated in *S. cerevisiae* by shuttle mutagenesis. DNA containing the transposon was then excised from the plasmid and transformed into yeast, where it replaced the chromosomal locus by homologous recombination.

With both mTn-3xHA/lacZ and mTn-3xHA/GFP, about 10% of transformants were identified as producing β-galactosidase protein (441,500 mTn-3xHA/lacZ transformants for *SP24* and 92,000 mTn-3xHA/GFP and 62,000 mTn-3xHA/lacZ transformants for *ARP100*, *SER1*, and *SER2*, respectively). Strains expressing the reporter genes were used for further analysis. The appropriate position of the transposon insertion in these strains was determined by size analysis of PCR products obtained from their genomic DNA (Materials and Methods). In some instances PCR products were sequenced, enabling exact identification of insertion points (Fig. 2).

### Efficiency of Cre-Mediated *loxP-loxP* Recombination, Although Efficient Cre-Mediated Recombination between *loxP*

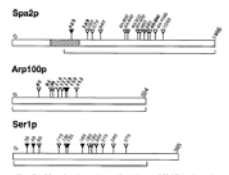


Fig. 2. Map showing amino acid positions of HAT tag insertion in the yeast proteins Sp24p, Arp100p, and Ser1p. Regions mutagenized are indicated by vertical bars. The amino acid positions are indicated by the numbers. The HAT tag consists of three or four copies of the HA epitope. For the HAT tag, the amino acid positions are indicated by the numbers.

