

CENG 465

Introduction to Bioinformatics

Spring 2010-2011

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

Course Web Page:

http://www.ceng.metu.edu.tr/~tcan/ceng465_s1011/

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 problems in bioinformatics
- Opportunity to interact with algorithms, tools, data in current practice

High level overview of the course

- 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
 - Profile HMMS
 - Multiple sequence alignment
 - Efficient pattern search: suffix trees
- Phylogenetic trees (1 week)

Course outline

- Protein structures (4 weeks)
 - Structure prediction (secondary, tertiary)
 - Analyze protein structures for clues regarding function
 - Structure alignment
- Microarray data analysis (2 weeks)
 - Correlations, clustering
- Gene/Protein networks, pathways (2 weeks)
 - Protein-protein, protein/DNA interactions
 - Construction and analysis of large scale networks

Grading

- Midterm exam - 30%
- Final exam - 40%
- Assignments (written/programming) - 30%

Miscellaneous

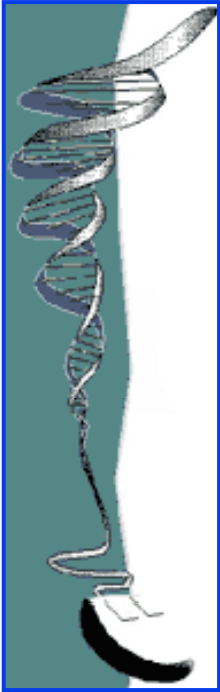
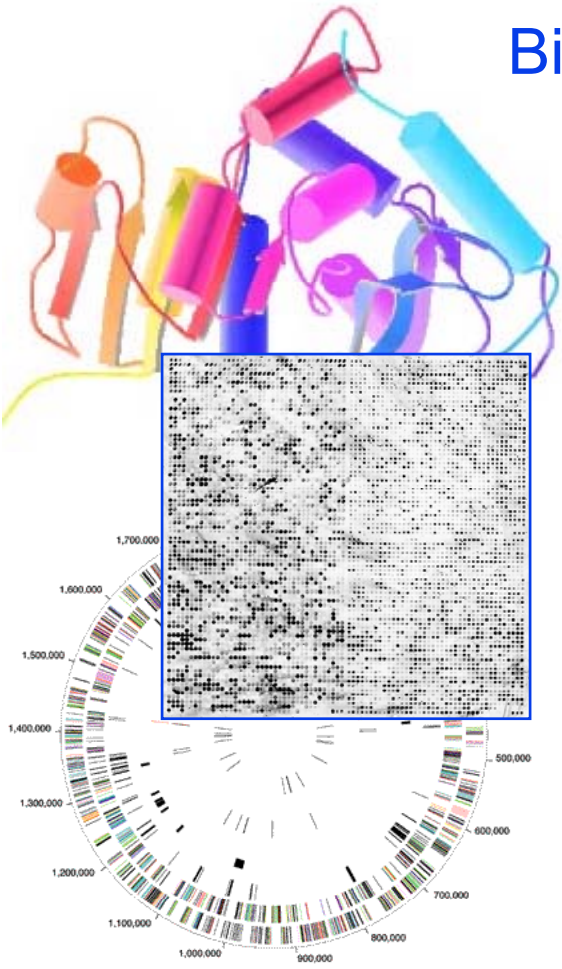
- Course webpage
 - http://www.ceng.metu.edu.tr/~tcan/ceng465_s1011/
 - Lecture slides and reading materials
 - Assignments
 - Teaching assistant: Sefa Kilic (sefa@ceng, B-204)
- Newsgroup
 - metu.ceng.course.465
 - You should follow the newsgroup for course related announcements
 - Students from other departments should get a CENG account for this semester (Room: A-210) in order to access the newsgroup

Bioinformatics: A simple view

Biological
Data

+

Computer
Calculations



What is Bioinformatics?

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

- One idea for a definition?

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

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

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

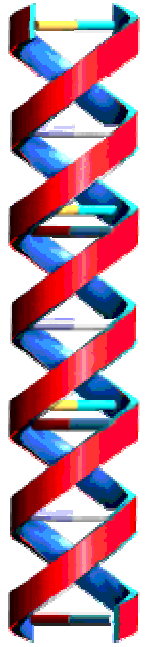
Computing *versus* Biology

looking into the future

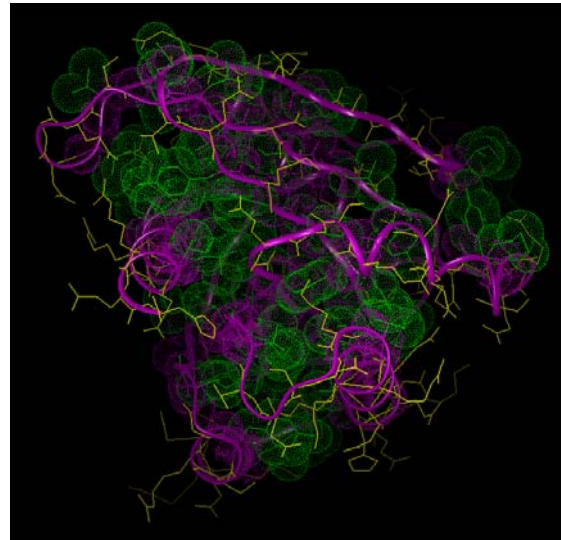
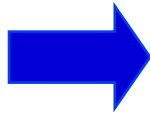
- *Like physics, where general rules and laws are taught at the start, biology will surely be presented to future generations of students as a set of basic systems duplicated and adapted to a very wide range of cellular and organismic functions, following basic evolutionary principles constrained by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

--Temple Smith, Current Topics in Computational Molecular
Biology

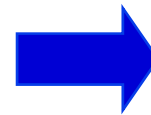
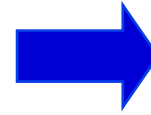
Introductory Biology



DNA
(Genotype)

