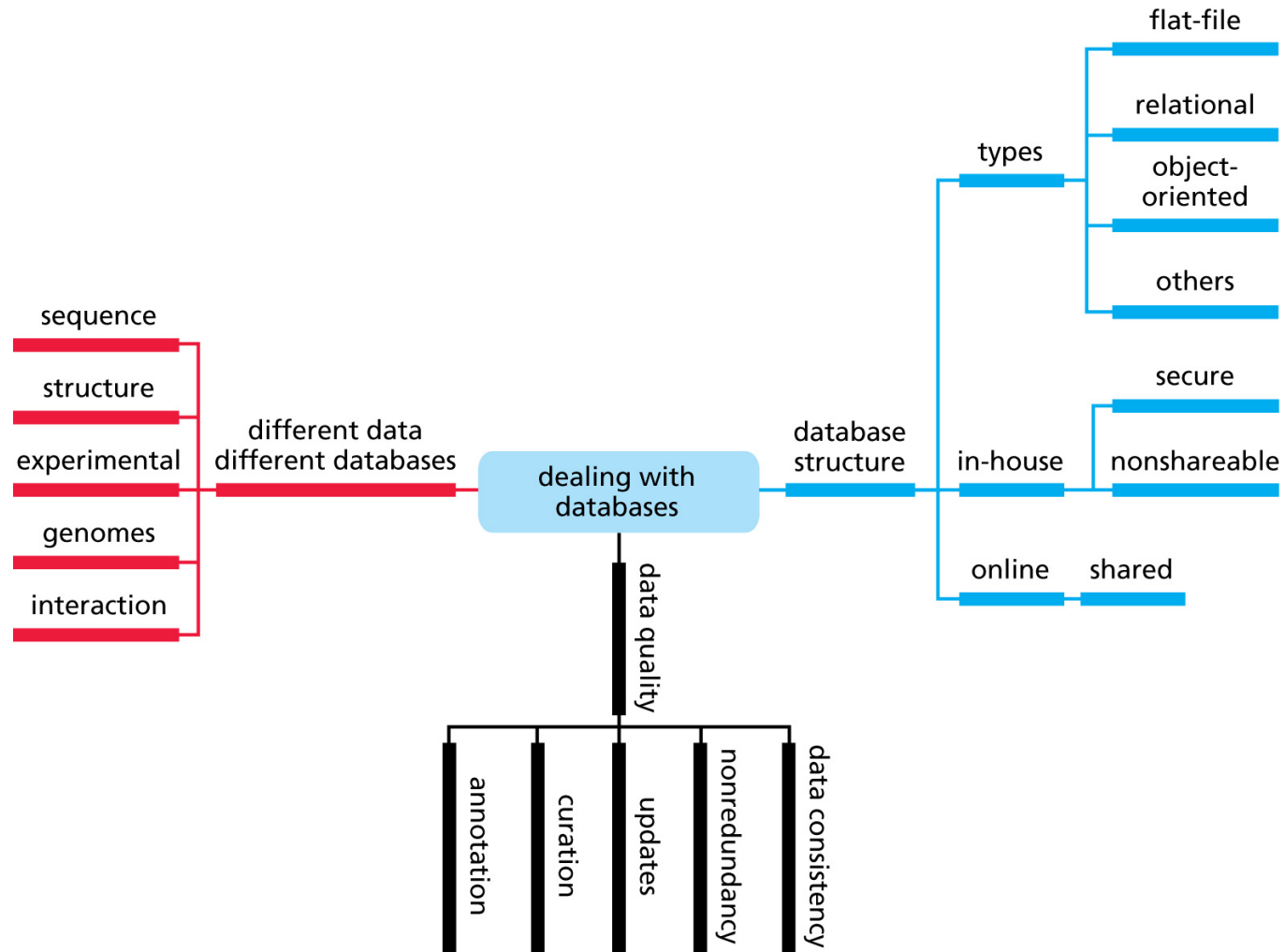
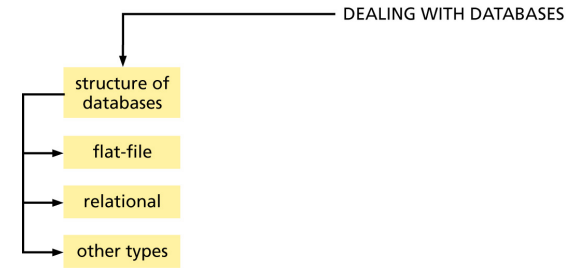


General aspects of databases



Structure of databases



Flat database: it is the simplest form of a database where collections of data (aminoacid sequence) are stored as a large txt file or more than one txt file.

Relational database: it stores the data within a number of tables, each consisting of records and fields. Each table will be linked to at least one other by a shared field called a KEY.

protab1			
Protein-code	Protein-name	Length	Species-origin
P1001	Hemoglobin	145	Bovine
P1002	Hemoglobin	136	Ovine
P1003	Eye Lens Protein	234	Human
.....			

protab2	
Protein-code	Protein-sequence
P1001	MDRTHGFDLKLSPRTVNQWLMLALFFGHS...
P1002	MDKTSHGFEIKLLTPKKLQQWLMIAIYFGHT...
P1003	SRTHEEEGKLMQWPPRPLYIALFTEPPYP...
.....	

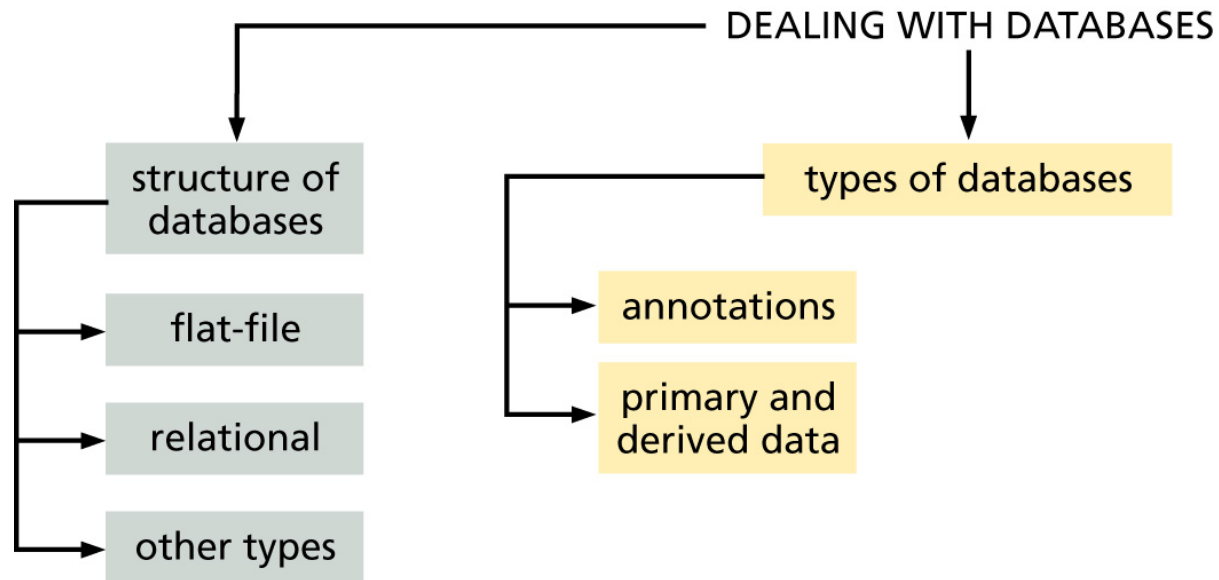
Type of databases

Data: it is the minimal content of a database including data's identity (for example protein name and source) and the author/submitter responsible for the entry.

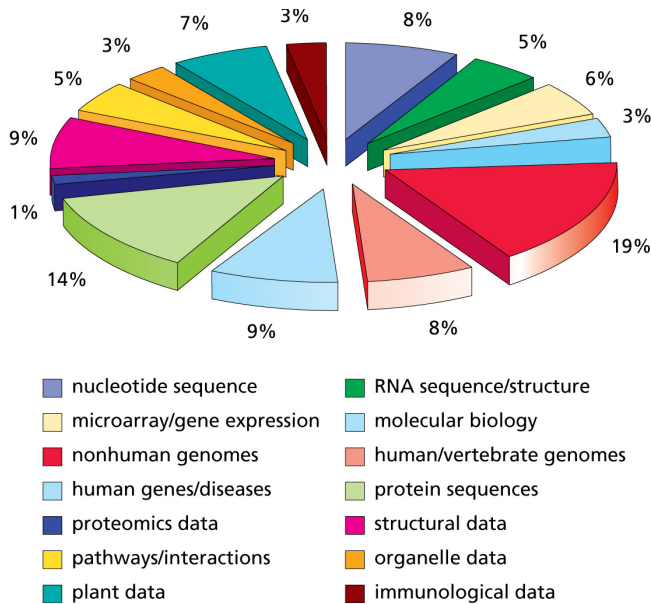
Annotation: provide more information to the data (published papers, lists of entries in other databases, gene structure)

Primary data: they include the raw experimental results.

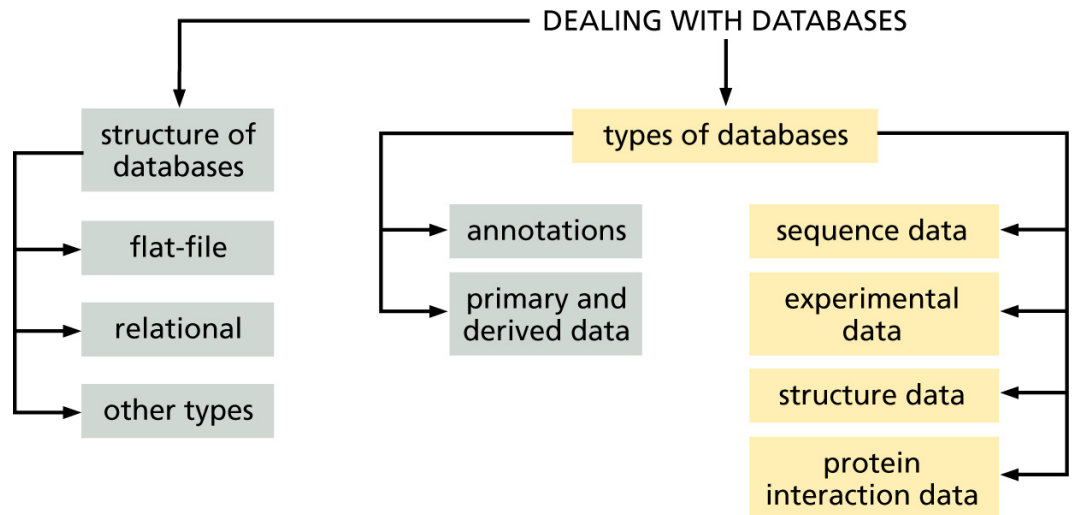
Derived data: based on the data existing at the time (example: conserved protein sequence motifs).



Looking for databases



Distribution of the type of databases as classified at the Nucleic Acid Research (NRA) Molecular Biology Database Collection Web site. In 2006 there were 858 databases listed, classified into 14 main categories.



Sequence database

1. DNA sequences:

- Raw genomic sequence (chromosomal DNA)
- cDNA (from mRNA)
- Expressed sequence tags (ESTs). Partial cDNA sequence.

