

Analysis of the tertiary protein structure for Protein Misfolding Diseases

Panagiotis Vlamos

vlamosp@gmail.com Professor, Department of Head of the Director of Bioinformatics and Human Electrophysiology Lab (BiHELab) Chair of Hellenic Initiative Against Alzheimer's (HIAAD)

1



Panagiotis Vlamos

Panagiotis Vlamos is a Full **Professor,** Chairman of the University Research Center of the Ionian University and Director of the Laboratory of Bioinformatics and Human Electrophysiology **(BiHELab)** of the Ionian University.

He completed his undergraduate studies in Mathematics at the University of Athens and obtained his Doctorate in Applied Mathematics from the National Technical University of Athens. He has received many awards and has authored more than 250 papers in scientific journals and conference proceedings, as well as 16 educational books. His research goal is to help bridge the translation gap from data to models and from models to drug discovery and personalized therapy, developing novel approaches to biological and clinical problems, using high-performance computational methods.

Professor Vlamos is a pioneer in the field of Computational Biomarkers and he was the first to introduce the notion of Metabiomarkers. In addition, Professor Vlamos is president of the Hellenic Initiative against Alzheimer's (HIAAD), a joint collaboration between researchers from Johns Hopkins University in the USA and scientists from BiHELab who combine their knowledge and expertise in order to tackle the "epidemic" of Alzheimer's Disease and related disorders in Greece. He is the founder and Head of the Program Committee of the biannual World Congress on "Genetics, Geriatrics and Research in Neurodegenerative Diseases" **(GeNeDis).** Professor Vlamos is one of the three editors of the Handbook of Computational Neurodegeneration to be published by Springer International Publishing and the founder of the MSc Program "Bioinformatics & Neuroinformatics".

BiHeLab's Main Interests:

- **BiHELab's** goal is to help bridge the translational gap from data to models and from models to drug discovery and personalized therapy by fostering collaborations and developing original quantitative approaches to biological and clinical problems.
- The *Bioinformatics and Human Electrophysiology Lab (BiHELab)* focuses on recent advances in geriatrics and neurodegeneration, ranging from basic science to clinical and pharmaceutical developments.
- The **main objective** of the laboratory is the creation of new and effective protocols for the diagnosis of various types of dementia and neurological disorders through identification, mapping, biological analysis as well as mathematical modelling and simulation of all factors associated with these diseases, in order to improve existing techniques and to design new targeted treatments.



Protein Misfolding Diseases

Human diseases caused by defects in protein folding, stability and aggregation

Disease	Protein affected	Description
Cystic fibrosis	Oystic fibrosis transmembrane conductance regulator (CFTR)	The ∆Phe508 mutant has wild-type activity, but impaired folding in the endoplasmic reticulum leads to degradation.
α1 Antitrypsin deficiency	α1 Antitrypsin (also known as SERPINA1)	80% of Glu342Lys mutants misfold and are degraded. Pathology is due to aggregation in patients with a reduced degradation rate.
SCAD deficiency	Short-chain acyl-CoA dehydrogenase (SCAD)	Impaired folding of Arg22Trp mutants leads to rapid degradation.
Alzheimer disease	Presenilin, γ-secretase	Mutations cause incorrect cleavage by the γ-secretase protease to produce the amyloid β-peptide; this aggregates into extracellular amyloid plaques.
Parkinson disease	α-Synuclein	Oxidative damage causes misfolding and aggregation. Hereditary forms are linked to deficiency in ubiquitin-mediated degradation.
Huntington disease	Huntingtin	CAG expansions in the Huntingtin gene lead to an abundance of polyglutamine fragments that aggregate and associate non-specifically with other cellular proteins.
Sickle cell anaemia	Haemoglobin	The Glu6Val mutation leads to aggregation in red blood cells.

Nature Reviews Genetics 2005;6:678-687





protein tangles that grow and connections between nerve cells shrivel



Precision medicine

• Precision medicine (PM) is a medical model that proposes the customization of healthcare, with medical decisions, treatments, practices, or products being tailored to the individual patient, instead of a one-drug-fits-all model.

• Findings from basic, clinical, and social/behavioral research, data from digital health, 'omic technologies, imaging, and computational health sciences and ethical and legal guidelines are integrated into a knowledge network, which informs both science and care for individuals and populations.





Protein Structure - Overview



Secondary structure

How primary and tertiary differs

 Schematic representation of the sequence alignment (A) versus structural alignment (B) of chain A versus chain D from PDB ID 1vr4. The two chains are 100% identical in sequence. The aligned parts are colored green (chain A) and cyan (chain D), while the unaligned parts are colored orange and magenta, respectively.



Kosloff, M., & Kolodny, R. (2008). Sequence-similar, structure-dissimilar protein pairs in the PDB. Proteins: Structure, Function, and Bioinformatics, 71(2), 891-902.

How we use databases and online tools in building analytical pipelines

Nucleic Acids Research

Issues Section browse
Advance articles Submit
Purchase

All Nucleic Acids Re



Volume 49, Issue D1 8 January 2021

Article Contents

Abstract

NEW AND UPDATED DATABASES

NAR ONLINE MOLECULAR
BIOLOGY DATABASE
COLLECTION

ACKNOWLEDGEMENTS

FUNDING

REFERENCES

Comments (0)

Next >

The 2021 Nucleic Acids Research database issue and the online molecular biology database collection ∂ Daniel J Rigden ☎, Xosé M Fernández

About

Nucleic Acids Research, Volume 49, Issue D1, 8 January 2021, Pages D1–D9, https://doi.org/10.1093/nar/gkaa1216 Published: 23 December 2020

📕 PDF 🛛 💵 Split View 🛛 🌜 Cite 🎤 Permissions 🛛 🔩 Share 🔻

Abstract

The 2021 Nucleic Acids Research database Issue contains 189 papers spanning a wide range of biological fields and investigation. It includes 89 papers reporting on new databases and 90 covering recent changes to resources previously published in the Issue. A further ten are updates on databases most recently published elsewhere. Seven new databases focus on COVID-19 and SARS-CoV-2 and many others offer resources for studying the virus. Major returning nucleic acid databases include NONCODE, Rfam and RNAcentral. Protein family and domain databases include COG, Pfam, SMART and Panther. Protein structures are covered by RCSB PDB and dispersed proteins by PED and MobiDB. In metabolism and signalling, STRING, KEGG and WikiPathways are featured, along with returning KLIFS and new DKK and KinaseMD, all focused on kinases. IMG/M and IMG/VR update in the microbial and viral genome resources section, while human and model organism genomics resources include Flybase, Ensembl and UCSC Genome Browser. Cancer studies are covered by updates from canSAR and PINA, as well as newcomers CNCdatabase and Oncovar for cancer drivers. Plant comparative genomics is catered for by updates from Gramene and GreenPhylDB. The entire Database Issue is freely available online on the Nucleic Acids Research website (https://academic.oup.com/nar). The NAR

The NAR online Molecular Biology Database Collection has been substantially updated, revisiting nearly 1000 entries, adding 90 new resources and eliminating 86 obsolete databases, bringing the current total to 1641 databases. It is available at

https://www.oxfordjournals.org/nar/database/c/.

OXFORD Journals

You are here: NAR Journal Home » Database Summary Paper Categories

NAR Database Summary Paper Category List

Nucleotide Sequence Databases RNA sequence databases Protein sequence databases Structure Databases Genomics Databases (non-vertebrate) Metabolic and Signaling Pathways Human and other Vertebrate Genomes Human Genes and Diseases Microarray Data and other Gene Expression Databases Proteomics Resources Other Molecular Biology Databases Organelle databases Plant databases Immunological databases Cell biology

Compilation Paper
Category List
Alphabetical List
Category/Paper List
Search Summary Papers

Compilation Paper
Category List
Alphabetical List
Category/Paper List
Search Summary Papers

Oxford University Press is not responsible for the content of external internet sites

Why to archive research data

Accessibility

- One-stop shop
- Uniform data representation/annotation

Persistence

- Typical websites have a half-life of 2 years
- Professional management

Context

- Comparisons against all other entries
- Validation
- Integration with other resources

Facilitate further analysis

• Database-wide studies and data-mining

The archive cost is ~1% of the cost of generating the data Immeasurable scientific value

Ensembl

Ensembl is a genome browser for vertebrate genomes that supports research in comparative genomics, evolution, sequence variation and transcriptional regulation.

Components:

- genome browser
- tools
- databases
- APIs

Ensembl to provides a centralized resource for geneticists, molecular biologists and other researchers studying the genomes

CEnsembl ₪	AST/BLAT VEP Tools BioMar	t Downloads Help & Docs Blog
Human (GRCh38.p13	3) ▼	
Location: 13:32,315,086-32,400,266	Gene: BRCA2	
Gene-based displays	Gene: BRCA2 ENSG000	00139618
Splice variants Transcript comparison Gene alleles	Description	BRCA2 DNA repair associated [Source:HGNC Symbol;Acc:HGNC:1101@]
E Sequence	Gene Synonyms Location	BRCC2, FACD, FAD1, FAD1, FANCD1, FANCD1, XRCC11 <u>Chromosome 13: 32,315,086-32,400,266</u> forward strand.
Comparative Genomics Genomic alignments Gene tree	About this gene	GRCh38:CM000675.2 This gene has 11 transcripts (<u>splice variants), 242 orthologues,</u> is a member of <u>1 Ensembl protein family</u> and is associated with <u>172 phenotypes</u> .
 Gene gain/loss tree Orthologues Paralogues 	Transcripts	Show transcript table
Ensembl protein families	Summary @	
- GO: Cellular component	Name	BRCA2 @ (HGNC Symbol)
- Phenotypes	CCDS UniProtKB	This gene is a member of the Human CCDS set: CCDS9344.1 @
Genetic Variation	RefSeg	This Ensembl/Gencode gene does not contain any transcripts for which we have selected identical model(s) in RefSeg. If there are other RefSeg transcripts available they will be in the External references table
 Variant image Structural variants 	LRG	LRG_293 provides a stable genomic reference framework for describing sequence variants for this gene
 Gene expression Pathway 	Ensembl version	ENSG0000139618.15
 Regulation External references 	Other assemblies	This gene maps to <u>32,889,223-32,974,403</u> in GRCh37 coordinates. View this locus in the GRCh37 archive: ENSG00000139618@
Supporting evidence ID History	Gene type	Protein coding
└─ Gene history	Annotation method	Annotation for this gene includes both automatic annotation from Ensembl and Havana manual curation, see article.
🔅 Configure this page	Annotation Attributes	overlapping locus [Definitions]
2 Custom tracks	Go to Region in I	Detail for more tracks and navigation options (e.g. zooming)
🛃 Export data		
Share this page	\$ 24 < ⊞ 🛛 🗞 🗟	
♣ Bookmark this page	Genes (Comprehensive set	32.31Mb 32.32Mb 32.35Mb 32.35Mb <t< th=""></t<>
	2	

The Universal Protein Resource (UniProt) is a comprehensive resource for protein sequence and annotation data. The UniProt databases are the UniProt Knowledgebase (UniProtKB), the UniProt Reference Clusters (UniRef), and the UniProt Archive (UniParc).

