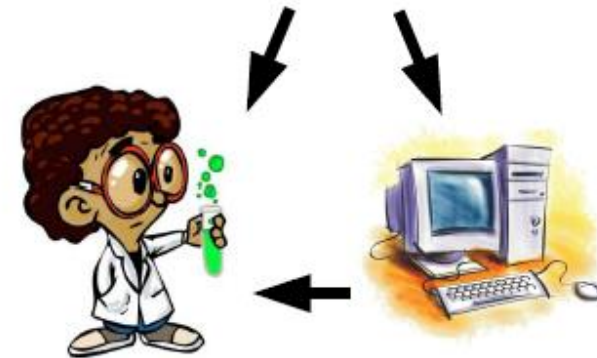# Lecture 2:
# Biological Databases

# NCBI  DataBases

# Database

- A **database** is an organized collection of data .

- The data is typically organized to model relevant aspects of reality (for example, the availability of rooms in hotels), in a way that supports processes requiring this information  (for example, finding a hotel with vacancies).

# Biological databases

- **Make biological data available** ...

  1. ... **to scientists**.

  2. ... **in computer-readable form**.

     - Analysis (computer based)

     - Handle and share large volumes of data

     - Interface for computer based systems (Algorithms, Web interfaces)

- Store data

  - Defined formats

  - Automated storage and retrieval of experimental data
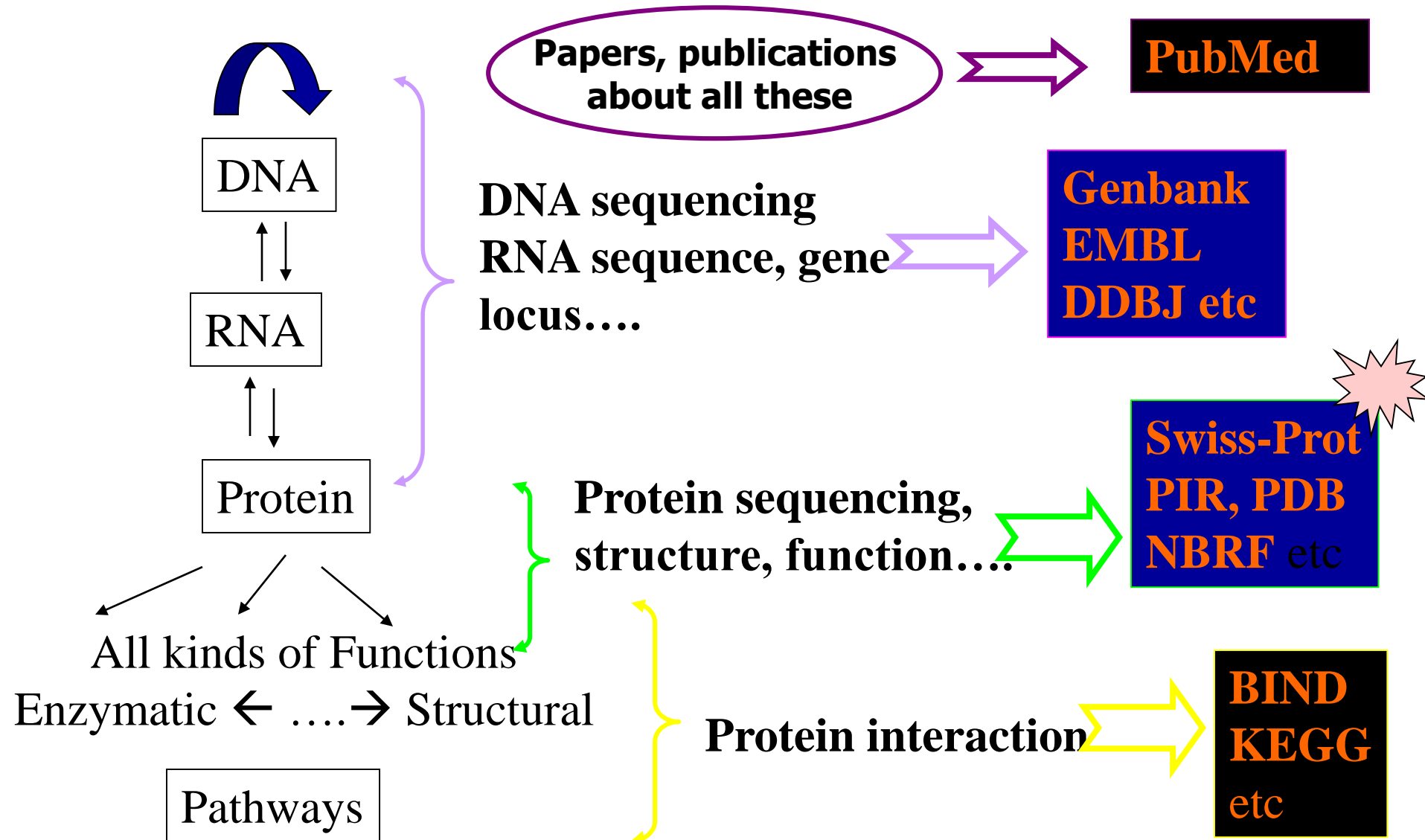
- Link knowledge with external resources

