

~~An introduction to biological databases~~

Marie-Claude.Blatter@isb-sib.ch



EMBnet

MCB, feb 2005



What is a database ?

A collection of

- structured

- searchable (index)

→ table of contents

- updated periodically (release)

→ new edition

- cross-referenced (hyperlinks)

→ links with other db

data

- Includes also associated tools (software)
necessary for db access/query, db updating, db information insertion, db information deletion....



EMBnet

MCB, feb 2005

Why biological databases ?

- Exponential growth in biological data.

- Data (genomic sequences, 3D structures, 2D gel analysis, MS analysis, Microarrays....) are no longer published in a conventional manner, but directly submitted to databases.



EMBnet

MCB, feb 2005

Distribution of databases

- Books, articles 1968 -> 1985
- Computer tapes 1982 -> 1992
- Floppy disks 1984 -> 1990
- CD-ROM 1989 -> ?
- FTP 1989 -> ?
- On-line services 1982 -> 1994
- • WWW 1993 -> ?
- DVD 2001 -> ?



EMBnet

MCB, feb 2005

Some statistics and remarks

- More than 1000 different 'biological' databases
- Variable size: <100Kb to >10Gb
 - DNA: > 10 Gb
 - Protein: 1 Gb
 - 3D structure: 5 Gb
 - Other: smaller
- Update frequency: daily to annually
- How to find them ?
 - Amos' links: www.expasy.org/alinks.html
 - Biohunt: <http://www.expasy.org/BioHunt/>
 - Google: <http://www.google.com/>



EMBnet

MCB, feb 2005

Google Web Images Groups News News Froogle more >

protein database Search Advanced Search Preferences

Web

The RCSB Protein Data Bank
... 25-Jan-2005 Winter Newsletter Published: A Look at 2004 Depositions, Early Protein Structures, Beta Site Features, and More The Winter 2005 issue of the RCSB ...
www.rcsb.org/pdb/ - 17k - 31 Jan 2005 - [Cached](#) - [Similar pages](#)

ExPASy Proteomics Server
... SWISS-MODEL Repository - Automatically generated protein models. GermOnLine - Knowledgebase on germ cell differentiation, Ashbya Genome Database; Links to many ...
www.expasy.org/ - 17k - 31 Jan 2005 - [Cached](#) - [Similar pages](#)

ExPASy - Swiss-Prot and TrEMBL
The UniProt Knowledgebase consists of Swiss-Prot; a curated protein sequence database which strives to provide a high level of annotation (such as the ...
www.expasy.org/sprot/ - 12k - 31 Jan 2005 - [Cached](#) - [Similar pages](#)
[[More results from www.expasy.org](#)]

Entrez Protein
The Entrez Protein database is a collection of protein sequences compiled from several sources, including SwissProt, PRF, PDR, and translations from ...
www.ncbi.nlm.nih.gov/entrez/query.fcgi?db=Protein - [Similar pages](#)

www.pdb.bnl.gov Redirection Page
The www.pdb.bnl.gov web site has changed location. The new site is available at <http://www.rcsb.org>. Please enter this URL in the ...
www.pdb.bnl.gov/ - 2k - [Cached](#) - [Similar pages](#)

Protein Information Resource Home Page - Protein Sequence Database ...
... in collaboration with the Munich Information Center for Protein Sequences(MIPS) and the Japanese International Protein Sequence Database (JIPID) maintains the ...
pir.georgetown.edu/ - 18k - [Cached](#) - [Similar pages](#)



EMBnet

MCB, feb 2005

~~✓ Protein Sequences is different than Nucleotide Sequence~~

The ten important bioinformatics databases *

GenBank/DDJB/EMBL	www.ncbi.nlm.nih.gov	Nucleotide sequences
Ensembl	www.ensembl.org	Human/mouse genome
PubMed	www.ncbi.nlm.nih.gov	Literature references
NR	www.ncbi.nlm.nih.gov	Protein sequences
Swiss-Prot	www.expasy.org	Protein sequences
InterPro	www.ebi.ac.uk	Protein domains
OMIM	www.ncbi.nlm.nih.gov	Genetic diseases
Enzymes	www.expasy.org	Enzymes
PDB	www.rcsb.org/pdb/	Protein structures
KEGG	www.genome.ad.jp	Metabolic pathways

*according to the « Bioinformatics for dummies »



EMBnet

MCB, feb 2005

Categories of databases for Life Sciences

Sequences (DNA, protein)

- Genomics
- Mutation/polymorphism
- Protein domain/family (----> tools)
- Proteomics (2D gel, Mass Spectrometry)
- 3D structure
- Metabolism
- Bibliography
- 'Others' (Microarrays, Protein protein interaction...)



EMBnet

MCB, feb 2005



Categories of databases for Life Sciences

- Sequences (DNA, protein)
 - Genomics
 - Mutation/polymorphism
 - Protein domain/family (----> tools)
 - Proteomics (2D gel, Mass Spectrometry)
 - 3D structure
 - Metabolism
 - Bibliography
 - 'Others' (Microarrays, Protein protein interaction...)

Ideal minimal content of a sequence database entry

- 1 : Sequences !!
- 2 : Accession number (AC) (unique identifier)
- 3 : Taxonomic data
- 4 : References
- 5 : ANNOTATION/CURATION Imre
- 6 : Keywords
- 7 : Cross-references a
- 8 : Documentation



EMBnet

MCB, feb 2005

Sequence Databases: some « technical » definitions

Data storage management:

- flat file: text file, human readable
- relational database (e.g., Oracle, Postgres)
- object oriented database

Sequence format (for BLAST, prediction tools...)

- Fasta, RAW
- GCG
- NBRF/PIR
- MSF....
- standardized format ?



EMBnet

MCB, feb 2005

Sequence database : format

SWISS-PROT (protein db) (flat file)

✓ Accession number

```
ID EPO_HUMAN      STANDARD;      PRT;    193 AA.  
AC P01588; Q9UHA0; Q9UEZ5; Q9UDZ0;  
DT 21-JUL-1986 (Rel. 01, Created)  
DT 21-JUL-1986 (Rel. 01, Last sequence update)  
DT 20-AUG-2001 (Rel. 40, Last annotation update)  
DE Erythropoietin precursor.  
GN EPO.  
OS Homo sapiens (Human).  
OC Eukaryota; Metazoa; Chordata; Craniata; Vertebrata; Euteleostomi;  
OC Mammalia; Eutheria; Primates; Catarrhini; Hominidae; Homo.  
OX NCBI_TaxID=9606;  
RN [1]  
RP SEQUENCE FROM N.A.  
RX MEDLINE=85137899; PubMed=3838366;  
RA Jacobs K., Shoemaker C., Rudersdorf R., Neill S.D., Kaufman R.J.,  
RA Mufson A., Seehra J., Jones S.S., Hewick R., Fritsch E.F.,  
RA Kawakita M., Shimizu T., Miyake T.;  
RT "Isolation and characterization of genomic and cDNA clones of human  
RT erythropoietin.";  
RL Nature 313:806-810(1985).  
...  
CC -!- FUNCTION: ERYTHROPOIETIN IS THE PRINCIPAL HORMONE INVOLVED IN THE  
CC REGULATION OF ERYTHROCYTE DIFFERENTIATION AND THE MAINTENANCE OF A  
CC PHYSIOLOGICAL LEVEL OF CIRCULATING ERYTHROCYTE MASS.  
CC -!- SUBCELLULAR LOCATION: SECRETED.  
CC -!- TISSUE SPECIFICITY: PRODUCED BY KIDNEY OR LIVER OF ADULT MAMMALS.  
CC AND BY LIVER OF FETAL OR NEONATAL MAMMALS.  
CC -!- PHARMACEUTICAL: Available under the names EpoGen (Amgen) and  
CC Procrit (Ortho Biotech).  
...  
DR EMBL; X02158; CAA26095.1; -.  
DR EMBL; X02157; CAA26094.1; -.  
DR EMBL; M11319; AAA52400.1; -.  
DR EMBL; AF053356; AAC78791.1; -.  
DR EMBL; AF202308; AAF23132.1; -.  
DR EMBL; AF202306; AAF23132.1; JOINED.  
...  
KW Erythrocyte maturation; Glycoprotein; Hormone; Signal; Pharmaceutical.
```

✓ Taxonomy

✓ Reference

✓ Annotations (comments)

✓ Cross-references

MeSH

Keywords

✓ Another DB
to query
list

Sequence database: format

Annotations (features)

```
FT SIGNAL      1      27  
FT CHAIN      28     193      ERYTHROPOIETIN.  
FT PROPEP     190     193      MAY BE REMOVED IN PROCESSED PROTEIN.  
FT DISULFID    34     188  
FT DISULFID    56      60  
FT CARBOHYD   51      51      N-LINKED (GLCNAC...).  
FT CARBOHYD   65      65      N-LINKED (GLCNAC...).  
FT CARBOHYD  110     110      N-LINKED (GLCNAC...).  
FT CARBOHYD  153     153      O-LINKED (GALNAC...).  
FT VARIANT    131     132      SL -> NF (IN AN HEPATOCELLULAR  
FT          CARCINOMA).  
FT          /FTId=VAR_009870.  
FT VARIANT    149     149      P -> Q (IN AN HEPATOCELLULAR CARCINOMA).  
FT          /FTId=VAR_009871.  
FT CONFLICT    40      40      E -> Q (IN REF. 1; CAA26095).  
FT CONFLICT    85      85      Q -> QQ (IN REF. 5).  
FT CONFLICT   140     140      G -> R (IN REF. 1; CAA26095).  
**  
** ##### INTERNAL SECTION #####  
**CL 7q22;  
SQ SEQUENCE 193 AA; 21306 MW; C91F0E4C26A52033 CRC64;  
MGVHECPAWL WLLLSLLSLP LGLPVLGAPP RLICDSRVLE RYLLEAKEAE NITTGCAEH  
SLNENITVPD TKVNFYAWKR MEVGQQAVEV WQGLALLSEA VLRGQALLVN SSQPWEPLQL  
HVDKAVSGLR SLTTLLRALG AQKEAISPPD AASAAPLRTI TADTFRKLF R VYSNFLRGKL  
KLYTGEACRT GDR  
//
```

Sequence



' I < you
(hash) comment

Sequence database: format

1. ~~list~~

...The fasta format:

> My_Sequence_Name

MGVHECPAWLWLLSLLSLPLGLPVLGAPPRLICDSRVLEAKAEAE
NITTGCAEHCSLNENITVPDTKVNFYAWKRMEVGQQAVEVWQGLALLSEA
VILRGQALLVNNSQPWEPLQLHVDKAVSGRLSLTLLRALGAQKEAISPPD
AASAAPLRTITADTFRKLFRVYSNFLRGKLKLYTGEACRTGDR

> sequence
right angle bracket



EMBnet

MCB, feb 2005

✓

Database 1a: nucleotide sequences

- The 3 main public nucleic acid sequence databases are EMBL (Europe) / GenBank (USA) / DDBJ (Japan)
« different views of the same data set » within 2 to 3 days (since 1990)
- EMBL: since 1982
- Specialized databases for the different types of RNAs (i.e. tRNA, rRNA, tm RNA, uRNA, etc...)
- 3D structure (DNA and RNA) → PDB **(Protein Data Bank)**
- Others: Aberrant splicing db; Eukaryotic promoter db (EPD); RNA editing sites, Multimedia Telomere Resource



EMBnet

MCB, feb 2005

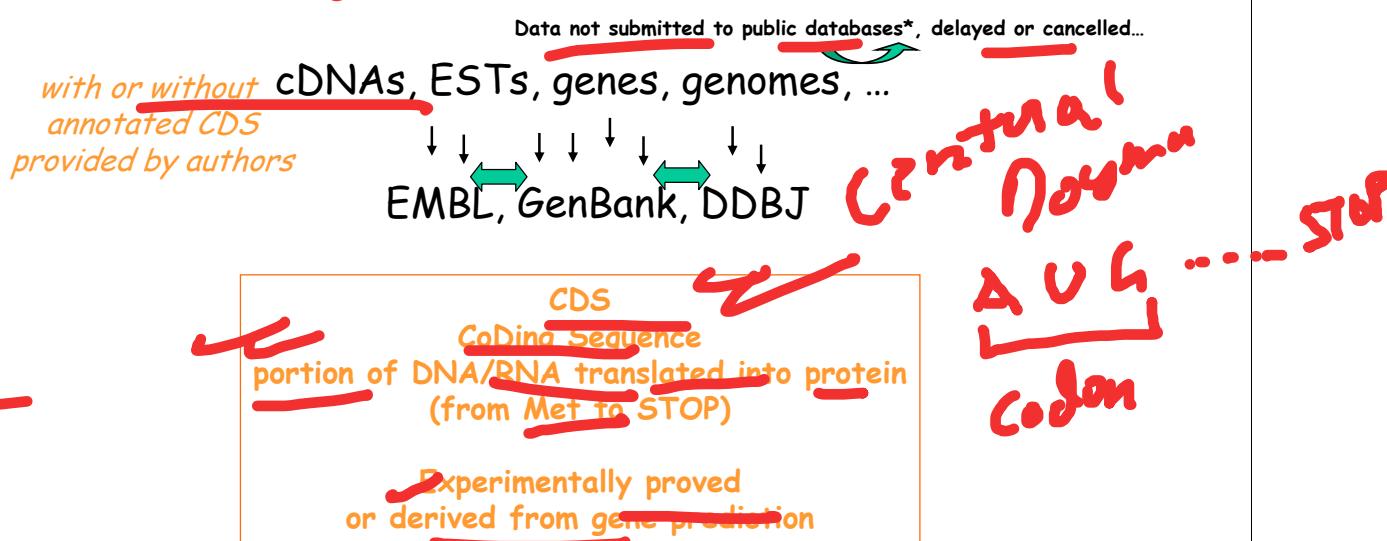
Nucleotide and associated topics databases

- [EMBL](#) - EMBL Nucleotide sequence db (EBI)
- [Genbank](#) - GenBank Nucleotide Sequence db (NCBI)
- [DDBJ](#) - DNA Data Bank of Japan
- [dbEST](#) - dbEST (Expressed Sequence Tags) db (NCBI)
- [dbSTS](#) - dbSTS (Sequence Tagged Sites) db (NCBI)
- [NDB](#) - Nucleic Acid Databank (3D structures)
- [BNAASDB](#) - Nucleic acid structure db from University of Pune
- [AsDb](#) - Aberrant Splicing db
- [ACUTS](#) - Ancient conserved untranslated DNA sequences db
- [Codon Usage Db](#)
- [EPD](#) - Eukaryotic Promoter db
- [HOVERGEN](#) - Homologous Vertebrate Genes db
- [IMGT](#) - ImMunoGeneTics db [Mirror at [EBI](#)]
- [ISIS](#) - Intron Sequence and Information System
- [RDP](#) - Ribosomal db Project
- [gRNAs db](#) - Guide RNA db
- [PLACE](#) - Plant cis-acting regulatory DNA elements db
- [PlantCARE](#) - Plant cis-acting regulatory DNA elements db
- [sRNA db](#) - Small RNA db
- [ssu rRNA](#) - Small ribosomal subunit db
- [lsu rRNA](#) - Large ribosomal subunit db
- [5S rRNA](#) - 5S ribosomal RNA db
- [tmRNA Website](#)
- [tmRDB](#) - tmRNA db
- [tRNA](#) - tRNA compilation from the University of Bayreuth
- [uRNADB](#) - uRNA db
- [RNA editing](#) - RNA editing site
- [RNAmod db](#) - RNA modification db
- [SOS-DGDB](#) - Db of Drosophila DNA sequences annotated with regulatory binding sites
- [TelDB](#) - Multimedia Telomere Resource
- [TRADAT](#) - TRAnscription Databases and Analysis Tools
- [Subviral RNA db](#) - Small circular RNAs db (viroid and viroid-like)
- [MPDB](#) - Molecular probe db
- [OPD](#) - Oligonucleotide probe db
- [VectorDB](#) - Vector sequence db (seems dead!)

Amos'links

<http://www.expasy.org/alinks.html#DNA>

Real life of a sequence ...



~~EMBL/GenBank/DDBJ~~

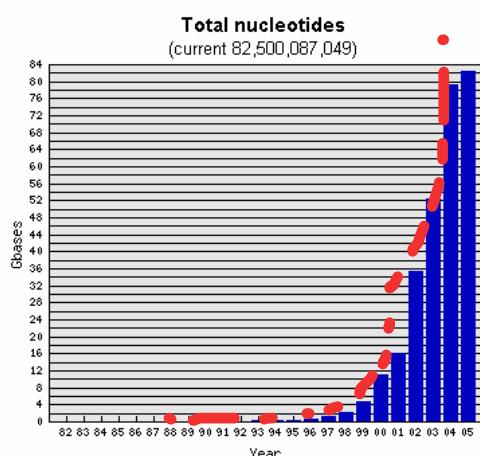
- Serve as archives
- Contain all public sequences derived from:
 - Genome projects (> 80 % of entries)
 - Sequencing centers (cDNAs, ESTs...)
 - Individual scientists (15 % of entries)
 - Patent offices (i.e. European Patent Office, EPO)
- Currently: 46×10^6 sequences, $\sim 80 \times 10^9$ bp;
- Sequences from > 80'000 different species;
- Contribution: EMBL 10 %; GenBank 73 %; DDBJ 17 %



EMBnet

MCB, feb 2005

The tremendous increase in nucleotide sequences



1980: 80 genes fully sequenced !



EMBnet

MCB, feb 2005

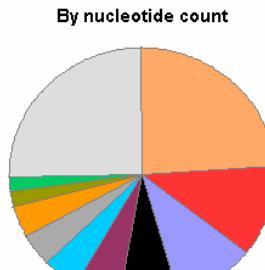
EMBL STATISTICS

This morning the EMBL Database contained **82,500,087,049** nucleotides in
...**239,522** entries.

Breakdown by entry type:

Entry Type	Entries	Nucleotides
Standard	42,213,672	46,346,877,856
Constructed (CON)	403,272	n/a
Third Party Annotation (TPA)	4,549	332,038,478
Whole Genome Shotgun (WGS)	4,618,027	35,770,720,283

New projects:
Environmental sequences
(no taxonomic information)



More than 80'000 species, but...
Human/Mouse/Dar!
Organisms with the highest redundancy!



