

Computer Aided Drug Design

—Sequence Analysis

Qin Xu

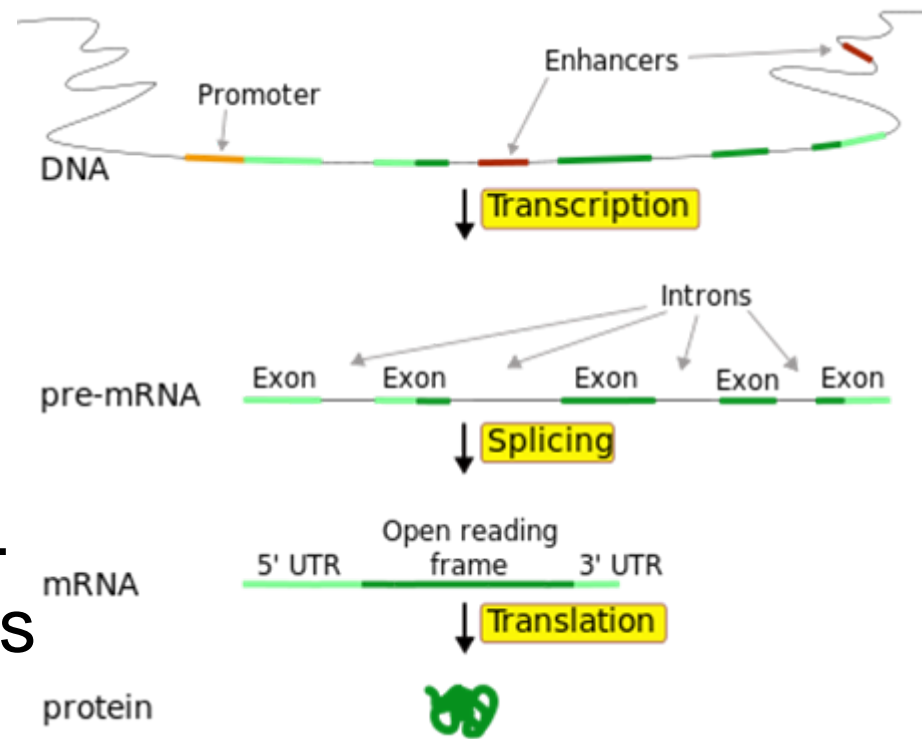
<http://cbb.sjtu.edu.cn/~qinxu/CADD.htm>

Course Outline

- Introduction and Case Study
- Drug Targets
 - Sequence analysis
 - Protein structure prediction
 - Molecular simulation
- Drug Design
 - Combinatorial library
 - 3D-QSAR
 - Statistical methods
- Molecular Docking

Basic problems in sequence analysis

- **Sequence alignment**
- **Genome construction**
 - genes → proteins
 - the non-coding sequences
- **Functional element finding**
 - Exons, introns, promoters, ...
- **Regulatory network analysis**
 - Transcription factors



Targets

Sequence Alignment

- DNA/RNA alignment
 - ACGT,U
- Protein alignment
 - 20 amino acids
- Global alignment and local alignment
 - Lower score

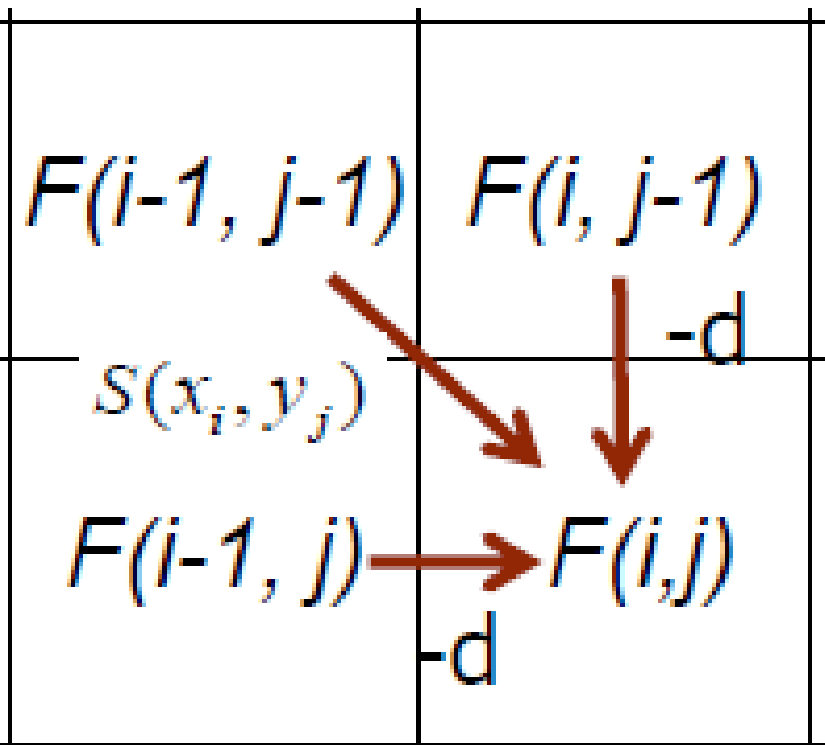
Dynamic programming sequence alignment algorithms

- Needleman/Wunsch global alignment
- Smith/Waterman local alignment
- Linear and affine gap penalties

Needleman/Wunsch global alignment

- Needleman, SB and Wunsch, CD. A general method applicable to the search for similarities in the amino acid sequence of two proteins”, *J. Mol. Biol.* 48:443-453, 1970.
- Two sequences $X = x_1 \dots x_n$ and $Y = y_1 \dots y_m$
- Let $F(i, j)$ be the optimal alignment score of $X_{1 \dots i}$ of X up to x_i and $Y_{1 \dots j}$ of Y up to Y_j ($0 \leq i \leq n, 0 \leq j \leq m$), then we have

Needleman/Wunsch global alignment



$$F(0,0) = 0$$

$$F(i, j) = \max \begin{cases} F(i-1, j-1) + s(x_i, y_j) \\ F(i-1, j) - d \\ F(i, j-1) - d \end{cases}$$