- **A database can be thought of as a large table, where the rows represent records and the columns represent fields.**

| Field  Record | Name | Length | Sequence | Enzyme |
|---|---|---|---|---|
| QA001 | MTGA | 243 | MYQWI… | yes |
| QA002 | Ribosomal protein L9 | 267 | MAAPV… | no |
| QA003 | Flagellin | 374 | GSSIL… | no |
| QA004 | GDPMH | 157 | MFLRQ… | yes |

4

# The "perfect" database

1. Comprehensive, but easy to search.
2. Annotated, but not "too annotated".
3. A simple, easy to understand structure.
4. Cross-referenced.
5. Minimum redundancy.
6. Easy retrieval of data.

# Bioinformatics Databases

DNA → RNA → Protein

Papers, publications about all these ⟹ **PubMed**

DNA sequencing RNA sequence, gene locus…. ⟹ **Genbank EMBL DDBJ** etc

Protein sequencing, structure, function…. ⟹ **Swiss-Prot PIR, PDB NBRF** etc

All kinds of Functions
Enzymatic ← …→ Structural

Pathways

Protein interaction ⟹ **BIND KEGG** etc

6

# 生物信息学数据库简介

**Part 1 outline:**

## 1. Biological information and databases

- **Overview and definition, types of biological databases**

## 2. Popular databases, records, data format

- **PubMed, Genbank, SwissProt, OMIM, PDB, KEGG, BIND, Pfam, PROSITE,**

## 3. Accessing biological databases, retrieval systems

- **Entrez, SRS**

## 4. Searching biological databases

- **Data quality, coverage, redundancy, errors**

# Biological Information

**Nucleic acids:**

- **DNA sequence, genes, gene products (proteins), mutation, gene coding, distribution patterns, motifs**
- **Genomics: genome, gene structure and expression, genetic map, genetic disorder**
- **RNA sequence, secondary structure, 3D structure, interactions**

**Proteins:**

- **Protein sequence, corresponding gene, secondary structure, 3D structure, function, motifs, homology, interactions**
- **Proteomics: expression profile, proteins in disease processes etc.**
- **Ligands and drugs (inhibitors, activators, substrates, metabolites)**

# Biological Information

## Pathways:

- Molecular networks, biological chain events, regulation, feedback, kinetic data

## Function:

- Binding sites, interactions, molecular action (binding, chemical reaction, etc.)
- Biological effect (signaling, transport, feedback, regulation, modification, etc.)
- Functional relationship, protein families, motifs, and homologs

# Biological databases

Lists of biological databases

- INFOBIOGEN Catalog of Databases
  http://www.infobiogen.fr/services/dbcat/

- Nucleic Acids Research Database Listing
  http://www3.oup.co.uk/nar/database/c/

  - These serve as starting point of biological databases.
  - More than 500 databases have been catalogued to date and those from the two listings satisfy minimal criteria for the content, access, and quality.
  - Other sites as a starting point.

