Database Systems in Bioinformatics



Several Exercises

- Retrieve Protein Sequences by SWISS-PROT
- Compare Protein Sequences by **BLAST**
- Use Entrez for Searching
- Search in the **Tree of Life**

ExPASy and SWISS-PROT

- Created and managed by one of the pioneers of protein informatics: **Prof. Amos Bairoch**
- ExPASy server: world-leading resource of protein information
- **SWISS-PROT**: name of the database
- Numerous analysis tools

Example 1

- 1. Point your favorite Internet browser to the website www.expasy.org
- 2. Inform yourself about the abbreviation **ExPASy**

Search (Swisshort RANK) in the Case of the Swisshort RANK is specific to the Swisshort Rank of Bioinformatics (SB) is backado to be analysis of protein Rank is specific to the Swisshort Rank of Bioinformatics (SB) is backado to be analysis of protein sequences and forwards as well as 2 O FACE (Disclammer Rank oncol).

General Case of the Swisshort Rank of the Swisshort Rank of Bioinformatics (SB) is backado to be analysis of protein sequences and forwards as well as 2 O FACE (Disclammer Rank oncol).

General Case of the Swisshort Rank of the Swisshort Rank of Bioinformatics (SB) is backado to be analysis of protein sequences and forwards as well as 2 O FACE (Disclammer Rank oncol).

General Case of the Swisshort Rank of the Swisshort Rank

- Protein inoviedgebase PROSITE - Protein families and domains SWISS-20PAGE - Two-dimensional polyacrylamide ge electrophoresis ENZYME - Enzyme nomenclature SWISS-MODEL Repository - Automatically generated
- SWISS-MODEL Repository Automatically generate protein models





You get a large amount of information about the dUTPase protein of E. Coli

Search in UniProt Knowledgebase (Swiss-Prot and TrEMBL) for: dUTPase coli

UniProtKB/Swiss-Prot Release 54.0 of 24-Jul-2007 UniProtKB/TrEMBL Release 37.0 of 24-Jul-2007

Number of sequences found in UnProt Knowledgebase (Swiss-Proc_m and TrEMEL)_{pg} 3
 Note that the selected sequences can be saved to a file to be later reminied to do so, go to the bottom of this page
 For more directed searches, you can use the Sequence Retrieval System SRS.

Search in UniProtKB/Swiss-Prot: There are matches to 3 out of 276256 entries

DUT_EC057 (P64007) Deexyndrae 5-septosphate nucleotiddhydislase (EC 3 6 1 23) (dJTPase) (dJTP pyrophosphatase). (GENE Nameed.g. Conterest_Curalitationes/2504, Ec94515) - Escherichia ciki 0157147

-	
7 Click e.g.	on DUT_ECOLI (P06968)
Note: most headings are clickable, even	If they don't appear as links. They link to the user menual or other documents.
Entry information	
Entry name	DUT_ECOLI
Primary accession number	P06968
Secondary accession number	Q2M7V4
Integrated into Swiss-Prot on	April 1, 1988
Sequence was last modified on	April 1, 1968 (Sequence version 1)
Name and origin of the protein	July 24, 2007 (Entry Version 76)
Protein name	Desvauidine 5° trinhosphate publicatide/publicate
Support	EC 3.6.1.23
Ognorigina	dUTPase
	dUTP pyrophosphatase
Gene name	Name: dut
	Synonyms: dnaS, sof
From	Escharichia.coli [TavD: 562194AMAP restacma]
Taxonomy	Bacteria, Proteobacteria, Gammaproteobacteria, Enterobacteriales, Enterobacteriaceae;
Drotsin evictorice	Eschenchia.
Proven existence	 Exitinguida de broteini liavae

Example 1

Top section = general information:

- Entry name
- Unique identifier (this code is worth writing down because it is used to cross-reference related entries in other databases)

Further entry information

Entry name	DUT_ECOLI
Primary accession number	P06968
Secondary accession number	Q2M7V4
Integrated into Swiss-Prot on	April 1, 1988
Sequence was last modified on	April 1, 1988 (Sequence version 1)
Annotations were last modified on	July 24, 2007 (Entry version 78)

Example 1

Top-middle section = biochemical description

- Standard name
- International Enzyme Committee number
- Synonyms
- List of bibliographic references

