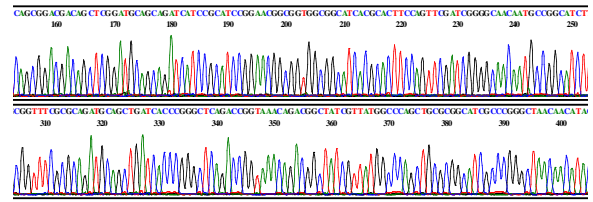


Molekulárně biologické databáze

Pro zajímavost, nebude součástí zkoušky...

Důležité, pravděpodobně bude u zkoušky...

Molekulárně biologická data



GATAGCGTAATGATCGGCTGGCTCCGCATTTGATGTTGGTCCCAAGAAAATAACCGCTCACGGTCCCATCACGATCCGACACCGGAAAATCGCGG
 TACAGTGGTCCGCCCCCGCCAGCACATCCCTGCCCCAATAAGATCTTTAGCGGACGACGCTCGGATGAGCAGATCATCCGCATCCGGAACGGC
 AGCGCGTTCCGCGAGATCGACGATCACCGGGCTCAGACCGGTAACAGACGGTATGTTATGGTCCAGTGGCGGATGCCCGGGTAAACAATA
 CATAACGGTGGGACCAATCAATCACGGTCCGGCCGGCGGATCACCGCTGGTCTCCGGTAGGGCTCCAGCAGGGTAAACCGCATCCAGAAATCACCGCAT

Molekulárně biologická data

MALDI-TOF → Identifikace proteinů

↓
 Sekvence proteinů

MDRNGNFGSLPPNTAFKAIIFYANAADRQDLK
 LPIDDAPEPAATFVGNSEDEGVRLFTLNSKG
 GKIRIEASANGRQSATDARLAPLSAGDTVW
 LGWLGAEDEGADADYNDGIVLQWPIT

Molekulárně biologická data

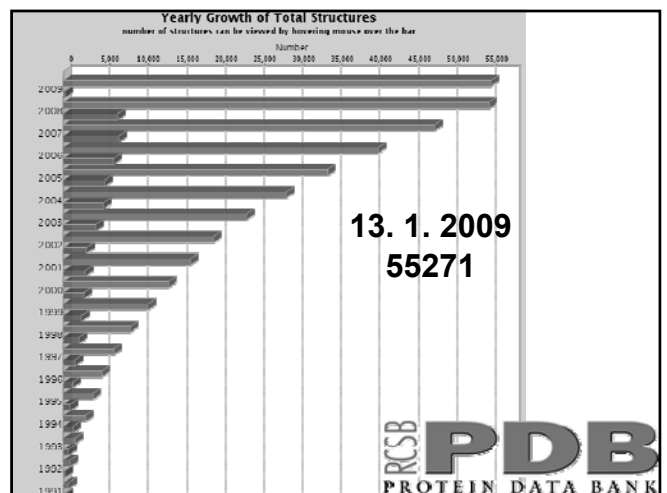
Proteinová krystalografie

NMR spektroskopie

Molekulárně biologická data


- Výkonné technologie:
 Automatické sekvencování
 MALDI-TOF
 NMR spektroskopie
 Proteinová krystalografie

Výrazný nárůst množství biologických dat.



Éra reverzní genetiky

Klasická genetik



↓

```
GATAGCGTAAATGATCGGCTGGCTGCCCATTTCC
TACAGGTGGTCCGCCGCCGCCAGCACATCGC
GGTGGCGGCATCACGCCACTTCCAGTGGATCGG
AGCGCGGTTTCCGCCGAGATCGAGTGAATCACCC
```

Fenotyp → **Genotyp**


Reverzní genetik
Automatické DNA sekvencování

↓

Produkcce velkého množství dat

```
GATAGCGTAAATGATCGGCTGGCTGCCCATTTCC
TACAGGTGGTCCGCCGCCGCCAGCACATCGC
```

↓



Genotyp → **Fenotyp**

Bi7201 Základy genomiky, podzimní semestr

Molekulárně biologická data

- Nutnost organizovaného ukládání a skladování dat.
- Nutnost prohlížení a analyzování uložených dat.

Databáze je určitá uspořádaná množina informací (dat) uložená na paměťovém médiu.

V širším smyslu jsou součástí databáze i softwarové prostředky, které umožňují manipulaci s uloženými daty a přístup k nim.

