

Public Databases in Health and Life Sciences

“the potential to translate big data into big discovery”



Academic year 2023-2024

Public Databases in Health and Life Sciences



Topic 4. Secondary and derived protein databases.

- Databases containing data derived from the analysis or treatment of primary data.
- Sequence based prediction of protein function.
- Database of protein motif, domains, families and functional sites.
- Conserved domain databases and domain classification.
- Databases for protein structural domains.

Practical sessions #4 and #5

Primary vs. Secondary Databases

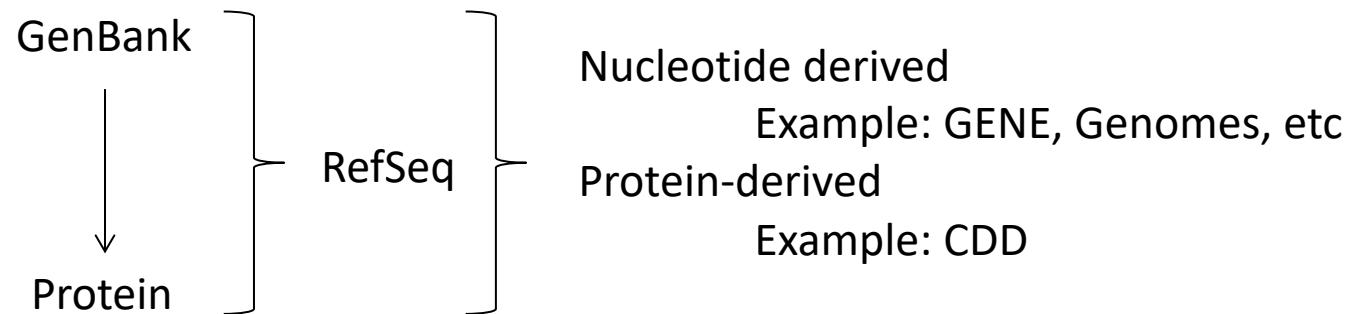
Primary db

- Contains experimentally-derived data such as nucleotide sequences, protein sequences or macromolecular structures.
- Experimental results are submitted directly into the database by researchers, and the data are essentially archival in nature.

Secondary/Derived db

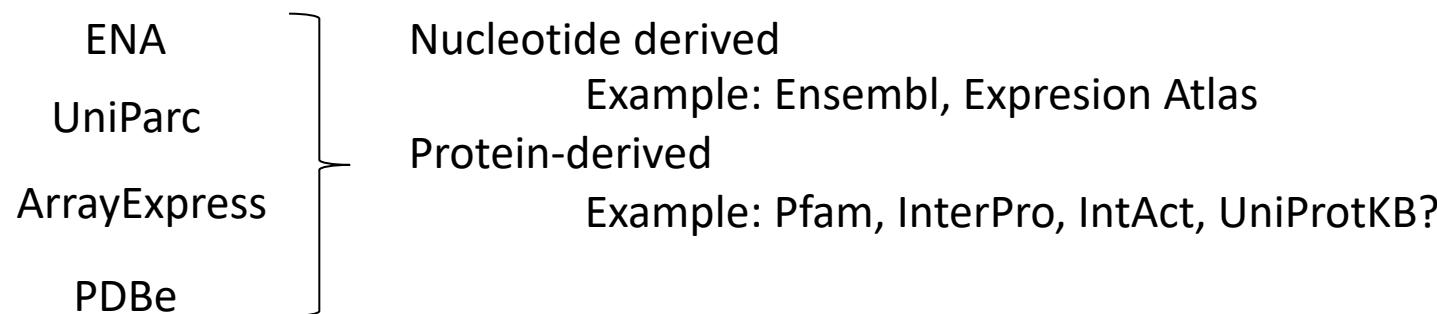
- Data derived from the analysis or treatment of primary data.

NCBI: Derivative Databases



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

EMBL-EBI: Derivative Databases

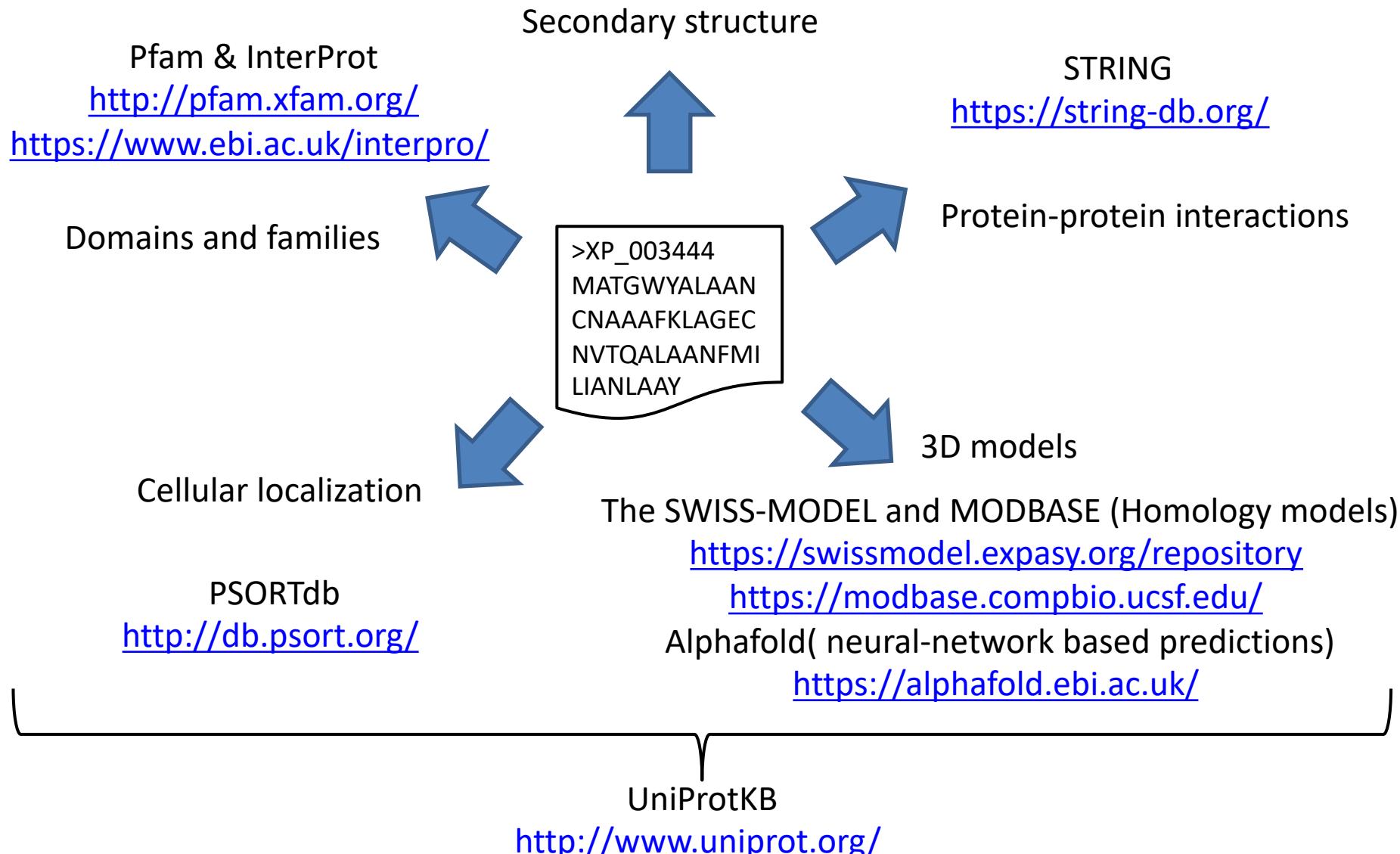


<https://www.ebi.ac.uk/>

Determining protein function

- Experimentally determining function (time consuming and expensive)
 - Functional assays
 - 3D structure
 - Co-localization
 - Protein-protein interaction (experimental approaches)
 - Microarray experiments
- **Predictive methods**
 - Transfer of function based on homology
 - BLAST
 - Based on the amino acid sequence (predicted protein secondary structure, transmembrane helices, subcellular localization, and posttranslational modifications)
 - 3D modeling
 - Motif and Domains
- Sequence Based Predictions at <http://www.expasy.org/tools/>
- Combination of all above

Protein function: predictive methods and associated databases



Protein classification

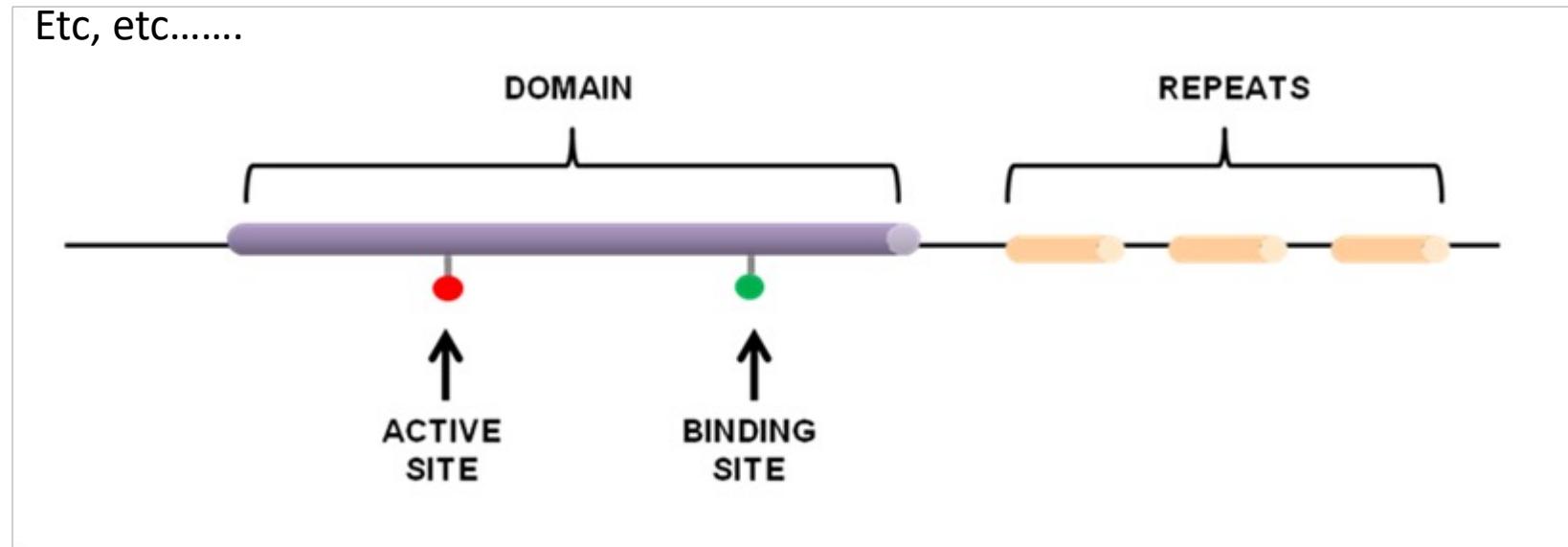
Proteins can be classified into different groups based on:

- the **SEQUENCE FEATURES** they possess
 - Groups of amino acids that confer certain characteristics to a protein
- the **DOMAINS** and **MOTIFS** they contain
 - Distinct functional and/or structural units in a protein.
- the **FAMILIES** to which they belong
 - A group of proteins that share a common evolutionary origin.

