

Slide 1

BIOINFORMATIKA- UVOD
Podatkovne zbirke v biokemiji in molekularni biologiji
31.3.-1.4.2005

1. Internet, strežniki z biokemijsko in molekularno-biološko vsebino
2. Nacionalni center za biotehnoško informacijo (NCBI), Expert Protein Analysis System (ExPASy)
3. Proteinske podatkovne zbirke, UniProt, nukleotidne podatkovne zbirke, GenBank
4. Bibliografske podatkovne zbirke, Pubmed, OMIM, Bookshelf

VIRI
<http://rcul.uni-lj.si/~bfbiochem/bioinfo/bioinfo.htm>; Spletna stran o bioinformatiki
Literatura- Popoln seznam na zgornjih spletnih straneh
Attwood TK, Parry-Smith DJ (1999) Introduction to Bioinformatics. Prentice Hall, Harlow, United Kingdom

Slide 2

BIOINFORMATIKA

Vmesnik med biologijo in računalništvom.

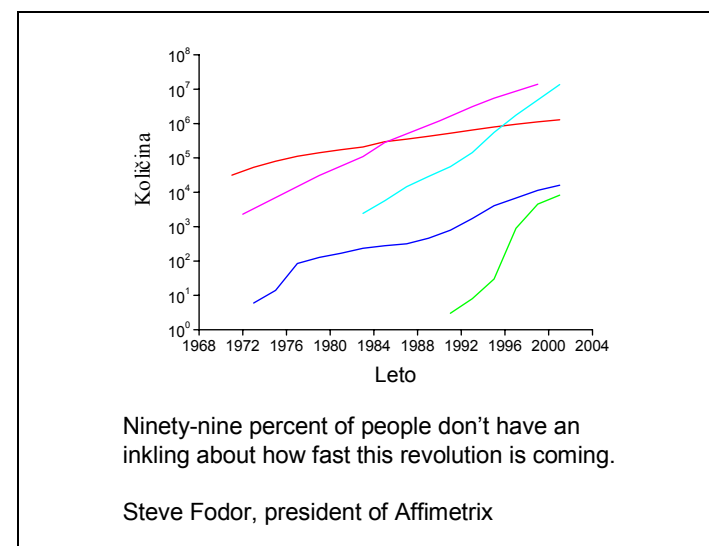
Matematične, statistične in računalniške metode za reševanje bioloških problemov z uporabo DNA in proteinskih zaporedij in povezane informacije.

Organiziranje, analiza in distribucija biološke informacije pri opisovanju in reševanju bioloških problemov.

Zbiranje, organiziranje, shranjevanje in iskanje biološke informacije v podatkovnih zbirkah.

NI ISKANJE ČLANKOV IN BRSKANJE ZA REVIJAMI!!!

Slide 3



Slide 4

```
>gi|122615|sp|P02023|HBB_HUMAN HEMOGLOBIN BETA CHAIN  
MVHLTPEEKSAVTALWGKVNVDVGGGALGRLLVVYFWTQRFESFGDLSTP  
DAVMGNPVKAHGKVLGAFSDGLAHLNLRGTFATLSLHCDKLVDPENFR  
LLGNVLCVLAHFGKEFTPPVQAAYQKVVAGVANALAHKYH
```

Vihinen M (2001) *Biomol. Eng.* 18:241-248

Slide 5

STREŽNIKI Z BIOKEMIJSKO IN MOL.- BIOL. INFORMACIJO

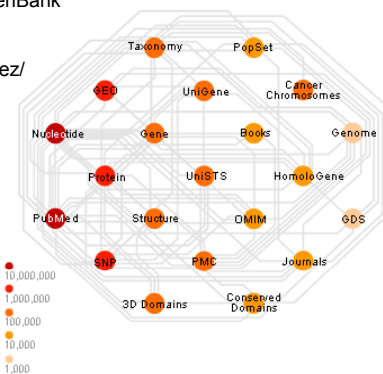
SRS	Sequence Retrieval System
NCBI	The National Center for Biotechnology Information
ExpASY	Expert Protein Analysis System
KEGG	Kyoto Encyclopedia of Genes and Genomes