Analytické nástroje

- Vyhledávací software
Nutnost snadného, rychlého a specifického vyhledání informací.
- Srovnávání dat (sekvencí)
Sequence alignment – „seřazení“ sekvencí.

```

LFFNTAKAIFYANAADRDEKLFIDAFEAATFVQNSDGVLL--PFLNKGKIKIIE
LFFNTDRALFFANAARQDHLKLFIDGSEPEAAHYKLTTRDGPFE--ATLNSGNKINFE
LFFNIXGVTAIHAANDQIIDLVIDDPPKAAATPKGAGAQDQNLGKVLDDGNHVVVI
LFFNIAEGVTALVHSSAPQTIIEWPVDHNRKAAATPQAGTQDANLNTQIVNKGKIKVIVV
LFFN-aFg---lanesad-QtikifidD-pFAATfkgag-----l-t-clnDgnokiRve

ASANGRCATDALLAPLSAGD-----TVNLEWLGASDASABVNDGIVILQWFIT
VSVNGKPRATDALLAPINOKKSDGSPFTVNFQIVVSDHHSBYNDGIVVILQWFIT
VMAHGRPRRLGSRQVDIFKKE-----YFSLGSRDHAHDVNDGIVVFNWFLG
VFAHGRPKIKISGQVDIFKKE-----YFSLGSRDHAHDVNDGIVVFNWFLG
VFAHGRPKIKISGQVDIFKKE-----YFSLGSRDHAHDVNDGIVVFNWFLG
vsaNGrpSat--R---ifkks-----tvyfGivgsEDGaDaDYNDGIViLqNFIg
    
```

Rozdělení molekulárně biologických databází

- Databáze:
 - Primární
 - Sekundární
 - Strukturální

```


EDRPIKFSTEGATSQSYKQFIEALRERLRGLLHDIPVLPDPTTLQERNRYIT
VELNSDTESEIEVGLDVNAYVAYRAGTQSYFLRDAPSSASDYLFTGTDQHS
LFFYGTGDLERWAHQSRQQLPLGLQALTHGISFFRSGDNEEKARTLIVII
QMVAAARFRYISNRVVSIGTAFQPDAAIISLENNWNLRSRGVQESVQDT
FFNQVLTINRNEPVIIVDSLSHPTVAVLALMLFVCPNPNIVEKSKICSRYP
TVRIIGRRDGMVDVYDNGYHNGNRIIMMKCKDRLEENQLWTLKSDKTIKRSNGK
    
```

↓

Ribosome-inactivating protein, subdomain 1

Ribosome-inactivating protein, subdomain 2

Ricin B-like lectins



Rozdělení molekulárně biologických databází

- Databáze:
 - Primární
 - Sekundární
 - Strukturální

Primární databáze obsahují anotované sekvence NA nebo proteinů.

Rozdělení molekulárně biologických databází

- Databáze:
 - Primární
 - Sekundární
 - Strukturální

Sekundární databáze obsahují informace odvozené z primárních databází ve formě charakteristických vzorů sekvencí, tj. funkčních nebo strukturálních motivů získaných srovnáním primárních dat (sekvencí).

Rozdělení molekulárně biologických databází

- Databáze:
- Primární
- Sekundární
- Strukturní

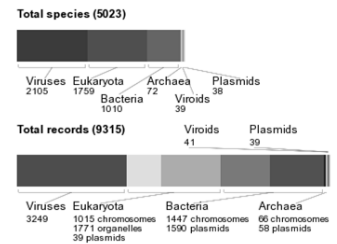


Obsahují struktury proteinů (nukleových kyselin) a jejich anotace.



Rozdělení molekulárně biologických databází

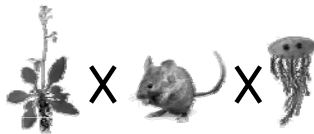
- Databáze:
- Primární
- Sekundární
- Strukturní



Genomové zdroje

Rozdělení molekulárně biologických databází

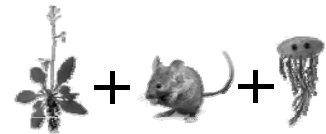
- Databáze:
- Specializované
- Univerzální



Specializované databáze obsahují informace o určité proteinové rodině nebo skupině proteinů, případně o určitém organismu.

Rozdělení molekulárně biologických databází

- Databáze:
- Specializované
- Univerzální



Univerzální databáze obsahují informace o proteinech (NA) ze všech organismů.

Rozdělení univerzálních proteinových databází

- Univerzální databáze:
- „Skladiště“ sekvencí – sequence repository
- „Manuálně“ spravovaná – curated database

Rozdělení univerzálních proteinových databází

- „Skladiště“ sekvencí – sequence repository
- Kromě sekvencí obsahují málo nebo žádné dodatečné informace.
- Záznamy generovány automaticky.
- Proteiny mohou být zastoupeny několika různými záznamy (sekvencemi) = „nadbytečnost“ (redundance) sekvencí.

Rozdělení univerzálních proteinových databází

- **Manuálně spravované – curated databases**
Záznamy obsahují dodatečné informace.
Informace jsou před vložením do databáze validovány experty.
Všechny záznamy o stejné proteinové sekvenci jsou sdružovány do jediného = non-redundant dataset.

