

CENG 465

Introduction to Bioinformatics

Fall 2013-2014

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

Course web site:

http://www.ceng.metu.edu.tr/~tcan/ceng465_f1314

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)
 - Sequence alignment by Dynamic Programming
 - Statistical significance of sequence alignments
 - Suffix trees
 - Profile hidden Markov models
 - Multiple sequence alignment
- Phylogenetic trees, hierarchical clustering methods (1 week)

Course outline

- Protein structures (3 weeks)
 - Structure prediction (secondary, tertiary)
 - Structural alignment
- Microarray data analysis (2 weeks)
 - Correlations, clustering
- Gene/Protein networks, pathways (3 weeks)
 - Protein-protein, protein/DNA interactions
 - Construction and analysis of large scale networks
 - Computational systems biology

Teaching assistant

- Hilal Kilic Arslan
 - will be grading your homework assignments
- Contact info:
 - hkilic@ceng.metu.edu.tr
 - Tel: +90(312)210-5522
 - Office: B202

Grading

- Midterm exam - 40%
- Final exam - 40%
- Assignments - 20%

Online materials

- Course web site
 - http://www.ceng.metu.edu.tr/~tcan/ceng465_f1314
 - Lecture slides and reading materials
 - Assignments
- Newsgroup
 - metu.ceng.course.465
 - Accesible via cow.ceng.metu.edu.tr/News
 - Students from other departments should get a temporary ceng account from A-210 (admins) to be able to post
- COW (Ceng On the Web)
 - Students from other departments should get a temporary ceng account from A-210 (admins) to access
 - Homeworks will be submitted via cow

What is Bioinformatics?

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

- One idea for a definition?

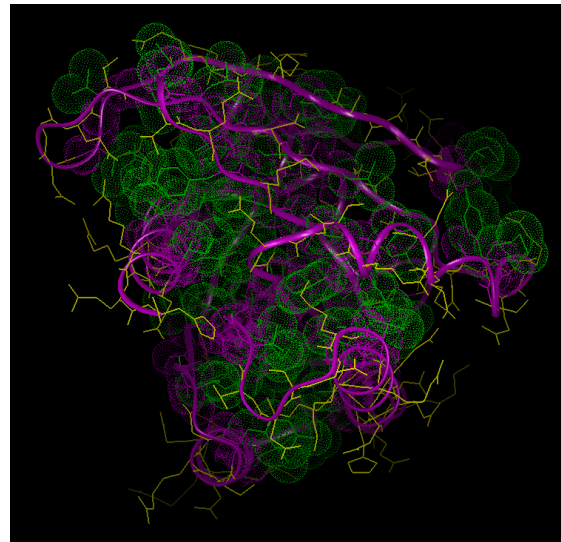
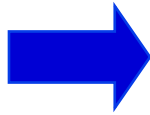
Bioinformatics is conceptualizing **biology in terms of molecules** (in the sense of physical-chemistry) and then applying **“informatics” techniques** (derived from disciplines such as applied math, CS, and statistics) to understand and **organize the information associated** with these molecules, **on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

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

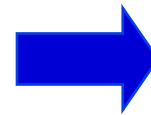
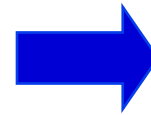
Introductory Biology



DNA
(Genotype)