Slide 6

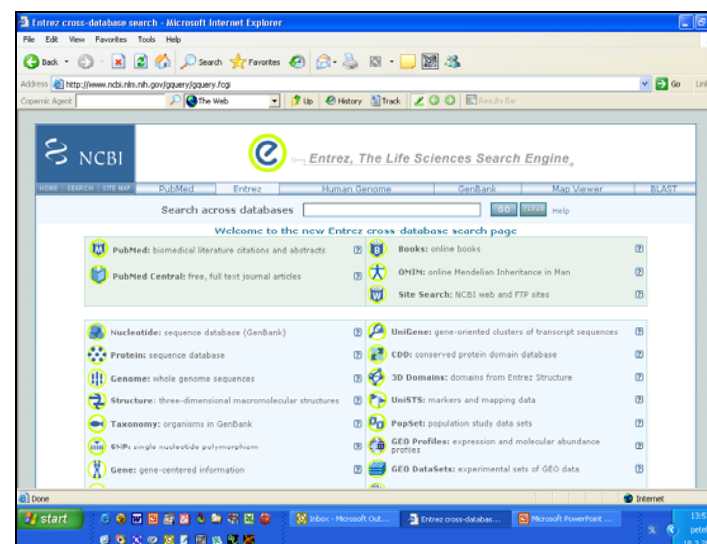
NCBI- The National Center for Biotechnology Information
1988 Kot del National Library of Medicine
Namen: *razvoj novih informacijskih tehnologij, ki pomagajo razumeti molekularne in genetske procese*

1992 Upravljanje z GenBank

Entrez
<http://www.ncbi.nlm.nih.gov/Entrez/>
Brskalnik za iskanje po bioloških podatkovnih zbirkah, ki jih ureja NCBI (proteinska zaporedja, nukleotidna zaporedja, gensko mapiranje, OMIM, 3D strukture iz PDB, PubMed)



Slide 7



Slide 8

ExpASY Expert Protein Analysis System

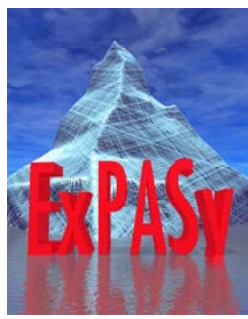
<http://www.expasy.org/>
Swiss Institute of bioinformatics (SIB)

The ExpASY (Expert Protein Analysis System) proteomics server from the Swiss Institute of Bioinformatics (SIB) is dedicated to molecular biology with an emphasis on data relevant to **proteins**.

Appel R.D., Bairoch A., Hochstrasser D.F.
A new generation of information retrieval tools for biologists: the example of the ExpASY WWW server.
Trends Biochem. Sci. 19:258-260(1994).

Podatkovne zbirke: **UniProt**
"SwissProt"
TrEMBL
PROSITE
SWISS 2D-PAGE

Programi: Orodja za proteomiko
Melanie (analiza 2D gelov)
Biochemical Pathways



Slide 9

PODATKOVNE ZBIRKE


PRIMARNE	("archival", arhivske)	
SEKUNDARNE	("curated", izpeljane)	
Nucleic Acids Research (2003) 31, number 1		
Primarne	nukleotidne (GenBank, EMBL, DDBJ), proteinske (UniProt), strukturne (Protein Data Bank)	
Bibliografske	PubMed, OMIM	

ATGGTGCACCTGAC
TCCTGAGGAGAAG...

GenBank, EMBL, DDBJ

MVHLTPEEK...


UniProt



Protein Data Bank

Slide 10

PROTEINSKE PODATKOVNE ZBIRKE



the universal protein resource

UniProt <http://www.expasy.uniprot.org/index.shtml>

UniProt (Universal Protein Resource) **is the world's most comprehensive catalog of information on proteins.** It is a central repository of protein sequence and function created by joining the information contained in Swiss-Prot, TrEMBL, and PIR.

