

6.

Protein classifications

Common Motifs

(a) Helix-loop-helix



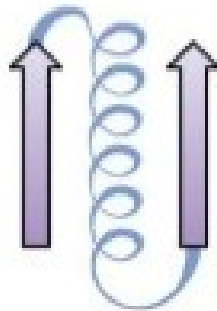
(b) Coiled coil



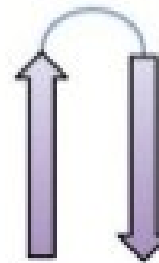
(c) Helix bundle



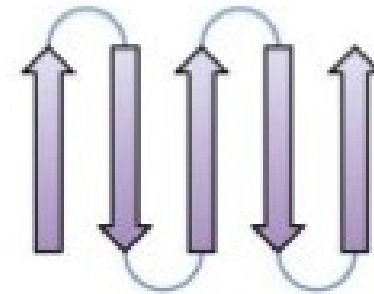
(d) $\beta\alpha\beta$ unit



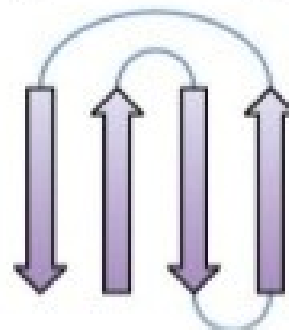
(e) Hairpin



(f) β meander



(g) Greek key



(h) β -sandwich



Motifs Combine to form Domains

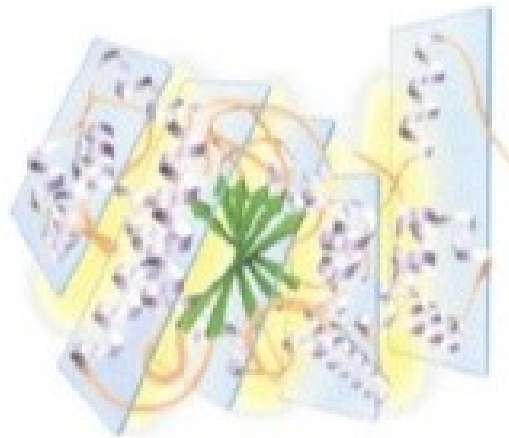


(a) Cytochrome *c*'



Parallel twisted sheet

(b) Phosphoglycerate kinase
(Domain 2)



(c) Phosphatase
(Domain 2)



Alpha/beta barrel

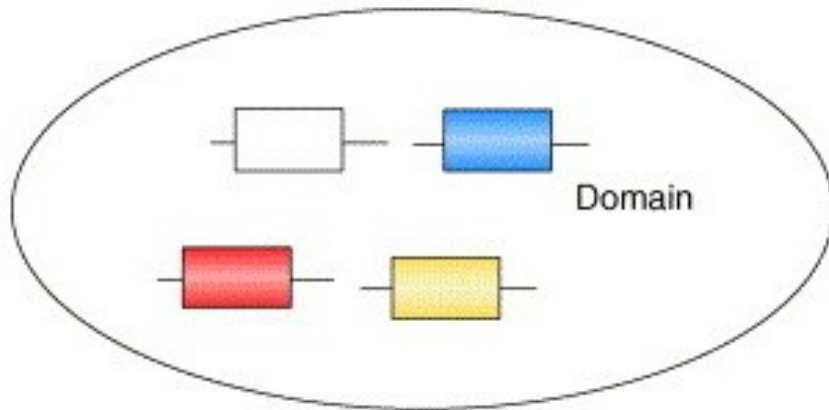
(d) Triose phosphate isomerase

- Domains are independent folding units in a 3^o structure of a protein
- Individual domains have specific function

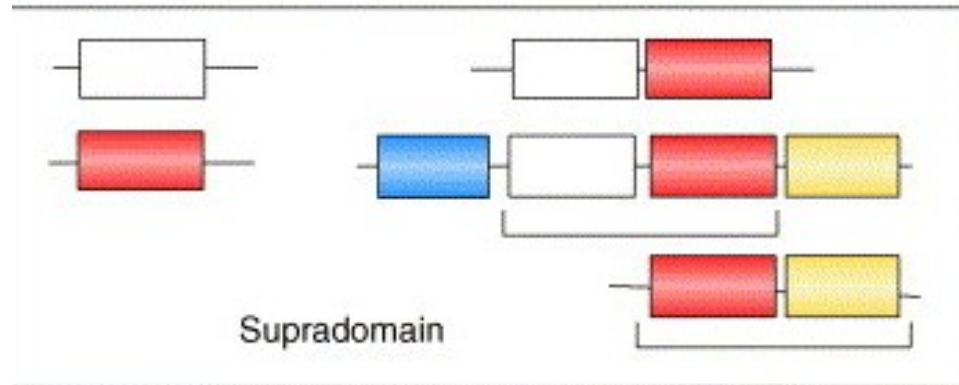
- Hydrophobic interactions are the major driving force in folding domains

Domains are reused and reinvented

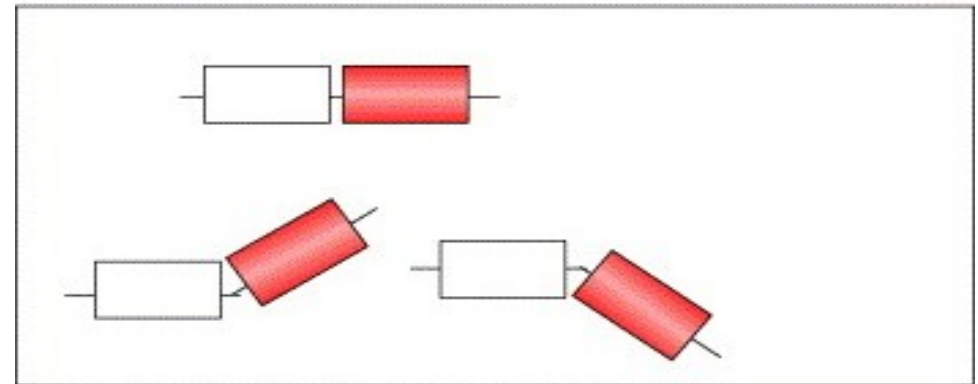
The repertoire of domain superfamilies...



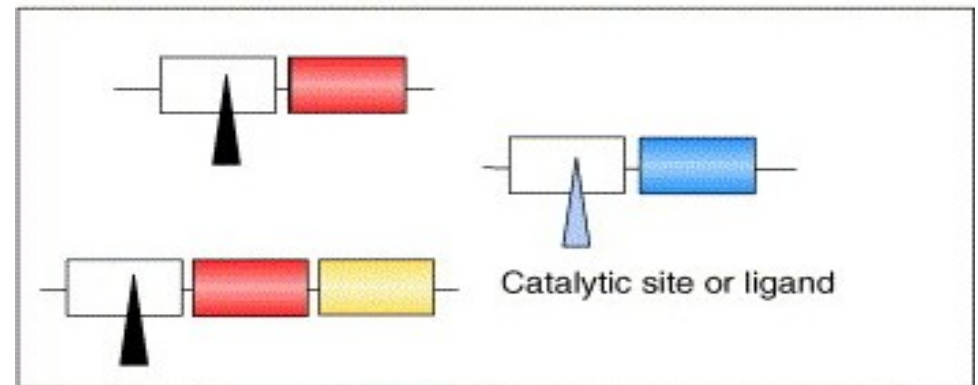
...duplicates and recombines to form single and multi-domain proteins.



The same combination can adopt different geometries...



...and/or different functions.