Rozdělení molekulárně biologických databází

- **Databáze:**
Primární
Sekundární
Strukturní

Genomové zdroje
Složené databáze


Složené databáze

- **Složené (composite) databáze:**
Slučují data z několika primárních databází.
Eliminace redundantních dat.
Různá priorita zdrojových databází podle kvality validace a anotace (eliminace redundantních dat z databáze s nižší prioritou).

Molekulárně biologické databáze

Nucleic Acids Research

<http://www3.oup.co.uk/nar/database/a/>



Nucleic Acids Research
1999-2009

Nucleotide Sequence Databases
International Nucleotide Sequence Database Collaboration
Coding and non-coding DNA
Gene structure, introns and exons, splice sites
Transcriptional regulator sites and transcription factors
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

1170 databází

EBI/NCBI/CIB

Instituce zabývající se shromažďováním, správou a poskytováním dat a informací a vývojem analytických nástrojů.

EBI

Evropský institut pro bioinformatiku



European Bioinformatics Institute

<http://www.ebi.ac.uk/>

NCBI

Národní centrum pro biotechnologické informace



National Center for Biotechnology Information

<http://www.ncbi.nlm.nih.gov/>

CIB

Centrum pro informační biologii



Center for Information Biology

<http://www.cib.nig.ac.jp/>

EBI – Evropský institut pro bioinformatiku



European Bioinformatics Institute

- Založen roku 1992 jako součást European Molecular Biology Laboratory - EMBL.
- Sídlo v Hinxtonu ve Velké Británii.

Welcome to the EBI

The European Bioinformatics Institute (EBI) is a non-profit academic organisation that forms part of the European Molecular Biology Laboratory (EMBL).

The EBI is a centre for research and services in bioinformatics. The Institute manages databases of biological data including nucleic acid, protein sequences and macromolecular structures.



Our Mission

- To provide freely available data and bioinformatics services to all facets of the scientific community in ways that promote scientific progress
- To contribute to the advancement of biology through basic investigator-driven research in bioinformatics
- To provide advanced bioinformatics training to scientists at all levels, from PhD students to independent investigators
- To help disseminate cutting-edge technologies to industry

NCBI - Národní centrum pro biotechnologické informace

National Center for Biotechnology Information
National Library of Medicine National Institutes of Health

- Založeno v roce 1988 jako oddělení Národní lékařské knihovny (National Library of Medicine – NLM) v USA.
- Součást National Institutes of Health – NIH.


What does NCBI do?
Established in 1988 as a national resource for molecular biology information, NCBI creates public databases, conducts research in computational biology, develops software tools for analyzing genome data, and disseminates biomedical information - all for the better understanding of molecular processes affecting human health and disease.



CIB – Centrum pro informační biologii


- Založeno jako oddělení Národního genetického institutu (国立遺伝学研究所, NIG) v Japonsku.

Research Organization of Information and Systems
National Institute of Genetics <http://www.nig.ac.jp/>



Primární databáze NA

- EMBL** - Evropský institut pro bioinformatiku
- GenBank** - Národní centrum pro biotechnologické informace
- DDBJ** - Národní genetický institut (NIG)

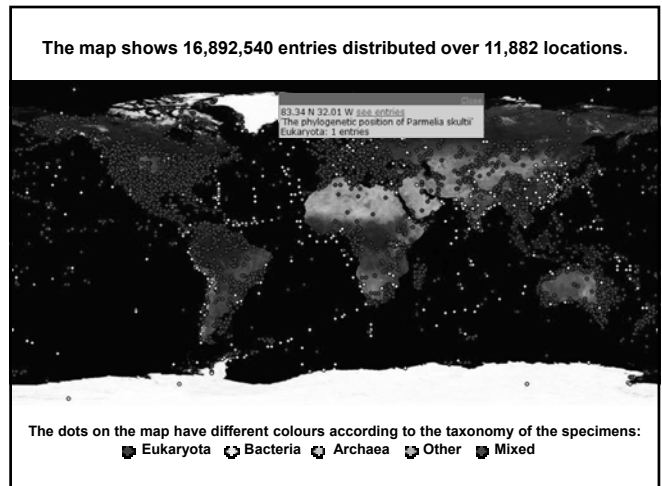
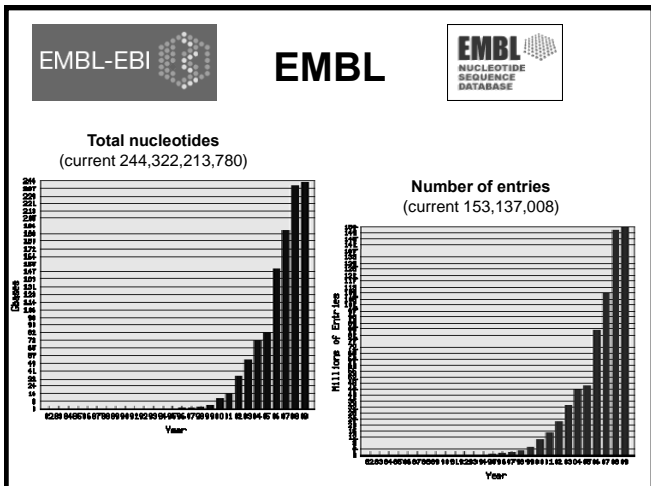


EMBL

- EMBL Nucleotide Sequence Database (EMBL-Bank)** byla založena roku 1980 jako první databáze nukleotidových sekvencí.
- Obsahuje sekvence RNA a DNA.
- Zdroje sekvencí: vloženy přímo autory, genomové projekty, patenty

This morning the EMBL Database contained
244,322,213,780 nucleotides
in 153,137,008 entries.

