

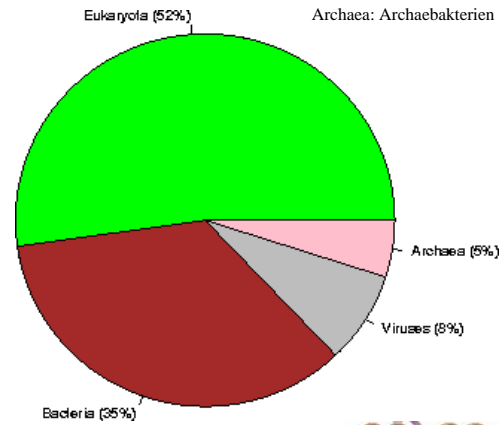
Protein-Datenbanken

- Sequenz-Datenbanken (Vertreter: Swiss-Prot)
- Domain/Familien-Datenbanken (Vertreter: InterPro)
- Struktur-Datenbanken (Vertreter: PDB)
- Vorsicht: Die Grenzen zwischen diesen Datenbank-Typen sind unscharf!



Swiss-Prot

- <http://www.ebi.ac.uk/swissprot/>
- Repository aller bekannten Proteinsequenzen
- Basiert auf Submission, Übersetzung und aktiver Suche, intensive (manuelle) Datenpflege
 - > 30 "Scientific Database Curators"
 - Redundanzfreiheit
 - Vierteljährliche Releases
- Tools für Protein-Analyse (z.B. Homologie-Modellierung)



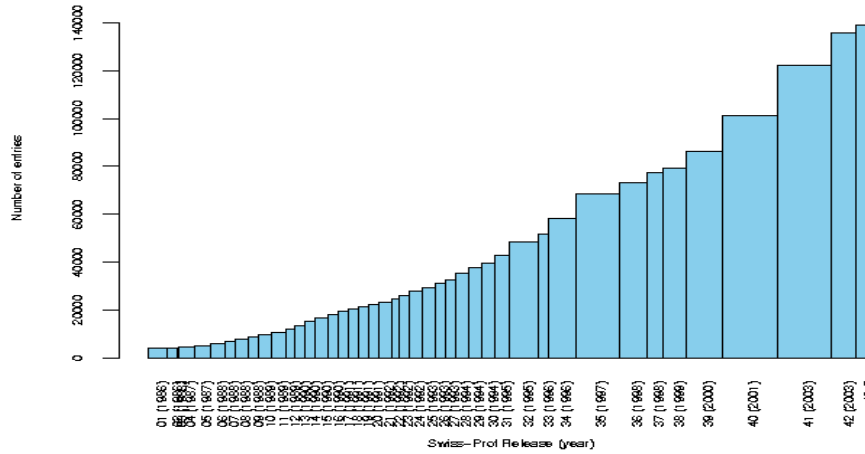
Swiss-Prot: Wachstum

Release 42.5 of 21-Nov-2003 of Swiss-Prot contains 138922 sequence entries, comprising 51131444 amino acids abstracted from 110725 references.

3077 sequences have been added since release 42, the sequence data of 201 existing entries has been updated and the annotations of 8149 entries have been revised. This represents an increase of 2%.

The growth of the database is summarized below.

Size of the Swiss-Prot database



Swiss-Prot: Daten

- "Flaches" Datenmodell (Entry-basiertes Modell), sehr ähnlich zu EMBL
 - Autor, Datum, Länge, Methode, letzte Änderung
 - Organismus
 - Proteinsequenz (z.B. im FASTA-Format)
 - Links zu anderen Datenquellen, Literaturreferenzen
- Oracle-Dumps verfügbar (ca. 140 Tabellen)
- XML-Export
- Keine Änderungsübersicht!
- TrEMBL (Translations of EMBL)
 - Supplement zu Swiss-Prot
 - Enthält alle automatisch in AS-Sequenzen übersetzte CDS-Sequenzen aus EMBL
 - Keine Überschneidung mit (manuell) eingebrachten Swiss-Prot-AS-Sequenzen
 - SP-TrEMBL: Geplanter Nachfolger von Swiss-Prot



Swiss-Prot: Beispieleintrag

```

ID GUNB_CLOTM STANDARD; PRT; 563 AA.
AC P04956;
DT 13-AUG-1987 (Rel. 05, Created)
DT 13-AUG-1987 (Rel. 05, Last sequence update)
DT 01-FEB-1995 (Rel. 31, Last annotation update)
DE ENDOGLUCANASE B PRECURSOR (EC 3.2.1.4) (EGB) (ENDO
DE (CELLULASE B).
GN CE1B.
OS Clostridium thermoCELLUM.
OC Bacteria; Firmicutes; Bacillus/Clostridium group;
OC Clostridium.
OX NCBI_TaxID=1515;
RN [1]
RP SEQUENCE FROM M.A.
RC STRAIN=NCIB 10682;
RX MEDLINE=86148508; PubMed=3453102;
RA Grepinet O., Beguin P.;
RT "Sequence of the cellulase gene of Clostridium the
RT endoglucanase B.";
RL Nucleic Acids Res. 14:1791-1799(1986).
CC -!- FUNCTION: THIS ENZYME CATALYZES THE ENDOHYDROL
CC GLUCOSIDIC LINKAGES IN CELLULOSE, LICHENIN AND
CC GLUCANS.
CC -!- CATALYTIC ACTIVITY: ENDOHYDROLYSIS OF 1,4-BETA
CC LINKAGES IN CELLULOSE.
CC -!- DOMAIN: A 24 RESIDUES DOMAIN IS REPEATED TWICE
CC WELL AS IN OTHER C.THERMOCELLUM CELLULOSOME EN
CC MAY FUNCTION AS THE BINDING LIGAND FOR THE SL
CC -!- SIMILARITY: BELONGS TO CELLULASE FAMILY A (FAM
CC HYDROLASES).
DR EMBL; X03592; CAA27266.1; -.
DR PIR; A23512; CZCLBM.
DR HSP; P54583; ICEE.
DR InterPro; IPR002105; Dockerin_1.
DR InterPro; IPR002048; EF-hand.
DR InterPro; IPR001547; Glyco_hydro_F5.
DR Pfam; PF00150; cellulase; 1.
DR Pfam; PF00404; Dockerin_1; 2.
DR PROSITE; PS00018; EF_HAND; UNKNOWN_1.
DR PROSITE; PS00448; CLOS_CELLULOSOME_RPT; 2.
DR PROSITE; PS00659; GLYCOSYL_HYDROL_F5; 1.
KW Cellulose degradation; Hydrolase; Glycosidase; Repeat; Signal.
FT SIGNAL 1 27 OR 31.
FT CHAIN 28 563 ENDOGLUCANASE B.
FT ACT_SITE 204 204 PROTON DONOR (BY SIMILARITY).
FT ACT_SITE 363 363 NUCLEOPHILE (BY SIMILARITY).
FT DOMAIN 502 557 2 X 24 AA APPROXIMATE REPEATS.
FT REPEAT 502 526 1.
FT REPEAT 534 557 2.
SQ SEQUENCE 563 AA; 63929 MW; 866FE55704A1DE4B CRC64;
MKKFLVLLIA LIMIATLLVV PGVQTSAEGS YADLAEPPDD WLHVEGTNIV DKYGNKVVIT
GANWFGMCR ERMLLDSYHS DIADIELVA DKGIVNVVRMP IATDLLYAGS QGIYPPSDT
SYMNALAGL NSYELFMFL ENFKRVGIRV ILDVHSPETD NQGHNYPLVY NTLITEIIFK
KAVVVAERY KNDDTIIGFD LKMEPHNTG TMKIKAQSAI WDDSNHPNHW KRVAEETALA
ILEVHFNVLV FVEGVEMYPK DGIWDETFD TSPVTGNNDY YGNWGGNLR GVKDYFIMLG
KYQSLVYSP HDYGPVYEQ DWFKGFITA NDEQAKRLY EQCWRDWMAY IMEEGISPLL
LGEWGMTEG GHPLLDLNLK YLRMRDFIL ENKYKRLHTF WCINIDSADT GGLFTRDEGT
PFPGRDLKW MDMKYDMYLY PVLWKTEDGK FIGLDHKIPL GRNGISISQL SMYTPSVTPS
PSATPSPTTI TAPPTDVTY GDVNGDGRVN SSDVALLKRY LGLVENVINK EAADVNVSGT
VNSDLDLAIMK RYVLRISISEL PYK
//

