# The Importance of Biological Databases in Biological Discovery

助理教授: 張學偉

高雄醫學大學　生物醫學暨環境生物學系
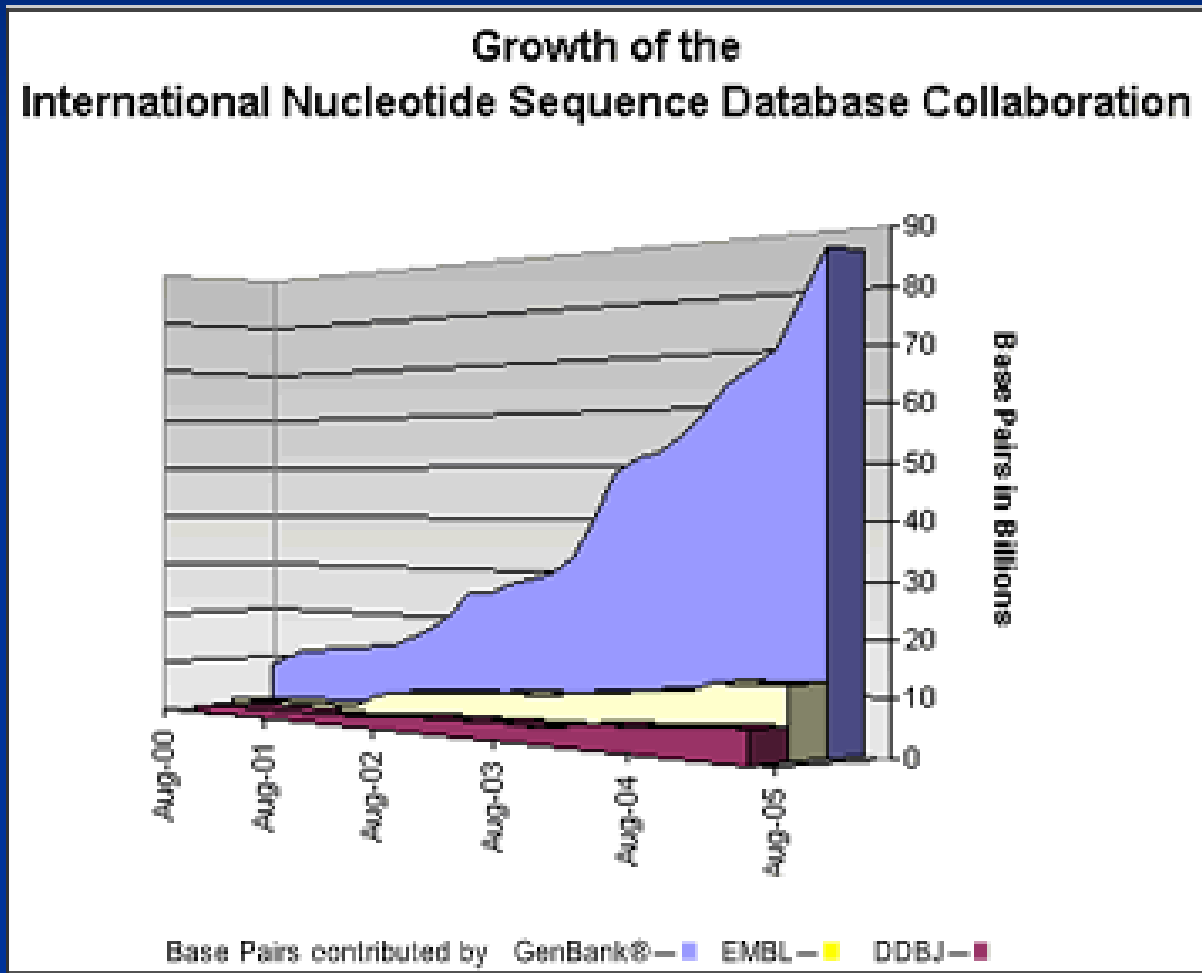
# The database that most biologists are familiar with is GenBank.

■ **GenBank**

-is the annotated collection of all publicly available *DNA* and *protein* sequences.

-is maintained by NCBI (National Center for Biotechnology) at NIH (National Institute of Health).

- -represents a collaborative effort between NCBI, <u>EMBL (European Molecular Biology Laboratory)</u> and <u>DDBJ (DNA Data Bank of Japan)</u>



http://www.ncbi.nlm.nih.gov/Genbank/

- GenBank® is the NIH genetic sequence database, an annotated collection of all publicly available DNA sequences (*Nucleic Acids Research 2007 Jan 1;35(Database issue):D16-20*).