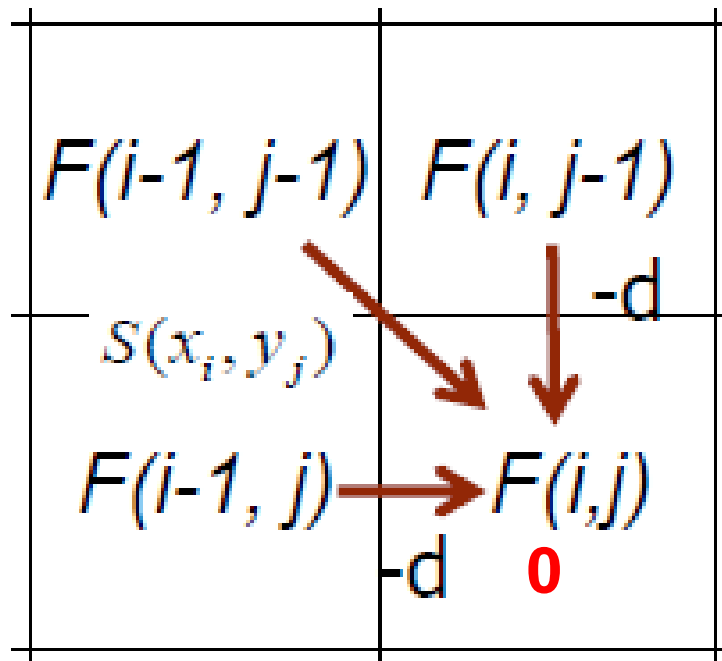
		G	C	C	C	T	A	G	C	G
	0	-2	-4	-6	-8	-10	-12	-14	-16	-18
G	-2	1	-1	-3	-5	-7	-9	-11	-13	-15
C	-4	-1	2	0	-2	-4	-6	-8	-10	-12
G	-6	-3	0	1	-1	-3	-5	-5	-7	-9
C	-8	-5	-2	1	2	0	-2	-4	-4	-6
A	-10	-7	-4	-1	0	1	1	-1	-3	-5
A	-12	-9	-6	-3	-2	-1	2	0	-2	-4
T	-14	-11	-8	-5	-4	-1	0	1	-1	-3
G	-16	-13	-10	-7	-6	-3	-2	1	0	0

图 2-1 Needleman-Wunsch 算法二维表

$d=2; s(x_i, y_j) = 1$ (same) | -1 (different)

Smith/Waterman local alignment

Smith, TF and Waterman, MS. Identification of Common Molecular Subsequences”, *J. Mol. Biol.* 147: 195-197, 1981



$$F(0,0) = 0$$

$$F(i, j) = \max \begin{cases} 0 \\ F(i-1, j-1) + s(x_i, y_j) \\ F(i-1, j) - d \\ F(i, j-1) - d \end{cases}$$

		G	C	C	C	T	A	G	C	G
	0	0	0	0	0	0	0	0	0	0
G	0	1	0	0	0	0	0	1	0	1
C	0	0	2	1	1	0	0	0	2	0
G	0	1	0	1	0	0	0	1	0	3
C	0	0	2	1	2	0	0	0	2	1
A	0	0	0	1	0	1	1	0	0	1
A	0	0	0	0	0	0	2	0	0	0
T	0	0	0	0	0	1	0	1	0	0
G	0	1	0	0	0	0	0	1	0	1

图 2-2 Smith-Waterman 算法二维表

$d=2; s(x_i, y_j) = 1$ (same) | -1 (different)

Linear and affine gap penalties

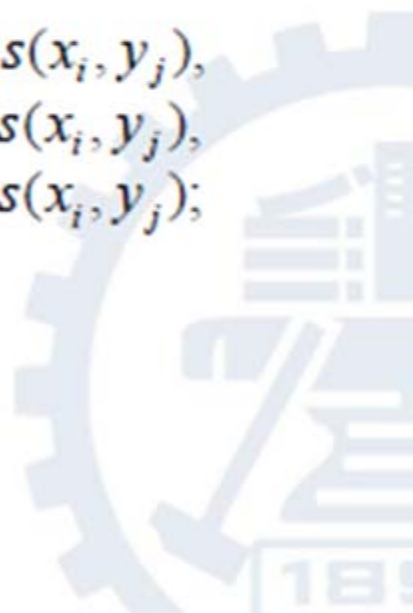
- Linear: $w(k) = k d$
- Affine: $w(k) = d + (k-1) e$
- Let $M(i,j)$, $I_x(i,j)$, $I_y(i,j)$ be the best scores up to (i,j) :
 - $M(i,j)$: x_i is aligned to y_j ;
 - $I_x(i,j)$: x_i is aligned to a gap;
 - $I_y(i,j)$ y_j is aligned to a gap

then we have

$$M(i, j) = \max \begin{cases} M(i-1, j-1) + s(x_i, y_j), \\ I_x(i-1, j-1) + s(x_i, y_j), \\ I_y(i-1, j-1) + s(x_i, y_j); \end{cases}$$

$$I_x(i, j) = \max \begin{cases} M(i-1, j) - d, \\ I_x(i-1, j) - e, \end{cases}$$

$$I_y(i, j) = \max \begin{cases} M(i, j-1) - d, \\ I_y(i, j-1) - e. \end{cases}$$



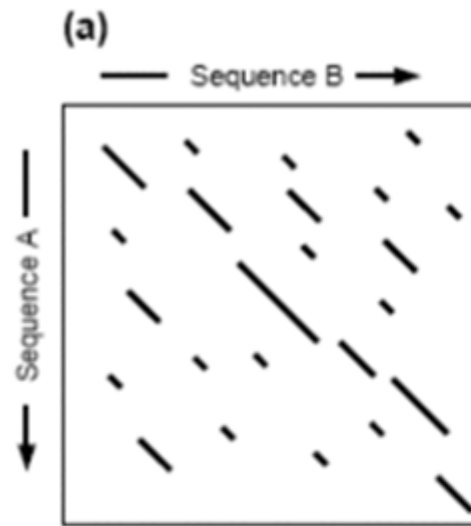
Improved algorithms

- Gotoh, O. An improved algorithm for matching biological sequences, *J. Mol. Biol.* 162:705-708, 1982
 - More efficient Needleman/Wunsch and Smith/Waterman algorithms.
- Myers, E. W. and Miller, W. Optimal alignment in linear space”, *CABIOS* 4: 11-17, 1988.
 - A method to reduce the memory cost from $O(n^2)$ to $O(n)$
- Henikoff, S. and Henikoff, J.G. Amino Acid Substitution Matrices from Protein Blocks. *PNAS* 89 (22): 10915–10919, 1992.
 - BLOSUM, gap penalty matrix for sequence alignment of proteins.

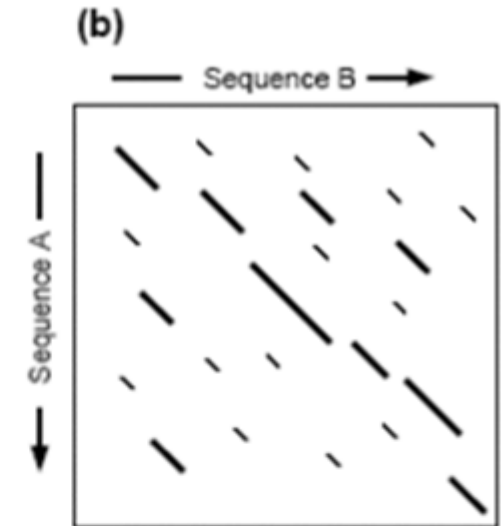
Heuristic algorithms

- Trading optimality, completeness, accuracy, or precision for speed
 - Much faster
 - May completely miss the optimal alignment
- How to save work?
 - Seed: find significant gap-less matches and then extend them
 - Searching in the diagonal band of the matrix
- Two important algorithms for alignment
 - FASTA
 - BLAST

