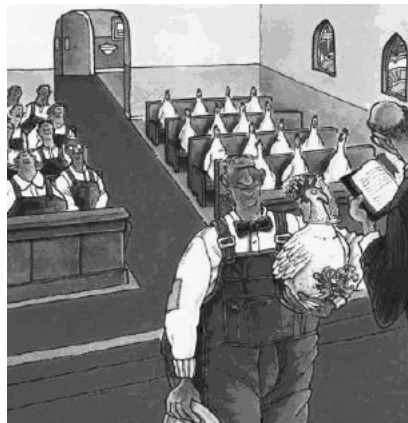# Introduction to BIOINFORMATICS

---

## What is Bioinformatics?

### The marriage between biology and informatics

**Understanding nature's mute elegant language of living cells is the question of modern molecular biology.**

**From an alphabet of only four letters representing the chemical subunits of DNA, emerges a syntax of life processes whose most complex expression is man.**
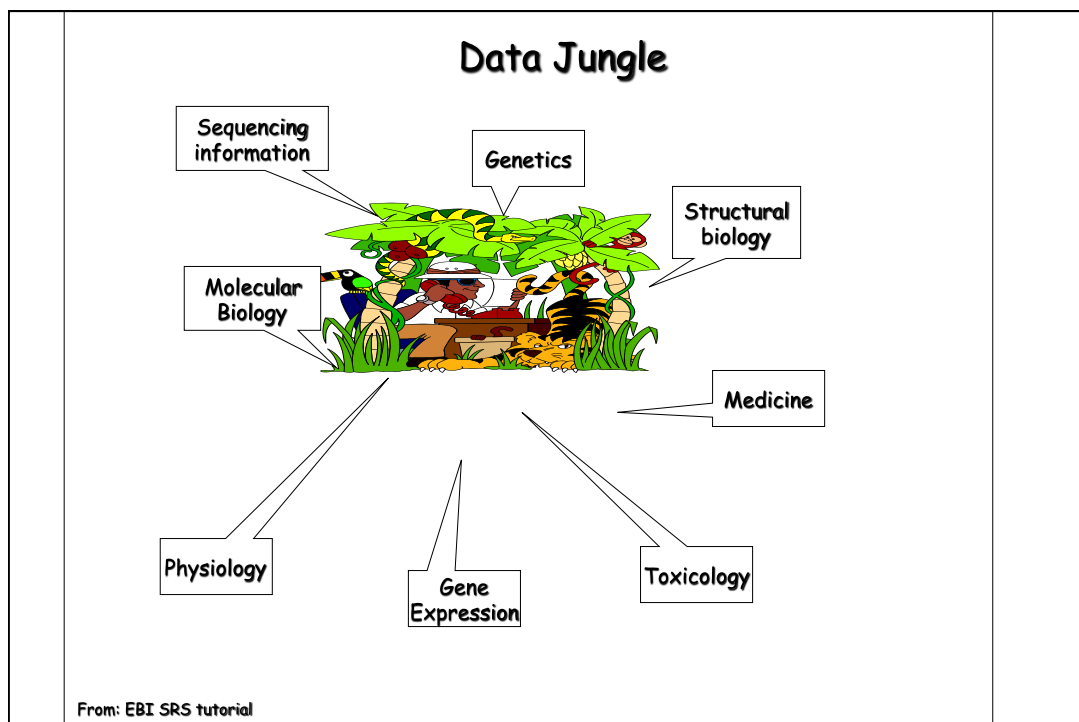
*From the National Centre for Biotechnology Information (NCBI), http://www.ncbi.nlm.nih.gov*

---

**The challenge is in finding new approaches to deal with the volume and complexity of data, and in providing researchers with better access to analysis and computing tools in order to advance understanding of our genetic legacy and its role in health and diseases.**

*From the National Centre for Biotechnology Information (NCBI), http://www.ncbi.nlm.nih.gov*

# Essentially, Bioinformatics has three components

■ **The creation of databases allowing the storage and management of large biological data sets.**

■ **The development of algorithms and statistics to determine relationships among members of large data sets.**