- There are approximately 65,369,091,950 bases in 61,132,599 sequence records in the *traditional GenBank divisions* and 80,369,977,826 bases in 17,960,667 sequence records in the *WGS division* as of August 2006.

# The growth of GenBank

- Human Genome Project
- Other systematic sequencing projects
  → accumulation of sequence data

# NCBI & Boolean Search

- ## Entrez homepage of NCBI

is not a database itself but rather the interface through which all of its component databases can be accessed and traversed.

**Boolean logic**
**From Wikipedia, the free encyclopedia**

# KEGG (Kyoto Encyclopedoa of Genes and Genomes)

- a widely-used compendium of biochemical pathways that has practical use in the modeling of expression data and in understanding higher-order cellular processes.

## 3. Useful links & tools- Glossary

- **BioInformatics Glossary**;

- **Genome Glossary NCBI**;

- **-Omes and -omics glossary**

# Biological databases

- Like any other database
  - Data organization for optimal analysis

- Data is of different types
  - Raw data (DNA, RNA, protein sequences)
  - Curated data (DNA, RNA and protein annotated sequences and structures, expression data)

# Characteristics of biological data

- **Complex** ➡ Thoughtful data modeling

  Data types range from sequences, 3-dimensional structures, pathways, images, text, and a wide variety of annotation.

- **Heterogeneous** ➡ Universal schema

  storage format, management, and access vary widely

- **Dynamic** ➡ Flexible designing

  contents and schema change routinely and rapidly (twice/year)

- **Inconsistent** ➡ Ontology

  lack standards at the ontology level

  - Controlled vocabulary for consistent naming for biomedical terms within and between databases
  - Data models for modeling or abstraction of biological system and processes

# The growth of public domain bio-databases



**(The Molecular Biology Database Collection from *Nucleic Acids Research*)**

# Three main tendencies of bio-databases

- **Database proliferation**
  - Dozens to hundreds at the moment

- **In the next 5 years biological data analysis will be trifurcated**
  - Bio-webs : remote data analysis and mining
  - Bio-grids : transparent high-end computing
  - Bio-semantic webs : biological knowledge

- **More and more scientific discoveries result from inter-database analysis and mining**

# Cross-references

Nucleotide seq db
EMBL

PTM
GlycoSuiteDB
PhosSite

Families, domains, sites
HAMAP
InterPro
Pfam
PRINTS
ProDom
PROSITE
SMART
TIGRFAMs

Organism-specific dbs
DictyDb
dbSNP
EcoGene
FlyBase
GeneDB_Spombe
Genew
Gramene
HIV
Leproma
ListiList
MaizeDB
MGD
MypuList
OMIM
SagaList
SGD
StyGene
SubtiList
TIGR
TubercuList
WormPep
ZFIN

**Swiss-Prot**

Explicit links

2D-gel electrophoresis
ANU-2DPAGE
Aarhus/Ghent-2DPAGE
COMPLUYEAST-2DPAGE
ECO2DBASE
HSC-2DPAGE
MAIZE-2DPAGE
PHCI-2DPAGE
PMMA-2DPAGE
Siena-2DPAGE
SWISS-2DPAGE

Miscellaneous dbs
MEROPS
PIR
REBASE
TRANSFAC

Structural dbs
HSSP
PDB

swissprot

# Gene Ontology database-1

"The Gene Ontology (GO) project seeks to provide a set of *structured vocabularies* for *specific biological domains* that can be used to describe gene products in any organism."