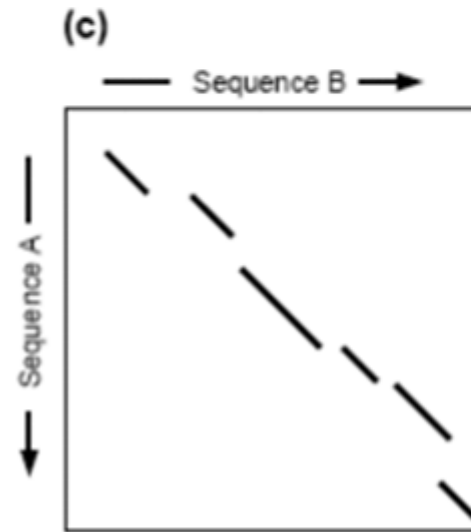
FASTA Algorithm



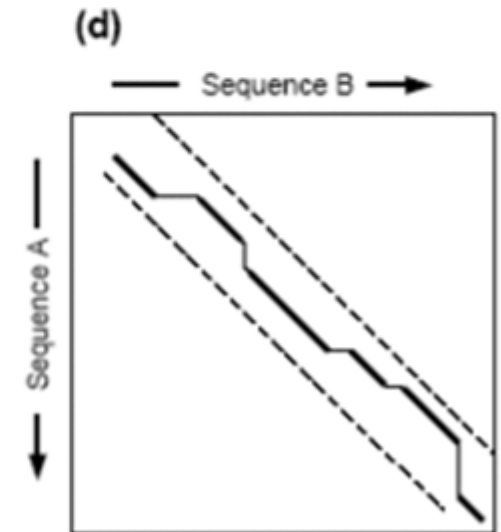
Find runs of identities



Re-score using PAM matrix
Keep top scoring segments.



Apply "joining threshold"
to eliminate segments that
are unlikely to be part of the alignment
that includes highest scoring segment.



Use dynamic programming
to optimise the alignment in a
narrow band that encompasses
the top scoring segments.

FASTA

- FASTA = "FAST-All" = "FAST-P" (protein) + "FAST-N" (nucleotide)
- FASTA file format

```
>1SIO:A|PDBID|CHAIN|SEQUENCE
```

```
AAPTAYTPLDVAQAYQFPEGLDGGQGCIAIIELGGGYD  
EASLAQYFASLGVPAPQVVSVDGASNQPTGDPSGP  
DGEVELDIEVAGALAPGAKFAVYFAPNTDAGFLDAITTA  
IHDPTLKPSVVSISWGGPEDSWTSAAIAAMNRAFLDAA  
ALGVTVLAAAGDSGSTDGEQDGLYHVDFPAASPYVLA  
CGGTRLVASGGRIAQETVWNDGPDGGATGGGVSRIF  
PLPAWQEHANVPPSANPGASSGRGVPDLAGNADPAT  
GYEVVIDGEATVIGGTSAVAPLFAALVARINQKLGKAVG  
YLNPTLYQLPADVFHDITEGNNDIANRAQIYQAGPGWD  
PCTGLGSPIGVRLQALLPSASQPQP
```


Basic Local Alignment Search Tool (BLAST)

- **Publications:**

- **Ungapped BLAST** – Altschul et al., 1990
- **Gapped BLAST, PSI-BLAST** - Altschul et al., 1997

- **Input:**

- **Query (target) sequence:** DNA, RNA or Protein
- **Scoring Scheme** – gap penalties, substitution matrix for proteins, identity/mismatch scores for DNA/RNA
- **Word length W** – typical is $W=3$ for proteins and $W=11$ for DNA/RNA

- **Output:**

- Statistically significant matches

Sequencing by BLAST

Different programs are available according to the type of query

Program	Query	Database
blastp	protein	protein
blastn	nucleotide	nucleotide
blastx	nucleotide protein	protein
tblastn	protein	nucleotide protein
tblastx	nucleotide protein	nucleotide protein

Standard Protein BLAST

BLAST Basic Local Alignment Search Tool

Home Recent Results Saved Strategies Help

NCBI/ BLAST/ blastp suite **Standard Protein BLAST**

blastn blastp **blastx** tblastn tblastx

BLASTP programs search protein databases using a protein query. [more...](#)

Enter Query Sequence

Enter accession number(s), gi(s), or FASTA sequence(s) [Clear](#) Query subrange [Query subrange](#)

From

To

Or, upload file [Browse...](#)

Job Title

Enter a descriptive title for your BLAST search

Align two or more sequences

Choose Search Set

Database [Help](#)

Organism Exclude [+](#)

Enter organism common name, binomial, or tax id. Only 20 top taxa will be shown.

Exclude Models (XM/XP) Uncultured/environmental sample sequences

Entrez Query

Enter an Entrez query to limit search

Program Selection

Algorithm

blastp (protein-protein BLAST)

PSI-BLAST (Position-Specific Iterated BLAST)

PHI-BLAST (Pattern Hit Initiated BLAST)

DELTA-BLAST (Domain Enhanced Lookup Time Accelerated BLAST)

Choose a BLAST algorithm

BLAST Search database Non-redundant protein sequences (nr) using Blastp (protein-protein BLAST)

Show results in a new window

- Non-redundant protein sequences (nr)
- Reference proteins (refseq_protein)
- UniProtKB/Swiss-Prot (swissprot)
- Patented protein sequences (pat)
- Protein Data Bank proteins (pdb)
- Metagenomic proteins (env_nr)
- Transcriptome Shotgun Assembly proteins (tsa_nr)

Results of BLAST

[Edit and Resubmit](#) [Save Search Strategies](#) [Formatting options](#) [Download](#) [YouTube](#) [How to read this page](#) [Blast report description](#)

1sio

Query ID	lcl 376901	Database Name	nr
Description	None	Description	All non-redundant GenBank CDS translations+PDB+SwissProt+PIR+PRF excluding environmental samples from WGS projects
Molecule type	amino acid	Program	BLASTP 2.2.28+ Citation
Query Length	364		

Other reports: [Search Summary](#) [Taxonomy reports](#) [Distance tree of results](#) [Related Structures](#) [Multiple alignment](#)

Graphic Summary

[Show Conserved Domains](#)

Putative conserved domains have been detected, click on the image below for detailed results.

The diagram shows a horizontal bar representing the query sequence from position 1 to 364. Key features are marked with red triangles: 'active site' at approximately position 100, 'catalytic triad' at approximately position 110, and 'calcium binding site' at approximately position 300. Below the bar, a red bar indicates a 'Specific hit' for 'Peptidases_553', and a yellow bar indicates a 'Superfamily' for 'Peptidases_S8_553 superfamily'. The 'Multi-domains' section is empty.