What are sequence features?

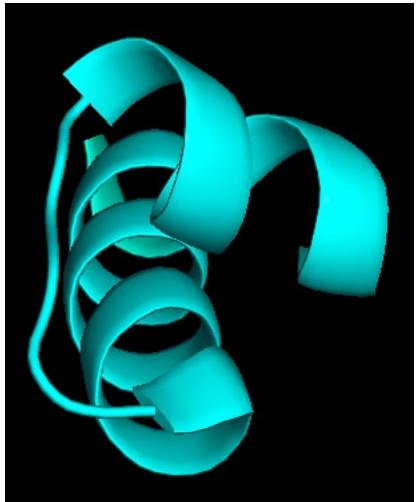
Sequence features are groups of amino acids that confer certain characteristics upon a protein and may be important for its overall function. Such features include:

- Active sites
- Binding sites
- Post-translational modification sites
- Short amino acid repeats
- Signal peptides
- Etc, etc.....

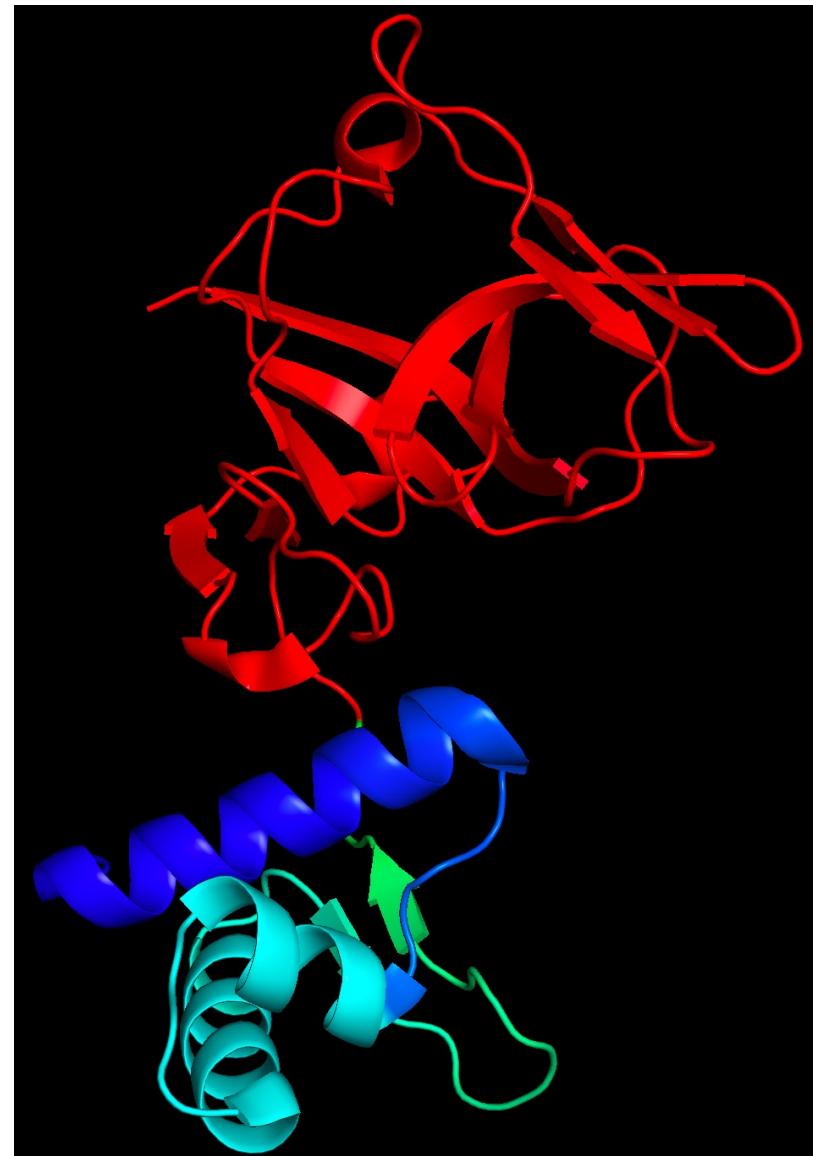
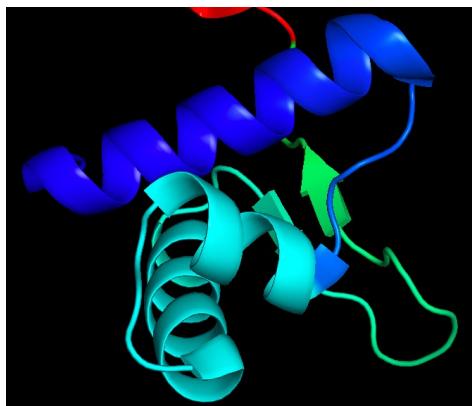


Motifs (AKA: fold) and Domains are units within the protein

Motif



Domains



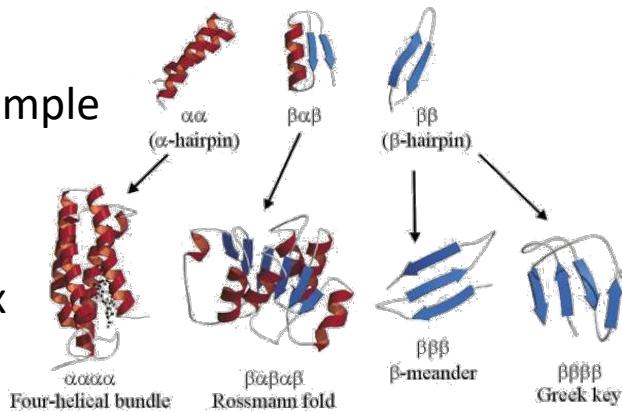
Motif AKA Fold (and rarely Supersecondary Structure):

Motif or fold is a recognizable folding pattern involving two or more elements of secondary structure and the connection(s) between them.

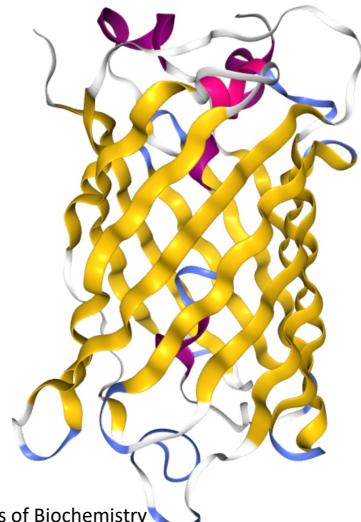
Generally speaking, we use the word “motif” when we talk about simple elements and the word fold for more complex elements

<https://qph.fs.quoracdn.net/main-qimg-028bcdcfaba956d3df3cea5316688054-c>

They can be very simple



A bit more complex

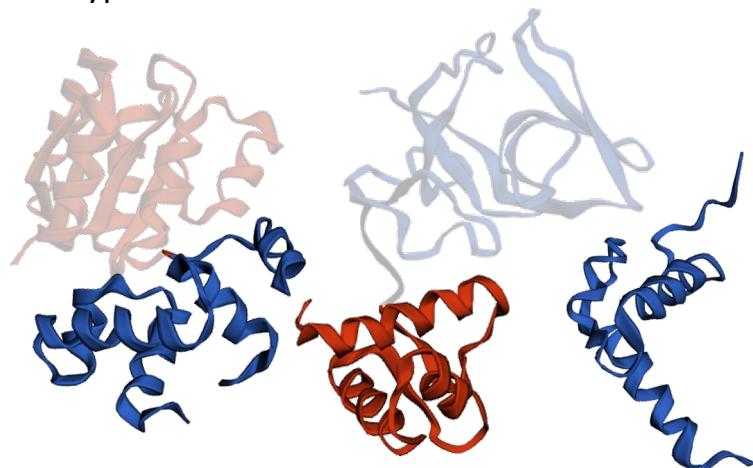


Much more complex

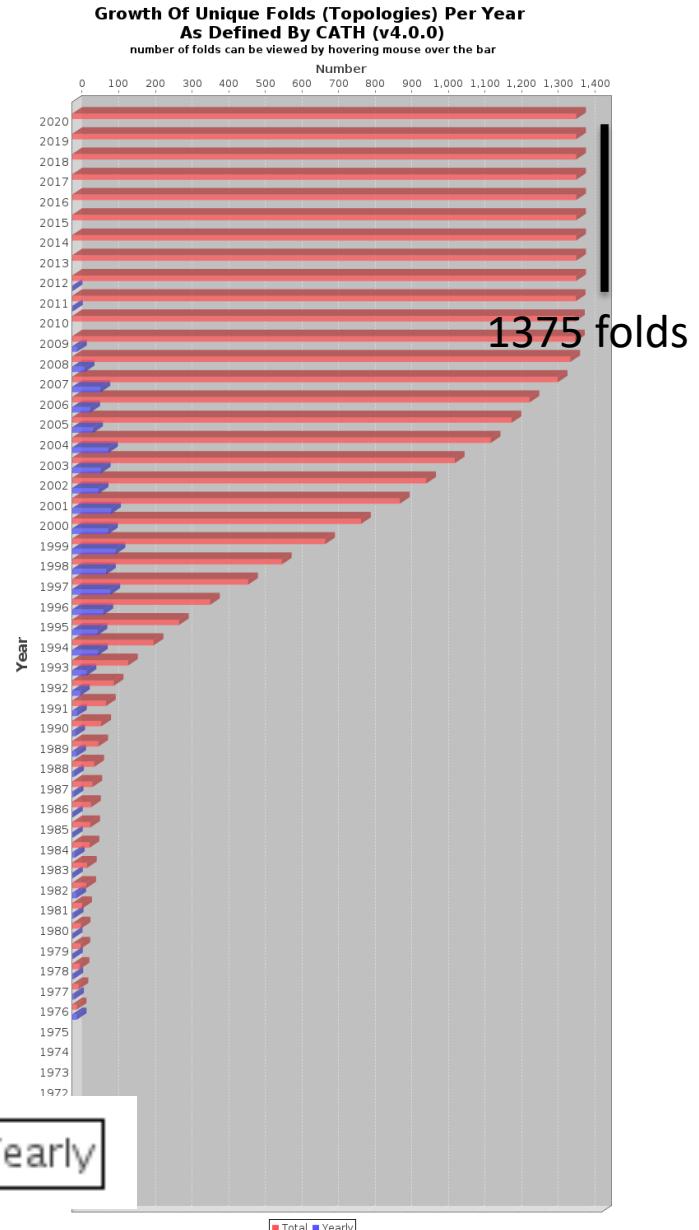
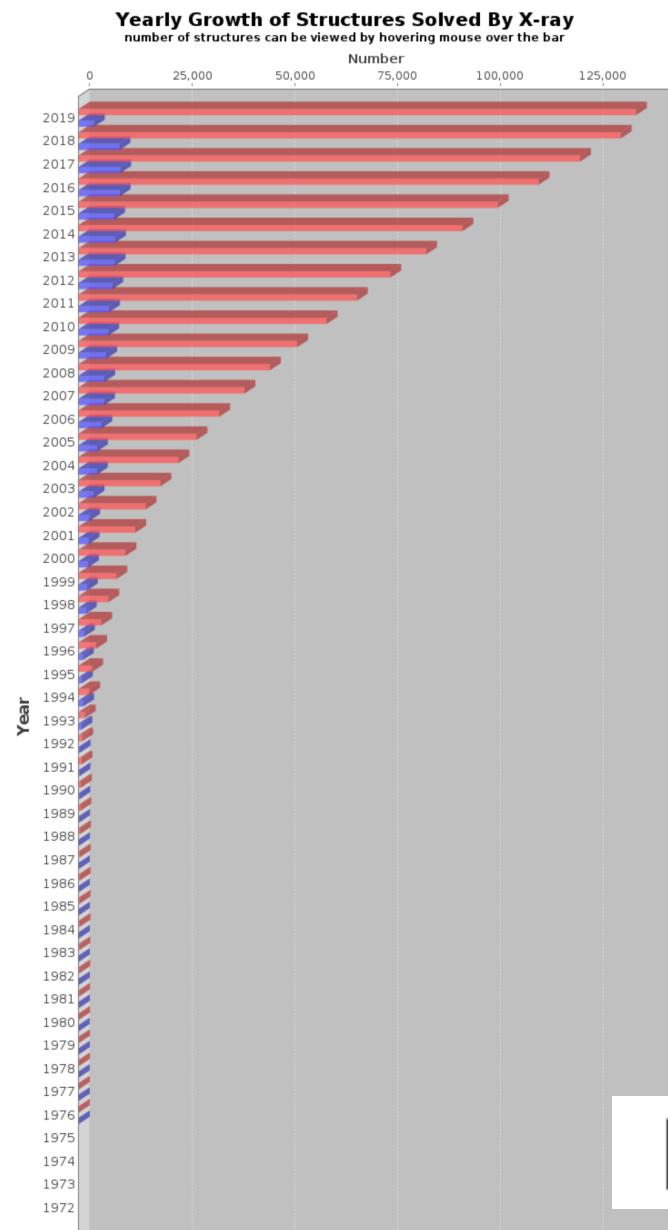
A motif may or may not be independently stable and it may or may not be related to a specific function.

