

Digital information cycle

BINF 630: Bioinformatics Methods

Iosif Vaisman

Email: ivaisman@gmu.edu

Creation and capture
Storage and management
Rights management
Search and access
Distribution

Electronic publishing

Quality (peer review, retrospective evaluation)
Reliability (stability of serves, control over alterations, proper archiving and mirroring)

Hypertext Functionality in Electronic Publishing

Active references
Forwarding references
Dynamic publishing

Ethical, Legal, and Economical Issues of Electronic Publishing

Intellectual property rights
Ownership of information
Information as a commodity

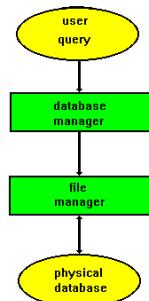
Database



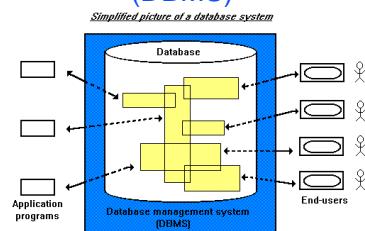
Database

database	a collection of related structured information about entities
file	a collection of records
record	a set of fields
field	a single characteristic of an entity
character	a symbol used in data field

Database Organization



Database Management System (DBMS)



Four major components of DBMS:

Data * Hardware * Software * Users

Data Model

- A named logical unit (record type, data item)
 - Relationships among logical units

Relationships among logical units

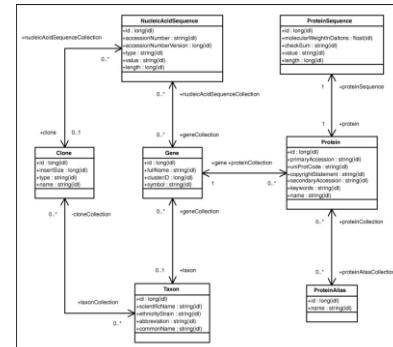
- one to one
 - one to many
 - many to one

Relational Database Model

protab1			
Protein-code	Protein-name	Length	Species-origin
P1001	Hemoglobin	145	Bovine
P1002	Hemoglobin	136	Ovine
P1003	Eye Lens Protein	234	Human
.....			

protab2	
Protein-code	Protein-sequence
P1001	MDRTHGFDLKLSPRTVNQWLMALFFGH...
P1002	MDKTSHGFEIKLLTPKKLQQWLMAIYFGHT...
P1003	SRTHEEEGKLMQWPPLRPLYIALFTPEPPY...
.....	

Object-oriented Database Model

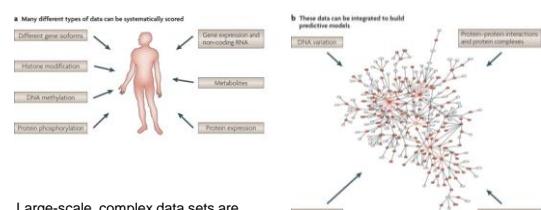


G A Komatsoulis et al., 2008

Database administration

- Redundancy eliminated
 - Inconsistency avoided
 - Data shared
 - Standards enforced
 - Security applied
 - Integrity maintained
 - Requirements balanced

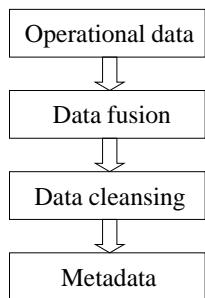
Data integration



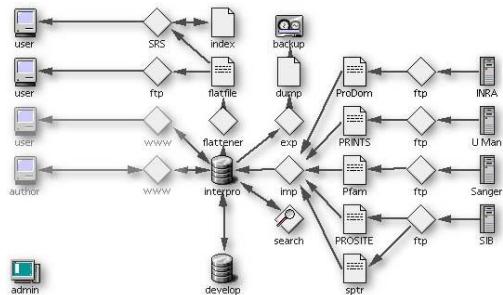
Large-scale, complex data sets are shown as a network in which the nodes represent variables of biological interest, such as DNA variation, RNA variation, protein levels, protein states, metabolite levels and disease-associated traits, and the edges between these nodes represent causal relationships between the variables.

