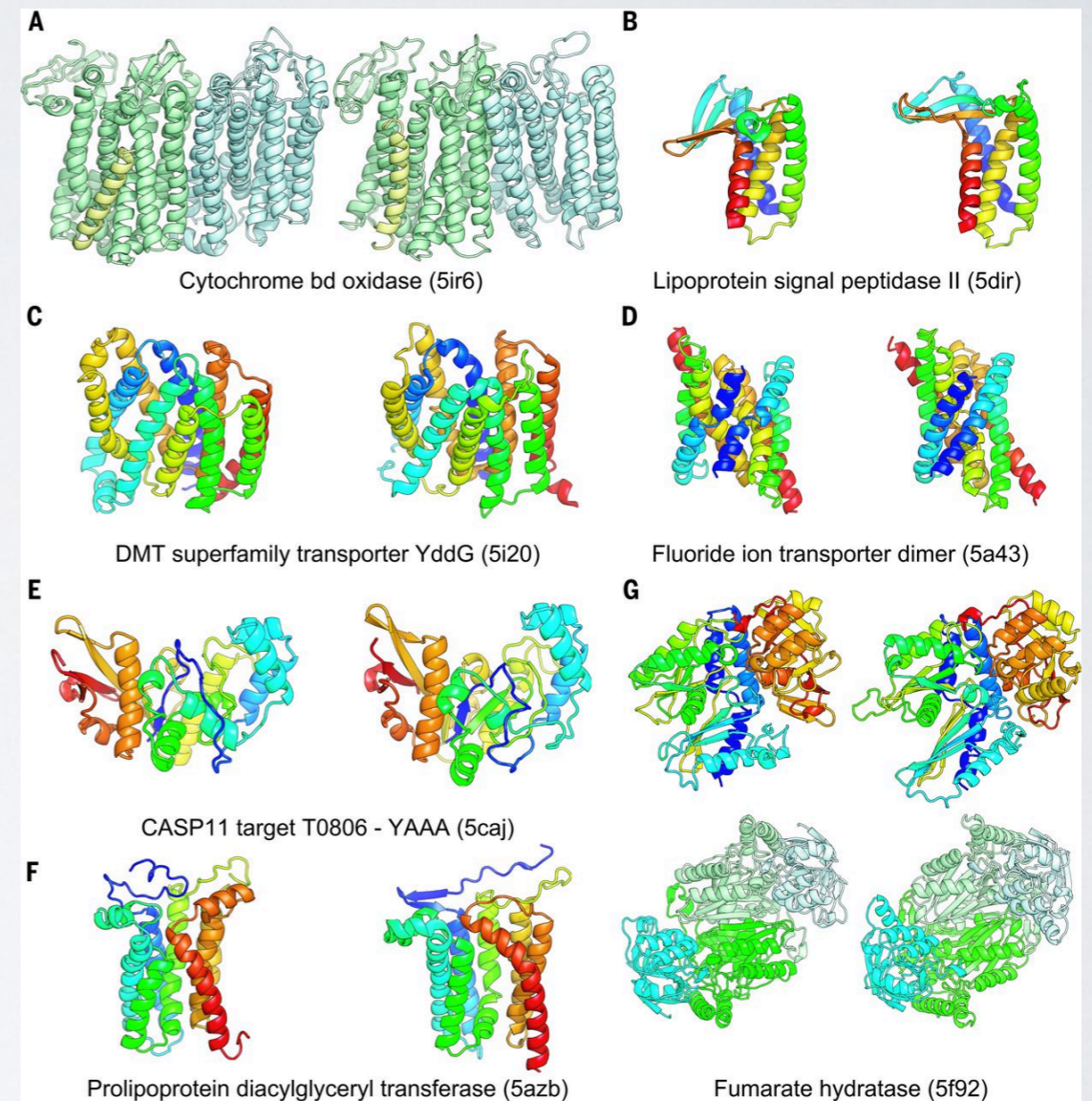
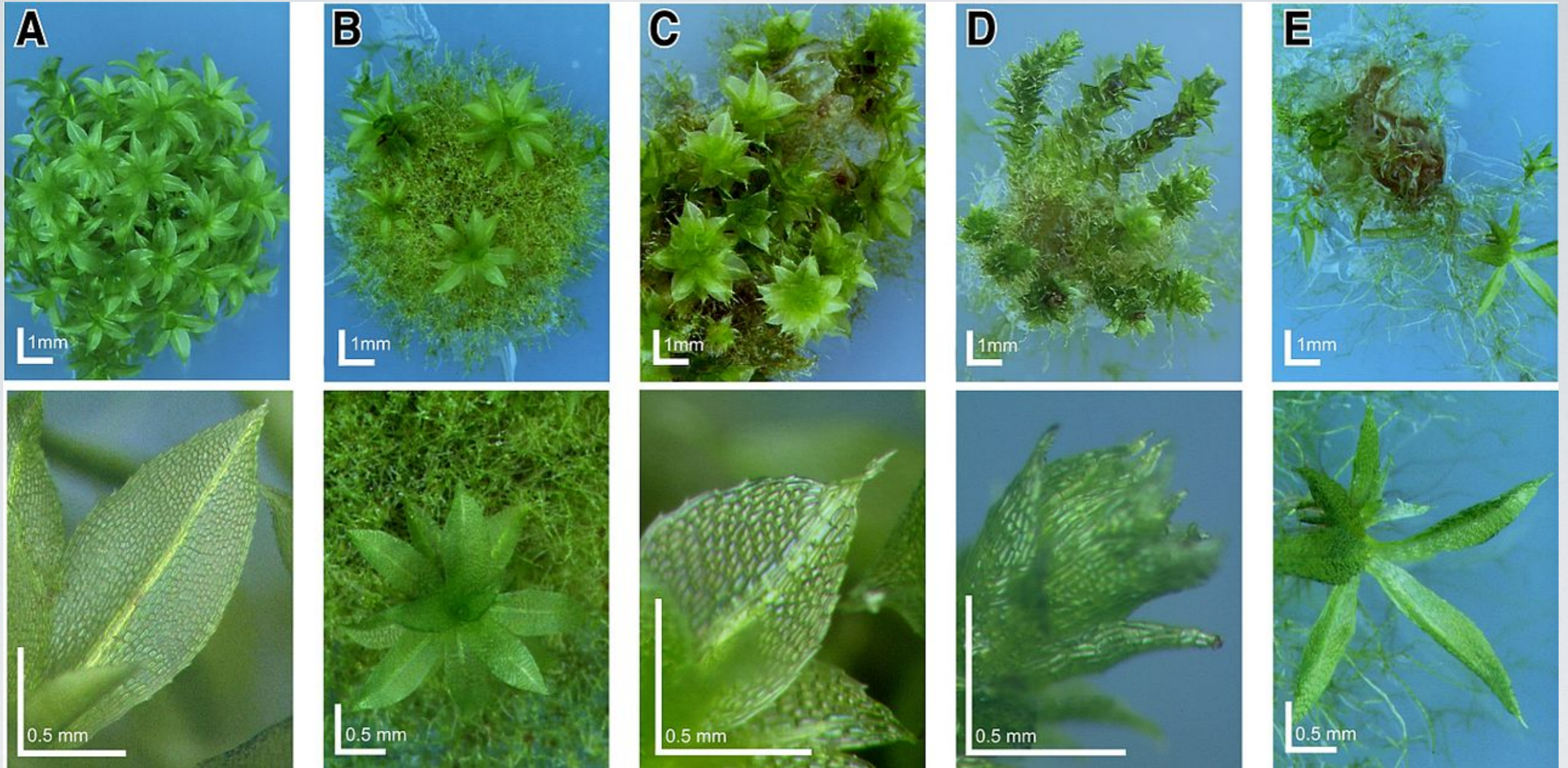