A motif is **not** a hierarchical structural element falling between secondary and tertiary structure.

The same motif can have several types (which do not necessarily coincide with their classification in Databases). See below three types of helix-turn-helix fold



The number of folds is finite



■ Total ■ Yearly

■ Total ■ Yearly

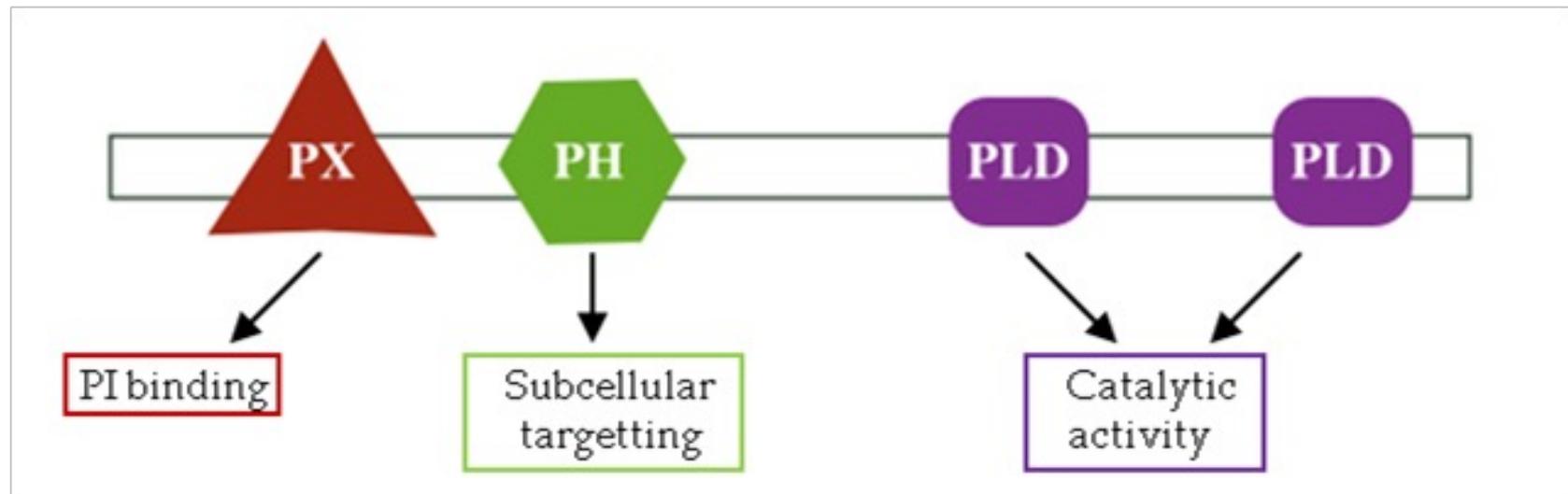
1375 folds

Motifs (AKA: fold) and Domains are units within the protein

- Motifs (AKA: “fold”) are used to describe what can be called a supersecondary structure of the protein. This is the combination of alpha-helices and beta-structures in the protein. It may not necessarily have a biological functions nor be independently stable.
- Domains are distinct functional and/or structural units in a protein. Usually, they are responsible for a particular function or interaction, contributing to the overall roll of the protein. They can fold by themselves and are stable once they are folded. They are supposed to be individual proteins in origin that while evolving they have joined others and are expressed together in order to be more efficient.
- Domains are a sequence of amino acids that are conserved between many different organisms.
- They are often used to hypothesize the overall function of uncharacterized proteins because the sequence of amino acids usually carries a similar function.
- They are especially important in classifying protein families.

Domains are distinct functional and/or structural units in a protein

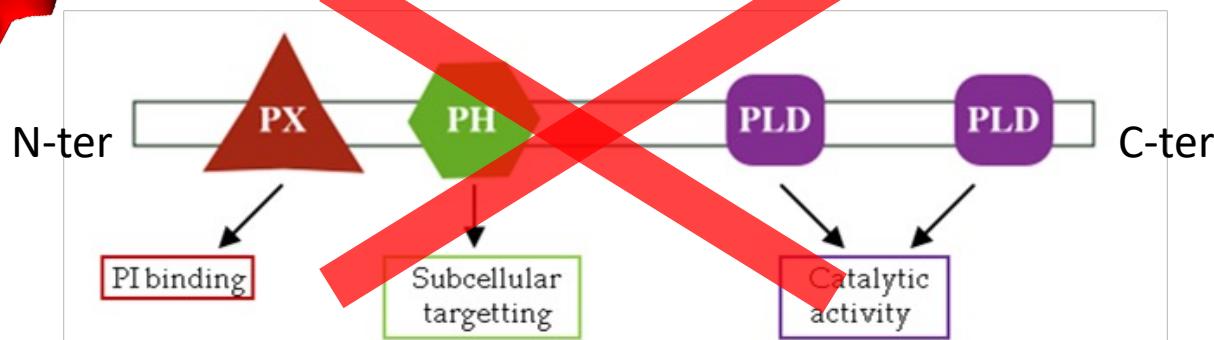
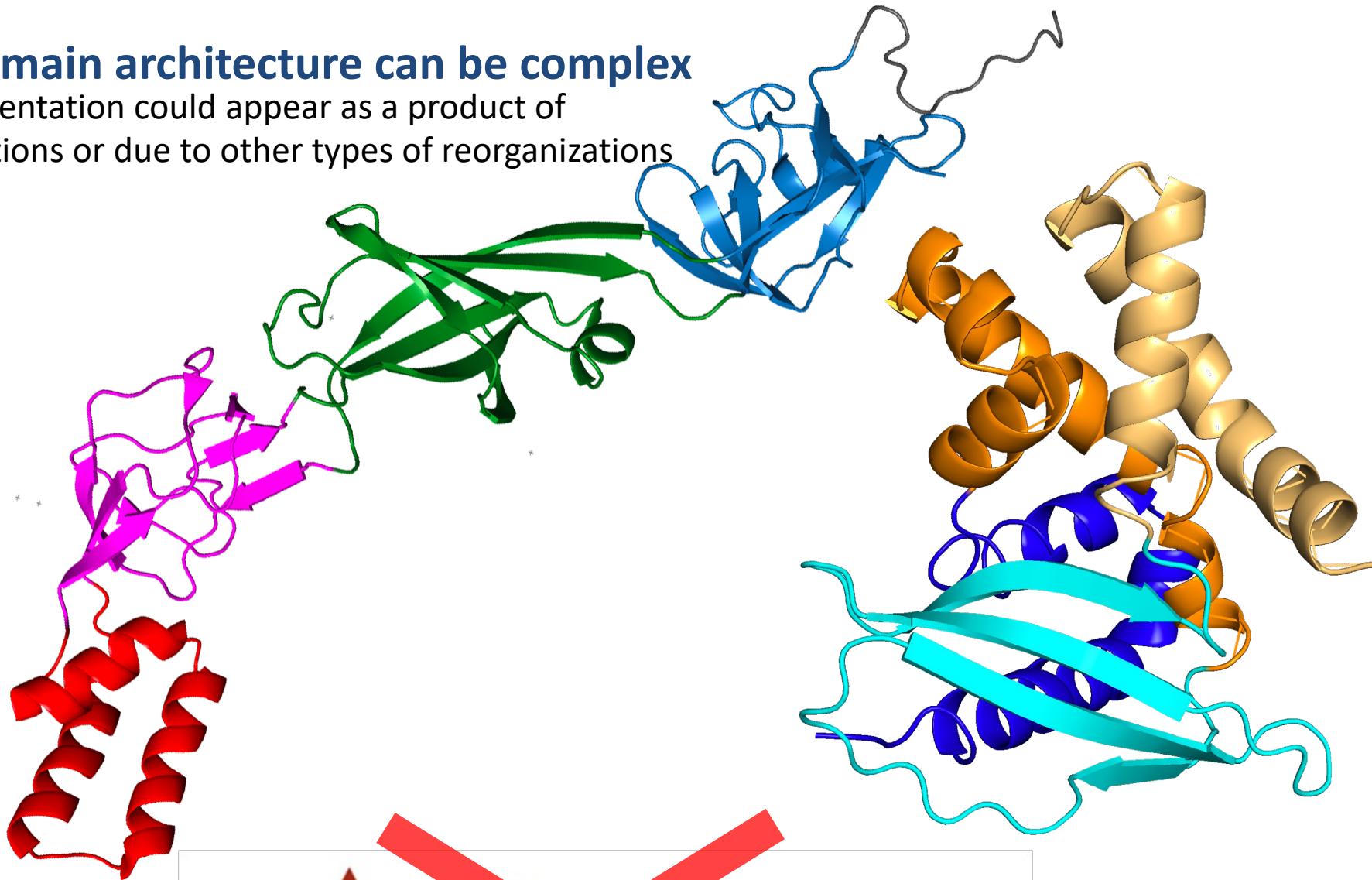
Proteins can be composed of several domains.



Phospholipase D1 Domains

Domain architecture can be complex

Segmentation could appear as a product of insertions or due to other types of reorganizations



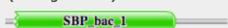
Sequence vs structural domains:

Periplasmic oligogalacturonide binding protein from *Yersinia enterocolitica*

Sequence search results

Show the detailed description of this results page.

