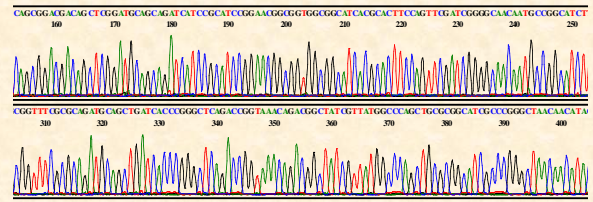


Molekulárně biologické databáze

Pro zajímavost, nebude součástí zkoušky...

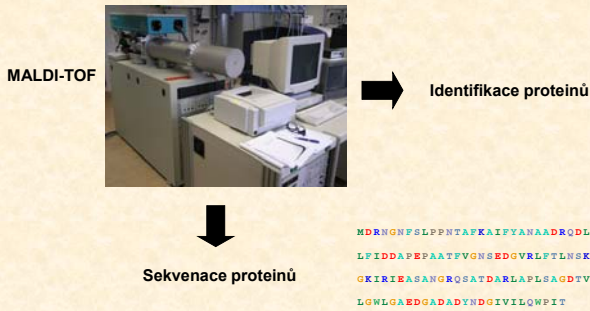
Důležité, pravděpodobně bude u zkoušky...

Molekulárně biologická data

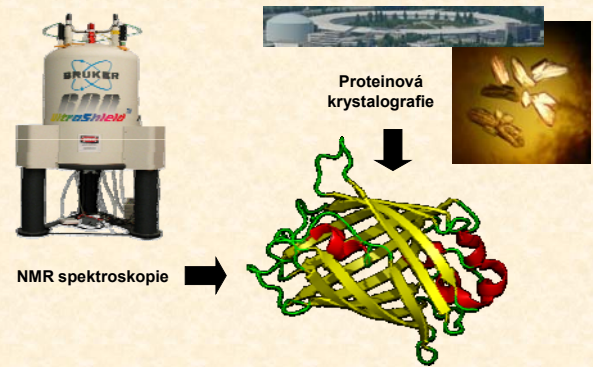


GATAGCGTAAATGATCGGCTGGCTGGCCATTTCACTGTTGGTTCCCAAGAAAATAACCGCTCACGGTCCATCACGATCGCACACCGAAAATCGGCG
 TACAGTGGTCCGCCCCCGCCAGCACATCGCTGGCCATAAATGATCTTTCAGCGACGACGATCGGATGAGCAGATCATCCGATCCGGAACGGC
 GGTGGCGCATCACCGCATTCGATGGATCGGGCAACATGCCCGCATCTTCAGGGCAAAAGCGAATAACACACCGCCACTTCGCGCGGACGACCGC
 AGCGCGTTTCGCGCAGATCGAGTGCATCACCCGGGCTCAGACCGGTAACAGACGGTATCGTTATGGCCAGCTCGCGCGCATCGCGGCTAACCA
 CATAACGGTGGCGACCATCAATCACCGTCCGGCGCGCGGATCACCGCTGGTTCGGGATAGGGCTCAGCAGGGTAAACCGCATCCAGAAATCACCGCAT

