

# Introduction to Bioinformatics

Sami Khuri

Department of Computer Science

San José State University

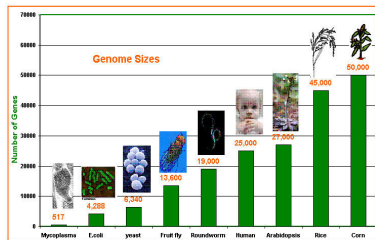
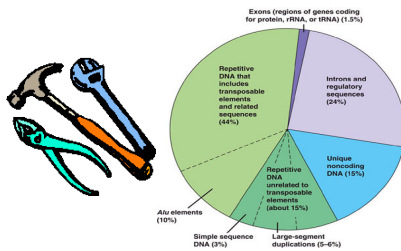
San José, California, USA

khuri@cs.sjsu.edu

www.cs.sjsu.edu/faculty/khuri

@2010 Sami Khuri

## What is Bioinformatics?



- The Human Genome Project (HGP)
- Mapping
- Model Organisms
- Types of Databases
- Applications of Bioinformatics
- Genome Research



@2010 Sami Khuri



## The Human Genome Project

- The **HGP** is a multinational effort, begun by the USA in 1988, whose aim is to produce a complete physical map of all human chromosomes, as well as the entire human DNA sequence.
  - As part of the project, genomes of other organisms such as bacteria, yeast, flies and mice are also being studied.
- The primary goal of the project is to make a series of descriptive diagrams (called **maps**) of each human chromosome at increasingly finer resolutions.

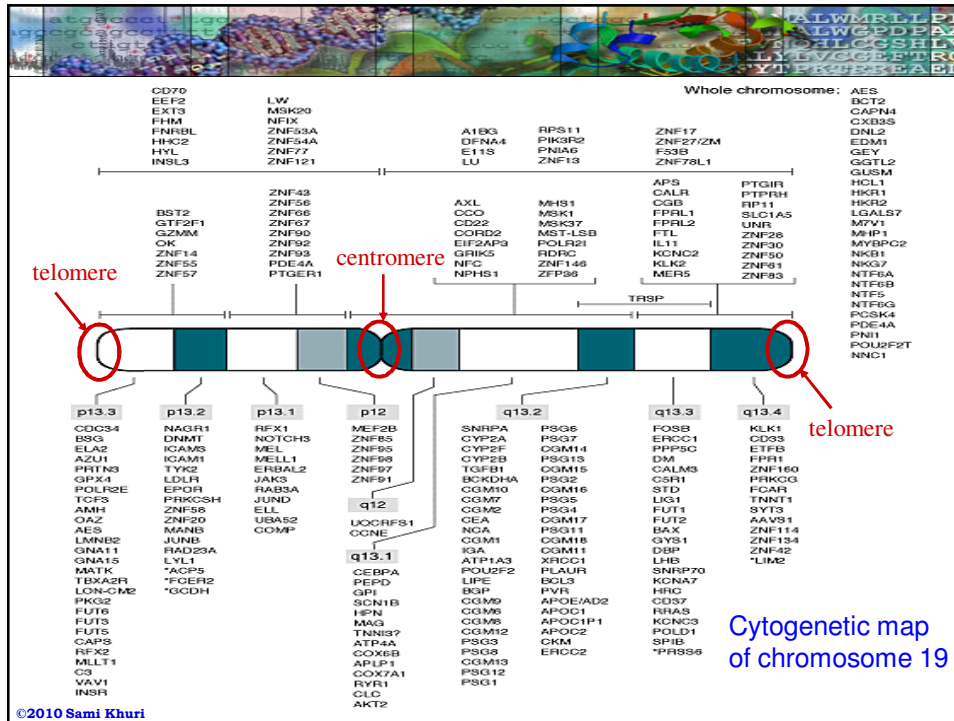
©2010 Sami Khuri



## The HGP Goal

- The ultimate goal of genome research is to find all the **genes** in the **DNA sequence** and to develop tools for using this information in the study of **human biology** and **medicine**.
- **Mapping** involves:
  - dividing the chromosomes into smaller fragments that can be propagated and characterized
  - ordering (mapping) them to correspond to their respective locations on the chromosomes.

©2010 Sami Khuri

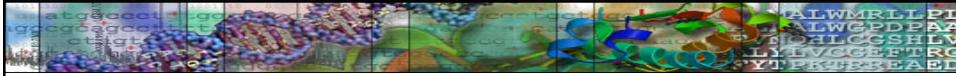


Cytogenetic map of chromosome 19

## Goals of the HGP

- To *identify* all the approximately 20,000-25,000 genes in human DNA,
- To *determine* the sequences of the 3.2 billion chemical base pairs that make up human DNA,
- To *store* this information in databases,
- To *improve* tools for data analysis,
- To *address* the ethical, legal, and social issues (ELSI) that may arise from the project.

©2010 Sami Khuri



## HGP Finished Before Deadline

- In 1991, the USA Congress was told that the HGP could be done by 2005 for \$3 billion.
- It ended in 2003 for \$2.7 billion, because of efficient computational methods.

