# Hydrophobicity profile

**Hydrophobicity profile** is simply the plot of the hydrophobicity indices of the residues against their sequence numbers.

**E.g. SAMPLEDATAWITHHYDINDICES**

**-1, 1, 1, -1, 2, -2, -2, 1, -1, 1, 2, 2, -1, -2, -2, 1, -2, 2, -1, -2, 2, 1, -2, -1**

>1a91_
A MENLNMDLLY MAAAVMMGLA AIGAAIGIGI
LGGKFLEGAA RQPDLIPLLR TQFFIVMGLV
DAIPMIAVGL GLYVMFAVA

A, C, G, M, Y: 1
F, I, L, V, W: 2
D, E, H, K, R: -2
N, P, Q, S, T: -1

# Sample data

**Nozaki-Tanford-Jones (Ht)**

A: 0.87 D: 0.66 C:1.52 E: 0.67

F: 2.87 G: 0.10 H: 0.87 I: 3.15

K: 1.64 L: 2.17 M: 1.67 N: 0.09

P: 2.77 Q: .00 R: 0.85 S: 0.07

T: 0.07 V: 1.87 W: 3.77 Y: 2.67

**Ponnuswamy-Gromiha (Hgm)**

A: 13.85 D: 11.61 C: 15.37 E: 11.38

F: 13.93 G: 13.34 H: 13.82 I: 15.28

K: 11.58 L: 14.13 M: 13.86 N: 13.02

P: 12.35 Q: 12.61 R: 13.10 S: 13.39

T: 12.70 V: 14.56 W: 15.48 Y: 13.88

Chain ID

Amino Ac

q31: size (choth
q32: relative mu
q33: aa compos
q34: pk-n (sobe
q35: pk-c (sobe
q36: melting po
q37: specific ro
q38: dihedral ar
q39: point muta
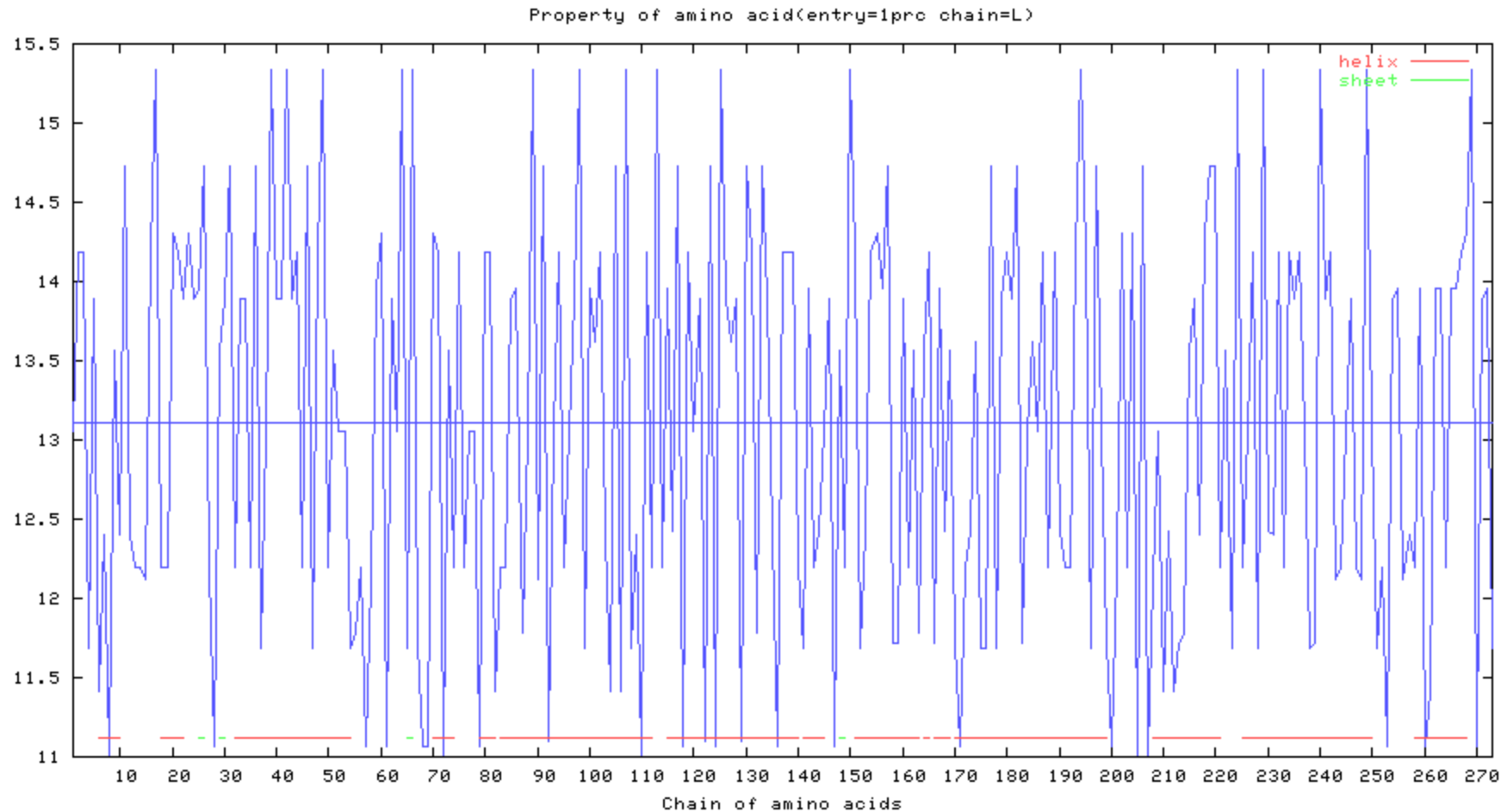q40: residue ac
q41: av access
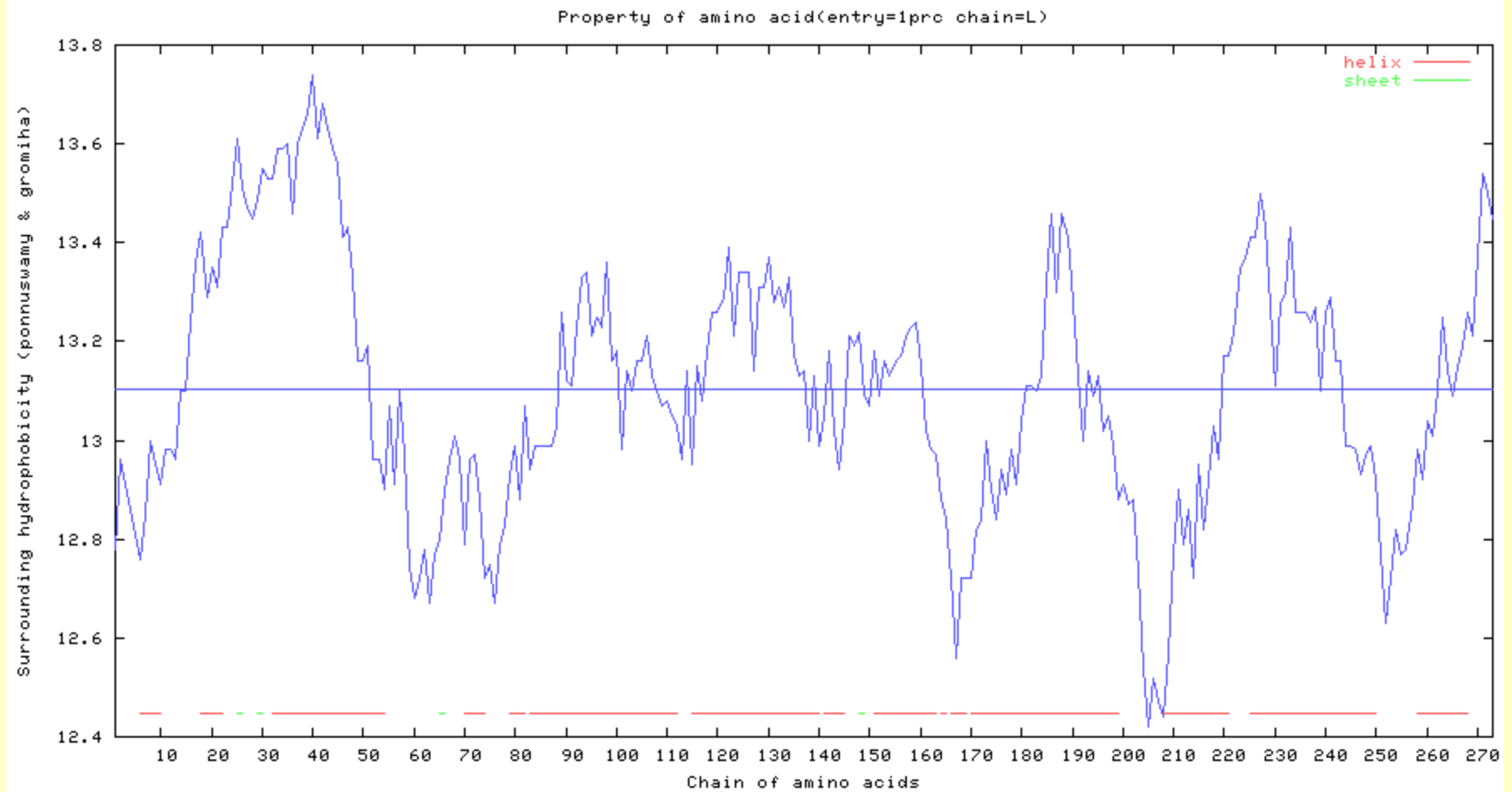q42: hydrophob

