All things protein...

Marylyn D Ritchie, PhD Professor, Biochemistry and Molecular Biology Director, Center for Systems Genomics The Pennsylvania State University





Amino Acid Sequence

ORIGIN

/translation="MANFLLPRGTSSFRRFTRESLAAIEKRMAEKQA GLPEEEAPRPQLDLQASKKLPDLYGNPPQELIGEPLEDLDPFYSTQK FRFSATNALYVLSPFHPIRRAAVKILVHSLFNMLIMCTILTNCVFMA EYTFTAIYTFESLVKILARGFCLHAFTFLRDPWNWLDFSVIIMAYTT RTFRVLRALKTISVISGLKTIVGALIQSVKKLADVMVLTVFCLSVFA RHKCVRNFTALNGTNGSVEADGLVWESLDLYLSDPENYLLKNGTSDV CPEGYRCLKAGENPDHGYTSFDSFAWAFLALFRLMTQDCWERLYQQT FMLVIFLGSFYLVNLILAVVAMAYEEQNQATIAETEEKEKRFQEAME RGVDTVSRSSLEMSPLAPVNSHERRSKRRKRMSSGTEECGEDRLPKS LSLTRGLSRTSMKPRSSRGSIFTFRRRDLGSEADFADDENSTAGESE LRRTSAQGQPSPGTSAPGHALHGKKNSTVDCNGVVSLLGAGDPEATS EHPPDTTTPSEEPGGPQMLTSQAPCVDGFEEPGARQRALSAVSVLTS CPPCWNRLAQRYLIWECCPLWMSIKQGVKLVVMDPFTDLTITMCIVL MTSEFEEMLQVGNLVFTGIFTAEMTFKIIALDPYYYFQQGWNIFDSI SRMSNLSVLRSFRLLRVFKLAKSWPTLNTLIKIIGNSVGALGNLTLV GMQLFGKNYSELRDSDSGLLPRWHMMDFFHAFLIIFRILCGEWIETM CLLVFLLVMVIGNLVVLNLFLALLLSSFSADNLTAPDEDREMNNLQL VKRTTWDFCCGLLRQRPQKPAALAAQGQLPSCIATPYSPPPPETEKV GEQPGQGTPGDPEPVCVPIAVAESDTDDQEEDEENSLGTEEESSKQQ PPDSRTWSQVSATASSEAEASASQADWRQQWKAEPQAPGCGETPEDS TAELLEQIPDLGQDVKDPEDCFTEGCVRRCPCCAVDTTQAPGKVWWR SWFETFIIFMILLSSGALAFEDIYLEERKTIKVLLEYADKMFTYVFV FKKYFTNAWCWLDFLIVDVSLVSLVANTLGFAEMGPIKSLRTLRALR RVVVNALVGAIPSIMNVLLVCLIFWLIFSIMGVNLFAGKFGRCINQT NNKSQCESLNLTGELYWTKVKVNFDNVGAGYLALLQVATFKGWMDIM QPQWEYNLYMYIYFVIFIIFGSFFTLNLFIGVIIDNFNQQKKKLGGQ YNAMKKLGSKKPQKPIPRPLNKYQGFIFDIVTKQAFDVTIMFLICLN SPEKINILAKINLLFVAIFTGECIVKLAALRHYYFTNSWNIFDFVVV IIOKYFFSPTLFRVIRLARIGRILRLIRGAKGIRTLLFALMMSLPAL FIYSIFGMANFAYVKWEAGIDDMFNFQTFANSMLCLFQITTSAGWDG YCDPTLPNSNGSRGDCGSPAVGILFFTTYIIISFLIVVNMYIAIILE PLSEDDFDMFYEIWEKFDPEATQFIEYSVLSDFADALSEPLRIAKPN VSGDRIHCMDILFAFTKRVLGESGEMDALKIQMEEKFMAANPSKISY EEVSAMVIQRAFRRHLLQRSLKHASFLFRQQAGSGLSEEDAPEREGL PLGPPSSSSISSTSFPPSYDSVTRATSDNLQVRGSDYSHSEDLADFP