2. Protein sequences (UniProtKB, Swiss-Prot, NCBI Protein Database)

```
LOCUS       RM_005358             7235 bp  mRNA    linear   PRI 02-AUG-2006
DEFINITION   Homo sapiens LIM domain 7 (LMO7), mRNA.
ACCESSION   RM_005358
VERSION     RM_005358.4  GI:111119012
KEYWORDS     -
SOURCE      Homo sapiens (human)
```

```
gene
1..7235
/ gene="LMO7"
/ locus="LMO7"
/ function="Synonyms: LOMF, FBX20, FEZO20, KIA0858"
/ db_xref="GeneID:1008"
/ db_xref="MIM:604104"
/ db_xref="EMBL:1008"
/ db_xref="FRC045078"
/ db_xref="MIM:604104"
CDS
1..1111
/ gene="LMO7"
```

```
ORIGIN
1 ggaagaaat ggaataaata ggaactatg gtggagtag gtagagagg atttcaaca
61 ttaatpqa taaagatga caagctctg taastpqa atctcgatc acagctctg
121 tgaagatga gatgtgatg ttcaacaag cactcaaga attatgaag tgttgagcc
181 agpctcaat taaataaga tgaactatg taactatga agactatg tggtagaga
241 aaatacaaa tccatgatg gatgtttag atctatcag atataactt taagaaga
301 aaataatg tccatgatg taactatga gctttccaa gattatggt tccactaga
361 agpctgatg tgaagatg cttaactatg taactatg tctgttatg ataatatg
421 tttgttcca caagcttag tgaagatg gatattaga tatpqaat aatatatg
481 gctttcaaa aaaaataaa ttaattgat tctcaacta tatytaaa taattgttg
541 taagtaga tcaactatg gttatgagt ctgtttctg tctgttatg taactatg
601 tttgttagg aataataa cttaagaaa aatatcttca tgcacttt acatgaag
661 taataaga gttactaag tccagcagg caaatatg tcttaataa ctactgtct
721 aaactatg taactatg tgaactatg ctcaactatg gttcaacta tcaactatg
781 caactatg tcaactatg taaataag agatgatg ataatatg atttgttct
841 caagcagg gaaactatg atattctt tttgtttag aaagtatg attttttaa
901 atcaactat agatgatg agpctatg caactatg agatgatg gaactatg
961 ctctatga atactatg acaactatg actctatg tcaactatg tttatgta
1021 atctatga ctggagag agatgatg tttctctt tgaactatg tttatgta
1081 tcaactat ctggagag tttctctt tcaactatg tttatgta tttatgta
1141 tttatgta tttatgta gatgaatg aatgatg taagtatg taagtatg agpctatg
1201 aaatgatg aatgatg taagtatg taagtatg taagtatg taagtatg taagtatg
1261 agpctatg taagtatg taagtatg taagtatg taagtatg taagtatg taagtatg
1321 ctctatga taagtatg taagtatg taagtatg taagtatg taagtatg taagtatg
1381 gttatgta gttatgta gttatgta gttatgta gttatgta gttatgta gttatgta
1441 atctatga atcaactatg taactatg atcaactatg taactatg taactatg
1501 aatgatga caaactatg taactatg atcaactatg taactatg taactatg
1561 atctatga taactatg taactatg atcaactatg taactatg taactatg
1621 gttatgta atcaactatg taactatg atcaactatg taactatg taactatg
1681 gttatgta atcaactatg taactatg atcaactatg taactatg taactatg
1741 gttatgta atcaactatg taactatg atcaactatg taactatg taactatg
1801 caactatg taactatg taactatg atcaactatg taactatg taactatg
1861 gttatgta atcaactatg taactatg atcaactatg taactatg taactatg
1921 atcaactatg taactatg taactatg atcaactatg taactatg taactatg
1981 caactatg taactatg taactatg atcaactatg taactatg taactatg
2041 gttatgta atcaactatg taactatg atcaactatg taactatg taactatg
2101 gttatgta atcaactatg taactatg atcaactatg taactatg taactatg
2161 gttatgta atcaactatg taactatg atcaactatg taactatg taactatg
2221 caactatg taactatg taactatg atcaactatg taactatg taactatg
2281 atcaactatg taactatg taactatg atcaactatg taactatg taactatg
2341 atcaactatg taactatg taactatg atcaactatg taactatg taactatg
2401 atcaactatg taactatg taactatg atcaactatg taactatg taactatg
2461 atcaactatg taactatg taactatg atcaactatg taactatg taactatg
2521 gttatgta atcaactatg taactatg atcaactatg taactatg taactatg
2581 atcaactatg taactatg taactatg atcaactatg taactatg taactatg
2641 atcaactatg taactatg taactatg atcaactatg taactatg taactatg
2701 caactatg taactatg taactatg atcaactatg taactatg taactatg
2761 gttatgta atcaactatg taactatg atcaactatg taactatg taactatg
2821 gttatgta atcaactatg taactatg atcaactatg taactatg taactatg
2881 atcaactatg taactatg taactatg atcaactatg taactatg taactatg
2941 atcaactatg taactatg taactatg atcaactatg taactatg taactatg
3001 atcaactatg taactatg taactatg atcaactatg taactatg taactatg
3061 gttatgta atcaactatg taactatg atcaactatg taactatg taactatg
3121 agpctatg taactatg taactatg atcaactatg taactatg taactatg
3181 gttatgta atcaactatg taactatg atcaactatg taactatg taactatg
3241 tttcttcc agpctatg tttcttcc agpctatg tttcttcc agpctatg
3301 gatgtatg tttcttcc agpctatg tttcttcc agpctatg tttcttcc agpctatg
3361 caactatg taactatg taactatg atcaactatg taactatg taactatg
3421 tttgttga caaataag gttatgta gttatgta gttatgta gttatgta
3481 agpctatg taactatg taactatg atcaactatg taactatg taactatg
3541 atcaactatg taactatg taactatg atcaactatg taactatg taactatg
3601 gttatgta atcaactatg taactatg atcaactatg taactatg taactatg
```

```
3661 atgtatgta ctctgatg ttaactatg gaaactatg caactatg tgaactatg
3721 gatcttcc agpctatg gatctatg atgtatga agpctatg tgaactatg
3781 gttatgta atgttctga atcaactatg tcaactatg ttttttga aaactatg
3841 agpctatg taactatg taactatg atcaactatg tcaactatg ttttttga
3901 caactatg tcaactatg ttttttga atgtatga gttatgta agpctatg
3961 gttatgta agpctatg gttatgta atgtatga atgtatga agpctatg
4021 aaactatg agpctatg agpctatg atgtatga atgtatga agpctatg
4081 ttttttga agpctatg atgtatga atgtatga atgtatga agpctatg
4141 ttttttga atgtatga atgtatga atgtatga atgtatga agpctatg
4201 atgtatga atgtatga atgtatga atgtatga atgtatga agpctatg
4261 gttatgta agpctatg gttatgta atgtatga atgtatga agpctatg
4321 agpctatg agpctatg gttatgta atgtatga atgtatga agpctatg
4381 gttatgta agpctatg agpctatg atgtatga atgtatga agpctatg
4441 agpctatg ttttttga ttttttga atgtatga atgtatga agpctatg
4501 atgtatga gttatgta atgtatga atgtatga atgtatga agpctatg
4561 atgtatga atgtatga atgtatga atgtatga atgtatga agpctatg
4621 agpctatg atgtatga atgtatga atgtatga atgtatga agpctatg
4681 agpctatg caactatg atgtatga atgtatga atgtatga agpctatg
4741 aaactatg ctctgatg ttttttga atgtatga atgtatga agpctatg
4801 gatcttcc atgtatga atgtatga atgtatga atgtatga agpctatg
4861 gatcttcc atgtatga atgtatga atgtatga atgtatga agpctatg
4921 agpctatg atgtatga atgtatga atgtatga atgtatga agpctatg
4981 agpctatg caactatg gttatgta atgtatga atgtatga agpctatg
5041 atgtatga atgtatga atgtatga atgtatga atgtatga agpctatg
5101 ttttttga atgtatga atgtatga atgtatga atgtatga agpctatg
5161 ttttttga atgtatga atgtatga atgtatga atgtatga agpctatg
5221 atgtatga atgtatga atgtatga atgtatga atgtatga agpctatg
5281 ttttttga atgtatga atgtatga atgtatga atgtatga agpctatg
5341 atgtatga atgtatga atgtatga atgtatga atgtatga agpctatg
5401 atgtatga atgtatga atgtatga atgtatga atgtatga agpctatg
5461 ttttttga atgtatga atgtatga atgtatga atgtatga agpctatg
5521 atgtatga atgtatga atgtatga atgtatga atgtatga agpctatg
5581 atgtatga atgtatga atgtatga atgtatga atgtatga agpctatg
5641 gttatgta atgtatga atgtatga atgtatga atgtatga agpctatg
5701 atcaactatg taactatg taactatg atcaactatg taactatg taactatg
5761 atcaactatg taactatg taactatg atcaactatg taactatg taactatg
5821 atcaactatg taactatg taactatg atcaactatg taactatg taactatg
5881 atcaactatg taactatg taactatg atcaactatg taactatg taactatg
5941 atcaactatg taactatg taactatg atcaactatg taactatg taactatg
6001 ttttttga atgtatga atgtatga atgtatga atgtatga agpctatg
6061 atgtatga atgtatga atgtatga atgtatga atgtatga agpctatg
6121 atgtatga atgtatga atgtatga atgtatga atgtatga agpctatg
6181 atgtatga atgtatga atgtatga atgtatga atgtatga agpctatg
6241 atgtatga atgtatga atgtatga atgtatga atgtatga agpctatg
6301 atgtatga atgtatga atgtatga atgtatga atgtatga agpctatg
6361 atcaactatg taactatg taactatg atcaactatg taactatg taactatg
6421 atgtatga atgtatga atgtatga atgtatga atgtatga agpctatg
6481 atgtatga atgtatga atgtatga atgtatga atgtatga agpctatg
6541 ttttttga atgtatga atgtatga atgtatga atgtatga agpctatg
6601 gttatgta atgtatga atgtatga atgtatga atgtatga agpctatg
6661 gttatgta atgtatga atgtatga atgtatga atgtatga agpctatg
6721 agpctatg atgtatga atgtatga atgtatga atgtatga agpctatg
6781 gatcttcc atgtatga atgtatga atgtatga atgtatga agpctatg
6841 gatcttcc atgtatga atgtatga atgtatga atgtatga agpctatg
6901 atcaactatg taactatg taactatg atcaactatg taactatg taactatg
6961 atcaactatg taactatg taactatg atcaactatg taactatg taactatg
7021 atcaactatg taactatg taactatg atcaactatg taactatg taactatg
7081 ttttttga atgtatga atgtatga atgtatga atgtatga agpctatg
7141 ttttttga atgtatga atgtatga atgtatga atgtatga agpctatg
7201 ttttttga atgtatga atgtatga atgtatga atgtatga agpctatg
```

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	LMO7_HUMAN
Primary accession number	Q8W11
Secondary accession numbers	O15462 O95346 Q9UKC1 Q9UQM5 Q9Y6A7
Integrated into Swiss-Prot on	March 15, 2004
Sequence was last modified on	March 15, 2004 (Sequence version 2)
Annotations were last modified on	July 25, 2006 (Entry version 39)

Name and origin of the protein	
Protein name	LIM domain only protein 7
Synonyms	LOMP F-box only protein 20
Gene name	Name: LMO7 Synonyms: FBX20, FEZO20, KIA0858
From	Homo sapiens (Human) [TaxID: 9606]
Taxonomy	Eukaryota; Metazoa; Chordata; Craniata; Vertebrata; Eumetazoa; Mammalia; Eutheria; Euarchontales; Primates; Haplorhina; Catarrhini; Homidae; Homo.

References	
[1]	NUCLEOTIDE SEQUENCE [mRNA] (ISOFORM 3), AND TISSUE SPECIFICITY. TISSUE=Brain, and Peripheral blood leukocyte; DOI=10.1007/s00439-001-0646-6; PubMed=11935316 [NCBI, EMBASE, EBI, Israel, Japan] Rozenblum E., Vaheriisto P., Sandberg T., Berthelsen J.T., Syrjakoski K., Weaver D., Haraldsson K., Johannsdottir H.K., Wehman P., Nigam S., Goltsov N., Robbins C., Pak E., Dutra A., Gilliland E., Stephan D.A., Bailey-Wilson J., Joo S.-H.H., Kainu T., Kallioniemi O.-P.; "A genomic map of a 6-Mb region at 13q21-q22 implicated in cancer development: identification and characterization of candidate genes."; Hum. Genet. 110:111-121(2002).

Key	From	To	Length	Description	FTId	
CHAIN	1	1683	1683	LIM domain only protein 7.	PRO_0000075824	
DOMAIN	54	168	115	CH.		
DOMAIN	1042	1128	87	PDZ.		
DOMAIN	1612	1678	67	LIM zinc-binding.		
	10	20	30	40	50	60
MKKIRICHIF	TFYSWMSYDV	LFQTELGAL	EIWRQLCAH	VCICVGLYL	RDRCVSKDI	
	70	80	90	100	110	120
ILRTEQNSGR	TILIKAVTEK	NFETKDFRS	LENGVLLCD	INKLPGVIK	KINRLSTPIA	
	130	140	150	160	170	180
GLDNINFLK	ACEQIGLKEA	QLFHPGLQD	LSNRVTWKE	ETDRRVKNVL	ITLYLWRKKA	

Structural database

They contain information about the structure of small molecules, proteins, DNA and RNA sequences, carbohydrates.

RCSB PDB
PROTEIN DATA BANK



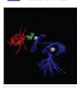
A MEMBER OF THE **PDB**
An Information Portal to Biological Macromolecular Structures

Contact Us | Help | Print Page | PDB ID or keyword | Author | SEARCH | Advanced Search

Home Search Results Queries

91 Structure Hits | 127 Web Page Hits | 1 Unreleased Structure

1 2 3 4 5 .. 10

- 1X62  **Solution structure of the LIM domain of carboxyl terminal LIM domain protein 1**
Release Date: 17-Nov-2005 Exp. Method: NMR 20 Structures
Structural Protein
Mol. Id: 1 Molecule: C Terminal Lim Domain Protein 1 Fragment: Lim Domain
Authors: Qin, X.R., Nagashima, T., Hayashi, F., Yokoyama, S.
- 1X4K  **Solution structure of LIM domain in LIM-protein 3**
Release Date: 14-Nov-2005 Exp. Method: NMR 20 Structures
Metal Binding Protein
Mol. Id: 1 Molecule: Skeletal Muscle Lim Protein 3 Fragment: Lim Domain
Authors: He, F., Muto, Y., Inoue, M., Kigawa, T., Shirouzu, M., Terada, T., Yokoyama,
- 1X4L  **Solution structure of LIM domain in Four and a half LIM domains protein 2**
Release Date: 14-Nov-2005 Exp. Method: NMR 20 Structures
Metal Binding Protein
Mol. Id: 1 Molecule: Skeletal Muscle Lim Protein 3 Fragment: Lim Domain
Authors: He, F., Muto, Y., Inoue, M., Kigawa, T., Shirouzu, M., Terada, T., Yokoyama,

Protein folds have also been classified according to the conservation of the fold. They include CATH and SCOP.

Structural Classification of Proteins



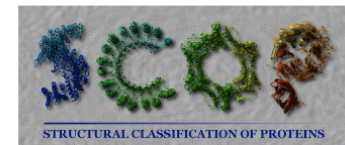
Welcome to **SCOP**: Structural Classification of Proteins.
1.75 release (June 2009)

38221 PDB Entries. 1 Literature Reference. 110800 Domains. (excluding nucleic acids and theoretical models).
Folds, superfamilies, and families [statistics here](#).
[New folds superfamilies families](#).
[List of obsolete entries and their replacements](#).

Authors. Alexey G. Murzin, John-Marc Chandonia, Antonina Andreeva, Dave Howorth, Loredana Lo Conte, Bartlett G. Ailey, Steven E. Brenner, Tim J. P. Hubbard, and Cyrus Chothia. scop@mrc-lmb.cam.ac.uk

Reference: Murzin A. G., Brenner S. E., Hubbard T., Chothia C. (1995). SCOP: a structural classification of proteins database for the investigation of sequences and structures. *J. Mol. Biol.* 247, 536-540. [PDF]

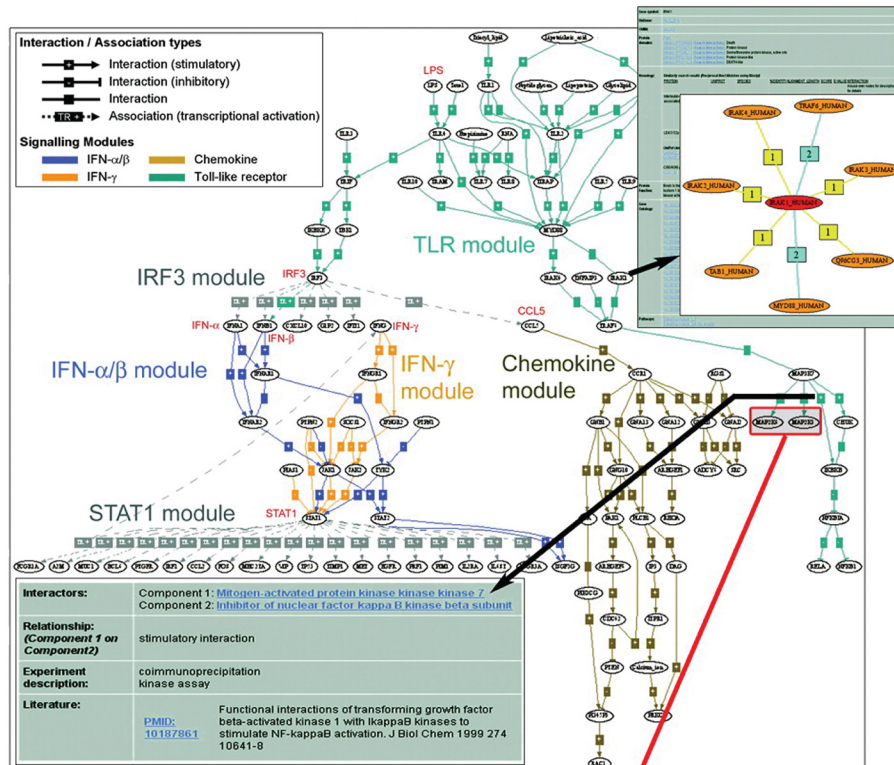
Recent changes are described in: Lo Conte L., Brenner S. E., Hubbard T.J.P., Chothia C., Murzin A. (2002). SCOP database in 2002: refinements accommodate structural genomics. *Nucl. Acid Res.* 30(1), 264-267. [PDF], Andreeva A., Howorth D., Brenner S.E., Hubbard T.J.P., Chothia C., Murzin A.G. (2004). SCOP database in 2004: refinements integrate structure and sequence family data. *Nucl. Acid Res.* 32:D226-D229. [PDF], and Andreeva A., Howorth D., Chandonia J.-M., Brenner S.E., Hubbard T.J.P., Chothia C., Murzin A.G. (2007). Data growth and its impact on the SCOP database: new developments. *Nucl. Acids Res.* 2008 36: D419-D425: doi:10.1093/nar/gkm993 [PDF].