EMBnet

MCB, feb 2005

an EMBL entry

DNA (genomic)
or
RNA

ID HSERPG standard; genomic DNA; HUM; 3398 BP.
XX
AC X02158;
XX
SV X02158.1
XX
DT 13-JUN-1985 (Rel. 06, Created)
DT 22-JUN-1993 (Rel. 36, Last updated, Version 2)
XX
DE Human gene for erythropoietin
XX
KW erythropoietin; glycoprotein hormone; hormone; signal peptide
XX
OS Homo sapiens (human)
OC Eukaryota; Metazoa; Chordata; Craniata; Vertebrata; Euteleostomi; Mammalia;
OC Eutheria; Primates; Catarrhini; Hominidae; Homo.
XX
RN [1]
RP 1-3398
RX MEDLINE; 85137899.
RA Jacobs K., Shoemaker C., Rudersdorf R., Neill S.D., Kaufman R.J.,
RA Mufson A., Seehra J., Jones S.S., Hewick R., Fritsch E.F., Kawakita M.,
RA Shimizu T., Miyake T.;
RT Isolation and characterization of genomic and cDNA clones of human
erythropoietin;
RL Nature 313:806-810(1985).
XX
DR GDB; 119110; EPO.
DR GDB; 119615; TIMP1.
DR Swiss-Prot; P01588; EPO_HUMAN.
XX
...

standard; RNA; HUM; 1342 BP.
XX
AC X02157;
XX
SV X02157.1
XX
DT 13-JUN-1985 (Rel. 06, Created)
DT 30-MAR-1995 (Rel. 43, Last updated, Version 3)
XX
DE Human mRNA for fetal erythropoietin

keyword
taxonomy
references
Cross-references



EMBnet

MCB, feb 2005

CC Data kindly reviewed (24-FEB-1986) by K. Jacobs

FH Key Location/Qualifiers

FH

FT source 1..3398

FT /db_xref=taxon:9606

FT /organism=Homo sapiens

FT mRNA join(397..627,1194..1339,1596..1682,2294..2473,2608..3327)

FT CDS join(615..627,1194..1339,1596..1682,2294..2473,2608..2763)

FT /db_xref=SWISS-PROT:P01588

FT /product=erythropoietin

FT /protein_id=CAA26095.1

FT /translation=MGVHECPAWLWLLSLLSLPLGLPVLGAPPRLLCDSRVLQRYLLE

FT AKEAEINTTGCAEHCSLNENITVPDTKVNFYAWKRMEVGQQAVEVWQGLALLSEAVLRG

FT QALLVNSSQPWEPLQLHVDKAVSGLRSLLTLLRALGAQKEAISPPDAASAAPLRTITAD

FT TFRKLFRVVSNFLRGKLKYTGACRTGDR

FT mat_peptide join(1262..1339,1596..1682,2294..2473,2608..2763)

FT /product=erythropoietin

FT sig_peptide join(615..627,1194..1261)

FT exon 397..627

FT /number=1

FT intron 628..1193

FT /number=1

FT exon 1194..1339

FT /number=2

FT intron 1340..1595

FT /number=2

FT exon 1596..1682

FT /number=3

FT intron 1683..2293

FT /number=3

FT exon 2294..2473

FT /number=4

FT intron 2474..2607

FT /number=4

FT exon 2608..3327

FT /note=3' untranslated region

FT /number=5

XX

SQ Sequence 3398 BP; 698 A; 1034 C; 991 G; 675 T; 0 other;

agcttctggg cttccagacc cagctacttt ggggaactca gcaaccagg catctctgg

60

sequence

tctccgcca agacgggtt gccccccagg aggtgtccgg gagcccccgg tttcccgat 120

CDS
Coding Sequence
(proposed by submitters)

Annotation

(Prediction or
experimentally determined)



EMBnet

MCB, feb 2005

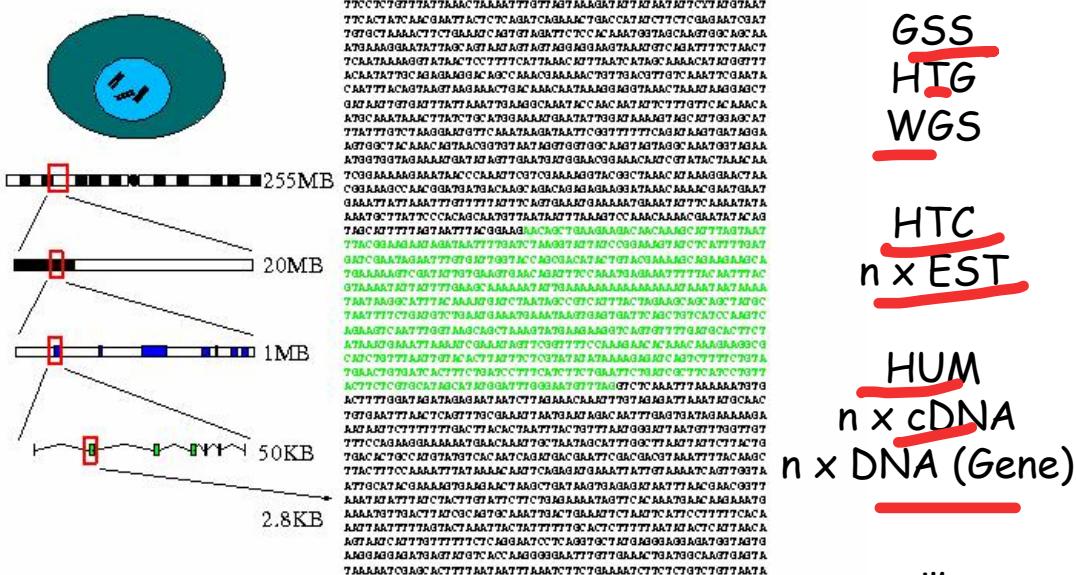
□ 1: NC_000001, Homo sapiens chro...[gi:51511461]

[Click here to see all features and the sequence of this contig record.](#)

LOCUS NC_000001 245522847 bp DNA linear CON 24-Aug-2004
DEFINITION Homo sapiens chromosome 1, complete sequence.
ACCESSION NC_000001
VERSION NC_000001.8 GI:51511461
KEYWORDS HTG.
SOURCE Homo sapiens (human)
ORGANISM Homo sapiens
Eukaryota; Metazoa; Chordata; Craniata; Vertebrata; Euteleostomi;
Mammalia; Eutheria; Primates; Catarrhini; Hominoidea; Homo.
REFERENCE 1 (bases 1 to 245522847)
AUTHORS International Human Genome Sequencing Consortium.
TITLE Finishing the euchromatic sequence of the human genome
JOURNAL Unpublished (2004)
COMMENT GENOME ANNOTATION [REFSEQ](#): Features on this sequence have been produced for build 35 version 1 of the NCBI's genome annotation [see [documentation](#)].
On Aug 24, 2004 this sequence version replaced gi:[42406218](#).
The DNA sequence is part of the third release of the finished human reference genome. It was assembled from individual clone sequences by the Human Genome Sequencing Consortium in consultation with NCBI staff.
FEATURES Location/Qualifiers
source 1..245522847
/organism="Homo sapiens"
/mol_type="genomic DNA"
/db_xref="taxon:9606"
/chromosome="1"
CONTIG join([NT_077402.1](#):1..167280,gap(50000),[NT_077911.1](#):1..40302,
gap(100000),[NT_077912.1](#):1..153649,gap(50000),[NT_077913.3](#):1..798564,
gap(50000),[NT_077914.2](#):1..192022,gap(50000),
[NT_004350.17](#):1..1014565,gap(50000),[NT_004321.16](#):1..1130358,
gap(60000),[NT_004547.16](#):1..1430159,gap(60000),
[NT_021937.17](#):1..7616319,gap(50000),[NT_077382.3](#):1..237250,
gap(50000),[NT_004873.16](#):1..3518496,gap(50000),
[NT_004610.17](#):1..8601969,gap(50000),[NT_037485.3](#):1..2854727,
gap(200000),[NT_004538.16](#):1..1064515,gap(50000),
[NT_004511.17](#):1..16145453,gap(50000),[NT_032977.7](#):1..57683299,
gap(50000),[NT_019273.17](#):1..16783271,gap(50000),
[NT_086586.1](#):1..189539,gap(50000),[NT_077389.3](#):1..398739,
gap(20290000),[NT_079485.2](#):1..164360,gap(300000),
[NT_077932.2](#):1..127263,gap(50000),[NT_077933.1](#):1..170669,gap(100000),
[NT_004434.17](#):1..1002899,gap(50000),[NT_054398.4](#):1..281532,
gap(50000),[NT_054400.3](#):1..1065823,gap(50000),[NT_077936.2](#):1..340928,
gap(150000),[NT_079483.1](#):1..172789,gap(50000),[NT_034401.5](#):1..220313,
gap(50000),[NT_034033.3](#):1..309240,gap(50000),[NT_086599.1](#):1..129426,
gap(50000),[NT_004487.17](#):1..56331646,gap(50000),
[NT_086602.1](#):1..170924,gap(50000),[NT_021877.17](#):1..16682800,
gap(50000),[NT_077939.1](#):1..422453,gap(50000),
[NT_004559.12](#):1..11387711,gap(50000),[NT_004836.16](#):1..13649709,
gap(50000),[NT_032968.7](#):1..174886,gap(80000))

//

~~The big problem = the redundancy~~



 swissprot



EMBnet

MCB, feb 2005

~~EMBL /GenBank/ DDBJ~~

~~Sets of sequence museum where sequences are preserved for eternity as they were determined, interpreted and published originally by their authors~~

(primary sequence repository)

~~The authors have full authority over the content of the entries they submit!~~

(editorial control of the content belongs to the authors)

(exception: TPA, since January 2003)



Submission: FTP email, Webin, etc...

The preferred submission tool is webin

- Check sequence for vector contamination (help)
- Webin - WWW Submission System (help)
- Updating EMBL Entries (help)
- Webin-Align (help)

Documentation

- Information for Submitters (including alignments, genome projects and bulk submissions)
- The DDBJ/EMBL/GenBank Feature Table Definition
- Sequin (including instructions for downloading)

MRF2_HUMAN AAH36831	KLGGYETITARRQWKHIYDELGGNPGSTSAACTRRHYERLILPYERFIKGEEKPLPPI
MRF2_HUMAN AAH36831	KPRKQENSSQENENKTVKSGTKRIKHEIPSKKEKENAPKPQDAEVSSFQEKEQETLIS
MRF2_HUMAN AAH36831	QKSIPEPLPAADMKKIEGYQEFSAKPLASRVDPEKDNETDQGSNSEKVAEEAGEKGPTP
MRF2_HUMAN AAH36831	PLPSAPLAPEKDSALVPGASKQPLTSPSALVDSKQESKLCCFTESPESEPQEASFPSFPT
MRF2_HUMAN	TQPPLANQNETEDDKLPAMADYIANCTVKVDQLGSDDIHNALKTPKVLVVQSFDMFKDK

Protein sequence derived from the traduction of a vector contamination

EMBL/GenBank/DDBJ

- Unexpected information you can find in these db:

```
FT source          1..124
FT              /db_xref="taxon:4097"
FT              /organelle="plastid:chloroplast"
FT              /organism="Nicotiana tabacum"
FT              /isolate="Cuban cahibo cigar, gift from
FT              President Fidel Castro"
```

✓✓

- Or:

```
FT source          1..17084
FT              /chromosome="complete mitochondrial genome"
FT              /db_xref="taxon:9267"
FT              /organelle="mitochondrion"
FT              /organism="Didelphis virginiana"
FT              /dev_stage="adult"
FT              /isolate="fresh road killed individual"
FT              /tissue_type="liver"
```

✓✓

```
FT CDS      complement(45959..47332)
FT          /db_xref="SPTREMBL:Q9UZ71"
FT          /note="PAB2386"
FT          /transl_table=11
FT          /product="4-AMINOBUTYRATE qui se dilate AMINOTRANSFERASE
(EC 2.6.1.19)"
FT          /protein_id="CAB50188.1"
FT          /translation="MDYPRIVVNPPGPKAKELEIREKRLSTGIGVKLPLVPKRGFGP
FIEDVDGNVFIDFLAGAAAAGTGYSHPKLVAKVEQVELIQHSMIGYTHSERAIRVAEK
LVKISPIKNSKVLFGLSGSDAVDMAIKVSFKSTRPWI LAFIGAYHGQTLGATSVASFQ
VSQKRGYSPLMPNVFWVPYPNPYRNPGINGYEEPQELVNRVVEYLEDYVFSHVVPPDE
VAAFFAEPIQGDAGIVVPPENFFKELKLLDEHGILLVMDEVQTGIGRTGKWFASEWFE
VKPDMIIFGKGVASGMGLSGVIGREDIMDITSGSALLTPAANPVISAADATLEIIEEE
FT          NiceProt View of TrEMBL: Q9UZ71
```

Printer

[Entry info] [Name and origin] [References] [Comments] [Cross-references] [Keywords] [Features]

Note: most headings are clickable, even if they don't appear as links. They link to the user manual or other documents.

Entry Information

Entry name	Q9UZ71
Primary accession number	Q9UZ71
Secondary accession numbers	None
Entered in TrEMBL in	Release 13, May 2000
Sequence was last modified in	Release 13, May 2000
Annotations were last modified in	Release 26, March 2004

Name and origin of the protein

Protein name	4-aminobutyrate QUI SE DILATE aminotransferase
Synonym	EC 2.6.1.19
Gene name	PYRAB12830 or PAB2386
From	Pyrococcus abyssi [TaxID: 29292]
Taxonomy	Archaea; Euryarchaeota; Thermococci; Thermococcales; Thermococcaceae; Pyrococcus.

The second generation of nucleotide sequence databases

Gene-centric databases

All the sequence information relevant to a given gene
is made accessible at once

i.e. Locus Link/RefSeq

Genome-centric databases

Information about gene sequence, relative position,
strand orientation, biochemical functions...

Information management systems that are able to connect
specialized sequence collection and browsing tools

i.e. Ensembl, TIGR

Gene-centric databases

The screenshot shows the LocusLink interface for the EPOR gene. At the top, there's a sequence logo for the EPOR gene. Below it, a search bar has "LocusLink" selected. The main content area displays "Click to Display mRNA-Genomic Alignments (spanning 6544 bps)" with a link to "Homo sapiens Official Gene Symbol and Name (HGNC)". Below this, the EPOR gene symbol is shown as "EPOR: erythropoietin receptor". A "LocusID: 2057" identifier is present. The "Overview" section contains a "RefSeq Summary" detailing the EPOR receptor's role in erythropoiesis. The "Proteome Summary" notes its membership in the cytokine receptor family. The "Locus Type" is listed as a gene with a protein product, and the "Product" is identified as the "erythropoietin receptor precursor". Navigation links on the left include "LocusLink Home", "EPOR Index", "Top of Page", "Nomenclature", "Overview", "Function", "Relationships", "Map", "RefSeq", "GenBank", and "Links".

New: Replaced by « Entrez Gene »
on March 1, 2005

Collaborators	Download	FAQ	Submit GeneRIF	(All Pubs)	?
Collaborators	Download	FAQ	Submit GeneRIF	(All Pubs)	?
GeneRIF: Gene References into Function:					
Collaborators	Download	FAQ	12118093	• functional significance of expression in breast cancer	
Collaborators	Download	FAQ	12027890	• evidence for p142;p-Y431 being a new high affinity binding site for SOCS-3 on the EpoR	
Collaborators	Download	FAQ	11884148	• Amino acid determinants of beta-hairpin conformation in erythropoietin receptor agonist peptides derived from a phage display library	
Collaborators	Download	FAQ	12021194	• The extracellular binding site for ERF is now characterized. The site is located in the membrane proximal, extracellular part of the receptor. ERF binds to a region on the EPOR that contains the same sequence as ERF	
Gene Ontology™:					
Term	Evidence	Source	Pub		
signal transduction	NR	Proteome			
erythropoietin receptor	E	Proteome	pm		
integral plasma membrane protein	E	Proteome	pm		
Other Ontologies:					
Term	Evidence	Source	Pub		
Differentiation	NR	Proteome	pm		
Dendrite	E	Proteome	pm		
Other development	NR	Proteome	pm		
Cell fate specification	NR	Proteome	pm		
Integral membrane	NR	Proteome	pm		
Cell body (soma)	E	Proteome	pm		
Receptor (signalling)	NR	Proteome	pm		



FMRnet

MGR feb 2005

~~LocusLink is tightly linked to RefSeq
("independent curated resources")~~

«Reference Sequences»

- for RNA (NM_)

- for genomic (NT_)
 - for protein (NP_)

Links to all the sequences found in
EMBL/GenBank/DDBJ
corresponding to this gene

1: NM_000121. Homo sapiens eryt...[gi:4557561]

LOCUS	NM_000121	1849 bp	mRNA	linear	PRI 23-AUG-2004
DEFINITION	Homo sapiens erythropoietin receptor (EPOR), mRNA.				
ACCESSION	NM_000121				
VERSION	NM_000121.2	GI:4557561			
KEYWORDS	.				
SOURCE	Homo sapiens (human)				
ORGANISM	Homo sapiens				
Eukaryota; Metazoa; Chordata; Craniata; Vertebrata; Euteleostomi; Mammalia; Eutheria; Primates; Catarrhini; Hominidae; Homo.					
REFERENCE	1 (bases 1 to 1849)				
AUTHORS	AcS, G., Xu, X., Chu,C., AcS,P. and Verma,A.				
TITLE	Prognostic significance of erythropoietin expression in human endometrial carcinoma				
JOURNAL	Cancer 100 (11), 2376-2386 (2004)				
PUBLMED	15160341				
REMARK	GeneRIF: Increased erythropoietin receptor expression is associated with advanced-stage disease, lymphovascular invasion, lymph node metastasis of endometrial carcinoma				
REFERENCE	2 (bases 1 to 1849)				
AUTHORS	Eccles,T.G., Patel,A., Verma,A., Nicholson,D., Lukes,Y., Tuttle,R.M. and Francis,G.L.				
TITLE	Erythropoietin and the erythropoietin receptor are expressed by papillary thyroid carcinoma from children and adolescents. Expression of erythropoietin receptor might be a favorable prognostic indicator				
JOURNAL	Ann. Clin. Lab. Sci. 33 (4), 411-422 (2003)				
PUBLMED	14584755				
REMARK	GeneRIF: Majority of papillary thyroid carcinomas (PTC) from children and adolescents express EPO-R, a finding associated with favorable prognostic indicators and a lower risk of recurrence.				
REFERENCE	3 (bases 1 to 1849)				
AUTHORS	Arcasoy,M.O., Jiang,X. and Haroon,Z.A.				
TITLE	Expression of erythropoietin receptor splice variants in human cancer				
JOURNAL	Biochem. Biophys. Res. Commun. 307 (4), 999-1007 (2003)				
PUBLMED	12878211				
REMARK	GeneRIF: Expression of erythropoietin receptor splice variants in tumor cell lines.				
REFERENCE	4 (bases 1 to 1849)				
AUTHORS	Matsuhashi,S., Canis,M., Yokomizo,R., Yaegashi,N., Bruhat,M.A. and Okamura,K.				
TITLE	Expression of erythropoietin and erythropoietin receptor in peritoneal endometriosis				
	Hum. Reprod. 18 (1), 152-156 (2003)				
	12525458				
	GeneRIF: Epo and Epo-R localized within glandular epithelial cells in both peritoneal endometriosis and eutopic endometrium. Epo-R expression lower in black peritoneal lesions.				
	5 (bases 1 to 1849)				

The corresponding RefSeq entry for the mRNA

NCBI Reference Sequence
<http://www.ncbi.nlm.nih.gov/RefSeq/>

RefSeq

Working with whole genome databases:

Genome-centric databases

« Browsing resources »

Remark: Genome-centric databases give usually access to several genomes, but some are « specialized » in particular organisms, i.e. TIGR: bacteria and plants

The screenshot shows the Ensembl Genome Browser homepage. At the top, there are logos for Swissprot, SIB, EMBnet, and The Wellcome Trust Sanger Institute, along with the EBI logo. The date 'MCB, feb 2005' is visible. A large red circle highlights the 'Ensembl' logo and the 'Species - Ensembl v19' section. A red bracket on the right side covers the text 'Ensembl provides a bioinformatics framework to organise biology around the sequences of large genomes.' Below the main header, there's a search bar and a 'About Ensembl' section. The 'Species - Ensembl v19' section lists various species with their assembly details. To the right, there's a 'Have you tried ...?' section featuring the 'Ensembl Chicken Preview Browser'. At the bottom, there's a 'Help and documentation' section with links to news, DAS, Apollo, Help Desk, and Documentation.

Ensembl provides a bioinformatics framework to organise biology around the sequences of large genomes.

Available now are:
human, mouse, rat, fugu,
zebrafish, mosquito,
Drosophila, C. elegans, and
C. briggsae, chicken...