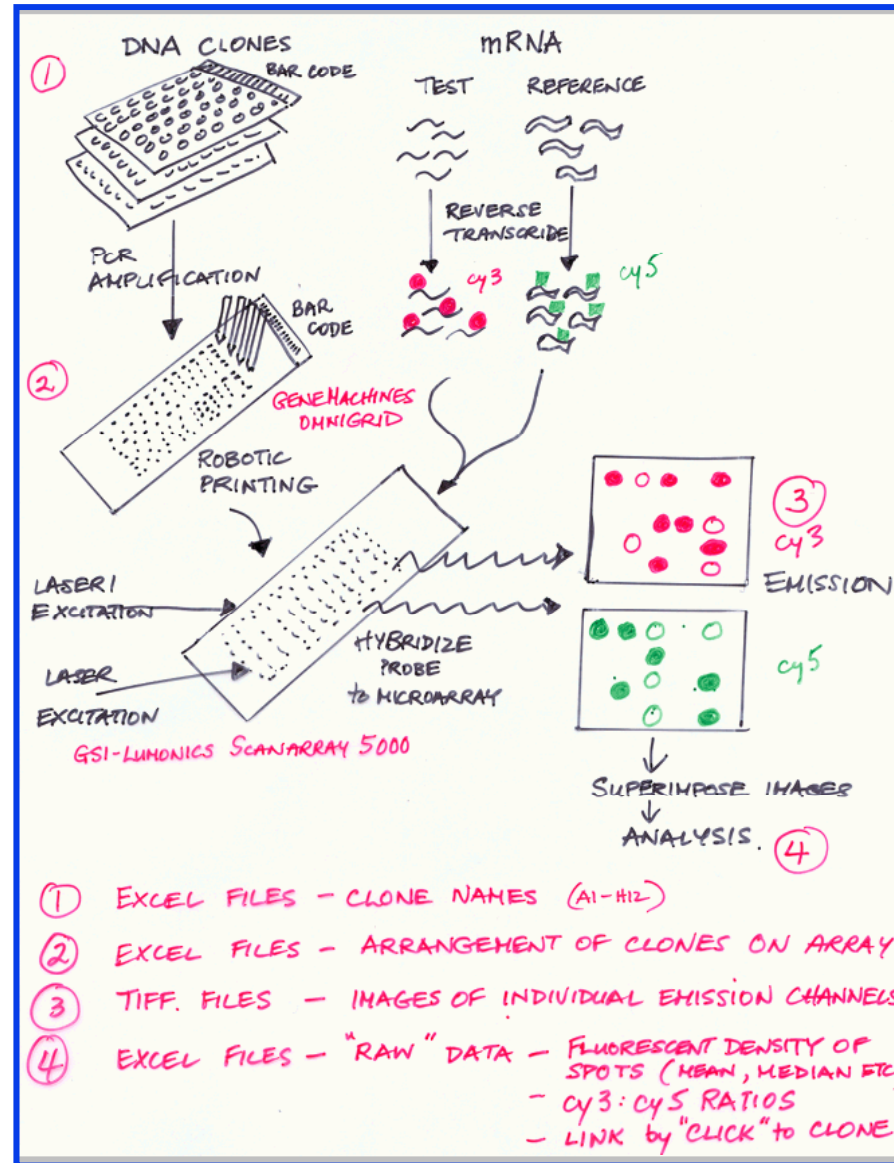
# Array Data

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

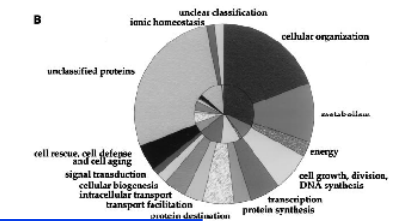
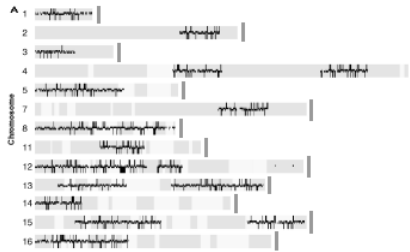
Elizabeth A. Winzeler,<sup>1\*</sup> Daniel D. Shoemaker,<sup>2\*</sup> Anna Astromoff,<sup>1\*</sup> Hong Liang,<sup>1\*</sup> Keith Anderson,<sup>1</sup> Bruno Andre,<sup>3</sup> Rhonda Bangham,<sup>4</sup> Rocio Benito,<sup>5</sup> Jef D. Boeke,<sup>6</sup> H. Carla Connelly,<sup>6</sup> Karen Davis,<sup>1</sup> Mohamed El Bakkoury,<sup>3</sup> Françoise Erik Gentalen,<sup>11</sup> Guri Giaever,<sup>1</sup> Ted Jones,<sup>1</sup> Michael Laub,<sup>1</sup> Howard David J. Lockhart,<sup>11</sup> Anca Lu Nasilha M'Rabet,<sup>3</sup> Patrice M. Chai Pai,<sup>1</sup> Corinne Rebschung,<sup>8</sup> Christopher J. Roberts,<sup>2</sup> Petra R. Michael Snyder,<sup>4</sup> Sharon Sookha Steeve Veronneau,<sup>7</sup> Marleer Teresa R. Ward,<sup>2</sup> Robert Wysocki Katja Zimmermann, Mark Johnston,<sup>13</sup>

The functions of many open reading sequencing projects are unknown. New, to systematically determine their function *S. cerevisiae* strains were constructed, by a precise deletion of one of 2026 ORFs (genome). Of the deleted ORFs, 17 per medium. The phenotypes of more than parallel. Of the deletion strains, 40 per in either rich or minimal medium.

The budding yeast *S. cerevisiae* serves as an important experimental organism for revealing gene function. In addition to carrying out all the

other essential genes (60% of M:1.0, R:3, and *ade1* (0.56, M:1.0). In addition, the *gpf1* (0.78, M:0.99, R) deletion showed a minimal medium-specific growth defect (15). *GTP1* (*YOR070C*) is a GTPase

ic analysis of the deletion ar those whose cognate ential to life, is a formidable of many genes will likely under very specialized n, necessitating the exami- ferent conditions. Previous ed that the barcodes allowed ndance of their respective arded when 12 strains were vely for many generations ng scheme thus has the ponate the phenotypic analysis rans by allowing the growth ms to be assayed simulta- 558 homozygous deletion ed were pooled (12) and d minimal media for about During this time, aliquots en the two pools. The tags and hybridized to high-den- ining the tag complements The hybridization data were the relative growth rates for tant in the population (14), that the growth rate for each independently with the UP- NTAG signals would agree, hich both the UPTAG and



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 construct deletion strains for all annotated

essential genes (60% of M:1.0, R:3, and *ade1* (0.56, M:1.0). In addition, the *gpf1* (0.78, M:0.99, R) deletion showed a minimal medium-specific growth defect (15). *GTP1* (*YOR070C*) is a GTPase

ic analysis of the deletion ar those whose cognate ential to life, is a formidable of many genes will likely under very specialized n, necessitating the exami- ferent conditions. Previous ed that the barcodes allowed ndance of their respective arded when 12 strains were vely for many generations ng scheme thus has the ponate the phenotypic analysis rans by allowing the growth ms to be assayed simulta- 558 homozygous deletion ed were pooled (12) and d minimal media for about During this time, aliquots en the two pools. The tags and hybridized to high-den- ining the tag complements The hybridization data were the relative growth rates for tant in the population (14), that the growth rate for each independently with the UP- NTAG signals would agree, hich both the UPTAG and

ial (short black bars) and 356 essential genes (all cutive groups on multiple chromosomes. A lighter posional duplication blocks (23). For 15 of the 356 previously described. These inconsistencies may be the conditions used for germination of spores. For

# Other Whole-Genome Experiments

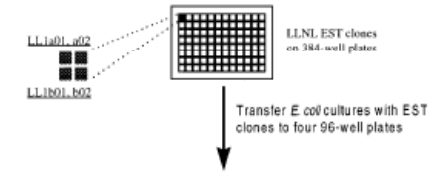


## Construction of a modular yeast two-hybrid cDNA library from human EST clones for the human genome protein linkage map

Shao-bing Hua<sup>1\*</sup>, Ying Luo<sup>1,2</sup>, Mengsheng Qiu<sup>3</sup>, Eva Chan<sup>2</sup>, Helen Zhou<sup>4</sup>, Li Zhu  
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 29 April 1998; Received by E.Y. Chen