©2010 Sami Khuri

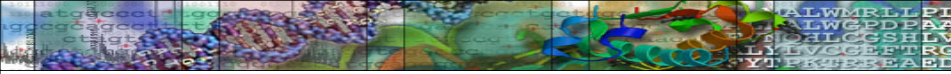


## Other Species

As part of the HGP, genomes of other organisms, such as bacteria, yeast, flies and mice are also being studied.

 <p>p53 gene pax6 gene</p>	 <p>C. elegans Diabetes</p>	 <p>Baker's yeast DNA repair Cell division</p>
	 <p>Chimps are infected with SIV Very rarely progress to AIDS</p>	

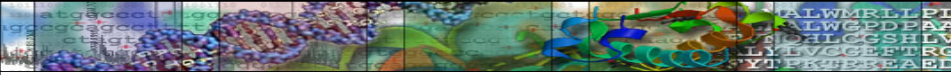
©2010 Sami Khuri


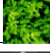
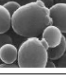







## Model Organisms

- A **model organism** is an organism that is extensively studied to understand particular biological phenomena.
- **Why have model organisms?** The hope is that discoveries made in model organisms will provide insight into the workings of other organisms.
- **Why is this possible?** This works because evolution reuses fundamental biological principles and conserves metabolic, regulatory, and developmental pathways.

©2010 Sami Khuri



Name	Genome BP	Genes	Chromosomes
HSV1 (Herpes virus) 	1.5x10 <sup>5</sup>	70	1
Escherichia Coli 	4.6x10 <sup>6</sup>	4,300	1
Saccharomyces cerevisiae 	1.2x10 <sup>7</sup>	5,900	16
Caenorhabditis Elegans 	1.0x10 <sup>8</sup>	19,100	6
Drosophila melanogaster 	1.8x10 <sup>8</sup>	13,600	6
Arabidopsis Thaliana 	1.2x10 <sup>8</sup>	25,500	5
Mus Musculus 	2.5x10 <sup>9</sup>	~30,000	20+X/Y
Homo sapiens 	2.9x10 <sup>9</sup>	~30,000	22+X/Y

David Gilbert

©2010 Sami Khuri






## Sequencing the Chimpanzee




Almost human...

Sequencing the chimpanzee has emerged as a top genomic priority. David Cyranoski asks the chimp's champions what they hope to gain from studying the genome of our closest living relative. *Nature*, vol 418, August 2002

©2010 Sami Khuri



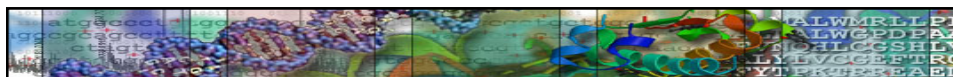


## Studying Diseases

<b>Comparison of disease susceptibility between chimps and humans</b>			
	<i>Condition</i>	<i>Human</i>	<i>Chimp</i>
<i>Definite differences</i>			
	HIV progression to AIDS	Common	Very rare
	Influenza A symptoms	Moderate/severe	Mild
	Hepatitis B/C late complications	Moderate/severe	Mild
	<i>Plasmodium falciparum</i> malaria	Susceptible	Resistant
	Menopause	Universal	Rare
<i>Probable differences</i>			
	<i>Escherichia coli</i> K99 gastroenteritis	Resistant	Sensitive?
	Alzheimer's disease pathology	Complete	Incomplete
	Coronary atherosclerosis	Common	Uncommon
	Epithelial cancers	Common	Rare

Source: Olson, M. V. et al. A white paper advocating complete sequencing of the genome of the common chimpanzee, Pan troglodytes. (2002).

©2010 Sami Khuri



## Studying Human Diseases

Organism	Human Diseases
<i>E. coli</i>	DNA repair; colon cancer and other cancers
Yeast	Cell cycle; cancer, Werner syndrome
<i>Drosophila</i>	Cell signaling; cancer
<i>C. elegans</i>	Cell signaling; diabetes
Zebrafish	Developmental pathways; cardiovascular disease
Mouse	Gene expression; Lesch-Nyhan disease, cystic fibrosis, fragile-X syndrome, and many other diseases

Copyright © 2006 Pearson Prentice Hall, Inc.