1 agacggcggc ggcgcccgta ggatgcaggg atcgctcccc cggggccgct gagcctgcgc 61 ccagtgcccc gagccccgcg ccgagccgag tccgcgccaa gcagcagccg cccaccccgg 121 ggcccggccg ggggaccagc agcttcccca caggcaacgt gaggagagcc tgtgcccaga 181 agcaggatga gaagatggca aacttcctat tacctcgggg caccagcagc ttccgcaggt 241 tcacacggga gtccctggca gccatcgaga agcgcatggc agagaagcaa gcccgcggct 301 caaccacctt gcaggagagc cgagaggggc tgcccgagga ggaggctccc cggccccagc 361 tggacctgca ggcctccaaa aagctgccag atctctatgg caatccaccc caagagctca 421 tcggagagcc cctggaggac ctggacccct tctatagcac ccaaaagact ttcatcgtac 481 tgaataaagg caagaccatc ttccggttca gtgccaccaa cgccttgtat gtcctcagtc 541 ccttccaccc catccggaga gcggctgtga agattctggt tcactcgctc ttcaacatgc 601 tcatcatgtg caccatcctc accaactgcg tgttcatggc ccagcacgac cctccaccct 661 ggaccaagta tgtcgagtac accttcaccg ccatttacac ctttgagtct ctggtcaaga 721 ttctggctcg aggcttctgc ctgcacgcgt tcactttcct tcgggaccca tggaactggc 781 tggactttag tgtgattatc atggcataca caactgaatt tgtggacctg ggcaatgtct 841 cagcettacg cacettecga gteetecggg ceetgaaaae tatateagte attteagge 901 tgaagaccat cgtgggggcc ctgatccagt ctgtgaagaa gctggctgat gtgatggtcc 961 tcacagtctt ctgcctcagc gtctttgccc tcatcggcct gcagctcttc atgggcaacc 1021 taaggcacaa gtgcgtgcgc aacttcacag cgctcaacgg caccaacggc tccgtggagg 1081 ccgacggctt ggtctgggaa tccctggacc tttacctcag tgatccagaa aattacctgc 1141 tcaagaacgg cacctctgat gtgttactgt gtgggaacag ctctgacgct gggacatgtc 1201 cggagggcta ccggtgccta aaggcaggcg agaaccccga ccacggctac accagcttcg 1261 atteettige etgggeettt ettgeactet teegeetgat gaegeaggae tgetgggage 1321 gcctctatca gcagaccctc aggtccgcag ggaagatcta catgatcttc ttcatgcttg 1381 tcatcttcct ggggtccttc tacctggtga acctgatcct ggccgtggtc gcaatggcct 1441 atgaggagca aaaccaagcc accatcgctg agaccgagga gaaggaaaag cgcttccagg 1501 aggccatgga aatgctcaag aaagaacacg aggccctcac catcaggggt gtggataccg 1561 tgtcccgtag ctccttggag atgtcccctt tggccccagt aaacagccat gagagaagaa 1621 gcaagaggag aaaacggatg tcttcaggaa ctgaggagtg tggggaggac aggctcccca 1681 agtctgactc agaagatggt cccagagcaa tgaatcatct cagcctcacc cgtggcctca 1741 gcaggacttc tatgaagcca cgttccagcc gcgggagcat tttcaccttt cgcaggcgag 1801 acctgggttc tgaagcagat tttgcagatg atgaaaacag cacagcgggg gagagcgaga 1861 gccaccacac atcactgctg gtgccctggc ccctgcgccg gaccagtgcc cagggacagc 1921 ccagtcccgg aacctcggct cctggccacg ccctccatgg caaaaagaac agcactgtgg 1981 actgcaatgg ggtggtctca ttactggggg caggcgaccc agaggccaca tccccaggaa 2041 gccacctcct ccgccctgtg atgctagagc acccgccaga cacgaccacg ccatcggagg

2101 agccaggcgg gccccagatg ctgacctccc aggctccgtg tgtagatggc ttcgaggagc 2161 caggagcacg gcagcgggcc ctcagcgcag tcagcgtcct caccagcgca ctggaagagt 2221 tagaggagtc tcgccacaag tgtccaccat gctggaaccg tctcgcccag cgctacctga