Distribution of 100 Blast Hits on the Query Sequence

Mouse-over to show define and scores, click to show alignments

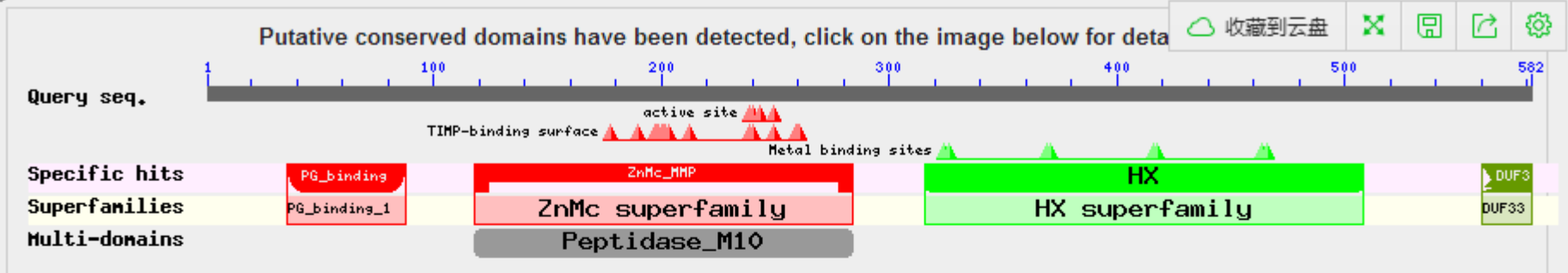
Color key for alignment scores	
Score Range	Color
<40	Black
40-50	Blue
50-80	Green
80-200	Purple
>=200	Red

The heatmap shows the distribution of 100 Blast Hits across the query sequence. The x-axis is labeled 'Query' and has markers at positions 1, 70, 140, 210, 280, and 350. The y-axis represents 100 individual hits. The heatmap is predominantly red, indicating high alignment scores (>=200) across most of the sequence.

Results of BLASTp on PDB

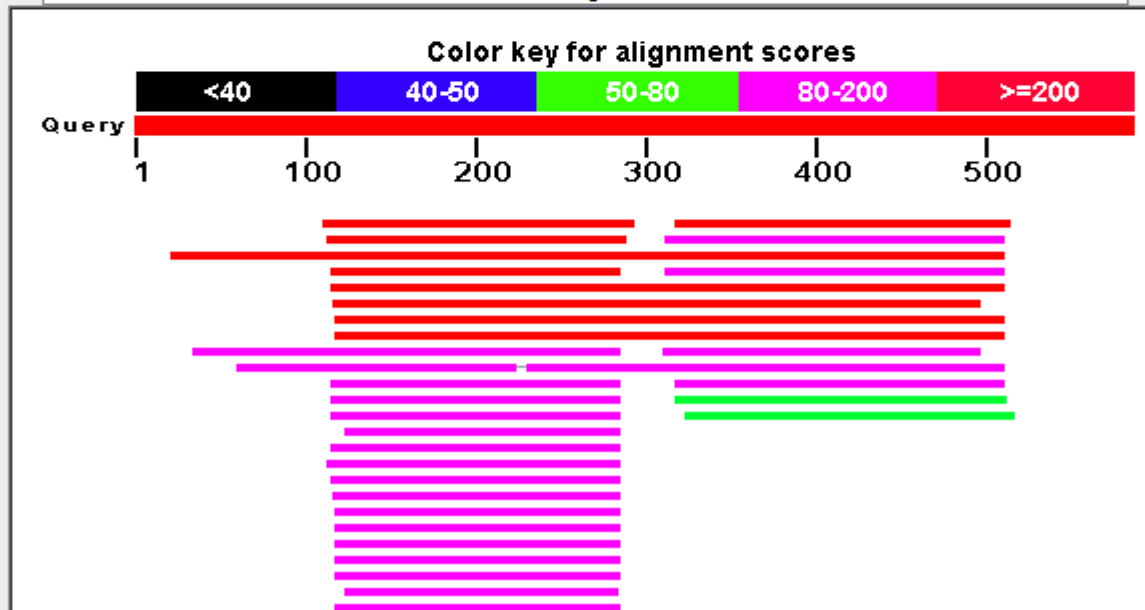
Graphic Summary

Show Conserved Domains



Distribution of 104 Blast Hits on the Query Sequence

Mouse over to see the define, click to show alignments



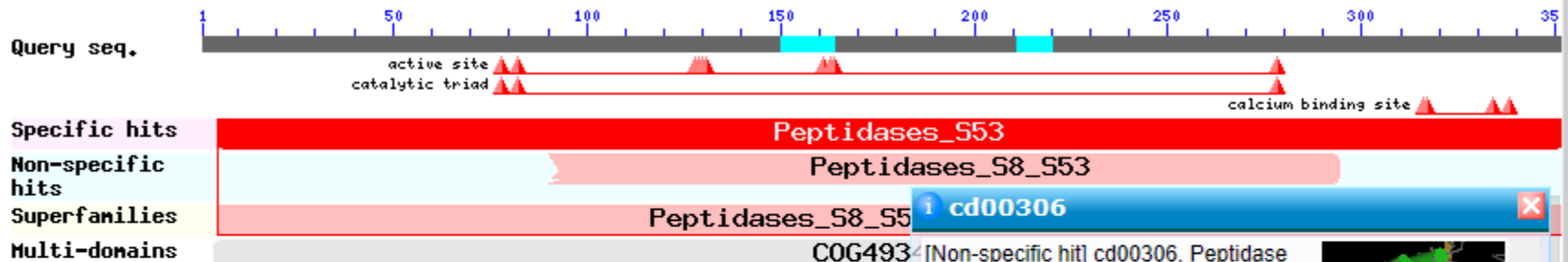
Results about conserved domain

Conserved domains on [gi|47169143|pdb|1SIOA|]

[View concise](#)

Chain A, Structure Of Kumamolisin-As Complexed With A Covalently-Bound Inhibitor, Acipf

Graphical summary [show options >](#)



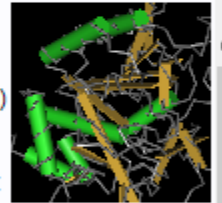
[Search for similar domain architecture](#)

List of domain hits

	Description
[+]Peptidases_S53[cd04056], Peptidase domain in the S53 family; Members of the peptidases S53 (sedolisin) family include endopeptidases and exopeptidases. The S53 family contains a catalytic triad Glu/Asp/Ser with an additional acidic residue Asp in the oxyanion hole, similar to that of subtilisin. The serine residue here is the nucleophilic equivalent of the serine residue in the S8 family, while glutamic acid has the same role here as the histidine base	
[+]Peptidases_S8_S53[cd00306], Peptidase domain in the S8 and S53 families; Members of the peptidases S8 (subtilisin and kexin) and S53 (sedolisin) family include endopeptidases and exopeptidases. The S8 family has an Asp/His/Ser catalytic triad similar to that found in trypsin-like proteases, but do not share their three-dimensional structure and are not homologous to trypsin. Serine acts as a nucleophile, aspartate as an electrophile, and histidine as a base. The S53 family contains a catalytic triad Glu/Asp/Ser with an additional acidic residue Asp in the oxyanion hole, similar to that of subtilisin. The serine residue here is the nucleophilic equivalent of the serine residue in the S8 family, while glutamic acid has the same role here as the histidine base	
[+]COG4934[COG4934], Predicted protease [Posttranslational modification, protein turnover, chaperone]	