Upravljata jo SIB (Swiss Institute of Bioinformatics) in EBI/EMBL (The European Bioinformatics Institute, UK). Vsebuje proteinska zaporedja

UniProt **Različica 4.3** vsebuje: Swiss-Prot različico 46.3 of 15-Mar-2005: **176.469** zapisov iz **8.998** vrst

Zelo kvalitetna **določitev** biološke funkcije (**ANOTACIJA, "annotation"**) → zelo uporabna podatkovna zbirka

1996 TrEMBL (SP-TrEMBL, REM-TrEMBL). Prevedena nukleotidna zaporedja vseh kodirajočih zaporedij (*coding sequences, CDS*) iz EMBL. TrEMBL različica 29.3 15-Mar-2005 vsebuje **1.640.768** zapisov

Slide 11

SwissProt

Nucleic Acids Research (2003)
31, 365-370

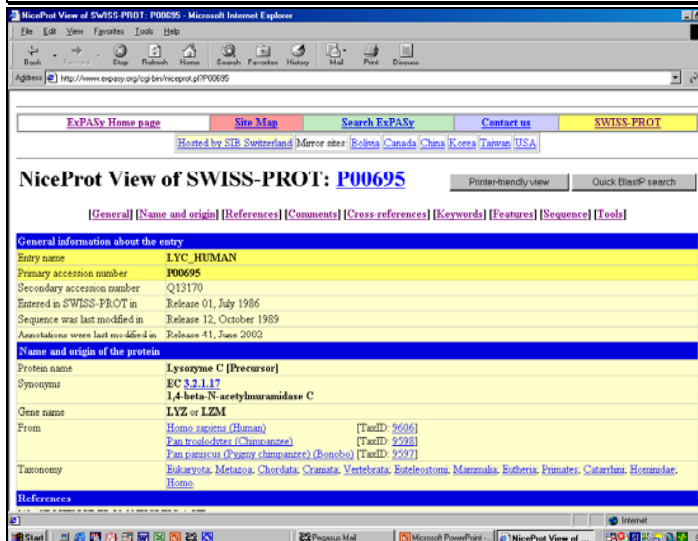
Core data aminokislinsko zaporedje
ime proteina
taksonomski podatki
citati

Določitev ("annotation")

- funkcija proteina
- podatki o encimski aktivnosti (katalitska aktivnost, kofaktorji, metabolne poti...)
- biološko pomembna mesta na molekuli (domene), posttranslacijske modifikacije
- M določena z masno spektroskopijo
- ktivno specifična izražanje
- razvojno specifično izražanje
- sekundarna struktura
- kvartarna struktura
- polimorfizmi
- podobnost z drugimi proteini
- uporaba v biotehnoških procesih
- napake (konflikti) v zaporedjih

Vpis "na roke" iz primarnih publikacij, preglednih člankov.
Pomoč zunanjih ekspertov za posamezne družine proteinov.

Slide 12



The screenshot shows the UniProt NiceProt View for the protein P00695 (LYC_HUMAN). The page includes navigation links like 'ExPASy Home page', 'Site Map', 'Search ExPASy', 'Contact us', and 'SWISS-PROT'. The main content area displays 'General information about the entry' with fields for Entry name, Primary accession number, Secondary accession number, Entered in SWISS-PROT in, Sequence was last modified in, and Annotations were last modified in. Below this is the 'Name and origin of the protein' section, which lists the protein name (Lysozyme C [Precursor]), Synonyms (EC 3.2.1.17, 1,4-beta-N-acetylmuramidase C), Gene name (LYZ or LZM), and From (Homo sapiens (Human), Pan troglodytes (Chimpanzee), Pan paniscus (Pygmy chimpanzee) (Bonobo)). The Taxonomy section shows the classification: Bacteria, Firmicutes, Bacilli, Lactobacillales, Lactobacillaceae, Lactobacillus.

Slide 17

STRUKTURA GenBank ZAPISA

GBFF- GenBank flat file oblika zapisa v GenBank. Zapis, ki si ga izmenjujejo podatkovne zbirke.

