



Protein homology and resources for protein families

Marco Punta

Centre for Evolution and Cancer
Institute of Cancer Research
London, UK

Outline:

Now:

- Homology (definition, implications, how to detect)
- Protein domains and protein families

To follow:

- Practical: Profile-based sequence searches for protein function annotation

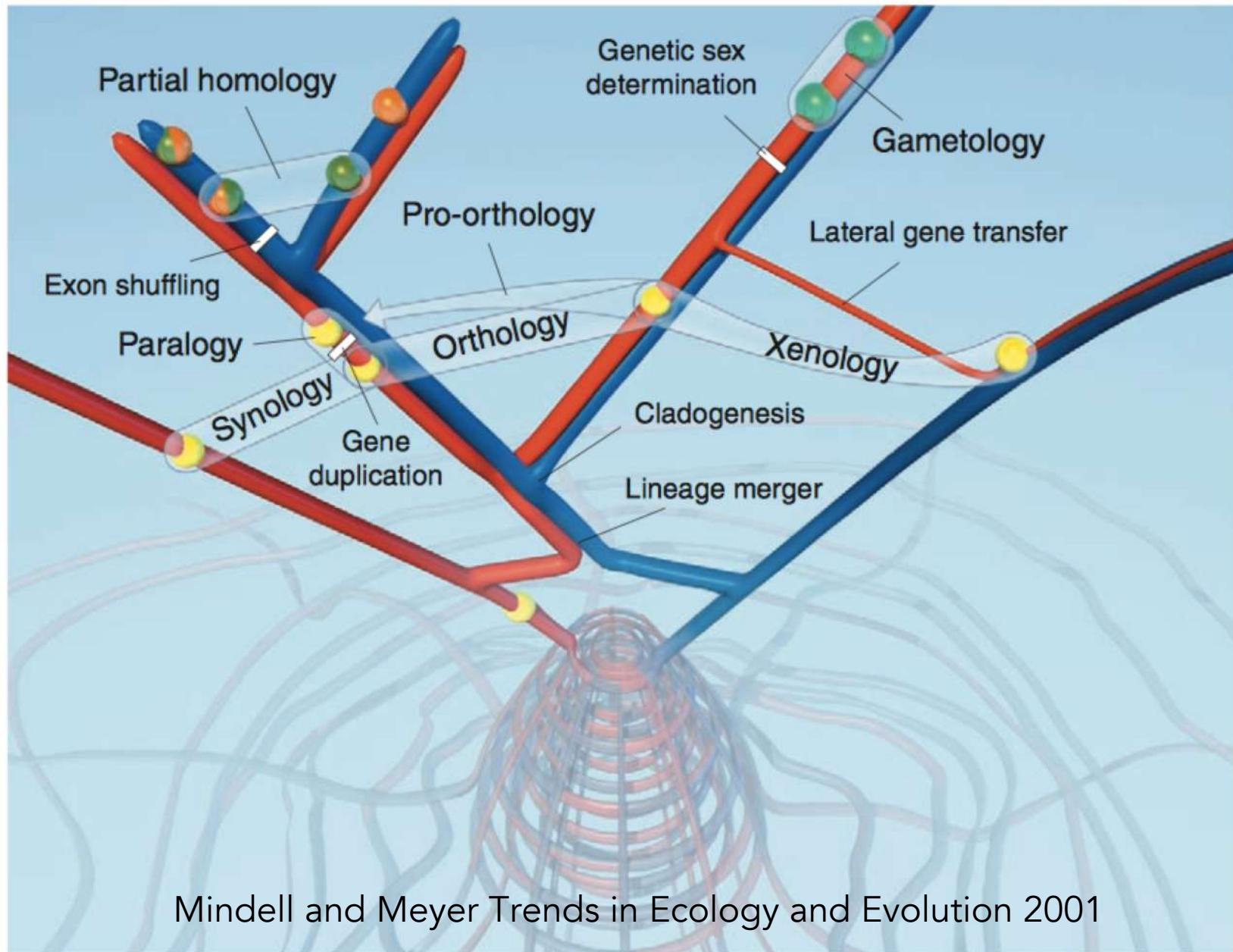
Homology

Definition:

Two proteins are **homologous** if they share a common ancestor, i.e. they are evolutionary related

Origin of homology in proteins

- Speciation (orthology)
- Gene duplication (paralogy)
- Horizontal gene transfer (xenology)
- Whole genome duplication (ohnology)
- etc. etc.

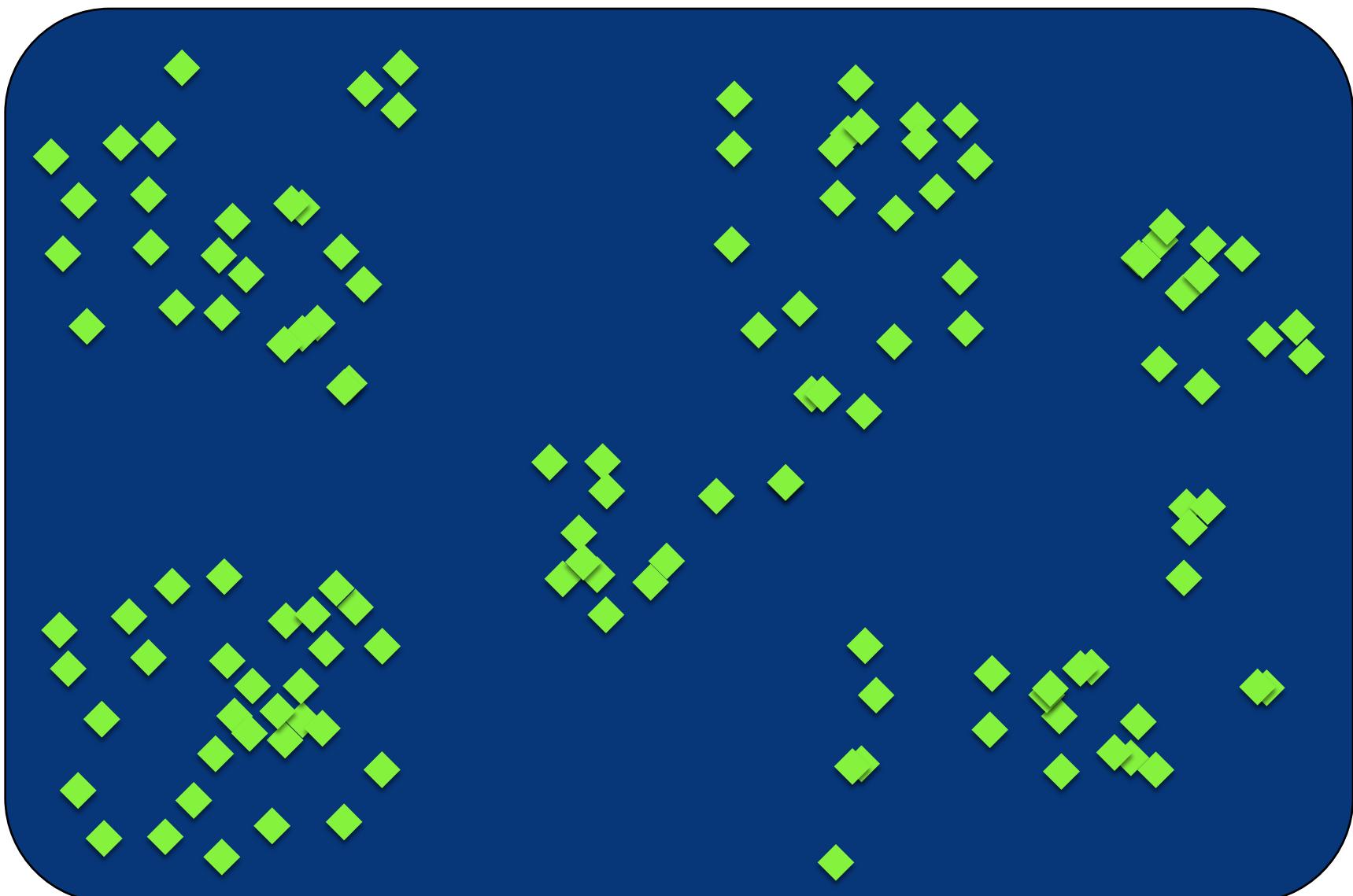


Protein Families

Definition:

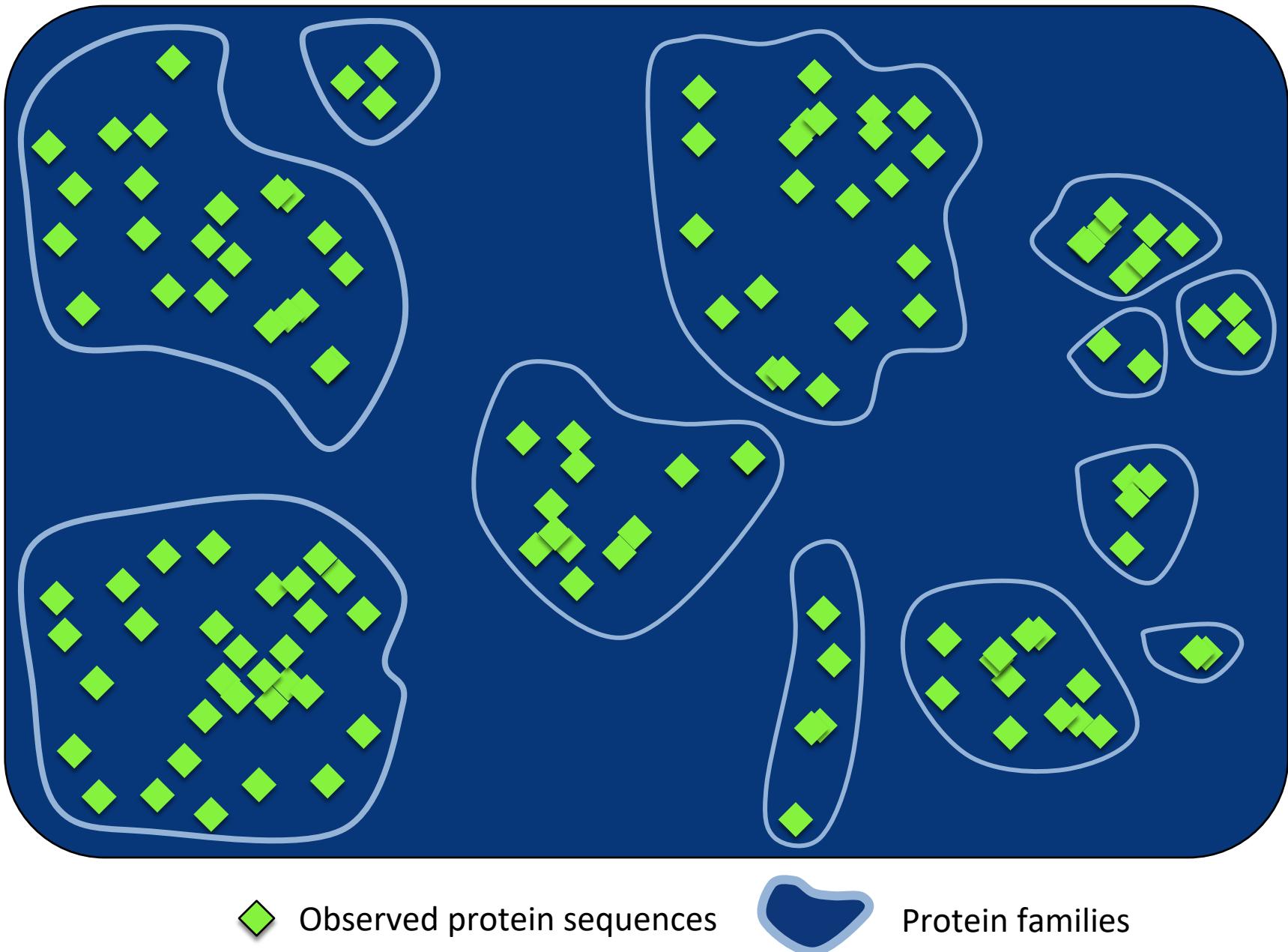
We call 'family' a group of evolutionary related proteins and/or protein regions

The sequence space



◆ Observed protein sequences

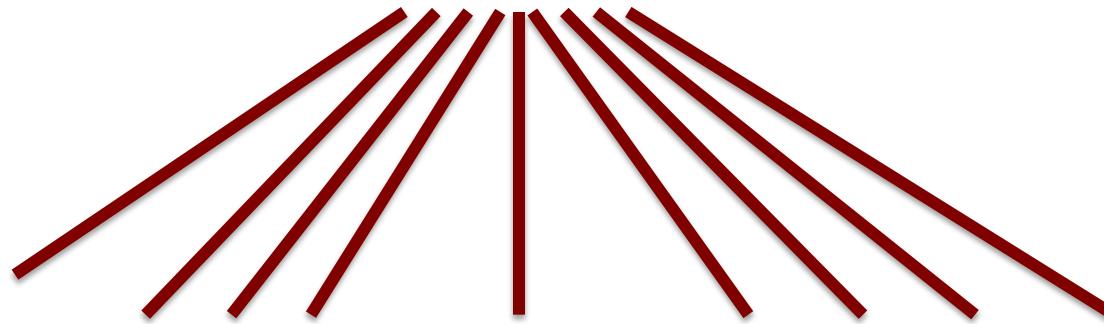
The sequence space



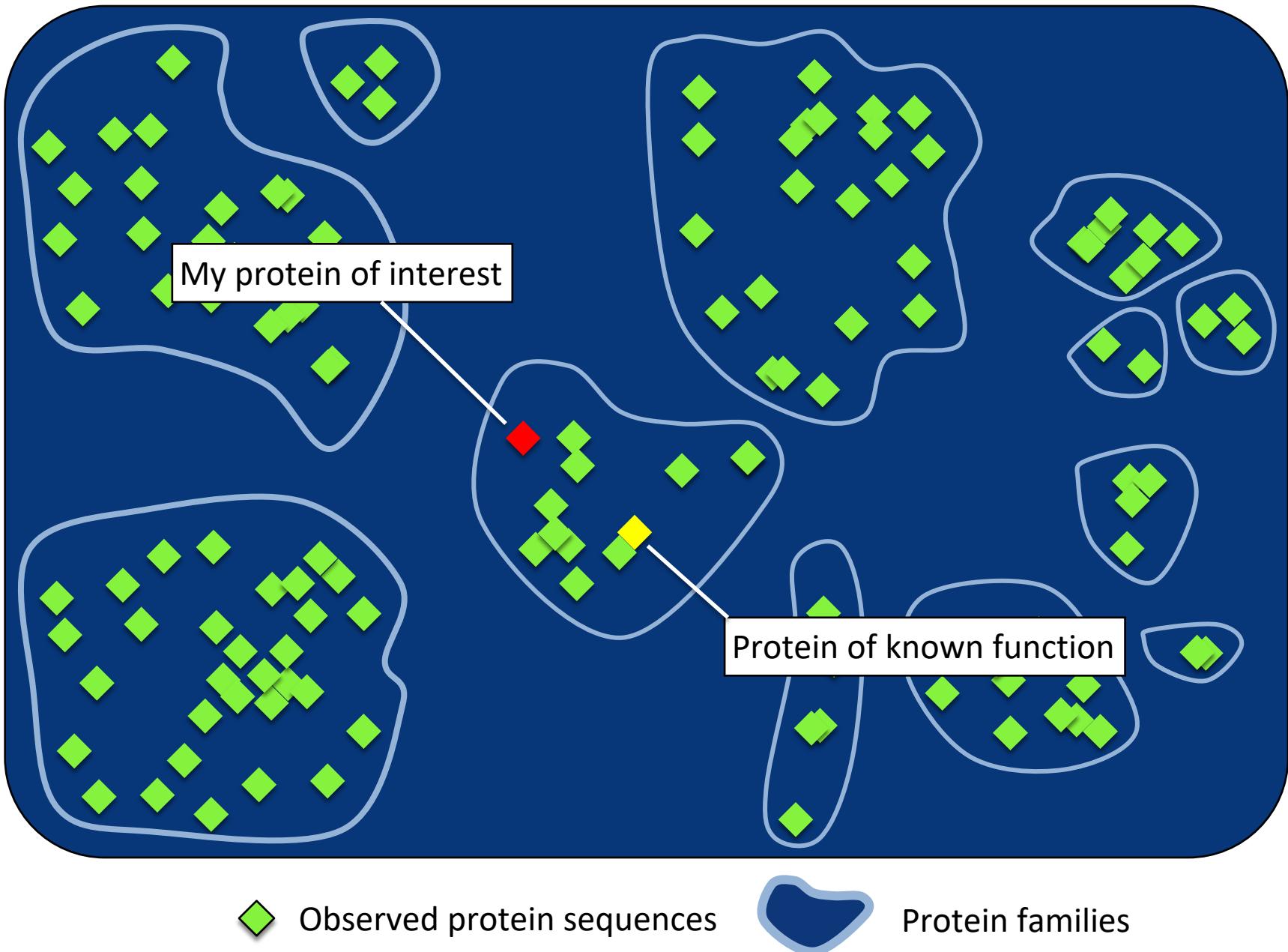
Homology: why interesting?