Operation

TEXT : ○

GRAPH :

1PRC

Property of amino acid(entry=1prc chain=L)



Start    Clea

**1PRC**

Property of amino acid(entry=1prc chain=L)

Open
TEX
GRA

Start

# Amphipathicity

Amphipathic character of amino acid residues is the **periodicities in the ploar/nonpolar character of the amino acid sequence** in a protein.

This has been examined by assigning a numerical hydrophobicity to each residue and searching for periodicity in the resulting one-dimensional function.

# Amphipathicity: $\alpha$-helices

The residues of an $\alpha$-helical segment are considered on four adjacent edges along the direction of the helical axis. The average hydrophobicity of the residues constituting the edge i (i = 1,4) is given by

$$\alpha_i = (\Sigma h_{i+j})/n,$$
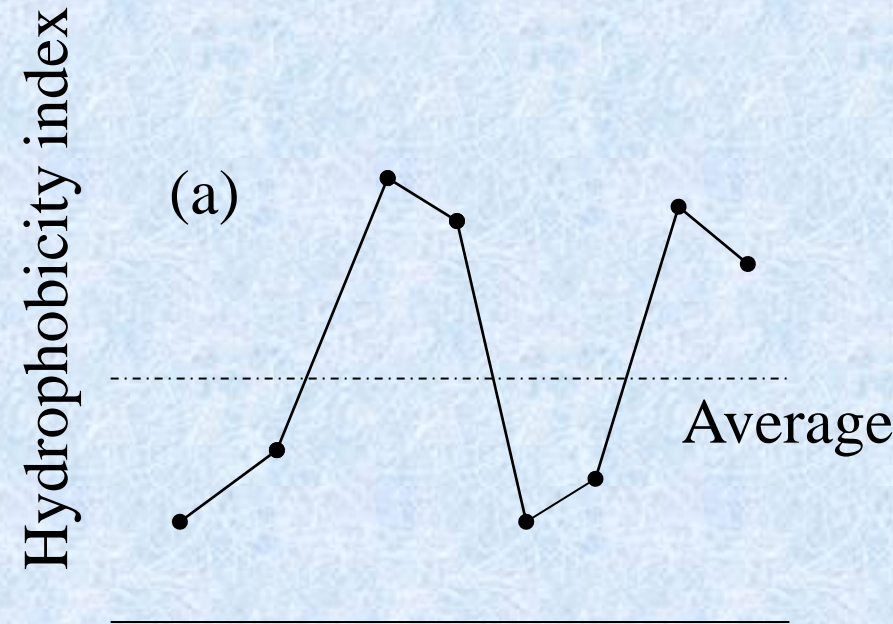
where n is the total number of residues in the edge,

j increases at an interval of 4 from 0 to m, m being the number of residues in the helix;

h is the hydrophobic index of the residue.

The power of amphipathicity of a helix is taken to be

$$A_\alpha = |(a_1+a_2) - (a_3+a_4)| \text{ or } |(a_1+a_4) - (a_2+a_3)|.$$

It has been reported that 75% of the helical segments in known structures are amphipathic in nature.