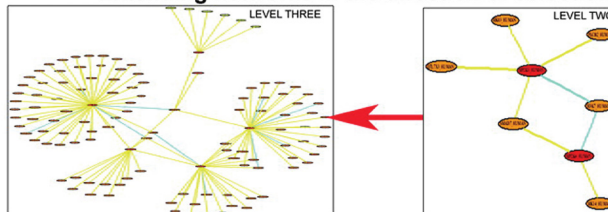
Protein interaction databases

They provide information about the interactions of proteins with other molecules, including other proteins.

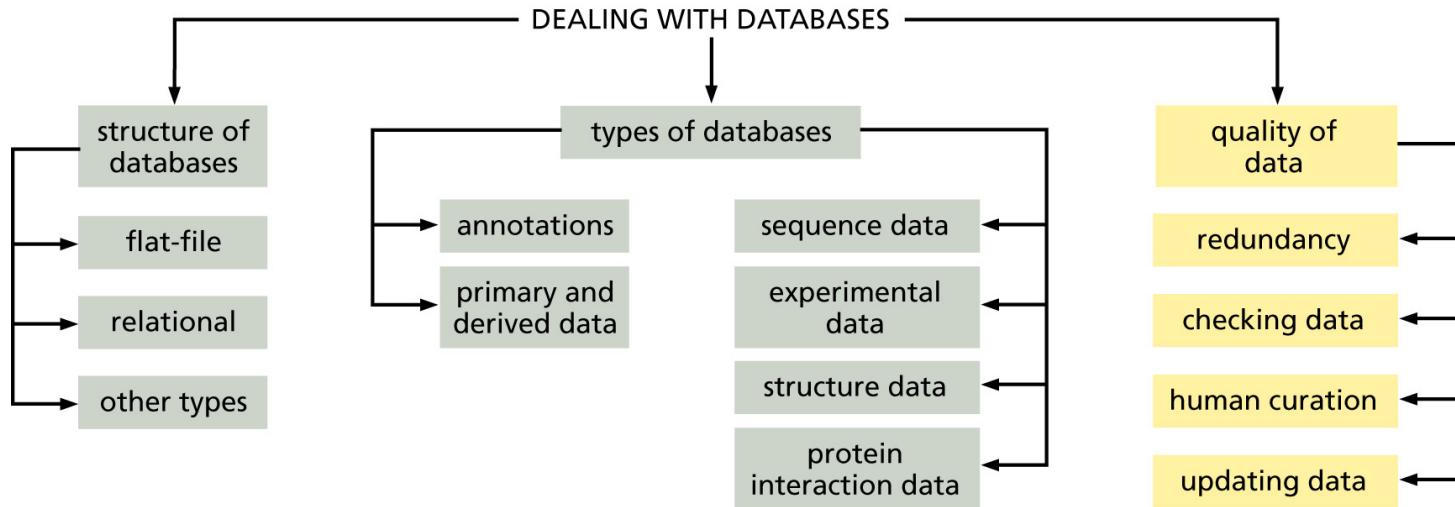
They include: the Database of Interacting Proteins (DIP) and the Molecular INTeraction Database (MINT).



Extending networks in the direction of selection



Quality of databases



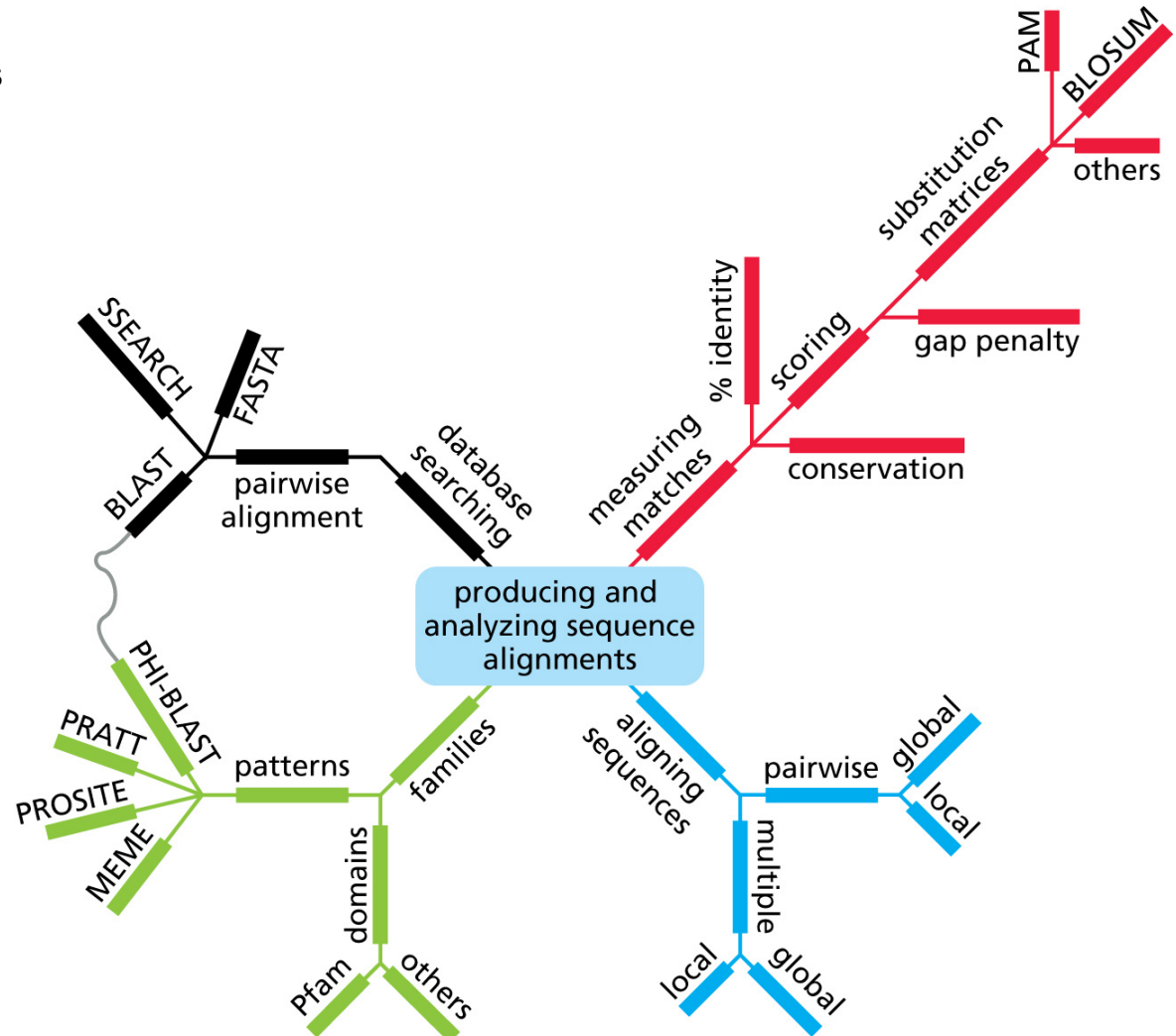
Non redundant databases: they include all the experimental data (from different labs) in one entry.

Checking data: a DNA sequence must contain only A, C, G, T. A protein sequence must correspond to a certain molecular weight according to the amino acids present.

Sequence alignments

Useful for:

- comparing an unknown sequence to all the sequences contained in a database;
- prediction of a protein structure
- construction of phylogenetic trees



Sequence alignments

Alignment is the task of locating equivalent regions of two or more sequences to maximize their similarity.

THIS SEQUENCE
THAT SEQUENCE

The differences in length between two or more sequences can be compensated by the introduction of **GAPS**.

THISISA- SEQUENCE
TH - - - -AT SEQUENCE

Gap penalty: each time a gap is introduced, the penalty is subtracted from the score, decreasing the overall score of the alignment.

(A)

```

Bovine PI-3Kinase p110a      LNWENPDMISEL L FGNNE I IFKNGDDL RQD ML TLQIIRIMENIWQNGQLDLRMLPYGCLSIGDCVGLIEV V RNSHTIMQIQCKGGLK GAL
cAMP-dependent protein kinase --WENPAQNTAHL D QFER I KTLGTGSFGRV ML VKHMETGNHYAMKILDQKQVVKLQIEHTLNEKRILQAV NFPFLVKLEFSFKD NSNLY

Bovine PI-3Kinase p110a      QFNSHTLHQWLKDNKNGEYDAAIDLFRSCAGYCVATFILGIGDRHNSNIMVKDDGQLFHIDFGHFLDHK K KKKFGYKRERVPFVLTQDF
cAMP-dependent protein kinase MVMEYVPGGEMFSLRRRIGRFSEPHARFYAAQIVLTFEYHLSDLIYRDLKPENLLIDQQGYIQVTDGFGFA KRVKGRTWXLCGTP EY LAP

Bovine PI-3Kinase p110a      L I V I S K G A Q E C T K T R E F E R F Q E M C Y K A Y L A I R Q H A N L F I N L F S M M L G S G M P E L Q S F D I A I Y I R K T L A L D K T E Q E A L E Y F M K Q M N D A H H G G
cAMP-dependent protein kinase E I I L S K G Y N K A V D W W A L G V L I Y E M A A G Y P P F F A D Q P I Q I Y E K I V S G K V R F P S H F S S D L K D L L R N L L Q V D L T K R F G N L K N G V N D I K N H K W F

Bovine PI-3Kinase p110a      W T T K M D W I F H T I K Q H A L N -----
cAMP-dependent protein kinase A T T D W I A I Y Q R K V E A P F I P K F K G P G D T S N F D D Y E E E E I R V X I N E K C G K E F S E F
    
```

A) An alignment where the gap penalty has been set very high.

B) An alignment with a very long gap penalty. Many more gaps have been introduced.

(B)

```

Bovine PI-3Kinase p110a      LNWENPDMISEL L FGNNE I IFKNGDDL RQD ML TLQIIRIMENIWQNGQLDLRMLPYGCLSIGDCVGLIEV V RNSHTIMQIQCKGGLK GAL
cAMP-dependent protein kinase ?-WENPAQNTAHL D QFER I KTLGTGSFGRV ML VKHM--ETGNHYAMK I LDKQKV-VKLGQIEHTLNEKRILQAVNFPFLVKLEFSFKD N-

Bovine PI-3Kinase p110a      QFNSHTLHQWLKDNKNGEYDAAIDLFRSCAGYCVATFILGIGDRHNSNIMVKD-DGQLFHIDFGHFLDHKKKKFGYKRERVPFVL--T
cAMP-dependent protein kinase -SNLYMMEYVPGGEMFSLRR-IGRFSEPHARFYAAQIVLTFEYHLSDLIYRDLKPENLLIDQQGYIQVTDGFGFAKRVKGRTWXLCGT

Bovine PI-3Kinase p110a      QDFL---I V I S K G A Q E C T K T R E F E R F -Q E M C --Y K A Y L A I R Q H A N L F I N L F S M M L G S G M P E L Q S F D I A I Y I R K T L A L D K T E Q E A L E Y F M K
cAMP-dependent protein kinase P E Y L A P E I I L S K G Y N K A V D W W A L G V L I Y E M A A G Y P P F F A -D Q P I Q I Y E K I V S G K V R F --P S H F S S D L K D L L R N L L Q V D L T K R --F G N L K N

Bovine PI-3Kinase p110a      Q M N D A H H G G W T T K M D W I -----F H T I K Q H A L -----N-----
cAMP-dependent protein kinase G V N D I K N H K W F A T T D W I A I Y Q R K V E A P F I P K F K G P G D T S N F D D Y E E E E I R V X I N E K C G K E F S E F
    
```

Sequence alignments

Similarity: the sequences show some degree of match.

Homology: similarity in sequence or structure due to descent from a common ancestor.

Mutation and selection over millions of years can result in considerable divergence between present-day sequences derived from the same ancestral gene.

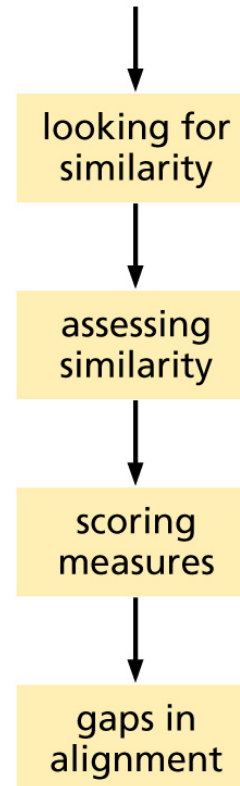
Bases at originally same position can change as a result of:

- Mutations
- Insertions
- Deletions
- Gene fusions

Homology \Rightarrow common ancestor \Rightarrow common structure or function?

Not always.....

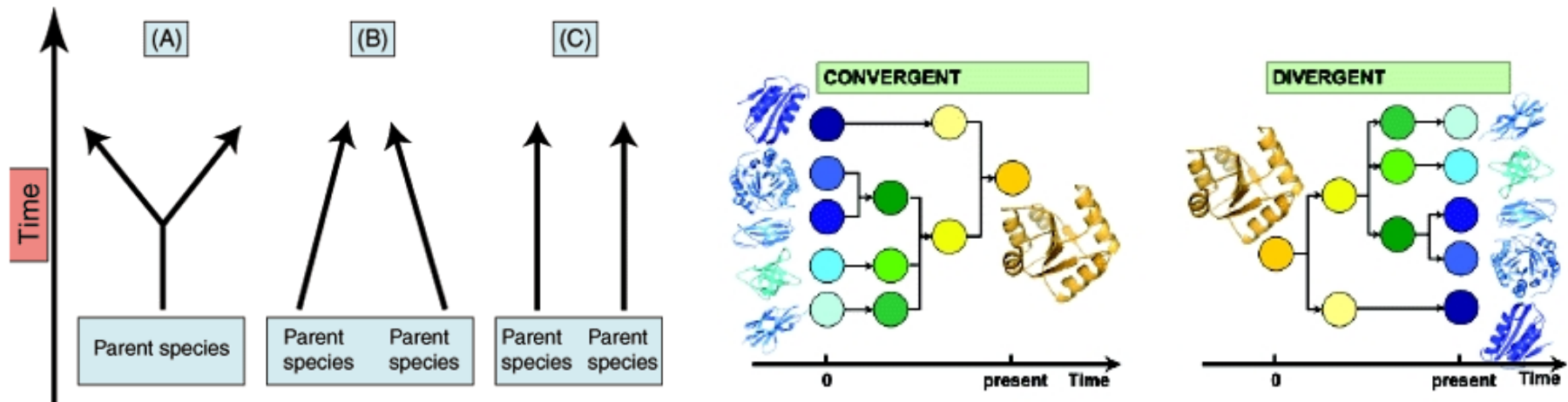
PRODUCING AND ANALYZING SEQUENCE ALIGNMENTS



Sequence alignments

Divergent evolution: mutation and selection can generate proteins with new functions but relatively little changes in sequence. Therefore, sequence similarity does not always imply a common function.