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

Ensembl/martview: example of queries

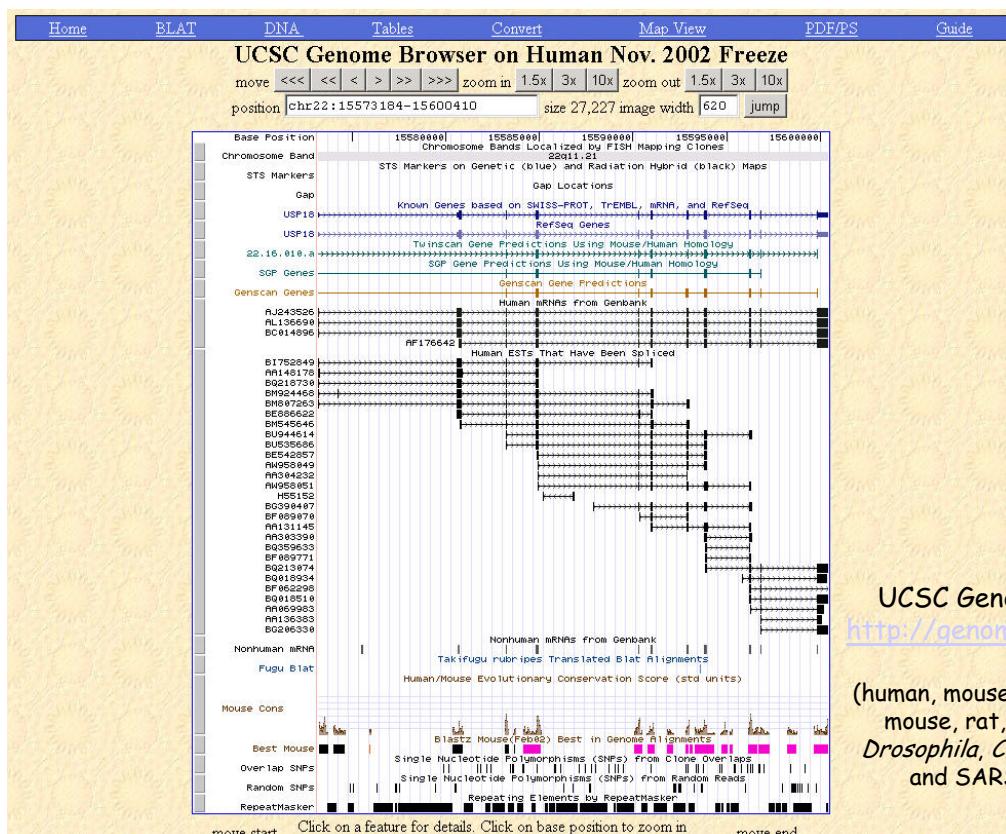
- Retrieve all mouse homologues of human disease genes containing transmembrane domains located between 1p22 and 1q22
- Retrieve the sequences 5kb upstream of all human « known » genes from chromosome 6

....



EMBnet

MCB, feb 2005



UCSC Genome Browser:
<http://genome.cse.ucsc.edu/>

(human, mouse, rat, chimpanzee,
mouse, rat, chicken, *Fugu*,
Drosophila, *C. briggsae*, yeast,
and SARS genomes.)

→

...and plants

<http://www.tigr.org/tdb/>

Database 1b: protein sequences

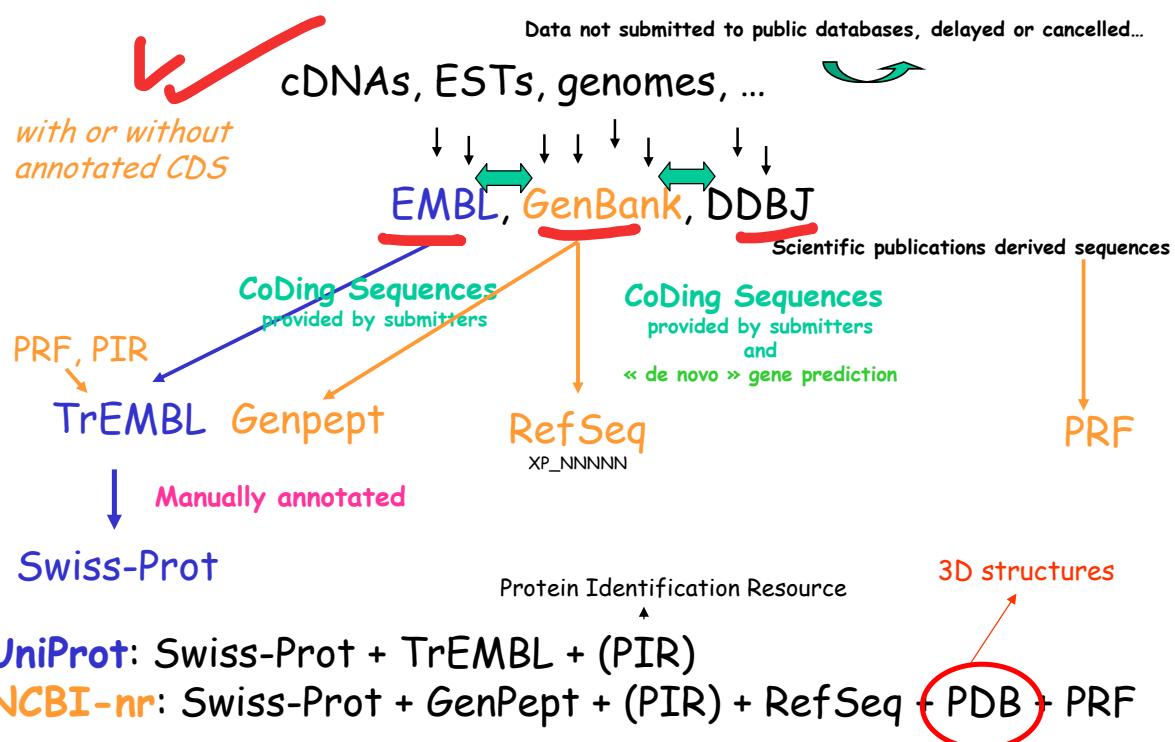
- ~~SWISS-PROT~~: created in 1986 (A.Bairoch) <http://www.expasy.org/sprot/>
- ~~TrEMBL~~: created in 1996; complement to SWISS-PROT; derived from ~~EMBL~~ CDS translations (« proteomic » version of EMBL)
- **(PIR-PSD: Protein Information Resources)**
<http://pir.georgetown.edu/>

- **Genpept: « proteomic » version of GenBank (~TrEMBL)**
- **RefSeq (NP_)**
- **PRF**

- **Many specialized protein databases for specific families or groups of proteins.**

Examples: **AMSDb** (antibacterial peptides), **GPCRDB** (7 TM receptors), **IMGT** (immune system) **YPD** (Yeast) etc.

Real life of a protein sequence ...



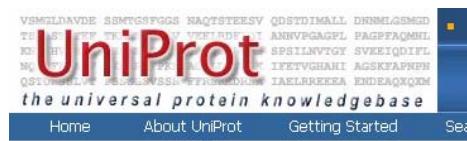
Protein sequence databases

The UniProt pathway

a central ressource for protein sequences
and function...

Since december 15, 2003

Swiss-Prot and TrEMBL constitute the



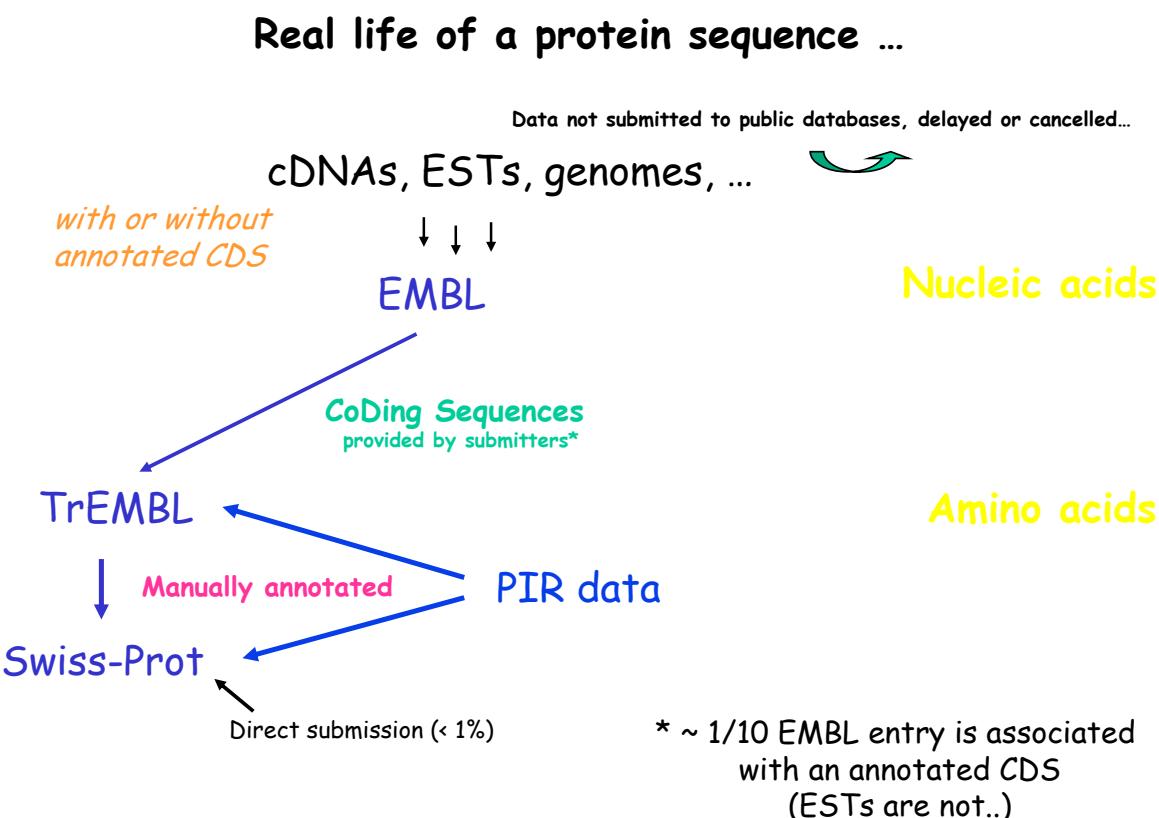
Knowledgebase

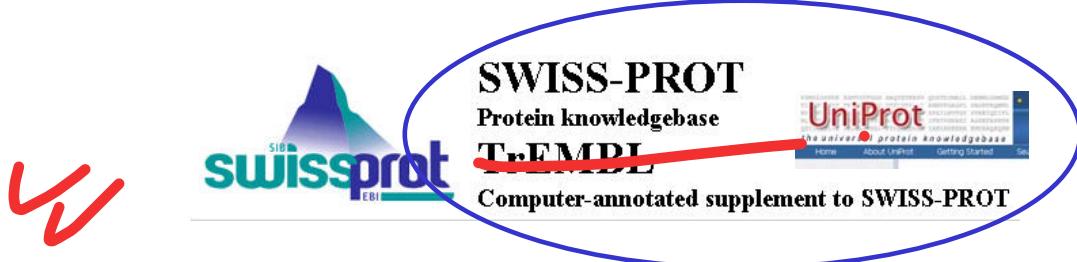
(integration of the PIR data)
(Protein Information Ressource)



EMBnet

MCB, feb 2005





-> give access to

all known* protein sequences

* submitted to the public databases (EMBL, GenBank, DDJB, Swiss-Prot)



EMBnet

MCB, feb 2005

Screen shot of EMBL database showing a protein sequence page. A red circle highlights a specific region of the sequence, and a blue arrow points from this region to the TrEMBL page below.



CDS

Screen shot of TrEMBL database showing a protein sequence page for Q9MMY8. A blue arrow points from the CDS icon on the EMBL page to this TrEMBL page. An illustration of a human brain is shown on the right side of the page.

TrEMBL

Screen shot of Swiss-Prot database showing a protein sequence page for P01508. A blue arrow points from the brain illustration on the TrEMBL page to this Swiss-Prot page.

EMBL

Swiss-Prot

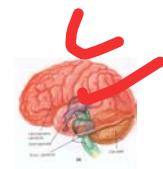
EMBL

```
Q9MMY8 protein (Q9MMY8) (Q9MMY8)
```



TrEMBL: Q9MMY8

Swiss-Prot: P01508



TrEMBL
Once in Swiss-Prot,
no more in TrEMBL
-> Minimal redundancy



Annotation of sequence differences
(conflicts, variants, splicing...)

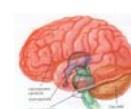
Average of 4.2 independent sequence reports
for each human protein

Swiss-Prot

```
P01508 protein (P01508) (P01508)
```

Swiss-Prot: P01508

TrEMBL: Q9MMY8



TrEMBL



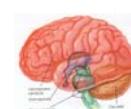
EMBL

```
Q9MMY8 protein (Q9MMY8) (Q9MMY8)
```



TrEMBL: Q9MMY8

Swiss-Prot: P01508



Swiss-Prot

~~Up-to-date sources:~~

✓ Swiss-Prot -> ExPASy
Since 1986 (www.expasy.org):

TrEMBL -> EBI (European Bioinformatics Institute)
Since 1996 (www.ebi.ac.uk/trembl/).



ExPASy Molecular Biology Server - Mozilla

File Edit View Go Bookmarks Tools Window Help

Back Forward Reload Stop http://www.expasy.org/ Search Print

Home Bookmarks SwissDoc SP internal page SRS Sanger SRS EBI PubMed OMIM LocusLink Genew MGD Ensembl Human Genome Google Wise2

ExPASyBar Search Swiss-Prot/TrEMBL by descr. for

PubMed Menu Search PubMed for

Site Map Search ExPASy Contact us

Hosted by SIB Switzerland Mirror sites: Australia Bolivia Canada China Korea Taiwan USA

Search Swiss-Prot/TrEMBL for Go Clear

ExPASy Molecular Biology The ExPASy (Expert Protein Analysis System) structures as well as 2-D PAGE (Disclaimer / Re

Swiss Institute of Bioinformatics (SIB) is dedicated to the analysis of protein sequences and

[Announcements] [Job opening] [Mirror Sites]

Databases

- [Swiss-Prot and TrEMBL](#) - Protein knowledgebase
- [PROSITE](#) - Protein families and domains
- [SWISS-2DPAGE](#) - Two-dimensional polyacrylamide gel electrophoresis
- [ENZYME](#) - Enzyme nomenclature
- [SWISS-3DIMAGE](#) - 3D images of proteins and other biological macromolecules
- [SWISS-MODEL Repository](#) - Automatically generated protein models
- [CD40lbase](#) - CD40 ligand defects
- [SeqAnalRef](#) - Sequence analysis bibliographic references
- [Links to many other molecular biology databases](#)

Tools and software packages

- [Proteomics and sequence analysis tools](#)
 - [Proteomics](#) [PeptideIdent, PeptideMass, ...]
 - [DNA > Protein](#) [Translate]
 - [Similarity searches](#) [BLAST]
 - [Pattern and profile searches](#) [ScanProsite]
 - [Post-translational modification and topology prediction](#)
 - [Primary structure analysis](#) [Protparam, pI/MW, ProScale]
 - [Secondary and tertiary structure prediction](#) [SWISS-MODEL, Swiss-PdbViewer]
 - [Alignment](#) [T-COFFEE, SIM]
 - [Biological text analysis](#)
- [Melanie 4](#) - Software for 2-D PAGE analysis
- [Roche Applied Science's Biochemical Pathways](#)

Education and services

- [The ExPASy FTP server](#)
- [SwissShop](#) - automatically obtain (by email) new sequence entries relevant to your field(s) of interest
- [Masters Degree in Bioinformatics](#)
- [Proteomics courses](#) - two courses covering Separation Sciences & Mass Spectrometry for Proteomics

Documentation

- [What's New on ExPASy](#)
- [SWISS-FLASH](#) electronic bulletins
- [Swiss-Prot documents](#)
- [How to create HTML links to ExPASy](#)
- [Complete table of available documents](#)

ExPASy

EBI

NCBI

In a Swiss-Prot entry, you can expect to find:

- All the names of a given protein (and of its gene);
- Its biological origin with links to the taxonomic databases;
- A selection of references;
- A summary of what is known about the protein: function, alternative products, PTM, tissue expression, disease, etc....;
- Numerous cross-references;
- Selected keywords;
- A description of important sequence features: domains, PTMs, variations, etc.;
- A (often corrected) protein sequence and the description of various isoforms/variants.

ExPASy Home page Site Map Search ExPASy Contact us SWISS-PROT
 Hosted by the HUCC US Mirror site Andhra Ganga Chaitanya Deemed University

NiceProt View of SWISS-PROT: O75144 Protein Summary Quick BLAST search

General Information Basic information Cross-references (CrossRef) (CrossView) (CrossSearch) (CrossList)

Accession number ID, AC, DT lines

Names and taxonomy DE, GN, OC, OS, OG lines

References RN, RP, RC, RX, RA, RL lines

Comments CC lines

Cross-references DR lines

Keywords KW lines

Features FT lines

Sequence SQ lines

View « by default » on the ExPASy server

MCB, feb 2005



Sequence quality

-> 1 gene / 1 species = 1 Swiss-Prot entry

For human: ~ 4,2 different independent sequence reports /gene

SEQ1_HUMAN	MNPRESSIFIEDYLKYFQDQVSREMLLQLLTDEAWNGFILAAELPRDEADELRKALNKLA
SEQ2_HUMAN	-----MLLRLLTDEAWNGFV-AAELPRDEADELRKALNKLA
SEQ1_HUMAN	SHMVMKDKNR-----SKEFPRLKRELEDHIRKLRALAEEVEUVHRGTTIANVVS
SEQ2_HUMAN	SHMVMKDKNRSKDQHROWHSKEFPRLKRELEDHIRKLRALAEEVE-----

Alternative splicing ?

Same gene ?

Polymorphisms ?

Alternative initiation ?

RNA editing ?

Usage of an alternative promoter ?

Selenocystein ?

Fragment ?

Sequencing errors ?

-> Identification and annotation of all sequence differences

Annotation (Comment lines)

- Function(s) and role(s); enzymes:
 - a. Catalytic activity (if EC number)
 - b. Cofactor
 - c. Enzyme regulation
 - d. Pathway
- Subunit (Protein/protein interactions)
- Subcellular location
- Alternative products (alt. splicing, alt. initiation, RNA editing)
- Tissue specificity (Northern and Western results)
- Developmental stage
- Induction
- Domain
- Post-translational modifications (PTM)
- Mass spectrometry
- Polymorphisms
- Disease
- Pharmaceutical
- Miscellaneous
- Similarities
- Caution
- Database (specialized cross-references)



EMBnet

MCB, feb 2005

Annotation/Curation (Comment lines)

Information is derived from:

- Publications; currently Swiss-Prot cites 1'500 different journals. 106 journals are cited more than 100 times.
- Databases;
- Personal communication;
- Prediction;
- Brain storming...



EMBnet

MCB, feb 2005

experimental qualifiers:

« - »: experimentally proved;

« By similarity »: experimentally proved in an ortholog or in another member of the family;

« Probable »: not proved, but realistic;

« Potential »: predicted (bioinformatic tools).



hob

1bV1

Comments

ICOL_HUMAN, Q75144

- **FUNCTION**: LIGAND FOR THE T-CELL-SPECIFIC CELL SURFACE RECEPTOR ICOS. ACTS AS A COSTIMULATORY SIGNAL FOR T-CELL PROLIFERATION AND CYTOKINE SECRETION, INDUCES ALSO B-CELL PROLIFERATION AND DIFFERENTIATION INTO PLASMA CELLS. COULD PLAY AN IMPORTANT ROLE IN MEDIATING LOCAL TISSUE RESPONSES TO INFLAMMATORY CONDITIONS, AS WELL AS IN MODULATING THE SECONDARY IMMUNE RESPONSE BY CO-STIMULATING MEMORY T-CELL FUNCTION (By similarity)
- **SUBCELLULAR LOCATION**: Type I membrane protein (By similarity)
- **ALTERNATIVE PRODUCTS**:
 - Alternative splicing [2 named forms]

Comment: Additional isoforms seem to exist.

Name	1
Isoform ID	Q75144-1
This is the isoform sequence displayed in this entry .	

Name	2
Isoform ID	Q75144-2
Features which should be applied to build the isoform sequence: VSP_002520	
- **TISSUE SPECIFICITY**: ISOFORM 1 IS WIDELY EXPRESSED (BRAIN, HEART, KIDNEY, LIVER, LUNG, PANCREAS, PLACENTA, SKELETAL MUSCLE, BONE MARROW, COLON, OVARY, PROSTATE, TESTIS, LYMPH NODES, LEUKOCYTES, SPLEEN, THYMUS AND TONSIL), WHILE [ISOFORM 2](#) IS DETECTED ONLY IN LYMPH NODES, LEUKOCYTES AND SPLEEN
- **INDUCTION**: CONSTITUTIVE EXPRESSION IS FURTHER ENHANCED BY TREATMENT WITH TNF-ALPHA IN PERIPHERAL BLOOD B-CELLS AND MONOCYTES, WHILE IT IS DECREASED IN DENDRITIC CELLS.
- **SIMILARITY**: BELONGS TO THE IMMUNOGLOBULIN SUPERFAMILY. BTN/MOG SUBFAMILY.
- **SIMILARITY**: Contains 1 immunoglobulin-like V-type domain.
- **SIMILARITY**: Contains 1 immunoglobulin-like C2-type domain.
- **CAUTION**: Ref 4 sequence differs from that shown in position 300 onward for an unknown reason.