Glava (header)
Lastnosti
Zaporedje

```

LOCUS       NM_000239             1487 bp    mRNA    linear    PRI 18-DEC-2001
DEFINITION Homo sapiens lysozyme (renal amyloidosis) (LYZ), mRNA.
ACCESSION  NM_000239
VERSION    NM_000239.1  GI:4557893
KEYWORDS   .
SOURCE     human.
  ORGANISM Homo sapiens
            Eukaryota; Metazoa; Chordata; Craniata; Vertebrata; Euteleostomi;
            Mammalia; Eutheria; Primates; Catarrhini; Hominidae; Homo.
REFERENCE  1 (bases 1 to 1487)
AUTHORS    Reitamo,S., Klockars,M., Adinolfi,M. and Osserman,E.F.
TITLE      Human lysozyme (origin and distribution in health and disease)
JOURNAL    Ric Clin Lab 8 (4), 211-231 (1978)
MEDLINE    79097291
PUBMED     366724
    
```

Slide 19

```

DEFINITION Homo sapiens lysozyme (renal amyloidosis) (LYZ), mRNA.
Bioški opis zaporedja (izpisan v FASTA formatu)
ACCESSION  NM_000239             NT_123456 constructed genomic contigs
Edinstvena koda, se ne spreminja, podana v publikacijah  NM_123456 mRNAs
                                                    NP_123456 proteins
                                                    NC_123456 chromosomes
VERSION    NM_000239.1  GI:4557893
Katera verzija zaporedja je. Se spreminja. GI (geneInfo identifier) edinstveno določena koda za zaporedje (enaka med različnimi podatkovnimi zbirkami).
KEYWORDS   .
SOURCE     human.
  ORGANISM Homo sapiens
            Eukaryota; Metazoa; Chordata; Craniata; Vertebrata; Euteleostomi;
            Mammalia; Eutheria; Primates; Catarrhini; Hominidae; Homo.
Izvor zaporedja.
REFERENCE  1 (bases 1 to 1487)
AUTHORS    Reitamo,S., Klockars,M., Adinolfi,M. and Osserman,E.F.
TITLE      Human lysozyme (origin and distribution in health and disease)
JOURNAL    Ric Clin Lab 8 (4), 211-231 (1978)
MEDLINE    79097291
PUBMED     366724
Bibliografske enote povezane z zaporedjem (kdaj sklonirano, kdaj določena sturktura gena...)
    
```

Slide 18

Glava

```

LOCUS       NM_000239             1487 bp    mRNA    linear    PRI 18-DEC-2001
    
```

kratka oznaka zapisa
 včasih ime lokusa
 Npr. HUMHBB človeški β-globinski lokus

dolžina
 Manj kot 50 bp
 ne sprejmejo, več kot 350 kb tudi ne

tip zaporedja
 DNA, tRNA, rRNA, mRNA

Odsek GenBank
 PRI primati
 ROD glodalci
 MAM ostali sesalci
 VRT ostali vretenčarji
 INV nevretenčarji
 PLN rastline
 BCT bakterije
 VRL virusi
 PHG bakteriofag
 SYN sintetične
 EST oznaka izraženega zaporedja
 PAT patent
 STS sequence tagged sites
 GSS genome survey sequences
 HTG "high throughput" genomski zaporedja
 HTC nedokončane "high throughput" EST zap.

datum, ko je zapis postal javen

Slide 20

```

FEATURES             Location/Qualifiers
     source            1..1487
                        /organism="Homo sapiens"
                        /db_xref="taxon:9606"
                        /chromosome="12"
                        /map="12q15"
                        /clone="pHL-2"
                        /cell_line="PMA treated U937"
     gene              1..1487
                        /gene="LYZ"
                        /db_xref="LocusID:1061"
                        /db_xref="MIM:113110"
     CDS               26..469
                        /gene="LYZ"
                        /EC_number="3.2.1.17"
                        /codon_start=1
                        /db_xref="LocusID:1061"
                        /db_xref="MIM:113110"
                        /product="lysozyme precursor"
                        /protein_id="AA022222"
                        /db_xref="GI:4557894"
                        /translation="MSALIVLGLVLLSVTVQGRVFCERLARTLRLGLMDGYRGISLA
                        NMKLAWESSINFRATNINAGDRSDTIGIFQINSRIWCDGKTPGAVNACHLSCAL
                        LQINIDAVAKKRVVSDPQIFAWVAVRRCQNRIVQVQVGVCCV"
     misc_peptide      26..469
                        /product="lysozyme"
     misc_feature      90..463
                        /note="Lys; Region: C-type lysozyme/alpha-lactalbumin
                        family. Alpha-lactalbumin is the regulatory subunit of
                        lactose synthase"
     misc_feature      90..466
                        /note="LYZ1; Region: Alpha-lactalbumin / lysozyme C"
     variation         288
                        /allele="C"
                        /allele="A"
                        /db_xref="dbSNP:130951"
     repeat_region     740..1054
                        /note="Alu repeat"
    
```