Schmidt et al. 2010

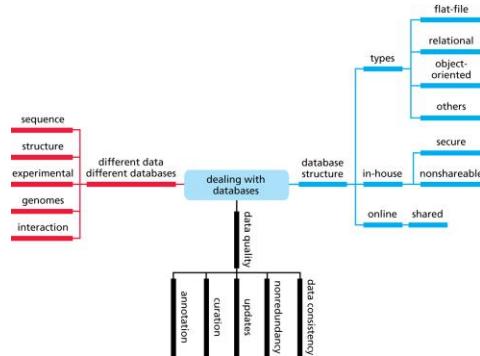
Data Warehouse



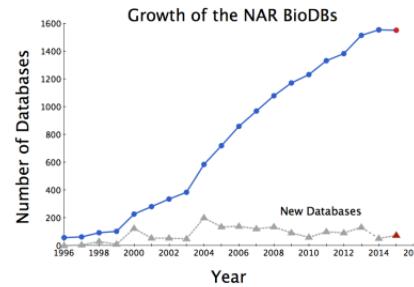
InterPro Database



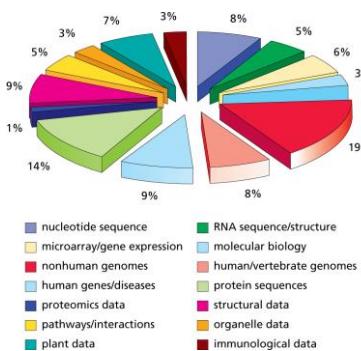
Databases in Bioinformatics



NAR Molecular Databases

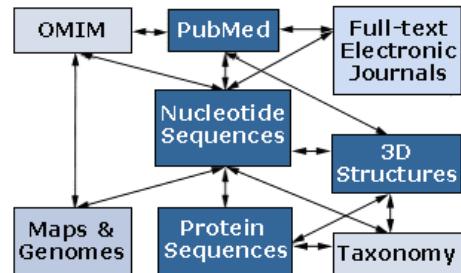


Molecular Databases



NCBI integrated search and retrieval system

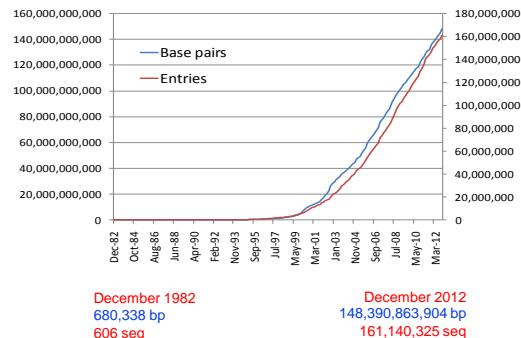
Entrez



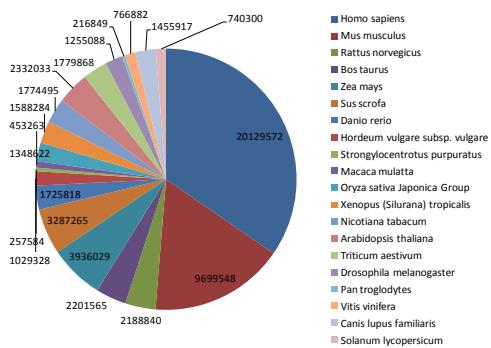
NCBI Databases

- nr** - All non-redundant GenBank CDS translations+PDB+SwissProt+PIR
- month** - All new or revised GenBank CDS released in the last 30 days
- swissprot** - the last major release of the SWISS-PROT protein sequence database (no updates)
- yeast** - Yeast (*Saccharomyces cerevisiae*) protein sequences.
- E. coli** - *E. coli* genomic CDS translations
- pdb** - Sequences derived from the 3-dimensional structure Brookhaven Protein Data Bank
- kabat** - Kabat's database of sequences of immunological interest

Growth of GenBank



GenBank Selected Per-Organism Statistics



NCBI Databases

Table 1. The Entrez databases (as of 3 September 2013)

