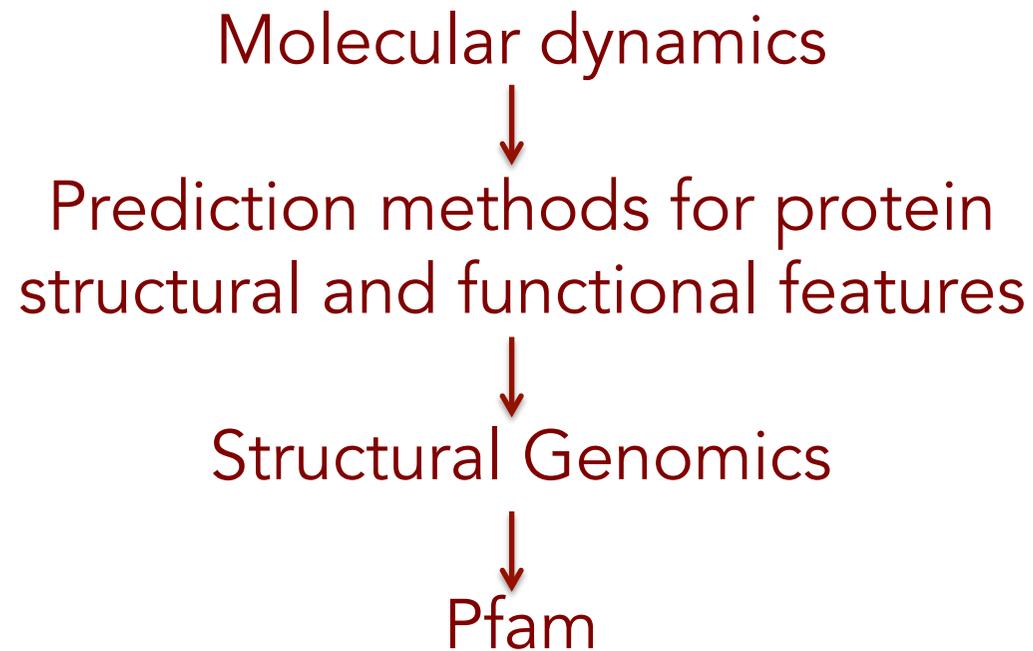


Marco Punta

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



Goals

- “Understand” annotation transfer by homology
- Know what protein family databases are and why they are useful

Outline

- Homology (definition, implications, how to detect)
- Exercise 1 and 2: homology-based function annotation transfer
- Protein domains
- Protein families
- Exercise 3: how to build a new (Pfam) protein family

Homology

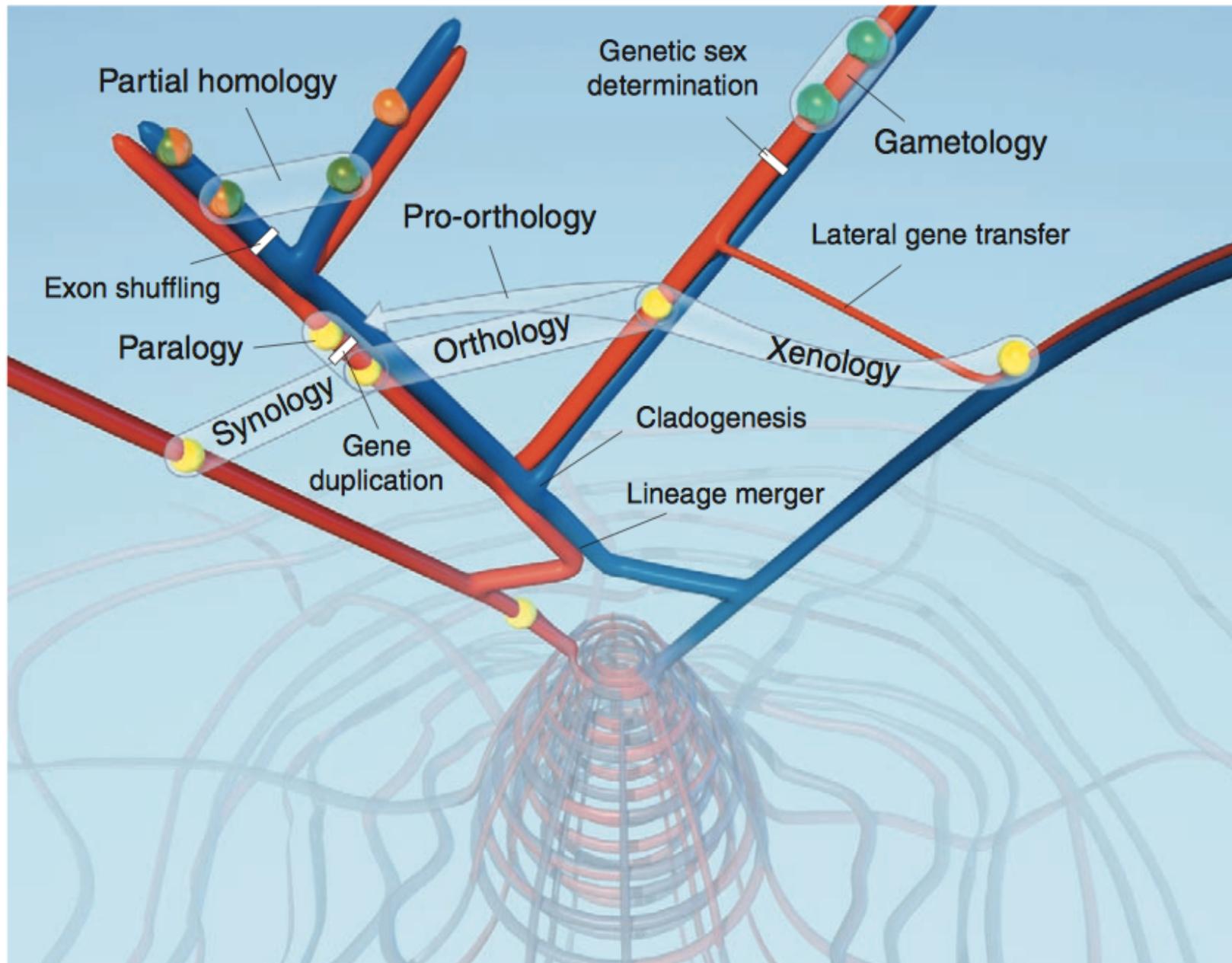
Definition:

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

Origins of homology in proteins

Origin of homology in proteins

- Speciation (orthology)
- Gene duplication (paralogy)
- Horizontal gene transfer (xenology)
- Whole genome duplication (ohnology)
- Gametology, Synology



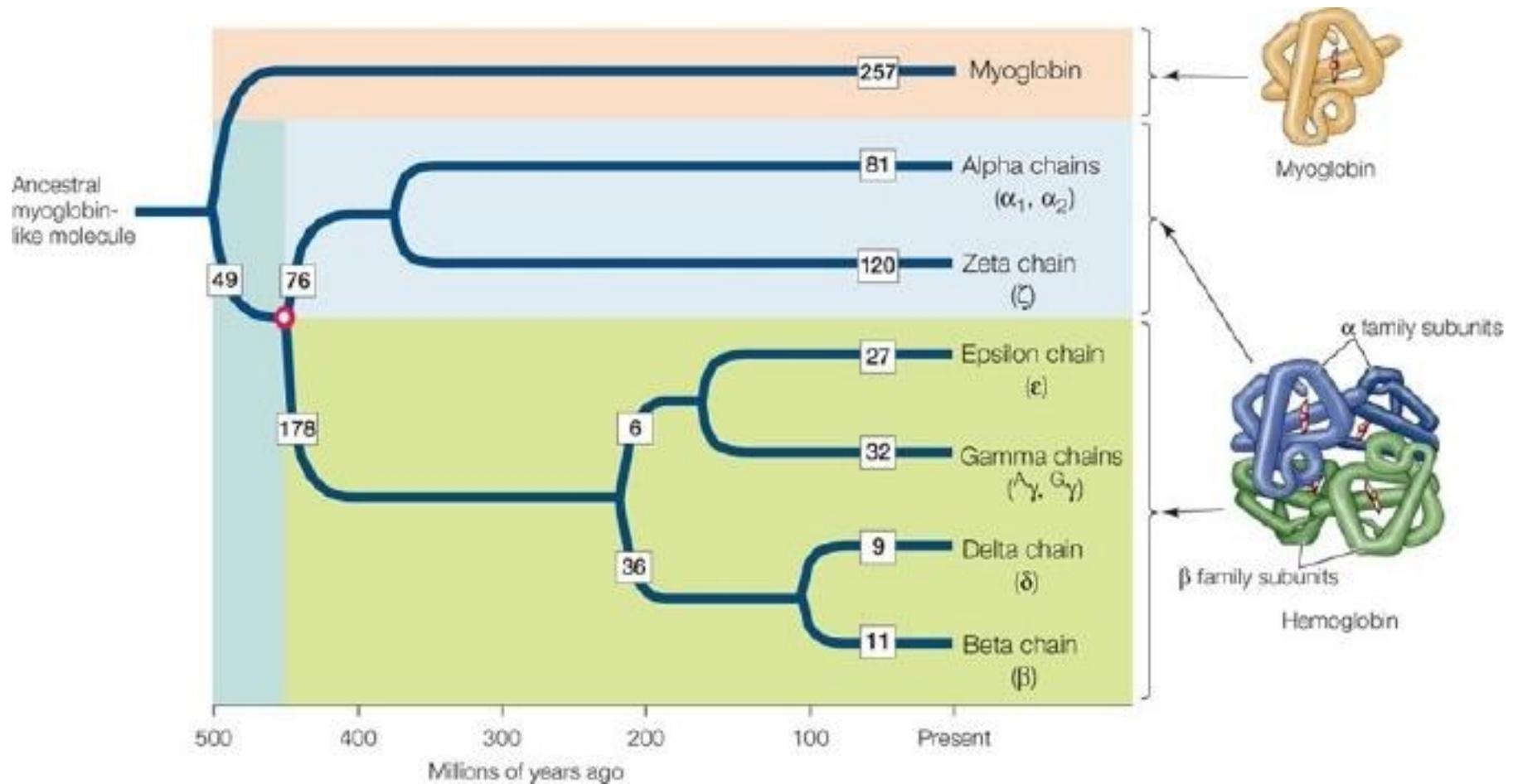
Mindell and Meyer Trends in Ecology and Evolution 2001

Protein Families

Definition:

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

Globins in Human



Homology: why bother?

Mind the gap!

Marco Punta

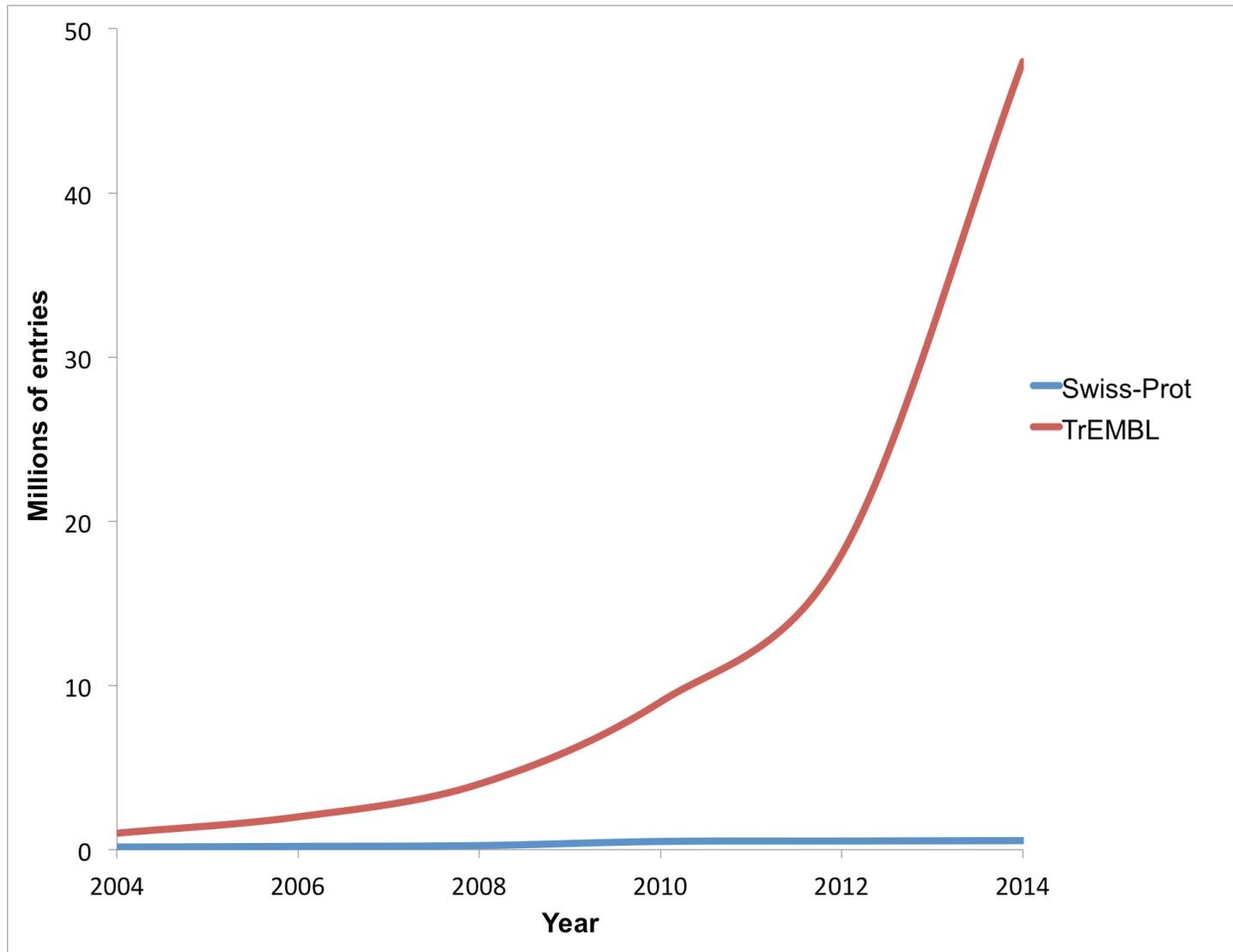


Figure courtesy of Alex Mitchell (EMBL-EBI)

EMBO Workshop, Budapest, 2016

Mind the gap!

Marco Punta

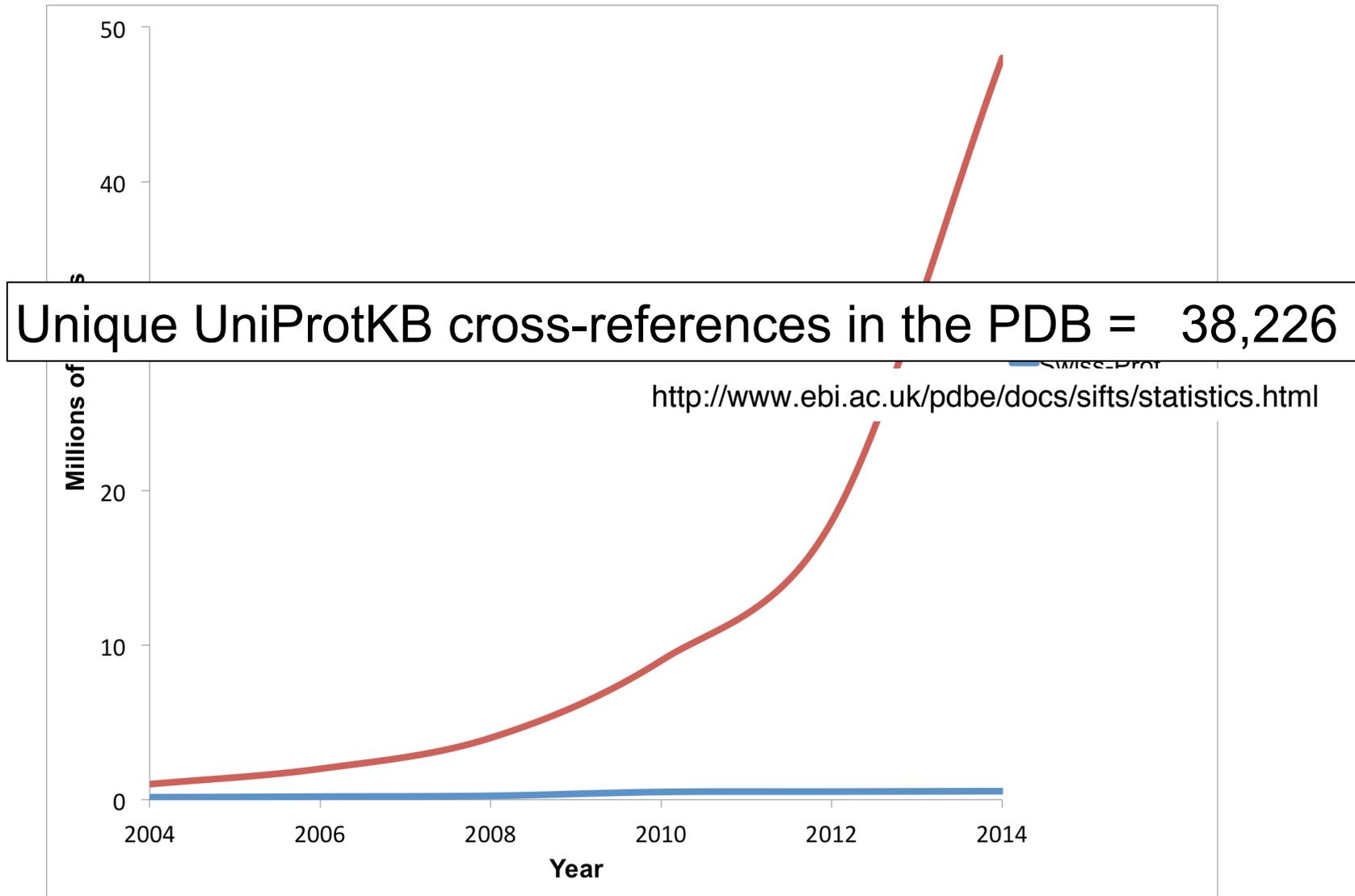


Figure courtesy of Alex Mitchell (EMBL-EBI)

EMBO Workshop, Budapest, 2016

Mind the gap!

Marco Punta

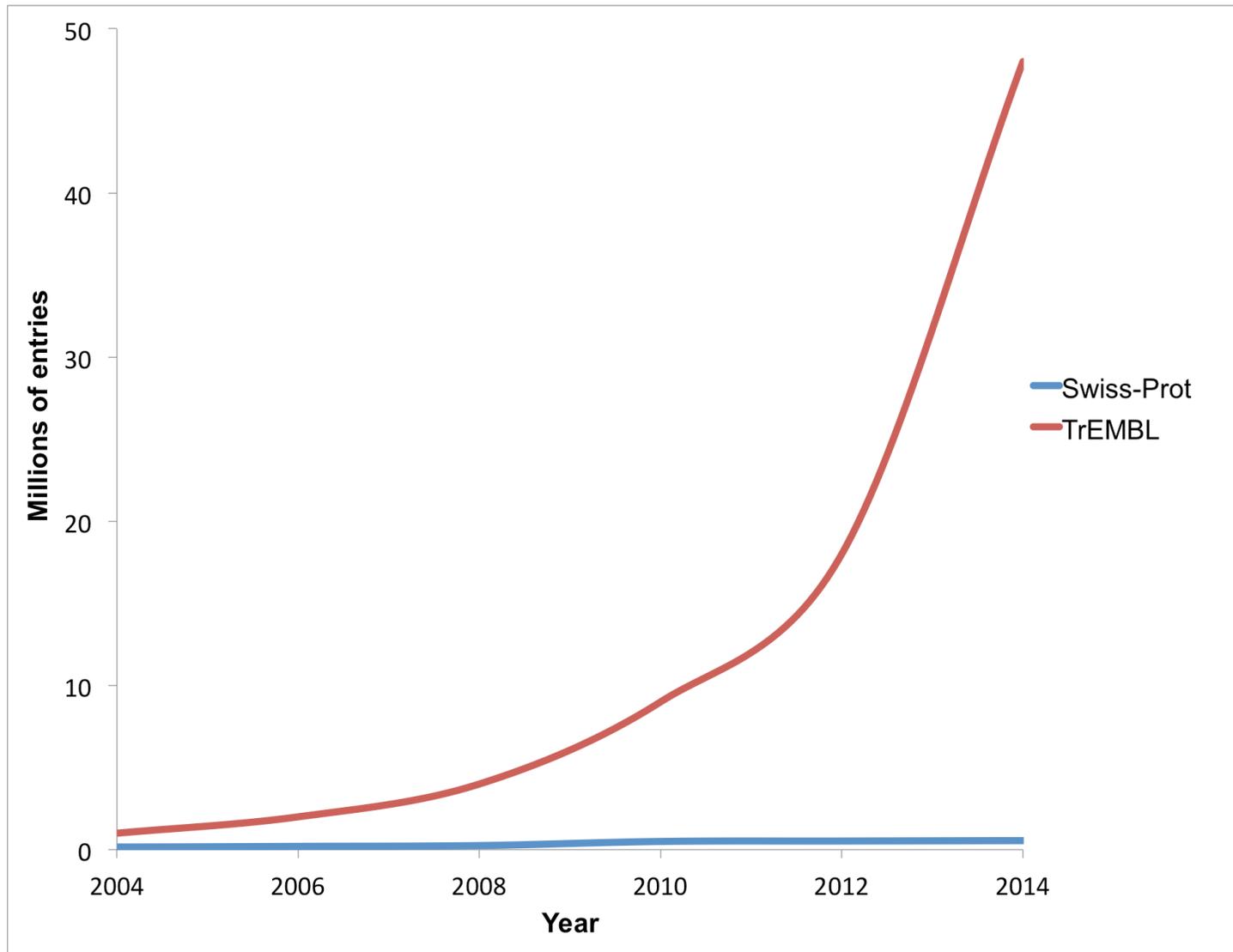


Figure courtesy of Alex Mitchell (EMBL-EBI)

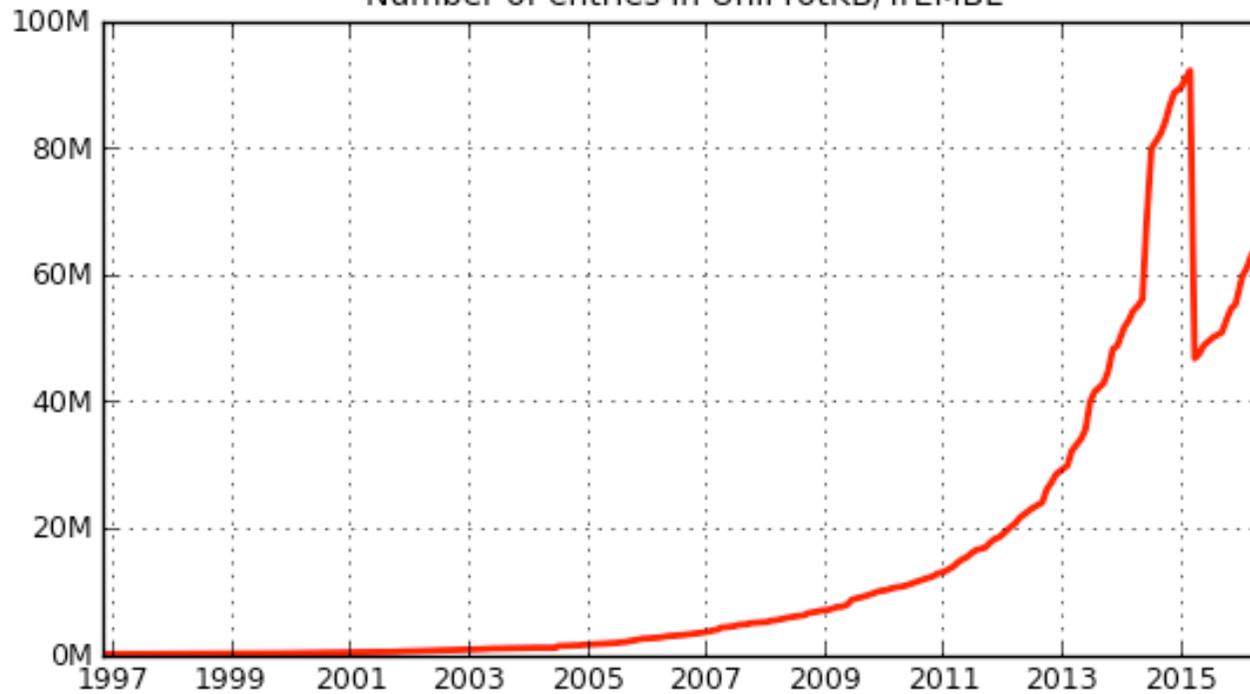
EMBO Workshop, Budapest, 2016

Mind the gap!

Marco Punta

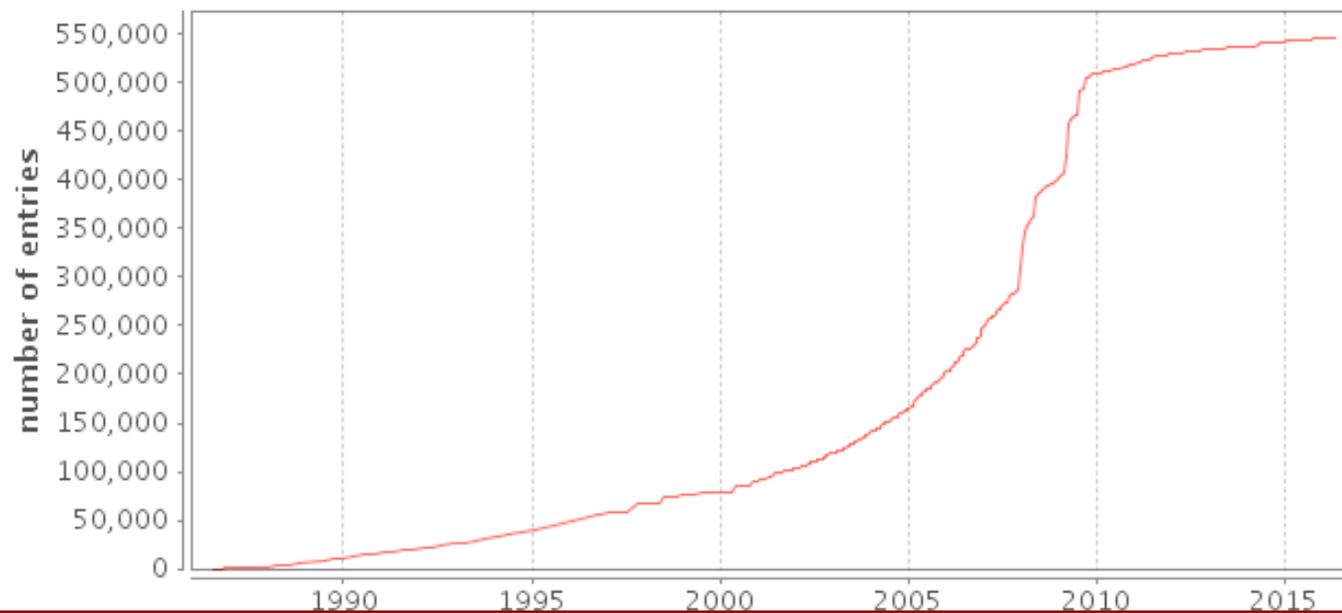
	Protein Existence (PE)	Number of entries
1	Evidence at protein level	92,536
2	Evidence at transcript level	57,757
3	Inferred from homology	387,589
4	Predicted	11,358
5	Uncertain	1,953

Number of entries in UniProtKB/TrEMBL



Marco Punta

Number of entries in UniProtKB/Swiss-Prot over time



Homology: why bother?

Homologous proteins come from a common ancestor, that is, a single ancestral protein with given sequence, structure and functions

Homology \Leftrightarrow similar sequence?

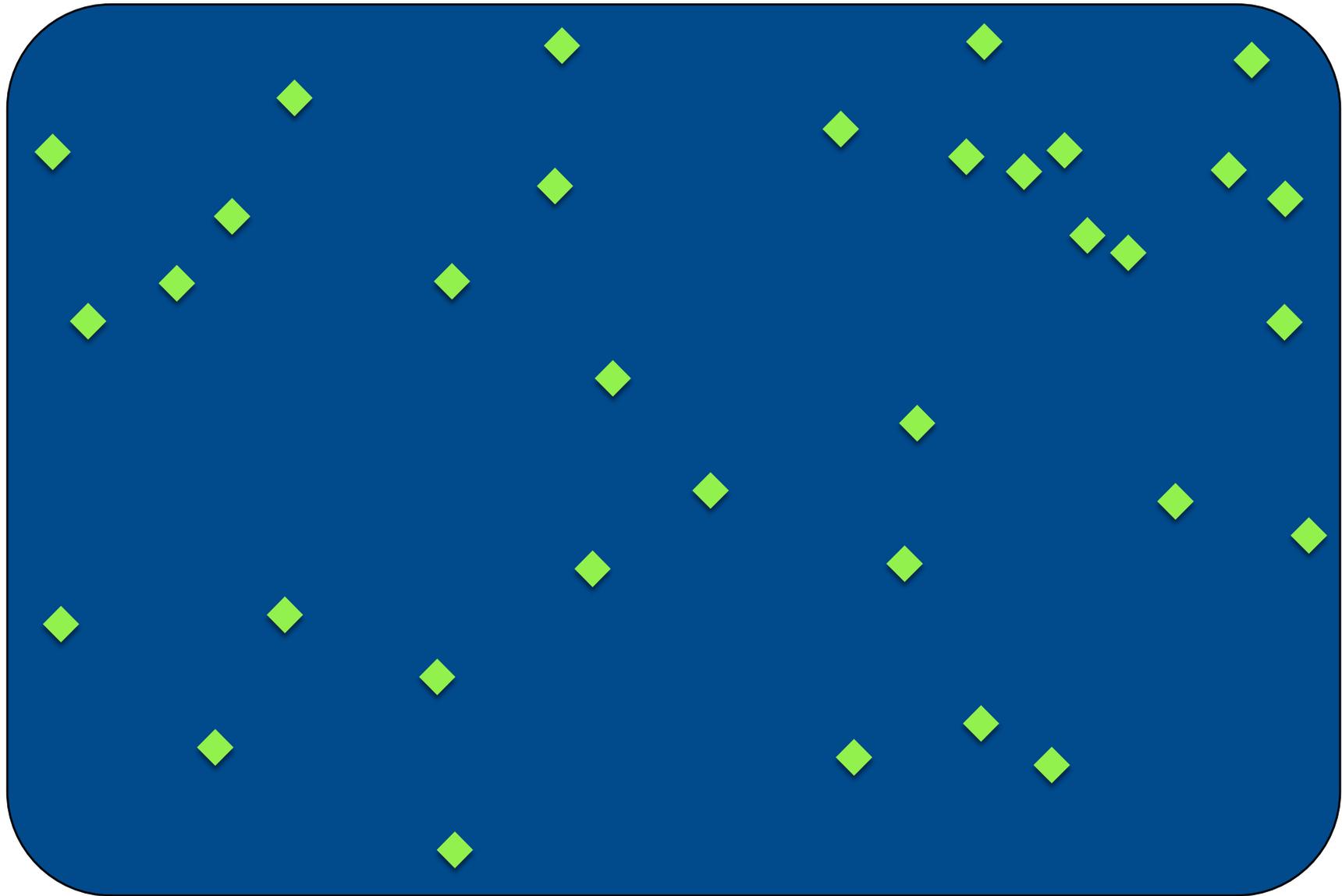
Homology \Leftrightarrow similar structure?

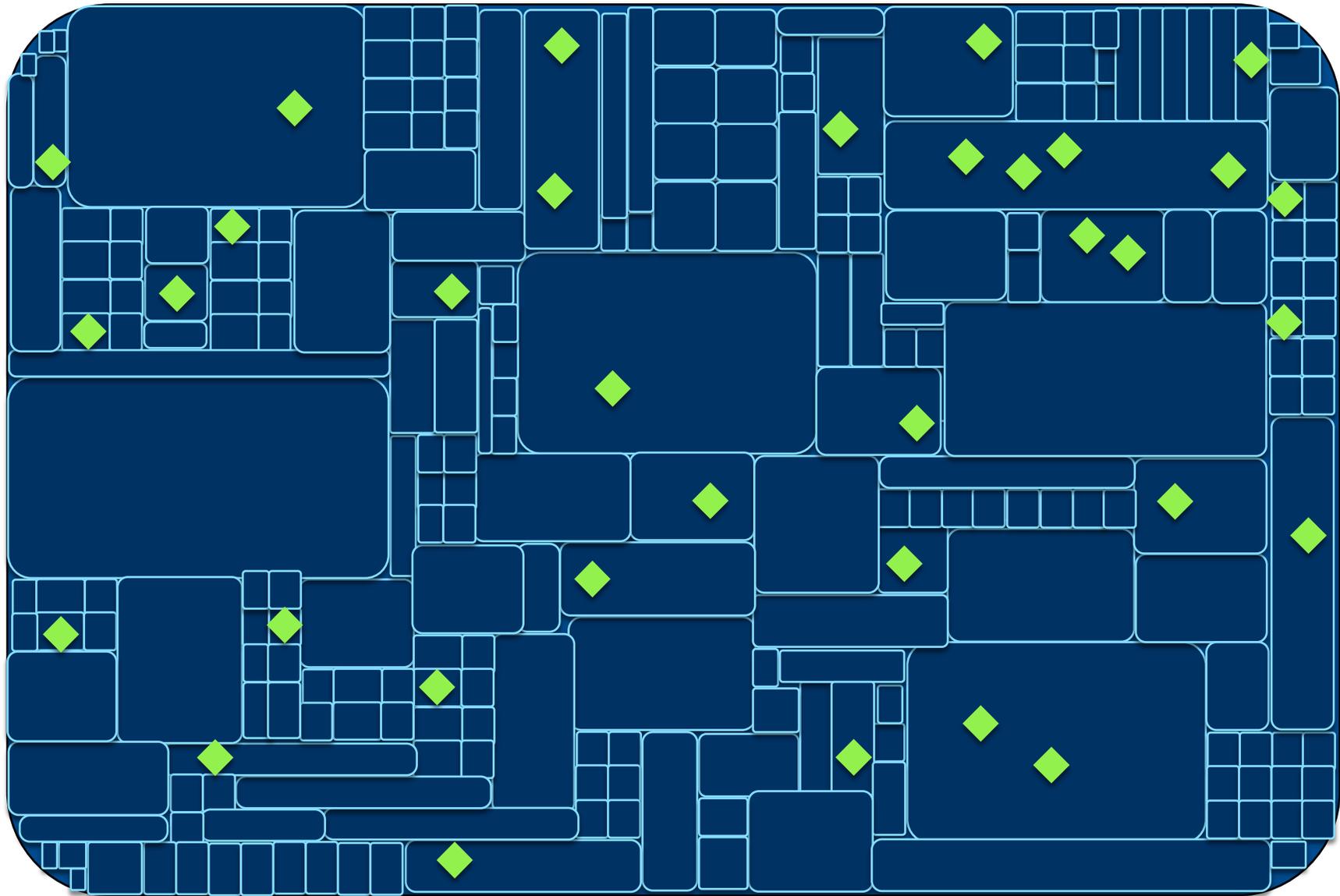
Homology \Leftrightarrow similar function?

Homologous protein regions have a similar (core) structure!

Chothia and Lesk *EMBO J* (1986)

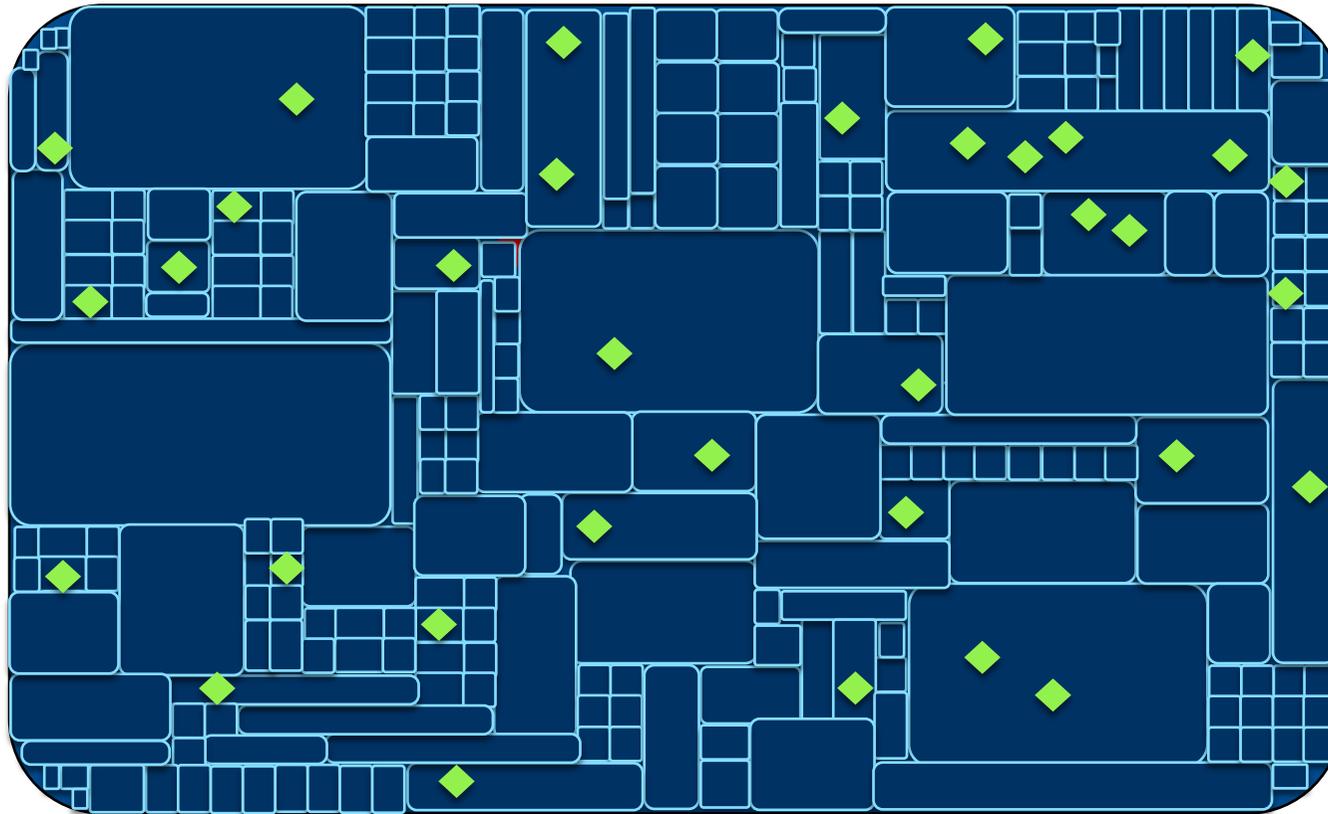






Homology Modelling

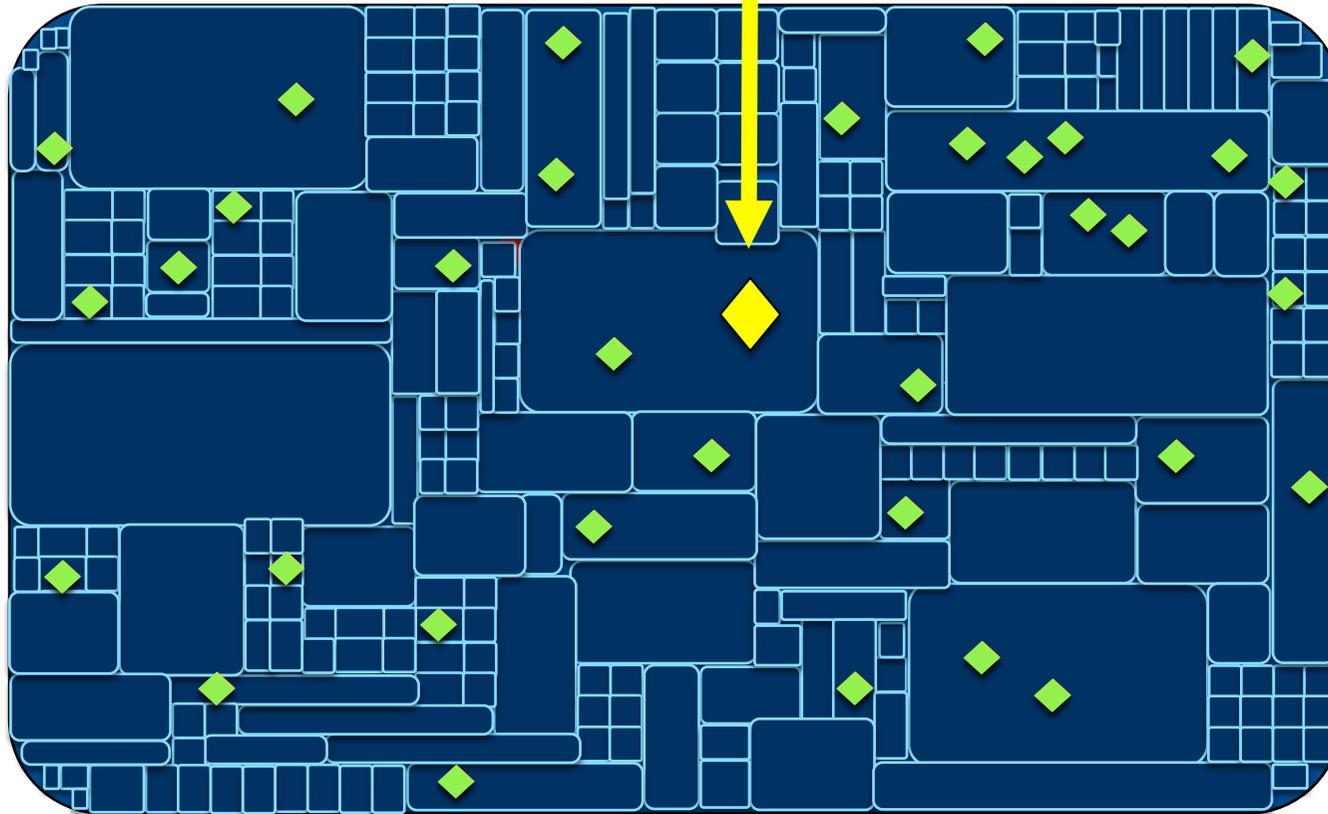
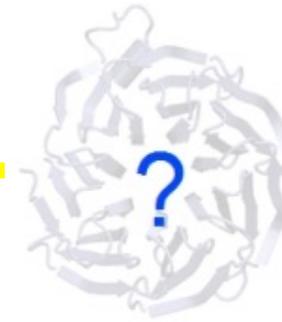
Unknown structure



<http://www.unil.ch/files/live/sites/pmf/files/shared/Technologies/Homology-modeling-450.jpg>

Homology Modelling

Unknown structure



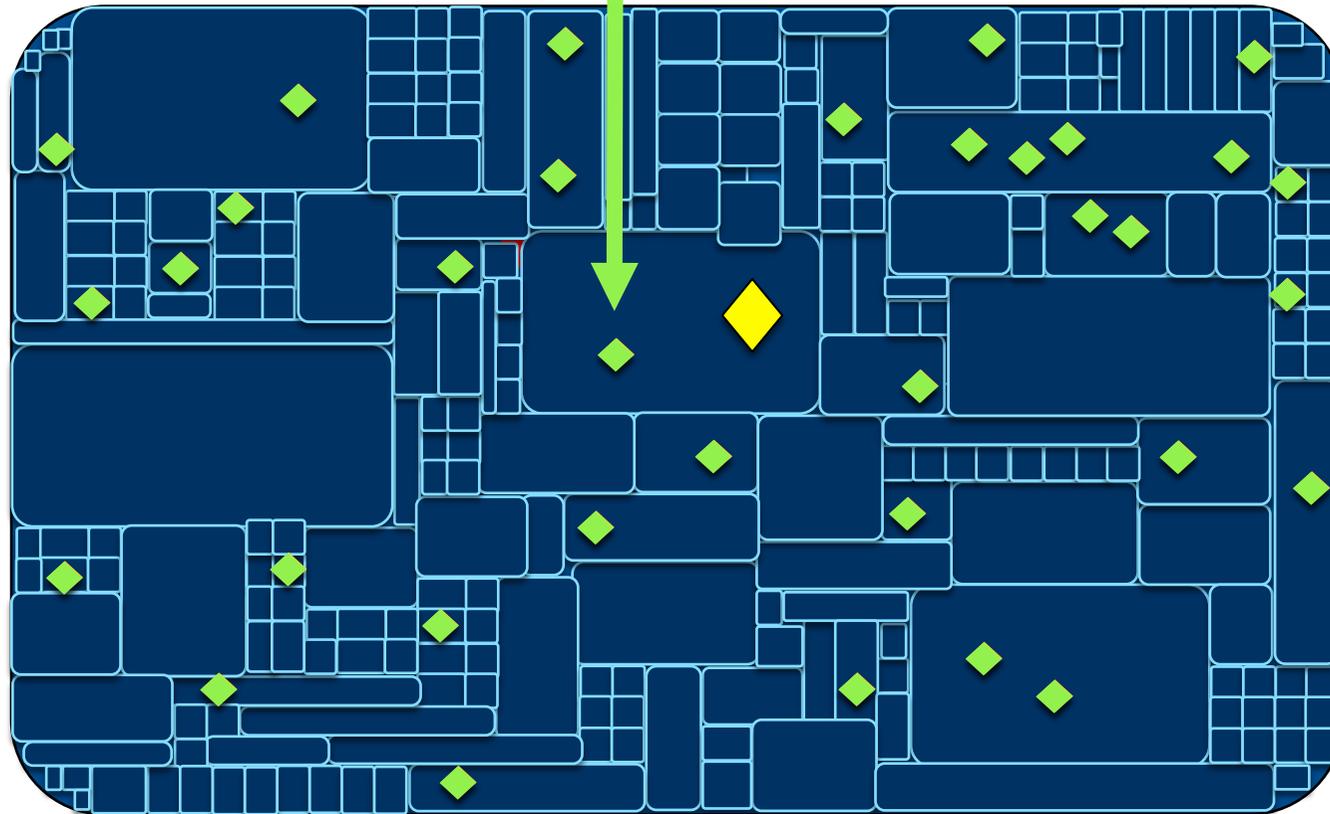
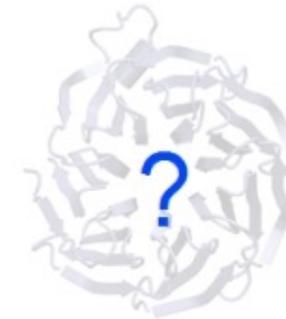
<http://www.unil.ch/files/live/sites/pmf/files/shared/Technologies/Homology-modeling-450.jpg>