UniPro Advanced - Q Search BLAST Align Retrieve/ID mapping Peptide search SPARQL Help Contact 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. New UniProt portal for the latest SARS-CoV-2 coronavirus UniParc UniRef Proteomes UniProtKB protein entries and receptors, updated independent of the general UniProt release cycle. UniProt Knowledgebase The UniProt Reference Clusters UniParc is a comprehensive and non-A proteome is the set of proteins view SARS-CoV-2 Proteins and Receptors (UniRef) provide clustered sets of redundant database that contains hought to be expressed by an Swiss-Prot (562,755) sequences from the UniProt most of the publicly available protein organism. UniProt provides proteomes Knowledgebase (including isoforms) Manually annotated and sequences in the world. for species with completely sequenced and selected UniParc records. genomes. reviewed. News Records with Information extracted from literature and curator-Forthcoming changes evaluated computational analysis. Planned changes for UniProt Supporting data TrEMBL (184,998,855) UniProt release 2020 03 Automatically annotated and not Mitochondrial call for help | Cross-referencences to IDEAL and BioGRID-ORCS | reviewed. Changes to ABCD, BioGrid and MycoCLAP cross-references | Un... Records that await full manual Taxonomy Subcellular locations Literature citations annotation. 44 Ð UniProt release 2020_02 Genome integrity maintenance by HMCES | Change of annotation topic Cross-ref. databases Diseases Keywords 'Interaction' | Cross-references to Antibodypedia, MetOSite and PHI-ba... 豊い EA XXX News archive Getting started UniProt data Protein spotlight You Tube Q Text search Integrity ▲ Download latest release

Our basic text search allows you to search all the resources available

Get the UniProt data

UniProtKB - P46	736 (BRCC	3_HUI	MAN)							🏦 Basket 👻			
Display	SBLAST ≡ Align	Format 🔒 Add t	o basket 🔇 Hist	itory				Help video	Add a publication	📢 Feedback			
Entry	Protein Lys-63-s	specific deub	quitinase BR	RCC36									
Publications	Gene BRCC3												
Feature viewer	Organism Homo sap	iens (Human)											
Feature table	Status 🛃 Review	wed - Annotatio	n score: 🔍	• Experimental evidence at properties of the	otein level ⁱ								
None	Function												
Function													
Vames & Taxonomy	Metalloprotease that sp Component of the BRCA	ecifically cleave: \1-A complex, a	s 'Lys-63'-linked complex that sp	l polyubiquitin chains (PubMed:19 pecifically recognizes 'Lys-63'-linl	9214193, PubMed:20656690, ked ubiquitinated histones H2	, PubMed:24075985, PubMed:263 2A and H2AX at DNA lesions sites	344097). Does not have act , leading to target the BRC/	tivity toward 'Ly A1-BARD1 hete	/s-48'-linked polyubiquiti rodimer to sites of DNA (n chains. damage at			
Subcellular location	double-strand breaks (E	OSBs). In the BR	CA1-A complex	, it specifically removes 'Lys-63'-	linked ubiquitin on histones I	H2A and H2AX, antagonizing the I	RNF8-dependent ubiquitina	tion at double-	strand breaks (DSBs)	-			
Pathology & Biotech	PubMed:26195665). Me	Med:20656690). Catalytic subunit of the BRISC complex, a multiprotein complex that specifically cleaves 'Lys-63'-linked ubiquitin in various substrates (PubMed:20656690, PubMed:24075985, PubMed:26344097, Med:26195665). Mediates the specific 'Lys-63'-specific deubiquitination associated with the COP9 signalosome complex (CSN), via the interaction of the BRISC complex with the CSN complex (PubMed:19214193). The BRISC											
PTM / Processing	complex is required for deubiguitination of the i	plex is required for normal mitotic spindle assembly and microtubule attachment to kinetochores via its role in deubiquitination NUMA1 (PubMed:26195665). Plays a role in interferon signaling via its role in the biquitination of the interferon receptor IFNAR1; deubiquitination increases IFNAR1 activity by enhancing its stability and cell surface expression (PubMed:24075985, PubMed:26344097). Down-regulates the response to											
Expression	bacterial lipopolysaccha	ride (LPS) via it	s role in IFNAR1	1 deubiquitination (PubMed:2407	5985). 🗣 12 Publications 👻	,							
Interaction	Cofactor ⁱ												
Structure	Zn ²⁺ By similarity - <u>Note</u> : Binds 1 zinc ion p	Curated Oer subunit.	By similarity 👻										
Family & Domains	Sites												
Sequences (5+)	Feature key	Position(s)	Description					Actions Gra	phical view	Length			
Similar proteins	Metal binding ⁱ	122	Zinc; catalytic (PROSITE-ProRule annotation -						1			
Cross-references	Metal binding ¹	124	Zinc; catalytic	PROSITE-ProRule annotation -						1			
Entry information	Metal binding ¹	135	Zinc; catalytic	PROSITE-ProRule annotation						1			
Miscellaneous	GO - Molecular functio	n ⁱ											
⊾Тор	 enzyme regulator a 	activity 🔮 Sour	ce: MGI 👻										
	 Lys63-specific deub metal ion binding 	oiquitinase activi	ty 🗣 Source: R	Reactome									

÷



Domains and Repeats

entre entre respected					
Feature key	Position(s)	Description	Actions	Graphical view	Length
) omain ⁱ	12 - 179	MPN 🛷 PROSITE-ProRule annotation 👻	BLAST		168

D	isplay	PROSITE ⁱ View protein in PROSITE PS50249 MPN, 1 hit	
E	ntry	Sequences (5+) ⁱ	
P	ublications	Sequence status ⁱ : Complete.	
Fe	eature viewer	Sequence processing ¹ : The displayed sequence is further processed into a mature form.	
Fe	eature table	This entry describes 5 isoforms ¹ produced by alternative splicing. Align Add to basket	
		This entry has 5 described isoforms and 6 potential isoforms that are computationally mapped. Show all 🔤 Align All	
_	None		
~	Function		Length: 316
	Names & Taxonomy	This isoform bas been chosen as the canonical sequence. All positional information in this entry refers to it. This is also the sequence that appears in the	Mass (Da): 36,072
	Subcellular location	downloadable versions of the entry.	Last modified: May 10, 2002 - v2 Checksum: ¹ 5720358C1A2F7421
•	Pathology & Biotech	< Hide	BLAST V GO
	PTM / Processing		
	Expression	10 20 30 40 50 MAVQVVQAVQ AVHLESDAFL VCLNHALSTE KEEVMGLCIG ELNDDTRSDS	
	Interaction	60 70 80 90 100	
	Structure	KFAYTGTEMR TVAEKVDAVR IVHIHSVIIL RRSDKRKDRV EISPEQLSAA 110 120 130 140 150	
	Family & Domains	STEAERLAEL TGRPMRVVGW YHSHPHITVW PSHVDVRTQA MYQMMDQGFV	
	Sequences (5+)	160 170 180 190 200 GLIFSCFIED KNTKTGRVLY TCFQSIQAQK SSESLHGPRD FWSSSQHISI	
	Similar proteins	210 220 230 240 250	
~	Cross-references	EGQKEEEKYE KIEIPIHIVP HVIIGKVCLE SAVELPKILC QEEQDAYKRI 260 270 280 290 300	
	Entry information	HSLTHLDSVT KIHNGSVFTK NLCSQMSAVS GPLLQWLEDR LEQNQQHLQE	
		310	
2	Miscellaneous	LQQEKEELMQ ELSSLE	
•	Тор		

PDB – PROTEIN DATA BANK

PDB – Protein Data Bank

- Archives atomistic experimental models (X-ray, NMR, EM) plus supporting experimental data since 1971
- > 134,000 entries (Oct-2017)
- EMBL-EBI involved since 1999
- Managed by wwPDB partners since 2003
- Collaboration on all aspects of the PDB archive
- Policies, procedures, formats, validation standards, ligands, journals, etc.
- Friendly competition on "data-out"
 - Serving PDB data with added-value; PDBbased services; other services, resources and activities





Other experimental 3D model archives

- NDB Nucleic acid Database
 - Operated by RCSB (1995)
 - >9100 structures (Oct-17)
 - Most structures also in PDB
 - http://ndbserver.rutgers.edu/
- CSD Cambridge Structural Database
 - Operated by CCDC (1965)
 - Crystal structures of "small molecules"
 - >875,000 structures (Oct-17)
 - http://www.ccdc.cam.ac.uk/



PDB – PROTEIN DATA BANK



and interoperability of these motocoles in maconnolecular structures, a contrology



https://www.rcsb.org/

sequence patterns in PDB structures

https://www.ebi.ac.uk/pdbe/node/1



pdbe.org/2yi7

- Some FAQs:
- Where is the structure published?
- Is it any good?
- How do I find out who has cited the structure?
- What is the assembly?
- What sequence / structural domains are present?
- Where is the ligand binding?

Primary publication:

Co-crystalization and in vitro biological characterization of 5-aryl-4-(5-substituted-2-4-dihydroxyphenyl)-1,2,3-thiadiazole hsp90 inhibitors. Sharp SY, Roe SM, Kazlauskas E, Cikotienė I, Workman P, Matulis D, Prodromou C

Details

Details

A PLoS ONE 7 e44642 (2012)

PMID: 22984537 🗹

Function and Biology

Biochemical function: Biological process: • response to stress 🗹

• ATP binding

1 distinct polypeptide molecule

Cellular component: o not assigned

Sequence domains:

Structure analysis

Entry contents:

Macromolecule:

- Heat shock protein Hsp90, N-terminal
- Histidine kinase-like ATPase, C-terminal domain 🗹
- Heat shock protein Hsp90 family 🗹
- Heat shock protein Hsp90, conserved site 🗹

Assembly composition: homo dimer (preferred)

2 bound ligands: Mg ⁺²

Ligands and Environments

Experiments and Validation

Metri	ic P	ercentile Ranks	Value
Rfree	e	0	0.191
Clashscore	e 199		9
Ramachandran outlier	s		0
Sidechain outlier	s		2.3%
RSRZ outlier	s 💶 🛛 🕻		4.8%
X-ray source:	Precentile relative to all X-ray etc DIAMOND BEAN	natures break of einstate resolution MLINE IO2	211
Spacegroup:	1222		
Unit cell:	a: 65.23Å a: 90°	b: 88.705Å β: 90°	c: 99.625Å γ: 90°
R-values:	R 0.176	R _{work} 0.175	R free 0.196

Expression system: Escherichia coli



2 citation in other articles

Volume of Hsp90 ligand binding and the unfolding phase diagram as a function of pressure and temperature. Petrauskas et al. (2013)

I more

PDB REDO

Details

The sliders below show the change in model quality between original PDB entry and the PDB_REDO entry

Model Geometry Fit model/data

PDB_REDO

Heat shock protein HSP 90-alpha

Molecule details : Chain: A Length: 229 amino acids Theoretical weight: 25.86 KDa Source organism: Homo sapiens Expression system: Escherichia coli UniProt: • P07900 @ (Residues: 1-229; Coverage: 31%) Gene names: HSP90A, HSP90AA1, HSPCA, HSPC1 Sequence domains: o Hsp90 protein 🗹 • Histidine kinase-, DNA gyrase B-, and HSP90-like ATPase

Structure domains: Heat Shock Protein 90

$1 \times BZ8$ $1 \times MG$ No modified residues













Validation issues detailed on sequence



UniProt coverage viewer



Protein Folds and Families, CATH, SCOP and Pfam

Proteins are generally composed of one or more functional regions, commonly termed **domains**. Different combinations of domains give rise to the diverse range of proteins found in nature. The identification of domains that occur within proteins can therefore provide insights into their function.



SCOP/CATH -Classification and evolution of domain structure superfamilies





cath



~450,000 domain structures

~6000 superfamilies

C A T H Home Search Browse Download About Support

Search CATH by keywords or ID



CATH v4.3 is nearly here! This has taken longer than anticipated due to issues relating to the recent lockdown, however we are still working hard to generate all the associated data for this upcoming release. In the meantime, get access to the very latest classification information in our daily updates. The core classification files for CATH v4.2 are available to download.