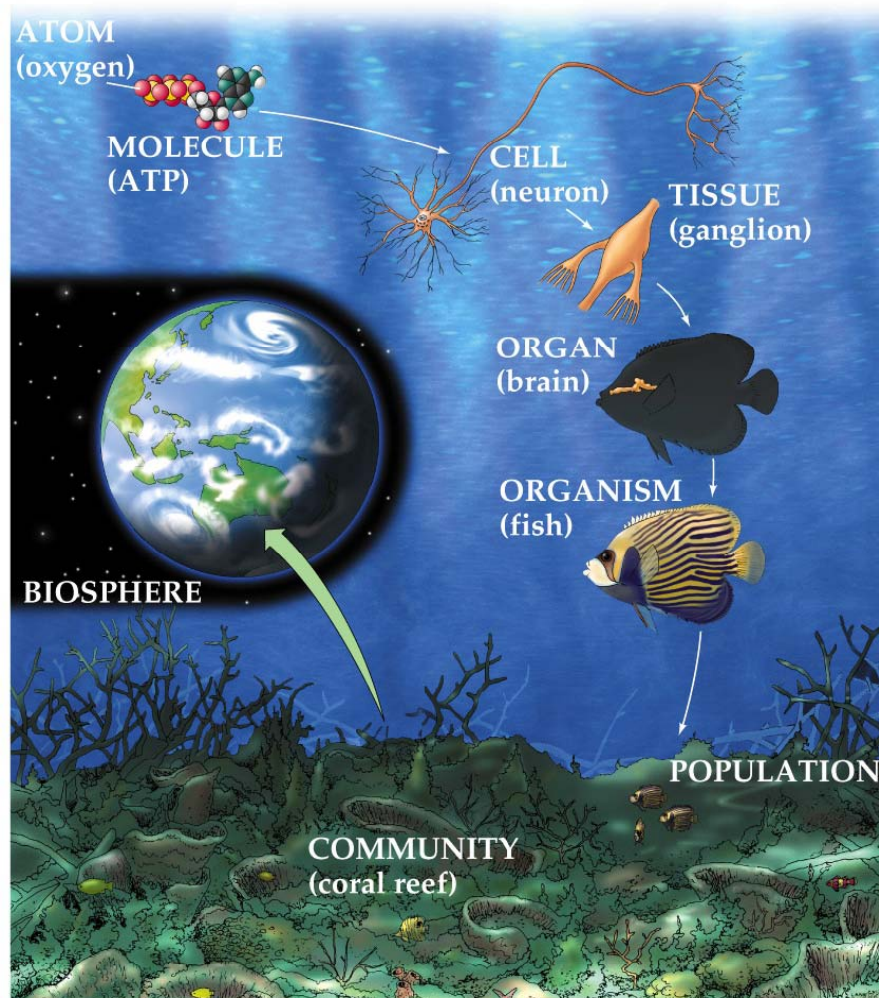


Protein



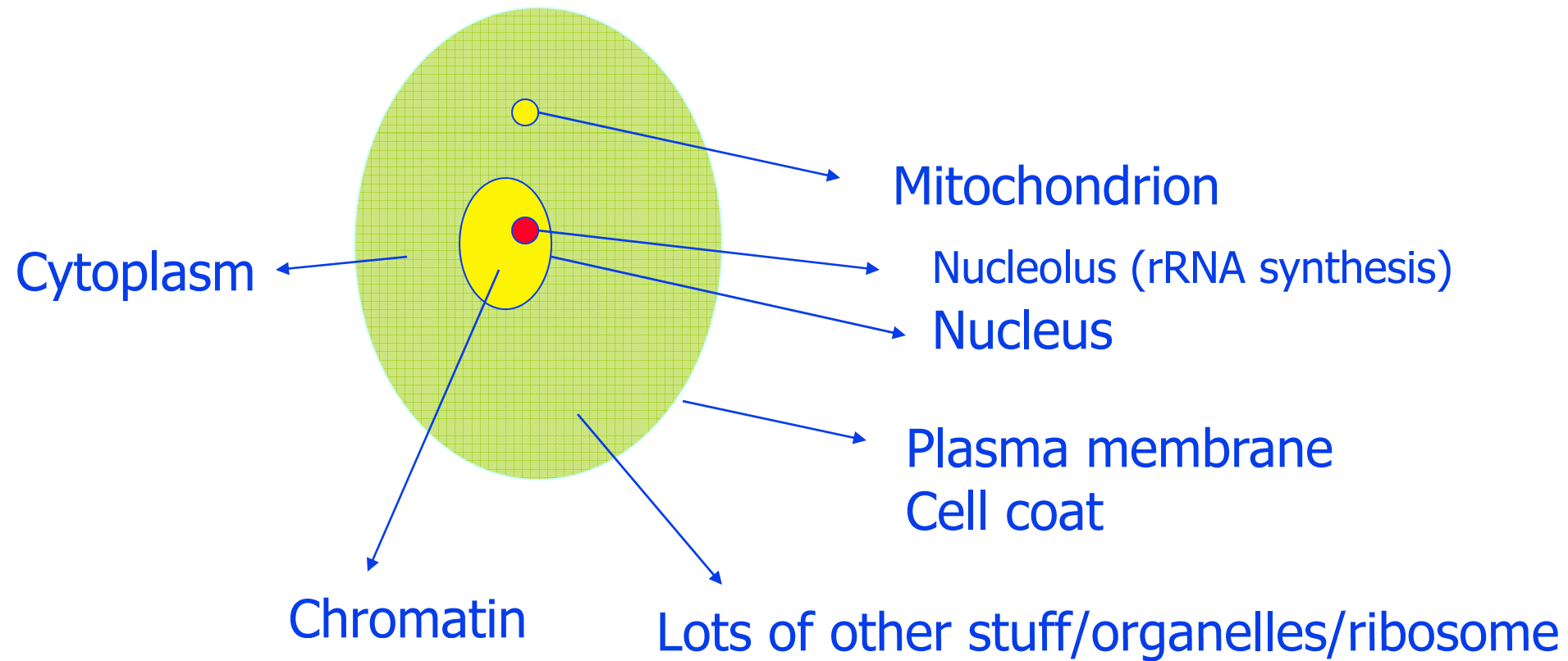
Phenotype

Scales of life

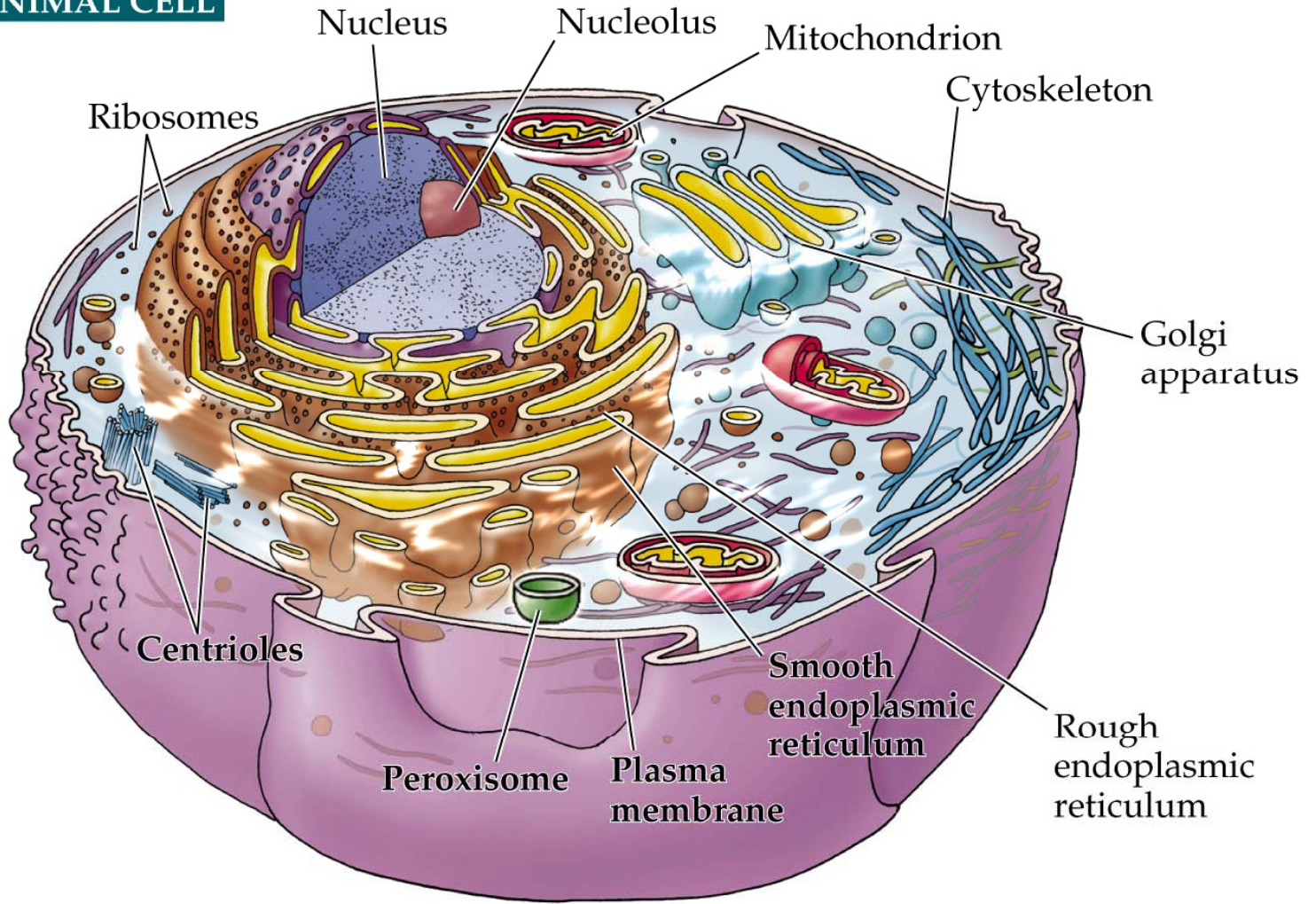


© 2001 Sinauer Associates, Inc.

Animal Cell



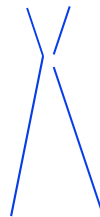
AN ANIMAL CELL



© 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

```
atggcaattaaaattgggtatcaatgggttttggcgtatcggccgtatcgtattccgtgca
gcacaacaccgtgatgacattgaagttgtaggtattaacgacttaatcgacggttgaatac
atggccttataatggtgaaatatgattcaactcacggtcgtttcgacggcactggtgaagtg
aaagatggtaacttagtgggttaatggtaaaactatccgtgtaactgcagaacgtgatcca
gcaaacttaaactgggggtgcaatcgggtggtgatatcgctggtgaagcgaactggtttattc
ttaactgatgaaactgctcgtaaacatatcactgcagggcgcaaaaaaagttgtattaact
ggcccatctaaagatgcaaccctatgttcgttcgtgggtgtaaaacttcaacgcatacgca
ggtcaagatatcgttttctaacgcactcttgtacaacaaactgttagctccttttagcacgt
gttggtcatgaaactttcgggtatcaaagatgggtttaaagaccactgttcacgcaacgact
gcaactcaaaaaactgtggatgggtccatcagctaaagactggcgcgggcgccgcggtgca