Homology Modelling

Template

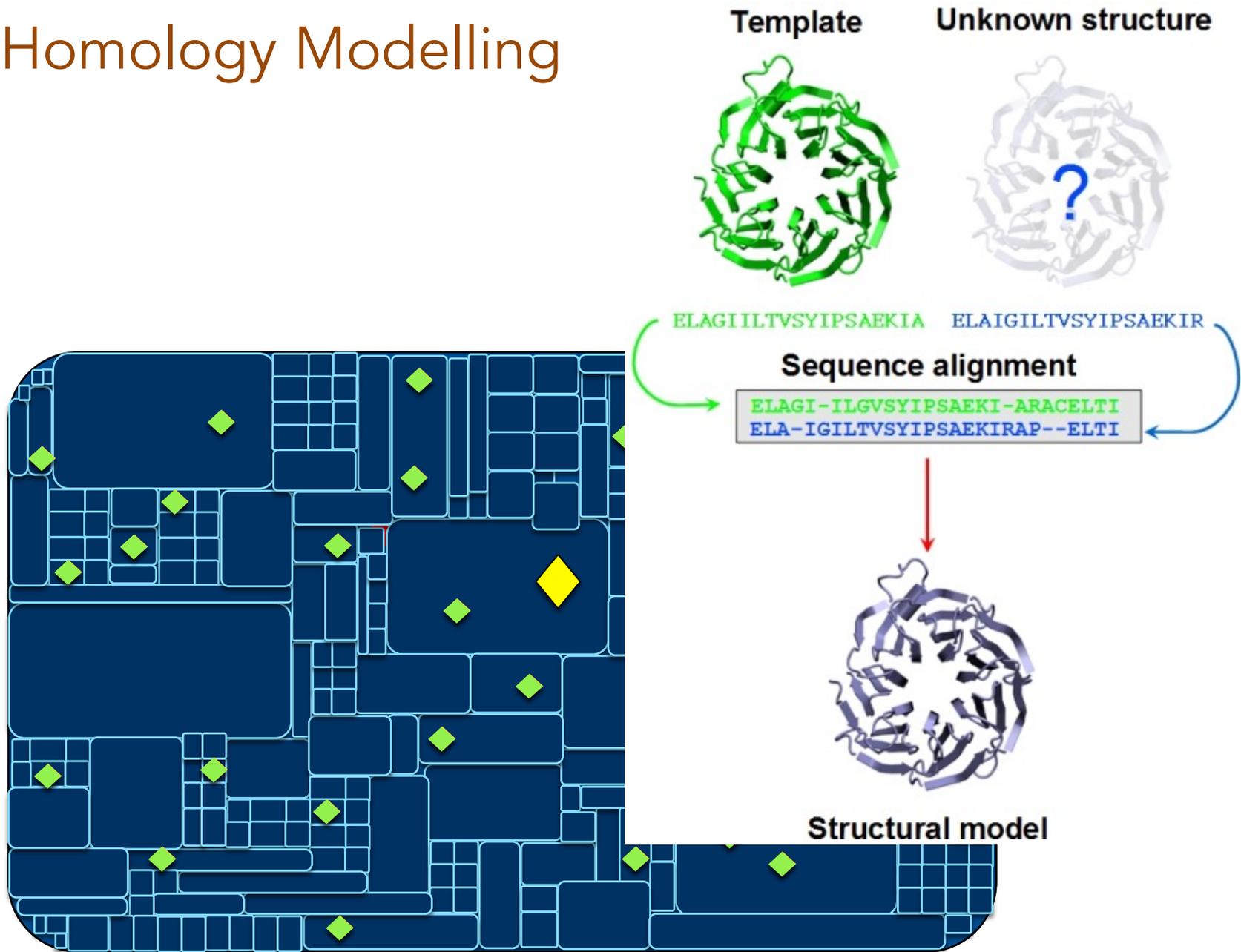


Unknown structure

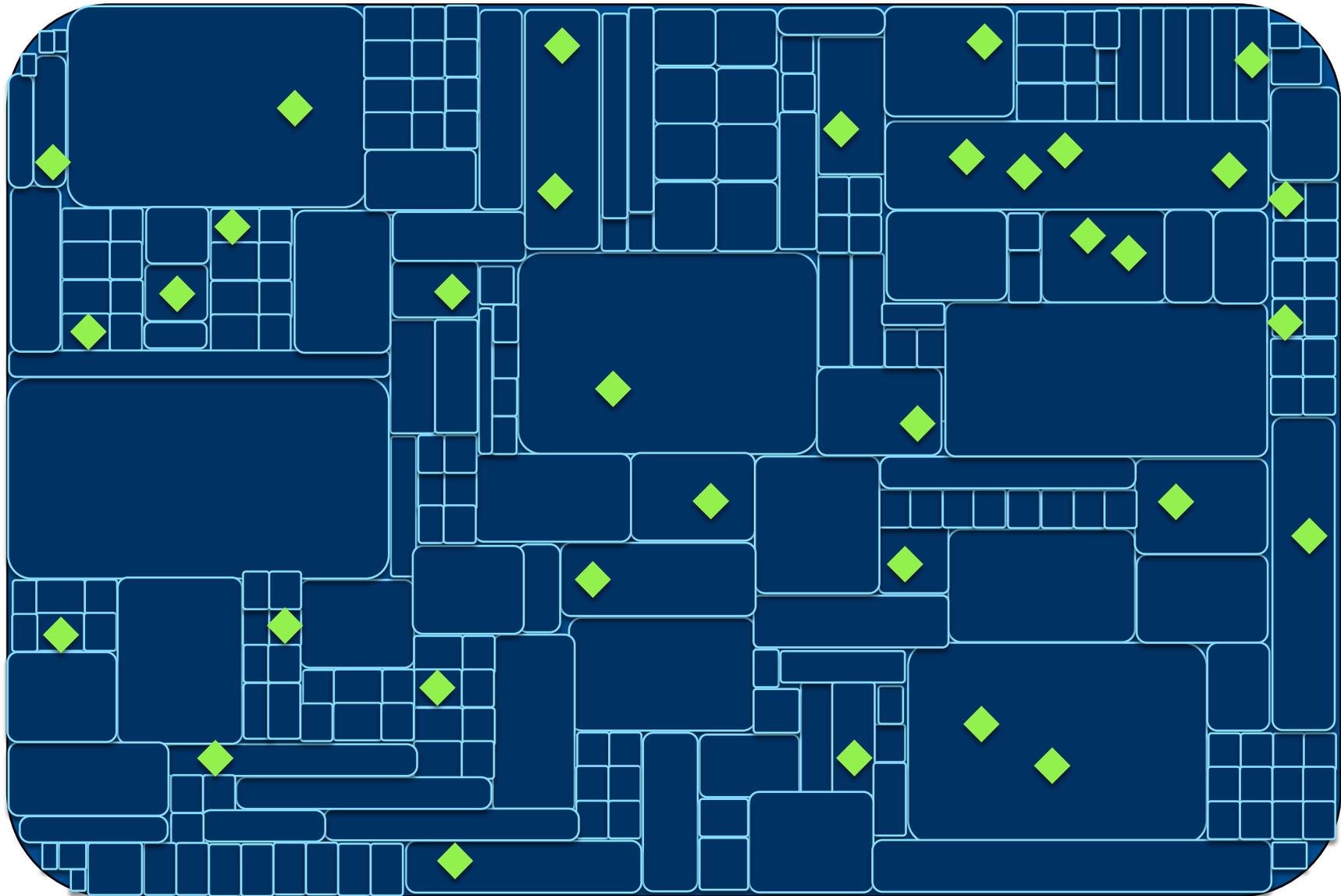


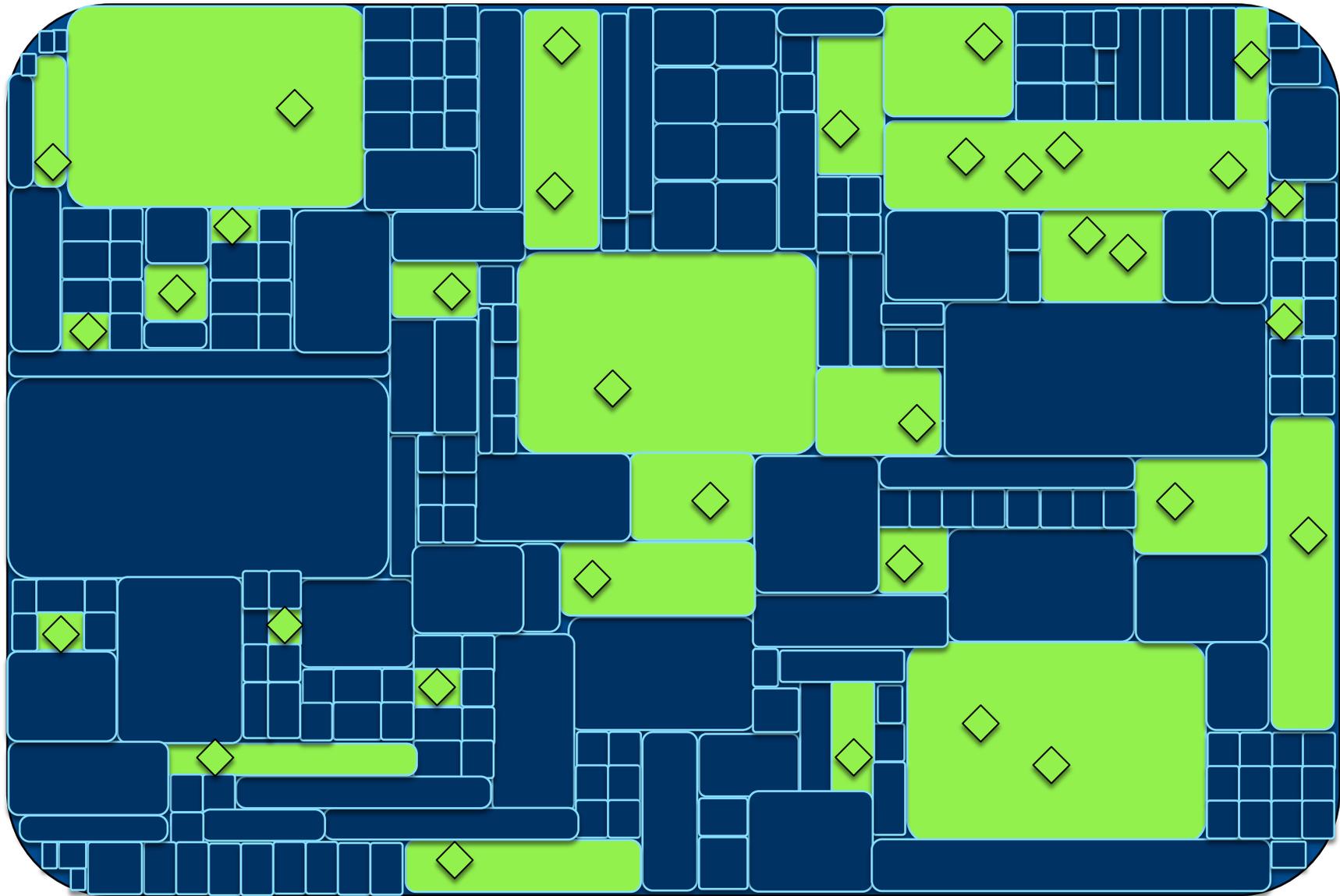
<http://www.unil.ch/files/live/sites/pmf/files/shared/Technologies/Homology-modeling-450.jpg>

Homology Modelling

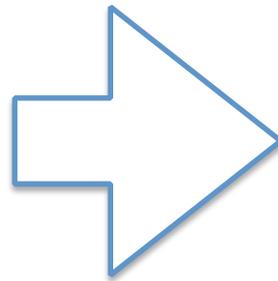


<http://www.unil.ch/files/live/sites/pmf/files/shared/Technologies/Homology-modeling-450.jpg>



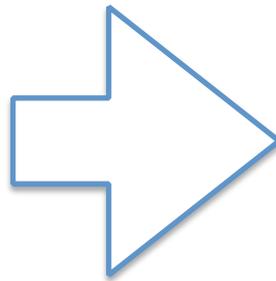


Homology



Structural similarity
(same fold)

Different fold



Not Homologous

Detecting homology

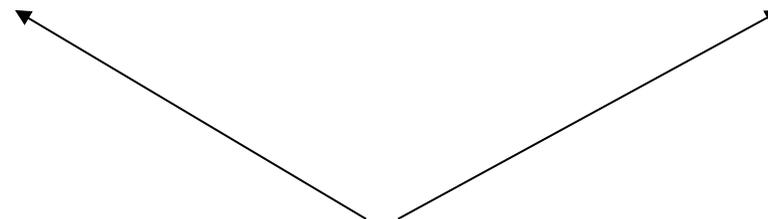
From sequence

Sequences of homologous proteins are related by an evolutionary process, they diverged from a common ancestor.

Modern day homologous proteins have evolved from the same ancestral sequence via a number of events (mutations, insertions, deletions, truncations,...)

ALHWRAALAAATVLLVIVLLAGSWLAVLAE

ALHWKAAGAATVLLVIVLLAGSYLAVLAE



ALHWRAAGAATVLLVIVLLAGSYLAVLAE

```

Human: 1  MGLSDGEWQLVLNVWGKVEADIPGHGQEVLIIRLFKGHPEKFDKFKHLKSEDEMKA 60
          MGLSDGEWQLVLNVWGKVEAD  GHGQEVLI  LFK HPETL  KFDKFK  LKSE  MK SE
Mouse: 1  MGLSDGEWQLVLNVWGKVEADLAGHGQEVLIIGLFKTHPETLDKFDKFKNLKSEEDMKGSE 60

Human: 61  DLKKHGATVLTALGGILKKKGHHEAEIKPLAQSHATKHKIPVKYLEFISECIIQVLQSKH 120
          DLKKHG  TVLTALG  ILKKKG  H  AEI  PLAQSHATKHKIPVKYLEFISE  II  VL  H
Mouse: 61  DLKKHGCTVLTALGTILKKKGQHAAEQPLAQSHATKHKIPVKYLEFISEIIIEVLKRRH 120

Human: 121 PGDFGADAQGAMKALELFRKDMASNYKELGFQG 154
          GDFGADAQGAM  KALELFR  D  A  YKELGFQG
Mouse: 121 SGDFGADAQGAMSKALELFRNDIAAKYKELGFQG 154
    
```

84% identities!

```

Human: 1  MGLSDGEWQLVLNVWGKVEADIPGHGQEVLIIRLFKGGHPETLEKFDKFKHLKSEDEMKA 60
          MGLSDGEWQLVLNVWGKVEAD  GHGQEVLI  LFK HPETL  KFDKFK  LKSE  MK SE
Mouse: 1  MGLSDGEWQLVLNVWGKVEADLAGHGQEVLIIGLFKTHPETLDKFDKFKNLKSEEDMKGSE 60

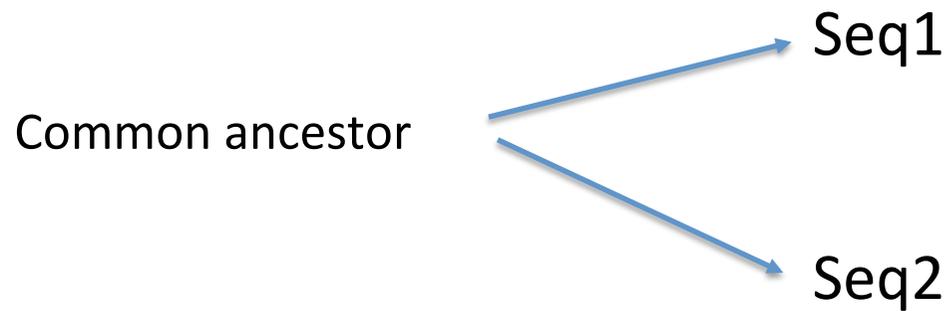
Human: 61  DLKKHGATVLTALGGILKKKGHHEAEIKPLAQSHATKHKIPVKYLEFISECIIQVLQSKH 120
          DLKKHG  TVLTALG  ILKKKG  H  AEI  PLAQSHATKHKIPVKYLEFISE  II  VL  H
Mouse: 61  DLKKHGCTVLTALGTILKKKGQHAAEQPLAQSHATKHKIPVKYLEFISEIIIEVLKRRH 120

Human: 121 PGDFGADAQGAMNKALELFRKDMASNYKELGFQG 154
          GDFGADAQGAM  KALELFR  D  A  YKELGFQG
Mouse: 121 SGDFGADAQGAMSKALELFRNDIAAKYKELGFQG 154
    
```

84% identities!

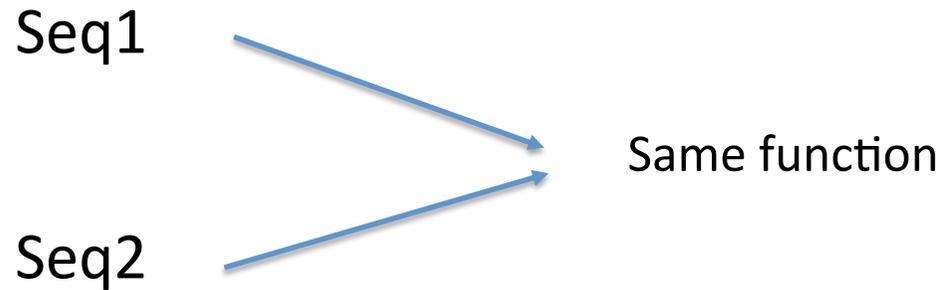
How can we explain “excess sequence similarity”* with respect to what we expect at random?

Homology: divergent evolution

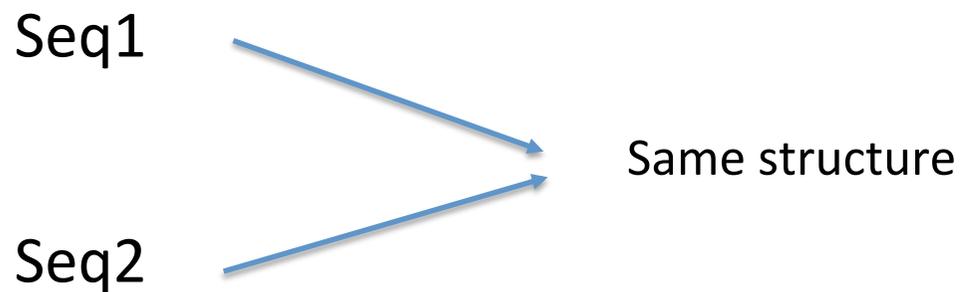


Stevens J Mol Recogn (2007)
Elias and Tawfik *J Biol Chem* (2012)

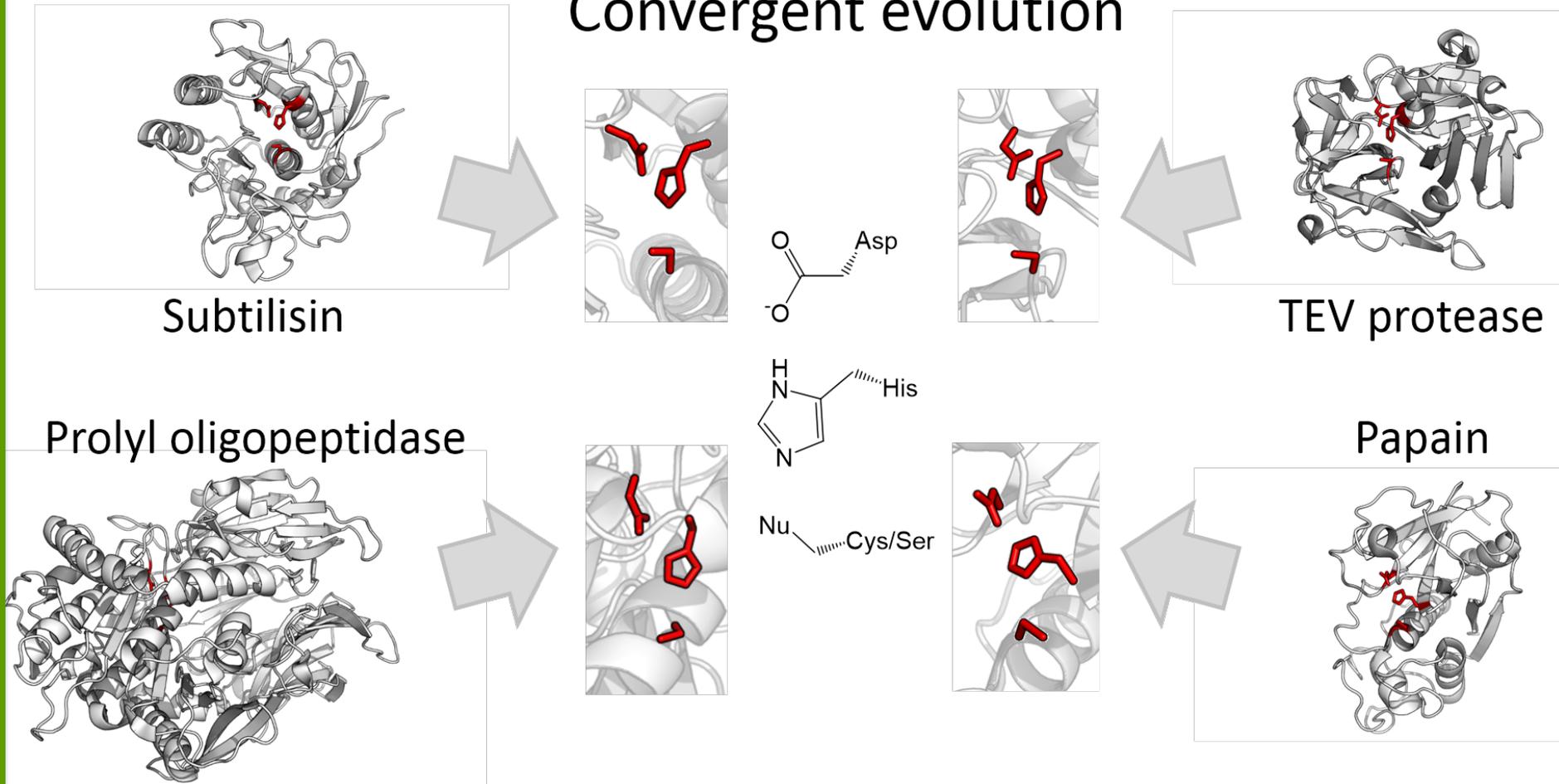
Analogy I: convergent evolution



Analogy II: concurrence

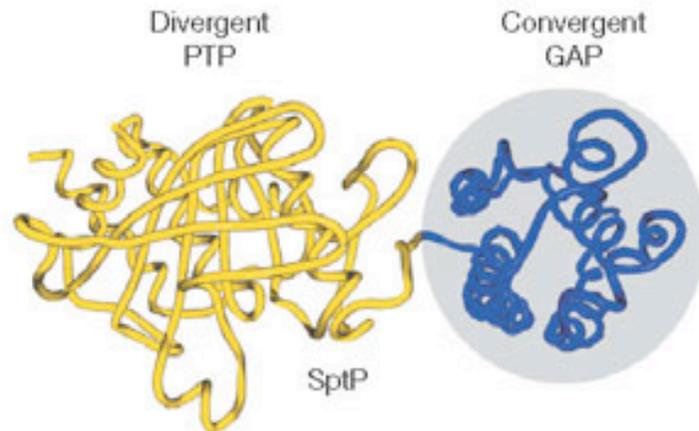


Convergent evolution



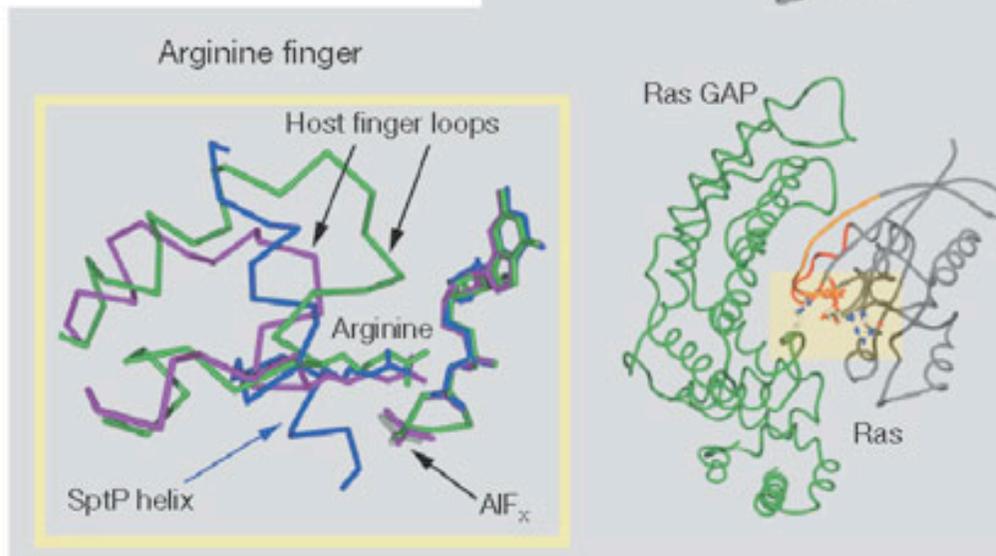
"Triad Convergence" by Thomas Shafee - Own work. Licensed under CC BY-SA 3.0 via Wikimedia Commons - https://commons.wikimedia.org/wiki/File:Triad_Convergence.png#/media/File:Triad_Convergence.png

Tyrosine phosphatase Rho GTPase activating domain



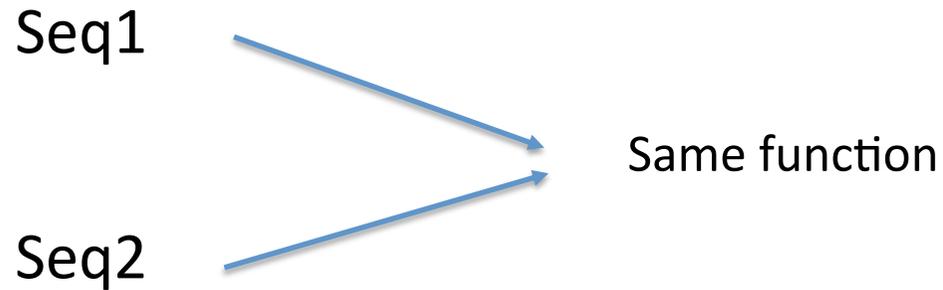
Salmonella virulence factor SptP

Activation of GTPases by *Salmonella* GAP leads to profuse rearrangements of the actin cytoskeleton and subsequent bacterial internalization into intestinal epithelial cells.

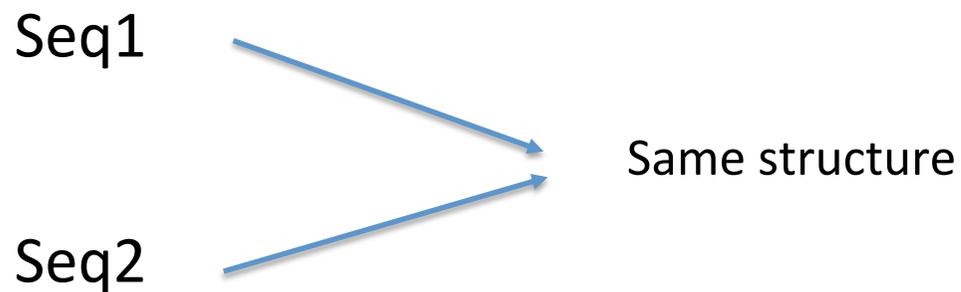


Stebbins and Galán Nature 412 (2001)

Analogy I: convergent evolution



Analogy II: concurrence




```

Human: 1  MGLSDGEWQLVLNVWGKVEADIPGHGQEVLIIRLFKGGHPETLEKFDKFKHLKSEDEMKA 60
          MGLSDGEWQLVLNVWGKVEAD  GHGQEVLI  LFK HPETL  KFDKFK  LKSE  MK SE
Mouse: 1  MGLSDGEWQLVLNVWGKVEADLAGHGQEVLIIGLFKTHPETLDKFDKFKNLKSEEDMKGSE 60

Human: 61  DLKKHGATVLTALGGILKKKGHHEAEIKPLAQSHATKHKIPVKYLEFISECIIQVLQSKH 120
          DLKKHG  TVLTALG  ILKKKG  H  AEI  PLAQSHATKHKIPVKYLEFISE  II  VL  H
Mouse: 61  DLKKHGCTVLTALGTILKKKGQHAAEQPLAQSHATKHKIPVKYLEFISEIIIEVLKRRH 120

Human: 121 PGDFGADAQGAMNKALELFRKDMASNYKELGFQG 154
          GDFGADAQGAM  KALELFR  D  A  YKELGFQG
Mouse: 121 SGDFGADAQGAMSKALELFRNDIAAKYKELGFQG 154
    
```

84% identities!

How can we explain “excess sequence similarity”* with respect to what we expect at random?

“[...]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)



Sequence alignment, 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, ...

Human: 1 MGLSDGEWQLVLNVWGKVEADIPGHGQEVLIIRLFKGHPEKFDKFKHLKSEDEMKA 60
 MGLSDGEWQLVLNVWGKVEAD GHGQEVLI LFK HPETL KFDKFK LKSE MK SE
 Mouse: 1 MGLSDGEWQLVLNVWGKVEADLAGHGQEVLIIGLFKTHPETLDKFDKFKNLKSEEDMKGSE 60

Human: 61 DLKKGATVLTALGGILKKKGHHEAEIKPLAQSHATKHKIPVKYLEFISECIIQVLQSKH 120
 DLKKG TVLTALG ILKKG H AEI PLAQSHATKHKIPVKYLEFISE II VL H
 Mouse: 61 DLKKGCTVLTALGTILKKGQHA AEIQPLAQSHATKHKIPVKYLEFISEIIIEVLKGRH 120

Human: 121 PGDFGADAQGAMKALELFRKDMASNYKELGFQG 154
 GDFGADAQGAM KALELFR D A YKELGFQG
 Mouse: 121 SGDFGADAQGAMSKALELFRNDIAAKYKELGFQG 154

84% identities!

Sequence alignment, 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, ...

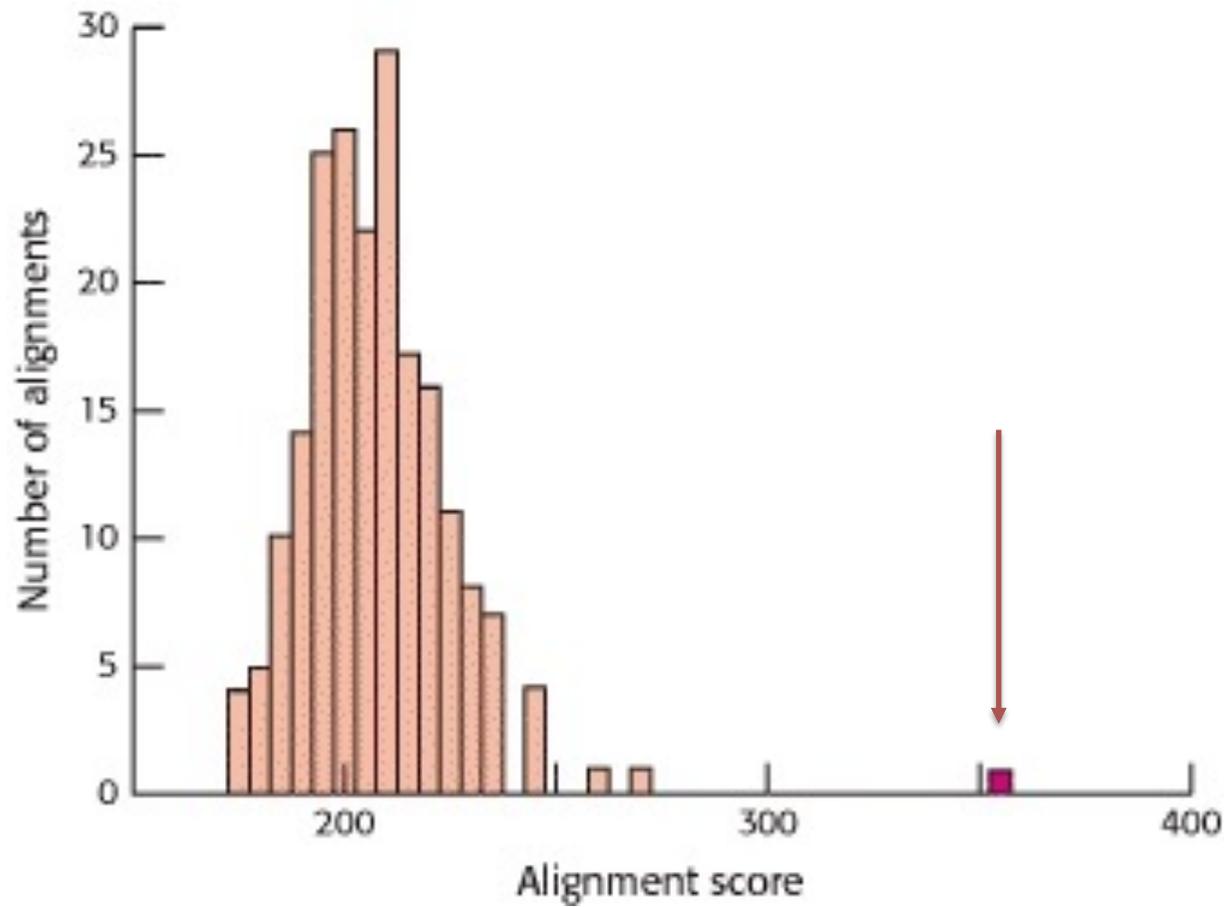


Image credits: <http://www.ncbi.nlm.nih.gov/books/NBK22456/figure/A937/?report=objectonly>

```

MYG_HUMAN      1 MGLSDGEWQLVLNVWGKVEADIPGHGQEVLIIRLFKGHPETLEKFDKFKHLKSEDEMKASE 60
                M LS      V   WGKV A      G E L R F   P T   F F      D   S
E9M4D4_HUMAN  1 MVLSPADKTNVKAAWGKVGAGAHAGEYGAEALERMFLSFPTTKTYFPHF-----DLSHGSA 54

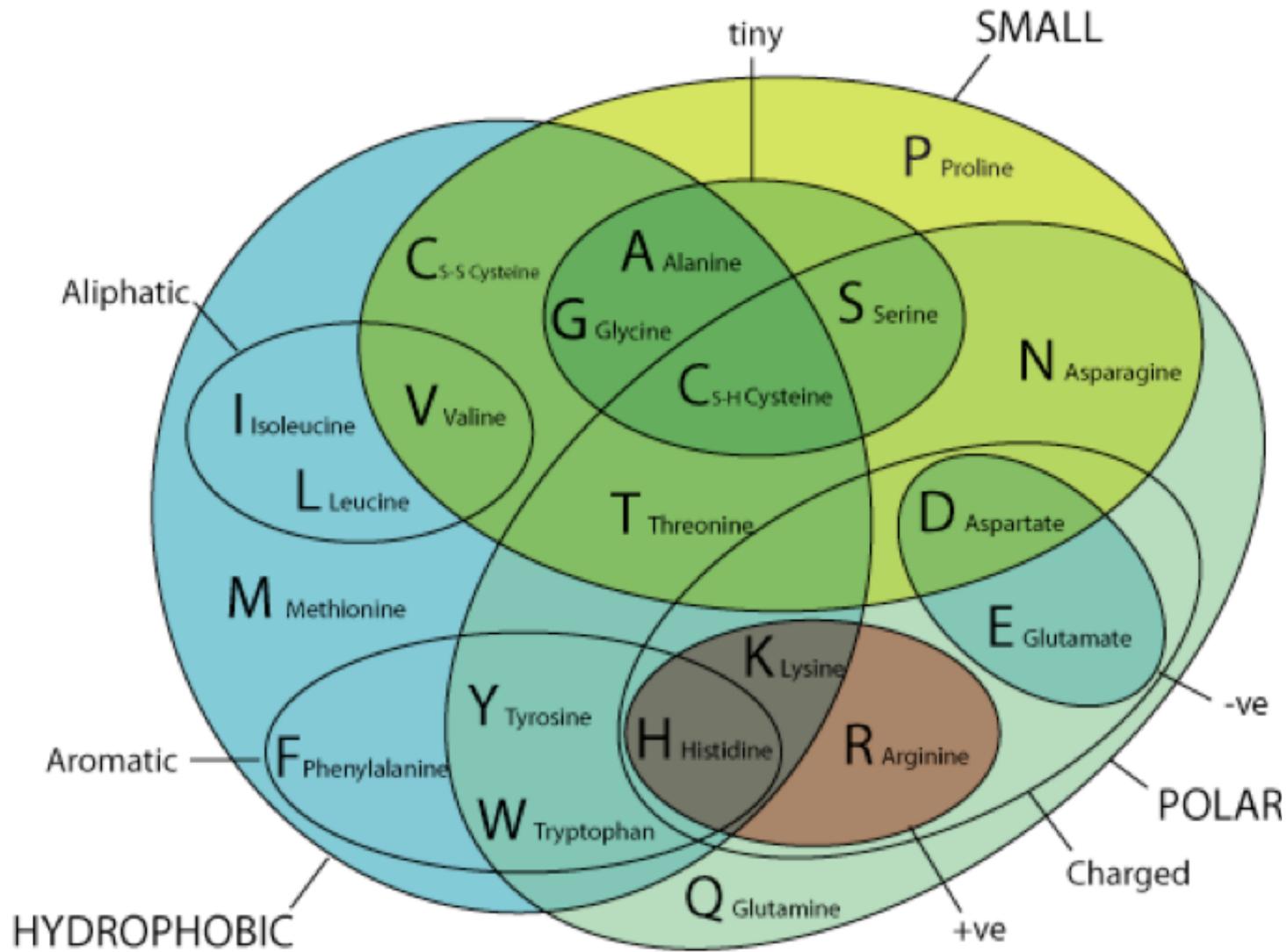
MYG_HUMAN      61 DLKKHGATVLTALGGILKKKGHHEAEIKPLAQSHATKHKI-PVKY 104
                K H   V AL                L   HA K   PV
E9M4D4_HUMAN  55 QVKGHSKKVADALTNVAHVDDMPNALSALSDLHAHKLRVDPVNF 99
    
```

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

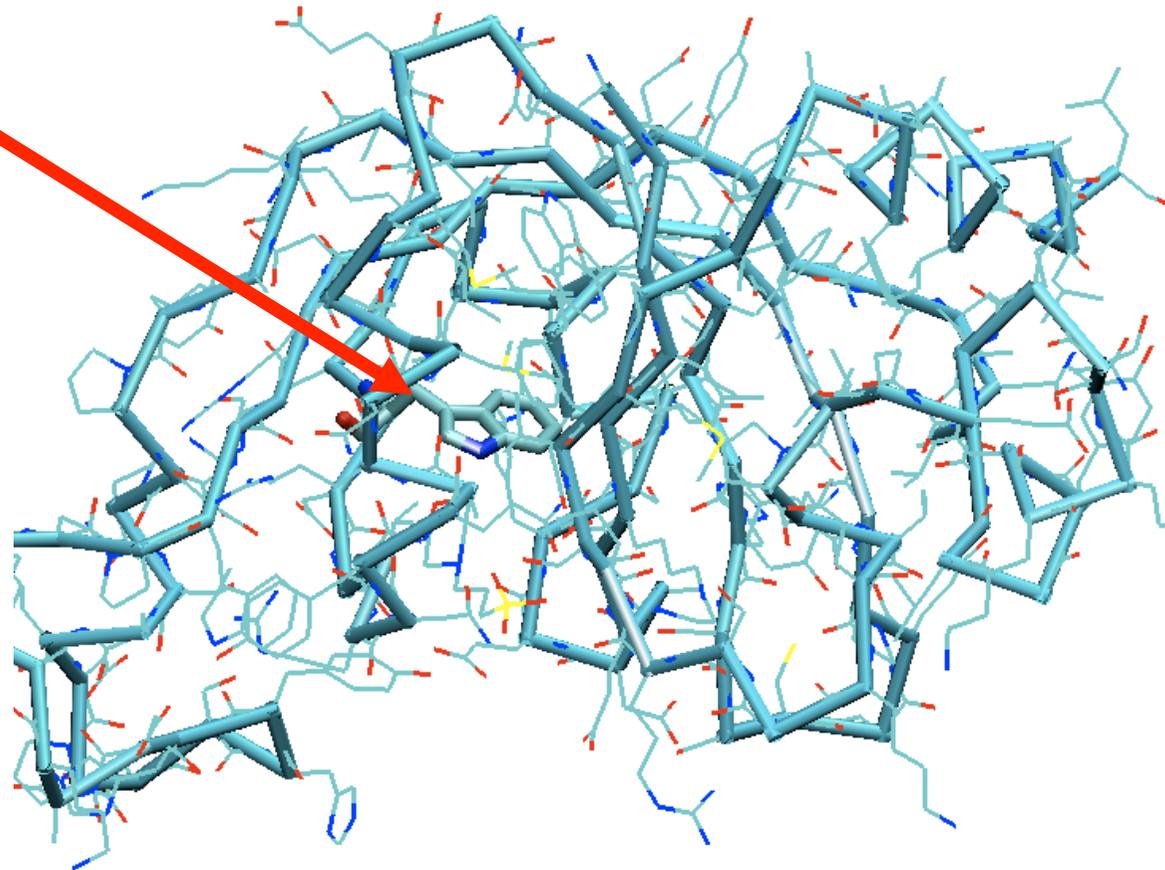
aa physico-chemical properties

Marco Punta



Protein structural and functional constraints

Trp (W)



An Introduction to Sequence Similarity (“Homology”) Searching

[William R. Pearson](#)¹

[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, orthology, paralogy, sequence alignment, multiple alignment, sequence evolution

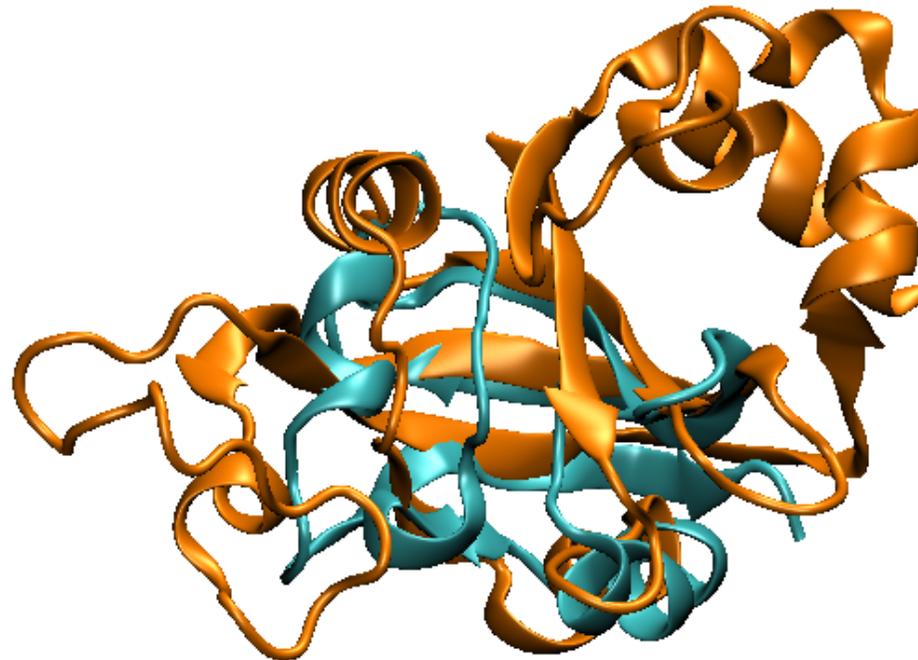
Subject: Bioinformatics, Bioinformatics Fundamentals, Finding Similarities and Inferring Homologies

From structure

Marco Punta

2EVE

1J2B



Z-score = 2.0

RMSD = 3.2

Lali = 54

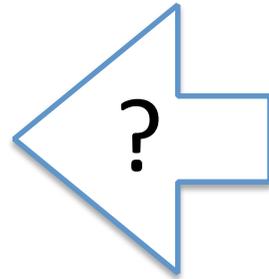
%id = 6

DALI: http://ekhidna.biocenter.helsinki.fi/dali_lite/start

CE, VAST, FATCAT, PDBeFold, ...