©2010 Sami Khuri

F	W	Y	Neurological	F	W	Y	Renal
+			Alzheimer-PS1	+			Diabetes Insipidus 2-AQP2
+			Alzheimer-APP	-			Polycystic Kidney 1-PKD1
-			Creutzfeldt-Jakob-PRNP	+			Polycystic Kidney 2-PKD2
+			Deafness, Hereditary-MYO15				
+			Dementia, Multi-Infarct-NOTCH3	F	W	Y	Endocrine
+			Duchenne MD*2-DMD	+			Diabetes-INS
+			Fragile-X-FRAXA	+			Diabetes-INSR
+			Huntington-HD	+			Hyperinsulinism-ABCC8
+			Limb Girdle MD*2A-CAPN3	+			Hyperinsulinism-KCNJ11
+			Limb Girdle MD*2B-YSF	+			Obesity-LEP
-			Limb Girdle MD*2E-BSG	-			Obesity-LEPR
+			Myotonic Dystrophy-DM1	-			Vitamin-D Resis. Rickets-VDR
+			Myotubular Myopathy 1-MTM1				
-			Parkinson-SNCA	F	W	Y	Metabolic
+			Parkinson-PARK2	+			Cystinuria, Type 1-SLC3A1
+			Parkinson-UCHL1	-			Hypercalcemia-CASR
+			Tay-Sachs-HEXA	+			Niemann-Pick C-NPC1
				+			SCID**-ADA
F	W	Y	Immune				
+			Bruton Agammaglobulin-BTK	F	W	Y	Other
+			Chronic Granulom.-CYBB	+			Cystic Fibrosis-ABCC7
+			Immunodeficiency-DNA Ligase 1	+			Hereditary Pancreatitis-PRSS1
-			Immunodeficiency-CD3G	+			Juvenile Glaucoma-GLC1A
+			SCID**-JAK3	+			Wolfram-WFS1
-			SCID**-RAG1				
-			SCID**-RAG2	F	W	Y	Cardiovascular
+			SCID**-ZAP70	+			Fam. Cardiac Myopathy-MYH7
				+			HDL Deficiency 1-ABCA1
F	W	Y	Cardiovascular				
+			Fam. Cardiac Myopathy-MYH7	F	W	Y	Birth Defects
+			HDL Deficiency 1-ABCA1	+			Holoprosencephaly 3-SHH
				+			Holoprosencephaly-SIX3
				+			Zellweger-PEX1

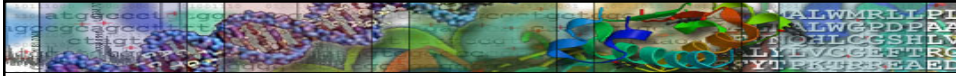
Flies have **orthologs** to humans disease-causing genes in categories such as:

- neurological
- renal
- immunological
- endocrine
- cardiovascular
- metabolic
- blood-vessel and
- cancerous disorders

Flies can provide insights into human disease at the **systems level**, revealing how different genes interact in vivo


Discovering Genomics, Campbell, 2007

©2010 Sami Khuri




## What is Bioinformatics? Set of Tools

- The use of computers to collect, analyze, and interpret biological information at the molecular level.
- A set of software tools for molecular sequence analysis



©2010 Sami Khuri



## What is Bioinformatics? A Discipline

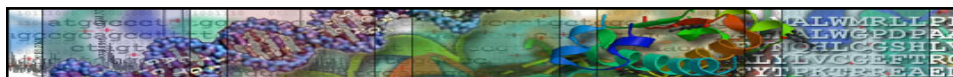
- The field of science, in which **biology**, **computer science**, and **information technology** merge into a single discipline.

*Definition of NCBI (National Center for Biotechnology Information)*

- The ultimate goal of **bioinformatics** is to enable the discovery of new biological insights and to create a global perspective from which unifying principles in biology can be discerned.

©2010 Sami Khuri

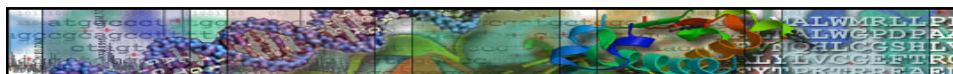




## Bioinformatics and the Internet

- The enormous increase in biological data has made it necessary to use **computer information technology** to collect, organize, maintain, access, and analyze the data.
- Computer speed, memory, and exchange of information over the Internet has greatly facilitated **bioinformatics**.
- The **bioinformatics** tools available over the Internet are accessible, generally well developed, fairly comprehensive, and relatively easy to use.

©2010 Sami Khuri



## What do Bioinformaticians do?

- Analyze and interpret data
- Develop and implement algorithms
- Design user interface
- Design database
- Automate genome analysis
- Assist molecular biologists in data analysis and experimental design.