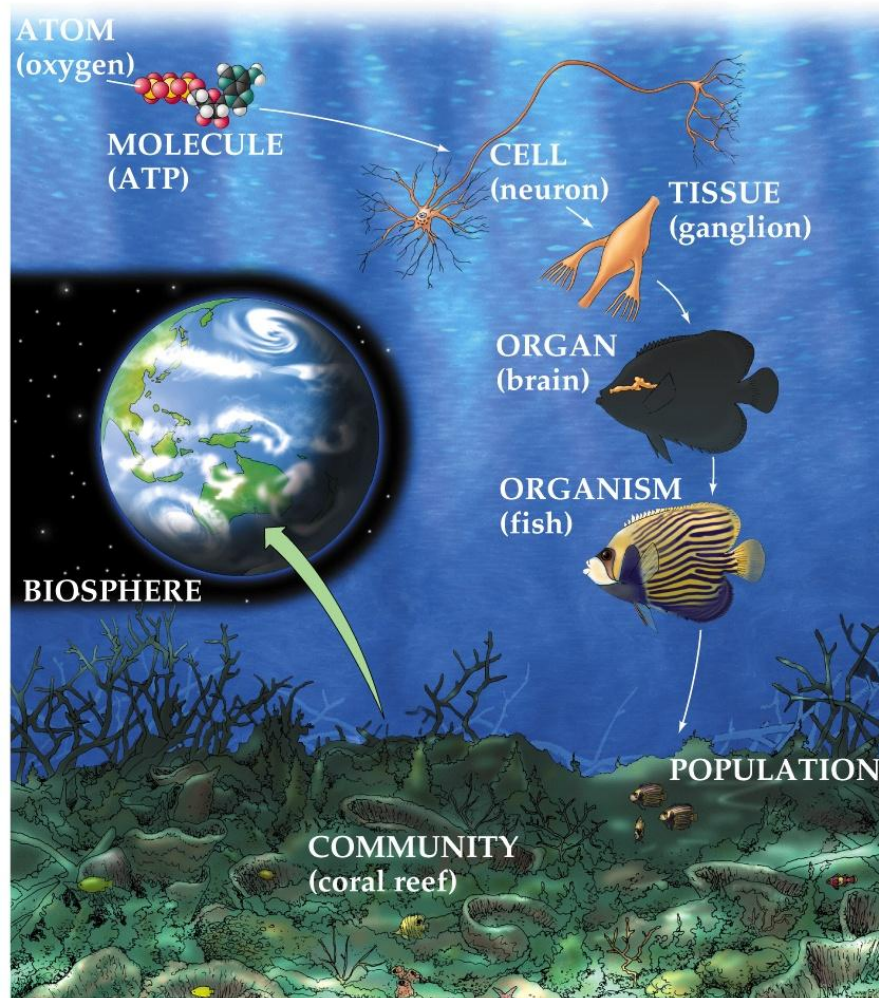


Protein



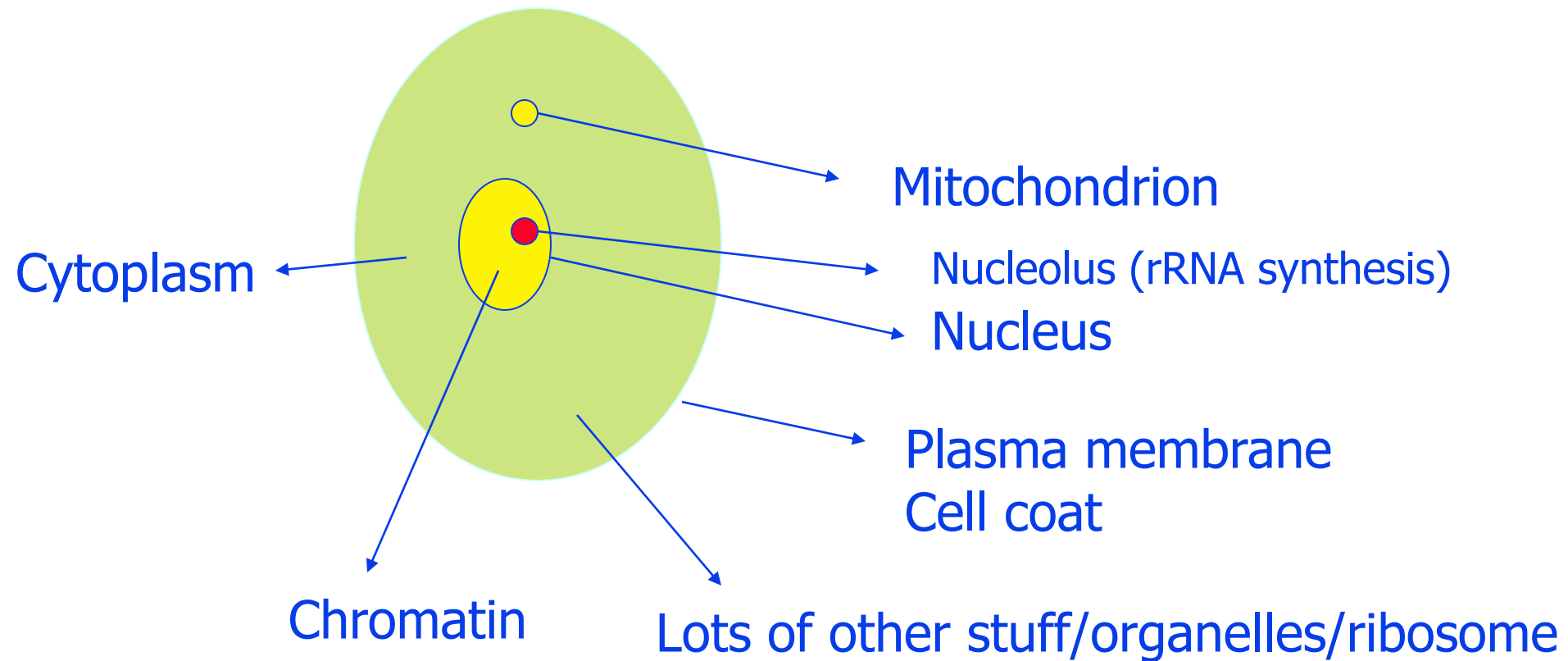
Phenotype

Scales of life

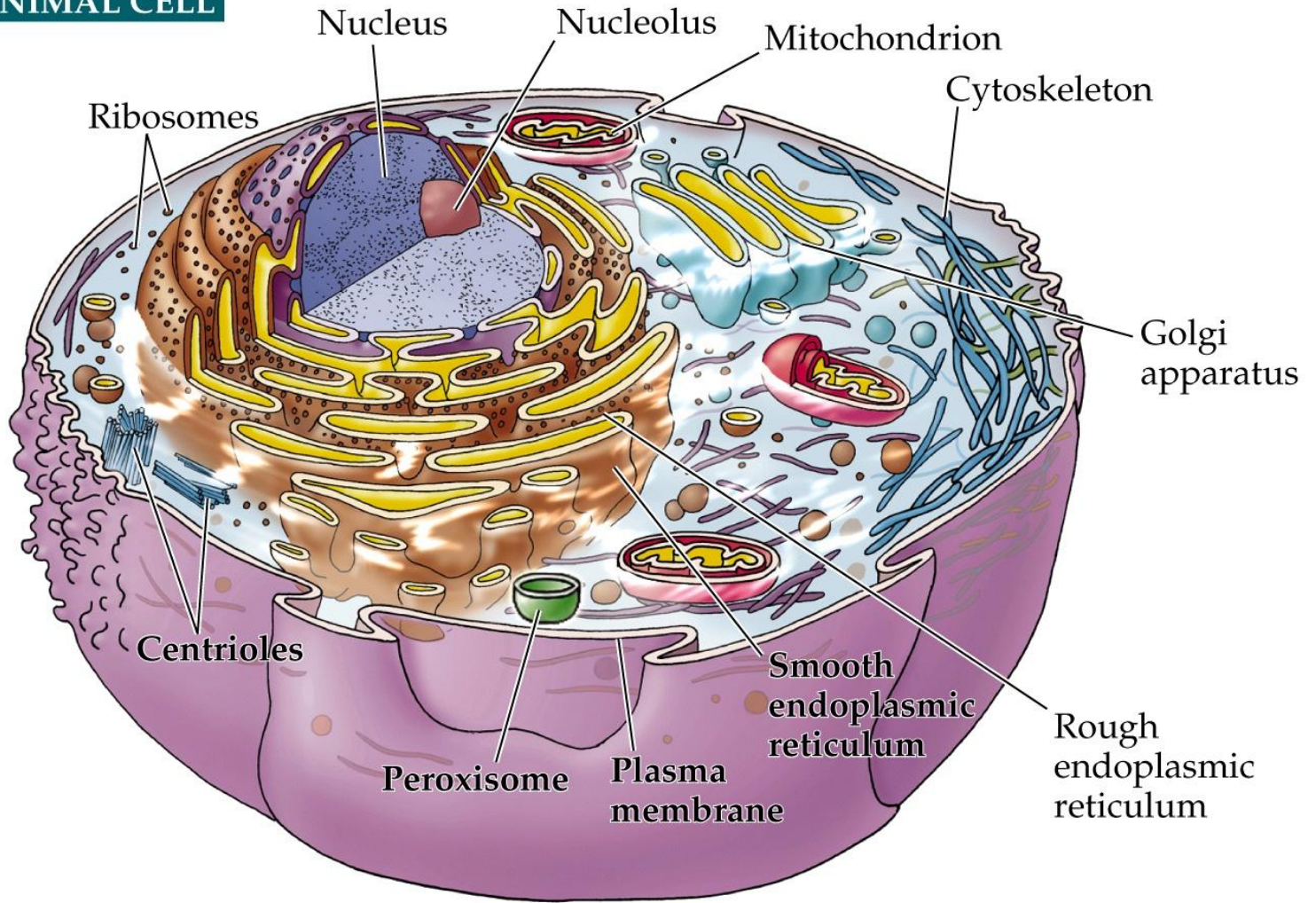


© 2001 Sinauer Associates, Inc.

Animal Cell



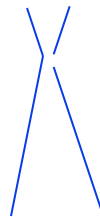
AN ANIMAL CELL



© 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

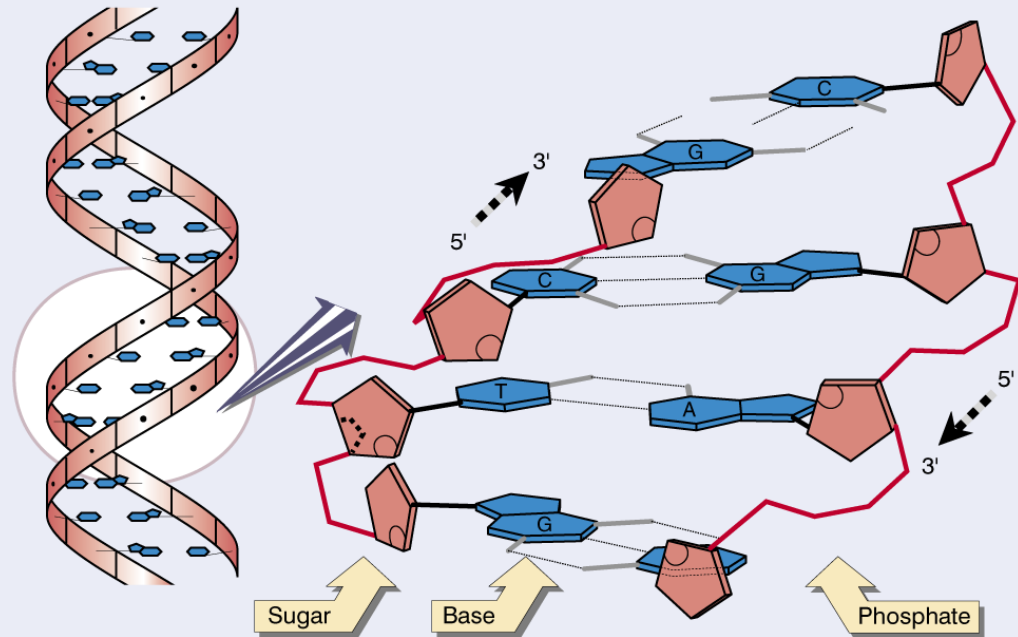
- Coding or Not?
- Parse into genes?
- 4 bases: AGCT
- ~1 Kb in a gene, ~2 Mb in genome
- ~3 Gb Human

```
atggcaattaaaattggtatcaatggTTTTGGTCGTATCGGCCGTATCGTATTCGGTGCA
gcacaacaccgtgatgacattgaagttgtaggtattaacgacttaatcgacggtgaatac
atggcTTATATGTTGAAATATGATTCAACTCACGGTCGTTTCGACGGCACTGTTGAAGTG
aaagatggtaacttagtggTTAATGGTAAAACATCCGTGTAACGTCAGAACGTGATCCA
gCAAACTTAACTGGGGTGCAATCGGTGTTGATATCGCTGTTGAAGCGACTGGTTTATTC
TTAACTGATGAACTGCTCGTAAACATATCACTGCAGGCGCAAAAAAGTTGTATTAAC
GGCCCATCTAAAGATGCAACCCCTATGTTTCGTTTCGTGGTGAAACTTCAACGCATACGCA
GGTCAAGATATCGTTTCTAACGCATCTTGTACAACAAACTGTTTAGCTCCTTTAGCACGT
GTTGTTTCATGAACTTTTCGGTATCAAAGATGGTTTAAATGACCACTGTTCAACGCAACGACT
GCAACTCAAAAACTGTGGATGGTCCATCAGCTAAAGACTGGCGCGGGCGGCGGGTGCA
TCACAAAACATCATTCCATCTTCAACAGGTGCAGCGAAAGCAGTAGGTAAGTATTACCT
GCATTAACGGTAAATTAACGTTATGGCTTTCCGTGTTCCAACGCCAAACGTATCTGTT
GTTGATTTAACAGTTAATCTTGAAAACCAGCTTCTTATGATGCAATCAAACAAGCAATC
AAAGATGCAGCGGAAGGTA AACGTTCAATGGCGAATTA AAAGGCGTATTAGGTTACT
GAAGATGCTGTTGTTTCTACTGACTTCAACGGTTGTGCTTTAACTTCTGTATTTGATGCA
GACGCTGGTATCGCATTAACTGATTCTTTTCGTTAAATTGGTATC . . .
```

```
. . . caaaaatagggttaatatgaatctcgatctccatTTTGTTCATCGTATTCAA
caacaagccaaaactcgtacaaatatgaccgcactTCGCTATAAAGAACACGGCTTGTGG
CGAGATATCTCTTGGAAAACTTCAAGAGCAACTCAATCAACTTCTCGAGCATTGCTT
GCTCACAATATTGACGTACAAGATAAAATCGCCATTTTGGCCATAATATGGAACGTTGG
GTTGTTTCATGAACTTTTCGGTATCAAAGATGGTTTAAATGACCACTGTTCAACGCAACGACT
ACAATCGTTGACATTGCGACCTTACAAATTCGAGCAATCACAGTGCCTATTTACGCAACC
AATACAGCCCAGCAAGCAGAATTTATCCTAAATCACGCCGATGTA AAAATTCTCTTCGTC
GGCGATCAAGAGCAATACGATCAAACATTGGAAATTGCTCATCATTGTCCAAAATTACAA
AAAATTGTAGCAATGAAATCCACCATTCAATTACAACAAGATCCTCTTTCTTGCCTTGG
```