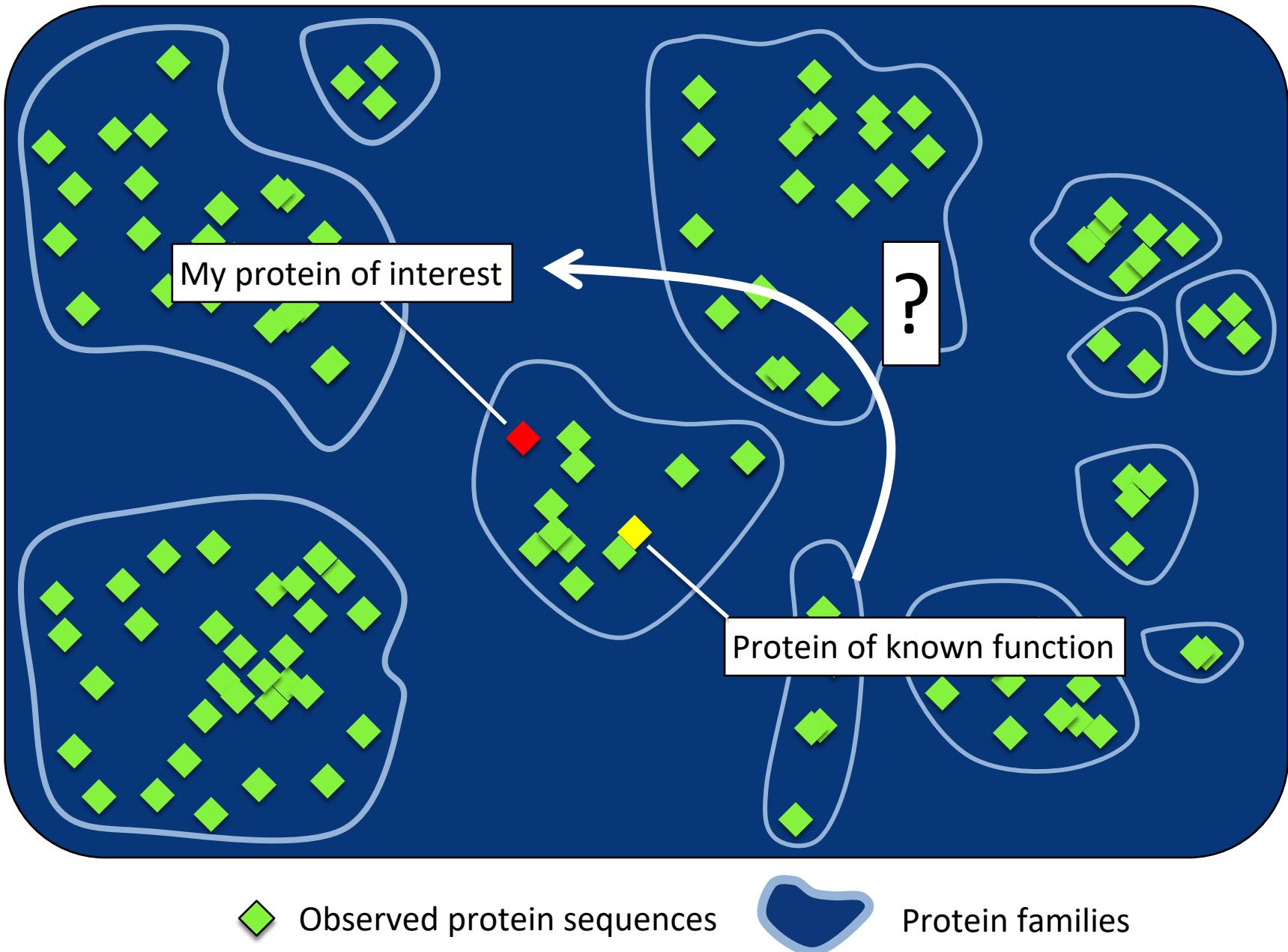
Homology: why interesting?

Common ancestor: one or more functions



Present day homologous proteins (family):
Share similar function(s)???



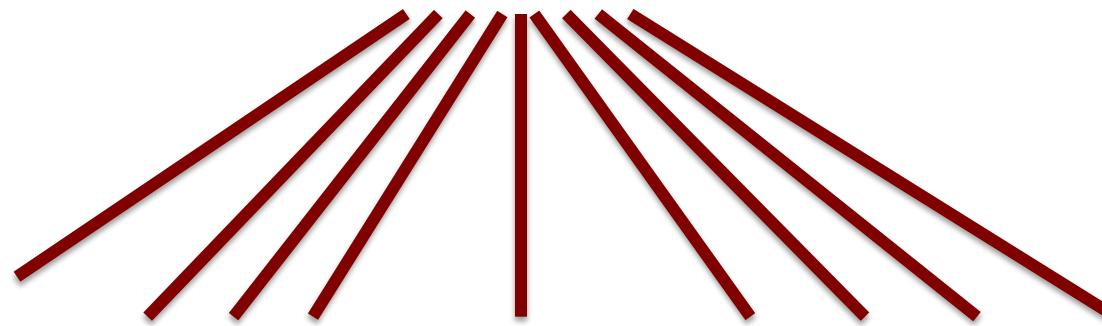


Detecting homology

How do we know two or more proteins
are homologous?

Detecting homology

Common ancestor: one sequence, one structure



Present day homologous proteins (family):
Similar sequence, similar structure???



“[...]we are justified to conclude that whenever **statistically significant sequence** or structural **similarity between proteins** or protein domains is observed, this is an indication of their divergent evolution from a common ancestor or, in other words, **evidence of homology**.”

Koonin and Galperin (2003)

Given two protein sequences how do we know if they are significantly similar (how likely that the observed similarity is random?)

Protein_1: 1 MGLSDGEWQLVLNWGKVEADIPGHGQEVLIRLFKGHPETLEKFDKFKHLKSEDEMKA 60
MGLSDGEWQLVLNWGKVEAD GHGQEVL I LFK HPETL KFDKFK LKSE MK SE

Protein_2: 1 MGLSDGEWQLVLNWGKVEADLAGHGQEVLIGLFKTHPETLDKFDKFKNLKSEEDMKG 60

Protein_1: 61 DLKKHGATVLTALGGILKKKGHHEAEIKPLAQSHATKHKIPVKYLEFISECIIQVLQSKH 120
DLKKHG TVLTALG ILKKKG H AEI PLAQSHATKHKIPVKYLEFISE II VL H

Protein_2: 61 DLKKHGCTVLTALGTILKKKGQHAAEIQPLAQSHATKHKIPVKYLEFISEIIIIEVLKKRH 120

Protein_1: 121 PGDFGADAQGAMNKALELFRKDMASNYKELGFQG 154
GDFGADAQGAM KALELFR D A YKELGFQG

Protein_2: 121 SGDFGADAQGAMS KALELFRNDIAAKYKELGFQG 154

Given two protein sequences how do we know if they are significantly similar (how likely that the observed similarity is random?)

Protein 1 1 MGLSDGEWQLVLNWGKVEADIPGHGQEVLIRLFKGHPETLEKFDKFKHLKSEDEMKA 60
M LS V WGKV A G E L R F P T F F D S

Protein 2 1 MVLSPADKTNVKAAWGKVGAGAEALERMFLSFPTTKTYFPHF-----DLSHGSA 54

Protein 1 61 DLKKHGATVLTALGGILKKKGHEAEIKPLAQSHATKHKI-PVKY 104
K H V AL L H A K PV

Protein 2 55 QVKGHSKKVADALTNAVAHVDDMPNALSALSDLHAHKLRVDPVN 99

What we need:

- Scoring system => empirically derived substitution matrices (PAMs, BLOSUMs,...)
- Efficient way to find highest scoring alignments => dynamic programming (Needleman-Wunsch, Smith-Waterman,...)
- Way to decide whether top score is high enough to infer homology (significance) => E-value, ...

BLOSUM62 matrix

Ala	4																			
Arg	-1	5																		
Asn	-2	0	6																	
Asp	-2	-2	1	6																
Cys	0	-3	-3	-3	9															
Gln	-1	1	0	0	-3	5														
Glu	-1	0	0	2	-4	2	5													
Gly	0	-2	0	-1	-3	-2	-2	6												
His	-2	0	1	-1	-3	0	0	-2	8											
Ile	-1	-3	-3	-3	-1	-3	-3	-4	-3	4										
Leu	-1	-2	-3	-4	-1	-2	-3	-4	-3	2	4									
Lys	-1	2	0	-1	-3	1	1	-2	-1	-3	-2	5								
Met	-1	-1	-2	-3	-1	0	-2	-3	-2	1	2	-1	5							
Phe	-2	-3	-3	-3	-2	-3	-3	-3	-1	0	0	-3	0	6						
Pro	-1	-2	-2	-1	-3	-1	-1	-2	-2	-3	-3	-1	-2	-4	7					
Ser	1	-1	1	0	-1	0	0	0	-1	-2	-2	0	-1	-2	-1	4				
Thr	0	-1	0	-1	-1	-1	-1	-2	-2	-1	-1	-1	-1	-2	-1	1	5			
Trp	-3	-3	-4	-4	-2	-2	-3	-2	-2	-3	-2	-3	-1	1	-4	-3	-2	11		
Tyr	-2	-2	-2	-3	-2	-1	-2	-3	2	-1	-1	-2	-1	3	-3	-2	-2	2	7	
Val	0	-3	-3	-3	-1	-2	-2	-3	-3	3	1	-2	1	-1	-2	-2	0	-3	-1	4
	Ala	Arg	Asn	Asp	Cys	Gln	Glu	Gly	His	Ile	Leu	Lys	Met	Phe	Pro	Ser	Thr	Trp	Tyr	Val

What we need:

- Scoring system => empirically derived substitution matrices (PAMs, BLOSUMs,...)
- Efficient way to find highest scoring alignments => dynamic programming (Needleman-Wunsch, Smith-Waterman,...)
- Way to decide whether top score is high enough to infer homology (significance) => E-value, ...

Statistical significance: E-values

My alignment score = S_0

The E-value of the alignment tells me how many alignments between unrelated sequences I can expect to find that have $S \geq S_0$ when searching the same database with my query sequence

If E-value = 10, I expect to find 10 such alignments

If E-value 10^{-5} , it is estimated that I would have to run 100,000 such searches before I find one such score between unrelated sequences

Excess sequence similarity → Homology

There are alternative possible explanations for excess sequence similarity (i.e. analogy due to functional or structural convergence)

In general, we are guided by observation (what is known today) and “the principles of parsimony (Occams’ razor) and likelihood”*

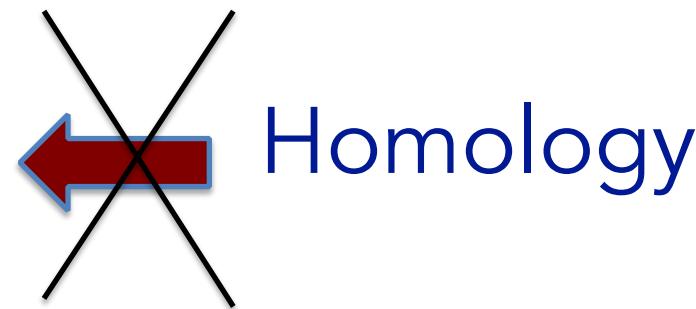
*Computational Structural Biology: Methods and Applications Torsten Schwede (2008)

Excess sequence
similarity → Homology

Cases where significance may misguide: low complexity
regions (homopolymers) (filter out), coiled-coils
(see e.g. Mistry et al NAR 2013)

The reverse statement is not true (plenty of examples)

Excess sequence similarity



Suggested reading

Curr Protoc Bioinformatics. Author manuscript; available in PMC 2014 Jun 1.

Published in final edited form as:

[Curr Protoc Bioinformatics. 2013 Jun; 0 3: 10.1002/0471250953.bi0301s42.](#)

doi: [10.1002/0471250953.bi0301s42](https://doi.org/10.1002/0471250953.bi0301s42)

PMCID: PMC3820096

NIHMSID: NIHMS519883

An Introduction to Sequence Similarity (“Homology”) Searching

[William R. Pearson¹](#)

[Author information ▶](#) [Copyright and License information ▶](#)

The publisher's final edited version of this article is available at [Curr Protoc Bioinformatics](#)

See other articles in PMC that [cite](#) the published article.

Abstract

[Go to: ▾](#)

Sequence similarity searching, typically with BLAST (units 3.3, 3.4), is the most widely used, and most reliable, strategy for characterizing newly determined sequences. Sequence similarity searches can identify “homologous” proteins or genes by detecting excess similarity – statistically significant similarity that reflects common ancestry. This unit provides an overview of the inference of homology from significant similarity, and introduces other units in this chapter that provide more details on effective strategies for identifying homologs.

Keywords: sequence similarity, homology, orthlogy, paralogy, sequence alignment, multiple alignment, sequence evolution

Subject: Bioinformatics, Bioinformatics Fundamentals, Finding Similarities and Inferring Homologies



“[...]we are justified to conclude that whenever **statistically significant** sequence or **structural similarity between proteins** or protein domains is observed, this is an indication of their divergent evolution from a common ancestor or, in other words, **evidence of homology**.”

Koonin and Galperin (2003)



Structure more conserved than sequence

Proteins. 2009 Nov 15;77(3):499-508. doi: 10.1002/prot.22458.

Structure is three to ten times more conserved than sequence--a study of structural response in protein cores.

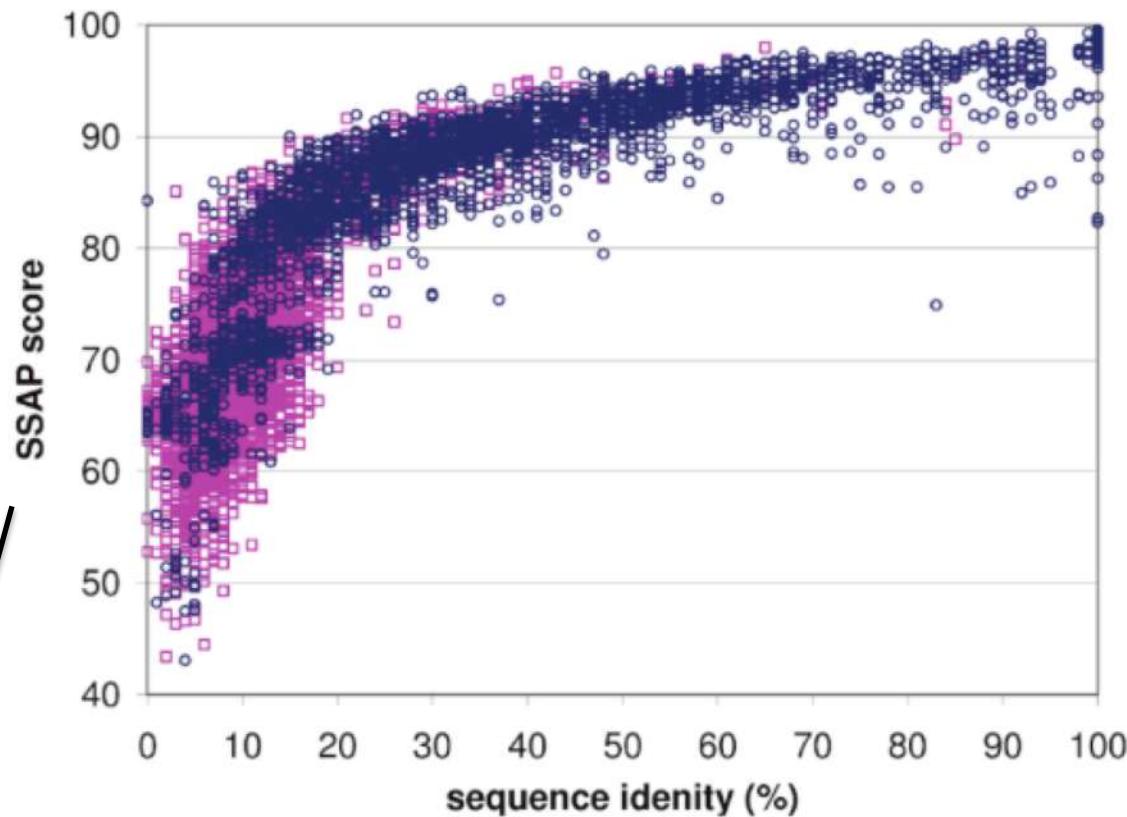
Illergård K¹, Ardell DH, Elofsson A.