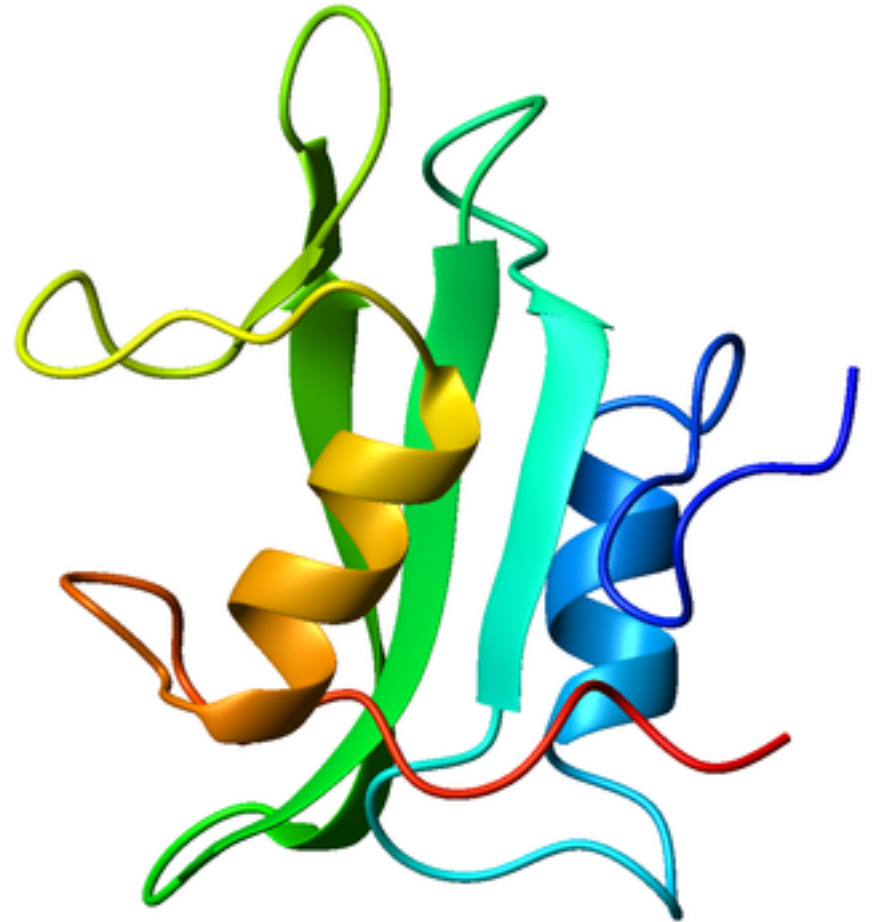


SH2 domain

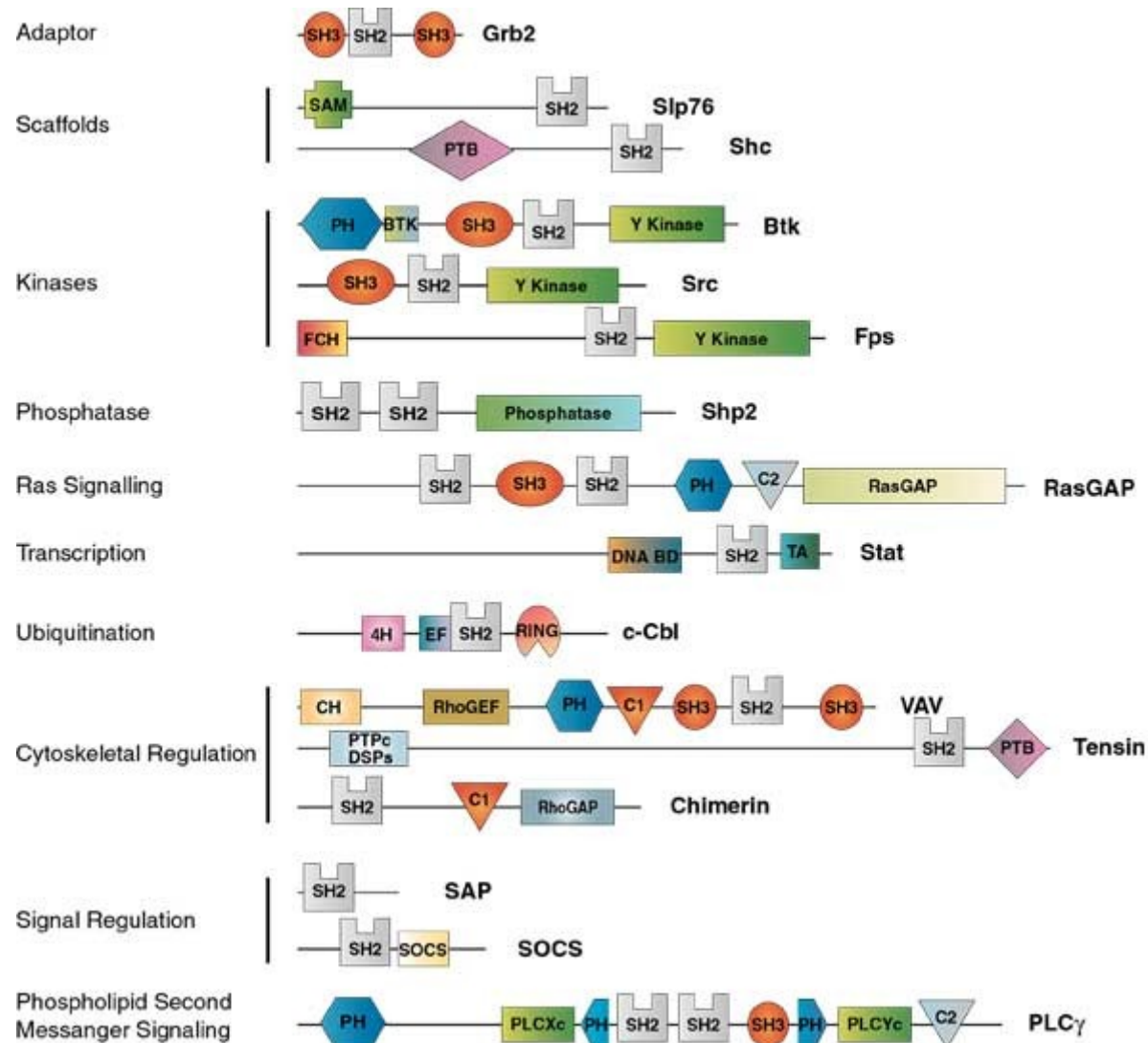
Involved in PPIs

Involved in signal
transduction

Binds phosphorylated
Tyr residues



SH2 domain



Protein classification

- Structural domains
- Protein families
- Sequence signatures

CATH

manually-curated hierarchical classification of protein domain structures.

The screenshot shows the top navigation bar of the CATH website with links for Home, Search, Browse, Download, About, and Support. A search bar on the right contains the text '2DHHA'. Below the navigation bar is a large grey box with the title 'CATH / Gene3D' and the text '26 million protein domains classified into 2,738 superfamilies'. At the bottom of this box are four buttons: 'Browse »' (red), 'Search »' (green), 'Download »' (blue), and 'Take the Tour »' (orange).

What is CATH?

CATH is a classification of protein structures downloaded from the Protein Data Bank. We group protein domains into superfamilies when there is sufficient evidence they have diverged from a common ancestor.

- [Search CATH by text, ID or keyword](#)
- [Browse CATH Hierarchy](#)
- [CATH Release Notes](#)