We found 1 Pfam-A match to your search sequence (all significant)



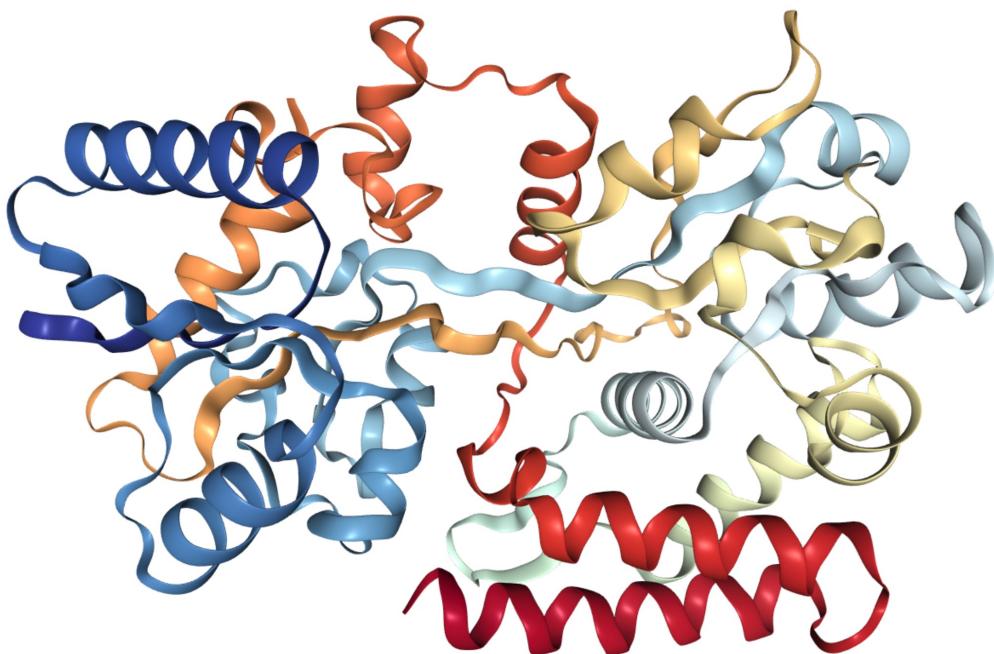
Show the search options and sequence that you submitted.

Return to the search form to look for Pfam domains on a new sequence.

Significant Pfam-A Matches

Show or hide all alignments.

Family	Description	Entry type	Clan	Envelope		Alignment		HMM		HMM length	Bit score	E-value	Predicted active sites	Show/hide alignment
				Start	End	Start	End	From	To					
SBP_bac_1	Bacterial extracellular solute-binding p...	Family	CL0177	19	314	25	311	13	312	315	87.5	1.5e-24	n/a	Show



PDB Information

PDB

Method

Host Organism

Gene Source

Primary Citation

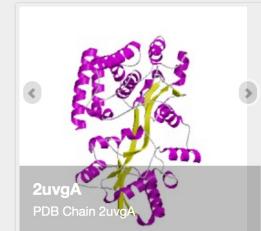
Header Sugar-Binding Protein

Released 2007-03-10

Resolution 2.200

CATH Insert Date 24 May, 2007

PDB Images (4)



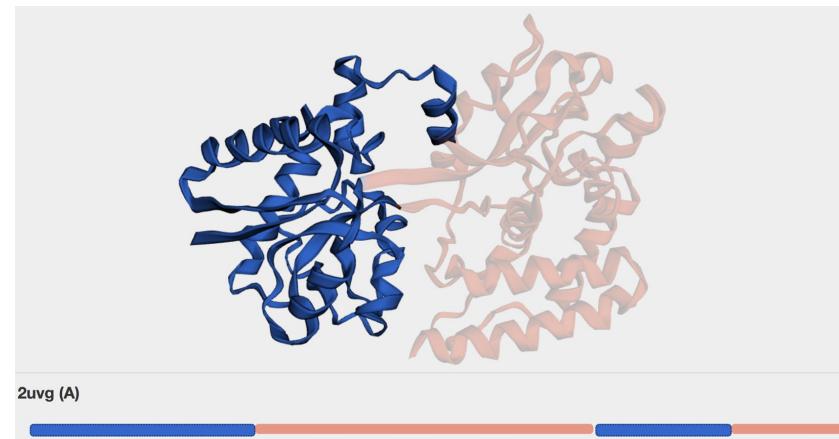
PDB Prints

PDB Chains (1)

Chain ID	Date inserted into CATH	CATH Status
A	24 May, 2007	Chopped ⓘ

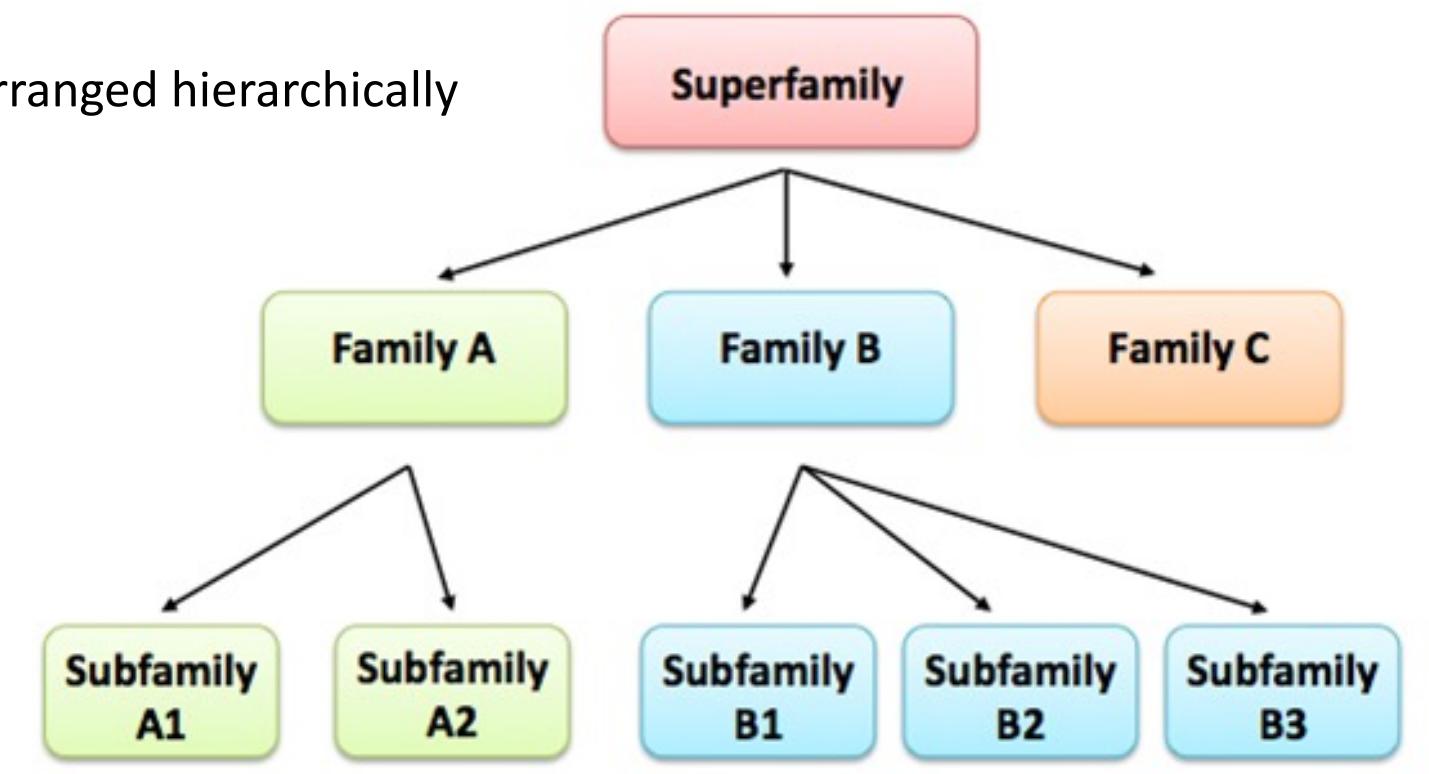
CATH Domains (2)

Domain ID	Date inserted into CATH	Superfamily	CATH Status
2uvgA01	09 Sep, 2013	3.40.190.10	Assigned ⓘ
2uvgA02	09 Sep, 2013	3.40.190.10	Assigned ⓘ



Protein families

- A group of proteins that share a common evolutionary origin.
- Proteins are sometimes assigned to families by virtue of the domain(s) they contain.
- They are arranged hierarchically



Domain Classification

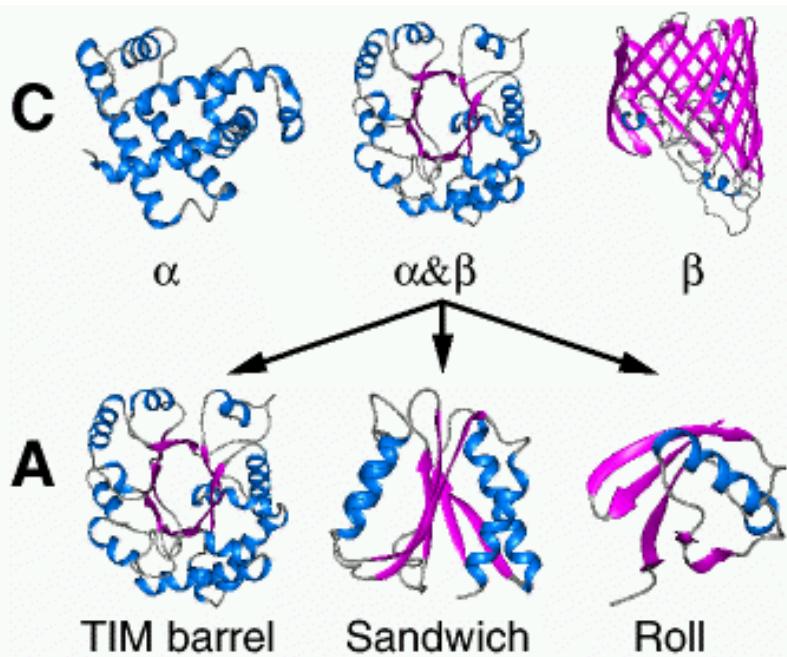
Sequence based (alignments)