Author information

Abstract

Protein structures change during evolution in response to mutations. Here, we analyze the mapping between sequence and structure in a set of structurally aligned protein domains. To avoid artifacts, we restricted our attention only to the core components of these structures. We found that on average, using different measures of structural change, protein cores evolve linearly with evolutionary distance (amino acid substitutions per site). This is true irrespective of which measure of structural change we used, whether RMSD or discrete structural descriptors for secondary structure, accessibility, or contacts. This linear response allows us to quantify the claim that structure is more conserved than sequence. Using structural alphabets of similar cardinality to the sequence alphabet, structural cores evolve three to ten times slower than sequences. Although we observed an average linear response, we found a wide variance. Different domain families varied fivefold in structural response to evolution. An attempt to categorically analyze this variance among subgroups by structural and functional category revealed only one statistically significant trend. This trend can be explained by the fact that beta-sheets change faster than alpha-helices, most likely due to that they are shorter and that change occurs at the ends of the secondary structure elements.

Structure vs Sequence similarity of homologous proteins



Structural similarity

Issues in structural comparison

- Statistical framework for structural alignments less solid than the one for sequence alignments:
- How should we define the alignment score?
- How do we find the optimal structural alignment?
- How do we define significance?
- Different methods may give different answers...

In practice

- Structural similarity important for suggesting homology in protein regions sharing insignificant sequence similarity, however:
- we generally have to look for additional signs of homology for validation (e.g. conservation structural and/or functional residues)
- Most protein families with no known function have no member of known structure

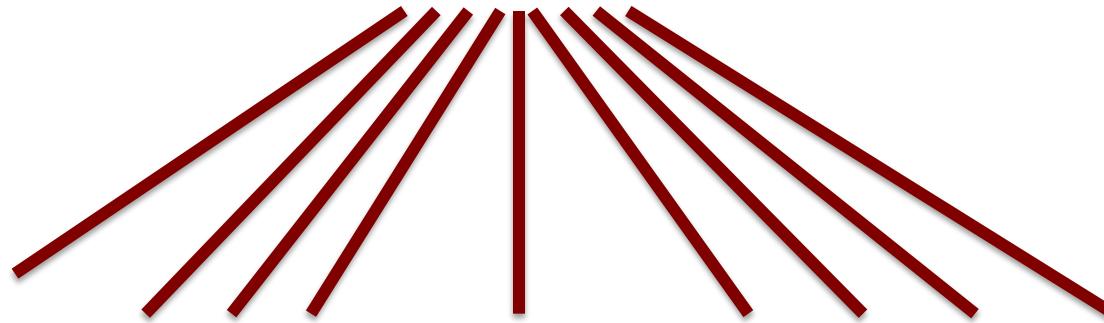
Here, we will focus on homology detection via
excess sequence similarity

Homology modeling

Same structural core ← Homology

Homology: why interesting?

Common ancestor: one or more functions



Present day homologous proteins (family):
Share similar function(s)???

How to compare function?



How to compare function?

CbiF

Catalytic activity



Methyltransferase

precorrin-4
methyltransferase activity**CbiA**

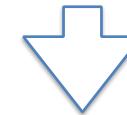
Catalytic activity



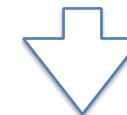
Ligase activity

cobyrinic acid a,c-diamide
synthase activity**CbiJ**

Catalytic activity



Reductase activity



precorrin-6A reductase activity

How to compare function?

CbiF

Catalytic activity



Methyltransferase

precorrin-4
methyltransferase activity**CbiA**

Catalytic activity



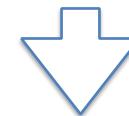
Ligase activity

cobyrinic acid a,c-diamide
synthase activity**CbiJ**

Catalytic activity



Reductase activity



precorrin-6A reductase activity

Molecular function



How to compare function?

CbiF

CbiA

CbiJ

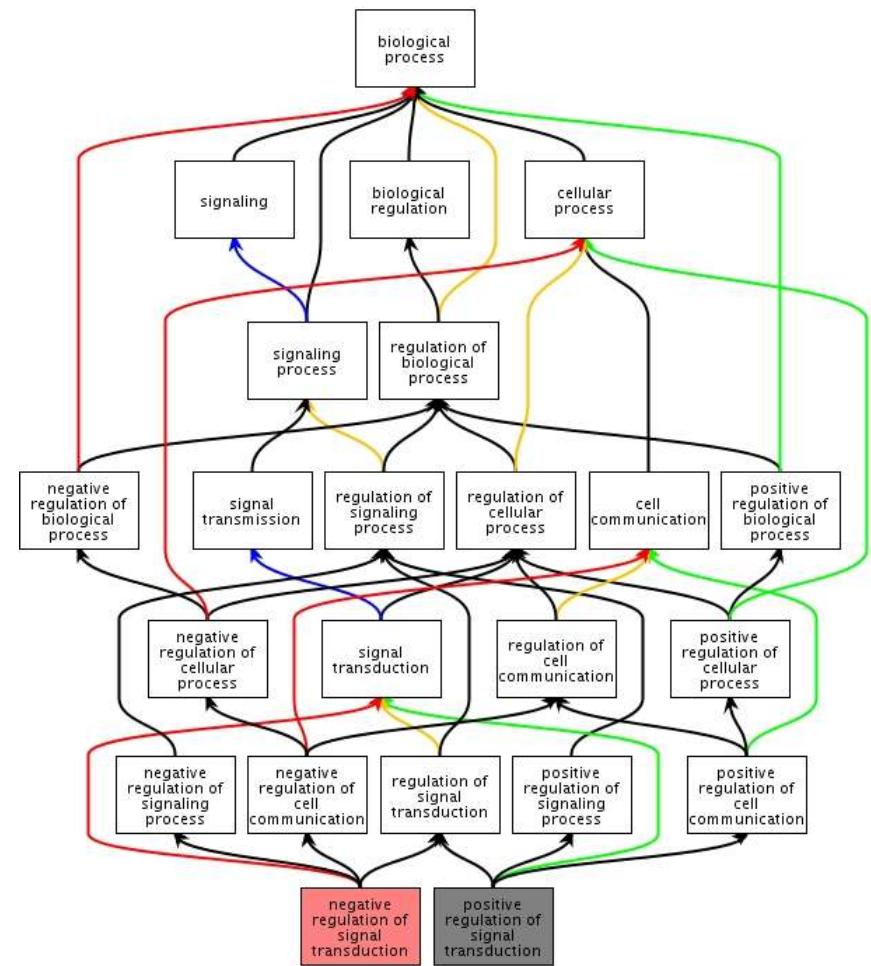


Cobalamin
biosynthetic process

Biological process

The Gene Ontology (GO)

- A way to capture biological knowledge in a written and computable form
- A set of concepts and their relationships to each other

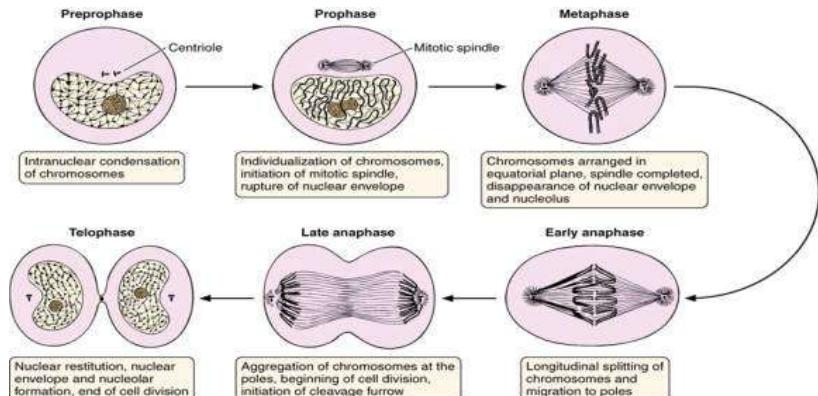


www.ebi.ac.uk/QuickGO

GO: 3 ontologies in 1

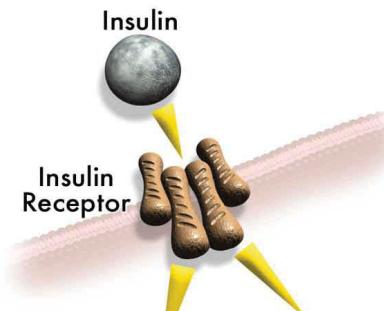
1. Molecular Function

An elemental activity or task or job



3. Cellular Component

Where a gene product is located

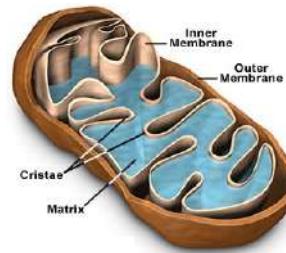


- protein kinase activity
- insulin receptor activity

2. Biological Process

A commonly recognised series of events

- cell division



- mitochondrion
- mitochondrial matrix
- mitochondrial inner membrane

Top-level EC numbers^[5]

Group	Reaction catalyzed	Typical reaction	Enzyme example(s) with trivial name
EC 1 <i>Oxidoreductases</i>	To catalyze oxidation/reduction reactions; transfer of H and O atoms or electrons from one substance to another	$AH + B \rightarrow A + BH$ (reduced) $A + O \rightarrow AO$ (oxidized)	Dehydrogenase, oxidase
EC 2 <i>Transferases</i>	Transfer of a functional group from one substance to another. The group may be methyl-, acyl-, amino- or phosphate group	$AB + C \rightarrow A + BC$	Transaminase, kinase
EC 3 <i>Hydrolases</i>	Formation of two products from a substrate by hydrolysis	$AB + H_2O \rightarrow AOH + BH$	Lipase, amylase, peptidase
EC 4 <i>Lyases</i>	Non-hydrolytic addition or removal of groups from substrates. C-C, C-N, C-O or C-S bonds may be cleaved	$RCOCOOH \rightarrow RCOH + CO_2$ or $[X-A-B-Y] \rightarrow [A=B + X-Y]$	Decarboxylase
EC 5 <i>Isomerases</i>	Intramolecule rearrangement, i.e. isomerization changes within a single molecule	$ABC \rightarrow BCA$	Isomerase, mutase
EC 6 <i>Ligases</i>	Join together two molecules by synthesis of new C-O, C-S, C-N or C-C bonds with simultaneous breakdown of ATP	$X + Y + ATP \rightarrow XY + ADP + Pi$	Synthetase

Structure vs Sequence similarity of homologous proteins

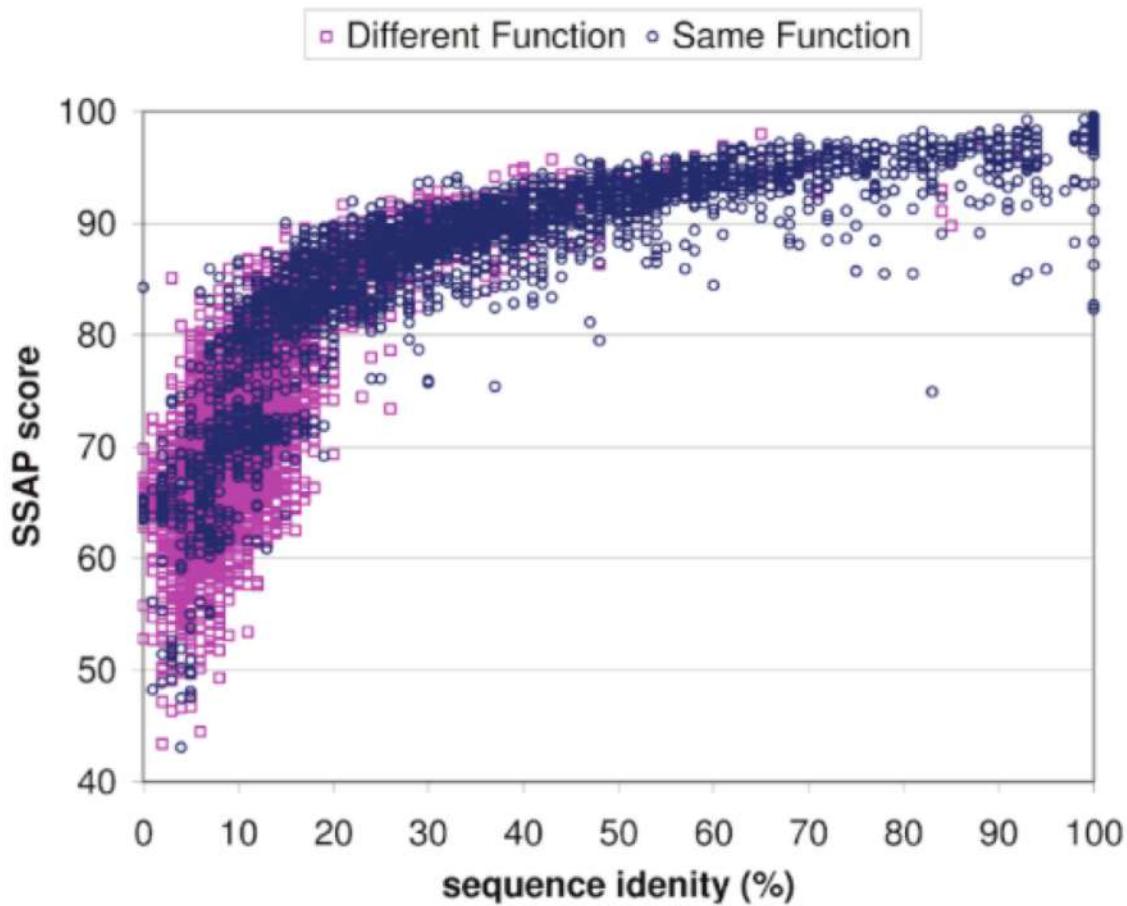


Fig. 7.8 Scatter plot showing the relationship between sequence, structure, and function of all homologues in enzyme superfamilies. Relatives having the same EC classification number are shown in blue. Those with different EC numbers are shown in pink.

Conservation of function in homologs

- Correlates with sequence/structural similarity
- No safe threshold
- Even single mutations can induce important changes in function (e.g. oncogenic and resistance point mutations in cancer)
- Protein may share only some of their functions
- Always look for additional evidence related to function (functional residues, functional motifs etc.)