PROTEIN DOMAINS, HMMS & MOTIFS

CLASSIFYING PROTEINS BY FUNCTION

- Important to be able to classify proteins as to what functions they perform
- This information is taken from experimental studies
- Genes have function determined from mutant





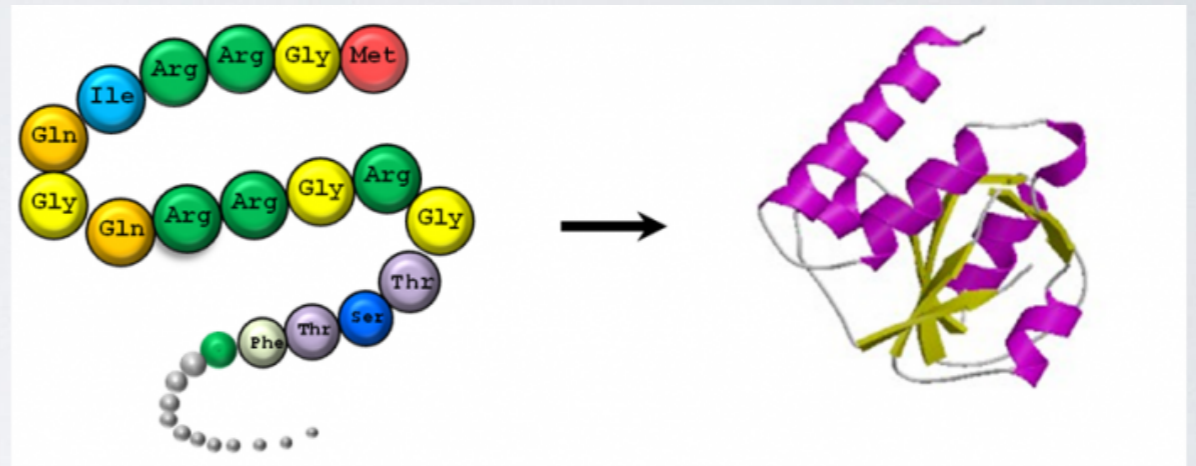
https://en.wikipedia.org/wiki/Reverse_genetics