🗟 NCBI 🛛 Resources 🗹 How To 🖂 All Databases 🔹 Search Biotechnology Information Proteins NCBI Home **Quick Links** BioProject (formerly Genome Resource List (A-Z) All Databases Downloads Submissions Tools How To Project) All Resources BioSystems Databases Chemicals & Bioassays Conserved Domain Database **BioProject (formerly Genome Project)** Data & Software (CDD) A collection of genomics, functional genomics, and genetics studies and links to their resulting datasets. This resource Protein Clusters DNA & RNA describes project scope, material, and objectives and provides a mechanism to retrieve datasets that are often difficult to find due to inconsistent annotation, multiple independent submissions, and the varied nature of diverse data types which Protein Database **Domains & Structures** are often stored in different databases. Reference Sequence (RefSeq) Genes & Expression BLAST (Stand-alone) **BioSystems** Genetics & Medicine Database that groups biomedical literature, small molecules, and sequence data in terms of biological relationships. BLAST Link (BLink) Genomes & Maps Basic Local Alignment Search Conserved Domain Database (CDD) Homology Tool (BLAST) A collection of sequence alignments and profiles representing protein domains conserved in molecular evolution. It also Literature includes alignments of the domains to known 3-dimensional protein structures in the MMDB database. Cn3D Proteins Conserved Domain Search HIV-1, Human Protein Interaction Database Service (CD Search) Sequence Analysis A database of known interactions of HIV-1 proteins with proteins from human hosts. It provides annotated bibliographies E-Utilities of published reports of protein interactions, with links to the corresponding PubMed records and sequence data. Taxonomy ProSplign Training & Tutorials Protein Clusters A collection of related protein sequences (clusters), consisting of Reference Sequence proteins encoded by complete Variation prokaryotic and organelle plasmids and genomes. The database provides easy access to annotation information, publications, domains, structures, external links, and analysis tools.

Protein Database

A database that includes protein sequence records from a variety of sources, including GenPept, RefSeq, Swiss-Prot, PIR, PRF, and PDB.

Sign in to NCBI

S NCBI Resources ⊙	How To 🗵		Sign in to NCBI
Protein	Protein Advance	ed	Search
QDIVEQI	RKETEKE	RT Protein	
DVGKKA VKPRVT	REGITTKR	The Protein database is a collection regions in GenBank, RefSeq and The fundamental determinants of bio	of sequences from several sources, including translations from annotated coding PA, as well as records from SwissProt, PIR, PRF, and PDB. Protein sequences are plogical structure and function.
Using Protein		Protein Tools	Other Resources
Quick Start Guide		BLAST	GenBank Home
FAQ		LinkOut	RefSeq Home
<u>Help</u>		<u>E-Utilities</u>	CDD
GenBank FTP		Blink	Structure
RefSeq FTP		Batch Entrez	

S NCBI R	Resources 🕑 How To 🕑		<u>Sign</u>	in to NCBI
Protein	Protein Advanced		Search	Help
Display Settin	ngs:	<u>Send to:</u>	ange region shown	•
GenBank: AA <u>FASTA</u> Gra	Al44622.1 uphics	Cu	stomize view	•
Go to: LOCUS DEFINITION ACCESSION VERSION DBSOURCE KEYWORDS SOURCE	AAI44622 1983 aa linear PRI 18-MAR-2009 SCN5A protein [Homo sapiens]. AAI44622 AAI44622.1 GI:219521582 accession <u>BC144621.1</u> MGC. Homo sapiens (human)	An: Rur Ider Hig Fin	alyze this sequence BLAST ntify Conserved Domains hlight Sequence Features d in this Sequence	
REFERENCE AUTHORS	Homo sapiens Eukaryota; Metazoa; Chordata; Craniata; Vertebrata; Euteleostomi; Mammalia; Eutheria; Euarchontoglires; Primates; Haplorrhini; Catarrhini; Hominidae; Homo. 1 (residues 1 to 1983) Strausberg,R.L., Feingold,E.A., Grouse,L.H., Derge,J.G., Klausner,R.D., Collins,F.S., Wagner,L., Shenmen,C.M., Schuler,G.D., Altschul,S.F., Zeeberg,B., Buetow,K.H., Schaefer,C.F., Bhat,N.K., Hopkins,R.F., Jordan,H., Moore,T., Max,S.I., Wang,J., Hsieh,F., Diatchenko,L., Marusina,K., Farmer,A.A., Rubin,G.M., Hong,L., Stapleton,M., Soares,M.B., Bonaldo,M.F., Casavant,T.L., Scheetz,T.E., Brownstein,M.J., Usdin,T.B., Toshiyuki,S.,	Pro	rtein 3D Structure Voltage-gated Sod Channel 1.5 C-terr Domain In Comple PDB: 4JQ0 Source: Homo sa Method: X-Ray D Resolution: 3.84 See all 4 s	ium ninal x With apiens biffraction Å