# Biological Database List

- **Database Categories List ()**
- [Nucleotide Sequence Databases](#)
- [RNA sequence databases](#)
- [Protein sequence databases](#)
- [Structure Databases](#)
- [Genomics Databases (non-vertebrate)](#)
- [Metabolic and Signaling Pathways](#)
- [Human and other Vertebrate Genomes](#)
- [Human Genes and Diseases](#)
- [Micro−array Data and other Gene Expression Databases](#)
- [Proteomics Resources](#)
- [Other Molecular Biology Databases](#)
- [Organelle databases](#)
- [Plant databases](#)
- [Immunological databases](#)

# 生物信息资源简介: NCBI

# National Center for Biotechnology Information

- Established in 1988

- Creates public databases,

- Conducts research in computational biology,

- Develops software tools for analyzing genome data,

- Disseminates biomedical information.

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

➢Advances science and health by providing access to biomedical and genomic information.

**http://www.ncbi.nlm.nih.gov/Sitemap/AlphaList.html**

## NCBI — Alphabetical Quicklinks Table

| PubMed | Entrez | BLAST | OMIM | Taxonomy | Structure |
|---|---|---|---|---|---|

### ALPHABETICAL QUICKLINKS TABLE

*(To view resource descriptions and a complete list of services, see the NCBI Resource Guide.*
*To view resources by category, see the graphical Site Map.)*

| | | | |
|---|---|---|---|
| About NCBI | Education | Map Viewer | Science Primer |
| Announcements | e-PCR | MeSH | Seminars |
| ASN.1 | Entrez | MGC | Sequin |
| BankIt | Entrez Utilities | Microbial Genomes | Site Search |
| BLAST | Expression | MMDB | SKY/M-FISH & CGH Database |
| BLink | FTP | Model Maker | Software Engineering |
| Books | GenBank | Mutation Databases (external) | Splign |
| Cancer Chromosomes | GenBank sample record | My NCBI (help, tutorial) | Statistics |
| CCDS | Genes | NCBI Home | Structures |
| CDART | Genes and Disease | NCBI News | Submit Data |
| CDD | Genomes (data, projects, submissions) | Nucleotide Sequences (Entrez) | Taxonomy |
| CGAP | GENSAT | OMIM | Tools |
| Clones | GEO (Expression) | OMSSA | TPA |
| Cn3D | Glossary | ORF Finder | Trace Archive |
| Coffee Break | Handbook | Plant Genomes | UniGene |
| COGs | HIV Interactions | Protein Sequences (Entrez) | UniSTS |
| Computational Biology Branch | HTGs | PubChem | VAST |
| Data Submissions | HomoloGene | PubMed | VecScreen |
| dbEST | Human Genome Resources | PubMed Central | Viruses |
| dbGSS | Human-Mouse Homology Maps | RefSeq | WGS |
| dbMHC | Journals | Research at NCBI | What's New |
| dbSNP | LinkOut | Retroviruses | |
| dbSTS | Malaria | SAGEmap | |

**NEW** *indicates a resource which has become available in the last 12 months.*

14

Nucleotides
GenBank
RefSeq (Reference Sequences)
dbEST (Expressed Sequence Tags)
dbGSS (Genome Survey Sequences)
dbMHC (Major Histocompatibility Complex)
dbSNP (Single Nucleotide Polymorphisms)
dbSTS (Sequence Tagged Sites)
TPA (Third Party Annotation Database)
Trace Archive
UniSTS (Sequence Tagged Sites)
PopSet (Evolutionary Relatedness)
UniVec (Vector Sequences)
WGS (Whole Genome Shotgun Sequences)