Experimental qualifiers:

« - »: experimentally proved;

« By similarity »: experimentally proved in an ortholog or in another member of the family;

« Probable »: not proved but realistic;

« Potential »: predicted (bioinformatic tools).

Comments

BRH2_HUMAN, Q9NY43

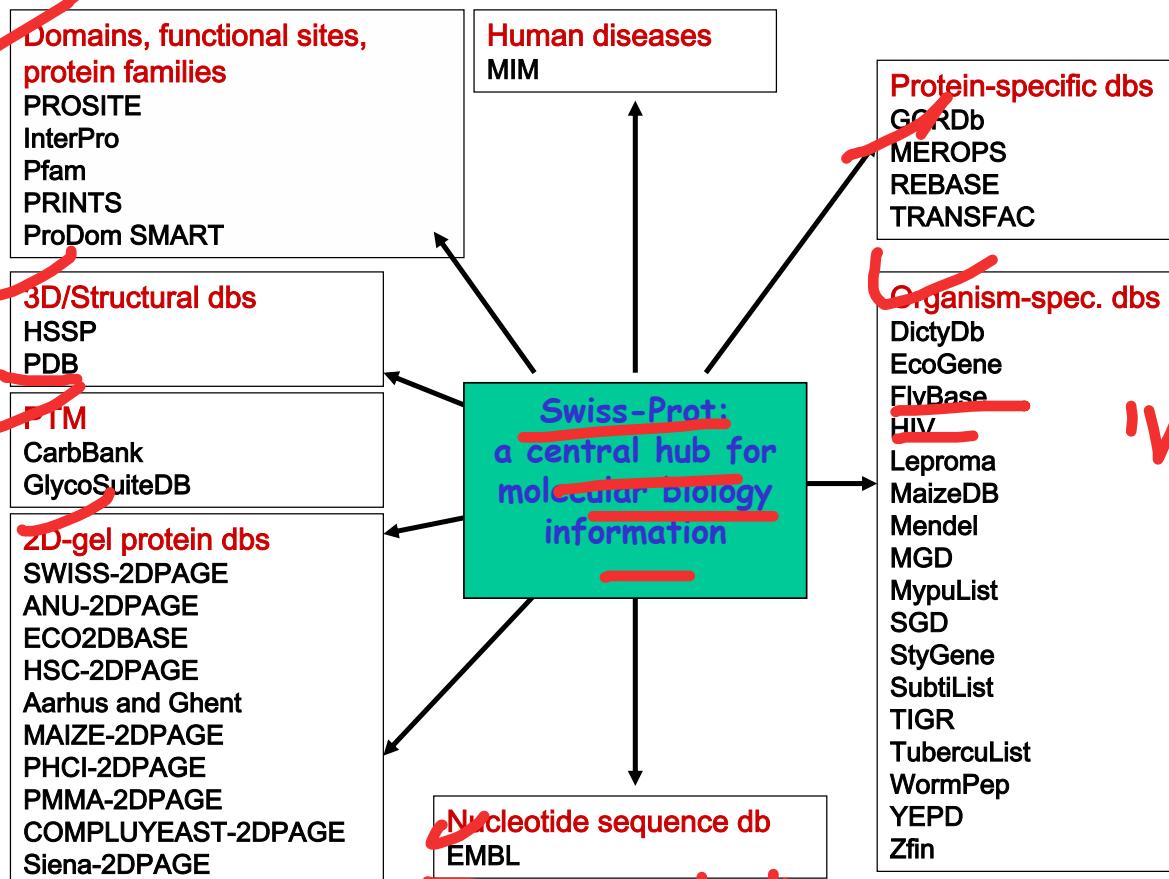
- **FUNCTION**: Potential regulator of neural basic helix-loop-helix genes (By similarity).
- **SUBCELLULAR LOCATION**: Nuclear (Probable)
- **SIMILARITY**: BELONGS TO THE E47/HOMEOBOX FAMILY.

Cross-references

- Explicit links to about 50 databases;
- Implicit X-references to 30 additional db added by the ExPASy servers on the WWW (such as GenBank, Ensembl, ...)
- => links to more than 80 databases from the ExPASy servers
- Currently 1.5×10^6 cross-references in Swiss-Prot

-> Connected with practically all the databases indexed under SRS.

Gasteiger et al., Curr. Issues Mol. Biol. (2001), 3(3): 47-55



IMPLICIT



Cross-references

ADN
(Index of low redundancy)

3D
genomic

Cross references	
X98172, CAA66854 1. -	[EMBL / GenBank / DDBJ] [CoDmgSequence]
X98173, CAA66854 1. -	[EMBL / GenBank / DDBJ] [CoDmgSequence]
X98174, CAA66855 1. -	[EMBL / GenBank / DDBJ] [CoDmgSequence]
X98175, CAA66856 1. -	[EMBL / GenBank / DDBJ] [CoDmgSequence]
X98176, CAA66857 1. -	[EMBL / GenBank / DDBJ] [CoDmgSequence]
X98177, CAA66858 1. ALT_Seq	[EMBL / GenBank / DDBJ] [CoDmgSequence]
X98178, CAA66859 1. ALT_Seq	[EMBL / GenBank / DDBJ] [CoDmgSequence]
U58143, AAC50602 1. -	[EMBL / GenBank / DDBJ] [CoDmgSequence]
U60520, AAC50645 1. -	[EMBL / GenBank / DDBJ] [CoDmgSequence]
AF102146, AAD24962 1. -	[EMBL / GenBank / DDBJ] [CoDmgSequence]
AF102139, AAD24962 1. JOINED	[EMBL / GenBank / DDBJ] [CoDmgSequence]
AF102140, AAD24962 1. JOINED	[EMBL / GenBank / DDBJ] [CoDmgSequence]
AF102141, AAD24962 1. JOINED	[EMBL / GenBank / DDBJ] [CoDmgSequence]
AF102142, AAD24962 1. JOINED	[EMBL / GenBank / DDBJ] [CoDmgSequence]
AF102143, AAD24962 1. JOINED	[EMBL / GenBank / DDBJ] [CoDmgSequence]
AF102144, AAD24962 1. JOINED	[EMBL / GenBank / DDBJ] [CoDmgSequence]
AF102145, AAD24962 1. JOINED	[EMBL / GenBank / DDBJ] [CoDmgSequence]
AF009620, AAB70913 1. -	[EMBL / GenBank / DDBJ] [CoDmgSequence]
AB038985, BAB32555 1. -	[EMBL / GenBank / DDBJ] [CoDmgSequence]
AB038982, BAB32555 1. JOINED	[EMBL / GenBank / DDBJ] [CoDmgSequence]
AB038983, BAB32555 1. JOINED	[EMBL / GenBank / DDBJ] [CoDmgSequence]
AB038984, BAB32555 1. JOINED	[EMBL / GenBank / DDBJ] [CoDmgSequence]
AF380342, AAK57437 1. -	[EMBL / GenBank / DDBJ] [CoDmgSequence]
AF422925, AAL87628 1. -	[EMBL / GenBank / DDBJ] [CoDmgSequence]
AF422926, AAL87629 1. -	[EMBL / GenBank / DDBJ] [CoDmgSequence]
AF422927, AAL87630 1. -	[EMBL / GenBank / DDBJ] [CoDmgSequence]
AF422928, AAL87631 1. -	[EMBL / GenBank / DDBJ] [CoDmgSequence]
AF422929, AAL87632 1. -	[EMBL / GenBank / DDBJ] [CoDmgSequence]
BC028223, AAH28223 1. -	[EMBL / GenBank / DDBJ] [CoDmgSequence]
1QDU; 10-JUL-00 [ExPASy / RCSB]	
C14_002 -	
HGNC_1509, CASP8.	
HGNC_1509, CASP8	
601763 [NCBI / EBI]	
607271 [NCBI / EBI]	
CASP8	
CASP8_Homo_sapiens	
CASP8_Homo_sapiens	
Q14790; Homo sapiens [Entry / Config view]	



EMBnet

MCB, feb 2005

FT lines ~~Feature table = Sequence description~~



- ~~Data derived from:~~
- ~~• Publications;~~
 - ~~• Databases;~~
 - ~~• Personal communication;~~
 - ~~• Prediction.~~

Features				
Key	From	To	Length	Description
SIGNAL	1	18	18	POTENTIAL.
CHAIN	19	302	284	100% LIGAND.
DOMAIN	19	256	238	EXTRACELLULAR (POTENTIAL).
TRANSMEM	257	277	21	POTENTIAL.
DOMAIN	278	302	25	CYTOPLASMIC (POTENTIAL).
DOMAIN	30	120	91	IG-LIKE V-TYPE DOMAIN.
DOMAIN	151	223	73	IG-LIKE C2-TYPE DOMAIN.
DISULFID	37	113		POTENTIAL.
DISULFID	158	216		POTENTIAL.
CARBOHYD	70	70		N-LINKED (GLCNAC...) (POTENTIAL).
CARBOHYD	137	137		N-LINKED (GLCNAC...) (POTENTIAL).
CARBOHYD	173	173		N-LINKED (GLCNAC...) (POTENTIAL).
CARBOHYD	186	186		N-LINKED (GLCNAC...) (POTENTIAL).
CARBOHYD	225	225		N-LINKED (GLCNAC...) (POTENTIAL).
VARSPLIC	300	302		GHV -> ESWNLNLLLS (IN ISOFORM 2).

General topology



Feature aligner



Feature table viewer

~~Sequence description:~~



~~Derived from:~~

- ~~• Publications;~~
- ~~• Databases;~~
- ~~• Personal communication;~~
- ~~• Prediction~~

Features					ICOL_HUMAN, O75144	
Key	From	To	Length	Description		
SIGNAL	1	18	18	POTENTIAL.		
CHAIN	19	302	284	ICOS LIGAND.		
DOMAIN	19	256	238	EXTRACELLULAR (POTENTIAL).		
TRANSMEM	257	277	21	POTENTIAL.		
DOMAIN	278	302	25	CYTOPLASMIC (POTENTIAL).		
DOMAIN	30	120	91	IG-LIKE V-TYPE DOMAIN.		
DOMAIN	151	223	73	IG-LIKE C2-TYPE DOMAIN.		
DISULFID	37	113	76	POTENTIAL.		
DISULFID	158	216	58	POTENTIAL.		
CARBOHYD	70	70	0	N-LINKED (GLCNAC...) (POTENTIAL).		
CARBOHYD	137	137	0	N-LINKED (GLCNAC...) (POTENTIAL).		
CARBOHYD	173	173	0	N-LINKED (GLCNAC...) (POTENTIAL).		
CARBOHYD	186	186	0	N-LINKED (GLCNAC...) (POTENTIAL).		
CARBOHYD	225	225	0	N-LINKED (GLCNAC...) (POTENTIAL).		
VARSPLIC	300	302	2	GHV -> ESWNLLLLS (IN ISOFORM 2).		

[Feature aligner](#)[Feature table viewer](#)

PTM

~~BRCA2_HUMAN, P51587~~

VARIANT	2466	2466	* A -> V (POLYMORPHISM; WAS ORIGINALLY THOUGHT TO BE LINKED TO OVARIAN CANCER).	VAR 008785
VARIANT	2480	2480	* L -> V.	VAR 008786
VARIANT	2490	2490	* I -> T.	VAR 008787
VARIANT	2502	2502	* R -> H (IN OVARIAN CANCER; COULD BE A POLYMORPHISM).	VAR 008788
VARIANT	2515	2515	* T -> I (IN BREAST CANCER; COULD BE A POLYMORPHISM).	VAR 008789
VARIANT	2787	2787	* R -> H (IN OVARIAN CANCER; SOMATIC MUTATION).	VAR 008790
VARIANT	2944	2944	* I -> F.	VAR 008791
VARIANT	2951	2951	* A -> T.	VAR 008792
VARIANT	2969	2969	* V -> M.	VAR 008793
VARIANT	3095	3095	* D -> E.	VAR 005108
VARIANT	3098	3098	* Y -> H (IN OVARIAN CANCER; COULD BE A POLYMORPHISM).	VAR 008794
VARIANT	3103	3103	* I -> M (IN MELANOMA).	VAR 005109
VARIANT	3118	3118	* M -> T (IN BREAST CANCER).	VAR 005110
VARIANT	3257	3257	* K -> R.	VAR 008795
VARIANT	3276	3276	* R -> S.	VAR 008796
VARIANT	3357	3357	* T -> R (IN BREAST CANCER).	VAR 005111
VARIANT	3412	3412	* I -> V (POLYMORPHISM; WAS ORIGINALLY THOUGHT TO BE ASSOCIATED WITH BREAST CANCER; dbSNP:1801426) [NCBI/Ensembl].	VAR 005112
CONFLICT	599	599	S -> F (IN REF. 2).	
CONFLICT	758	758	S -> N (IN REF. 1).	
CONFLICT	1761	1762	GY -> RI (IN REF. 1).	
CONFLICT	1767	1767	K -> N (IN REF. 1).	
CONFLICT	2536	2536	S -> P (IN REF. 3).	
CONFLICT	3216	3216	L -> LVS (IN REF. 3).	

~~Polymorphisms~~

~~Differences between the sequence shown and other submitted sequences~~

Sequence description:



Derived from:

- Publications;
- Databases;
- Personal communication;
- Prediction.

IUPAC

Features					ICOL_HUMAN, O75144
Key	From	To	Length	Description	
SIGNAL	1	18	18	POTENTIAL.	
CHAIN	19	302	284	ICOS LIGAND.	
DOMAIN	19	256	238	EXTRACELLULAR (POTENTIAL).	
TRANSMEM	257	277	21	POTENTIAL.	
DOMAIN	278	302	25	CYTOPLASMIC (POTENTIAL).	
DOMAIN	30	120	91	IG-LIKE V-TYPE DOMAIN.	
DOMAIN	151	223	73	IG-LIKE C2-TYPE DOMAIN.	
DISULFID	37	113	76	POTENTIAL.	
DISULFID	158	216	58	POTENTIAL.	
CARBOHYD	70	70	0	N-LINKED (GLCNAC...) (POTENTIAL).	
CARBOHYD	137	137	0	N-LINKED (GLCNAC...) (POTENTIAL).	
CARBOHYD	173	173	0	N-LINKED (GLCNAC...) (POTENTIAL).	
CARBOHYD	186	186	0	N-LINKED (GLCNAC...) (POTENTIAL).	
CARBOHYD	225	225	0	N-LINKED (GLCNAC...) (POTENTIAL).	
VARSPLIC	300	302	2	GHV -> ESWNLLLLS (IN ISOFORM 2).	Alternative splicing



[Feature aligner](#)



[Feature table viewer](#)

Alternative Splicing of Swiss-Prot entry: [O75144](#)

Name: 2
Isoform ID: O75144-2
Features which should be applied to build the isoform sequence: [V1](#)

All the alternatively spliced sequences are available for BLAST searches and proteomic tools at the ExPASy server

Sequence information

Length: 309 AA

10	20	30	40	50	60
MRLOSPGLLF	LLFSSSLRADT	QEKEVRAAVG	SDVVELSCACR	EKGSRFDLNDV	YVVWQTSKSK
70	80	90	100	110	120
TIVTYTHIPQN	ESLEENVDSRY	RNRAALMSPGC	MRLGDFSLBL	FNVTPQDEQK	FHCLVLVLSQSL
130	140	150	160	170	180
GFOEQLSVEV	TLIHVIAANFSV	PVVSAPHSPS	QDELTTFTCTG	INGYPRPRNVY	WINKTINSLL
190	200	210	220	230	240
DQALQNDTVF	LNNERGLYDVV	SVLRVARTPS	VNIQCCINNV	LLOQNLTVGS	QTONDICERD
250	260	270	280	290	300
KITENPVGSTG	EKNAATWSIL	AVLCLLVLVVA	VAIGHVCRDR	CLQHSYAGAW	AVSPETELTE
309					
					SWNLLLLS



BLAST [BLAST submission on ExPASy/SIB](#) or at [NCBI \(USA\)](#)



Sequence analysis tools: [ProtParam](#), [ProtScale](#), [Compute pI/Mw](#), [PeptideMass](#), [PeptideCutter](#), [Dotter](#) (Java)

[SeanProsite](#), [MotifScan](#)



[Direct Submission to SWISS-MODEL](#)

[ExPASy Home page](#)

[Site Map](#)

[Search ExPASy](#)

[Contact us](#)

[Proteomics tools](#)

[Swiss-Prot](#)

Hosted by SIB Switzerland | Mirror sites: [Australia](#) | [Bolivia](#) | [Canada](#) | [China](#) | [Korea](#) | [Taiwan](#) | [USA](#)

Swiss-Prot & TrEMBL introduce a new arithmetical concept !

$$170'000 + 1'600'000 \approx 1'200'000$$

Redundancy in TrEMBL

&

Redundancy between TrEMBL and Swiss-Prot

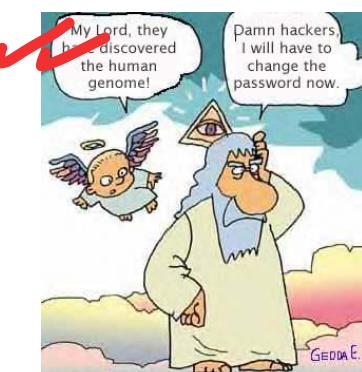
- In 2 years....more than 2'000'000 protein sequences
- But, in the future: redundancy is going to decrease:
« new » genome sequencing -> « new » proteins
(AB, sept 2002)

In the case of human proteins, the redundancy is still very high:

~~$$11'900 + 45'000 \approx \text{about } 22'000^*$$~~

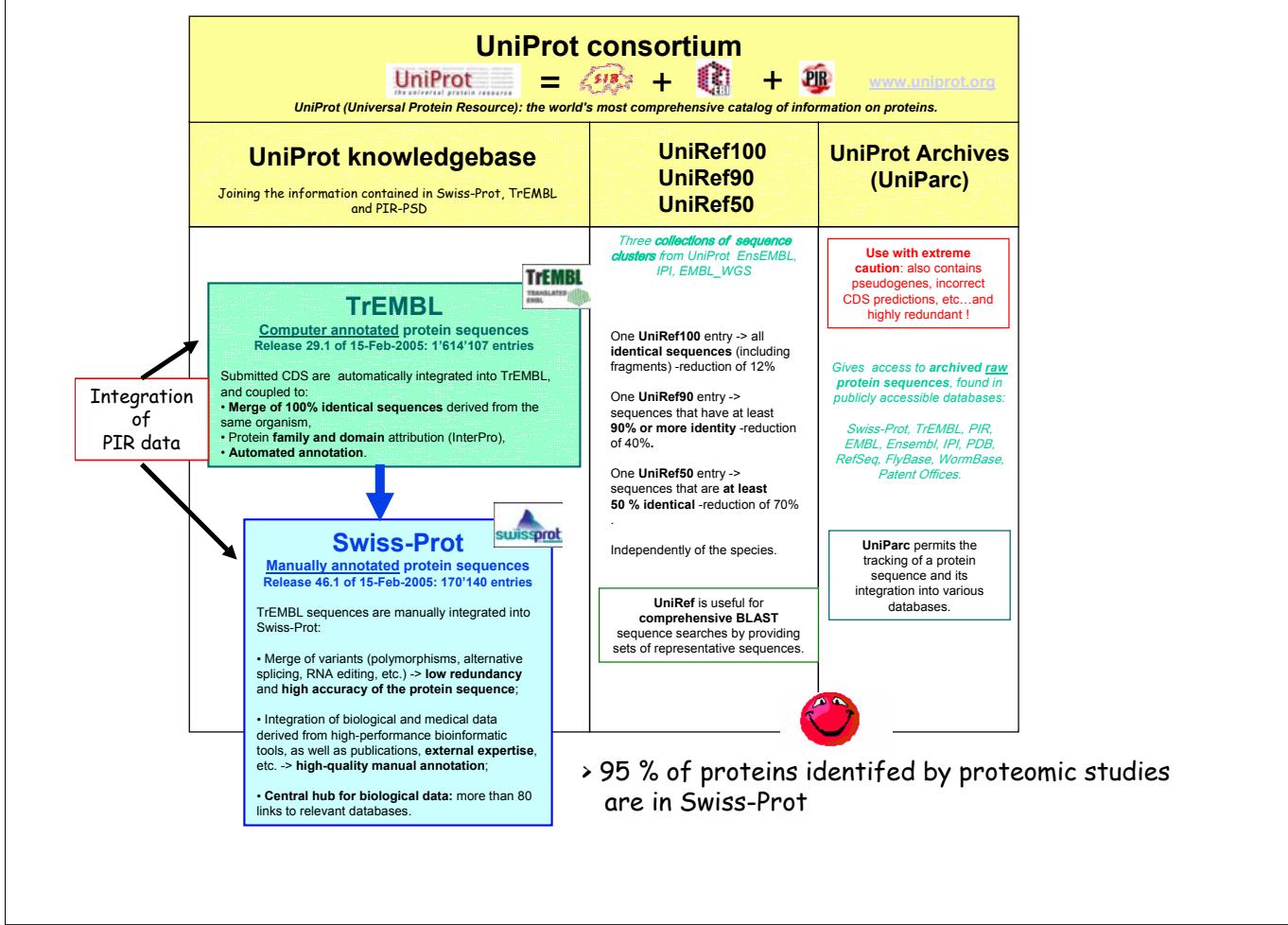
* human gene number estimation:
25'000-35'000

MS proteomics has verified more than 10% of human genes products, but has not identified significant numbers of unpredicted proteins (Southan C, Proteomics, 2004)



Are missing:

- Sequences not submitted to EMBL/GenBank/DDJB (and PIR)
- Not yet predicted or known genes (« no CDS provided by the submitters» or no DNA sequence)
- Confidential data (Patent application sequences)
- Immunoglobulins, T-cell receptors (-> UniParc)
- ...