Convergent evolution: proteins with very little sequence similarity to each other but in which a common protein fold and function are preserved.

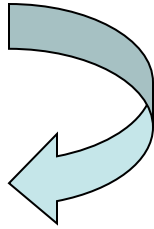


It is easier to compare to detect homology when comparing protein sequence than when comparing nucleic acid sequences.

1. There are only 4 letters to compare in the DNA alphabet compared to the 20 letters in the protein one
2. The genetic code is redundant
3. The 3D structure of a protein and hence its function, is determined by the amino acid sequence

Scoring alignments

The quality of an alignment is measured by giving it a quantitative score



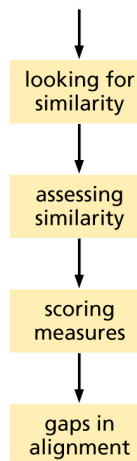
Percent identity: obtained by dividing the number of identical matches by the total length of the aligned region and multiplying by 100.

A good percentage of identity depends on the length of the sequence.



Substitution matrices: the score is assigned to each aligned pair of amino acids by a matrix that defines values for all possible pairs of residues.

PRODUCING AND ANALYZING
SEQUENCE ALIGNMENTS

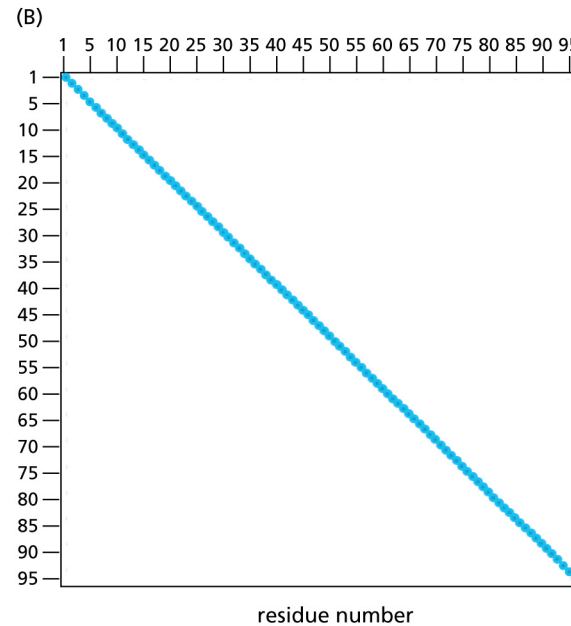
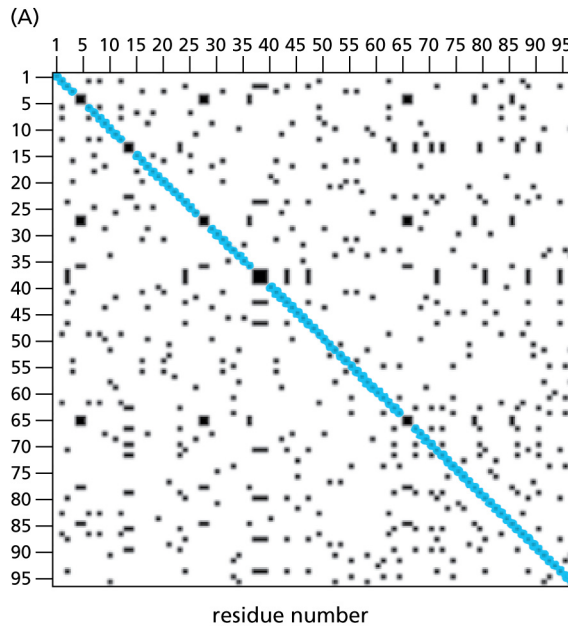
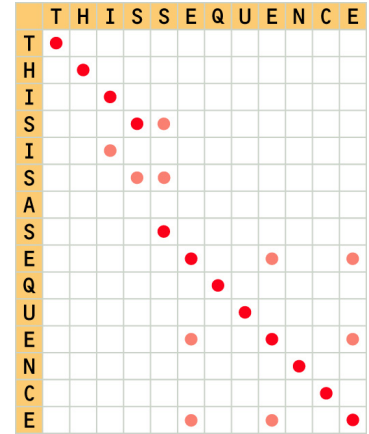


Scoring alignments: identity percentage and similarity percentage

Dot-plots: it is the simplest way to compare sequence similarities.

Use of filters:

- Window size allows to overlap fixed-length windows
- Minimum identity score: it is the minimum identity score fixed for the window previously set.



Two views of dot-plot representations of an SH2 sequence compared to itself. A) Unfiltered dot-plot. The identity is shown by the unbroken diagonale. There is some background noise. B) Dot-plot of the same sequence comparison with a window of 10 residues and a minimum identity score within the window set to 3.

Scoring alignments: identity percentage and similarity percentage

Similarity percentage: it takes into account the so-called conservative substitution

THISISA-SEQUENCE
TH---ATSEQUENCE

THISISA-SEQUENCE
THAT---SEQUENCE

```
gi|66361410|pdb|1ZBM|A -----MGHHHHHSHKIRVAHTPDADD 22
gi|154175534|ref|YP_001409022. -----MKNIKHIDVAHSPDADD 17
gi|6647837|sp|O28098.1|SUCD2_A -----MAIIVDERTKVVVQGITGQGK 22
gi|1711576|sp|P53598.1|SUCA_YE MLRSTVSKASLKICRHFHRESIPYDKTIKNLLLPKDTKVI FQGF TKGQGT 50
. : : : .

gi|66361410|pdb|1ZBM|A AFXFYAXTHGKVDT-WLEIEHVIEDIETLNKRAFNAEYEVTAISAHAYAL 71
gi|154175534|ref|YP_001409022. IFMYMAIKFGWVGSKNLSFTNTALDIQTLNEEALKSTYTATAISFALYPL 67
gi|6647837|sp|O28098.1|SUCD2_A FHTE RMLNYGTKIVAGVTPGKGGTEVLGVPVYDSVKEAVREADANASVIF 72
gi|1711576|sp|P53598.1|SUCA_YE FHASISQEYGTNVVGGTNPKKAQTHLGQPVFASVKDAIKETGATASAI F 100
. * : : : :

gi|66361410|pdb|1ZBM|A LDDKYRILSAGASVGDGYGPVVAKSEISLD-GKRIAVPGRYTTANLLK 120
gi|154175534|ref|YP_001409022. ISDDYALLRCAVSFGEYGPGLIKKRGVNLKRNFKVALSGAHTTNALLFR 117
gi|6647837|sp|O28098.1|SUCD2_A VPAPFAADAVMEAADAGIKVIVCITEGIPVHDELKMYWRVKEAGAT-LIG 121
gi|1711576|sp|P53598.1|SUCA_YE VPPPIAAAIAKESIEAEIPLAVCITEGIPQHDMLYIAEMLQTQDKTRLVG 150
: : : . : . : *

gi|66361410|pdb|1ZBM|A LAVE-DFEPVEXPFDRIIQAVLDEEVDAGLLIHEGQITYADYGLKCVL DL 169
gi|154175534|ref|YP_001409022. AAYP-EARIVYKNFLEIENAVLSGEVDAGVLIHESILGFSS-ELEVEREI 165
gi|6647837|sp|O28098.1|SUCD2_A PNCPGIISPG-KTHLGIMPVQIFKPGNVGIVRSGLTLYQIAYNLTKLGL 170
gi|1711576|sp|P53598.1|SUCA_YE PNCPGIINPATKVRIGIQPPKIFQAGKIGIISRSGLTYEAVQQTTKTDL 200
* : . *:: : . : :

gi|66361410|pdb|1ZBM|A WDWVSEQV--KLPLPLGLNAIRRDLSVEVQEEFLRAXRESIAFAIEN-PD 216
gi|154175534|ref|YP_001409022. WDVWCELAGENLPLPLGGMALRRSLPLTDAIECERVLTKAVATAAHKPF 215
gi|6647837|sp|O28098.1|SUCD2_A GQSTVVVGLGGDRIIGTDFVEVLRLEFDDKETKAVVLVGEIGGRDEEVAAE 220
gi|1711576|sp|P53598.1|SUCA_YE GQSLVIGMGGDAFPGTDFIDALKLFLFLEDETEGIIMLGEIGGKAEIEAAQ 250
: : : : :

gi|66361410|pdb|1ZBM|A EAIEYAX-----KYSRGLDREAKRFAXXVNDYTYNXPESVDAAL 257
gi|154175534|ref|YP_001409022. LSHMLME-----RNLRIDKEKPKIYLNLYANKDSISMNQTLKAL 256
gi|6647837|sp|O28098.1|SUCD2_A FIREMS-----KPVVGVVAGLTAPPK--RMGHAGAIIEGGVGTAE SKI 262
gi|1711576|sp|P53598.1|SUCA_YE FLKEYNFSRSKMPVASFIAGTIVAGQMKGVRMGHSGAIVEGSGTDAESK 300
* : . : .

gi|66361410|pdb|1ZBM|A KKLYEX-----AEAKGLIKMPKLDILRL-- 280
gi|154175534|ref|YP_001409022. NRLF EIGYDQGFYQPIDAHDYLIPT EYNDARFS- 290
gi|6647837|sp|O28098.1|SUCD2_A KALEAAG-----ARVGT PMEVAELVAEIL---- 287
gi|1711576|sp|P53598.1|SUCA_YE QALRDVG-----VAVVESPGYLGQALLDQFAFK 329
: * : : :
```

Scoring alignments: substitution matrices

(A)

C	9																			
S	-1	4																		
T	-1	1	5																	
P	-3	-1	-1	7																
A	0	1	0	-1	4															
G	-3	0	-2	-2	0	6														
N	-3	1	0	-2	-2	0	6													
D	-3	0	-1	-1	-2	-1	1	6												
E	-4	0	-1	-1	-1	-2	0	2	5											
Q	-3	0	-1	-1	-1	-2	0	0	2	5										
H	-3	-1	-2	-2	-2	-2	1	-1	0	0	8									
R	-3	-1	-1	-2	-1	-2	0	-2	0	1	0	5								
K	-3	0	-1	-1	-1	-2	0	-1	1	1	-1	2	5							
M	-1	-1	-1	-2	-1	-3	-2	-3	-2	0	-2	-1	-1	5						
I	-1	-2	-1	-3	-1	-4	-3	-3	-3	-3	-3	-3	1	4						
L	-1	-2	-1	-3	-1	-4	-3	-4	-3	-2	-3	-2	2	2	4					
V	-1	-2	0	-2	0	-3	-3	-3	-2	-2	-3	3	2	1	3	1	4			
F	-2	-2	-2	-4	-2	-3	-3	-3	-3	-3	-1	-3	-3	0	0	0	-1	6		
Y	-2	-2	-2	-3	-2	-3	-2	-3	-2	-1	2	-2	-2	-1	-1	-1	-1	3	7	
W	-2	-3	-2	-4	-3	-2	-4	-4	-3	-2	-2	-3	-3	-1	-3	-2	-3	1	2	11
	C	S	T	P	A	G	N	D	E	Q	H	R	K	M	I	L	V	F	Y	W

(B)

C	9																			
S	-1	3																		
T	-3	2	4																	
P	-3	1	-1	6																
A	-3	1	1	1	3															
G	-5	1	-1	-2	1	5														
N	-5	1	0	-2	0	0	4													
D	-7	0	-1	-2	0	0	2	5												
E	-7	-1	-2	-1	0	-1	1	3	5											
Q	-7	-2	-2	0	-1	-3	0	1	2	6										
H	-4	-2	-3	-1	-3	-4	2	0	-1	3	7									
R	-4	-1	-2	-1	-3	-4	-1	-3	-3	1	1	6								
K	-7	-1	-1	-2	-2	-3	1	-1	-1	0	-2	2	5							
M	-6	-2	-1	-3	-2	-4	-3	-4	-4	-1	-4	-1	0	8						
I	-3	-2	0	-3	-1	-4	-2	-3	-3	-3	-4	-2	-2	1	6					
L	-7	-4	-3	-3	-3	-5	-4	-5	-4	-2	-3	-4	-4	3	1	5				
V	-2	-2	0	-2	0	-2	-3	-3	-3	-3	-3	-3	-4	1	3	1	5			
F	-6	-3	-4	-5	-4	-5	-4	-7	-6	-6	-2	-4	-6	-1	0	0	-3	8		
Y	-1	-3	-3	-6	-4	-6	-2	-5	-4	-5	-1	-6	-6	-4	-2	-3	-3	4	8	
W	-8	-2	-6	-7	-7	-8	-5	-8	-8	-6	-5	1	-5	-7	-7	-5	-8	-1	-1	12
	C	S	T	P	A	G	N	D	E	Q	H	R	K	M	I	L	V	F	Y	W

Expectation value (E-value): the probability of different alignments with scores equivalent to or better than S that are expected to occur in a database search by chance. The lower the E value, the more significant the score.

It indicates the number of sequences that would be expected to have that score (or more) if the query sequence were compared against a database containing no sequences related to the query sequence. Thus, a lower E-value indicates that the sequences are more likely to be related than if the comparison had a higher E-value. An E-value of 0.00001 or less (also sometimes written as $1e-5$, which is shorthand for 1.0×10^{-5}) is often used as good initial evidence that a query and database sequence are related, although further investigation should always be carried out to obtain additional support for such a hypothesis.

Amino acids substitution scoring matrices. A) The BLOSUM-62 matrix and B) the PAM120 matrix. The colored shading indicates different physicochemical properties of the residues.