### Abstract

Identification of all human protein important information for functional studying protein-protein interactions construct two-hybrid cDNA libraries we have constructed a modular human Quality analysis of this library indicates human EST clones is feasible, and su first time that a comprehensive two- EST clones. © 1998 Elsevier Science



Keywords: Functional genomics, yeast two-hybrid, human EST clones, protein linkage map

### 1. Introduction

The Human Genome Project has produced a tremendous amount of DNA sequence data. Over 50 000 UniGenes have been identified (Schuler, 1995; Miller et al., 1996). Approximately 50% of these genes are expressed in human cells (Rowen et al., 1996). A minority of these UniGenes

## 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

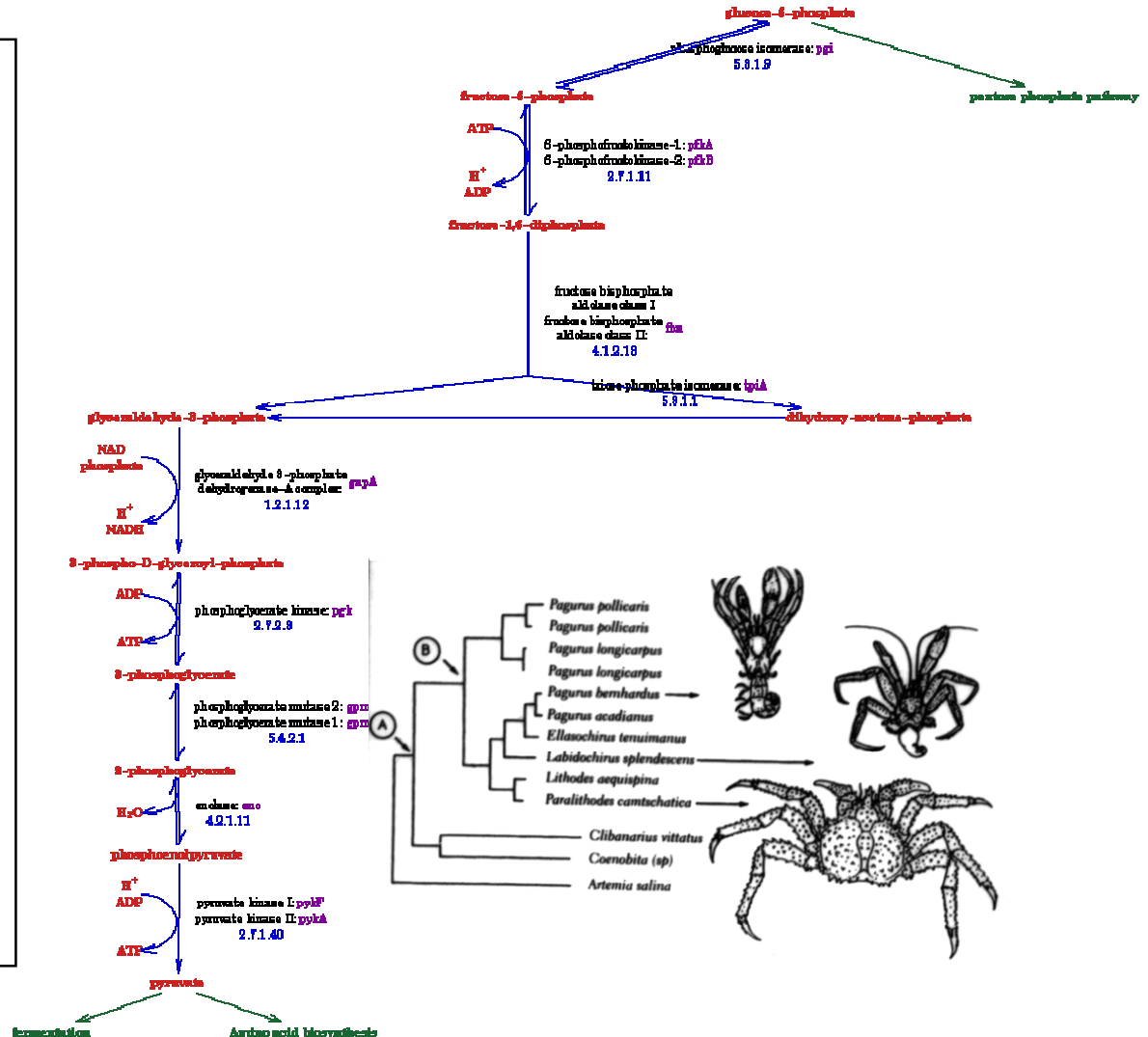
For yeast:  
6000 x 6000 / 2  
~ 18M interactions