Molekulárně biologická data



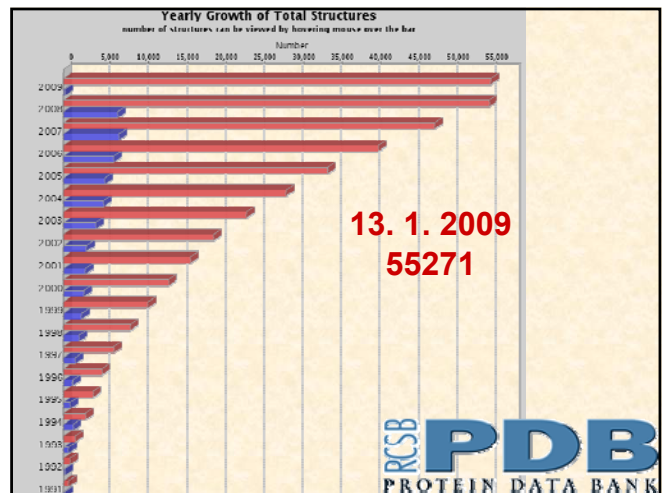
Molekulárně biologická data

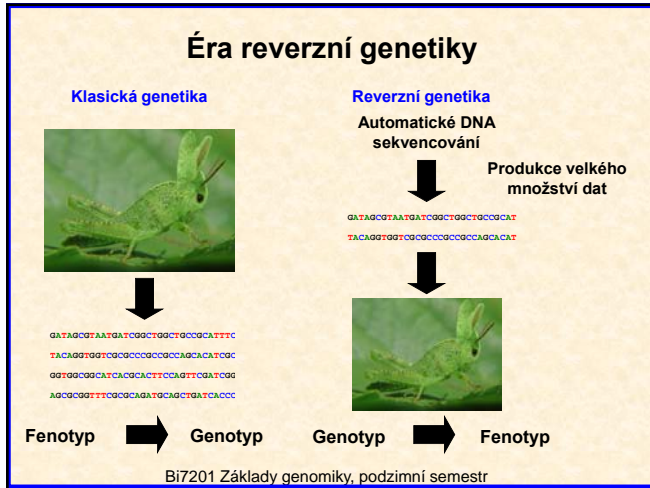


Molekulárně biologická data

- **Výkonné technologie:**
 Automatické sekvencování
 MALDI-TOF
 NMR spektroskopie
 Proteínová krystalografie

Výrazný nárůst množství biologických dat.





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.

```

LFFNTAKAIFYANAADRDLKLFIDAFEPAAATFVQNSRSDGVL--PFLNKGKXIIIE
LFFNTDRALFFANAABQDHKLFIDGSEPEAAHYKLLTRDGPFE--ATLNGSNKINFE
LFFNIXGVTALHAANDSTIDYIDDDPKFAATPKGAGAQDQNLGKVLDDGMRVSVI
LFFNIAEGVTALVHSSAPSTIEWVVDHNPKEAATPQAGTQDANLNTQIVNKGKRVVV
LFFN-afg---lanaad-QtikifidD-pFAATfkgag-----l-t-clngnkiRve

ASANGRCATDALLAPLSAGD-----TVNLEWLGAEDEABADYNDGIVILQWFIT
VSVNGKPRATDALLAPINOKKSDGSPFTVNFQIVVSRDHSBYNDGIVVILQWFIG
VMAHGPRRLGSRQVDIFKKE-----YFGLIGSEDAHDDYNDGIVVFNWFLG
VFAHSEPKIGSRQVDIFKKE-----YFGLIGSEDAHDDYNDGIVVFNWFLG
VFAHSEPKIGSRQVDIFKKE-----YFGLIGSEDAHDDYNDGIVVFNWFLG
vsaNGrpSat--R--ifkks-----tvyfGivgsEDGaDaDYNDGIVILqNFIg
    
```

Rozdělení molekulárně biologických databází

- Databáze:**
 - Primární
 - Sekundární
 - Strukturální


```
EDRPIKFSTEGATSQSYKQFIEALRERLRGLLHDIPVLPDPTTLQERNRYIT
VELSNDTSESIEVGLDVNAYVAYRAGTQSYFLRDAPSSASDYLFTGDQHS
LFFYGTGDLERWAHQSRQQLPLGLQALTHGISFFRSGGDNNEKARTLIVII
QMVAAARFRYISNRVSVISQGTAFQPDAAIISLENNWNLRSRGVQESVQDT
FFNQVLTINRNEPVIIVDSLHPTVAVLALMLFVCPNPNIVEKSKICSRYP
TVRIGGRDGMVDVYDNGYHNGNRIIMKCKDRLEENQLWTLKSDKTIIRSNKG
```

↓

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.

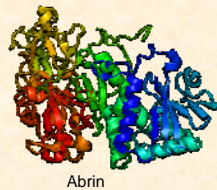
Rozdělení molekulárně biologických databází

- 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

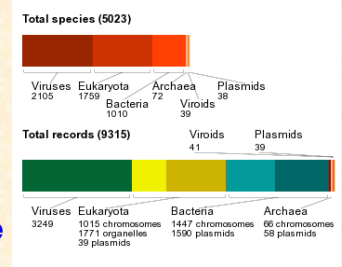


Obsahují struktury proteinů (nukleových kyselin) a jejich anotace.



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.

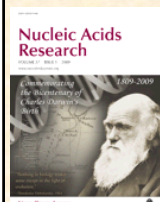
Eliminace redundantních dat.

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/>



[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

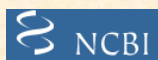
Evropský institut pro bioinformatiku



European Bioinformatics Institute

NCBI

Národní centrum pro biotechnologické informace



National Center for Biotechnology Information

CIB

Centrum pro informační biologii



Center for Information Biology

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

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

<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

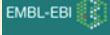


- Založeno jako oddělení Národního genetického institutu (国立遺伝学研究所, NIG) v Japonsku.