**A few key points:**
GO is a "structured" vocabulary, which is really a specialized type of a "controlled" vocabulary.
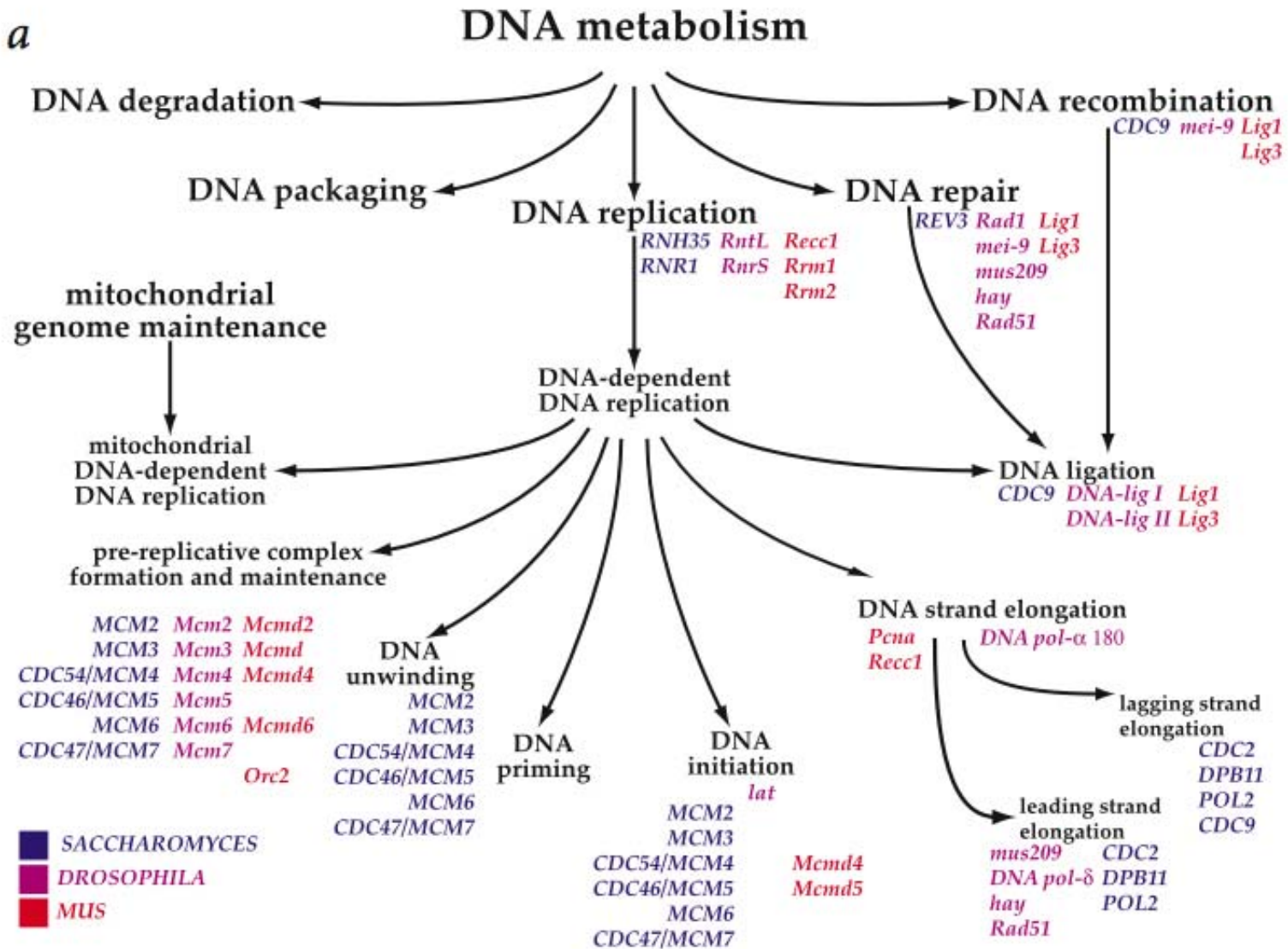
# Gene Ontology database-2

The ontologies in GO are intended to describe three biological areas, "molecular function", "biological processes" and "cellular components".

GO was *originally developed* through the collaboration of the members of *three model organism projects*: SGD, the *Saccharomyces* Genome database; FlyBase, the *Drosophila* genome database; and MGD/GXD, the Mouse Genome Informatics databases.