Lastnosti zaporedja "feature table".

Svojo lastno strukturo in povezave z drugimi podatkovnimi zbirkami

Slide 21

Zaporedje

```

BASE COUNT      435 a   306 c   308 g   438 t
ORIGIN
1 ctgacctct gacctagcag tcaacatgaa ggccttcatt gttctggggc ttgtctctct
61 ttctgttacg gtcccaggca aggtctttga aaggtgtgag ttggccagaa ctctgaaag
121 attgggaatg gatgctaca ggggaatcag cctagcaaac tggatgtgtt tggccaastg
181 ggaagatggt tacaacacac gagctacaaa ctcaactgpt ggaacagaa gcactgatta
241 tgggatattt cagatcaata gcgcctactg gtgtaatgat gcaaaaaccc caggagcagt
301 taatgctgtt catttatcct gcagtccttt gctgcaaat acatcctcgt atgctgtagc
361 ttgtgaaag aggtgtgtcc gtagtccaca aggcattaga gcaatgggtg catggagaaa
421 tctgttcaaa aacagagatg tccgtcagta tgcacaagpt tgtggaggtt aactccagaa
481 ttttctctct tcaactcatt ttgtctctct cacattaagg gagttagaat taagtgaag
541 gtcacactac cattatttcc ctttcaaaa aataatattt ttacagaagc aggagcaaaa
601 tatggccttt cttctaagag atataatggt cactaatgtg gttattttac attaagccta
661 caacattttt cagtttgaaa atagaactaa taccgttgaa aattttacta aaactctgtt
721 tacaataac atctccagta cattcctcct ttttttttt ttgagacagt ctgctctgt
781 cgcacagctt ggaatgcagt ggcgcaatct cggctcaact caactccac ctcccgggtt
841 caagcaatc tctgtctca gcttcccgag tagctgggat tacggggccc cgcaccacag
901 ccccgcaaat ttttgtattt tttagtagag acagggcttc accgtgttag ccaaggtggt
961 ctgatctctc tgaccttggt atccaccac ctggcctccc caaagtctg gattacac
1021 cgtgagccac tggccocggc cacattcagt ttttatcaaa gaataaccc agacttaatc
1081 ttgatgata cgtatagcc caatattaag taasaatat aagaaaagt tatcttaaat
1141 agactctagg caaataacca gctgatgaag gcactctgag ccttcaatcg ttcagctcac
1201 tccaaaacca gtaaaaataa ccaacttttg ttggccaata tgaattttt aaagagtag
1261 aatacaaat gatagaacca gactgctgta attgagaatt ttgatttctt aaagtgtgtt
1321 ttttccaaa ttgtgttccc ttaattgat taaattaat catgtaatat gattaaatct
1381 gggcagatg agctcaag tattgaaata attactaatt aatcaaaat gtaagttat
1441 gcatgagta aaaaataca acattcaat taaggcttt gcaacac
    
```