Latest Release Statistics

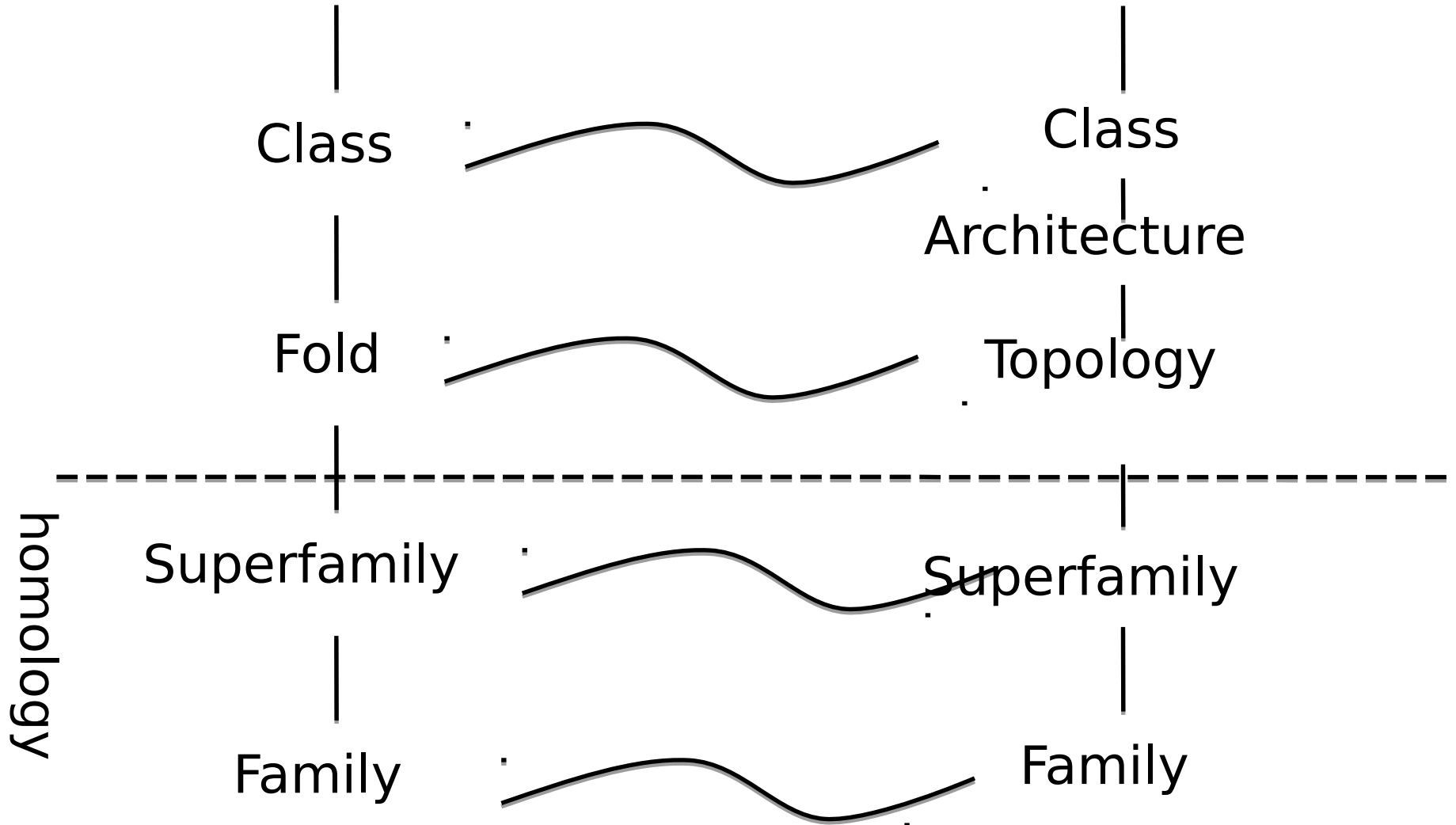
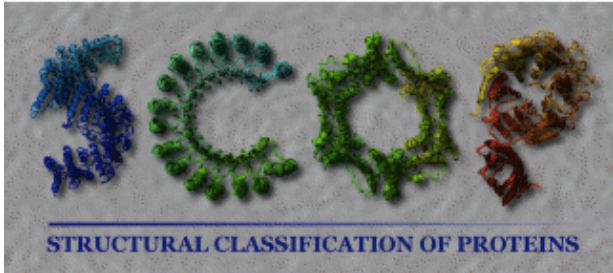
CATH v4.0 based on PDB dated March 26, 2013	
235,858	CATH Domains
2,738	CATH Superfamilies

Class: secondary structure content (e.g. mainly-alpha, mainly-beta, mixed alpha/beta or 'few secondary structures');

Architecture: general arrangement of the secondary structures

Topology (fold) takes into account the connectivity of secondary structures in the chain;

Homologous Superfamily: domains that are believed to be related by a common ancestor.



Protein sequence families

Protein families


- Defined based on MSA
- Identification of functional amino acids
- Diversity -> detecting remote homologues
- Identification of parts of the sequence space which have no functional annotation

Pfam database

Pfam: Family: Avidin (PF01382) - Mozilla Firefox

Pfam: Family: Avidin (P...

pfam.xfam.org/family/PF01382

EMBL-EBI  HOME | SEARCH | BROWSE | FTP | HELP | ABOUT

Pfam
keyword search **Go**

Family: Avidin (PF01382)

18 architectures 142 sequences 1 interaction 81 species 454 structures

- Summary
- Domain organisation
- Clan
- Alignments
- HMM logo
- Trees
- Curation & model
- Species
- Interactions
- Structures

Jump to...
enter ID/acc **Go**

Summary: Avidin family

Pfam includes annotations and additional family information from a range of different sources. These sources can be accessed via the tabs below.

Wikipedia: Avidin **Pfam** InterPro

This tab holds the annotation information that is stored in the Pfam database. As we move to using Wikipedia as our main source of annotation, the contents of this tab will be gradually replaced by the Wikipedia tab.

Avidin family **Provide feedback**

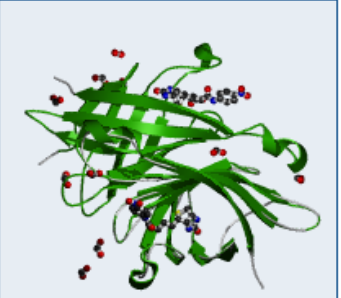
No Pfam abstract.

Internal database links

SCOOP: [DUF1143](#)

External database links

HOMSTRAD:	avi
PANDIT:	PF01382
PROSITE:	PDOC00499
Pseudofam:	PF01382
SCOP:	1s1f
SYSTEMS:	Avidin



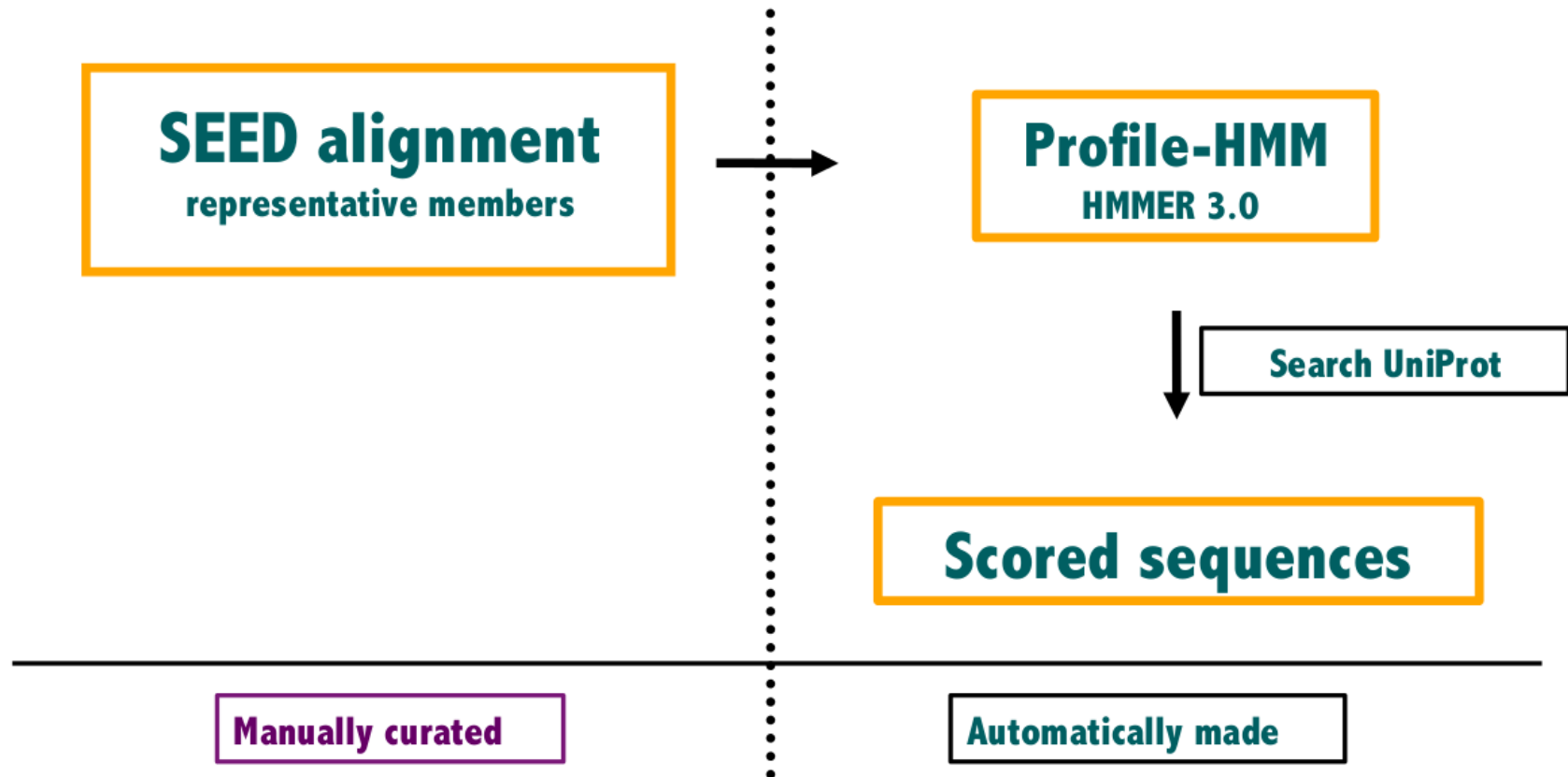
Example structure
PDB entry 2OF8: Crystal structure of AVR4 (D39A/C122S)-BNA complex
View a different structure:

Pfam sequence families

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

```
A6VZD9_MARMS/3-149 NY...WLMK.....SEPKDA.....FSIDDL...KRLKH...SPWDGVRNYQARN..FMK.EMNEGDLVFFYHSS.
THYN1_HUMAN/55-219 SH...WLMK.....SEPESRLEKGVVDVKFSIEDLK.AQPKQT..TCWDGVRNYQARN..FLR.AMKLGEEAFFYHSN.
A3ZLP5_9PLAN/16-165 RY...WLLK.....TEPES.....YSIDDLA.NEKKQT..TFWSGVRNYQARN..FMRDDMKVGDEVFFYHSN.
Q312I8_DESDG/2-150 RY...WLMK.....SEPEC.....FSLEDLV.NAPEQT..TPWDGVRNYQARN..FMRDEMPPGDKVLFYHSG.
Q1MR11_LAWIP/2-150 QY...WLFK.....SDTDC.....YSINDLQ.SAPNQT..TSWDGVRNYQARN..FMRDEMIRIGDLGFFYHSG.
Q729D5_DESVH/2-149 NY...WLFK.....SETDC.....FSVDDLA.ASPDAT..SSWDGVRNYQARN..FMR.TMRKGDVGFFYHSG.
Q3A4Z6_PELCD/6-155 RY...WLMK.....SEPGC.....FSIDDLK.DCPDGI..SPWDGVRNYQARN..LLRDEIKAGDGVLFYHSN.
Q74AS6_GEOSL/6-154 RY...WLFK.....SEPSC.....FSFDDLQ.SRPNGT..EHWDGVRNFQARN..LLRDEIKPGDGVLFYHSN.
B3E7D4_GEOLS/2-150 RY...WLFK.....SEPGC.....FSFQDLQ.ARPNAT..EQWDGVRNFQARN..FLRDEIKPGDRVLFYHSS.
A1ATY5_PELPD/2-150 NY...WLFK.....TEPGC.....FSFDNLK.NRPNMT..EPWDGVRNFQARN..YLRDTVKVGDLVLFYHSN.
Q054Y1_LEPBL/2-151 KH...WLFK.....TEPDV.....FSIDDLY.KAPSRI..APWEGVRNYQARN..FLRDSIQKGDLLIFYHSR.
Q5KKF6_CRYNE/7-162 EF...WQTA.....NG.....KFCRV..SPWDGVRNHEAKK..IMREKMKLGDKVLFYHSN.
Q7UMR5_RHOBA/2-156 KY...WLMK.....TEPNT.....FSIDDLA.EQPEQI..TCWEGVRNYQARN..LLRDEIEEGDQVLFYHSA.
Q60BW0_METCA/2-151 RY...WLMK.....TEPGE.....FGIDDLA.ARPAQT..EPWDGVRNYQARN..MMRDEMVKVGDVLFYHSN.
A4BS86_9GAMM/2-151 SY...WVMK.....SEPSV.....YGIDDLA.AQPSQT..DHWEGVRNYQARN..MLRDQMRPGDLALLYHSN.
A1WW93_HALHL/2-151 NR...WVMK.....SEPDV.....FGIDDLA.AAPQGT..DRWDGVRNYQVRN..MIRDHMRPGDAFFYHSN.
Q1K2S0_DESAC/2-151 NY...WLMK.....SEPEA.....FGIDDLQ.QMPEQT..EHWDGVRNYQARN..MMRDDMKIGDLAFFYHSN.
Q0A5L0_ALHEH/2-151 NY...WLMK.....SEPDE.....FGIEDLK.QRPDQI..EPWDGVRNYQARN..MMRDQMKVGDVLAFFYHSN.
Q83BN8_COXBU/2-151 NY...WLLK.....SEPTS.....YSIDDLF.REKNKI..TRWDGVRNYQARN..FMRDGMKKGDVLAFFYHSN.
Q7S515_NEUCR/116-280 QY...WLLK.....AEPLPRLNGYDVHFSIDDL..AARTSP..EPWDGIRNYSARN..NLR.SMRVGDVLAFFYHSN.
A6SM74_BOTFB/112-267 QY...WLMK.....AEPESRIEKGHDIKFSIDDL..AAKTEP..EPWD.....ARN..NLR.AMKKGDVLAFFYHSS.
Q0CLU5_ASPTN/109-262 SY...WLMK.....AEPESRIEKGVDVKFSIDDL..RERTKP..EPWD.....ARN..NMR.EMKKGDYAFFYHSN.
Q0U623_PHANO/118-277 VF...WLLK.....AEPLFRYENGVNVAFSIDDL..AACTVP..EPWGGVRNPQARN..NMQ.AMRKGDVGFFYHSN.
```

PFAM generálása



A description of the family, includes thresholds you to create the full alignment

Rules – No false positives. A family is not allowed to overlap with any other family