Blast search parameters

Data Source: Precalculated data, version = cdd.v.3.10
 Preset Options: Database: cdsearch/cdd Low complexity filter: yes E-value threshold: 0.01



Description of conserved domain

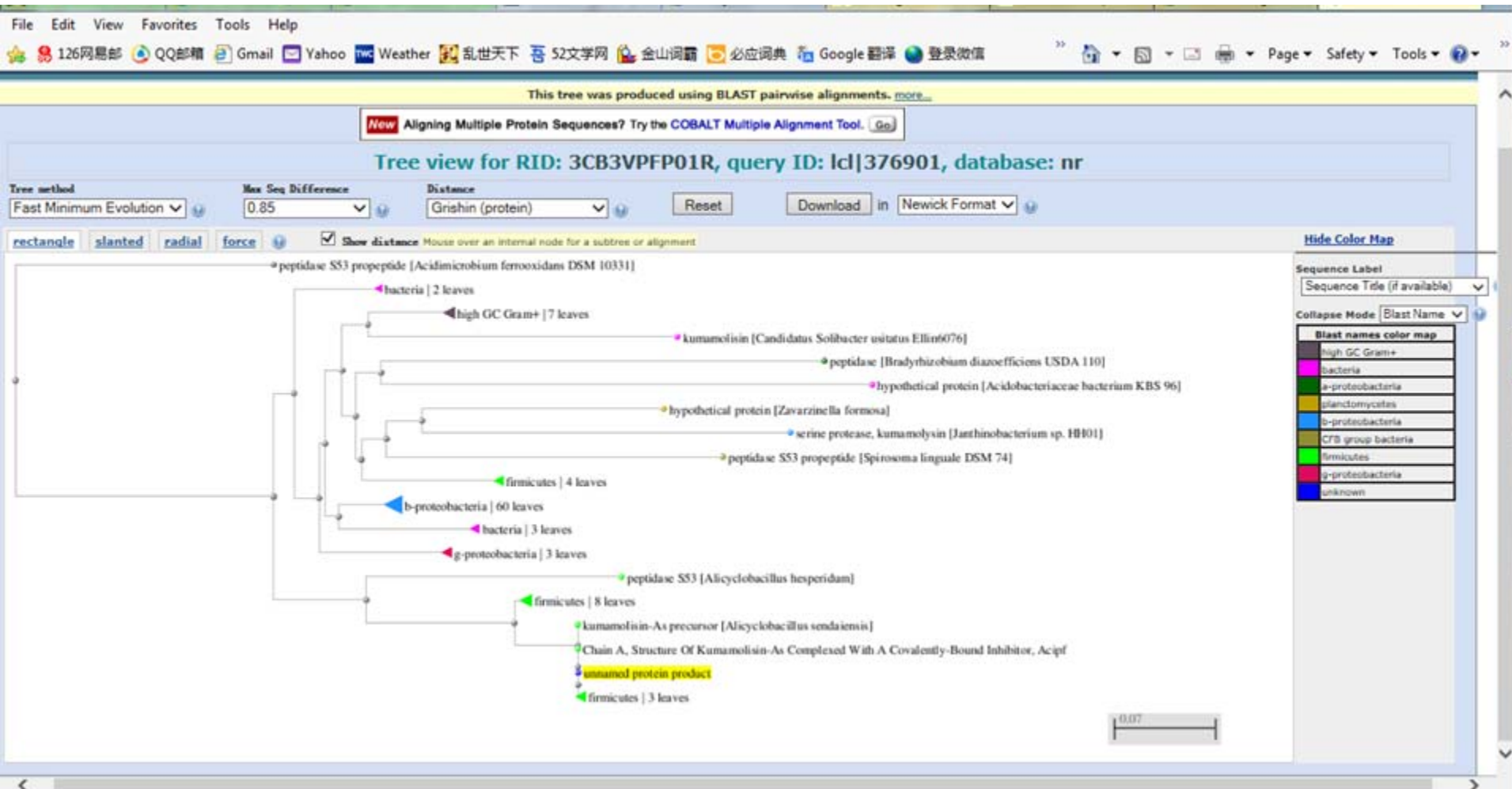
List of domain hits

Description	PssmId	Multi-
[-] Peptidases_S53[cd04056], Peptidase domain in the S53 family; Members of the peptidases S53 (sedolisin) family include endopeptidases and exo	173788	yes
<p>Peptidase domain in the S53 family; Members of the peptidases S53 (sedolisin) family include endopeptidases and exopeptidases sedolisin, kumamolysin, and (PSCF insensitive Carboxyl Proteinase. The S53 family contains a catalytic triad Glu/Asp/Ser with an additional acidic residue Asp in the oxyanion hole, similar to that of Asn stability of these enzymes may be enhanced by calcium, some members have been shown to bind up to 4 ions via binding sites with different affinity. Some members contain disulfide bonds. These enzymes can be intra- and extracellular, some function at extreme temperatures and pH values. Characterized sedolisins include Kum; extracellular calcium-dependent thermostable endopeptidase from Bacillus. The enzyme is synthesized with a 188 amino acid N-terminal preprotein region which is cleaved into the extracellular space with low pH. One kumamolysin paralog, kumamolysin-As, is believed to be a collagenase. TPP1 is a serine protease that functions as a tripeptidyl exopeptidase as well as an endopeptidase. Less is known about PSCP from Pseudomonas which is thought to be an aspartic proteinase.</p>		