DNA structure

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



Molecular Biology Information: Protein Sequence

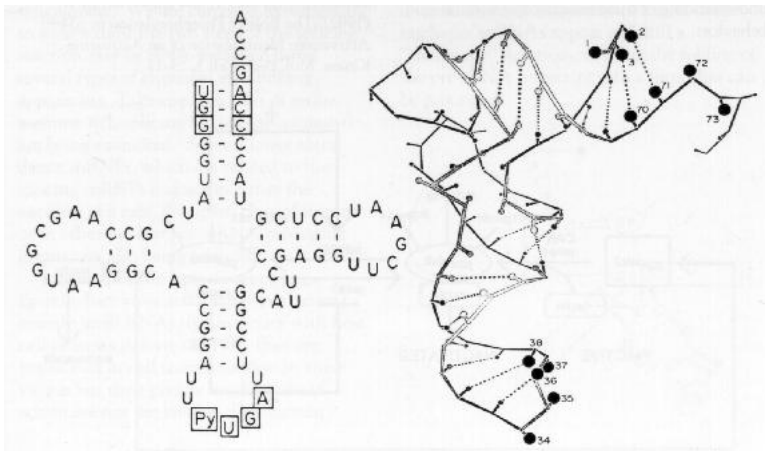
- 20 letter alphabet
 - ACDEFGHIKLMNPQRSTVWY but not BJOUXZ
- Strings of ~300 aa in an average protein (in bacteria),
 - ~200 aa in a domain
- ~1M known protein sequences

```
d1dhfa_ LNCIVAVSQNMGIGKNGDLPWPPPLRNEFRYFQRMTTTSSVEGKQ-NLVIMGKKTWFSI
d8dfr_  LNSIVAVCQNMGIGKDG NLPWPPPLRNEYKYFQRMSTSHVEGKQ-NAVIMGKKTWFSI
d4dfra_ ISLIAALAVDRVIGMENAMPWN-LPADLAWFKRNTL-----NKPVIMGRHTWESI
d3dfr_  TAFLWAQDRDGLIGKDGHLPW-HLPDDLHYFRAQTV-----GKIMVGRRTYESF
```

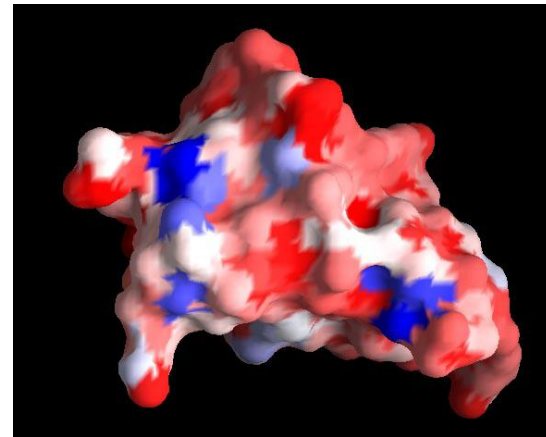
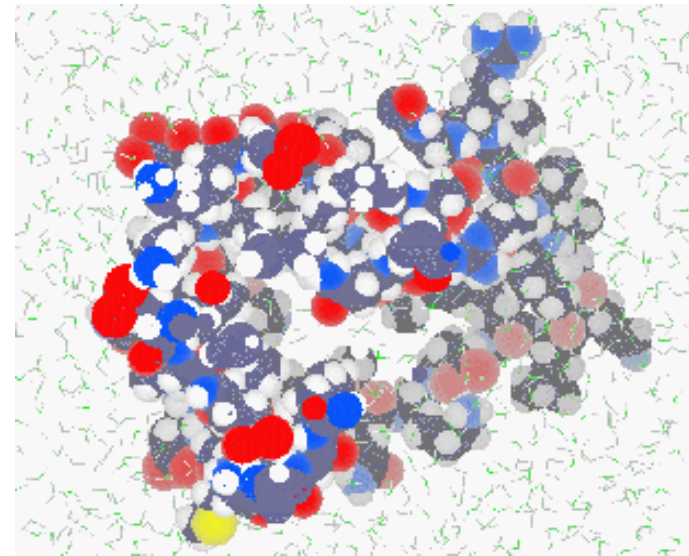
```
d1dhfa_ LNCIVAVSQNMGIGKNGDLPWPPPLRNEFRYFQRMTTTSSVEGKQ-NLVIMGKKTWFSI
d8dfr_  LNSIVAVCQNMGIGKDG NLPWPPPLRNEYKYFQRMSTSHVEGKQ-NAVIMGKKTWFSI
d4dfra_ ISLIAALAVDRVIGMENAMPWN-NLPADLAWFKRNTLD-----KPVIMGRHTWESI
d3dfr_  TAFLWAQDRNGLIGKDGHLPW-HLPDDLHYFRAQTVG-----KIMVGRRTYESF
```

Molecular Biology Information: Macromolecular Structure

- DNA/RNA/Protein
 - Almost all protein

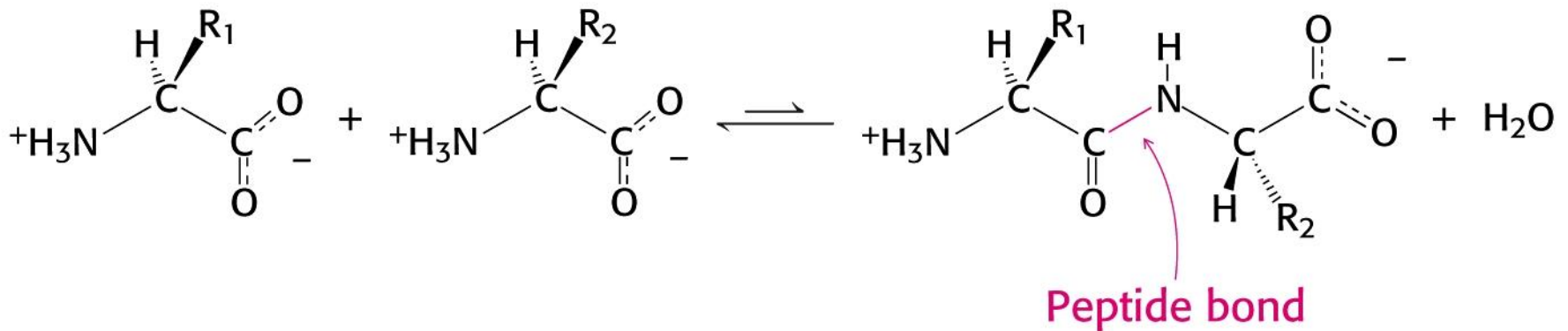


'Identity elements' in *Escherichia coli* glutamine tRNA.



More on Macromolecular Structure

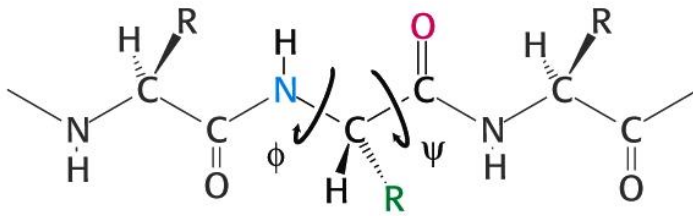
- Primary structure of proteins
 - Linear polymers linked by peptide bonds
 - Sense of direction



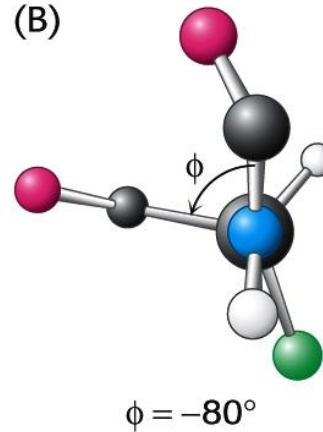
Secondary Structure

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

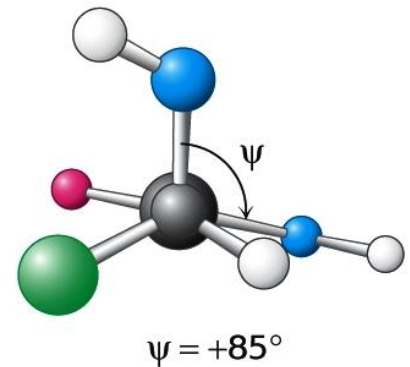
(A)



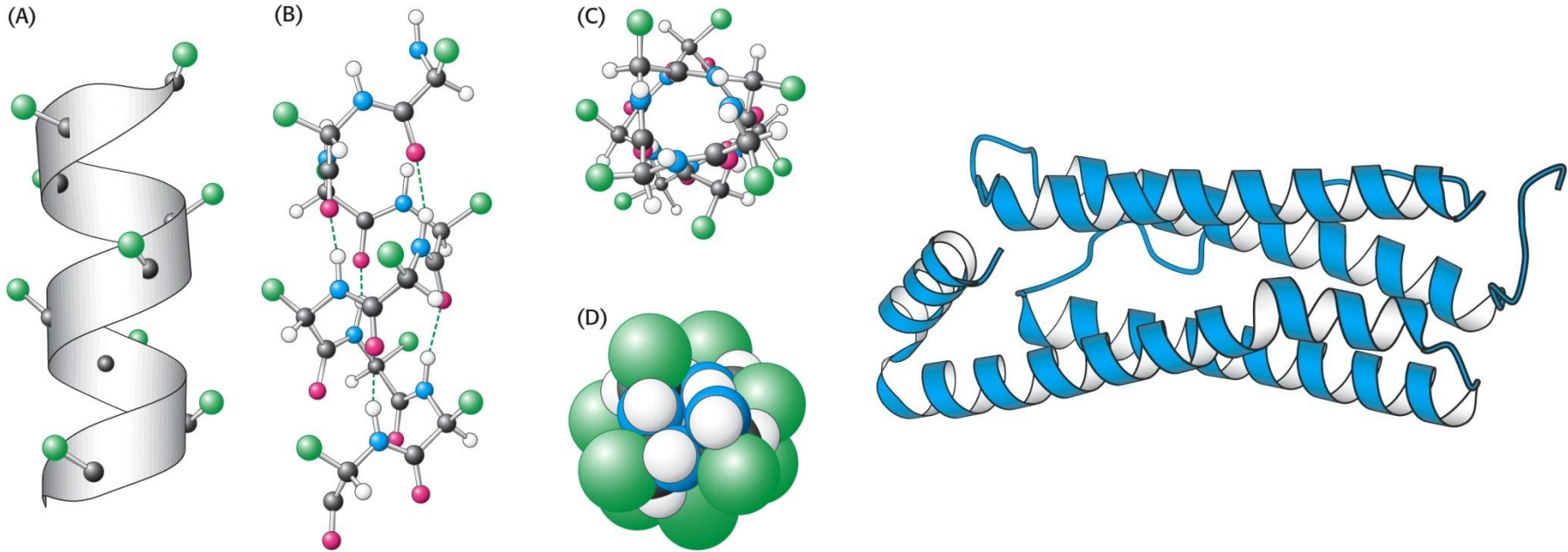
(B)