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)


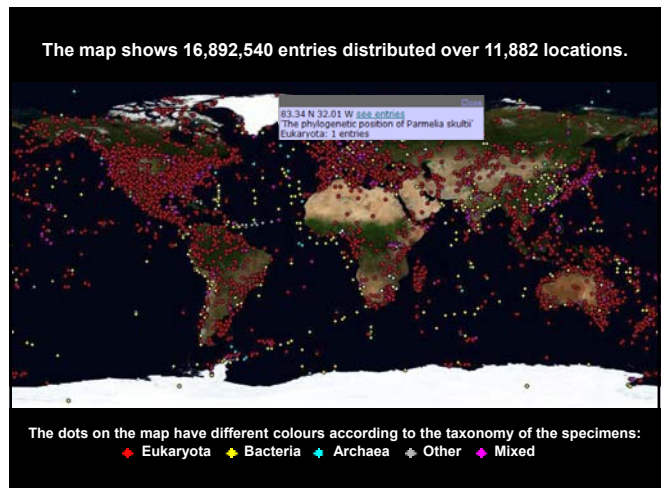
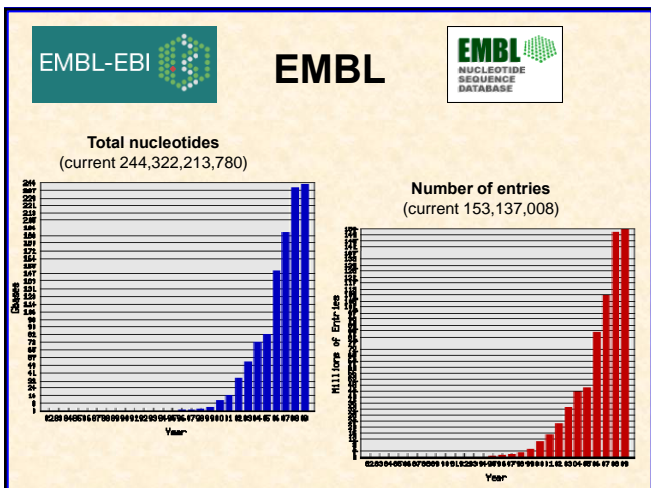





EMBL

- 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.

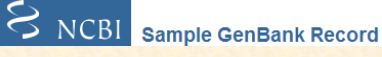
This morning = 21.1.2009

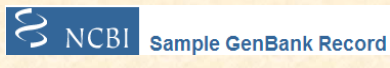


• Založena roku 1982 v rámci institutu NCBI.

GenBank[®] is the NIH genetic sequence database, an annotated collection of all publicly available DNA sequences (*Nucleic Acids Research*, 2009, Jan 26(Database issue), D26-50). There are approximately 85,759,566,764 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.



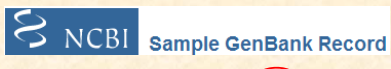
<http://www.ncbi.nlm.nih.gov/Sitemap/samplerecord.html>



```

LOCUS      SC049845      5028 bp      DNA           FLN           21-JUN-1999
DEFINITION Saccharomyces cerevisiae TCP1-beta gene, partial cds, and Ax12p
            (AX12) and Rev7p (REV7) genes, complete cds.
ACCESSION  U99645
VERSION   049845.1      GI:1293613
KEYWORDS  .
SOURCE    Saccharomyces cerevisiae (baker's yeast)
ORGANISM  Saccharomyces cerevisiae
            Eukaryota; Fungi; Ascomycota; Saccharomycotina; Saccharomycetes;
            Saccharomycetales; Saccharomycetaceae; Saccharomycetes.
REFERENCE  1 (bases 1 to 5028)
AUTHORS   Torpey,L.E., Gibbs,P.E., Helson,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
            Yeast 10 (11), 1503-1509 (1994)
JOURNAL   PUBMED 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
            Genes Dev. 10 (7), 777-793 (1996)
JOURNAL   PUBMED 8546915
REFERENCE  3 (bases 1 to 5028)
AUTHORS   Roemer,T.
TITLE     Direct Submission
            Submitted (22-FEB-1996) Terry Roemer, Biology, Yale University, New
            Haven, CT, USA
    
```