Database	Records	Section within this article	Data source
NCBI Web Site	21 929	Introduction	N
PubMed	23 952 796	Literature	C
PMC	2 836 592	Literature	D, C
NLM Catalog	1 483 089	Literature	C, N
MeSH	243 770	Literature	N
Books	225 312	Literature	C, N
Bioencyclopedia ^a	1 113 795	Taxonomy	C, N
Nucleotide ^a	101 599 766	DNA and RNA	D (GenBank), C, N
EST ^a	74 911 096	DNA and RNA	D (GenBank)
GS ^a	36 959 049	DNA and RNA	D (GenBank)
BioSample	2 100 817	DNA and RNA	N
ERA ^a	475 312	DNA and RNA	N
PopSet ^a	183 110	DNA and RNA	D (GenBank)
Protein ^a	94 102 424	Proteins	C, N
Protein Clusters ^a	382 691	Proteins	N
GEO Profiles ^a	91 193 791	Genes and expression	D
Gene ^a	31 367 498	Genes and expression	D
Gene ^a	14 167 800	Genes and expression	C, N
UniGene ^a	6 467 085	Genes and expression	N
Gene3D ^a	1 044 344	Genes and expression	N
Proteomes ^a	522 730	Genes and expression	C
Homologene ^a	133 548	Genes and expression	N
Clone ^a	33 135 797	Genomes	D, N
UnSTS ^a	545 913	Genomes	D (dhSTS)
BioProject ^a	98 358	Genomes	D
Assembly	17 707	Genomes	C, N
Genome ^a	10 929	Genomes	C, N

NCBI Databases (cont.)

McGenEpigenomics ^a	10 811	Genomics	D
SNP ^a	300 258 943	Genetics and medicine	D (dbSNP), N
dbVar ^a	3 584 019	Genetics and medicine	D
MedGen ^a	169 433	Genetics and medicine	C, N
dbGap ^a	154 971	Genetics and medicine	D
ChrVar ^a	49 000	Genetics and medicine	D, N
PubMed Health	41 262	Genetics and medicine	C
GTR ^a	29 212	Genetics and medicine	D
OMIA	2844	Genetics and medicine	C
PubChem Substance ^a	11 983 13 846	Chemicals and bioassays	D
PubChem Compound ^a	47 757 896	Chemicals and bioassays	N
PubChem BioAssay ^a	717 429	Chemicals and bioassays	D
Structure ^a	92 993	Domains and structures	C, N
CDD ^a	48 034	Domains and structures	C, N

^aIndicates that the data in this resource are available by FTP. D, direct submission; C, collaboration/agreement; N, internal NCBI/NLM curation.

Derivative databases

Systems	Year	Major features
Ranking search results		
ReMed	2010	Featuring multi-level relevance feedback for ranking
Quarter	2009	Allowing multiple term concepts in queries
MedlineRanker	2009	Ranking relevant documents through classification
MISearch	2009	Using implicit feedback for improving ranking
Haka	2008	Powered by Haka's proprietary semantic technology
SemanticMEDLINE	2008	Powered by cognition's proprietary search technology
MISeARCH	2008	MISeARCH
eBlast	2008	Finding relevant documents through classification
PubCrouse	2007	Ranking documents similar to input terms
Tweeze	2006	Sorting by impact factor and citation volume
Genome	2005	Query expansion with relevance ranking technique