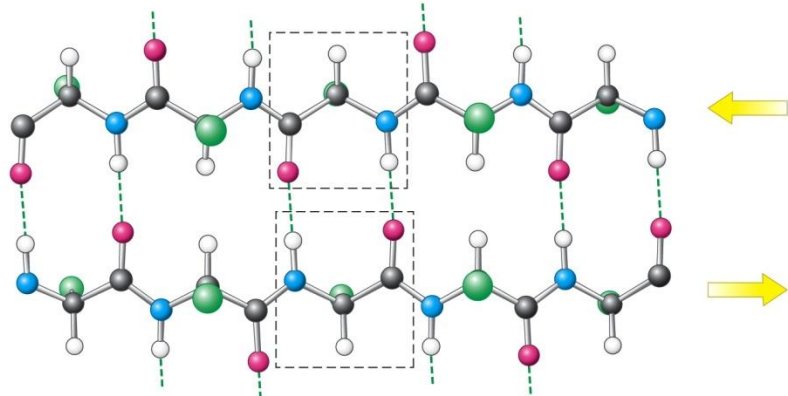
(C)



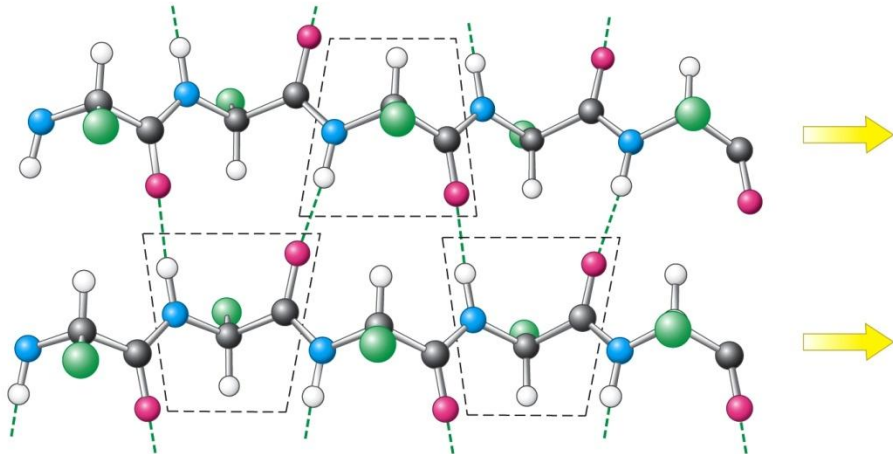
Alpha helix



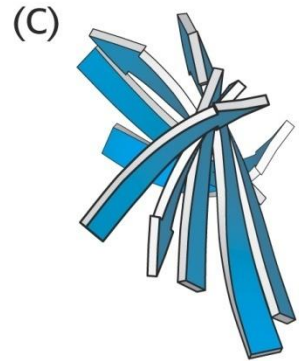
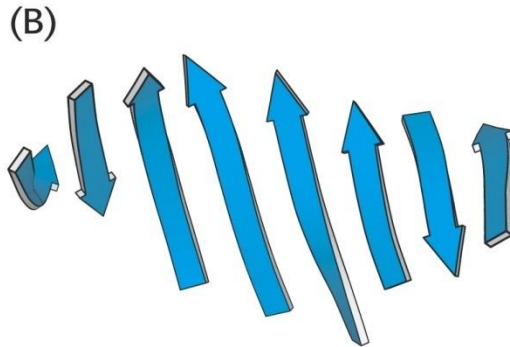
Beta sheet



anti-parallel



parallel

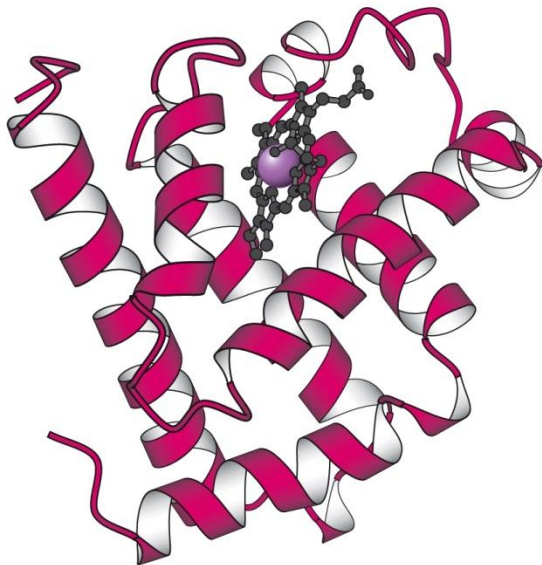


schematic

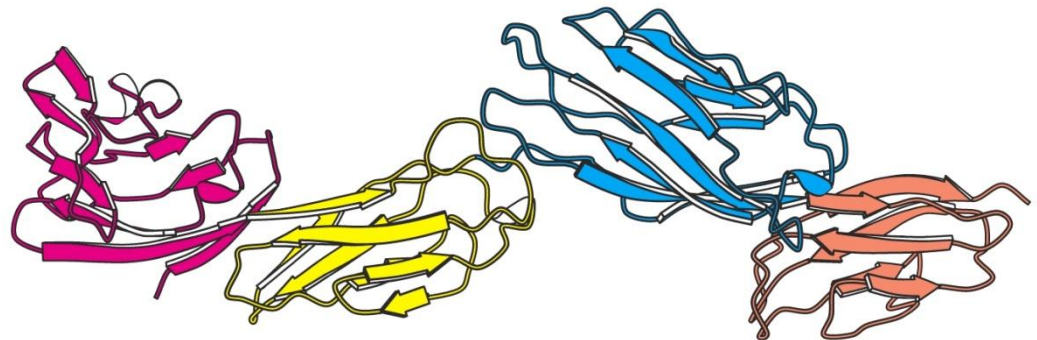
Tertiary Structure

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

(B)



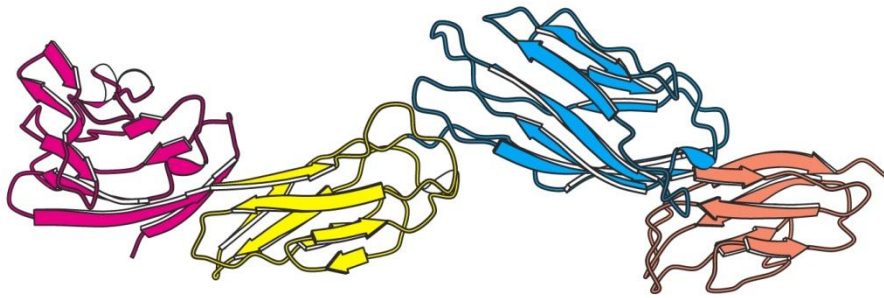
tertiary structure of
myoglobin



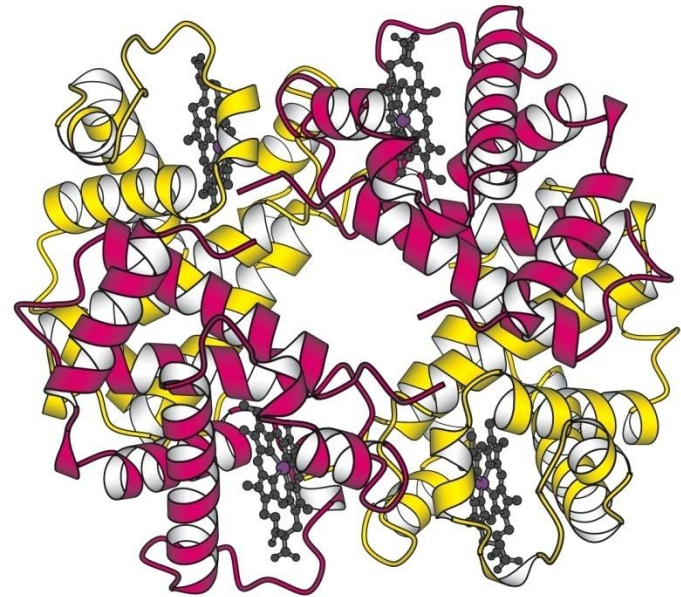
domains of CD4

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

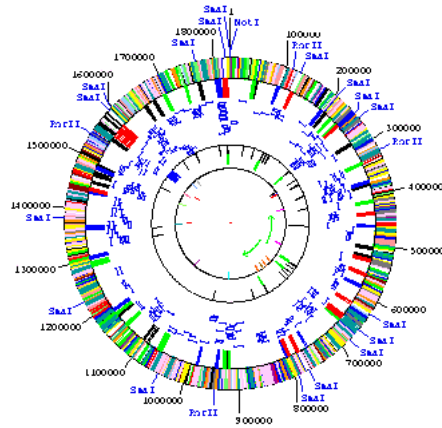
- One gene encodes one* protein.
- Like a program, it starts with start codon (e.g. ATG), then each three code one amino acid. Then a stop codon (e.g. TGA) signifies end of the gene.
- Sometimes, in the middle of a (eukaryotic) gene, there are introns that are spliced out (as junk) during transcription.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

A.A. Coding Table

Glycine (GLY)	GG*
Alanine(ALA)	GC*
Valine (VAL)	GT*
Leucine (LEU)	CT*
Isoleucine (ILE)	AT(*-G)
Serine (SER)	AGT, AGC
Threonine (THR)	AC*
Aspartic Acid (ASP)	GAT,GAC
Glutamic Acid(GLU)	GAA,GAG
Lysine (LYS)	AAA, AAG
Start: ATG, CTG, GTG	