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



```


LOCUS      SC049845      5028 bp      DNA           FLN           21-JUN-1999
DEFINITION Saccharomyces cerevisiae TCP1-beta gene, partial cds, and Ax12p
            (AX12) and Rev7p (REV7) genes, complete cds.
ACCESSION  U99645
VERSION   049845.1      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
TITLE     Cloning and sequence
            DNA damage-induced m
            Yeast 10 (11), 1503-
            PUBMED 7871890
REFERENCE  2 (bases 1 to 5028)
AUTHORS   Roemer,T., Madden,K.
TITLE     Selection of axial g
            plasma membrane glyo
            Genes Dev. 10 (7), 7
            PUBMED 8546915
REFERENCE  3 (bases 1 to 5028)
AUTHORS   Roemer,T.
TITLE     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 FLN.

The GenBank database is divided into 18 divisions:

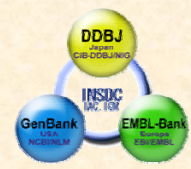
- 1 PRI - primate sequences
- 2 ROD - rodent sequences
- 3 MAM - other mammalian sequences
- 4 VRT - other vertebrate sequences
- 5 INV - invertebrate sequences
- 6 PLN - plant, fungal, and algal sequences
- 7 BCT - bacterial sequences
- 8 VIR - viral sequences
- 9 PHG - bacteriophage sequences
- 10 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>


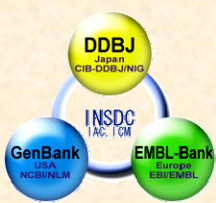


The DNA Data Bank of Japan

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



<http://www.ddbj.ac.jp/>

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

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

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 → vkládáno více sekvencí než je možno manuálně anotovat a zařadit do databáze.
- **TrEMBL** – Počítačově anotovaná data, odvozená z kódujících úseku 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**

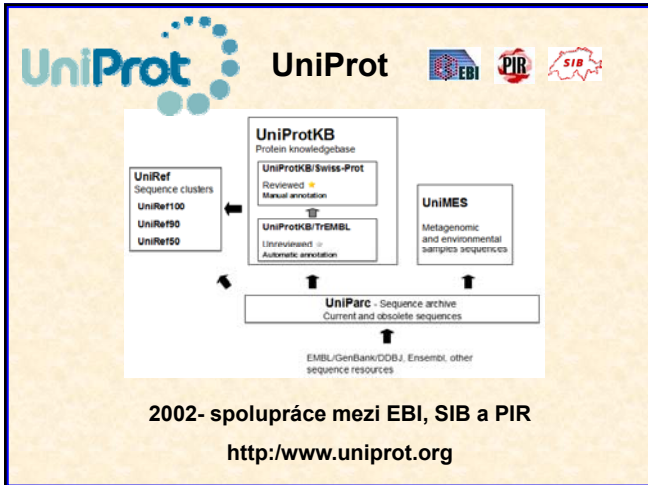
Složené databáze

- **Složené (composite) databáze:**
 - 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).



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 R D L K L F I D
L 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 O T I E V P V D D
  
```

[RQC] -G-P-R-Q-G-(RQC)-

This pattern is conserved in: [P] in or [C]q -any-Vai-any-amp-any-4-mp-famp-1-4-1-6iv 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-B00AB341813AD2EE LENGTH: 382 aa

InterPro	Ricin B lectin		
PF00027X			
Domain	PF00052		Ricin_B_lectin
InterPro			
Domain	S100458		RICIN
InterPro			
Domain	P56021		RICIN_B_LECTIN
InterPro	Ribosome-inactivating protein		
PF001874			
Family	PF00181		RP
InterPro			
Family	SGF56371		Ribosome_inactivat_prot
InterPro	Ricin B-related lectin		
94000997			
Domain	SF750170		RicinB_like
InterPro	Ribosome-inactivating protein, subdomain 2		
PF016138			
Domain	G305A.4.10.470.10		Ribosome_inactivat_prot_sub2
InterPro	Ribosome-inactivating protein subgroup		
PF01788			
Family	PF00396		SHIGARICIN

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

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
- 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

ABOUT THIS JOURNAL CONTACT THIS JOURNAL SUBSCRIPTIONS CURRENT ISSUE ARCHIVE SEARCH

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
- Protein structure
- 3D-Genomics
- 3D - 3D interacting domains
- ArchDB
- AS-GLIPS
- ASTRAL
- AutoPSI
- BANKING
- BioMagResBank
- CADB - Conformational Angles Database of Proteins
- CATN
- CC+
- CE
- CEC Central
- ColIGN
- Columb
- ConSurf-DB
- CPSB
- CSA - Catalytic Site Atlas
- DisProt - Database of Protein Disorder
- DRAPS
- Drugground
- DomInt - Database of Domain Interactions
- DSDBase - Disulfide Database
- DSM - a Database of Simulated Molecular Motions
- E-MSD - EBI-Macromolecular Structure Database
- el-site - Electrostatic surface of Functional site
- ESUICB
- FireDB
- FRN
- Gene3D
- Genomic Threading Database
- GTOP - Genomes To Protein structures
- HOMSTRAD - Homologous Structure Alignment Database
- HOPIent
- IMGT/3Dstructure-DB
- IMGTdb
- JAIN
- Jenabib: Jena Library of Biological Macromolecules
- KineticDB
- LPFC
- MALISAH
- MegaNatiBase
- MIMB
- MolBase
- MolMovDB - Database of Macromolecular Movements
- MOBI
- POB
- PROSPRDB
- ProteinAtlas
- PROSUM