©2010 Sami Khuri

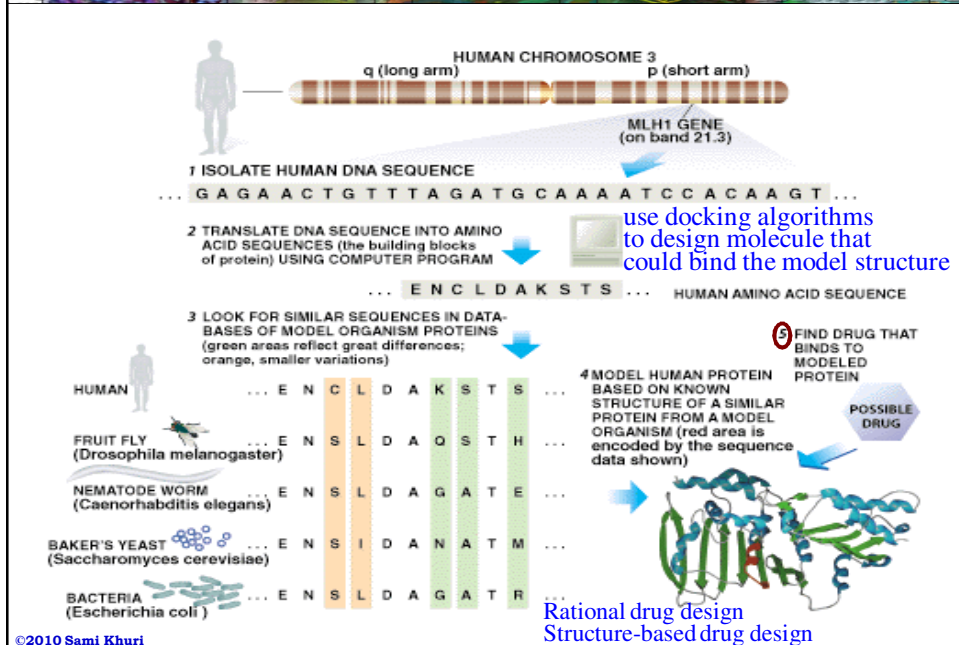


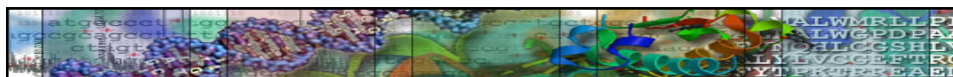
## Why Study Bioinformatics?

- Bioinformatics is intrinsically interesting
- Bioinformatics offers the prospect of finding better drug targets earlier in the drug development process.
  - By looking for genes in model organisms that are similar to a given human gene, researchers can learn about protein the human gene encodes and search for drugs to block it.



©2010 Sami Khuri

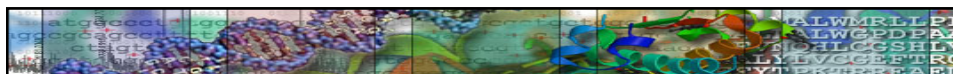




## Databases for Storage and Analysis

- Databases store data that need to be analyzed
- By comparing sequences, we discover:
  - How organisms are related to one another
  - How proteins function
  - How populations vary
  - How diseases occur
- The improvement of sequencing methods generated a lot of data that need to be:
  - stored
  - organized
  - curated
  - annotated
  - managed
  - networked
  - accessed
  - assessed

©2010 Sami Khuri



## Types of Databases

- **Sequence**
  - Genbank, SwissProt, 3D structure, carbohydrates, organism specific, phylogenetic, sequence patterns
- **Literature**
  - Medline, OMIM, Patents, eJournals
- **Graphical**
  - Swiss2D-Page
- **Expression Analysis Databases**
  - Microarrays
- **Protein Interaction Databases**
  - Pathways

©2010 Sami Khuri




## Three Major Databases



- **GenBank** from the NCBI (National Center of Biotechnology Information), National Library of Medicine <http://www.ncbi.nlm.nih.gov>
- **EBI** (European Bioinformatics Institute) from the European Molecular Biology Library <http://www.ebi.ac.uk>
- **DDBJ** (DNA DataBank of Japan) <http://www.ddbj.nig.ac.jp>

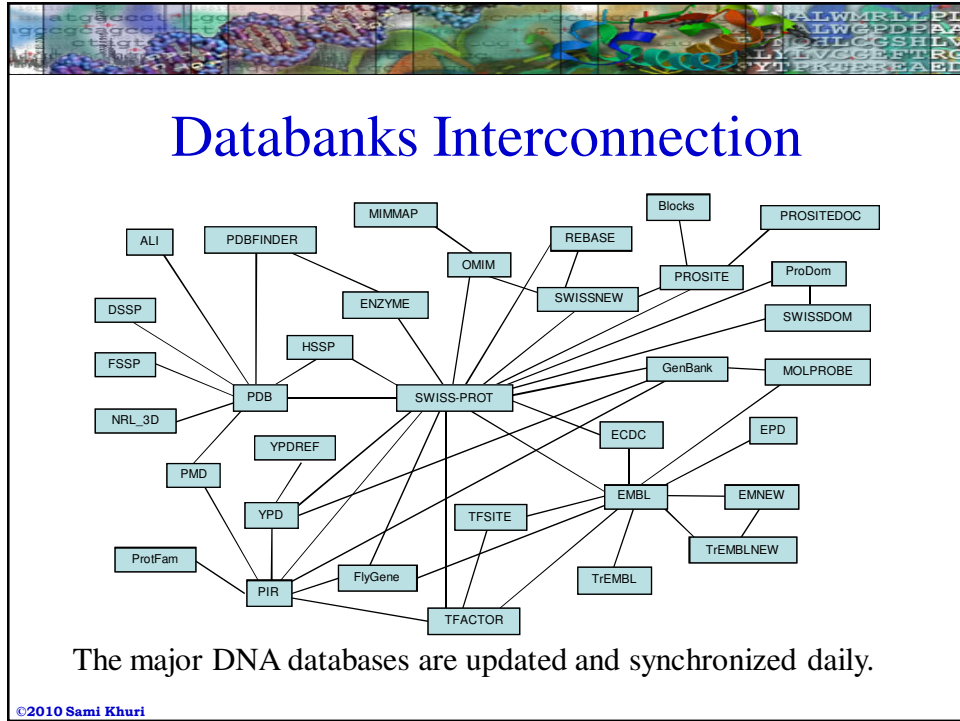
©2010 Sami Khuri



## GenBank Taxonomic Sampling

Homo sapiens	62.1%
Mus musculus	7.7%
Drosophila melanogaster	6.1%
Caenorhabditis elegans	3.3%
Arabidopsis thaliana	2.9%
Oryza sativa	1.3%
Rattus norvegicus	0.8%
Danio rerio	0.6%
Saccharomyces cerevisiae	0.6%

©2010 Sami Khuri

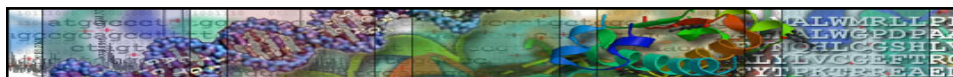


**The Annotation Challenge**

- One of the biggest challenges for maintainers of biological databases is the **annotation**:
  - putting sufficient **information** in the database such that there is no question of what the gene is,
  - creating the proper **links** between that information and the gene sequence and serial number.
- Correct **annotation** of genomic data is an active area of research.