tcacaaaacatcattccatcttcaacaggtgcagcgaagcagtaggtaaaagtattacct
gcattaaacggtaaaattaactgggtatggcctttccgtgttccaacgccaacgatatcgtt
gttgatttaacagttaatcttgaaaaaccagcttcttatgatgcaatcaacaagcaatc
aaagatgcagcgggaaggtaaaacggttcaatggcgaattaaaaggcgtattagggttacct
gaagatgctgttgttttctactgacttcaacggttgtgctttaacttctgtatttgatgca
gacgctggtatcgcattaactgattcttttcgttaaattgggtatc . . .
```

```
. . . caaaaatagggttaatatgaatctcgatctccattttgttcatcgtattcaa
caacaagccaaaactcgtacaaatatgaccgcacttcgctataaagaacacggcttgtgg
cgagatatctcttggaaaaactttcaagagcaactcaatcaactttctcgagcattgctt
gctcacaatattgacgtacaagataaaatcgccatttttggccataatatggaacgttgg
gttggtcatgaaactttcgggtatcaaagatgggtttaaagaccactgttcacgcaacgact
acaatcgttgacattgacgaccttacaattcagagcaatcacagtgacctatttacgcaacc
aatacagcccagcaagcagaatttatcctaaatcacgcccgatgtaaaaaattctcttcgtc
ggcgatcaagagcaatacgatcaaacattggaaattgctcatcattgtccaaaattacaa
aaaattgtagcaatgaaatccaccattcaattacaacaagatcctctttcttgcacttgg
```