GENETICS TO FUNCTION

By seeing which mutations break a protein can determine what functional role the protein plays

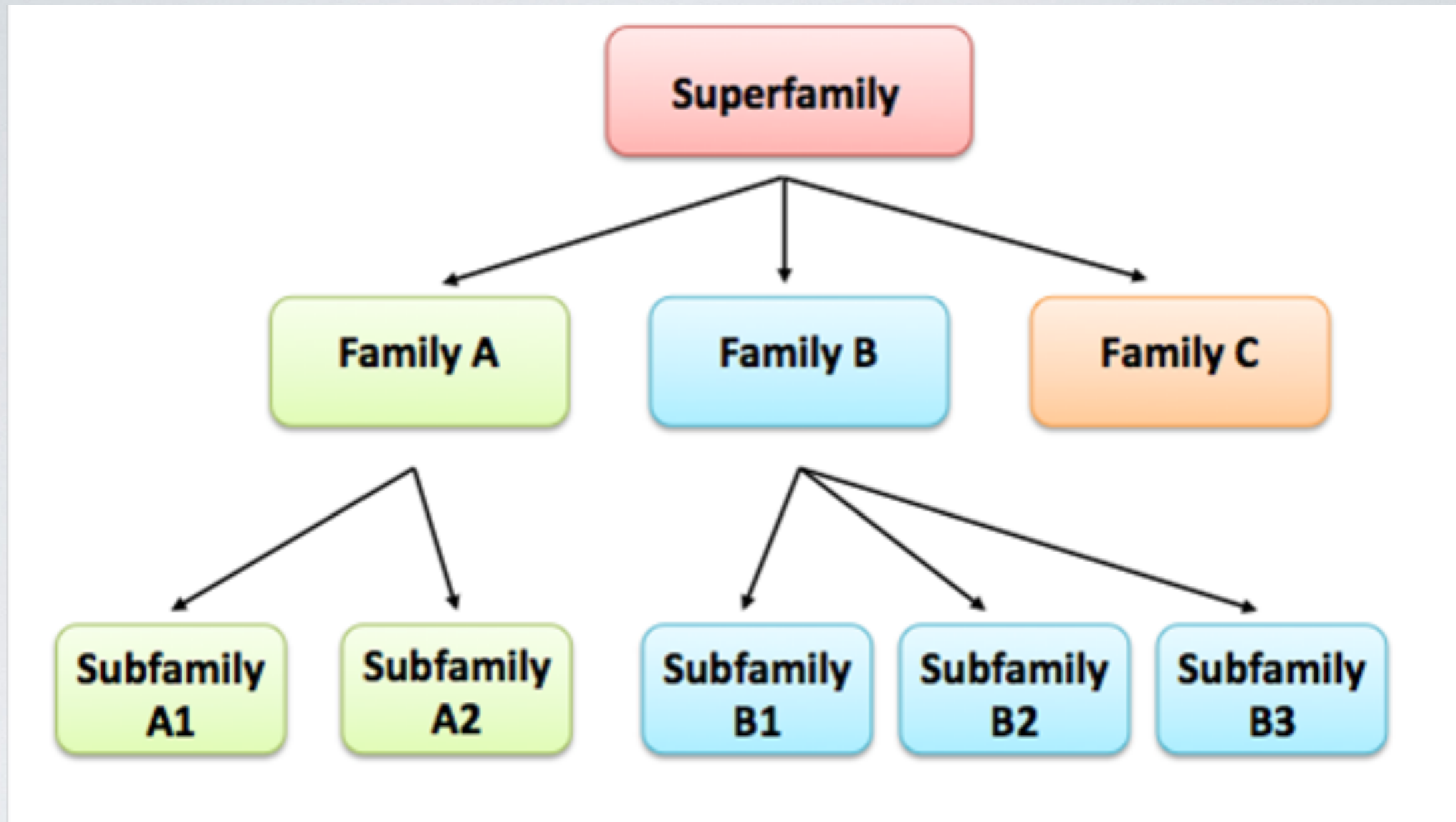
PROTEIN CLASSIFICATION

- Many many (!) proteins if consider all the types found in all organisms
- Proteins can be classified into **Families**
- Families can be classified into **Domains**
 - This can be discrete (eg DNA binding) or part of an enzyme
- Sequence can have **Features**