Clustering results into topics

Medline^a, OtaTe

2008 Clustering by important words, topics, journals, authors, etc.

2007 Clustering by Medline or LMSK concepts

2005 Clustering by Medline or GO terms

2004 Clustering by Medline, title/abstract, author, affiliation, or date

2001 Clustering by extracted keywords from abstracts

2005 Clustering by extracted keywords from abstracts

2008 Clustering by important words, topics, journals, authors, etc.

2007 Clustering by Medline or LMSK concepts

2005 Clustering by Medline or GO terms

2004 Clustering by Medline, title/abstract, author, affiliation, or date

2001 Clustering by extracted keywords from abstracts

2005 Clustering by extracted keywords from abstracts

2008 Clustering by important words, topics, journals, authors, etc.

2007 Clustering by Medline or LMSK concepts

2005 Clustering by Medline or GO terms

2004 Clustering by Medline, title/abstract, author, affiliation, or date

2001 Clustering by extracted keywords from abstracts

2005 Clustering by extracted keywords from abstracts

2008 Clustering by important words, topics, journals, authors, etc.

2007 Clustering by Medline or LMSK concepts

2005 Clustering by Medline or GO terms

2004 Clustering by Medline, title/abstract, author, affiliation, or date

2001 Clustering by extracted keywords from abstracts

2005 Clustering by extracted keywords from abstracts

2008 Clustering by important words, topics, journals, authors, etc.

2007 Clustering by Medline or LMSK concepts

2005 Clustering by Medline or GO terms

2004 Clustering by Medline, title/abstract, author, affiliation, or date

2001 Clustering by extracted keywords from abstracts

2005 Clustering by extracted keywords from abstracts

2008 Clustering by important words, topics, journals, authors, etc.

2007 Clustering by Medline or LMSK concepts

2005 Clustering by Medline or GO terms

2004 Clustering by Medline, title/abstract, author, affiliation, or date

2001 Clustering by extracted keywords from abstracts

2005 Clustering by extracted keywords from abstracts

2008 Clustering by important words, topics, journals, authors, etc.

2007 Clustering by Medline or LMSK concepts

2005 Clustering by Medline or GO terms

2004 Clustering by Medline, title/abstract, author, affiliation, or date

2001 Clustering by extracted keywords from abstracts

2005 Clustering by extracted keywords from abstracts

2008 Clustering by important words, topics, journals, authors, etc.

2007 Clustering by Medline or LMSK concepts

2005 Clustering by Medline or GO terms

2004 Clustering by Medline, title/abstract, author, affiliation, or date

2001 Clustering by extracted keywords from abstracts

2005 Clustering by extracted keywords from abstracts

2008 Clustering by important words, topics, journals, authors, etc.

2007 Clustering by Medline or LMSK concepts

2005 Clustering by Medline or GO terms

2004 Clustering by Medline, title/abstract, author, affiliation, or date

2001 Clustering by extracted keywords from abstracts

2005 Clustering by extracted keywords from abstracts

2008 Clustering by important words, topics, journals, authors, etc.

2007 Clustering by Medline or LMSK concepts

2005 Clustering by Medline or GO terms

2004 Clustering by Medline, title/abstract, author, affiliation, or date

2001 Clustering by extracted keywords from abstracts

2005 Clustering by extracted keywords from abstracts

2008 Clustering by important words, topics, journals, authors, etc.

2007 Clustering by Medline or LMSK concepts

2005 Clustering by Medline or GO terms

2004 Clustering by Medline, title/abstract, author, affiliation, or date

2001 Clustering by extracted keywords from abstracts

2005 Clustering by extracted keywords from abstracts

2008 Clustering by important words, topics, journals, authors, etc.

2007 Clustering by Medline or LMSK concepts

2005 Clustering by Medline or GO terms

2004 Clustering by Medline, title/abstract, author, affiliation, or date

2001 Clustering by extracted keywords from abstracts

2005 Clustering by extracted keywords from abstracts

2008 Clustering by important words, topics, journals, authors, etc.

2007 Clustering by Medline or LMSK concepts

2005 Clustering by Medline or GO terms

2004 Clustering by Medline, title/abstract, author, affiliation, or date

2001 Clustering by extracted keywords from abstracts

2005 Clustering by extracted keywords from abstracts

2008 Clustering by important words, topics, journals, authors, etc.

2007 Clustering by Medline or LMSK concepts

2005 Clustering by Medline or GO terms