Molecular Biology Information: Protein Sequence

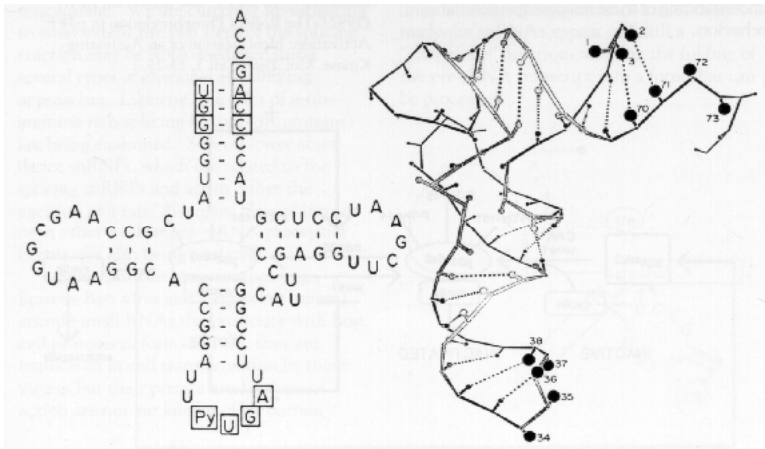
- 20 letter alphabet
 - ACDEFGHIKLMNPQRSTVWY but not BJOUXZ
- Strings of ~300 aa in an average protein (in bacteria),
 - ~200 aa in a domain
- ~1M known protein sequences

```
d1dhfa_ LNCIVAVSQNMGIGKNGDLPWPPLRNEFRYFQRMTTSSVEGKQ-NLVIMGKKTWFSI
d8dfr_  LNSIVAVCQNMGIGKDGNLPWPPLRNEYKYFQRMSTSHVEGKQ-NAVIMGKKTWFSI
d4dfra_ ISLIAALAVDRVIGMENAMPWN-LPADLAWFKRNTL-----NKPVIMGRHTWESI
d3dfr_  TAFLWAQDRDGLIGKDGHLPWH-LPDDLHYFRAQTV-----GKIMVGRRTYESF
```

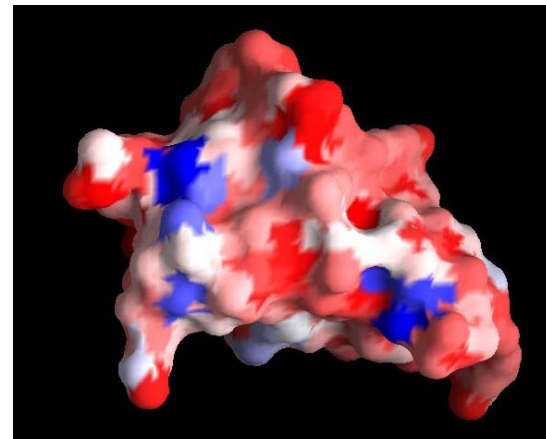
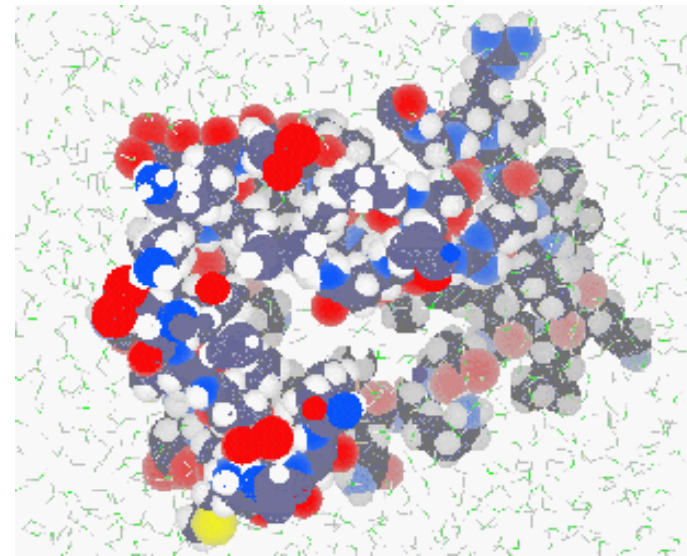
```
d1dhfa_ LNCIVAVSQNMGIGKNGDLPWPPLRNEFRYFQRMTTSSVEGKQ-NLVIMGKKTWFSI
d8dfr_  LNSIVAVCQNMGIGKDGNLPWPPLRNEYKYFQRMSTSHVEGKQ-NAVIMGKKTWFSI
d4dfra_ ISLIAALAVDRVIGMENAMPWN-NLPADLAWFKRNTLD-----KPVIMGRHTWESI
d3dfr_  TAFLWAQDRNGLIGKDGHLPWH-HLPDDLHYFRAQTVG-----KIMVGRRTYESF
```

Molecular Biology Information: Macromolecular Structure

- DNA/RNA/Protein
 - Almost all protein



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



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

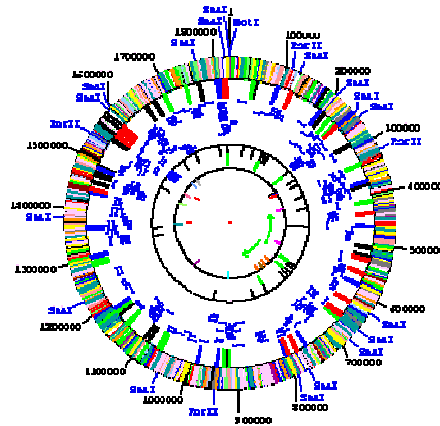
- One gene encodes one* protein.
- Like a program, it starts with start codon (e.g. ATG), then each three code one amino acid. Then a stop codon (e.g. TGA) signifies end of the gene.
- Sometimes, in the middle of a (eukaryotic) gene, there are introns that are spliced out (as junk) during transcription.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

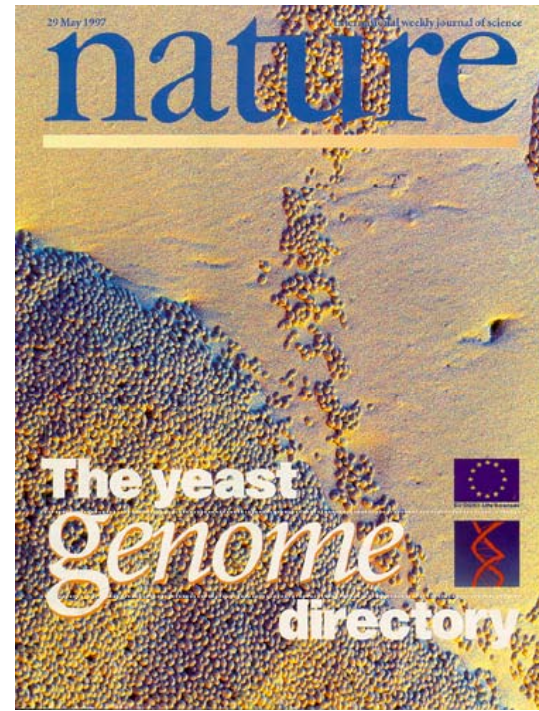


Genome sequences now accumulate so quickly that, in less than a week, a single laboratory can produce more bits of data than Shakespeare managed in a lifetime, although the latter make better reading.