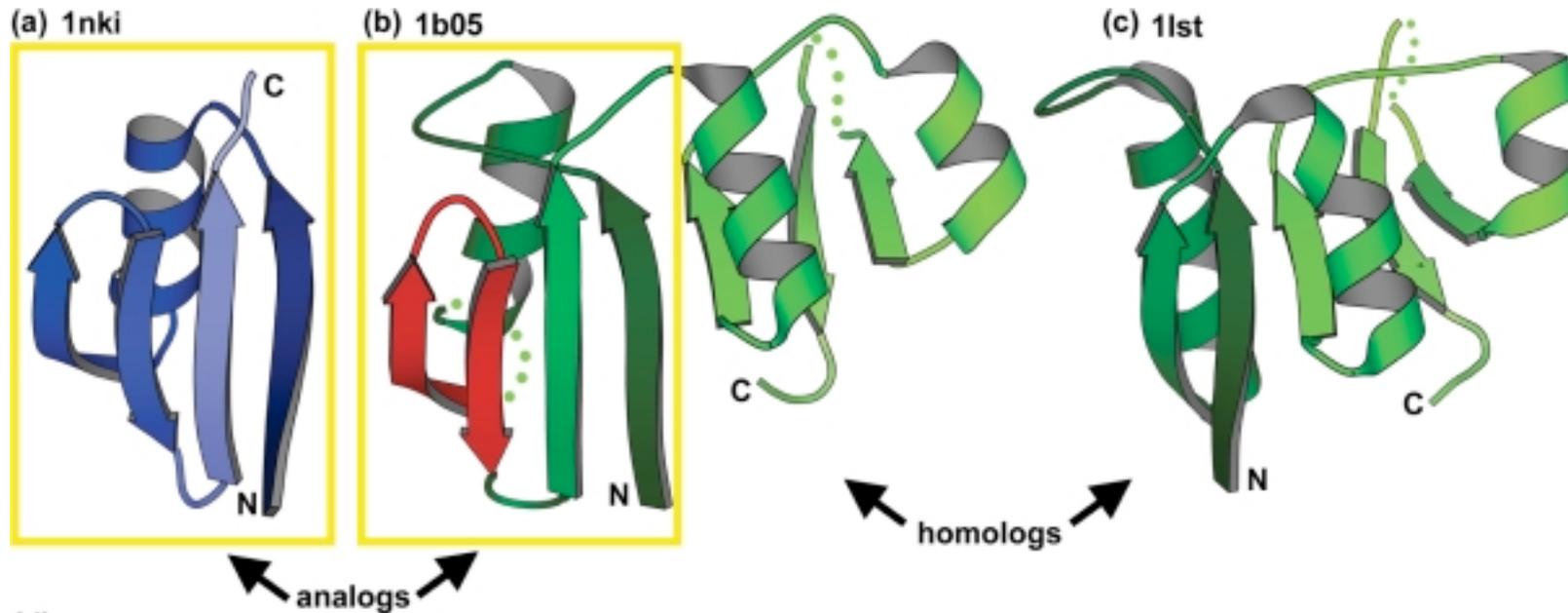
Bertonati, Punta et al. Proteins 2009

Homology



Structural similarity
(same fold)

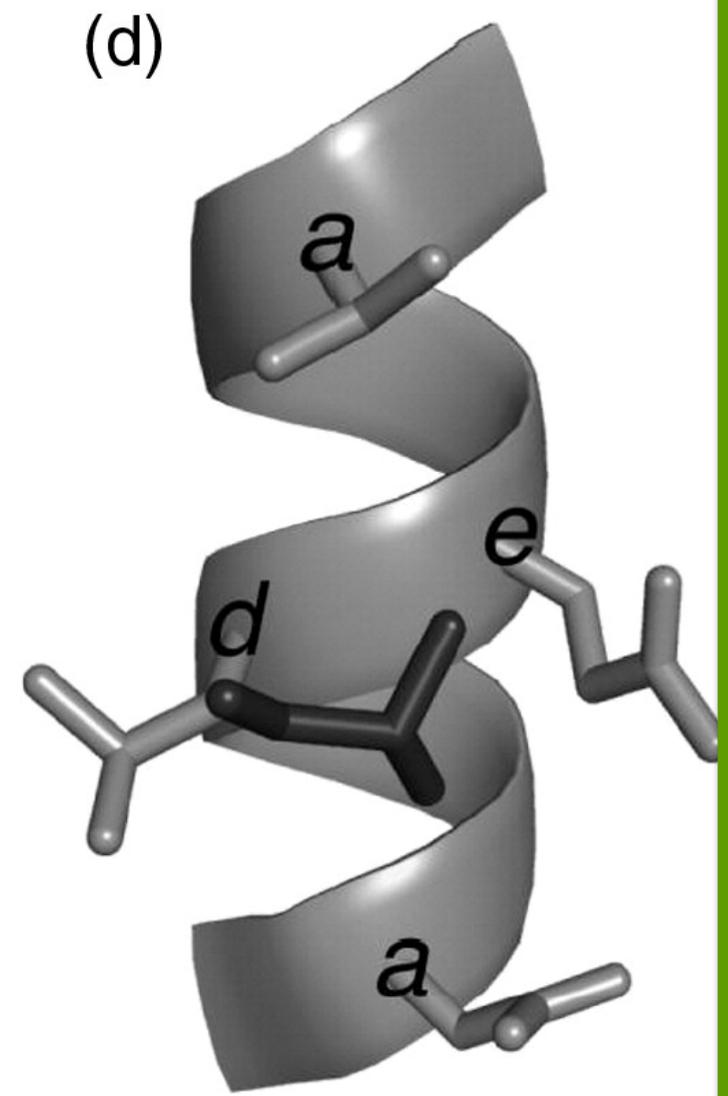
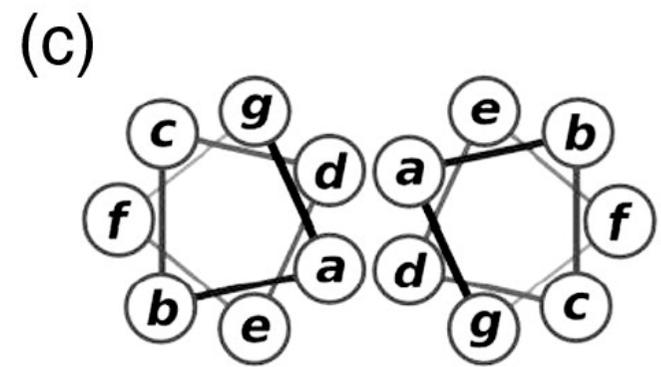
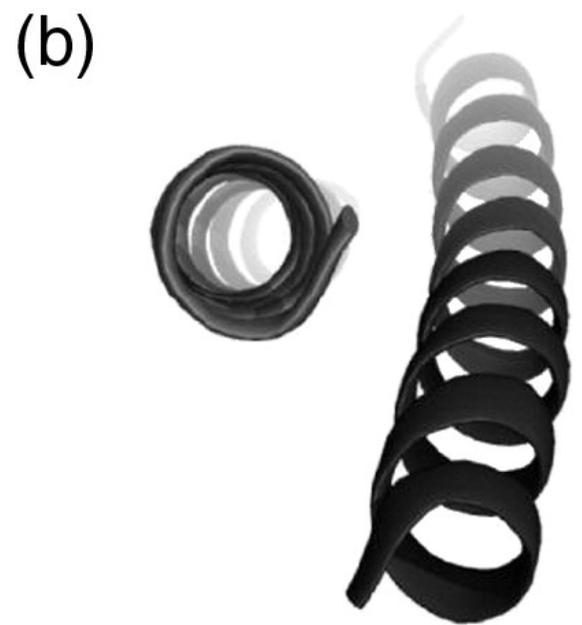
<http://prodata.swmed.edu/malisam/>



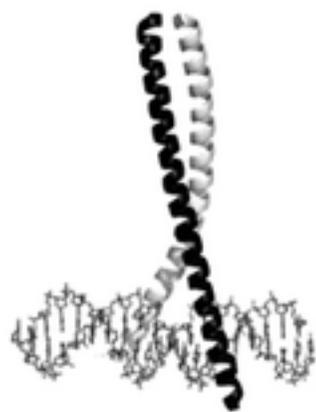
(d)

1nki_A	7	H	TLAVA	----	DLPASIAFYRDL	---	LGFRLEARMD	-QGAYLELGSL	--	WLCLSREP	-----	53								
1b05_A	14	T	LVRNNG	[12]	GVPESNVSRDLFE	[145]	YKLNWVNVNERIVLERN	[12]	QV	TYLPISSEVTDVNRYSGEIDMTYN	[6]	FQ	KLKKEIPNEVRV	[221]	ARLVKFWVGG	497				
1l1st_A	5	T	VRIGTD	[16]	GFDIDLGNEMCKRM	-----	-----	-----	-----	QVKC	TWVAS	-DFDALIPSLKAKKIDAIIS	[4]	TD	KRQQEI	--	AFSD	[110]	GVGLRKDDYE	205

Cheng et al. Nucleic Acids Res. 2008



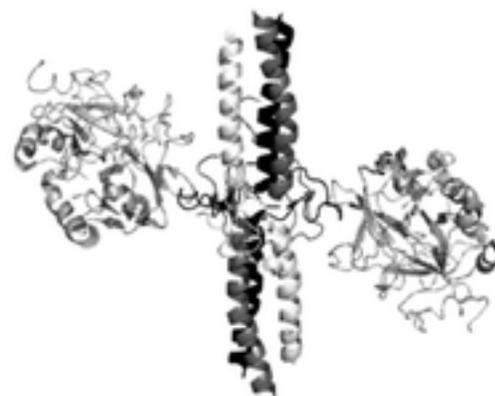
Rackham et al. J Mol Biol. 2010



Fos/Jun
1FOS



MADS box transcription factor
1EGW



Fibrinogen beta
2A45



HIV-1 gp41
1AIK



Tyrosine hydrolase
1TOH



Transcription
repressor
1LBH



Tetrabrachion
1FE6

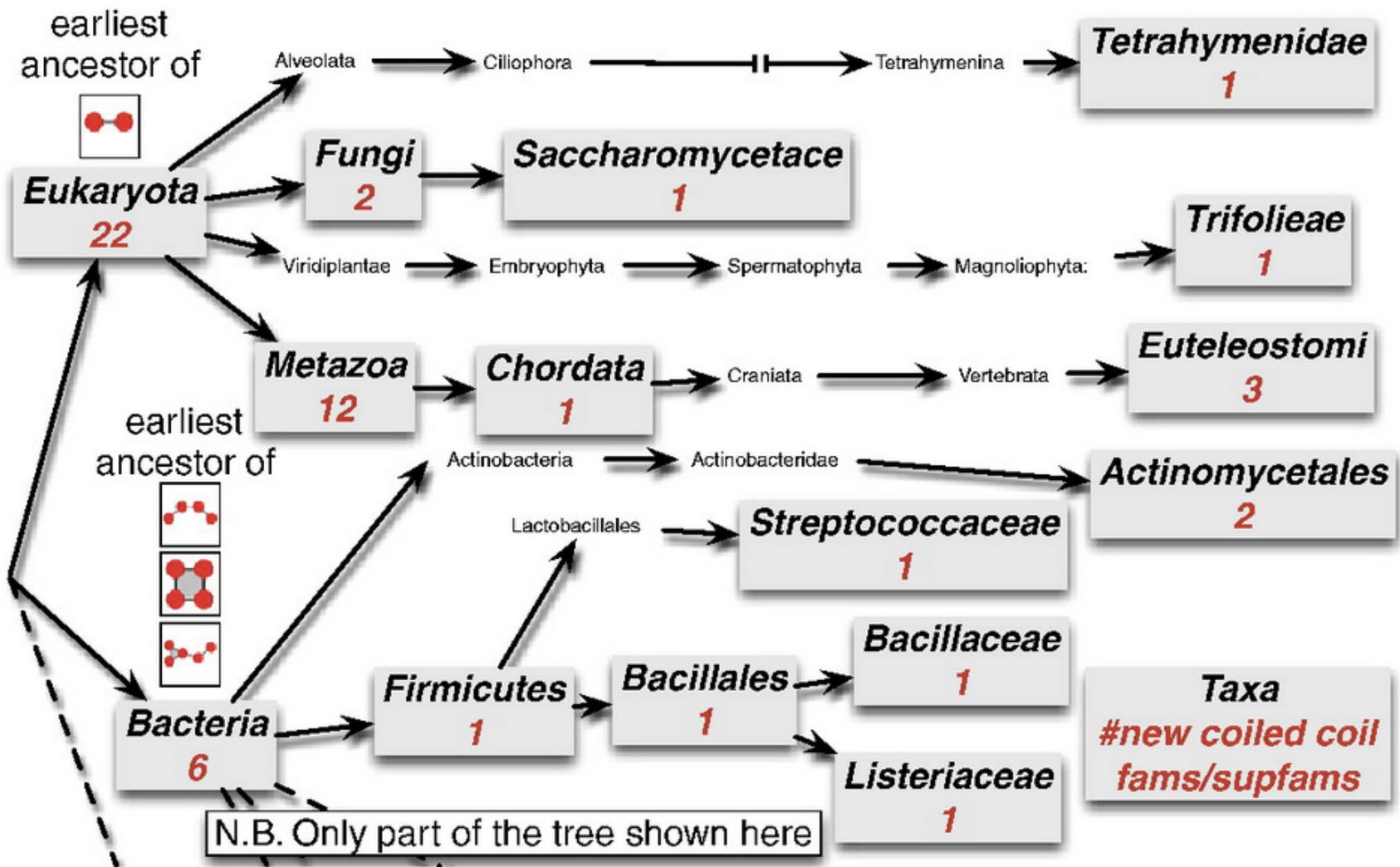


Phospholamban
1ZLL



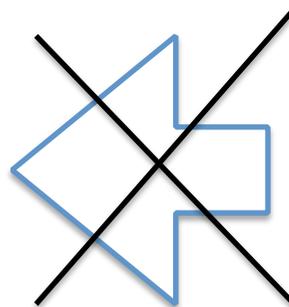
Cobalamin
Adenosyltransferase
2NT8

Coiled-coil evolution



Rackham et al. J Mol Biol. 2010

Homology



Structural similarity

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.

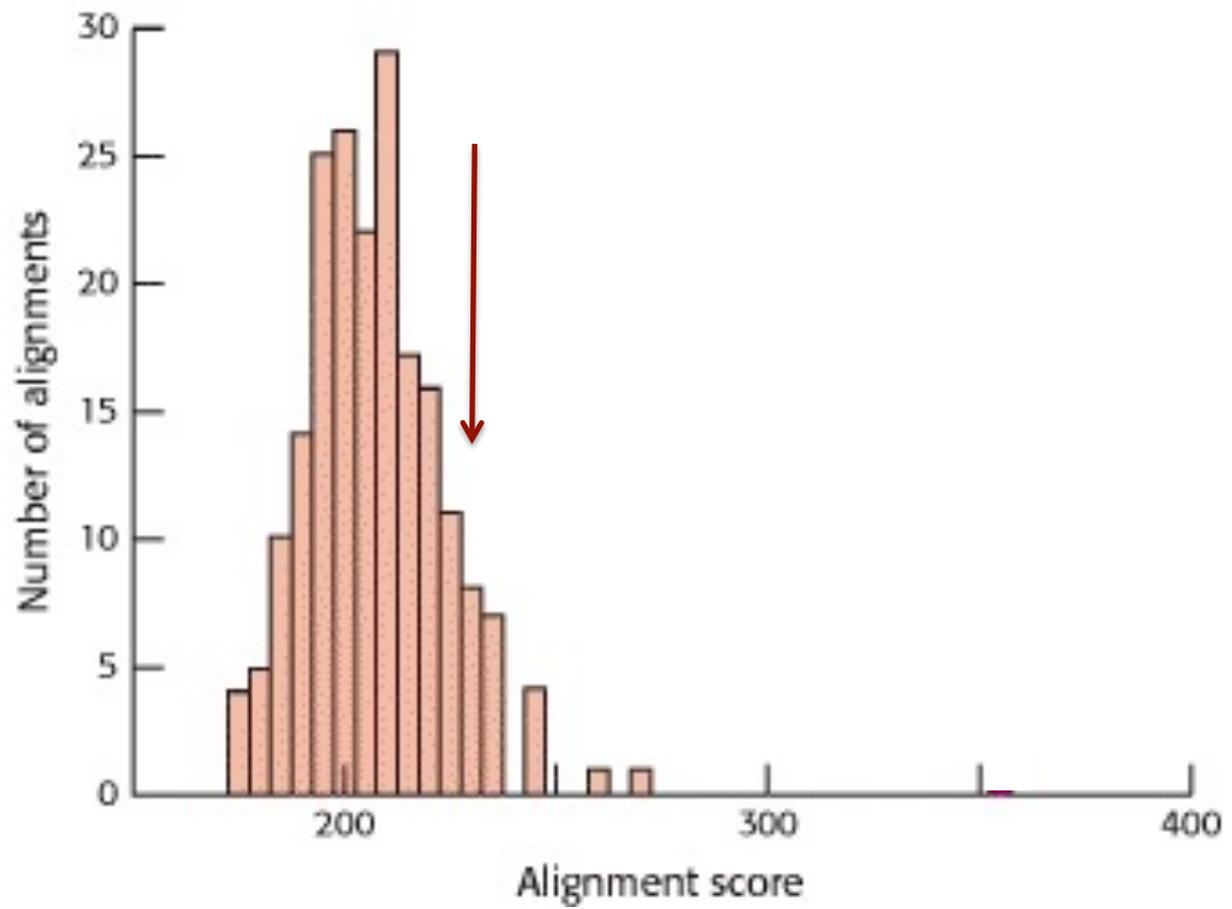
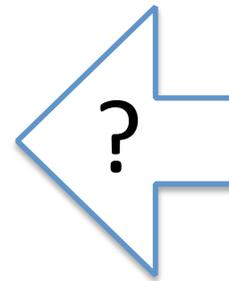


Image credits: <http://www.ncbi.nlm.nih.gov/books/NBK22456/figure/A937/?report=objectonly>

Excess sequence similarity

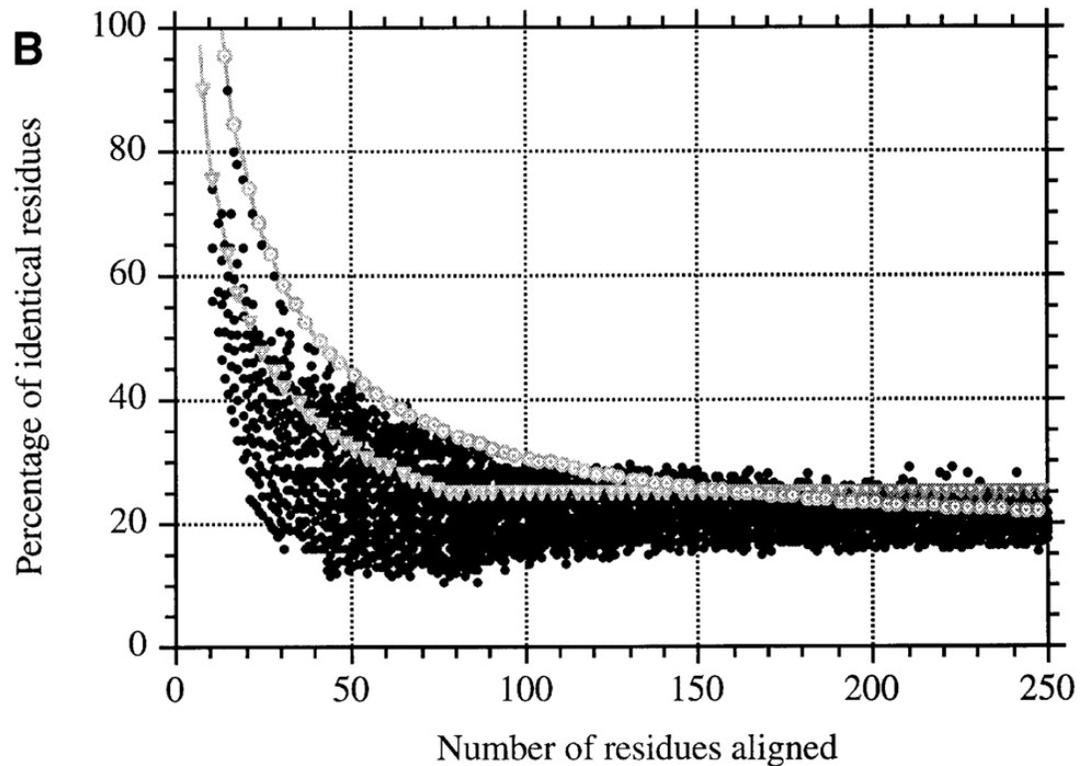


Homology

```

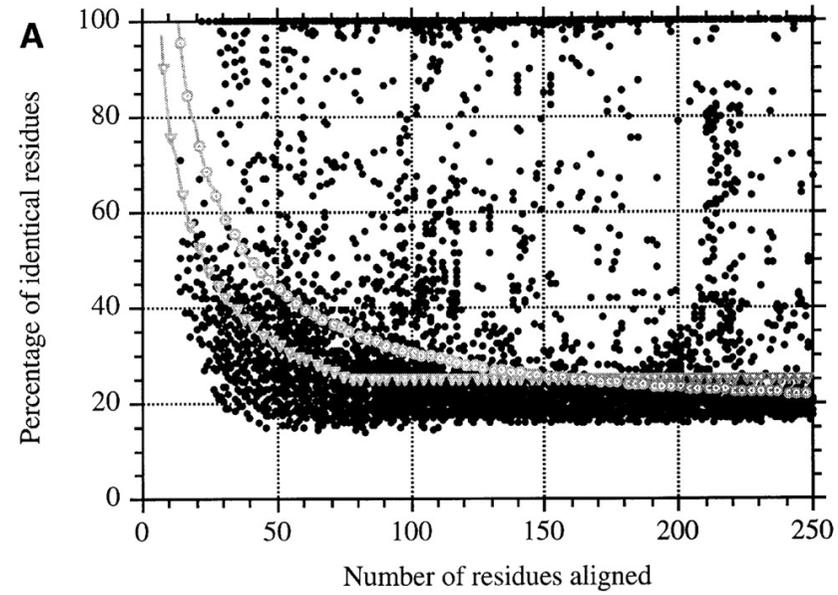
1 MGLSDGEWQLVLNVWGKVEADIPGHGQEVLI R L R L F K G H P E T L E K F D K F K H L K S E D E M K A S E 60
  M L S      V   W G K V A      G E L R F   P T   F F      D   S
1 MVLSPADKTNVKA AWGKVG AHAGEYGA EALERMFLSFPTTKTYFPHF-----DLSHGSA 54

61 DLKKHGATVLTALGGILKKKGHHEAEIKPLAQSHATKHKI-PVKY 104
   K H   V AL              L   HA K   PV
55 QVKGHSKKVADALTNAVAHVDDMPNALSALS D L H A H K L R V D P V N F 99
    
```

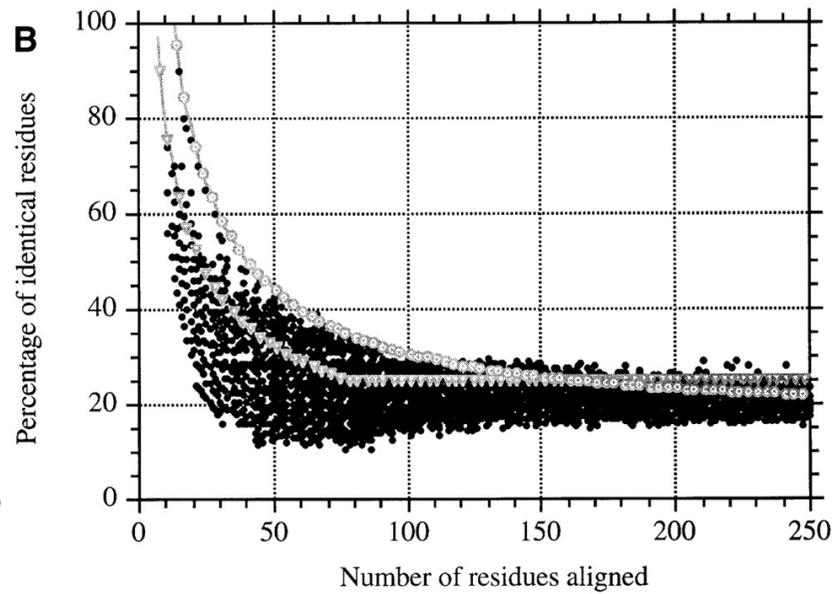


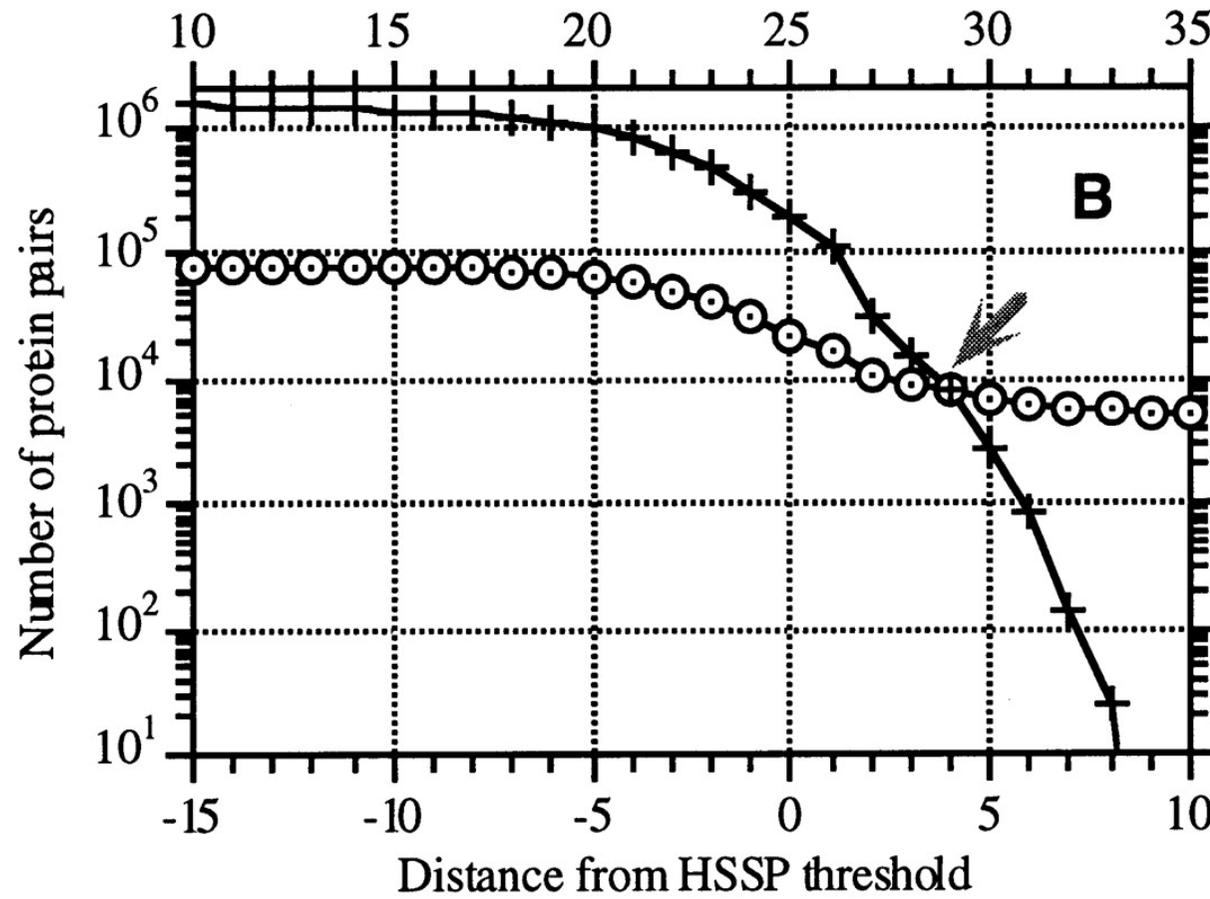
Unrelated proteins

Related proteins

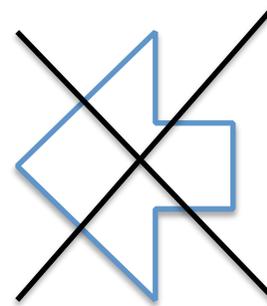


Unrelated proteins





Excess sequence similarity



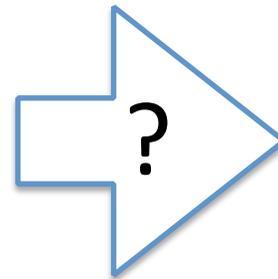
Homology

Homology \Leftrightarrow similar sequence?

Homology \Leftrightarrow similar structure?

Homology \Leftrightarrow similar function?

Homology



Similar function

The X-ray structure of a cobalamin biosynthetic enzyme, cobalt-precorrin-4 methyltransferase

Heidi L. Schubert¹, Keith S. Wilson¹, Evelyne Raux², Sarah C. Woodcock² and Martin J. Warren²

Biosynthesis of the corrin ring of vitamin B₁₂ requires the action of six S-adenosyl-L-methionine (AdoMet) dependent transmethyases, closely related in sequence. The first X-ray structure of one of these, cobalt-precorrin-4 transmethylase, CbiF, from *Bacillus megaterium* has been determined to a resolution of 2.4 Å. CbiF contains two α/β domains forming a trough in which S-adenosyl-L-homocysteine (AdoHcy) binds. The location of AdoHcy and a number of conserved residues, helps define the precorrin binding site. A second crystal form determined at 3.1 Å resolution highlights the flexibility of two loops around this site. CbiF employs a unique mode of AdoHcy binding and represents a new class of transmethylase.

The X-ray structure of a cobalamin biosynthetic enzyme, cobalt-precorrin-4 methyltransferase

Heidi L. Schubert¹, Keith S. Wilson¹, Evelyne Raux², Sarah C. Woodcock² and Martin J. Warren²

Biosynthesis of the corrin ring of vitamin B₁₂ requires the action of six S-adenosyl-L-methionine (AdoMet) dependent transmethyases, closely related in sequence. The first X-ray structure of one of these, cobalt-precorrin-4 transmethylase, CbiF, from *Bacillus megaterium* has been determined to a resolution of 2.4 Å. CbiF contains two α/β domains forming a trough in which S-adenosyl-L-homocysteine (AdoHcy) binds. The location of AdoHcy and a number of conserved residues, helps define the precorrin binding site. A second crystal form determined at 3.1 Å resolution highlights the flexibility of two loops around this site. CbiF employs a unique mode of AdoHcy binding and represents a new class of transmethylase.

The X-ray structure of a cobalamin biosynthetic enzyme, cobalt-precorrin-4 methyltransferase

Heidi L. Schubert¹, Keith S. Wilson¹, Evelyne Raux², Sarah C. Woodcock² and Martin J. Warren²

Biosynthesis of the corrin ring of vitamin B₁₂ requires the action of six S-adenosyl-L-methionine (AdoMet) dependent transmethyases, closely related in sequence. The first X-ray structure of one of these, cobalt-precorrin-4 transmethylase, CbiF, from *Bacillus megaterium* has been determined to a resolution of 2.4 Å. CbiF contains two α/β domains forming a trough in which S-adenosyl-L-homocysteine (AdoHcy) binds. The location of AdoHcy and a number of conserved residues, helps define the precorrin binding site. A second crystal form determined at 3.1 Å resolution highlights the flexibility of two loops around this site. CbiF employs a unique mode of AdoHcy binding and represents a new class of transmethylase.

The X-ray structure of a cobalamin biosynthetic enzyme, cobalt-precorrin-4 methyltransferase

Heidi L. Schubert¹, Keith S. Wilson¹, Evelyne Raux², Sarah C. Woodcock² and Martin J. Warren²

Biosynthesis of the corrin ring of vitamin B₁₂ requires the action of six S-adenosyl-L-methionine (AdoMet) dependent transmethyases, closely related in sequence. The first X-ray structure of one of these, cobalt-precorrin-4 transmethylase, CbiF, from *Bacillus megaterium* has been determined to a resolution of 2.4 Å. CbiF contains two α/β domains forming a trough in which S-adenosyl-L-homocysteine (AdoHcy) binds. The location of AdoHcy and a number of conserved residues, helps define the precorrin binding site. A second crystal form determined at 3.1 Å resolution highlights the flexibility of two loops around this site. CbiF employs a unique mode of AdoHcy binding and represents a new class of transmethylase.

Protein function(s)

Marco Punta

CbiF

Catalytic activity



methyltransferase



precorrin-4
methyltransferase activity

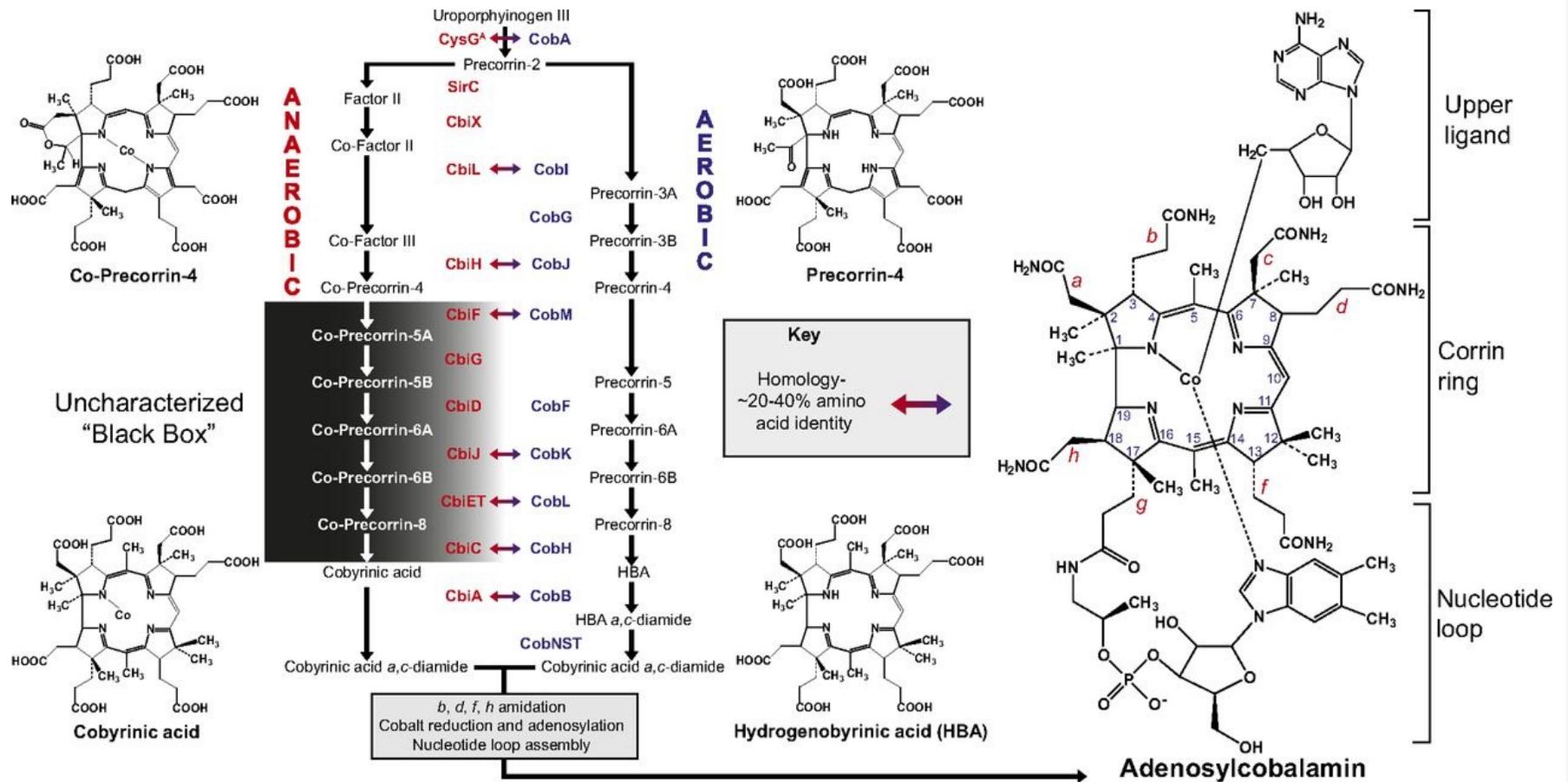
The X-ray structure of a cobalamin biosynthetic enzyme, cobalt-precorrin-4 methyltransferase

Heidi L. Schubert¹, Keith S. Wilson¹, Evelyne Raux², Sarah C. Woodcock² and Martin J. Warren²

Biosynthesis of the corrin ring of vitamin B₁₂ requires the action of six S-adenosyl-L-methionine (AdoMet) dependent transmethyases, closely related in sequence. The first X-ray structure of one of these, cobalt-precorrin-4 transmethylase, CbiF, from *Bacillus megaterium* has been determined to a resolution of 2.4 Å. CbiF contains two α/β domains forming a trough in which S-adenosyl-L-homocysteine (AdoHcy) binds. The location of AdoHcy and a number of conserved residues, helps define the precorrin binding site. A second crystal form determined at 3.1 Å resolution highlights the flexibility of two loops around this site. CbiF employs a unique mode of AdoHcy binding and represents a new class of transmethylase.

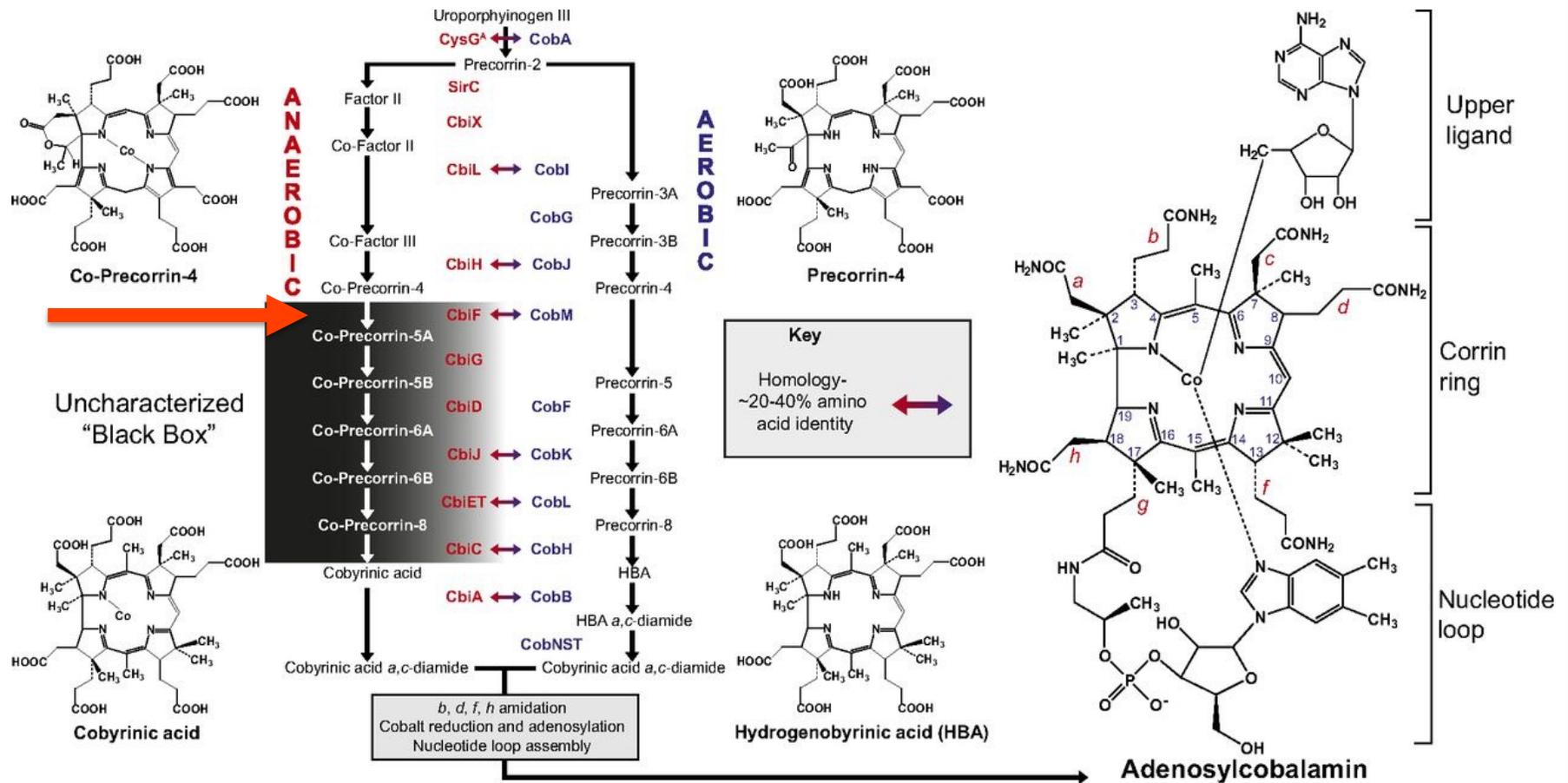
Cobalamin biosynthetic pathways

Marco Punta



Cobalamin biosynthetic pathways

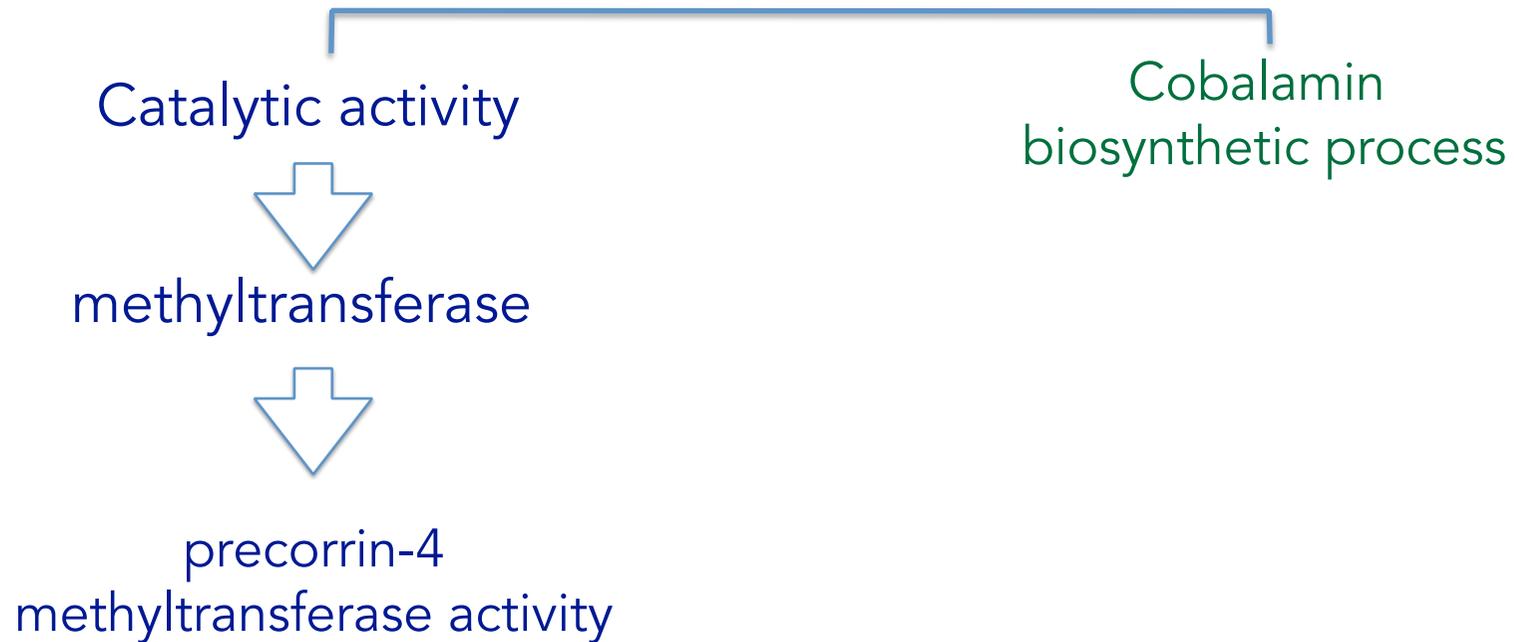
Marco Punta



Protein function(s)

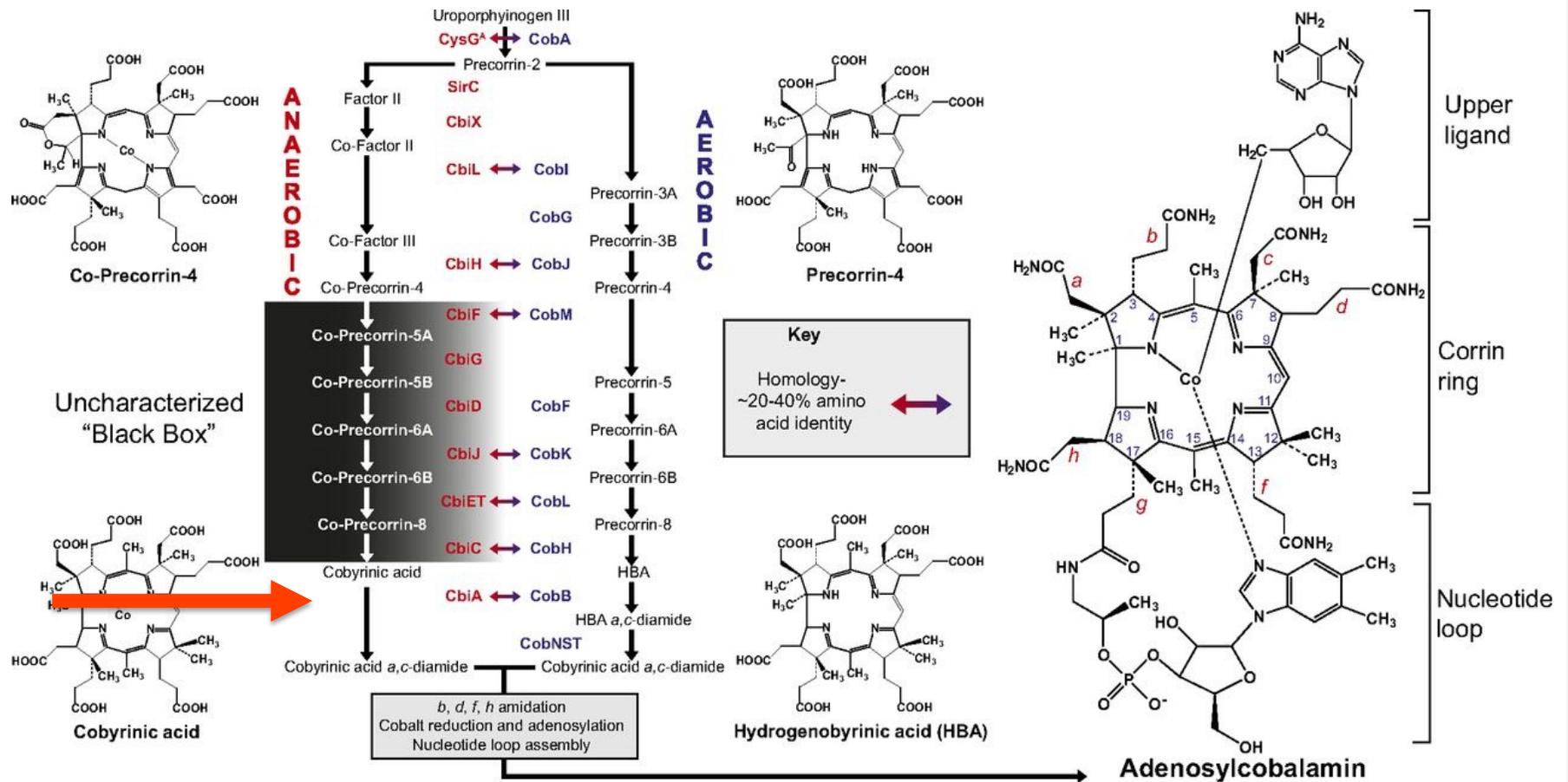
Marco Punta

CbiF



Cobalamin biosynthetic pathways

Marco Punta



Protein function(s)

CbiF

Catalytic activity



methyltransferase



precorrin-4
methyltransferase activity

Cobalamin
biosynthetic process

CbiA

Catalytic activity



hydrolase activity



cobyric acid a,c-diamide
synthase activity

Protein function(s)

Molecular function

CbiF

Biological process

Catalytic activity



methyltransferase



precorrin-4
methyltransferase activity

Cobalamin
biosynthetic process

CbiA

Catalytic activity



hydrolase activity



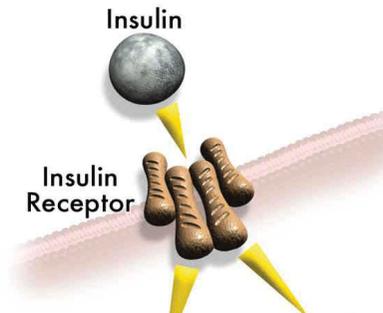
cobyric acid a,c-diamide
synthase activity

Molecular function

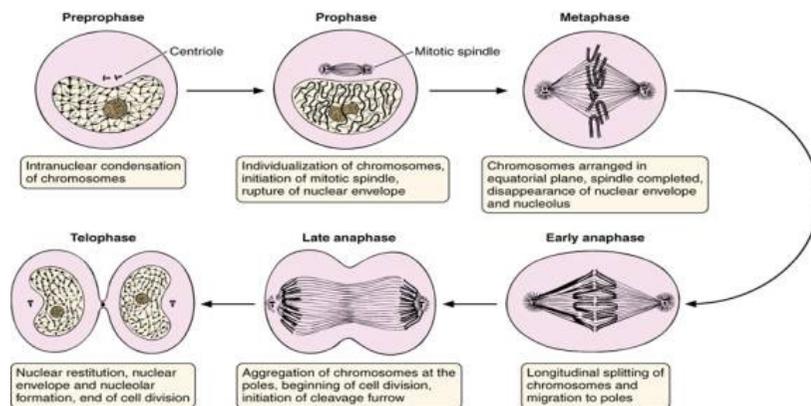
GO: 3 ontologies in 1

1. Molecular Function

An elemental activity or task or job



- protein kinase activity
- insulin receptor activity



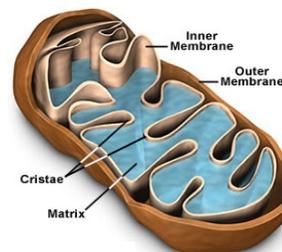
2. Biological Process

A commonly recognised series of events

- cell division

3. Cellular Component

Where a gene product is located



- mitochondrion
- mitochondrial matrix
- mitochondrial inner membrane

Bacillus megaterium CbiF GO annotation

Marco Punta

Database	Gene Product ID	Symbol	Qualifier	GO Identifier	GO Term Name	Aspect	Evidence	Reference	With	Taxon	Date	Assigned By	Pro Fo ID
Process													
UniProtKB	O87696	cbiF		GO:0006779	porphyrin-containing compound biosynthetic process	P	IEA	InterPro2GO	InterPro:IPR003043	1404	20151010	InterPro	
UniProtKB	O87696	cbiF		GO:0008152	metabolic process	P	IEA	InterPro2GO	InterPro:IPR000878 InterPro:IPR014776 InterPro:IPR014777	1404	20151010	InterPro	
UniProtKB	O87696	cbiF		GO:0009236	cobalamin biosynthetic process	P	IEA	InterPro2GO	InterPro:IPR006362	1404	20151010	InterPro	
UniProtKB	O87696	cbiF		GO:0009236	cobalamin biosynthetic process	P	IEA	UniProt Keywords2GO (UniProtKB/Swiss-Prot entries)	UniProtKB-KW:KW-0169	1404	20151010	UniProt	
UniProtKB	O87696	cbiF		GO:0009236	cobalamin biosynthetic process	P	IEA	UniPathway2GO	UniPathway:UPA00148	1404	20151010	UniProt	
UniProtKB	O87696	cbiF		GO:0032259	methylation	P	IEA	UniProt Keywords2GO (UniProtKB/Swiss-Prot entries)	UniProtKB-KW:KW-0489	1404	20151010	UniProt	
UniProtKB	O87696	cbiF		GO:0055114	oxidation-reduction process	P	IEA	InterPro2GO	InterPro:IPR003043	1404	20151010	InterPro	
Function													
UniProtKB	O87696	cbiF		GO:0008168	methyltransferase activity	F	IEA	InterPro2GO	InterPro:IPR000878 InterPro:IPR003043 InterPro:IPR014776 InterPro:IPR014777	1404	20151010	InterPro	
UniProtKB	O87696	cbiF		GO:0008168	methyltransferase activity	F	IEA	UniProt Keywords2GO (UniProtKB/Swiss-Prot entries)	UniProtKB-KW:KW-0489	1404	20151010	UniProt	
UniProtKB	O87696	cbiF		GO:0016740	transferase activity	F	IEA	UniProt Keywords2GO (UniProtKB/Swiss-Prot entries)	UniProtKB-KW:KW-0808	1404	20151010	UniProt	
UniProtKB	O87696	cbiF		GO:0043115	precorrin-2 dehydrogenase activity	F	IEA	InterPro2GO	InterPro:IPR003043	1404	20151010	InterPro	
UniProtKB	O87696	cbiF		GO:0046026	precorrin-4 C11-methyltransferase activity	F	IEA	InterPro2GO	InterPro:IPR006362	1404	20151010	InterPro	

Top-level EC numbers^[9]