CBI/ BLAST Home			
BLAST finds regions of similarity between biological sequences. mor	e		
New DELTA-BLAST, a r	nore sensitive protein-	protein search 💿	
BLAST Assembled Genomes			
Find Genomic BLAST pages:	□ <u>Human</u>	□ <u>Rabbit</u>	Zebrafish Clawod from
Enter organism name or idcompletions will be suggested		□ Guinea pig	Arabidopsi
	□ <u>Cow</u>	Fruit fly	□ <u>Rice</u>
	□ <u>Piq</u>	Honey bee	Yeast
	Dog	Chicken	□ <u>Microbes</u>
Basic BLAST			

	Algorithms: blastn, megablast, discontiguous megablast
<u>protein blast</u>	Search protein database using a protein query Algorithms: blastp, psi-blast, phi-blast, delta-blast
<u>blastx</u>	Search protein database using a translated nucleotide query
<u>tblastn</u>	Search translated nucleotide database using a protein query
<u>tblastx</u>	Search translated nucleotide database using a translated nucleotide query



About UniProt

The mission of UniProt is to provide the scientific community with a comprehensive, high quality and freely accessible resource of protein sequence and functional information.

UniProt is comprised of four components, each optimised for different uses:

1) The **UniProt Knowledgebase (UniProtKB)** is the central access point for extensive curated protein information, including function, classification, and cross-reference.

It consists of two sections:

- · UniProtKB/Swiss-Prot which is manually annotated and is reviewed and
- UniProtKB/TrEMBL which is automatically annotated and is not reviewed.

Links

Download Centre Release Statistics UniProt DAS Server QuickGO Posters

Submissions (SPIN) DAS Server BLAST ClustalW2 ID mapping UniSave

http://www.ebi.ac.uk/uniprot

- Mission of UniProt is to provide comprehensive, high quality and freely accessible resource of protein sequence and functional information
- UniProt Knowledgebase (UniProtKB)
 - central access point for extensive curated protein information, including function, classification, and crossreference
 - UniProtKB/Swiss-Prot which is manually annotated and is reviewed
 - UniProtKB/TrEMBL which is automatically annotated and is not reviewed

- UniProt Reference Clusters (UniRef) databases
 - clustered sets of sequences from the UniProtKB and selected UniProt Archive records to obtain complete coverage of sequence space at several resolutions while hiding redundant sequences
- UniProt Archive (UniParc)
 - comprehensive repository, used to keep track of sequences and their identifiers
- UniProt Metagenomic and Environmental Sequences (UniMES)
 - repository specifically developed for metagenomic and environmental data.

UniProt	UniProtKB 🗸	Advanced - Q
BLAST Align Retrieve/ID Mapping		Help Contact
How to use this tool The Basic Local Alignment Search Tool (BLAST) finds regions of between sequences, which can be used to infer functional and e relationships between sequences as well as help identify memb families.	f local similarity evolutionary ers of gene	 Enter either a protein or nucleotide sequence or a UniProt identifier (e.g.P00750 or A4_HUMAN or UPI000000001) into the form field. Optionally, change the program parameters with the dropdown menus under the form. Click the <i>Run BLAST</i> button.
		😮 Help 🛛 Tutorials and Videos 👌 Downloads
BLAST		
Protein sequence, Nucleotide sequence or UniProt identifier		
Target database E-Threshold Matrix	Filtering	Gapped Hits
UniProtKB ID Auto	▼ None	▼ yes ▼ 250 ▼
Run Blast in a separate window.		🔧 Run BLAST Clear