Protein name	Deoxyuridine 5'-triphosphate nucleotidohydrolase
Synonyms	EC 3.6.1.23 dUTPase dUTP pyrophosphatase
Gene name	Name: dut Synonyms: dnaS. sof OrderedLocusNames: b3840, JW3615
From	Escherichia coli [TaxID: 562] [HAMAP proteome]
Taxonomy	Bacteria, Proteobacteria, Gammaproteobacteria, Enterobacteriales, Enterobacteriaceae, Escherichia.
Protein existence	1: Evidence at protein level:

Bottom-middle sectionSeries of links to various functional classification schemes

Example 1

Sequence	databases
EMBL	X01714; CAA25859 1; Genomic_DNA_EMBL/GenBank/DDBJ[CoDingSequence] V01578; CAA24897 1; Genomic_DNA_EMBL/GenBank/DDBJ[CoDingSequence] L10228; AAA61993 1; Genomic_DNA_EMBL/GenBank/DDBJ[CoDingSequence] U00096; AAC76664 1; Genomic_DNA_EMBL/GenBank/DDBJ[CoDingSequence] AP009048; BAE77652 1; Genomic_DNA_EMBL/GenBank/DDBJ[CoDingSequence]
PIR	A30388; WPECDU.



8. To get **FASTA** format, click the FASTA format button (right on the bottom)

>P06968|DUT_ECOLI Decxyuridine 5'-triphosphate nucleotidohydrolase - Escherichia coli. HKKIDKWILDPFVGKEFEPLPTATSGSAGLDLARCLMDAVELAFODTTLVPFGLATHIAD PSLAMMEPSEGGEMEGUTLGMUCLIEDSTVGCOLMISVUNRGODSFTIOPGEFIACMI FVPVVQAEFNLVEDFDATDRGEGGFGHSGRQ

 For further studies (in our next example) copy the amino-acid sequence and paste it into the simple text file *dutpasecoli.txt*.

Exercise 1

Find the Human Insulin Protein

BLAST

- = Basic Local Alignment Search Tool
- Great sequence comparison tool
- Information can be used for: prediction of protein function, 3-D-structure, identification of homologues in other organisms
- Server: National Center for Biotechnologie Information = NCBI

Example 2

1. Point your browser to the website www.ncbi.nih.gov/BLAST/

BUBLAST Home	larity habusan historical sequences man	News
ULAS 1 man regions of unitary services accounces the services the services the services of the service of the s		Old BLAST Web Pages to be deleted Jame 19th 2007 As perviously announced access to the old pages will be
Chose a spacies genome to search, or list <u>all genemic.BLAST.databases</u> Bisman Bisman Bisman Bisman Bisman Bisman Bisman Bisman Bisman Bisman Bisman Bisman Bisman Bisman Bisman Bisman Bisman B	2007-06-01 12:15:00	
 Arabidopsis thaliana 	n - anno	Tip of the Day
Basic BLAST		Now to Search Custom Databases in Web Illiast Italias Entras







Scroll down until you reach a list of	of sequ	ences.
Sequences producing significant alignments:	(Bits)	Value
<pre>refiNP 418097.11 deoxyuridinetriphosphatase [Escherichia col</pre>	1 300	1e-00 G
pdb/1DUP/A Chain A, Deoxyuridine 5'-Triphosphate Nucleotido	H 300	1e-80 🗳
pdb/1EUW/A Chain A, Atomic Resolution Structure Of E. Coli I	utpa 300	1e-80 🖻
<pre>pdb/1RNJ & Chain &, Crystal Structure Of Inactive Butant Dut</pre>	p 299	3e-00 🗳
<pre>ref(NP_290220.1) deoxyuridine 5'-triphosphate nucleotidohydr</pre>	o <u>298</u>	7e-80 C
<pre>ref[NP_756326.1] deoxyuridine 5'-triphosphate nucleotidohydr</pre>	0 298	9e-80
<pre>ref(YP_543143.1) deoxyuridine 5'-triphosphate nucleotidohydr</pre>	297	2e-79
ref[YP_312545.1] deoxyuridinetriphosphatase [Shigella sonnei	296	2e=79
ref YP_691203.1] deoxyuridinetriphosphatase [Shigella flexne	r 296	3e-79
ref NP 709419.1] deoxyuridine 5'-triphosphate nucleotidohydr	296	3e-79
rering 450195.11 decoupriding \$1 bricksphate nucleotidohydr	203	20-75
ref VP 001174040 11 decoverridine 51-triphosphate nucleotidonydr	270	36-75
refive coll337632.11 decovaridine S'-triphosphate nucleotidol	278	94-74
ref 2P 01537436.11 deoxyuridine 5'-triphosphate nucleotidohy ref 2P 00822491.11 C000756: dUTPase [Yersinia bercovieri ATC	d 266 C 43 264	3e-70 1e-69
ref:YP_001004455.11 deoxyuridine 5'-triphosphate nucleotidob ref:ZP_00528007.11 C000756: dUTPase [Tersinia frederiksenii	ATCC 263	1e-69 2e-69
ref[2P_00825806.1] COG0756: dUTPase [Tersinia mollaretii ATC	C 43 263	2e-69
<pre>ref[NP_667437.1] deoxyuridine 5'-triphosphate nucleotidohydr</pre>	262	4e-69
refinp_932021.11 deoxyuridine 5'-triphosphate nucleotidohydr	:0 <u>262</u>	5e=69