Database	Gene Product ID	Symbol	Qualifier	GO Identifier	GO Term Name	Aspect	Evidence	Reference	With	Taxon	Date	Assigned By	Product Form ID
Process													
UniProtKB	P02144	MB		GO:0001666	response to hypoxia	P	IEA	Ensembl Compara	Ensembl:ENSMUSP00000125995	9606	20140913	Ensembl	
UniProtKB	P02144	MB		GO:0006810	transport	P	IEA	UniProt Keywords2GO (UniProtKB/Swiss-Prot entries)	UniProtKB-KW:KW-0813	9606	20140913	UniProt	
Function													
UniProtKB	P02144	MB		GO:0007507	heart development	P	IEA	Ensembl Compara	Ensembl:ENSMUSP00000125995	9606	20140913	Ensembl	
UniProtKB	P02144	MB		GO:0009725	response to hormone	P	IEA	Ensembl Compara	Ensembl:ENSRNOP0000006184	9606	20140913	Ensembl	
UniProtKB	P02144	MB		GO:0015671	oxygen transport	P	IEA	InterPro2GO	InterPro:IPR002335 InterPro:IPR012292	9606	20140913	InterPro	
UniProtKB	P02144	MB		GO:0015671	oxygen transport	P	IEA	Ensembl Compara	Ensembl:ENSRNOP0000006184	9606	20140913	Ensembl	
UniProtKB	P02144	MB		GO:0015671	oxygen transport	P	IEA	UniProt Keywords2GO (UniProtKB/Swiss-Prot entries)	UniProtKB-KW:KW-0561	9606	20140913	UniProt	
UniProtKB	P02144	MB		GO:0031444	slow-twitch skeletal muscle fiber contraction	P	IEA	Ensembl Compara	Ensembl:ENSRNOP0000006184	9606	20140913	Ensembl	
UniProtKB	P02144	MB		GO:0042542	response to hydrogen peroxide	P	IEA	Ensembl Compara	Ensembl:ENSRNOP0000006184	9606	20140913	Ensembl	
UniProtKB	P02144	MB		GO:0043353	enucleate erythrocyte differentiation	P	IEA	Ensembl Compara	Ensembl:ENSMUSP00000125995	9606	20140913	Ensembl	
UniProtKB	P02144	MB		GO:0050873	brown fat cell differentiation	P	IEA	Ensembl Compara	Ensembl:ENSMUSP00000125995	9606	20140913	Ensembl	
Component													
UniProtKB	P02144	MB		GO:0070062	extracellular vesicular exosome	C	IDA	PMID:23533145		9606	20140714	UniProt	
Database	Gene Product ID	Symbol	Qualifier	GO Identifier	GO Term Name	Aspect	Evidence	Reference	With	Taxon	Date	Assigned By	Product Form ID
Process													
UniProtKB	P02008	HBZ		GO:0000122	negative regulation of transcription from RNA polymerase II promoter	P	IEA	Ensembl Compara	Ensembl:ENSMUSP0000020531	9606	20140913	Ensembl	
UniProtKB	P02008	HBZ		GO:0006810	transport	P	IEA	UniProt Keywords2GO (UniProtKB/Swiss-Prot entries)	UniProtKB-KW:KW-0813	9606	20140913	UniProt	
UniProtKB	P02008	HBZ		GO:0015671	oxygen transport	P	IEA	InterPro2GO	InterPro:IPR002338 InterPro:IPR002340 InterPro:IPR012292	9606	20140913	InterPro	
UniProtKB	P02008	HBZ		GO:0015671	oxygen transport	P	IEA	UniProt Keywords2GO (UniProtKB/Swiss-Prot entries)	UniProtKB-KW:KW-0561	9606	20140913	UniProt	
UniProtKB	P02008	HBZ		GO:0043249	erythrocyte maturation	P	IEA	Ensembl Compara	Ensembl:ENSMUSP0000020531	9606	20140913	Ensembl	
UniProtKB	P02008	HBZ		GO:0005344	oxygen transporter activity	F	IEA	UniProt Keywords2GO (UniProtKB/Swiss-Prot entries)	UniProtKB-KW:KW-0561	9606	20140913	UniProt	
UniProtKB	P02008	HBZ		GO:0005344	oxygen transporter activity	F	TAS	PMID:7555018		9606	20030904	PINC	
UniProtKB	P02008	HBZ		GO:0005506	iron ion binding	F	IEA	InterPro2GO	InterPro:IPR000971 InterPro:IPR002338 InterPro:IPR002340 InterPro:IPR012292	9606	20140913	InterPro	
UniProtKB	P02008	HBZ		GO:0005515	protein binding	F	IPI	PMID:11159543	UniProtKB:P68871	9606	20140914	IntAct	
UniProtKB	P02008	HBZ		GO:0005515	protein binding	F	IPI	PMID:6683087	UniProtKB:P68871	9606	20140914	IntAct	
UniProtKB	P02008	HBZ		GO:0019825	oxygen binding	F	IEA	InterPro2GO	InterPro:IPR002338 InterPro:IPR002340 InterPro:IPR012292	9606	20140913	InterPro	
UniProtKB	P02008	HBZ		GO:0020037	heme binding	F	IEA	InterPro2GO	InterPro:IPR000971 InterPro:IPR002338 InterPro:IPR002340 InterPro:IPR012292	9606	20140913	InterPro	
UniProtKB	P02008	HBZ		GO:0046872	metal ion binding	F	IEA	UniProt Keywords2GO (UniProtKB/Swiss-Prot entries)	UniProtKB-KW:KW-0479	9606	20140913	UniProt	
Component													
UniProtKB	P02008	HBZ		GO:0005833	hemoglobin complex	C	IEA	InterPro2GO	InterPro:IPR002338 InterPro:IPR002340	9606	20140913	InterPro	
UniProtKB	P02008	HBZ		GO:0005833	hemoglobin complex	C	TAS	PMID:7555018		9606	20030904	PINC	
UniProtKB	P02008	HBZ		GO:0070062	extracellular vesicular exosome	C	IDA	PMID:23533145		9606	20140714	UniProt	

Database	Gene Product ID	Symbol	Qualifier	GO Identifier	GO Term Name	Aspect	Evidence	Reference	With	Taxon	Date	Assigned By	Product Form ID
Process													
UniProtKB	P02144	MB	GO:0001666	response to hypoxia	P	IEA	Ensembl Compara	Ensembl:ENSMUSP00000125995		9606	20140913	Ensembl	
UniProtKB	P02144	MB	GO:0006810	transport	P	IEA	UniProt Keywords2GO (UniProtKB/Swiss-Prot entries)	UniProtKB-KW:KW-0813		9606	20140913	UniProt	
UniProtKB	P02144	MB	GO:0007507	heart development	P	IEA	Ensembl Compara	Ensembl:ENSMUSP00000125995		9606	20140913	Ensembl	
UniProtKB	P02144	MB	GO:0009725	response to hormone	P	IEA	Ensembl Compara	Ensembl:ENSRNOP0000006184		9606	20140913	Ensembl	
UniProtKB	P02144	MB	GO:0015671	oxygen transport	P	IEA	InterPro2GO	InterPro:IPR002335 InterPro:IPR012292		9606	20140913	InterPro	
UniProtKB	P02144	MB	GO:0015671	oxygen transport	P	IEA	Ensembl Compara	Ensembl:ENSRNOP0000006184		9606	20140913	Ensembl	
UniProtKB	P02144	MB	GO:0015671	oxygen transport	P	IEA	UniProt Keywords2GO (UniProtKB/Swiss-Prot entries)	UniProtKB-KW:KW-0561		9606	20140913	UniProt	
UniProtKB	P02144	MB	GO:0031444	slow-twitch skeletal muscle fiber contraction	P	IEA	Ensembl Compara	Ensembl:ENSRNOP0000006184		9606	20140913	Ensembl	
UniProtKB	P02144	MB	GO:0042542	response to hydrogen peroxide	P	IEA	Ensembl Compara	Ensembl:ENSRNOP0000006184		9606	20140913	Ensembl	
UniProtKB	P02144	MB	GO:0043353	enucleate erythrocyte differentiation	P	IEA	Ensembl Compara	Ensembl:ENSMUSP00000125995		9606	20140913	Ensembl	
UniProtKB	P02144	MB	GO:0050873	brown fat cell differentiation	P	IEA	Ensembl Compara	Ensembl:ENSMUSP00000125995		9606	20140913	Ensembl	
Function													
UniProtKB	P02144	MB	GO:0005344	oxygen transporter activity	F	IEA	UniProt Keywords2GO (UniProtKB/Swiss-Prot entries)	UniProtKB-KW:KW-0561		9606	20140913	UniProt	
UniProtKB	P02144	MB	GO:0005506	iron ion binding	F	IEA	InterPro2GO	InterPro:IPR000971 InterPro:IPR012292		9606	20140913	InterPro	
UniProtKB	P02144	MB	GO:0019825	oxygen binding	F	IEA	InterPro2GO	InterPro:IPR002335 InterPro:IPR012292		9606	20140913	InterPro	
UniProtKB	P02144	MB	GO:0019825	oxygen binding	F	IEA	Ensembl Compara	Ensembl:ENSRNOP0000006184		9606	20140913	Ensembl	
UniProtKB	P02144	MB	GO:0020037	heme binding	F	IEA	InterPro2GO	InterPro:IPR000971 InterPro:IPR002335 InterPro:IPR012292	9606	20140913	InterPro		
UniProtKB	P02144	MB	GO:0046872	metal ion binding	F	IEA	UniProt Keywords2GO (UniProtKB/Swiss-Prot entries)	UniProtKB-KW:KW-0479		9606	20140913	UniProt	
Component													
UniProtKB	P02144	MB	GO:0070062	extracellular vesicular exosome	C	IDA	PMID:23533145			9606	20140714	UniProt	

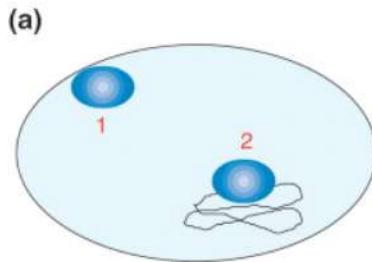
Database	Gene Product ID	Symbol	Qualifier	GO Identifier	GO Term Name	Aspect	Evidence	Reference	With	Taxon	Date	Assigned By	Product Form ID
Process													
UniProtKB	P02008	HBZ	GO:0000122	negative regulation of transcription from RNA polymerase II promoter	P	IEA	Ensembl Compara	Ensembl:ENSMUSP0000020531		9606	20140913	Ensembl	
UniProtKB	P02008	HBZ	GO:0006810	transport	P	IEA	UniProt Keywords2GO (UniProtKB/Swiss-Prot entries)	UniProtKB-KW:KW-0813		9606	20140913	UniProt	
UniProtKB	P02008	HBZ	GO:0015671	oxygen transport	P	IEA	InterPro2GO	InterPro:IPR002338 InterPro:IPR002340 InterPro:IPR012292		9606	20140913	InterPro	
UniProtKB	P02008	HBZ	GO:0015671	oxygen transport	P	IEA	UniProt Keywords2GO (UniProtKB/Swiss-Prot entries)	UniProtKB-KW:KW-0561		9606	20140913	UniProt	
UniProtKB	P02008	HBZ	GO:0043249	erythrocyte maturation	P	IEA	Ensembl Compara	Ensembl:ENSMUSP0000020531		9606	20140913	Ensembl	
UniProtKB	P02008	HBZ	GO:0005344	oxygen transporter activity	F	IEA	UniProt Keywords2GO (UniProtKB/Swiss-Prot entries)	UniProtKB-KW:KW-0561		9606	20140913	UniProt	
UniProtKB	P02008	HBZ	GO:0005344	oxygen transporter activity	F	TAS	PMID:7555018			9606	20030904	PINC	
UniProtKB	P02008	HBZ	GO:0005506	iron ion binding	F	IEA	InterPro2GO	InterPro:IPR000971 InterPro:IPR002338 InterPro:IPR002340 InterPro:IPR012292	9606	20140913	InterPro		
UniProtKB	P02008	HBZ	GO:0005515	protein binding	F	IPI	PMID:115871	InterPro:IPR000971 InterPro:IPR002338 InterPro:IPR002340 InterPro:IPR012292	9606	20140914	IntAct		
UniProtKB	P02008	HBZ	GO:0005515	protein binding	F	IPI	PMID:6	InterPro:IPR000971 InterPro:IPR002338 InterPro:IPR002340 InterPro:IPR012292	9606	20140914	IntAct		
UniProtKB	P02008	HBZ	GO:0019825	oxygen binding	F	IEA	InterPro2GO	InterPro:IPR002338 InterPro:IPR002340 InterPro:IPR012292	9606	20140913	InterPro		
UniProtKB	P02008	HBZ	GO:0020037	heme binding	F	IEA	InterPro2GO	InterPro:IPR000971 InterPro:IPR002338 InterPro:IPR002340 InterPro:IPR012292	9606	20140913	InterPro		
UniProtKB	P02008	HBZ	GO:0046872	metal ion binding	F	IEA	UniProt Keywords2GO (UniProtKB/Swiss-Prot entries)	UniProtKB-KW:KW-0479		9606	20140913	UniProt	
Component													
UniProtKB	P02008	HBZ	GO:0005833	hemoglobin complex	C	IEA	InterPro2GO	InterPro:IPR002338 InterPro:IPR002340		9606	20140913	InterPro	
UniProtKB	P02008	HBZ	GO:0005833	hemoglobin complex	C	TAS	PMID:7555018			9606	20030904	PINC	
UniProtKB	P02008	HBZ	GO:0070062	extracellular vesicular exosome	C	IDA	PMID:23533145			9606	20140714	UniProt	

SAME

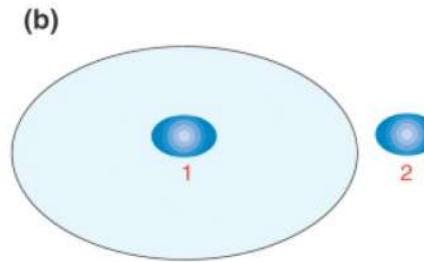
SAME

SAME

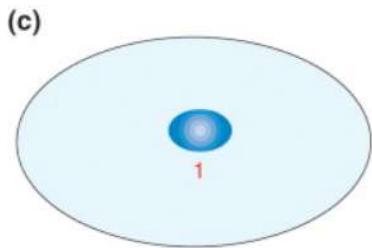
Moonlighting proteins



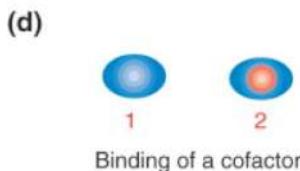
Different locations within the cell



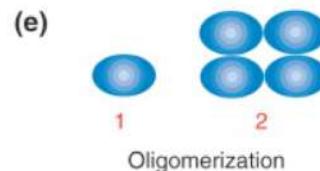
Inside and outside the cell



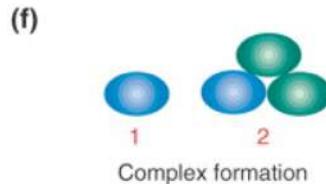
Expression by different cell types



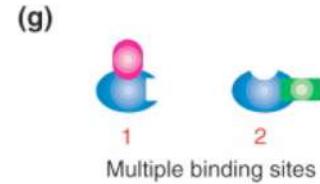
Binding of a cofactor



Oligomerization



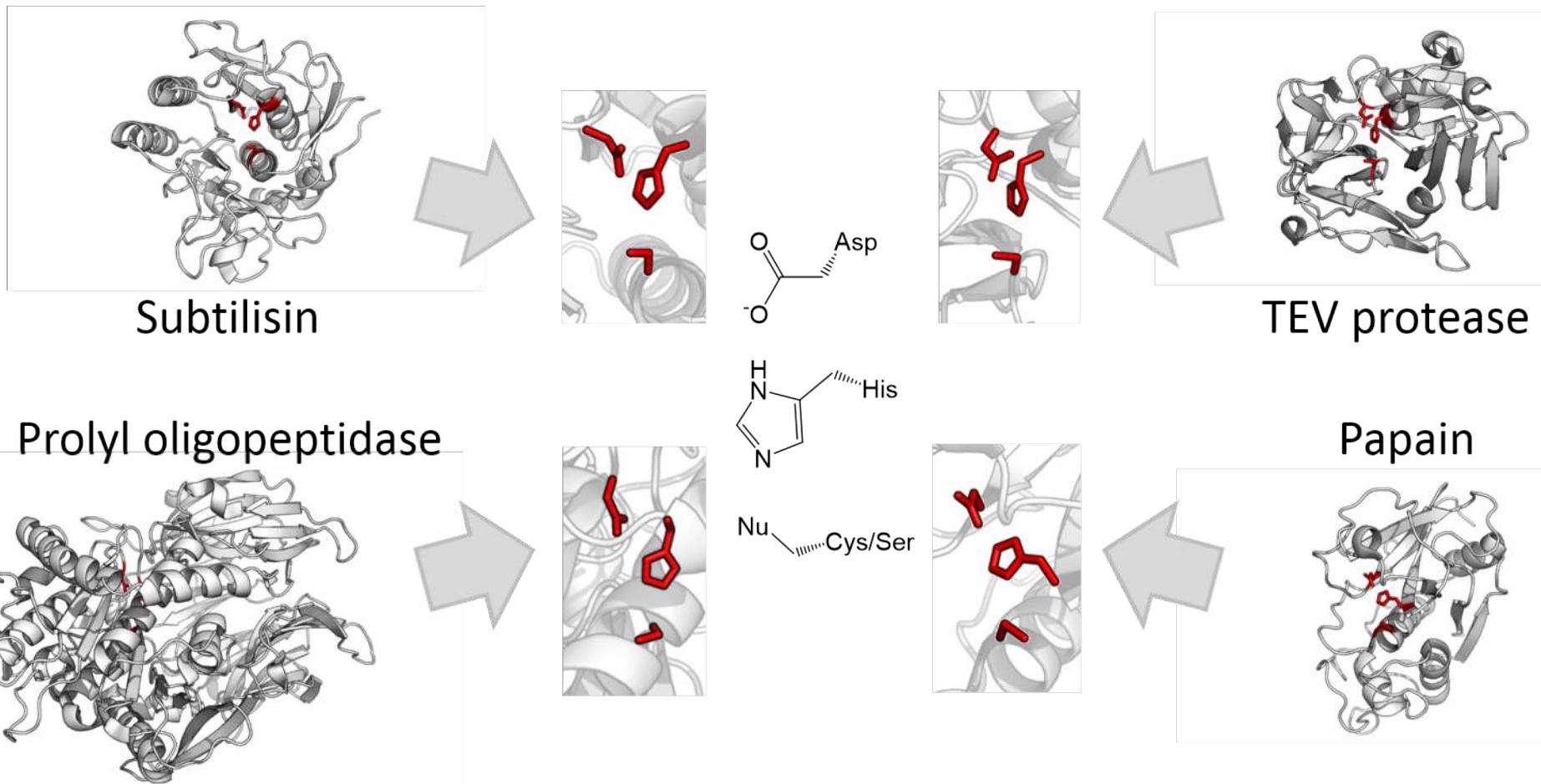
Complex formation



Multiple binding sites

Still...