Families and
superfamilies

Structure based

The arrange of
secondary structures

Q5E940_BOVIN	MPREDRATWKSNEYLKIIILDDYPKCFIVGADNVGSKMOIIRMSLRGK-AVVLMGKNTMMRKAIIGHLENN--PALE	76
RLAO_HUMAN	MPREDRATWKSNEYLKIIILDDYPKCFIVGADNVGSKMOIIRMSLRGK-AVVLMGKNTMMRKAIIGHLENN--PALE	76
RLAO_MOUSE	MPREDRATWKSNEYLKIIILDDYPKCFIVGADNVGSKMOIIRMSLRGK-AVVLMGKNTMMRKAIIGHLENN--PALE	76
RLAO_RAT	MPREDRATWKSNEYLKIIILDDYPKCFIVGADNVGSKMOIIRMSLRGK-AVVLMGKNTMMRKAIIGHLENN--PALE	76
RLAO_CHICK	MPREDRATWKSNEYMKIIILDDYPKCFIVGADNVGSKMOIIRMSLRGK-AVVLMGKNTMMRKAIIGHLENN--PALE	76
RLAO_RANSY	MPREDRATWKSNEYMKIIILDDYPKCFIVGADNVGSKMOIIRMSLRGK-AVVLMGKNTMMRKAIIGHLENN--PALE	76
Q7ZUG3_BRARE	MPREDRATWKSNEYLKIIILDDYPKCFIVGADNVGSKMOIIRSLRGK-AVVLMGKNTMMRKAIIGHLENN--PALE	76
RLAO_ICTPU	MPREDRATWKSNEYLKIIILDDYPKCFIVGADNVGSKMOIIRSLRGK-AVVLMGKNTMMRKAIIGHLENN--PALE	76
RLAO_DROME	MVRENKRAWKQYEIKVVEFLDFEEPKCFIVGADNVGSKMOIIRTSLRGK-AVVLMGKNTMMRKAIIGHLENN--PALE	76
RLAO_DCDC1	MSCAG-SKRKKFLIEKATKLFETTDKMIVAEADFGVSQLOKIRKSTGIGAVLGMKNTMIRKVIRDLADSK--PELD	75
Q54LP0_DCDC1	MSCAG-SKRKKFLIEKATKLFETTDKMIVAEADFGVSQLOKIRKSTGIGAVLGMKNTMIRKVIRDLADSK--PELD	75
RLAO_PLAF8	MAKLSQQKQOMYILEKLSSLIQQNSLLLVHDNVGSQMASVRSLSLRGK-AEILMGKNTIRITALKKN1QAV--PQLE	76
RLAO_SULAC	MIGAVITKKLAKWVSKWLEELTQKLTWKLTKHIIIANIEGPKADLHLERIRKLRGK-ADTVVKN1NLFLAKNAG--YDK	79
RLAO_SULTO	MIRMAVITQEKRKLAWKWIEEVKELELQKREHMIIIANIEGPKADLHLERIRKLRGK-ADTVVKN1NLFLAKNAG--LDVS	80
RLAO_SULSO	MKRKALALQKQRKAWSKLLEEVKELTEIIKNSNTIIILNQLEHFLDNLHGLR-ADTVVKN1NLFLAKNAG--IDLE	80
RLAO_AERPE	MSVVSIVGQMYKRFKEPDEWKTLMRLEFSKRVVLFADLTQFVVVQVRVKKLWKW-PMMWAKRKLILRMKAAGLE--LDDN	86
RLAO_PYRAE	MLMLAIGKRYVYRTROQPKCKVIVSEATELQQPKPYVYFLRILHEYRYJRRY-GVVIKILQPKFLFKIAFTKVTYGC--IPAE	85
RLAO_METAC	MAEERHTHEH1POWKDDEIENIKELQSKVKGFMVWIEGLLATTMOKIRRDLKDW-AVILVBRNTLLEERALNQLC--ETTP	78
RLAO_METMA	MAEERHTHEH1POWKDDEIENIKELQSKVKGFMVWIEGLLATTMOKIRRDLKDW-AVILVBRNTLLEERALNQLC--ESYL	78
RLAO_ARCFU	MAAVAS--PPPEYKVRAVEEIKRMSISKPVVAVLSERFNVPAQDMOKIRREFRGRK-AEIKVVKNTLLERALDALG--GDYL	75
RLAO_METKA	MAVKAKSQPPSGYEPKWAEWKRREEVKELELMDENENVGLVIANLDEIPAPDLOEIRK1LKERD1IIRMRBN1LMLRALEEK1LDER--PELE	88
RLAO_METTH	MAHYAEWKKKEVQEELHDLIKGEVVGIVLAAADIPARADLOKMRQTLRDS-ALIEMKRN1LISIALEAKREL--ENV	74
RLAO_METTL	MITAASEHK1APWPKIEFWNLKEELLKNGQIVALVDMMEVPAROLOEIRDTR-CGTM1KMBRN1LIERRAKEVAAETGNEPEA	82
RLAO_METVA	MIDAKSEHK1APWPKIEFWNLKEELLKNSAN1ALIDDMEVFPVALOEIRDTR-DQML1KMBRN1LIKRAVEEVAETGNEPEA	81
RLAO_METJA	METTKVKAHYAPWPKIEFWNLKEELLKNGQIVALVDMMEVPAROLOEIRDTR-DQML1KMBRN1LIRKRAVEEVAETGNEPEA	81
RLAO_PYRAB	MAHYAEWKKKEVEELANL1KSPVVALVDVSSMPAYPSLSMRRLLIRENGGLLIRVBRN1LIEIAKKAALKEGLGPPEL	77
RLAO_PYRHO	MAHYAEWKKKEVEELANL1KSPVVALVDVSSMPAYPSLSMRRLLIRENGGLLIRVBRN1LIEIAKKAALKEGLGPPEL	77
RLAO_PYRFU	MAHYAEWKKKEVEELANL1KSPVVALVDVSSMPAYPSLSMRRLLIRENGGLLIRVBRN1LIEIAKKAALKEGLGPPEL	77
RLAO_PYRKO	MAHYAEWKKKEVEELANL1KSPVVALVDVSSMPAYPSLSMRRLLIRENGGLLIRVBRN1LIEIAKKAALKEGLGPPEL	76
RLAO_HALMA	MSA1SERKETET1PEWKEQEVDA1VEMIES1ESVGVVNN1IAG1P1R1Q1D1M1R1L1HGT-AEL1VBRNTLLE1A1KRAA1ELGQPELE	79
RLAO_HALVO	MSESEVRQTEV1PQWQKREEVDELVD1IES1ESVGVVNN1IAG1P1R1Q1D1M1R1L1HGT-AE1VBRNTLLE1A1KRAA1ELGQPELE	79
RLAO_HALSA	MSA1EQRTTEEV1PEWKEQEV1ELV1LLET1DSVGVVNN1T1P1K1Q1D1M1R1L1HGT-AE1VBRNTLLE1A1KRAA1ELGQPELE	79
RLAO_THEAC	MKEV1SQQKE1VNE1T1R1KAKS1V1A1D1A1G1R1Q1D1R1G1KRN1G1-1N1K1V1K11L1F1K1A1L1G1D1E1K1S	72
RLAO_THEVO	MRK1N1K1K1E1V1S1E1A1D1T1K1S1A1V1D1K1G1V1R1Q1D1R1M1D1R1K1N1D1K1S1D1E1K1T	72
RLAO_PICTO	MTEP1Q1K1D1F1X1L1E1N1S1K1R1N1S1E1F1K1R1N1S1D1K1S1D1E1K1T	72



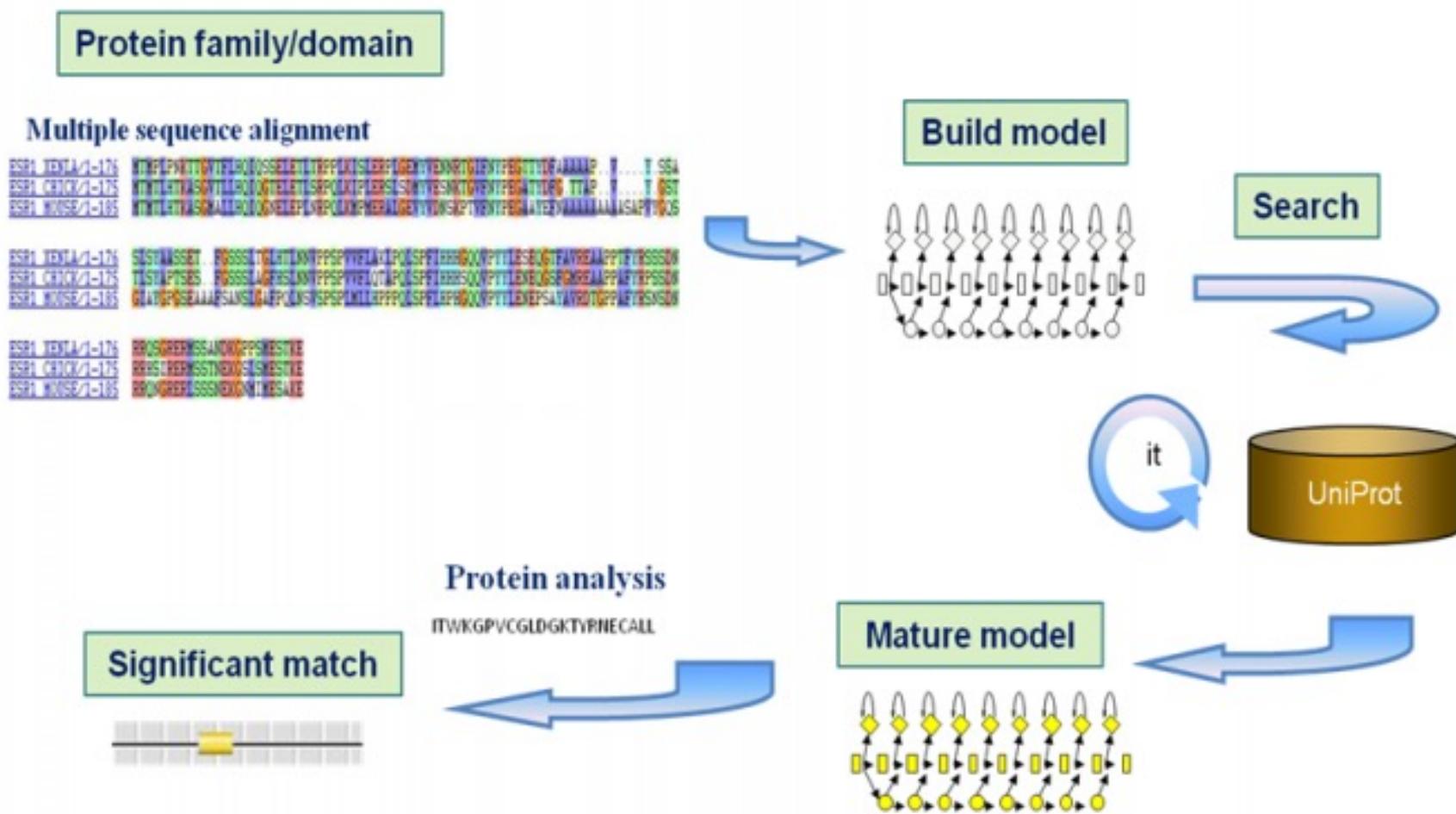
What are protein signatures?

Signature sequences are contiguous patterns of amino acids that are associated with a particular structure or function in proteins.