(a)

Hydrophobicity index

Average

# Amphipathicity: β-strands

A β-strand segment is considered to have two faces and the average hydrophobicity of residues constituting the face i (i = 1, 2) is given by
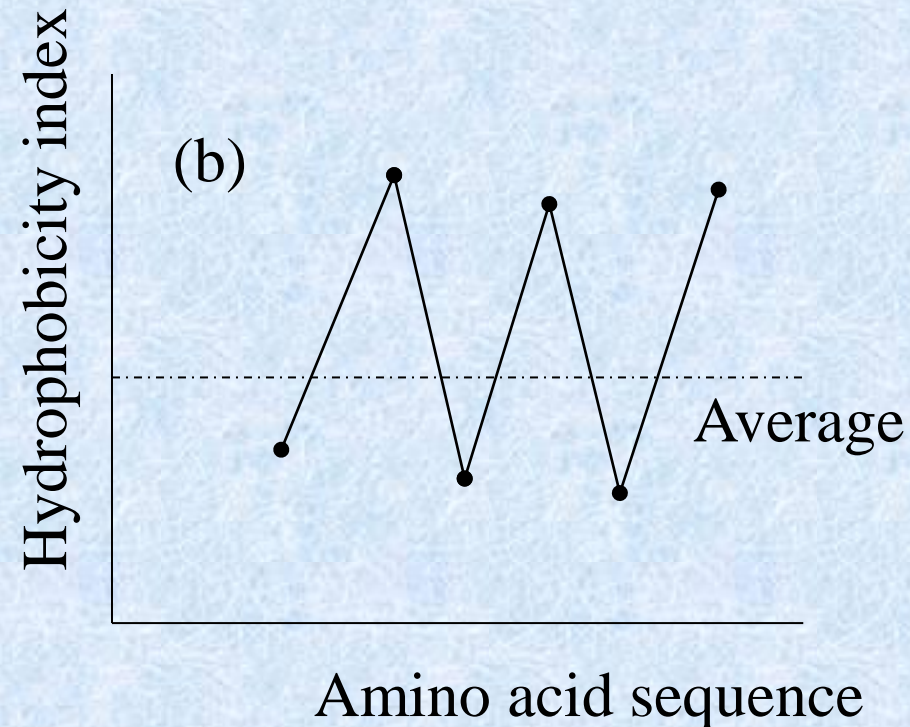
$$\beta_i = (\Sigma h_{i+j})/n,$$

where n is the total number of residues in the face, j increases at an interval of 2 from 0 to m, m being the number of residues in the strand;

The amphipathicity index of a strand is computed using the equation,
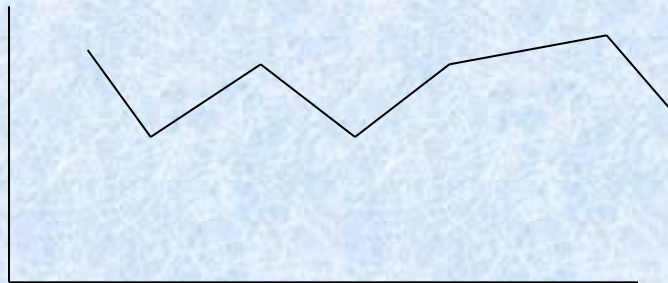
$$A_\beta = |\beta_1 - \beta_2|.$$

The structural analysis showed that about 65% of the β-strands possess amphipathic character.



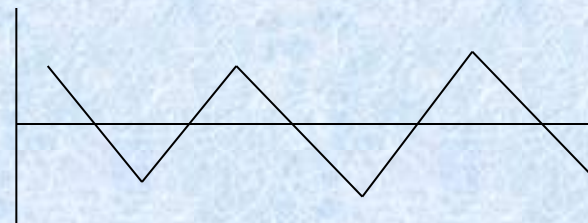(b)

Hydrophobicity index

Average

Amino acid sequence

# Patterns

- Identify the pattern of hydrophobic residues for membrane spanning helical proteins

- Amphiphathic character of β-strands by alternative hydrophobic-hydrophilic residues

E.g. AILVGYWFFVVA

AKINIHVTFKIKLP

# Pattern Definition

[LIVM]-[VIC]-x(2) -G-[DENQTA]-x-[GAC]-x(2)-[LIVMFY](4)-x (2)-G

1. Use capital letters for amino acid residues

2. Use "[…]" for a choice of multiple amino acids in a particular position. [LIVM] means that L, I, V, or M can be in the first position

3. Use "{…}" to exclude amino acids. {CF} means C and F should not be in that particular position

4. Use "x" or "X" for a position that can be any amino acid.

5. Use "(n)", where n is a number, for multiple positions; x(3) is the same as "xxx"

# PIR: Pattern Definition

[LIVM]-[VIC]-x(2) -G-[DENQTA]-x-[GAC]-x(2)-[LIVMFY](4)-x (2)-G

Illustrates a 17 amino acid peptide that has:

L, I, V, or M at position 1;

V, I, or C at position 2;

any residue at positions 3 and 4;

G at position 5 and so on ….

# PIR: Pattern search

# [LIVM]-[VIC]-x(2) -G-[DENQTA]-x-[GAC]-x(2) -[LIVMFY](4)-x(2)-G



Pattern Search Result (UniProtKB)

Query Pattern On ✓   Help ?

725 proteins | 37 pages | 20 / page |  |◄ «  ‹   22 | 23 | **24** | 25 | 26   › »►|

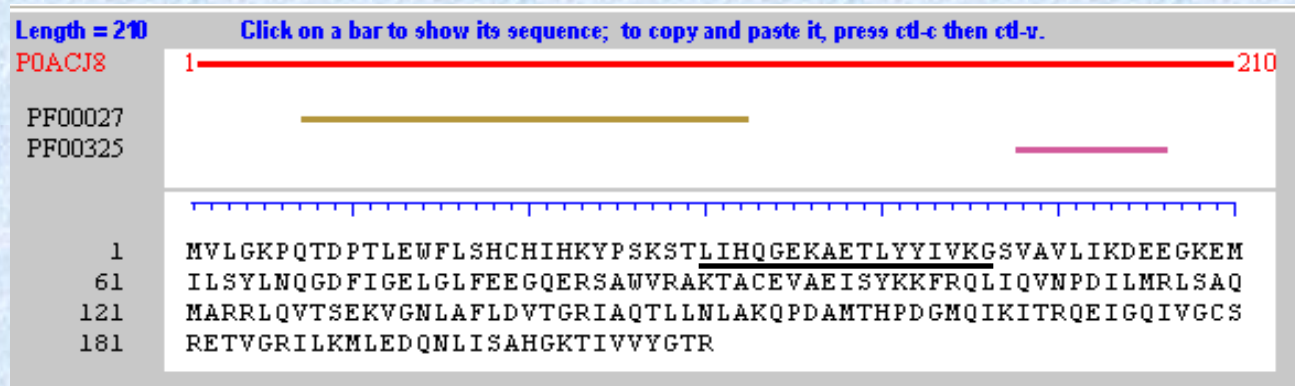Save Result As: ↘TABLE | ↘FASTA

1 selected  (show)