G	3D Structure Find out what 3D structure your protein adopts	Protein Evolution	Protein Function Investigate the function of your protein
	Find out more Go	Find out more	Find out more Go
	Conserved Sites Look at protein sites that are highly conserved and implicated in function	Download Data Download data files and query CATH via webservices	Pind out how CATH is created and maintained, how to link to CATH and more
ļ	Find out more Go	> Go	> Go

CATH: Protein Structure Classification



Cath architectures





2-layer ($\alpha\beta$) sandwich



3-layer ($\beta \alpha \beta$) sandwich



 $\alpha\beta$ -box



3-layer ($\alpha\beta\alpha$) sandwich



4-layer ($\alpha\beta\beta\alpha$) sandwich



 $\alpha\beta$ -horseshoe

Domain Structure Classification

Structure based methods for recognising fold similarities



Pfam database

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

Pfam also generates higher-level groupings of related entries, known as **clans**. A clan is a collection of Pfam entries which are related by similarity of sequence, structure or profile-HMM.

The data presented for each entry is based on the <u>UniProt Reference Proteomes</u> but information on individual UniProtKB sequences can still be found by entering the protein accession. Pfam *full* alignments are available from searching a variety of databases, either to provide different accessions (e.g. all UniProt and NCBI GI) or different levels of redundancy.

Pfam classification



Sequence based methods for detecting protein homologues





sequence-profile comparison PSI-BLAST, HMMer, SAM-T

AN	1 1	C	N	G	Т	L	Е	S	I	R	S	Б	Т	IJ	S	Ŋ	F	A I	J	E.	C	N	G	т	Е	Е	S	I	R	S	L	Т	Ы	S		N
GN	I X	L	Ŋ	G	S	ĸ	Е	М	V	V	D	L	L	IJ	G	LŊ	C	3 1	J		L	N	G	S	R	Е	м	v	v	D	L	L	IJ	G	L	Ŋ
AN		С	N	G	S	Q	Q	S	L	S	Е	Ь	1	D	L	N	Z	11	J		С	N	G	S	0	0	s	E.	s	Е	L	I	D	L		N
GN	P X	Ŀ	N	G	S	R	Q	s	Ι	ĸ	Е	Ι	V	Е	R	L N	C	3 1	N		L	N	G	S	8	Q	S	I	R	Е	Ι	v	Е	R	L	N
GN		И	N	G	D	8	#L	S	D	G	Е	ь	I	Ħ	Т	L N	C	3 1	J		М	N	G	D		£	5	E,	G	Е	Б	I	Æ	т	Ь	N
GN		Ь	Ŋ	G	S	R	H	М	V	Ħ	Е	L	W	S	N	LE	C	3 1	J		L	N	G	S	R		М	V	H	Е	Ь	W	s	N	L	в
GN		M	N	G	R	R	Q	S	L	G	Е	Ь	I	G	Т	L N	C	3 1	J.		И	N	G	12		0	S	6	G	Е	L	I	G	Т	L	N
GN		И	H	R	т	L	A	E	A	V	Q		V	Е	D	VK	C	3 1	IJ		И	H	R.	т	L	A	Е	A	V	0		v	Е	D	V	R
AN	香	C	Ŋ	G	Т	Т	A	S	I	Е	R	L	W	Q	V	Ŋ	7	1	J		С	N	G	т	Т	A	S	I	Е	R	Ь	v	Q	V		N

profile-profile comparison, HMM-HMM comparison PRC, COMPASS, HHpred

sequence-sequence comparison

BLAST, NW, SSEARCH
Blast

BLAST [®] » blas	tp suite		Home Recent Results Saved Strategies Help									
0	COVID-1 Get the latest public he Get the latest re Find NCBI SARS-CoV-2 literature, so	19 is an emerging, rapidly evolving situation. ealth information from CDC: <u>https://www.coronavirus.gov.</u> search from NIH: <u>https://www.nih.gov/coronavirus</u> . equence, and clinical content: <u>https://www.ncbi.nlm.nih.gov/sars-cov-2/</u> .										
Standard Protein BLAST												
blastn blastp blastx	tblastn tblastx											
Estas Ourse Os	BLASTP programs s	search protein databases using a protein query. more	Reset page Bookmark									
Or, upload file Job Title	Imber(s), gi(s), or FASTA sequence(s) Imber(s), gi(s), or FASTA sequence(s) From From To To Enter a descriptive title for your BLAST search re sequences		BLAST results will be displayed in a new format by default You can always switch back to the Traditional Results page.									
Choose Search	Set											
Database Organism Optional Exclude Optional	Non-redundant protein sequences (nr) Image: Completion of the suggested in the suggested in the suggested in the suggested in the suggested is the suggest of the suggest											
Program Select Algorithm	 ○ Quick BLASTP (Accelerated protein-protein BLAST) ● 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 	<u>https://blast.ncbi.nlm.nih.gov/Blast.cg</u> <u>E=BlastHome</u>	gi?CMD=Web&PAGE_TYP									
BLAST	Search database nr using Blastp (protein-protein BLAST)											

HMMER: biosequence analysis using profile hidden Markov models

HMMER is used for searching sequence databases for sequence homologs, and for making sequence alignments. It implements methods using probabilistic models called profile hidden Markov models (profile HMMs).

HMMER is often used together with a profile database, such as <u>Pfam</u>. But HMMER can also work with query *sequences*, not just profiles, just like BLAST. For example, you can search a protein query sequence against a database with **phmmer**, or do an iterative search with **jackhmmer**.

HMMER is designed to detect remote homologs as sensitively as possible, relying on the strength of its underlying probability models. In the past, this strength came at significant computational expense, but as of the new HMMER3 project, HMMER is now essentially as fast as BLAST.



VALIDATION OF STRUCTURAL SIMILARITY

C A T H

Based on CATH (Protein Structure Classification database) to Validate homologues at least 2 out of 3 of the following conditions should be met:

Significant structural similarity Significant sequence similarity

Functional similarity

Structure alignment

Structure alignment may be defined as identification of residues occupying "equivalent" geometrical positions

Unlike in sequence alignment, residue type is neglected

- Used for
 - measuring the structural similarity
 - protein classification and functional analysis
 - database searches

Multiple structure alignment





PDBeFold is an interactive service that allows you to **identify structures** that are **similar** to that of your reference protein. It is a very powerful **structure alignment** tool which can perform both pairwise and multiple three-dimensional alignment. In addition, PDBeFold gives you various options to sort the results of your structural alignment query.

	OSubmission Fo						
	3D alignment	Om	ultiple				
G	Query		Tarç	get			
Source: PDB entr PDB code 1sar Select chains Chains: *(all)	y ▼ view Find chains	Source:	Whole PDB PDB entry SCOP entry Coordinate Whole PDB List of PDB All SCOP 1. SCOP 1.73 List of SCO File set	archive file archive codes .73 archiv subset P 1.73 co	▼ e des		
Lowest acceptable r	natch (%) 70	Lowest ac	cceptable ma	tch (%)	70		
 match individua match connect 	al chains ivity	✓ best r	matches only e matches or	nly			
✓ if no matches within limits of acceptability are found, show close ones							

PDBeFold output

- **Table of matched Secondary Structure Elements**
- **Table of matched backbone** C_{α} **-atoms with distances between them at best structure** superposition
- Rotation-translation matrix of best structure superposition
- Visualisation in Jmol and Rasmol
- \Box *r.m.s.d.* of C_{α}-alignment
- \Box Length of C_{α} -alignment N_{align}
- **\Box** Number of gaps in C_{α}-alignment
- Quality score *Q*
- □ Statistical significance scores *P*(*S*), *Z*
- Sequence identity

The Results Page For Pairwise Alignment

				Dunad	Ν.	A ST	N 9/							
##	Q	Ρ	z	Rmsa	Nalgn	Ng	⁷⁰ seq	%sse	Match	%sse	Nres	×	Title	
1	1.00	26.8	15.5	0.00	184	0	100	100	2mjp:A	100	184		STRUCTURE-BASED IDENTIFICATION OF THE BIOCHEMICAL FUNCTION OF A HYPOTHETICAL PROTEIN FROM METHANOCOCCUS JANNASCHII:MJ0226	
2	1.00	25.4	15.0	0.20	184	0	100	100	1b78:A	100	184		STRUCTURE-BASED IDENTIFICATION OF THE BIOCHEMICAL FUNCTION OF A HYPOTHETICAL PROTEIN FROM METHANOCOCCUS JANNASCHII:MJ0226	
3	0.95	22.3	14.1	0.49	182	1	99	93	1b78:B	93	184		STRUCTURE-BASED IDENTIFICATION OF THE BIOCHEMICAL FUNCTION OF A HYPOTHETICAL PROTEIN FROM METHANOCOCCUS JANNASCHII:MJ0226	
4	0.95	20.0	13.3	0.55	182	1	99	86	2mjp:B	100	184 STRUCTURE-BASED IDENTIFICATION OF THE BIOCHEMICAL FUNCTION OF A HYPOTHETICAL PROTEIN FROM METHANOCOCCUS JANNASCHILMJ0226		STRUCTURE-BASED IDENTIFICATION OF THE BIOCHEMICAL FUNCTION OF A HYPOTHETICAL PROTEIN FROM METHANOCOCCUS JANNASCHILMJ0226	
5	0.85	15.3	11.6	1.14	182	2	50	93	2dvn:B	100	186	6 D STRUCTURE OF PH1917 PROTEIN WITH THE COMPLEX OF IMP FROM PYROCOCC		
<u>6</u>	0.68	9.1	8.9	1.92	181	3	49	93	1v7r:A	100	186		STRUCTURE OF NUCLEOTIDE TRIPHOSPHATE PYROPHOSPHATASE FROM PYROCOCCUS HORIKOSHII 0T3	
<u>7</u>	0.67	8.5	8.6	1.96	181	3	49	93	2dvo:A	100	185		STRUCTURE OF PH1917 PROTEIN WITH THE COMPLEX OF ITP FROM PYROCOCCUS HORIKOSHII	
8	0.67	9.0	8.9	1.95	181	3	49	93	2dvn:A	100	186		STRUCTURE OF PH1917 PROTEIN WITH THE COMPLEX OF IMP FROM PYROCOCCUS HORIKOSHII	
9	0.66	6.7	7.7	1.93	177	5	50	79	2dvp:A	92	184		STRUCTURE OF NTPASE FROM PYROCCOUS HORIKOSHII	
<u>10</u>	0.65	7.9	8.3	2.08	181	3	49	93	2e5x:A	100	185		STRUCTURE OF NUCLEOTIDE TRIPHOSPHATE PYROPHOSPHATASE FROM PYROCOCCUS HORIKOSHII OT3	
<u>11</u>	0.65	8.7	8.8	1.42	168	7	35	71	2car:B	83	194		CRYSTAL STRUCTURE OF HUMAN INOSINE TRIPHOSPHATASE	
<u>12</u>	0.64	8.6	8.8	1.48	170	7	35	71	2car:A	83	196		CRYSTAL STRUCTURE OF HUMAN INOSINE TRIPHOSPHATASE	
<u>13</u>	0.64	8.8	8.8	1.41	168	7	35	71	2i5d:A	83	195		CRYSTAL STRUCTURE OF HUMAN INOSINE TRIPHOSPHATE PYROPHOSPHATASE	
<u>14</u>	0.63	5.4	7.0	1.81	173	6	35	71	1vp2:A	91	189		CRYSTAL STRUCTURE OF PUTATIVE XANTHOSINE TRIPHOSPHATE PYROPHOSPHATASE/HAM1 PROTEIN HOMOLOG (TM0159) FROM THERMOTOGA MARITIN AT 1.78 A RESOLUTION	