Take home message

- Be aware of the differences between TrEMBL and Swiss-Prot.
 - Always cite the Accession number, not the ID.
 - We need your feedback!
- swiss-prot@expasy.org



Lin

10,000

Righting the wrongs

"Sequences are rarely deposited in a "mature" state; as with all scientific research, DNA and protein annotation is a continual process of learning, revision and corrections."

"Sequencing error rates: ~1 base in 10'000"

"Making people aware of errors is good and great; making people aware that they're responsible also for correcting errors is even greater"

C. Hardley, EMBO reports, 4(9), 2003.



EMBnet

MCB, feb 2005

Protein sequence databases

The NCBI-nr pathway
(Entrez protein)



EMBnet

MCB, feb 2005

~~Real life of a protein sequence ...~~

Data not submitted to public databases, delayed or cancelled...

cDNAs, ESTs, genomes, ...

↓ ↓ ↓ ↓ ↓ ↓
EMBL, GenBank, DDBJ

Coding Sequences
provided by submitter

Genpept

Coding Sequences
provided by submitter
and
<< de novo >> gene prediction

RefSeq
XP_NNNNN

Scientific publications derived sequences

PRF

~~NCBI nr~~: Swiss-Prot + GenPept + (PIR) + RefSeq + PDB + PRF

Lar
• ↗

(
gent contue
dans le PRF)

Protein sequences: « NR database » Entrez protein

<http://www.ncbi.nlm.nih.gov/entrez/query.fcgi?db=Protein>

The protein entries in the Entrez search and retrieval system have been compiled from a variety of sources, including SwissProt, PIR, PRF, PDB, and translations from annotated coding regions in GenBank and RefSeq.

Draft Human Genome
Explore [human genome resources](#) or browse the human genome sequence using the [Map Viewer](#).

~~NCBI - au~~

*departs I you
you back*

Scientific publications
derived sequences
« Journal scan »
(integrated into TrEMBL)

NCBI-nr: Swiss-Prot + GenPept + (PIR) + RefSeq + PDB + PRF

All PIR data have been
integrated into Swiss-Prot
and TrEMBL (UniProt)

derived
from GenBank/EMBL/DDBJ sequences
which have a CDS annotated on them
- equivalent to TrEMBL,
except that it is
redundant
with Swiss-Prot

3D structure database:
all the protein sequences
which have been crystallized
(Swiss-Prot/TrEMBL are
crosslinked to PDB)



EMBnet

MCB, feb 2005

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

- The RefSeq collection, which is tightly linked to LocusLink contains:
genomic DNA, transcript (RNA), and protein products

- RefSeq provides a non-redundant set of sequences, derived from GenBank,
the literature and gene prediction

- Release 3 includes over 800'000 proteins from 2218 organisms (including
1100 viruses and 150 bacteria).



EMBnet

MCB, feb 2005

□ 1: NP_057603. ameloblastin [Hom..[gi:7706067]

LOCUS NP_057603 447 aa linear PRI 05-OCT-2003
DEFINITION ameloblastin [Homo sapiens].
ACCESSION NP_057603
VERSION NP_057603.1 GI:7706067
DBSOURCE REFSEQ: accession NM_016519.2
KEYWORDS .
SOURCE Homo sapiens (human)
ORGANISM Homo sapiens
Eukaryota; Metazoa; Chordata; Craniata; Vertebrata; Euteleostomi;
Mammalia; Eutheria; Primates; Catarrhini; Hominidae; Homo.
REFERENCE 1 (residues 1 to 447)
AUTHORS Toyosawa,S., Fujiwara,T., Ooshima,T., Shintani,S., Sato,A.,
Ogawa,Y., Sobue,S. and Ijuhin,N.
TITLE Cloning and characterization of the human ameloblastin gene
JOURNAL Gene 256 (1-2), 1-11 (2000)
MEDLINE 20509985
PUBMED 11054529
REFERENCE 2 (residues 1 to 447)
AUTHORS MacDougall,M., Simmons,D., Gu,T.T., Forsman-Semb,K., Mardh,C.K.,
Mesbah,M., Forest,N., Krebsbach,P.H., Yamada,Y. and Berdal,A.
TITLE Cloning, characterization and immunolocalization of human
ameloblastin
JOURNAL Eur. J. Oral Sci. 108 (4), 303-310 (2000)
MEDLINE 20401566
PUBMED 10946765
REFERENCE 3 (residues 1 to 447)
AUTHORS Bartlett,J.D. and Simmer,J.P.
TITLE Proteinases in developing dental enamel
JOURNAL Crit. Rev. Oral Biol. Med. 10 (4), 425-441 (1999)
MEDLINE 20098047
PUBMED 10634581
REFERENCE 4 (residues 1 to 447)
AUTHORS MacDougall,M., DuPont,B.R., Simmons,D., Reus,B., Krebsbach,P.,
Karrman,C., Holmgren,G., Leach,R.J. and Forsman,K.
TITLE Ameloblastin gene (AMB) maps within the critical region for
autosomal dominant amelogenesis imperfecta at chromosome 4q21
JOURNAL Genomics 41 (1), 115-118 (1997)
MEDLINE 97271565
PUBMED 9126491
REFERENCE 5 (residues 1 to 447)
AUTHORS Lee,S.K., Krebsbach,P.H., Matsuki,Y., Nanci,A., Yamada,K.M. and
Yamada,Y.
TITLE Ameloblastin expression in rat incisors and human tooth germs
JOURNAL Int. J. Dev. Biol. 40 (6), 1141-1150 (1996)
MEDLINE 97184059
PUBMED 9033019
COMMENT REVIEWED [REFSEQ](#): This record has been curated by NCBI staff. The
reference sequence was derived from AF219994.2 and AF263464.1.

RefSeq/Protein

AC

KW

Taxonomy

References

GenBank source

GenBank source

As for the nucleic acid sequence, RefSeq chooses a protein Reference Sequence:
they do not annotate the sequence differences.

- If there is an alternative splicing event, there will be several entries for a same gene

Related Sequences

?

Nucleotide	Type	Protein	
AY009124	g	AAG27036	BL
AF209780	m	AAG35772	BL
→ AF219994	m	AAF37355	BL
→ AF263464	m	AAF73048	BL
AJ537437	m	CAD60938	BL
None	p	Q9NP70	BL

Annotation

Summary: Ameloblastin is thought to represent an unique ameloblast-specific gene product that may be important in enamel matrix formation and mineralization. The gene is located on chromosome 4 near other genes associated with mineralized tissues: osteopontin, bone sialoprotein, and bone morphogenetic protein 3. Based on its cytogenetic location, this gene is a candidate gene for one form of the disorder, dentinogenesis imperfecta, and/or the disorder, autosomal dominant amylogenesis imperfecta.

FEATURES Location/Qualifiers
source 1..447
/organism="Homo sapiens"
/db_xref="taxon:9606"
/chromosome="4"
/map="4q21"

Protein 1..447
/product="ameloblastin"

variation 255
/allele="V"
/allele="A"
/db_xref="dbSNP:[7439186](#)"
1..447

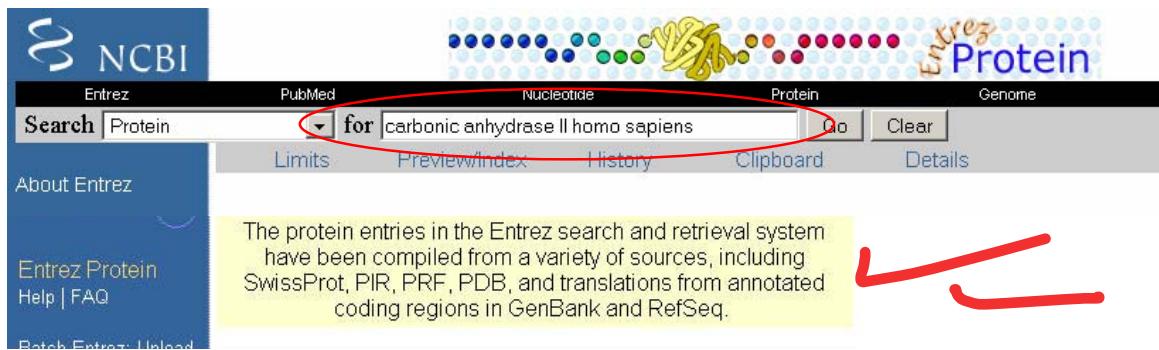
CDS /gene="AMBN"
/coded_by="NM_016519.2:76..1419"
/notes="go_component: extracellular matrix [goid 0005578]
[evidence NAS];
go_function: structural constituent of tooth enamel [goid 0030345] [evidence NAS];
go_process: bone mineralization [goid 0030282] [evidence NAS]"
/db_xref="GeneID:258"
/db_xref="LocusID:258"
/db_xref="MIM:[601259](#)"

Cross references

OPEN IN
1 msaskiplfk mkkdililcl lemsfafvpif pqqsgtpgma slsletmrql gslqrnlts
61 qysryigfks fnslwmgll phhssipwmr prehetqqye yslpvhpppl pscpslkpqg
121 pgikpflqsa aattnqatal kealqppihl ghlpigegel plvqqgvaps dkppkpelpg
181 vdfadpqgps lpgmdfpdpq gpsslgldfa dpqgstifqi arlishgpmp qnkqsplypg
241 mlyvpfgand lnaparlgim sseevaggre dpmaygamfp qfgggmrpgfe qmphpnpangg
301 dftlefdspv aatkpgenee ggaqgspmpe anpdnlenpa fltelepaph agllalpkdd
361 ipglprspsg kmkgklpsvtp aaadplmtpe ladvyrtyna dmttsvdfqe eatmdttmap
421 ns1qtsmpgn kaqepemmhhd awhfqep
//

Query at Entrez protein

<http://www.ncbi.nlm.nih.gov/entrez/query.fcgi?db=Protein>



The protein entries in the Entrez search and retrieval system have been compiled from a variety of sources, including SwissProt, PIR, PRF, PDB, and translations from annotated coding regions in GenBank and RefSeq.

Typical result of
a query at
« Entrez protein »

□ 3: NP_000058
carbonic anhydrase II; carbonic dehydratase II; carbonic dehydratase; carbonic anhydrase B [Homo sapiens]
gi|4557395|ref|NP_000058.1|[4557395]

□ 4: NP_000333
solute carrier family 4, anion exchanger, member 1 (erythrocyte membrane protein band 3, Diego blood group);
[Homo sapiens]
gi|4507021|ref|NP_000333.1|[4507021]

□ 5: NP_002334
lactotransferrin; lactoferrin [Homo sapiens]
gi|4505043|ref|NP_002334.1|[4505043]

□ 6: NP_003606
solute carrier family 4, sodium bicarbonate cotransporter, member 7; sodium bicarbonate cotransporter 2; solute
member 6 [Homo sapiens]
gi|19923176|ref|NP_003606.2|[19923176]

□ 7: P00918
Carbonic anhydrase II (Carbonate dehydratase II) (CA-II) (Carbonic anhydrase C)
gi|115456|sp|P00918|CAH2_HUMAN|[115456]

□ 8: AAH11949
Carbonic anhydrase II [Homo sapiens]
gi|15080386|gb|AAH11949.1|[15080386]

□ 9: 1QQ5A
Chain A, Carbonic Anhydrase II In Complex With Nanomolar Inhibitor
gi|47168469|pdb|1QQ5|A[47168469]



EMBnet

MCB, feb 2005

RefSeq

Swiss-Prot

Genpept

(gb/embl/ddbj)

DDB

The AC number jungle

Type of record	Sample Accession Format
GenBank/EMBL/DDBJ	One letter followed by five digits: e.g. U12345 Two letters followed by 6 digits: e.g. AF123456
Swiss-Prot/TrEMBL	One letter and five digits/letters: e.g. P12345
RefSeq nucleotide	Two letters, underscore bar and six digit: e.g. mRNA NM_000492 e.g. genomic NT_000907
RefSeq protein	e.g. NP_00483
RefSeq prediction	e.g. XM_000483 e.g. XP_000467
PDB (protein structure)	One digit followed by three letters: e.g. 1TUP

Categories of databases for Life Sciences

- Sequences (DNA, protein)
- Genomics
- Mutation/polymorphism
- Protein domain/family (----> tools)
- Proteomics (2D gel, Mass Spectrometry)
- 3D structure
- Metabolism
- Bibliography
- 'Others' (Microarrays, Protein protein interaction...)



EMBnet

MCB, feb 2005

~~Databases 2: 'genomics'~~

- Contain informations on gene chromosomal location (mapping) and nomenclature, and provide links to sequence databases; has usually no sequence; ↗
- Exist for most organisms important in life science research; usually species specific.
- Examples: MIM, GDB (human), MGD (mouse), FlyBase (Drosophila), SGD (yeast), MaizeDB (maize), SubtiList (*B.subtilis*), etc.; ↗
- Generally relational db (Oracle, SyBase or AceDb). ↗



EMBnet

MCB, feb 2005

Mendelian

Genetic units

MTM / OMIM

OMIM™: Online Mendelian Inheritance in Man

- catalog of human genes and genetic disorders
- contains a summary of literature and reference information. It also contains links to publications and sequence information.



EMBnet

MCB, feb 2005

*13371 ERYTHROPOEITIN (EPO)

alternative symbols

TABLE OF CONTENTS

- DESCRIPTION
- CLONING
- MAPPING
- ANIMAL MODEL
- REFERENCES
- SIL ALIOTS
- CORRECTIONS
- CREATION DATE
- EHT HISTORY

Database Links

[MEDLINE](#) [Protein](#) [DNA](#) [Genome](#) [LocusLink](#) [Gene Map](#) [GBD](#) [NCB](#) [Nomenclature](#)

Gene Map Locus 7q21

Note: pressing the symbol will find the citations in MEDLINE whose text most closely matches the text of the preceding OMIM paragraph, using the Entrez MEDLINE neighboring function.

TEXT

DESCRIPTION

Human erythropoietin is an acidic glycoprotein hormone with a molecular mass of 34 kD. As the prime regulator of red cell production, its major functions are to promote erythroid differentiation and to initiate hemoglobin synthesis.

CLONING

Lee-Huang (1984) cloned human erythropoietin cDNA in E. coli. McDonald et al. (1986) and Shuemaker and Matsick (1986) cloned the mouse gene and the latter workers showed that coding DNA and amino acid sequence are about 80% conserved between man and mouse. This is a much higher order of conservation than for various interleukins, interleukin-2, and GM-CSF. Shuemaker and Shiu (1986) described a human renal carcinoma cell line that continuously produces erythropoietin. Etzkorn et al. (1992) demonstrated the effectiveness of recombinant human erythropoietin in treating the anemia of end-stage renal disease.

Romanowski and Sykora (1994) reviewed the molecular structure of human erythropoietin in historical perspective. The EPO gene has 5 exons that code for a 193-amino acid prohormone. A 27-amino acid leader sequence is cleaved off the amino terminus of the prohormone, yielding the functional 166-amino acid protein. However, recombinant human EPO expressed in Chinese hamster ovary cells contains only 165 amino acids, having lost arg166. The mechanism for this was undefined, and whether EPO circulates in the plasma as a truncated arg166 was not known. Both the nucleotide and amino acid sequences of EPO are highly conserved among mammals.

In the central nervous system, neurons express EPO receptor (EPOR) and astrocytes produce EPO. EPO has been shown to protect primary cultured neurons from TNF α -mediated glutamate toxicity. Sakurada et al. (1990) reported in vivo evidence that EPO protects neurons against ischemia-induced cell death. They presented findings suggesting that EPO may exert neuroprotective effect by reducing the tonic osteo-mediated formation of free radicals or antagonizing their toxicity. Semenza et al. (2001) presented data suggesting that inhibition of neuronal apoptosis under short latency protective effects of EPO after cerebral ischemia and other brain injuries. They suggested that evaluation of EPO, a compound established as clinically safe, as neuroprotective therapy in acute brain injury is required.

Noval erythropoietin-stimulating protein (NEST) stimulates erythropoiesis in the same manner as human recombinant EPO. NEST is distinct from EPO in that it has additional basic which has been shown to confer an increased terminal half-life in animal models, patients with chronic renal failure, and cancer patients receiving multiple cycles of chemotherapy (Macdonald et al., 1999). In studies of 89 patients with nonmelanoma malignancies, Smith et al. (2001) found that NEST was well tolerated, with response rates ranging from 61 to 83%, depending on dosage.

MAPRED

Law et al. (1986) assigned EPO to chromosome 7 by Southern blot analysis of DNA from human/Chinese hamster cell hybrids with a cDNA clone for the entire coding region of the gene. Further localization to 7q11-q22 was achieved by *in situ* hybridization. They found a RFLP with a frequency of about 20% in a Chinese population. By hybridization analysis (dot-blot) of DNA from human chromosomes isolated by high resolution dual laser sorting, Powell et al. (1989) also located EPO on chromosome 7. They found a RFLP with a frequency of about 20% in a Chinese population. The chromosomal band of 7q, closely linked to COL1A2 (120169) and to DNA markers linked to CF (G19700). Because of the close linkage of EPO to COL1A2 and markers linked to CF, it is probably justified to narrow the assignment of EPO to 7q21-q22.

By *in situ* hybridization and by genomic analysis using RFLPs in interspecific mouse backcross DNAs, Lacombe et al. (1988) demonstrated that EPO is located on chromosome 5 in the mouse.

In addition to its role as a kidney cyto-kine regulating hemopoiesis, EPO is also produced in the brain after oxidative or inflammatory stress. The transcription factor HIF-1 (602340) upregulates EPO following hypoxic stimuli. Dangond et al. (2001) demonstrated that preconditioning with EPO protects neurons in models of ischemic and degenerative damage due to excitotoxicity and consequent apoptosis, including nitric oxide. Activation of neuronal EPO receptors (133171) prevents apoptosis induced by NM23 or nitric oxide by triggering cross talk between signaling pathways JAK2 (47796) and NFkB (see 164011). Pacewicz and Letton (2001) demonstrated that EPO receptor-mediated activation of JAK2 leads to phosphorylation of the inhibitor of NFkB (rel-B kappa-B-alpha) (164009), subsequent nuclear translocation of the transcription factor NFkB, and subsequent upregulation of the expression of a neurotrophic factor, ciliary neurotrophic factor, with a dominant interfering form of JAK2 or an I-kappa-B-alpha superrepressor blocks EPO-induced prevention of neuronal apoptosis. Thus, neuronal EPO receptors activate a neuroprotective pathway that is distinct from previously well characterized JAK and NFkB functions. Moreover, this EPO effect may underlie neuroprotection mediated by hypoxic-ischemic preconditioning.