©2010 Sami Khuri





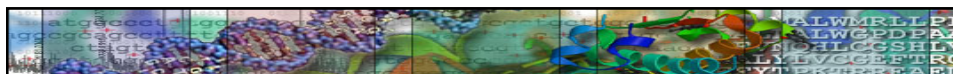
## Data Formats

- Most databases are relational DB. The format can be changed to export the data:
  - GenBank data is stored by NCBI in a Sybase database, but made public in a flatfile format.
- Each database and each sequence analysis program store/accept data in a different format.
  - The most commonly used data format is FASTA

>gil37222328|gblAY350722.1|Pan troglodytes masticatory myosin heavy chain (MYH16)

```
GAGCAGCTGAACAAGCTGATGACCACCCTCCATAGCACCCGCACCCCATTTTG  
TCCGCTGTATTATCCCCAATGAGTTTAAGCAATCGG
```

©2010 Sami Khuri

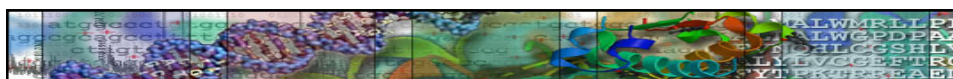


## Reasons for Searching Databases

**Searching a database** can answer the following questions:

- A researcher has just sequenced a gene.  
**Has someone already found it?**
- A researcher has a sequence of unknown function.  
**Is there a homology with another sequence that has a known function?**
- A researcher has found a new protein in a lower organism.  
**Is there a homology in a higher species?**

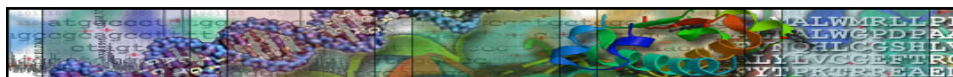
©2010 Sami Khuri



## Protein Databases

- **GenPept** from NCBI.
- **ExPASy** (Expert Protein Analysis System)  
<http://www.expasy.ch>
- **SwissProt, TrEMBL**  
<http://www.ebi.ac.uk>
- **PIR** (Protein Identification Resource)  
<http://www-nbrf.georgetown.edu/pirwww>
- **DISC** - DNA Information and Stock Center, Japan  
<http://www.dna.affrc.go.jp>

©2010 Sami Khuri



## SwissProt and SRS

- The **SwissProt** protein sequence database is at ISREC (Swiss Institute for Experimental Cancer Research) in Epalinges, Lausanne.
- The **Sequence Retrieval System** (SRS) at the European Bioinformatics Institute:
  - allows both simple and complex concurrent searches of one or more sequence databases,
  - may also be used on a local machine to assist in the preparation of local sequence databases.

©2010 Sami Khuri



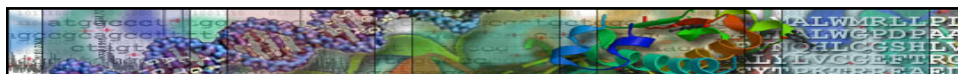
## What does NCBI do?

**NCBI:** established in 1988 as a national resource for molecular biology information.

- it creates public databases,
- it conducts research in computational biology,
- it develops software tools for analyzing genome data, and
- it disseminates biomedical information,

all for the better understanding of molecular processes affecting human health and disease.

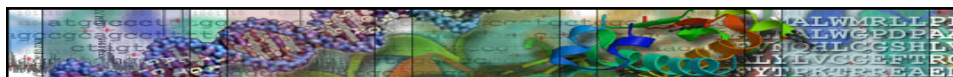
©2010 Sami Khuri



## GenBank

GenBank is the NIH genetic sequence database of all publicly available DNA and derived protein sequences, with annotations describing the biological information these records contain.

©2010 Sami Khuri

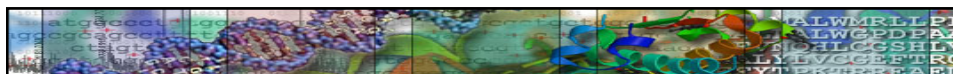


## GenBank DNA Sequence Entry I

The main important fields of an entry are:

- **Locus:** name of locus, length and type of sequence, classification of organism, date of entry. Not maintained among other databases
- **Definition:** description of entry
- **Accession:** accession number of original source. A citable entity; does not change when record is updated.