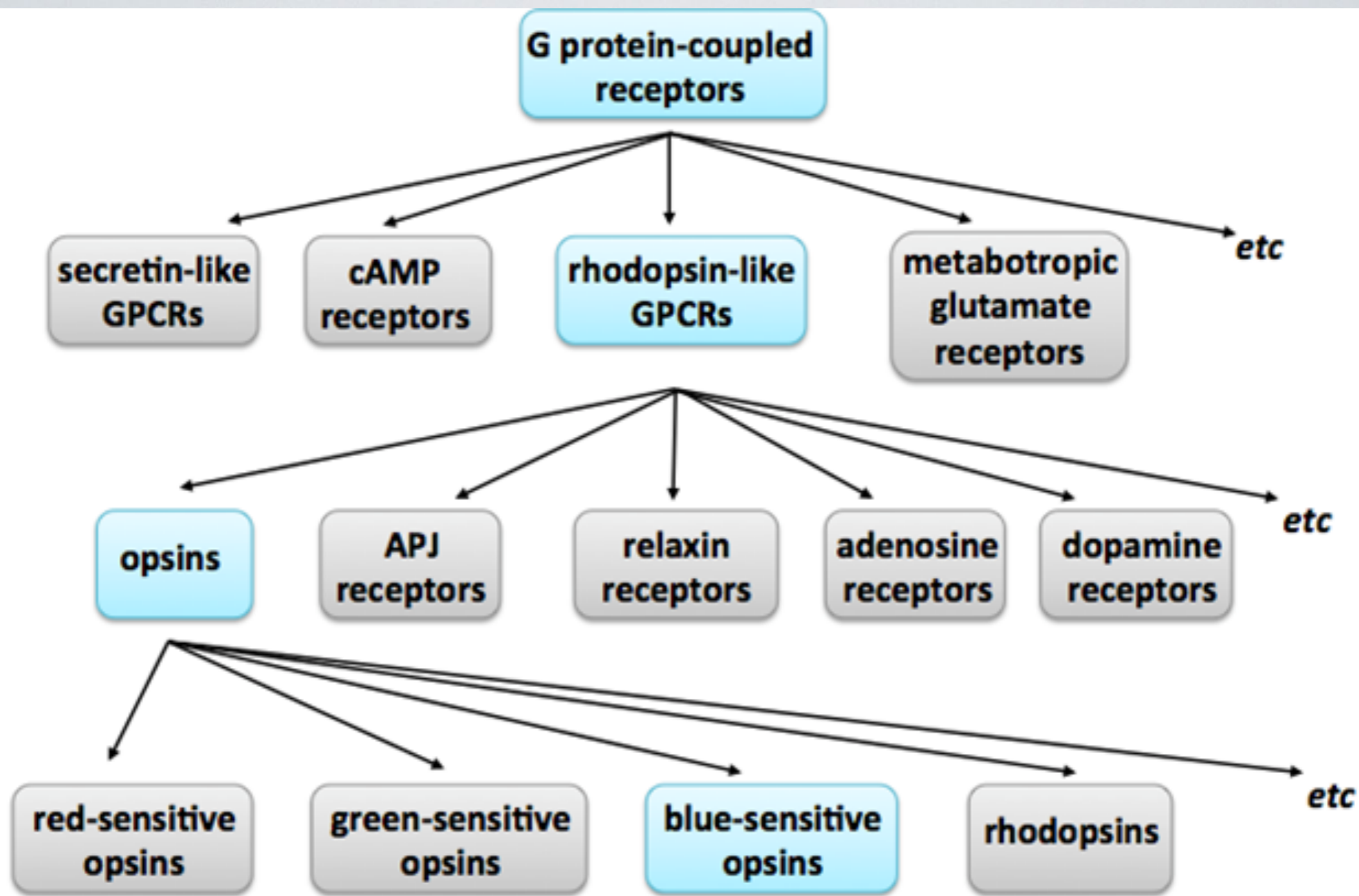
Protein sequence

Protein Structure and
Folds



PROTEIN CLASSIFICATION

Classification of domains requires recognition of regions in protein which are evolutionary conserved and function as a unit



GPCR SUPERFAMILY

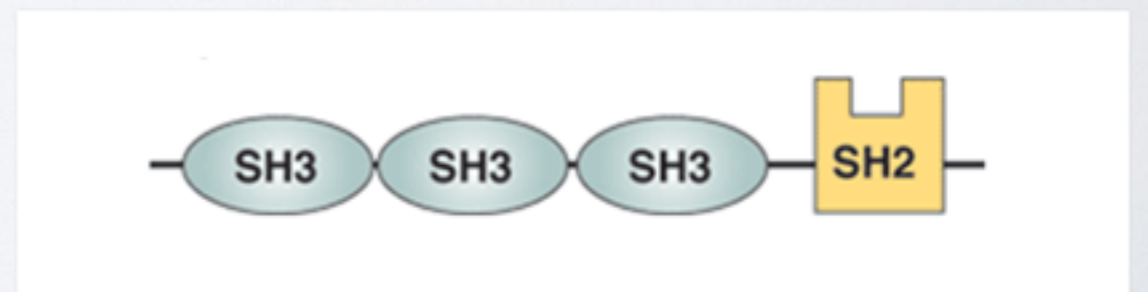
Sub-types of a family - top to bottom this is a classification that is general to specific