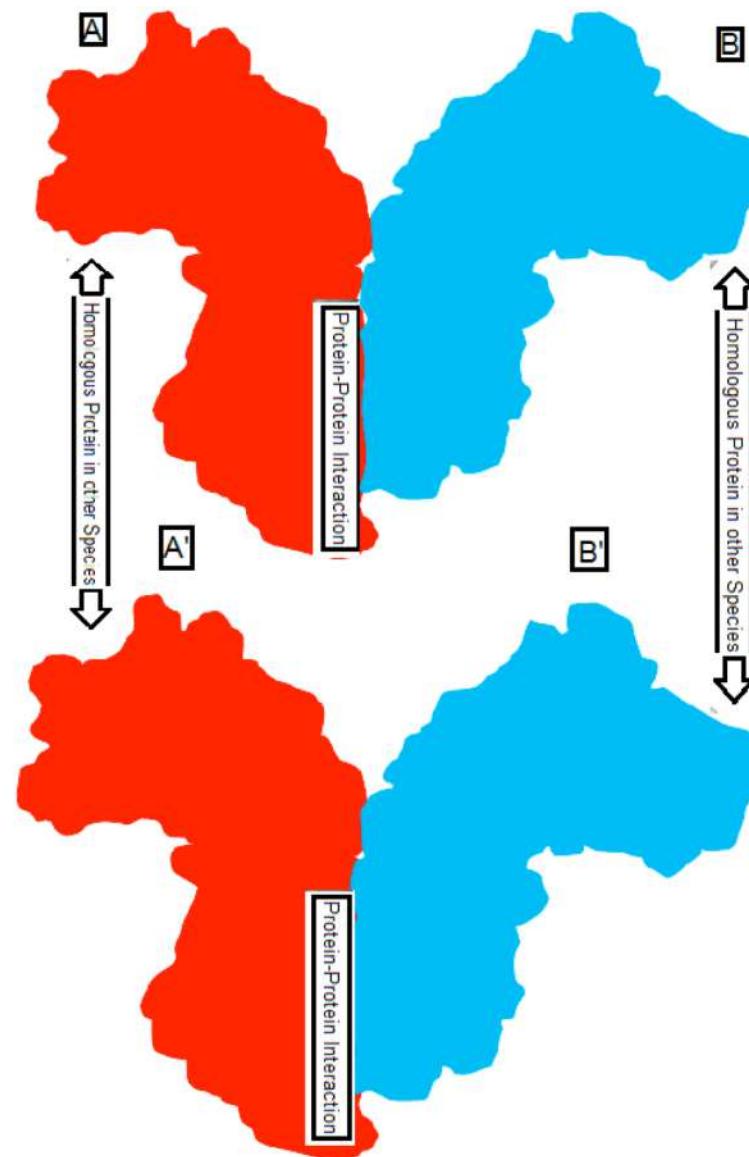
“Detecting sequence similarity in order to uncover homologous relationships between proteins remains the single most powerful tool for function prediction” Ochoa....Singh PLOS CB (2015)



"Triad Convergence" by Thomas Shafee - Own work. Licensed under CC BY-SA 3.0 via Wikimedia Commons - https://commons.wikimedia.org/wiki/File:Triad_Convergence.png#/media/File:Triad_Convergence.png

How about homology and
protein-protein interactions?

Interologs



Interologs

What Evidence Is There for the Homology of Protein-Protein Interactions?

Anna C. F. Lewis, Nick S. Jones, Mason A. Porter, Charlotte M. Deane 

Published: September 20, 2012 • <http://dx.doi.org/10.1371/journal.pcbi.1002645>

Article	Authors	Metrics	Comments	Related Content
▼				

Abstract

Author Summary

Introduction

Results/Discussion

Materials and Methods

Supporting Information

Acknowledgments

Author Contributions

References

Reader Comments (1)

Media Coverage (0)

Figures

Abstract

The notion that sequence homology implies functional similarity underlies much of computational biology. In the case of protein-protein interactions, an interaction can be inferred between two proteins on the basis that sequence-similar proteins have been observed to interact. The use of transferred interactions is common, but the legitimacy of such inferred interactions is not clear. Here we investigate transferred interactions and whether data incompleteness explains the lack of evidence found for them. Using definitions of homology associated with functional annotation transfer, we estimate that conservation rates of interactions are low even after taking interactome incompleteness into account. For example, at a blastp *E*-value threshold of 10^{-70} , we estimate the conservation rate to be about 11% between *S. cerevisiae* and *H. sapiens*. Our method also produces estimates of interactome sizes (which are similar to those previously proposed). Using our estimates of interaction conservation we estimate the rate at which protein-protein interactions are lost across species. To our knowledge, this is the first such study based on large-scale data. Previous work has suggested that interactions transferred within species are more reliable than interactions transferred across species. By controlling for factors that are specific to within-species interaction prediction, we propose that the transfer of interactions within species might be less reliable than transfers between species. Protein-protein interactions appear to be very rarely conserved unless very high sequence similarity is observed. Consequently, inferred interactions should be used with care.

Rpb4-Rpb7 complex crystallized in both *H. sapiens* (pdb code:2c35) and *S. cerevisiae* (pdb code:1y14).

Control interactive selection and zoom in the structure panel

Basic swap between positions 31 and 35

Residues Focus: R31A, N35A, E41A, F31B, E35B, A178

Invariant position at E35

Multiple sequence Alignment:

gi/Index	R31A	N35A	E41A	F31B	E35B	A178
Anopheles_gambiae	H	-	-	-	-	-
Aedes_aegypti	H	-	-	-	-	-
Drosophila_melanogaster	H	-	-	-	-	-
Ixodes_scapularis	H	-	-	-	-	-
Acyrtosiphon_pisum	H	-	-	-	-	-
Schistosoma_mansoni	-	-	-	-	-	-
Trichoplax_adhaerens	-	-	-	-	-	-
Loa_loa	S	-	-	-	-	-
Dictyostelium_discoidium	S	-	-	-	-	-
Malassezia_globosa	S	-	-	-	-	-
Ustilago_maydis	S	-	-	-	-	-
Laccaria_bicolor	S	-	-	-	-	-
Schizophyllum_comune	-	-	-	-	-	-
Schizosaccharomyces_pombe	-	-	-	-	-	-
Schizosaccharomyces_japonicus	-	-	-	-	-	-
Tuber_melanoporum	-	-	-	-	-	-
Magnaporthe_grisea	-	-	-	-	-	-
Uncinocarpus_reesii	-	-	-	-	-	-
Coccidioides_posadasii	-	-	-	-	-	-
Trichophyton_verrucosum	-	-	-	-	-	-
Aspergillus_nidulans	-	-	-	-	-	-
Aspergillus_fumigatus	-	-	-	-	-	-
Talaromyces_stipitatus	-	-	-	-	-	-
Verticillium_albo-atrum	-	-	-	-	-	-
Gibberella_zaeae	-	-	-	-	-	-
Hectria_haematoceca	-	-	-	-	-	-
Botryotinia fuckeliana	-	-	-	-	-	-
Podospora_anserina	-	-	-	-	-	-
Neurospora_crassa	-	-	-	-	-	-

Chain A (2c35)

Chain B (2c35)

Chain A (2c35)
Chain A (1y14)

Chain B (2c35)
Chain B (1y14)

OPEN

Assessment of the key regulatory genes and their Interologs for Turner Syndrome employing network approach

Received: 13 March 2018

Accepted: 15 June 2018

Published online: 04 July 2018

Anam Farooqui, Safia Tazyeen, Mohd. Murshad Ahmed, Aftab Alam, Shahnawaz Ali^{ID}, Md. Zubair Malik, Sher Ali & Romana Ishrat

Turner Syndrome (TS) is a condition where several genes are affected but the molecular mechanism remains unknown. Identifying the genes that regulate the TS network is one of the main challenges in understanding its aetiology. Here, we studied the regulatory network from manually curated genes reported in the literature and identified essential proteins involved in TS. The power-law distribution analysis showed that TS network carries scale-free hierarchical fractal attributes. This organization of the network maintained the self-ruled constitution of nodes at various levels without having centrality-lethality control systems. Out of twenty-seven genes culminating into leading hubs in the network, we identified two key regulators (KRs) i.e. KDM6A and BDNF. These KRs serve as the backbone for all the

INTEROLOG FINDER

Start New Analysis Page

navigation

Pages

[Start New Analysis](#)

[Downloads](#)

[About Us](#)

[Help and FAQ](#)

Protein or Proteins of interest:

Paste identifiers, comma separated or list format:

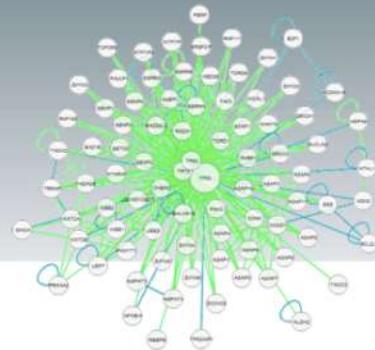
NCBI IDs are preferred, but gene names, Ensembl IDs,
and several other identifiers are translated
example: 672, TP53, ENSG00000107331

Species:

- Homo sapiens*
- Mus musculus*
- Drosophila melanogaster*
- Caenorhabditis elegans*
- Saccharomyces cerevisiae*

[Get interologs](#)

On the next page results and possible synonyms to your input proteins will be displayed; chose the correct ID from the following synonym list and click extend.
If you need further help with identifiers, please visit [Ensembl](#) or [NCBI](#). Click on the button to add selected genes.



Home

Home

Sample 1

Sample 2

Help

Contact us

Protein-Protein Interaction Search

Input an interacting protein pair as a query to search its homologous interactions across multiple species

Press the ? to obtain more information on that specific field.

Query protein pair (sequences in FASTA format or [UniProt ID](#)) :

Input sequences in FASTA format

Interacting partner 1:

```
>sp|P61967|AP1S1_MOUSE  
MMRFMLLFSRQGKRLQKWLATSDKERKKMVRLEMQVVLARKPKMCSFLEWRDLKVYYK  
RYASLYFCCAIEGQDNELTITLELIHYVELLDKYFGSVCELDIIFNFEKAYFILDEFMLMG  
GDVQDTSKKSVLKAIEQADLLQEEDESPRSVLEEMGLA
```

Interacting partner 2:

```
>sp|P22892|AP1G1_MOUSE  
MPAPIRLRELIRLIRTARTQAEEREMIQKECAAIRSSFREEDNTYRCRNVAKLLYMHMLG  
YPAHFGQLECLKLIAQSQKFTDKRIGYLGMALLDERQDVHLLMTNCIKNDLNHSTQFVQG  
LALCTLGCMGSSEMCRDLAGEVEKLLKTSNSYLRKKAALCAHVIRKPELMEMFLPATK  
NLLNEKNHGVLHTSVVLLTEMCERSPDMLAHFRKLVPQLVRILKNLIMSGYSPEHDVSGI
```

Input UniProt ID (Ex: AP1S1_MOUSE)

Interacting partner 1: AP1S1_MOUSE

Interacting partner 2: AP1G1_MOUSE

[Search](#)

[Clear](#)

Options:

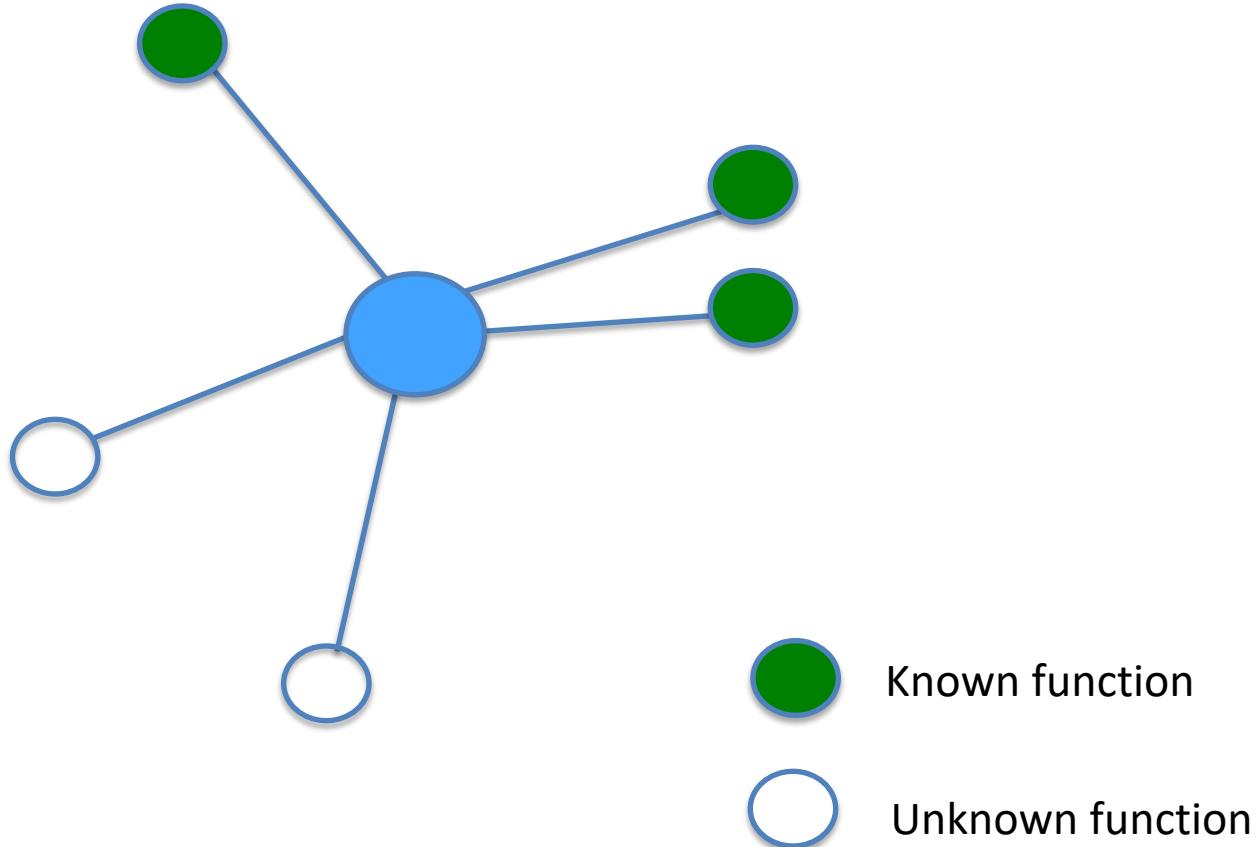
E-value cut-off threshold for homolog searching [?](#)