● BLAST  ● FASTA  ● Pattern Match  ● Multiple Alignment  ● Domain Display

| Protein AC/ID | Protein Name | Length | Organism Name | PIRSF ID | Match Range |
|---|---|---|---|---|---|
| ☐ O96777/O96777_HELVI /ProClass  UniProtKB/TrEMBL | Cyclic nucleotide and voltage-activated ion channel | 678 | Heliothis virescens (Tobacco budworm moth) | | 477-493; |
| ☐ O97119/O97119_LIMPO /ProClass  UniProtKB/TrEMBL | Cyclic nucleotide-gated ion channel LCNG1 | 900 | Limulus polyphemus (Atlantic horseshoe crab) | | 532-548; |
| ☐ P0A1L7/SPAQ_SALTY /ProClass  UniProtKB/Swiss-Prot | Surface presentation of antigens protein spaQ | 86 | Salmonella typhimurium | PIRSF004669 | 4-20; |
| ☐ P0A1L8/SPAQ_SALTI /ProClass  UniProtKB/Swiss-Prot | Surface presentation of antigens protein spaQ | 86 | Salmonella typhi | PIRSF004669 | 4-20; |
| ☐ P0A1M2/SPAQ_SALSE /ProClass  UniProtKB/Swiss-Prot | Surface presentation of antigens protein spaQ | 86 | Salmonella senftenberg | PIRSF004669 | 4-20; |
| ☐ P0A1M3/SPAQ_SALTP /ProClass  UniProtKB/Swiss-Prot | Surface presentation of antigens protein spaQ | 86 | Salmonella typhisuis | PIRSF004669 | 4-20; |
| ☑ P0ACJ8/CRP_ECOLI /ProClass  UniProtKB/Swiss-Prot | Catabolite gene activator; (AltName: Full=cAMP receptor protein; AltName: Full=cAMP regulatory protein) | 210 | Escherichia coli (strain K12) | PIRSF003151 | 30-46; |
| ☐ P0ACJ9/CRP_ECOL6 /ProClass  UniProtKB/Swiss-Prot | Catabolite gene activator; (AltName: Full=cAMP receptor protein; AltName: Full=cAMP regulatory protein) | 210 | Escherichia coli O6 | PIRSF003151 | 30-46; |
| ☐ P0ACK0/CRP_ECO57 /ProClass  UniProtKB/Swiss-Prot | Catabolite gene activator; (AltName: Full=cAMP receptor protein; AltName: Full=cAMP regulatory protein) | 210 | Escherichia coli O157:H7 | PIRSF003151 | 30-46; |
| ☐ P29973/CNGA1_HUMAN /ProClass  UniProtKB/Swiss-Prot | cGMP-gated cation channel alpha-1; (AltName: Full=CNG channel alpha-1; Short=CNG-1; Short=CNG1; AltName: Full=Cyclic nucleotide-gated channel alpha-1; AltName: Full=Cyclic nucleotide-gated channel, photoreceptor; AltName: Full=Cyclic nucleotide-gated cation channel 1; AltName: Full=Rod photoreceptor cGMP-gated channel subunit alpha) | 690 | Homo sapiens (Human) | PIRSF002403 | 506-522; |
| ☐ P29974/CNGA1_MOUSE /ProClass  UniProtKB/Swiss-Prot | cGMP-gated cation channel alpha-1; (AltName: Full=CNG channel alpha-1; Short=CNG-1; Short=CNG1; AltName: Full=Cyclic nucleotide-gated channel alpha-1; AltName: Full=Cyclic nucleotide-gated channel, photoreceptor; AltName: Full=Cyclic nucleotide-gated cation channel 1; AltName: Full=Rod photoreceptor cGMP-gated channel subunit alpha) | 684 | Mus musculus (Mouse) | PIRSF002403 | 498-514; |
| ☐ P36600/KAPR_SCHPO /ProClass  UniProtKB/Swiss-Prot | cAMP-dependent protein kinase regulatory subunit; ( Short=PKA regulatory subunit) | 412 | Schizosaccharomyces pombe (Fission yeast) | PIRSF000548 | 171-187;305-321; |
| ☐ P49605/KAPR_USTMA /ProClass  UniProtKB/Swiss-Prot | cAMP-dependent protein kinase regulatory subunit; ( Short=PKA regulatory subunit) | 525 | Ustilago maydis (Smut fungus) | PIRSF000548 | 241-257;375-391; |