ANIMAL MODEL

Synthesis of erythropoietin in the kidney and liver in response to hypoxia depends on both protein synthesis and heme synthesis. Gotoh et al. (1989) proposed a model in which a ligand-dependent conformational change in a heme protein accounts for the hypoxia signal by which hypoxia acts as cobalt and nickel stimulates the production of erythropoietin. Semenza et al. (1989) generated transgenic mice containing the human erythropoietin gene and found increased erythropoietin mRNA expression not only in liver and kidney but in all other transgenic tissues analyzed. The mice were polythemic, with increased erythropoietin precursors in hemato-poietic tissues and increased erythrocyte indices in peripheral blood. From further studies in these transgenic mice, Semenza et al. (1989) concluded that different DNA sequences flanking the EPO gene control liver versus kidney expression of the gene and that some of these sequences are located 3-prime to the gene.

Nathkaik et al. (1995) examined whether the secretion of erythropoietin from genetically modified cells could represent an alternative to repeated injections of the recombinant hormone for treating chronic anemias responsive to EPO. Primary mouse skin fibroblasts were transduced with a retroviral vector in which the murine cDNA was expressed under the control of the murine phosphoglycerate kinase promoter. "Neo-egress" containing the genetically modified fibroblasts embedded into collagen lattices were implanted into the peritoneal cavity of mice. Increased hematopoiesis and elevated serum EPO concentration were observed in recipient animals over a 10-month observation period. The approach was considered applicable to the treatment of human anemias.

Chen et al. (1991) investigated the rate of expression and biologic effects of transplanting autologous vascular smooth muscle cells transduced with a retroviral vector encoding rat erythropoietin cDNA. Vector-derived Epo secretion caused increases in reticulocytes followed by clinically significant increases in hematocrit and hemoglobin for up to 11 weeks. There were no significant differences between control and treated animals in the number of white blood cells and platelets. Kidney and to a lesser extent liver are specific organs that synthesize Epo in response to tissue oxygenation. In the treated animals, endogenous Epo mRNA was largely downregulated in liver and absent from liver. These results indicate to the authors that autologous smooth muscle cells can be genetically modified to provide treatments of anemias due to Epo deficiency and suggest that this cell type may be targeted in the treatment of other diseases requiring systemic therapeutic protein delivery.

Similar experiments were performed by Kessler et al. (1996), who demonstrated that, following a single intramuscular administration of a recombinant adenovirus-associated virus (AAV) vector containing the beta-galactosidase gene into adult mice, protein expression was detected in myoblasts for at least 22 weeks. Furthermore, a single intramuscular administration of an AAV vector containing a gene for human erythropoietin gene resulted in dose-dependent secretion of erythropoietin and corresponding increases in red blood cell production that persisted for up to 40 weeks. Primary human myocytes transduced in vitro with the AAV-Epo vector also showed dose-dependent production of Epo.

SEE ALSO

Taylor et al. (1985) · Ian et al. (1985) · Semenza et al. (1989)

REFERENCES

Gene Lynx

A portal to the human genome

Release 1.0 beta
26 Dec 2001
32657 records

Text search BLAST search GeneLynx guide GeneLynx info

View GeneLynx record
Enter a GeneLynx ID:
ID: Go

GeneLynx Home Text Search BLAST search Linking to GeneLynx Resource submission GeneLynx guide GeneLynx info

GeneLynx is a portal to a collection of hyperlinks for each human gene. It is implemented as an easily extensible relational database with a straightforward user interface.

You can access the information about a particular human gene by providing any reasonable identifier - just type a keyword, ANY accession number or ID below, or submit a related protein or nucleotide sequence on the BLAST search page. You can also perform a more refined keyword search on the Text search page.

Parts of GeneLynx were out of function January 11-13, 2002 due to server misconfiguration.
We apologize for the inconvenience.

Quick Search
Enter one or more terms separated by spaces.
Combine terms with: AND OR
 Exclude low-scoring hits

Send comments and questions to Boris Lenhard

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

Gene Lynx

GeneLynx #5230

Release 1.1
07 May 2002
32226 records

View GeneLynx record
Enter a GeneLynx ID:
ID: Go

GeneLynx Home Text Search BLAST search Linking to GeneLynx Resource submission GeneLynx guide GeneLynx info

Summary pages

LocusLink	2056
GeneCards	EPO
Unigene	Hs.2303
Swiss_Prot	EPO_HUMAN
KEGG gene	2056
EGAD	3760
euGenes	HUGN0002056
MIPS	771
HumanPSD	EPO

Genomic resources

Genomic sequences	NC_000010 AF053356 NC_000010 X02158	NC_000010 G20209	NC_000010 M11319
GDB	119110		
GenAtlas	EPO		
Ensembl gene	ENSG00000087083	ENSG00000130427	

Collections of hyperlinks for each human gene

Categories of databases for Life Sciences

- ✓ Sequences (DNA, protein)
- ✗ Genomics
- ✗ Mutation/polymorphism
 - Protein domain/family (----> tools)
 - Proteomics (2D gel, Mass Spectrometry)
 - 3D structure
 - Metabolism
 - Bibliography
 - 'Others' (Microarrays, Protein protein interaction...)



EMBnet

MCB, feb 2005

✗ Mutation/polymorphism: definitions

- ✗ SNPs: single nucleotide polymorphisms; occur approximately once every 100 to 300 bases
(distinction between sequencing error and polymorphism !)
- ✗ c-SNPs: coding single nucleotide polymorphisms
(Single Nucleotide Polymorphisms within cDNA sequences)
- ✗ SAPs: single amino-acid polymorphisms
- ✗ Missense mutation: -> SAP
- ✗ Nonsense mutation: -> STOP
- Insertion/deletion of nucleotides -> frameshift...



EMBnet

MCB, feb 2005

Databases 3: mutation/polymorphism

- Contain informations on sequence variations linked or not to genetic diseases;
- Mainly human but: OMIA - Online Mendelian Inheritance in Animals
- General db:
 - ~~OMIM~~
 - HGMD - Human Gene Mutation db
 - SVD - Sequence variation db
 - HGBASE - Human Genic Bi-Allelic Sequences db
 - dbSNP - Human single nucleotide polymorphism (SNP) db
- Disease-specific db: most of these databases are either linked to a single gene or to a single disease;
 - p53 mutation db
 - ADB - Albinism db (Mutations in human genes causing albinism)
 - Asthma and Allergy gene db
 -

For human (Amos'link)

- [HGMD](#) - Human Gene Mutation db
- [SVD](#) - EBI Sequence variation db
- [HGBASE](#) - Human Genic Bi-Allelic Sequences db
- [The SNP consortium](#)
- • [dbSNP](#) - Human single nucleotide polymorphism (SNP) db
 - [ALFRED](#) - Allele Frequency Db
 - [SeattleSNPs](#) - UW-FHCRC Variation Discovery Resource
 - [PicSNP](#) - Catalog of non-synonymous SNP
- [List of mutation databases from OMIM](#)
- [List of mutation databases from IMT \(Finland\)](#)

Organism	Number of Submissions (ss#'s)	Number of RefSNP Clusters (rs#'s) (# validated)	Number of (ss#'s) with frequency	Number of (ss#'s) with genotype
Homo sapiens / geneReport	21,581,724	10,054,521 (5,054,875)	488,391	2,727,888
Gallus gallus / geneReport	3,841,959	3,291,672 (3,281,479)		3,624,831
Oryza sativa	3,931,108	3,899,916 (22,057)		
Anopheles gambiae	1,368,805	1,136,268 (0)		
Mus musculus / geneReport	600,576	581,577 (535,237)		545,500
Canis familiaris	975,417	957,828 (12,315)	17	
Rattus norvegicus	45,360	43,229 (702)		669
Saccharum hybrid cultivar	42,853	42,853 (0)		
Danio rerio	2,031	2,025 (1,903)		2,031
Sus scrofa	5,450	1,521 (24)		
Bos indicus x Bos taurus	2,425	2,425 (0)		
Bos taurus	2,058	2,058 (45)	54	
Caenorhabditis elegans	1,065	1,065 (0)		



EMBnet

MCB, feb 2005

Mutation/polymorphism

- No single source for all SNPs (~100 SNPs db) !
- Generally modest size; lack of coordination and format standards in these databases making it difficult to access the data.
- ! Numbering of the mutated amino acid depends on the db (aa no 1 is not necessarily the initiator Met !)
- There are initiatives to unify these databases (politic/founding problems)
Mutation Database Initiative (4th July 1996).

→ SVD - Sequence Variation Database project at EBI (HMutDB)
<http://www.ebi.ac.uk/mutations/central/>

→ HUGO Mutation Database Initiative (MDI).
Human Genome Variation Society
<http://www.genomic.unimelb.edu.au/mdi/dblist/dblist.html>



EMBnet

MCB, feb 2005

Categories of databases for Life Sciences

1)

- Sequences (DNA, protein)
- Genomics
- Mutation/polymorphism
- Protein domain/family (----> tools)
- Proteomics (2D gel, Mass Spectrometry)
- 3D structure
- Metabolism
- Bibliography
- 'Others' (Microarrays, Protein protein interaction...)



EMBnet

MCB, feb 2005

Protein domain/family: some definitions

- Most proteins have « modular » structures
- Estimation: ~ 3 domains / protein



EMBnet

MCB, feb 2005

Some statistics

Proteome Analysis @EBI		
15 most common domains for <i>H. sapiens</i> [help]		
InterPro	Matches per genome (Proteins matched)	Name
IPR007110	6473(983)	Immunoglobulin-like
IPR007087	29392(949)	Zn-finger, C2H2 type
IPR001841	1182(391)	Zn-finger, RING
IPR008938	879(391)	ARM repeat fold
IPR001849	1078(354)	Pleckstrin-like
IPR009057	406(342)	Homeodomain-like
IPR006209	3019(322)	EGF-like domain
IPR000504	1540(296)	RNA-binding region RNP-1 (RNA recognition motif)
IPR008941	690(291)	TPR-like
IPR001452	2206(290)	SH3
IPR008957	934(258)	Fibronectin, type III-like fold
IPR009058	216(209)	Winged helix DNA-binding
IPR001478	1335(208)	PDZ/DHR/GLGF domain
IPR008973	318(200)	C2 calcium/lipid-binding domain, CaLB
IPR000210	495(196)	BTB/POZ domain

<http://www.ebi.ac.uk/proteome/HUMAN/interpro/top15d.html>



EMBnet

MCB, feb 2005

Protein domain/family: some definitions

- Domains (conserved sequences or structures) are identified by multiple sequence alignments

Sequence ID	start	end	weight	Sequence
3 EPO_HUMAN	2.41			APPRLICDSRVLERYLLEAKEAENVTGCGSEHCSLNENITVPTKVNFYAWKRM
2 EPO_RAT	2.61			APPRLICDSRVLERYLLEAKEAENVTGCAEGPRLSENITVPTKVNFYAWKRM
3 EPO_FELCA	2.99			APPRLICDSRVLERYLLEAKEAENATMGCAGCSFSENITVPTKVNFYAWKRM
8 Consensus	8.01			APPRLICDSRVLERYLLEAKEAENVTGCAEGCSLNENITVPTKVNFYAWKRM
1 PROSITE				

- Domains can be defined by different methods:

- Pattern (regular expression); used for very conserved domains
- Profiles (weighted matrices): two-dimensional tables of position specific match-, gap-, and insertion-scores, derived from aligned sequence families; used for less conserved domains
- Hidden Markov Model (HMM); probabilistic models; an other method to generate profiles.

Pattern-Profile

HPT1_HUMAN : NLTTGATLINEQNLTTAKNA
 ACRO_RABIT : YHACGGVLLNAHVVLTAAHCS
 KLKE_HUMAN : RFLCGGALLSGQNVITATHCL
 MCT3_SHEEP : SYICGGFLVREDPVLTAAHCP
 TRB2_HUMAN : MHFCGGSLIHPQNVLTAAHCE
 PRTC_HUMAN : KLACCGAVLIHPSNVLTAAHCA
 EL2_MOUSE : RHNCGGSLVANNNVLTAAHCH
 HPT_CANFA : NLTSQAPLINEQNLMTTAKNV
 VSP3_TRIFL : GALCGGTLLINQEVVLTASHCL
 TMS3_HUMAN : YHLCGGSVITPLNLTAAHCA
 TRY2_RAT : YHFCGGSLINDQNVVSAAHCF
 MCT2_RAT : RVICGGELISRQFVLTAAHCF
 HPT_MUSSA : GLTTCAHLISDQNLTTAKNN
 TRY4_LUCCU : SHSCGGSVYNSRIIVTAAHCY
 PLMN_MACMU : MHFCGGTLISPEVVLTAGHON

- Pattern [LIVM]-[ST]-A-[STAG]-H-C

→ Yes or no

- Profile:

```
ID TRYPSIN_DOM; MATRIX.
AC PS50240;
DT DEC-2001 (CREATED); DEC-2001 (DATA UPDATE); DEC-2001 (INFO UPDATE).
DE Serine proteases, trypsin domain profile.
MA /GENERAL_SPEC: ALPHABET="ABCDEFGHIKLMNPQRSTUVWXYZ"; LENGTH=234;
MA /DISJOINT: DEFINITION=PROTECT; N1=6; N2=229;
MA /NORMALIZATION: MODE=-1; FUNCTION=LINEAR; R1=0.0169; R2=0.00836256; TEXT='-'LogE';
MA /CUT_OFF: LEVEL=0; SCORE=1134; N_SCORE=9.5; MODE=-1; TEXT='!';
MA /CUT_OFF: LEVEL=-1; SCORE=775; N_SCORE=6.5; MODE=-1; TEXT='?';
MA /DEFAULT: M0=-9; D=-20; I=-20; BI=-60; EI=-60; MI=-105; MD=-105; DM=-105;
MA /I: BI=0; BI=-105; BD=-105;
MA      A B D E F G H I K L M N P Q R S T V W Y
MA /M: SY='I'; M=-8,-29,-34,-26, 3,-34,24, 34,-26,19, 15,-24,-21,-21,-24,-19,-8,-25,-19, 3;
MA /M: SY='N'; M= 0, 14, 10, 1,-22,-1, 6,-23,-4,-26,-17,20,-14,-1,-6, 13, 2,-20,-34,-15;
MA /M: SY='E'; M=-4, 4, 7, 14,-26,-13,-7,-23, 3,-22,-16, 2, 7, 3,-3, 2,-2,-21,-30,-18;
MA /M: SY='R'; M=-12, 5, 5, 7,-23,-17, 3,-24, 8,-20,-12, 7,-16, 10, 12,-2,-6,-21,-27,-9;
MA /M: SY='W'; M=-16,-33,-35,-27, 13,-22,-24,-11,-18,-13,-13,-31,-27,-20,-18,-30,-21,-18, 97, 25;
MA /M: SY='V'; M= 1,-29,-31,-28,-1,-30,-29, 31,-22, 13, 11,-27,-27,-26,-22,-12,-2, 41,-27,-8;
MA /M: SY='L'; M=-8,-29,-31,-22, 9,-30,-21, 23,-27, 37, 20,-28,-28,-21,-20,-25,-8, 17,-20,-1;
MA /M: SY='T'; M= 2,-1,-9,-9,-11,-17,-19,-10,-10,-13,-11, 1,-11,-9,-10, 23, 43, 0,-32,-12;
MA /M: SY='A'; M= 45,-9,-19,-10,-20,-2,-15,-11,-10,-11,-10,-9,-11,-9,-19, 10, 1,-1,-21,-18;
MA /M: SY='P'; M= 40,-9,-17,-8,-21, 5,-18,-14,-9,-13,-12,-8,-11,-9,-16, 9,-2,-5,-21,-21;
MA /M: SY='H'; M= 18, 0, 0, 1,-21,-19, 89,-29,-8,-21,-1, 9,-19, 11, 0,-7,-17,-29,-30, 16;
MA /M: SY='C'; M=-9,-18,-28,-29,-20,-29,-29,-29,-20,-19,-18,-39,-29,-29,-9,-9,-49,-29;
MA /I: EI=0; IE=-105; DE=-105;
//
```

→ score/threshold

Protein domain/family databases

- Contains biologically significant « pattern / profiles / HMM » formulated in such a way that, with appropriate computational tools, it can rapidly and reliably determine to which known family of proteins (if any) a new sequence belongs to
- Used as a tool to identify the function of uncharacterized proteins translated from genomic or cDNA sequences (« functional diagnostic »)
- Either manually curated (i.e. PROSITE, PfamA, PRINTS, SMART, TIGRFAM etc.) or automatically generated (i.e. PfamB, ProDom, DOMO)

Protein domain/family db

		Interpro
PROSITE	Patterns / Profiles	
ProDom	Aligned motifs (PSI-BLAST) (Pfam B)	
PRINTS	Aligned motifs	
Pfam	HMM (Hidden Markov Models)	
SMART	HMM	
TIGRfam	HMM	

DOMO Aligned motifs
BLOCKS Aligned motifs (PSI-BLAST)
CDD Pfam and SMART

→ A Conserved Domain Database and Search Service



EMBnet

MCB, feb 2005

Prosite <http://www.expasy.org/prosite/>

- Created in 1988 (SIB)
- Contains functional domains fully annotated, based on two methods: patterns and profiles
- Entries are deposited in PROSITE in two distinct files:
 - Pattern/profiles with the list of all matches in SWISS-PROT
 - Documentation

15-Aug-2004: contains 1277 documentation entries that describe 1736 different patterns, rules and profiles/matrices.



EMBnet

MCB, feb 2005

NiceSite View of PROSITE: [PS00817](#)

General information about the entry	
Entry name	EPO_TPO
Accession number	PS00817
Entry type	PATTERN
Date	OCT-1993 (CREATED); NOV-1995 (DATA UPDATE); JUL-1998 (INFO UPDATE).
PROSITE documentation	PDOC00644
Name and characterization of the entry	
Description	Erythropoietin / thrombopoietin signature.
Pattern	P-x(O)-C-D-x-R-[LIVM](2)-x-[KR]-x(14)-C.
Numerical results	
<ul style="list-style-type: none"> SWISS-PROT release number: 38, total number of sequence entries in that release: 80000. Total number of hits in SWISS-PROT: 14 hits in 14 different sequences Number of hits on proteins that are known to belong to the set under consideration: 14 hits in 14 different sequences Number of hits on proteins that could potentially belong to the set under consideration: 0 hits in 0 different sequences Number of false hits (on unrelated proteins): 0 hits in 0 different sequences Number of known missed hits: 0 Number of partial sequences which belong to the set under consideration, but which are not hit by the pattern or profile because they are partial (fragment) sequences: 1 Precision (true hits / (true hits + false positives)): 100.00 % Recall (true hits / (true hits + false negatives)): 100.00 % 	
Comments	
<ul style="list-style-type: none"> Taxonomic range: Eukaryotes Maximum known number of repetitions of the pattern in a single protein: 1 'Interesting' site in the pattern: 3,disulfide 'Interesting' site in the pattern: 11,disulfide 	
Cross-references	
SWISS-PROT	True positive hits:
	<p>EPO_BOVIN (P48617), EPO_CANFA (P33707), EPO_FEELCA (P33708), EPO_HUMAN (P01588), EPO_MACFA (P07865), EPO_MACMU (Q28513), EPO_MOUSE (P07321), EPO_PIG (P49157), EPO_RAT (P29576), EPO_SHEEP (P33709), TPO_CANFA (P42705), TPO_HUMAN (P40225), TPO_MOUSE (P40226), TPO_RAT (P49745)</p> <p>Potential' hits (sequences that belong to the set under consideration, but which were not picked up because the region(s) that are used as a 'fingerprint' (pattern or profile) is not yet available in the data bank (partial sequence)):</p> <p>TPO_PIG (P42706)</p>

NiceSite View of PROSITE: PS50097

ProDom

<http://protein.toulouse.inra.fr/prodom/current/html/home.php>

- ProDom is a comprehensive set of protein domain families automatically generated from the SWISS-PROT and TrEMBL sequence databases
- consists of an automated compilation of homologous domain alignment.
- 2004.1: ProDom families were generated automatically using PSI-BLAST. built from non fragmentary sequences from SWISS-PROT + TREMBL - Sept, 2003



EMBnet

MCB, feb 2005



Home



Contact



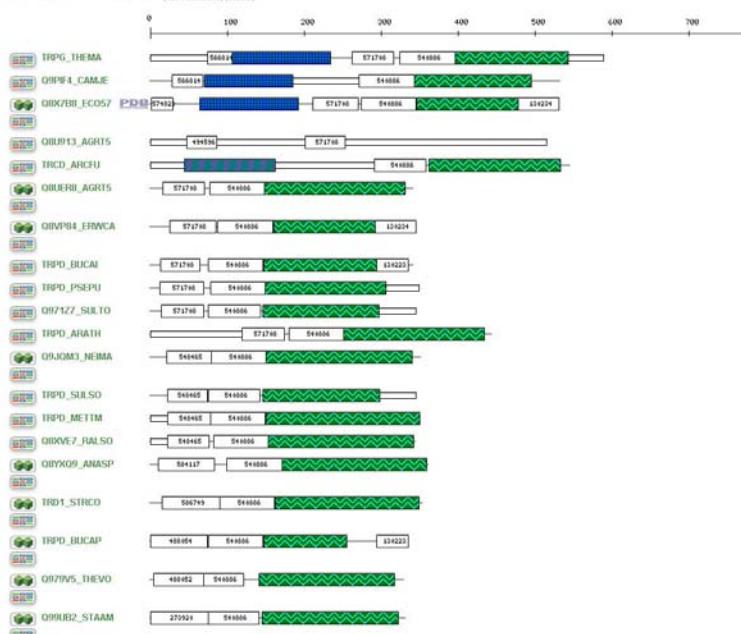
Release2002.1

Proteins sharing a prodom "Domain" with TRPD_BUCSC (Q44602)



SEARCH

[◀◀◀◀◀▶▶▶▶▶] [new window] [close]



Top

InterPro

www.ebi.ac.uk/interpro

- Search simultaneously many domain databases.
- Single set of documents linked to the various methods.
- InterPro release 8.1 contains 11330 entries representing 2933 domains, 8126 families, 222 repeats, 27 active sites, 21 binding sites and 20 post-translational modification sites.