Cd Length: 361 Bit Score: 321.19 E-value: 4.01e-107

	10	20	30	40	50	60	70	80		
1SIOA	5	AYTFLDVAQAYQFPE	GLDGGQCIAIIE	GGGYDEASLAQYFASL	GVPAPQVVS	VSDGASNqPTGDP	SGPDGEVELD	82		
Cdd:cd04056	1	GYTPADLAALYNIPP	1GYTGGQTIGIIE	FGGGYNPSDLQTF	FFQLfGLPAPT	VFIVVIGGGN	APGTSSGWGGEASLD	79		
1SIOA	83	IEVAGALAPGAKFAV	YFAPNT-DAGFLDA	ITTAIHDP	TLKPSVVSISWGG	PEDSWTSA	AIAAMNRAFLD	AAALGVTVLAA	161	
Cdd:cd04056	80	VEYAGAIAPGANITL	YFAPGTvTNGPL	LAFLAAVLDNPN	LPVSVISISYGE	PEQSLPPA	YAQRVCNLF	AAAAQGITVLAA	159	
1SIOA	162	AGDSGSTDGEQDG	---LYHVDFP	AASPYVLA	CGGTRLVA	SGGRIAQET	VWVNDGPD	GGATGGGVS	RIFPLPAWQ	EHA-NV 236
Cdd:cd04056	160	SGDSGAGCGCGDG	agtGFSV	SFPASSPYVT	AVGGTTL	LYTgGTGSSA	ESTVWS	SEGGWGG	SGGGSNY	FPRPSYQSGAvLG 239
1SIOA	237	PPSANPGASSGRG	VPDLAGNADP	ATGYEVV	IDGEATV	IGGTS	AVAPLFAAL	VARINQ	KLKGA---	VGYLNPTLYQLPAD 312
Cdd:cd04056	240	LPPSGLYNGS	GRGVPDVAAN	ADPGTYL	VVWNGQ	WYLVGGT	SAAAPL	FAGLIAL	INQARLAA	gkppLGFLNPLLYQLAAT 319
1SIOA	313	---VFHDITE	GNNDIANr	AQIQAG	PWDPCTGL	GSPIGV	RLL	352		
Cdd:cd04056	320	apsAFNDITS	GNNGCG	GAGY	PAGPWDP	VTGLG	TPNFAKLL	361		

Distance tree



Descriptions of searched results

Descriptions

Sequences producing significant alignments:

Select: [All](#) [None](#) Selected:0

[Alignments](#) [Download](#) [GenPept](#) [Graphics](#) [Distance tree of results](#) [Multiple alignment](#)

	Description	Max score	Total score	Query cover	E value	Ident	Accession
<input type="checkbox"/>	kumamolisin-As precursor [Allycycobacillus sendanensis]	725	725	100%	0.0	100%	BAC41257.1
<input type="checkbox"/>	Chain A, Structure Of Kumamolisin-As Complexed With A Covalently-Bound Inhibitor, Acplf -pdb 1SIO B Chain B, Structure Of Kumamolisin-As C	717	717	100%	0.0	100%	1SIO_A
<input type="checkbox"/>	Chain A, Structure Of Kumamolisin-As Mutant, D164n	716	716	100%	0.0	99%	1ZVJ_A
<input type="checkbox"/>	Chain A, Kumamolisin-As E78h Mutant	715	715	100%	0.0	99%	1SIU_A
<input type="checkbox"/>	Chain A, Structure Of Double Mutant, D164n, E78h Of Kumamolisin-as -pdb 1ZVK B Chain B, Structure Of Double Mutant, D164n, E78h Of Kum	703	703	98%	0.0	99%	1ZVK_A
<input type="checkbox"/>	Peptidase S53 propeptide [Allycycobacillus acidocaldarius subsp. acidocaldarius DSM 445] -ref WP_012810159.1 peptidase S53 [Allycycobacillus	687	687	100%	0.0	94%	YP_003184192.1
<input type="checkbox"/>	Pro- Kumamolisin [Allycycobacillus acidocaldarius subsp. acidocaldarius Tc-4-1] -ref WP_014463596.1 peptidase S53 [Allycycobacillus acidocald	659	659	100%	0.0	93%	YP_005517214.1
<input type="checkbox"/>	kumamolisin precursor [Bacillus sp. MN-32]	659	659	100%	0.0	93%	BA885637.2
<input type="checkbox"/>	Chain A, High Resolution Crystal Structure Of Mutant E23a Of Kumamolisin, A Sedolisin Type Proteinase (Previously Called Kumamolysin Or Ksc	650	650	100%	0.0	92%	1T1G_A
<input type="checkbox"/>	Chain A, High Resolution Crystal Structure Of The Intact Pro- Kumamolisin, A Sedolisin Type Proteinase (Previously Called Kumamolysin Or Ksc	656	656	100%	0.0	92%	1T1E_A
<input type="checkbox"/>	Chain A, High Resolution Crystal Structure Of Mutant W129a Of Kumamolisin, A Sedolisin Type Proteinase (Previously Called Kumamolysin Or K	647	647	100%	0.0	92%	1T1I_A
<input type="checkbox"/>	peptidase S53 [Allycycobacillus acidocaldarius] -qb EED05251.1 Peptidase S53 propeptide [Allycycobacillus acidocaldarius LAA1]	650	650	98%	0.0	92%	WP_008341072.1
<input type="checkbox"/>	Chain 1, High Resolution Crystal Structure Of A Thermostable Serine-Carboxyl Type Proteinase, Kumamolisin (Kscp) -pdb 1GT9 2 Chain 2, High	640	640	98%	0.0	93%	1GT9_1
<input type="checkbox"/>	peptidase S53 [Allycycobacillus hesperidum] -qb EJY56072.1 Peptidase S53 propeptide [Allycycobacillus hesperidum URH17-3-68]	489	489	96%	2e-165	72%	WP_006446582.1
<input type="checkbox"/>	peptidase S53 [Paenibacillus elqii]	443	443	97%	9e-149	65%	WP_010499475.1
<input type="checkbox"/>	hypothetical protein [Bacillus sp. 105MF]	420	420	98%	4e-140	59%	WP_018767626.1
<input type="checkbox"/>	hypothetical protein [Zavarzinella formosa]	413	413	96%	1e-138	61%	WP_020470275.1
<input type="checkbox"/>	peptidase S53 [Burkholderia sp. SJ98] -qb EKS66971.1 peptidase S53 propeptide [Burkholderia sp. SJ98]	410	410	98%	2e-136	62%	WP_008353804.1
<input type="checkbox"/>	hypothetical protein [Bacillus cereus] -qb EOO25493.1 hypothetical protein ICC_04848 [Bacillus cereus BAG1X1-1] -qb EOO43642.1 hypotheti	409	409	97%	1e-135	59%	WP_016082867.1

Alignment

Alignments

Download [GenPept](#) [Graphics](#)

kumamolisin-As precursor [Alicyclobacillus sendaiensis]
 Sequence ID: [dbj|BAC41257.1](#) Length: 553 Number of Matches: 1

Range 1: 190 to 553 [GenPept](#) [Graphics](#)

[Next Match](#) [Previous Match](#)