PROTEIN DOMAINS

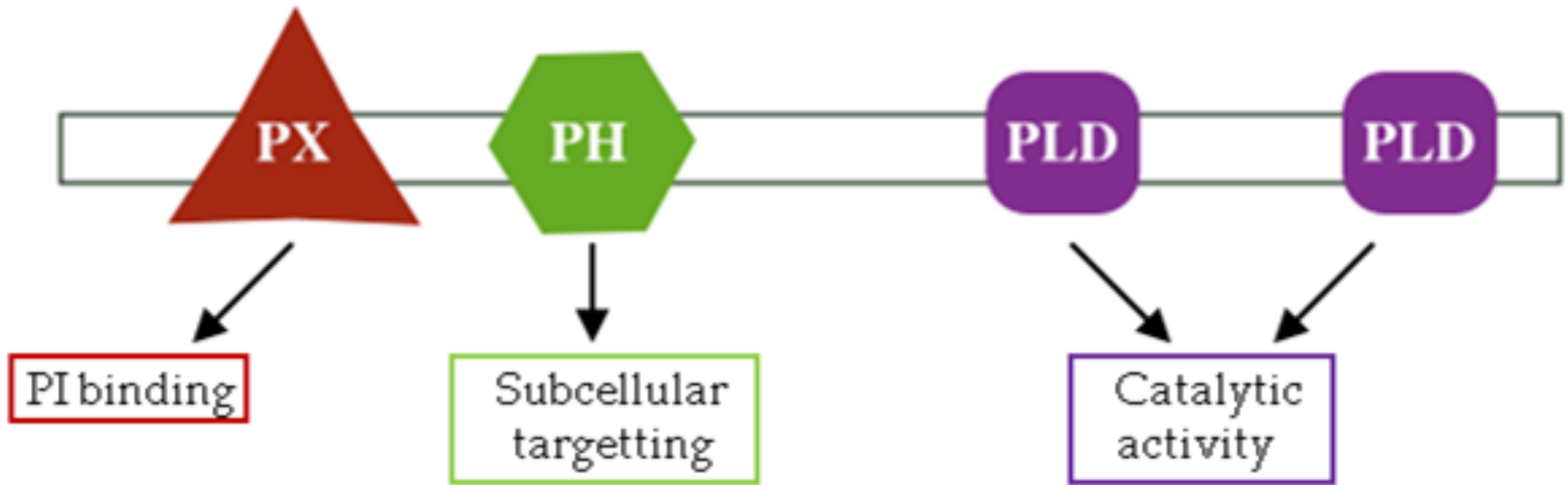
- Distinct functional or structural units of proteins
- SH3 structure shows the 3D folds of the protein when modeled
- Multiple domains can be found within a protein