## Links to iProClass and UniProtKB

## Link to NCBI taxonomy

## Link to PIRSF report

Pattern: LIHQGEKAETLYYIVKG

[LIVM]-[VIC]-x(2) -G-[DENQTA]-x-[GAC]-x(2)-[LIVMFY](4)-x (2)-G
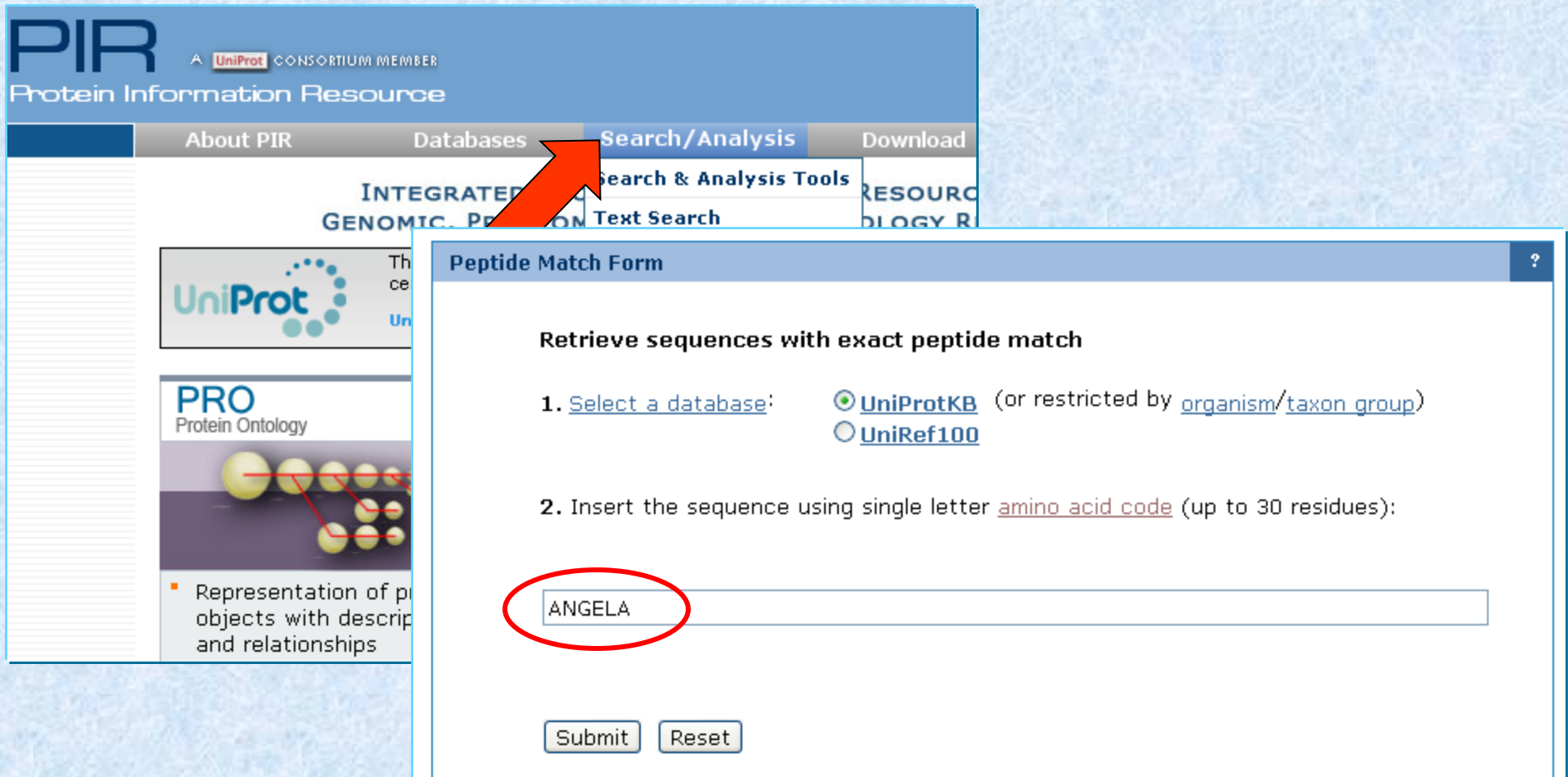
New β-signal motif : $P_o x G h_y x H_y x H_y$
[K,R,H,Q,N,S,T].G [I,V,L,F,M,Y,W,A,C].[I,V,L,F,M,Y,W].[I,V,L,F,M,Y,W]

**Algorithm**

**K. Imai, M.M. Gromiha and P. Horton (2008) Cell**

# PIR: Specific Peptide search

# PIR: Specific Peptide search

# Position specific scoring matrices (Profiles)

Position specific scoring matrices (PSSM) or profiles express the patterns inherent in a multiple sequence alignment of a set of homologous sequences.

The basic idea to use profiles is to match the query sequences from the database against the sequences in the alignment table, giving higher weight to positions that are conserved than to those are variable.

These profiles are obtained with a set of probability scores for each amino acid (or gap) at each position of the alignment.

# Profiles: Applications

(i) they permit greater accuracy in alignments of distantly-related sequences,

(ii) the conservation patterns facilitate identification of other homologous sequences,

(iii) patterns from the sequences are useful in classifying subfamilies within a set of homologues,

(iv) most structure prediction methods are reliable if based on multiple sequence alignment rather than on a single sequence etc.

## a) Alignment Matrix

```
A  A  T  T  G  A
A  G  G  T  C  C
A  G  G  A  T  G
A  G  G  C  G  T
```

|   | 1 | 2 | 3 | 4 | 5 | 6 |
|---|---|---|---|---|---|---|
| A | 4 | 1 | 0 | 1 | 0 | 1 |
| C | 0 | 0 | 0 | 1 | 1 | 1 |
| G | 0 | 3 | 3 | 0 | 2 | 1 |
| T | 0 | 0 | 1 | 2 | 1 | 1 |

consensus: **A G G T G N**

$$\ln \frac{(n_{i,j} + p_i)/(N+1)}{p_i} \approx \ln \frac{f_{i,j}}{p_i}$$

## b) Weight Matrix

|   | 1 | 2 | 3 | 4 | 5 | 6 |
|---|---|---|---|---|---|---|
| A | 1.2 | 0 | -1.6 | 0 | -1.6 | 0 |
| C | -1.6 | -1.6 | -1.6 | 0 | 0 | 0 |
| G | -1.6 | .96 | .96 | -1.6 | .59 | 0 |
| T | -1.6 | -1.6 | 0 | .59 | 0 | 0 |

test sequence: **A G G T G C**

**Fig. 1.** Examples of the simple matrix model for summarizing a DNA alignment. (**a**) An alignment matrix describing the alignment of the four 6-mers on top. The matrix contains the number of times, $n_{i,j}$, that letter $i$ is observed at position $j$ of this alignment. Below the matrix is the consensus sequence corresponding to the alignment (N indicates that there is no nucleotide preference). (**b**) A weight matrix derived from the alignment in (a). The formula used for transforming the alignment matrix to a weight matrix is shown above the arrow. In this formula, $N$ is the total number of sequences (four in this example), $p_i$ is the *a priori* probability of letter $i$ (0.25 for all the bases in this example) and $f_{i,j} = n_{i,j}/N$ is the frequency of letter $i$ at position $j$. The numbers enclosed in blocks are summed to give the overall score of the test sequence. The overall score is 4.3, which is also the maximum possible score with this weight matrix.

# Method



Given amino acid sequence get PSI-BLAST results

**Normalize it**

**X-min/(max-min)**

**E.g. 5, 6, 9, 2, 1**

**(5-1)/(9-1) = 4/8 = 0.5**

**(9-1)/(9-1) = 8/8/ = 1.0**

**(1-1)/(9-1) = 0/8 = 0.0**

# PSSM-400



**Value of LA=** Σ value of L in column A
(shown in bold)

**Value of EC=** Σ value of E in column C
(shown in bold and italics)

EC: 0.08, 0.10, 0.05

LA: 0.24, 0.21, 0.24

QD ?