```



Swiss-Prot: FASTA-Format

■ Alternatives Format für AS-Sequenzen

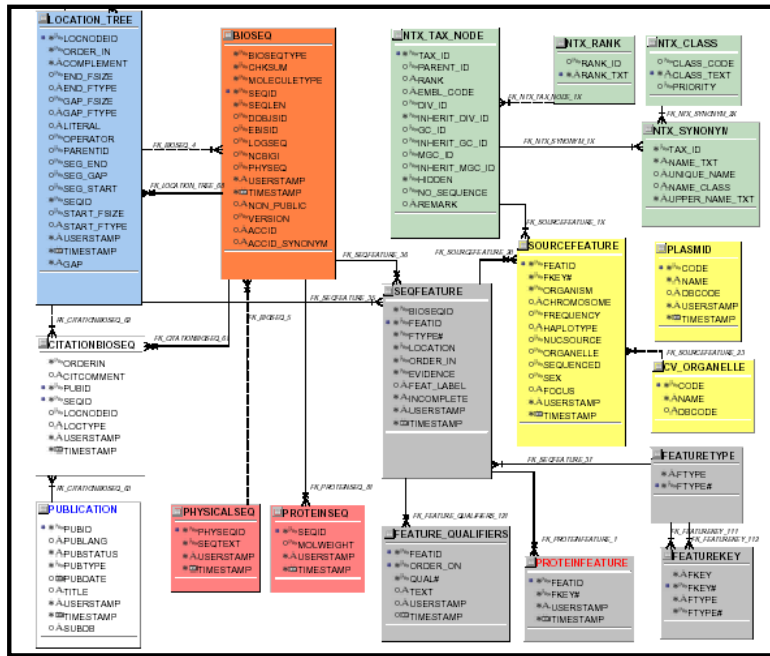
- begins with a description line indicated by a “>” sign
- followed by amino acid seq. in capital letters,
- no numbers, no blocks
- line length usually 80 characters

Example:

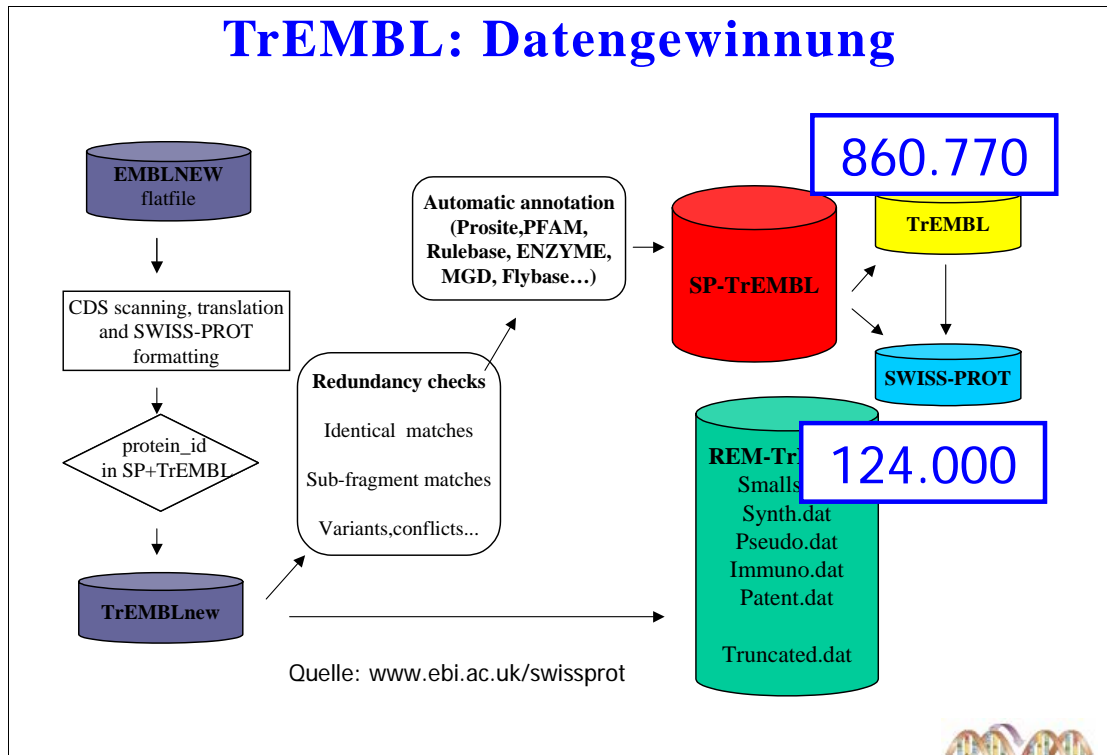
```
>gi|532319|pir|TVFV2E|TVFV2E envelope protein
ELRLRYCAPAGPALLKCNADADYDGFKTNCSNVSVVHCTNLMNTTVTGLLNGSYSENRT
QIWQKHRTSNDSALILLNKHYNLTVTCKRPGNKTVLPVTIMAGLVFHSQKYNLRLQAWC
HFPSNWKGAWEVKEEIVNLPKERYRGTNDPKRIFFRQWGDPEANLWFNCHGEFFYCK
MDWFLNLYLNNLTVDADHNECKNTSGTKSGNKRAPGPCVQRTYVACHIRSVIIWLETISKK
TYAPPRECHLECTSTVTGMTVELNYIPKNRTNVTLSFQIESIWAELDRYKLVETPIGF
APTEVRRYTGGERQKRVFVXXXXXXXXXXXXXXXXXXXXXXXXXVQSQHLLAGILQQQKNL
LAAVEAQQQMLKLTIWGVK
```



Swiss-Prot: Relationales Schema



TrEMBL: Datengewinnung



Swiss-Prot Web Interface

NiceProt View of SWISS-PROT: P29358

P29358

Printer-friendly view Quick BlastP search

[\[General\]](#) [\[Name and origin\]](#) [\[References\]](#) [\[Comments\]](#) [\[Cross-references\]](#) [\[Keywords\]](#) [\[Features\]](#) [\[Sequence\]](#) [\[Tools\]](#)

General information about the entry

Entry name	143B_BOVIN
Primary accession number	P29358
Secondary accession numbers	None
Entered in SWISS-PROT in	Release 24, December 1992
Sequence was last modified in	Release 33, February 1996
Annotations were last modified in	Release 41, June 2002

Name and origin of the protein

Protein name	14-3-3 protein beta/alpha
Synonyms	Protein kinase C inhibitor protein-1 KCIP-1
Gene name	YWHAB
From	Bos taurus (Bovine) [TaxID: 9913] Ovis aries (Sheep) [TaxID: 9940]
Taxonomy	Eukaryota ; Metazoa ; Chordata ; Craniata ; Vertebrata ; Euteleostomi ; Mammalia ; Eutheria ; Cetartiodactyla ; Ruminantia ; Pecora ; Bovoidea ; Bovidae ; Bovinae ; Bos

References

[1] SEQUENCE
SPECIES=Bovine;
MEDLINE=91108808; PubMed=1671102; [NCBI, ExPASy, EBI, Israel, Japan]

Dokument: Done (2.273 Sek.)