There are different types of signatures, built using different computational approaches. These include:

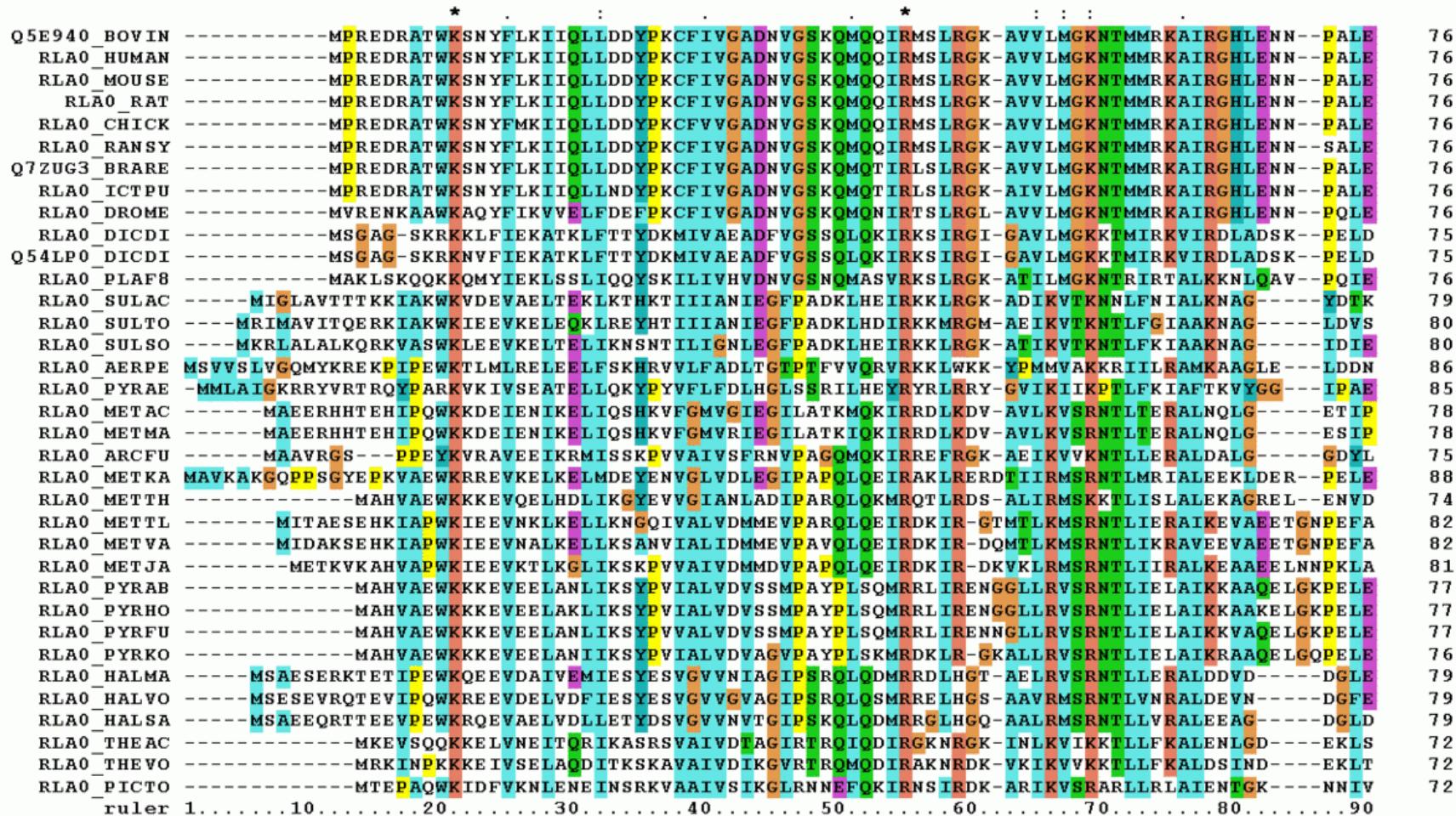
- patterns (motifs)
- profiles (pfam is a profile database)
- fingerprints
- hidden Markov models (HMMs) (a type of profiles)

The process of building a protein signature



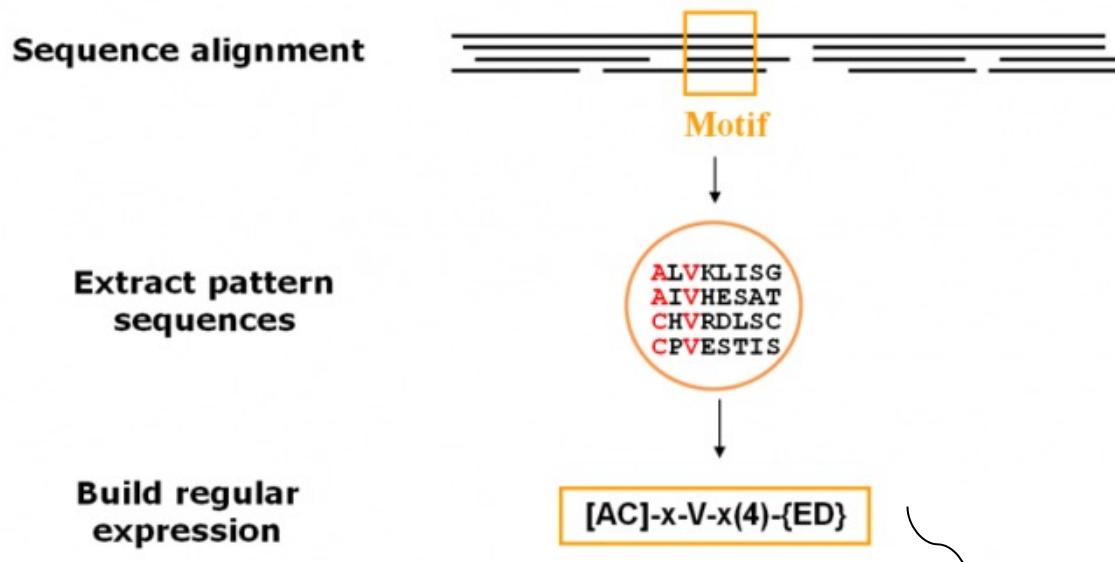
Protein classification: An introduction to EMBL-EBI resources
At <https://www.ebi.ac.uk/training/online/>

Multiple sequence alignments can provide us with valuable information for protein classification since they allow us to identify the (often few) amino acid residues that are conserved in distantly related proteins.



What are patterns (AKA: motifs => BEWARE!!! Same name different concept)?

Many important sequence features, such as binding sites or the active sites of enzymes, consist of only a few amino acids that are essential for protein function. When creating patterns, a conserved motif is used to build a regular expression.



The pattern illustrated here is translated as:
[Ala or Cys]-any-Val-any-any-any-any-{any but Glu or Asp}.

Databases that use patterns:

Prosite

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

Pattern syntax at PROSITE for pattern (motif) searching

Pattern	Explanation
[AC]-x-V-x(4)-{ED}	[Ala or Cys]-any-Val-any-any-any-any-{any but Glu or Asp}
<A-x-[ST](2)-x(0,1)-V	Ala-any-[Ser or Thr]-[Ser or Thr]-(any or none)-Val at the N-terminal of the sequence
<{C}*>	No Cys from the N-terminal to the C-terminal i.e. All sequences that do not contain any Cys.
IIRIFHLRNI	Ile-Ile-Arg-Ils-Phe-His-Leu-Arg-Asn-Ile



ScanProsite Results Viewer

Output format: Graphical view - this view shows ScanProsite results together with ProRule-based predicted intra-domain features [[help](#)].

[include splice variants \(Swiss-Prot\)](#)

Hits for **USERPAT1{P-x(2)-G-E-S-G(2)-[AS]}** motif on all UniProtKB/Swiss-Prot (release 2017_09 of 27-Sep-17: 555594 entries) database sequences :

found: 11 hits in 11 sequences

What are profiles?

They are built by converting multiple sequence alignments (MSAs) into position-specific scoring systems. Amino acids at each position in the alignment are scored according to the frequency with which they occur.

PWMs are often derived from a set of aligned sequences that are thought to be functionally related and have become an important part of many software tools for computational motif discovery.

A position weight matrix (PWM), also known as a position-specific weight matrix (PSWM) or position-specific scoring matrix (PSSM), is a commonly used representation of motifs (patterns) in biological sequences.

Databases that use profiles to classify proteins:

CDD <https://www.ncbi.nlm.nih.gov/Structure/cdd/cdd.shtml>

HAMAP <http://hamap.expasy.org/>

PROSITE <http://prosite.expasy.org/>

PRODOM <http://prodom.prabi.fr/prodom/current/html/home.php>

http://www.ncbi.nlm.nih.gov/Class/Structure/pssm/pssm_viewer.cgi

P - consensus sequence position C - consensus sequence residue

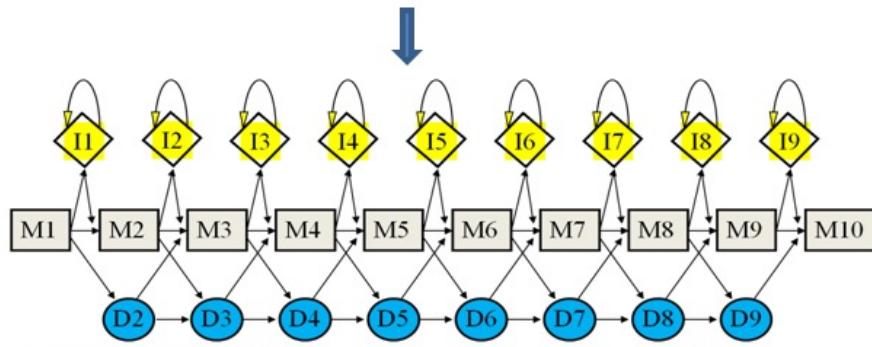
P	C	A	G	I	L	V	M	F	W	P	C	S	T	Y	N	Q	H	K	R	D	E
<u>1</u>	M	-1	-3	1	2	1	6	0	-2	-3	-1	-2	-1	-1	-2	0	-2	-1	-1	-3	-2
<u>2</u>	I	-1	-4	4	1	3	3	-1	-3	-3	-1	-2	-1	-1	-3	-2	-3	-3	-3	-3	-3
<u>3</u>	K	-1	-2	-3	-3	-2	-2	-3	-4	-2	-3	0	2	-3	0	0	-1	2	3	4	0
<u>4</u>	N	-2	1	-3	-3	-3	-2	1	-3	-2	-3	-1	-1	0	3	3	5	-1	-1	2	0
<u>5</u>	N	-2	-2	0	0	-1	1	-2	-4	-2	-3	0	0	-3	3	-1	-1	-1	-2	4	1
<u>6</u>	L	1	-2	1	1	1	0	-2	-3	4	-2	0	0	-2	-3	-2	-3	-2	-3	-3	-2
<u>7</u>	G	-1	5	-4	-4	-4	-3	-4	-3	-3	2	1	-2	-3	-1	4	-2	-1	-2	-2	-1
<u>8</u>	V	-2	-4	4	-1	4	-1	-2	-4	-4	-2	-2	2	-2	-3	-3	0	-3	-1	-4	-3
<u>9</u>	A	-1	-4	3	2	3	1	-2	-3	-3	-2	-2	2	-2	-3	-3	-4	-3	-3	-4	-3
<u>10</u>	V	0	-3	-1	-1	1	0	-2	-3	0	-3	-1	2	0	0	2	2	0	0	1	0
P	C	A	G	I	L	V	M	F	W	P	C	S	T	Y	N	Q	H	K	R	D	E
<u>11</u>	I	-1	-5	6	2	3	0	-1	-3	-4	-3	-3	-2	1	-5	-4	-4	-4	-4	-5	-4
<u>12</u>	G	0	3	-4	-3	-3	-3	-4	-4	-3	-4	-1	-2	-3	4	1	2	-1	-2	3	1
<u>13</u>	S	-2	4	-5	-4	-4	-3	-4	-4	-3	-4	0	-2	-4	0	-1	-2	3	3	-1	-2
<u>14</u>	K	-2	-3	-2	-3	-2	-2	-3	-4	-3	-4	-1	0	-1	1	1	4	4	3	-2	0
<u>15</u>	Q	-2	-1	-3	-3	-2	-2	-3	-3	0	-3	0	2	0	-1	2	-1	3	3	0	2
<u>16</u>	Y	-4	-5	-3	-1	-3	-2	6	3	-5	-4	-4	-4	7	-5	-4	-1	-4	-4	-5	-5
<u>17</u>	A	4	-3	0	-2	4	-2	-3	-4	-3	-2	-1	-2	-3	-4	-3	-4	-3	-4	-4	-3
<u>18</u>	V	1	-1	1	-1	3	-1	-3	-4	-3	5	2	1	-3	-3	-3	-3	-3	-3	-3	-3
<u>19</u>	N	-1	5	-5	-5	-4	-4	-5	-5	-3	-4	2	-2	-4	4	-2	-2	-2	-3	0	-2
<u>20</u>	L	-3	-6	1	6	-1	0	-1	-4	-5	-3	-4	-3	-3	-5	-4	-5	-4	-4	-6	-5

What are HMMs?