# Applications

Protein secondary structure prediction

Discrimination of proteins belonging to different classes, types
    etc.

Identifying the binding sites, functionally important residues etc.

# Large scale analysis

**Non redundant sequences**

    No two protein sequences have the sequence identity of more than a specific cutoff (say, 40%).

    Redundancy cause a bias in any analysis.

    E.g. Consider two sequences

    ADIKLAAIKL and KILASDPQWE: Average A is 4/20 = 0.20

    If one of these sequences appear twice:

    ADIKLAAIKL, ADIKLAAIKL and KILASDPQWE: Average A is 7/30= 0.23

                                          **(over-represented)**

    ADIKLAAIKL, KILASDPQWE and KILASDPQWE: Average A is 5/30 = 0.17

                                         **(under-represented)**

# Programs

CD-HIT: **C**luster **D**atabase at **H**igh **I**dentity with **T**olerance.

The program takes a **fasta format sequence** database as input and produces a set of '**non-redundant**' (nr) representative sequences as output.

It uses clustering algorithm and eliminates the redundant sequences.

The main advantages of this program are given below:

   (i) it can handle huge datasets,

   (ii) it is easy to download and

   (iii) the results can be obtained quickly.

CD-HIT can be used to create the non-redundant dataset of less than 40% sequence identity.

CD-HIT

Blastclust

PISCES

**http://cd-hit.org/**

# Algorithm

**Greedy incremental algorithm**: selects representative protein sequence sets

Sequences with the identity of more than the threshold will be discarded.

Longest sequences, the first and proceed with shorter ones.

Sequence identity is the number of identical residues divided by the length of the shorter sequence

# Short-word filtering system

Explicit alignment is time consuming

Algorithm without aligning.

Sequences with >90% sequence identity

Decapeptides: query and database (at least 1)

Pentapepides: 85%

Tetrapeptides: 80%

Tripeptides: 75%

Dipeptides: 65% -> efficiency decreases

Compare word size and number of same words with sequence identity

# Clustering methods

**_k_-means clustering** is a method of cluster analysis which aims to partition _n_ **observations** into _k_ **clusters** in which each observation belongs to the cluster with the nearest mean.

observations ($\mathbf{x}_1, \mathbf{x}_2, \ldots, \mathbf{x}_n$), _k_-means clustering aims to partition the _n_ observations into _k_ sets ($k \le n$)$\mathbf{S} = \{S_1, S_2, \ldots, S_k\}$ so as to minimize the within-cluster sum of squares (WCSS):

$$\arg\min_{\mathbf{S}} \sum_{i=1}^{k} \sum_{\mathbf{x}_j \in S_i} \|\mathbf{x}_j - \boldsymbol{\mu}_i\|^2$$

# Clustering methods based on composition

**Hamming distance**

$$D^H = \Sigma \, |\text{Comp}(1)_i - \text{comp}(2)_i|, \; i=1,20$$

**Euclidean distance**

$$D^E = \{\Sigma \, [\text{Comp}(1)_i - \text{comp}(2)_i]^2\}^{1/2}$$

# CD-HIT Installation

## Installation

**Most CD-HIT programs were written in C++.**

**Download** current CD-HIT at http://bioinformatics.org/cd-hit/

   Example

        cd-hit-v4.5.4-2011-03-07.tgz

**Unpack** the file with

   "**tar xvf** cd-hit-v4.5.4-2011-03-07.tgz **--gunzip**"