SH3-structure

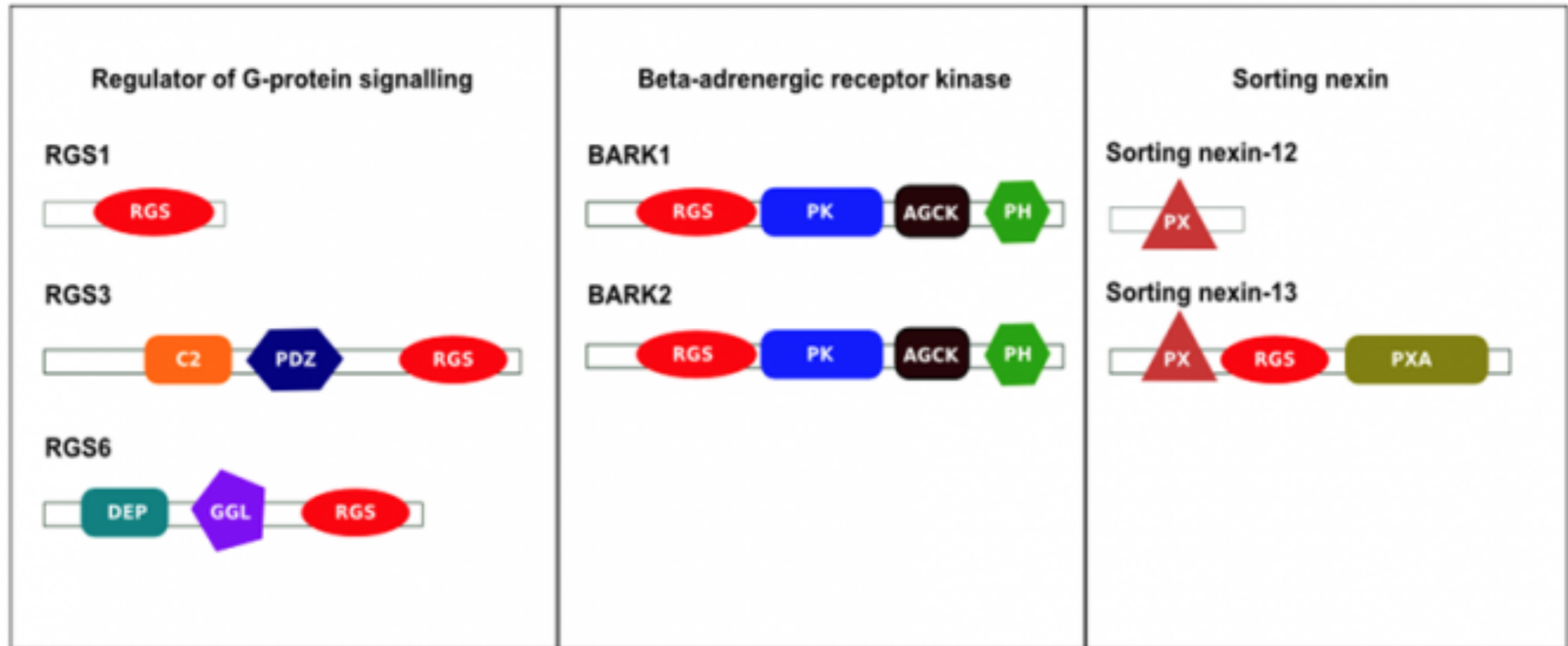


Multidomain protein schematic



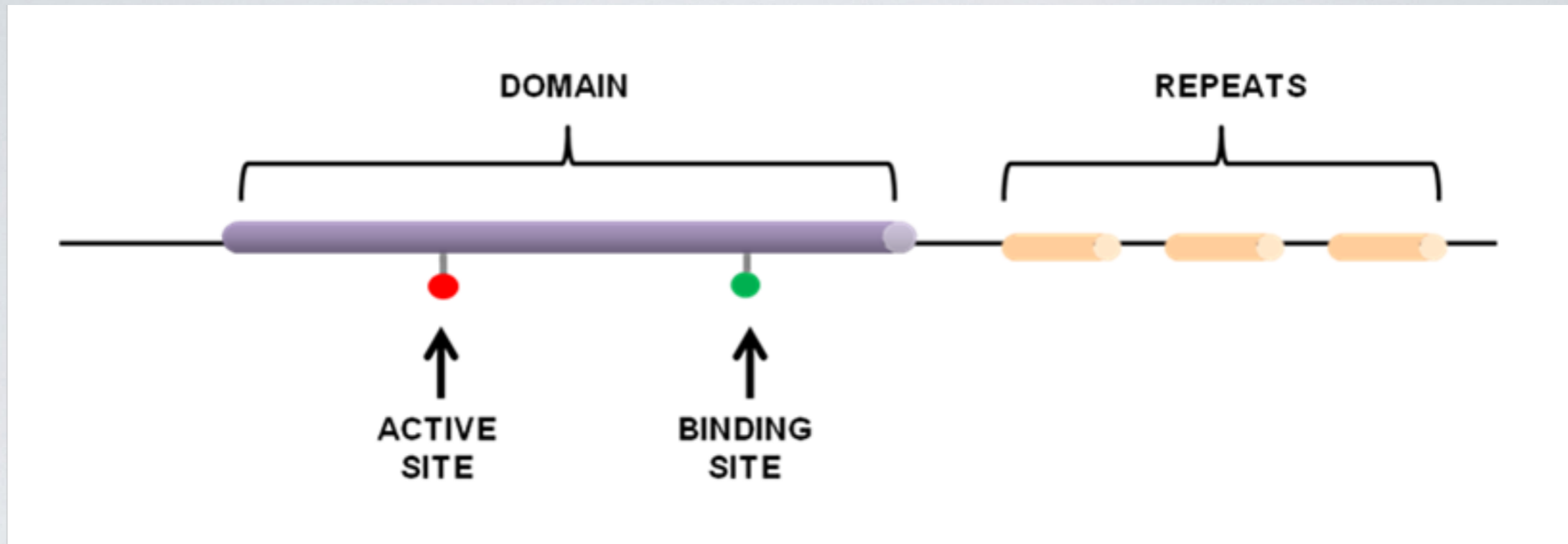
PROTEIN DOMAINS

These domains can have specific functions based



SIGNATURES OF DOMAINS CAN BE CHARACTERISTIC

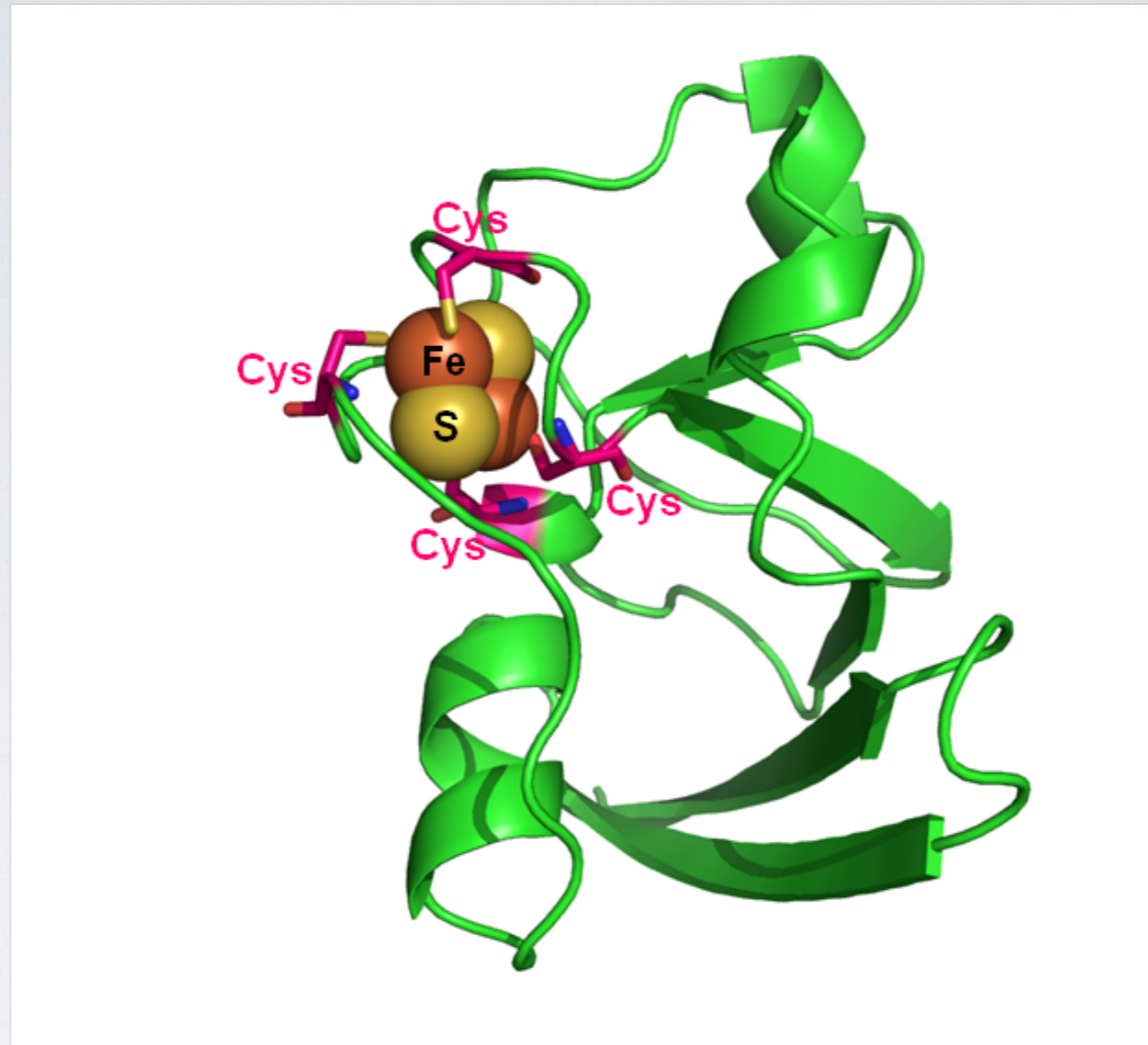
the order and domain content of a protein can be a signature of the type of function



SEQUENCE FEATURES

Sequence features are groups of amino acids which confer certain characteristics

- Could be active site with particular function in enzyme
- binding site for protein-DNA, protein-RNA, protein-protein interactions
- post translational modification site
- repeats within a protein (eg short motifs that repeat)



PROTEINS CLASSIFIED BY FEATURES

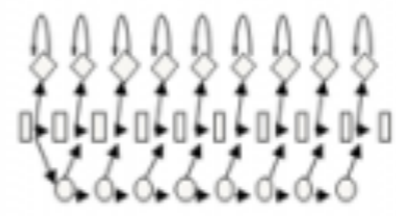
Sequence features like the type of iron- and sulfur-binding residues are used to classify a protein - this is a 2Fe-2S ferredoxin