Proteins
RefSeq (Reference Sequences)
CDD (Conserved Domain Database)

MMDB (Molecular Modeling DataBase)
3D Domains
PubChem BioAssay
PubChem Compound
PubChem Substance

Gene
UniGene
HomoloGene
CCDS (Consensus CoDing Sequence)

GEO (Gene Expression Omnibus)
Entrez GEO Profiles
Entrez GEO DataSets
GENSAT

TaxBrowser
Entrez Taxonomy

Literature Databases
PubMed
PubMed Central
OMIM
Books

Molecular Databases
Nucleotide Sequences
Protein Sequences
Structures
Genes
Gene Expression
Taxonomy

**Databases**

Genomes
Entrez Genome
Entrez Genome Project
Map Viewer
Cancer Chromosomes
SKY/M-FISH & CGH Database

Quer
Entre
Data

NCBI

**Data Submissions**

16

# 这么多数据库如何组织在一起?

- **Entrez** - provides integrated access to nucleotide and protein sequence data from differnt organisms, along with 3D protein structures, genomic mapping information, PubMed MEDLINE, and more.

# ENTREZ

# ENTREZ results

# 1. PubMed
## PubMed is one of the literature databases in the NCBI family.



**In the NCBI family with Nucleotide, Protein, Genome, Structure, OMIM, PMC, Journals, Books and more**

1. includes millions of citations from MEDLINE and other life science journals for biomedical articles back to the 1950s.
2. **PubMed** includes links to full text articles and other related resources.

**Search terms may be topics, authors or journals.**

20

# CONTROLLING CANCER

Cancer is caused by defects that occur when healthy cells divide and die off, a process that occurs millions of times in a lifetime. ONCOGENES direct cells to divide and multiply and TUMOR SUPPRESSOR GENES order old cells to die. CLEANUP GENES correct mistakes made during the cycle. The longer one lives, the likelier these defects will pile up until one or more of the three types of genes goes haywire, and cells multiply out of control. Eventually, they can spread, or metastasize, throughout the body.

**Scientists are pursuing four main avenues to rein in out-of-control cells...**

## BLOCK GROWTH FACTORS

Drugs stop proteins emitted by oncogenes that encourage cells to proliferate wildly.

Examples:
- Iressa (AstraZeneca)
- Erbitux (ImClone)
- Tarceva (Genentech/OSI)

## ENTICE CELLS TO DIE

Drugs reactivate the tumor suppressor genes or block enzymes that protect cancer cells from dying.

Examples:
- Velcade (Millennium)
- Advexin (Introgen)

## MOBILIZE THE IMMUNE SYSTEM

Vaccines and other immunotherapies train the body's immune system to go after cancer cells.

Examples:
- GVAX (Cell Genesys)
- Provenge (Dendreon)
- Oncophage (Antigenics)

## STARVE THE TUMOR

Cancer cells need blood to thrive. So-called anti-angio-genesis drugs stop blood vessels from reaching them.

Examples:
- Avastin (Genentech)
- Neovastat (Aeterna)
- Thalidomide (Celgene)

TUMOR

PROTEINS

DYING CANCER CELLS

METASTASES

ANTIBODY

BLOOD VESSEL

22

# PubMed Search

Cancer treatment by targeting blood supply:

Cancer growth depends on blood supply (why?) and thus requires the growth of new blood vessels – angiogenesis

Proteins involved in angiogenesis may be potential anticancer targets

You can find some of these targets by searching Pubmed

| Key Word | No. of Entries |
|---|---|
| Cancer | |
| Cancer Blood supply | |
| Cancer Blood supply Protein | |
| Cancer Blood supply Enzyme | |
| Cancer Blood supply Enzyme Drug | |

**PubMed**.gov
US National Library of Medicine
National Institutes of Health

[ PubMed ▾ ]   [                                        ]  **Search**

Advanced

Display Settings: ⊙ Abstract                                    Send to: ⊙

## Combined targeting of HER2 and VEGFR2 for effective treatment of HER2-amplified breast cancer brain metastases.

Kodack DP, Chung E, Yamashita H, Incio J, Duyverman AM, Song Y, Farrar CT, Huang Y, Ager E, Kamoun W, Goel S, Snuderl M, Lussiez A, Hiddingh L, Mahmood S, Tannous BA, Eichler AF, Fukumura D, Engelman JA, Jain RK.