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

https://en.wikipedia.org/wiki/Enzyme_Commission_number

EC 1 **Oxidoreductases**
EC 1.3 **Acting on the CH-CH Group of Donors**
EC 1.3.1 **With NAD⁺ or NADP⁺ as acceptor**
EC 1.3.1.21 **7-dehydrocholesterol reductase**

IUBMB Enzyme Nomenclature

EC 1.3.1.21

Database	Gene Product ID	Symbol	Qualifier	GO Identifier	GO Term Name	Aspect	Evidence	Reference	With	Taxon	Date	Assigned By	Product Form ID
Process													
JniProtKB	P02144	MB		GO:0001666	response to hypoxia	P	IEA	Ensembl Compara	Ensembl:ENSMUSP00000125995	9606	20140913	Ensembl	
JniProtKB	P02144	MB		GO:0006810	transport	P	IEA	UniProt Keywords2GO (UniProtKB/Swiss-Prot entries)	UniProtKB-KW:KW-0813	9606	20140913	UniProt	
JniProtKB	P02144	MB		GO:0007507	heart development	P	IEA	Ensembl Compara	Ensembl:ENSMUSP00000125995	9606	20140913	Ensembl	
JniProtKB	P02144	MB		GO:0009725	response to hormone	P	IEA	Ensembl Compara	Ensembl:ENSRNOP00000006184	9606	20140913	Ensembl	
JniProtKB	P02144	MB		GO:0015671	oxygen transport	P	IEA	InterPro2GO	InterPro:IPR002335 InterPro:IPR012292	9606	20140913	InterPro	
JniProtKB	P02144	MB		GO:0015671	oxygen transport	P	IEA	Ensembl Compara	Ensembl:ENSRNOP00000006184	9606	20140913	Ensembl	
JniProtKB	P02144	MB		GO:0015671	oxygen transport	P	IEA	UniProt Keywords2GO (UniProtKB/Swiss-Prot entries)	UniProtKB-KW:KW-0561	9606	20140913	UniProt	
JniProtKB	P02144	MB		GO:0031444	slow-twitch skeletal muscle fiber contraction	P	IEA	Ensembl Compara	Ensembl:ENSRNOP00000006184	9606	20140913	Ensembl	
JniProtKB	P02144	MB		GO:0042542	response to hydrogen peroxide	P	IEA	Ensembl Compara	Ensembl:ENSRNOP00000006184	9606	20140913	Ensembl	
JniProtKB	P02144	MB		GO:0043353	enucleate erythrocyte differentiation	P	IEA	Ensembl Compara	Ensembl:ENSMUSP00000125995	9606	20140913	Ensembl	
JniProtKB	P02144	MB		GO:0050873	brown fat cell differentiation	P	IEA	Ensembl Compara	Ensembl:ENSMUSP00000125995	9606	20140913	Ensembl	
Function													
JniProtKB	P02144	MB		GO:0005344	oxygen transporter activity	F	IEA	UniProt Keywords2GO (UniProtKB/Swiss-Prot entries)	UniProtKB-KW:KW-0561	9606	20140913	UniProt	
JniProtKB	P02144	MB		GO:0005506	iron ion binding	F	IEA	InterPro2GO	InterPro:IPR000971 InterPro:IPR012292	9606	20140913	InterPro	
JniProtKB	P02144	MB		GO:0019825	oxygen binding	F	IEA	InterPro2GO	InterPro:IPR002335 InterPro:IPR012292	9606	20140913	InterPro	
JniProtKB	P02144	MB		GO:0019825	oxygen binding	F	IEA	Ensembl Compara	Ensembl:ENSRNOP00000006184	9606	20140913	Ensembl	
JniProtKB	P02144	MB		GO:0020037	heme binding	F	IEA	InterPro2GO	InterPro:IPR000971 InterPro:IPR002335 InterPro:IPR012292	9606	20140913	InterPro	
JniProtKB	P02144	MB		GO:0046872	metal ion binding	F	IEA	UniProt Keywords2GO (UniProtKB/Swiss-Prot entries)	UniProtKB-KW:KW-0479	9606	20140913	UniProt	
Component													
JniProtKB	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													
JniProtKB	P02008	HBZ		GO:000122	negative regulation of transcription from RNA polymerase II promoter	P	IEA	Ensembl Compara	Ensembl:ENSMUSP00000020531	9606	20140913	Ensembl	
JniProtKB	P02008	HBZ		GO:0006810	transport	P	IEA	UniProt Keywords2GO (UniProtKB/Swiss-Prot entries)	UniProtKB-KW:KW-0813	9606	20140913	UniProt	
JniProtKB	P02008	HBZ		GO:0015671	oxygen transport	P	IEA	InterPro2GO	InterPro:IPR002338 InterPro:IPR002340 InterPro:IPR012292	9606	20140913	InterPro	
JniProtKB	P02008	HBZ		GO:0015671	oxygen transport	P	IEA	UniProt Keywords2GO (UniProtKB/Swiss-Prot entries)	UniProtKB-KW:KW-0561	9606	20140913	UniProt	
JniProtKB	P02008	HBZ		GO:0043249	erythrocyte maturation	P	IEA	Ensembl Compara	Ensembl:ENSMUSP00000020531	9606	20140913	Ensembl	
Function													
JniProtKB	P02008	HBZ		GO:0005344	oxygen transporter activity	F	IEA	UniProt Keywords2GO (UniProtKB/Swiss-Prot entries)	UniProtKB-KW:KW-0561	9606	20140913	UniProt	
JniProtKB	P02008	HBZ		GO:0005344	oxygen transporter activity	F	TAS	PMID:7555018		9606	20030904	PINC	
JniProtKB	P02008	HBZ		GO:0005506	iron ion binding	F	IEA	InterPro2GO	InterPro:IPR000971 InterPro:IPR002338 InterPro:IPR002340 InterPro:IPR012292	9606	20140913	InterPro	
JniProtKB	P02008	HBZ		GO:0005515	protein binding	F	IPI	PMID:11159543	UniProtKB:P68871	9606	20140914	IntAct	
JniProtKB	P02008	HBZ		GO:0005515	protein binding	F	IPI	PMID:6683087	UniProtKB:P68871	9606	20140914	IntAct	
JniProtKB	P02008	HBZ		GO:0019825	oxygen binding	F	IEA	InterPro2GO	InterPro:IPR002338 InterPro:IPR002340 InterPro:IPR012292	9606	20140913	InterPro	
JniProtKB	P02008	HBZ		GO:0020037	heme binding	F	IEA	InterPro2GO	InterPro:IPR000971 InterPro:IPR002338 InterPro:IPR002340 InterPro:IPR012292	9606	20140913	InterPro	
JniProtKB	P02008	HBZ		GO:0046872	metal ion binding	F	IEA	UniProt Keywords2GO (UniProtKB/Swiss-Prot entries)	UniProtKB-KW:KW-0479	9606	20140913	UniProt	
Component													
JniProtKB	P02008	HBZ		GO:0005833	hemoglobin complex	C	IEA	InterPro2GO	InterPro:IPR002338 InterPro:IPR002340	9606	20140913	InterPro	
JniProtKB	P02008	HBZ		GO:0005833	hemoglobin complex	C	TAS	PMID:7555018		9606	20030904	PINC	
JniProtKB	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													
JniProtKB	P02144	MB		GO:0001666	response to hypoxia	P	IEA	Ensembl Compara	Ensembl:ENSMUSP00000125995	9606	20140913	Ensembl	
JniProtKB	P02144	MB		GO:0006810	transport	P	IEA	UniProt Keywords2GO (UniProtKB/Swiss-Prot entries)	UniProtKB-KW:KW-0813	9606	20140913	UniProt	
JniProtKB	P02144	MB		GO:0007507	heart development	P	IEA	Ensembl Compara	Ensembl:ENSMUSP00000125995	9606	20140913	Ensembl	
JniProtKB	P02144	MB		GO:0009725	response to hormone	P	IEA	Ensembl Compara	Ensembl:ENSRNOP00000006184	9606	20140913	Ensembl	
JniProtKB	P02144	MB		GO:0015671	oxygen transport	P	IEA	InterPro2GO	InterPro:IPR002335 InterPro:IPR012292	9606	20140913	InterPro	
JniProtKB	P02144	MB		GO:0015671	oxygen transport	P	IEA	Ensembl Compara	Ensembl:ENSRNOP00000006184	9606	20140913	Ensembl	
JniProtKB	P02144	MB		GO:0015671	oxygen transport	P	IEA	UniProt Keywords2GO (UniProtKB/Swiss-Prot entries)	UniProtKB-KW:KW-0561	9606	20140913	UniProt	
JniProtKB	P02144	MB		GO:0031444	slow-twitch skeletal muscle fiber contraction	P	IEA	Ensembl Compara	Ensembl:ENSRNOP00000006184	9606	20140913	Ensembl	
JniProtKB	P02144	MB		GO:0042542	response to hydrogen peroxide	P	IEA	Ensembl Compara	Ensembl:ENSRNOP00000006184	9606	20140913	Ensembl	
JniProtKB	P02144	MB		GO:0043353	enucleate erythrocyte differentiation	P	IEA	Ensembl Compara	Ensembl:ENSMUSP00000125995	9606	20140913	Ensembl	
JniProtKB	P02144	MB		GO:0050873	brown fat cell differentiation	P	IEA	Ensembl Compara	Ensembl:ENSMUSP00000125995	9606	20140913	Ensembl	
Function													
JniProtKB	P02144	MB		GO:0005344	oxygen transporter activity	F	IEA	UniProt Keywords2GO (UniProtKB/Swiss-Prot entries)	UniProtKB-KW:KW-0561	9606	20140913	UniProt	
JniProtKB	P02144	MB		GO:0005506	iron ion binding	F	IEA	InterPro2GO	InterPro:IPR000971 InterPro:IPR012292	9606	20140913	InterPro	
JniProtKB	P02144	MB		GO:0019825	oxygen binding	F	IEA	InterPro2GO	InterPro:IPR002335 InterPro:IPR012292	9606	20140913	InterPro	
JniProtKB	P02144	MB		GO:0019825	oxygen binding	F	IEA	Ensembl Compara	Ensembl:ENSRNOP00000006184	9606	20140913	Ensembl	
JniProtKB	P02144	MB		GO:0020037	heme binding	F	IEA	InterPro2GO	InterPro:IPR000971 InterPro:IPR002335 InterPro:IPR012292	9606	20140913	InterPro	
JniProtKB	P02144	MB		GO:0046872	metal ion binding	F	IEA	UniProt Keywords2GO (UniProtKB/Swiss-Prot entries)	UniProtKB-KW:KW-0479	9606	20140913	UniProt	
Component													
JniProtKB	P02144	MB		GO:0070062	extracellular vesicular exosome	C	IDA	PMID:23533145		9606	20140714	UniProt	

SAME

SAME

Database	Gene Product ID	Symbol	Qualifier	GO Identifier	GO Term Name	Aspect	Evidence	Reference	With	Taxon	Date	Assigned By	Product Form ID
Process													
JniProtKB	P02008	HBZ		GO:000122	negative regulation of transcription from RNA polymerase II promoter	P	IEA	Ensembl Compara	Ensembl:ENSMUSP00000020531	9606	20140913	Ensembl	
JniProtKB	P02008	HBZ		GO:0006810	transport	P	IEA	UniProt Keywords2GO (UniProtKB/Swiss-Prot entries)	UniProtKB-KW:KW-0813	9606	20140913	UniProt	
JniProtKB	P02008	HBZ		GO:0015671	oxygen transport	P	IEA	InterPro2GO	InterPro:IPR002338 InterPro:IPR002340 InterPro:IPR012292	9606	20140913	InterPro	
JniProtKB	P02008	HBZ		GO:0015671	oxygen transport	P	IEA	UniProt Keywords2GO (UniProtKB/Swiss-Prot entries)	UniProtKB-KW:KW-0561	9606	20140913	UniProt	
JniProtKB	P02008	HBZ		GO:0043249	erythrocyte maturation	P	IEA	Ensembl Compara	Ensembl:ENSMUSP00000020531	9606	20140913	Ensembl	
Function													
JniProtKB	P02008	HBZ		GO:0005344	oxygen transporter activity	F	IEA	UniProt Keywords2GO (UniProtKB/Swiss-Prot entries)	UniProtKB-KW:KW-0561	9606	20140913	UniProt	
JniProtKB	P02008	HBZ		GO:0005344	oxygen transporter activity	F	TAS	PMID:7555018		9606	20030904	PINC	
JniProtKB	P02008	HBZ		GO:0005506	iron ion binding	F	IEA	InterPro2GO	InterPro:IPR000971 InterPro:IPR002338 InterPro:IPR002340 InterPro:IPR012292	9606	20140913	InterPro	
JniProtKB	P02008	HBZ		GO:0005515	protein binding	F	IPI	PMID:11159	UniProtKB:P68871	9606	20140914	IntAct	
JniProtKB	P02008	HBZ		GO:0005515	protein binding	F	IPI	PMID:11159	UniProtKB:P68871	9606	20140914	IntAct	
JniProtKB	P02008	HBZ		GO:0019825	oxygen binding	F	IEA	InterPro2GO	InterPro:IPR002338 InterPro:IPR002340 InterPro:IPR012292	9606	20140913	InterPro	
JniProtKB	P02008	HBZ		GO:0020037	heme binding	F	IEA	InterPro2GO	InterPro:IPR000971 InterPro:IPR002338 InterPro:IPR002340 InterPro:IPR012292	9606	20140913	InterPro	
JniProtKB	P02008	HBZ		GO:0046872	metal ion binding	F	IEA	UniProt Keywords2GO (UniProtKB/Swiss-Prot entries)	UniProtKB-KW:KW-0479	9606	20140913	UniProt	
Component													
JniProtKB	P02008	HBZ		GO:0005833	hemoglobin complex	C	IEA	InterPro2GO	InterPro:IPR002338 InterPro:IPR002340	9606	20140913	InterPro	
JniProtKB	P02008	HBZ		GO:0005833	hemoglobin complex	C	TAS	PMID:7555018		9606	20030904	PINC	
JniProtKB	P02008	HBZ		GO:0070062	extracellular vesicular exosome	C	IDA	PMID:23533145		9606	20140714	UniProt	

SAME

SAME

Functional property to be conserved	Sequence identity	Conservation rate	Reference
Non-enzyme	50%	98%*	[88]
All 4 EC numbers	70%**	90%	[89]
All 4 EC numbers	40%**	70%	[89]
First 3 EC numbers	50%**	90%	[89]
First 3 EC numbers	30%**	70%	[89]
All 4 EC numbers	50%	30%	[90]
First 3 EC numbers	25%	70%	[91]
SWISS-PROT keywords	40%	70%	[92]
Subcellular localization (11 classes)	70%	90%	[93]

*98% of non enzymes that have at least one enzyme homolog.

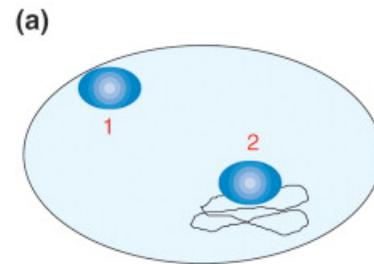
**Global identity, defined in [89].

Note: different estimates for the same functional aspects reflect the different methods, procedures, and datasets used to assess sequence similarity by the various groups.

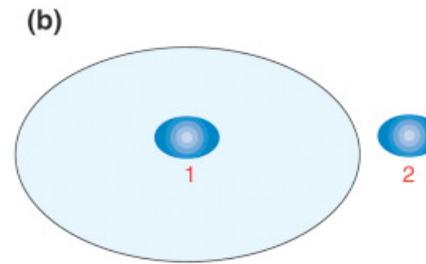
doi:10.1371/journal.pcbi.1000160.t001

Moonlighting proteins

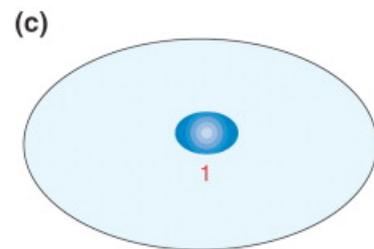
Marco Punta



Different locations within the cell



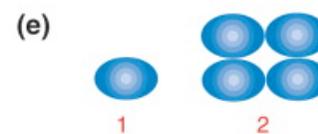
Inside and outside the cell



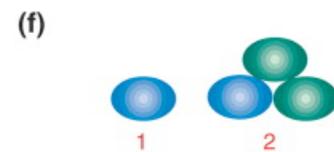
Expression by different cell types



Binding of a cofactor



Oligomerization



Complex formation



Multiple binding sites

Do homologous protein regions perform a similar function?

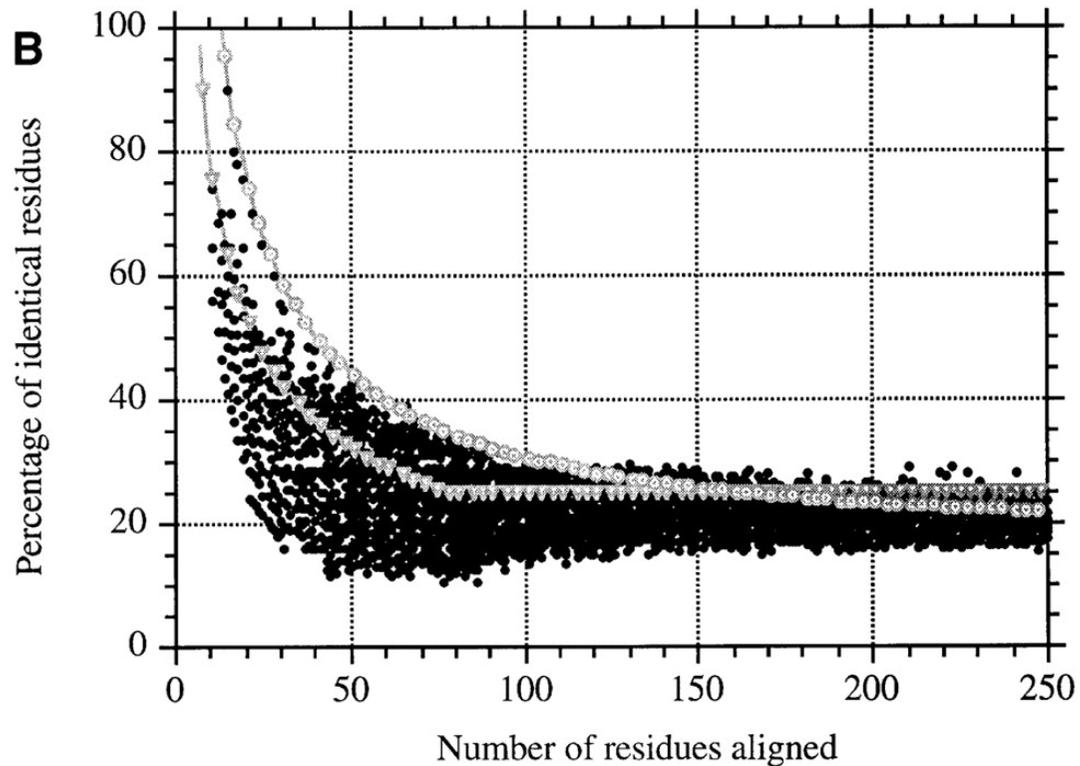
Homologous proteins may share a number of functional features, however:

- functional drift can lead to different functions or aspects of function
- while functional similarity generally correlates with evolutionary distance, no distance is safe for inferring function (very closely related proteins can have slightly to radically different functions)

```

1 MGLSDGEWQLVLNVWGKVEADIPGHGQEVLIIRLFKGGHPETLEKFDKFKHLKSEDEMKASE 60
  M LS      V  WGKV A      G E L R F  P T  F F      D  S
1 MVLSPADKTNVKAAGKVGGAHAGEYGAEALERMFLSFPTTKTYFPHF-----DLSHGSA 54

61 DLKKHGATVLTALGGILKKKGHHEAEIKPLAQSHATKHKI-PVKY 104
   K H  V AL                      L  HA K  PV
55 QVKGHSKKVADALTNAVAHVDDMPNALSALSADLHAHKLRVDPVNF 99
    
```



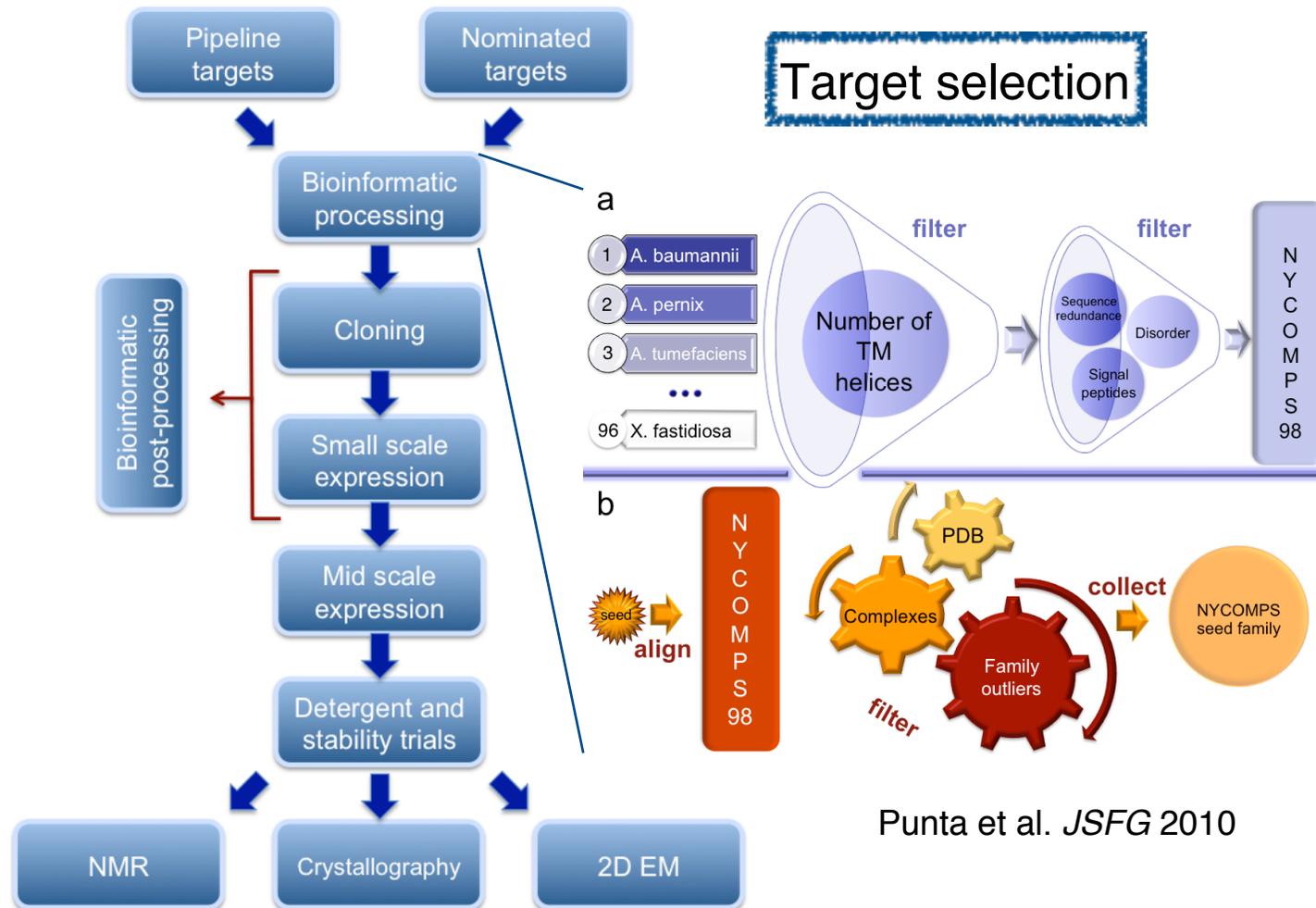
Unrelated proteins

Exercise

Homology-based function annotation transfer #1

NYCOMPS pipeline

Marco Punta

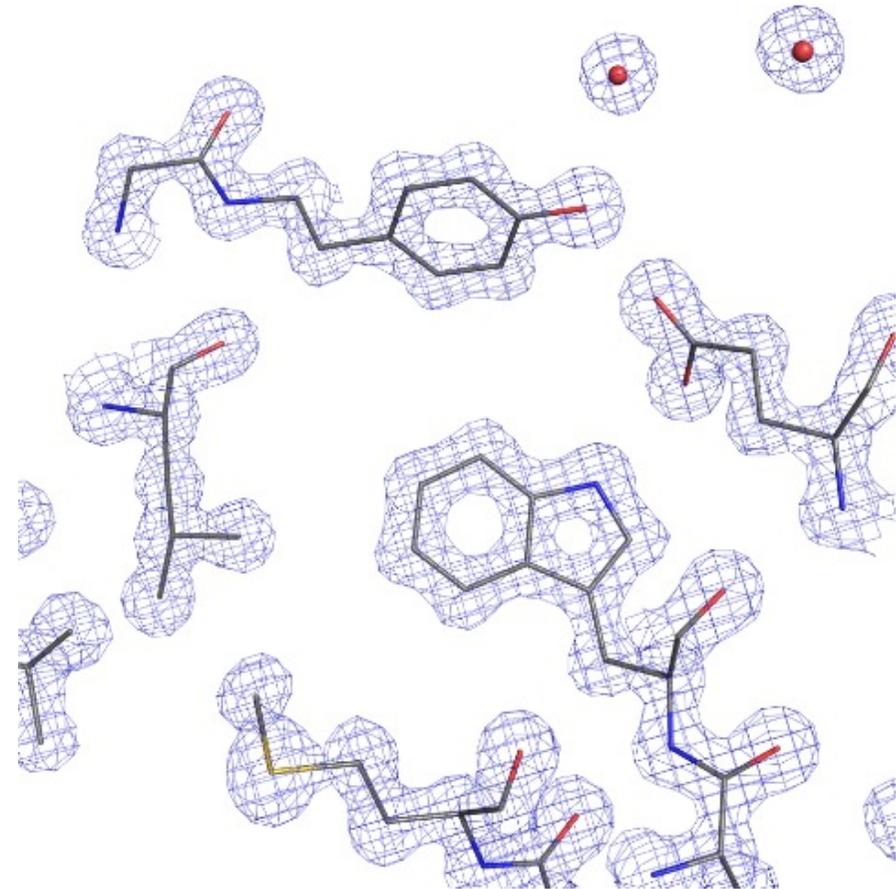
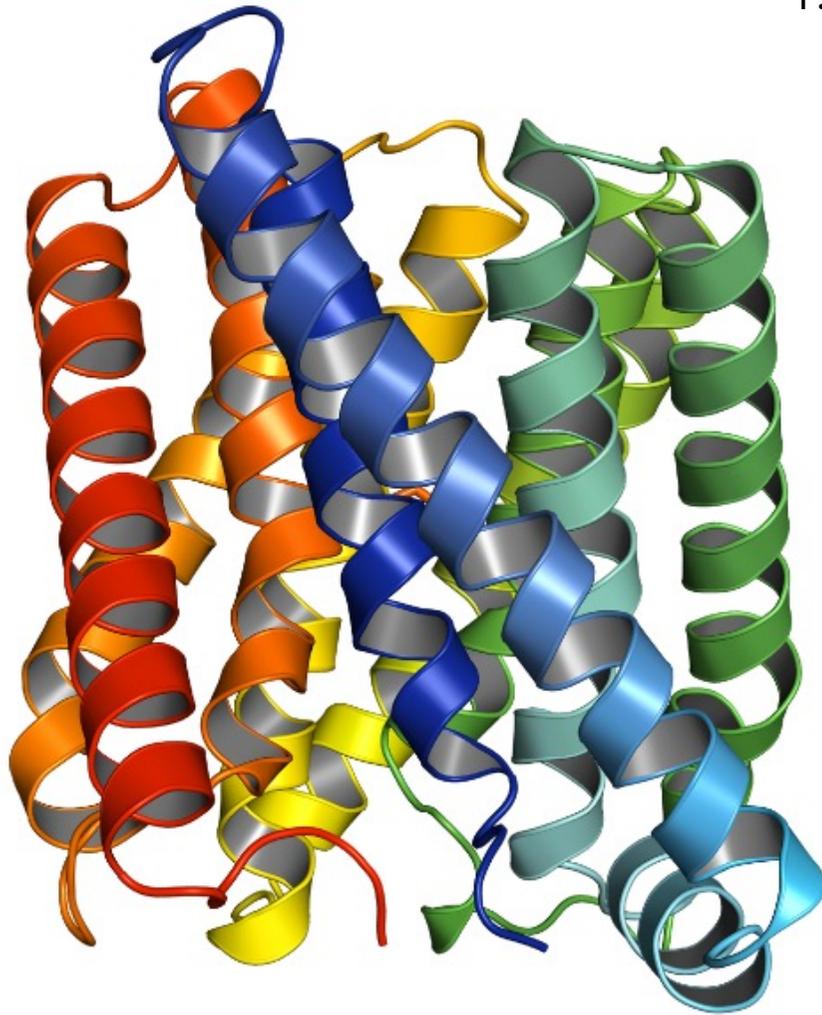


Love et al. *JSFG* 2010

Punta et al. *JSFG* 2010

H. influenzae protein [3M71] ← PDB id
1.20 Å

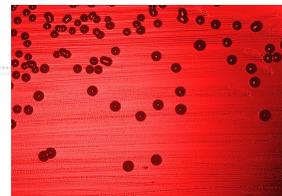
Marco Punta



Chen et al. *Nature* 467 (2010)

Alignment

Q9LD83	SLAC1_ARATH - Guard cell S-type anion channel SLA... - Arabidopsis thal...		
E-value: 3e-10	Positives : 41.0%		
Score: 160	Query Length: 328		
Ident.: 22.0%	Match Length: 556		
P44741	20	PFPL--PTGYFGIPLGLAALS LawFHLE-----NLFPARMVSDVLGIVASAVWILFILM	72
		PF L P G FGI LGL++ ++ W L N +++ V+ + + V +	
Q9LD83	183	PFLFRFPIGCFGICLGLSSQAVLWLALAKSPATNFLHITPLINLVVWLFSLVVLVSVSFT	242
P44741	73	YAYKLRYYFEEVRAEYHSPVRFSFIALIPITTMLVG---DILYRWNPLIAEVLWIGTIG	129
		Y K +YFE V+ EY PVR +F + M + ++ N IW +G	
Q9LD83	243	YILKCIFYFEAVKREYFHPVRVNFVFFAPWVCMFLAISVPPMFSNPKYLHPAIWCVFMG	302
P44741	130	QLLFSTLRVSELWQGGVFEQ--KSTHPSFYLPVAANFTSASSLALLGYHDLGYLFFGAG	187
		F L++ W G + K +PS +L +V NF A + +G+ ++ + G	
Q9LD83	303	PYFFLELKIYGQWLSGGKRRRLCKVANPSSHL-SVVGNFVGAILASKVGWDEVAKFLWAVG	361
P44741	188	MIAWIIFEPVLLQHLRISSLEPQFRATMGIVLAPAFVCSAYLSINHGEVDTLAKILWGY	247
		+++ L Q L S P+ + + A S + +G+ D ++ +	
Q9LD83	362	FAHYLVVFVTLYQRLPTSEALPKELHPVYSMFIAAPSAASIAWNTIYGQFDGCSRTCFFI	421
P44741	248	GFLQLFLLRLFPWIVEKGLNIGLWAFSFC LASMANSATAFY----HGNVLQGVSI FAFV	303
		L+ + ++ W++ F + + A+ AT Y G + +++	
Q9LD83	422	ALFLYISLVARINFFTGFKFSVAWWSYTFMTT-ASVATIKYAEAVPGYPSRALALTL SF	480
P44741	304	FSNVMIGLLVLMTI 317	
		S M+ +L + T+	
Q9LD83	481	ISTAMVCVLFVSTL 494	

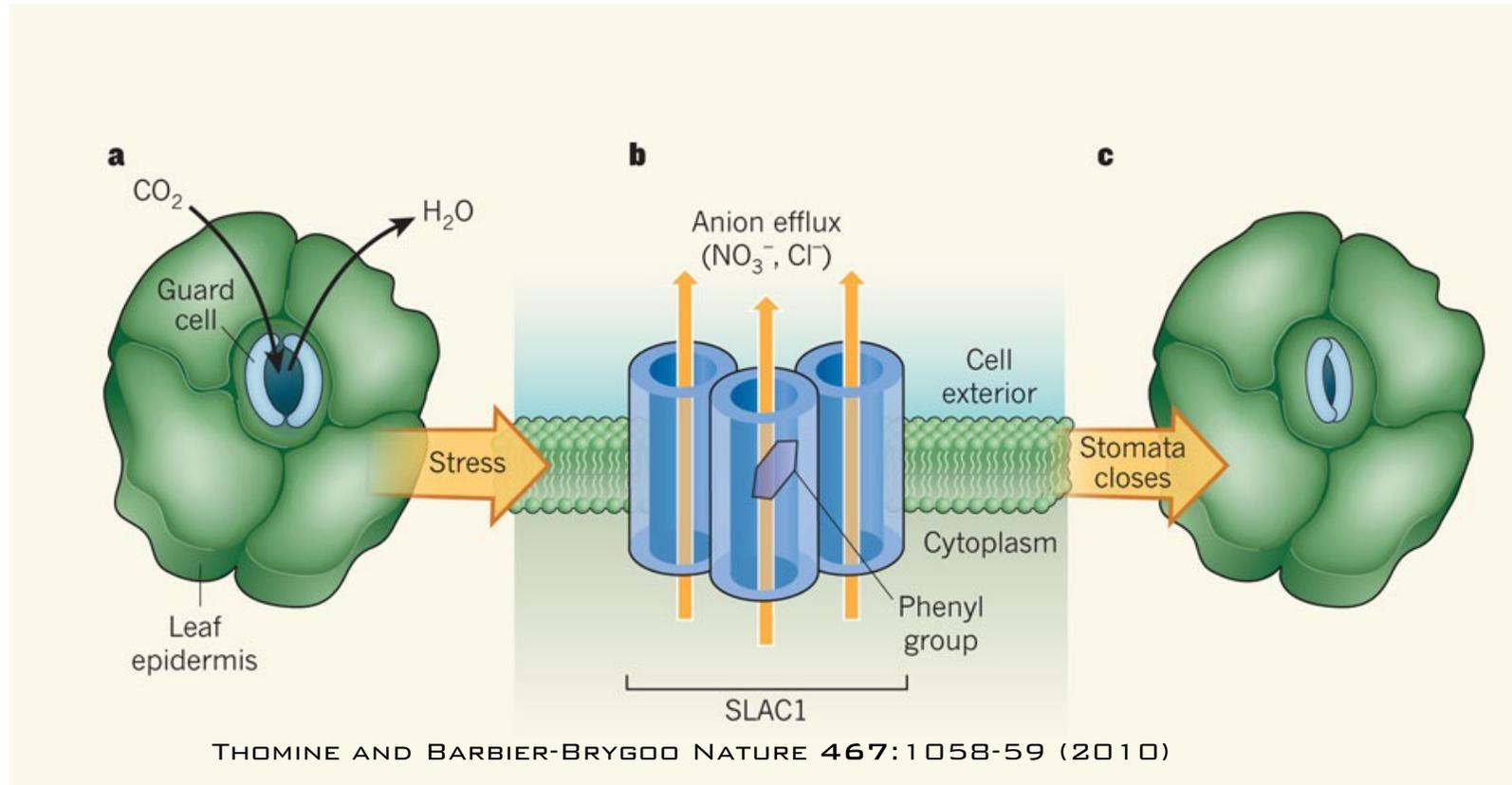


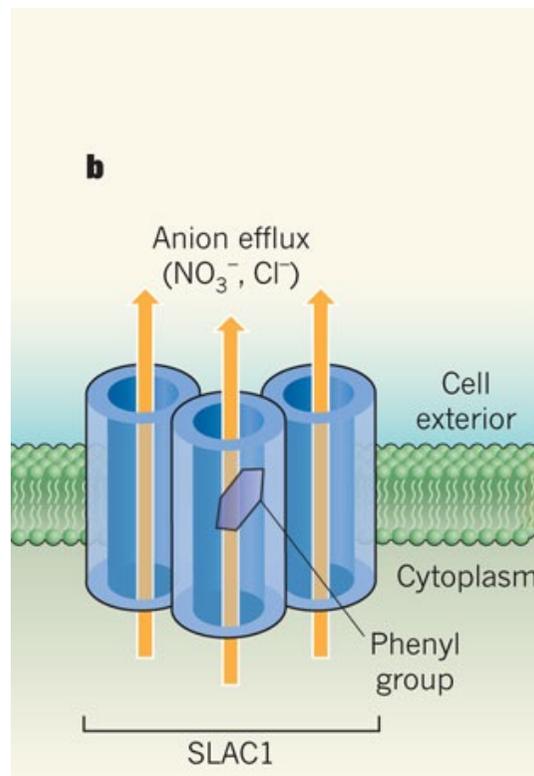
E-value is the number of matches with a given score (or higher) that we expect to occur by chance.

This depends on database size!

For an alignment with score S and $E\text{-value}=1$, we expect to have by chance 1 match with the same or higher score.

For an alignment with score S and $E\text{-value}=1$, we expect to have by chance 1 match with the same or higher score. If $E\text{-value}$ is 0.001 then we expect by chance 0.001 matches with the same or higher score.





Protein Families

aidanbudd.github.io/ppisnd/trainingMaterial/marcoPunta/

EMBO Budapest excellence in life sciences

UNIVERSITÀ DI BOLOGNA

EMBO Practical Course

Course Program
Introduction To Linux Command-line
Introduction To PPI Networks
STRING
Network Visualization With Cytoscape
Chimera
Unseminar
Structure And Interfaces Of PPIs
Peptide And Protein Docking
MSA & Jalview
Protein Families
Protein Feature Prediction
Repeats And Low Complexity Regions
Intrinsically Disordered Proteins
Short Linear Motifs
REST Services
Molecular Dynamics

PROTEIN FAMILIES

by Marco Punta

EXERCISE 1

- 3m71.pdb
- 3m71_del.pdb
- PF03595_seed_mod.txt.aln

EXERCISE 2

- P29973.fasta
- mystery-protein.fasta

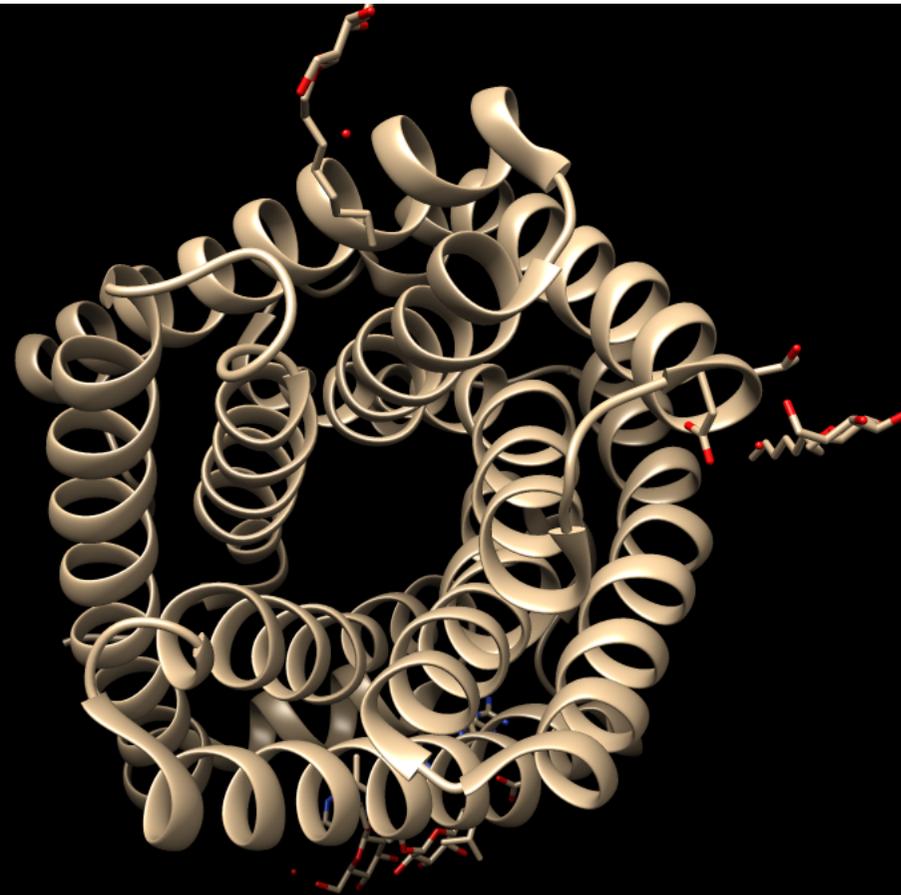
EXERCISE 3

- 2lhu.pdb
- family-building-exercise.fasta
- hmmer-ali.fasta
- jalview-ali.fasta