**Change directory** by "**cd** cd-hit-v4.5.4-2011-03-07"

**Compile** the programs by "**make**"

**Run** the program

# Run CD-HIT

**./cd-hit -i db -o db90 -c 0.9 -n 5**

    **db: input file name**

    **db90:output file name**

    **0.9, means 90% identity (clustering threshold)**

    **5 is the size of word**

**Choice of word size:**

    **-n 5 for thresholds 0.7 ~ 1.0**

    **-n 4 for thresholds 0.6 ~ 0.7**

    **-n 3 for thresholds 0.5 ~ 0.6**

    **-n 2 for thresholds 0.4 ~ 0.5**

# Example

## ./cd-hit -i hemoglobin_fasta -o db85 -c 0.85 -n 5

```
>sp|P69905|HBA_HUMAN Hemoglobin subunit alpha OS=Homo sapiens GN=HBA1 PE=1 SV=2
MVLSPADKTNVKAAWGKVGAHAGEYGAEALERMFLSFPTTKTYFPHFDLSHGSAQVKGHG
KKVADALTNAVAHVDDMPNALSALSDLHAHKLRVDPVNFKLLSHCLLVTLAAHLPAEFTP
AVHASLDKFLASVSTVLTSKYR
>sp|P01946|HBA_RAT Hemoglobin subunit alpha-1/2 OS=Rattus norvegicus GN=Hba1 PE=1 SV=3
MVLSADDKTNIKNCWGKIGGHGGEYGEEALQRMFAAFPTTKTYFSHIDVSPGSAQVKAHG
KKVADALAKAADHVEDLPGALSTLSDLHAHKLRVDPVNFKFLSHCLLVTLACHHPGDFTP
AMHASLDKFLASVSTVLTSKYR
>sp|P01942|HBA_MOUSE Hemoglobin subunit alpha OS=Mus musculus GN=Hba PE=1 SV=2
MVLSGEDKSNIKAAWGKIGGHGAEYGAEALERMFASFPTTKTYFPHFDVSHGSAQVKGHG
KKVADALASAAGHLDDLPGALSALSDLHAHKLRVDPVNFKLLSHCLLVTLASHHPADFTP
AVHASLDKFLASVSTVLTSKYR
>sp|P01966|HBA_BOVIN Hemoglobin subunit alpha OS=Bos taurus GN=HBA PE=1 SV=2
MVLSAADKGNVKAAWGKVGGHAAEYGAEALERMFLSFPTTKTYFPHFDLSHGSAQVKGHG
AKVAAALTKAVEHLDDLPGALSELSDLHAHKLRVDPVNFKLLSHSLLVTLASHLPSDFTP
AVHASLDKFLANVSTVLTSKYR
>sp|P01958|HBA_HORSE Hemoglobin subunit alpha OS=Equus caballus GN=HBA PE=1 SV=2
MVLSAADKTNVKAAWSKVGGHAGEYGAEALERMFLGFPTTKTYFPHFDLSHGSAQVKAHG
KKVGDALTLAVGHLDDLPGALSNLSDLHAHKLRVDPVNFKLLSHCLLSTLAVHLPNDFTP
AVHASLDKFLSSVSTVLTSKYR
>sp|P69907|HBA_PANTR Hemoglobin subunit alpha OS=Pan troglodyte
MVLSPADKTNVKAAWGKVGAHAGEYGAEALERMFLSFPTTKTYFPHFDLSHGSAQVKGHG
KKVADALTNAVAHVDDMPNALSALSDLHAHKLRVDPVNFKLLSHCLLVTLAAHLPAEFTP
AVHASLDKFLASVSTVLTSKYR
>sp|P01959|HBA_EQUAS Hemoglobin subunit alpha OS=Equus asinus G
MVLSAADKTNVKAAWSKVGGNAGEFGAEALERMFLGFPTTKTYFPHFDLSHGSAQVKAHG
KKVGDALTLAVGHLDDLPGALSNLSDLHAHKLRVDPVNFKLLSHCLLSTLAVHLPNDFTP
AVHASLDKFLSTVSTVLTSKYR
>sp|P01965|HBA_PIG Hemoglobin subunit alpha OS=Sus scrofa GN=HB
VLSAADKANVKAAWGKVGGQAGAHGAEALERMFLGFPTTKTYFPHFNLSHGSDQVKAHGQ
KVADALTKAVGHLDDLPGALSALSDLHAHKLRVDPVNFKLLSHCLLVTLAAHHPDDFNPS
VHASLDKFLANVSTVLTSKYR
>sp|P06635|HBA_PONPY Hemoglobin subunit alpha OS=Pongo pygmaeus
MVLSPADKTNVKTAWGKVGAHAGDYGAEALERMFLSFPTTKTYFPHFDLSHGSAQVKDHG
KKVADALTNAVAHVDDMPNALSALSDLHAHKLRVDPVNFKLLSHCLLVTLAAHLPAEFTP
AVHASLDKFLASVSTVLTSKYR
>sp|P60529|HBA_CANFA Hemoglobin subunit alpha OS=Canis familiar
VLSPADKTNIKSTWDKIGGHAGDYGGEALDRTFQSFPTTKTYFPHFDLSPGSAQVKAHGK
KVADALTTAVAHLDDLPGALSALSDLHAYKLRVDPVNFKLLSHCLLVTLACHHPTEFTPA
VHASLDKFFAAVSTVLTSKYR
```

```
[gromiha@INSIGHT1 cd-hit-v4.5.4-2011-03-07]$ more db85
>sp|P69905|HBA_HUMAN Hemoglobin subunit alpha OS=Homo sapiens GN=HBA1 PE=1 SV=2
MVLSPADKTNVKAAWGKVGAHAGEYGAEALERMFLSFPTTKTYFPHFDLSHGSAQVKGHG
KKVADALTNAVAHVDDMPNALSALSDLHAHKLRVDPVNFKLLSHCLLVTLAAHLPAEFTP
AVHASLDKFLASVSTVLTSKYR
>sp|P01946|HBA_RAT Hemoglobin subunit alpha-1/2 OS=Rattus norvegicus GN=Hba1 PE=1 SV=3
MVLSADDKTNIKNCWGKIGGHGGEYGEEALQRMFAAFPTTKTYFSHIDVSPGSAQVKAHG
KKVADALAKAADHVEDLPGALSTLSDLHAHKLRVDPVNFKFLSHCLLVTLACHHPGDFTP
AMHASLDKFLASVSTVLTSKYR
>sp|P01965|HBA_PIG Hemoglobin subunit alpha OS=Sus scrofa GN=HBA PE=1 SV=1
VLSAADKANVKAAWGKVGGQAGAHGAEALERMFLGFPTTKTYFPHFNLSHGSDQVKAHGQ
KVADALTKAVGHLDDLPGALSALSDLHAHKLRVDPVNFKLLSHCLLVTLAAHHPDDFNPS
VHASLDKFLANVSTVLTSKYR
>sp|P60529|HBA_CANFA Hemoglobin subunit alpha OS=Canis familiaris GN=HBA PE=1 SV=1
VLSPADKTNIKSTWDKIGGHAGDYGGEALDRTFQSFPTTKTYFPHFDLSPGSAQVKAHGK
KVADALTTAVAHLDDLPGALSALSDLHAYKLRVDPVNFKLLSHCLLVTLACHHPTEFTPA
VHASLDKFFAAVSTVLTSKYR
[gromiha@INSIGHT1 cd-hit-v4.5.4-2011-03-07]$
```

Connected to 10.93.219.140                                    SSH2 - aes12

# Blastclust

**Blastclust** is a program within the standalone BLAST package used to cluster either protein or nucleotide sequences.

The program begins with pairwise matches and places a sequence in a cluster if the sequence matches at least one sequence already in the cluster.

 In the case of proteins, the blastp algorithm is used to compute the pairwise matches.

The general command to create a set of non-redundant set of protein sequences is *blastclust -i infile -o outfile -p T -L .9 -b T -S 95*,

where <u>infile</u> and <u>outfile</u> are input and output files, respectively.

T stands for protein;

the coverage of the length and sequence identity cutoff are 90% (-L .9) and 95% (-S 95), respectively.

# PISCES

PISCES is a protein sequence culling server to produce subsets of non-redundant sequences using Protein Data Bank entries or Uniprot sequences in FASTA format.

Sequence identities for PDB sequences are determined by the combination of Combinatorial Extension structural alignment and PSI-BLAST alignment.

non-PDB sequences are culled with sequence identities from PSI-BLAST. PISCES does not search the non-redundant sequence database, but rather use the user's input sequences as the database.

This server will usually be used to cull a related set of sequences, for instance those from a PSI-BLAST search.