# Biological Process Ontology
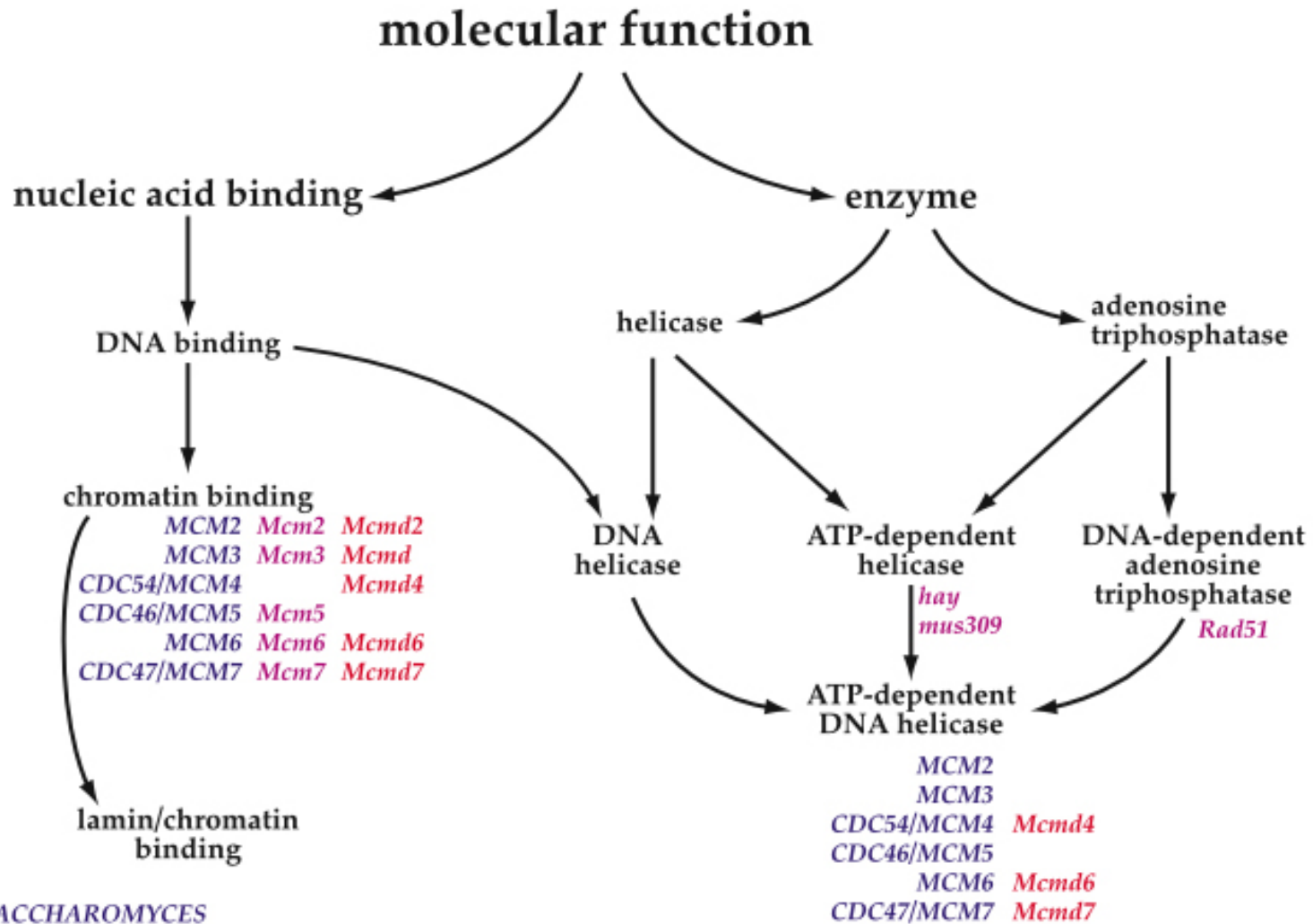