Score	Expect	Method	Identities	Positives	Gaps
725 bits(1871)	0.0	Compositional matrix adjust.	364/364(100%)	364/364(100%)	0/364(0%)
Query 1		AAPTAYTPLDVAQAYQFPEGLDGGQCCIAI IELGGGYDEASLAQYFASLGVPAPQVVSVS	60		
Sbjct 190		AAPTAYTPLDVAQAYQFPEGLDGGQCCIAI IELGGGYDEASLAQYFASLGVPAPQVVSVS	249		
Query 61		VDGASNQPTGDPSPGDGEVELDIEVAGALAPGAKFAVYFAPNTDAGFLDAITTAIHDPTL	120		
Sbjct 250		VDGASNQPTGDPSPGDGEVELDIEVAGALAPGAKFAVYFAPNTDAGFLDAITTAIHDPTL	309		
Query 121		KPSVVISISWGGPEDSWTSAIIAAMNRAFLDAAALGVTVLAAAGDSGSTDGEQDGLYHVDF	180		
Sbjct 310		KPSVVISISWGGPEDSWTSAIIAAMNRAFLDAAALGVTVLAAAGDSGSTDGEQDGLYHVDF	369		
Query 181		PAASPYVLACGGTRLVASGGRIAQETVWNDGPDGGATGGGVSRIFFLPWQEHANVPPSA	240		
Sbjct 370		PAASPYVLACGGTRLVASGGRIAQETVWNDGPDGGATGGGVSRIFFLPWQEHANVPPSA	429		
Query 241		NFGASSGRGVFDLGNADPATGYEVVIDGEATVIGGTSAVAPLFAALVARINQKLGKAVG	300		
Sbjct 430		NFGASSGRGVFDLGNADPATGYEVVIDGEATVIGGTSAVAPLFAALVARINQKLGKAVG	489		
Query 301		YLNFTLYQLPADVFHDITEGNNDIANRAQIYQAGPGWDPCTGLGSPIGVRLQALLPSAS	360		
Sbjct 490		YLNFTLYQLPADVFHDITEGNNDIANRAQIYQAGPGWDPCTGLGSPIGVRLQALLPSAS	549		
Query 361		QPQP 364			
Sbjct 550		QPQP 553			

CLUSTAL

- Multiple sequence alignment
 - Do a pairwise alignment
 - Create a (or use a user-defined tree)
 - Use the guide tree to carry out a multiple alignment
- Three main version
 - ClustalW: command line interface
 - ClustalX: graphical user interface
 - Clustal Omega: the latest online version

CLUSTAL



Clustal: Multiple Sequence Alignment

Multiple alignment of nucleic acid and protein sequences



Clustal Omega

- Latest version of Clustal - fast and scalable (can align hundreds of thousands of sequences in hours), greater accuracy due to new HMM alignment engine
- Command line/web server only (GUI public beta available soon)



ClustalW/ClustalX

- "Classic Clustal"
- GUI (ClustalX), command line (ClustalW), web server versions available

Multiple Sequence Alignment

Clustal Omega is a new multiple sequence alignment program that uses seeded guide trees and HMM profile-profile techniques to generate alignments.

STEP 1 - Enter your input sequences

Enter or paste a set of sequences in any supported format:

```
AAPTAYTPLDVAQAYQFPEGLDGGQGCIAIIELGGGYDEASLAQYFASLGVPAPQVVSVDGASNQPTGDPSGPDGEVE
LDIEVAGALAPGAKFAVYFAPNTDAGFLDAITTAIHDTLKPSSVVSISWGGPEDSWTSAIIAAMNRAFLDAAALGVTVLA
AAGDSGSTDGEQDGLYHVDFFAASPYVLACGGTRLVASGGRIAQETVWNDGPDGGATGGGVSRIFPLPAWQEHANVPPSA
NPGASSGRGVPDLGNADPATGYEVVIDGEATVIGGTSAVAPLFAALVARINQKLGKAVGYLNPTLYQLPADVFHDITEG
NNDIANRAQIQAGPGWDPCTGLGSPIGVRLQALLPSASQPQP
```

Or, [upload](#) a file:

STEP 2 - Set your parameters

OUTPUT FORMAT

The default settings will fulfill the needs of most users and, for that reason, are not visible.

(Click here, if you want to view or change the default settings.)

STEP 3 - Submit your job

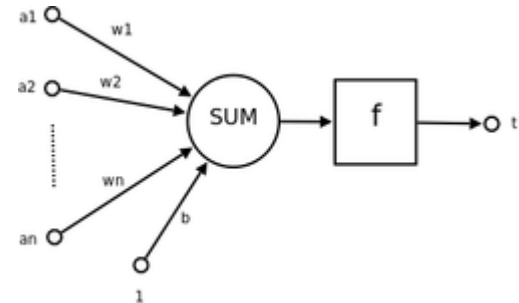
Be notified by email *(Tick this box if you want to be notified by email when the results are available)*

Course Outline

- Introduction and Case Study
- Drug Targets
 - Sequence analysis
 - Protein structure prediction
 - Molecular simulation
- Drug Design
 - Combinatorial library
 - 3D-QSAR
 - Statistical methods
- Molecular Docking

Classes of Secondary Structural Prediction

- Chou-Fasman
 - P_{α} , P_{β} , P_{turn}
 - Overlap
- Amino Acid Hydrophobicity
- k-Nearest Neighbor algorithm (kNN, 最邻近法, 相似片段法)
- Artificial Neural Networks (ANNs, 人工神经网络)



Amino Acid Hydrophobicity

AA	Abrr.	Hydrophobicity	AA	Abrr.	Hydrophobicity
Ala	A	1.8	Leu	L	3.8
Arg	R	-4.5	Lys	K	-3.9
Asn	N	-3.5	Met	M	1.9
Asp	D	-3.5	Phe	F	2.8
Cys	C	2.5	Pro	P	-1.6
Gln	Q	-3.5	Ser	S	-0.8
Glu	E	-3.5	Thr	T	-0.7
Gly	G	-0.4	Trp	W	-0.9
His	H	-3.2	Tyr	Y	-1.3
Ile	I	4.5	Val	V	4.2

Methods of Secondary Structural Prediction

- DPM（双重预测方法）
- DSC
- PHDsec
- SOMPA
- MLRC
- Jpred
- Predict Protein

Jpred

<http://www.compbio.dundee.ac.uk/~www-jpred/>



Jpred 3
Incorporating Jnet

A Secondary Structure Prediction Server

Update (19 Aug 2014): JNet v.2.3.1. [Please see news page for details.](#)

Sequence:

[Help](#)
[Advanced](#)

Make Prediction

Clear

[The Barton Group - The University of Dundee](#)

Citation: Cole C, Barber JD & Barton GJ. *Nucleic Acids Res.* 2008. 35 (suppl. 2) W197-W201 [[link](#)]

[More citations](#)



Jpred results

```
      : 1-----11-----21-----31-----41-----51-----61-----71-----81-----91-----101-----111-----
OrigSeq : MSPAPRPPRCLLLPLLTLG TALASLGS AQSSSF SPEAWLQQYGYLPPGDLRHTHTQRSPQSL SAAIAAMQKFYGLQVTGKADATMKAMRRPRCGVPDKFGAEIKANVRRKRYAI
Jnet    : -----HHHHHHHHHHHHHH-----HHHHHHHHHH-----HHHHHHHHHHHHHH-----HHHHHHHH-----EE-----
jhmm    : -----HHHHHHHHHHHHHHHH-----HHHHHHHHHH-----HHHHHHHHHHHHHH-----HHHHHHHH-----EE-----
jpssm   : -----HHHHHHHHHHHHHH-----HHHHHHHHHH-----HHHHHHHHHHHHHH-----HHHHHHHH-----E-----
Lupas 14 : -----
Lupas 21 : -----
Lupas 28 : -----
Jnet_25 : -----BBBBBBBBBBBBBBBB-----B-B-BB-BB-B-----B-BB-BB-BB-B-BBB-B--BB-BB--B-BBB-----B--B--BBB-----
Jnet_5  : -----BBBBBBBBBBBB-----B-BB-----B-B-B-----B-----
Jnet_0  : -----B--B-----B--B-----B-----
Jnet Rel : 999312378999999998603678877775278999998623677777777777642899999999985167888766278999986236777777777776510124357
```