PDB - Protein Data Bank

RCB PDB A MEMBER OF THE **CPDB** 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

	Method	Molecule Type				Total
		Proteins	Nucleic Acids	Protein/NA Complexes	Other	
	X-ray	44706	1116	2060	24	47906
	NMR	6726	636	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 using it has been developed.

Documents describing the PDB file format is available from the website at <http://www.rcsb.org/pdb/files.html>.

Where a copy of the PDB file format specification (PDB) is not available.

- File Formats & File Formats
- File Exchange Dictionary
- Chemical Component
- 3D Format
- File Format
- File Format
- File Format
- File Format

- 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://ftp.rcsb.org/pub/wwpdb/data/ccdictionaries>).

Record Format

column	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	seqRes	Residue sequence number.
25	Alpha	code	Code for insertion of residues.
31 - 35	Real(1,3)	x	Orthogonal coordinate for X in Angstroms.
36 - 40	Real(1,3)	y	Ortho. ATOM 2 CA GLU A 1 64.373 11.709 60.583 1.00 79.99 C
41 - 44	Real(1,3)	z	Ortho. ATOM 3 C GLU A 1 63.512 10.438 60.597 1.00 79.31 C
45 - 48	Real(1,3)	occupancy	Ortho. ATOM 4 O GLU A 1 63.540 9.485 61.574 1.00 79.23 O
49 - 50	Real(1,3)	tempFactor	Ortho. ATOM 5 CB GLU A 1 63.005 12.794 59.603 1.00 79.36 C
61 - 64	Real(1,3)	tempFactor	Temp. ATOM 6 CG GLU A 1 62.880 13.819 60.228 1.00 78.52 C
65 - 68	Real(1,3)	tempFactor	Temp. ATOM 7 CD GLU A 1 61.525 13.275 60.476 1.00 78.50 C
69 - 72	String(4)	element	Elem. ATOM 8 OE1 GLU A 1 60.915 12.482 59.923 1.00 77.14 O
73 - 76	String(4)	charge	Charg. ATOM 9 OE2 GLU A 1 61.064 13.659 61.776 1.00 77.48 O
77 - 80	String(4)	charge	Charg. ATOM 10 H1 GLU A 1 66.076 10.680 60.918 1.00 20.00 H
81 - 84	String(4)	charge	Charg. ATOM 11 H2 GLU A 1 65.774 10.893 59.265 1.00 20.00 H
85 - 88	String(4)	charge	Charg. ATOM 12 H3 GLU A 1 66.387 12.177 60.222 1.00 20.00 H

PyMOL Viewer

Color: by element

By Secondary Structure: by chain

Helix Sheet Loop by ss

Helix Sheet Loop by e1, n

Auto

Front

Reverse

Close

Yellow

Magenta

Cyan

Orange

Pink

Gray

Mouse Mode 3-Button Viewing

Buttons: L H A W

Keys: Rotate Move Move2 Slab

Ctrl + Box - Box D11e Move

Ctrl +/- PKA PK1 HUS

Ctrl + S Sale D11e Move

Single Click +/- Ctrl Menu

Double Click Menu - PKR

Right-click Residues

Frame 1 / 11 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
Nucleic acid structure
GROFIT
GRSDB
ITSD
MARNA
NDB - Non-Canonical Interactions in RNA
QuaBase
Rfam
RNA FRABASE
RNA SSTRAND
RNAjunction
SARS-CoV RNA 3SS
SCOR - Structural Classification of RNA
Viterbi-DB
Protein structure

NDB - Nucleic Acid Database

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

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

NDB ID: AD0001

Biological Unit 1

Other Views
Asymmetric Unit
Crystal Packing
Entire Biological Unit 1

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

ATLAS

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)

EBI, NCBI – genomové databáze

MGD - Mouse Genome Database
TIGR Gene Indices
Genome annotation terms, ontologies and nomenclature
Taxonomy and identification
General genomics databases
Viral genomics databases
Prokaryotic genome databases
Unicellular eukaryotes genome databases
Fungal genome databases
Invertebrate genome databases

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>