Family overlaps

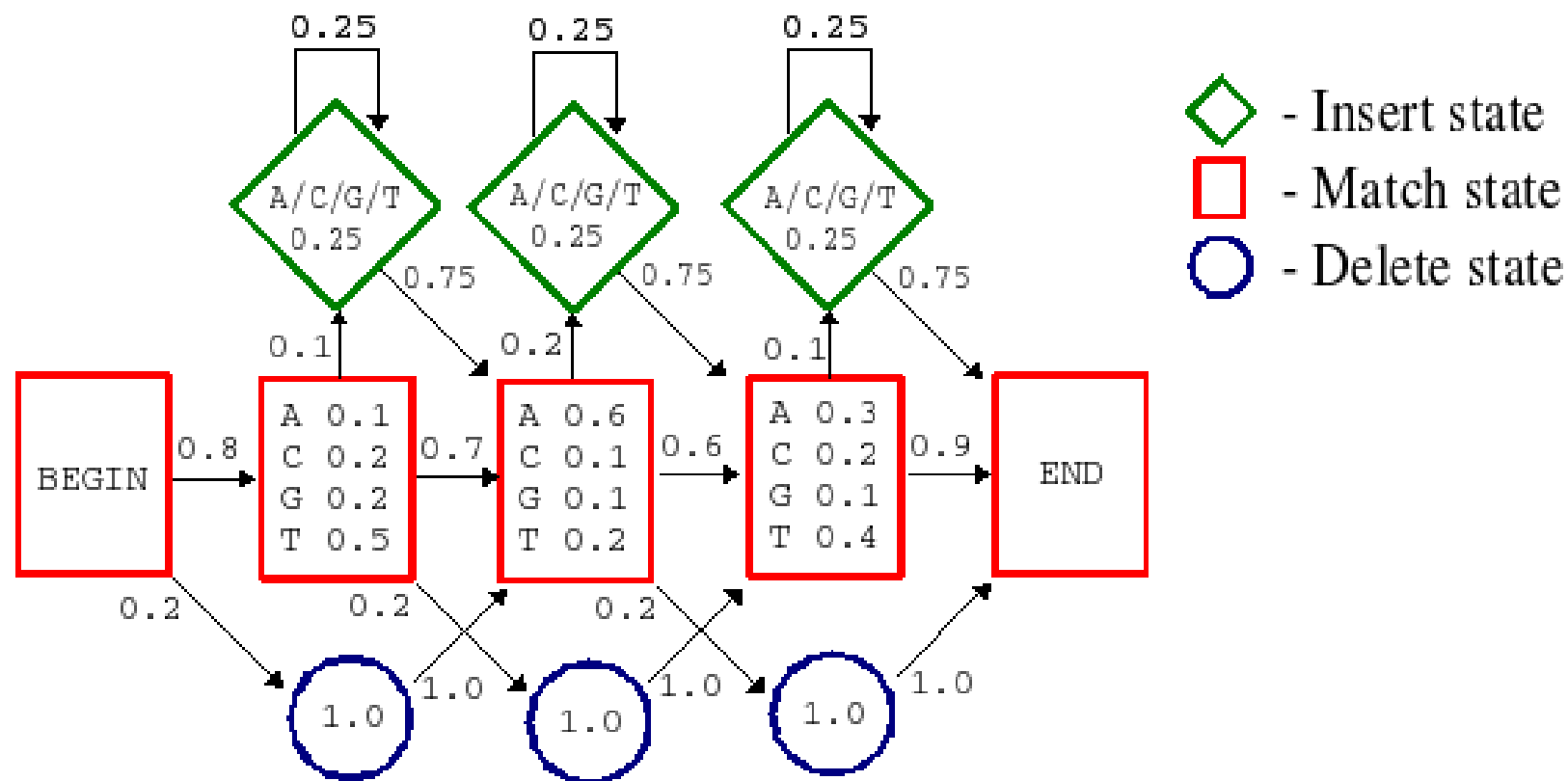
Old Family



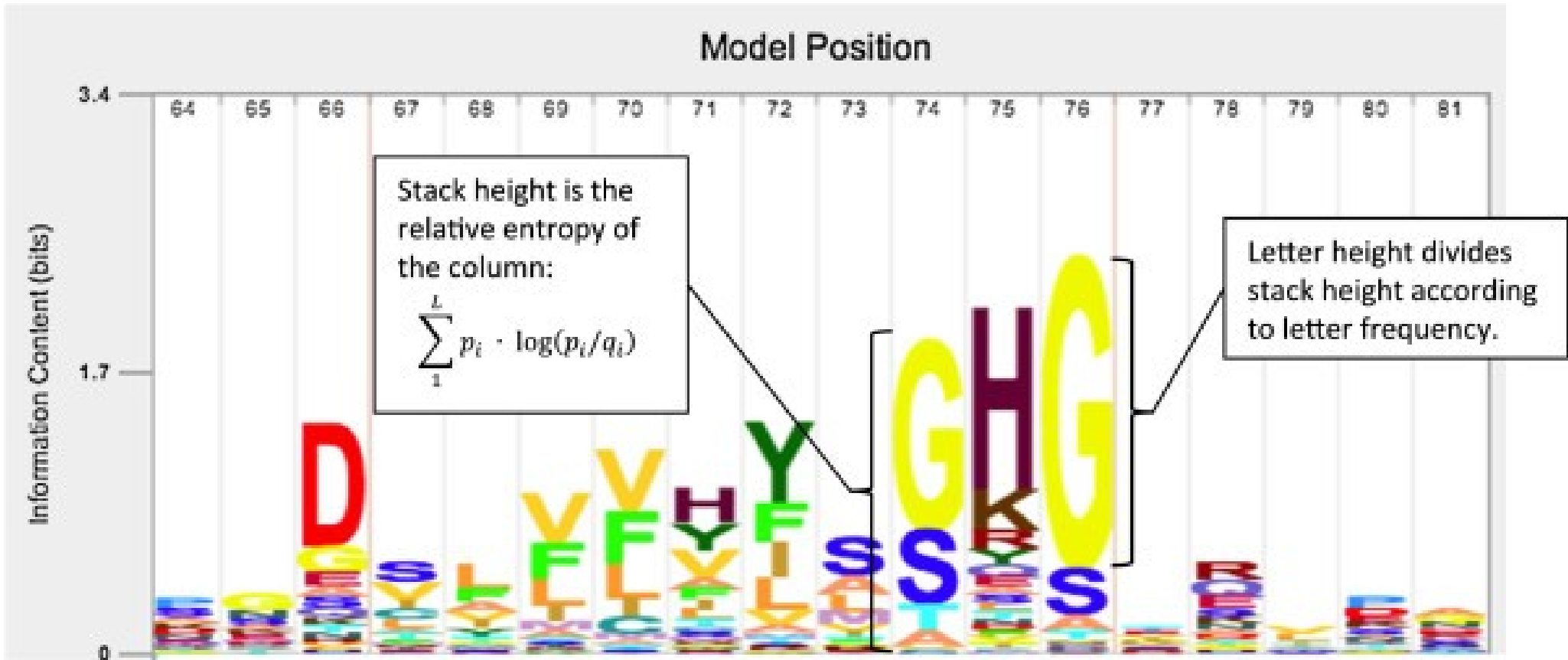
New Family



Profile Hidden Markov Models - Encapsulate diversity



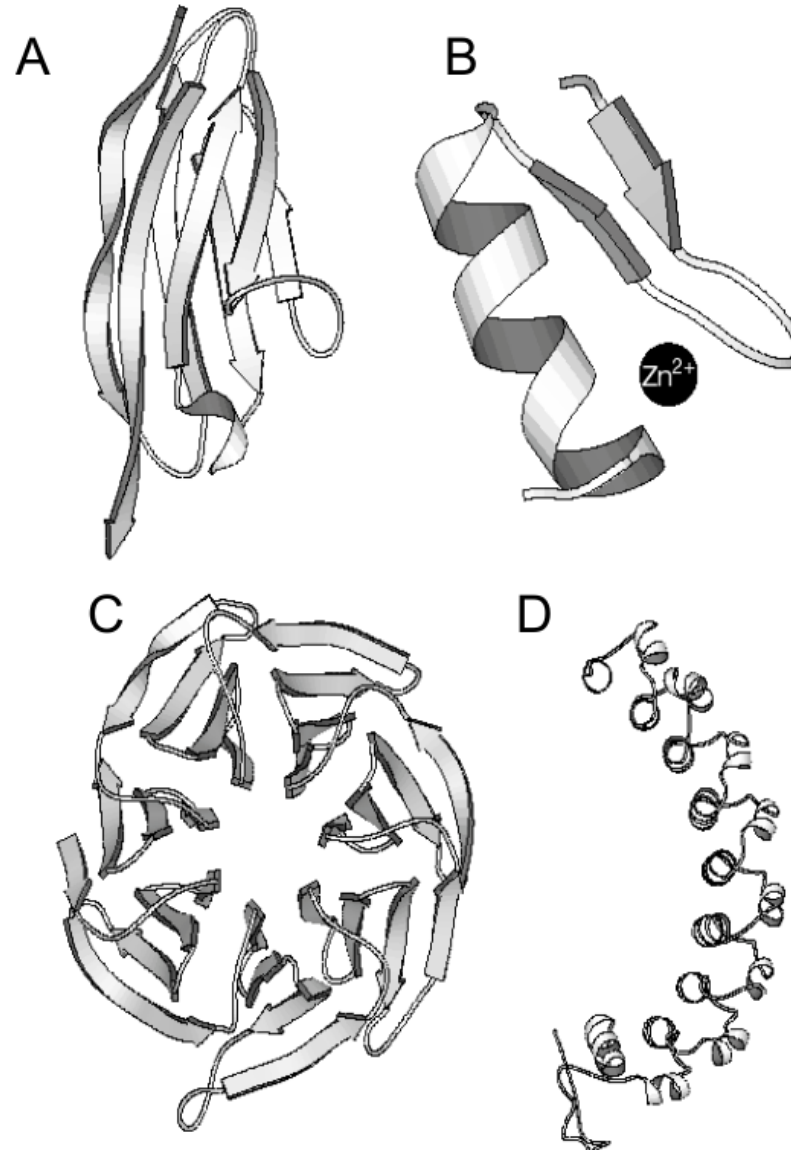