Swiss-Prot: Annotationen

- CC-Felder für Kommentare
 - Unterteilt in Topics
 - Beispiele: Caution, Disease, Function, Regulation, ...
- FT: Feature Table
 - Modifikationen, Sequenzabschnitte, Sekundärstruktur
- KW: Keywords
 - Ca. 800 verschiedene Keywords
- Einträge oft Mischung aus Controlled Vocabularies und Freitext
- Seit kurzem: Evidence Codes für alle Annotationen (Curator, Opinion, By Similarity, Experiment, ...)



Swiss-Prot: Versionierung / Identifikation

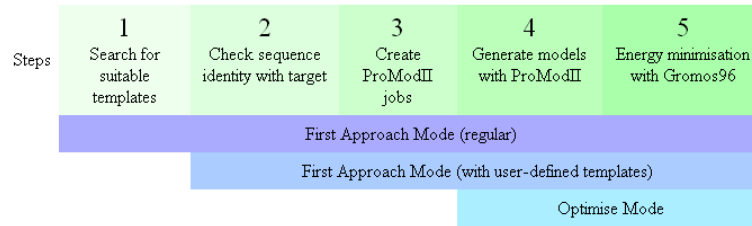
- Swiss-Prot Release ca. alle 3 Monate
- ID und AC Line
 - ID: X_Y; X: "Name" des Proteins; Y: "Name" der Spezies
 - Keine Standards für Proteinnamen
 - Spezies mit wissenschaftlichen oder umgangssprachlichen Namen
 - AC: Accession Number
 - Primäre ID
 - Kann mehrere Einträge enthalten (Merged Entries)
- Keine Versionen von Einträgen
 - Last Update
 - Keine Änderungsübersichten



Swiss-Model: 3D-Strukturbestimmung

■ Ausgangsproblematik: Nach derzeitigem biochemischem Kenntnisstand ist es nur in Ausnahmefällen möglich, von der AS-Sequenz auf die 3-Struktur zu schließen

■ Ausweg: Vergleich mit ähnlichen Sequenzen und deren Struktur (falls bekannt)



Step	Program/Method	Database	Action
1	BLASTP2	ExNRL-3D	Will find all similarities of target sequence with sequences of known structure.
2	SIM	-	Will select all templates with sequence identities above 25% and projected model size larger than 20 residues. Furthermore, this step will detect domains which can be modelled based on unrelated templates
3	-	-	Generate ProModII input files
4	ProModII	ExpDB	Generate all models
5	Gromos96	-	Energy minimisation of all models



BLAST 2.0

- **BLAST**
 - verschiedene Varianten für AS und Nuk. Sequenzen
 - schneller als dyn. Programmierung
 - aber weniger sensitiv und berücksichtigt keine Lücken
 - findet lokales statt globales Alignment
 - nutzt AS Austauschmatrizen, PAM, BLOSUM
- **Ansatz zur Beschleunigung**
 - Anfrage-Sequenz wird in Wörter der Länge W ($W=3$) zerlegt
 - Wortliste wird um ähnliche Wörter erweitert
 - Nur Worte mit Score $\geq T$ werden in DB gesucht
 - Wort-Treffer werden nach links und rechts erweitert



BLAST 2.0

- Beispiel: (W=2, T=8) Anfrage: qlnfsagw

Initiales Wort Erweiterte Liste

ql ql, qm, hl, zl

ln ln, lb

nf nf, af, ny, df, qf, ef, gf, hf, kf, sf, tf, bf, zf

fs fs, fa, fn, fd, fg, fp, ft, fb, ys

sa nothing scores 8 or higher

ag ag

gw gw, aw, rw, nw, dw, qw, ew, hw, iw, kw, mw, pw, sw, tw, vw, bw, zw, xw

- Invertierte Liste als Index nutzen
 - bei W=3 und 20AS nur 8000 verschiedene Worte möglich
- Probleme
 - ABCDEGH und ABCDEEFGH



Swiss-Model

The screenshot displays the Swiss-PdbViewer 3.7 interface. The main window shows a ribbon representation of a protein structure. The 'Alignment' panel at the bottom left compares the 'TARGET' sequence with three reference sequences (TRP35, TRP36, TRP37). The 'Control Panel' on the right lists various residues and their visibility and movement status. The 'Ramachandran Plot' on the right shows the distribution of phi and psi angles for the protein structure.

Alignment

?	TARGET	D	W	H	V	D	V	M	D	G	R	F	V	P	N	I	T	I	G	A	P	V
1	TRP35	D	W	H	V	D	V	M	D	G	R	F	V	P	N	I	T	I	G	A	P	V
2	TRP36	D	W	H	V	D	V	M	D	G	R	F	V	P	N	I	T	I	G	A	P	V
3	TRP37	D	W	H	V	D	V	M	D	G	R	F	V	P	N	I	T	I	G	A	P	V

Control Panel

TARGET	visible	?	can move	checkbox
s TRP35	<input checked="" type="checkbox"/>		<input checked="" type="checkbox"/>	<input type="checkbox"/>
s LEU36	<input checked="" type="checkbox"/>		<input checked="" type="checkbox"/>	<input type="checkbox"/>
s HIS37	<input checked="" type="checkbox"/>		<input checked="" type="checkbox"/>	<input type="checkbox"/>
s MET38	<input checked="" type="checkbox"/>		<input checked="" type="checkbox"/>	<input type="checkbox"/>
s ASP39	<input checked="" type="checkbox"/>		<input checked="" type="checkbox"/>	<input type="checkbox"/>
s ILE40	<input checked="" type="checkbox"/>		<input checked="" type="checkbox"/>	<input type="checkbox"/>
s MET41	<input checked="" type="checkbox"/>		<input checked="" type="checkbox"/>	<input type="checkbox"/>
s ASP42	<input checked="" type="checkbox"/>		<input checked="" type="checkbox"/>	<input type="checkbox"/>
s GLY43	<input checked="" type="checkbox"/>		<input checked="" type="checkbox"/>	<input type="checkbox"/>
s HIS44	<input checked="" type="checkbox"/>		<input checked="" type="checkbox"/>	<input type="checkbox"/>
s PHE45	<input checked="" type="checkbox"/>		<input checked="" type="checkbox"/>	<input type="checkbox"/>
s VAL46	<input checked="" type="checkbox"/>		<input checked="" type="checkbox"/>	<input type="checkbox"/>
s PRO47	<input checked="" type="checkbox"/>		<input checked="" type="checkbox"/>	<input type="checkbox"/>

Ramachandran Plot

The Ramachandran plot shows the distribution of phi (φ) and psi (ψ) angles. The x-axis represents phi (φ) and the y-axis represents psi (ψ), both ranging from -180 to 180 degrees. The plot shows several peaks, indicating regions of high conformational stability.