2004 Clustering by Medline, title/abstract, author, affiliation, or date

2001 Clustering by extracted keywords from abstracts

2005 Clustering by extracted keywords from abstracts

2008 Clustering by important words, topics, journals, authors, etc.

2007 Clustering by Medline or LMSK concepts

2005 Clustering by Medline or GO terms

2004 Clustering by Medline, title/abstract, author, affiliation, or date

2001 Clustering by extracted keywords from abstracts

2005 Clustering by extracted keywords from abstracts

2008 Clustering by important words, topics, journals, authors, etc.

2007 Clustering by Medline or LMSK concepts

2005 Clustering by Medline or GO terms

2004 Clustering by Medline, title/abstract, author, affiliation, or date

2001 Clustering by extracted keywords from abstracts

2005 Clustering by extracted keywords from abstracts

2008 Clustering by important words, topics, journals, authors, etc.

2007 Clustering by Medline or LMSK concepts

2005 Clustering by Medline or GO terms

2004 Clustering by Medline, title/abstract, author, affiliation, or date

2001 Clustering by extracted keywords from abstracts

2005 Clustering by extracted keywords from abstracts

2008 Clustering by important words, topics, journals, authors, etc.

2007 Clustering by Medline or LMSK concepts

2005 Clustering by Medline or GO terms

2004 Clustering by Medline, title/abstract, author, affiliation, or date

2001 Clustering by extracted keywords from abstracts

2005 Clustering by extracted keywords from abstracts

2008 Clustering by important words, topics, journals, authors, etc.

2007 Clustering by Medline or LMSK concepts

2005 Clustering by Medline or GO terms

2004 Clustering by Medline, title/abstract, author, affiliation, or date

2001 Clustering by extracted keywords from abstracts

2005 Clustering by extracted keywords from abstracts

2008 Clustering by important words, topics, journals, authors, etc.

2007 Clustering by Medline or LMSK concepts

2005 Clustering by Medline or GO terms

2004 Clustering by Medline, title/abstract, author, affiliation, or date

2001 Clustering by extracted keywords from abstracts

2005 Clustering by extracted keywords from abstracts

2008 Clustering by important words, topics, journals, authors, etc.

2007 Clustering by Medline or LMSK concepts

2005 Clustering by Medline or GO terms

2004 Clustering by Medline, title/abstract, author, affiliation, or date

2001 Clustering by extracted keywords from abstracts

2005 Clustering by extracted keywords from abstracts

2008 Clustering by important words, topics, journals, authors, etc.

2007 Clustering by Medline or LMSK concepts

2005 Clustering by Medline or GO terms

2004 Clustering by Medline, title/abstract, author, affiliation, or date

2001 Clustering by extracted keywords from abstracts

2005 Clustering by extracted keywords from abstracts

2008 Clustering by important words, topics, journals, authors, etc.

2007 Clustering by Medline or LMSK concepts

2005 Clustering by Medline or GO terms

2004 Clustering by Medline, title/abstract, author, affiliation, or date

2001 Clustering by extracted keywords from abstracts

2005 Clustering by extracted keywords from abstracts

2008 Clustering by important words, topics, journals, authors, etc.

2007 Clustering by Medline or LMSK concepts

2005 Clustering by Medline or GO terms

2004 Clustering by Medline, title/abstract, author, affiliation, or date

2001 Clustering by extracted keywords from abstracts

2005 Clustering by extracted keywords from abstracts

2008 Clustering by important words, topics, journals, authors, etc.

Example of a Genbank entry

```
LOCUS    VIBHALUXA    3141 bp    DNA      BCT      15-FEB-1996
DEFINITION V.harveyi luciferase alpha and beta subunit (luxA and luxB) genes, complete cds.
ACCESSION M10961 M13494
NID      g155174
KEYWORDS luciferase.
SOURCE   Vibrio harveyi DNA.
ORGANISM Vibrio harveyi
          Eubacteria; Proteobacteria; gamma subdivision; Vibrionaceae; Vibrion.
REFERENCE 1 (bases 1 to 1838)
AUTHORS Cohn,D.H., Mileham,A.J., Simon,M.I., Nealson,K.H., Rausch,S.K., Bonam,D. and Baldwin,T.O.
TITLE    Nucleotide sequence of the luxA gene of Vibrio harveyi and the complete amino acid sequence of the alpha subunit of bacterial luciferase
JOURNAL  J. Biol. Chem. 260 (10), 6139-6146 (1985)
MEDLINE  85207595
REFERENCE 2 (bases 1745 to 3141)
AUTHORS Johnston,T.C., Thompson,R.B. and Baldwin,T.O.
TITLE    Nucleotide sequence of the luxB gene of Vibrio harveyi and the complete amino acid sequence of the beta subunit of bacterial luciferase
JOURNAL  J. Biol. Chem. 261 (11), 4805-4811 (1986)
MEDLINE  86168191
```