Edwin L. Steele Laboratory for Tumor Biology, Department of Radiation Oncology, Massachusetts General Hospital and Harvard Medical School, Boston, MA 02114, USA.

## Abstract
Brain metastases are a serious obstacle in the treatment of patients with human epidermal growth factor receptor-2 (HER2)-amplified breast cancer. Although extracranial disease is controlled with HER2 inhibitors in the majority of patients, brain metastases often develop. Because these brain metastases do not respond to therapy, they are frequently the reason for treatment failure. We developed a mouse model of HER2-amplified breast cancer brain metastasis using an orthotopic xenograft of BT474 cells. As seen in patients, the HER2 inhibitors trastuzumab and lapatinib controlled tumor progression in the breast but failed to contain tumor growth in the brain. We observed that the combination of a HER2 inhibitor with an anti-VEGF receptor-2 (VEGFR2) antibody significantly slows tumor growth in the brain, resulting in a striking survival benefit. This benefit appears largely due to an enhanced antiangiogenic effect: Combination therapy reduced both the total and functional microvascular density in the brain xenografts. In addition, the combination therapy led to a marked increase in necrosis of the brain lesions. Moreover, we observed even better antitumor activity after combining both trastuzumab and lapatinib with the anti-VEGFR2 antibody. This triple-drug combination prolonged the median overall survival fivefold compared with the control-treated group and twofold compared with either two-drug regimen. These findings support the clinical development of this three-drug regimen for the treatment of HER2-amplified breast cancer brain metastases.

**Save items**

☆ Add to Favorites  ▾

**Related citations in PubMed**

Effect of lapatinib on the outgrowth of metastatic brea [J Natl Cancer Inst. 2

# PubMed searching skills

- keyword
- 连词的应用 "and" "or" "not"
- limits
- 同义词...

# My NCBI

- **My NCBI** is a central place to customize NCBI Web services. To use it, you must first register, and your browser must accept cookies.
- You can use **My NCBI** to:
  - Save searches
  - Set up e-mail alerts for new content
  - Display links to Web resources (LinkOut)
  - Choose filters that group search results

# PubMed Central (PMC)

**PubMed Central (PMC)** is the U.S. National Institutes of Health (NIH) free digital archive of biomedical and life sciences journal literature.

| | Find Articles | Advanced search |

**Browse PMC journals:** [A-B] [C-H] [I-M] [N-S] [T-Z] [Full List] [New Journals]

Receive notice of new journals and other major updates to PMC: join the **PMC News mail list** or subscribe to the PMC News **RSS feed** .

All the articles in PMC are free (sometimes on a delayed basis). Some journals go beyond free, to **Open Access**. Find out what that means.

PMC's **utilities** include an OAI service that provides XML of the full-text of some articles, functions for scripting PMC searches and linking to specific PMC articles from your site, and more ...

Looking for a modern journal article DTD? Take a look at NLM's **Journal Publishing XML DTD and schema**.

It's about preservation and access: **digitizing the complete run of back issues** of many of the journals in PMC.

The **PMC journal list** comprises journals that deposit material in PMC on a routine basis and generally make all their published articles available here. Find out how to **include your journal** in PMC.

PMC also has the **author manuscripts** of articles published by NIH-funded researchers in various non-PMC journals. Increasing free access to these articles is the goal of the **NIH Public Access** policy. Similar manuscripts from researchers funded by the Wellcome Trust are available in PMC as well.

Eligible researchers should use the **NIH Manuscript Submission** system to deposit manuscripts.