C:\Program Files\DeepView\spdbv\temp\proslst2.txt

```

#PSOSITE
PS00001; N-glycosylation site.
PS00004; cAMP- and cGMP-dependent protein kinase
PS00005; Protein kinase C phosphorylation site.
PS00006; Casein kinase II phosphorylation site.
PS00007; Tyrosine kinase phosphorylation site.
PS00008; N-myristoylation site.
PS00342; Microbodies C-terminal targeting signal
PS01085; Ribulose-phosphate 3-epimerase family s
PS01086; Ribulose-phosphate 3-epimerase family s
    
```



InterPro

■ Sekundärdatenbank zu Proteinsequenzen (Schwerpunkt: Protein-Domains)

■ Motivation

- Bestimmte Sequenzabschnitte (Motifs) bestimmen Funktion des Proteins
- Datenbanken zur Beschreibung interessanter Domänen (Proteinfamilien) nötig
- Untersuchung neuer Sequenzen auf Vorhandensein bekannter Domänen – Rückschlüsse auf Funktion

■ InterPro: Integrierte Datenbank von Proteindomänen-Datenbanken

The InterPro consortium:

- Co-ordinated by EBI (R. Apweiler & team) 
- PROSITE (A. Bairoch, P. Bucher, N. Hulo, C. Sigrist, L. Cerutti, M. Pagni, L. Falquet) 
- PRINTS (T. Attwood, P. Bradley) 
- PFAM (R. Durbin, A. Bateman, S. Griffiths-Jones) 
- PRODOM (D. Kahn, F. Servant) 
- SMART (C. Ponting, R. Copley, N. Dickens)  
- TIGRFAMs (D. Haft, O. White) 

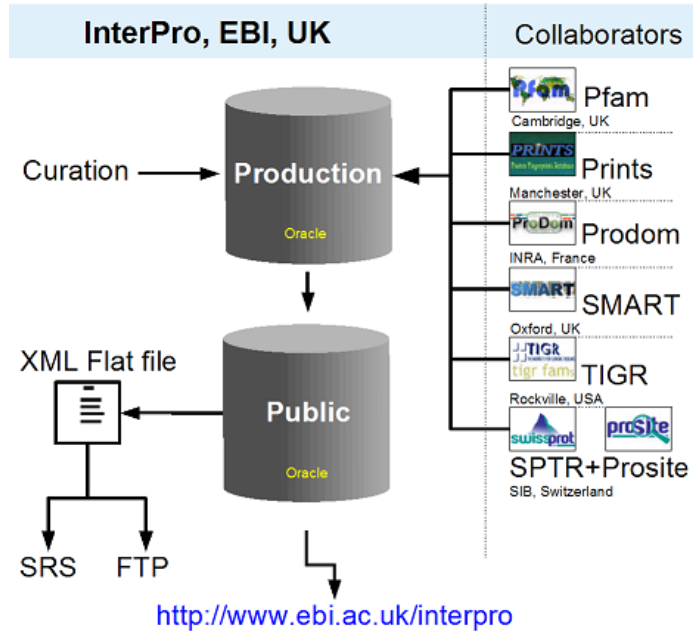


InterPro: Biologischer Fokus

- **Family** - group of evolutionarily related proteins, that share one or more domains/repeats in common.
- **Domain** - independent structural unit which can be found alone or in conjunction with other domains or repeats.
- **Repeat** - region occurring more than once that is not expected to fold into a globular domain on its own.
- **PTM** (post-translational modification) -The sequence motif is defined by the molecular recognition of this region in a cell.

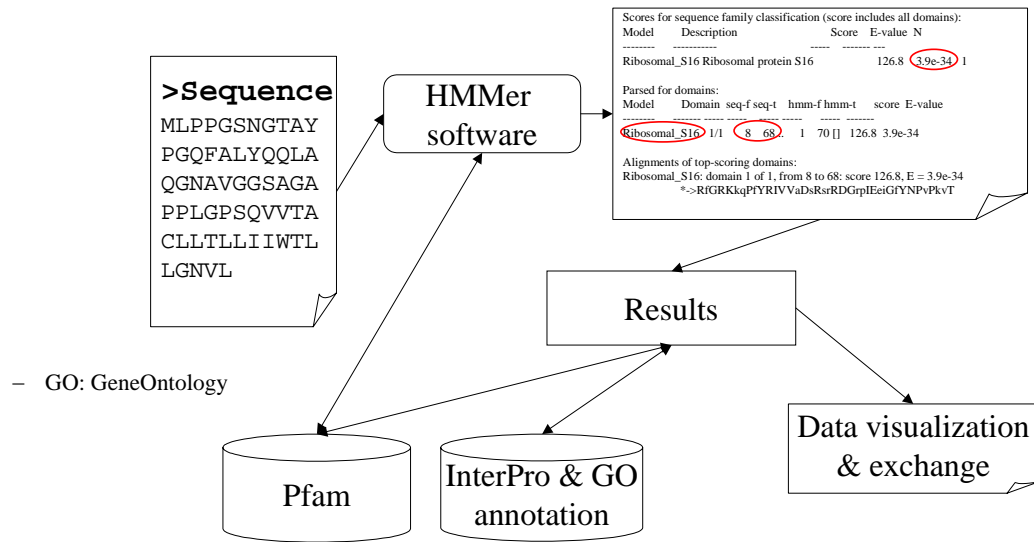


InterPro: Datengewinnung



InterPro: Datengewinnung (2)

- Beispiel: Pfam (Protein families database of alignments and HMMs; Multiple sequence alignments and hidden Markov models of common protein domains)



InterPro: Datengewinnung (3)

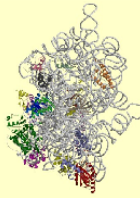


Figure 1: 1pns

Ribosome

Crystal structure of a streptomycin dependent ribosome from *e. coli*, 30s subunit of 70s ribosome. this file, 1pns, contains the 30s subunit, two trnas, and one mma molecule. the 50s ribosomal subunit is in file 1pnu.