Sequence alignments: BLAST



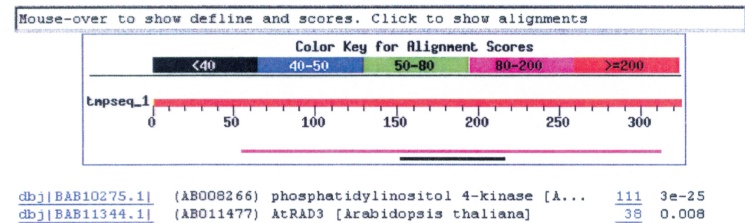
results of **BLAST**

Description	BLAST equivalent
Protein compared to protein database or DNA to DNA database. For protein, ktup = 2 by default (ktup = 1 is more sensitive); default for DNA is 6; 4 or 3 is more sensitive. 1 should be used for short DNA stretches.	blastp/blastn
Uses Smith-Waterman algorithm. Can search protein to protein or DNA to DNA. Can be more sensitive than fasta with protein sequences.	
DNA compared to protein database. DNA translated into all three frames. faster than fastx but better. Used to see if DNA encodes a protein.	blastx
Protein compared to DNA database. Mainly used to identify EST sequences. This is preferred over fastx as protein comparison is more sensitive than DNA.	tblastn (tblastx compares translated DNA to translated DNA database)
Mixed peptide sequence (such as obtained by Edman degradation) compared to protein database.	
Mixed peptide sequence compared to DNA database.	

(A)

sp P32871 P11A	BOVIN	PHOSPHATIDYLINOSITOL 3-KINASE CATALYTI...	680	0.0
sp P42336 P11A	HUMAN	PHOSPHATIDYLINOSITOL 3-KINASE CATALYT...	676	0.0
sp P42337 P11A	MOUSE	PHOSPHATIDYLINOSITOL 3-KINASE CATALYT...	674	0.0
sp P42338 P11B	HUMAN	PHOSPHATIDYLINOSITOL 3-KINASE CATALYT...	338	9e-93
sp O35904 P11D	MOUSE	PHOSPHATIDYLINOSITOL 3-KINASE CATALYT...	332	7e-91
sp O00329 P11D	HUMAN	PHOSPHATIDYLINOSITOL 3-KINASE CATALYT...	331	2e-90
↓				
sp P47473 R1R1	MYCGE	RIBONUCLEOSIDE-DIPHOSPHATE REDUCTASE A...	34	0.59

(B) **Distribution of 2 Blast Hits on the Query Sequence**



(C) ... This CD alignment includes 3D structure. To display structure, download [Cn3D v3.00!](#)

Mouse-over boxes to display more information

Sequences producing significant alignments:

	Score	E
	(bits)	value
• gnl Smart PI3Kc	Phosphoinositide 3-kinase, catalytic domain, Phosphoinositide ...	301 3e-83
• gnl Pfam pfam00454	PI3_P14_kinase, Phosphatidylinositol 3- and 4-kinases	263 9e-72

• [gnl|Smart|PI3Kc](#), Phosphoinositide 3-kinase, catalytic domain, Phosphoinositide 3-kinase isoforms participate in a variety of processes, including cell motility, the Ras pathway, vesicle trafficking and secretion, and apoptosis. These homologues may be either lipid kinases and/or protein kinases: the former phosphorylate the 3-position in the inositol ring of inositol phospholipids. The ataxia telangiectesia-mutated gene produced, the targets of rapamycin (TOR) and the DNA-dependent kinase have not been found to possess lipid kinase activity. Some of this family possess PI-4 kinase activities.

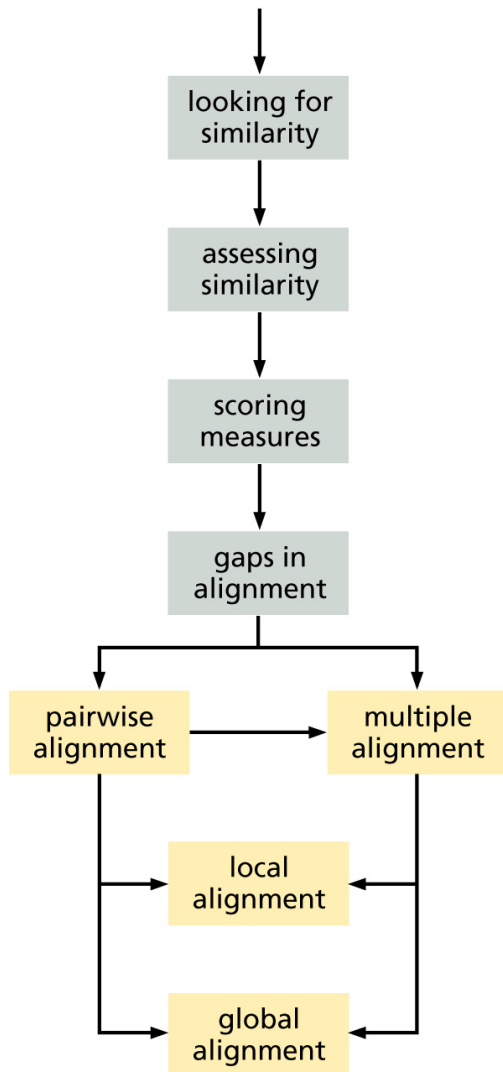
Add query to multiple alignment, display sequences

Length = 265
Score = 301 bits (763), Expect = 3e-83

Query: 19 IIFKNGDDLRODMLTLQIRIHENIWCNQLDLRMLPYGCLSGDCVGLIEUVNRSHTIM 78
Sbjct: 2 IIFKNGDDLRODMLILQILRIMESIVETESLDLCLLPYGCISTGDKIGNIEIVKDATTTIA 61

Types of alignments

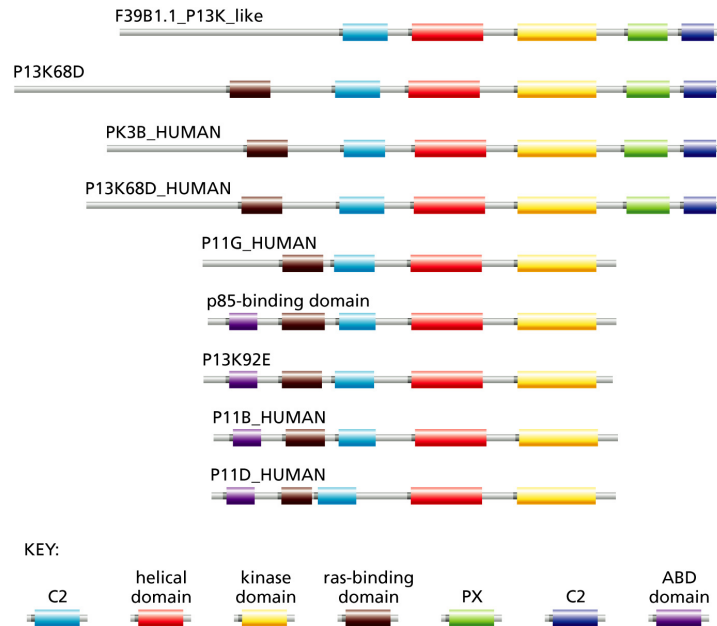
PRODUCING AND ANALYZING SEQUENCE ALIGNMENTS



Global alignment: it is used to find or compare closely related sequences that are similar over their whole sequence.

Local alignment: can reveal that parts of sequences are related.

It is useful in multidomain proteins.



P13-kinase is a multidomain protein. Output from Pfam.

Multiple alignments

They can be constructed by different techniques.

(A) structural/functional alignment from BALiBase

```
1 csy SHEKMPWFHGKISRREESEQIVLIGSKTNGKFLIRARD--NNGSYALCCLLHEGKVLHYRIDKDKTGKLSIPEGK-KFDTLWQLVEHYSYK-----DGLLRVL-TVPCQK
1 gri EMKPHPWFFGKIPRAKAEEML-SKQRHDGAFIIRSESA-APGDFSLSVKFGNDVQHFVKVLRDAGAGYFL-WVV-KFNSLNELVDYHRSTS-VSRNQIFLRDIEQVPPQQ-
1 aya ---MRRWFHPNITGVEAENLLLRGVDGSLFARPSKSN-PPGDFTLVRRNGAVTHIKIQNTGDYDLYGGGEKFA-TLAELVQYMEHHGQLKEKNGDVIEL-KYPLN-
2 pna -LQDAEWYWGDISREEVNEKLRD--ADGTFVLVRDASTKMHGDYTLTLRKGKGNKLIKIFHRD-GKYGFSDDL-TFNSVVELINHYRNES-LAQYNPKLDVKL-LYPVS-
1 bfi HHDEKTNVVGSSNRNKAENLLRGR--RDGTFVLVRES--KQGCYACSVVVDGVEVKHCVINKTATG-YGFAEPYNYLSSLKELVLHYQHTS-LVQHNDSLNVTLA-AYPVYA
```

(B) DIALIGN multiple sequence alignment

```
1 csy SHEKMPWFHGKISRREESEQIVLIGSKT-NGKFLIRAR-DN--NGSYALCCLLHEGKVLHYRIDKDKTGKLSIPEGK-K-FDTLWQLVEHYSYK-----DGLLRVL-TVPCQK
1 gri EMKPHPWFFGKIPRAKAEEML--SKQRHDGAFIIRSESA--PGDFSLSVKFGNDVQHFVKVLRDAGAGYFLWVV-K-FNSLNELVDYHRST--SVSRNQIFLRDIEQVPPQQ-
1 aya M---MRRWFHPNITGVEAENLLLRGVDGSLFARPSKSN--PGDFTLVRRNGAVTHIKIQNTGDYDLYG-GEK-FATLAELVQYMEHHGQLKEKNGDV-IELK-YPLN-
2 pna LQDAE-WYWGDISREEVNEKL--RDTA-DGTFVLVRDA-STKMHGDYTLTLRKGKGNKLIKIFHRD-GKYGFSDDL-TFNSVVELINHYRNE--SLAQYNPKLDVKL-LYPVS-
1 bfi HHDEKTNVVGSSNRNKAENLL--RGR-DGTFVLVRES-SK--QGCYACSVVVDGVEVKHCVINKTATGYGFAE-PYNYLSSLKELVLHYQHT--SLVQHNDSLNVTLA-AYPVYA
```

(C) ClustalW multiple sequence alignment

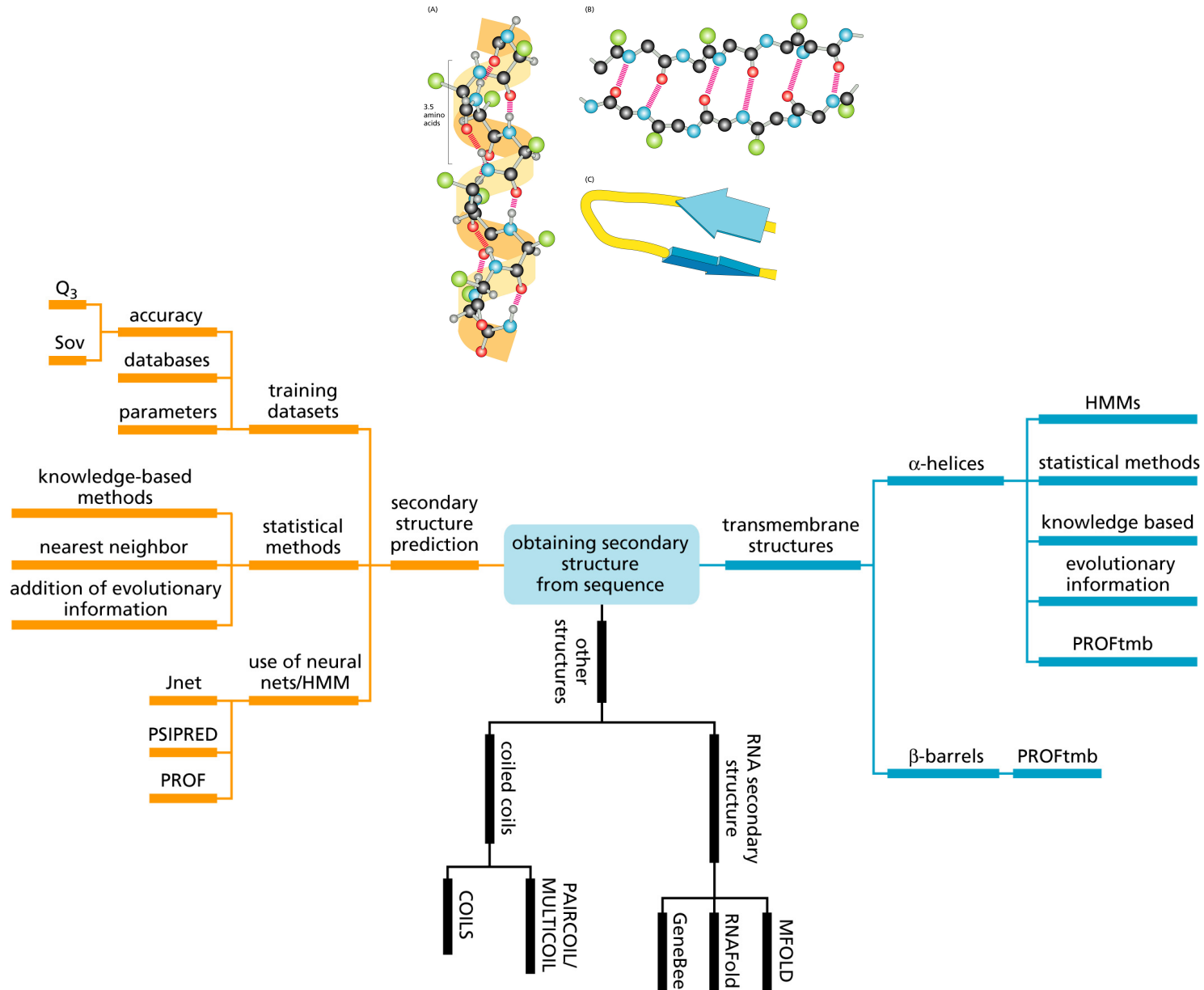
```
1 csy SHEKMPWFHGKISRREESEQIVLIGSKTNGKFLIRARDN--NNGSYALCCLLHEGKVLHYRIDKDKTGKLSIPEGK-KFD-TLWQLVEHYSYK-----ADGLLRVLTVPCQK
1 gri EMKPHPWFFGKIPRAKAEEMLSKQRHDGAFIIRSESA-APGDFSLSVKFGNDVQHFVKVLRDAGAGY-FLWVVK-FNSLNELVDYHRSTS-VSRNQIFLRDIEQVPPQQ-
1 aya ---MRRWFHPNITGVEAEN-LLLRGVDGSLFARPSKSN-PPGDFTLVRRNGAVTHIKIQNTGDYDLYGGGEKFA-TLAELVQYMEHHGQLKEKNGDVIELKYPLN-
2 pna -LQDAEWYWGDISREEVNEKL--EKL-RDADGTFVLVRDASTKMHGDYTLTLRKGKGNKLIKIFHRD-GKYGFSDDL-TFNSVVELINHYRNES-LAQYNPKLDVKL-LYPVS-
1 bfi HHDEKTNVVGSSNRNKAEM--NLLRGRDGTFLVRESK--QGCYACSVVVDGVEVKHCVINKT-ATGYGFAEPYNYLSSLKELVLHYQHTS-LVQHNDSLNVTLA-AYPVYA
```