10 10^{-1} 10^{-10} (Default) Other: (Ex: -50 = 10^{-50})

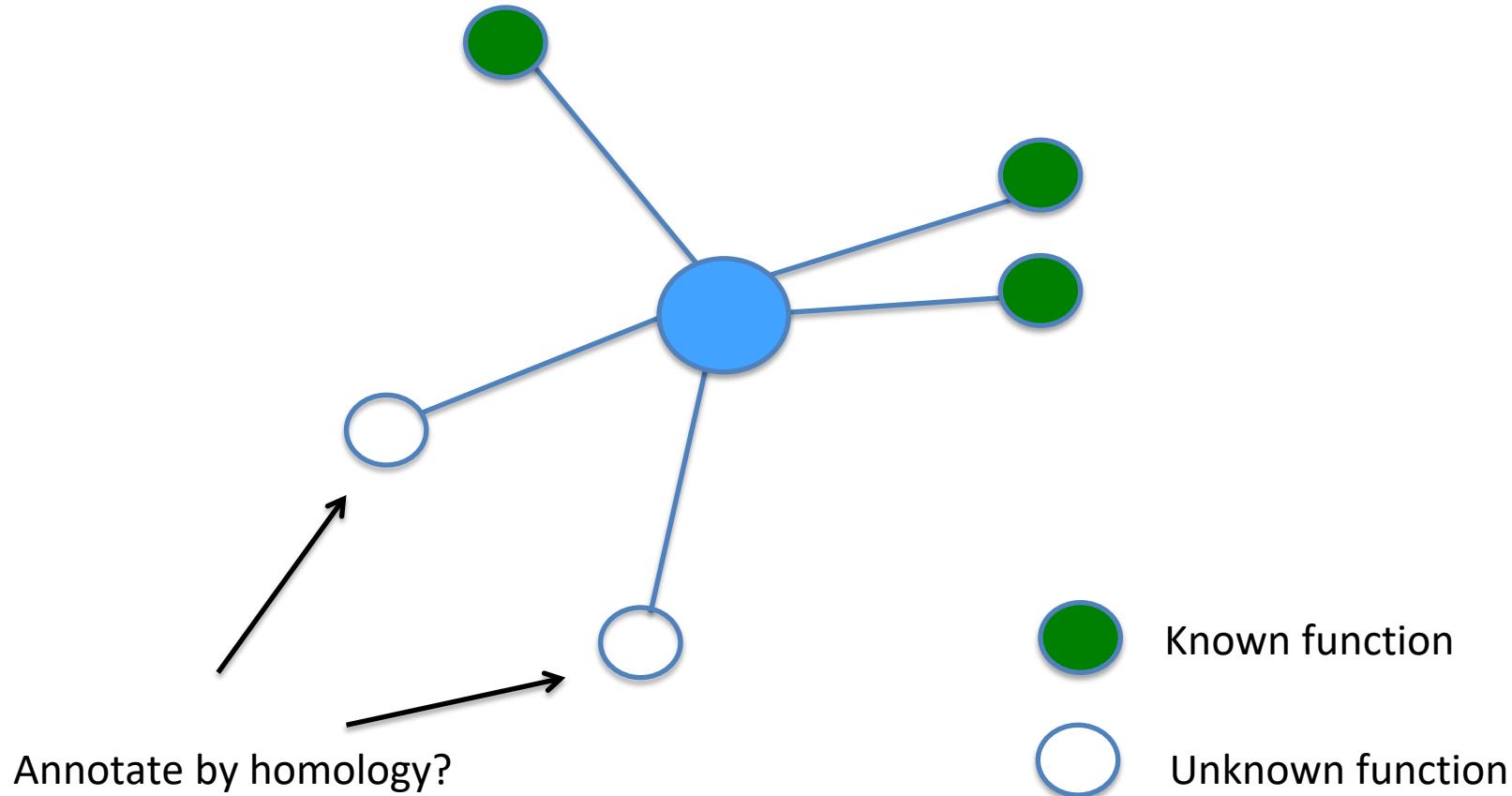
Joint E-value [?](#)

10^{-100} 10^{-40} (Default) 10^{-10} Other: (Ex: -50 = 10^{-50})

Homology for annotating interactors in known P-P interaction networks



Homology for annotating interactors in known P-P interaction networks



Here, we will focus on homology detection via
excess sequence similarity

Exercise

We are going to align two protein sequences.

The first one is:

P29973 | CNGA1_HUMAN cGMP-gated cation channel
alpha-1

the second one is a mystery protein...

Let's visit the BLAST website:

https://blast.ncbi.nlm.nih.gov/Blast.cgi?PROGRAM=blastp&PAGE_TYPE=BlastSearch&LINK_LOC=blasthome

See file: blast_pairwise_link.txt



BLAST® » blastp suite

Align Sequences Protein BLAST

[blastn](#) [blastp](#) [blastx](#) [tblastn](#) [tblastx](#)

Enter Query Sequence

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

Or, upload file

[Choose File](#) no file selected [?](#)

Job Title

Enter a descriptive title for your BLAST search [?](#) Align two or more sequences [?](#)

Enter Subject Sequence

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

Or, upload file

[Choose File](#) no file selected [?](#)From To

Program Selection

Algorithm

 blastp (protein-protein BLAST)Choose a BLAST algorithm [?](#)**BLAST**

Search protein sequence using Blastp (protein-protein BLAST)

 Show results in a new window[+ Algorithm parameters](#)

BLAST® » blastp suite

Align Sequences Protein BLAST

blastn blastp blastx tblastn tblastx

Enter Query Sequence

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

Enter accession number(s), gi(s), or FASTA sequence(s) [?](#)

Clear

Query subrange [?](#)

Upload file **mystery_protein.fasta**

Click



Choose File mystery_protein.fasta [?](#)

Title

Enter a descriptive title for your BLAST search [?](#)

Align two or more sequences [?](#)

Enter Subject Sequence

Enter accession number(s), gi(s), or FASTA sequence(s) [?](#)

Clear

Subject subrange [?](#)

Upload file **CNGA1_human.fasta**

Click



Choose File CNGA1_human.fasta [?](#)

Program Selection

Algorithm

blastp (protein-protein BLAST)

[Choose a BLAST algorithm](#) [?](#)

BLAST

Search protein sequence using **Blastp (protein-protein BLAST)**

Show results in a new window

[Algorithm parameters](#)



BLAST® » blastp suite

Align Sequences Protein BLAST

[blastn](#) [blastp](#) [blastx](#) [tblastn](#) [tblastx](#)

Enter Query Sequence

BLASTP programs search protein subjects using a protein query. [more...](#)Enter accession number(s), gi(s), or FASTA sequence(s) [?](#)[Clear](#)Query subrange [?](#)

```
>mystery_protein
MGNGSVKPKHSKHPDGHSGNLTTDALRNKVTELERELRRKDAEIQEREYHLKELREQLSK
QTVAIAEELTEELQNKCQLNKLQDVHMQGGSPHQASPDVKPLEVHRKTSGLVLHSRRG
AKAGVSAEPTTRTYDLNKPPEFSFEKARVRKDSSEKKLTDALNKQFLKRLDPQQIKDM
VECMYGRNYQQGSYIJKQGEPGNHIFVLAEGRLEVWKLSSIPMWTTFGELAILYNC
```

Or, upload file [Choose File](#) no file selected

Job Title

 mystery_proteinEnter a descriptive title for your BLAST search [?](#) Align two or more sequences [?](#)**or Copy and Paste Sequences**

Enter Subject Sequence

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

```
>sp|P29973|CNGA1_HUMAN cGMP-gated cation channel alpha-1 OS=Homo sapiens
OX=9606 GN=CNGA1 PE=1 SV=3
MKLSMKNNIINTQQSFVTMPNVIVPDIEKEIRRMEENGACSSFSEDDDSASTSEESENNP
HARGSF SYKSLRKGGPSQREQYLPGAIALFVNNSNNKDQEPEEKKKKKKEKKSKSDDKN
ENKNDPEKKKKKKDKKEKKKEEKS KDKEEEKKEVVVIDPSGN TYYNWLF CILPV MYNW
```

Or, upload file [Choose File](#) no file selected

Program Selection

Algorithm

 blastp (protein-protein BLAST)[Choose a BLAST algorithm](#) [?](#)**BLAST**Search protein sequence using **Blastp (protein-protein BLAST)** Show results in a new window[Algorithm parameters](#)



BLAST® » blastp suite

Align Sequences Protein BLAST

[blastn](#) [blastp](#) [blastx](#) [tblastn](#) [tblastx](#)

Enter Query Sequence

BLASTP programs search protein subjects using a protein query. [more...](#)Enter accession number(s), gi(s), or FASTA sequence(s) [?](#)[Clear](#)

```
>mystery_protein
MGNGSVKPKHSKHPDGHSGNLTTDALRNKVTELERELRRKDAEIQEREYHLKELREQLSK
QTVAIAEELTEELQNKCQLNKLQDVHMQGGSPHQASPDVKPLEVHRKTSGLVLHSRRG
AKAGVSAEPTTRTYDLNKPPEFSFEKARVRKDSSEKKLTDALKNQFLKRLDPQQIKDM
VECMYGRNYQQGSYIJKQGEPGNHIFVLAEGRLEVFGKEKLSSIPMWTTFGEYLINYC
```

Or, upload file

[Choose File](#) no file selected [?](#)

Job Title

mystery_protein

Enter a descriptive title for your BLAST search [?](#) Align two or more sequences [?](#)Query subrange [?](#)

From

To

Enter Subject Sequence

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

```
>sp|P29973|CNGA1_HUMAN cGMP-gated cation channel alpha-1 OS=Homo sapiens
OX=9606 GN=CNGA1 PE=1 SV=3
MKLSMKNNIINTQQSFVTMPNVIVPDIEKEIRRMEENGACSSFSEDDDSASTSEESENENP
HARGSFSYKSLRKGGPSQREQYLPGAIALFNVNNSSNKDQEPEEKKKKKKEKKSKSDDKN
ENKNDPEKKKKKKDKKEKKKEEKSKDKEEEKKEVVVIDPSGNTYYNWLFCTLPVMYNW
```

Or, upload file

[Choose File](#) no file selected [?](#)Subject subrange [?](#)

From

To

Program Selection

Algorithm

 blastp (protein-protein BLAST)[Choose a BLAST algorithm](#) [?](#)[BLAST](#)Search protein sequence [?](#) Show results in a new window**Then BLAST!**[Algorithm parameters](#)

BLAST Results

Edit and Resubmit Save Search Strategies > Formatting options > Download

YouTube How to read this page Blast report des...

Job title: mystery_protein (762 letters)

RID XXGSCDW8114 (Expires on 11-05 16:59 pm)

Query ID Icl|Query_208167

Description mystery_protein

Molecule type amino acid

Query Length 762

Subject ID Icl|Query_208169

Description sp|P29973|CNGA1_HUMAN cGMP-gated cation channel alpha-1 OS=Homo sapiens OX=9606 GN=CNGA1 PE=1 SV=3

Molecule type amino acid

Subject Length 690

Program BLASTP 2.8.1+ > Citation

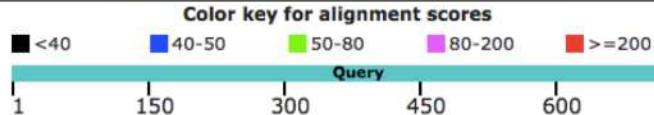
Other reports: > Search Summary [Multiple alignment] [MSA viewer]

New Analyze your query with SmartBLAST

Graphic Summary

Distribution of the top 5 Blast Hits on 1 subject sequences ⓘ

Mouse over to see the title, click to show alignments



Dot Matrix View



Descriptions

Sequences producing significant alignments:

Select: All None Selected:0

All Alignments Download Graphics Multiple alignment



	Description	Max score	Total score	Query cover	E value	Ident	Accession
<input type="checkbox"/>	sp P29973 CNGA1_HUMAN cGMP-gated cation channel alpha-1 OS=Homo sapiens OX=9606 GN=CNGA1 PE=1 SV=3	41.6	136	35%	1e-07	27%	Query_208169

BLAST Results

Edit and Resubmit Save Search Strategies > Formatting options > Download

YouTube How to read this page Blast report des...

Job title: mystery_protein (762 letters)

RID XXGSCDW8114 (Expires on 11-05 16:59 pm)

Query ID Icl|Query_208167

Description mystery_protein

Molecule type amino acid

Query Length 762

Subject ID Icl|Query_208169

Description sp|P29973|CNGA1_HUMAN cGMP-gated cation channel alpha-1 OS=Homo sapiens OX=9606 GN=CNGA1 PE=1 SV=3

Molecule type amino acid

Subject Length 690

Program BLASTP 2.8.1+ > Citation

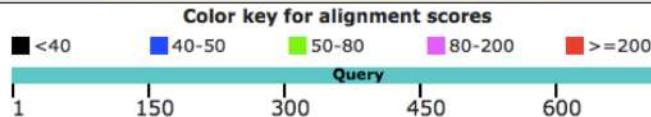
Other reports: > Search Summary [Multiple alignment] [MSA viewer]

New Analyze your query with SmartBLAST

Graphic Summary

Distribution of the top 5 Blast Hits on 1 subject sequences ⓘ

Mouse over to see the title, click to show alignments



Dot Matrix View

Descriptions

Sequences producing significant alignments:

Select: All None Selected:0

All Alignments Download Graphics Multiple alignment



	Description	Max score	Total score	Query cover	E value	Ident	Accession
<input type="checkbox"/>	sp P29973 CNGA1_HUMAN cGMP-gated cation channel alpha-1 OS=Homo sapiens OX=9606 GN=CNGA1 PE=1 SV=3	41.6	138	65%	1e-07	27%	Query_208169

41.6 138 65% 1e-07 27% Query_208169

E-value=10⁻⁷

Is our mystery protein a cGMP-gated cation channel?