<http://aidanbudd.github.io/ppisnd/trainingMaterial/marcoPunta/>

1. OPEN Chimera

2. File -> Open "3M71.pdb"

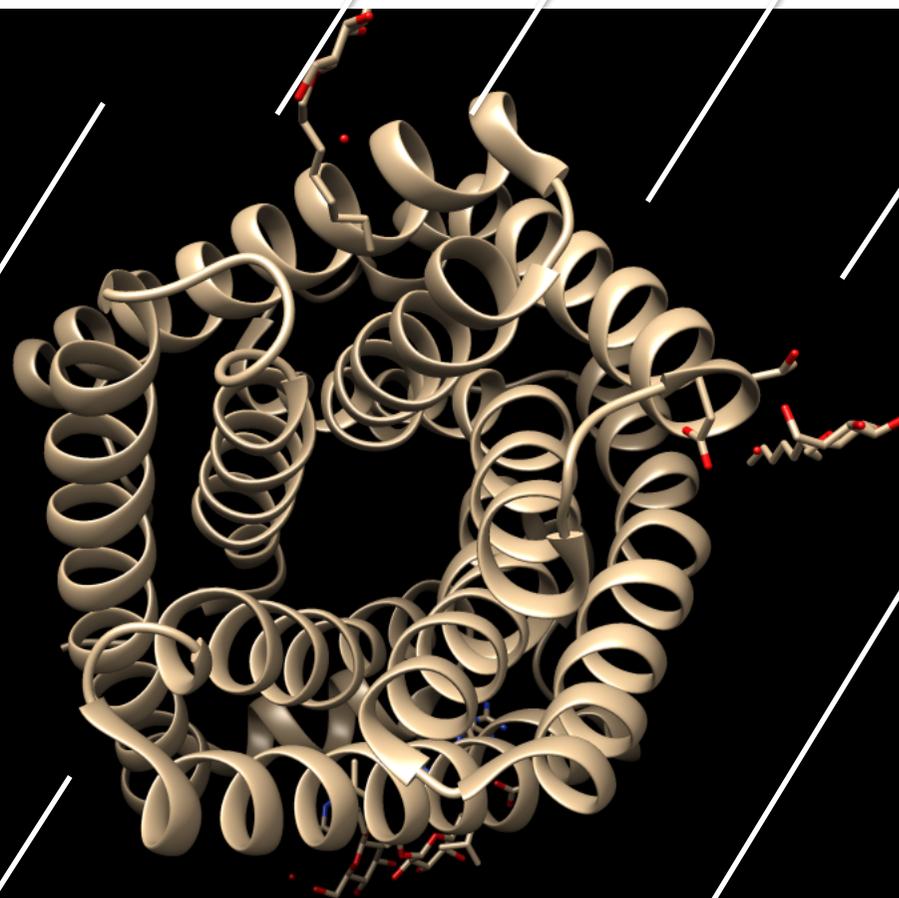


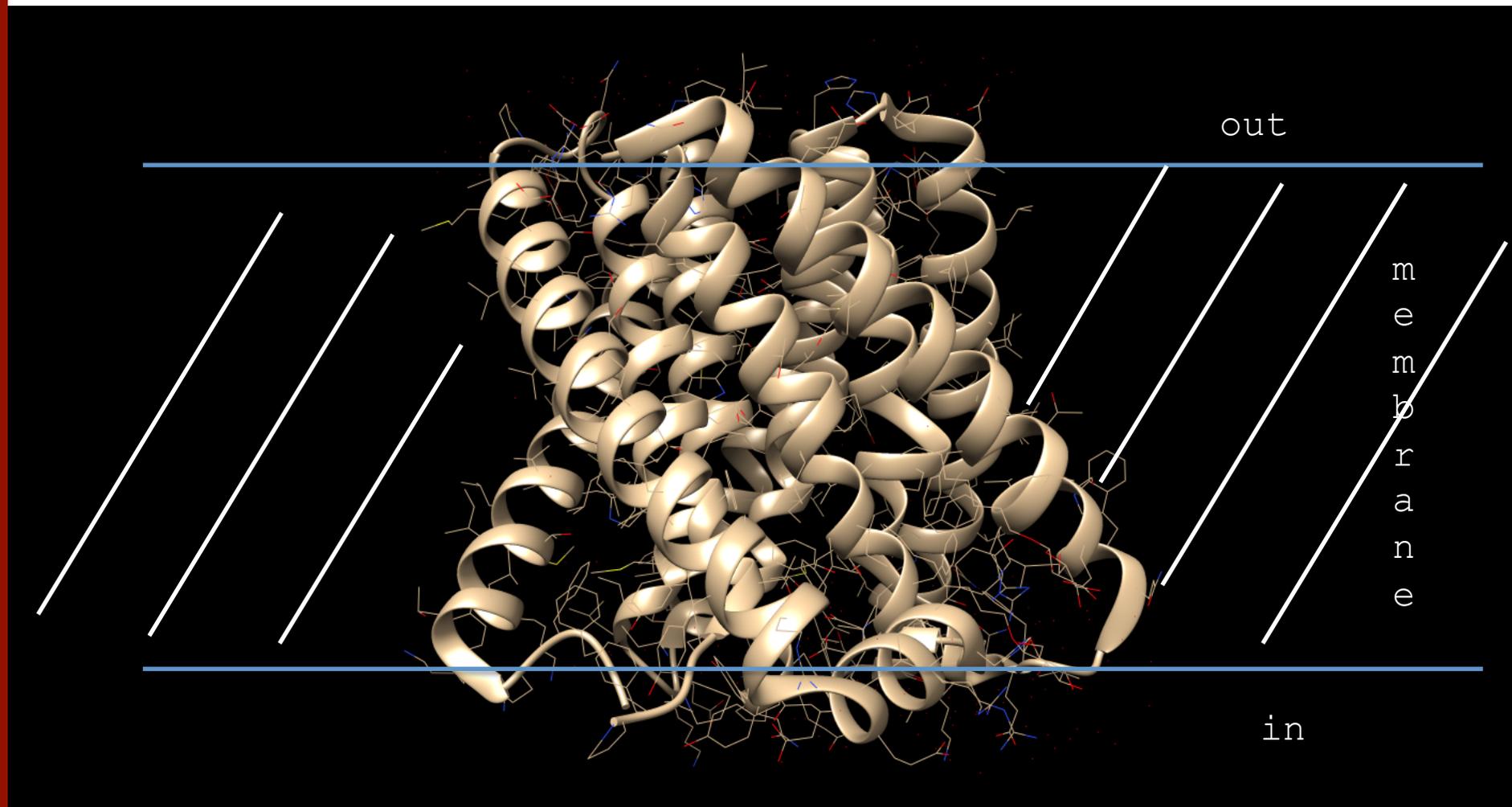
out

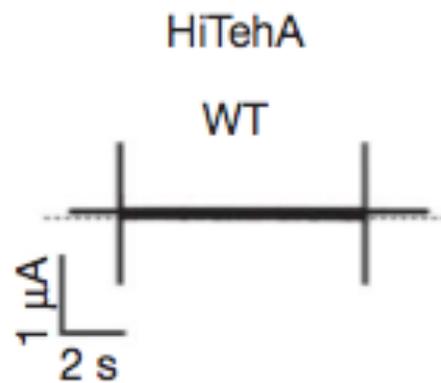
out

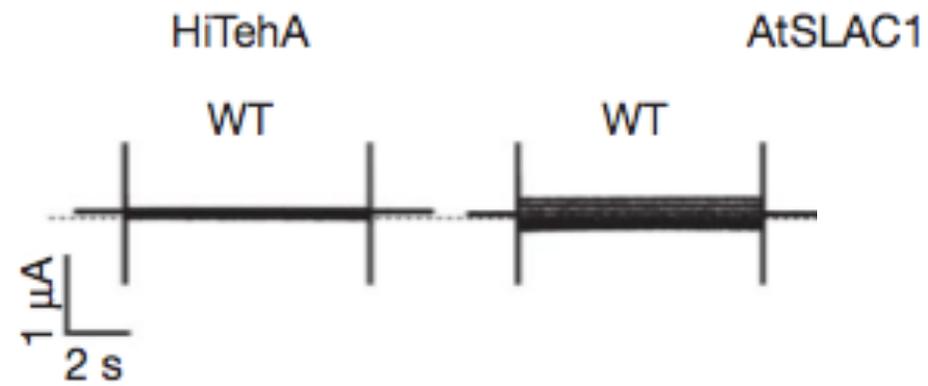
out

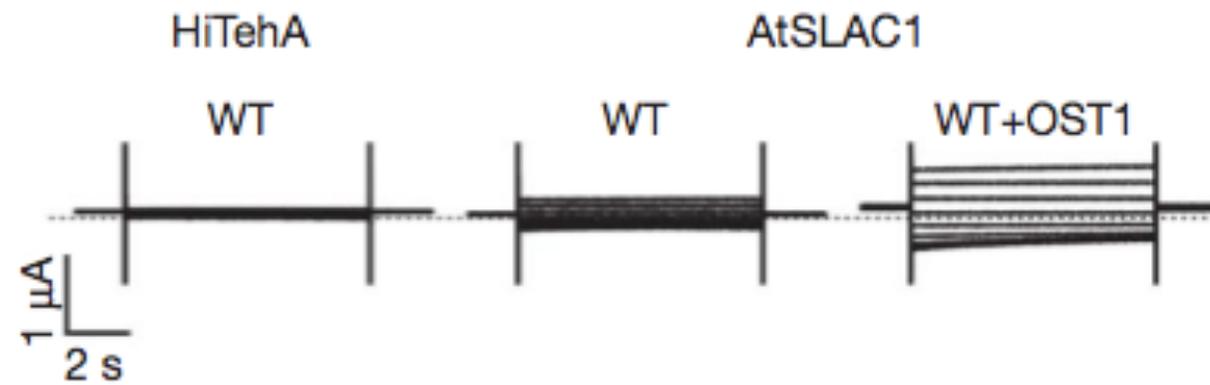
out





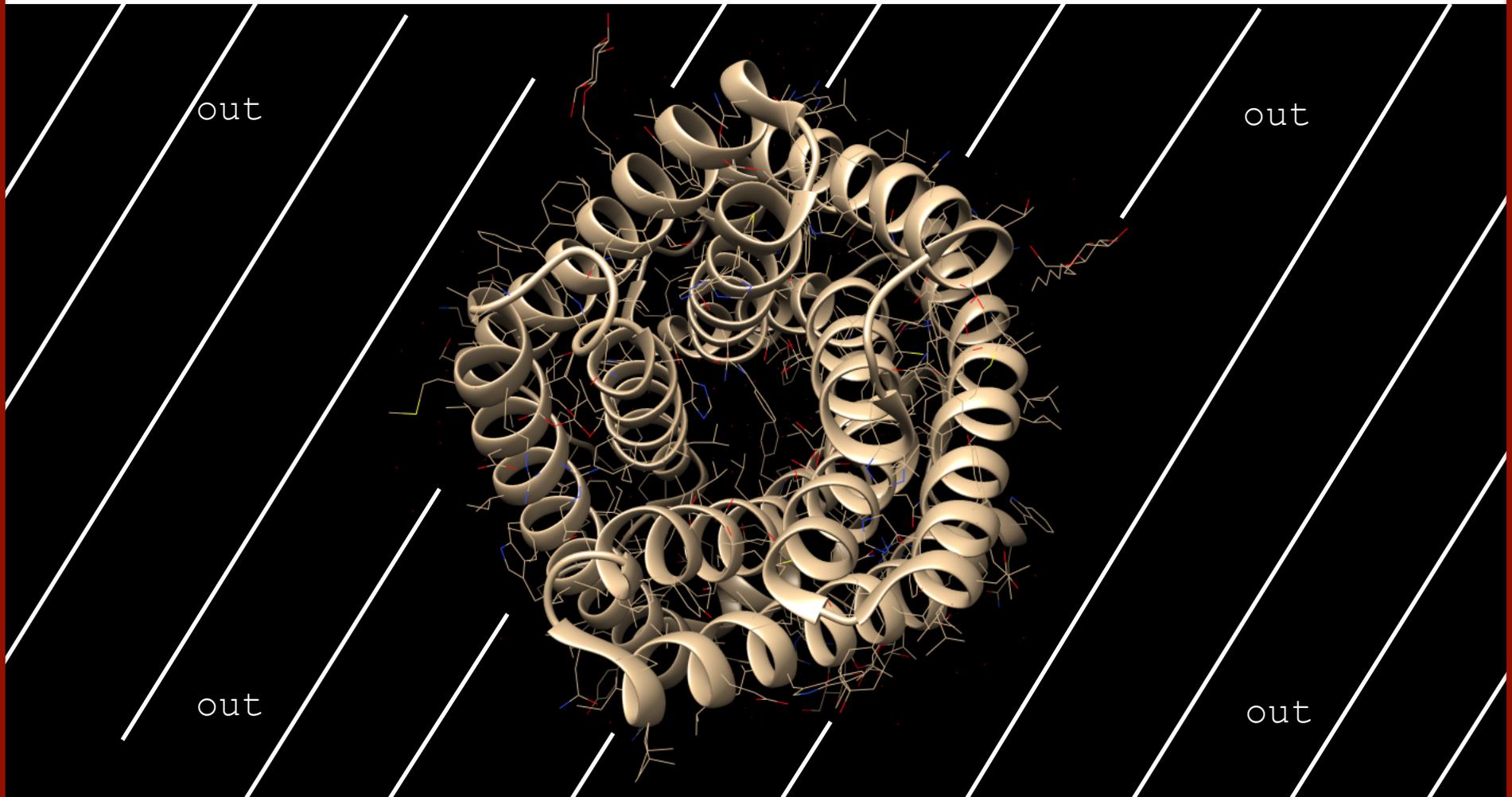






1. Actions -> Atoms/Bonds -> wire
2. Actions -> Atoms/Bonds -> show

1. Actions -> Atoms/Bonds -> wire
2. Actions -> Atoms/Bonds -> show



1. File -> Fetch by ID

Database	ID	Example
<input type="radio"/> NDB		pde024
<input type="radio"/> PDB		1yti
<input type="radio"/> PDB (mmCIF)		1yti
<input type="radio"/> PDB (biounit)		1hho
<input type="radio"/> CATH		1cukA01
<input type="radio"/> SCOP		d1g0sa_
<input type="radio"/> cellPACK		HIV-1_0.1.6
<input type="radio"/> PubChem		12123
<input checked="" type="radio"/> CASTp	3m71	1www
<input type="radio"/> EDS (2fo-fc)		1a0m
<input type="radio"/> EDS (fo-fc)		1a0m
<input type="radio"/> EMDB		5625
<input type="radio"/> EMDB & fit PDBs		1048
<input type="radio"/> PQS		2cwj
<input type="radio"/> ModBase		P04848
<input type="radio"/> VIPERdb		1ej6
<input type="radio"/> UniProt		P01138 NGF_HUMAN

Buttons: Set download directory, Ignore any cached data, Keep dialog up after Fetch, Fetch, Web Page, Close, Help

EMBO-Course-Budapest-2016-DDT.pptx

CASTp: 3m71

Columns

ID ▼	MS volume	pocket MS area	# openings	mouth MS area
33	881.8	728.0	2	162.3
32	667.5	511.1	2	84.8
31	313.7	243.6	1	46.1
30	318.7	217.3	1	94.5
29	343.9	189.7	1	122.4
28	194.8	173.5	1	37.0
27	151.3	154.8	1	45.0
26	238.1	165.4	1	75.0
25	148.3	119.0	1	51.0
24	76.3	69.6	1	26.9

Color name: #8787cecebeb

Treatment of Chosen Pocket Atoms

- Select
- Color (and color all other atoms No)
- Surface (colored by hydrophobicity)
- Zoom in on
- Exclude mouth atoms

Quit Hide Help

The screenshot displays a protein structure visualization with a data table and a settings panel. The table lists cavity metrics for various residues, and the settings panel shows options for selecting and coloring atoms.

ID	MS volume	pocket MS area	# openings	mouth MS area
33	881.8	728.0	2	162.3
32	667.5	511.1	2	84.8
31	313.7	243.6	1	46.1
30	318.7	217.3	1	94.5
29	343.9	189.7	1	122.4
28	194.8	173.5	1	37.0
27	151.3	154.8	1	45.0
26	238.1	165.4	1	75.0
25	148.3	119.0	1	51.0
24	76.3	69.6	1	26.9

Treatment of Chosen Pocket Atoms

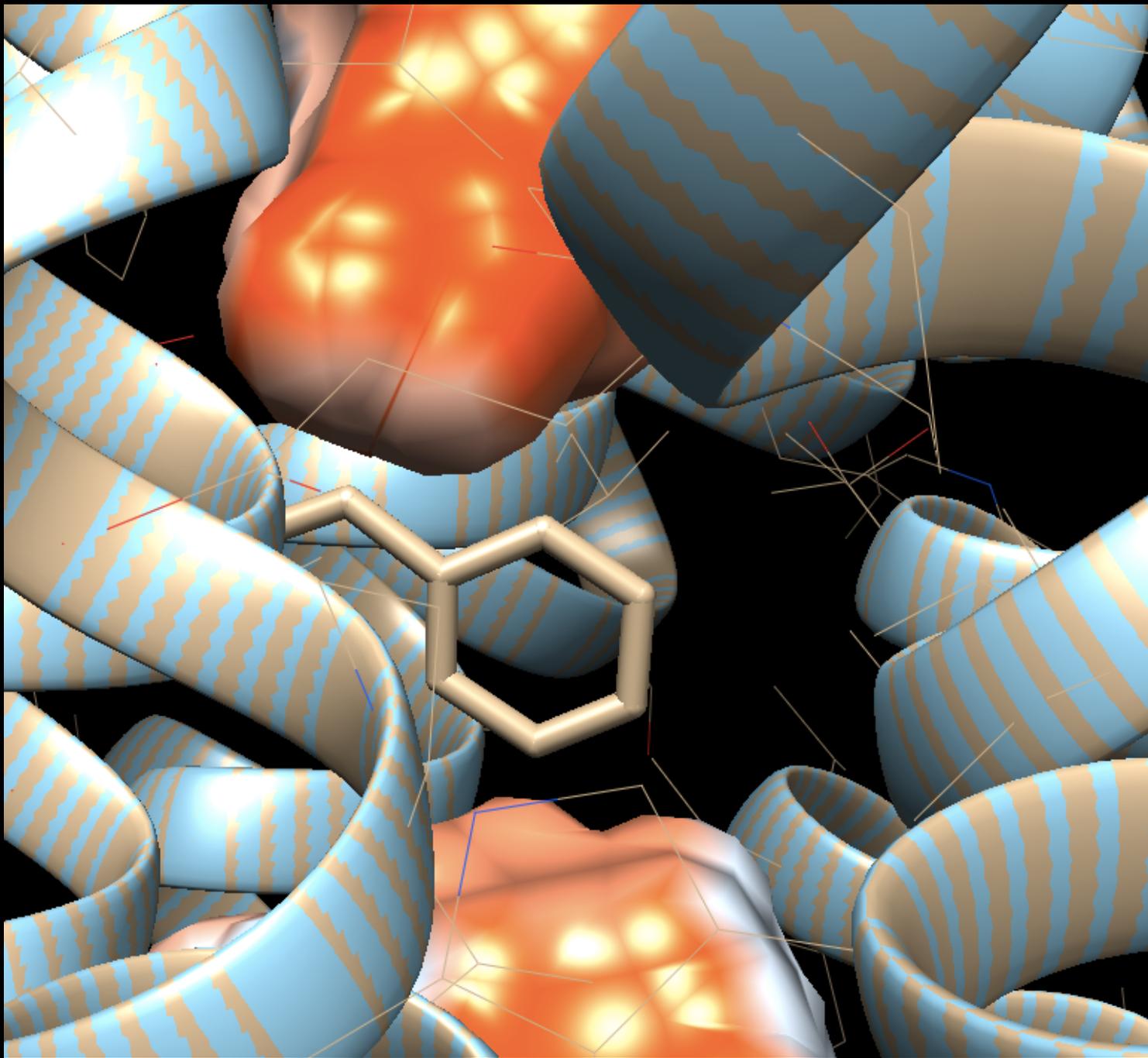
- Select
- Color (and color all other atoms)
- Surface (colored by hydrophobicity)
- Zoom in on
- Exclude mouth atoms

Buttons: Quit, Hide, Help

Hold down
'Shift'
key to select
multiple
cavities

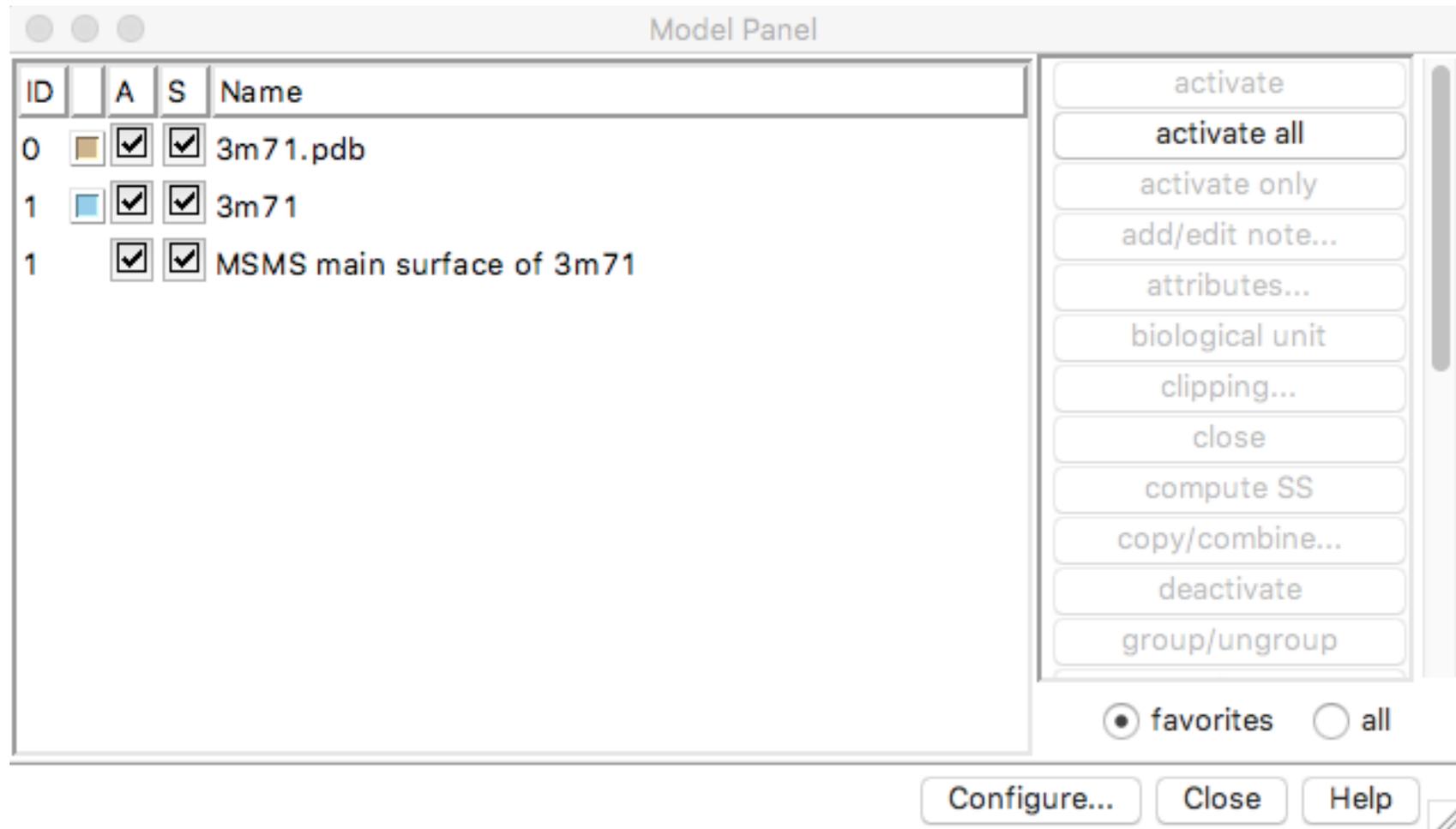




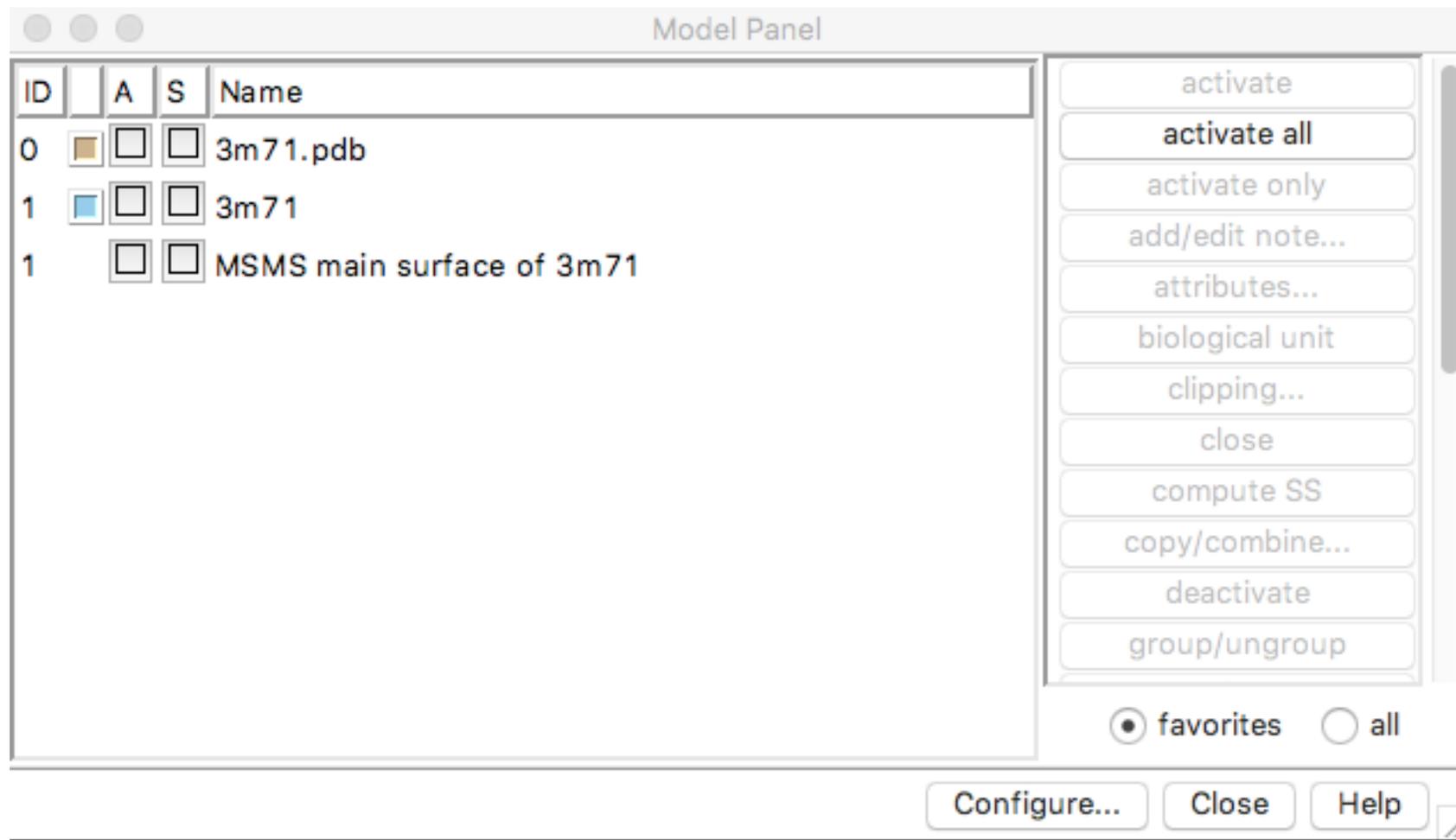


1. Favorites -> Model panel

1. Favorites -> Model panel

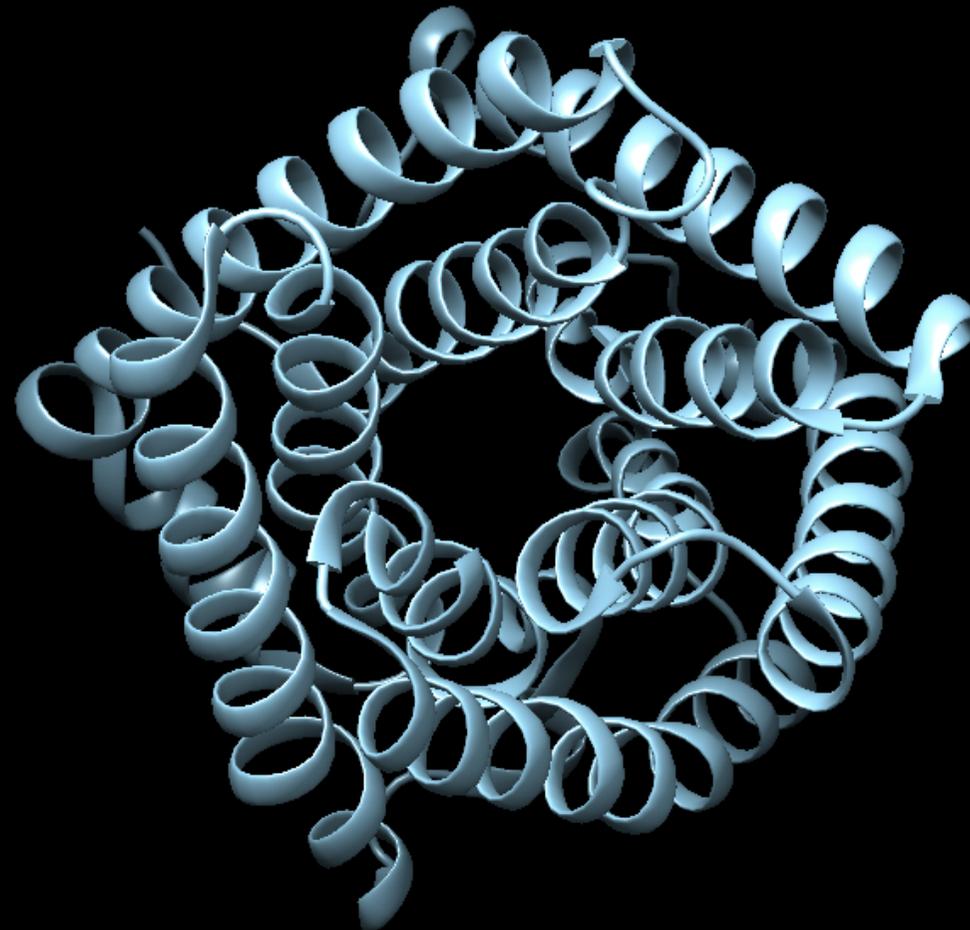


1. Favorites -> Model panel



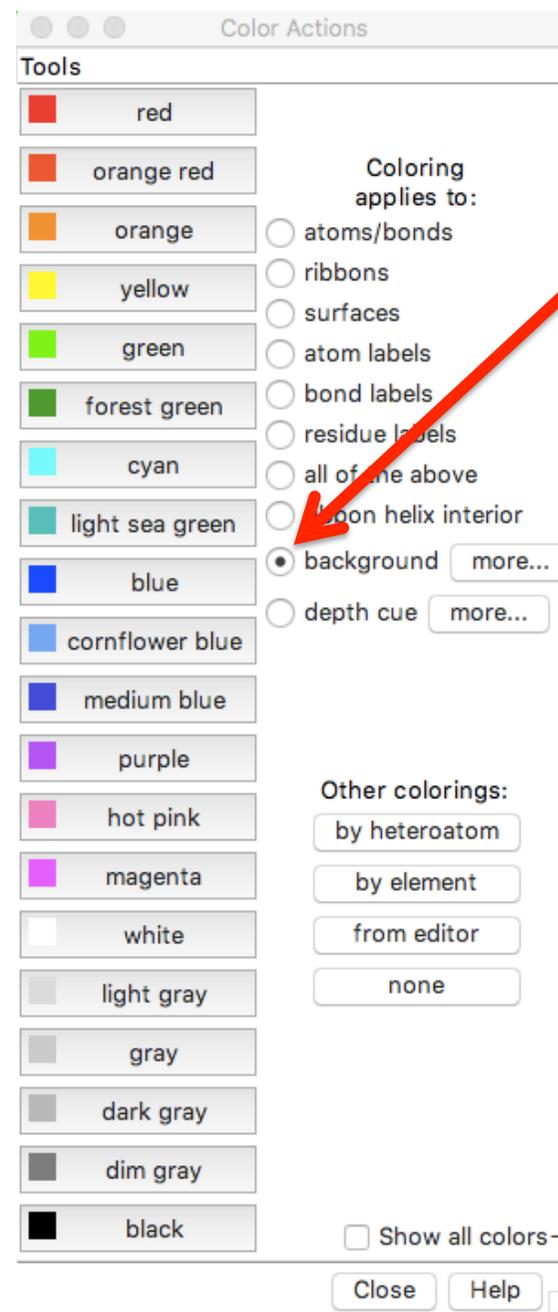
1. File -> Open "3M71_del.pdb"

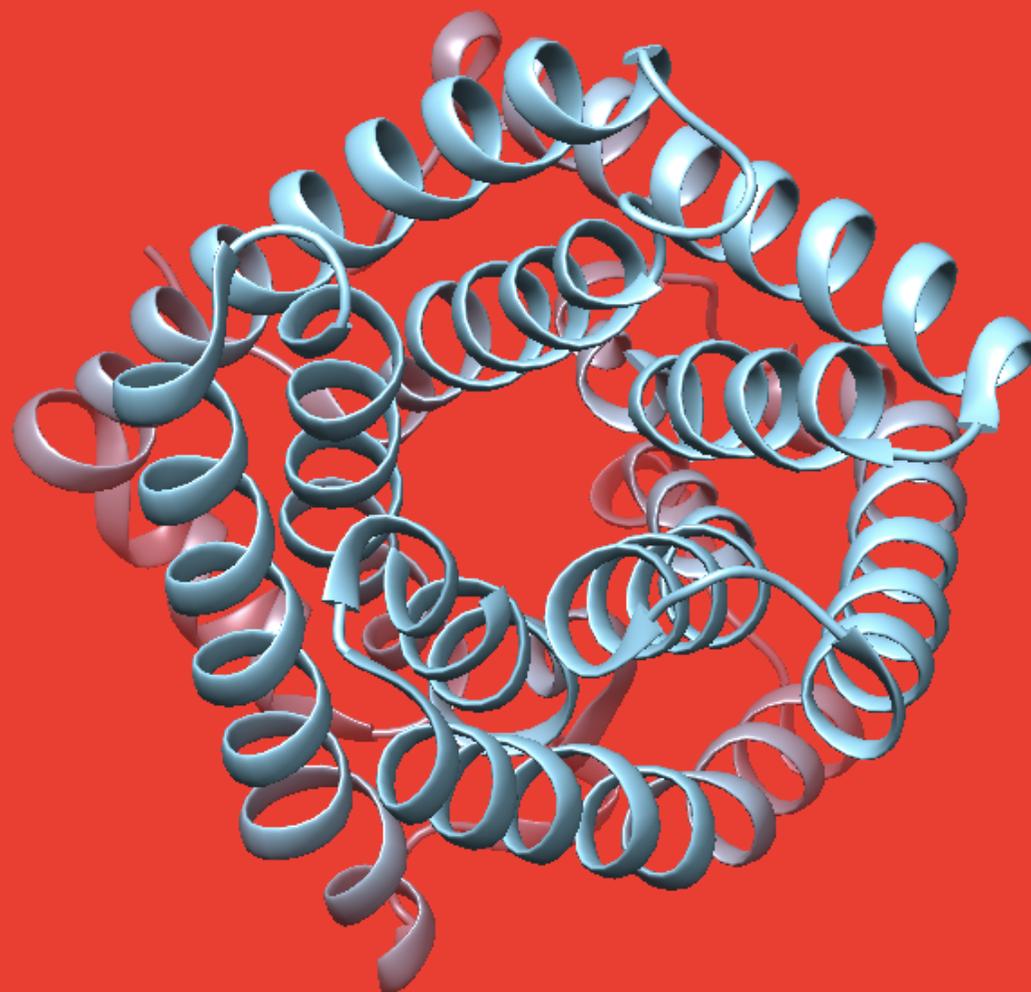
1. File -> Open "3M71_del.pdb"



1. File -> Open "3M71_del.pdb"
2. Actions -> Color -> all options

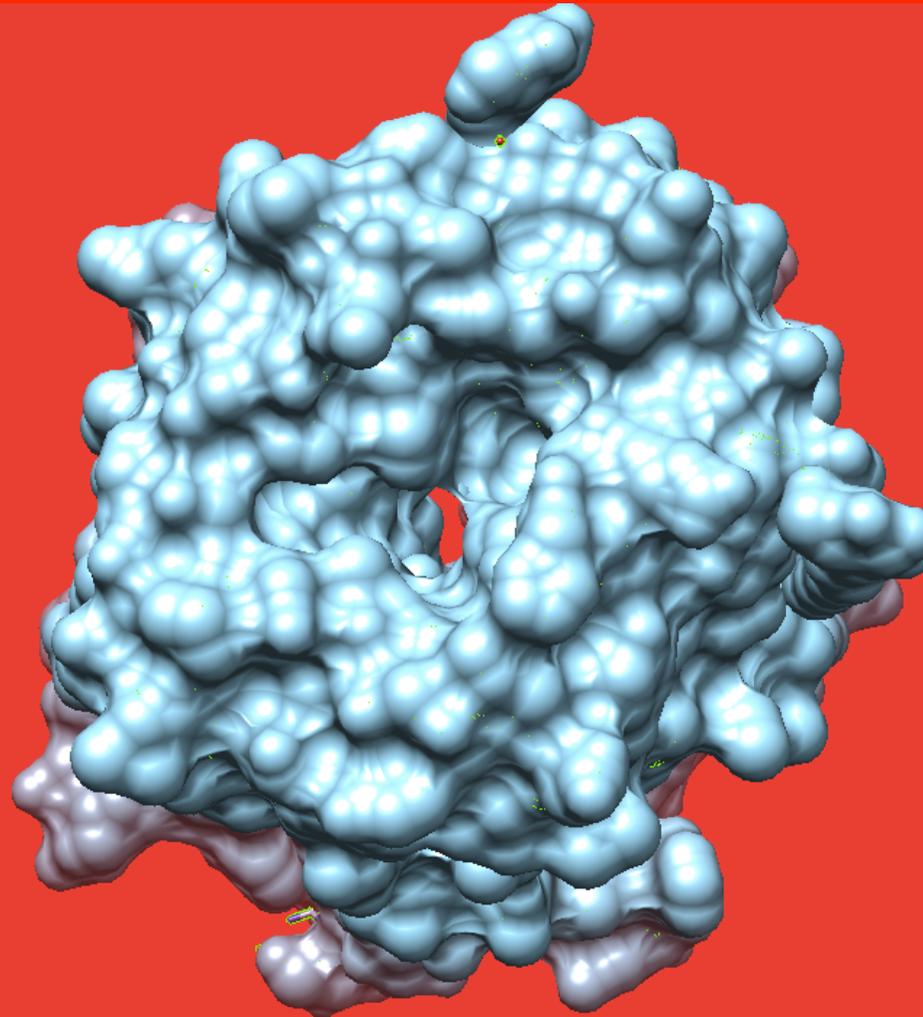
1. File -> Open "3M71_del.pdb"
2. Actions -> Color -> all options
3. Click on 'background' and select red



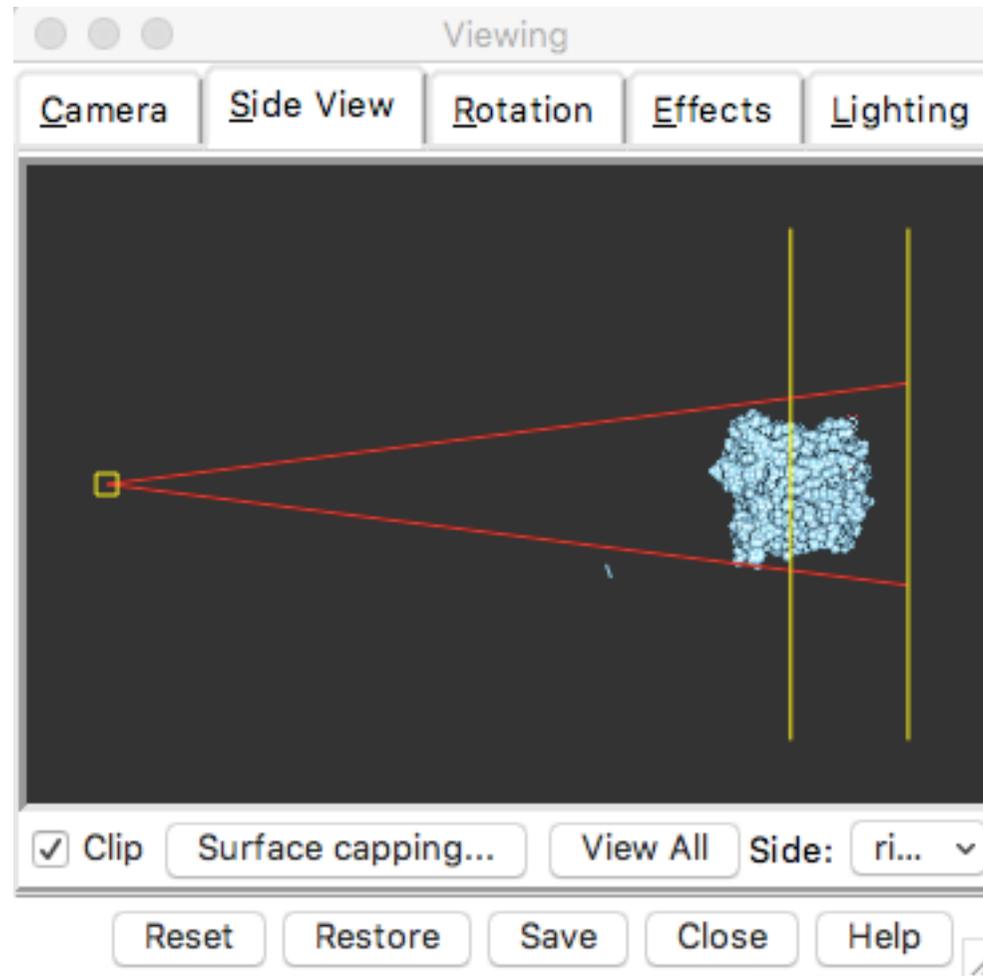


1. Select -> Chain -> A -> 3m71_del.pdb (#1)
2. Actions-> Surface -> show

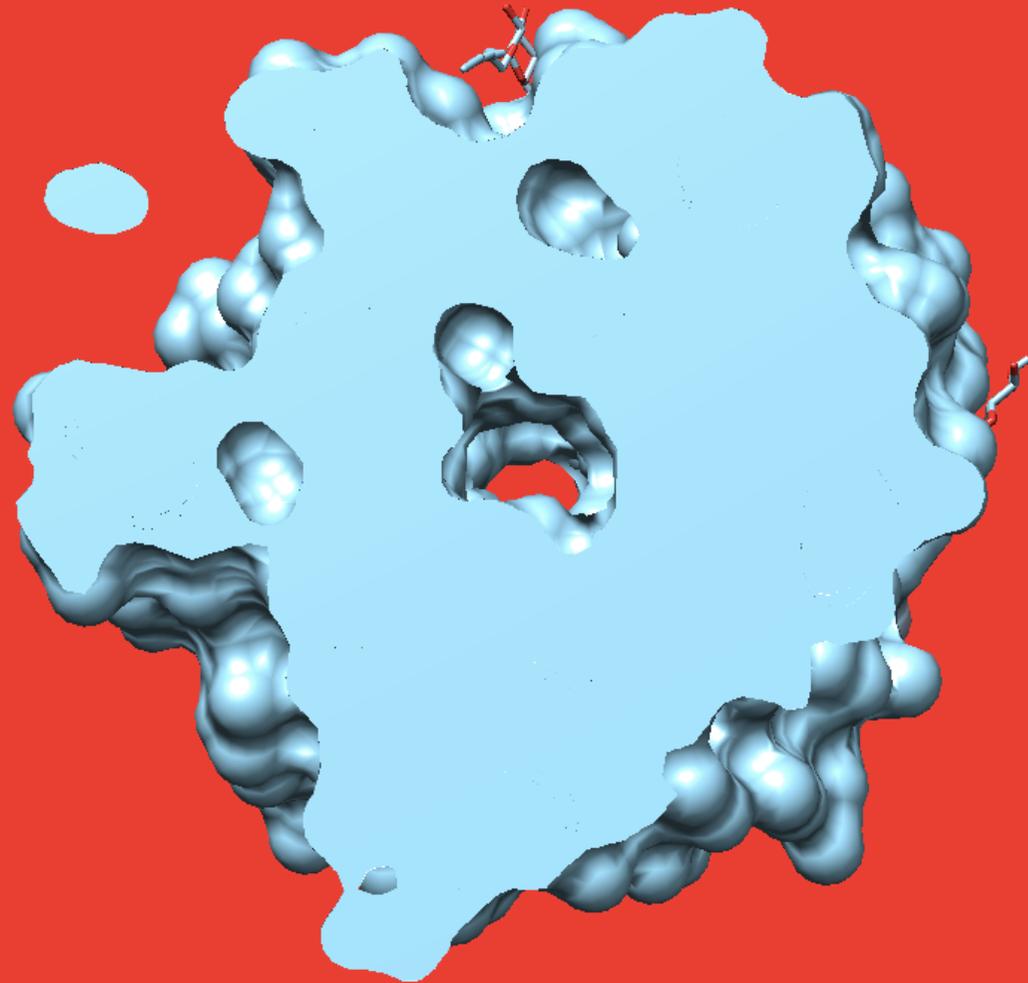
1. Select -> Chain -> A -> 3m71_del.pdb (#1)
2. Actions-> Surface -> show



1. Favorites -> Side View

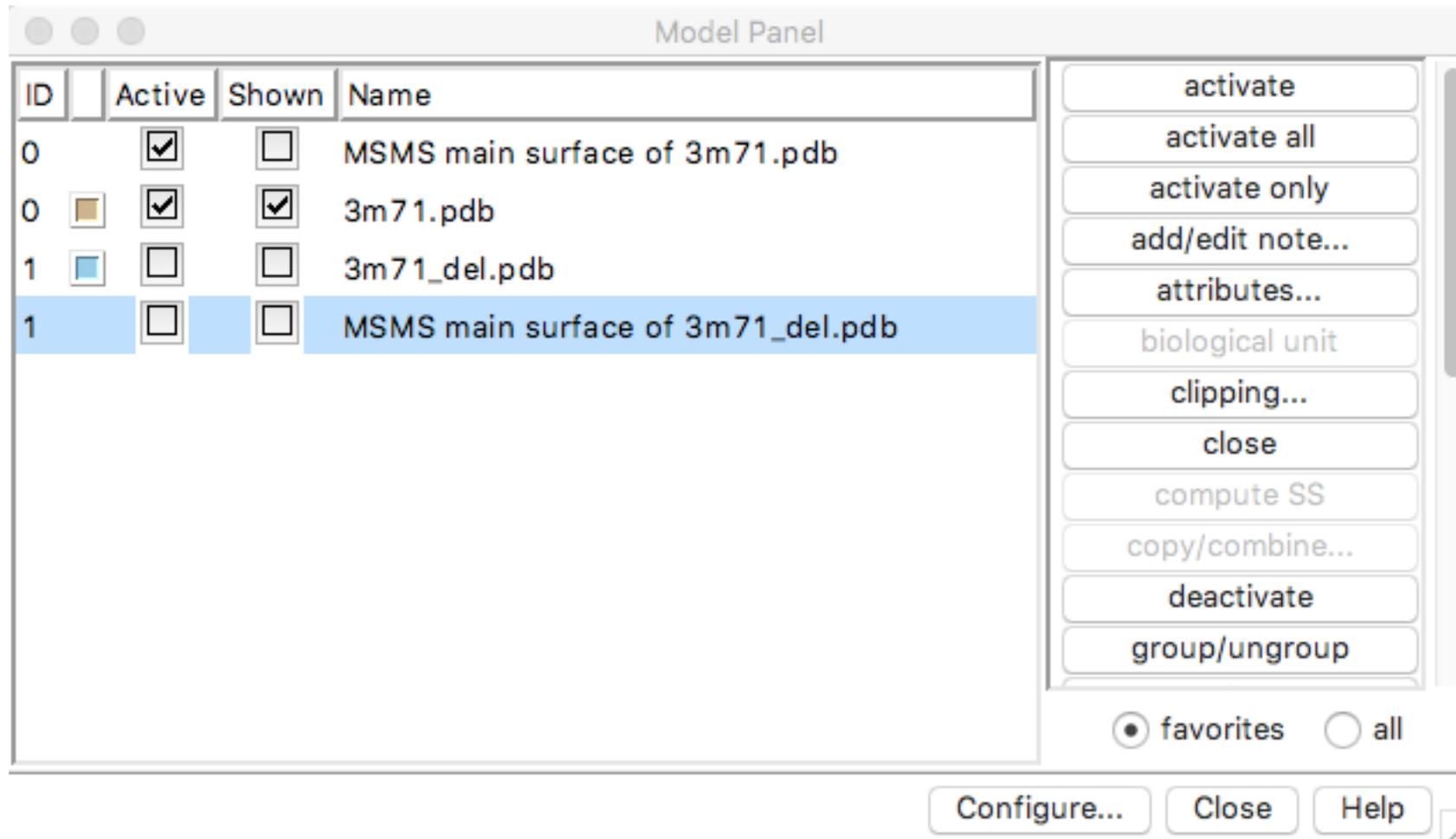


1. Favorites -> Side View

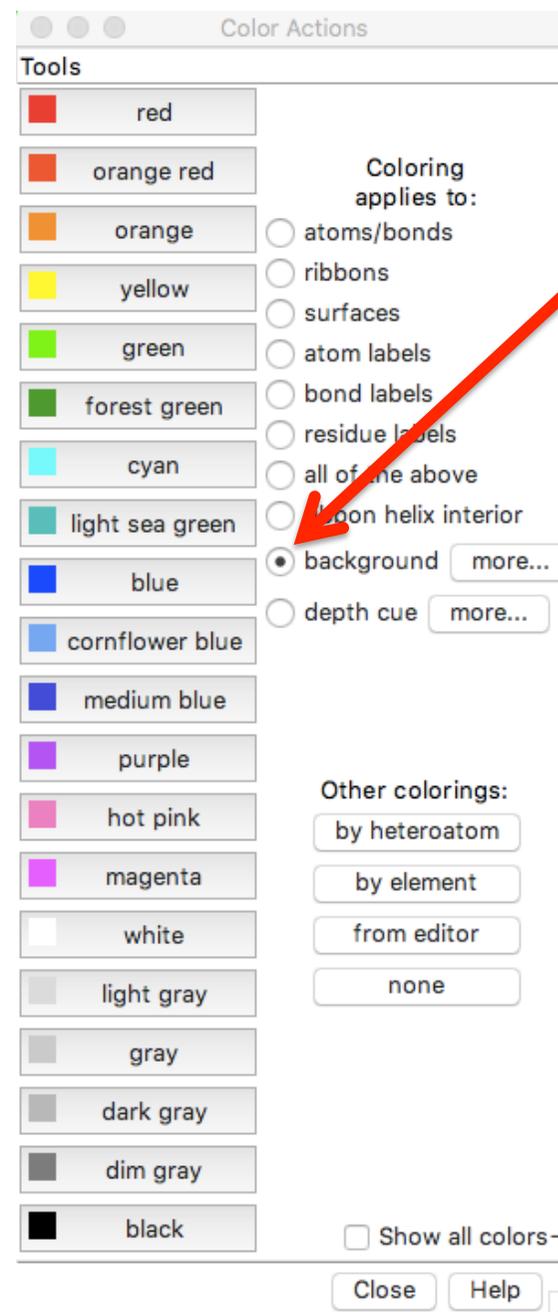


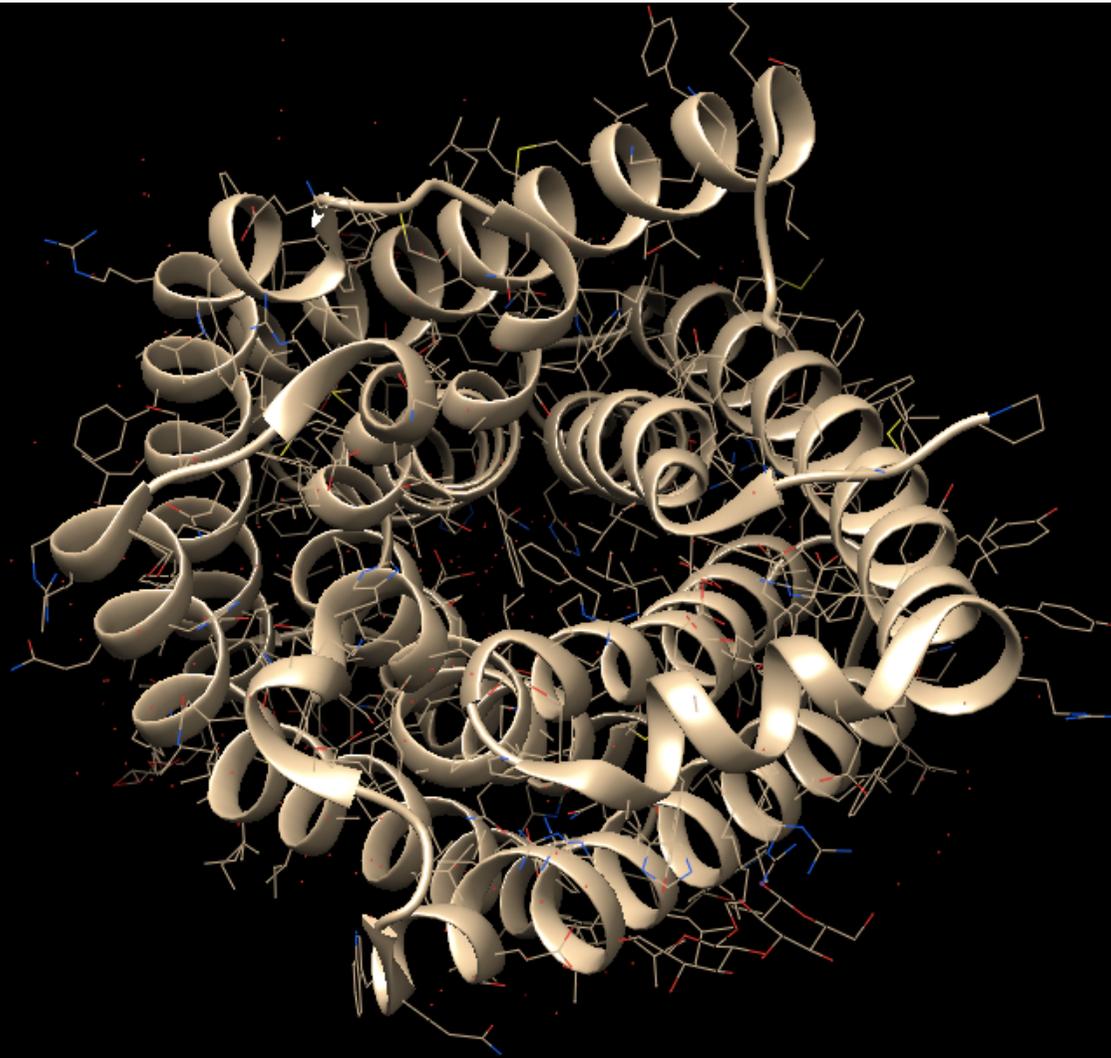


1. Favorites -> Model panel



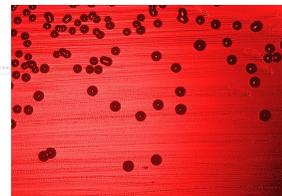
1. File -> Open "3M71_del.pdb"
2. Actions -> Color -> all options
3. Click on 'background' and select black