## 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

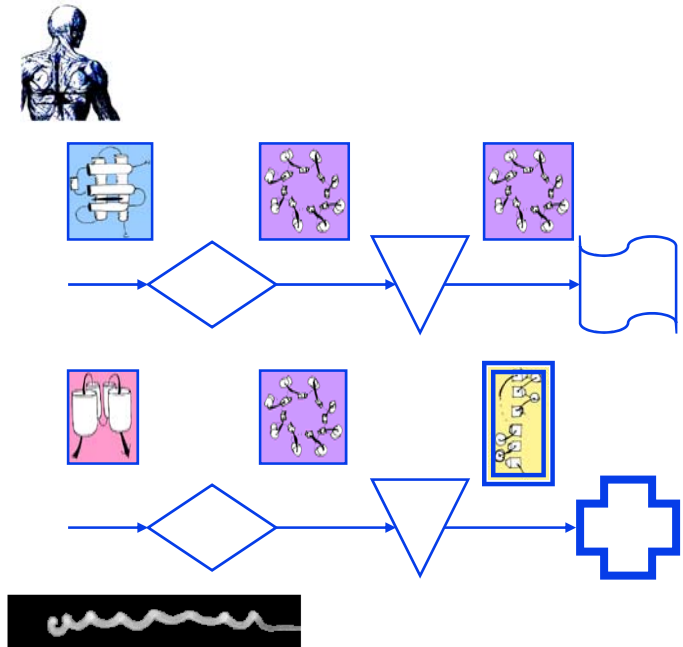
# Molecular Biology Information: Other Integrative Data

- Information to understand genomes
  - Metabolic Pathways (glycolysis), traditional biochemistry
  - Regulatory Networks
  - Whole Organisms Phylogeny, traditional zoology
  - Environments, Habitats, ecology
  - The Literature (MEDLINE)
- The Future....



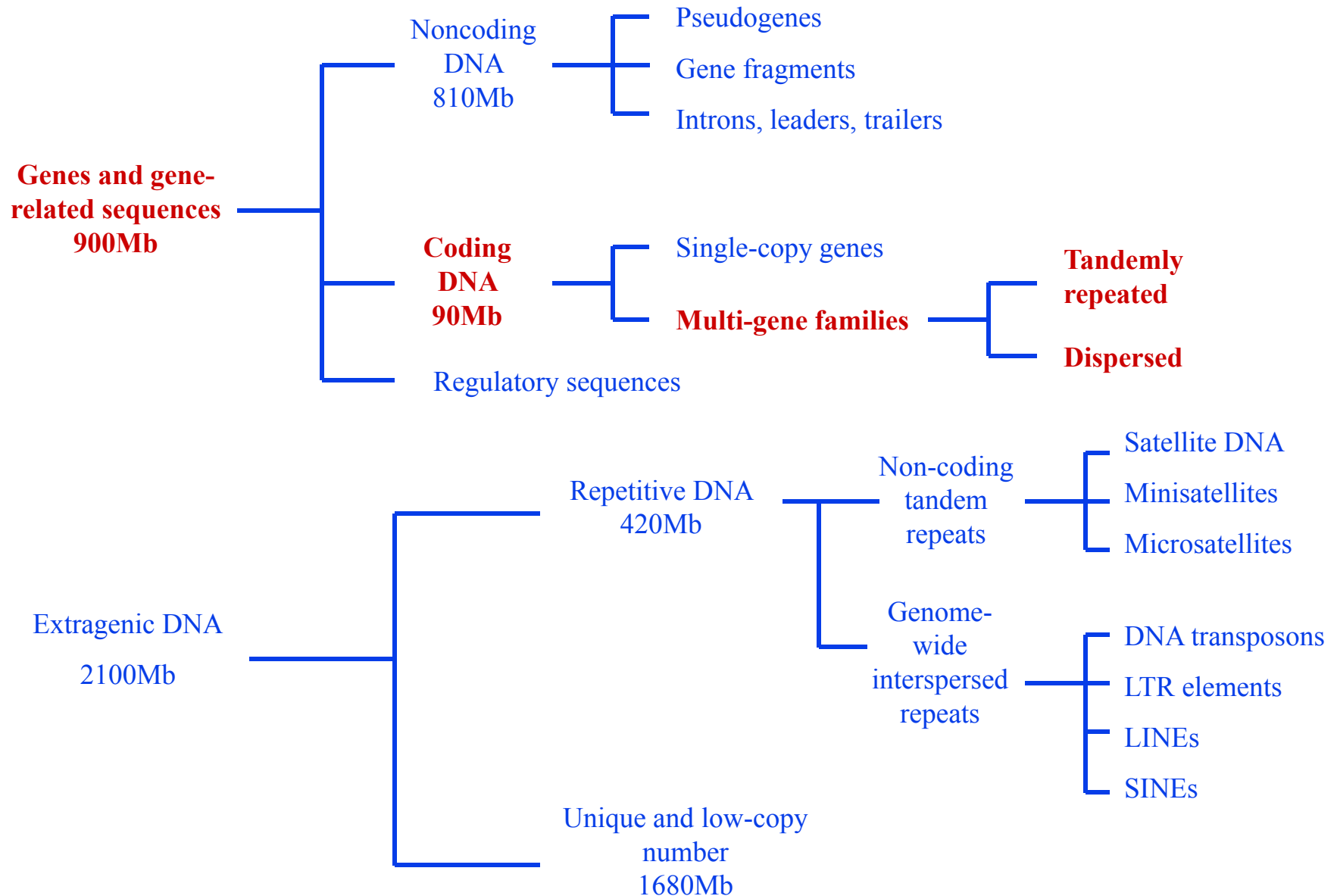
# Organizing Molecular Biology Information: Redundancy and Multiplicity

- 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 ↔ ....