Arginine (ARG)	CG*
Asparagine (ASN)	AAT, AAC
Glutamine (GLN)	CAA, CAG
Cysteine (CYS)	TGT, TGC
Methionine (MET)	ATG
Phenylalanine (PHE)	TTT,TTC
Tyrosine (TYR)	TAT, TAC
Tryptophan (TRP)	TGG
Histidine (HIS)	CAT, CAC
Proline (PRO)	CC*
Stop	TGA, TAA, TAG

Molecular Biology Information: Whole Genomes

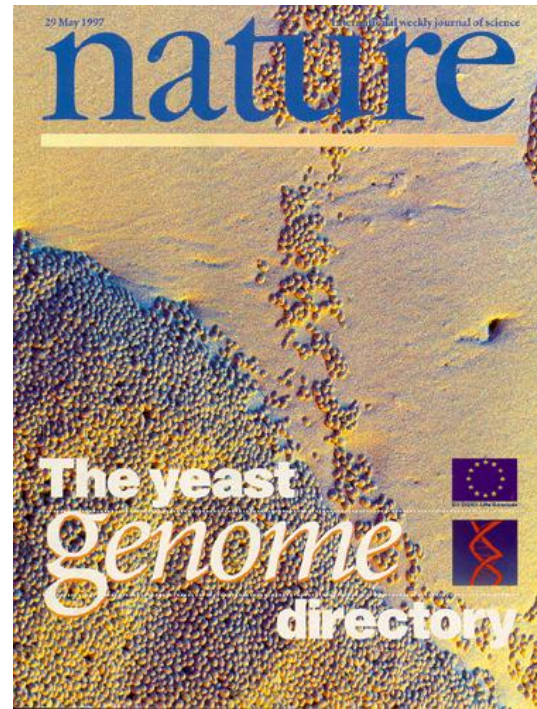


Genome sequences now accumulate so quickly that, in less than a week, a single laboratory can produce more bits of data than Shakespeare managed in a lifetime, although the latter make better reading.

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

1995

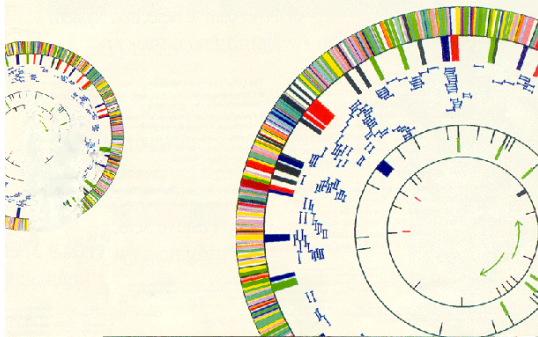
Bacteria,
1.6 Mb,
~1600 genes
[*Science* 269: 496]



Genomes highlight the Finiteness of the “Parts” in Biology

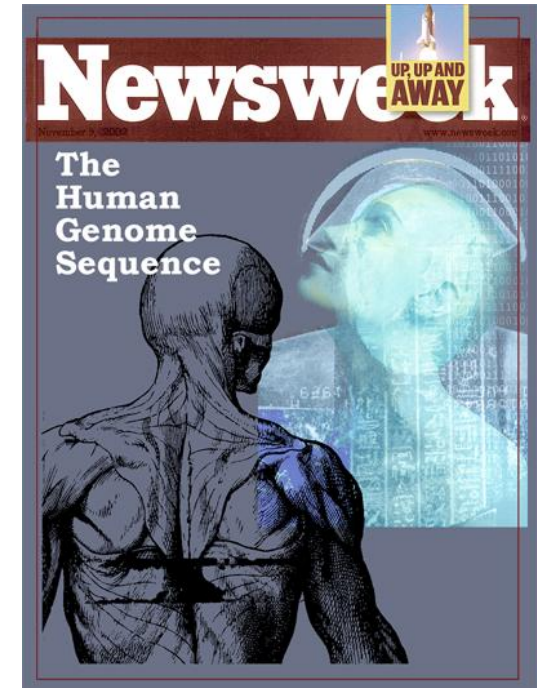
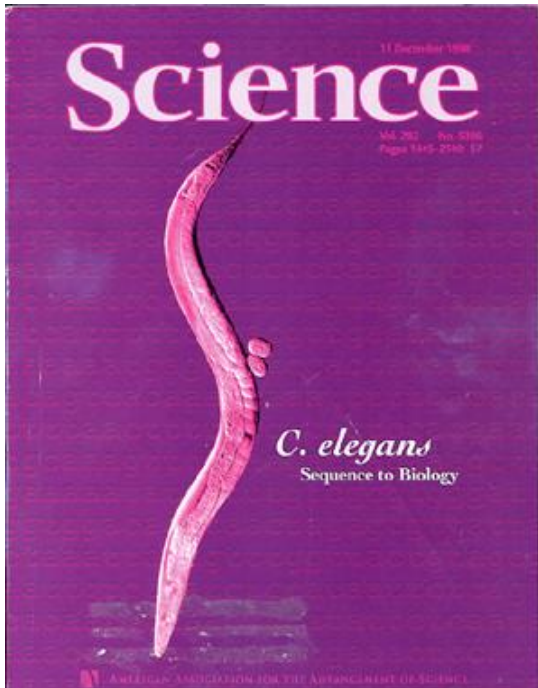
1997

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



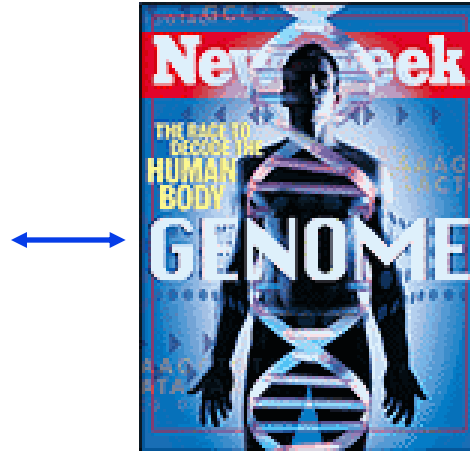
1998

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***

Dissecting the Regulatory Circuitry of a Eukaryotic Genome

Frank C. P. Holstege,* Ezra G. Jennings,*¹ John J. Wyrick,*¹ Tong Ihn Lee,*¹ Christopher J. Hengartner,*¹ Michael R. Green,¹ Todd R. Golub,*⁵ Eric S. Lander,*¹ and Richard A. Young*^{1||}
¹Whitehead Institute for Biomedical Research, Cambridge, Massachusetts 02142
²Department of Biology, Massachusetts Institute of Technology, Cambridge, Massachusetts 02139
³Howard Hughes Medical Institute, Program in Molecular Medicine, University of Massachusetts Medical Center, Worcester, Massachusetts 01605
⁴Dana-Farber Cancer Institute and Harvard Medical School, Boston, Massachusetts 02115

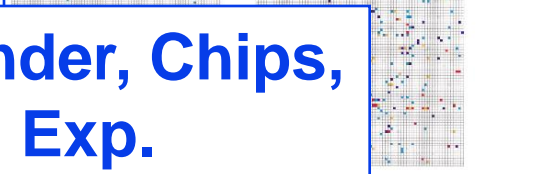
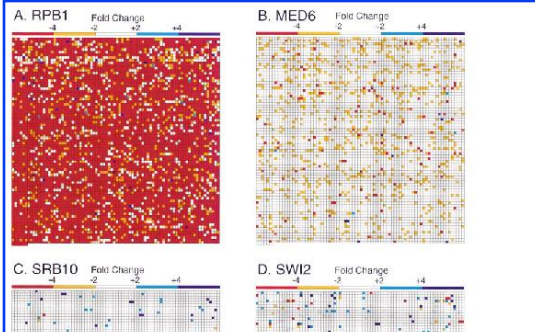


Fig. 3. Schematic representation of the mTn construct and the derived HAT tag elements. Each mTn contains the coding regions for the tetracycline efflux (Tet) and Ura3 proteins, and the mTn element from Tn3. mTn-3HA/lacZ and mTn-3HA/lacZ contain a truncated lacZ gene. mTn-3HA/GFP contains the coding region for GFP mutant (G111). In each case, these are flanked by the sites and Tn3 terminal repeats (TR). Between lacp and the right TR, each transposon contains a sequence encoding three tandem copies of the HA epitope. Between the left TR and lacZ is a sequence encoding either an additional copy of the HA epitope (mTn-6HA/lacZ) or the factor Xa protease cleavage site (mTn-3HA/lacZ, mTn-3HA/GFP). Exposure of these transposons to Cre recombinase catalyzes the formation of a smaller element encoding the HAT tag, shown in the right. The site of the lac recombinase insertion is indicated by a triangle. (Not drawn to scale.)

Specific transcription factors, a novel mechanism for the regulation of specific sets of genes
 Figure 2. Genome-wide Expression Data for Selected Components of the RNA Polymerase II Holoenzyme
 Changes 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, then III, etc., until the last gene on the right arm of chromosome XVI is reached. The results are shown for (A) Rpb1, (B) Med6, (C) Srb10, and (D) Swi2.

use II with that obtained by its inactivation. Comparison of the two data sets reveals that the transcriptional program in the *trb1* mutant (see Technology, Protocol 1) is similar to that of the wild-type, but that some other gene's altered mRNA levels. The 506 genes we have identified that require Med6 function to the same extent as Rpb1 function are those at which promoter-associated transcriptional regulators are most abundant through interactions with Med6. The function of the Srb/mediator complex also not known (Thompson et al., 1994; Koleske and Young, 1994; Henry-Meyers et al., 1998). To determine dependence of gene expression on Srb10, we compared the wild-type expression of an *SRB5* gene and its wild-type expression (see the web site for details). The results indicate that 16% of all function for their expression. With rain and other constitutive mutants.