Logos



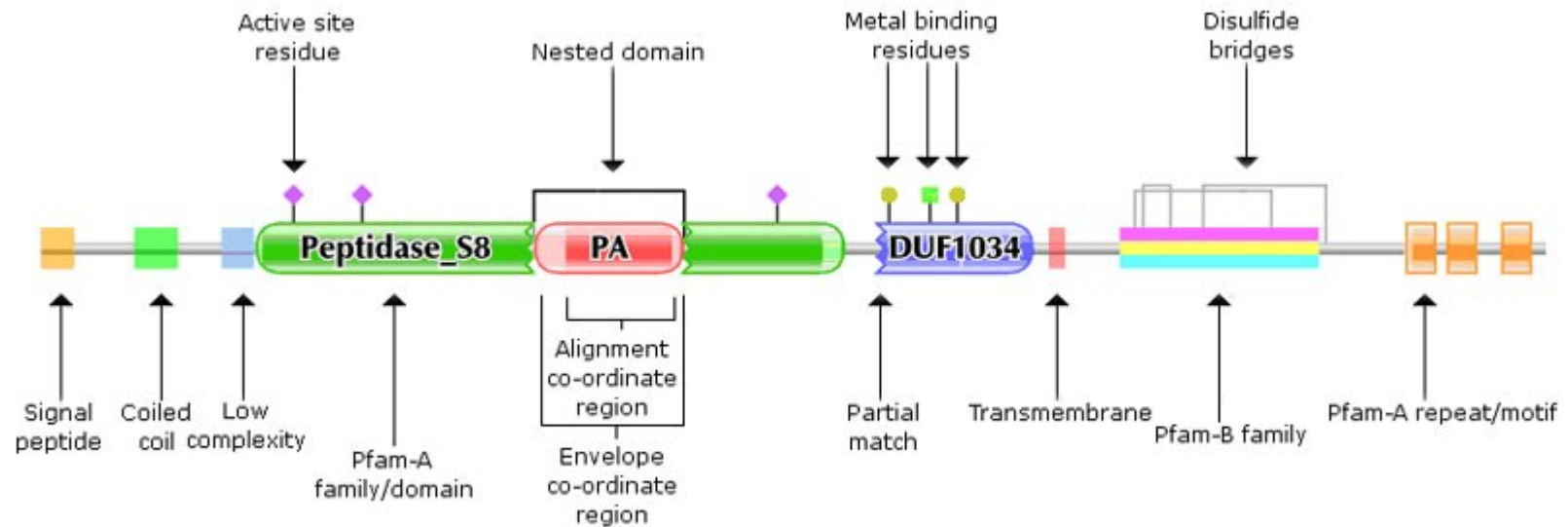
Types of PFAM families

• PFAM

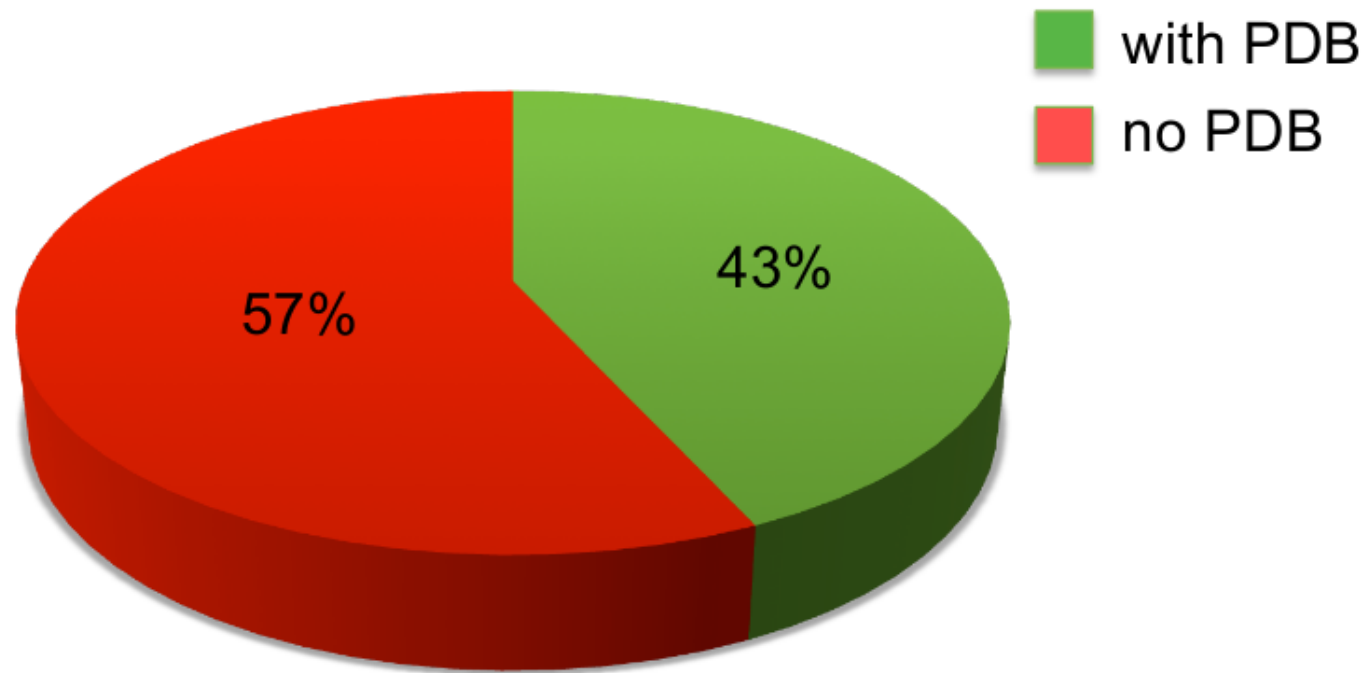
- Domains – structural information
- Families
- Repeats
- Motifs