Key:

Domain	Chain	Start Residue	End Residue
Ribosomal_S2	B	10	226
Ribosomal_S3_C	C	119	202
Ribosomal_S3_N	C	2	62
KH	C	65	112
Ribosomal_S4	D	3	98
S4	D	99	146
Ribosomal_S5	E	5	71
Ribosomal_S5_C	E	80	153
Ribosomal_S6	F	2	93

Accession number: PF00886

Ribosomal protein S16

[Add Annotation](#)

This family forms **structural complexes** with other Pfam families, to view them click [here](#)

INTERPRO description (entry IPR000307)

Ribosomes are the particles that catalyze mRNA-directed protein synthesis in all organisms. The codons of the mRNA are exposed on the ribosome to allow tRNA binding. This leads to the incorporation of amino acids into the growing polypeptide chain in accordance with the genetic information. Incoming amino acid monomers enter the ribosomal A site in the form of aminoacyl-tRNAs complexed with elongation factor Tu (EF-Tu) and GTP. The growing polypeptide chain, situated in the P site as peptidyl-tRNA, is then transferred to aminoacyl-tRNA and the new peptidyl-tRNA, extended by one residue, is translocated to the P site with the aid of the elongation factor G (EF-G) and GTP as the deacylated tRNA is released from the ribosome through one or more exit sites ([MEDLINE:21198157](#)), ([MEDLINE:21185928](#)). About 2/3 of the mass of the ribosome consists of RNA and 1/3 of protein. The proteins are named in accordance with the subunit of the ribosome which they belong to - the small (S1 to S31) and the large (L1 to L44). Usually they decorate the rRNA cores of the subunits.

Many of ribosomal proteins, particularly those of the large subunit, are composed of a globular, surfaced-exposed domain with long finger-like projections that extend into the rRNA core to stabilize its structure. Most of the proteins interact with multiple RNA elements, often from different domains. In the large subunit, about 1/3 of the 23S rRNA nucleotides are at least in van der Waal's contact with protein, and L22 interacts with all six domains of the 23S rRNA. Proteins S4 and S7, which initiate assembly of the 16S rRNA, are located at junctions of five and four RNA helices, respectively. In this way proteins serve to organize and stabilize the rRNA tertiary structure. While the crucial activities of decoding and peptide transfer are RNA based, proteins play an active role in functions that may have evolved to streamline the process of protein synthesis. In addition to their function in the ribosome, many ribosomal proteins have some function 'outside' the ribosome ([MEDLINE:21185928](#)), ([MEDLINE:20566949](#)).

Ribosomal protein S16 is one of the proteins from the small ribosomal subunit. It belongs to a family of ribosomal proteins which, on the basis of sequence similarities PUB00005070, groups:

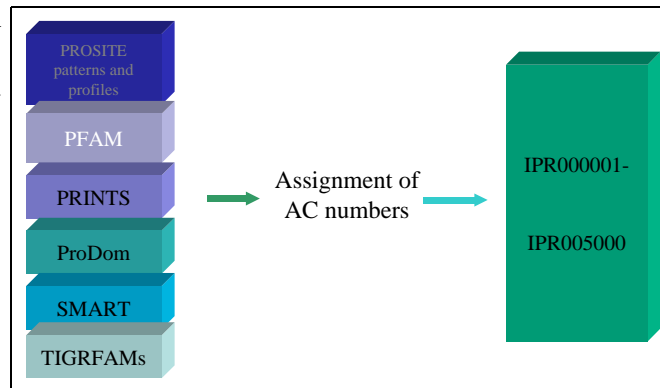
- Eubacterial S16.
- Algal and plant chloroplast S16.
- Cyanelle S16.
- *Neurospora crassa* mitochondrial S24 (cyt-21).

S16 proteins have about 100 amino-acid residues.

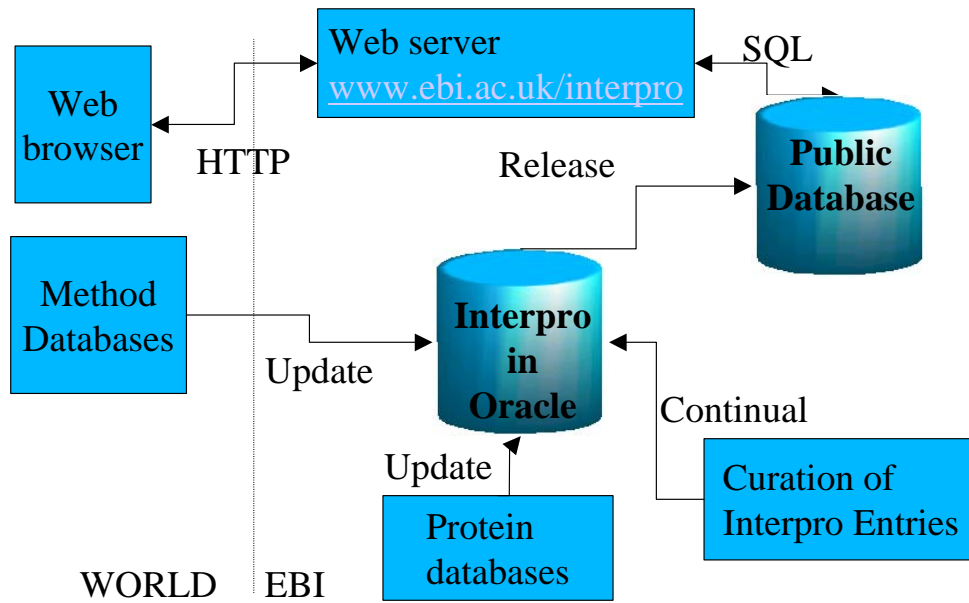


InterPro: Datengewinnung (4)

- Quellen bleiben eigenständig
- Regelmäßige Aktualisierungen
- Jeder Entry der Quelle wird Entry in InterPro
 - Aber: Zusammenhänge bleiben erhalten (Verifizierbarkeit!)
- Größtenteils manuelles Verfahren
 - Redundante Einträge
 - Sub/Superdomänen-Relationen zwischen Entries



InterPro: Architektur

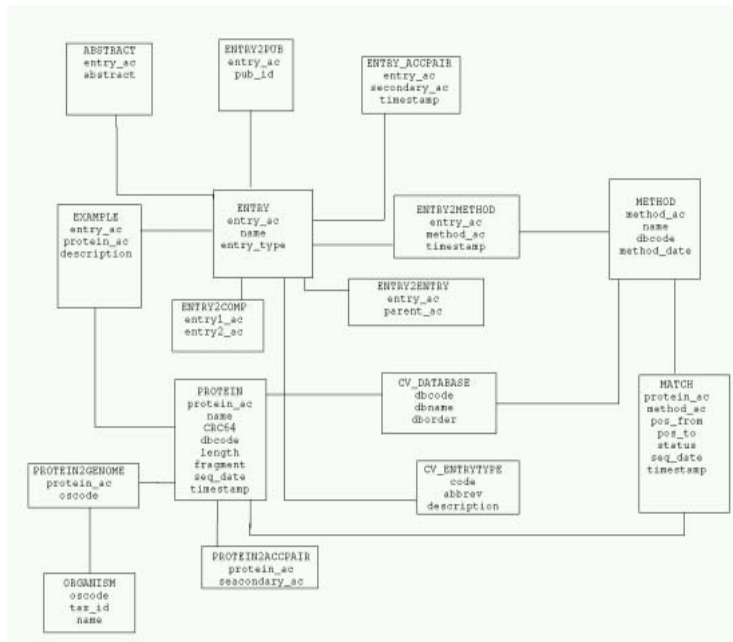


InterPro: Datenarten

- Basic Data
 - InterPro Entries (ENTRY)
 - Proteins (PROTEIN)
 - Methods (METHOD)
 - Annotation
 - Abstracts (ABSTRACT)
 - Publications (PUB, AUTHOR, BOOK ...)
 - Examples (EXAMPLE)
 - Cross References
 - Hierarchical Relationships (ENTRY2ENTRY, ENTRY2COMP)
 - Methods Mapping (ENTRY2METHOD)
 - Matches (MATCH)
 - Supporting Data
 - Secondary AC numbers (ENTRY_ACCPAIR)
 - Proteome Analysis Data (PROTEIN2GENOME, ORGANISM)
 - Audit Tables
- Methods - match domains and families
 - Eg: PF00001: 7 transmembrane receptor (rhodopsin family)
 - Proteins
 - Eg: O00155: PROBABLE G PROTEIN-COUPLED RECEPTOR GPR25.
 - Matches – precomputed
 - Eg: PF00001: matches O00155 at amino acids 56-306
 - Entries – logical groupings of Methods
 - Eg: IPR000276: Rhodopsin-like GPCR superfamily