©2010 Sami Khuri



## GenBank DNA Sequence Entry II

- **Keywords:** keywords for cross referencing this entry
- **Source:** source organism of DNA
- **Organism:** description of organism
- **Comment:** biological function or database information
- **Features:** information about sequence by base position or range of positions

©2010 Sami Khuri



## GenBank DNA Sequence Entry III

- **Base count:** count of A, C, G, T and other symbols
- **Origin:** text indicating the start of the sequence.

The sequence entry is assumed by all computer programs to lie between the identifiers **ORIGIN** and **//**.

©2010 Sami Khuri

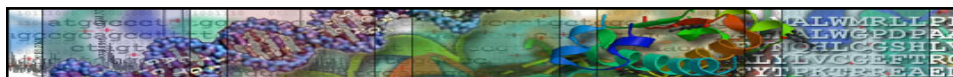


## Division of Organisms

<b>BCT</b> Bacterial	<b>PLN</b> Plant
<b>FUN</b> Fungal	<b>PRI</b> Primate
<b>HUM</b> Homo sapiens	<b>PRO</b> Prokaryotic
<b>INV</b> Invertebrate	<b>ROD</b> Rodent
<b>MAM</b> Other mammalian	<b>SYN</b> Synthetic & chimeric
<b>ORG</b> Organelle	<b>VRL</b> Viral
<b>PHG</b> Phage	<b>VRT</b> Other vertebrate

©2010 Sami Khuri

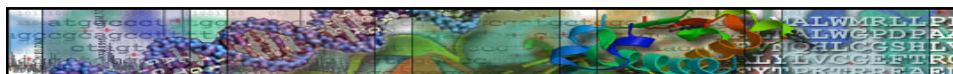




## Interesting Databases

- UCSC Human Genome Browser
  - <http://genome.ucsc.edu/>
- Organism specific information:
  - Yeast: <http://genome-www.stanford.edu/Saccharomyces/>
  - Arabidopsis: <http://www.tair.org/>
  - Mouse: <http://www.jax.org/>
  - Fruit fly: <http://www.fruitfly.org/>
  - Nematode: <http://www.wormbase.org/>

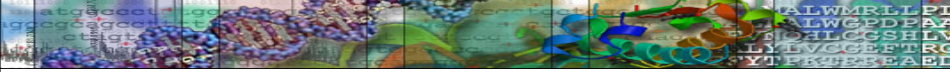
©2010 Sami Khuri



## European Molecular Biology Laboratory

- The **European Molecular Biology Laboratory (EMBL)** was established in 1974.
- It is supported by sixteen countries.
- EMBL consists of five facilities:
  - The main Laboratory in Heidelberg (Germany),
  - Outstations in Hamburg (Germany), Grenoble (France) and Hinxton (the U.K.), and an external Research Programme in Monterotondo (Italy).

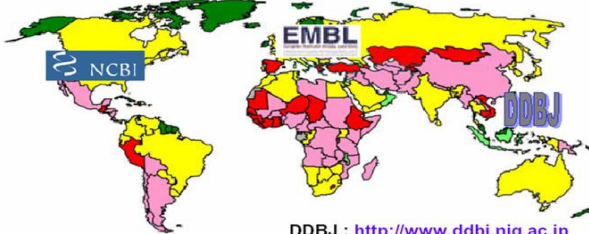
©2010 Sami Khuri



## NCBI – EMBL - DDJB

NCBI : <http://www.ncbi.nlm.nih.gov/>  
NCBI, at the NIH campus, USA

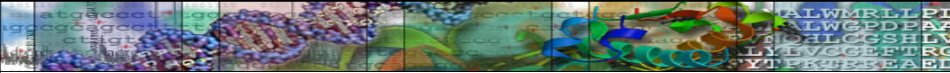
EMBL : <http://www.embl-heidelberg.de/>  
European Molecular Biology Laboratory, UK



DDBJ : <http://www.ddbj.nig.ac.jp>  
DNA Databank of Japan

### Nucleic acid Databases

©2010 Sami Khuri



## Applications of Genome Research

Current and potential applications of Genome Research include:

- Molecular Medicine
- Microbial Genomics
- Risk Assessment
- Bioarcheology, Anthropology, Evolution and Human Migration
- DNA Identification
- Agriculture, Livestock Breeding and Bioprocessing

©2010 Sami Khuri



## Molecular Medicine

- Improve the **diagnosis** of disease
- Detect genetic **predispositions** to disease
- Create drugs **based on molecular information**
- Use **gene therapy** and control systems as drugs
- Design **custom drugs** on individual genetic profiles.

©2010 Sami Khuri



## Microbial Genomics

- Swift detection and treatment in clinics of disease-causing microbes: pathogens
- Development of new energy sources: biofuels
- Monitoring of the environment to detect chemical warfare
- Protection of citizens from biological and chemical warfare
- Efficient and safe clean up of toxic waste.