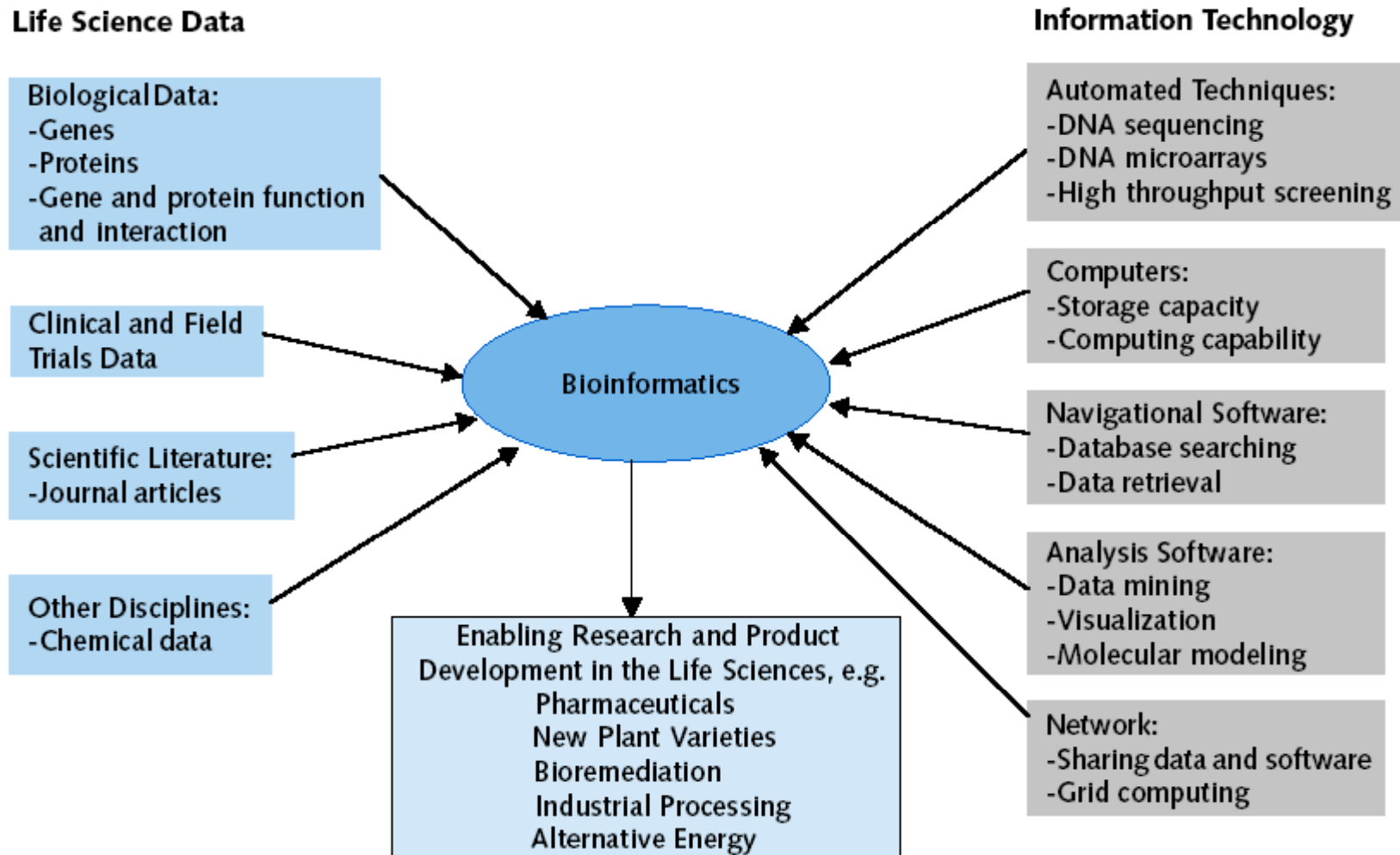
# Human genome



# Where to get data?

- GenBank
  - <http://www.ncbi.nlm.nih.gov>
- Protein Databases
  - SWISS-PROT: <http://www.expasy.ch/sprot>
  - PDB: <http://www.pdb.bnl.gov/>
- And many others

**Figure 6.1. Bioinformatics Uses Information Technology to Manage and Analyze Information Generated by the Life Sciences**