- List contains the sequences with significant similarity
- Ranked by decreasing score values
- Using a protein sequence taken from a database (like in this example) the best matching protein is the one that you started with
- Score value depends on the length of the most similar segments between 2 sequences







Entrez

- The Life Sciences Search Engine
- 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

<text><text><image><image><text>





Example 3

First Section of the flat file: Header

LOCUS DEFINITION	X01714 1609 bp DNA linear BCT 12-SEP-1993 E. coli dut gene for dUTPase (E 3.6.1.23) (deoxyuridine 5'-triphosphate nucleotidohydrolase).			
ACCESSION	X01714			
VERSION	X01714.1 GI:41296			
KEYWORDS	dUTPase; unidentified reading frame.			
SOURCE	Escherichia coli			
ORGANISM	Escherichia coli			
	Bacteria; Proteobacteria; Gammaproteobacteria; Enterobacteriales; Enterobacteriaceae; Escherichia.			
REFERENCE	1 (bases 1 to 1609)			
AUTHORS	Lundberg,L.G., Thoresson, H.O., Karlstrom, O.H. and Nyman, P.O.			
TITLE	Nucleotide sequence of the structural gene for dUTPase of			
	Escherichia coli K-12			
JOURNAL	EMBO J. 2 (6), 967-971 (1983)			
PUBMED	6139280			
COMMENT	Data kindly reviewed (25-NOV-1985) by L. Lundberg.			

Example 3

- LOCUS: locus name, size of the Nucleotide sequence in base pairs, nature of the molecule (here DNA), its topology (linear or circular)
- **DEFINITION**: short definition, here *E. coli dut gene*
- ACCESSION: unique identifier
- VERSION: synonymous or past numbers
- **KEYWORDS**: terms that characterize the entry

- SOURCE: common name of the relevant organism to which the sequence belongs
- **ORGANISM**: more complete identification of the organism
- **REFERENCE**: authors, titles, journals, ...
- **COMMENT**: info that doesn't fit in the previous sections

Example 3

Middle Section: Feature table

FEATURES	Location/Qualifiers	
source	11609	
	/organism="Escherichia coli"	
	/mol_type="genomic DNA"	
	/db_xref="taxon:562"	
promoter	286291 /note="-35 region"	
promoter	310316 /note="-10 region"	
misc_feature	322324 /note="put. transcription start region"	
RBS	330333	
	/note="put. rRNA binding site"	
CDS	343798	
	/note="unnamed protein product; dUTP-ase (aa 1-151)"	
	/codon_start=1	
	/transl_table=11	
	/protein_id="CAA25859.1"	
	/db_xref="GI:41297"	
	/db xref="GOA:P06968"	
	/db_xref="InterPro:IPR008180"	
		22/2

Example 3

- **source**: indicates the origin of specific regions of the sequence
- **promoter**: shows the precise coordinates of a promoter element, in X01714 a -35 box is indicated from position 286 to 291
- **Misc feature**: (miscellaneous feature) indicates the putative location of the transcription start (mRNA synthesis)
- **RBS**: (Ribosome Binding Site) indicates the location of the last upstream element
- **CDS**: (CoDing Segment) describes the gene open reading frame

Exercise 3

Find the complete genome for the mitochondrion of the blue whale





Example 4: Tree of Life

Find the **path from the origin to human** by clicking

Eukaryotes Animals

.

- Synapsida
- Bilateria
- Deuterostomia
- Chordata
- Craniata Vertebrata
- Primates
- Catarrhini Hominidae

Amniota

 Therapsida Mammalia

Eutheria

- Gnathostomata
- Sarcopterygii . Terrestrial Vertebrates .
- Homo Homo sapiens

Literature

■ Claverie, J-M, Notredame, C.: *Bioinformatics* for Dummies. Wiley Publishing, Inc. 2003

example	pages
1	46 – 49
2	63 – 67
3	78 - 82