EMBnet

MCB, feb 2005

InterPro

InterPro Entry IPR001323

Erythropoietin/thrombopoietin

Database	InterPro
Accession	IPR001323; EPO_TPO (matches 18 proteins)
Name	Erythropoietin/thrombopoietin
Type	Family
Dates	08-OCT-1999 (created) 23-NOV-2000 (last modified)
Signatures	PS00817, EPO_TPO (19 proteins) PF00728, EPO_TPO (18 proteins)
Children	IPR003013, Erythropoietin signature (12 proteins) IPR003973, Thrombopoietin (5 proteins)
Parents	
Function	glycoprotein hormone (GO:0005181)
Abstract	Erythropoietin, a plasma glycoprotein, is the primary physiological mediator of erythropoiesis [1]. It is involved in the regulation of the level of peripheral erythrocytes by stimulating the differentiation of erythroid progenitor cells, found in the spleen and bone marrow, into mature erythrocytes [2]. It is primarily produced in adult kidneys and fetal liver, acting by attachment to specific binding sites on erythroid progenitor cells, stimulating their differentiation [3]. Severe kidney dysfunction causes reduction in the plasma levels of erythropoietin, resulting in chronic anaemia - injection of purified erythropoietin into the blood stream can help to relieve this type of anaemia. Levels of erythropoietin in plasma fluctuate with varying oxygen tension of the blood, but androgens and prostaglandins also modulate the levels to some extent [3]. Erythropoietin glycoprotein sequences are well conserved, a consequence of which is that the hormones are cross-reactive among mammals, i.e. that from one species, say human, can stimulate erythropoiesis in other species, say mouse or rat [4].
Examples	• E49752 TPO_RAT • 233701 EPO_SHEEP • 233700 EPO_FELCA View examples
References	<ol style="list-style-type: none">Shoemaker C.B., Matscock L.D. <i>Murine erythropoietin gene - Cloning, expression, and human gene homology.</i> Mol Cell Biol 6 849-858(1986) [MEDLINE:87039105] [PUB00003648]Takeuchi M., Takanaki S., Miyazaki H., Kato T., Hozai S., Kobata N., Kobata A. <i>Comparative study of the asparagine-linked sugar chains of human erythropoietin purified from urine and the culture medium of recombinant Chinese hamster ovary cell.</i> J Biol Chem 263 3657-3663(1988) [MEDLINE:88153657] [PUB00002479]Lin F.R., Lin C.H., Lui P.H., Brown J.K., Egrie J.C., Smalling R., Fox G.M., Chen K.K., Castro M., Suggs S. <i>Monkey erythropoietin gene - Cloning, expression and comparison with the human erythropoietin gene.</i> Gene 44 201-209(1986) [MEDLINE:87055236] [PUB00001756]Nagao M., Suga H., Okano M., Matsuda S., Narta H., Ikura K., Sanaki R. <i>Nucleotide sequence of rat erythropoietin.</i> Biochim Biophys Acta 1171 99-102(1992) [MEDLINE:93042015] [PUB00000639]
Database links	Blocks: IPR001323 PROSITE doc: PSOC00644
Matches	Table all Graphical all Condensed graphical view

~~From a Swiss-Prot entry:~~




InterPro	IPR007219 ; Fungal_trans. IPR001138 ; Fungi_Trscrp_N. IPR005600 ; Gal4_dimer. Graphical view of domain structure.
Pfam	PF04082 ; Fungal_trans; 1. PF03902 ; Gal4_dimer; 1. PF00172 ; Zn_clus; 1. Pfam graphical view of domain structure.
PRINTS	PR00054 ; FUNGALZNCYS.
SMART	SM00066 ; GAL4; 1.
PROSITE	PS00463 ; ZN2_CY6_FUNGAL_1; 1. PS50048 ; ZN2_CY6_FUNGAL_2; 1. PROSITE graphical view of domain structure.
ProDom	[Domain structure / List of seq. sharing at least 1 domain]
BLOCKS	P04386 .



EMBnet

MCB, feb 2005

Example: GAL4_YEAST

InterPro domain architecture:



InterPro Entry	Method accession	Graphical match [?]
IPR001138;	PF00172	
IPR001138;	PR00054	
IPR001138;	PS00463	
IPR001138;	PS50048	
IPR001138;	SM00066	
IPR005600;	PF03902	
IPR007219;	PF04082	
Classification	PDB Chain/Domain ID & View 3D	PDB Chain/Structural Domains [?]
1d66	1d66a	
1d66	1d66b	
4.10.240.10.1	4.10.240.10.1	
g.38.1.1	g.38.1.1	
h.1.3.1	h.1.3.1	



EMBnet

MCB, feb 2005

Categories of databases for Life Sciences

- Sequences (DNA, protein)
- Genomics
- Mutation/polymorphism
- Protein domain/family (----> tools)
- **Proteomics** (2D gel, Mass Spectrometry)
 - 3D structure
 - Metabolism
 - Bibliography
 - 'Others' (Microarrays, Protein protein interaction...)



EMBnet

MCB, feb 2005

Databases 5: proteomics

Contain informations obtained by 2D-PAGE: images of master gels and description of identified proteins

- Examples: SWISS-2DPAGE, ECO2DBASE, Maize-2DPAGE, Sub2D, Cyano2DBase, etc.

- Composed of image and text files

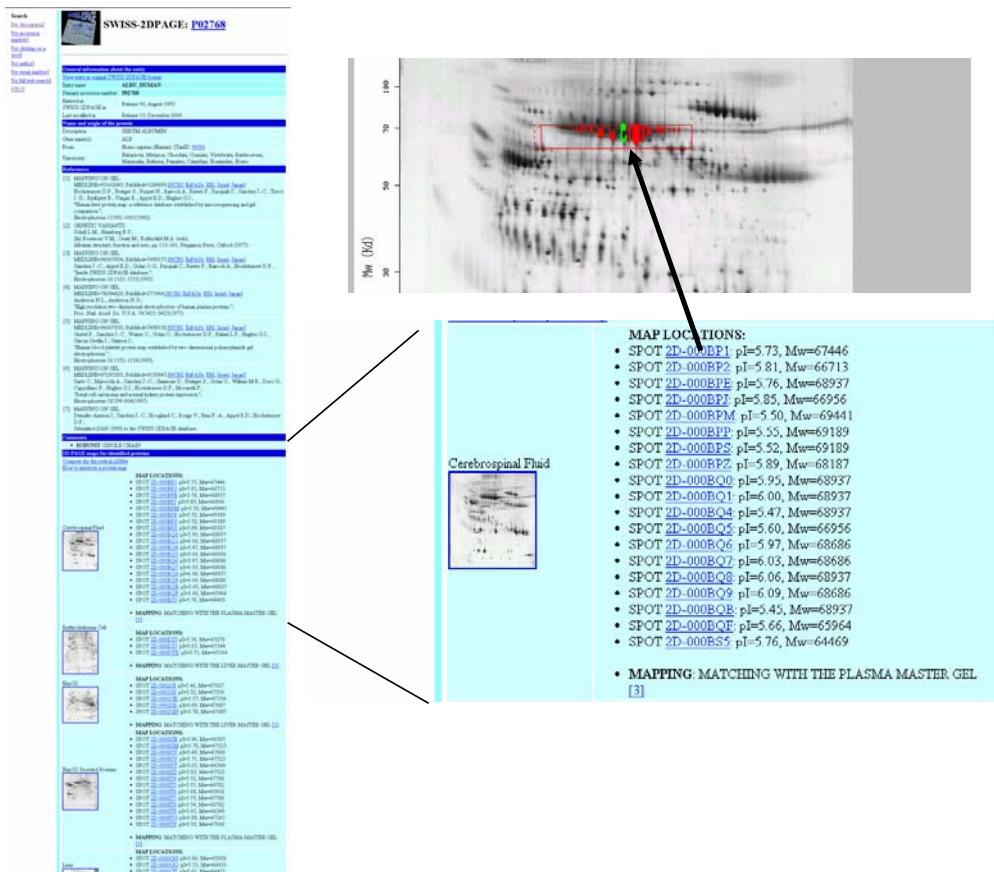
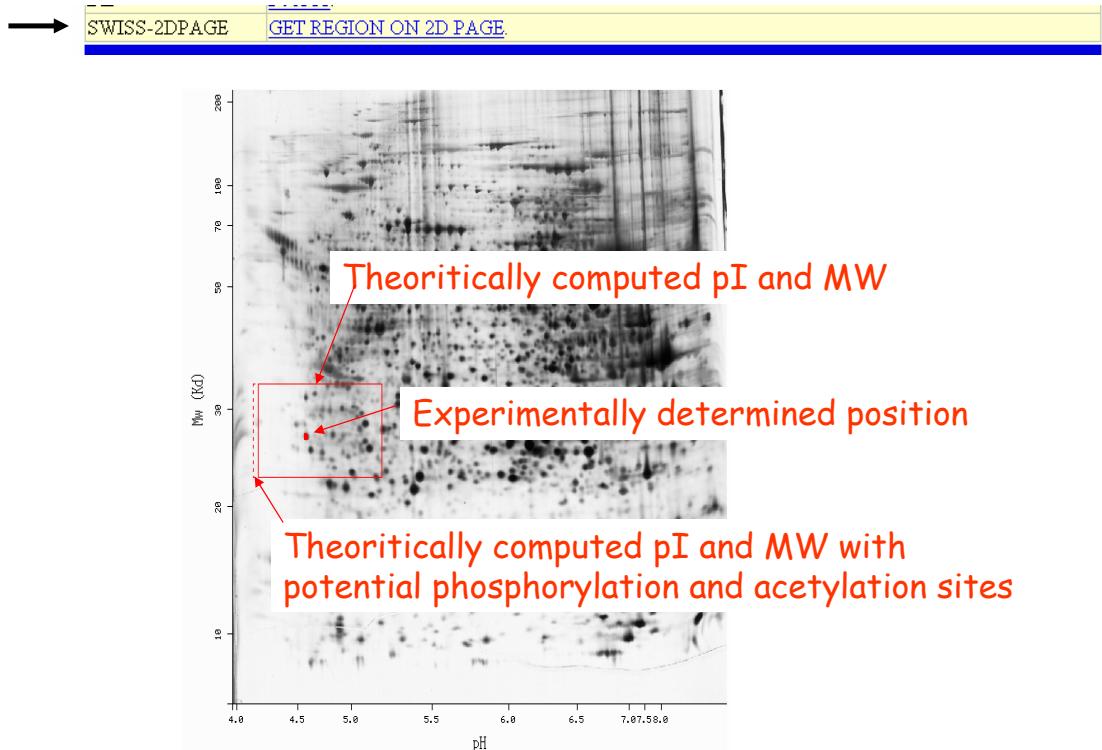
- There is currently no protein Mass Spectrometry (MS) database (not for long...)



EMBnet

MCB, feb 2005

This protein does not exist in the current release of SWISS-2DPAGE.



Categories of databases for Life Sciences

- 1) • Sequences (DNA, protein)
- 2) • Genomics → *domains, glycosylation, tools*
- 3) • Mutation/polymorphism *variants*
- 4) • Protein domain/family (----> tools) *PASS, SMART*
- 5) • Proteomics (2D gel, Mass Spectrometry) *SWISS-2DPAGE, MassSpecDB*
- 6) • 3D structure
- 7) • Metabolism
- 8) • Bibliography
- 'Others' (Microarrays, Protein protein interaction...)



EMBnet

MCB, feb 2005

Databases 6: 3D structure

- Contain the spatial coordinates of macromolecules whose 3D structure has been obtained by X-ray or NMR studies
- Proteins represent more than 90% of available structures (others are DNA, RNA, sugars, viruses, protein/DNA complexes...)
- Only one: PDB (Protein Data Bank),

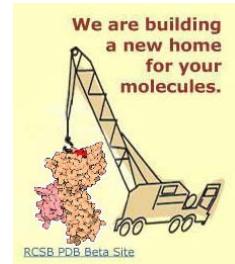


EMBnet

MCB, feb 2005

PDB: Protein Data Bank

www.rcsb.org/pdb/



- Managed by Research Collaboratory for Structural Bioinformatics (RCSB) (USA).
- Contains structure data on proteins, nucleic acids, protein-nucleic acid complexes, and viruses
- Associated with specialized programs allow the visualization of the corresponding 3D structure (e.g., SwissPDB viewer, Chimera, RasMol).
- Currently there are ~29'500 structural data for about 8'000 different proteins, but far less protein family (highly redundant)!



EMBnet

MCB, feb 2005

PDB: example

```
HEADER LYASE(OXO-ACID)          01-OCT-91 12CA 12CA 2
COMPND CARBONIC ANHYDRASE /II (CARBONATE DEHYDRATASE) (/HCA II) 12CA 3
COMPND 2 (E.C.4.2.1.1) MUTANT WITH VAL 121 REPLACED BY ALA (/V121A) 12CA 4
SOURCE HUMAN (HOMO SAPIENS) RECOMBINANT PROTEIN           12CA 5
AUTHOR S.KNAIR,D.W.CHRISTIANSON                         12CA 6
REVDAT 1 15-OCT-92 12CA 0                               12CA 7
JRNL  AUTH S.KNAIR,T.L.CALDERONE,D.W.CHRISTIANSON,C.A.FIERKE 12CA 8
JRNL  TITL ALTERING THE MOUTH OF A HYDROPHOBIC POCKET. 12CA 9
JRNL  TITL 2 STRUCTURE AND KINETICS OF HUMAN CARBONIC ANHYDRASE 12CA 10
JRNL  TITL 3 /II$ MUTANTS AT RESIDUE VAL-121        12CA 11
JRNL  REF J.BIOL.CHEM. V. 266 17320 1991   12CA 12
JRNL  REFN ASTM JBCHA3 US ISSN 0021-9258            071 12CA 13
REMARK 1                                         12CA 14
REMARK 2                                         12CA 15
REMARK 2 RESOLUTION 2.4 ANGSTROMS.             12CA 16
REMARK 3                                         12CA 17
REMARK 3 REFINEMENT.                          12CA 18
REMARK 3 PROGRAM      PROLSQ                 12CA 19
REMARK 3 AUTHORS      HENDRICKSON,KONNERT    12CA 20
REMARK 3 R VALUE      0.170                  12CA 21
REMARK 3 RMSD BOND DISTANCES 0.011 ANGSTROMS 12CA 22
REMARK 3 RMSD BOND ANGLES   1.3 DEGREES     12CA 23
REMARK 4                                         12CA 24
REMARK 4 N-TERMINAL RESIDUES SER 2, HIS 3, HIS 4 AND C-TERMINAL 12CA 25
REMARK 4 RESIDUE LYS 260 WERE NOT LOCATED IN THE DENSITY MAPS AND, 12CA 26
REMARK 4 THEREFORE, NO COORDINATES ARE INCLUDED FOR THESE RESIDUES. 12CA 27
.....
```



EMBnet

MCB, feb 2005

PDB (cont.)

```
SHEET 3 S10 PHE 66 PHE 70-1 O ASN 67 N LEU 60 12CA 68
SHEET 4 S10 TYR 88 TRP 97-1 O PHE 93 N VAL 68 12CA 69
SHEET 5 S10 ALA 116 ASN 124-1 O HIS 119 N HIS 94 12CA 70
SHEET 6 S10 LEU 141 VAL 150-1 O LEU 144 N LEU 120 12CA 71
SHEET 7 S10 VAL 207 LEU 212 1 O ILE 210 N GLY 145 12CA 72
SHEET 8 S10 TYR 191 GLY 196-1 O TRP 192 N VAL 211 12CA 73
SHEET 9 S10 LYS 257 ALA 258-1 O LYS 257 N THR 193 12CA 74
SHEET 10 S10 LYS 39 TYR 40 1 O LYS 39 N ALA 258 12CA 75
TURN 1 T1 GLN 28 VAL 31 TYPE VIB (CIS-PRO 30) 12CA 76
TURN 2 T2 GLY 81 LEU 84 TYPE II(PRIME)(GLY 82) 12CA 77
TURN 3 T3 ALA 134 GLN 137 TYPE I (GLN 136) 12CA 78
TURN 4 T4 GLN 137 GLY 140 TYPE I (ASP 139) 12CA 79
TURN 5 T5 THR 200 LEU 203 TYPE VIA (CIS-PRO 202) 12CA 80
TURN 6 T6 GLY 233 GLU 236 TYPE II (GLY 235) 12CA 81
CRYST1 42.700 41.700 73.000 90.00 104.60 90.00 P 21 2 12CA 82
ORIGX1 1.000000 0.000000 0.000000 0.00000 12CA 83
ORIGX2 0.000000 1.000000 0.000000 0.00000 12CA 84
ORIGX3 0.000000 0.000000 1.000000 0.00000 12CA 85
SCALE1 0.023419 0.000000 0.006100 0.00000 12CA 86
SCALE2 0.000000 0.023981 0.000000 0.00000 12CA 87
SCALE3 0.000000 0.000000 0.014156 0.00000 12CA 88
ATOM 1 N TRP 5 8.519 -0.751 10.738 1.00 13.37 12CA 89
ATOM 2 CA TRP 5 7.743 -1.668 11.585 1.00 13.42 12CA 90
ATOM 3 C TRP 5 6.786 -2.502 10.667 1.00 13.47 12CA 91
ATOM 4 O TRP 5 6.422 -2.085 9.607 1.00 13.57 12CA 92
ATOM 5 CB TRP 5 6.997 -0.917 12.645 1.00 13.34 12CA 93
ATOM 6 CG TRP 5 5.784 -0.209 12.221 1.00 13.40 12CA 94
ATOM 7 CD1 TRP 5 5.681 1.084 11.797 1.00 13.29 12CA 95
ATOM 8 CD2 TRP 5 4.417 -0.667 12.221 1.00 13.34 12CA 96
ATOM 9 NE1 TRP 5 4.388 1.418 11.515 1.00 13.30 12CA 97
ATOM 10 CE2 TRP 5 3.588 0.375 11.797 1.00 13.35 12CA 98
ATOM 11 CE3 TRP 5 3.837 -1.877 12.645 1.00 13.39 12CA 99
ATOM 12 CZ2 TRP 5 2.216 0.208 11.656 1.00 13.39 12CA 100
ATOM 13 CZ3 TRP 5 2.465 -2.043 12.504 1.00 13.33 12CA 101
ATOM 14 CH2 TRP 5 1.654 -1.001 12.009 1.00 13.34 12CA 102
```

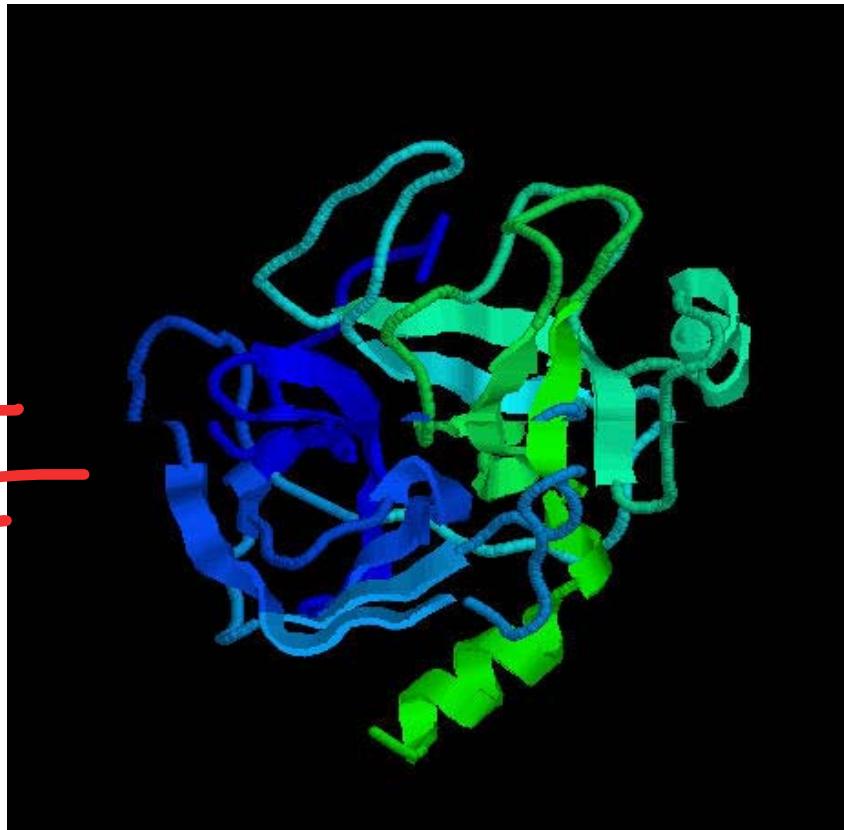
Coordinates of each atom



EMBnet

MCB, feb 2005

The same
PDB entry
"visualized"
with Chimera



Industry of databases around PDB

- HSSP: Homology-derived secondary structure of proteins.
<http://www.sander.ebi.ac.uk/hssp/>

- Structure classification

-CATH

-SCOP

-...