Notes

Key:

- Colour code for alignment:
- Blue - Complete identity at a position
- Shades of red - The more red a position is, the higher the level of conservation of chemical properties of the amino acids
- Jnet - Final secondary structure prediction for query
- jalign - Jnet alignment prediction
- jhmm - Jnet hmm profile prediction
- jpssm - Jnet PSIBLAST pssm profile prediction
- Lupas - Lupas Coil prediction (window size of 14, 21 and 28)
- Note on coiled coil predictions
 - = less than 50% probability
 - c = between 50% and 90% probability
 - C = greater than 90% probability
- Jnet_25 - Jnet prediction of burial, less than 25% solvent accessibility
- Jnet_5 - Jnet prediction of burial, less than 5% exposure
- Jnet_0 - Jnet prediction of burial, 0% exposure
- Jnet Rel - Jnet reliability of prediction accuracy, ranges from 0 to 9, bigger is better.

SOMPA



Pôle BioInformatique Lyonnais
Network Protein Sequence Analysis

NPS@ is the [IBCP](#) contribution to [PBIL](#) in Lyon, France

[\[HOME\]](#) [\[NPS@\]](#) [\[HELP\]](#) [\[REFERENCES\]](#) [\[NEWS\]](#) [\[MPSA\]](#) [\[ANTHEPROT\]](#) [\[Geno3D\]](#) [\[SuMo\]](#) [\[Positions\]](#) [\[PBIL\]](#)

Monday, August 26th 2013: added support of BCL2DB database ([see news](#))

SOPMA SECONDARY STRUCTURE PREDICTION METHOD

[\[Abstract\]](#) [\[NPS@ help\]](#) [\[Original server\]](#)

Sequence name (optional) :

Paste a protein sequence below : [help](#)

Output width :

Parameters

Number of conformational states :

Similarity threshold :

Window width :

PredictProtein

<https://www.predictprotein.org/>



Welcome Guest | [Log in](#) | [Register](#) | [Help](#) [Tutorials](#)

```
MSPAPRPPRCLLLPLLLTLGTALASLGSAQSSSF SPEAWLQQYGYLPPGDLRHTHTQRSPQSLSAIAIAMQK
FYGLQVTGKADADTMKAMRRPRCGVPDKFGAEIKANVRRKRYAIQGLKWQHNEITFCIQNYTPKVGEYAT
YEAIRKAFRVWESATPLRFREVYAYIREGHEKQADIMIFFAEGFHGDSTPFDDGEGGFLAHAYFPGPNIG
GDTHFDSAEPWTVRNEIDLNGNDIFLVAVHELGHALGLEHSSDPSAIMAPFYQWMDTENFVLPDDDDRRGIQ
QLYGGESGFPTKMPPQPRTTSRPSVPDKPKNPTYGPNICDGNFDTVAMLRGEMFVFKERWFWRVRNNQVM
DGYPMPIGQFWRGLPASINTAYERKDGKVFVFFKGDKHWWFDEASLEPGYPKHIKELGRGLPTDKIDAALF
WMPNGKTYFFRGNKYRFEELRAVDSEYPKNIKVWEGIPESPRGSFMGSDEVFTYFYKGNKYWKFNQK
LKVEPGYPKSALRDWMGCPSGGRPDEGTEETEVIIEVDEEGGGAVSAAAVLPLVLLLLLVLAVGLAVF
FFRRHGTPRRLLYCQRSLLDKV
```

Clear

PredictProtein

[Example Input 1

Example Input 2]

PredictProtein results

VIEWS

Dashboard >

STRUCTURE ANNOTATION

Secondary Structure and Solvent Accessibility >

Transmembrane Helices >

Protein Disorder and Flexibility >

Disulphide Bridges >

FUNCTION ANNOTATION

Effect of Point Mutations >

Gene Ontology Terms >

Subcellular Localization >

Binding Sites >

ADDITIONAL SERVICES

Literature Search >

HELP

Site Tutorial >

Secondary Structure and Solvent Accessibility Prediction for MMP14_HUMAN

Export

Visual

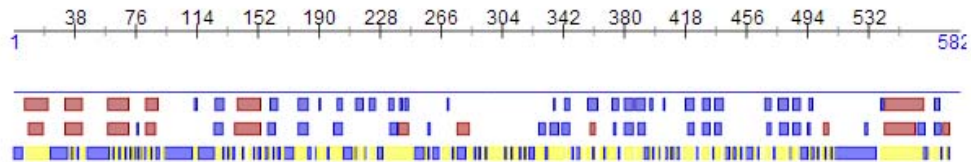
Help

What am I seeing Here? This viewer lays out predicted features that correspond to regions within the queried sequence. Mouse over the different colored boxes to learn more about the annotations. Note that this panel may show results from two prediction methods: RePROF (new and experimental) and PROFsec (veteran). The PROFsec method will be retired by the end of 2014. See references below and help sections for more information.

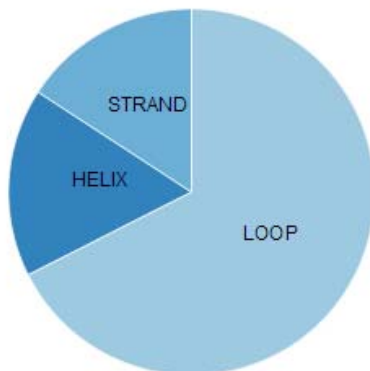
Zoom - Start:1, End:582



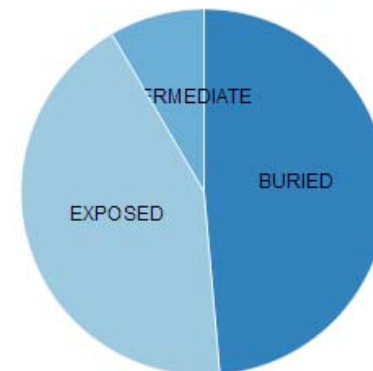
Export to image



Secondary Structure Composition



Solvent Accessibility



Course Outline

- Introduction and Case Study
- Drug Targets
 - Sequence analysis
 - Protein structure prediction
 - Molecular simulation
- Drug Design
 - Combinatorial library
 - 3D-QSAR
 - Statistical methods
- Molecular Docking