Alignments

[Download](#) [Graphics](#) Sort

sp|P29973|CNGA1_HUMAN cGMP-gated cation channel alpha-1 OS=Homo sapiens OX=9606 GN=CNGA1 PE=1 SV=3
Sequence ID: Query_89161 Length: 690 Number of Matches: 5

Range 1: 474 to 579 [Graphics](#)

▼ Next Match ▲ Previous Match

Score	Expect	Method	Identities	Positives	Gaps
41.6 bits(96)	1e-07	Compositional matrix adjust.	30/110(27%)	54/110(49%)	11/110(10%)
Query 281	LRSVSSLKKNLPEDKLTKIIDCLEVEYYDKGDYIIREGEEGSTFFILAKGVKVQSTEGLH	340			
	L+ V + + L + + Y GDYI ++G+ G +I+ +GK+ V				
Sbjct 474	LKKVRIFADCEAGLLVELVLKLQPQVYSPGDYICKKGDIGREMYIIKEGKLAVV---AD	529			
Query 341	DQPQLIKTLQKGEYFGEKALI-----SDDVRSANIIA-EENDVACLVID	383			
	D L G YFGE +++ + + R+ANI + +D+ CL D				
Sbjct 530	DGVTQFVVLSDGSYFGEISILNIKGSKAGNRRTANIKSIGYSDLFCLSKD	579			



E-value=10⁻⁷

Range 2: 472 to 579 [Graphics](#)

▼ Next Match ▲ Previous Match

▲ First Match

Score	Expect	Method	Identities	Positives	Gaps
37.7 bits(86)	1e-06	Compositional matrix adjust.	26/108(24%)	53/108(49%)	8/108(7%)
Query 161	DALNKNQFLKRLLDPQQIKDMVECMYGRNYQQGSYIJKQGEPEGNHIFVLAEGRLEVQGEK	220			
	D L K + + + +V + + Y G YI K+G+ G +++++ EG+L V +				
Sbjct 472	DTLKKVRIFADCEAGLLVELVLKLQPQVYSPGDYICKKGDIGREMYIIKEGKLAVVADDG	531			
Query 221	LLSSIPM--WTFGELAIL-----YNCTRTASVKAITNVKTWALDR	260			
	+ + + + FGE++IL RTA++K+I + L ++				
Sbjct 532	VTQFVVLSDGSYFGEISILNIKGSKAGNRRTANIKSIGYSDLFCLSKD	579			



E-value=10⁻⁶

Range 3: 317 to 356 [Graphics](#)

▼ Next Match ▲ Previous Match

▲ First Match

Score	Expect	Method	Identities	Positives	Gaps
21.9 bits(45)	0.11	Compositional matrix adjust.	14/44(32%)	23/44(52%)	5/44(11%)
Query 593	VDFGFAKKIGSGQKTWTCGTPEYVAPEV-ILNKGHDFSVDFWS	635			
	V + +K IG G TW + P+ PE L + + +S+ +WS				
Sbjct 317	VFYSISKAIGFGNDTWVY---PDINDPEFGRLARKYVYSL-YWS	356			

Range 4: 180 to 196 [Graphics](#)

▼ Next Match ▲ Previous Match

▲ First Match

Score	Expect	Method	Identities	Positives	Gaps
18.5 bits(36)	1.2	Compositional matrix adjust.	6/17(35%)	10/17(58%)	0/17(0%)
Query 539	WSILRDRGSFDEPTSKF	555			

Alignments

Download ▾ Graphics

Sort by: E value

sp|P29973|CNGA1_HUMAN cGMP-gated cation channel alpha-1 OS=Homo sapiens OX=9606 GN=CNGA1 PE=1 SV=3
Sequence ID: Query_89161 Length: 690 Number of Matches: 5

Range 1: 474 to 579 [Graphics](#)

▼ Next Match ▲ Previous Match

Score	Expect	Method	Identities	Positives	Gaps
41.6 bits(96)	1e-07	Compositional matrix adjust.	30/110(27%)	54/110(49%)	11/110(10%)

Query	281	LRSVSLKLKNLPEDKLTKIIDCLEVEYYDKGDYIIREGEESTFFILAKGKVKTQSTEGH	340
		L+ V + + L + + Y GDYI ++G+ G +I+ +GK+ V	
Sbjct	474	LKKVRIFADCEAGLLVELVLKLQPQVYSPGDYICKKGDIGREMYIIKEGKLAVV---AD	529
Query	341	DQPQLIKTLQKGEYFGEKALI-----SDDVRSANIIA-EENDVACLVID	383
		D L G YFGE +++ + + R+ANI + +D+ CL D	
Sbjct	530	DGVTQFVVLSDGSYFGEISILNIKGSKAGNRRTANIKSIGYSDLFCLSKD	579

Range 2: 472 to 579 [Graphics](#)

▼ Next Match ▲ Previous Match ▲ First Match

Score	Expect	Method	Identities	Positives	Gaps
37.7 bits(86)	1e-06	Compositional matrix adjust.	26/108(24%)	53/108(49%)	8/108(7%)

Query	161	DALKNQFLKRLDPQQIKDMVEC MYGRNYQQGSYI IKQGE PGNHIFVLAEGRLEV FQGEK	220
		D L K + + + ++V + + Y G YI K+G+ G +++++ EG+L V +	
Sbjct	472	DTLKKVRIFADCEAGLLVELVLKLQPQVYSPGDYICKKGDIGREMYIIKEGKLAVVADDG	531
Query	221	LLSSIPM--WTTFGELAIL-----YNCTRTASVKAITNVKTWALDRE	260
		+ + + + FGE++IL RTA++K+I + L ++	
Sbjct	532	VTQFVVLSDGSYFGEISILNIKGSKAGNRRTANIKSIGYSDLFCLSKD	579

Alignments

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

Sort by: E value

sp|P29973|CNGA1_HUMAN cGMP-gated cation channel alpha-1 OS=Homo sapiens OX=9606 GN=CNGA1 PE=1 SV=3
Sequence ID: Query_89161 Length: 690 Number of Matches: 5

Range 1: 474 to 579 [Graphics](#)

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

Score	Expect	Method	Identities	Positives	Gaps
41.6 bits(96)	1e-07	Compositional matrix adjust.	30/110(27%)	54/110(49%)	11/110(10%)

Query	281	LRSVSLKLKNLPEDKLTKIIDCLEVEYYDKGDYIIREGEESTFFILAKGKVKTQSTEGH	340
		L+ V + + L + + Y GDYI ++G+ G +I+ +GK+ V	
Sbjct	474	LKKVRIFADCEAGLLVELVLKLQPQVYSPGDYICKKGDIGREMYIIKEGKLAVV---AD	529
Query	341	DQPQLIKTLQKGEYFGEKALI-----SDDVRSANIIA-EENDVACLVID	383
		D L G YFGE +++ + + R+ANI + +D+ CL D	
Sbjct	530	DGVTQFVVLSDGSYFGEISILNIKGSKAGNRRTANIKSIGYSDLFCLSKD	579

Range 2: 472 to 579 [Graphics](#)

[▼ Next Match](#) [▲ Previous Match](#) [▲ First Match](#)

Score	Expect	Method	Identities	Positives	Gaps
37.7 bits(86)	1e-06	Compositional matrix adjust.	26/108(24%)	53/108(49%)	8/108(7%)

Query	161	DALKNQFLKRLDPQQIKDMVEC MYGRNYQQGSYIIKQGE PGNHIFVLAEGRLEV FQGEK	220
		D L K + + + ++V + + Y G YI K+G+ G +++++ EG+L V +	
Sbjct	472	DTLKKVRIFADCEAGLLVELVLKLQPQVYSPGDYICKKGDIGREMYIIKEGKLAVVADDG	531
Query	221	LLSSIPM--WTTFGELAIL-----YNCTRTASVKAITNVKTWALDRE	260
		+ + + + FGE++IL RTA++K+I + L ++	
Sbjct	532	VTQFVVLSDGSYFGEISILNIKGSKAGNRRTANIKSIGYSDLFCLSKD	579

cGMP-gated cation channel alpha-1

P29973 (CNGA1_HUMAN)

Marco Punta



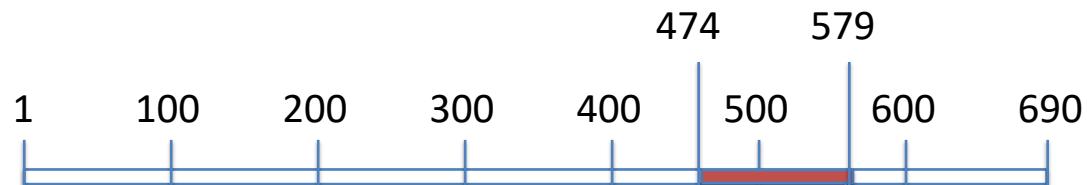
Mistery protein



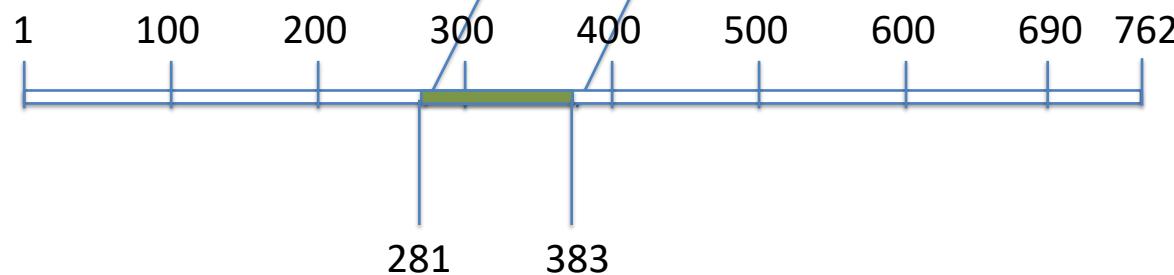
cGMP-gated cation channel alpha-1

P29973 (CNGA1_HUMAN)

Marco Punta



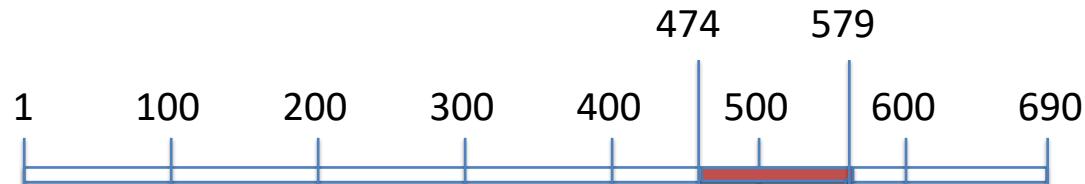
Mistery protein



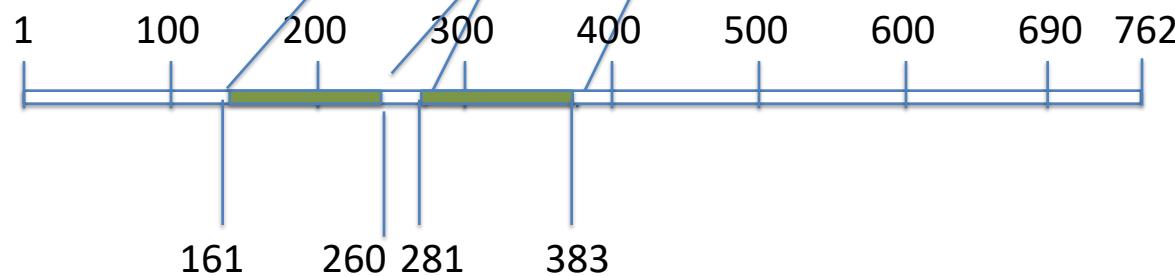
cGMP-gated cation channel alpha-1

P29973 (CNGA1_HUMAN)

Marco Punta



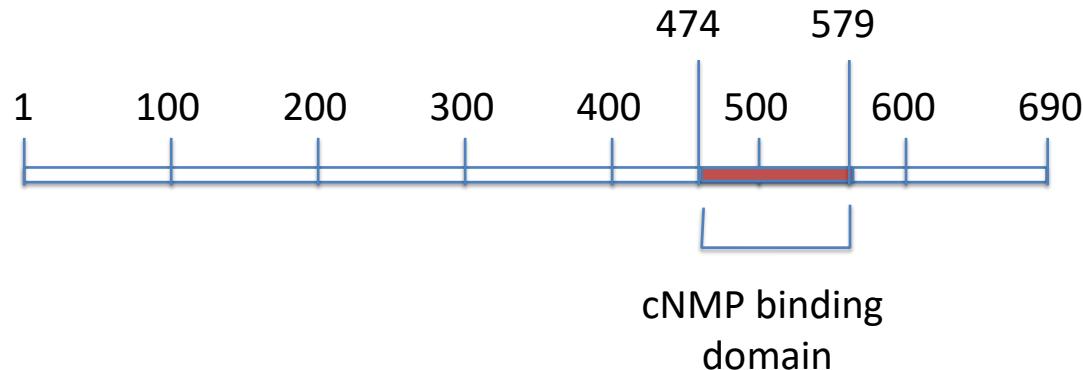
Mistery protein



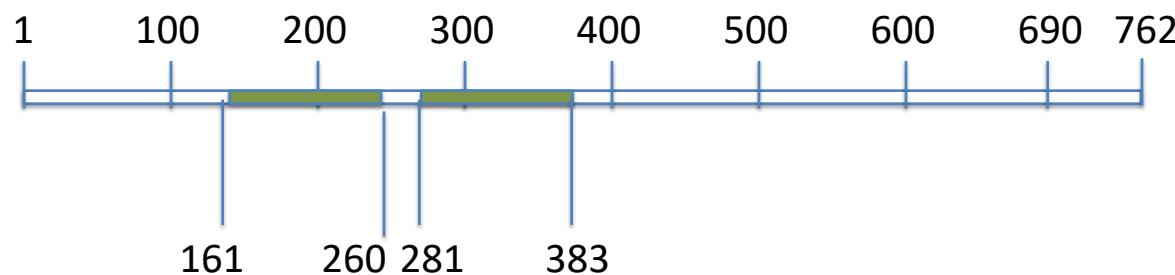
cGMP-gated cation channel alpha-1

P29973 (CNGA1_HUMAN)

Marco Punta



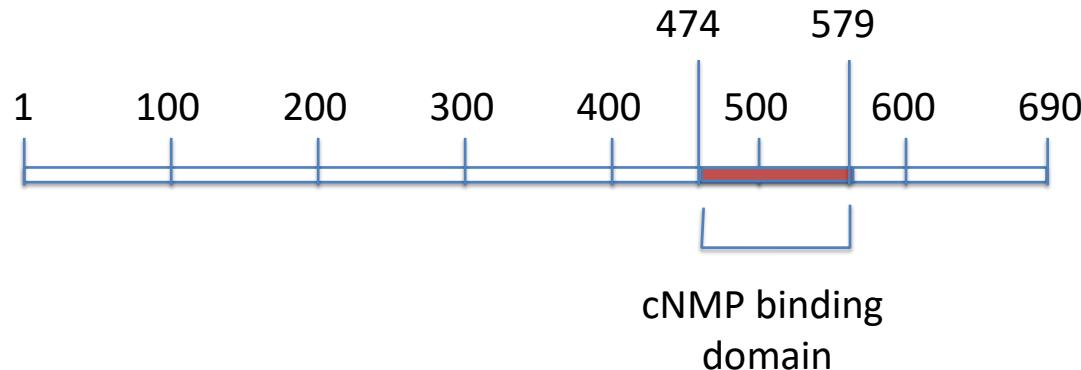
Mystery protein



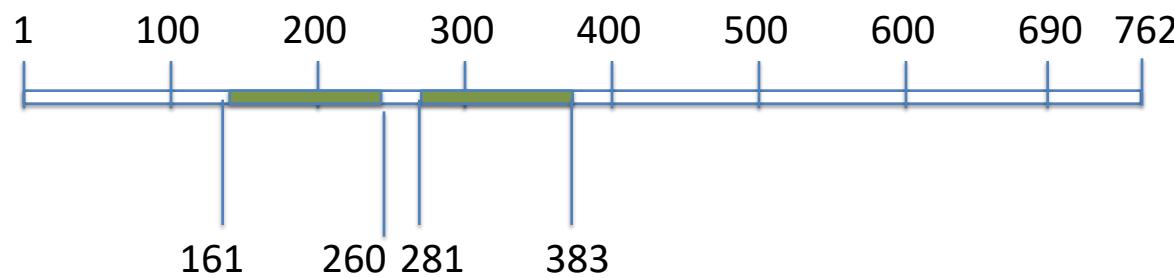
cGMP-gated cation channel alpha-1

P29973 (CNGA1_HUMAN)

Marco Punta



Mistery protein

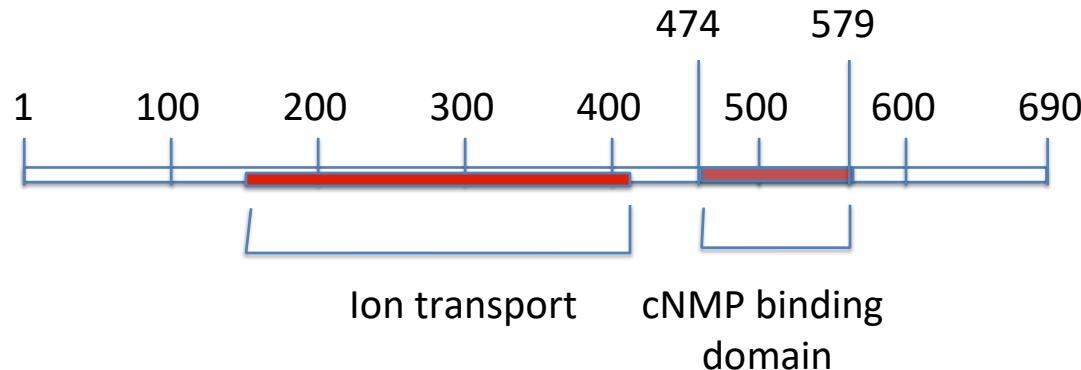


cNMP binding domains?