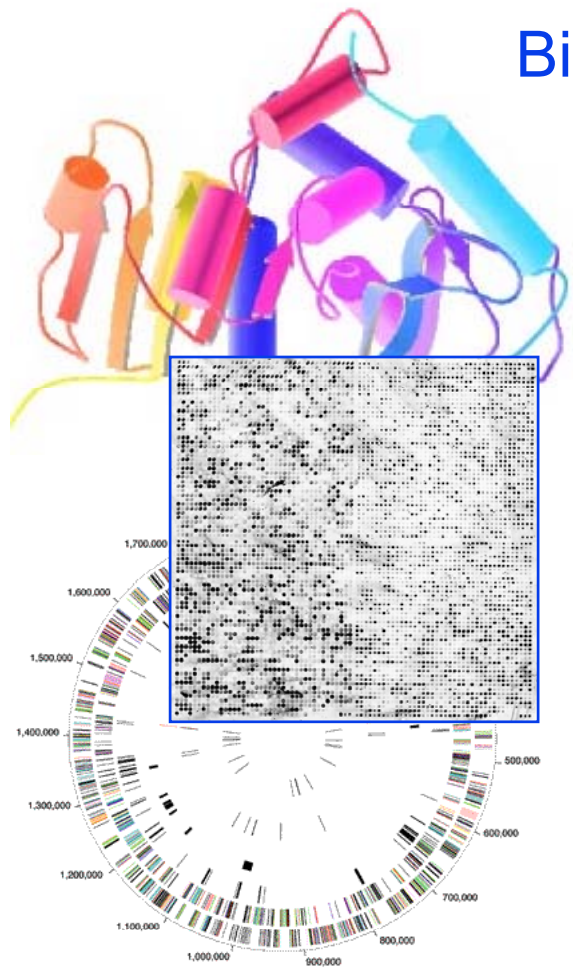


# Bioinformatics: A simple view

Biological  
Data

+

Computer  
Calculations



# Application domains

**Table 6.2. Number of Survey Respondents Indicating Bioinformatics Research Activities by Application, 2002**

Application	Number of firms in application	Conduct bioinformatics research
Human Health	780	247
Animal Health	144	37
Agricultural & Aquacultural/Marine	128	41
Marine & Terrestrial Microbial	41	19
Industrial and Agricultural-Derived Processing	132	45
Environmental Remediation and Natural Resource Recovery	41	12
Other <b>Bio-defense</b>	160	30

Note: The total number of firms that responded to the biotechnology survey was 1,031, and 304 of these firms indicated that they had some activity in bioinformatics. The number of firms by biotechnology application does not add up to the total number of firms that responded to the survey because firms were classified in an application if they indicated it as either a "primary" or "secondary" focus.

Source: Survey data from *Critical Technology Assessment of Biotechnology in U.S. Industry*, U.S. Department of Commerce, Technology Administration and Bureau of Industry and Security, August 2002.



# Kinds of activities

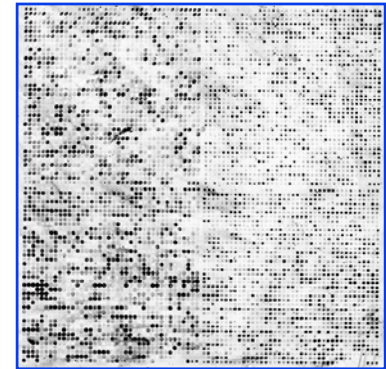
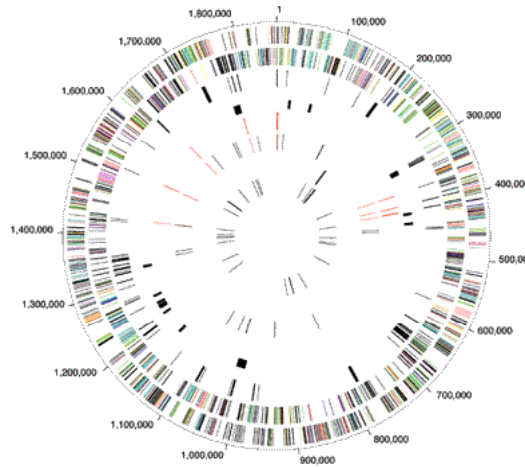
	Conduct research on/in	Approved, marketed, or in production		Total
		Product(s)	Process(es)	
<b>DNA-based</b>				
Bioinformatics	29	2	1	30
Genomics, pharmacogenetics	29	3	2	30
DNA sequencing/synthesis/ amplification, genetic engineering	39	5	3	43
<b>Biochemistry/Immunology</b>				
Drug design & delivery	33	4	2	38
Synthesis/sequencing of proteins and peptides	27	3	1	30
Combinatorial chemistry, 3-D molecular modeling	18	1	0	19

Note: The total number of responses to the biotechnology activity question was 1021. Percents do not add up to 100 percent because firms can have more than one activity.

Source: Survey data from *Critical Technology Assessment of Biotechnology in U.S. Industry*, U.S. Department of Commerce, Technology Administration and Bureau of Industry and Security, August 2002.