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

1995

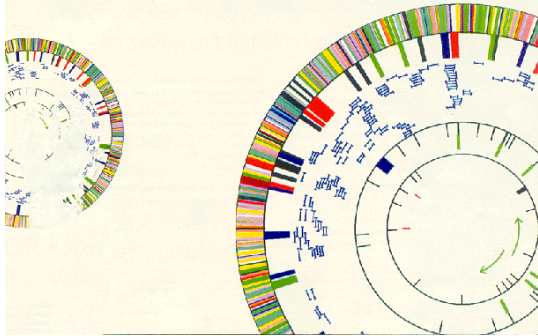
Bacteria,
1.6 Mb,
~1600 genes
[*Science* 269: 496]



Genomes highlight the Finiteness of the “Parts” in Biology

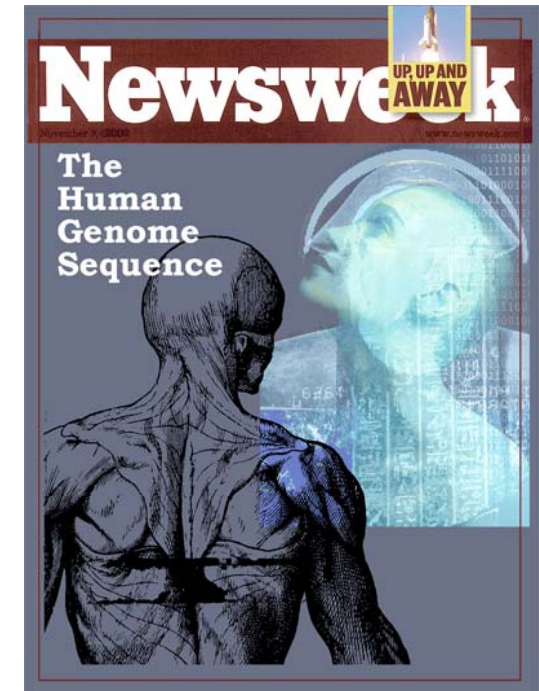
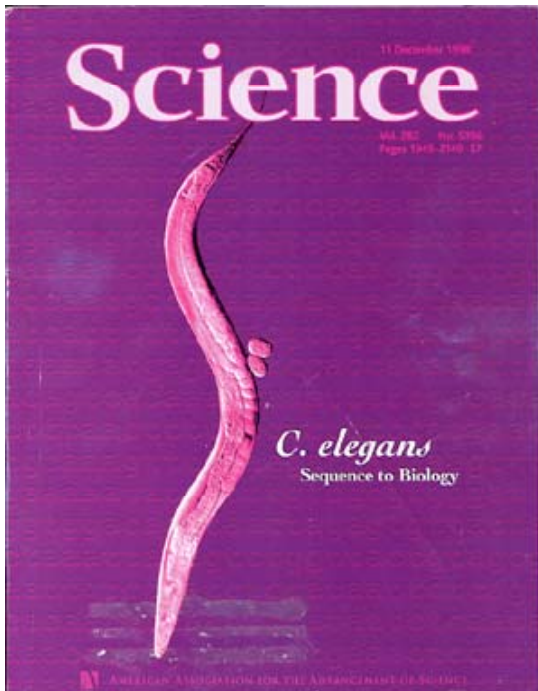
1997

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



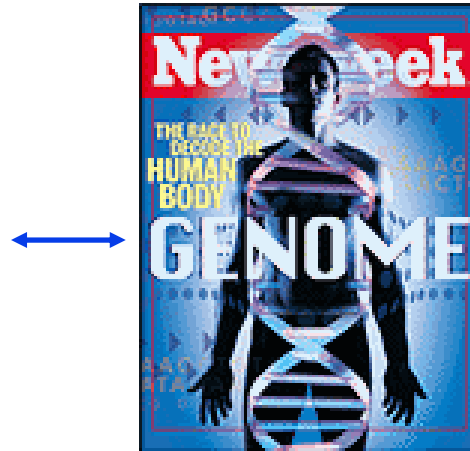
1998

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*

YGA 99-1133R

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

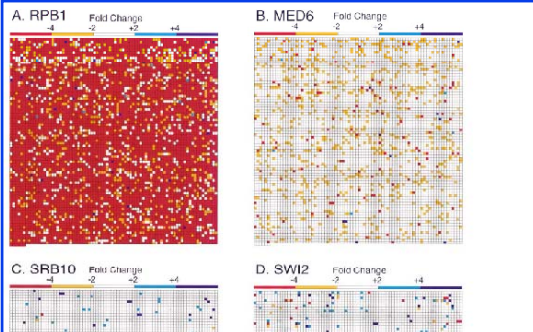


Figure 2. Genomewide Expression Data for Selected Components of the RNA Polymerase II Holoenzyme. Changes in mRNA levels when a mutant is compared to its isogenic wild-type counterpart is presented in a grid format. In the left grid square represents the left-most gene on chromosome I, and the squares to its right represent adjacent genes, in a fashion through chromosome I, then II, 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 inactivator. Comparison of the two data sets reveals that the expression of a subset of genes is altered by the inactivator. 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 is also not known (Thompson et al., 1994; Koleske and Young, 1994; Hen-Meyers et al., 1998). To determine dependence of gene expression on an *SRB5* gene and its wild-type compared (see the web site for detailed results indicate that 16% of all genes for their expression. With rain and other constitutive mutants

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, a 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-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 transposons by mutagenesis of the *SPA2*, *ARP100*, *SER1*, and *BDF1* genes.