Analyzing the result from a particular pairwise alignment



Analyzing the result from a particular pairwise alignment



Residue-by-Residue Structure alignment result

				notations					
PI	DB 2mj	p:A	SI	Dist. (Å)	PDB 2e5x:A				
					E	A:MET	1		
+	A:LYS	10	10×	2.20	s+	A:LYS	2		
s-	A:ILE	11		2.05	s-	A:ILE	3		
s٠	A:TYR	12	122	1.69	s-	A:PHE	4		
s-	A:PHE	13	225	1.35	s-	A:PHE	5		
s-	A:ALA	14	1	1.19	s-	A:ILE	6		
•	A:THR	15	225	1.83	•	A:THR	7		
-	A:GLY	16		2.95	•	A:SER	8		
+	A:ASN	17	:::·	3.82	H+	A:ASN	9		
H+	A:PRO	18	115	4.87	H+	A:PRO	10		
H+	A:ASN	19		5.09	H-	A:GLY	11		
H+	A:LYS	20	10×	3.70	H+	A:LYS	12		
H-	A:ILE	21	:::	3.77	H-	A:VAL	13		
H+	A:LYS	22	::	4.19	H+	A:ARG	14		
H+	A:GLU	23	121	3.13	H+	A:GLU	15		
H-	A:ALA	24	1	2.22	H-	A:VAL	16		
H+	A:ASN	25		2.84	H-	A:ALA	17		
н-	A:ILE	26		2.55	H+	A:ASN	18		

3D Structural alignment

- Structure alignment includes both secondary structure assignment and residue identity:
- For aligned residues:
 - Exact match
 - No match

Multiple 3D alignment using PDBefold

	Submission Form multiple
	3D alignment Opairwise
List of entries	PDB::3d5i
PDB::1sar:A PDB::1py3:C PDB::3d5i:A	Source: PDB entry ~ PDB code 3d5i view
	Select chains ~ Find chains
~	Chains: A
V	"iewer: Jmol ∨ Delete entry Update List New entry
	Home Submit your query

Results from multiple 3D alignment

3D Structural alignment



$PolyPhen-2 \ \ {\rm prediction \ of \ functional \ effects \ of \ human \ SNPs}$

PolyPhen-2 (**Poly**morphism **Phen**otyping v**2**) is a tool which predicts possible impact of an amino acid substitution on the structure and function of a human protein using straightforward physical and comparative considerations.

PolyPhen-2 prediction of functional effects of human nsSNPs	
Home About Help Downloads Batch query WHES	S.db
PolyPhen-2 (Polymorphism Phenotyping v2) is a tool which predicts possible impact of an amino acid substitution on the stru	cture and function of a human protein using straightforward physical and comparative considerations. Please, use the form below to submit your query.
Query Data	
Protein or SNP identifier	
Protein sequence in FASTA format	
Position	
Substitution	AA ₁ A R N D C E Q G H I L K M F P S T W Y V AA ₂ A R N D C E Q G H I L K M F P S T W Y V
Query description	
	Submit Query Clear Check Status
	Display advanced query options
Coffware 9 was auspart was adaption	Web design & development: h



3D structure prediction tools

Tertiary (or 3-D) structure prediction tools fall into two main methods: Ab initio, and comparative protein modeling.

Ab initio (or de novo) protein structure prediction methods attempt to predict tertiary structures from sequences based on general principles that govern protein folding energetics and/or statistical tendencies of conformational features that native structures acquire, without the use of explicit templates. Ab initio protein structure prediction thus requires vast amount of computational power and time to solve the native conformation of a protein, and remains one of the top challenges for modern science.

Most popular servers include <u>Robetta</u> (using the Rosetta software package), <u>SWISS-MODEL</u>, <u>PEPstr</u>, <u>QUARK</u>.

If a protein of known tertiary structure shares at least 30% of its sequence with a potential homolog of undetermined structure, comparative methods that overlay the putative unknown structure with the known can be utilized to predict the likely structure of the unknown. Homology modeling and protein threading are two main strategies that use prior information on other similar protein to propose a prediction of an unknown protein, based on its sequence.

Homology modeling and protein threading software include <u>RaptorX</u>, <u>FoldX</u>, <u>HHpred</u>, <u>I-TASSER</u>, and more.

$ALPHA\ FOLD\ (https://deepmind.com/blog/article/AlphaFold-Using-AI-for-scientific-discovery)$



AlphaFold placed first in the overall rankings of the 13th CASP In November 2020, DeepMind's new version, AlphaFold 2, won CASP14



Improved protein structure prediction using potentials from deep learning. <u>https://www.nature.com/articles/s41586-019-1923-</u> <u>7.epdf?author_access_token=Z_KaZKDqtKzbE7Wd5Htwl9RgN0jAjWel9jnR3ZoTv0MCcgAwHMgRx9mvLjNQdB2TlQQaa7l420UCtGo8vYQ39gg8IFWR9mAZtvsN_1Pr</u> ccXflbc6e-tGSgazNL_XdtQzn1PHfy21qdcxV7Pw-k3htw%3D%3D

https://github.com/deepmind/deepmind-research/tree/master/alphafold_casp13



(The server completed predictions for <u>555532 proteins</u> submitted by <u>132543 users</u> from <u>148 countries or regions</u>) (<u>The template library</u> was updated on <u>2020/07/11</u>)

I-TASSER (Iterative Threading ASSEmbly Refinement) is a hierarchical approach to protein structure and function prediction. It first identifies structural templates from the PDB by multiple threading approach LOMETS, with full-length atomic models constructed by iterative template-based fragment assembly simulations. Function insights of the target are then derived by re-threading the 3D models through protein function database <u>BioLiP</u>. I-TASSER (as 'Zhang-Server') was ranked as the No 1 server for protein structure prediction in recent community-wide <u>CASP7</u>, <u>CASP8</u>, <u>CASP9</u>, <u>CASP10</u>, <u>CASP11</u>, <u>CASP12</u>, and <u>CASP13</u> experiments. It was also ranked as the best for function prediction in <u>CASP9</u>. The server is in active development with the goal to provide the most accurate protein structure and function predictions using state-of-the-art algorithms. Please report problems and questions at <u>I-TASSER message board</u> and our developers will study and answer the questions accordingly. (>> More about the server ...)

Structure models for the 2019-nCov Coronavirus genome by C-I-TASSER

[Queue] [Forum] [Download] [Search] [Registration] [Statistics] [Remove] [Potential] [Decoys] [News] [Annotation] [About] [FAQ]

I-TASSER On-line Server (View an example of I-TASSER output):

Copy and paste your sequence below ([10, 1500] residues in FASTA format). Click here for a sample input:

Or upload the sequence from your local computer: Choose File No file chosen

Email: (mandatory, where results will be sent to)

Password: (mandatory, please click <u>here</u> if you do not have a password)

ID: (optional, your given name of the protein)

Option I: Assign additional restraints & templates to guide I-TASSER modeling.

Option II: Exclude some templates from I-TASSER template library.

Option III: Specify secondary structure for specific residues.

I-TASSER results for job id S555535

[Click on S555535_results.tar.bz2 to download the tarball file including all modeling results listed on this page]

(Click on Annotation of I-TASSER Output to read the instructions for how to interpret the results on this page)

Submitted Sequence in FASTA format >seq MAKSSFKISNPLEARMSESSRIREKYPDRIPVIVEKAGQSDVPDIDKKKYLVPADLTVGQ FVYVVRKRIKLGAEKAIFVFVKNTLPPTAALMSAIYEEHKDEDGFLYMTYSGENTFGSLT VA Predicted Secondary Structure 20 100 40 60 80 120 Sequence MAKSSFKISNPLEARMSESSRIREKYPDRIPVIVEKAGQSDVPDIDKKKYLVPADLTVGQFVVVVKRIKLGAEKAIFVFVKNTLPPTAALMSAIYEEHKDEDGFLYMTYSGENTFGSLTVA Conf.Score 986532331899999999999999878765157887636257852245237606997088999997643256988619999898516654759999998348778779996178656885659 H:Helix; S:Strand; C:Coil Predicted Solvent Accessibility 20 40 60 100 120 80 Sequence MAKSSFKISNPLEARMSESSRIREKYPDRIPVIVEKAGQSDVPDIDKKKYLVPADLTVGQFVYVVRKRIKLGAEKAIFVFVKNTLPPTAALMSAIYEEHKDEDGFLYMTYSGENTFGSLTVA Prediction 87645147633155036405412741473030103326756145153432113472323302310443171466310101024322464330340156344421000000132412233648

Values range from θ (buried residue) to 9 (highly exposed residue)

Predicted normalized B-factor

(B-factor is a value to indicate the extent of the inherent thermal mobility of residues/atoms in proteins. In I-TASSER, this value is deduced from threading template proteins from the PDB in combination with the sequence profiles derived from sequence databases. The reported B-factor profile in the figure below corresponds to the normalized B-factor of the target protein, defined by B=(B'-u)/s, where B' is the raw B-factor value, u and s are respectively the mean and standard deviation of the raw B-factors along the sequence. Click here to read more about predicted normalized B-factor)

Predicted normalized B-factor

(B-factor is a value to indicate the extent of the inherent thermal mobility of residues/atoms in proteins. In I-TASSER, this value is deduced from threading template proteins from the PDB in combination with the sequence profiles derived from sequence databases. The reported B-factor profile in the figure below corresponds to the normalized B-factor of the target protein, defined by B=(B'-u)/s, where B' is the raw B-factor value, u and s are respectively the mean and standard deviation of the raw B-factors along the sequence. <u>Click here to read more about predicted normalized B-factor</u>)



Top 10 threading templates used by I-TASSER

(I-TASSER modeling starts from the structure templates identified by LOMETS from the PDB library. LOMETS is a meta-server threading programs, where each threading programs, where each threading programs can generate tens of thousands of template alignments. I-TASSER only uses the templates of the highest significance in the threading alignments, the significance of which are measured by the Z-score, i.e. the difference between the raw and average scores in the unit of standard deviation. The templates in this section are the 10 best templates selected from the LOMETS threading programs. Usually, one template of the highest Z-score is selected from each threading program, where the threading programs are sorted by the average performance in the large-scale benchmark test experiments.)