	UniProtKB -	Advanced - Q			
BLAST Align Retrieve/ID Mapping			Help Contact		
How to use this tool Align two or more protein sequences with the Clustal Omega program (see also this FAQ) to view their characteristics alongside each other.		 Enter either protein sequences in FASTA format or UniProt identifiers into the form field, for example: TPA_HUMAN TPA_PIG 			
Align		2. Click the <i>kun Align</i> button.	Tutorials and Videos Downloads		

Protein sequences (FASTA) or UniProt identifiers



Aligning multiple sequences

- Highlights areas of similarity which may be associated with specific features that have been more highly conserved than other regions
 - These regions in turn can help classify sequences or to inform experiment design
- important step for phylogenetic analysis, which aims to model the substitutions that have occurred over evolution and derive the evolutionary relationships between sequences
- Clustal Omega improves on ClustalW in a number of ways - alignment accuracy and improved scaling to many sequences are the main results

- Single worldwide repository of information about the 3D structures of large biological molecules, including proteins and nucleic acids
- Understanding the shape of a molecule helps to understand how it works
- Can be used to help deduce a structure's role in human health and disease, and in drug development
- Structures in the archive range from tiny proteins and bits of DNA to complex molecular machines like the ribosome
- PDB archive is available at no cost to users
- PDB archive is updated each week at the target time of Wednesday 00:00 UTC (Coordinated Universal Time)
- Most recent release is timestamped and linked on every page in the top right header.





4 Structure Hits 3 Citations 3 Ligand Hits 7 Web Page Hit	5				
Query Parameters:				Query	Details Save Query to MyP[
Text Search for: scn5a Other search suggestions:					
Gene View Protein Feature View • SCN5A - sodium (5) • SCN5A (5)					
					clos
Query Refinements: Select an item or pie chart 😮					Hie
Organism	Experimental Method	X-ray Resolution	Release Date	Polymer Type	Protein Symmetry
 Homo sapiens only (4) Eukaryota only (4) 	X-ray (3)Solution NMR (1)	 less than 1.5 Å (1) 2.0 - 2.5 Å (1) 3.0 and more Å (1) more choices 	 2005 - 2010 (1) 2010 - today (3) this year (1) more choices 	• Protein (4)	Asymmetric (4)more choices
Protein Stoichiometry					

- Monomer (2)
- Heteromer (2)
 more choices...

UniProtKB: Q14524 🖉



38340 Structure Hits 14104 Citations	5605 Ligand Hits				
Query Parameters:					Query Details Save Quer
TaxonomyTree Search for Bacteria (eu	ibacteria)				
Query Refinements: Select an item	n or pie chart 😧				
Organism	E Taxonomy	Experimental Method	X-ray Resolution	Release Date	Polymer Type
 Escherichia coli (4888) Escherichia coli K-12 (2076) Bacillus subtilis (1037) Thermus thermophilus HB8 (995) Mycobacterium tuberculosis (925) Staphylococcus aureus (798) Pseudomonas aeruginosa (677) Other (27135) 	 Bacteria only (37229) Bacteria/Eukaryota (770) Bacteria/Other (182) Bacteria/Viruses (87) Bacteria/Archaea (32) Bacteria/Eukaryota/Viruses (14) Bacteria/Eukaryota/Other (12) Other (14) 	X-ray (36100) Solution NMR (1907) Electron Microscopy (270) Hybrid (28) Solid-State NMR (19) Neutron Diffraction (9) Electron Crystallography (6) Fiber Diffraction (1)	 less than 1.5 Å (2857) 1.5 - 2.0 Å (12776) 2.0 - 2.5 Å (12172) 2.5 - 3.0 Å (5737) 3.0 and more Å (2569) more choices 	 before 2000 (3123) 2000 - 2005 (6451) 2005 - 2010 (12673) 2010 - today (16093) this year (2752) this month (84) more choices 	 Protein (36115) Mixed (2102) RNA (120)
Enzyme Classification	SCOP Classification	Protein Symmetry	Protein Stoichio	metry 🌔 Membr	ane Proteins
 3: Hydrolases (7032) 2: Transferases (5586) 1: Oxidoreductases (4683) 4: Lyases (2077) 5: Isomerases (1313) 6: Ligases (960) 	 Alpha and beta proteins (a/b) (6190) Alpha and beta proteins (a+b) (4077) All beta proteins (3250) All alpha proteins (2485) Multi-domain proteins (alpha an (569) Small proteins (379) Membrane and cell surface prote (375) Other (338) 	 Asymmetric (18125) Cyclic (14374) Dihedral (4600) Tetrahedral (225) Octahedral (76) Helical (46) Tcosahedral (12) more choices 	 Homomer (18766) Monomer (15804) Heteromer (2888) more choices 	 ALPHA-HELIC BETA-BARRE MONOTOPIC 	CAL (1032) L (331) MEMBRANE PROTEINS (84)