It takes the amino acid sequence in FASTA format and sends the list of non-redundant protein sequences by e-mail.

http://dunbrack.fccc.edu/pisces/

# PISCES

```
> 7467903|Genbank|Outer membrane integral membrane
MSKFTITIFITTLLFTGSVIALDLEQALTEGYKNNEELKAAQIKFLNAIE
QFPQAFSGFMPNVGLQINRQNSKTKYNKKYVNRLGITPRETASTQGILTI
EQSLFNGGASIAALKAAQSGFRASRSEYYAGEQKVLLNLITAYLDCVESK
EKYDISESRVRTNIQQVKTVEEKLRLGEATAIDIAAARAGLAAAETNKLA
AYADFQGKKANFIKVFGIEANDITMPDLPDRLPISLDEFTRKAAKFNPDI
NSARHNVTVTKALEMVQKGKLLPQVSVKLLSGGTNYNPQEPVIQNINNRI
YTTTLSVNIPIYPEGGAQYSRIRSAKNQTRNSVVQLDSAIKQIKAGVVSV
WEGFETAKSRIVAANQGVEAAQISYNGIVQEEIVGSKTILDVLDAEQKLY
EAKITRVDAYKNSVLASYQMKLLTGELTAKSLKLKVKYFSPEEEFNNLKK
KMFIGF
> 11559475|Genbank|Outer membrane integral membrane
MTRNRFVMRRIATTLLVAGIIVSQAAYAQVTLNFVNADIDQVAKAIGAAT
GKTIIVDPRVKGQLNLVAERPVPEDQALKTLQSALRMQGFALVQDHGVLK
VVPEADAKLQGVPTYIGNAPQARGDQVITQVFELHNESANNLLPVLRPLI
SPNNTVTAYPANNTIVVTDYADNVRRIAQIISGVDSAAGAQVQVVPLRNA
NAIDLAAQLQKMLDPGAIGNSDATLKVSVTADPRTNALLLRASNASRLAA
AKRLVQQLDAPSAVPGNMHVVPLRNADAVKLAKTLRGMLGKGGNDSGSSA
SSNDANSFNQNGGSSASGNFSTGTSGTPPLPSGGLGGSSSSSYGGSGGSS
GGGLGTGGLLGGDKDKSGDDNQPGGMIQADSATNSLIITASDPVYRNLRS
VIDQLDARRAQVYIEALIVELNSTTQGNLGIQWQVASGQFLGGTNLAPTA
GNGLGNSIINLTAGGLTNAAGGITGGGLASNLGQLSQGLNIGWLHNMFGV
QGLGALLQYFAGVSDANVLSTPNLITLDNEEAKIVVGQNVPIATGSYSNL
TSGTTSNAFNTYDRRDVGLTLHVKPQITDGGILKLQLYTEDSAVVNGTTN
SQTGPTFTKRSIQSTILADNGEIIVLGGLMQDNYQVSNSKVPLLGDIPWI
GQLFRSESKVRAKTNLMVFLRPVIISDRSTAQEVTSNRYDYIQGVTGAYK
SDNNVIRDKDDPVVPPMPLGPSQGGTAAGNLFDLDKMRRQQLQRQVVPVP
AQPLPEATPAQPQGVPLQAVPQQPLTTAPGASQ
> 7469324|Genbank|Outer membrane integral membrane
MRSNSVKNFRFWLTTEIATCCLLALAPAQAETVSQSNTLDGDLRTAIAGD
SSRDWLQFEKSLEQSLKQKEEIDSWKPSLELMQAKSLVKPGQKLTNIELL
VQELEALSDPLALNFPEPNQTSVAQMAPPSRPMPPPPAGSGQVMFPNPEI
IIQQQGGVPQRGASPQVGNPSILSPAVPVAPVRSRAVPPPVGDLAISNIN
ASFDMIDLGQRGQVNVPSLVLREAPAREVLAVLTRYAGMNLIFTDNQNNE
GTPTPGTPPGGQVAPPQAQSTITLDIQNESVQDVFNYVLMASGLKASRRG
NTIFAGANLLPSARNIITRTIRLNQASAESVASTLASQGAEVNILFEGQE
DVQLAENAPPRVIKQPPTLVPLTVQKPANDSSVLILEGLVVSTDPRLNTV
TLVGEPRNVELASSMITQMDARRQVAVNVKIIDINLNNIQDYDSSFSFG
IGDSFFVQDSGSAVMRFGDTAPVQEIDINNNLGRITNPPAIVNPFQDGEI
FFDLNRITNIEVPLGPGTIPINFFTSGSGAVSNNPLFNGVTEFPIVEVDE
QGLLTITQPEFGLPSFYQYPKKFQAQIDAQIRSGNAKILTDPTLIVQEGE
AAQVKLTESVIASVDTQVDTQGDTAVRTITPVLEDVGLTLNVIVDRIDDN
GFITLRVNPIVASPAGTQVFDSGAGAINEITLINKRELTSGVVRLRDDQT
FILSGIISELQRSTTSKVPILGDLPVIGALFRQSTDTTDRSEVIILMTPK
IIHDSTEAQFGFRYNPDAATAEFLRQKGFPVQAQP
```

# Sequence id

| IDs | length |
| --- | --- |
| 7467903|Genbank|Outer | 456 |
| 11559475|Genbank|Outer | 783 |
| 7469324|Genbank|Outer | 785 |
| 15596906|Genbank|Outer | 295 |
| P26466|SwissProt|Outer | 452 |
| 15597487|Genbank|Outer | 452 |
| 5640161|Genbank|Outer | 889 |
| P13949|SwissProt|Outer | 201 |
| P10170|SwissProt|Outer | 260 |
| P16945|SwissProt|Outer | 292 |
| 7208425|Genbank|Outer | 560 |
| 15598604|Genbank|Outer | 891 |
| 13470835|Genbank|Outer | 794 |
| P19196|SwissProt|Outer | 835 |
| 12620518|Genbank|Outer | 230 |
| P31600|SwissProt|Outer | 990 |
| 15596468|Genbank|Outer | 616 |
| P15727|SwissProt|Outer | 482 |
| 7520765|Genbank|Outer | 778 |
| P16466|SwissProt|Outer | 1577 |
| 3228547|Genbank|Outer | 700 |
| P06970|SwissProt|Outer | 812 |
| P22340|SwissProt|Outer | 505 |
| P44601|SwissProt|Outer | 565 |
| P35077|SwissProt|Outer | 584 |
| 12721580|Genbank|Outer | 444 |
| 7470479|Genbank|Outer | 654 |
| P13794|SwissProt|Outer | 350 |
| P06111|SwissProt|Outer | 257 |
| P24126|SwissProt|Outer | 530 |
| P29041|SwissProt|Outer | 759 |
| 5759281|Genbank|Outer | 462 |