Rank	PDB Hit	lden1 lden2	2 Cov	Norm. Z-score	Download Align.		20 	40 	60 	80 	100 	120
						Sec.Str CCCCCCCC Seq MAKSSFKISN	CHHHHHHHHHHHHHHCCC IPLEARMSESSRIREKYPD	CCSSSSSSCCCCCCCCCCSS RIPVIVEKAGQSDVPDIDKKKY	SCCCCCCHHHHHHHHHHH LVPADLTVGQFVYVVRKRI	CCCCCCSSSSSSCCCCCCC KLGAEKAIFVFVKNTLPPTA	CCHHHHHHHHCCCCCCC	SSSSSSCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCC
1	<u>1eo6B</u>	0.60 0.57	0.96	3.40	Download	M-KWMFKEDH	ISLEHRCVESAKIRAKYPD	RVPVIVEKVSGSQIVDIDKRKY	(LVPSDITVAQFMWIIRKRI	QLPSEKAIFLFVDKTVPQSS	LTMGQLYEKEKDEDGF	LYVAYSGENTFGF
2	<u>3h9dA</u>	0.53 0.51	0.95	2.81	Download	- KDSKYKMSH	ITFESRQSDAAKVRERHPD	RLPIICEKVYNSDIGELDRCKF	LVPSDLTVGQFVSVLRKRV	QLEAESALFVYTNDTVLPSS	AQMADIYSKYKDEDGF	LYMKYSGEATFG
3	<u>4co7A</u>	0.58 0.57	0.97	3.22	Download	SMKWMFKEDH	ISLEHRCVESAKIRAKYPD	RVPVIVEKVSGSQIVDIDKRKY	(LVPSDITVAQFMWIIRKRI	QLPSEKAIFLFVDKTVPQSS	LTMGQLYEKEKDEDGF	LYVAYSGENTFGF
4	<u>51xiD</u>	0.56 0.53	0.96	3.11	Download	-MKFQYKEDH	IPFEYRKKEGEKIRKKYPD	RVPVIVEKAPKARVPDLDKRKY	/LVPSDLTVGQFYFLIRKRI	HLRPEDALFFFVNNTIPPTS	ATMGQLYEDNHEEDYF	LYVAYSDESVYGK
5	<u>51xiD</u>	0.56 0.53	0.96	2.28	Download	-MKFQYKEDH	IPFEYRKKEGEKIRKKYPD	RVPVIVEKAPKARVPDLDKRKY	/LVPSDLTVGQFYFLIRKRI	HLRPEDALFFFVNNTIPPTS	ATMGQLYEDNHEEDYF	LYVAYSDESVYGK
6	<u>3h9dA</u>	0.53 0.51	0.95	3.88	<u>Download</u>	-KDSKYKMSH	ITFESRQSDAAKVRERHPD	RLPIICEKVYNSDIGELDRCKF	LVPSDLTVGQFVSVLRKRV	QLEAESALFVYTNDTVLPSS	AQMADIYSKYKDEDGF	LYMKYSGEATFG
7	<u>51xiD</u>	0.56 0.53	0.96	3.45	<u>Download</u>	-MKFQYKEDH	IPFEYRKKEGEKIRKKYPD	RVPVIVEKAPKARVPDLDKRKY	/LVPSDLTVGQFYFLIRKRI	HLRPEDALFFFVNNTIPPTS	ATMGQLYEDNHEEDYF	LYVAYSDESVYGK
8	<u>1eo6A</u>	0.60 0.57	0.95	2.98	<u>Download</u>	M-KWMFKEDH	ISLEHRCVESAKIRAKYPD	RVPVIVEKVSGSQIVDIDKRKY	(LVPSDITVAQFMWIIRKRI	QLPSEKAIFLFVDKTVPQSS	LTMGQLYEKEKDEDGF	LYVAYSGENTFG
9	<u>3h9dA</u>	0.53 0.51	0.95	3.24	Download	-KDSKYKMSH	ITFESRQSDAAKVRERHPD	RLPIICEKVYNSDIGELDRCKF	LVPSDLTVGQFVSVLRKRV	QLEAESALFVYTNDTVLPSS	AQMADIYSKYKDEDGF	LYMKYSGEATFG
10	<u>1gnuA</u>	0.00 0.54	0.96	3.21	Download							

(a) All the residues are colored in black; however, those residues in template which are identical to the residue in the query sequence are highlighted in color. Coloring scheme is based on the property of amino acids, where polar are brightly coloured while non-polar residues are colored in dark shade. (more about the colors used)

(b) Rank of templates represents the top ten threading templates used by I-TASSER.

(c) Ident1 is the percentage sequence identity of the templates in the threading aligned region with the query sequence.

(d) Ident2 is the percentage sequence identity of the whole template chains with query sequence.

(e) Cov represents the coverage of the threading alignment and is equal to the number of aligned residues divided by the length of query protein.

(f) Norm. Z-score is the normalized Z-score of the threading alignments. Alignment with a Normalized Z-score >1 mean a good alignment and vice versa.

(g) Download Align. provides the 3D structure of the aligned regions of the threading templates.

(h) The top 10 alignments reported above (in order of their ranking) are from the following threading programs:

1: MUSTER 2: FFAS-3D 3: SPARKS-X 4: HHSEARCH2 5: HHSEARCH1 6: Neff-PPAS 7: HHSEARCH 8: pGenTHREADER 9: wdPPAS 10: cdPPAS

Top 5 final models predicted by I-TASSER

(For each target, I-TASSER simulations generate a large ensemble of structural conformations, called decoys. To select the final models, I-TASSER uses the SPICKER program to cluster all the decoys based on the pair-wise structure similarity, and reports up to five models which corresponds to the five largest structure clusters. The confidence of each model is quantitatively measured by C-score that is calculated based on the significance of threading template alignments and the convergence parameters of the structure assembly simulations. C-score is typically in the range of [-5, 2], where a C-score of a higher value signifies a model with a higher confidence and vice-versa. TM-score and RMSD are estimated based on C-score and protein length following the correlation observed between these qualities. Since the top 5 models are ranked by the cluster size, it is possible that the lower-rank models have a higher C-score in rare cases. Although the first model has a better quality in most cases, it is also possible that the lower-rank models have a better quality than the higher-rank models as seen in our benchmark tests. If the I-TASSER simulations converge, it is possible to have less than 5 clusters generated; this is usually an indication that the models have a good quality because of the converged simulations.)

- More about C-score
- Local structure accuracy profile of the top five models

(By right-click on the images, you can export image file or change the configurations, e.g. modifying the background color or stopping the spin of your models)



- Download Model 1
- C-score=0.51 (Read more about C-score)
- Estimated TM-score = 0.78±0.10
- Estimated RMSD = 3.4±2.3Å



Reset to initial orientation Spin On/Off

- <u>Download Model 2</u>
- C-score = -0.65



Reset to initial orientation Spin On/Off

<u>Download Model 3</u>
C-score = 0.47



Reset to initial orientation Spin On/Off

- <u>Download Model 4</u>
- C-score = -2.03

Proteins structurally close to the target in the PDB (as identified by TM-align)

(After the structure assembly simulation, I-TASSER uses the TM-align structural alignment program to match the first I-TASSER model to all structures in the PDB library. This section reports the top 10 proteins from the PDB that have the closest structural similarity, i.e. the highest <u>TM-score</u>, to the predicted I-TASSER model. Due to the structural similarity, these proteins often have similar function to the target. However, users are encouraged to use the data in the next section 'Predicted function using COACH' to infer the function of the target protein, since COACH has been extensively trained to derive biological functions from multi-source of sequence and structure features which has on average a higher accuracy than the function annotations derived only from the global structure comparison.)





TM-align is an algorithm for sequence independent protein structure comparisons. For two protein structures of unknown equivalence, TM-align first generates optimized residue-to-residue alignment based on structural similarity using heuristic dynamic programming iterations. An optimal superposition of the two structures built on the detected alignment, as well as the <u>TM-score</u> value which scales the structural similarity, will be returned. TM-score has the value in (0,1], where 1 indicates a perfect match between two structures. Following strict statistics of structures in the PDB, scores below 0.2 correspond to randomly chosen unrelated proteins while those higher than 0.5 assume generally the same fold in SCOP/CATH.

News: New TM-align now allows for input structure with either PDB or PDBx/mmCIF format. Meanwhile C++ version of TM-align is now officially released.

TM-align on-line (view an example of output)

Note: This server is only for pair-wise structure comparison. If you want to match one protein structure with all proteins in the PDB library, you can do it in COFACTOR Server.

· Input Structure 1 in PDB format or PDBx/mmCIF format (mandatory):

Please copy and paste your structure file here. Sample input

HEADER	OXYGEN TRANSPORT 13-DEC-97 101M	
TITLE	SPERM WHALE MYOGLOBIN F46V N-BUTYL ISOCYANIDE AT PH 9.0	
COMPND	MOL_ID: 1;	
COMPND	2 MOLECULE: MYOGLOBIN;	
COMPND	3 CHAIN: A;	
COMPND	4 ENGINEERED: YES;	
COMPND	5 MUTATION: YES	
SOURCE	MOL_ID: 1;	
SOURCE	2 ORGANISM_SCIENTIFIC: PHYSETER CATODON;	
SOURCE	3 ORGANISM COMMON: SPERM WHALE;	
1 4 I		- b

Or upload the stucture file:

Choose File No file chosen

Input Structure 2 in <u>PDB format</u> or <u>PDBx/mmCIF format</u> (mandatory):

Please copy and paste your structure file here. Sample input

HEADER	OXYGEN STORAGE 22-FEB-89 1MBA		*
TITLE	APLYSIA LIMACINA MYOGLOBIN. CRYSTALLOGRAPHIC ANALYSIS AT		
TITLE	2 1.6 ANGSTROMS RESOLUTION	- 1	
COMPND	MOL_ID: 1;		
COMPND	2 MOLECULE: MYOGLOBIN;		
COMPND	3 CHAIN: A;		
COMPND	4 ENGINEERED: YES		
SOURCE	MOL_ID: 1;		
SOURCE	2 ORGANISM_SCIENTIFIC: APLYSIA LIMACINA;		
SOURCE	3 ORGANISM COMMON: SLUG SEA HARE;		*
I ≤ 1		F	

Or upload the stucture file:

Choose File No file chosen

• Input Email: (optional, where results will be sent to)

TM-align (Version 20190822) An algorithm for protein structure alignment and comparison Based on statistics: 0.0 < TM-score < 0.30, random structural similarity 0.5 < TM-score < 1.00, in about the same fold Reference: Y Zhang and J Skolnick, Nucl Acids Res 33, 2302-9 (2005) Please email your comments and suggestions to: zhng@umich.edu

Name of Chain_1: A835798 Name of Chain_2: B835798 Length of Chain_1: 154 residues Length of Chain_2: 146 residues

Aligned length= 143, RMSD= 1.83, Seq_ID=n_identical/n_aligned= 0.245 TM-score= 0.81453 (if normalized by length of Chain_1) TM-score= 0.85377 (if normalized by length of Chain_2) (You should use TM-score normalized by length of the reference protein)

(":" denotes aligned residue pairs of d < 5.0 A, "." denotes other aligned residues)

MVLSEGEWQLVLHVWAKVEADVAGHGQDILIRLFKSHPETLEKFDRVKHLKTEAEMKASEDLKKHGVTVLTALGAILKKK--G-HHEAELKPLAQSHATKHKIPIKYLEFISEAIIHVLHSRHPGNFGADAQGAMNKALELFRKDIAAKYKELGYQG -SLSAAEADLAGKSWAPVFANKNANGLDFLVALFEKFPDSANFFADFKGKS-VADIKASPKLRDVSSRIFTRLNEFVNNAANAGKMSAMLSQFAKEHV-GFGVGSAQFENVRSMFPGFVASVAA--PPAGADAANTKLFGLIIDALKA-A-----GA

Visualization (Protein-1 in blue and Protein-2 in red)



•TMScore: TM-score is intended as a more accurate measure of the quality of full-length protein structures than the often used RMSD measure •RMSD — a different structure comparison measure

•GDT — a different structure comparison measure •Longest continuous segment (LCS)

 Global distance calculation (GDC_sc, GDC_all) — Structure comparison measures that use full-model information (not just α-carbon) to assess similarity
 Local global alignment (LGA) — Protein structure alignment program and structure comparison measure

USE CASES

Chronic Lymphocytic Leukemia CLL

- 19,635 IGHV-IGHD-IGHJ gene rearrangement sequences from CLL patients were examined. Based on IGHV gene usage, all cases were subgrouped into 3 phylogenetic clans namely
 - Clan I | IGHV1/5/7), n=5439;
 - Clan II | IGHV2/4/6, n=5092 and
 - Clan III | IGHV3, n=9104.
- A region weight factor was set based on the <u>physicochemical properties</u> of the amino acids and <u>the entropy</u> of the SHMs in the complementarity determining regions (CDRs) and the framework regions (FRs).
- The Euclidean distance between the sequences was defined as a combination of the Blosum Matrix and the region weight factor,
- DBSCAN clustering algorithm
- A force-directed graph drawing visualization the associations between SHMs and CLL groups was underlined using the VH CDR3 patterns as similarity metric.