(D) divide-and-conquer multiple sequence alignment

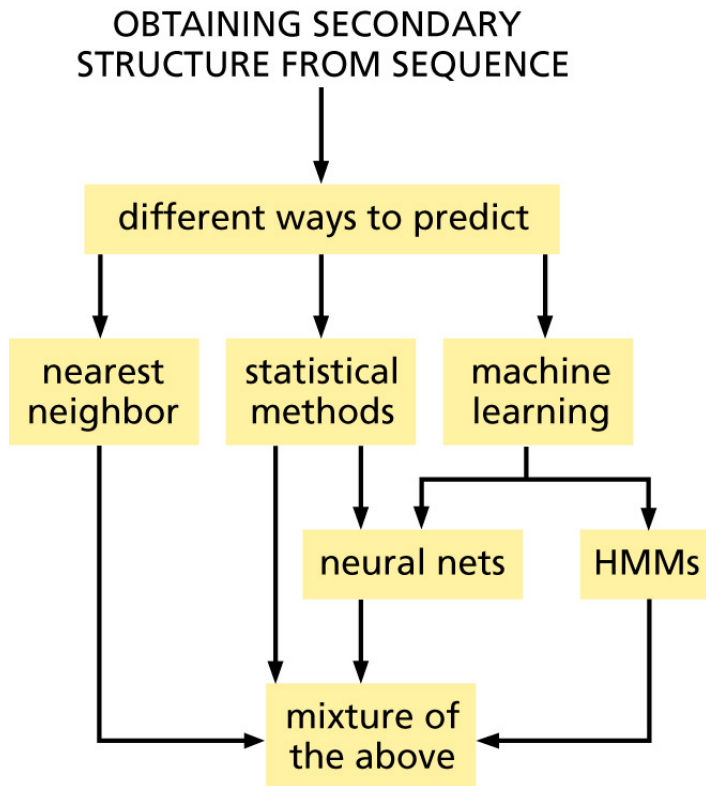
```
1 csy SHEKMPWFHGKISRREESEQIVLIGSKTNGKFLIRA-RDNN-GSYALCCLLHEGKVLHYRIDKDKTGKLSIPEGK-K-FDTLWQLVEHY-SY----KADGLLRV-L-TVPCQK
1 gri EMKPHPWFFGKIPRAKAEEMLS-KQRHDGAFIIRSESA-APGDFSLSVKFGNDVQHFVKVLRDAGAGY-FLWVVK-FNSLNELVDYH-RSTSVSRNQIFLRDIEQVPPQQ-
1 aya ---MRRWFHPNITGVEAENLL-TRGVDGSLFARP-SKSNPPGDFTLVRRNGAVTHIKIQNTGDY-DLYGGGEK-FATLAELVQYMEHHGQLKEKNGDVIEL-KYPLN-
2 pna -LQDAEWYWGDISREEVNEKL--RDTADGTFVLVRDASTKMHGDYTLTLRKGKGNKLIKIFHRD-GKYGFSDDL-TFNSVVELINHY-RNESLAQYNPKLDVKL-LYPVS-
1 bfi HHDEKTNVVGSSNRNKAENLL--RGRDGTFLVRESK-QGCYACSVVVDGVEVKHCVINKTATGY-GFAEPYNYLSSLKELVLHY-QHTSLVQHNDSLNVTLA-AYPVYA
```

Structural alignments: if the structure of one of the proteins is known, then the gap penalty can be increased for regions of known secondary structure such as helices and strands, as these regions are less likely to suffer insertions or deletions. This will mean that few or no gaps are introduced into these regions.

Protein secondary structure prediction



Types of secondary structure prediction



Statistical methods are based on rules that give the probability that a residue will form part of a particular secondary structure.

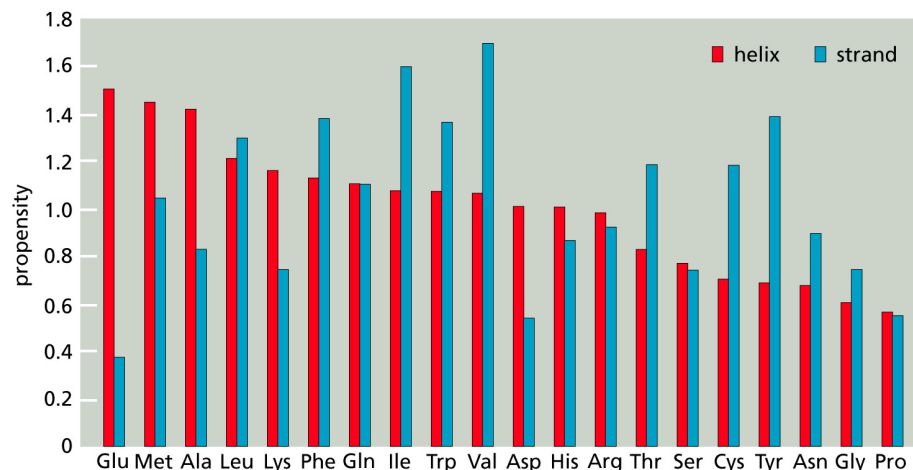
The probabilities are derived from analysing structure and sequence data from large sets of proteins of known structure.

Nearest neighbor methods are statistical methods that incorporate additional information about protein structure (shapes, sizes and physicochemical properties of the different amino acid residues).

Machine learning approaches train a neural net or other learning algorithms to acquire structure-sequence relationships which can then be applied to predict structure from a protein sequence.

Statistical and knowledge-based methods: Chou and Fasman

A.A.	P(a)	P(b)	P(turn)	f(i)	f(i+1)	f(i+2)	f(i+3)
Alanine	142	83	66	0.060	0.076	0.035	0.058
Arginine	98	93	95	0.070	0.106	0.099	0.085
Asparagine	67	89	156	0.161	0.083	0.191	0.091
Aspartic acid	101	54	146	0.147	0.110	0.179	0.081
Cysteine	70	119	119	0.149	0.050	0.117	0.128
Glutamic acid	151	37	74	0.056	0.060	0.077	0.064
Glutamine	111	110	98	0.074	0.098	0.037	0.098
Glycine	57	75	156	0.102	0.085	0.190	0.152
Histidine	100	87	95	0.140	0.047	0.093	0.054
Isoleucine	108	160	47	0.043	0.034	0.013	0.056
Leucine	121	130	59	0.061	0.025	0.036	0.070
Lysine	114	74	101	0.055	0.115	0.072	0.095
Methionine	145	105	60	0.068	0.082	0.014	0.055
Phenylalanine	113	138	60	0.059	0.041	0.065	0.065
Proline	57	55	152	0.102	0.301	0.034	0.068
Serine	77	75	143	0.120	0.139	0.125	0.106
Threonine	83	119	96	0.086	0.108	0.065	0.079
Tryptophan	108	137	96	0.077	0.013	0.064	0.167
Tyrosine	69	147	114	0.082	0.065	0.114	0.125
Valine	106	170	50	0.062	0.048	0.028	0.053



Chou-Fasman is one of most commonly used algorithms

- measured frequencies at which each amino acid appeared in particular types of secondary sequences in a set of proteins of known structure
- assigns the amino acids three conformational parameters based on the frequency at which they were observed in alpha helices, beta sheets and beta turns
 1. **P(a) = propensity to form alpha helices**
 2. **P(b) = propensity to form beta sheets**
 3. **P(turn) = propensity to form beta turns**
- also assigns 4 turn parameters based on frequency at which they were observed in the first, second, third or fourth position of a beta turn
 1. **f(i) = probability of being in position 1**
 2. **f(i+1) = probability of being in position 2**
 3. **f(i+2) = probability of being in position 3**
 4. **f(i+3) = probability of being in position 4**

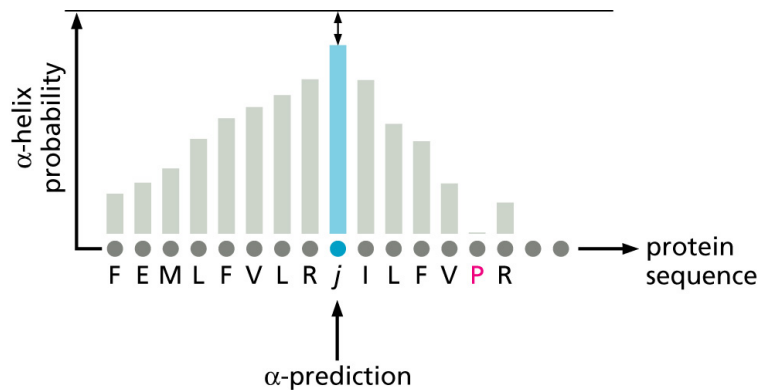
Statistical and knowledge-based methods: Chou and Fasman

identifies helix and sheet "nuclei", then applies a set of heuristic rules to determine if these clusters of amino acids are sufficient to nucleate a region of alpha-helix or beta-sheet.

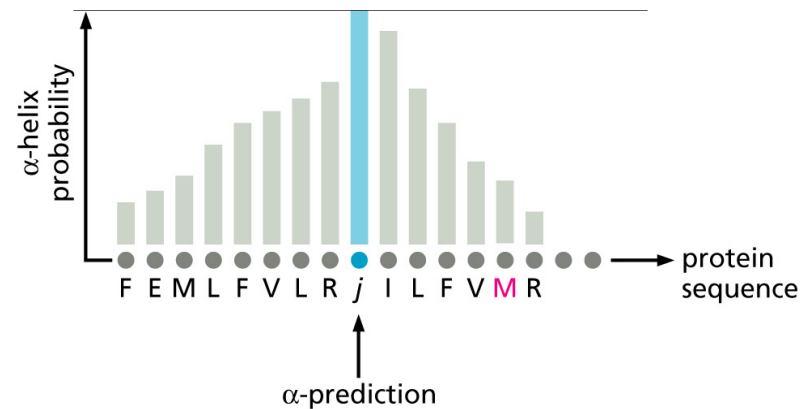
- **helix: 4 out of 6 amino acids with $P(a) > 100$**
 - extends the nucleus in each direction until reach four amino acids in a row with $P(a) < 100$
 - for each of these regions, add up all the $P(a)$ and all the $P(b)$ values.
 - If the total $P(a)$ is larger than the total of $P(b)$ and the run is more than 5 amino acids long, then it is predicted to be alpha helix
- **sheet: 4 out of 6 amino acids with $P(b) > 100$ (some people use 3 out of 5).**
 - extends the nucleus in each direction until reach four amino acids in a row with $P(b) < 100$
 - for each of these regions, add up all the $P(a)$ and all the $P(b)$ values.
 - If the total $P(b)$ is larger than the total of $P(a)$, the run is more than 5 amino acids long, and the average $P(b) > 100$ then it is predicted to be beta sheet.
- **If helices and sheets overlap then compare the total $P(a)$ and total $P(b)$ for the overlapping region. If the total $P(a)$ is larger than the total of $P(b)$ then it is predicted to be alpha helix (and vice-versa)**
- **beta turn**
 - calculate the likelihood of a turn $P(t)$ for amino acid at position i as the sum of $f(i)$ + the $f(i+1)$ value for the following amino acid + the $f(i+2)$ value for the next amino acid + the $f(i+3)$ value for the amino acid at the plus three position.
 - Predict a beta- turn at position i if the following criteria are met:
 - the calculated $P(t)$ is > 0.5
 - the average $P(\text{turn})$ for amino acids i to $i+3$ is > 100
 - the sum of the $P(\text{turn})$ values for amino acids i to $i+3$ is larger than the sum of the $P(a)$ and $P(b)$ values
- **Accuracy = 50-85%, depending on the protein**

Statistical and knowledge-based methods: GOR

It incorporates the effects of local interactions between amino acids residues by taking successive windows of 17 residues and considering the effect of residues from position $j-8$ to $j+8$ on the conformation of the residue at position j .



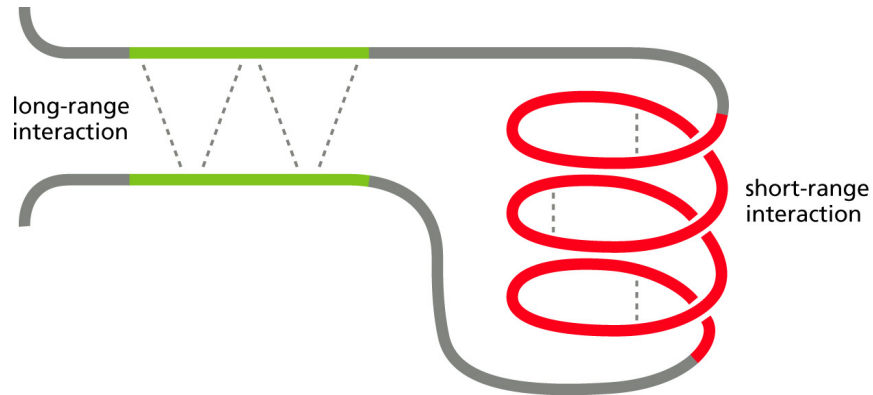
The effect of an helix breaker (Pro) at position $j+5$. The proline diminishes the overall additive propensity of residue j to form helix



The effect of a non helix breaker (Met) at position $j+5$. The methionine improves the overall additive propensity of residue j to form helix

Nearest neighbor methods

The formation of secondary structure in proteins does not only depend on local interactions (beta-sheets are made up of beta-strands that are separated from some distance in the polypeptide chain).



Neural networks methods

The algorithm will learn by iterative changes to its parameters until the predicted structure is as similar to the observed structure as possible.