Example of a Genbank entry

```
FEATURES          Location/Qualifiers
gene             707..1774
                 /gene="luxA"
CDS              707..1774
                 /gene="luxA"
                 /codon_start=1
                 /product="luciferase alpha subunit"
                 /db_xref="PID:g155175"
                 /transl_table=11
                 /translation="MKFGNFLTYQPPELSQTEVMKRLVNLGKASEGCGFDVTWLLHE
HFTFGLLGNPYVAAAHLLGATETLNVGTAAIVLPTAHHPVRQAEDVNLLQMSKGRFR
FGICRGLYIKDFRVFGTMDMNSRALMCWYILMKEGFNEGYIAADNEHIKFPKIQJNPF
SAYTQQGAPVYVAESEASSTTEWAERGLPMILSWIINTHEKAQLDLYNEVATEHGYD
VTKIDHCLSYISITSVHDSRANKDICRNFLGHWMYDSYVNATKIFDDSDQTGYDNKKQ
WRDFVLKGHKDTNRKRIDSYEINPVGTPPEECIAIIQODIDATGIDINICCGFEANGSEE
EIIASMKLFQSDVMVPPYLKEKQ"
BASE COUNT      883 a   665 c   741 g   852 t
ORIGIN          1 bp upstream of EcoRI site.
                 1 gatgttccac ttatccatgg ctttgtcactg tttatctcg gttgtcagacc
                 61 aaggccacac aaaaacattt ctggatggcg gcaatgttcct cagctgtgtt cgcctatggaa
                 121 gtatctttcg atttggatgt gtttttttcg attatcgccg ttggtggtgt gaacctgggt
                 181 gacacactag aaaaacgtcg tggtttgtat taectcagg tgcctatcg tgatgtaccaa
                 ...
...
```

Example of a Pubmed entry

```
PMDID- 15887224
GWID - N/A
STARCDLINE
DA - 20050718
DCOM- 20060522
IS - 1097-0134 (Electronic)
IS - 0887-3585 (Linking)
DP - 2005 Aug 15
TI - A descriptor for protein secondary structure assignment based on a simple topological descriptor.
PG - 513-24
AB - A simple, five-element descriptor, derived from the Delaunay tessellation of a protein structure in a single point per residue representation, can be assigned to each residue in protein. This descriptor characterizes mainly local topology and accessibility in the neighborhood of the residue and does not explicitly depend on putative hydrogen bonds or any geometric parameter, including bond length, angles, and areas. Rules based on this descriptor can be used for accurate, robust, and computationally efficient secondary structure assignment that correlates well with the existing methods.
AD - Laboratory for Structural Bioinformatics, School of Computational Sciences, George Mason University, Manassas, Virginia 20110, USA.
FAU - Taylor, Todd
AU - Taylor T
FAU - Rivera, Margarita
AU - Rivera M
FAU - Wilson, Glenda
AU - Wilson G
FAU - Vaisman, Iosif I
AU - Vaisman II
LA - eng
PT - Journal Article
PT - Research Support, U.S. Gov't, Non-P.H.S.
PL - United States
```

Example of a Pubmed entry

```
TA - Proteins
JT - Proteins
JID - 8700181
RN - 0 (Proteins)
SP -
MH - Algorithms
MH - Amino Acid Sequence
MH - Biophysical Phenomena
MH - Biophysics
MH - Computational Biology/*methods
MH - Computer Simulation
MH - Hydrogen Bonding
MH - Models, Chemical
MH - Models, Molecular
MH - Molecular Sequence Data
MH - Protein Conformation
MH - Protein Folding
MH - Protein Structure, Secondary
MH - Protein Structure, Tertiary
MH - Proteins/chemistry
MH - Proteins/physiology/*methods
MH - Sensitivity and Specificity
MH - Software
MH - Structural Homology, Protein
EDAT- 2005/05/12 09:00
MHDA- 2006/05/23 09:00
CRDT- 2005/05/12 09:00
```