Fig. 3. Map showing amino acid positions of HAT tag insertions in the yeast proteins Spz2, Arp100p, and Ser1p. Regions highlighted are the amino acid positions of the HAT tag insertions. The HAT tag is a 262-bp DNA element containing a sequence encoding three tandem copies of the HA epitope. For the HAT tag, the amino acid positions are indicated by black and grey circles.

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 transposon creates insertion mutations in the target gene, allowing phenotypic analysis. The transposon can be reduced by *cre-loxP*-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 *SP2*, *ARP100*, *SER1*, and *BDF1* genes.

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). Determining

antibody into a protein of interest, the time and expense of generating specific antibodies and associated reagents is avoided

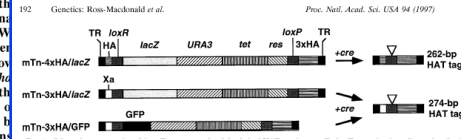


Fig. 1. Schematic representation of the mTn construct and the derived HAT tag elements. Each mTn contains the coding regions for the tetracycline efflux (Tet) and Ura3 proteins, and the mTn element from Tn3. mTn-3HA/lacZ and mTn-3HA/lacZ contain a truncated lacZ gene. mTn-3HA/GFP contains the coding region for GFP mutant (G111). In each case, these are flanked by the sites and Tn3 terminal repeats (TR). Between lacp and the right TR, each transposon contains a sequence encoding three tandem copies of the HA epitope. Between the left TR and lacZ is a sequence encoding either an additional copy of the HA epitope (mTn-6HA/lacZ) or the factor Xa protease cleavage site (mTn-3HA/lacZ, mTn-3HA/GFP). Exposure of these transposons to Cre recombinase catalyzes the formation of a smaller element encoding the HAT tag, shown in the right. The site of the lac recombinase insertion is indicated by a triangle. (Not drawn to scale.)

transposon, and contain the Tn3 *res* site for resolution of transposon conjugates. Tn3-encoded enzymes catalyze transposon transposition and resolution are provided to mTn. All three transposons contain the *URA3* and *tet* genes for selection in *S. cerevisiae*. *E. coli* respectively. Transposon mTn-3HA/lacZ and mTn-6HA/lacZ contain a *lacZ* gene that lacks an initiator methionine, while transposon mTn-3HA/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 use 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 (e.g., refs. 4, 23, and 26).

A *loxP* element lies at one end of the transposon and a *lacp* 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 *lox* sites are internal to sequences encoding multiple copies of an epitope from the influenza virus hemagglutinin protein (the HA epitope; ref. 25). The mTn-3HA transposon also contains a factor Xa protease cleavage site (19) in the region external to the *loxP* site. Expression of the Cre recombinase induces recombination between the *lox* 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-3HA) or 262-bp (mTn-6HA) element. This element consists of a single *loxP* site and sequences encoding three or four copies of the HA epitope, flanked by the Tn3 terminal repeats (Fig. 1). The mTn-3HA-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-3HA) or 89 amino acids (mTn-6HA) into the protein. We designate these insertions HAT tags.

Mutagenesis of Yeast Genes. Transposon mTn-3HA/lacZ and mTn-6HA/lacZ were tested by mutagenesis of the yeast *SP2* gene. *SP2* encodes a nonessential protein that localizes to sites of polarized growth; *spz2* mutants exhibit defects in

Fig. 3. Map showing amino acid positions of HAT tag insertions in the yeast proteins Spz2, Arp100p, and Ser1p. Regions highlighted are the amino acid positions of the HAT tag insertions. The HAT tag is a 262-bp DNA element containing a sequence encoding three tandem copies of the HA epitope. For the HAT tag, the amino acid positions are indicated by black and grey circles.

Young/Lander, Chips, Abs. Exp.

The Brown Lab
 Stanford University Department of Biochemistry

The MGuide
 The Complete Guide to MicroArrays
 Build your own arrayer and scanner

The transcriptional program in the response of human fibroblasts to serum
 The web supplement to Iyer V.R. et al. (1999) Science 283:83-87

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 transcriptome

Brown, μ array, Rel. Exp. over Timecourse

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

Snyder, Transposons, Protein Exp.

Array Data

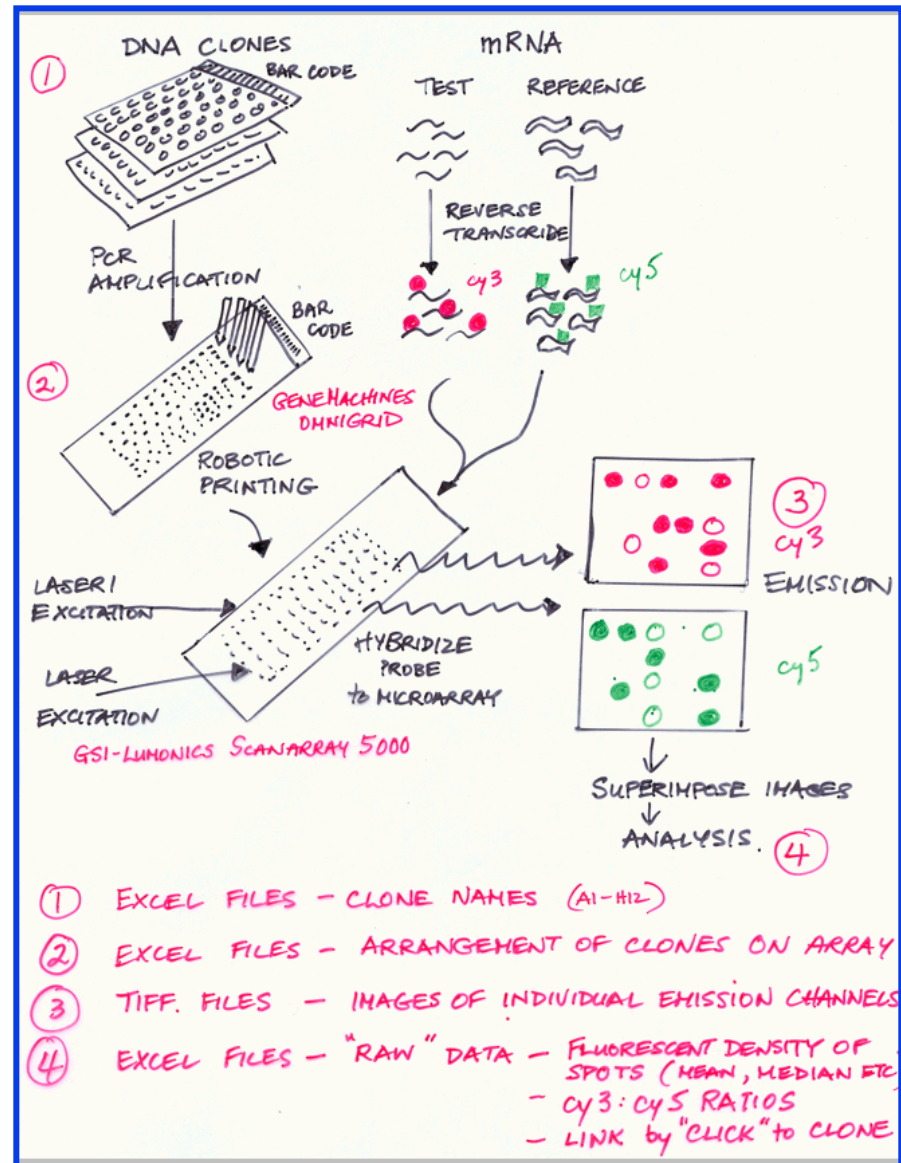
Yeast Expression Data in Academia:

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 x 10 = 60K floats

telling signal from background



(courtesy of J Hager)

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

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

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 coordinate deletion strains for all annotated other essential genes (66% of those within 5 kb of another essential gene) whereas 47% of nonessential essential genes were generally 50 kb of the telomeres (Fig. 1). In addition, the *gpr1* (0.78, M; 0.99, R), *del1* (0.83, R; 0.98, M), *overlaps ribosomal protein rpl26a*; *enc1* (0.83, R; 0.97, M) and *ym013w* showed a minimal medium-specific growth defect (15). *GTP1* (*YOR070C*) is a GTPase. Altogether, almost 40% of the deletants



The functions of many open reading sequencing projects are unknown. Now, to systematically determine their *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

the phenotypic analysis strains by allowing the growth to be assayed simultaneously. 558 homozygous deletion were pooled (12) and minimal media for about During this time, aliquots from the two pools. The tags and hybridized to high-density the hybridization data were the relative growth rates for mutant in the population (14) that the growth rate for each independently with the UP-NTAG signals would agree, such both the UPFAG and

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

Other Whole-Genome Experiments



Gene 215 (1998) 143-152

GENE
AN INTERNATIONAL JOURNAL ON GENES AND GENOMES

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

Shao-bing Hua 1*, Ying Luo1,2, Mengsheng Qiu1,3, Eva Chan 2, Helen Zhou 4, Li Zhu

GeneNet Group, CLONTECH Laboratories Inc., 1020 East Meadow Circle, Palo Alto, CA 94303, USA

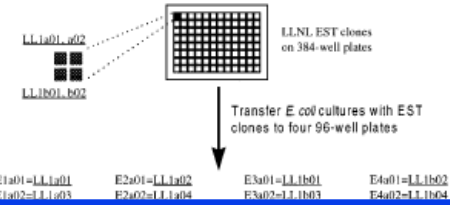
Received 1 February 1998; received in revised form 28 April 1998; accepted 29 April 1998; Received by E.Y. Chen

Abstract

Identification of all human protein-protein interactions 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 so first time that a comprehensive two-hybrid EST clones. © 1998 Elsevier Science

148

S. b. Hua et al. / Gene 215 (1998) 143-152



Keywords: Functional genomics

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 a minority of these UniGenes.