This morning = 21.1.2009





GenBank

- Založena roku 1982 v rámci institutu NCBI.

GenBank[®] is the NIH genetic sequence database, an annotated collection of all publicly available DNA sequences (*Nucleic Acids Research*, 2009 Jan 26 [Database issue] 37:25-30). There are approximately 65,759,566,764 bases in 82,863,685 sequence records in the traditional GenBank divisions and 108,635,736,141 bases in 27,439,206 sequence records in the WGS division as of February 2009.



Sample GenBank Record

<http://www.ncbi.nlm.nih.gov/Sitemap/samplerecord.html>



Sample GenBank Record

```

LOCUS      SC049845      5028 bp      DNA           PLN           21-JUN-1999
DEFINITION Saccharomyces cerevisiae TCP1-beta gene, partial cds, and Ax12p
            (AXL2) and Rev7p (REV7) genes, complete cds.
ACCESSION  U98845
VERSION    U98845.1      GI:1293613
KEYWORDS   .
SOURCE     Saccharomyces cerevisiae (baker's yeast)
ORGANISM   Saccharomyces cerevisiae
            Eukaryota; Fungi; Ascomycota; Saccharomycotina; Saccharomycetes;
            Saccharomycetales; Saccharomycetaceae; Saccharomyces.
REFERENCE  1 (bases 1 to 5028)
AUTHORS   Torpey,L.E., Gibbs,P.E., Helsen,J. and Lawrence,C.W.
TITLE     Cloning and sequence of REV7, a gene whose function is required for
            DNA damage-induced mutagenesis in Saccharomyces cerevisiae
JOURNAL    Yeast 10 (11), 1803-1809 (1994)
FUNSD     7871890
REFERENCE  2 (bases 1 to 5028)
AUTHORS   Roemer,T., Madden,K., Chang,J. and Snyder,M.
TITLE     Selection of axial growth sites in yeast requires Ax12p, a novel
            plasma membrane glycoprotein
JOURNAL    Genes Dev. 10 (7), 777-793 (1996)
FUNSD     0546915
REFERENCE  3 (bases 1 to 5028)
AUTHORS   Roemer,T.
TITLE     Direct Submission
JOURNAL    Submitted (22-FEB-1996) Terry Roemer, Biology, Yale University, New
            Haven, CT, USA
  
```

<http://www.ncbi.nlm.nih.gov/Sitemap/samplerecord.html>



Sample GenBank Record

```

LOCUS      SC049845      5028 bp      DNA           PLN           21-JUN-1999
DEFINITION Saccharomyces cerevisiae TCP1-beta gene, partial cds, and Ax12p
            (AXL2) and Rev7p (REV7) genes, complete cds.
ACCESSION  U98845
VERSION    U98845.1      GI:1293613
KEYWORDS   .
SOURCE     Saccharomyces cerevi
ORGANISM   Saccharomyces cerevi
            Eukaryota; Fungi; As
            Saccharomycetales; 3
REFERENCE  1 (bases 1 to 5028)
AUTHORS   Torpey,L.E., Gibbs,P
            Cloning and sequence
            DNA damage-induced m
            Yeast 10 (11), 1803-
            7871890
REFERENCE  2 (bases 1 to 5028)
AUTHORS   Roemer,T., Madden,K.
            Selection of axial g
            plasma membrane glyo
            Genes Dev. 10 (7), 7
            0546915
REFERENCE  3 (bases 1 to 5028)
AUTHORS   Roemer,T.
            Direct Submission
            Submitted (22-FEB-19
            Haven, CT, USA
  
```

GenBank Division The GenBank division to which a record belongs is indicated with a three letter abbreviation. In this example, GenBank division is PLN.