InterPro: Oracle-Schema (Auszug)



■ Insgesamt 41 Tabellen (ohne Beziehungstabellen)



InterPro: Oracle-Anfragen

- How many short (<100 aa) Drosophila proteins have C2H2 zinc fingers (IPR000822) ?
 - select count(p.protein_ac) from protein p, entry2method e, match m
 - where m.protein_ac = p.protein_ac
 - and m.method_ac = e.method_ac
 - and p.len < 100
 - and e.entry_ac = 'IPR000822'
 - and p.protein_ac in (select protein_ac from protein2genome
 - where oscode = 'DROME');
- 10
 - Which InterPro entries containing only Pfam signatures are common for Human and *V. cholerae* proteomes (SwissProt proteins only) ?
 - 1 select e.entry_ac, count(e.entry_ac) from entry2method e,protein2genome g, match m, protein p
 - where g.oscode = 'HUMAN' and m.protein_ac = g.protein_ac
 - and g.protein_ac = p.protein_ac and p.dbcode = 'S' and m.method_ac = e.method_ac
 - and m.dbcode = 'H' having count(e.entry_ac) = 1 group by e.entry_ac
 - intersect
 - select e.entry_ac, count(e.entry_ac) from entry2method e,protein2genome g, match m, protein p
 - where g.oscode = 'VIBCH' and m.protein_ac = g.protein_ac
 - and g.protein_ac = p.protein_ac and p.dbcode = 'S' and m.method_ac = e.method_ac
 - and m.dbcode = 'H' having count(e.entry_ac) = 1 group by e.entry_ac
 - | | |
|-------------|---|
| • IPR000206 | 1 |
| • IPR000307 | 1 |
| • IPR000398 | 1 |
| • ... | |
| • ... | |
| • IPR002930 | 1 |
| • IPR003156 | 1 |
 - (28 entries)



InterPro: Web Interface

InterProScan by Evgenii.Zdobnov@ebi.ac.uk - Microsoft Internet Explorer

File Edit View Favorites Tools Help

Inter

InterProScan [InterPro Home](#) [Tools@EBI](#) [EBI Home](#)
[README](#) [FAQs](#) [ACC query](#)

Enter or paste a **protein Sequence** in **FASTA** format:

Upload a file

HMMProfam HMMTigr ProfileScan FingerPRINTScan HMMSmart BlastProDom
 Coil Seg IMHMM

e-mail

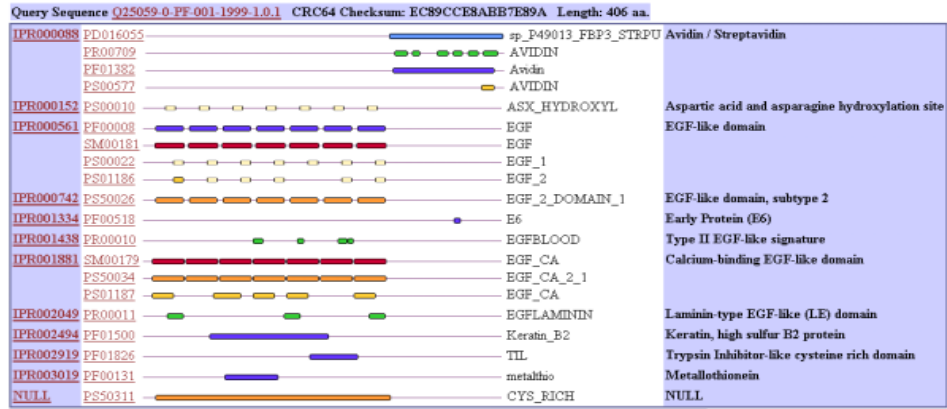
interactive job do InterPro lookup show GO terms report scores

in case of translation required use the following transcript length treshold
and [codon table](#)

Internet



InterPro: Web Interface (2)



InterPro	Results of EPrintScan vs. PRINTS	Results of HMMSmart vs. SMART	Results of HMMPfam vs. PFAM-A	Results of ScanProsite vs. PROSITE	Results of ProfileScan vs. PROFILES	Results of BlastProDom vs. PRODOM
IPR000088 Avidin / Streptavidin	PF00709 [283-297]T [303-311]T [331-343]T [349-359]T [366-377]T [383-399]T		PF01382 [281-395]T	PS00577 [381-395]T		PF016055 [277-405]T
IPR000152 Aspartic acid and asparagine hydroxylation site				PS00010 [27-38]T [65-76]T...		