■ **The use of these tools for the analysis and interpretation of various types of biological data, including DNA, RNA and protein sequences, protein structures, gene expression profiles, and biochemical pathways**

---

## Data Jungle



Sequencing information

Genetics

Structural biology

Molecular Biology

Medicine

Physiology

Gene Expression

Toxicology

From: EBI SRS tutorial

## Factors that made bioinformatics so important

- Rapid and cheap techniques for DNA sequencing

- The development of powerful computers

- Internet and the Wide World Web

## DATA SOURCES FOR DATABASES

- Direct scientific submission

- Genome sequencing labs and groups

- Scientific literature

- Patent applications

**DATABASE: a collection of data that has a regular structure and that is organized in such a way that a computer can easily find and retrieve information.**

**A database is generally a collection of RECORDS, available through specific entries, each of which contains one or more FIELD.**

RECORD

ID=Unique identifier

| ID | Locus | Sequence |

**FIELDS**

---

# Data Resources at NCBI

### Databases: Primary and Derivative

*Primary Databases*

- Archival submissions of experimental results
- Database staff organize but don't add additional information

| Genbank | dbEST | dbSNP | Probe |

*Derivative Databases*

- Curated/expert review
- Computationally derived
- Combinations

| Refseq | Genomes | UniGene | UniSTS | Homologene |

Primary vs. Derivative Sequence Databases

# ■ Bioinformatics Developers

## They develop tools for bioinformatics

- Experts in Mathematics, Statistics and Informatics
- Computational biologists

# ■ Bioinformatics Users

## They use the tools of bioinformatics

- Researchers (Biologists, Biotechnologist,…)

# Biological complexity

COMPLEXITY ↑

| Ecological processes & populations |
| Tissue & organism physiology |
| Cellular & developmental processes |
| Biochemical pathways & processes |
| Complete genomes |
| Genes, Proteins, RNA……. |

Introduction to Bioinformatics: www.bioinformatics.com/courses.com/bioinfom

---

Database Retrieval

Sequencing Project Management

Restriction Mapping

Nucleic Acid Sequences

DNA/RNA Folding

Nucleic Acid Sequence Analysis

Seeking Coding regions

Database Retrieval

Translation to amino acids

Protein Sequences

Database Similarity Searching

Protein Sequence analysis

Pair wise Sequence Comparison

Multiple Sequence Alignment

Prediction of Function

Phylogeny

Motifs and Patterns

Structure prediction

Structure analysis

**Sites where the integration among databases
and between databases and software
is developed**

- **USA**
  **NCBI (National Center for Biotechnology Information)**

- **Europe**
  **EBI (European Bioinformatics Institute,Hinxton, UK)**

- **Japan**
  **NIG (National Institute of Genetics)**

# Biological database history

**1965**

M. Dayhoff *et al.* published "Atlas of Protein Sequences and Strucures"

**1982**

EMBL started the DNA sequence collection

**1983**

Genbank started the DNA sequence collection

**1984**

DNA sequence databases of Japan

**1988**

Embl/GenBank/DDBJ agreed on common format for data elements

The International Sequence Database Collaboration

# NCBI: The National Center for Biotechnology Information



Bethesda, Maryland

## Created by US Congress in 1988, NCBI is part of the National Library of Medicine at National Institutes of Health.

www.nih.gov

---

## The National Center for Biotechnology Information (NCBI)

- Created as a part of the *National Library of Medicine* in 1988
  - Establish public databases
  - Research in computational biology
  - Develop software tools for sequence analysis
  - Disseminate biomedical information
- Tools: BLAST(1990), Entrez (1992)
- GenBank (1992)
- Free MEDLINE (PubMed, 1997)
- Other databases: dbEST, dbGSS, dbSTS, MMDB, OMIM, UniGene, GeneMap, Taxonomy, CGAP, SAGE, Gene, RefSeq

http://www.ncbi.nlm.nih.gov/Sitemap/index.html

# Databases connections

http://www.ncbi.nlm.nih.gov/Database/datamodel/index.html

# Entrez databases and the connections between them



http://www.ncbi.nlm.nih.gov/Database/datamodel/index.html