The GenBank database is divided into 18 divisions:

- 1 PRI - primate sequences
- 2 ROD - rodent sequences
- 3 MAM - other mammalian sequences
- 4 VRT - other vertebrate sequences
- 5 INV - invertebrate sequences
- 6 PLN - plant, fungal, and algal sequences
- 7 BCT - bacterial sequences
- 8 VIR - viral sequences
- 9 PHG - bacteriophage sequences
- 10 SYN - synthetic sequences
- 11 UNA - unannotated sequences
- 12 EST - EST sequences (expressed sequence tags)
- 13 PAT - patent sequences
- 14 STS - STS sequences (sequence tagged sites)
- 15 OSS - OSS sequences (genome survey sequences)
- 16 HTG - HTG sequences (high-throughput genomic sequences)
- 17 HTC - unfinished high-throughput cDNA sequencing
- 18 ENV - environmental sampling sequences

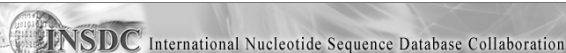
<http://www.ncbi.nlm.nih.gov/Sitemap/samplerecord.html>



DNA Data Bank of Japan

The DNA Data Bank of Japan

- Původně zahrnovala data především z japonských výzkumů.
- V současnosti úzká spolupráce s ostatními databázemi.



International Nucleotide Sequence Database Collaboration



<http://www.insdc.org/>

DDBJ: DNA Data Bank of Japan
CIB-DDBJ: Center for Information Biology and DNA Data Bank of Japan
NIIG: National Institute of Genetics

EBI: European Bioinformatics Institute
EMBL: European Molecular Biology Laboratory

NCBI: National Center for Biotechnology Information
NLM: National Library of Medicine

IAC: International Advisory Committee
ICM: International Collaborative Meeting

Primární databáze proteinů

- Univerzální databáze: „Skladiště“ sekvencí – sequence repository
- Manuálně spravovaná – curated database

Příklad: GenBank versus RefSeq



National Center for Biotechnology Information

National Library of Medicine

National Institutes of Health

Primární databáze proteinů

GenBank	RefSeq
Not curated	Curated
Author submits	NCBI creates from existing data
Only author can revise	NCBI revises as new data emerge
Multiple records for same loci common	Single records for each molecule of major organisms
Records can contradict each other	
No limit to species included	Limited to model organisms
Data exchanged among INSDC members	Exclusive NCBI database
Akin to primary literature	Akin to review articles
Proteins identified and linked	Proteins and transcripts identified and linked
Access via NCBI Nucleotide databases	Access via Nucleotide & Protein databases

GenPept - GenBank Gene Products Data Bank
RefSeq - Reference Sequence



Primární databáze proteinů

- PIR-PSD - Protein Information Resource Protein Sequence Database.
- Nejstarší univerzální „curated“ databáze proteinů.
- Komplexní, non-redundant data, rozčleněna podle proteinových rodin a nadrodin, doplněna funkčními, strukturními a bibliografickými daty.

<http://pir.georgetown.edu/>

swissprot Swiss-PROT + TrEMBL

- Swiss-Prot - „Curated“ databáze založená na Univerzitě v Ženevě v roce 1986. Spravovaná Švýcarským institutem pro bioinformatiku (SIB - Swiss Institute of Bioinformatics).
- Vysoká úroveň anotace \implies vkládáno více sekvencí než je možno manuálně anotovat a zařadit do databáze.
- TrEMBL – Počítačově anotovaná data, odvozená z kódujících úseků sekvencí v DDBJ/EMBL/GenBank, která ZATÍM nejsou zařazena v Swiss-Prot.



swissprot Swiss-PROT + TrEMBL

- Anotace:
 - Funkce
 - Katalytická aktivita
 - Podjednotky
 - Domény
 - Biotechnologické využití
 - Sekvenční homologie
 - Posttranslační modifikace
 - Reference atd.

<http://www.expasy.org/sprot/>

Složené databáze

- Databáze:
 - Primární
 - Sekundární
 - Strukturní

Genomové zdroje
Složené databáze

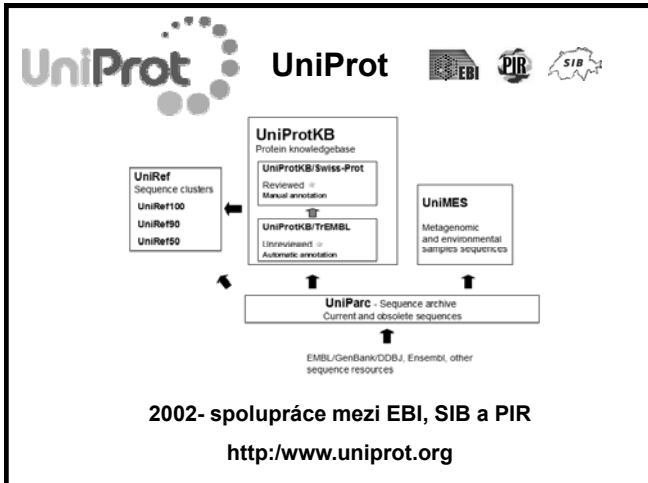
Složené databáze

- Složené (composite) databáze:
 - Stučují data z několika primárních databází.
 - Eliminace redundantních dat.
 - Různá priorita zdrojových databází podle kvality validace a anotace (eliminace redundantních dat z databáze s nižší prioritou).