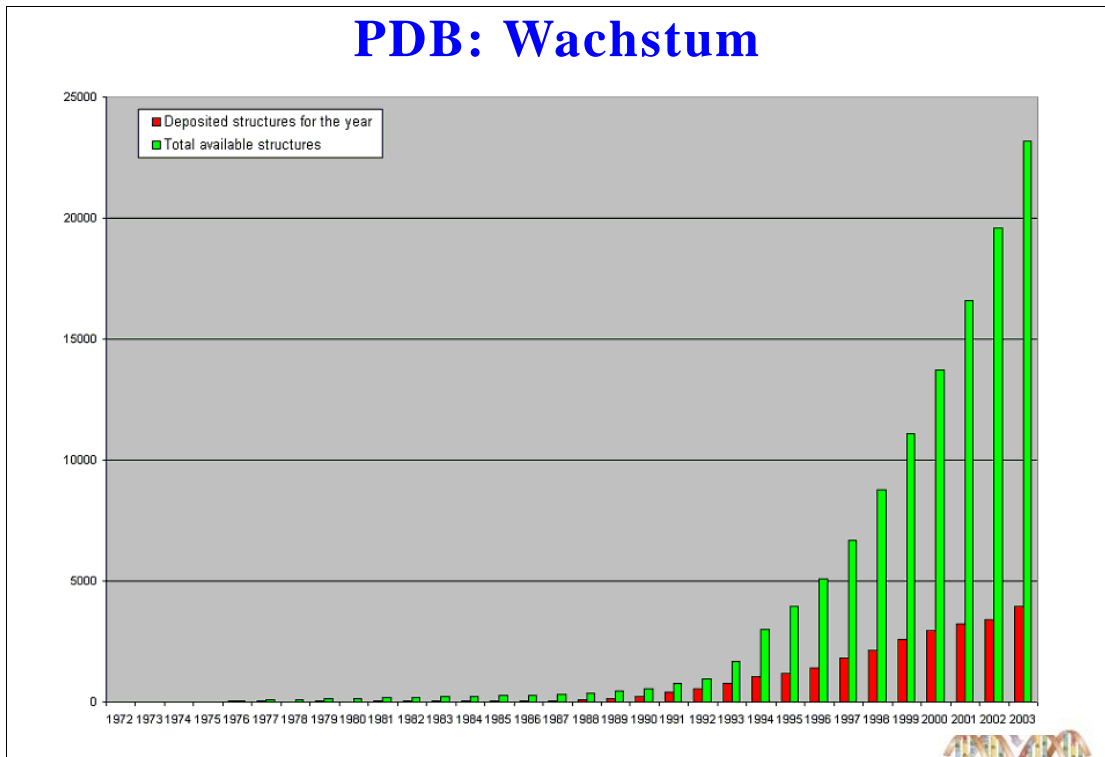


PDB: Protein Data Bank

- Protein-Struktur-Datenbank
- Motivation: Proteine falten sich in komplexe Strukturen, die entscheidend für die Funktion sind
- Strukturaufklärung
 - Röntgenkristallographie (seit 50'er Jahren), Massenspektrometrie, Nuclear Magnetic Resonance (NMR)
- Protein Data Bank
 - Repository aller (bekannten) Protein-3D-Strukturen
 - Seit 1971 in Brookhaven; seit 1999: Rutgers University
- Entry-Based Legacy Format; sehr komplexes 3D-Datenmodell
- Enge Kooperation mit OMG "Specification for Macromolecular Structure, v 1.0" (http://www.omg.org/technology/documents/formal/macro_molecular.htm)



PDB: Wachstum



PDB: Strukturabbildung (2)

- [ENTITY_SOURCE_TYPE](#)
- [ENTITY_SRC_GEN](#)
- [ENTITY_SRC_NAT](#)
- [ENTITY_SRC_SYN](#)
- [ENTRY](#)
- [ENZYME_CLASS](#)
- [ENZYME_CLASS_SYNC](#)
- [ENZYME_STRUCT](#)
- [EXPERIMENTAL_METHOD](#)
- [EXPTL](#)
- [EXPTL_CRYSTAL](#)
- [EXPTL_CRYSTAL_GROWTH](#)
- [EXPTL_CRYSTAL_GROWTH_CONDITIONS](#)
- [GEOM_ANGLE](#)
- [GEOM_BOND](#)
- [GEOM_CONTACT](#)
- [GEOM_TORSION](#)
- [GO_SORT_PDB](#)
- [GO_TERM](#)
- [GO_TERM2TERM](#)
- [GO_TERM_BAYES](#)
- [GO_TERM_PDB](#)
- [GO_TERM_SORT](#)
- [JOURNAL_ABBREVIATION](#)
- [JOURNAL_NAME](#)

Data items in the GEOM_ANGLE category record details about the molecular and crystal angles, as calculated from the contents of the ATOM, CELL, and SYMMETRY data.

Table Name : geom_angle

Note

DDL Code

```
CREATE TABLE GEOM_ANGLE(
  PARTITION_RANGE          CHAR(1)          NOT NULL,
  GEOM_ANGLE_ID            NUMBER(38, 0)      NOT NULL,
  ATOM_SITE_AUTH_ASYM_ID_1 VARCHAR2(10),
  ATOM_SITE_AUTH_ASYM_ID_2 VARCHAR2(10),
  ATOM_SITE_AUTH_ASYM_ID_3 VARCHAR2(10),
  ATOM_SITE_AUTH_ATOM_ID_1 VARCHAR2(10),
  ATOM_SITE_AUTH_ATOM_ID_2 VARCHAR2(10),
  ATOM_SITE_AUTH_ATOM_ID_3 VARCHAR2(10),
  ATOM_SITE_AUTH_COMP_ID_1 VARCHAR2(10),
  ATOM_SITE_AUTH_COMP_ID_2 VARCHAR2(10),
  ATOM_SITE_AUTH_COMP_ID_3 VARCHAR2(10),
  ATOM_SITE_AUTH_SEQ_ID_1  VARCHAR2(10),
  ATOM_SITE_AUTH_SEQ_ID_2  VARCHAR2(10),
  ATOM_SITE_AUTH_SEQ_ID_3  VARCHAR2(10),
  ATOM_SITE_ID_1           VARCHAR2(10)      NOT NULL,
  ATOM_SITE_ID_2           VARCHAR2(10)      NOT NULL,
  ATOM_SITE_ID_3           VARCHAR2(10)      NOT NULL,
  ATOM_SITE_LABEL_ALT_ID_1 VARCHAR2(10),
  ATOM_SITE_LABEL_ALT_ID_2 VARCHAR2(10),
  ATOM_SITE_LABEL_ALT_ID_3 VARCHAR2(10),
  ATOM_SITE_LABEL_ASYM_ID_1 VARCHAR2(10),
  ATOM_SITE_LABEL_ASYM_ID_2 VARCHAR2(10),
  ATOM_SITE_LABEL_ASYM_ID_3 VARCHAR2(10),
  ATOM_SITE_LABEL_ATOM_ID_1 VARCHAR2(10)
```



PDB: Web Interface

View Structure

Interactive 3D Display:


Choose from the following [display options](#) (asymmetric unit only):


- [VRML \(default options\)](#): Interactive immersive ribbon diagram
- [VRML \(custom options, full screen display\)](#): Interactive immersive ribbon or cylinder diagram with ligands
- [Rasmol](#)
- [Swiss-PdbViewer](#)
- [MICE - Molecular Interactive Collaborative Environment](#) (requires Java Plugin)
- [FirstGlance](#) (needs [Chime](#))
- [Protein Explorer](#) (needs [Chime](#))
- [Sting Millennium](#) (needs [Chime](#) and Java)
- Java (simple interactive sequence/structure/property backbone diagram):

[HELP](#)
Download Help
[VRML](#)
[Rasmol](#)
[Swiss-PdbViewer](#)
[Chime](#)
[MICE](#)

Still Images:

Asymmetric Unit **Assumed Biological Molecules**





PDB: Web Interface (2)

The screenshot displays the QuickPDB web interface. At the top, it says "QuickPDB" and provides instructions: "Sequence: drag or click to select residues | 3D: double click to select residue". Below this is a sequence viewer showing three chains of a protein (1CX2:A, 1CX2:B, 1CX2:C, 1CX2:D) with their respective amino acid sequences. The cursor is positioned at 1CX2:D 666. On the left, there is a "Polymers:" list containing 1CX2:A, 1CX2:B, 1CX2:C, and 1CX2:D. Below the list are controls for "Exposure" (set to "Hydrophobic"), a "Stereo" checkbox (unchecked), "Mouse" (set to "Rotate"), and "Color" (set to "Blue"). There are also "Reset" and "Close" buttons. The main area shows a 3D visualization of the protein structure, rendered as a blue and red ribbon. At the bottom, it says "Applet QuickPDB v1.1 (C) 1996-1998 SDSC, by Ilya Shindyalov & Phil Bourne" and "Java Applet Window".