2 hybrids, linkage maps

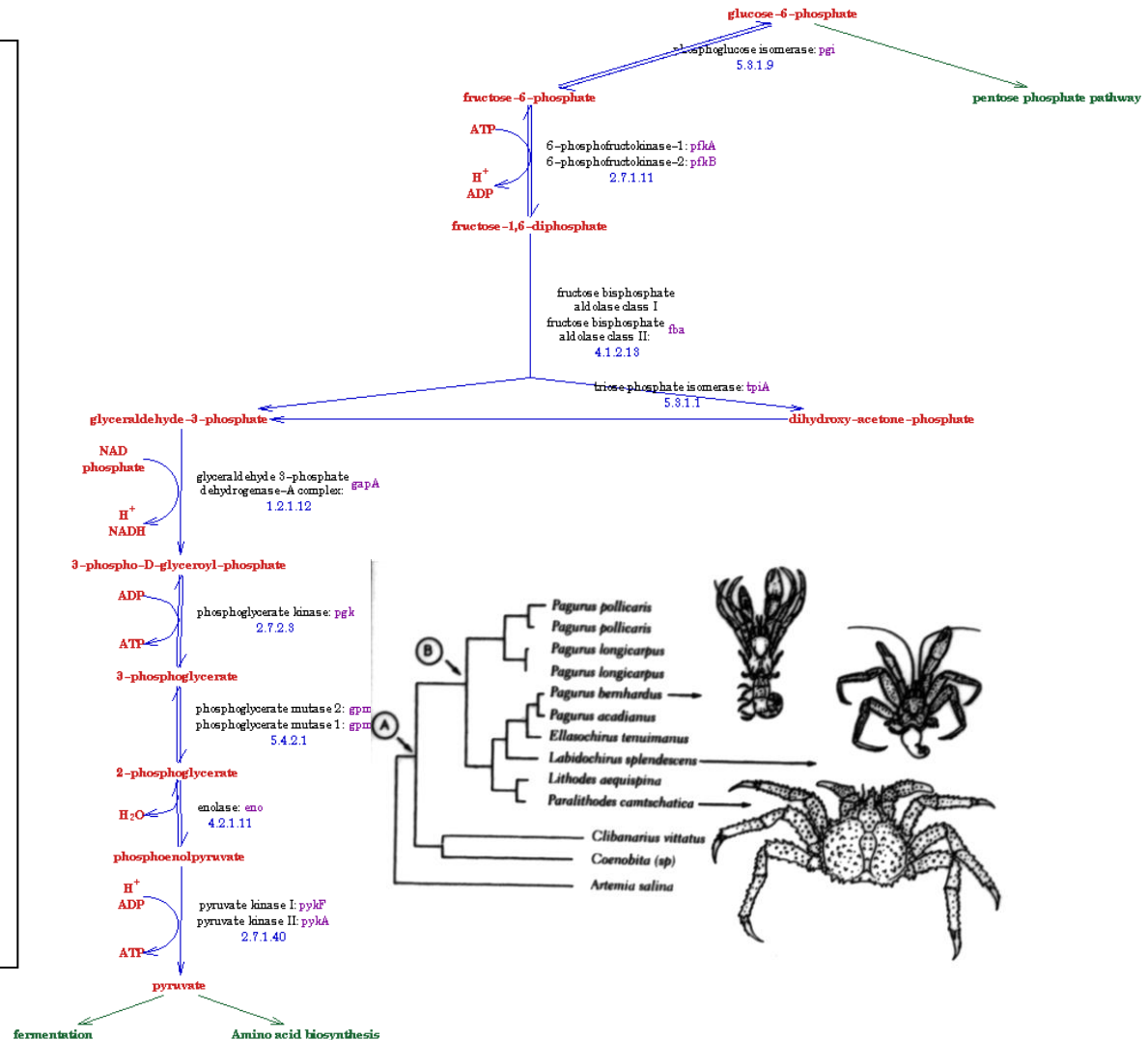
Hua, S. B., Luo, Y., Qiu, M., Chan, E., Zhou, H. & Zhu, L. (1998). Construction of a modular yeast two-hybrid cDNA library from human EST clones for the human genome protein linkage map. *Gene* **215**, 143-52

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

* Corresponding author. Tel: +1 650 927 2200; e-mail: sbhua@clontech.com
1 These authors contributed equally to this work.
2 Present address: Rigel, Inc., 94086 USA.
3 Present address: Department of Neurobiology, School of Medicine, University of Kentucky, KY 40292, USA.

Molecular Biology Information: Other Integrative Data

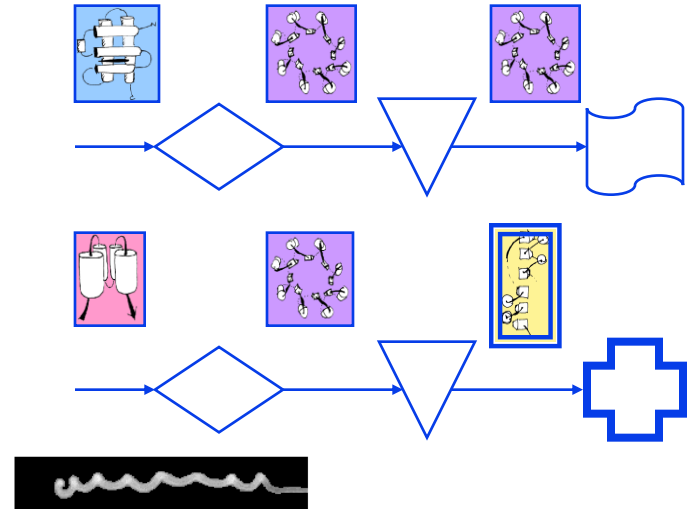
- Information to understand genomes
 - Metabolic Pathways (glycolysis), traditional biochemistry
 - Regulatory Networks
 - Whole Organisms Phylogeny, traditional zoology
 - Environments, Habitats, ecology
 - The Literature (MEDLINE)
- The Future....



Organizing

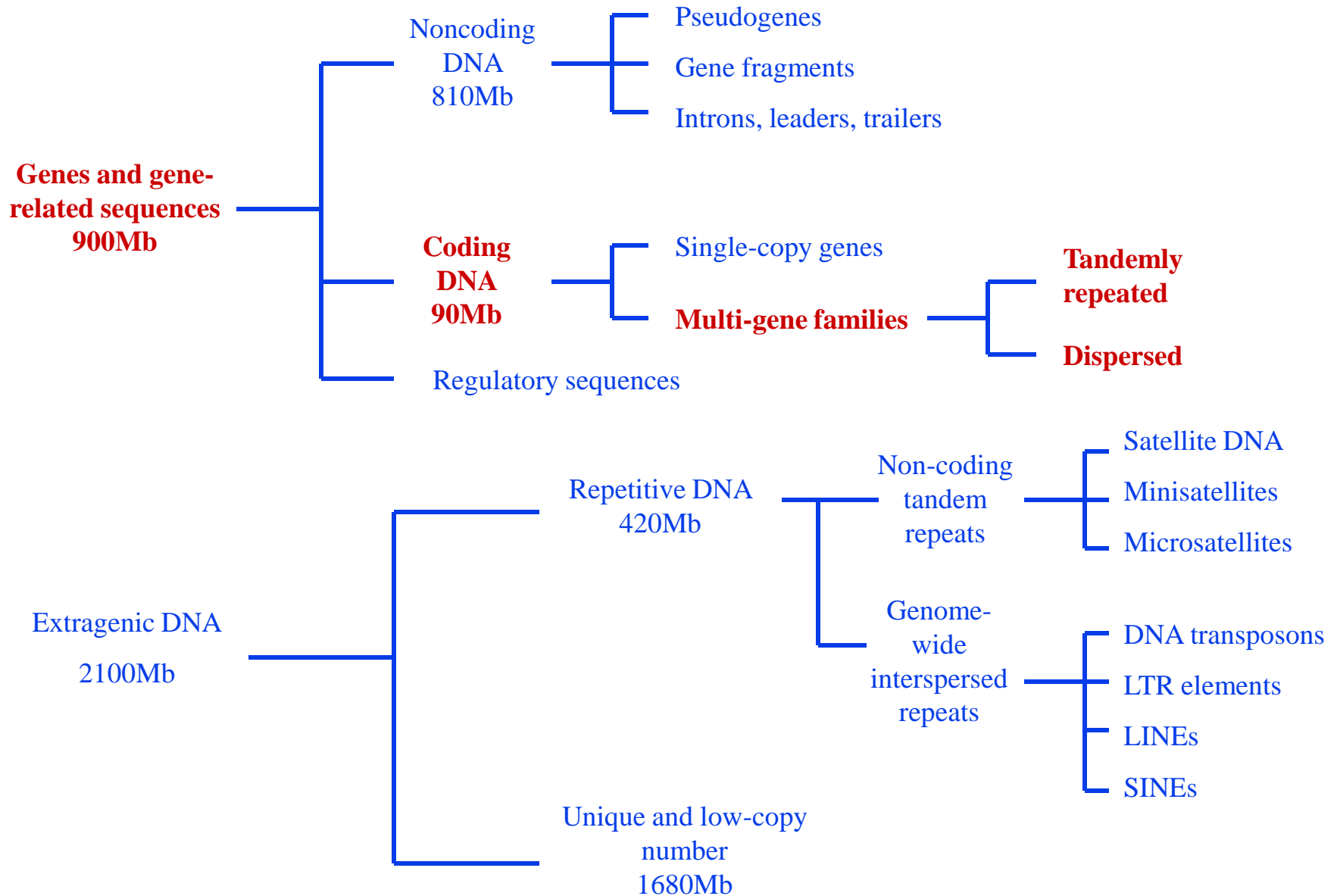
Molecular Biology Information: Redundancy and Multiplicity

- 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

Life Science Data

Biological Data:
-Genes
-Proteins
-Gene and protein function and interaction

Clinical and Field Trials Data

Scientific Literature:
-Journal articles

Other Disciplines:
-Chemical data

Bioinformatics

Information Technology

Automated Techniques:
-DNA sequencing
-DNA microarrays
-High throughput screening