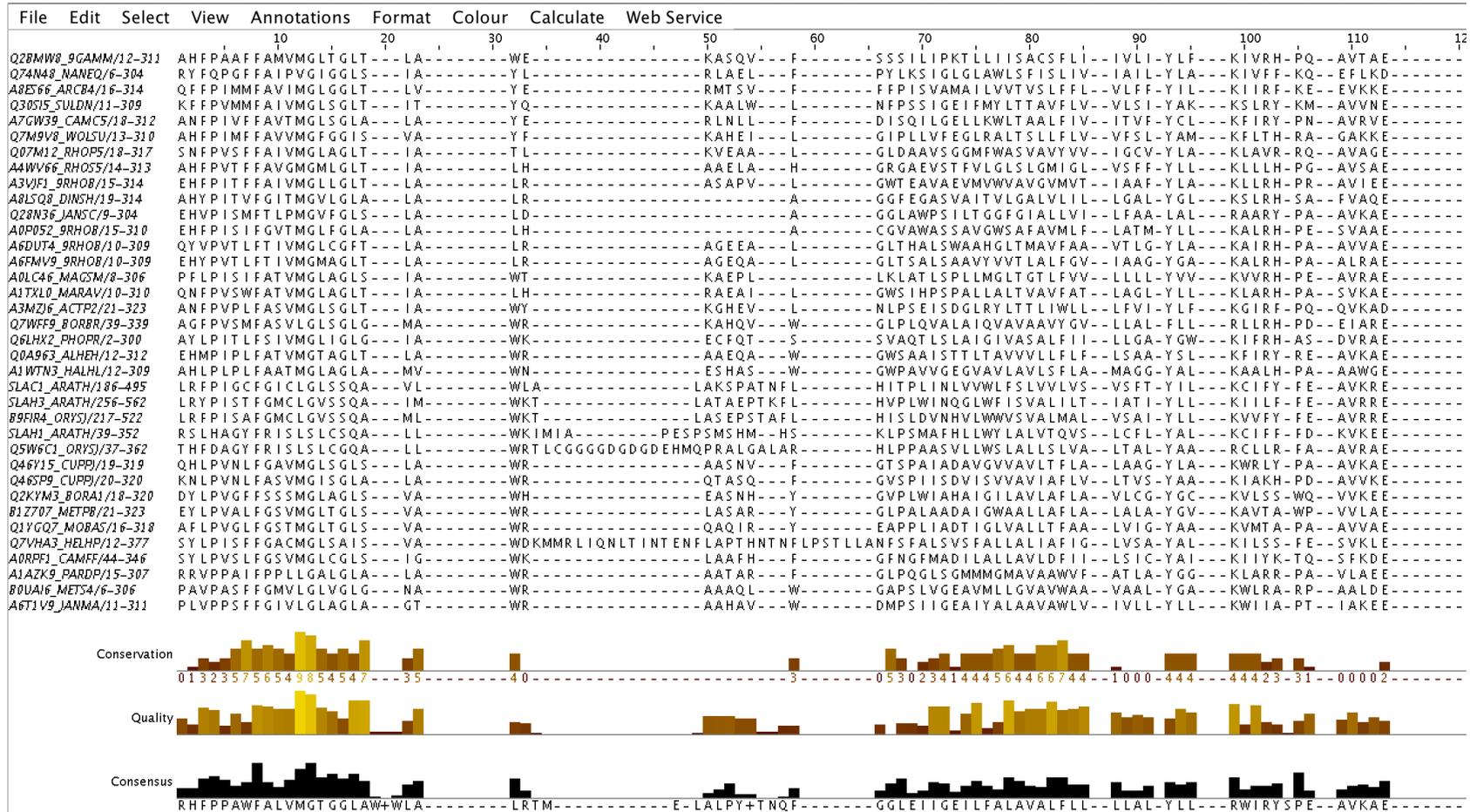
Alignment

Q9LD83	SLAC1_ARATH - Guard cell S-type anion channel SLA... - Arabidopsis thal...		
E-value: 3e-10	Positives : 41.0%		
Score: 160	Query Length: 328		
Ident.: 22.0%	Match Length: 556		
P44741	20	PFPL--PTGYFGIPLGLAALS LawFHLE-----NLFP AARMVSDVLGIVASAVWILFILM	72
		PF L P G FGI LGL++ ++ W L N +++ V+ + + V +	
Q9LD83	183	PFLLRFPICGFCGICLGLSSQAVLWLALAKSPATNFLHITPLINLVVWLFSLVVLVSVSFT	242
P44741	73	YAYKLRYYFEEVRAEYHSPVRF SFIALIPITTMLVG---DILYRWNPLIAEVLWIGTIG	129
		Y K +YFE V+ EY PVR +F + M + ++ N IW +G	
Q9LD83	243	YILKCIFYFEAVKREYFHPVRVNF FFAPWVCMFLAISVPPMFSPNRKYLHPAIWCVFMG	302
P44741	130	QLLFSTLRVSELWQGGVFEQ--KSTHPSFYLPVAANFTSASSLALLGYHDLGYLFFGAG	187
		F L++ W G + K +PS +L +V NF A + +G+ ++ + G	
Q9LD83	303	PYFFLELKIYGQWLSGGKRR LCKVANPSSHL-SVVGNFV GAILASKVGWDEVAKFLWAVG	361
P44741	188	MIAWIIFEPVLLQHLRISSLEPQFRATMGIVLAPAFVCSAYLSINHGEVDTLAKILWGY	247
		+++ L Q L S P+ + + A S + +G+ D ++ +	
Q9LD83	362	FAHYLVVFVTLYQRLPTSEALPKELHPVYSMFIAAPSAASIAWNTIYGQFDGCSRTCFFI	421
P44741	248	GFLQLFFLLRLFPWIVEKGLNIGLWAFSFC LASMANSATAFY----HGNVLQGVSI FAFV	303
		L+ + ++ W++ F + + A+ AT Y G + +++	
Q9LD83	422	ALFLYISLVARINF FTGFKFSVAWWSYTFE MTT-ASVATIKYAEAVPGYPSRALALTL SF	480
P44741	304	FSNVMIGLLV LMTI 317	
		S M+ +L + T+	
Q9LD83	481	ISTAMVCVLFVSTL 494	

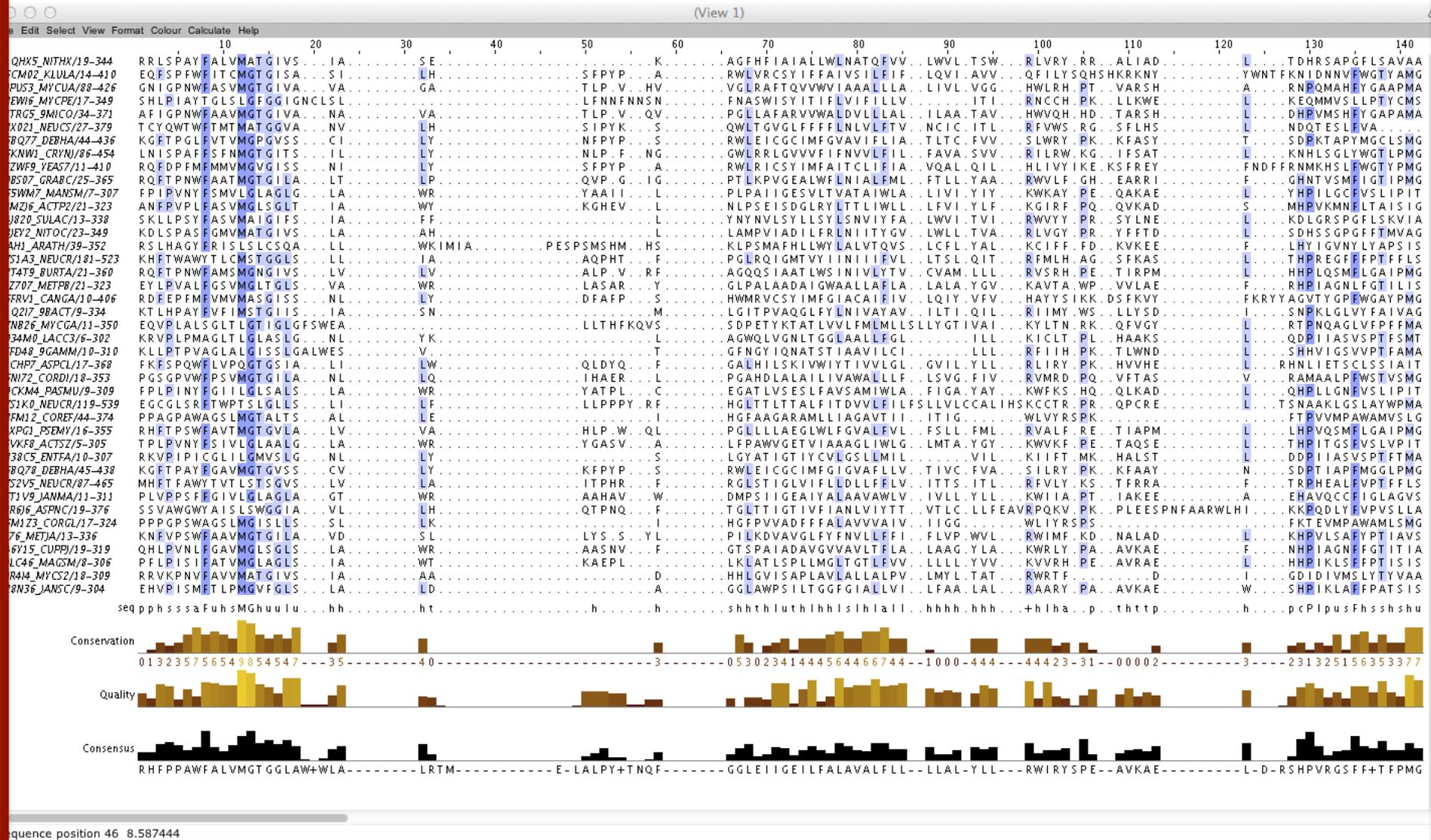


1. OPEN Jalview

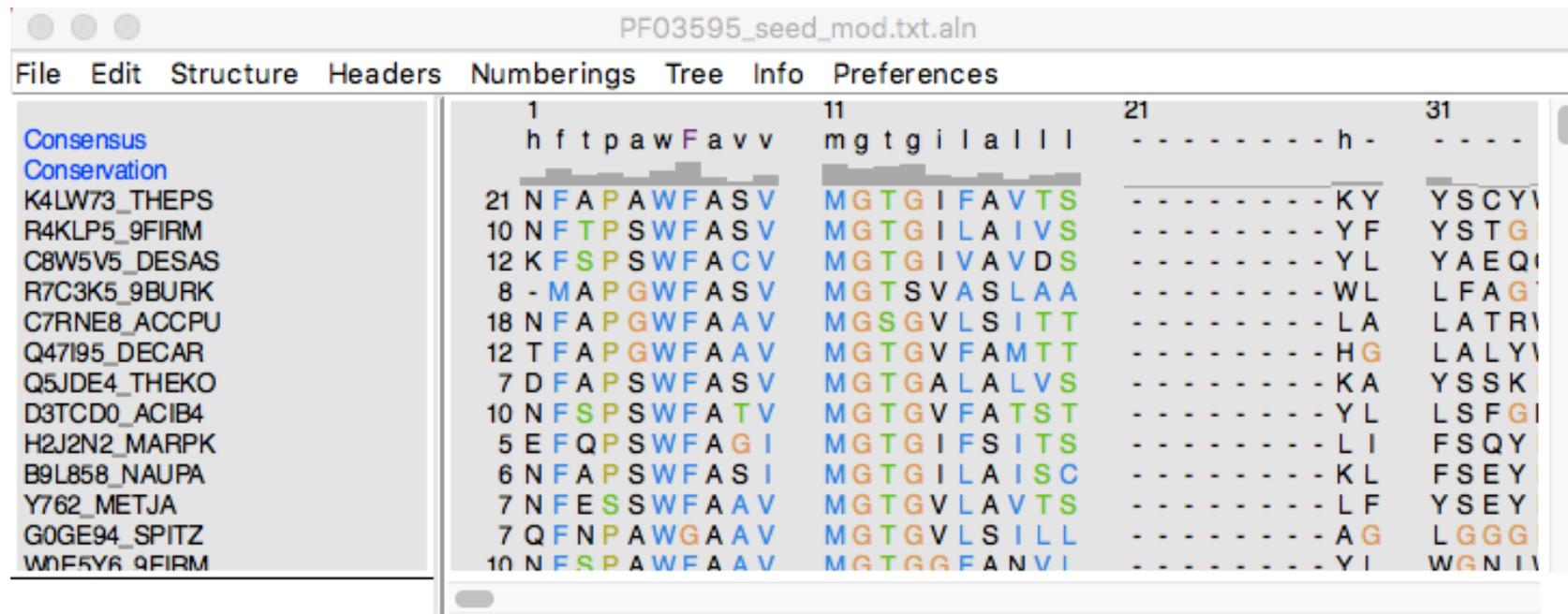
2. File -> Input Alignment -> From File "PF03595_seed.txt"



1. Colour -> BLOSUM62



1. Tools-> Sequence -> Multialign Viewer
2. Choose "PF03595_seed_mod.txt"
3. Select Aligned FASTA

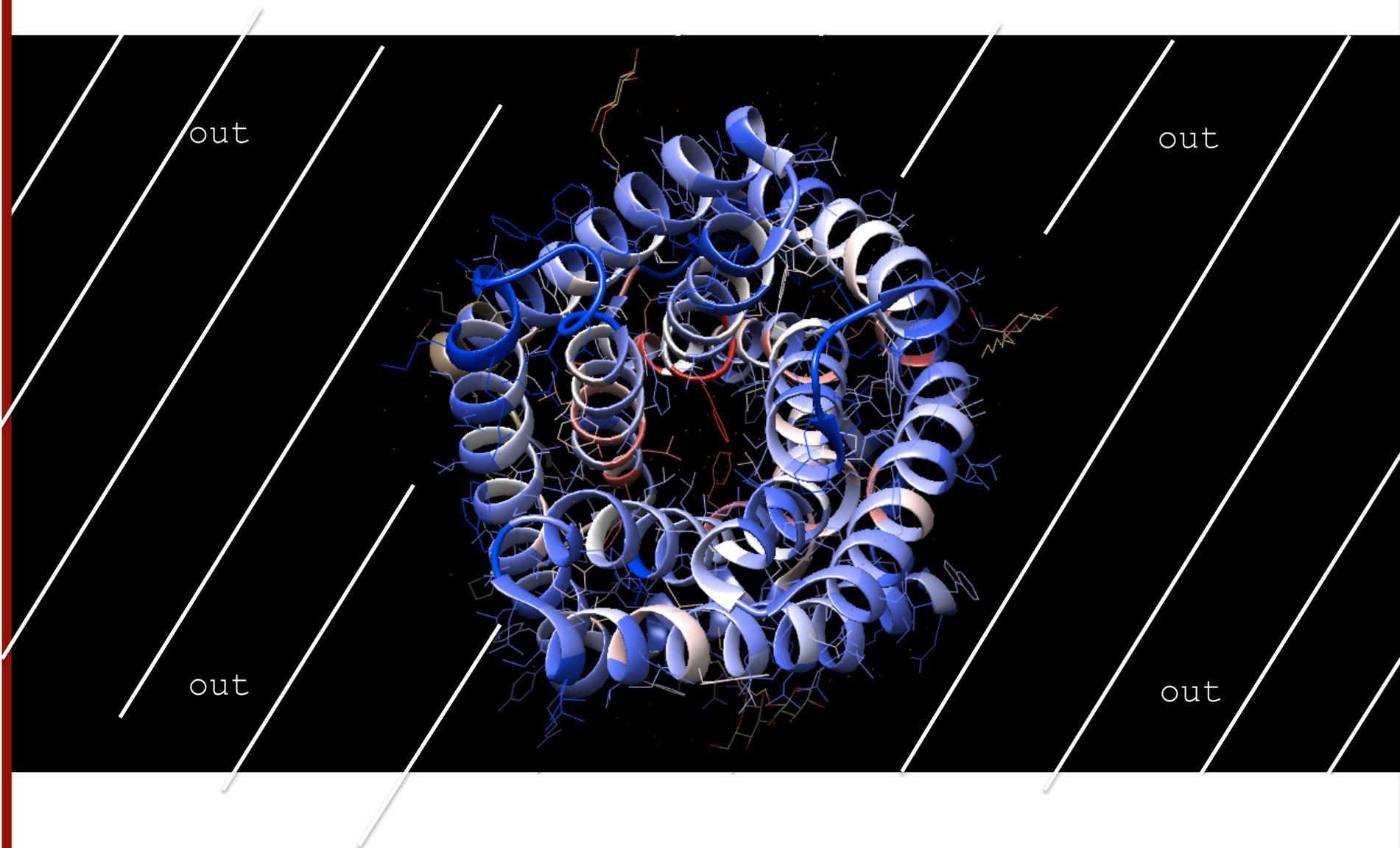


Right-click to focus on residue
 Right-shift-click to focus on region

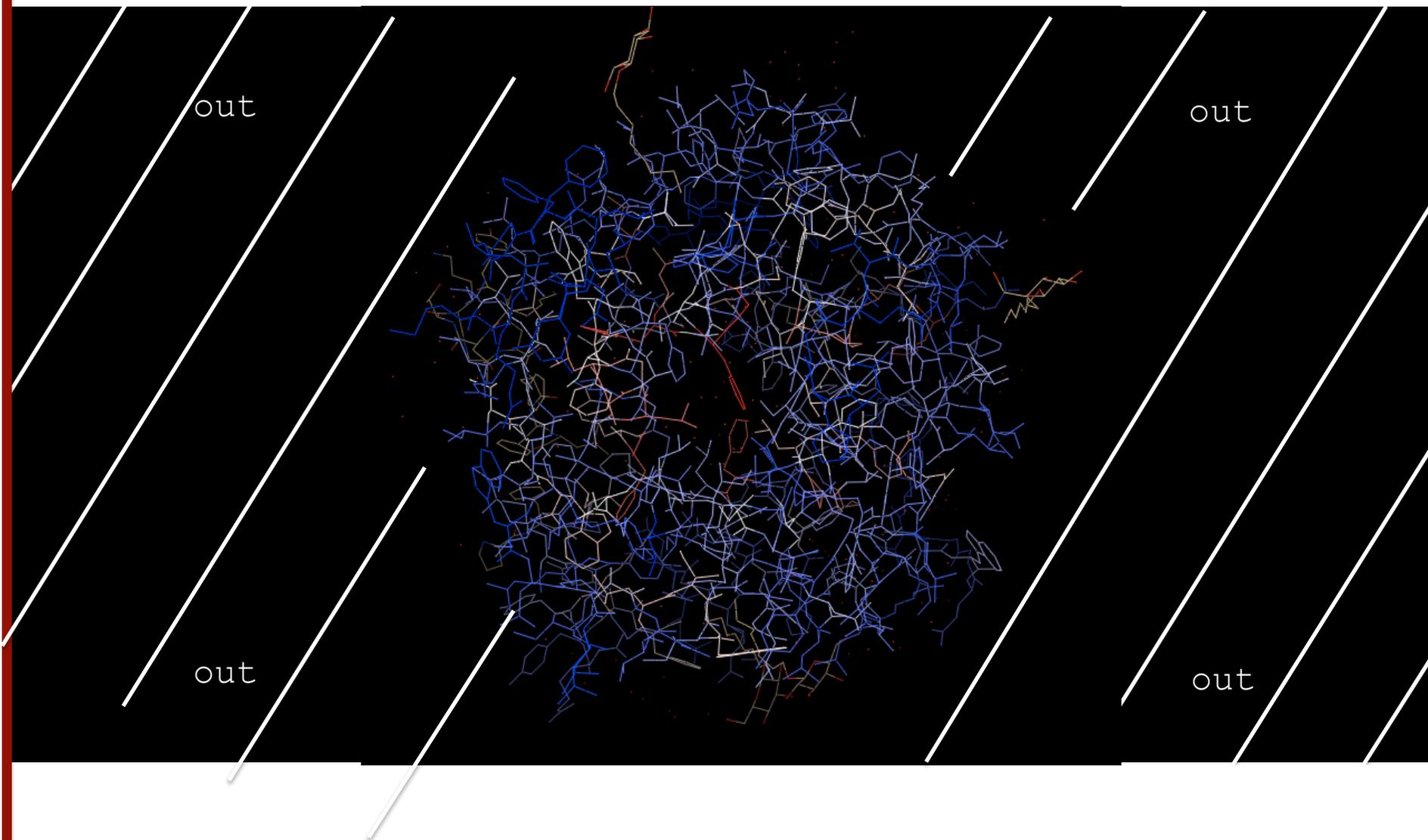
Quit Hide Help

1. Tools-> Sequence -> Multialign Viewer
2. Choose "PF03595_seed_mod.txt"
3. Select Aligned FASTA

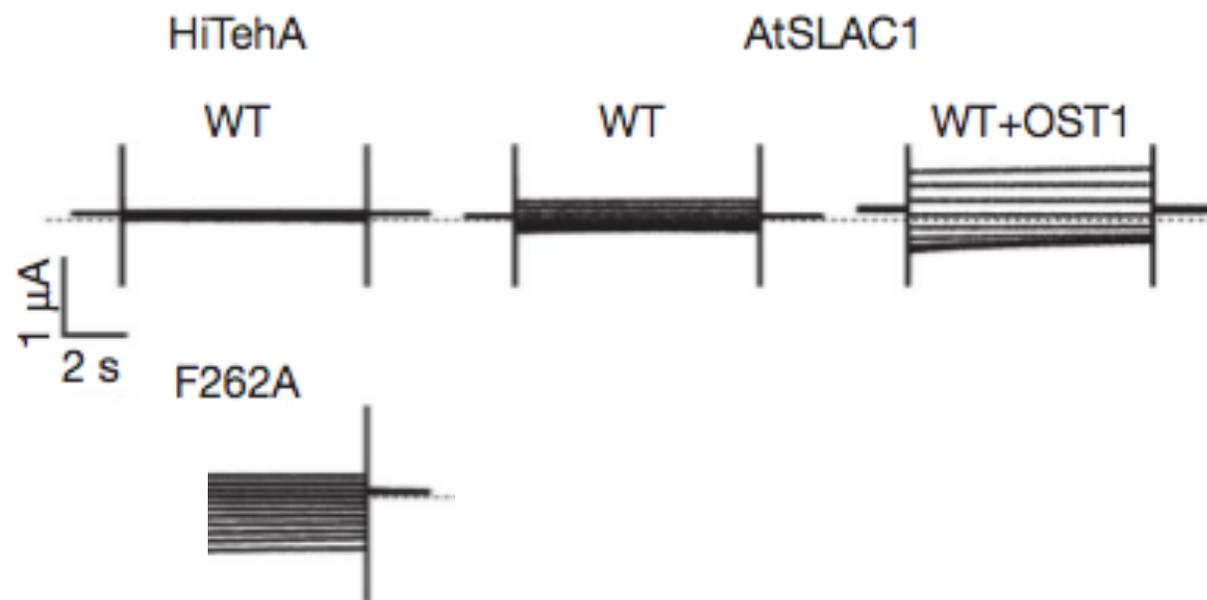
4. Structure -> Render by Conservation

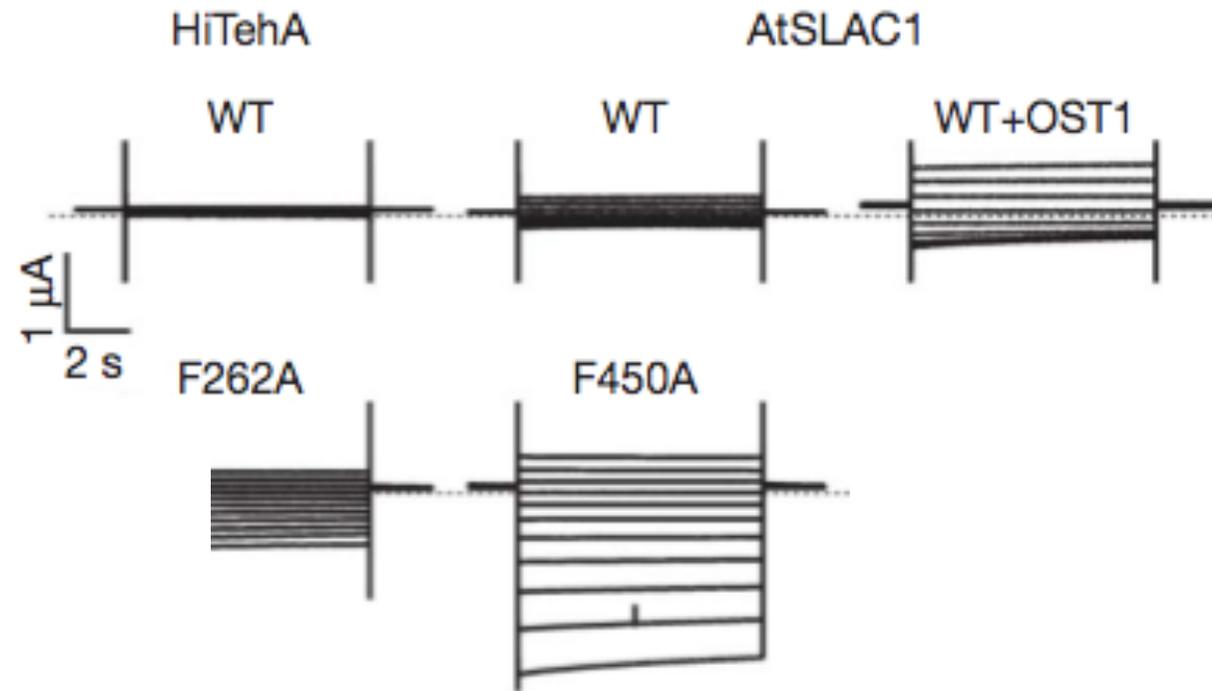


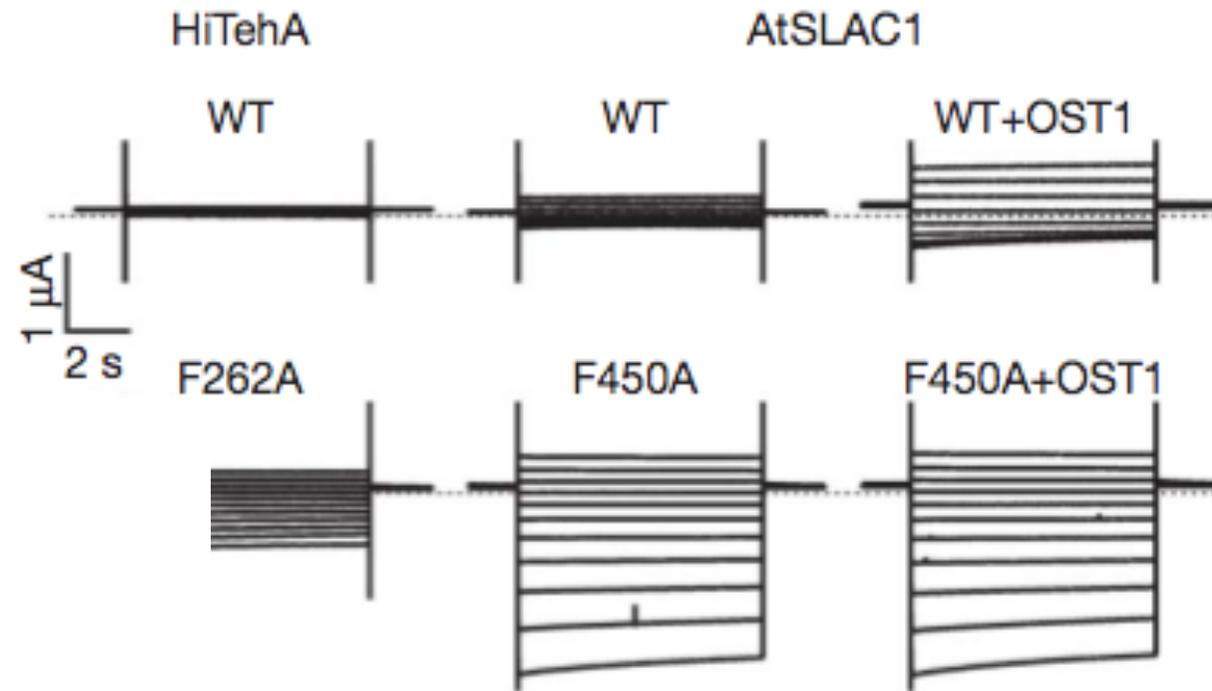
1. Actions-> Ribbon -> hide



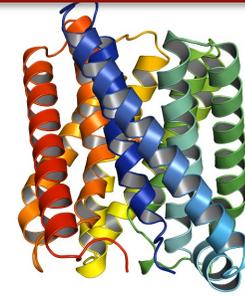




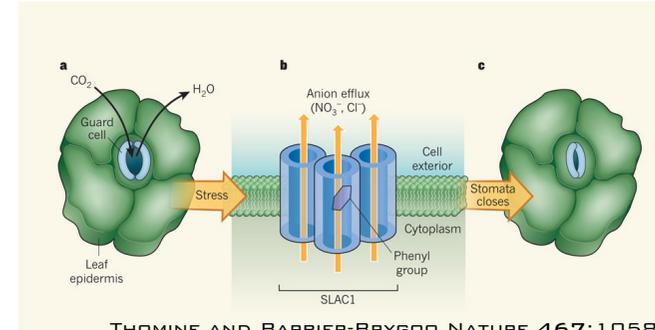




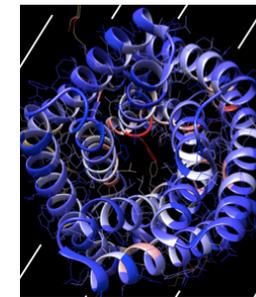
H. influenzae protein structure



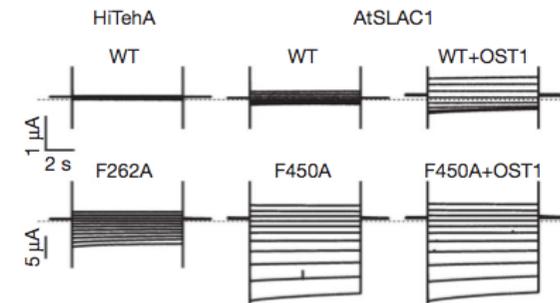
Functional hypothesis via homology to SLAC1

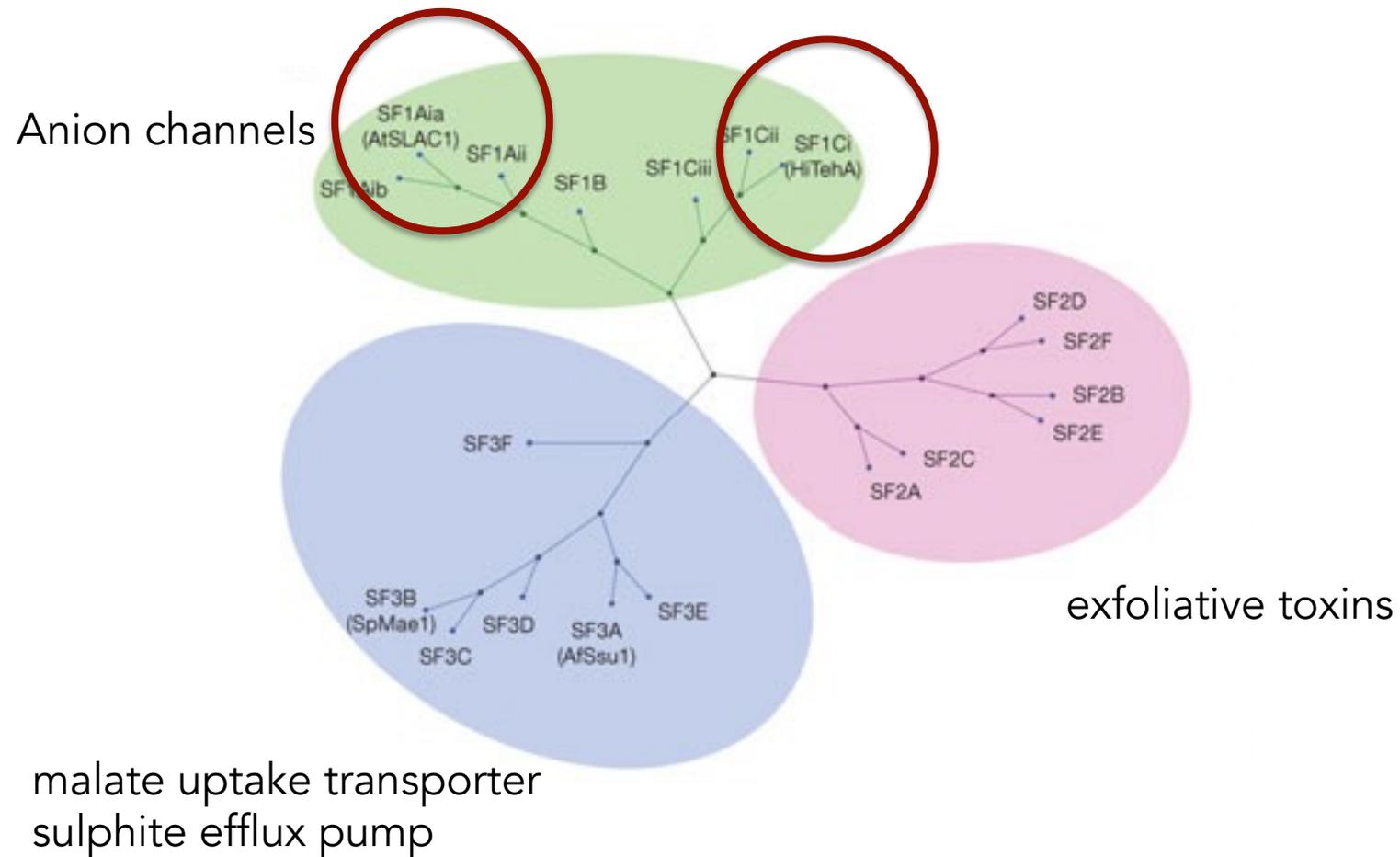


Identification potential functional residues using sequence conservation across the family and structural knowledge



Suggested experiments to test functional hypothesis





Exercise

Homology-based function annotation transfer #2

Protein Families

aidanbudd.github.io/ppisnd/trainingMaterial/marcoPunta/

EMBO Budapest excellence in life sciences

UNIVERSITÀ DI BOLOGNA

EMBO Practical Course

Course Program
Introduction To Linux Command-line
Introduction To PPI Networks
STRING
Network Visualization With Cytoscape
Chimera
Unseminar
Structure And Interfaces Of PPIs
Peptide And Protein Docking
MSA & Jalview
Protein Families
Protein Feature Prediction
Repeats And Low Complexity Regions
Intrinsically Disordered Proteins
Short Linear Motifs
REST Services
Molecular Dynamics

PROTEIN FAMILIES

by Marco Punta

EXERCISE 1

- 3m71.pdb
- 3m71_del.pdb
- PF03595_seed_mod.txt.aln

EXERCISE 2

- P29973.fasta
- mystery-protein.fasta

EXERCISE 3

- 2lhu.pdb
- family-building-exercise.fasta
- hmmer-ali.fasta
- jalview-ali.fasta

<http://aidanbudd.github.io/ppisnd/trainingMaterial/marcoPunta/>

>sp|P29973|CNGA1_HUMAN cGMP-gated cation channel alpha-1 OS=Homo sapiens
MKLSMKNNIINTQQSFVTMPNVIVPDIEKEIRRMENGACSSFSEDDDSASTSESEENENP
HARGSF SYKSLRKGGPSQREQYLPGAIALFNVNNSNKDQEPEEKKKKKKEKKS SDDKN
ENKNDPEKKKKKKDKEKKKKEEKSKDKKEEEKKEVVVIDPSGNTYYNWLFCITLPVMYNW
TMVIARACFDELQSDYLEYWLILDYVSDIVYLDIMFVRTRTGYLEQGLLVKEELKLINKY
KSNLQFKLDVLSLIPTDLLYFKLGWNYPEIRLNRLLRFSRMFEFFQRTETRTNYPNIFRI
SNLVMYIVIIHWNACVFYSISKAIGFGNDTWVYPDINDPEFGRLARKYVYSLYWSTLTL
TTIGETPPPVRDSEYVFVVVDFLIGVLIFATIVGNIGSMISNMNAARAEFQARIDAIKQY
MHFRNVSKDMEKRVIKWF DYLWTNKKTVDEKEVLKYL PDKLRAEIAINVHLDTLKKVRI F
ADCEAGLLVELVLKLQPQVYSPGDYICKKGDIGREMYIIKEGKLAVVADDGVTQFVVLSD
GSYFGEISILNIKGSKAGNRRTANIKSIGYSDLFCLSKDDLMEALTEY PDAKTMLEEK GK
QILMKDGLLDLNIANAGSDPKDLEEKVTRMEGSVDLLQTRFARILAEYESMQQK LKQRLT
KVEKFLKPLIDTEFSSIEGPGAESGPIDST

>mystery protein

MGNGSVKPKHSHKHPDGHSGNLTTDALRNKVTELERELRRKDAEIQEREYHLKELREQLSK
QTVAIAELTEELQNKCIQLNKLQDVVHMQGGSPLOASPDKVPLEVHRKTSGLVSLHSRRG
AKAGVSAEPTTRTYDLNKPPEFSFEKARVRKDSSEKKLITDALNKNQFLKRLDPQQIKDM
VECMYGRNYQQGSYIIKQGEPGNHIFVLAEGRLEVFQGEKLLSSIPMWTTFGELAILYNC
TRTASVKAITNVKTWALDREVFQONIMRRTAQARDEQYRNFLRSVSLKLNLPEDKLTKIID
CLEVEYYDKGDYIIREGEEGSTFFILAKGKVKVTQSTEGHDQPQLIKTLQKGEYFGEKAL
ISDDVRSANIIAEENDVACLVIDRETFNQTVGTFEELQKYLEGYVANLNRDDEKRHAKRS
MSNWKLSKALSLEMIQLKEKVARFSSSSPFQNLLEIIATLGVGGFGRVELVKVKNENVAF
MKCIRKKHIVDTKQQEHVYSEKRILEELCSPFIVKLYRTFKDNKYVYMLLEACLGGELWS
ILRDRGSFDEPTSKFCVACVTEAFDYLHRLGIIYRDLKPENLILDAEGYLKLVDFGFAKK
IGSGQKTWTFCGTPEYVAPEVILNKGHDFSVDVFWSLGILVYELLTGNPPFSGVDQMMTYN
LILKGIEKMDFPRKITRRPEDLIRRLCRQNPTERLGNLKNGINDIKKHRWLNGFNWEGLK
ARSLPSPLQRELKGPIDHSYFDKYPPEKGMPPDELSGWDKDF

<http://blast.ncbi.nlm.nih.gov/Blast.cgi?PAGE=Proteins>

Align Sequences Protein BLAST

[blastn](#)**[blastp](#)**[blastx](#)[tblastn](#)[tblastx](#)BLASTP programs search protein subjects using a protein query. [more...](#)

Enter Query Sequence

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

```
TTIGETPPPVRDSEYVFFVVDFLIGVLIFATIVGNIGSMISNMNAARAEFQARIDAIKQY
MHFRNVSKDMEKRVIKWFDYLWTNKKTVDEKEVLKYLDPDKLRAEIAINVHLDLTKKVRIF
ADCEAGLLVELVLKLPQVYSPGDYICKKGDIGREMYIIKEGKLAVVADDGVTQFVVLSD
GSYFGEISILNIKGSKAGNRRRTANIKSIGYSDLFCLSKDDLMEALTEYPDAKTMLEEKGK
QILMKDGLLDLNIANAGSDPKDLEEKVTRMEGSVDLLQTRFARILAEYESMQQKLRQLT
KVEKFLKPLIDTEFSSIEGPGAESGPIDST
```

From

To

Or, upload 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 

```
MSNWKLSKALSLEMIQLKEKVARFSSSSPFQNLIIATLGVGGFGRVELVKVKNENVAFA
MKCIRKKHIVDTKQQEHVYSEKRILEELCSPFIVKLYRTFKDNKYVYMLLEACLGGELWS
ILRDRGSFDEPTSKFCVACVTEAFDYLHRLGIIYRDLKPENLILDAEGYLKLVDFGFAK
IGSGQKTWTFCCGTPEYVAPEVILNKGHDFSVDFWSLILVYELLTGNPPFSGVDQMMTYN
LILKGIEKMDFFPKITRRPEDLIRRLCRQNPTERLGNLKNGINDIKKHRWLNGFNWEGLK
ARSLPSPLORELKGPIDHSYFDKYPPEKGMPPDELSGWDKDF
```

From

To

Or, upload file

 No file selected. 

Alignments

Download [Graphics](#) Sort by: E value

unnamed protein product

Sequence ID: lcl|Query_22995 Length: 762 Number of Matches: 3

Range 1: 281 to 383 [Graphics](#)

▼ Next Match ▲ Previous Match

Score	Expect	Method	Identities	Positives	Gaps
41.2 bits(95)	1e-07	Compositional matrix adjust.	30/110(27%)	54/110(49%)	11/110(10%)
Query 474	LKKVRIFADCEAGLLVELVLKLPQVSPGDYICKKGDIGREMYIIKEGKLAVV----	AD	529		
	L+ V + + L +++ L+ + Y GDYI ++G+ G +I+ +GK+ V				
Sbjct 281	LRSVSLKLNLPEDKLTKIIDCLEVEYYDKGDYIIREGEEGSTFFILAKGKVKVTQSTEGH		340		
Query 530	DGVTQFVVLSDGSYFGEISILNIKSGKAGNRRRTANIKSIGYSDLFCLSKD		579		
	D L G YFGE +++ + + R+ANI + +D+ CL D				
Sbjct 341	DQPQLIKTLQKGEYFGEKALI-----SDDVRSANIAIA-EENDVACLVID		383		

Range 2: 161 to 260 [Graphics](#)

▼ Next Match ▲ Previous Match ▲ First Match

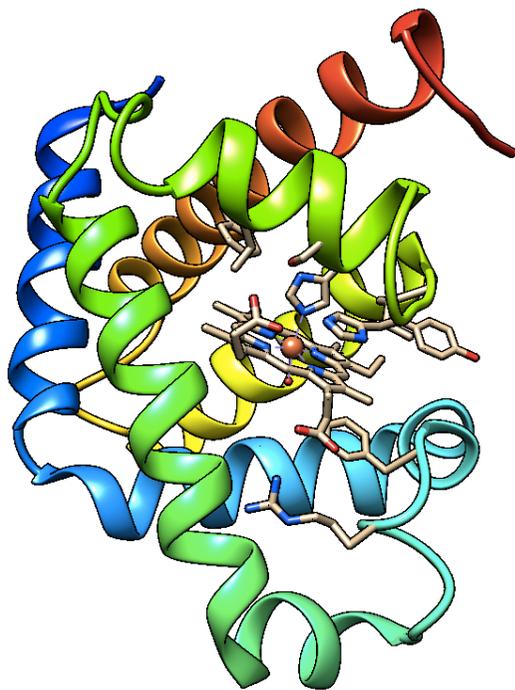
Score	Expect	Method	Identities	Positives	Gaps
38.1 bits(87)	1e-06	Compositional matrix adjust.	26/108(24%)	53/108(49%)	8/108(7%)
Query 472	DTLKKVRIFADCEAGLLVELVLKLPQVSPGDYICKKGDIGREMYIIKEGKLAVVADDG		531		
	D L K + + + ++V + + Y G YI K+G+ G ++++ EG+L V +				
Sbjct 161	DALNKNQFLKRLDPQQIKDMVECMYGRNYQQGSYIIKQGEPEGNHIFVLAEGRLEVFQGEK		220		
Query 532	VTQFVVLSDGSYFGEISILNIKSGKAGNRRRTANIKSIGYSDLFCLSKD		579		
	+ + + FGE++IL RTA++K+I + L ++				
Sbjct 221	LLSSIPM--WTFFGELAIL-----YNCRTASVKAITNVKTWALDRE		260		

Range 3: 593 to 649 [Graphics](#)

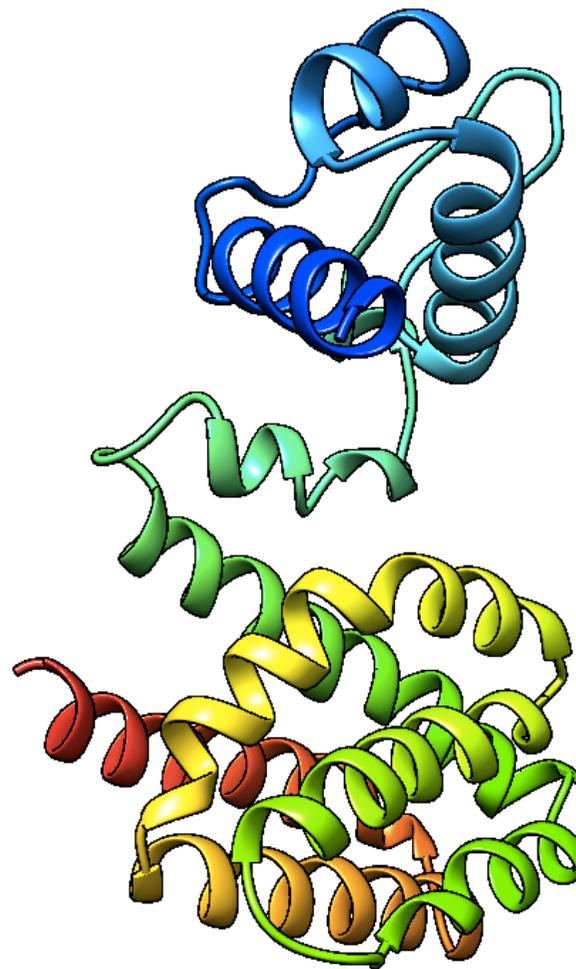
▼ Next Match ▲ Previous Match ▲ First Match

Score	Expect	Method	Identities	Positives	Gaps
22.3 bits(46)	0.081	Compositional matrix adjust.	17/58(29%)	27/58(46%)	7/58(12%)
Query 317	VFYSISKAIGFGNDTWVY---PDINDPEFGRLARKYVYSL-YWS--TLTLTTIGETPP		368		
	V + +K IG G TW + P+ PE L + + +S+ +WS L + PP				
Sbjct 593	VDFGFAKKIGSGQKTWTF CGTPEYVAPEV-ILNKGHDFSVDFWSLGLVYELLTGNPP		649		

Mystery protein is a cGMP-gated cation channel?