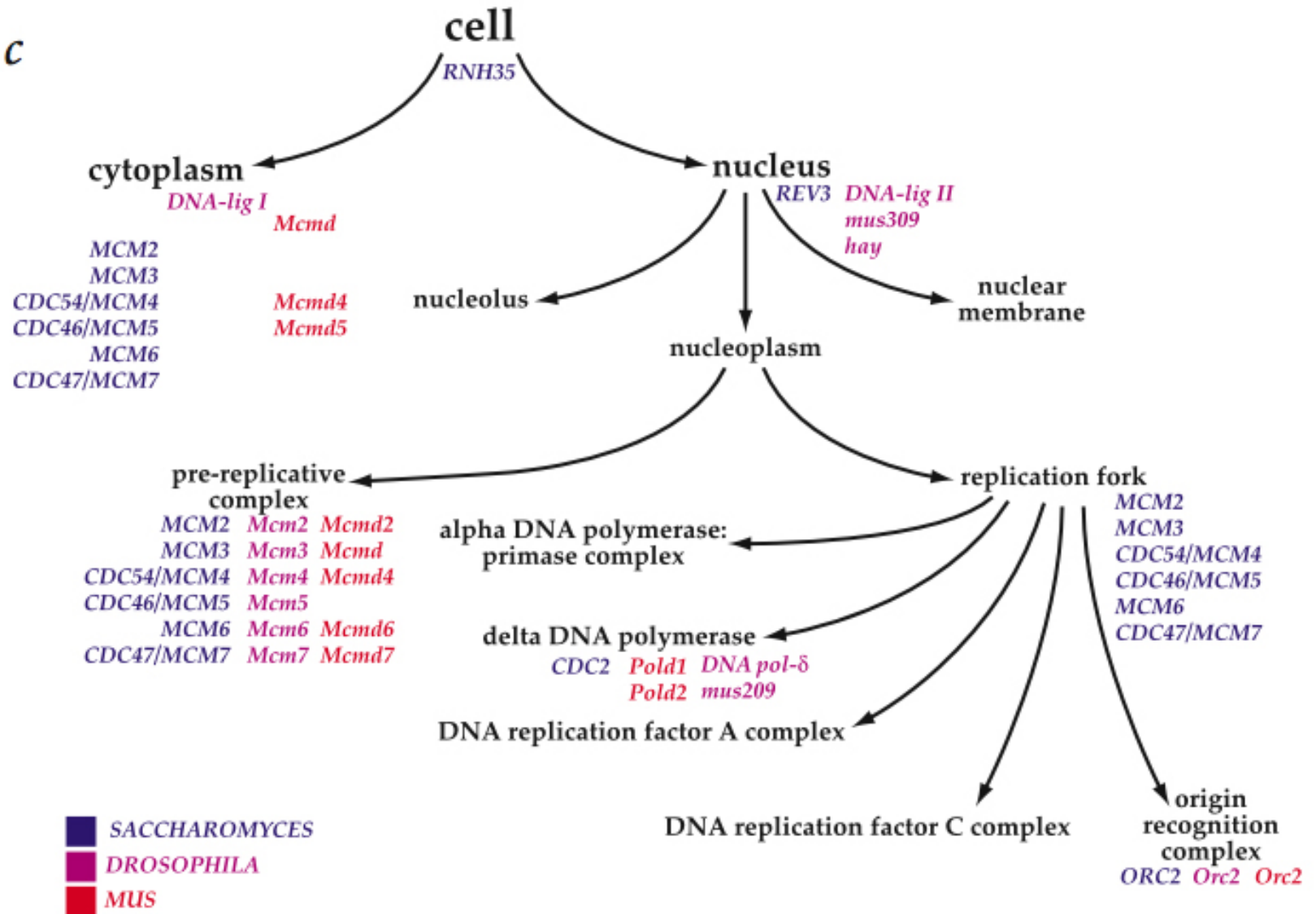
# Molecular Function Ontology