Swiss-PROT + TrEMBL

OWL (Swiss-PROT + PIR + Genbank + NRL-3D)



Sekundární databáze NA a proteinů

Sekundární databáze obsahují informace odvozené z primárních databází ve formě charakteristických vzorů sekvencí, tj. funkčních nebo strukturních motivů získaných srovnáním primárních dat (sekvencí).

- Vyhledávání „vzoru“ charakteristického pro určitou skupinu proteinů.
- Možnost predikce funkce proteinů.



Sekundární databáze NA a proteinů

- Databáze mohou obsahovat:
 Proteinové DOMÉNY odvozené ze známých struktur
 Proteinové sekvence seřazené do SEKVENČNÍCH RODIN
 CHARAKTERISTICKÉ MOTIVY odvozené z těchto sekvencních rodin.



```

L P N T A K A I P Y A N A A D R D L K L F I D
I P N T D R A I F F A N A A R Q H I K L P T G
L P H I K G V T A L T H A A N D T I D I V D D
L P N I A G V T A L V N S S A P Q T I E V P V D
  
```

[AAC] -G-P-G-G-(EED) -

This pattern is conserved in: (P)in or (C)ys-arg-Val-arg-arg-arg-4-arg-Arg-Asp/Val or (Asp)

Sekundární databáze NA a proteinů

- Sekundární proteinové databáze:
PROSITE, Pfam, PRINTS, ProDom, SMART, TIGRFAMS
 V současné době sdruženy do integrované klasifikační databáze proteinů InterPro

<http://www.ebi.ac.uk/InterProscan/>

Table View	Raw Output	XML Output	Original Sequences	SUBMIT ANOTHER JOB
SEQUENCE: Sequence_1 CRC64: D00AB341613AD2EE LENGTH: 382 aa				
InterPro	Ricin B lectin			
PF00027X	PF00052			Ricin_B_lectin
Domain	SM00458			RICN
InterPro	PS50231			RICN_B_LECTIN
InterPro	Ribosome-inactivating protein			
PF001874	PF00181			RP
Family	SGF56371			Ribosome_inactivat_prot
InterPro	Ricin B-related lectin			
SK000997	SKF50370			RicinB_like
Domain				
InterPro	Ribosome-inactivating protein, subdomain 2			
PF016139	G30CA.4.10.470.10			Ribosome_inactivat_prot_sub2
Domain				
InterPro	Ribosome-inactivating protein subgroup			
PF001788	PF00196			SHIGARICN
Family				

Sekundární databáze NA a proteinů

- Sekundární proteinové databáze:
PROSITE, Pfam, PRINTS, ProDom, SMART, TIGRFAMS
 V současné době sdruženy do integrované klasifikační databáze proteinů InterPro

<http://www.ebi.ac.uk/InterProscan/>

- Sekundární databáze NA
 TRANSFAC

Strukturní databáze

Nucleic Acids Research

2009 NAR Database Summary Papers Category List

- Nucleotide Sequence Databases
- Protein sequence databases
- Structure Databases
- Small molecules
- Carbohydrates
- Nucleic acid structure
- Protein structure
- 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://www3.oup.co.uk/nar/database/a/>

Strukturní databáze proteinů

Nucleic Acids Research

2009 NAR Database Summary Papers

- Nucleotide Sequence Databases
- RNA sequence databases
- Protein sequence databases
- Structure Databases
- Small molecules
- Carbohydrates
- Nucleic acid structure
- Protein structure
- 3D-Genomics
- 3DID - 3D interacting domains
- ArchDB
- AS-ALPS
- ASTRAL
- AutoPSI
- BANKING
- BioMagResBank
- CADB - Conformational Angles DataBase of Proteins
- CATN
- CC+
- CE
- CnC Central
- CoGNE
- Columbia
- ConSurf-DB
- CPSB
- CSA - Catalytic Site Atlas
- DisProt - Database of Protein Disorder
- DMAPS
- Dooground
- DomIna - Database of Domain Insertions
- DSDBASE - Disulfide Database
- DSIM - a Database of Simulated Molecular Motions
- E-HSD - EBI-Macromolecular Structure Database
- eF-site - Electrostatic surface of Functional site
- ESCATDB
- FireDB
- FRN
- Gene3D
- Genomic Threading Database
- GTOP - Genomes To Protein structures
- HOMSTRAD - Homologous Structure Alignment Database
- HOStront
- IMGT/3Dstructure-DB
- INCTDB
- JAIL
- Jenablib: Jena Library of Biological Macromolecules
- KineticDB
- LPFC
- MALISAH
- MegaMotifBase
- MIMDB
- MultiBase
- MolMovDB - Database of Macromolecular Movements
- PDB3D
- PDB
- SCOP-PRODB
- PROSITE
- PROSUM

PDB - Protein Data Bank

RCSB PDB PROTEIN DATA BANK

A MEMBER OF THE **CFDB** MyPDB: Login | Register

An Information Portal to Biological Macromolecular Structures

as of Tuesday Feb 17, 2009 there are 55941 Structures | PDB Statistics

- Databáze obsahuje experimentálně získané struktury proteinů, nukleových kyselin a komplexů informačních biomakromolekul.