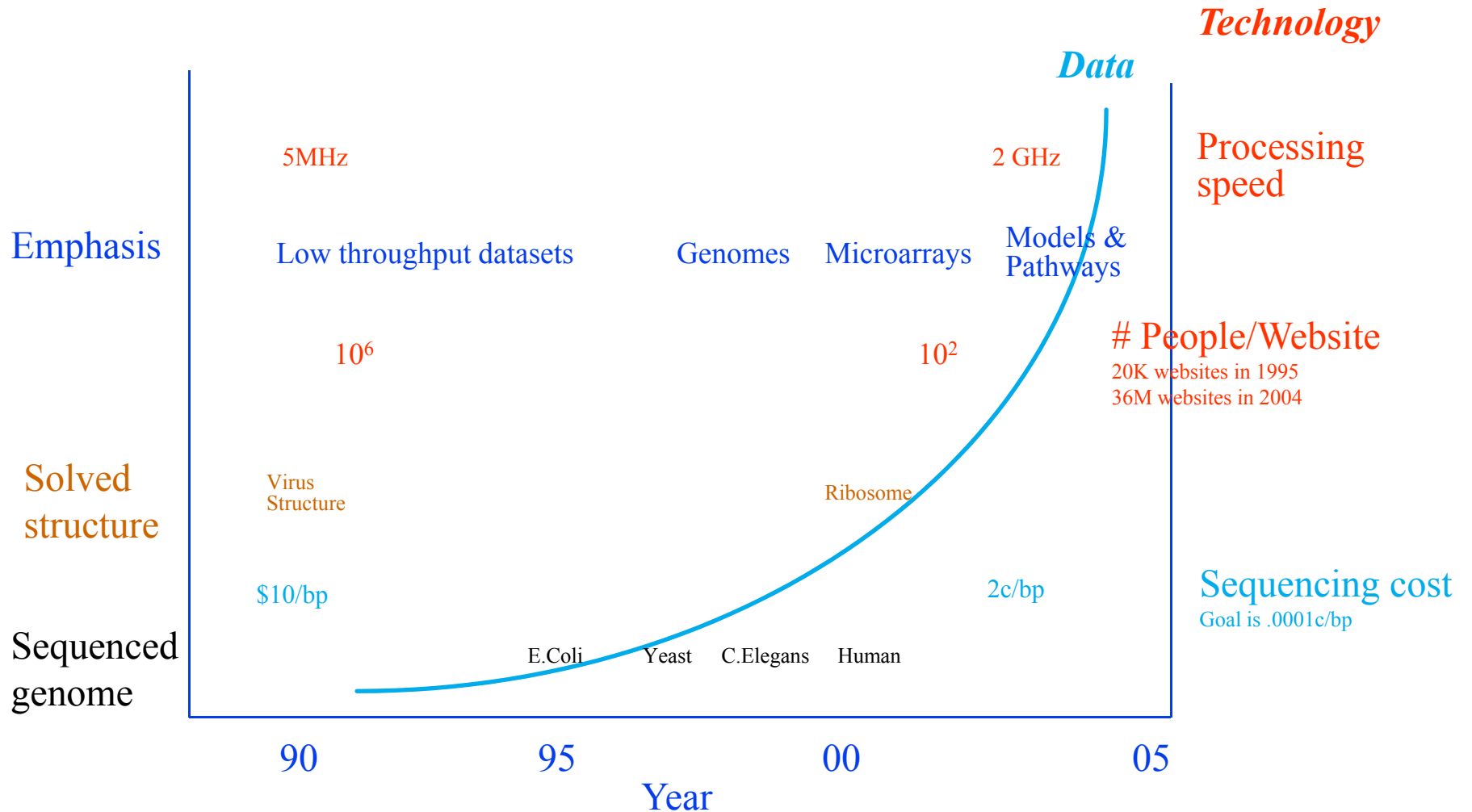
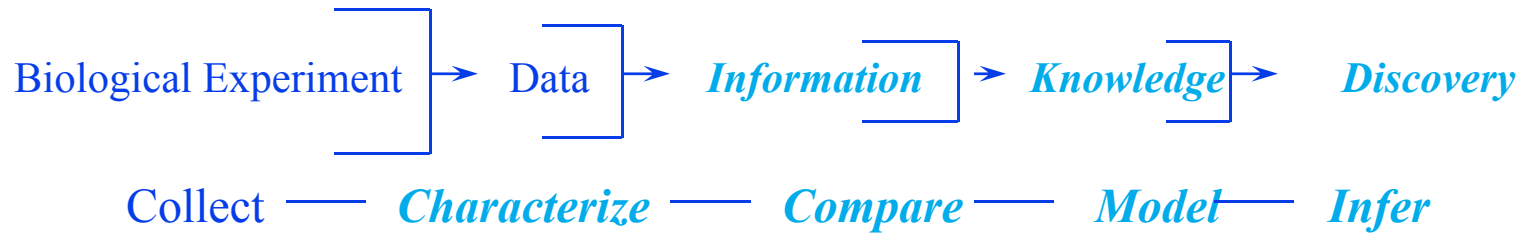
# Motivation