Slide 22

OBLIKA ZAPISA ZAPOREDIJ

Računalnik vs človek
GBFF
ASN.1
FASTA(Pearson) najbolj enostaven zapis

```

>P00695
MKALIVLGLVLLSVTVQGVFERCELARTLKRGLGMDGYRGISLANWMCLAKWESGYNTRATNYNAGDRS
TDYGFIFQINSRYWCDNGKTPGAVNACHLSCSALLQDNIADAVACAKRVVRDPQGIRAWVAWRNRCQNRD
VRQYVQGGCV

>gi|4557894|ref|NP_000230.1| lysozyme precursor [Homo sapiens]
MKALIVLGLVLLSVTVQGVFERCELARTLKRGLGMDGYRGISLANWMCLAKWESGYNTRATNYNAGDRS
TDYGFIFQINSRYWCDNGKTPGAVNACHLSCSALLQDNIADAVACAKRVVRDPQGIRAWVAWRNRCQNRD
VRQYVQGGCV

>gi|4557893|ref|NM_000239.1| Homo sapiens lysozyme (renal amyloidosis)
(LYZ), mRNA
CTAGCACTCTGACCTAGCAGTCAACATGAAGCCTCTCATTGTTCTGGGGCTTGCTCCTCTTCTGTTACG
GTCCAGGGCAAGGCTTTGAAAGGTGTGAGTTGGCCAGAACTCTGAAAAGATTGGGAATGGATGGCTACA
GGGGAATCAGCCTAGCAAACTGGATGTGTTTGGCCAAATGGGAGAGTGGTTCAACAACACACGAGCTACAAA
CTACAATGCTGGAGACAGAACACTGATATGGGATATTGAGATCAATAGCCGCTACTGGTGTATGAT
GGCAAAACCCAGGAGCAGTAAATGCCTGTCATTTATCTCGAGTCTTGTGTCGAAGATAACATCGCTG
    
```

Slide 23

BIBLIOGRAFSKE PODATKOVNE ZBIRKE



NCBI- tekstovne podatkovne zbirke

- PubMed
- OMIM
- BOOK-SHELF
- PubMed Central

PubMed
Na NCBI. Prosto dostopna podatkovna zbirka.
<http://www.ncbi.nlm.nih.gov/entrez/query.fcgi?db=PubMed>
Vsebuje citate iz MEDLINE in še dodatno informacijo (nekatero druge revije, ki jih v MEDLINE ni)
MEDLINE: vsebuje bibliografsko informacijo 4300 revij iz 70 držav; vsebuje preko 11 mio citatov od 1960 naprej. (OLDMEDLINE- 1958-1965)
PubMed Journal Browser- informacija o revijah, iz katerih so članki v PubMed; **JournalLinkOutProvider-** povezave na spletne strani založnikov, ki imajo revije na spletu; **MeSH browser (Medical Subject Heading)-** pojmi, ki jih PubMed uporablja za indeksiranje člankov

Slide 24

OMIM

On-line Mendelian Inheritance in Man
<http://www.ncbi.nlm.nih.gov/entrez/query.fcgi?db=OMIM>

Katalog človeških genov in bolezni s katerimi so povezani. (=fenotipski dodatek človeškemu genomu). Ureja dr. Victor A McKusick, John Hopkins University. Tekstovna informacija in reference, povezave na podatkovne zbirke znotraj NCBI in ostale.

Skupaj 14925 opisov (18.3.2005)

Iskanje s ključnimi besedami

	Autosomal	X-Linked	Y-Linked	Mitochondrial	Total
* Gene with known sequence	2561	430	48	37	10076
* Gene with known sequence and phenotype	364	37	0	0	401
# Phenotype description, molecular basis known	1527	138	2	27	1694
‡ Mendelian phenotype or locus, molecular basis unknown	1328	134	4	0	1466
Other, mainly phenotypes with suspected mendelian basis	2145	153	2	0	2300
Total	14925	892	56	64	15937

BOOK-SHELF



37 "on-line" knjig z vsebino mol biologije, biokemije, celične biologije...

Iskanje gesel

Npr. Molecular Biology of the Cell (Alberts), Biochemistry (Stryer), Molecular Cell Biology (Lodish)...

PUB-MED CENTRAL

<http://www.pubmedcentral.nih.gov/>



Arhiv revij (>100) s področja naravoslovja v polnem tekstu!!!

Npr. EMBO J

Infection and Immunity

Journal of Bacteriology

The Plant Cell

Plant Physiology....

PNAS