- Homology-derived 3D structure db:

Swiss-Model Redepository (SMR): feb 2005: 555'900 models.



EMBnet

MCB, feb 2005

SWISS-MODEL REPOSITORY

[Home](#) | [Advanced Search](#) | >>[Swiss-Model](#) | [HELP](#)

Swiss-Prot/TrEMBL AC or ID:

Welcome to the SWISS MODEL Repository

The SWISS-MODEL Repository is a database of annotated three-dimensional comparative protein structure models generated by the fully automated homology-modelling pipeline SWISS-MODEL. The repository is developed at the Biozentrum Basel within the Swiss Institute of Bioinformatics.



REPOSITORY STATUS

Models:	482043
UniProt:	2.5
PDB:	20.08.2004
Last update:	28.09.2004



NOTE: The SWISS-MODEL repository contains theoretically calculated models, which may contain significant errors.

Jürgen Kopp & Torsten Schwede

swissmodel.expasy.org/repository

[[Disclaimer](#)] [[HELP](#)] [[Contact](#)]

Annotated 3D comparative protein structure models generated by the fully automated homology-modelling pipeline SWISS-MODEL.

Precompute the 3D model of protein domains (~200 amino acids, biggest model: 1500 aa) which share about 40 % similarity with a 3D experimentally determined template.

<http://swissmodel.expasy.org/repository/>



EMBnet

MCB, feb 2005

Categories of databases for Life Sciences

- Sequences (DNA, protein)
- Genomics
- Mutation/polymorphism
- Protein domain/family (----> tools)
- Proteomics (2D gel, Mass Spectrometry)
- 3D structure
- Metabolism
- Bibliography
- 'Others' (Microarrays, Protein protein interaction...)



EMBnet

MCB, feb 2005

~~Databases 7: metabolic~~

- ✓ ~~Contain informations that describe enzymes, biochemical reactions and metabolic pathways;~~
- ~~ENZYME and BRENDA: nomenclature databases that store informations on enzyme names and reactions;~~
- Metabolic databases: EcoCyc (specialized on *Escherichia coli*), KEGG, EMP/WIT;
Usually these databases are tightly coupled with query software that allows the user to visualise reaction schemes.



EMBnet

MCB, feb 2005

- There are about 3750 "EC numbers"
 - ~ 1900 are linked to Swiss-Prot sequence
 - ~ 200 are linked to a TrEMBL sequence
 - ~ 1450 can not be linked to any sequence !



EMBnet

MCB, feb 2005

~~BRENDA~~
~~Useful to prepare
lab's experiments !~~

<http://www.brenda.uni-koeln.de/>

BRENDA
The Comprehensive Enzyme Information System

Complete Entry of EC-Number 1.2.3.4
[open printable version in a new window](#)

EC NUMBER	CODACTOR	SPECIFIC ACTIVITY [nmol/min/mg]	LOCATION
ORGANISM	KM VALUE [mM]	PURIFICATION	
SYSTEMATIC NAME	PH OPTIMUM	MOLECULAR WEIGHT	
RECOMMENDED NAME	PH RANGE	SUBUNITS	
SYNONYMS	TEMPERATURE OPTIMUM [°C]	PH STABILITY	
CAS REGISTRY NUMBER	METALS, IONS	TEMPERATURE STABILITY [°C]	
REACTION	INHIBITORS	STORAGE STABILITY	
REACTION TYPE	SOURCE/TISSUE	LINKS TO OTHER DATABASES	
SUBSTRATES/PRODUCTS		REFERENCES	

<...> = reference, #_#=organism

EC NUMBER
1.2.3.4

ORGANISM
#18 *Sorghum vulgare* (Variety CSH-5 <1>, CSH-1 <3>) <1, 3>
#28 *Pseudomonas* sp. (OCX-53) <2>
#38 *Hordeum vulgare* (cv.stripy) <4>
#44 *Tilapia contraversa* (parrotfish, fungus) <5>
#58 *Hylaeocodium splendens* <6>
#68 *Erythridialphus squamatus* <6>
#78 *Hylaeocodium ferrum* <6>

SYSTEMATIC NAME
Oxalate oxygen oxidoreductase

RECOMMENDED NAME
Oxalate oxidase

SYNONYMS
Oxidase, oxalate
Ferro-manganese dehydrogenase
Oxalic acid oxidase #48 <5>

CAS REGISTRY NUMBER
9031-79-2

REACTION
Oxalate + O₂ = 2 CO₂ + H₂O₂

REACTION TYPE
Redox reaction

SUBSTRATES/PRODUCTS
#1 Oxalate + O₂ #1-78 (almost specific for oxalate, #28 <2>; no other electron acceptor found, #38 <4>) <1-5>
#1 CO₂ + H₂O₂
#2 Additional information: #28 (oxidation at extreme low rate: glyoxylic acid, DL-malic acid, citric acid, not oxidized: succinic acid, formic acid, <2>)
#2 ?

SPECIFIC ACTIVITY [nmol/min/mg]

IntEnz = Enzyme + BRENDA + NC-IUBMB nomenclature
<http://www.ebi.ac.uk/intenz/index.html>



EMBL-EBI
European Bioinformatics Institute

EBI Home About EBI Research Services Toolbox Databases Downloads Submissions Integrated Enzyme Database (IntEnz)

IntEnz Home

News Advanced search Proposed entries Display information About IntEnz Future development

Exact phrase

[EC 1 Oxidoreductases](#)
[EC 2 Transferases](#)
[EC 3 Hydrolases](#)
[EC 4 Lyases](#)
[EC 5 Isomerases](#)
[EC 6 Ligases](#)

News

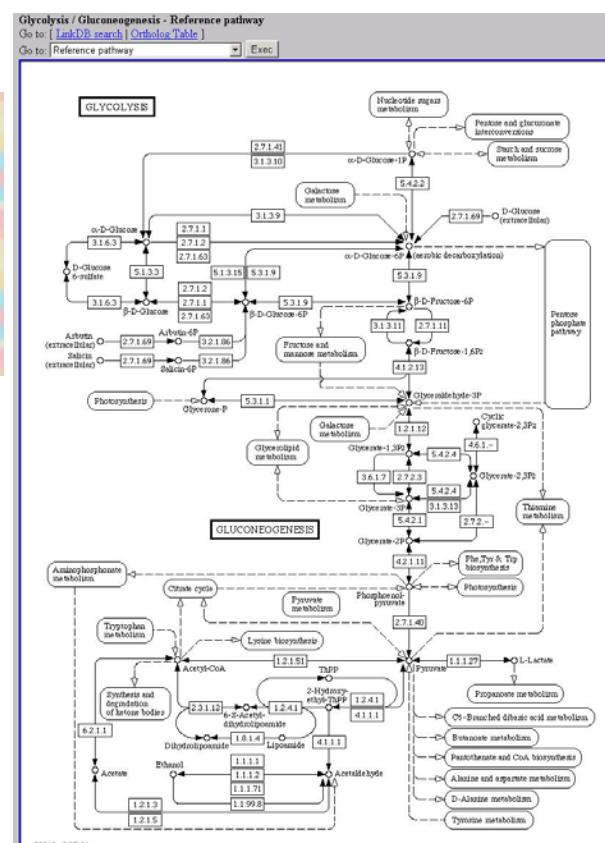
IntEnz NEWS
IntEnz release 3
... more

IntEnz is the name for the Integrated relational Enzyme database and is the most up-to-date version of the Enzyme Nomenclature. The Enzyme Nomenclature comprises the recommendations of the Nomenclature Committee of the International Union of Biochemistry and Molecular Biology (NC-IUBMB) on the Nomenclature and Classification of Enzyme-Catalysed Reactions.



FMRnet

MCB, feb 2005



<http://www.genome.ad.jp/kegg>

Categories of databases for Life Sciences

- Sequences (DNA, protein)
- Genomics
- Mutation/polymorphism
- Protein domain/family (----> tools)
- Proteomics (2D gel, Mass Spectrometry)
- 3D structure
- Metabolism
- **Bibliography**
- 'Others' (Microarrays, Protein protein interaction...)



EMBnet

MCB, feb 2005

Databases 8: bibliographic

- ✓ Bibliographic reference databases contain citations and abstract informations of published life science articles;
- ✓ Example: Medline
- Other more specialized databases also exist (i.e. Agricola <http://agricola.nal.usda.gov/>, EMBASE (not free)...).



EMBnet

MCB, feb 2005

~~Medline~~

- Comprehensive database of primary scientific literature in the biomedical area.
- More than 4,000 biomedical journals published in the United States and 70 other countries
- Contains over 15 million indexed citations since 1966 until now
- Citations prior to the mid-1960s are located in OLDMEDLINE.
- Contains links to biological db
 - Many papers not dealing with humans are not in Medline !
 - Before 1970, keeps only the first 10 authors !
 - Not all journals have citations since 1966 ! (they go back...)
 - Indexed by Google in 2004 !



EMBnet

MCB, feb 2005

~~PubMed~~

<http://www.ncbi.nlm.nih.gov/entrez/query.fcgi?db=PubMed>

- Maintained by the US National Library of Medicine.
- Allows access to the citations from MEDLINE and additional life science journals.
- Includes links to many sites providing full text articles and other related resources.
- Gives also access to :
 - In Process Citations
 - Publisher supplied citations: citations directly submitted to PubMed ([Record as supplied by publisher]).
- PMID (PubMed ID) UI (Medline ID)

New:

DOI (Digital Object Identifier) are names (characters and/or digits) assigned to objects of intellectual property such as electronic journal articles, images, learning objects, ebooks, any kind of content.

Server: <http://dx.doi.org>

-> biggest advance to track documents on the web !



EMBnet

MCB, feb 2005

Categories of databases for Life Sciences

- Sequences (DNA, protein)
- Genomics
- Mutation/polymorphism
- Protein domain/family (----> tools)
- Proteomics (2D gel, Mass Spectrometry)
- 3D structure
- Metabolism
- Bibliography
- 'Others' (Microarrays, Protein protein interaction...)



EMBnet

MCB, feb 2005

Databases 9: others

- There are many databases that cannot be classified in the categories listed previously;
- Examples: ReBase (restriction enzymes), TRANSFAC (transcription factors), CarbBank, GlycoSuiteDB (linked sugars), Protein-protein interactions db (Intact, BIND), Protease db (MEROPS), biotechnology patents db, etc.;
- As well as many other resources concerning any and new aspects of macromolecules and molecular biology (Microarrays).



EMBnet

MCB, feb 2005

Amos links: Microarrays

Microarray databases and resources

- [ArrayExpress](#) - MicroArray informatics at the EBI
- [ChipDB](#) - Gene expression db
- [ExpressDB](#) - Yeast and E.coli RNA expression db
- [GeneX](#) - Gene Expression project
- [GEO](#) - Gene Expression Omnibus
- [RAD](#) - RNA abundance db
- [SMD](#) - Stanford microarray db



EMBnet

MCB, feb 2005

Interactome

- Protein/protein interaction:
description from 1 to more than 20'000 interactions / publication
 - Several databases: Intact, BIND, DIP.
 - Proteomics standard initiative since 2005



EMBL-EBI
European Bioinformatics Institute

EBI Home About EBI Research Services Toolbox Databases Downloads HOME

IntAct Project

Search IntAct:

- [IntAct Home](#)
- [Search](#)
- [Project](#)
- [Members](#)
- [Documentation](#)
- [Software](#)
- [Data](#)
- [Statistics](#)
- [Availability](#)
- [Status](#)
- [Data submission](#)
- [Mailing Lists](#)
- [Publications](#)
- [Acknowledgements](#)

• Examples:

- Gene name: [lsm7](#)
- IntAct Ac: [EBI-141](#)
- SPTR Ac: [Q08162](#)
- SPTR Id: [m44_yeast](#)
- InterPro Ac: [IPR001900](#)
- GO Id: [GO:0000176](#)

Project description

- IntAct provides a freely available, open source database system and analysis tools for protein interaction data. All interactions are derived from literature curation or direct user submissions and are freely available.

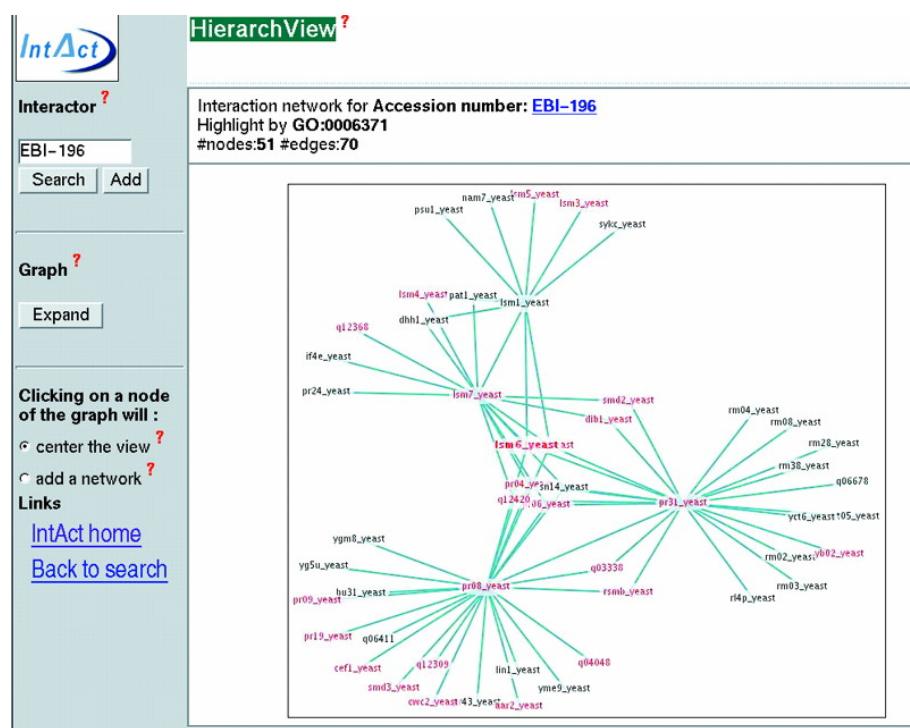
Members of the IntAct consortium

<http://www.ebi.ac.uk/intact/index.html>



EMRnet

MGR feb 2005



EMBnet

MCB, feb 2005

Gene Ontology (GO) database

The Gene Ontology (GO) project (<http://www.geneontology.org/>) provides structured, controlled vocabularies and classifications that cover several domains of molecular and cellular biology and are freely available for community use in the annotation of genes, gene products and sequences.

The three organizing principles of GO are **molecular function (MF)**, **biological process (BP)** and **cellular component (CC)**.



EMBnet

MCB, feb 2005

Proliferation of databases

- Which does contain the highest quality data ?
- Which is the more comprehensive ?
- Which is the more up-to-date ?
- Which is the less redundant ?
- Which is the more indexed (allows complex queries) ?
- Which Web server does respond most quickly ?
-??????



EMBnet

MCB, feb 2005

To benefit from the data stored in a database, we need:

- easy access to the information

→ a method for extracting only that information needed to answer a specific biological question

Examples: Entrez (NCBI), SRS (Europ), tools such as BLAST, Peptident...



EMBnet

MCB, feb 2005

Some important practical remarks

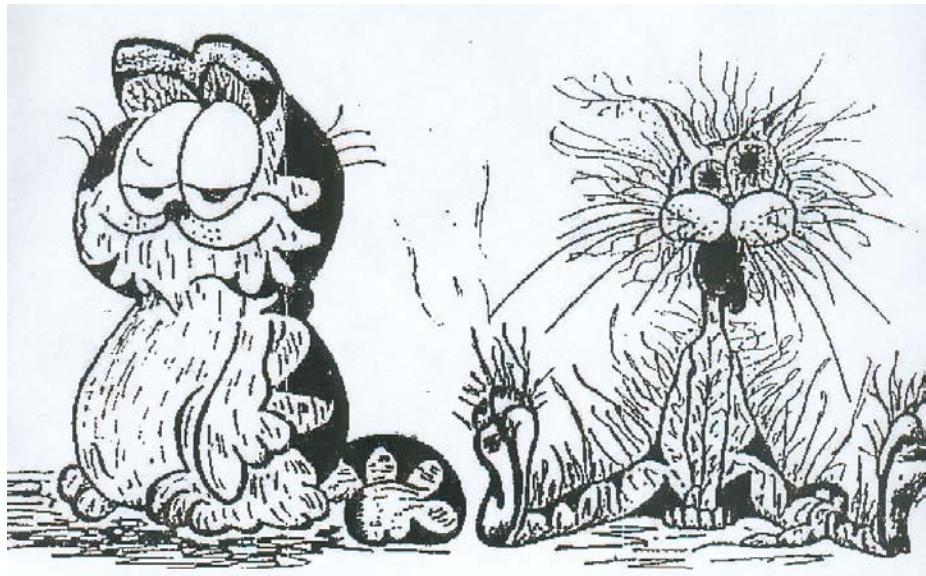
- Databases: many errors (automated annotation) !
- Not all db are available on all servers
- The update frequency is not the same for all servers;
- Some servers add automatically cross-references to an entry (implicit links) in addition to already existing links (explicit links)...different looks...



EMBnet

MCB, feb 2005

Before the introduction to databases...



~~After~~ the introduction to databases...



EMBnet

MCB, feb 2005