Hidden Markov models (HMMs) are used by many databases. Like profiles, they can be used to convert multiple sequence alignments into position-specific scoring systems. HMMs are adept at representing amino acid insertions and deletions, meaning that they can model entire alignments, including divergent regions. They are sophisticated and powerful statistical models, very well suited to searching databases for homologous sequences.

Multiple sequence alignment

Sequence 1:	F	K	L	L	S	H	C	L	L	V
Sequence 2:	F	K	A	F	G	Q	T	M	F	Q
Sequence 3:	Y	P	I	V	G	Q	E	L	L	G
Sequence 4:	F	P	V	V	K	E	A	I	L	K
Sequence 5:	F	K	V	L	A	A	V	I	A	D
Sequence 6:	L	E	F	I	S	E	C	I	I	Q
Sequence 7:	F	K	L	L	G	N	V	L	V	C



I = insert state

M = match state

D = delete state

Databases that use HMM to classify proteins:

Pfam

SMART

<http://pfam.xfam.org/>

<http://smart.embl-heidelberg.de/>

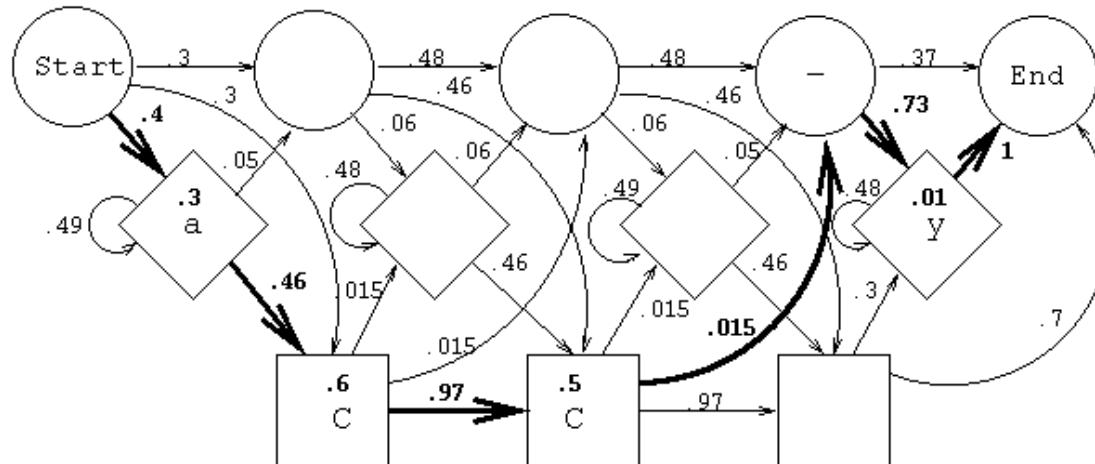
The Pfam database is a large collection of protein families

It stores profiles that contain all information inside a multiple alignment

LEXA_AERHH/1-65	MKPLTPRQAEVLELIKANMNETGMPPTRAEIAQKLGFKSANAAEEHLKALAKKGVIEIMPGTSRG
LEXA_SHIDS/1-65	MKALTARQQEVFDLIRDHISQTGMPPTRAEIAQRLGFRSPNAAAEEHLKALARKGVIEIVSGASRG
LEXA_HAEIN/1-65	MRPLTPRQEVLDLLKRHETTGMPPTRAEISRELGFKSANAAEEHLKALSRKGAIIEIPGASRG
LEXA_PSEAE/1-65	MQKLTTPRQAEIILSFIKRCLLEDHGFPPTRAEIAQELGFKSPNAAAEEHLKALARKGAIEMTPGASRG
LEXA_BACSU/1-65	MTKLSKRQLDLRFIKAEVKSKGYPPSVREIGEAVGGLASSSTVHGHHLARLETKGLIRRDPTKPRA
LEXA_MYCLE/20-84	DSGLTERQRTILNVRASVTSRQYPPSIREIADAVGTTSSSVAHQQLRTLERKGYLRRDPNRPRA
LEXA_STRC0/25-89	SSGLTDRQRRVIEVIRDSVQRRGYPPSMREIGQAVGSSSTSSVAHQMLALERKGFLRRDPHRPRA
LEXA_THEME/1-64	MKDTERQKVLLFIEEFIEKNQYPPSVREIARRFRIT.PRGALLHLIALEKKGYIERKNGKPKRA
LEXA_THEME/1-64 (SS)	XX---HHHHHHHHHHHHHHSS---HHHHHHHHHTS---HHHHHHHHHHHHHTSEE-GGG-CCGX
P73722_SYN3/1-65	MEPLTRAQKELFDWLVSYIDETOHAPSIRQMMRAMNLRSAPIQSRLERLRNKGYVDWTDGKART



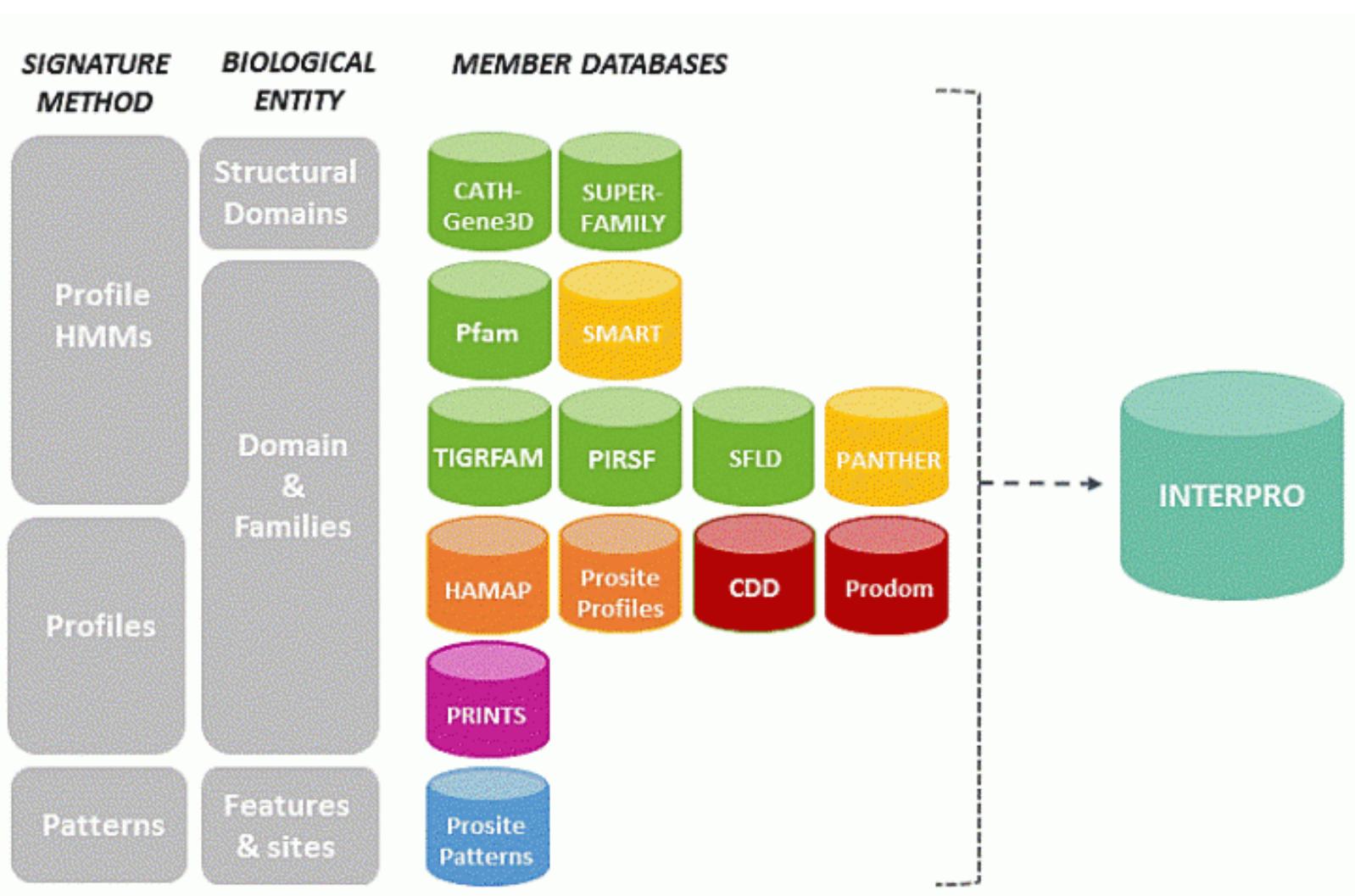
These profiles are stored in the form of a Hidden Markov model (HMM)



InterPro: Databases that combine patterns, profiles and HMMs to classify proteins

<https://www.ebi.ac.uk/interpro/>

InterPro is the main resource for protein classification at the EBI.