©2010 Sami Khuri



## DNA Identification I

- Identify potential suspects whose DNA may match evidence left at crime scenes
- Exonerate persons wrongly accused of crimes
- Establish paternity and other family relationships
- Match organ donors with recipients in transplant programs

©2010 Sami Khuri



## Louis XVII



**Louis XVII:** son of Louis XVI and Marie-Antoinette who died from tuberculosis in 1795 at the age of 12

©2010 Sami Khuri



## DNA Identification II

- Identify endangered and protected species as an aid to wildlife officials and also to prosecute poachers
- Detect bacteria and other organisms that may pollute air, water, soil, and food
- Determine pedigree for seed or livestock breeds
- Authenticate consumables such as wine and caviar

©2010 Sami Khuri

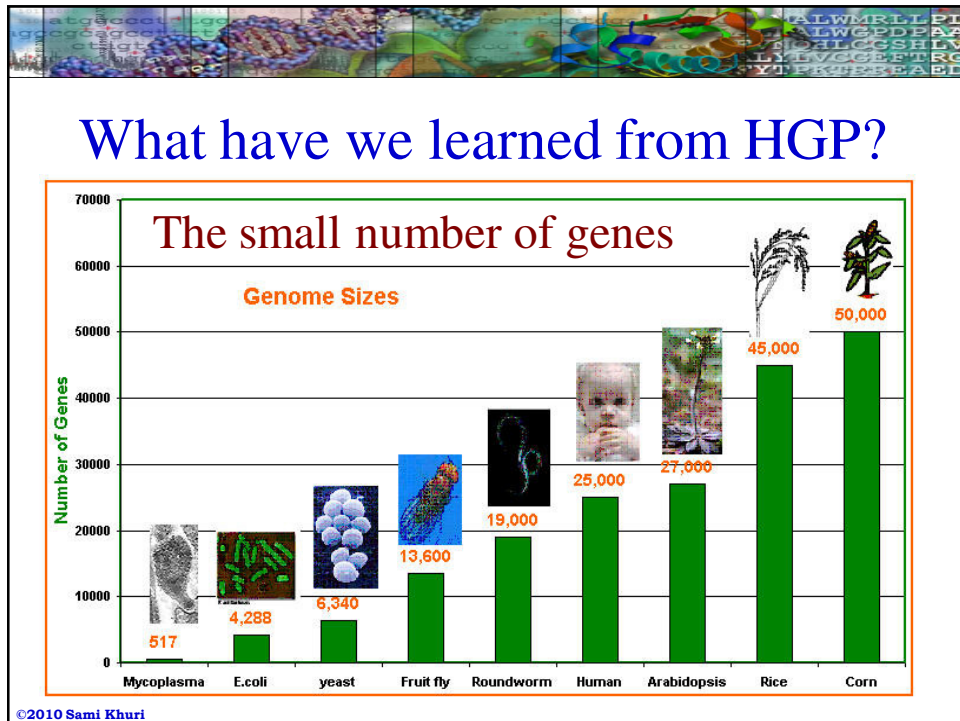
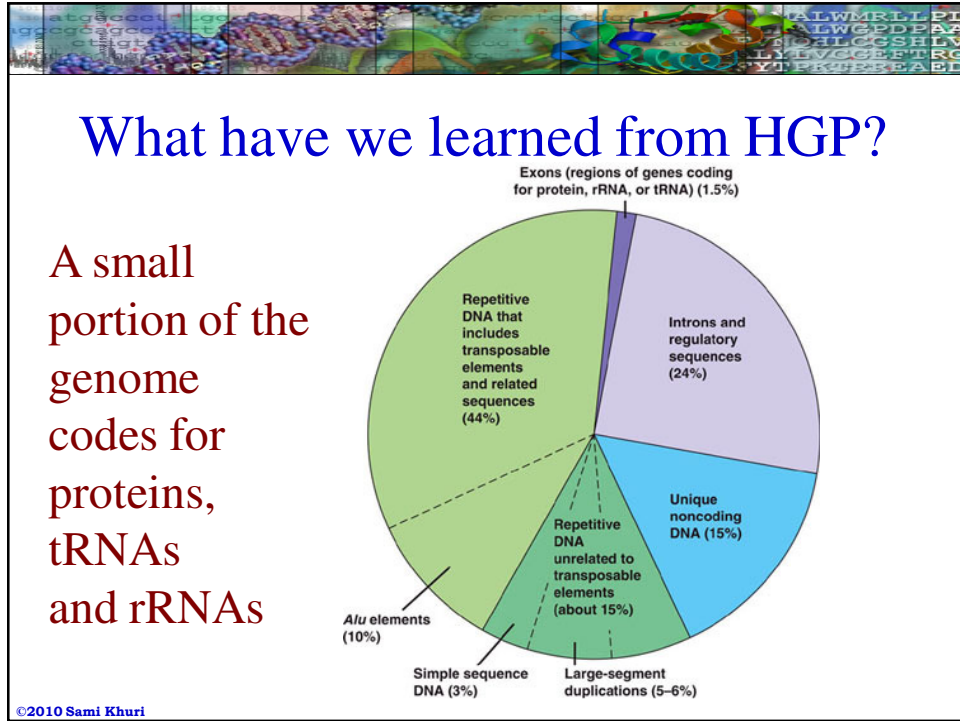