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

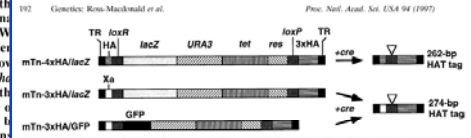


Fig. 1. Schematic representation of the mTn construct and the derived HAT tag element. Each mTn contains the coding region for the transposase (TR) and Ura3 protein, and the cre element from *Tet*. mTn-*lacZ* and mTn-*HA* contain a transposon flanked by two *loxP* sites. mTn-*HA* also contains a sequence 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-*HA*), or the factor Xa protease cleavage site (mTn-*HA*), or a GFP coding region (mTn-*GFP*). Expression of these transposons to *Cre* recombination results in the formation of a smaller element encoding the HAT tag element (to the right). The site of the *loxP* recombination reaction is indicated by a triangle. (Not drawn to scale.)

transposon, and contains the *Tet* *res* site for resolution of transposon conjugates. *Tet*-resoluted elements catalyze transposon transposition and resolution as provided in nature. All three transposons contain the *URA3* and *arg* genes for selection in *S. cerevisiae*. *E. coli* (respectively, Transposon mTn-*lacZ*, mTn-*HA*, and mTn-*HA*), while transposon mTn-*HA* also contains an initiator methionine, while transposon mTn-*HA* also contains an initiator methionine, while transposon mTn-*HA* also contains an initiator methionine. Levels of both activities can be measured quantitatively and have been shown to provide indices of gene expression (e.g., refs. 4, 23, and 24). 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* recombination are divergent from one another and undergo levels of spontaneous recombination. The *loxP* sites are internal to sequences encoding multiple copies of an epitope from the influenza virus hemagglutinin protein (the HA epitope; ref. 25). The mTn-*HA* transposon also contains a factor Xa protease cleavage site (19) in the region external to the *lacP* site. Expression of the *Cre* recombination 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-*HA*) or 262-bp (mTn-*HA*) element. This element consists of a single *loxP* site and sequences encoding three or four copies of the HA epitope, flanked by the *Tet* terminal repeats (Fig. 1). The mTn-*HA* derived element also contains a sequence encoding the factor Xa cleavage site. When the transposon 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-*HA*) or 89 amino acids (mTn-*HA*) into the protein. We designate these insertions HAT tags.

Mutagenesis of Yeast Genes. Transposons mTn-*HA* and mTn-*HA* were used to mutagenize the yeast *SPA2* gene. *SPA2* encodes a nonessential protein that is critical to sites of polarized growth; *spa2* mutants exhibit defects in

was mutagenized in *E. coli* 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-*HA* and mTn-*HA*, about 10% of transformants were identified as producing β -galactosidase protein (41,500 mTn-*HA* and 41,500 mTn-*HA* transformants for *SPA2* and *ARP100*, 25,200 and 62,500 mTn-*HA* transformants for *SPA2*, *ARP100*, and *SER1*, respectively). Strains expressing the reporter genes were used for further analysis. The approximate 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* Recombination. Although efficient Cre-mediated recombination between *loxP*

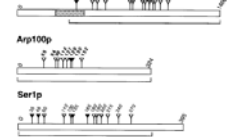


Fig. 2. Map showing amino acid positions of HAT tag insertion into the yeast proteins Spa2p, Arp100p, and Ser1p. Regions recognized as insertion sites are indicated by brackets. Insertion positions determined by sequencing the derived mTn-*HA* DNA. The HAT tag consists of the *loxP* site, the HA epitope, and the factor Xa cleavage site. For the mTn-*HA* derived element, the HAT tag also contains a sequence encoding the factor Xa cleavage site.

Young/Lander, Chips, Abs. Exp.

Specific transcription factors, a novel mechanism for the resolution of specific sets of genes

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 response of human fibroblasts to serum

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

Exploring the Genes of the Database

See the entire transcript



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

Snyder, Transposons, Protein Exp.

Brown, μ array, Rel. Exp. over Timecourse

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,² Patrice M. Chai Pai,¹ Corinne Rebschung,⁸ Christopher J. Roberts,² Petra R. Michael Snyder,⁴ Sharon Sookhai Steeve Véronneau,⁷ Marleer Teresa R. Ward,² Robert Wysocki Katja Zimmermann, Mark Johnston,¹³

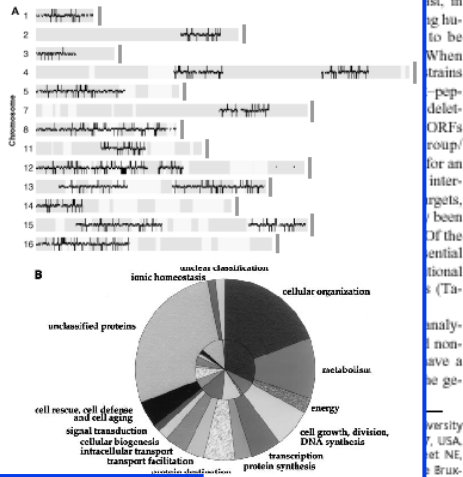
that serve as strain identifiers (6, 7). We show that these barcodes allow large numbers of deletion strains to be pooled and analyzed in parallel in competitive growth assays. This direct, simultaneous, competitive assay of fitness increases the sensitivity, accuracy and speed with which growth defects 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 generate deletion strains for all unannotated essential genes (69% of which were within 5 kb of another gene). Whereas 47% of nonessential essential genes were generally 50 kb of the telomeres (Fig. 1), whereas 90% of nonessential essential genes were 70% high-copy number of transcripts. The function of the essential genes is shown in Fig. 1.

The analysis of the deletion strains by allowing the growth of those whose cognate essential to life, is a formidable task for many genes will likely be under very specialized conditions, necessitating the examination of different conditions. Previous studies that the barcodes allowed an increase of their respective abundance when 12 strains were assayed in parallel.