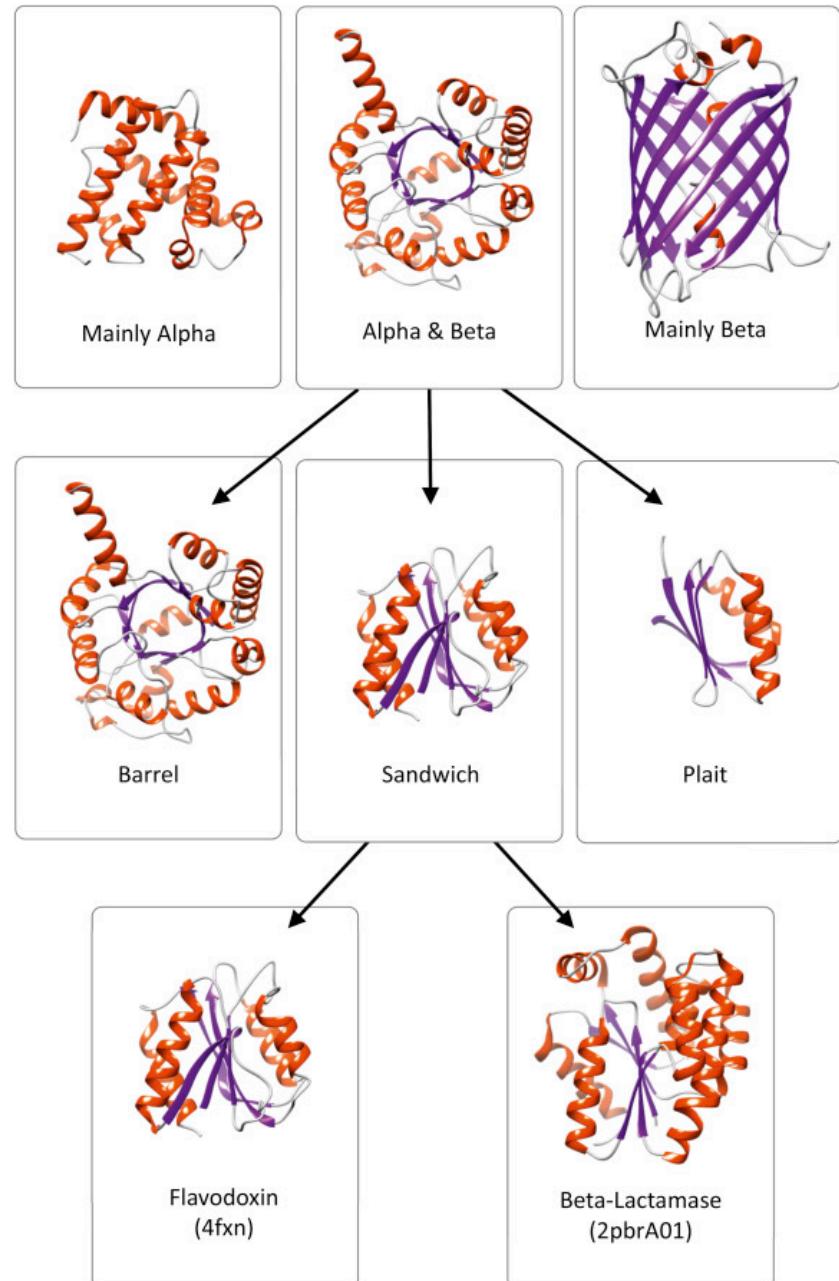


CATH: Protein Structure Classification Database

<http://www.cathdb.info/>

The first three levels of the CATH structure classification hierarchy:

1. Class (based on secondary structure content).
2. Architecture (based on gross spatial arrangement of secondary structures).
3. Topology or Fold (similar folding arrangement of secondary structures).

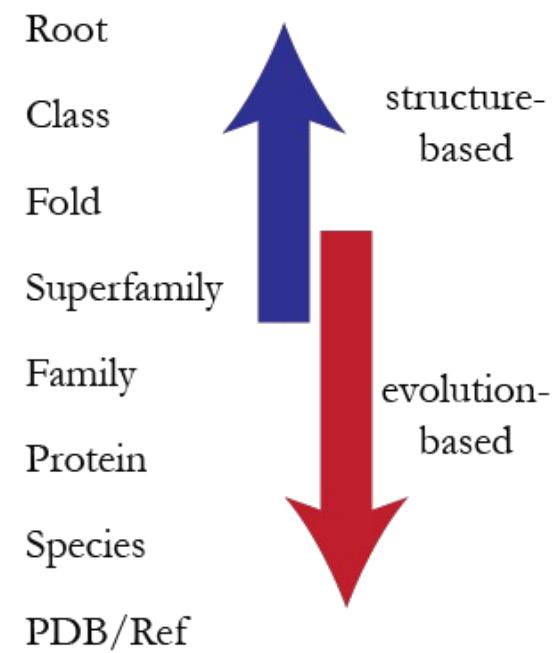
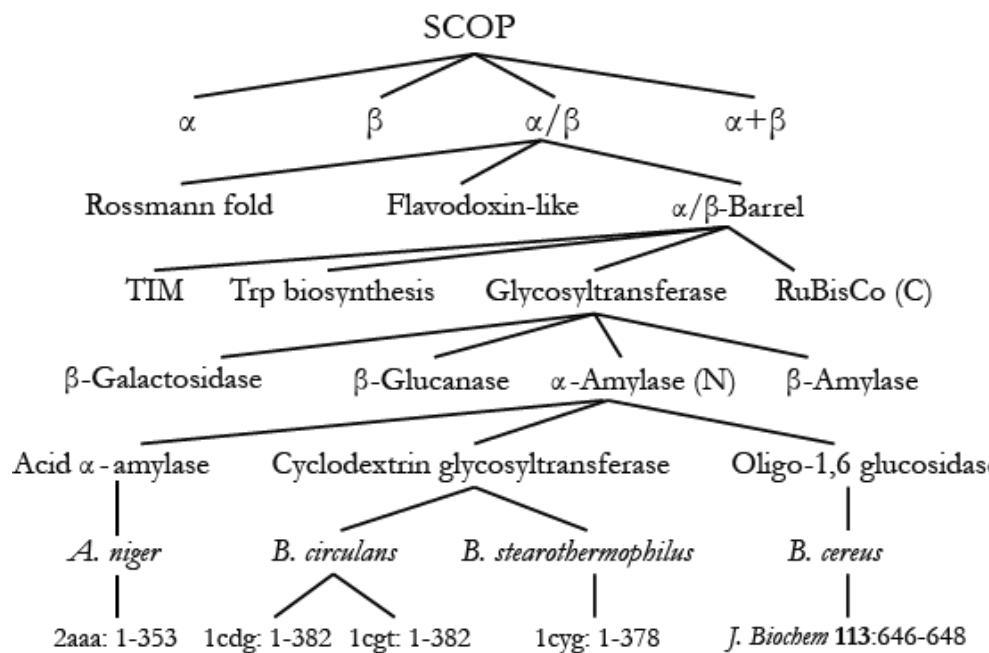


SCOP2: Structural Classification of Proteins database

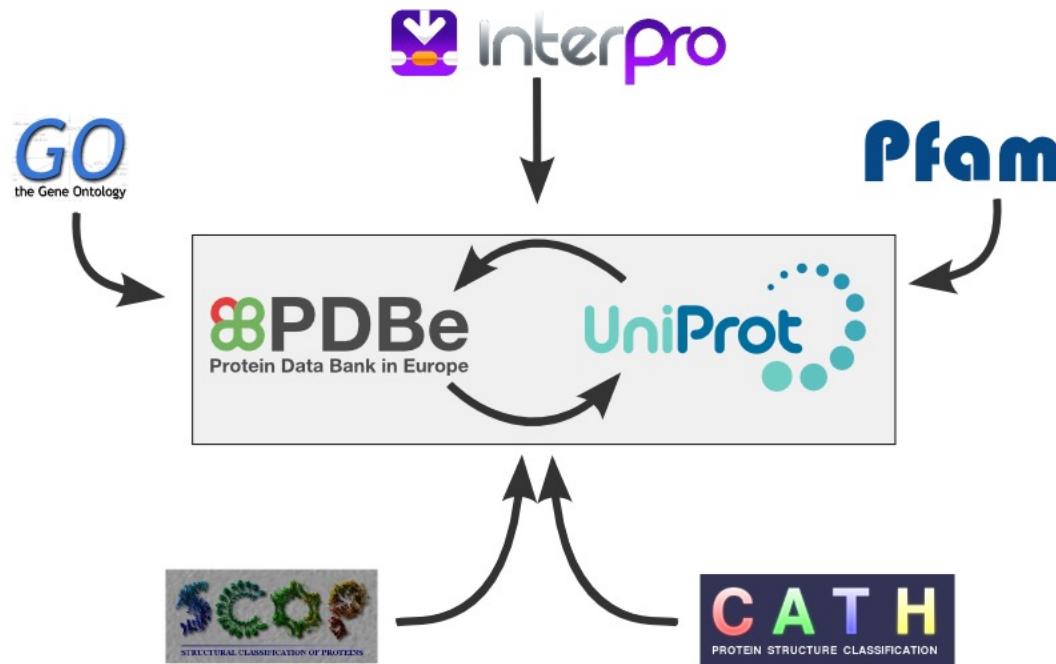
Manually (visually) classified

<http://scop2.mrc-lmb.cam.ac.uk/>

- All alpha
- All beta
- Alpha and beta (a/b) proteins with alternating alpha-helices and beta-strands
- Alpha + beta (a+b) proteins with segregated alpha-helices and beta-strands
- Small proteins little or no secondary structures



SIFTS: Structure integration with function, taxonomy and sequence



<https://www.ebi.ac.uk/pdbe/docs/sifts/index.html>

pdb_chain_uniprot.tsv.gz	A summary of the PDBe to UniProt residue level mapping, showing the start and end residues of the mapping using SEQRES, PDB sequence and UniProt numbering.
pdb_chain_taxonomy.tsv.gz	A summary of the NCBI tax_id(s),scientific_name(s) and chain type for each PDB chain that has been processed.
pdb_pubmed.tsv.gz	A summary of the Pubmed id(s) associated with each PDB entry, together with an ordinal number.
pdb_chain_enzyme.tsv.gz	A summary of the EC number(s) for each PDB chain that has been processed.
pdb_chain_go.tsv.gz	A summary of the GO identifier(s) for each PDB chain that has been processed.
pdb_chain_interpro.tsv.gz	A summary of the InterPro identifier(s) for each PDB chain that has been processed.
pdb_chain_pfam.tsv.gz	A summary of the Pfam domain identifier(s)(derived via the UniProt mapping) for each PDB chain that has been processed.
pdb_chain_cath_uniprot.tsv.gz	A summary of the CATH identifier(s) and UniProt primary accession number(s) for each PDB chain that has been processed.
pdb_chain_scop_uniprot.tsv.gz	A summary of the SCOP identifier(s) and UniProt primary accession number(s) for each PDB chain that has been processed.
uniprot_pdb.tsv.gz	A summary of the UniProt to PDB mappings showing the UniProt accession followed by a semicolon-separated list of PDB four letter codes.
pdb_chain_ensembl.tsv.gz	A summary of the Ensembl identifier(s) and UniProt primary accession number(s) for each PDB chain that has been processed.
pdb_chain_hmmer.tsv.gz	A summary of the Pfam domain identifier(s) (automatically calculated using HMMER) for each PDB chain that has been processed.
uniprot_segments_observed.tsv.gz	A summary of the UniProt to PDB residue level mapping (observed residues only), showing the start and end residues of the mapping using SEQRES, PDB sequence and UniProt numbering.
uniprot_segments_observed.csv.gz	

Specialized Sequence Databases in Health and Life Science

[Ensembl](#) automatic annotation databases for human, mouse, other vertebrate and eukaryote genomes.

[SILVA](#), a database of ribosomal RNAs

[VectorBase](#) The NIAID Bioinformatics Resource Center for Invertebrate Vectors of Human Pathogens

[Wormbase](#), genome of the model organism *Caenorhabditis elegans*

[PATRIC](#), the PathoSystems Resource Integration Center

[Flybase](#), genome of the model organism *Drosophila melanogaster*

[EcoCyc](#) genome and the biochemical machinery of the model organism *E. coli* K-12

[ENZYME](#) (Enzyme database at ExPASy, CH)

[BRENDA](#) (The comprehensive enzyme information system, DE)

[ReLiBase](#) (database system for analysing receptor/ligand complexes deposited in the Protein Data Bank at EBI, UK)

[HGMD disease-causing mutations](#) (HGMD Human Gene Mutation Database)

[HUGO](#) (Official Human Genome Database: HUGO Gene Nomenclature Committee)

[IEDB](#) Immune Epitope Database

[OrthoMCL](#) Ortholog Groups of Protein Sequences from Multiple Genomes.

[DrugBank](#) combines detailed drug data with comprehensive drug target information

[PRIDE](#) Archive - proteomics data repository

[EST](#) database is a collection of short single-read transcript sequences from GenBank.