# Cellular Component Ontology

# The CANCER GENOME ANATOMY PROJECT

- The <u>Gene Ontology Browser</u> (GO Browser) classifies human and mouse genes by molecular function, biological process, and cellular component.

# What GO is Not

1. GO is *not a way* to unify biological databases. Sharing nomenclature is a step toward unification, but is not, in itself, sufficient.

2. GO is *not* a dictated standard, mandating nomenclature across databases. Groups participate because of self-interest and cooperate to arrive at a consensus.

3. GO does *not* define homologies between gene products from different organisms. The use of the GO results in shared annotations for gene products from different organisms, and this may reflect an evolutionary relationship, but the shared annotation is in itself not sufficient for such a determination.
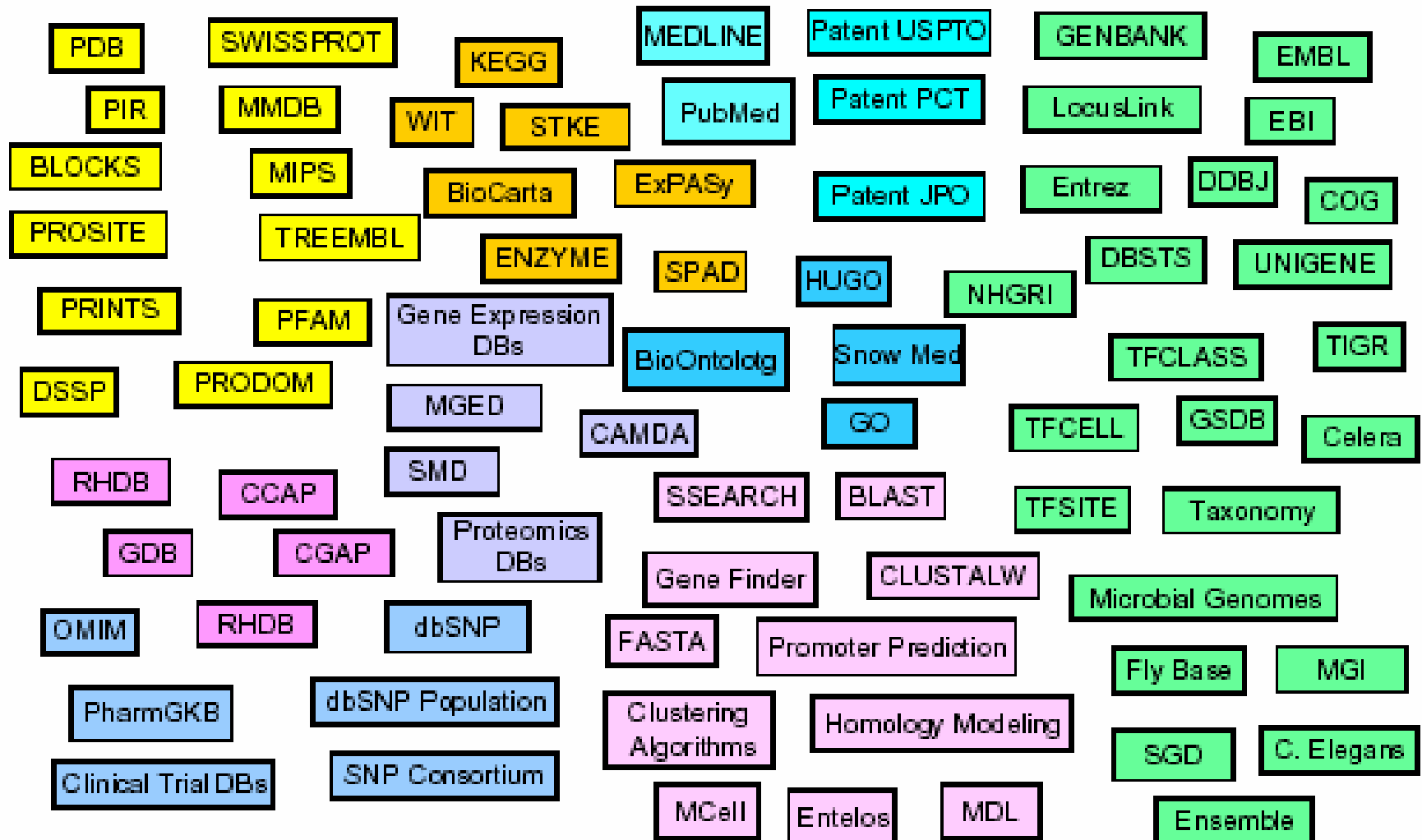
# Knowledge discovery in databases

- **Data mining** is a technique to discover hidden information in large databases. This information, e.g. trends and patterns, can be used to build predictive models.

- Example: extracting predictive information of gene expression from genome sequence databases.

# Bio-databases:
# a short word on problems

- Even today we face some key limitations
  - **There is no standard format**
    - Every database or program has its own format
  - **There is no standard nomenclature**
    - Every database has its own names
  - **Data is not fully optimized**
    - Some datasets have missing information without indications of it
  - **Data errors**
    - Data is sometimes of poor quality, erroneous, misspelled

# Swimming in Data Sources

# The Molecular Biology Database Collection: 2007 update

Michael Y. Galperin*

National Center for Biotechnology Information, National Library of Medicine,
National Institutes of Health, Bethesda, MD 20894, USA

## ABSTRACT

The NAR online Molecular Biology Database Collection is a public resource that contains links to the databases described in this issue of *Nucleic Acids Research*, previous NAR database issues, as well as a selection of other molecular biology databases that are freely available on the web and might be useful to the molecular biologist. The 2007 update includes 968 databases, 110 more than the previous one. Many databases that have been described in earlier issues of NAR come with updated summaries, which reflect recent progress and, in some instances, an expanded scope of these databases. The complete database list and summaries are available online on the *Nucleic Acids Research* web site http://nar.oxfordjournals.org/.

## NAR 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
Microarray Data and other Gene Expression Databases
Proteomics Resources
Other Molecular Biology Databases
Organelle databases
Plant databases
Immunological databases

http://nar.oxfordjournals.org/cgi/reprint/35/suppl_1/D3  Free download

# Nucleic Acids Research



Database & server issues

# 4. Bioinformatics Journal list

**BMC Bioinformatics**

**BMC Genomics**

**BMC Medical Informatics and Decision Making**

**ALGORITHMS FOR MOLECULAR BIOLOGY**