The functions of many open reading frame sequencing projects are unknown. Now, to systematically determine their functions, yeast strains were constructed, by a precise deletion of one of 2026 ORFs genome. Of the deleted ORFs, 17 per cent medium. The phenotypes of more than 1000 deletion strains, 40 per cent 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



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 Luo 1,2, Mengsheng Qiu 1,3, Eva Chan 2, Helen Zhou 4, Li Zhu

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

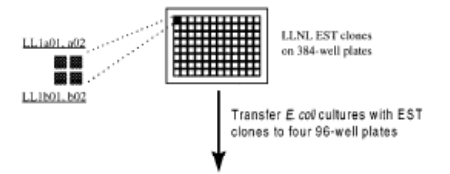
Received 1 February 1998; received in revised form 28 April 1998; accepted 29 April 1998; Received by E.Y. Chen

Abstract

Identification of all human protein-protein interactions is an important information for functional genomics. We have constructed a modular human two-hybrid cDNA libraries using human EST clones. Quality analysis of this library indicates that human EST clones is feasible, and so far the first time that a comprehensive two-hybrid cDNA library has been constructed.

148

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



Keywords: Functional genomics; Yeast two-hybrid; Human EST clones

1. Introduction

The Human Genome Project has produced a tremendous amount of DNA sequence information. Over 50 000 UniGenes have been identified (Schuler, 1995; Miller et al., 1996). Approximately 50% of the human genome is transcribed into RNA, but only a minority of these UniGenes are known to be expressed in any given tissue or cell type.

2 hybrids, linkage maps

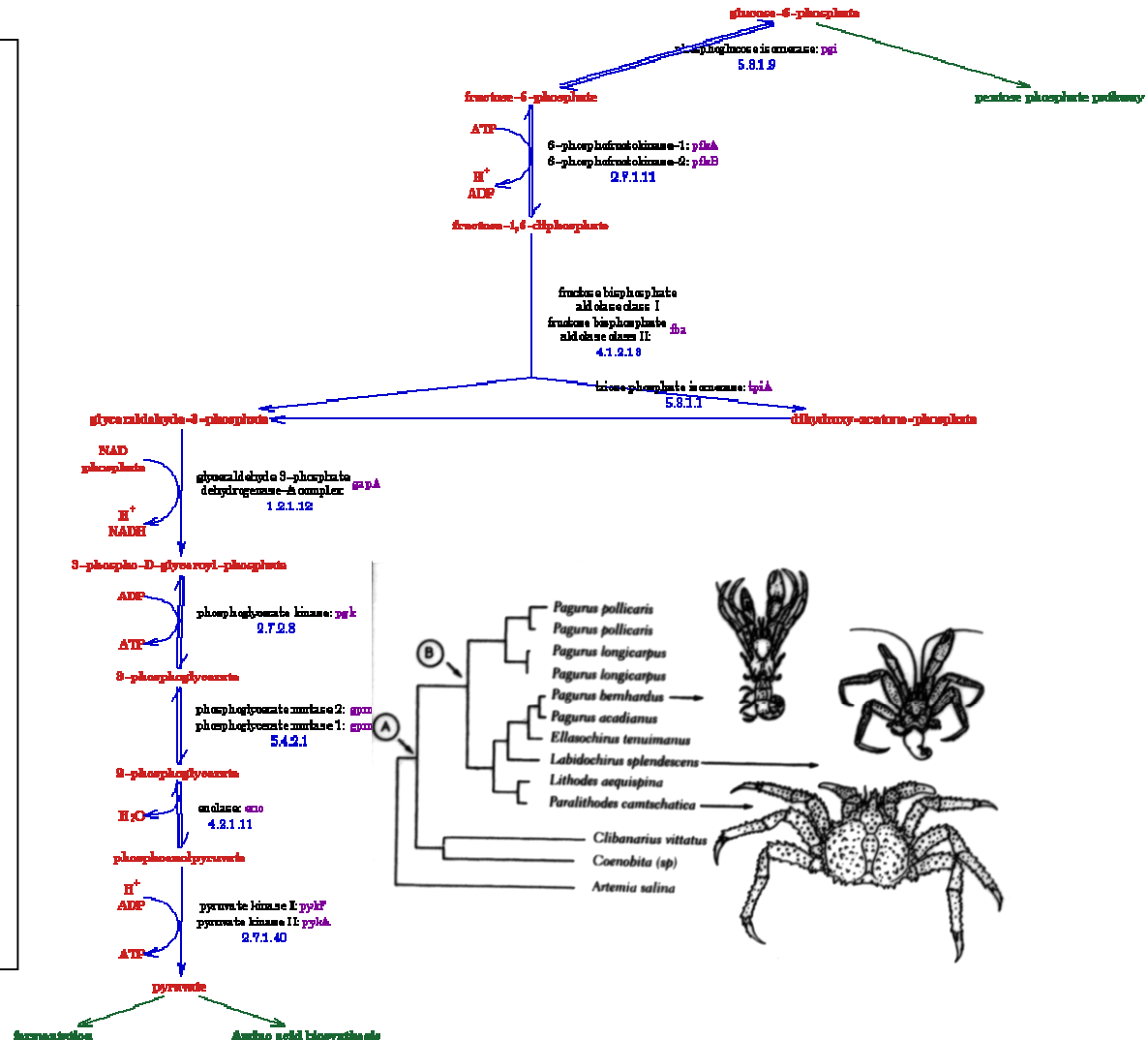
Hua, S. B., Luo, Y., Qiu, M., Chan, E., Zhou, H. & Zhu, L. (1998). Construction of a modular yeast two-hybrid cDNA library from human EST clones for the human genome protein linkage map. *Gene* **215**, 143-52

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

* Corresponding author. Tel: +1 650 924 6500. 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

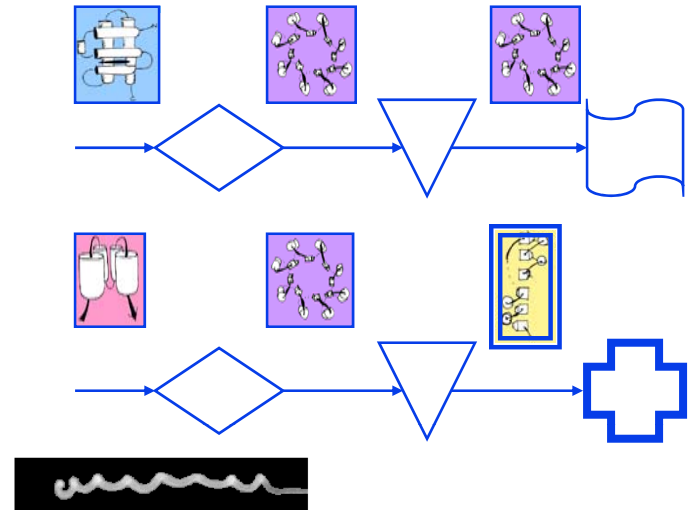
- Information to understand genomes
 - Metabolic Pathways (glycolysis), traditional biochemistry
 - Regulatory Networks
 - Whole Organisms
 - Phylogeny, traditional zoology
 - Environments, Habitats, ecology
 - The Literature (MEDLINE)
- The Future....