PSIPRED is a three stage method:

1. It generates a multiple sequence alignment
2. It generates an initial secondary structure
3. It filters the initial prediction

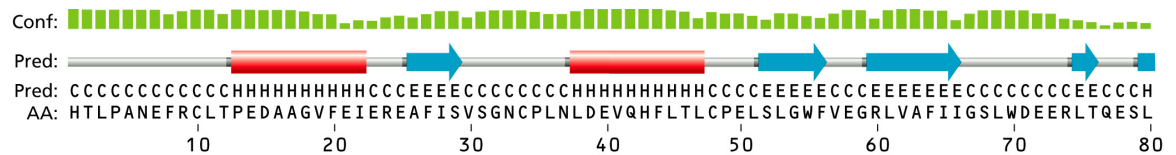
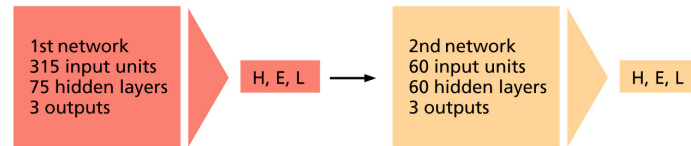
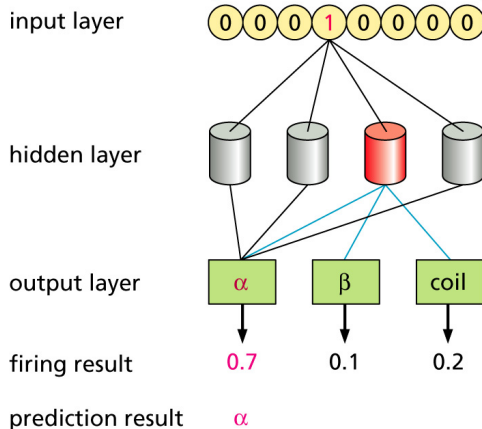
A	R	N	D	C	Q	E	G	H	I	L	K	M	F	P	S	T	W	Y	V
-3	-4	-4	-4	-3	-4	-4	-4	-2	-1	-1	-4	-1	8	-5	-3	-3	0	2	-2
0	-1	-1	3	-4	3	4	1	-1	-4	-4	0	-3	-4	-2	-1	-2	-4	-3	-3
0	-1	2	1	-3	4	0	-1	-2	-4	-3	1	-2	-4	-2	0	-4	-3	-3	-2
-2	-3	-4	-5	-2	-3	-4	-6	-4	0	6	0	0	-1	4	-2	2	0	-4	-2
0	-3	-1	-2	0	-2	4	-3	-3	0	-2	-2	-4	-3	3	1	-4	-4	-4	-3
0	2	0	4	-4	1	2	1	-2	-3	-4	0	-3	-4	-3	1	-2	-5	-4	-4
-1	5	3	-2	-4	-1	-1	1	-2	-1	-4	1	-3	-4	-3	1	-2	-5	-4	-4
-2	-3	-4	-5	-3	-3	-4	-5	-4	3	4	-1	1	2	-4	-3	-2	-3	-1	0
-2	3	2	-3	-4	2	1	-3	-2	-3	-3	1	1	-4	-3	2	1	-4	-3	-1
0	2	3	1	-4	0	0	0	-2	-4	-4	1	-2	0	-5	-4	0	0	-4	-4
5	-3	-3	-3	-2	-3	-3	-2	-3	1	-2	-3	-2	1	-3	0	1	-4	-2	0
-1	-4	-5	-5	-3	-4	-4	-5	-4	3	3	-4	2	3	-5	-3	-2	5	-1	2
0	3	3	0	-4	3	0	1	-2	-4	-4	1	-3	-4	-3	1	-1	-4	-3	-4
0	3	3	0	-4	3	0	1	-2	-4	-4	1	-3	-4	-3	1	-1	-4	-3	-3
-1	0	1	0	-4	1	-1	-2	-4	-3	5	-2	0	-3	0	-2	-1	0	0	-3
-1	1	3	-2	-4	0	-2	4	-2	-4	-4	0	-3	0	0	-3	0	-3	0	-4

window of 15 rows

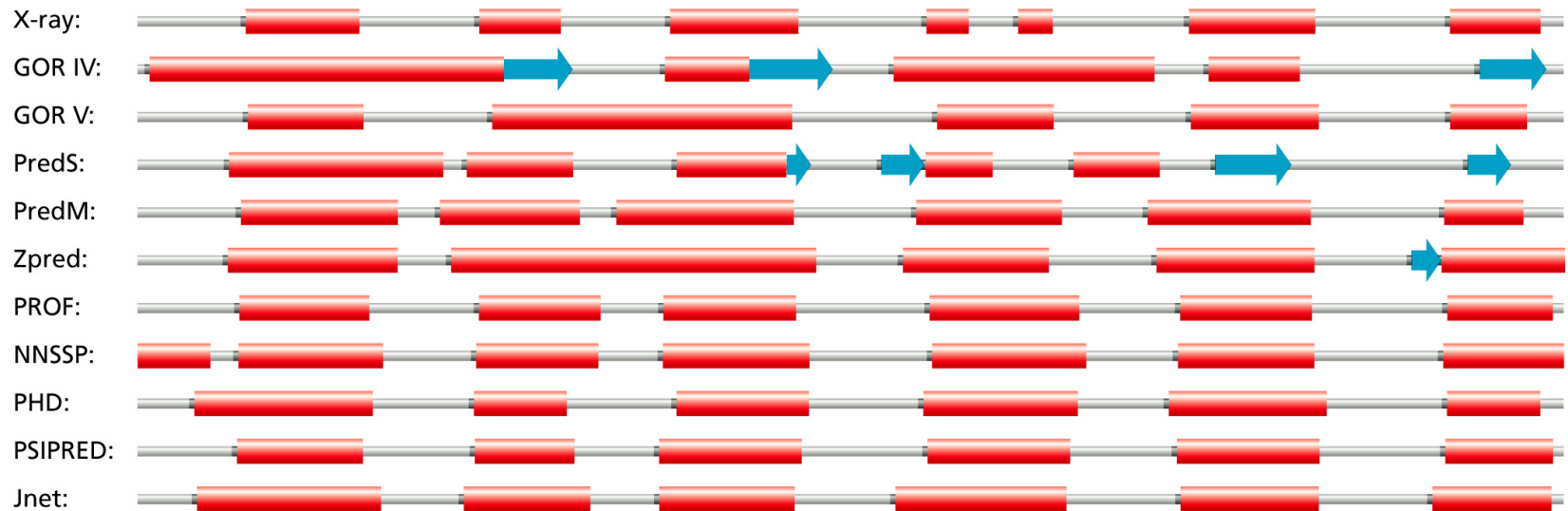
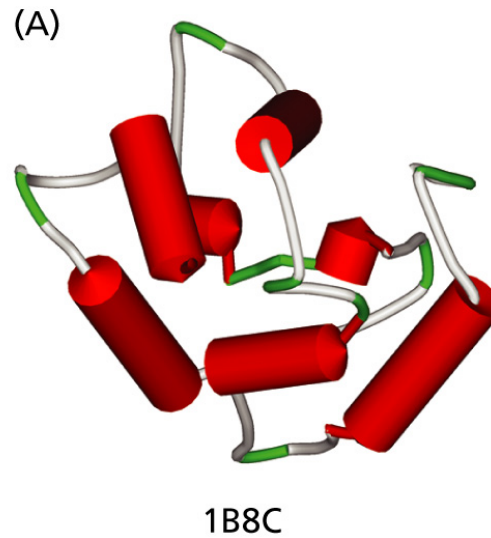
A	R	N	D	C	Q	E	G	H	I	L	K	M	F	P	S	T	W	Y	V
0.4	0.3	0.3	0.3	0.2	0.9	0.3	0.3	0.4	0.4	0.4	0.4	0.4	0.9	0.1	0.4	0.4	0.5	0.7	0.4
0.3	0.2	0.3	0.8	0.4	0.3	0.7	0.1	0.6	0.2	0.3	0.3	0.5	0.2	0.1	0.4	0.8	0.2	0.3	0.2
0.1	0.1	0.4	0.3	0.5	0.1	0.1	0.3	0.1	0.1	0.4	0.2	0.4	0.9	0.3	0.4	0.4	0.9	0.3	0.6
.
.
.

15 x 20 scaled inputs to 1st network

N-terminal...THISISAHIDDENMESSAGE...C-terminal

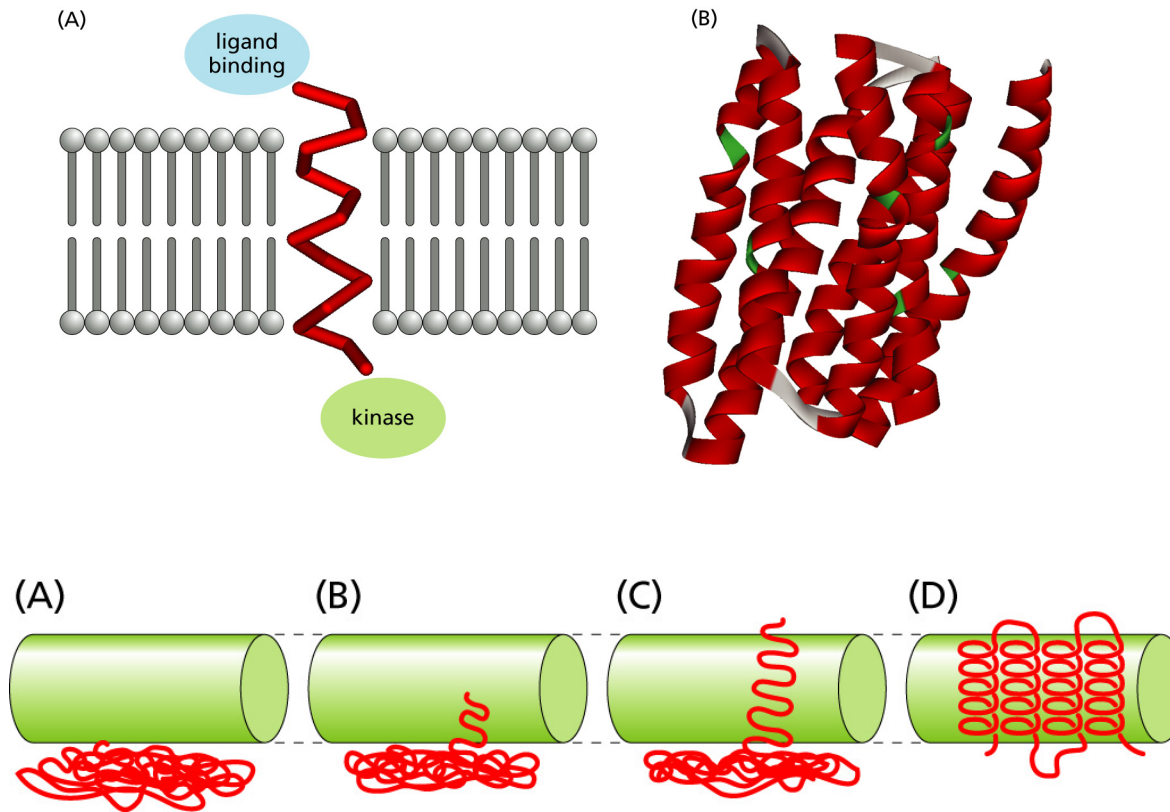


Secondary structure prediction methods



Transmembrane proteins

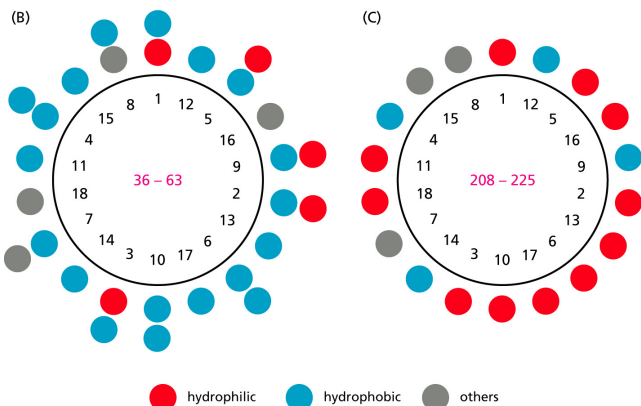
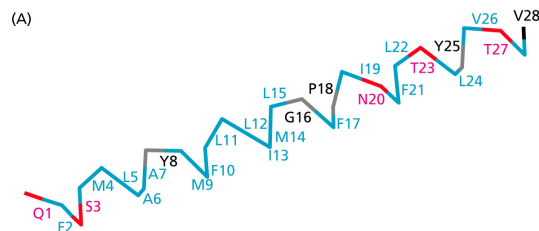
Membrane proteins are functionally important. For example, the receptors are formed by 1 or more helices spanning the membrane



The four main ways in which proteins may be attached to a membrane. A) Attachment by interactions between the protein and the cytosolic face of the lipid bilayer. B) Attachment via an anchor (lipidic or terminals of the protein) that are added post-translationally. C) Transmembrane proteins have part of the protein chain embedded in the lipid bilayer. D) Transmembrane proteins where the protein chain threads back and forth across the membrane multiple times.

Transmembrane proteins

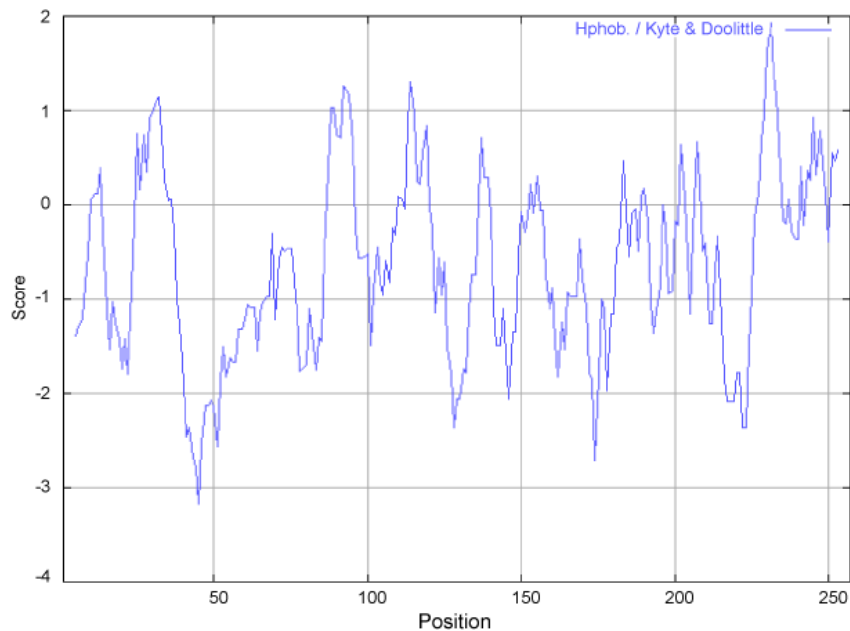
Helix wheel



Hydrophobicity diagram

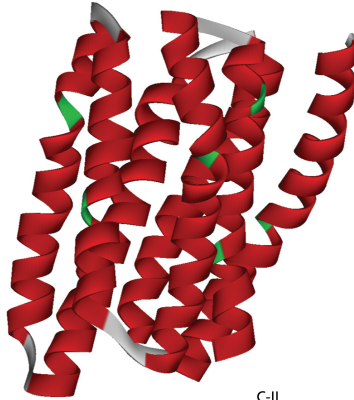
Using the scale **Hphob. / Kyte & Doolittle**, the individual values for the 20 amino acids are:

Ala: 1.800	Arg: -4.500	Asn: -3.500	Asp: -3.500	Cys: 2.500	Gln: -3.500
Glu: -3.500	Gly: -0.400	His: -3.200	Ile: 4.500	Leu: 3.800	Lys: -3.900
Met: 1.900	Phe: 2.800	Pro: -1.600	Ser: -0.800	Thr: -0.700	Trp: -0.900
Tyr: -1.300	Val: 4.200	: -3.500	: -3.500	: -0.490	