Distance between sequences:

- Position Entropy
- the Blosum Matrix
- The physicochemical similarity was defined according to the 11 "IMGT Physicochemical" classes, for the 20 amino acids"

```
Ala
Arg -1
        5
Asn
   -2
        0
    -2 -2
    0
       -3
           -3
               -3
            0
                0
        0
            0
                2
            0 -1 -3 -2
       -2
                         -2
Gly
        0
               ^{-1}
                  -3
                       0
                          0
His
            1
       -3 -3 -3 -1 -3 -3
                             -4
       -2
           -3 -4
                  -1 -2 -3
        2
            0
              ^{-1}
                  -3
                      1
                          1
                             -2
           -2 -3
                  ^{-1}
                      0 -2
                            -3
       -3 -3 -3 -2 -3 -3 -3
Pro
    -1 -2 -2 -1 -3 -1 -1 -2 -2 -3
Ser
       -1
                0
                  -1
                       0
                          0
                              0
                                           0
Thr
    0
       ^{-1}
           0 -1 -1 -1 -1 -2 -2 -1 -1 -1
Trp
    -3
       -3 -4 -4 -2 -2 -3 -2 -2 -3
                                      -2
                                          -3
    -2 -2 -2 -3 -2 -1 -2 -3 2 -1 -1 -2 -1
                                                  3
                                                     -3
                  -1 -2 -2 -3 -3 3
Val
              -3
                                       1 -2 1 -1 -2
   Ala Arg Asn Asp Cys Gln Glu Gly His Ile Leu Lys Met Phe Pro Ser Thr Trp Tyr Val
```

Volume' classes			Hydropathy' classes									
	in A ³	H	Iydrop	hobic		Neutral		Hydrophilic				
Very large	189-228	F		W		Y						
Large	162-174	Ι	L	М				K R				
Medium	138-154	v					Н		Е	Q		
Small	108-117			С	Ρ	Т			D	Ν		
Very small	60-90	Α			G	S						
		Al	phatic	Sulfur		Hydroxil	Bas	ic	Acidic	Amide		
					-	Uncharged	Charged			Uncharged		
			Non	polar				P	olar			

- Region weight factor was set based on the <u>physicochemical properties</u> of the amino acids and <u>the</u> <u>entropy</u> of the SHMs in the complementarity determining regions (CDRs) and the framework regions (FRs).
- The SHM variance of each position was calculated through the position Shannon entropy $-\sum_{i=1}^{\frac{m_i}{N}\log_2 \frac{m_i}{N}}$

where m_i / N is the frequency of the SHM *i* in every codon of the IGHV genes sequence.



The points a,b,c underline the maximum entropy in CDRs while d,e,f in FRs for ClanI, Clan II, ClanIII respectively.







	Cluster 1	Cluster 2	Cluster 3	Outliers
Gr. clan I	4496	766	n/a	218
Gr. clan II	3992	824	171	107
Gr. clan III	7671	1092	n/a	342



Group No	Jaccard Index	Subset	Associations	VH CDR3 pattern	Sequence Logo	Group E				1.1242	
Group A- Claul_Clusterl	0.841	#28A	5-M-78_V-A-100	ARXXXGXXYYYYYGMDx	AR-YSGSYYWWYG DVI	Clan2_Cluster1	0.330	#148B	S-T-31_N-D-59	(not officially defined)	
GroupB-	0.317	=]	R-W-55_V-A-100, G-D-36_V-A-100, G-D-36_R-W-55,	ARx[NQ]W[AVL IbcootFDx		Group F- Clan2_Clusterl	0.250	#14	P-A-46_Y-T-58	x[KRH]GGxWxFDx	
Chan_Crasteri		-	M-1-39_R-W-55, R-W-55_S-M-78 R-W-55_V-A-100, G-D-36_V-A-100.			Group G- Clan2_Clusterl	0.240	#14	S-T-40_I-M-78	x[KRH]GGxWxFDx	
Group B- Clan1_Cluster1	0.493	#28A	G-D-36_R-W-55, M-I-39_R-W-55, R-W-55_S-M-79	ARamGarYYYYYGMDx		Group H- Clan2 Clusterl	0.400	#77	P-S-16_C-Y-103	[AVLI]RGxxx[ST]GWxxxxx	
Group C- Clan1_Cluster1	0.384	#59	A-T-38_G-R-55, F-L-62_T-I-64	AzexDFWSGzex							
Group D- Clan2_Cluster1	0.067	#4	S-T-40_P-S-45	[AVLI]RGxxxxxxx[KRH]RYYYYGx[DE]x		Group I- Clan3_Clusterl	0.292	#31	S-E-38_S-G-62	ARDOGGGGGX XXXXYYYXMD x	

results

- Clan based patterns in terms of both quantity and quality of mutations.
- Some frequent SHMs were linked to specific subsets. And are now considered subset specific.
- Next level analysis will include associations between more than two SHMs.
- SHM patterns in the VH region will then be associated with CDR3 patterns leading to the identification of yet uncharacterized group of patients.

Datasets Description

Dataset	Protein structures	Predefined Subsets
D1	137	6 (D1.A ~ D1.F)
D2	925	6 (D1.A ~ D1.F) + Unknown

Test dataset

Ground Truth
Subset 1: 37 cases 28 cases: IGHV1, 8 cases: IGHV5, 1 case: IGHV7
Subset 2: 43 cases
Subset 4: 22 cases
Subset 6: 12 cases
Subset 7: 12 cases 2 cases: 21 aa CDR3, 1 case: 22 aa
CDR3, 4 cases: 23 aa CDR3, 2 cases: 24 aa CDR3, 2
cases: 25 aa CDR3, 1 case: 26 aa CDR3
Subset 8: 11 cases 7 cases: 19 aa CDR3 and 4 cases:
18 aa CDR3

D2

D1

Marcatili, P., Rosi, A., & Tramontano, A. (2008). PIGS: automatic prediction of antibody structures. Bioinformatics, 24(17), 1953-1954.