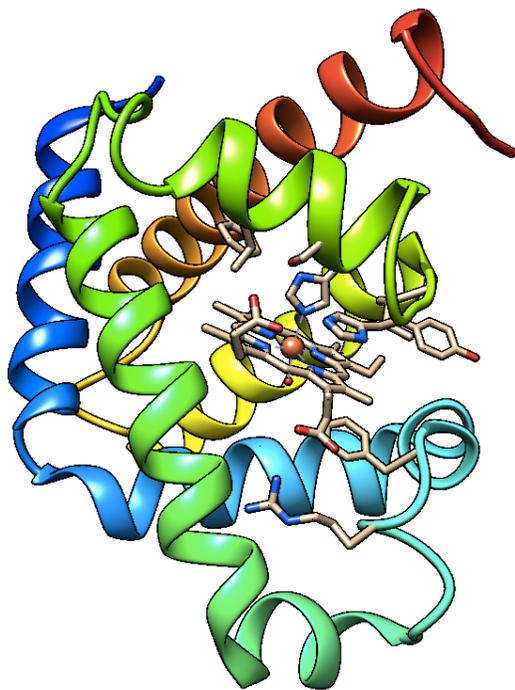
1MBN



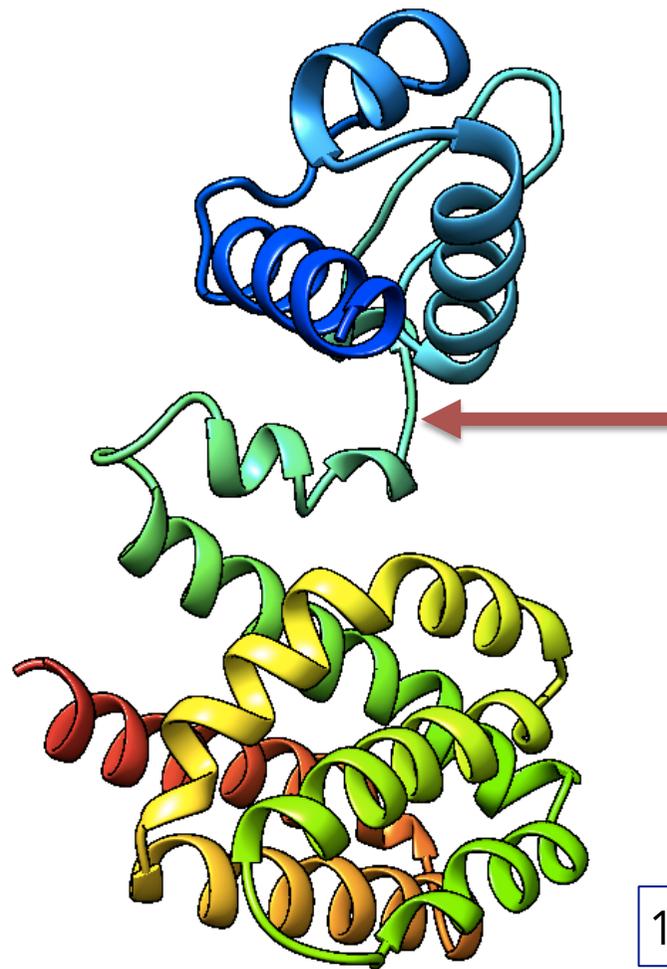
1HW2

Colour Scheme:





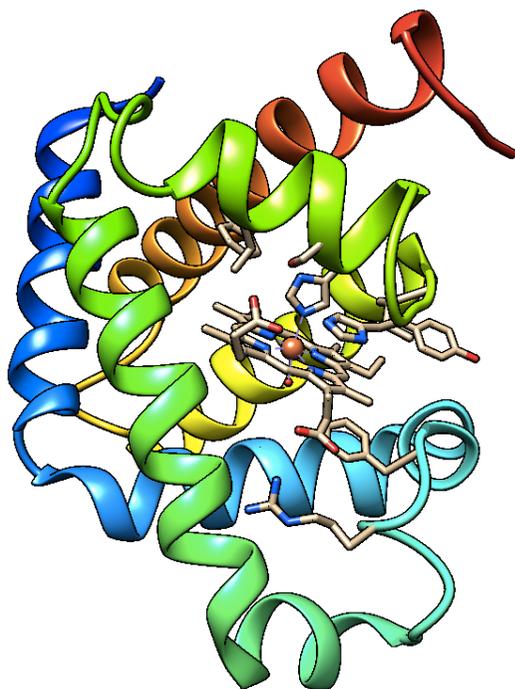
1MBN



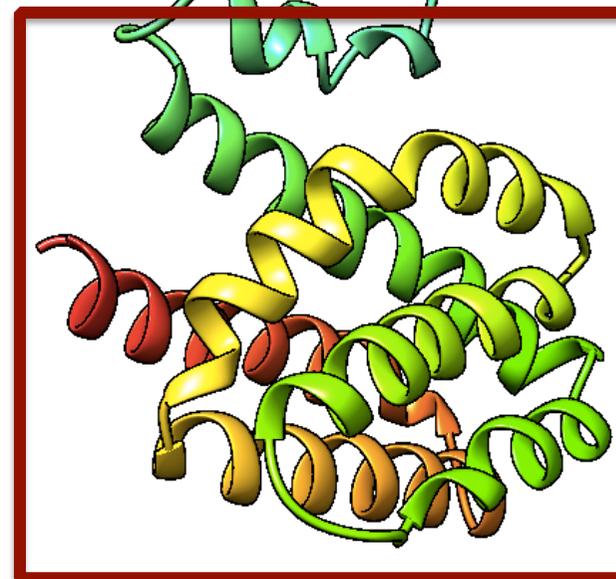
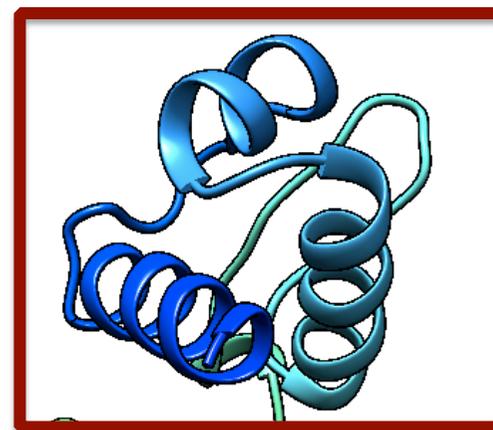
1HW2

Colour Scheme:





1MBN



1HW2

Colour Scheme:

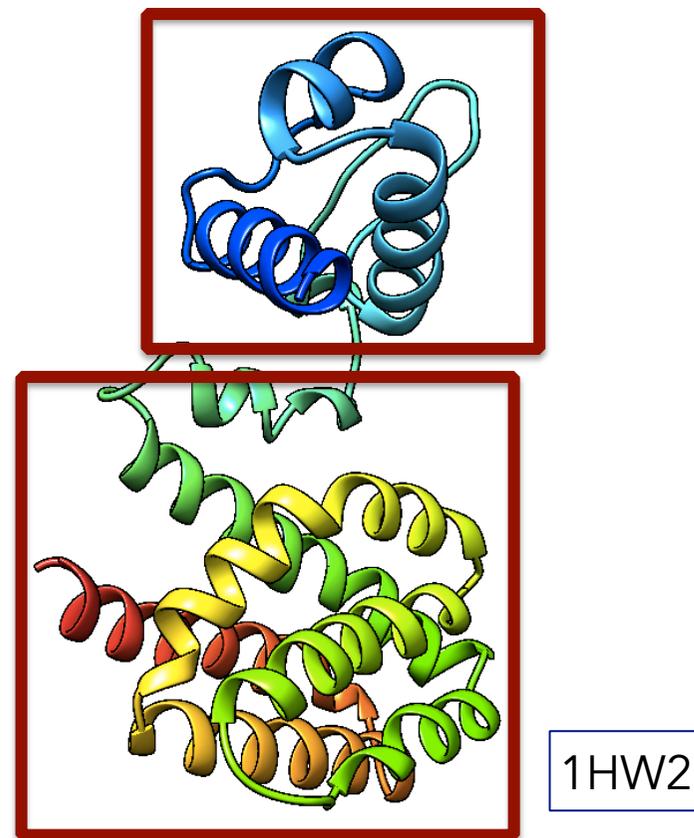


N'

C'

Definition (Wikipedia):

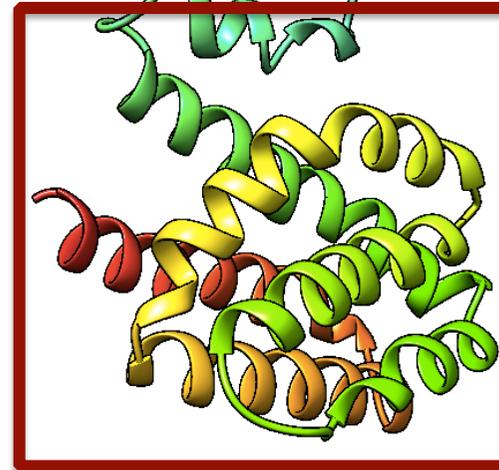
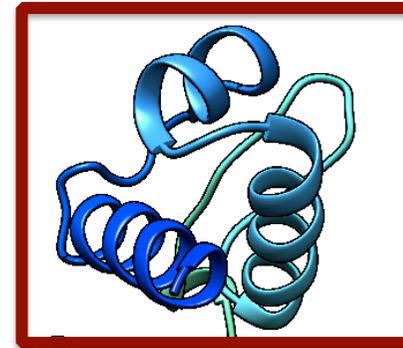
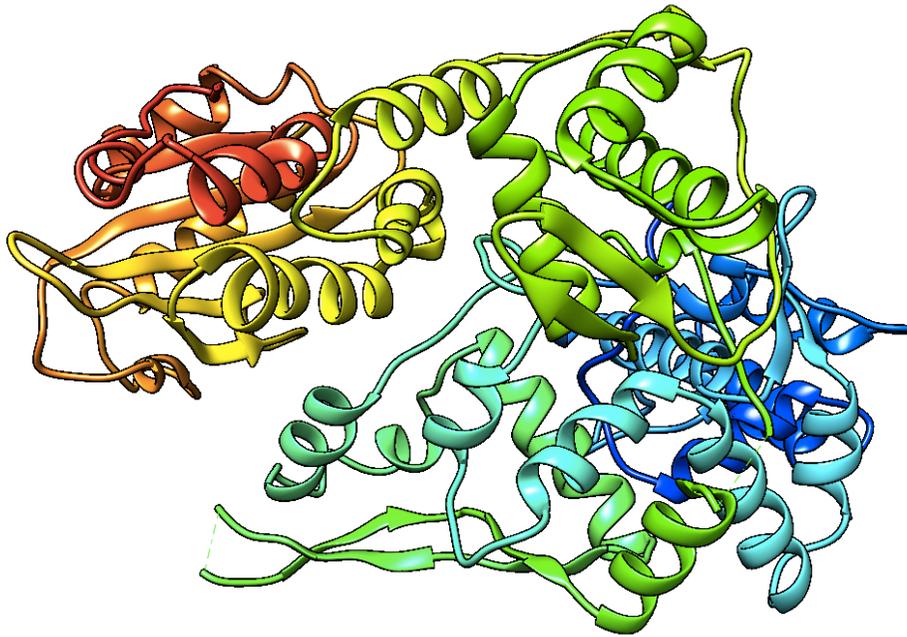
A protein domain is a conserved part of a given protein sequence and (tertiary) structure that can evolve, function, and exist independently of the rest of the protein chain. Each domain forms a compact three-dimensional structure and often can be independently stable and folded.



Colour Scheme:



1FOK

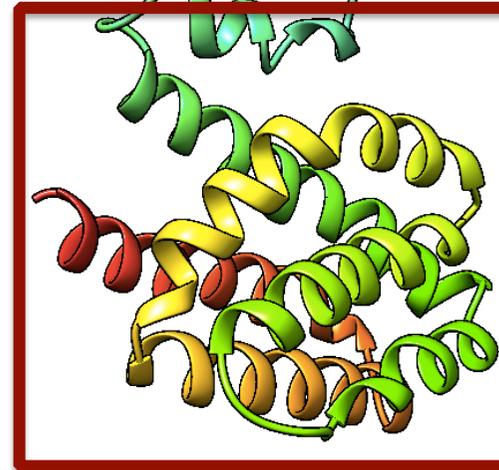
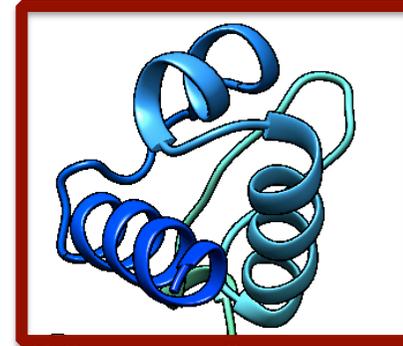
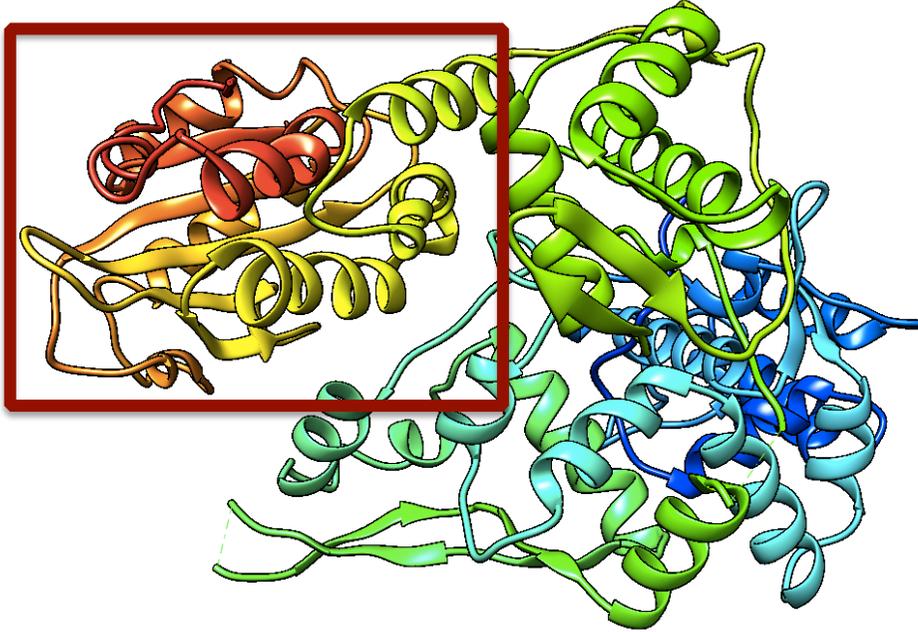


1HW2

Colour Scheme:



1FOK

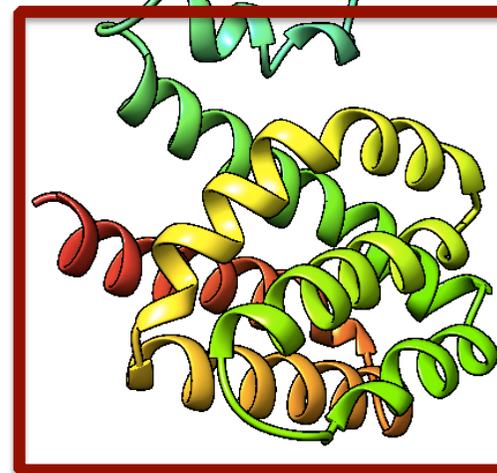
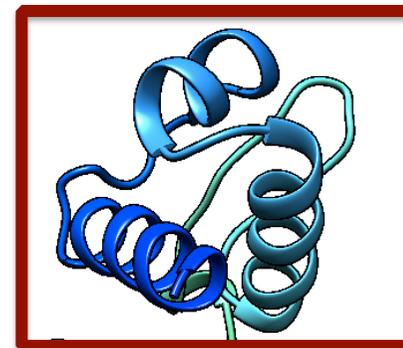
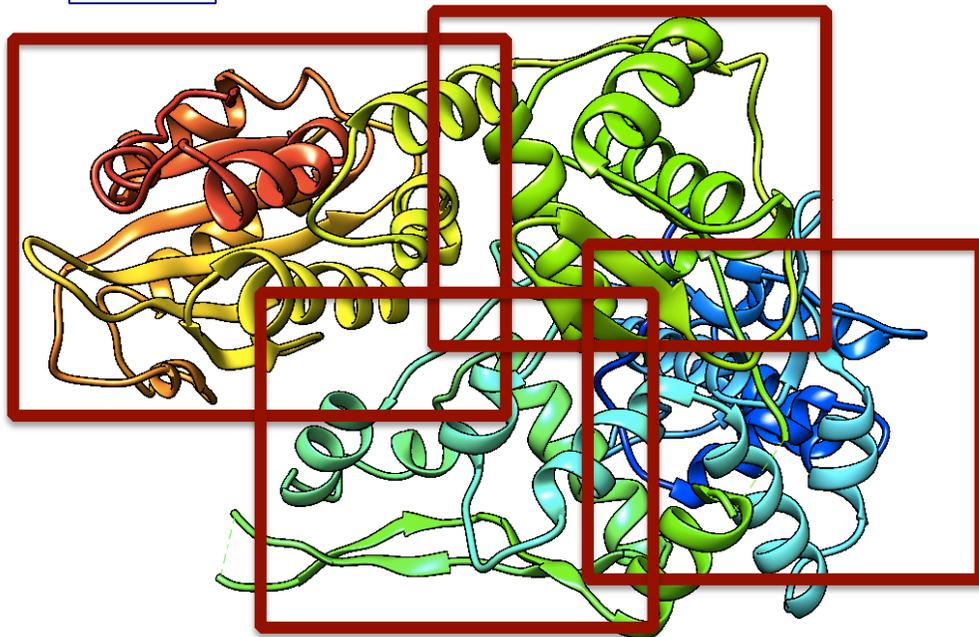


1HW2

Colour Scheme:



1FOK

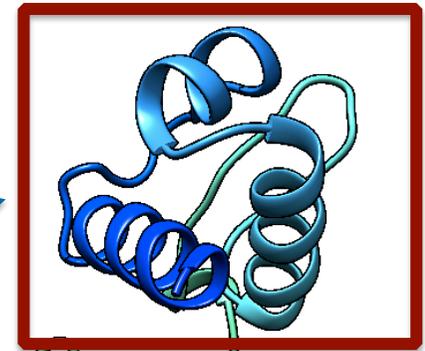
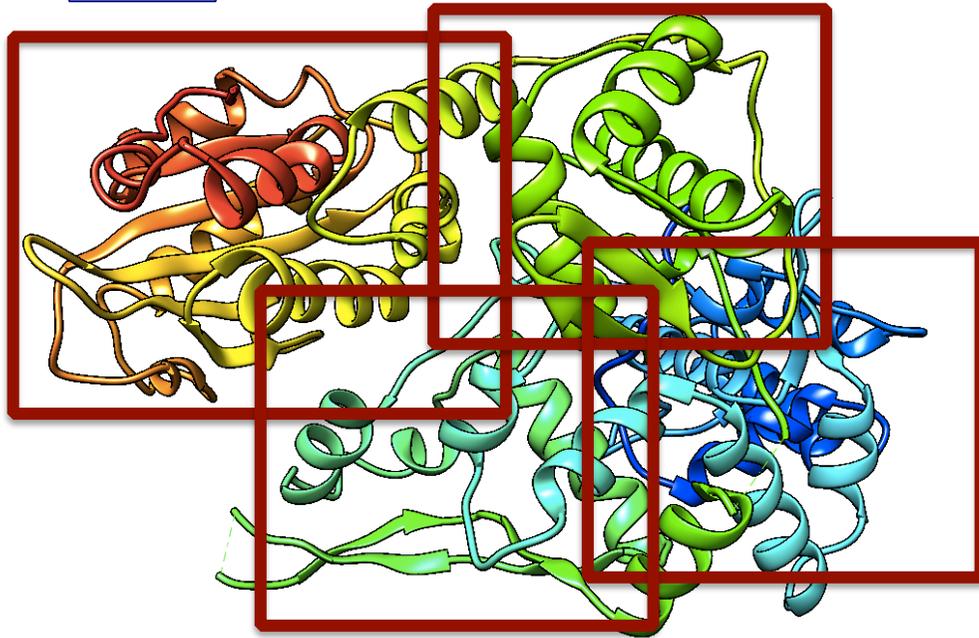


1HW2

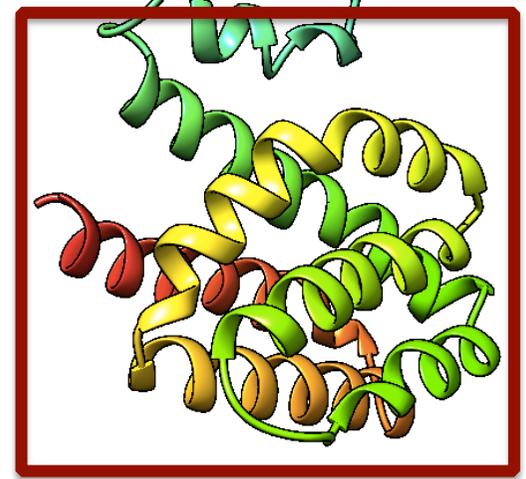
Colour Scheme:



1FOK



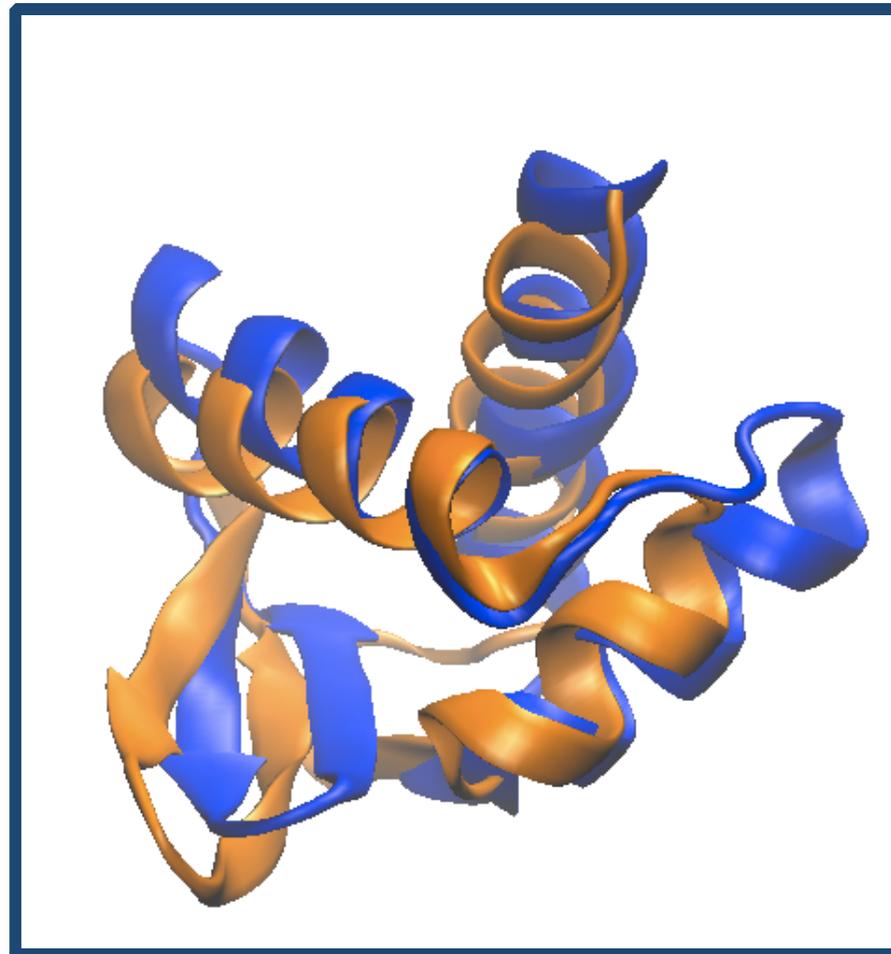
1HW2



Colour Scheme:



Winged helix domain (WHD)

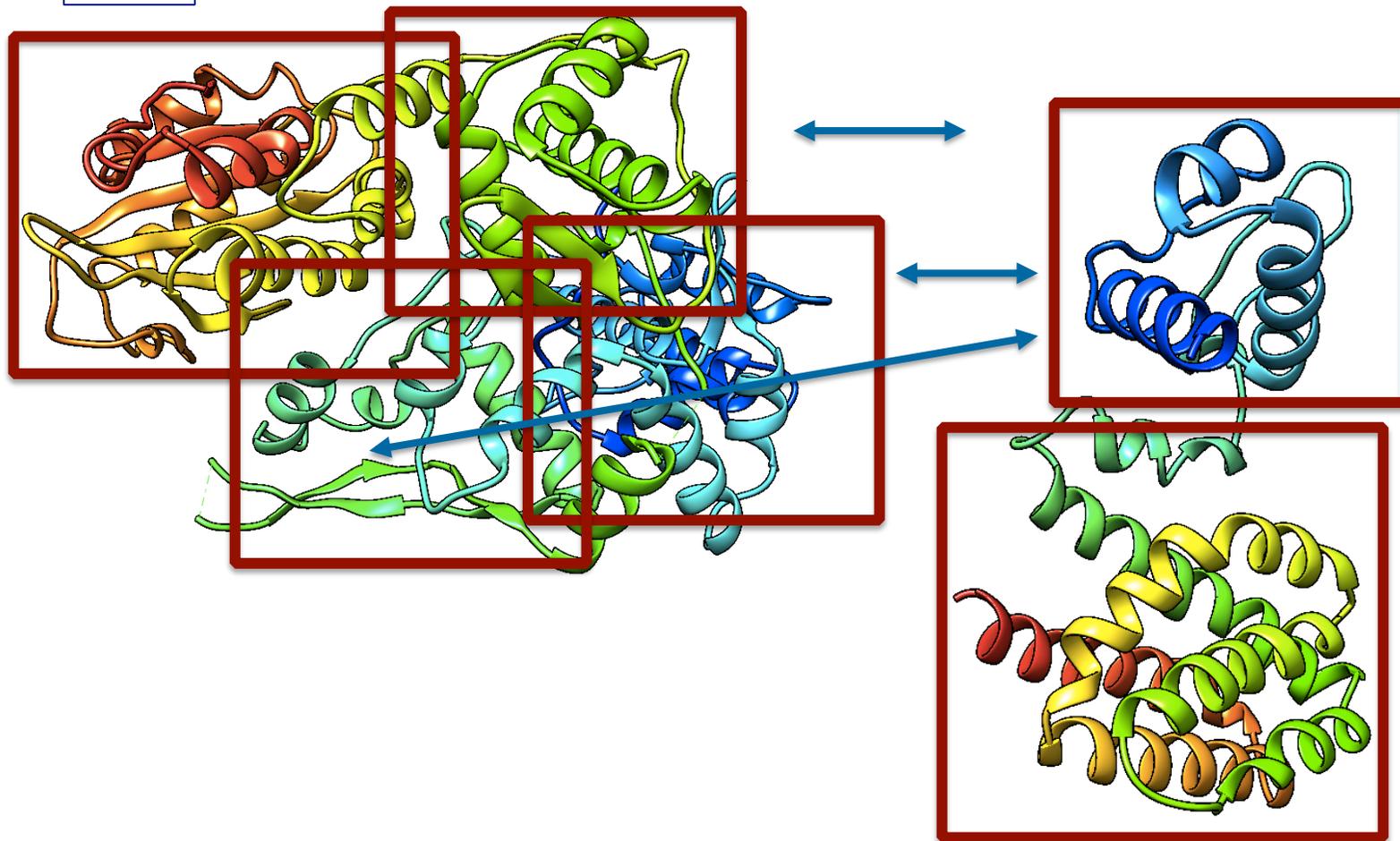


Z-score = 4.0

%id = 8%

RMSD = 2.7Å

1FOK

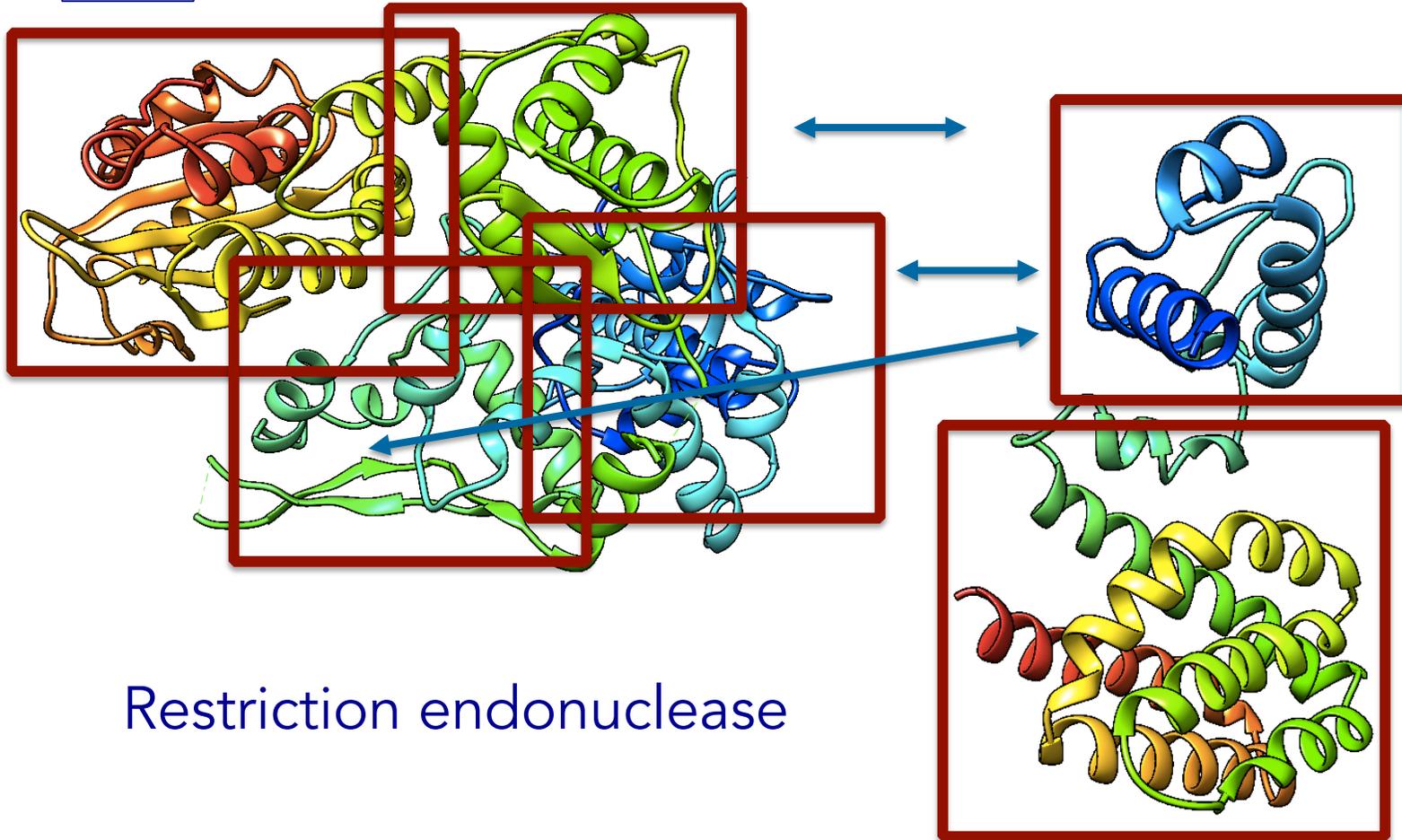


1HW2

Colour Scheme:



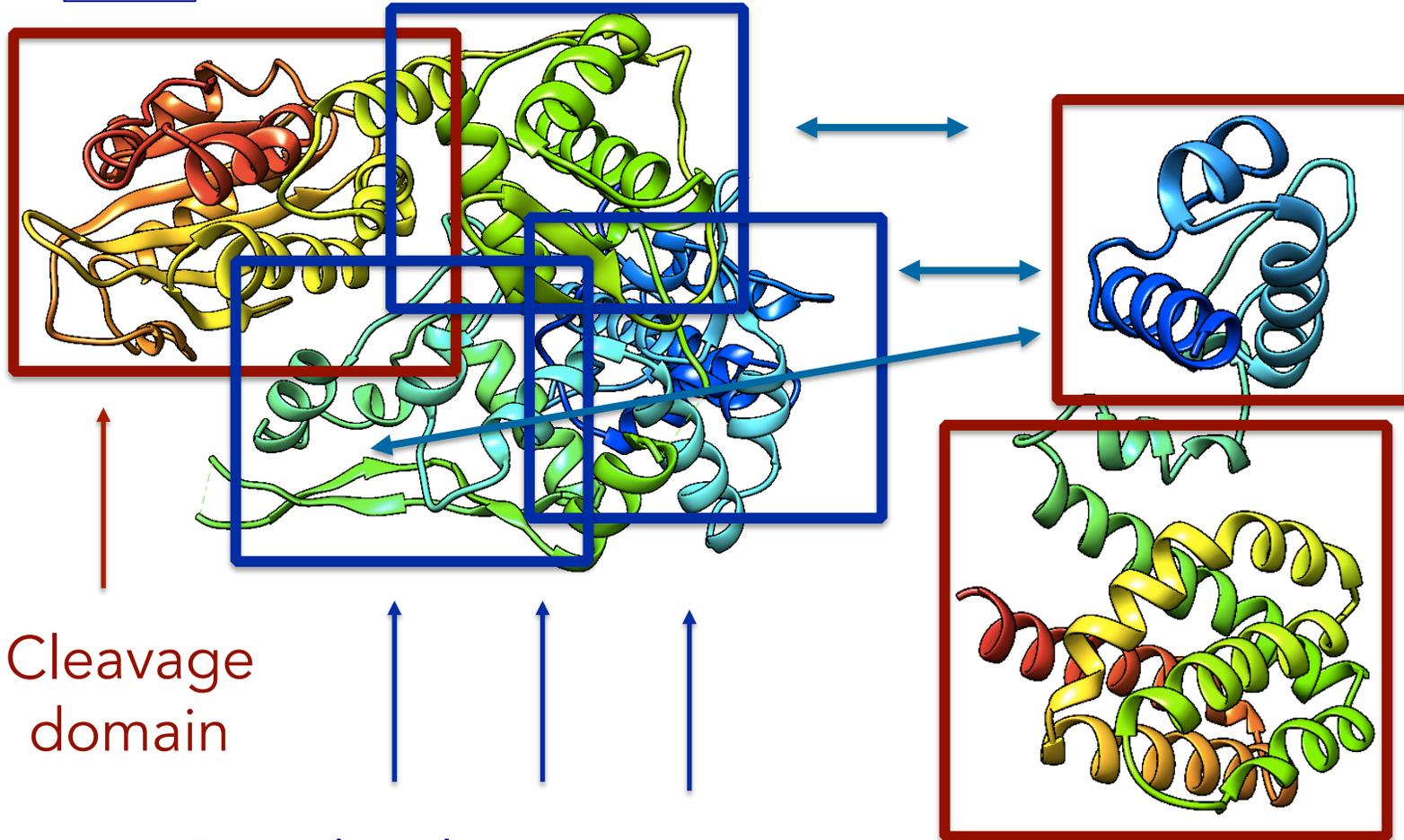
1FOK



Restriction endonuclease

1HW2

1FOK



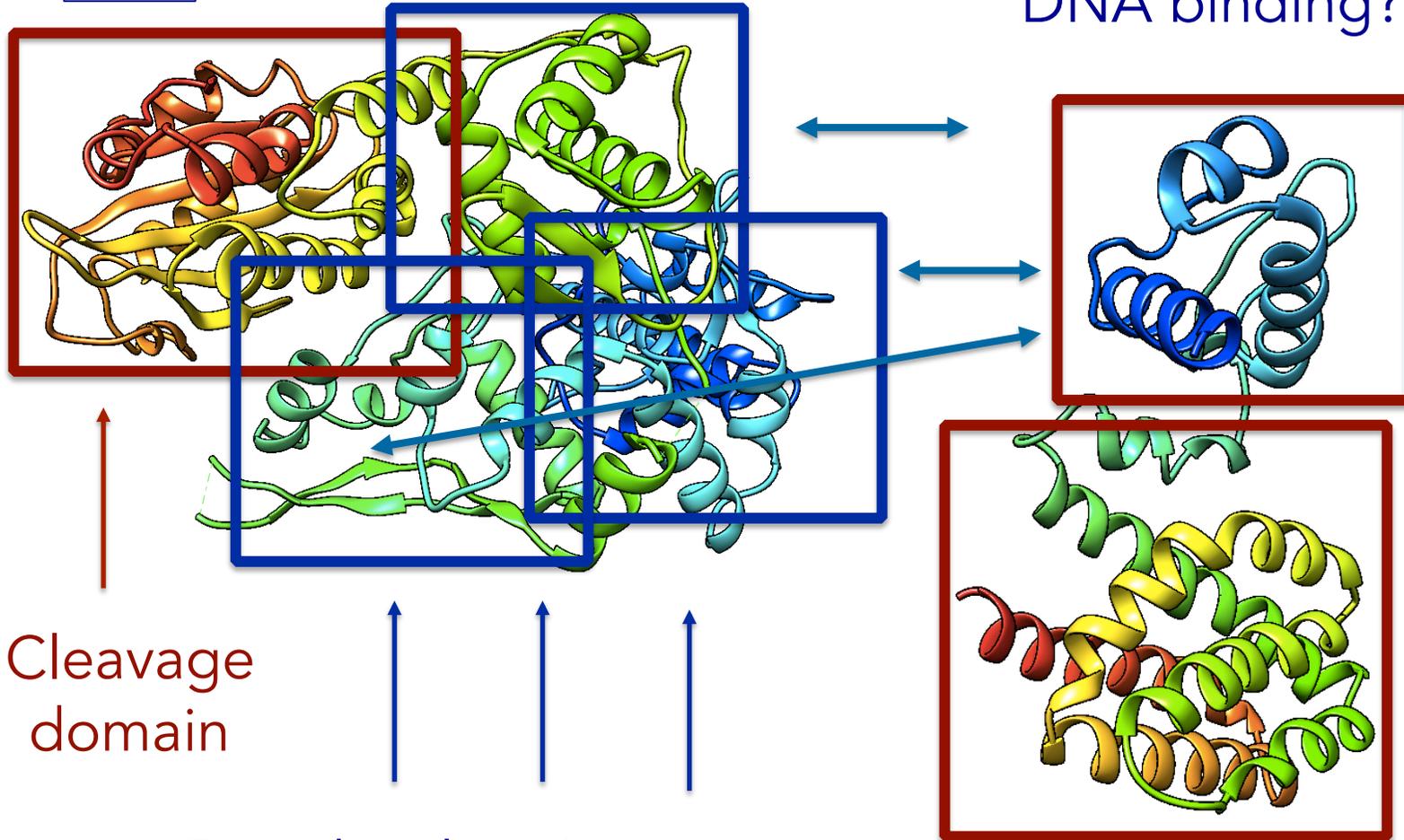
Cleavage domain

DNA binding (targeting to a specific DNA sequence)

1HW2

1FOK

DNA binding?



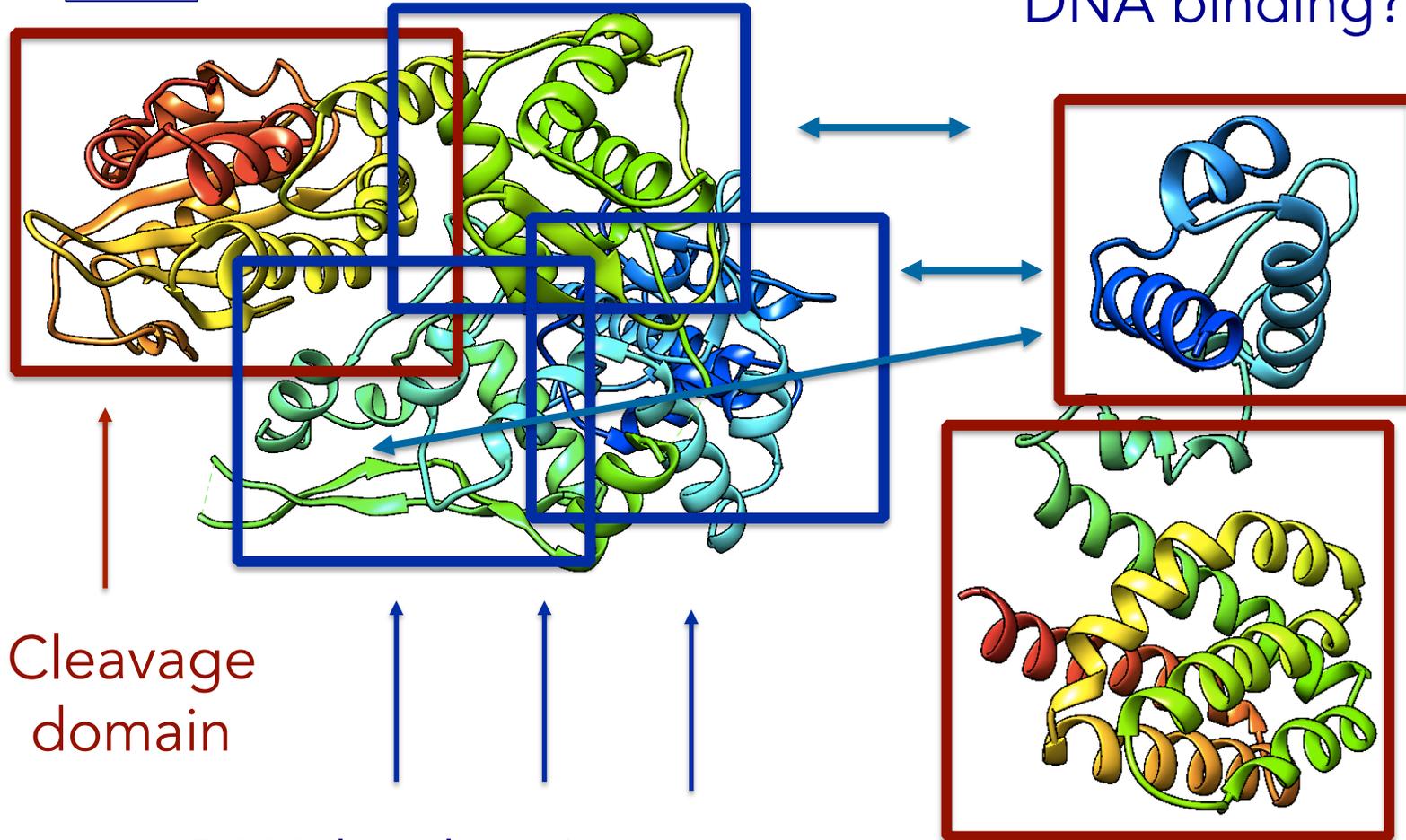
Cleavage domain

DNA binding (targeting to a specific DNA sequence)

1HW2

1FOK

DNA binding?



Cleavage domain

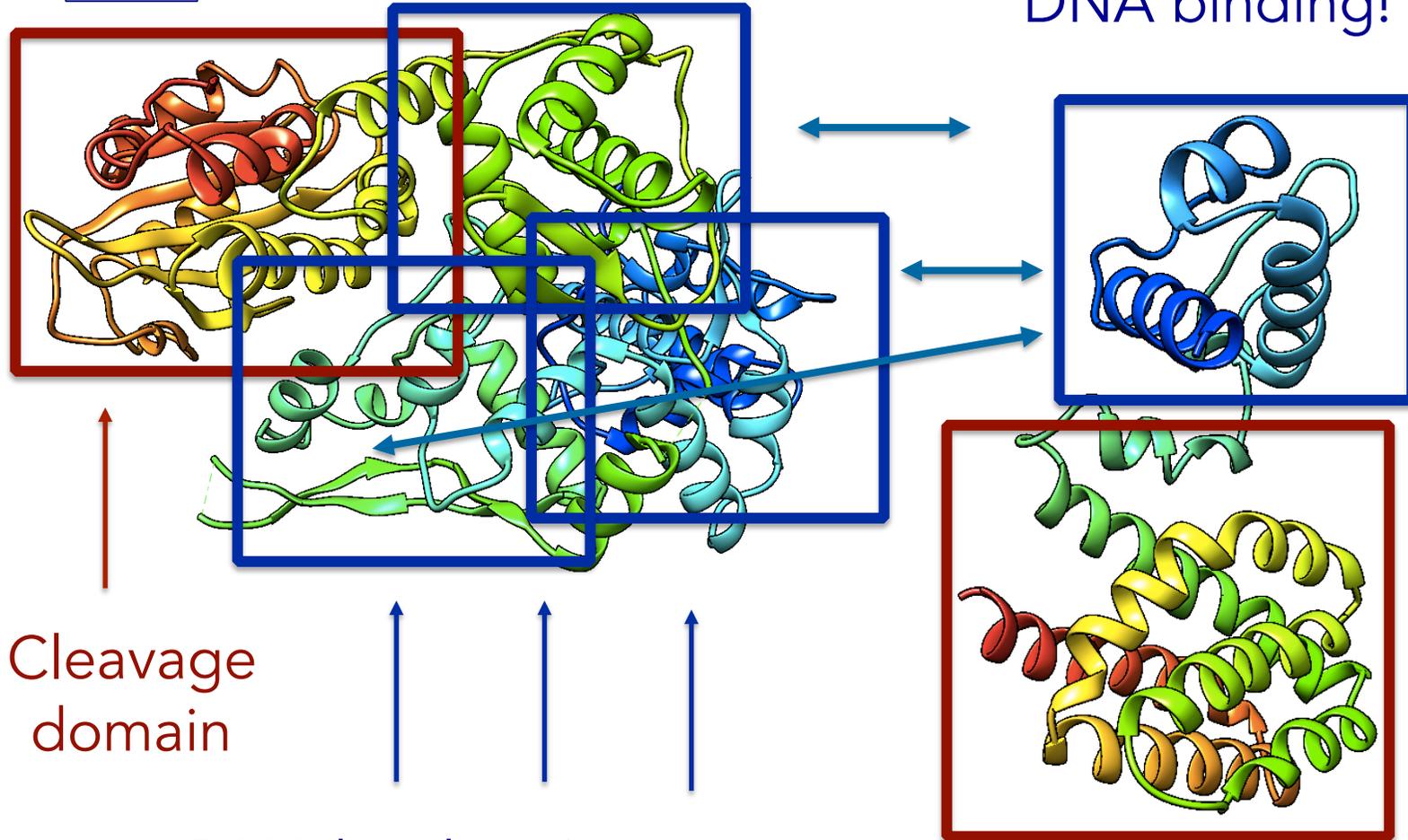
DNA binding (targeting to a specific DNA sequence)

1HW2

?

1FOK

DNA binding!