PDB Current Holdings Breakdown

Exp. Method	X-ray	Molecule Type				Total
		Proteins	Nucleic Acids	Protein/NA Complexes	Other	
		44706	1116	2060	24	47906
	NMR	6726	836	142	7	7711
	Electron Microscopy	146	16	55	0	217
	Other	96	5	4	2	107
Total		51674	1973	2261	33	55941

<http://www.rcsb.org/pdb/>

PDB formát

PDB File Format

The Protein Data Bank (PDB) format provides a standard representation for macromolecular structure data derived from X-ray diffraction and NMR studies. The representation was created in the 1970's and a large amount of software (e.g. Mol) has been written.

- PDB formát – původní formát databáze.
- 1997 – mmCIF (macromolecular Crystallographic Information File).
- Záznamy jsou v databázi uloženy v obou formátech a volně stažitelné.
- PDB formát – rozeznáván téměř všemi programy pro práci se strukturami.

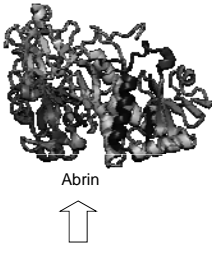
PDB formát

The ATOM records present the atomic coordinates for standard amino acids and nucleotides. They also present the occupancy and temperature factor for each atom. Non-polymer chemical coordinates use the HETATM record type. The element symbol is always present on each ATOM record; charge is optional.

Changes in ATOM/HETATM records result from the standardization atom and residue nomenclature. This nomenclature is described in the Chemical Component Dictionary (<http://www.rcsb.org/pdb/dictionaries/atomcomps>).

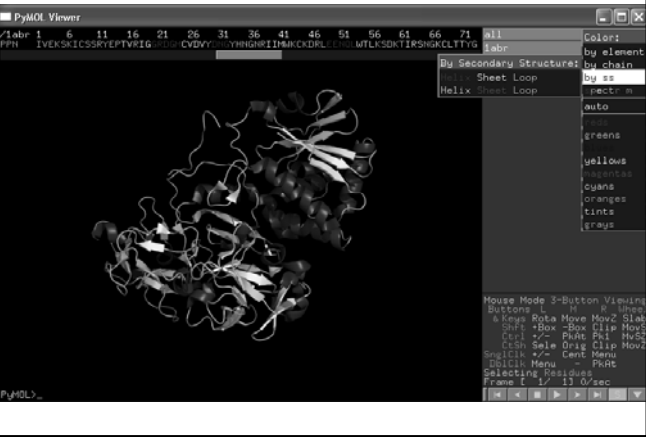
Record Format

columns	DATA TYPE	FIELD	DEFINITION
1 - 4	Record name	"ATOM"	
7 - 11	Integer	serial	Atom serial number.
12 - 14	Atom	name	Atom name.
17	Character	altLoc	Alternate location indicator.
18 - 20	Residue name	resName	Residue name.
22	Character	chainID	Chain identifier.
23 - 24	Integer	seqNum	Residue sequence number.
25	Atom	locId	Code for location of residues.
31 - 35	Real(float)	x	Orthogonal coordinates for 3-D Angstroms.
36 - 40	Real(float)	y	Ortho ATOM 2 CA GLU A 1
41 - 44	Real(float)	z	Ortho ATOM 3 C GLU A 1
45 - 50	Integer	occupancy	Ortho ATOM 4 O GLU A 1
51 - 55	Integer	TempFactor	Ortho ATOM 5 CB GLU A 1
56 - 60	Integer	TempFactor	Ortho ATOM 6 CG GLU A 1
61 - 65	Integer	TempFactor	Ortho ATOM 7 CD GLU A 1
66 - 70	Integer	TempFactor	Ortho ATOM 8 OE1 GLU A 1
71 - 75	Integer	TempFactor	Ortho ATOM 9 OE2 GLU A 1
76 - 80	Integer	TempFactor	Ortho ATOM 10 NI GLU A 1
			Ortho ATOM 11 N2 GLU A 1
			Ortho ATOM 12 N3 GLU A 1



Abrin

PyMOL Viewer



PyMOL Viewer interface showing the protein structure of Abrin. The command line at the bottom shows 'PDB01_...' and 'Color: by element'. The secondary structure menu is open, showing options like Sheet, Loop, Helix, and a 'Selecting Residues' dialog box with 'Frame 1 / 1 / 0/sec'.