Dataset Patients Subsets D1 106 6 (D1.A ~ D1.F) D2 894 N/A EXPDIA THEORETICAL MODEL, MODELLER 9.14 2016/01/20 18:57:51 REMARK 6 MODELLER OBJECTIVE FUNCTION: 297.9526 REMARK 6 MODELLER OBJECTIVE FUNCTION: 297.9526 ATOM 1 N VAL H 2 -10.312 3.544 19.698 1.00 0.00 ATOM 2 CA VAL H 2 -11.487 2.542 19.698 1.00 0.00 CA ATOM 3 CB VAL H 2 -11.487 1.658 20.956 1.00 0.00 CA ATOM 5 CG2 VAL H 2 -10.196 0.849 21.127 1.00 0.00 CA ATOM 5 CG2 VAL H 2 -10.196 0.849 21.127 1.00 0.00 CA ATOM 6 C VAL H 2 -12.770 3.081 19.445 1.00 0.00 CA ATOM 8 N GLN H 3 -14.986 3.636 2.044 1.00 0.00 CA
D1 106 6 (D1.A ~ D1.F) D2 894 N/A EXPDTA THEORETICAL MODEL, MODELLER 9.14 2016/01/20 18:57:51 REMARK 6 MODELLER OBJECTIVE FUNCTION: 297.9526 REMARK 6 MODELLER BEST TEMPLATE % SEQ ID: 100.000 1 ATOM 1 N VAL H 2 -11.473 1.658 20.956 1.00 0.00 ATOM 2 CA VAL H 2 -11.473 1.658 20.956 ATOM 3 CB VAL H 2 -11.473 1.658 20.956 1.00 0.00 ATOM 4 CG1 VAL H 2 -12.664 0.741 2.081 1.00 0.00 C0 ATOM 5 CG2 VAL H 2 -12.770 3.081 19.445 1.00 0.00 C0 ATOM 6 VAL H 2 -12.770 3.081 19.445 1.00 0.00 C0 ATOM 8 N GLN H 3 -13.662 3.226
D2 894 N/A EXPDTA THEORETICAL MODELL, MODELLER 9.14 2016/01/20 18:57:51 REMARK 6 MODELLER OBJECTIVE FUNCTION: 297.9526 REMARK 6 MODELLER BEST TEMPLATE % SEQ ID: 100.000 1 ATOM 1 N VAL H 2 -10.312 3.544 19.698 1.00 0.00 N ATOM 1 N VAL 4 2 -11.387 2.542 19.698 1.00 0.00 NC ATOM 2 CA VAL 4 2 -11.473 1.658 20.956 1.00 0.00 CC ATOM 4 CG1 VAL 4 2 -12.684 0.741 20.081 1.00 0.00 CC ATOM 5 CG2 VAL H 2 -12.770 3.081 19.445 1.00 0.00 CC ATOM 6 VAL H 2 -12.770 3.081 19.445 1.00 0.00 CC ATOM 7 0
EXPDTA THEORETICAL MODEL, MODELLER 9.14 2016/01/20 18:57:51 REMARK 6 MODELLER OBJECTIVE FUNCTION: 297.9526 REMARK 6 MODELLER BEST TEMPLATE % SEQ ID: 100.000 100.000 ATOM 1 N VAL H -10.312 3.544 19.814 1.00 0.00 N ATOM 2 CA VAL H -11.473 1.658 20.956 1.00 0.00 CD ATOM 3 CB VAL H -11.473 1.658 20.956 1.00 0.00 CD ATOM 4 CGI VAL H -10.196 0.849 21.127 1.00 0.00 CD ATOM 5 CG2 VAL H 2 -12.684 0.741 20.881 1.00 0.00 CD ATOM 4 CGI VAL H 2 -12.770 3.081 19.445 1.00 0.00 CD ATOM 7 0 VAL H 2 -13.074 3.395 18.296 1.00 0.00 CD ATOM 7
REMARK 6 HODELLER OBJECTIVE FUNCTION 120.000 REMARK 6 MODELLER BEST ITME FUNCTION 100.000 ATOM 1 N VAL H 2 -10.312 3.544 19.698 1.00 0.00 N ATOM 2 CA VAL H 2 -11.387 2.542 19.698 1.00 0.00 CC ATOM 2 CA VAL H 2 -11.387 2.542 19.698 1.00 0.00 CC ATOM 3 CB VAL H 2 -12.664 0.741 20.881 1.00 0.00 CC ATOM 5 CG2 VAL H 2 -12.770 3.081 19.445 1.00 0.00 CC ATOM 6 C VAL H 2 -13.074 3.395 18.296 1.00 0.00 CC ATOM 7 O VAL
ALTARK 6 HODELLER BEST HERELAR 5 SEQ DI: 100.000 ATOM 1 N VAL H 2 -10.312 3.544 19.814 1.00 0.00 N ATOM 2 CA VAL H 2 -11.387 2.542 19.698 1.00 0.00 CC ATOM 3 CE VAL H 2 -11.473 1.658 20.956 1.00 0.00 CC ATOM 4 CGI VAL H 2 -12.664 0.741 20.881 1.00 0.00 CC ATOM 5 CG2 VAL H 2 -12.770 3.081 19.445 1.00 0.00 CC ATOM 6 C VAL H 2 -13.074 3.395 18.296 1.00 0.00 NC ATOM 7 O VAL H 3 -14.986 3.636 20.044 1.00 0.00 NC ATOM 10 CB
AIOM 1 N VAL H 2 -10.512 3.544 19.618 1.00 0.00 N ATOM 2 CA VAL H 2 -11.473 1.658 20.956 1.00 0.00 CO ATOM 3 CB VAL H 2 -11.473 1.658 20.956 1.00 0.00 CO ATOM 4 CG1 VAL H 2 -12.684 0.741 20.881 1.00 0.00 CO ATOM 5 CG2 VAL H 2 -12.770 3.081 19.445 1.00 0.00 CO ATOM 6 C VAL H 2 -12.770 3.081 19.445 1.00 0.00 CO ATOM 6 C VAL H 2 -12.770 3.081 19.445 1.00 0.00 CO ATOM 7 0 VAL H 2 -13.074 3.395 18.296 1.00 0.00 CO ATOM 9 CA<
A10M 2 CA VAL H 2 -11.387 2.542 19.698 1.00 0.00 CO ATOM 3 CB VAL H 2 -11.473 1.658 2.9542 19.698 1.00 0.00 CO ATOM 4 CG1 VAL H 2 -12.664 0.741 20.881 1.00 0.00 CO ATOM 5 CG2 VAL H 2 -12.664 0.741 20.881 1.00 0.00 CO ATOM 6 C VAL H 2 -12.770 3.081 19.445 1.00 0.00 CO ATOM 7 O VAL H 2 -13.074 3.395 18.296 1.00 0.00 CO ATOM 8 N GLN H 3 -14.966 3.636 20.044 1.00 0.00 CO ATOM 9 CA GLN H 3 -15.890 2.416 19.853 1.00 0.00 CO ATOM 10 CB GLN H 3 -15.622 2.016 <t< td=""></t<>
A10M 3 CB VAL H 2 -11.473 1.658 20.956 1.00 0.00 CC ATOM 4 CGI VAL H 2 -12.684 0.741 20.881 1.00 0.00 CC ATOM 5 CG2 VAL H 2 -10.196 0.849 21.127 1.00 0.00 CC ATOM 6 C VAL H 2 -12.770 3.081 19.445 1.00 0.00 CC ATOM 6 C VAL H 2 -13.074 3.395 18.296 1.00 0.00 CC ATOM 7 O VAL H 2 -13.074 3.395 18.296 1.00 0.00 CO ATOM 8 N GLN H 3 -14.986 3.636 20.044 1.00 0.00 CC ATOM 10 CB GLN H 3 -15.420 1.460 18.769 1.00 0.00 CC ATOM 12 CD GLN H
AIOM 4 CGI VAL H 2 -12.884 0.741 20.881 1.00 0.00 CC ATOM 5 CG2 VAL H 2 -10.96 0.849 21.127 1.00 0.00 CC ATOM 6 C VAL H 2 -12.770 3.081 19.445 1.00 0.00 CC ATOM 6 C VAL H 2 -12.770 3.081 19.445 1.00 0.00 CC ATOM 7 0 VAL H 2 -13.074 3.395 18.296 1.00 0.00 CC ATOM 9 CA GLN H 3 -14.662 3.226 20.453 1.00 0.00 CC ATOM 9 CA GLN H 3 -15.820 2.416 19.853 1.00 0.00 CC ATOM 11 CG GLN H 3 -15.622 2.016 17.374 1.00 0.00 CC ATOM 12 CD GLN H 3 -16.525 2.819 17.140 <t< td=""></t<>
AIOM 5 CG2 VAL H 2 -10.196 0.849 21.127 1.00 0.00 C ATOM 6 C VAL H 2 -10.196 0.849 21.127 1.00 0.00 C ATOM 6 C VAL H 2 -13.074 3.395 18.296 1.00 0.00 C ATOM 8 N GLN H 3 -13.662 3.226 20.453 1.00 0.00 N ATOM 9 CA GLN H 3 -14.986 3.636 20.0453 1.00 0.00 C ATOM 9 CA GLN H 3 -15.890 2.416 19.853 1.00 0.00 C ATOM 11 CG GLN H 3 -15.622 2.016 17.374 1.00 0.00 C ATOM 13 OEI GLN H 3
ATOM 6 C VAL H 2 -12.770 3.081 19.425 1.00 0.00 CC ATOM 7 O VAL H 2 -13.074 3.395 18.296 1.00 0.00 CC ATOM 8 N GLN H 3 -13.662 3.226 20.453 1.00 0.00 N ATOM 9 CA GLN H 3 -14.986 3.636 20.044 1.00 0.00 CC ATOM 10 CB GLN H 3 -15.420 1.460 18.769 1.00 0.00 CC ATOM 11 CG GLN H 3 -15.420 1.460 18.769 1.00 0.00 CC ATOM 12 CD GLN H 3 -15.622 2.016 17.374 1.00 0.00 CC ATOM 12 CD GLN H 3 -16.525 2.819 17.140 1.00 0.00 CC ATOM
AIOM 7 0 VAL 2 -13.074 3.395 18.296 1.00 0.00 CO ATOM 8 N GLN H 3 -13.662 3.226 20.453 1.00 0.00 N ATOM 9 CA GLN H 3 -14.986 3.636 20.044 1.00 0.00 CO ATOM 10 CB GLN H 3 -15.890 2.416 19.853 1.00 0.00 CO ATOM 11 CG GLN H 3 -15.622 2.016 17.374 1.00 0.00 CO ATOM 12 CD GLN H 3 -15.622 2.016 17.374 1.00 0.00 CO ATOM 13 OEI GLN H 3 -14.778 1.590 16.441 1.00 0.00 CO ATOM 14 NE2 GLN H 3 -15.513 4.
ATOM 8 N GLN H 3 -13.662 3.226 20.453 1.00 0.00 N ATOM 9 CA GLN H 3 -14.986 3.636 20.044 1.00 0.00 CO ATOM 10 CB GLN H 3 -15.890 2.416 19.853 1.00 0.00 CO ATOM 11 CG GLN H 3 -15.420 1.460 18.769 1.00 0.00 CO ATOM 11 CG GLN H 3 -15.622 2.016 17.374 1.00 0.00 CO ATOM 13 OE1 GLN H 3 -16.525 2.819 17.140 1.00 0.00 CO ATOM 14 NE2 GLN H 3 -14.778 1.590 16.441 1.00 0.00 N ATOM 15 C GLN H
ATOM 9 CA GLN H 3 -14.986 3.636 20.044 1.00 0.00 CC ATOM 10 CB GLN H 3 -15.890 2.416 19.853 1.00 0.00 CC ATOM 11 CG GLN H 3 -15.420 1.460 18.769 1.00 0.00 CC ATOM 12 CD GLN H 3 -15.622 2.016 17.374 1.00 0.00 CC ATOM 13 OE1 GLN H 3 -16.525 2.819 17.140 1.00 0.00 CC ATOM 14 NE2 GLN H 3 -14.778 1.590 16.441 1.00 0.00 N ATOM 15 C GLN H 3 -15.613 4.581 21.036 1.00 0.00 N
ATOM 10 CB GLN H 3 -15.890 2.416 19.853 1.00 0.00 CC ATOM 11 CG GLN H 3 -15.420 1.460 18.769 1.00 0.00 CC ATOM 12 CD GLN H 3 -15.622 2.016 17.374 1.00 0.00 CC ATOM 13 OE1 GLN H 3 -16.525 2.819 17.140 1.00 0.00 CC ATOM 14 NE2 GLN H 3 -14.778 1.590 16.441 1.00 0.00 N ATOM 15 C GLN H 3 -15.613 4.581 21.036 1.00 0.00 N
ATOM 11 CG GLN H 3 -15.420 1.460 18.769 1.00 0.00 CC ATOM 12 CD GLN H 3 -15.622 2.016 17.374 1.00 0.00 CC ATOM 13 OEI GLN H 3 -16.525 2.819 17.140 1.00 0.00 CC ATOM 14 NE2 GLN H 3 -14.778 1.590 16.441 1.00 0.00 N ATOM 15 C GLN H 3 -15.613 4.581 21.036 1.00 0.00 C
ATOM 12 CD GLN H 3 -15.622 2.016 17.374 1.00 0.00 CC ATOM 13 OE1 GLN H 3 -16.525 2.819 17.140 1.00 0.00 CC ATOM 14 NE2 GLN H 3 -14.778 1.590 16.441 1.00 0.00 N ATOM 15 C GLN H 3 -15.613 4.581 21.036 1.00 0.00 C
ATOM 13 OEI GLN H 3 -16.525 2.819 17.140 1.00 0.00 C ATOM 14 NE2 GLN H 3 -14.778 1.590 16.441 1.00 0.00 N ATOM 15 C GLN H 3 -15.613 4.581 21.036 1.00 0.00 N
ATOM 14 NE2 GLN H 3 -14.778 1.590 16.441 1.00 0.00 N ATOM 15 C GLN H 3 -15.613 4.581 21.036 1.00 0.00 C
ATOM 15 C GLN H 3 -15.613 4.581 21.036 1.00 0.00 C
ATOM 16 O GLN H 3 -15.496 4.408 22.247 1.00 0.00 0
ATOM 17 N LEU H 4 -16.308 5.623 20.523 1.00 0.00 N
ATOM 18 CA LEU H 4 -17.040 6.541 21.354 1.00 0.00 C
ATOM 19 CB LEU H 4 -16.344 7.903 21.390 1.00 0.00 C
ATOM 20 CG LEU H 4 -17.047 9.003 22.188 1.00 0.00 C
ATOM 21 CD1 LEU H 4 -17.113 8.639 23.663 1.00 0.00 0
ATOM 22 CD2 LEU H 4 -16.340 10.337 22.003 1.00 0.00 0
ATOM 23 C LEU H 4 -18.425 6.616 20.781 1.00 0.00 C
ATOM 24 O LEU H 4 -18.597 6.919 19.600 1.00 0.00 0
ATOM 25 N VAL H 5 -19.460 6.339 21.604 1.00 0.00 N
ATOM 26 CA VAL H 5 -20.814 6.394 21.116 1.00 0.00 C
ATOM 27 CB VAL H 5 -21.490 5.011 21.169 1.00 0.00 C
ATOM 28 CG1 VAL H 5 -22.939 5.109 20.719 1.00 0.00 C
ATOM 29 CG2 VAL H 5 -20.728 4.013 20.310 1.00 0.00 0
ATOM 30 C VAL H 5 -21.547 7.404 21.945 1.00 0.00 C
ATOM 31 O VAL H 5 -21.410 7.435 23.167 1.00 0.00 0
ATOM 32 N GLN H 6 -22.369 8.249 21.291 1.00 0.00 N
ATOM 33 CA GLN H 6 -23.047 9.315 21.972 1.00 0.00 C
ATOM 34 CB GLN H 6 -22.757 10.655 21.293 1.00 0.00
ATOM 35 CG GLN H 6 -21.294 11.067 21.331 1.00 0.00 0
ATOM 36 CD GLN H 6 -21.049 12.407 20.666 1.00 0.00
ATOM 37 OE1 GLN H 6 -19.919 12.731 20.301 1.00 0.00
ATOM 38 NE2 GLN H 6 -22,109 13,190 20,507 1.00 0.00 N
ATOM 39 C GLN H 6 -24.522 9.037 22.000 1.00 0.00
ATOM 40 0 GLN H 6 -25,031 8,210 21,246 1,00 0,00
ATOM 41 N SER H 7 -25 246 9 732 22 904 1 00 0 00 N

1st Approach : Clustering D1 using Single descriptors





A high degree of

similarity was observed in all cases (85%-95% depending on the

structure).