Cleavage domain

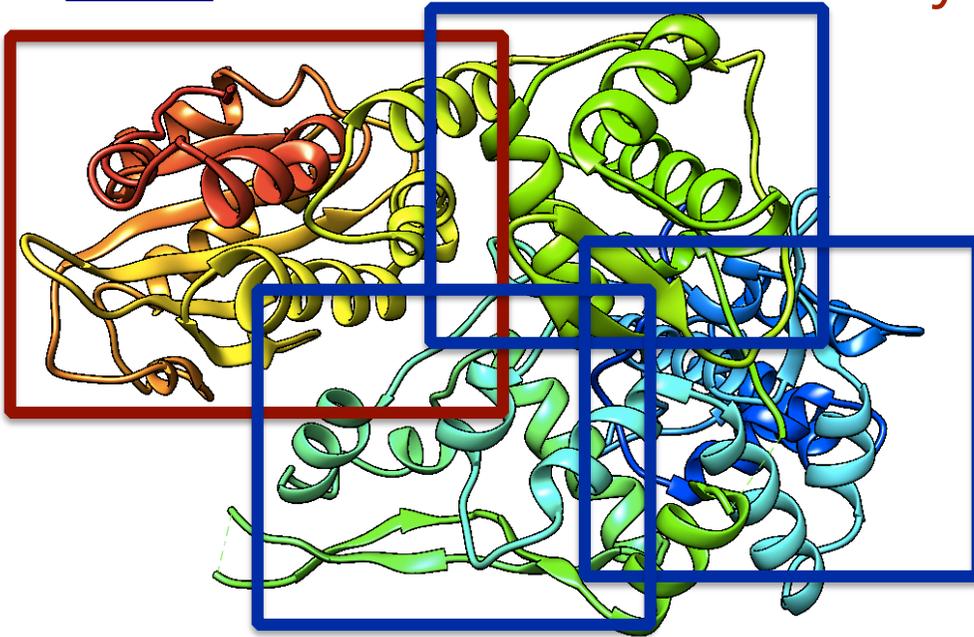
DNA binding (targeting to a specific DNA sequence)

acyl-CoA binding domain controls affinity

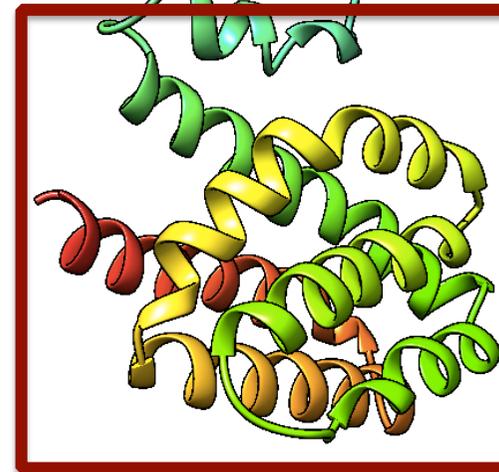
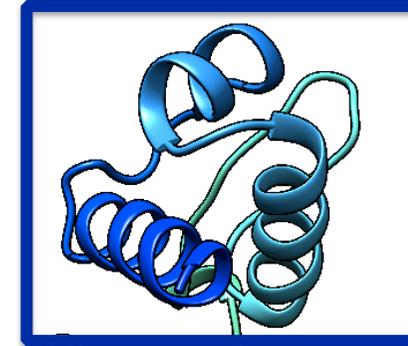
1HW2

1FOK

“syntactical change”



Restriction endonuclease



1HW2

Transcription factor

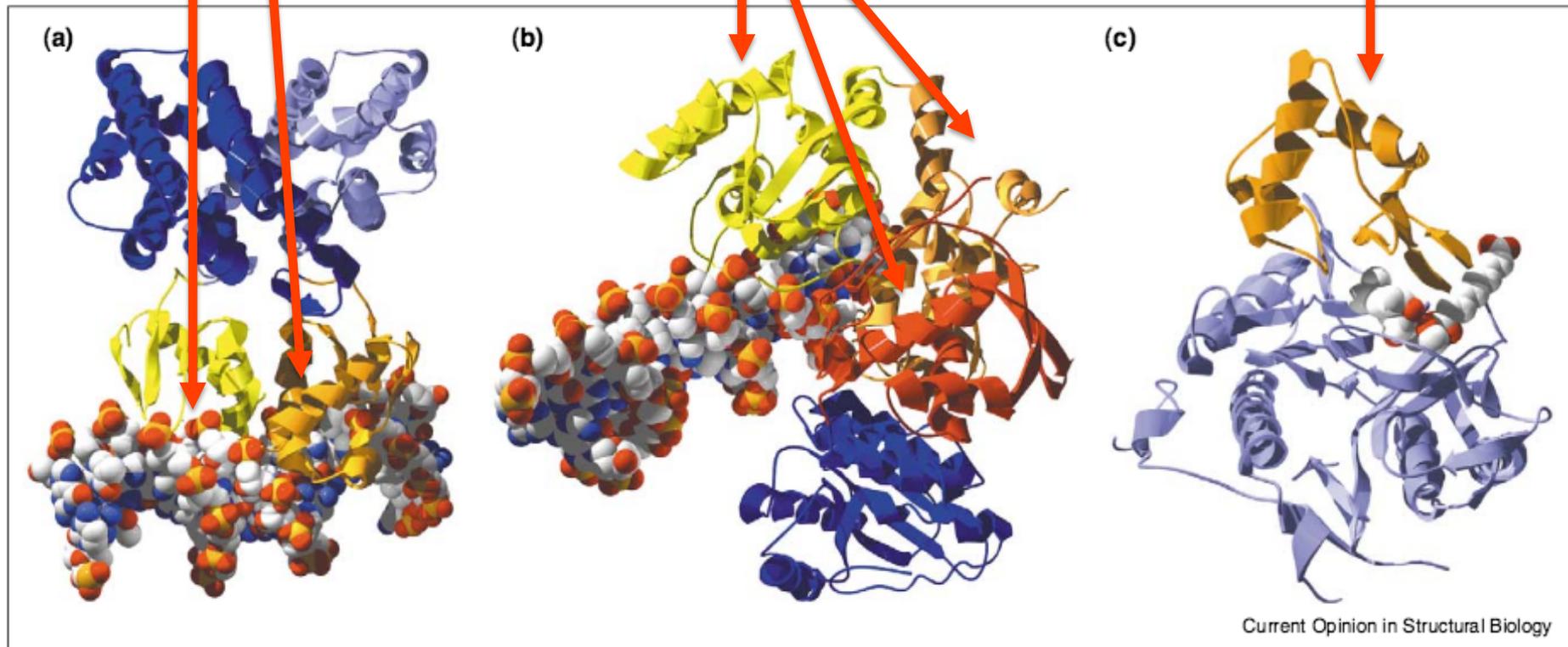
Semantic change

Marco Punta

DNA binding

DNA binding

substrate specificity pocket



Transcription
factor

Restriction
endonuclease

Human methionine
aminopeptidase 2

“syntactical change”

Marco Punta

DNA sequence
recognised



Restriction endonuclease

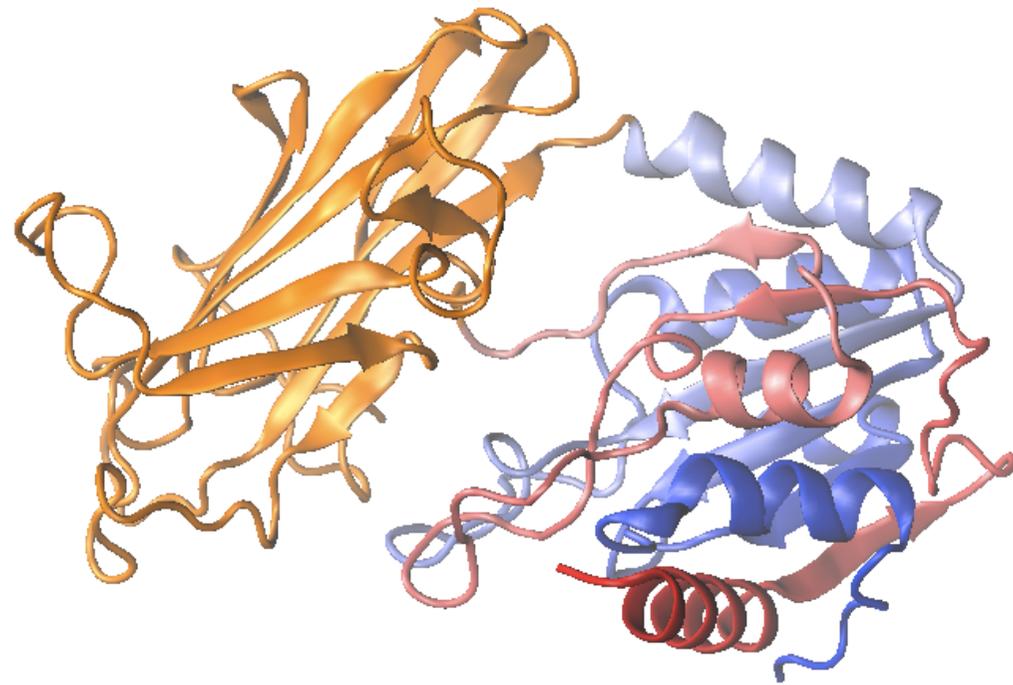
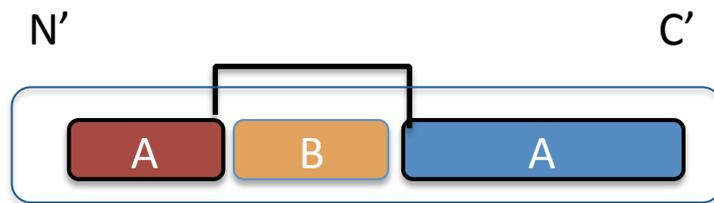
5'-GGATG-3'

Transcription factor

5'-TGGNNNNNCCA-3'

“Nested” domains

Marco Punta



3ABZ-truncated

Alignments

Download [Graphics](#) Sort by:

unnamed protein product

Sequence ID: lcl|Query_22995 Length: 762 Number of Matches: 3

Range 1: 281 to 383 [Graphics](#)

▼ Next Match ▲ Previous Match

Score	Expect	Method	Identities	Positives	Gaps
41.2 bits(95)	1e-07	Compositional matrix adjust.	30/110(27%)	54/110(49%)	11/110(10%)
Query 474	LKKVRIFADCEAGLLVELVLKLPQVSPGDYICKKGDIGREMYIIKEGKLAVV----	AD	529		
	L+ V + + L +++ L+ + Y GDYI ++G+ G +I+ +GK+ V				
Sbjct 281	LRSVSLKKNLPEDKLTKIIDCLEVEYYDKGDYIIREGEEGSTFFILAKGKVKVTQSTEGH		340		
Query 530	DGVTQFVVLSDGSYFGEISILNIKSGKAGNRRRTANIKSIGYSDLFCLSKD		579		
	D L G YFGE +++ + + R+ANI + +D+ CL D				
Sbjct 341	DQPQLIKTLQKGEYFGEKALI-----SDDVRSANIAIA-EENDVACLVID		383		

Range 2: 161 to 260 [Graphics](#)

▼ Next Match ▲ Previous Match ▲ First Match

Score	Expect	Method	Identities	Positives	Gaps
38.1 bits(87)	1e-06	Compositional matrix adjust.	26/108(24%)	53/108(49%)	8/108(7%)
Query 472	DTLKKVRIFADCEAGLLVELVLKLPQVSPGDYICKKGDIGREMYIIKEGKLAVVADDG		531		
	D L K + + + ++V + + Y G YI K+G+ G ++++ EG+L V +				
Sbjct 161	DALNKNQFLKRLDPQQIKDMVECMYGRNYQQGSYIIKQGEPGNHIFVLAEGRLEVFQGEK		220		
Query 532	VTQFVVLSDGSYFGEISILNIKSGKAGNRRRTANIKSIGYSDLFCLSKD		579		
	+ + + FGE++IL RTA++K+I + L ++				
Sbjct 221	LLSSIPM--WTFFGELAIL-----YNCRTASVKAITNVKTWALDRE		260		

Range 3: 593 to 649 [Graphics](#)

▼ Next Match ▲ Previous Match ▲ First Match

Score	Expect	Method	Identities	Positives	Gaps
22.3 bits(46)	0.081	Compositional matrix adjust.	17/58(29%)	27/58(46%)	7/58(12%)
Query 317	VFYSISKAIGFGNDTWVY---PDINDPEFGRLARKYVYSL-YWS--TLTLTTIGETPP		368		
	V + +K IG G TW + P+ PE L + + +S+ +WS L + PP				
Sbjct 593	VDFGFAKKIGSGQKTWTF CGTPEYVAPEV-ILNKGHDFSVDFWSLGLVYELLTGNPP		649		

cGMP-gated cation channel alpha-1

P29973 (CNGA1_HUMAN)

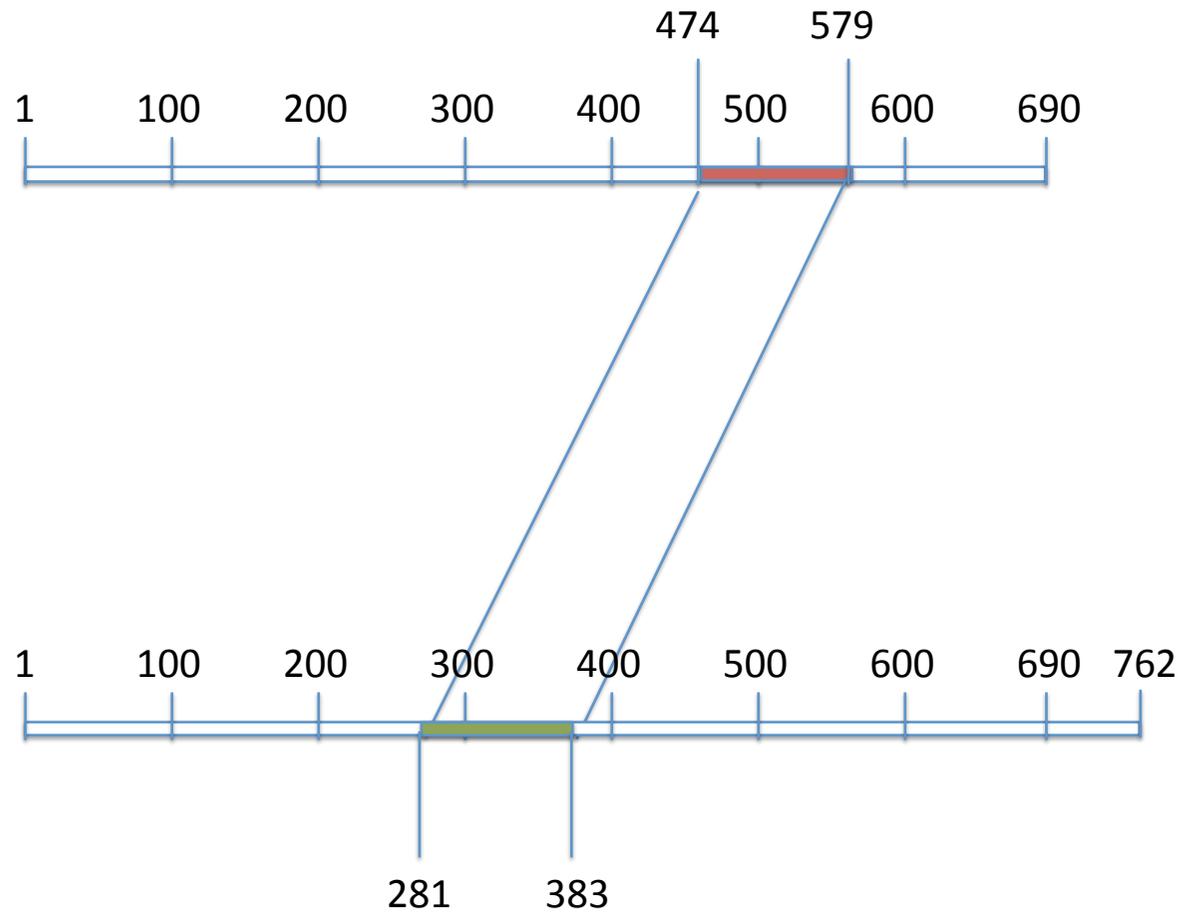
Marco Punta



cGMP-gated cation channel alpha-1

P29973 (CNGA1_HUMAN)

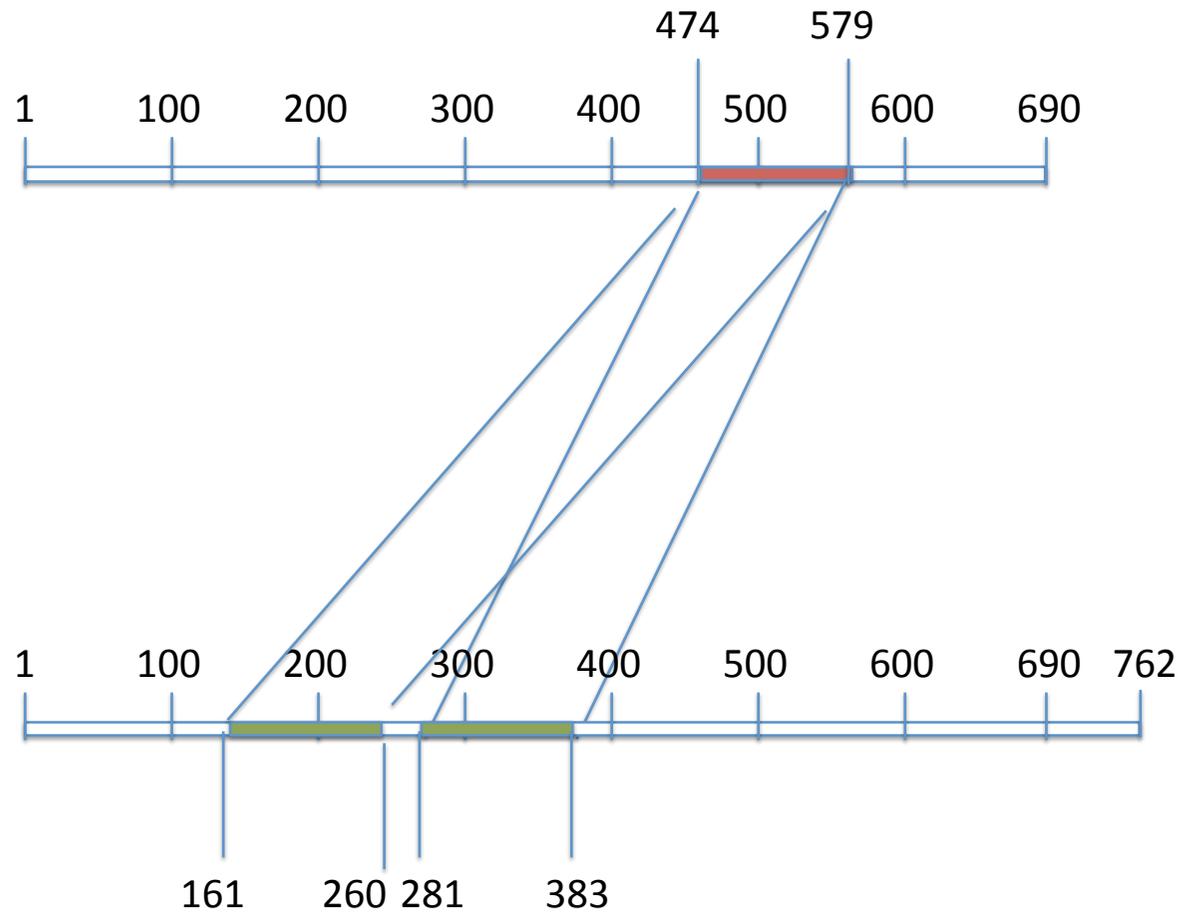
Marco Punta



cGMP-gated cation channel alpha-1

P29973 (CNGA1_HUMAN)

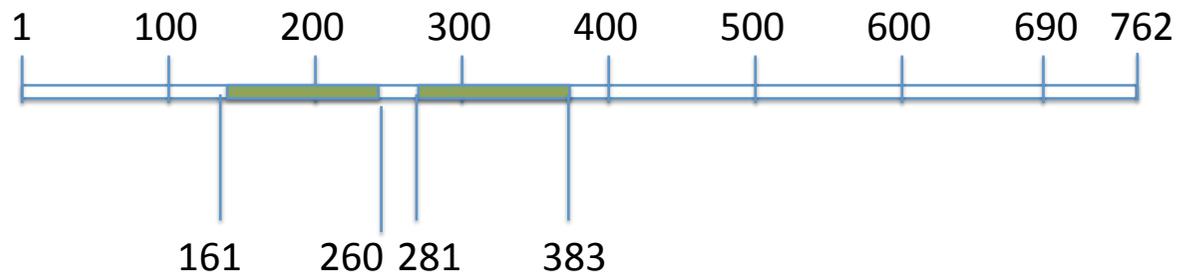
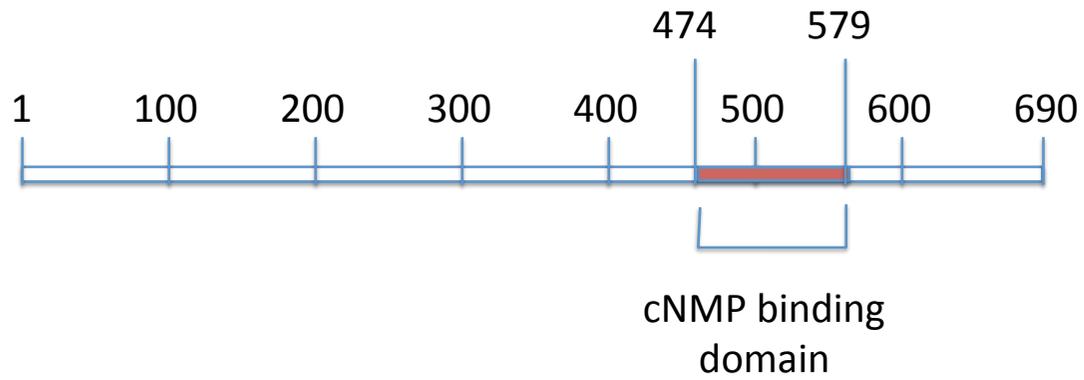
Marco Punta



cGMP-gated cation channel alpha-1

P29973 (CNGA1_HUMAN)

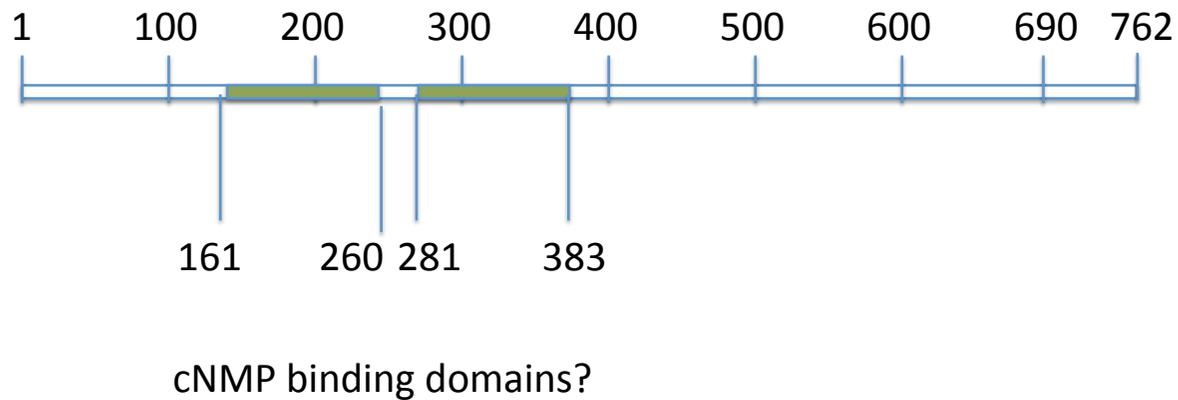
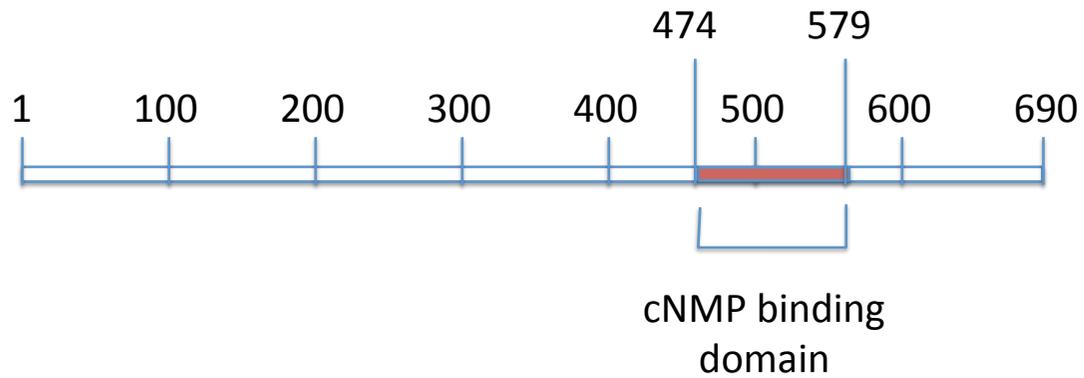
Marco Punta



cGMP-gated cation channel alpha-1

P29973 (CNGA1_HUMAN)

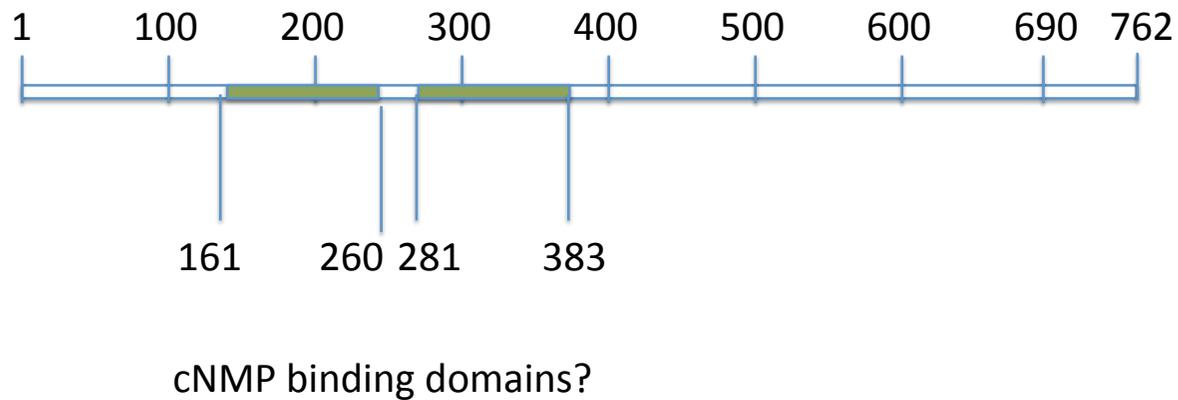
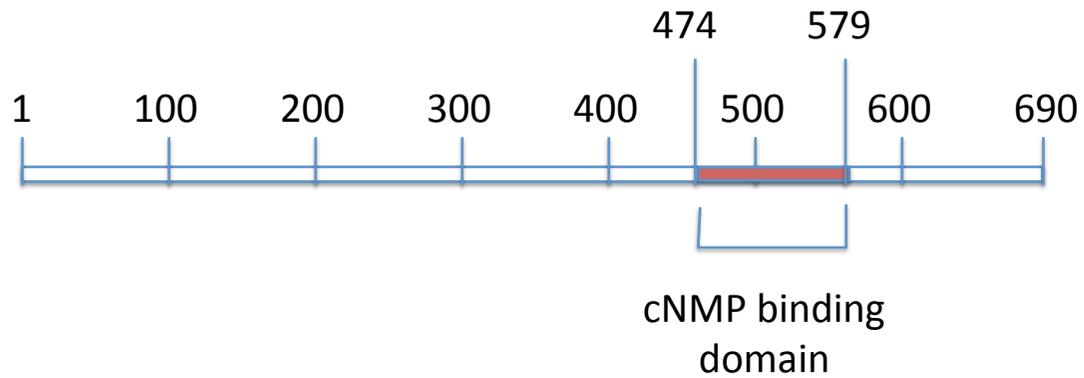
Marco Punta



cGMP-gated cation channel alpha-1

P29973 (CNGA1_HUMAN)

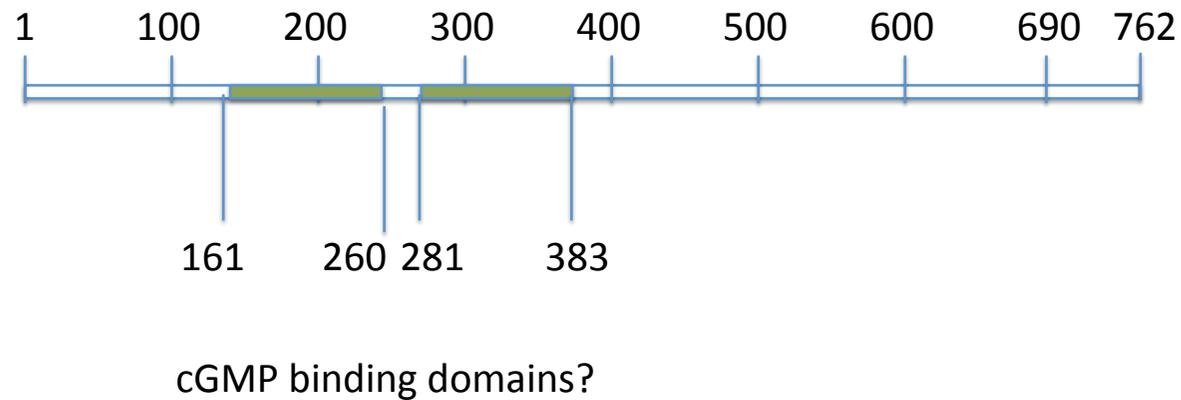
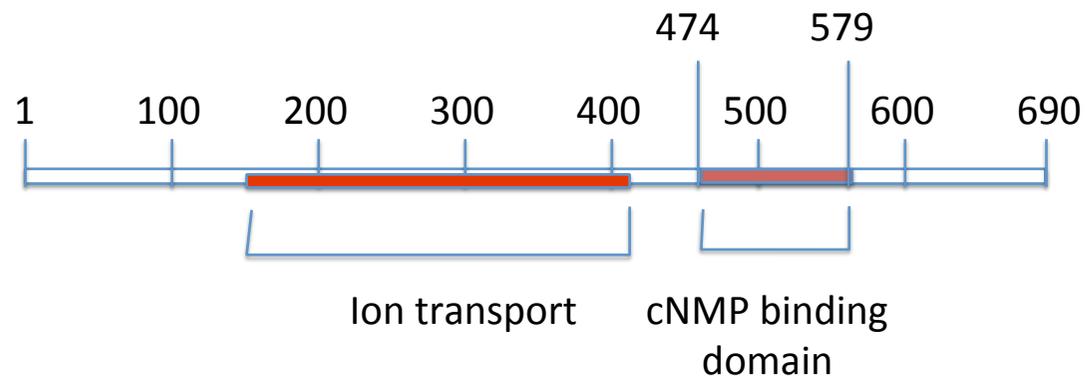
Marco Punta



cGMP-gated cation channel alpha-1

P29973 (CNGA1_HUMAN)

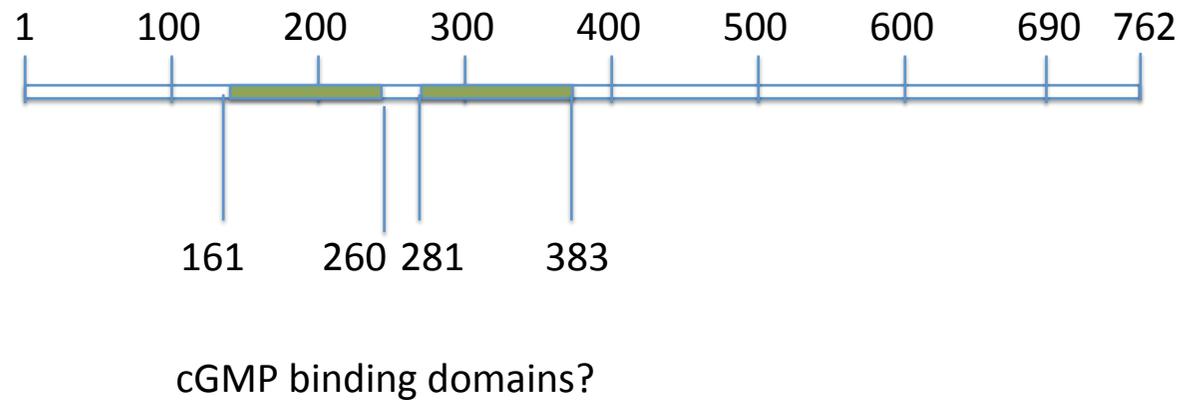
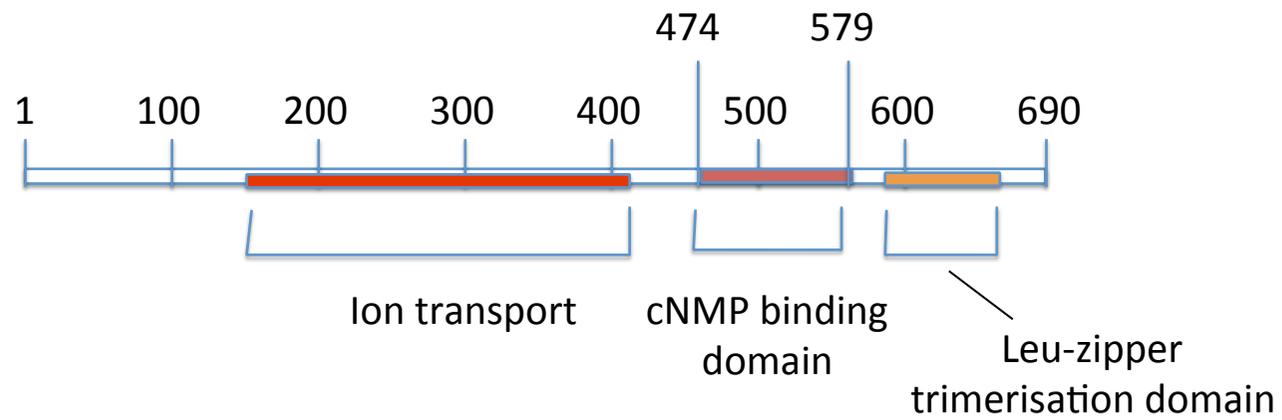
Marco Punta



cGMP-gated cation channel alpha-1

P29973 (CNGA1_HUMAN)

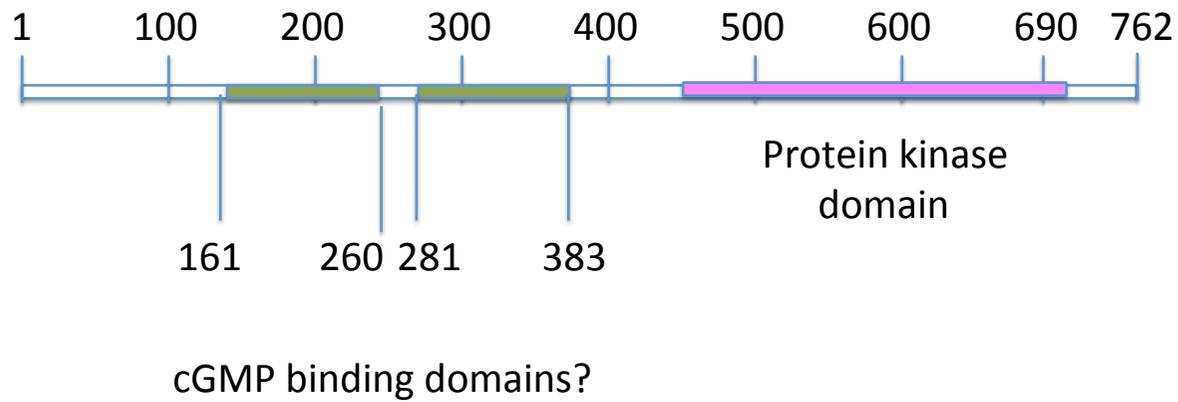
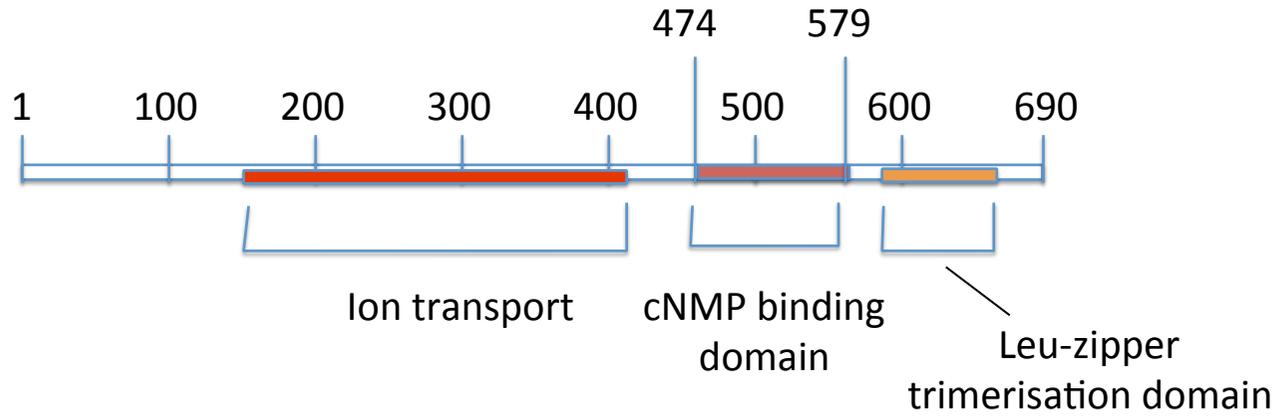
Marco Punta



cGMP-gated cation channel alpha-1

P29973 (CNGA1_HUMAN)

Marco Punta



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.

Function annotation transfer by homology

Homologous proteins can share a number of functional features, however:

- functional drift can lead to radically different functions
- while functional similarity correlates with function, no similarity threshold is safe for transfer
- Proteins may have multiple functions
- if more than one functional domain is present annotation transfer can be attempted only between domains that are homologous and NOT for the full-length protein function (but still not 'safe')

Globular proteins?

Water-soluble, globe-like shape.

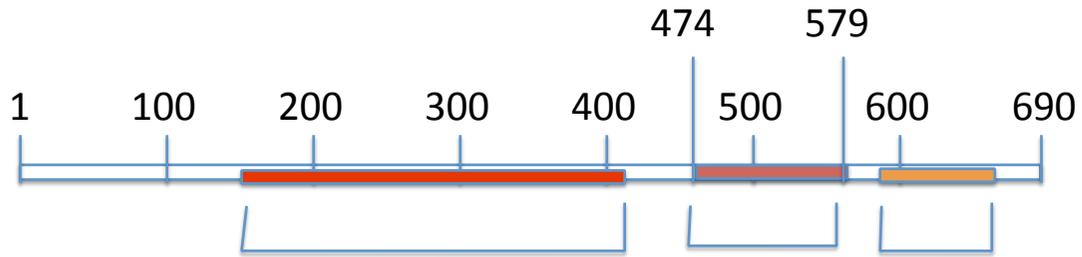
Non globular proteins?

Membrane proteins, fibrous proteins, disordered proteins

cGMP-gated cation channel alpha-1

P29973 (CNGA1_HUMAN)

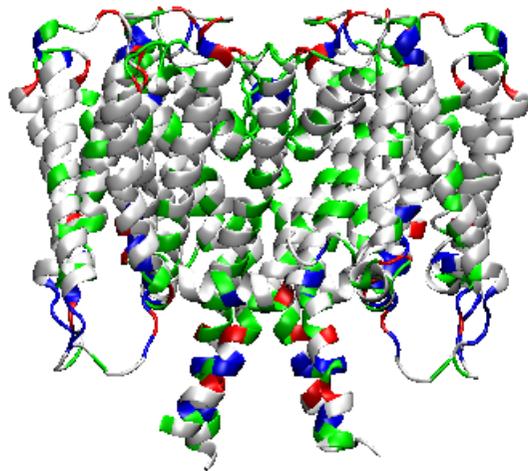
Marco Punta



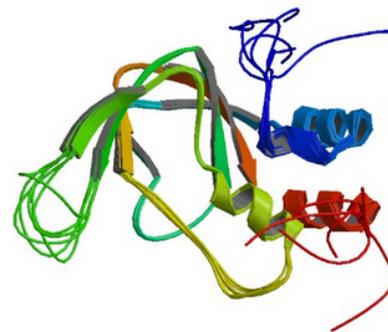
Ion transport

cNMP binding domain

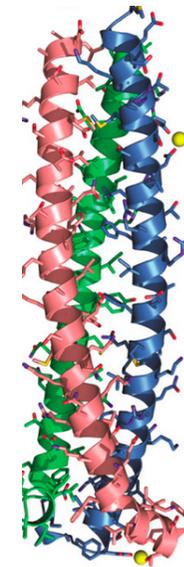
Leu-zipper trimerisation domain



Transmembrane



Globular

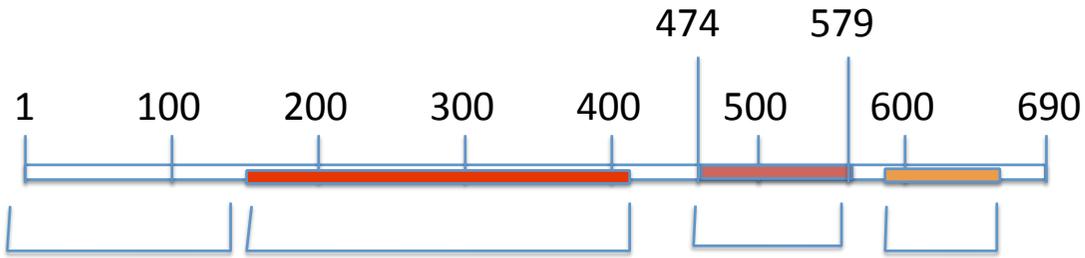


Coiled-coil

cGMP-gated cation channel alpha-1

P29973 (CNGA1_HUMAN)

Marco Punta

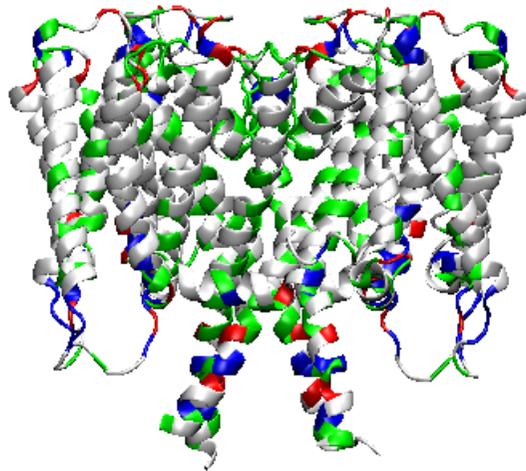


Predicted disordered

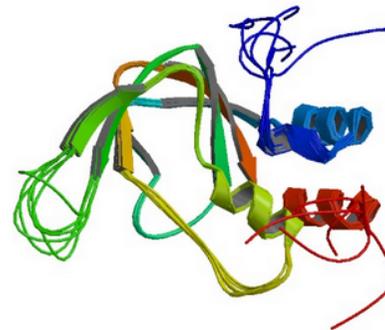
Ion transport

cNMP binding domain

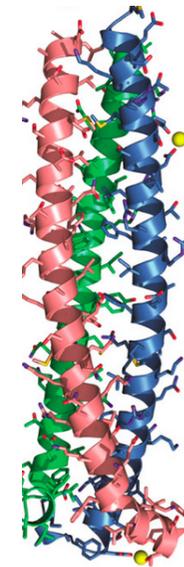
Leu-zipper trimerisation domain



Transmembrane



Globular

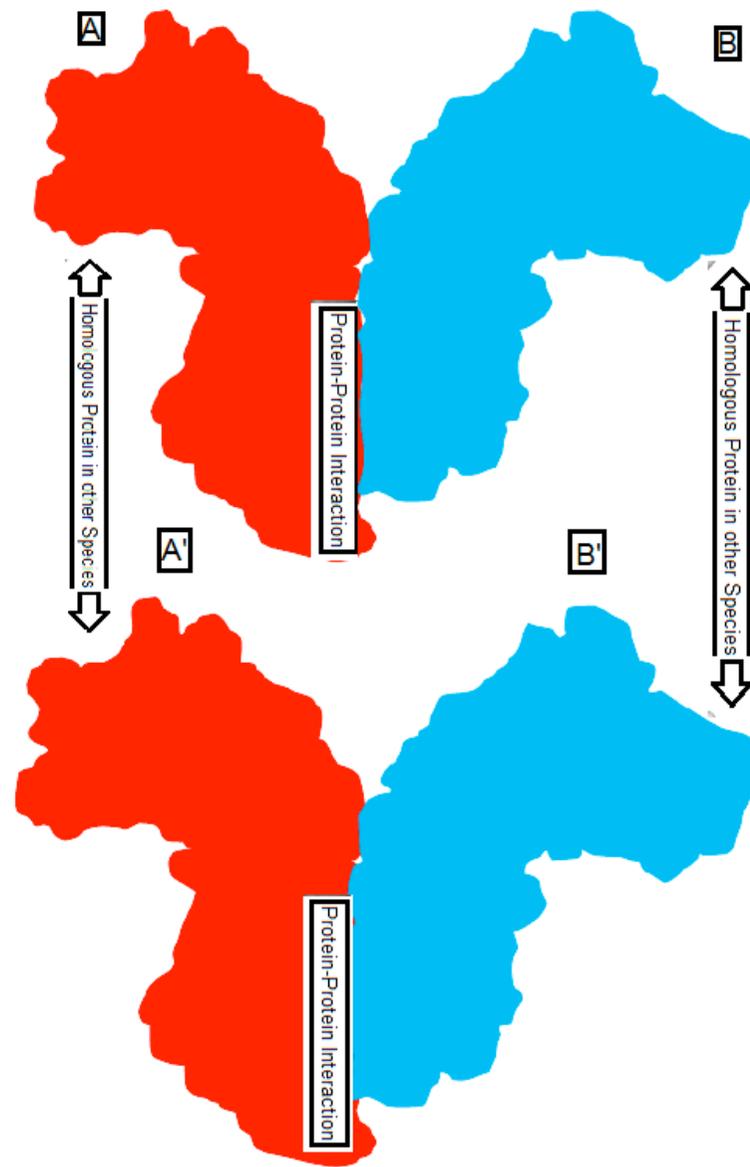


Coiled-coil

Homology \Leftrightarrow similar sequence?

Homology \Leftrightarrow similar structure?

Homology \Leftrightarrow similar function?



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)

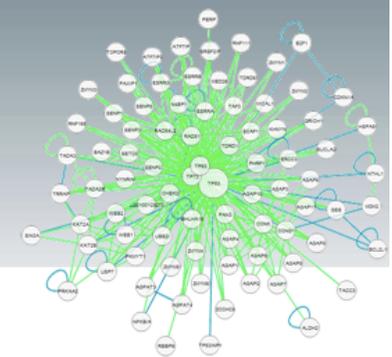
Figures

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.

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.

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
MMRFMLLFSRQGKLRLLQKWYLATSDKERKKMVRELMQVVLARKPKMCSFLEWRDLKVVYK
RYASLYFCCAIEGQDNELITLELIHRYVELLDKYFGSVCELDIIFNFEKAYFILDEFMLG
GDVQDTSKKSVLKAIEQADLLQEEDES PRSVLEEMGLA
```

Interacting partner 2:

```
>sp | P22892 | AP1G1_MOUSE
MPAPIRLRELIRTIRTARTQAEEREMIQECAAIRSSFREEDNTYRCRNVAKLLYMHMLG
YPAHFGQLECLKLIASQKFTDKRIGYLGAMLLDERQDVHLLMTNCIKNDLNHSTQFVQG
LALCTLGCMGSSEMCRDLAGEVEKLLKTSNSYLRKKAALCAVHVIRKVPPELMMFLPATK
NLLNEKNHGVLHTSVVLLTEMCERSPDMLAHFRKLVLPQLVRLKKNLIMSGYSPEHDVSGI
```

Input UniProt ID (Ex: AP1S1_MOUSE)

Interacting partner 1:

Interacting partner 2:

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)

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

InterEvol PyMol Visualization Tool

Structure Loader
Load PDB
Load the PDB structure of a complex (binary or multi-chain complex)

Alignment Loader
A(128aa) B(171aa)
1st sequence of every alignment should be the same as the sequence of every pdb chain.
Alignments can be generated at: <http://biodev.cea.fr/interevol/interevalign.aspx>

Alignment Visualization
Show Alignment (select)
To display the alignment:
- Directly select residues in the alignment
- OR write the selection in the command line

Residues Focus
R31A N35A E41A F31B E35B A47B
Residues are written as AminoAcid,Index,Chain.
Click on one or several buttons to zoom in the structure panel

Multiple sequence Alignment

gi/Index	R31A	N35A	E41A	F31B	E35B	A47B
Anopheles gambiae	N	N	N	N	N	N
Aedes aegypti	N	N	N	N	N	N
Drosophila melanogaster	N	N	N	N	N	N
Ixodes scapularis	N	N	N	N	N	N
Acyrtosiphon pisum	N	N	N	N	N	N
Schistosoma mansoni	N	N	N	N	N	N
Trichoplax adhaerens	N	N	N	N	N	N
Loa loa	N	N	N	N	N	N
Dictyostelium discoideum	N	N	N	N	N	N
Malassezia globosa	N	N	N	N	N	N
Ustilago maydis	N	N	N	N	N	N
Laccaria bicolor	N	N	N	N	N	N
Schizophyllum commune	N	N	N	N	N	N
Schizosaccharomyces pombe	N	N	N	N	N	N
Schizosaccharomyces japonicus	N	N	N	N	N	N
Tuber melanosporum	N	N	N	N	N	N
Magnaporthe oryzae	N	N	N	N	N	N
Uncinocarpus reesii	N	N	N	N	N	N
Coccidioides posadasii	N	N	N	N	N	N
Trichophyton verrucosum	N	N	N	N	N	N
Aspergillus nidulans	N	N	N	N	N	N
Aspergillus fumigatus	N	N	N	N	N	N
Talaromyces stipitatus	N	N	N	N	N	N
Verticillium albo-atrum	N	N	N	N	N	N
Gibberella zeae	N	N	N	N	N	N
Hectactia haematococca	N	N	N	N	N	N
Botryotinia fuckeliana	N	N	N	N	N	N
Podospora anserina	N	N	N	N	N	N
Neurospora crassa	N	N	N	N	N	N

Once selected, aligned positions are displayed

Basic swap between positions 31 and 35

Invariant position at E35

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