Pfam pictograms

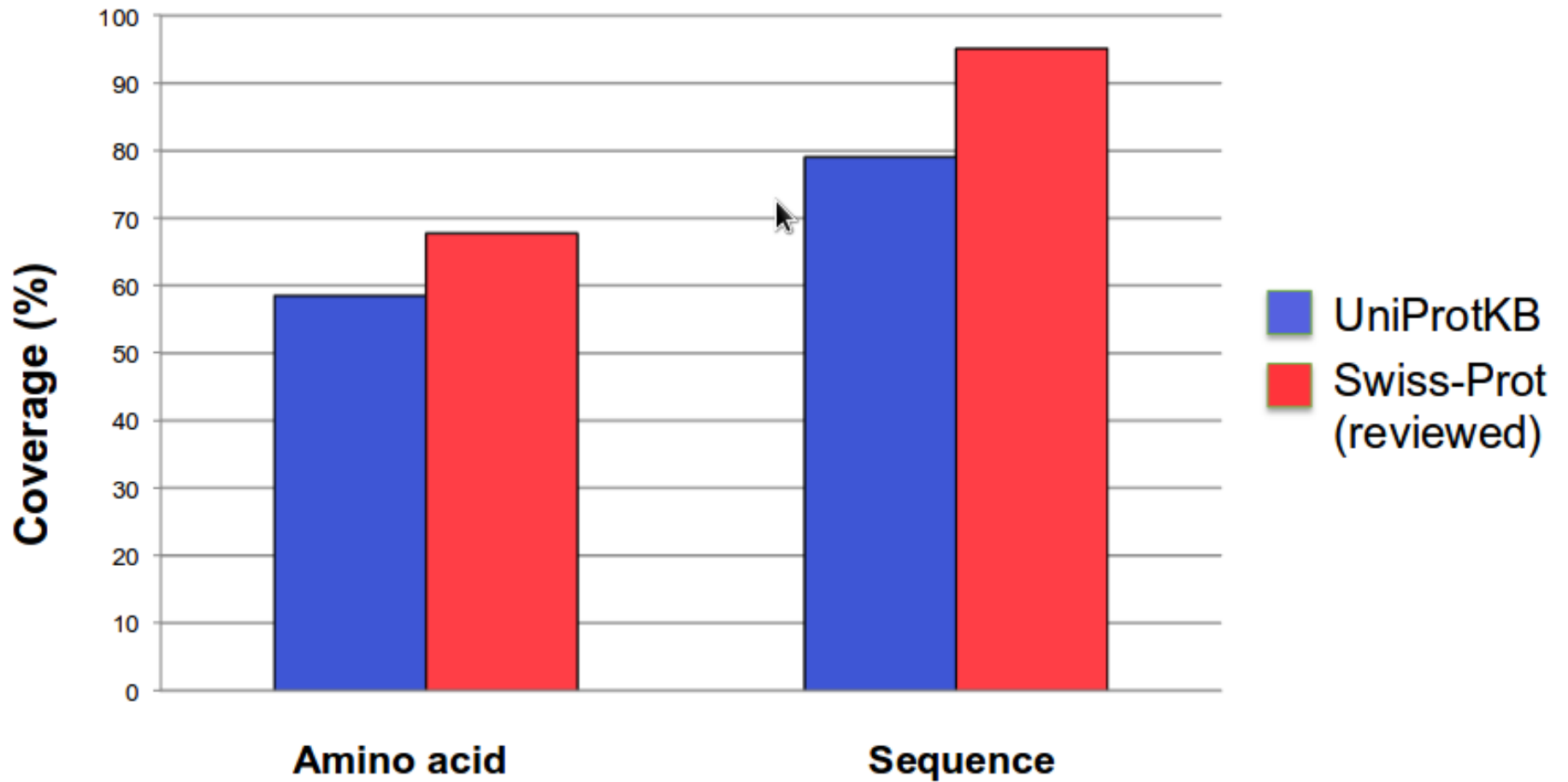


Pfam families with PDB structure

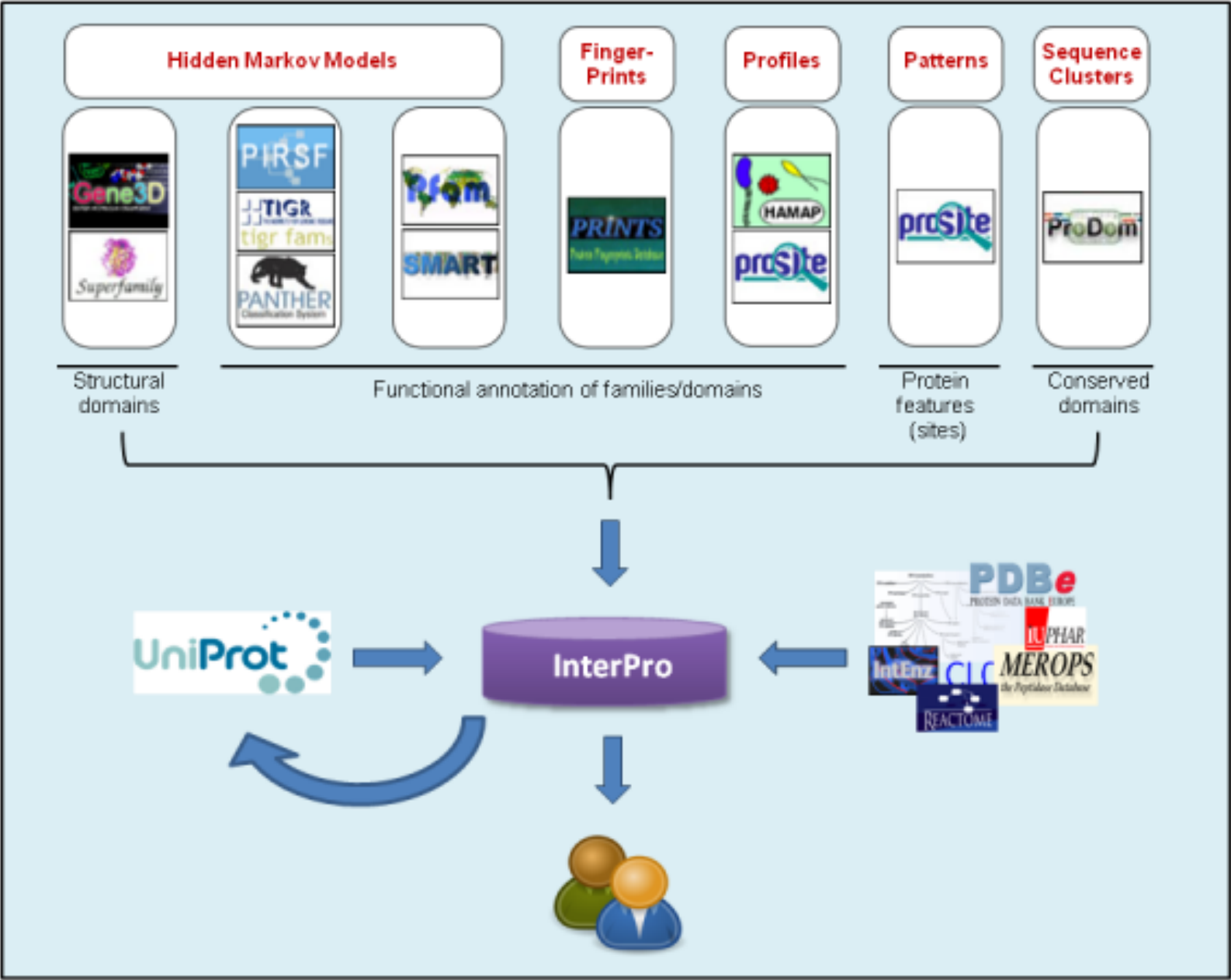


100%=all Pfam families

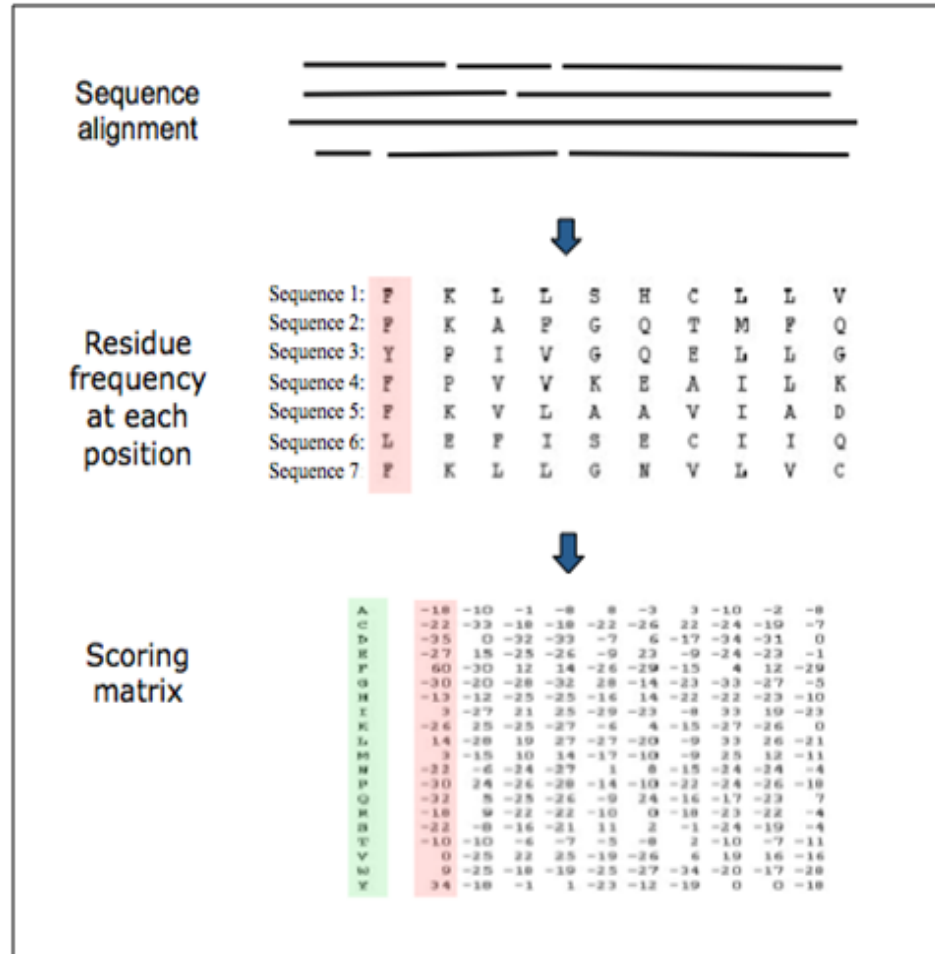
Pfam coverage



Interpro



Profiles

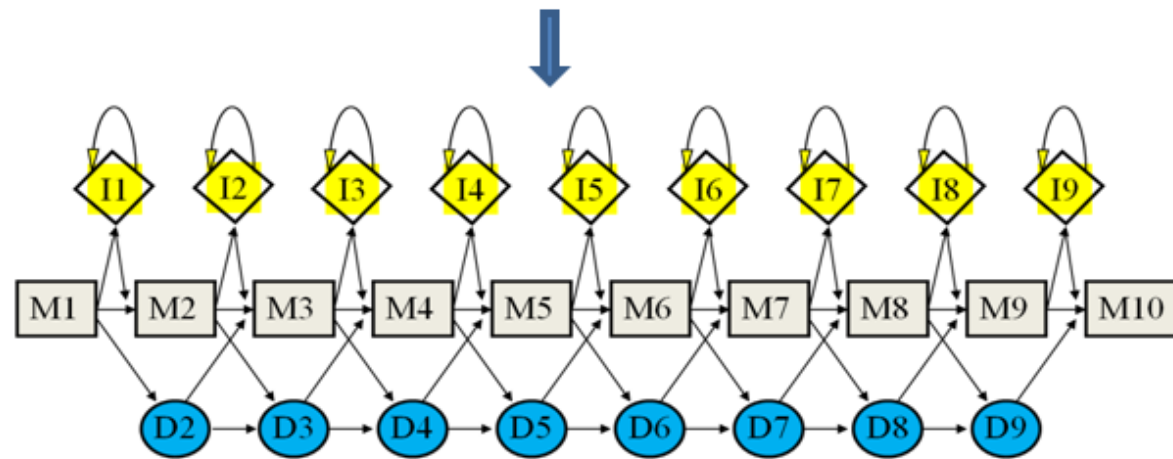


e.g. HAMAP, PROSITE, PRODOM

HMM

Multiple sequence alignment

Sequence 1:	F	K	L	L	S	H	C	L	L	V
Sequence 2:	F	K	A	F	G	Q	T	M	F	Q