Based on CATH (Protein Structure Classification database) to Validate homologues at least 2 out of 3 of the following conditions should be met:



To examine sequence similarity two new descriptors: Dayhoff & Physicochmical Similarity were examined

2nd Approach : Clustering based on 3D descriptors and Sequence Similarity



Rand Index

0.879669

0.890833

0.895234

0.937849

0.937849

Δ

2nd Approach : Clustering based on 3D descriptors and Sequence Similarity Clustering D1 using Descriptors Combination



2nd Approach : Clustering based on 3D descriptors and Sequence Similarity



				Clustering D2 using Descriptors Combin		
	3dsc	fpfh	rsd	Dayhoff	Physicochemical	Rand Index
11 clusters	0.5	0.2	0.1	0.2	0	0.993

RESULTS

Clusters were distinct ang highly homogeneous.

- This finding could be interpreted as important evidence that this novel approach enables robust clustering of patients with CLL.
- IG 3D models from cases with stereotyped subsets were clustered together strengthening the relevance of our findings.
- Altogether, clustering of CLL patients into discrete subgroups based on feature extraction from IG
 3D models may have highly relevant implications for refined patient sub-classification.

BMC Bioinformatics

METHODOLOGY





Automated shape-based clustering of 3D immunoglobulin protein structures in chronic lymphocytic leukemia

Eleftheria Polychronidou^{1*†}, Ilias Kalamaras^{1†}, Andreas Agathangelidis², Lesley-Ann Sutton⁵, Xiao-Jie Yan⁶, Vasilis Bikos⁷, Anna Vardi⁸, Konstantinos Mochament¹, Nicholas Chiorazzi⁶, Chrysoula Belessi⁹, Richard Rosenquist⁵, Paolo Ghia¹¹, Kostas Stamatopoulos², Panayiotis Vlamos¹⁰, Anna Chailyan³, Nanna Overby⁴, Paolo Marcatili⁴, Anastasia Hatzidimitriou² and Dimitrios Tzovaras¹

From 5th International Work-Conference on Bioinformatics and Biomedical Engineering Granada, Spain. 26-28 April 2017

Polychronidou E, et. al.Automated shape-based clustering of 3D immunoglobulin protein structures in chronic lymphocytic leukemia. BMC Bioinformatics. 2018 Nov 20;19(Suppl 14):414. doi: 10.1186/s12859-018-2381-1. PMID: 30453883; PMCID: PMC6245605.

Alzheimer's Disease

Alzheimer's Disease is a folding disorder

- Irreversible progressive dementia with long prodromal stages and up to now there is a lack of effective pharmacotherapy options evaluation of the prediction algorithms
- While the clinical symptoms of the disease are defined by cognitive impairment, the causes leading to memory decline are strongly tied to deposits of misfolded protein aggregates.

In Silico Analysis of the Apolipoprotein E and the Amyloid & Peptide Interaction: Misfolding Induced by Frustration of the Salt Bridge Network

Jinghui Luo¹, Jean-Didier Maréchal², Sebastian Wärmländer¹, Astrid Gräslund¹, Alex Perálvarez-Marín^{1,3*¤}

1 Department of Biochemistry and Biophysics, Stockholm University, Stockholm, Sweden, 2 Unitat de Química Física, Departament de Química, Universitat Autònoma de Barcelona, Bellaterra, Spain, 3 Unitat de Biofísica, Departament de Bioquímica i de Biologia Molecular i Centre d'Estudis en Biofísica, Universitat Autònoma de Barcelona, Bellaterra, Spain

Abstract

The relationship between Apolipoprotein E (ApoE) and the aggregation processes of the amyloid β (A β) peptide has been shown to be crucial for Alzheimer's disease (AD). The presence of the ApoE4 isoform is considered to be a contributing risk factor for AD. However, the detailed molecular properties of ApoE4 interacting with the A β peptide are unknown, although various mechanisms have been proposed to explain the physiological and pathological role of this relationship. Here, computer simulations have been used to investigate the process of A β interaction with the N-terminal domain of the human ApoE isoforms (ApoE2, ApoE3 and ApoE4). Molecular docking combined with molecular dynamics simulations have been undertaken to determine the A β peptide binding sites and the relative stability of binding to each of the ApoE isoforms. Our results show that from the several ApoE isoforms investigated, only ApoE4 presents a misfolded intermediate when bound to A β . Moreover, the initial α -helix used as the A β peptide model structure also becomes unstructured due to the interaction with ApoE4. These structural changes appear to be related to a rearrangement of the salt bridge network in ApoE4, for which we propose a model. It seems plausible that ApoE4 in its partially unfolded state is incapable of performing the clearance of A β , thereby promoting amyloid forming processes. Hence, the proposed model can be used to identify potential drug binding sites in the ApoE4-A β complex, where the interaction between the two molecules can be inhibited. Apolipoprotein E (ApoE) and the aggregation processes of the amyloid β (A β) peptide



Apolipoprotein E (ApoE) and the aggregation processes of the amyloid β (A β) peptide



Nat Rev Neurol. 2013 Feb; 9(2): 106–118.



EDBM103] Support for researchers with an emphasis on young researchers-cycle B '





What is the connection between the different isoforms of

well known proteins implicated in the pathogenesis and

progression of Alzheimer's Disease?

Could we find a correlation between their misfolded forms and the events occurring at the molecular level? • Computational simulation of the effect of an increasing electric field strength on the protein structure



J Proteomics. 2012 Oct 22; 75(18-2): 5533–5543.

• Computational simulation of the effect of an increasing electric field strength on the protein structure



Methodology

- study of the experimentally determined protein structures,
- evaluation of the prediction algorithms
- prediction of the AD related proteins with unknown structures
- evaluation of the mutation footprint of the established mutations in the tertiary structure.

Innovation

- ability to highlight broad ensembles of genes that drive key disease mechanisms
- These genes may hold strong potential for the identification of promising therapeutic targets.







Software	Method	Description	Category
MODELLER	Satisfaction of spatial restraints	Standalone program mainly in Fortran and Python	homology
SWISS-MODEL	Local similarity/fragment assembly	Automated webserver (based on ProModII)	homology
HHpred	Template detection, alignment, 3D modeling	Interactive webserver with help facility	threading
I-TASSER	Combination of ab initio folding and threading methods	Structural and function predictions	Combination
ROBETTA	Rosetta homology modeling and ab initio fragment assembly with Ginzu domain prediction	Webserver	Ab initio
QUARK	Monte Carlo fragment assembly	On-line server for protein modeling (best for ab initio folding in CASP9)	Ab initio









STRUCTURAL SIMILARITY and MUTATION analysis IN Alzheimer's disease



Summary of the 246 mutations of the Presenilin-1

>pdb|5a63|B

sp|P49768|PSN1_HUMAN Presenilin-1 OS=Homo sapiens OX=9606 GN=PSEN1 PE=1 SV=1

MTELPAPLSYFQNAQMSEDNHLSNTVRSQNDNRERQEHNDRRSLGHPEPLSNGRPQGN SRQVVEQDEEEDEELTLKYGAKHVIMLFVPVTLCMVVVVATIKSVSFYTRKDGQLIYTPFTE DTETVGQRALHSILNAAIMISVIVVMTILLVVLYKYRCYKVIHAWLIISSLLLLFFFSFIYLGEVF KTYNVAVDYITVALLIWNFGVVGMISIHWKGPLRLQQAYLIMISALMALVFIKYLPEWTAWLI LAVISVYDLVAVLCPKGPLRMLVETAQERNETLFPALIYSSTMVWLVNMAEGDPEAQRRVS KNSKYNAESTERESQDTVAENDDGGFSEEWEAQRDSHLGPHRSTPESRAAVQELSSSILA GEDPEERGVKLGLGDFIFYSVLVGKASATASGDWNTTIACFVAILIGLCLTLLLLAIFKKALP ALPISITFGLVFYFATDYLVQPFMDQLAFHQFYI

*Blue : AD - Unclear Pathogenicity, *Green : AD - Not Pathogenic, *Red: AD - Pathogenic

Presenilin 1 is a large multi-pass transmembrane protein which participates in the formation of γ -secretase complex and has a strong impact in the clinical course of the disease.

Mutations that are reported in PSEN1 protein. The mutations are grouped into three categories.

STRUCTURAL SIMILARITY and MUTATION analysis IN Alzheimer's disease

Dendrogram based on 3DSC Descriptors

Dendrogram based on FPFH Descriptors



Dendrogram based on RSD Descriptors



dist.mat2

holust (*, "complete")



Cluster Dendrogram



holest (*, "complete")



dist.mat2 hclust (*, "complete")

STRUCTURAL SIMILARITY and MUTATION analysis IN Alzheimer's disease



Higher Percentage of structure alteration due to mutation

Insights for the pathogenicity of the R35Q mutation

Summary

- Risk prediction standpoint,
- networks exhibiting coexistent genetic variation and biological perturbation would represent prime targets in the development of personalized, burden-based genetic susceptibility tests.
- Therapeutic strategies development

• pathways and networks displaying multi-omics relationships in AD would reduce the search space for rational drug design and may highlight "hub" genes for therapeutic cocktail approaches, such as in the polypharmacy strategies successfully employed for AIDS and various cancers.

> Adv Exp Med Biol. 2020;1195:227-236. doi: 10.1007/978-3-030-32633-3_31.

Alzheimer's Disease: The Role of Mutations in Protein Folding

Eleftheria Polychronidou¹, Antigoni Avramouli², Panayiotis Vlamos²

Affiliations + expand PMID: 32468481 DOI: 10.1007/978-3-030-32633-3_31

Abstract

Misfolded proteins result when a protein follows the wrong folding pathway. Accumulation of misfolded proteins can cause disorders, known as amyloid diseases. Unfortunately, some of them are very common. The most prevalent one is Alzheimer's disease. Alzheimer's disease is a neurodegenerative disorder and the commonest form of dementia. The current study aims to assess the impact of somatic mutations in PSEN1 gene. The said mutations are the most common cause of familial Alzheimer's disease. As protein functionality can be affected by mutations, the study of possible alterations in the tertiary structure of proteins may reveal new insights related to the relationship between mutations and protein functions. To examine the effect of mutations, the

Polychronidou E, Avramouli A, Vlamos P. Alzheimer's Disease: The Role of Mutations in Protein Folding. Adv Exp Med Biol. 2020;1195:227-236. doi: 10.1007/978-3-030-32633-3_31. PMID: 32468481.

5. Experimental design for studying protein misfolding

JPK Force Robot

- Atomic force microscope is a very sensitive tool for quantitatively measuring forces on a single molecule level
 - Molecular stretching experiments



JPK Force Robot



A Low entropy conformation – unlikely in solution

Contour length

B High entropy conformation – typical in solution



Typical dimensions << contour length

Fig. 2 Schematic diagram of two conformations of a long molecule.



Fig. 7 Schematic diagram of bacteriorhodopsin unfolding, with intermediate states as the alpha-helices are pulled out in pairs.



Thank you!

Computer Science is no more about computers as astronomy is about telescopes.

There should be no such thing as boring Mathematics!

Edsger Dijkstra