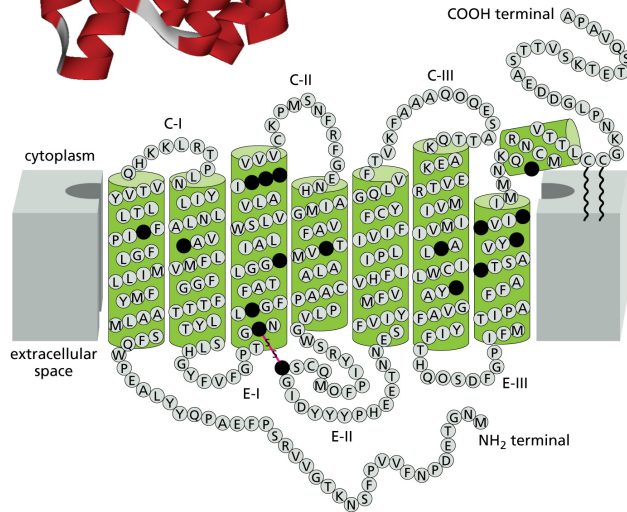


Transmembrane proteins

(A)



(B)



X-RAY	MNGTEGPNFY	VVFSNKTGVV	RSPFEAPQYY	LAEPWQFSML	AAYMFLLIIML	50
HMSTOP	MNGTEGPNFY	VVFSNKTGVV	RSPFEAPQYY	LAEPWQFSML	AAYMFLLIIML	
SOSUI	MNGTEGPNFY	VVFSNKTGVV	RSPFEAPQYY	LAEPWQFSML	AAYMFLLIIML	
DAS	mngtegnpny	vpfsnktgvv	rspfeapqyy	laepwqfsm	aaymflliiml	
TMHMM	MNGTEGPNFY	VVFSNKTGVV	RSPFEAPQYY	LAEPWQFSML	AAYMFLLIIML	
TMpred	MNGTEGPNFY	VVFSNKTGVV	RSPFEAPQYY	LAEPWQFSML	AAYMFLLIIML	
PHDhtm	MNGTEGPNFY	VVFSNKTGVV	RSPFEAPQYY	LAEPWQFSML	AAYMFLLIIML	
TMAP	MNGTEGPNFY	VVFSNKTGVV	RSPFEAPQYY	LAEPWQFSML	AAYMFLLIIML	

X-RAY	GFPINFLTLY	VTVQHKKLRT	PLNYILLNLA	VADLFMVVFGG	FTTTYLTSLH	100
HMSTOP	GFPINFLTLY	VTVQHKKLRT	PLNYILLNLA	VADLFMVVFGG	FTTTYLTSLH	
SOSUI	GFPINFLTLY	VTVQHKKLRT	PLNYILLNLA	VADLFMVVFGG	FTTTYLTSLH	
DAS	GFPINFLTLY	vtvqhkklrt	pLnyILLNLA	VADLFMVVFGG	FTTTYLtslh	
TMHMM	GFPINFLTLY	VTVQHKKLRT	PLNYILLNLA	VADLFMVVFGG	FTTTYLTSLH	
TMpred	GFPINFLTLY	VTVQHKKLRT	PLNYILLNLA	VADLFMVVFGG	FTTTYLTSLH	
PHDhtm	GFPINFLTLY	VTVQHKKLRT	PLNYILLNLA	VADLFMVVFGG	FTTTYLTSLH	
TMAP	GFPINFLTLY	VTVQHKKLRT	PLNYILLNLA	VADLFMVVFGG	FTTTYLTSLH	

X-RAY	GYFVFGPTGC	NLEGGFATLG	GEIALWSLVV	LAIERYVVVC	KPMSNFRFGE	150
HMSTOP	GYFVFGPTGC	NLEGGFATLG	GEIALWSLVV	LAIERYVVVC	KPMSNFRFGE	
SOSUI	GYFVFGPTGC	NLEGGFATLG	GEIALWSLVV	LAIERYVVVC	KPMSNFRFGE	
DAS	gyfvfgptgc	nlegffatlg	geIALWSLVV	LAIERYVvvc	kpmnsnfrfge	
TMHMM	GYFVFGPTGC	NLEGGFATLG	GEIALWSLVV	LAIERYVVVC	KPMSNFRFGE	
TMpred	GYFVFGPTGC	NLEGGFATLG	GEIALWSLVV	LAIERYVVVC	KPMSNFRFGE	
PHDhtm	GYFVFGPTGC	NLEGGFATLG	GEIALWSLVV	LAIERYVVVC	KPMSNFRFGE	
TMAP	GYFVFGPTGC	NLEGGFATLG	GEIALWSLVV	LAIERYVVVC	KPMSNFRFGE	

X-RAY	NHAIMGVAFT	WVMALACAAP	PLVGWSRYIP	EGMQCSCGID	YYTPHEETNN	200
HMSTOP	NHAIMGVAFT	WVMALACAAP	PLVGWSRYIP	EGMQCSCGID	YYTPHEETNN	
SOSUI	NHAIMGVAFT	WVMALACAAP	PLVGWSRYIP	EGMQCSCGID	YYTPHEETNN	
DAS	nhaimgvaft	wvmalacaap	plvgwsryip	egmqcscgid	yytphheetnn	
TMHMM	NHAIMGVAFT	WVMALACAAP	PLVGWSRYIP	EGMQCSCGID	YYTPHEETNN	
TMpred	NHAIMGVAFT	WVMALACAAP	PLVGWSRYIP	EGMQCSCGID	YYTPHEETNN	
PHDhtm	NHAIMGVAFT	WVMALACAAP	PLVGWSRYIP	EGMQCSCGID	YYTPHEETNN	
TMAP	NHAIMGVAFT	WVMALACAAP	PLVGWSRYIP	EGMQCSCGID	YYTPHEETNN	

X-RAY	ESFVIYMFVV	HFIIPLIVIF	FCYGQLVFTV	KEAAAQQQES	ATTQKAEKEV	250
HMSTOP	ESFVIYMFVV	HFIIPLIVIF	FCYGQLVFTV	KEAAAQQQES	ATTQKAEKEV	
SOSUI	ESFVIYMFVV	HFIIPLIVIF	FCYGQLVFTV	KEAAAQQQES	ATTQKAEKEV	
DAS	esfviymfvv	hfiiplivif	fcygqlvftv	keaaaqqqes	attqkaekv	
TMHMM	ESFVIYMFVV	HFIIPLIVIF	FCYGQLVFTV	KEAAAQQQES	ATTQKAEKEV	
TMpred	ESFVIYMFVV	HFIIPLIVIF	FCYGQLVFTV	KEAAAQQQES	ATTQKAEKEV	
PHDhtm	ESFVIYMFVV	HFIIPLIVIF	FCYGQLVFTV	KEAAAQQQES	ATTQKAEKEV	
TMAP	ESFVIYMFVV	HFIIPLIVIF	FCYGQLVFTV	KEAAAQQQES	ATTQKAEKEV	

X-RAY	TRMVIIMVIA	FLICWLPYAG	VAFYIFTHQG	SDFGPIFMTI	PAFFAKTSAV	300
HMSTOP	TRMVIIMVIA	FLICWLPYAG	VAFYIFTHQG	SDFGPIFMTI	PAFFAKTSAV	
SOSUI	TRMVIIMVIA	FLICWLPYAG	VAFYIFTHQG	SDFGPIFMTI	PAFFAKTSAV	
DAS	tRmviimvia	flicwlpYag	vafyifthqg	sdfgpifmti	paFFaktsav	
TMHMM	TRMVIIMVIA	FLICWLPYAG	VAFYIFTHQG	SDFGPIFMTI	PAFFAKTSAV	
TMpred	TRMVIIMVIA	FLICWLPYAG	VAFYIFTHQG	SDFGPIFMTI	PAFFAKTSAV	
PHDhtm	TRMVIIMVIA	FLICWLPYAG	VAFYIFTHQG	SDFGPIFMTI	PAFFAKTSAV	
TMAP	TRMVIIMVIA	FLICWLPYAG	VAFYIFTHQG	SDFGPIFMTI	PAFFAKTSAV	

X-RAY	YNPVIYIMMN	KQFRNCMVTT	LCCGKNPLGD	DEASTTVSKT	ETSQVAPA	348
HMSTOP	YNPVIYIMMN	KQFRNCMVTT	LCCGKNPLGD	DEASTTVSKT	ETSQVAPA	
SOSUI	YNPVIYIMMN	KQFRNCMVTT	LCCGKNPLGD	DEASTTVSKT	ETSQVAPA	
DAS	ynpviyimn	kqfrncmvt	lccgknplgd	deasttvskt	etsqvapa	
TMHMM	YNPVIYIMMN	KQFRNCMVTT	LCCGKNPLGD	DEASTTVSKT	ETSQVAPA	
TMpred	YNPVIYIMMN	KQFRNCMVTT	LCCGKNPLGD	DEASTTVSKT	ETSQVAPA	
PHDhtm	YNPVIYIMMN	KQFRNCMVTT	LCCGKNPLGD	DEASTTVSKT	ETSQVAPA	
TMAP	YNPVIYIMMN	KQFRNCMVTT	LCCGKNPLGD	DEASTTVSKT	ETSQVAPA	

Protein Sequence Motifs or Patterns

What is required is a method of searching for the occurrence of short sequence patterns, or motifs.

A motif, in general, is any conserved element of a sequence alignment (CONSENSUS), whether composed of a short sequence of contiguous residues or a more distributed pattern. Functionally related sequences will share similar distribution patterns of critical functional residues that are not necessarily contiguous.



Figure 4.15

Residues that contribute to one of the blocks returned by the BLOCKS database after submission of the PI3-kinase p100 α sequence.

(A) A block for four homologous sequences, and (B) for 31 homologous sequences. These representations are called logos, and are computed using a position-specific scoring matrix. This block contains the active-site amino acids and the DFG kinase motif. The size of the letters indicates the level of conservation and the colors indicate physicochemical properties of the residues: acidic, red; basic, blue; small and polar, white; asparagine and glutamine, green; sulfur-containing amino acids, yellow; hydrophobic, black; proline, purple; glycine, gray; aromatic, orange.

Protein Sequence Motifs or Patterns

The PROSITE database is a compilation of motifs and patterns extracted from protein sequences and compiled by inspection of protein families. This database can be searched with an unknown protein sequence to obtain a list of hits to possible patterns or protein signatures.



Database of protein domains, families and functional sites

[Home](#) [ScanProsite](#) [ProRule](#) [Documents](#) [Downloads](#) [Links](#) [Funding](#)

PROSITE consists of [documentation entries](#) describing protein domains, families and functional sites as well as associated [patterns](#) and [profiles](#) to identify them [[More details](#) / [References](#) / [Disclaimer](#) / [Commercial users](#)]. PROSITE is complemented by [ProRule](#), a collection of rules based on profiles and patterns, which increases the discriminatory power of profiles and patterns by providing additional information about functionally and/or structurally critical amino acids [[More details](#)].

Release 20.67, of 03-Nov-2010 (1598 documentation entries, 1308 patterns, 909 profiles and 898 ProRule)

PROSITE access

e.g: PDOC00022, PS50089, SH3, zinc finger

add wildcard '*'

Browse:

- [by documentation entry](#)
- [by ProRule description](#)
- [by taxonomic scope](#)
- [by number of positive hit](#)

PROSITE tools

Scan a sequence against PROSITE patterns and profiles - quick scan

(Output includes graphical view and feature detection)



Enter your sequence or a UniProtKB (Swiss-Prot or TrEMBL) ID or AC [[help](#)]:

- [ScanProsite](#) - advanced scan
- [PRATT](#) - allows to interactively generate conserved patterns from a series of unaligned proteins.
- [MyDomains - Image Creator](#) ^{new} - allows to generate custom domain figures.

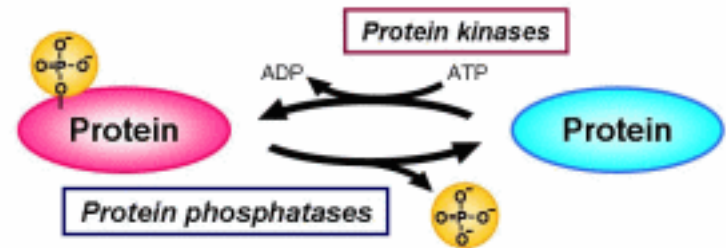
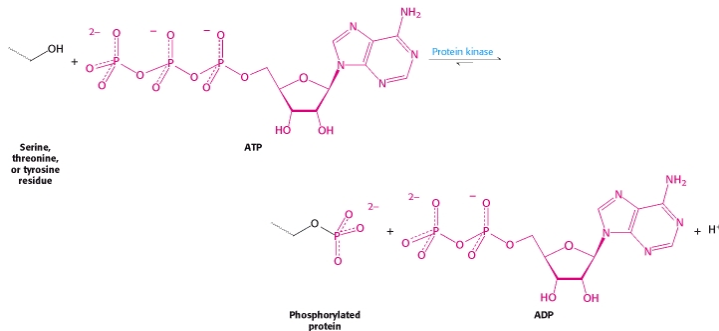


Protein Sequence Motifs or Patterns

Common covalent modifications of protein activity

Modification	Donor molecule	Example of modified protein	Protein function
Phosphorylation	ATP	Glycogen phosphorylase	Glucose homeostasis; energy transduction
Acetylation	Acetyl CoA	Histones	DNA packing; transcription
Myristoylation	Myristoyl CoA	Src	Signal transduction
ADP-ribosylation	NAD	RNA polymerase	Transcription
Farnesylation	Farnesyl pyrophosphate	Ras	Signal transduction
γ -Carboxylation	HCO_3^-	Thrombin	Blood clotting
Sulfation	3'-Phosphoadenosine-5'-phosphosulfate	Fibrinogen	Blood-clot formation
Ubiquitination	Ubiquitin	Cyclin	Control of cell cycle

Copyright © 2002, W. H. Freeman and Company



The consensus sequence recognized by protein kinase A is Arg-Arg-X-Ser-Z or Arg-Arg-X-Thr-Z, in which X is a small residue, Z is a large hydrophobic one, and Ser or Thr is the site of phosphorylation. It should be noted that this sequence is not absolutely required.

Protein Sequence Motifs or Patterns

NetPhos predicts phosphorylation sites in a protein sequence due to kinase acting post-translationally.

```
Name: test1          Length: 26 <-- Sequence name, length
QWERRRITYELVISLIVESYEAHYEAH <-- Submitted sequence
.....T.....SY...Y... <-- Assignments. S,T,Y indicates

                                predicted phosphorylation sites

Ser: 1 Thr: 1 Tyr: 2          <-- No. of predicted S,T,Y phosph. sites
```

Serine predictions

Name	Pos	Context	Score	Pred
test1	13	ELVISLIVE	0.017	.
test1	18	LIVESYEAH	0.942	*S*

Threonine predictions

Name	Pos	Context	Score	Pred
test1	7	ERRRITYELV	0.921	*T*

Tyrosine predictions

Name	Pos	Context	Score	Pred
test1	8	RRRITYELVI	0.056	.
test1	19	IVESYEAHY	0.502	*Y*
test1	23	YEAHYEAH-	0.885	*Y*