Strukturní databáze NA

Nucleic Acids Research

ABOUT THIS JOURNAL CONTACT THIS JOURNAL SUBSCRIPTIONS CURRENT ISSUE

Oxford Journals > Life Sciences > Nucleic Acids Research > Database Summary Paper

2009 NAR Database Summary Papers

Nucleotide Sequence Databases
RNA sequence databases
Protein sequence databases
Structure Databases
Small molecules
Carbohydrates
Nucleic acid structure
GREGG
GRSDB
ITSD
MARNA
NDB - Non-Canonical Interactions in RNA
QuaRBase
Rfam
RNA FRADASE
RNA SSTRAND
RNAJunction
SARS-CoV RNA SSS
SCOR - Structural Classification Of RNA
VIMOR DB
Protein structure


NDB - Nucleic Acid Database

Number of Released Structures:
4089 Structures
Last Update: 15-Jan-2009

WELCOME TO THE NUCLEIC ACID DATABASE
a repository of three-dimensional structural information about nucleic acids

NDB ID: AD0001 NDB Atlas X-Ray Atlas

Title:	STRUCTURE OF A DNA IN LOW SALT CONDITIONS D (GACCGGTC)
Molecular Description:	5'-D (GpApCpCpGpCpGpTpC) -3'
Structural Features:	A DOUBLE HELIX
Nucleic Acid Sequence:	Class: A1 (DB) (DA) (DC) (DD) (DE) (DF) (DG) (DH) (DI) (DJ) (DK)
Primary Citation:	Foley, J.B., Lee, M. X-ray crystal structures of half the human papilloma virus E2 binding site, 4GACCGGTC. Nucleic Acids Res. 26, pp 5719-5727, 1998.
Experimental Information:	X-RAY DIFFRACTION
Space Group:	P 6 ₃ 2 2
Cell Constants:	a = 38.444 b = 38.444 c = 88.175 (Angstroms) α = 90.00 β = 90.00 γ = 120.00 (degrees)
Crystallization Conditions:	Method: VAPOR DIFFUSION Drop: WATER, MPD, Spermine HCl, Na Cacodylate Reservoir: WATER, MPD
Refinement:	The structure was refined using the J-FLOOR program. The R-value is 21.9 for 179 reflections in the resolution range 5.000 to 2.200 Angstroms with Falls = 3.900 signal/Falls.



Biological Unit 1
Other Views
Asymmetric Unit
Crystal Packing
Enlarge Biological Unit 1

<http://ndbserver.rutgers.edu/>

Atlas Deposit Download Search Reports Education Standards Tools Links

ATLAS

The Plain Melody for BDF062, Strand A (C G C T G G)

of Nucleic Acid Containing Structures

X-Ray Atlas

- Gallery Index
- Index Listing [text only]

NMR Atlas

- Gallery Index
- Index Listing [text only]
- Sorted Galleries**
- Musical Atlas
- About this Atlas

The NDB Atlas provides summary information and images for each structure in the database. These images provide many looks at the varied structures of nucleic acids.

The Atlas is first divided by experimental type, and then by structure type. Features include:

- images of the asymmetric and biological units, and crystal packing pictures for nucleic acid structures from X-ray crystallographic experiments
- images of the average and ensemble structure from NMR experiments
- links to coordinate files, experimental data files
- tables of derived data, including torsion angles and hydrogen bonding classifications
- special features for RNA structures, including images of secondary and tertiary structure

A more detailed description of the NDB Atlas features is available at "About this Atlas"

Genomové zdroje

Nucleic Acids Research

ABOUT THIS JOURNAL CONTACT THIS JOURNAL SUBSCRIPTIONS CURRENT ISSUE

Oxford Journals > Life Sciences > Nucleic Acids Research > Database Summary Paper Categories


2009 NAR Database Summary Papers Category List

Nucleotide Sequence Databases
RNA sequence databases
Protein sequence databases
Structure Databases
Genomics Databases (non-vertebrate)
MGD - Mouse Genome Database
TIGR Gene Indices
Genome annotation terms, ontologies and nomenclature
Taxonomy and identification
General genomics databases
Viral genome databases
Prokaryotic genome databases
Unicellular eukaryotes genome databases
Fungal genome databases
Invertebrate genome databases

EBI, NCBI – genomové databáze

Vyhledávací systémy


- Nutnost organizovaného ukládání a skladování dat.
- Nutnost prohlížení a analyzování uložených dat.



Databáze je určitá uspořádaná množina informací (dat) uložená na paměťovém médiu.

V širším smyslu jsou součástí databáze i softwarové prostředky, které umožňují manipulaci s uloženými daty a přístup k nim.

Vyhledávací systémy



- Textové vyhledávání v databázích

NCBI – Entrez

<http://www.ncbi.nlm.nih.gov/Entrez/>

Entrez is the integrated, text-based search and retrieval system used at NCBI for the major databases, including PubMed, Nucleotide and Protein Sequences, Protein Structures, Complete Genomes, Taxonomy, and others.

<http://www.ncbi.nlm.nih.gov/Entrez/tutor.html>