Get **answers** to other questions about PubMed Central.

27

## 2. OMIM

**OMIM is another literature database in the NCBI family.**

**It is the online version of a catalog of human genes and genetic disorders.**

# Biological databases: OMIM
## Online Mendelian Inheritance in Man
### (http://www.ncbi.nlm.nih.gov/Omim/)

- The OMIM database contains abstracts and texts describing genetic disorders to support genomics efforts and clinical genetics. It provides gene maps, and known disorder maps in tabular listing formats. Contains keyword search.

Hamosh A. *et* al. Online Mendelian Inheritance in Man (OMIM), a knowledge base of human genes and genetic disorders *Nucleic Acids Res*. 2002 30: 52-55.

# Biological databases: OMIM web-page

**National Center for Biotechnology Information**

01101011

OMIM™
Online Mendelian Inheritance in Man

New!
Try
OMIM
in Entrez

## Home Page

Welcome to OMIM(TM), Online Mendelian Inheritance in Man. This database is a catalog of human genes and genetic disorders authored and edited by Dr. Victor A. McKusick and his colleagues at Johns Hopkins and elsewhere, and developed for the World Wide Web by NCBI, the National Center for Biotechnology Information. The database contains textual information, pictures, and reference information. It also contains copious links to NCBI's Entrez database of MEDLINE articles and sequence information.

**NEW** The OMIM Morbid Map, a catalog of genetic diseases and their cytogenetic map locations arranged alphabetically by disease, is now available.

## Browsing OMIM

- Search the OMIM Database
- Search the OMIM Gene Map
- Search the OMIM Morbid Map
- The OMIM numbering system
- View the OMIM Update Log
- OMIM Statistics
- Citing OMIM in the literature
- How to create WWW links to OMIM
- The OMIM Gene List

# Biological databases: OMIM search engine

# Search OMIM

# Nucleic Acids databases

**What info are in these databases:**

- DNA sequence, genes, gene products (proteins), mutation, gene coding, distribution patterns, motifs

- Genomics: genome, gene structure and expression, genetic map, genetic disorder

- RNA sequence, secondary structure, 3D structure, interactions

# GenBank

- The complete release notes for the current version of GenBank are available on the NCBI ftp site.

- A new release is made every two months.

- GenBank is part of the International Nucleotide Sequence Database Collaboration, which comprises the DNA DataBank of Japan (DDBJ), the European Molecular Biology Laboratory (EMBL), and GenBank at NCBI. These three organizations exchange data on a daily basis.

# GenBank

- GenBank® is the NIH genetic sequence database, an annotated collection of all publicly available DNA sequences.

DDBJ/EMBL/GenBank database growth

Note: CON division is not counted in statistics of DDBJ periodical releases.

# Where does large data come from? High-throughput techniques



Fred Sanger

• Nobel prize in chemistry in 1958
"for his work on the **structure** of proteins, especially that of insulin"

• Nobel prize in chemistry in 1980
"for their contributions concerning the determination of **base sequences** in nucleic acids"

**DNA databases:**

**GenBank**
**Web page**

# What might we want to know about a sequence?

- Is this sequence similar to any known genes? How close is the best match?  Significance?
- What do we know about that gene?
  - Genomic (chromosomal location, allelic information, regulatory regions, etc.)
  - Structural (known structure? structural domains? etc.)
  - Functional (molecular, cellular & disease)
- Evolutionary information:
  - Is this gene found in other organisms?
  - What is its taxonomic tree?