Sequence: drag or click to select residues | 3D: double click to select residue

```
1CX2:A 1  AWPCCMPCCMRGECMSTGFDQYKDCDCTRTGFGYGENCTTPEFLTRIKLLKPTPNTVRYILTMFRGVVNI
1CX2:A 71  VNNIFPLRSLIMKVLTSRVYLIDSPPTVQVHYGKSGEAFSNLSYYTRALPPVADDCTTFMGVKGKEL
1CX2:A 141  FDSKEVLEKULLRRETIPTDQGGSRHMFATFAQHFTHQFFKTDHKKRPPCTRCLEGRGVLDLHMIYGETLDRQ
1CX2:A 211  HKLLLEKDCLEKYQVICGEVYPTVQDTQVENIYPPHIFENLQTAUGCQEVFGLVPLMHVATIQLEDRQ
```

Cursor at 1CX2:D 666

Polymers:

- 1CX2:A
- 1CX2:B
- 1CX2:C
- 1CX2:D

Exposure: Hydrophobic

Stereo

Mouse: Rotate

Color: Blue

Reset Close

Applet QuickPDB v1.1 (C) 1996-1998 SDSC, by Ilya Shindyalov & Phil Bourne

Java Applet Window



Weitere Protein-Datenbanken



PIR Protein Information Resource

About PIR

Databases

Search and Retrieval

Download

Support

AN INTEGRATED PUBLIC RESOURCE OF PROTEIN INFORMATICS TO SUPPORT
GENOMIC AND PROTEOMIC RESEARCH AND SCIENTIFIC DISCOVERY

PIR produces the **Protein Sequence Database (PSD)** of functionally annotated protein sequences, which grew out of the *Atlas of Protein Sequence and Structure* (1965-1978) edited by Margaret Dayhoff and has been incorporated into an integrated knowledge base system of value-added databases and analytical tools.

ProClass, a central point for exploration of protein information, provides summary descriptions of protein family, function and structure for PIR-PSD, Swiss-Prot, and TrEMBL sequences, with links to over 50 biological databases. [Release 2.35, 24-Nov-2003, contains 1,169,177 entries.](#)

PIR-NREF, a comprehensive database for sequence searching and protein identification, contains non-redundant protein sequences from PIR-PSD, Swiss-Prot, TrEMBL, RefSeq, GenPept, and PDB. [Release 1.35, 24-Nov-2003, contains 1,397,398 entries.](#)

PIR has recently joined forces with [EBI](#) (European Bioinformatics Institute) and [SIB](#) (Swiss Institute of Bioinformatics) to establish the [UniProt](#) (United Protein Databases), the central resource of protein sequence and function.



PIR News Flash

Press Release: Protein Information Resource Adds
New Tools to Databases

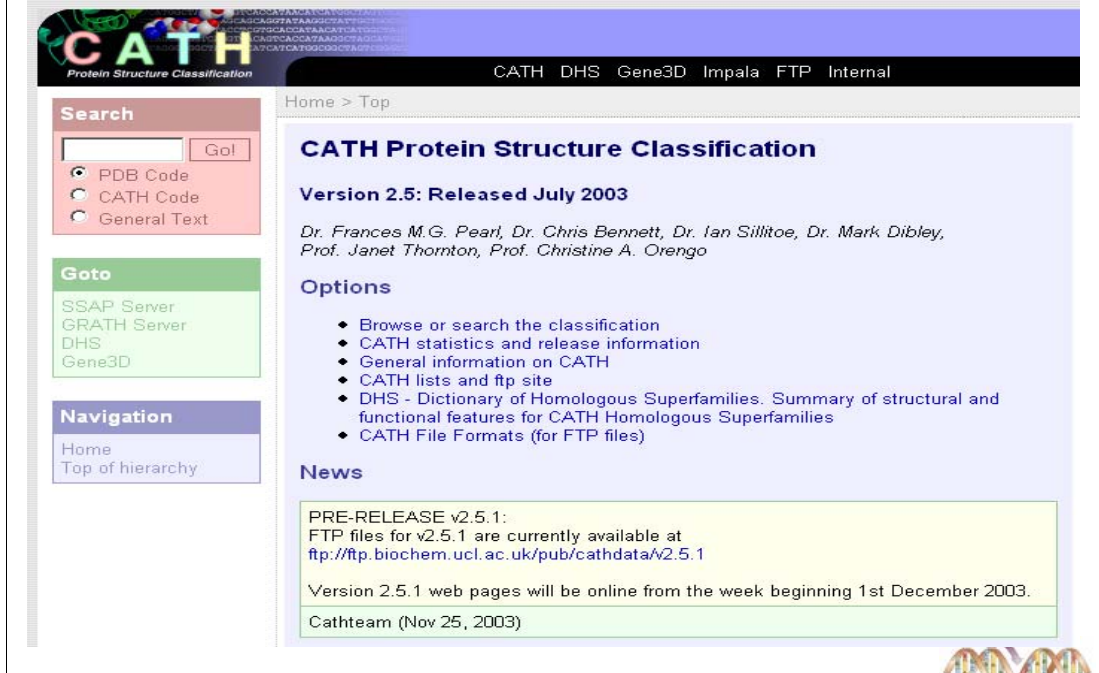
Text Search Protein Databases:

Find an Exact Peptide Match:

Type in a string of single letter amino acid
code (at least 3 letters)



Weitere Protein-Datenbanken (2)



The screenshot shows the CATH Protein Structure Classification website. At the top, there is a navigation bar with links for CATH, DHS, Gene3D, Impala, FTP, and Internal. Below this is a search bar with a 'Go!' button and radio buttons for 'PDB Code', 'CATH Code', and 'General Text'. To the left, there are sections for 'Goto' (SSAP Server, GRATH Server, DHS, Gene3D) and 'Navigation' (Home, Top of hierarchy). The main content area is titled 'CATH Protein Structure Classification' and 'Version 2.5: Released July 2003'. It lists the authors: Dr. Frances M.G. Pearl, Dr. Chris Bennett, Dr. Ian Sillitoe, Dr. Mark Dibley, Prof. Janet Thornton, and Prof. Christine A. Orengo. Under 'Options', there is a bulleted list of links: 'Browse or search the classification', 'CATH statistics and release information', 'General information on CATH', 'CATH lists and ftp site', 'DHS - Dictionary of Homologous Superfamilies. Summary of structural and functional features for CATH Homologous Superfamilies', and 'CATH File Formats (for FTP files)'. A 'News' section contains a yellow box with the text: 'PRE-RELEASE v2.5.1: FTP files for v2.5.1 are currently available at ftp://ftp.biochem.ucl.ac.uk/pub/cathdata/v2.5.1'. Below this is a green box with the text: 'Version 2.5.1 web pages will be online from the week beginning 1st December 2003.' and 'Cathteam (Nov 25, 2003)'. A small DNA double helix icon is located in the bottom right corner of the page.

Weitere Protein-Datenbanken (3)

■ UniProt

- Beinhaltet PIR, Swiss-Prot und TrEMBL
- Ablösung einer langen Parallelentwicklung
- Erster Release noch nicht verfügbar

■ OWL

- Nicht-redundante Sammlung von Proteinsequenzen
- Enthält: Swiss-Prot, PIR, GenBank

■ ... und viele mehr



Zusammenfassung

- Motivation und historische Entwicklung
- Proteomics
 - Datengewinnung
 - PEDRo-Projekt
- Protein-Datenbanken
 - Sequenz-Datenbanken (Swiss-Prot)
 - Domain/Familien-Datenbanken (InterPro)
 - Struktur-Datenbanken (PDB)
 - Weitere Protein-Datenbanken