5666 Structure Hits 260 Unreleased Structures 2653 Citations 1501 Ligand Hits 269 Web Page Hits

Query Parameters:

Text Search for: virus

Other search suggestions:

Molecule Name	Structural Domains	Membrane Proteins	Retrieve	Molecule of the Month	Author
 SINDBIS VIRUS CAPSID (12) Hepatitis Delta virus ribozyme (8) Herpes virus entry (7) Herpes virus entry (7) Epstein-Barr virus receptor (6) Virus-induced (6) More - Find all 	 Human immunodeficiency virus type (246) Vaccinia Virus protein (445) Multimerization sendai virus [SCOP] (1) Hepatitis C Virus Capsid (1) Multimerization sendai virus [SCOP] (1) Hepatitis C Virus Capsid (1) More 	• Virus Coat (25)	Virus Structures	 Adenovirus [virus] Tobacco Mosaic Virus Simian Virus 40 Dengue Virus 	• Virus, C. (1)
Organism • Human immunodeficiency virus 1 (14) • Influenza A virus (576) • Hepatitis C virus (332) • Vaccinia virus (97) • African fever virus Malawi (1) • Woodchuck hepatitis virus 8 (1) More	t59) Eukaryota (53213) Bacteria (38340) Viruses (6906) Archaea (3830) Unassigned (2909) Other (964)	t Journal • Virus Re: • Virus Re: • Viruses	s PF00073 s. PF00506 PF00506 PF00519 PF00606 PF00693 PF00695 More	- picorna virus capsid (162 - Influenza virus nucleoprote - Papilloma virus helicase (9) - Herpes virus Glycoprotein B herpes virus (38) hepadna virus) in (16) (11)
Ontology Terms					

Query Details | Save Query to M

- regulation ... response to virus by virus ... (544)
- virus tail ... (4)
- virus receptor ... (30)
- Aura virus [Genome ... (3)
 B04.820 ... Influenza A virus [MeSH ... (426)
- fusion of virus membrane ... (284)

Pfam



HOME | SEARCH | BROWSE | FTP | HELP | ABOUT

Pfam 27.0 (March 2013, 14831 families)



The Pfam database is a large collection of protein families, each represented by **multiple sequence alignments** and **hidden Markov models (HMMs)**. <u>More...</u>

QUICK LINKS	YOU CAN FIND DATA IN PFAM IN VARIOUS WAYS		
SEQUENCE SEARCH	Analyze your protein sequence for Pfam matches		
VIEW A PFAM FAMILY	View Pfam family annotation and alignments		
VIEW A CLAN	See groups of related families		
VIEW A SEQUENCE	Look at the domain organisation of a protein sequence		
VIEW A STRUCTURE	Find the domains on a PDB structure		
KEYWORD SEARCH	Query Pfam by keywords		
JUMP TO	enter any accession or ID Go Example Enter any type of accession or ID to jump to the page for a Pfam family or clan, UniProt sequence, PDB structure, etc.		

Or view the help pages for more information

Recent Pfam blog & posts Moving to xfam.org & (posted 1 May 2014)

⊠Hide this

http://pfam.xfam.org/

DIP



http://dip.doe-mbi.ucla.edu/dip/Main.cgi

$\mathsf{D}\mathsf{M}^2$



Worldwide PDB





PDB Archive Snanshots:

Questions???