- Diversity and size of information
  - Sequences, 3-D structures, microarrays, protein interaction networks, *in silico* models, bio-images



- Understand the relationship
  - Similar to complex software design

# Bioinformatics - A Revolution



# Computing versus Biology

- *what computer science is to molecular biology is like what mathematics has been to physics .....*

-- Larry Hunter, ISMB'94

- *molecular biology is (becoming) an information science .....*

-- Leroy Hood, RECOMB'00

- *bioinformatics ... is the research domain focused on linking the behavior of biomolecules, biological pathways, cells, organisms, and populations to the information encoded in the genomes*

--Temple Smith, Current

Topics in Computational Molecular Biology

# Computing *versus* Biology

looking into the future

- *Like physics, where general rules and laws are taught at the start, biology will surely be presented to future generations of students as a set of basic systems ..... duplicated and adapted to a very wide range of cellular and organismic functions, following basic evolutionary principles constrained by Earth's geological history.*

--Temple Smith, Current Topics in Computational Molecular  
Biology

# Scalability challenges

- Recent issue of NAR devoted to data collections contains 719 databases
  - Sequence
    - Genomes (more than 150), ESTs, Promoters, transcription factor binding sites, repeats, ..
  - Structure
    - Domains, motifs, classifications, ..
  - Others
    - Microarrays, subcellular localization, ontologies, pathways, SNPs, ..

# Challenges of working in bioinformatics

- Need to feel comfortable in interdisciplinary area
- Depend on others for primary data
- Need to address important biological *and* computer science problems

# Skill set

- Artificial intelligence
- Machine learning
- Statistics & probability
- Algorithms
- Databases
- Programming



# Bioinformatics Topics

## Genome Sequence

- Finding Genes in Genomic DNA
  - introns
  - exons
  - promoters
- Characterizing Repeats in Genomic DNA
  - Statistics
  - Patterns
- Duplications in the Genome
  - Large scale genomic alignment

- 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
- Evolutionary Issues
  - Rates of mutation and change

# Computationally challenging problems

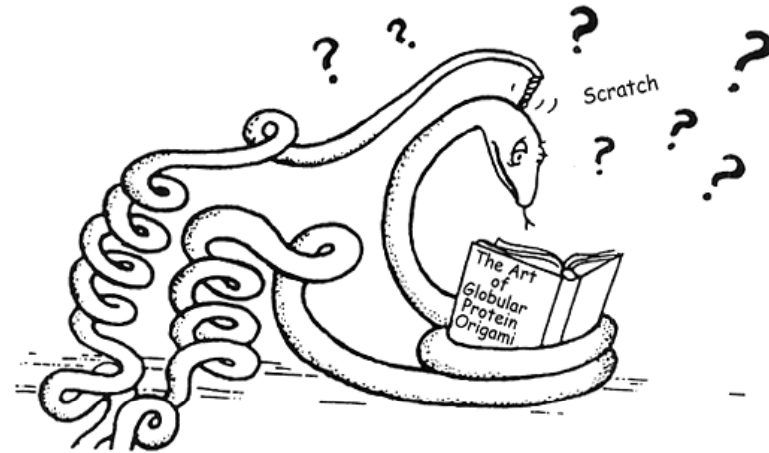
- More sensitive pairwise alignment
  - Dynamic programming is  $O(mn)$ 
    - $m$  is the length of the query
    - $n$  is the length of the database
- Scalable multiple alignment
  - Dynamic programming is exponential in number of sequences
  - Currently feasible for around 10 protein sequences of length around 1000
- Shotgun alignment
  - Current techniques will take over 200 days on a single machine to align the mouse genome

# Bioinformatics Topics

## Sequence / Structure

- Secondary Structure “Prediction”
  - via Propensities
  - Neural Networks, Genetic Alg.
  - Simple Statistics
  - TM-helix finding
  - Assessing Secondary Structure Prediction
- Structure Prediction: Protein and RNA

“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

- Tertiary Structure Prediction
  - Fold Recognition
  - Threading
  - Ab initio
- Function Prediction
  - Active site identification
- Relation of Sequence Similarity to Structural Similarity

# 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

# Computationally challenging problems

- Alignment against a database
  - Single comparison usually takes seconds.
  - Comparison against a database takes hours.
  - All-against-all comparison takes weeks.
- Multiple structure alignment and motifs
- Combined sequence and structure comparison
- Secondary and tertiary structure prediction

# Topics -- Databases

- Relational Database Concepts and how they interface with Biological Information
  - 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
- Visualization of Large Amounts of Information
- 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  $\rightarrow$  Energy  $\rightarrow$  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-Boltzmann Equation
- Lattice Models and Simplification

# General Types of “Informatics” techniques in Bioinformatics

- Databases
  - Building, querying
  - Schema design
  - Heterogeneous, distributed
- Similarity search
  - Sequence, structure
  - Significance statistics

- Finding Patterns
  - AI / Machine Learning
  - Clustering
  - Data mining
- Modeling & simulation
- Programming
  - Perl
  - Java/C/C++/..