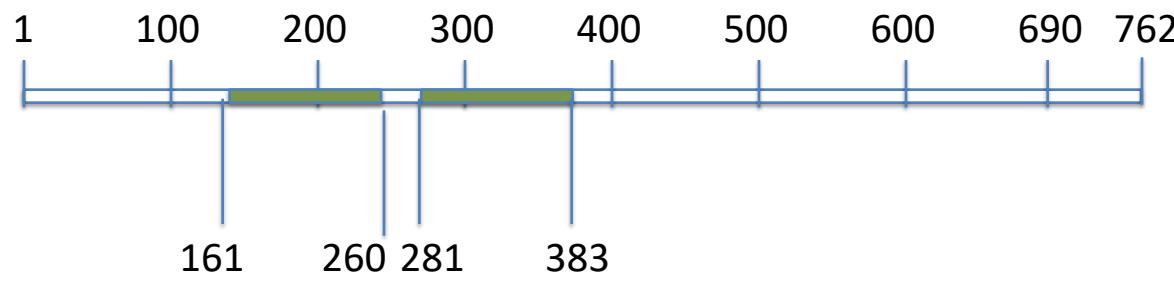
cGMP-gated cation channel alpha-1

P29973 (CNGA1_HUMAN)

Marco Punta



Mistery protein

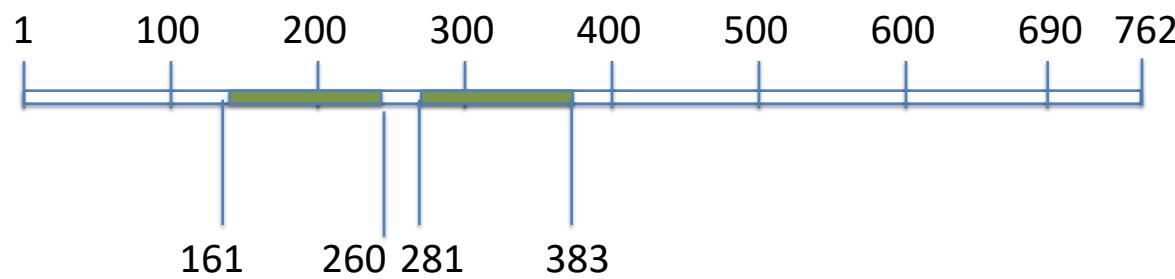
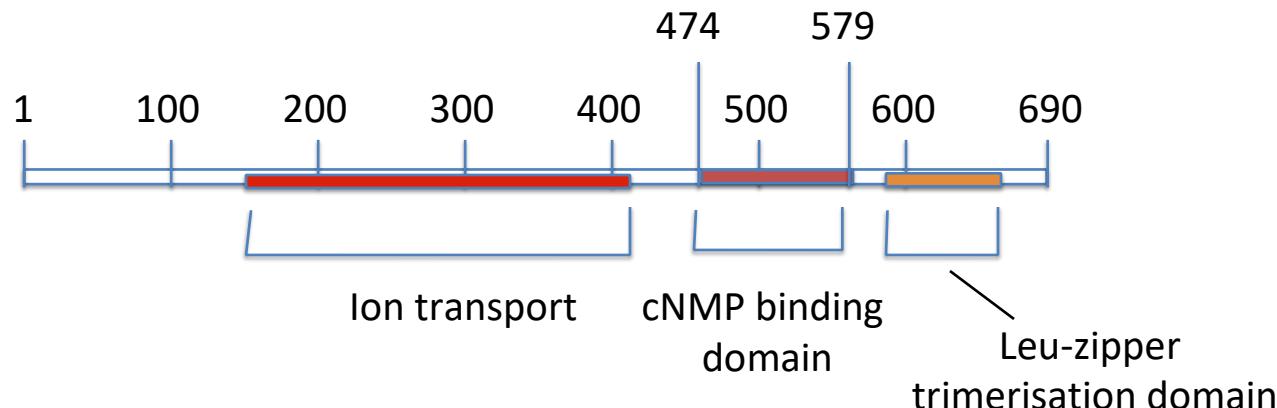


cNMP binding domains?

cGMP-gated cation channel alpha-1

P29973 (CNGA1_HUMAN)

Marco Punta
[color scale]

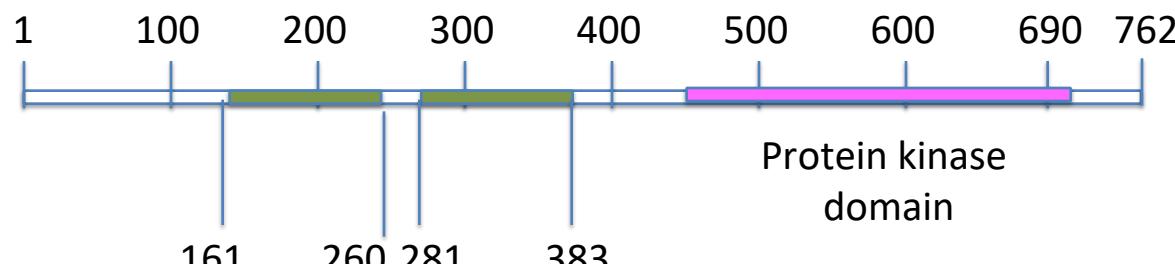
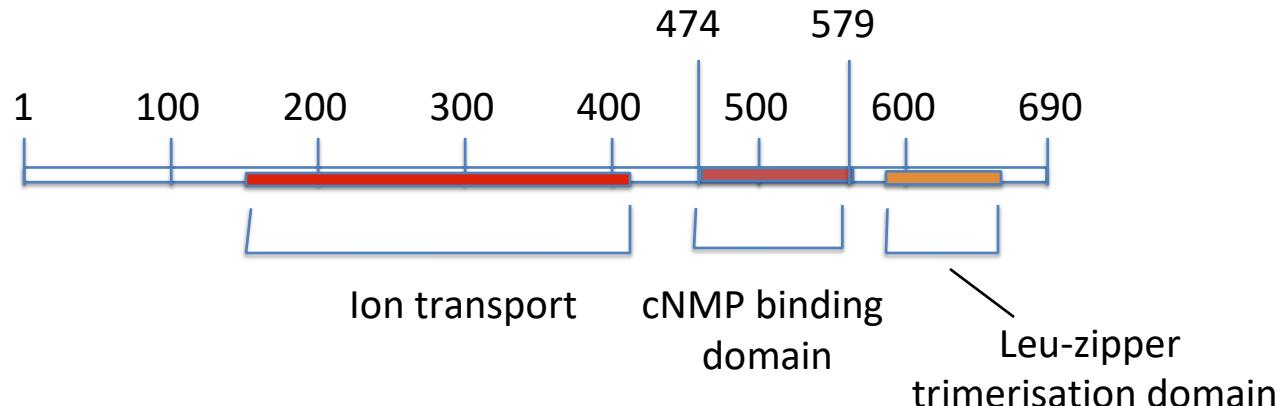


cNMP binding domains?

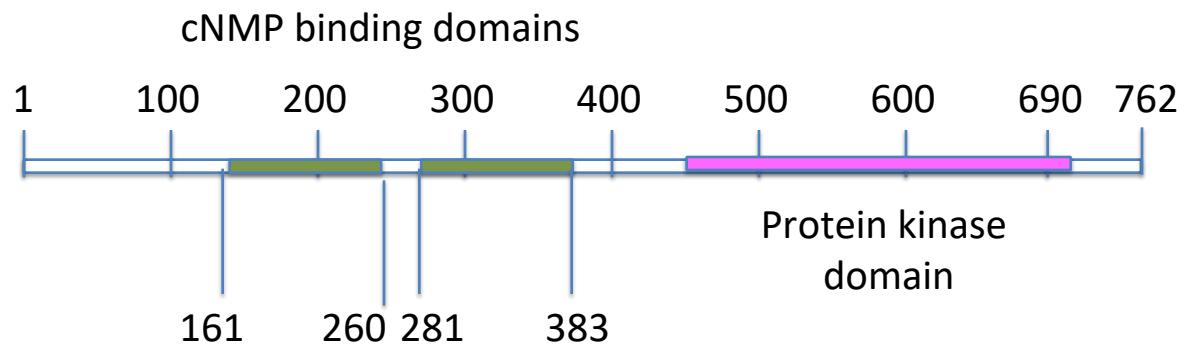
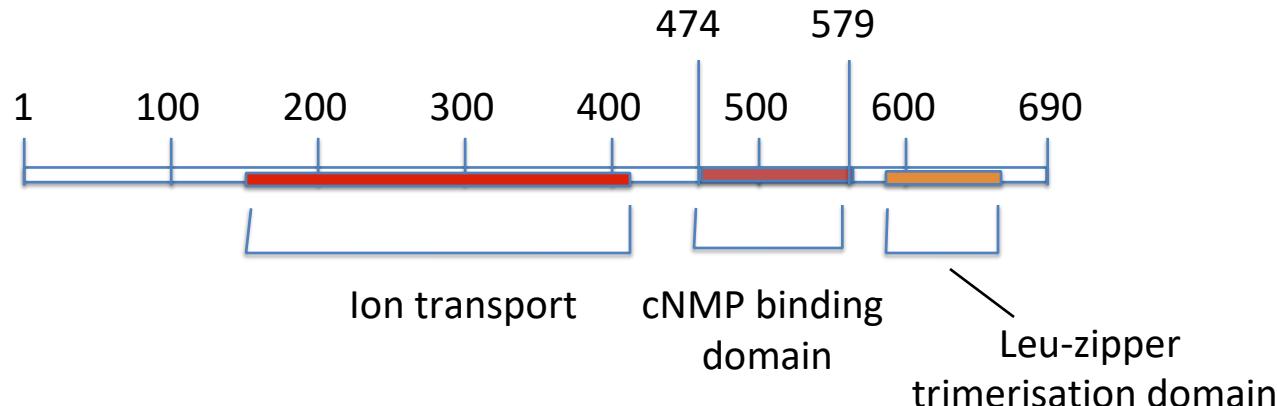
cGMP-gated cation channel alpha-1

P29973 (CNGA1_HUMAN)

Marco Punta



cNMP binding domains?



Mystery protein is a cGMP-dependent protein kinase 2
Q13237 (KGP2_HUMAN)

Definition (Wikipedia):

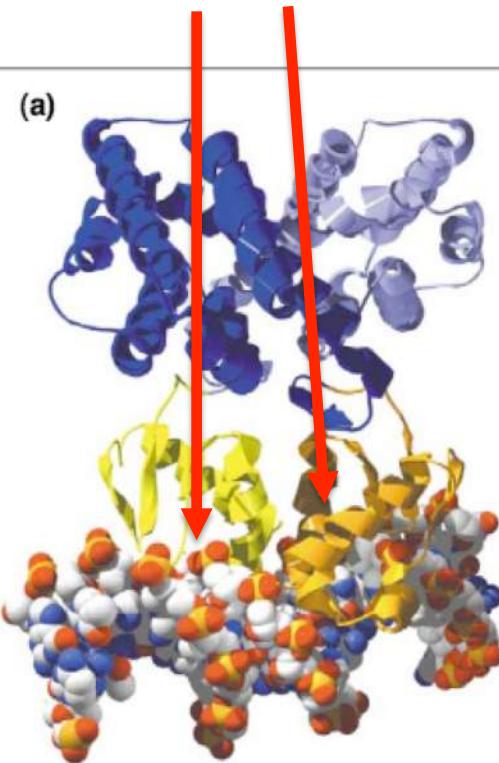
A protein domain is a conserved part of a given protein sequence and (tertiary) structure that can evolve, function, and exist independently of the rest of the protein chain. Each domain forms a compact three-dimensional structure and often can be independently stable and folded. Many proteins consist of several structural domains. One domain may appear in a variety of different proteins. Molecular evolution uses domains as building blocks and these may be recombined in different arrangements to create proteins with different functions.

Domains and function annotation

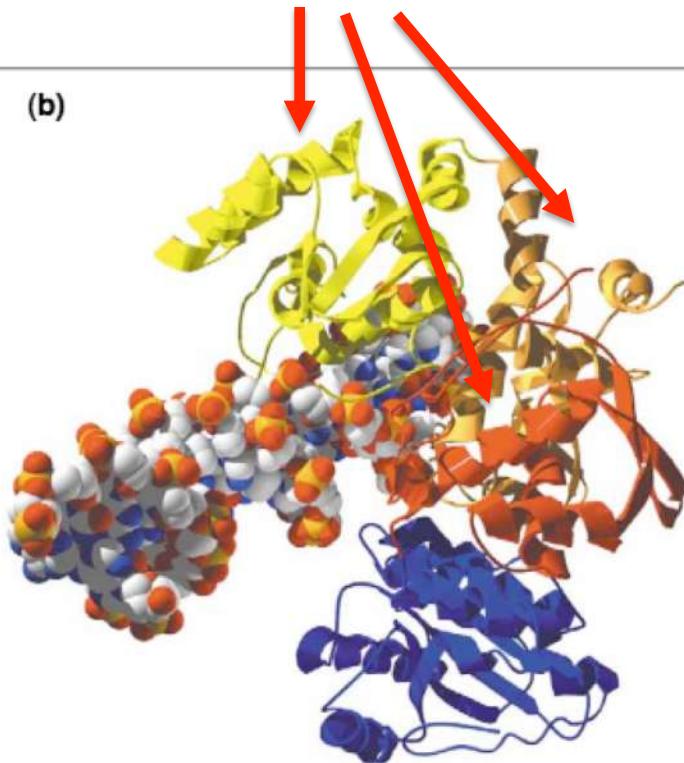
- Proteins may be homologous only in some regions (domains), this is especially true at longer evolutionary distance
- If so, function annotation transfer possible (still not safe) only between these regions

Winged helix domain

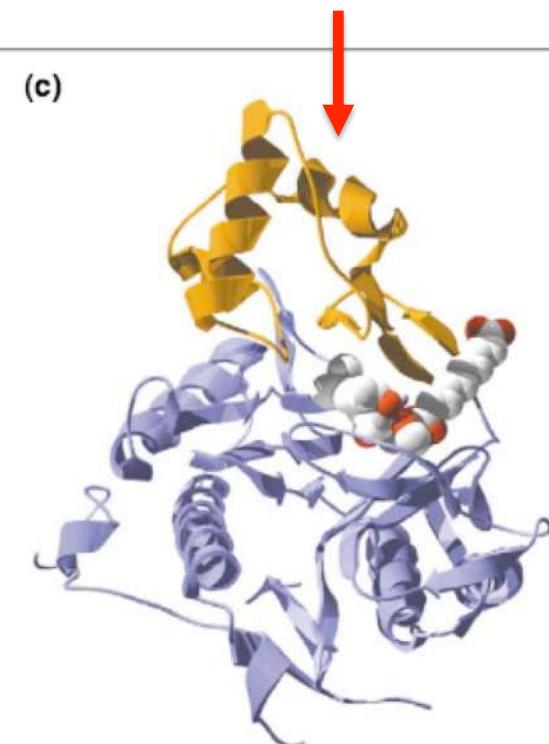
DNA binding



DNA binding



substrate specificity pocket



Current Opinion in Structural Biology

Transcription
factor

Restriction
endonuclease

Human methionine
aminopeptidase 2

Beyond just homology

RESEARCH

Open Access



An expanded evaluation of protein function prediction methods shows an improvement in accuracy

Yuxiang Jiang¹, Tal Ronnen Oron², Wyatt T. Clark³, Asma R. Bankapur⁴, Daniel D'Andrea⁵, Rosalba Lopore⁵, Christopher S. Funk⁶, Indika Kahanda⁷, Karin M. Verspoor^{8,9}, Asa Ben-Hur⁷, Da Chen Emily Koo¹⁰, Duncan Penfold-Brown^{11,12}, Dennis Shasha¹³, Noah Youngs^{12,13,14}, Richard Bonneau^{13,14,15}, Alexandra Lin¹⁶, Sayed M. E. Sahraeian¹⁷, Pier Luigi Martelli¹⁸, Giuseppe Profiti¹⁸, Rita Casadio¹⁸, Renzhi Cao¹⁹, Zhaolong Zhong¹⁹, Jianlin Cheng¹⁹, Adrian Altenhoff^{20,21}, Nives Skunca^{20,21}, Christophe Dessimoz^{22,87,88}, Tunca Dogan²³, Kai Hakala^{24,25}, Suwisa Kaewphan^{24,25,26}, Farrokh Mehryary^{24,25}, Tapio Salakoski^{24,26}, Filip Ginter²⁴, Hai Fang²⁷, Ben Smithers²⁷, Matt Oates²⁷, Julian

Next:

- Protein domain databases
- Protein profile searches
- Practical: Homology based annotation of interactors' function in P-P interaction networks