Sequence 3:	Y	P	I	V	G	Q	E	L	L	G
Sequence 4:	F	P	V	V	K	E	A	I	L	K
Sequence 5:	F	K	V	L	A	A	V	I	A	D
Sequence 6:	L	E	F	I	S	E	C	I	I	Q
Sequence 7:	F	K	L	L	G	N	V	L	V	C



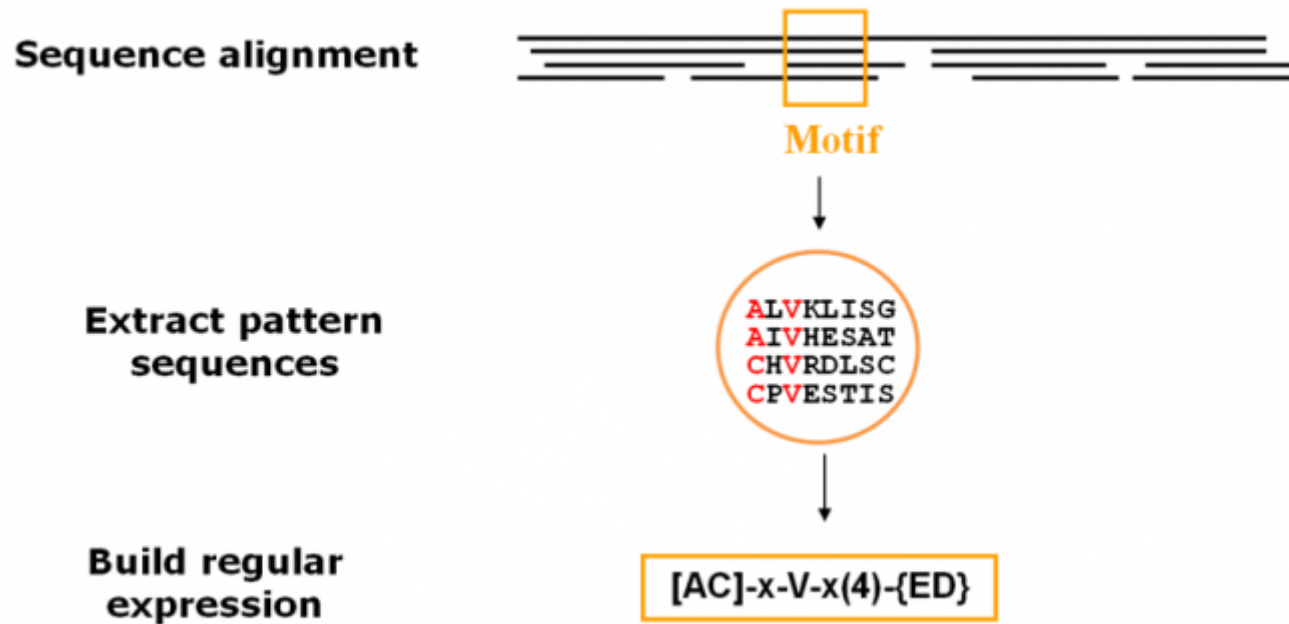
I = insert state

M = match state

D = delete state

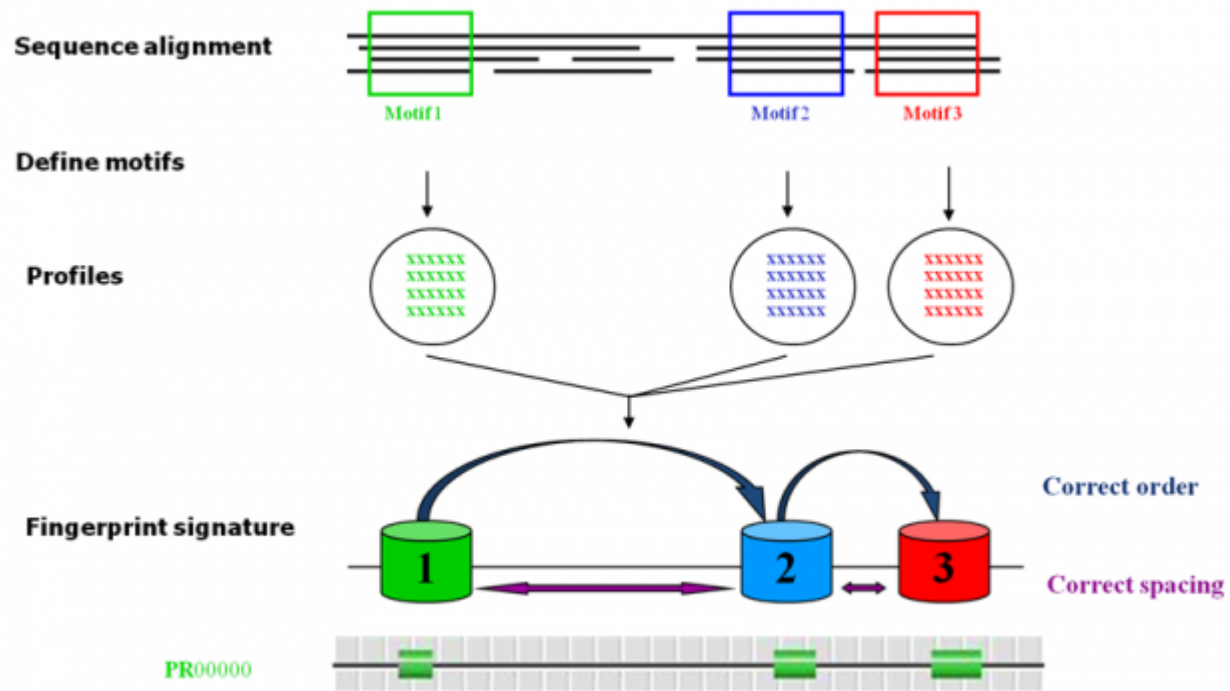
e.g. Pfam, SMART, TIGRFAM, PANTHER..

Patterns



e.g. PROSITE

Fingerprint



e.g. PRINTS

Domain architectures