# DNA databases

- An Example from GenBank– flat file

  - Human Alpha-Lactalbumin gene

    **This protein is a complex of 2 proteins A and B. In the absence of the B protein, the enzyme catalyzes the transfer of galactose from UDP-galactose to Nacetylglucosamine (cf. EC 2.4.1.90).**

# A GenBank entry – HEADER

# GenBank Entry – Links provided in the Header

- MapViewer – find the gene position in chromosome

- Related Sequences – other entries related to this gene (or sequence)

- OMIM– link to catalog of human genes and genetic disorders

- Protein – retrieve protein record from GenPept

- Medline and PubMed –literature abstracts related to this gene

- Taxonomy – Classification of organisms

- UniGene – Unified gene data

- UniSTS – Unified sequence tagged sites, marker and mapping data

- LinkOut – links to publishers, aggregators libraries, biological databases, sequence centers, and other Web resources

- REFSEQ – reference sequence standards

Note: These links are representative. Other links may also be found in GenBank entries.

# GenBank entry - FEATURES

```
FEATURES              Location/Qualifiers
     source           1..727
                      /organism="Homo sapiens"
                      /db_xref="taxon:9606"
                      /chromosome="12"
                      /map="12q13"
     gene             1..727
                      /gene="LALBA"
                      /db_xref="LocusID:3906"
                      /db_xref="MIM:149750"
     prim_transcript  1..727
                      /gene="LALBA"
     CDS              27..455
                      /gene="LALBA"
                      /EC_number="2.4.1.22"
                      /codon_start=1
                      /product="lactalbumin, alpha-"
                      /protein_id="NP_002280.1"
                      /db_xref="GI:4504947"
                      /db_xref="LocusID:3906"
                      /db_xref="MIM:149750"
                      /translation="MRFFVPLFLVGILFPAILAKQFTKCELSQLLKDIDGYGGIALPE
                      LICTMFHTSGYDTQAIVENNESTEYGLFQISNKLWCKSSQVPQSRNICDISCDKFLDD
                      DITDDIMCAKKILDIKGIDYWLAHKALCTEKLEQWLCEKL"
     sig_peptide      27..83
                      /gene="LALBA"
     misc_feature     84..440
                      /gene="LALBA"
                      /note="LYZ1; Region: Alpha-lactalbumin / lysozyme C"
                      /db_xref="CDD:LYZ1"
     misc_feature     84..440
```

3

# GenBank Entry– Links provided in the Feature section

LocusID – locus and display of genomic and mRNA sequences

MIM – Link to OMIM description, other entries for this sequence

EC_number – link to the corresponding cataloged enzymes

Protein_id – retrieve protein record from GenPept

CD– conserved protein domain (SMART),

CDD – conserved protein domain (Pfam).

# Biological databases: GenBank – SEQUENCE

```
BASE COUNT       175 a      183 c      163 g      206 t
ORIGIN
        1 atttcaggtt cttgggggta gccaaaatga ggttctttgt ccctctgttc ctggtgggca
       61 tcctgttccc tgccatcctg gccaagcaat tcacaaaatg tgagctgtcc cagctgctga
      121 aagacataga tggttatgga ggcatcgctt tgcctgaatt gatctgtacc atgtttcaca
      181 ccagtggtta tgacacacaa gccatagttg aaaacaatga aagcacggaa tatggactct
      241 tccagatcag taataagctt tggtgcaaga gcagccaggt ccctcagtca aggaacatct
      301 gtgacatctc ctgtgacaag ttcctggatg atgacattac tgatgacata atgtgtgcca
      361 agaagatcct ggatattaaa ggaattgact actggttggc ccataaagcc ctctgcactg
      421 agaagctgga acagtggctt tgtgagaagt tgtgagtgtc tgctgtcctt ggcacccctg
      481 cccactccac actcctggaa tacctcttcc ctaatgccac ctcagtttgt ttctttctgt
      541 tcccccaaag cttatctgtc tctgagcctt gggccctgta gtgacatcac cgaattcttg
      601 aagactattt tccagggatg cctgagtggt gcactgagct ctagacsctt actcagtgcc
      661 ttcgatggca ctttcactac agcacagatt tcacctctgt cttgaataaa ggtcccactt
      721 tgaagtc
//
```

# GenBank - NOTES

Majority of GenBank entries have similar form to our example.