Organizing

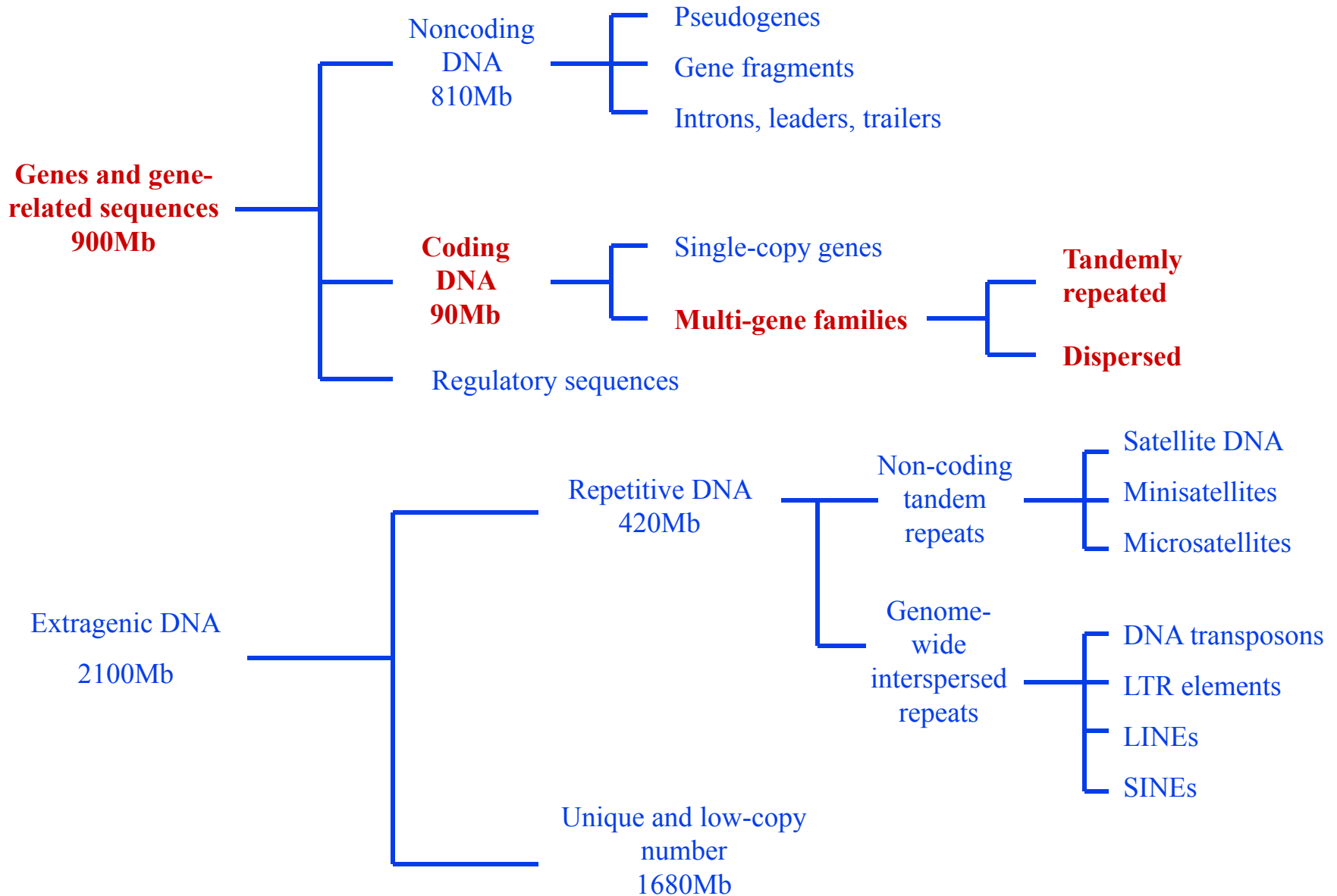
Molecular Biology Information: Redundancy and Multiplicity

- 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

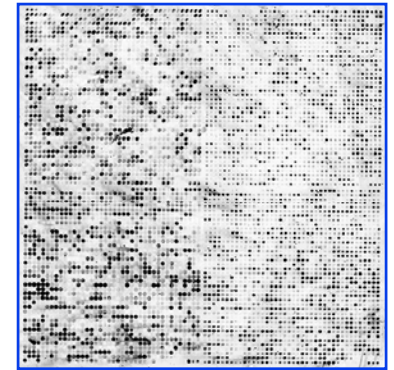
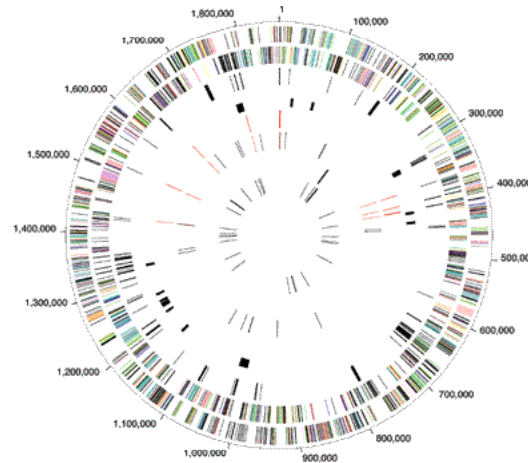


Where to get data?

- 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

Data

- 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

Scalability challenges

- As of December 2009, NAR online Database Collection, available at <http://www.oxfordjournals.org/nar/database/a/>, lists 1230 carefully selected databases covering various aspects of molecular and cell biology
 - 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

- Programming
- Algorithms
- Machine learning/Pattern recognition/AI
- Statistics & probability
- Mathematics