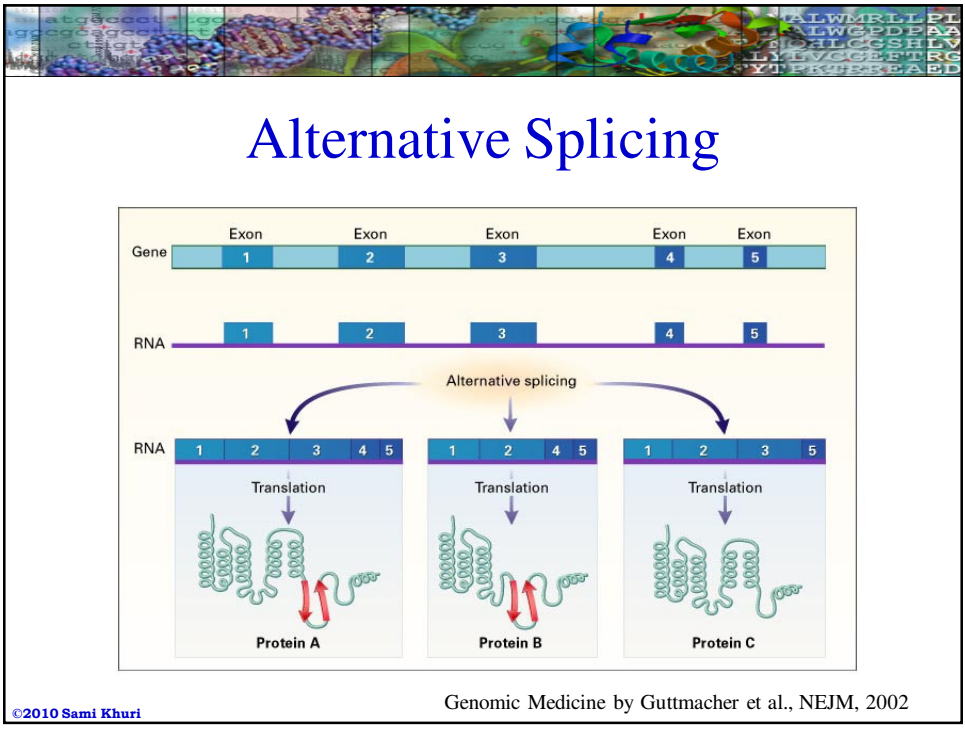
## Agriculture, Livestock Breeding and Bioprocessing

- Grow disease-resistant, insect-resistant, and drought-resistant crops
- Breed healthier, more productive, disease-resistant farm animals
- Grow more nutritious produce
- Develop biopesticides
- Incorporate edible vaccines into food products

©2010 Sami Khuri







Convert all this progress into real riches for science, society, and patients

©2010 Sami Khuri



## Objectives of Molecular Biology

- Extract the information in the genomes.
- Understand the structure of the genome.
- Apply this understanding to the diagnosis and treatment of genetic diseases.
- Explain the process of evolution by comparing genomes of related species.

©2010 Sami Khuri



## Goals of Modern Molecular Biology

- Read the entire genomes of living things
- Identify every gene
- Match each gene with the protein it encodes
- Determine the structure and function of each protein.

©2010 Sami Khuri



## Objectives of Bioinformatics

Development and use of **mathematical** and **computer science** techniques to help solving the problems in molecular biology.

©2010 Sami Khuri



## Bioinformatics Problems

- Reconstructing long DNA sequences from overlapping **string fragments**.
- Comparing two or more sequences for similarities.
- Storing, retrieving and comparing DNA **sequences** and **subsequences** in databases.
- Exploring frequently occurring patterns of nucleotides.
- Finding informative elements in protein and DNA sequences.
- Finding evolutionary relationships between organisms.

©2010 Sami Khuri



## Main Aim of the Problems

- The aim of these problems is to learn about the **functionality** and/or the **structure** of protein without actually having to physically construct the protein itself.
- The research is based on the assumption that similar sequences produce similar proteins.

©2010 Sami Khuri

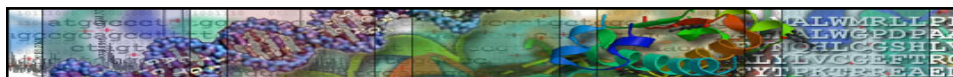


## Functional: Coding v/s Noncoding

	Coding Sequence (Genes)	Non-Coding Sequence
Identifying Computational Tools	Relatively Easy Improving Tools	Very Hard Poor predictive tools
Signals What to look for	We Have a Good Understanding	Very little is known
Complementary data we can use	Available – Ex. ESTs and cDNAs	Unavailable

©2010 Sami Khuri

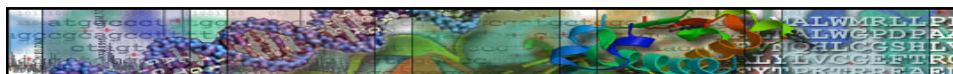




## Post Human Genome Project

- Major role for comparative sequence analysis will be the identification of functionally important, non-coding sequences.
- Need to study the relation between Sequence Conservation and Sequence Function.
- Focus on the interpretation of the human genome.
- Learn the functional landscape of the human genome.
- **Challenge:** go from sequence to function
  - i.e., define the role of each gene and understand how the genome functions as a whole.

©2010 Sami Khuri



## Impact on Other Fields of Science

- Mathematical Models
- Visualization/Animation of molecular structures
- Visualization/Animation of algorithms
- More sophisticated databases
- Data mining
- Statistical analysis of results
- Metabolic pathways
- Predictive algorithms

©2010 Sami Khuri