Computers:
-Storage capacity
-Computing capability

Navigational Software:
-Database searching
-Data retrieval

Analysis Software:
-Data mining
-Visualization
-Molecular modeling

Network:
-Sharing data and software
-Grid computing

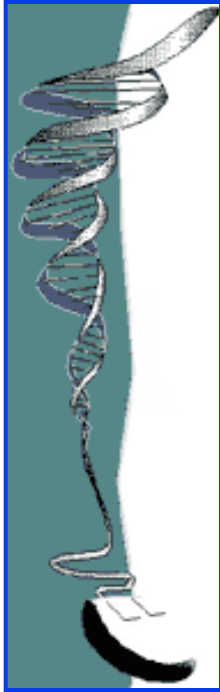
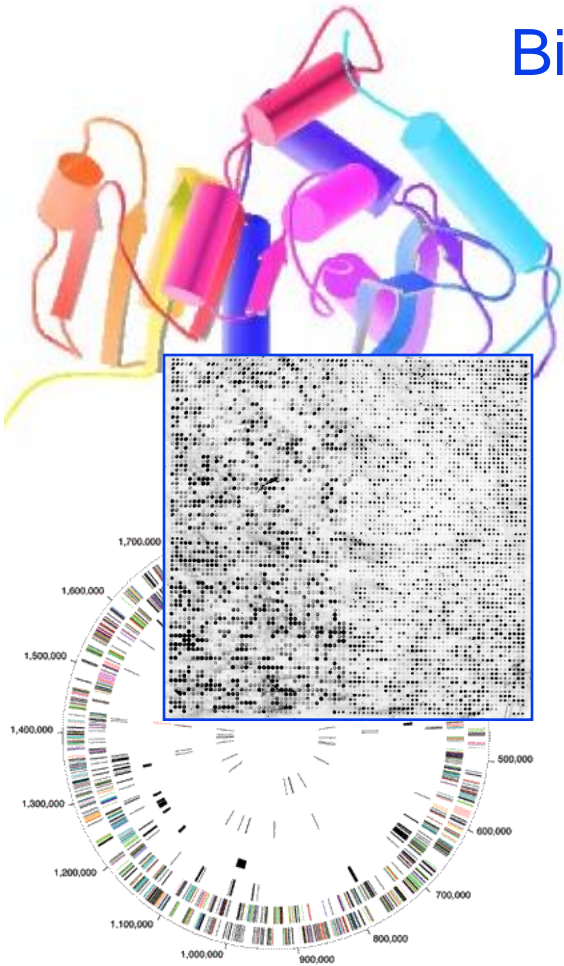
Enabling Research and Product Development in the Life Sciences, e.g.
Pharmaceuticals
New Plant Varieties
Bioremediation
Industrial Processing
Alternative Energy

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 Bio-defense	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

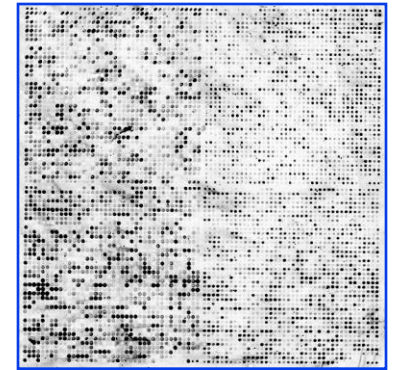
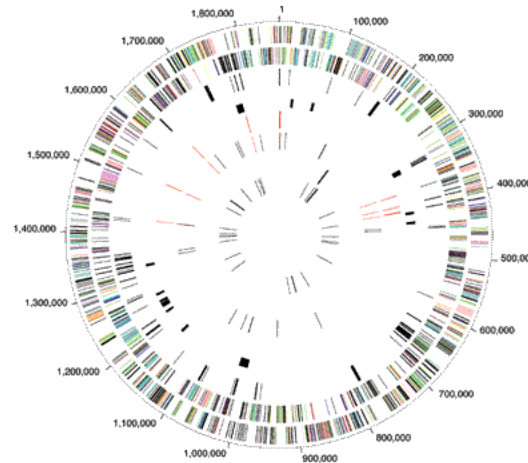
	Conduct research on/in	Approved, marketed, or in production		Total
		Product(s)	Process(es)	
DNA-based				
Bioinformatics	29	2	1	30
Genomics, pharmacogenetics	29	3	2	30
DNA sequencing/synthesis/ amplification, genetic engineering	39	5	3	43
Biochemistry/Immunology				
Drug design & delivery	33	4	2	38
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.

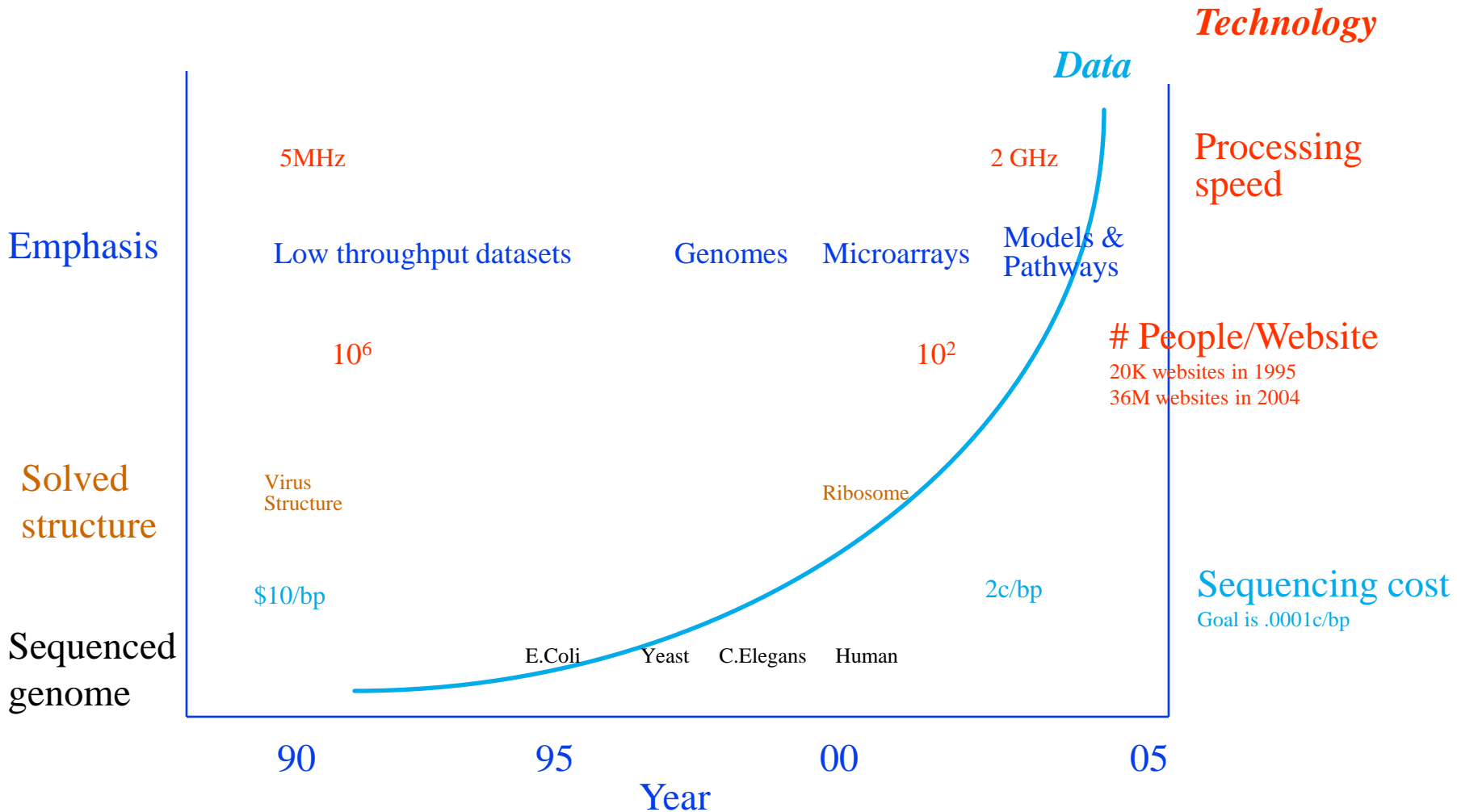
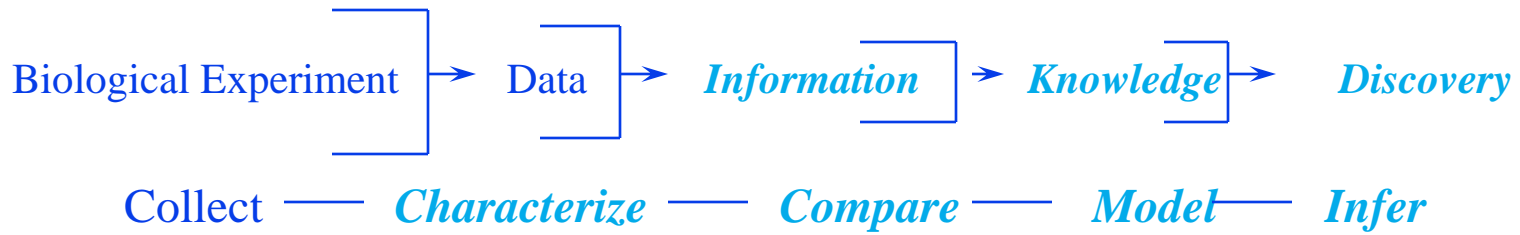
Motivation

- 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

Current problems

- Next generation sequencing
- Gene regulation
- Epigenetics and genetics of diseases, aging
 - SNPs, DNA methylation, histone modification
- Comparison of whole genomes
- Computational systems biology
 - Complexity, dynamics
- Structural bioinformatics, molecular dynamics simulations
- Text mining --- the BioCreative challenge
- and many more