When accessing the database, the following needs to be noticed:

- Some entries are huge, containing as much as 30,000 lines. (NT_021877 Homo sapiens chromosome 1 working draft sequence segment)

- Some entries have contig information instead of sequence information. (NT_021877 Homo sapiens chromosome 1 working draft sequence segment)

- Some entries are derived from cDNA sequences and thus represent putative genes/proteins. These should be used with caution.

- Some annotations are predicted using automated analysis. These should also be used with caution. (XM_131483 Mus musculus simi...[gi:20832685]).

# RefSeq

- RefSeq: sub-collection of NCBI databases with only non-redundant, highly annotated entries (genomic DNA, transcript (RNA), and protein products)



47

Search  All Databases  for  [            ]  Go

- [Brief Description](#)
- [Scope](#)
- [Announcements](#)
- [Access and Availability](#)
- [Distinguishing Features](#)
- [References](#)

## NCBI Reference Sequences

The Reference Sequence (RefSeq) collection aims to provide a comprehensive, integrated, non-redundant, well-annotated set of sequences, including genomic DNA, transcripts, and proteins. RefSeq is a foundation for medical, functional, and diversity studies; they provide a stable reference for genome annotation, gene identification and characterization, mutation and polymorphism analysis (especially RefSeqGene records), expression studies, and comparative analyses. [more...]

### ▶ Scope    ⬆

NCBI provides RefSeqs for taxonomically diverse organisms including eukaryotes, bacteria, and viruses. Additional records are added to the collection as data become publicly available.

> **May 12, 2011: RefSeq Release 47 available for FTP**
>
> This release includes:
>
> **Proteins:**    12,625,466
> **Organisms:**  12,000
> **Available at:** ftp://ftp.ncbi.nih.gov/refseq/release/
>
> To receive announcements of future RefSeq releases and incremental large updates please subscribe to NCBI's refseq-announce mail list: refseq-announce

### Site contents

**Information**

NCBI Handbook
Overview | FAQ ⍰
Accessions | Status |
Queries | Publications

**FTP**

RefSeq Release
Catalog | Notes
Genomes
BLAST databases

**Statistics**

Release Statistics

**Feedback**

NCBI Help Desk
Submit Updates
Submit GeneRIF

**Subscribe - eMail Lists**

RefSeq | Gene
Map Viewer | NCBI

**Related links**

Genomic Biology Home
Gene | Genome Project
Entrez Genomes Home
Map Viewer | UniGene

**Credits**

Collaborators
Microbial Providers
Viral Genome Advisors
NCBI Staff

# The RefSeq Accession number format and molecule types

| Accession | Molecule type |
|---|---|
| NC_xxxxxx | Complete genomic molecule |
| NG_xxxxxx | Genomic region |
| NM_xxxxxx | mRNA |
| NP_xxxxxx | Protein |
| NR_xxxxxx | RNA |
| NT_xxxxxx | computed Genomic contig |
| XM_xxxxxx | computed mRNA |
| XP_xxxxxx | computed Protein |

# Biological Databases

## Database Searching

1.  Most of the databases have a web-interface to search for data
2.  Databases must have methods for accessing and extracting data stored.
3.  The most basic search is keyword searching
    Keywords can be any word that occurs somewhere in the database
    records. It can be the name of the gene or protein (e.g. lactalbumin),
    species (e.g.*homo sapiens*, human), a taxonomy term
    (e.g.primates), or a word from the reference title (e.g. cancer)
4.  Others include: Entry Id number, sequence
5.   User can choose to view the data or save to your computer
6.  Databases typically have hyperlinks that help to navigate from one database to another easily

# Summary

- what is a biological database?

- Why need bioinformatics database?

- Different types of bioinformatics database

- NCBI database
  - Pubmed
  - OMIM
  - GenBank

# Homework

- What's the latest amount of data for PubMed, OMIM and GenBank database?

- Explore NCBI database, choose 2 other database you are interested to explore details. Give a summary of them.

# Gene