Protein family/domain

Multiple sequence alignment

ESFL_XENLA/1-176	NTWFLPVITTYVYFLRQIQSSELETITLTPPLKISLDFPIGKTYENNRIGIFNTPEITTYDFAAAAP...T SGA
ESFL_CRACK/1-175	NTWFLPVITTYVYFLRQIQSSELETITLTPPLKISLDFPIGKTYENNRIGIFNTPEITTYDFAAAAP...T SGT
ESFL_MOUSE/1-185	NTWFLPVITTYVYFLRQIQSSELETITLTPPLKISLDFPIGKTYENNRIGIFNTPEITTYDFAAAAP...T SGA
ESFL_XENLA/1-176	SLSYVASSSET...FSSSLTGLTUNVTPSPYVLAFLPQLSPFTIHRDQVPTTLESDGTFVREAAAPPTFSSSDW
ESFL_CRACK/1-175	TLSYVASSSET...FSSSLTGLTUNVTPSPYVLAFLPQLSPFTIHRDQVPTTLESDGTFVREAAAPPTFSSSDW
ESFL_MOUSE/1-185	GLAYVASSSET...FSSSLTGLTUNVTPSPYVLAFLPQLSPFTIHRDQVPTTLESDGTFVREAAAPPTFSSSDW
ESFL_XENLA/1-176	RRDGRERSSANDKPPSKESTIE
ESFL_CRACK/1-175	RRDGRERSSANDKPPSKESTIE
ESFL_MOUSE/1-185	RRDGRERSSANDKPPSKESTIE

Build model



Search



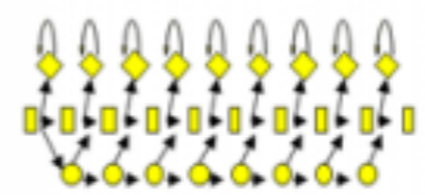
Protein analysis

ITWKGPVCGLDGKTYRNECALL

Significant match



Mature model



CLASSIFYING PROTEINS

Start with known proteins which are similar and determined to be homologous (BLASTP)

```

Q5E940_BOVIN -----MPREDRATWKSNYFLKIIQLLDDYPKCFIVGADNVGSKQMQQIRMSLRGK-AVVLMGKNTMMRKAIRGHLENN--PALE
RLA0_HUMAN -----MPREDRATWKSNYFLKIIQLLDDYPKCFIVGADNVGSKQMQQIRMSLRGK-AVVLMGKNTMMRKAIRGHLENN--PALE
RLA0_MOUSE -----MPREDRATWKSNYFLKIIQLLDDYPKCFIVGADNVGSKQMQQIRMSLRGK-AVVLMGKNTMMRKAIRGHLENN--PALE
RLA0_RAT -----MPREDRATWKSNYFLKIIQLLDDYPKCFIVGADNVGSKQMQQIRMSLRGK-AVVLMGKNTMMRKAIRGHLENN--PALE
RLA0_CHICK -----MPREDRATWKSNYFMKIIQLLDDYPKCFVVGADNVGSKQMQQIRMSLRGK-AVVLMGKNTMMRKAIRGHLENN--PALE
RLA0_RANSY -----MPREDRATWKSNYFLKIIQLLDDYPKCFIVGADNVGSKQMQQIRMSLRGK-AVVLMGKNTMMRKAIRGHLENN--SALE
Q7ZUG3_BRARE -----MPREDRATWKSNYFLKIIQLLDDYPKCFIVGADNVGSKQMQQIRMSLRGK-AVVLMGKNTMMRKAIRGHLENN--PALE
RLA0_ICTPU -----MPREDRATWKSNYFLKIIQLLNDYPKCFIVGADNVGSKQMQQIRMSLRGK-AIVLMGKNTMMRKAIRGHLENN--PALE
RLA0_DROME -----MVRENKAAWKAQYFIKVVVLFDEFKCFIVGADNVGSKQMQQIRMSLRGK-AVVLMGKNTMMRKAIRGHLENN--PQLE
RLA0_DICDI -----MSGAG-SKRKKLFIEKATKLFTTYDKMIVAEADFGVSSQLQKIRKSIIRGI-GAVLMGKNTMIRKVIIRDLADSK--PELD
Q54LP0_DICDI -----MSGAG-SKRKNVFIKATKLFTTYDKMIVAEADFGVSSQLQKIRKSIIRGI-GAVLMGKNTMIRKVIIRDLADSK--PELD
RLA0_PLAF8 -----MAKLSKQQKKQMYIEKLSLIQQYSKILIVHVDNVGSKQMQQIRMSLRGK-ATILMGKNTIRIRALKKNLQAV--PQIE
RLA0_SULAC -----MIGLAVTTTTKKIAKWVDEVAELTEKLTHTKTIIANIEGFPADKLHEIRKKLRGK-ADIKVTKNNLNFNIALKNAG-----YDTK
RLA0_SULTO -----MRIMAVITQERKIAKWVIEEVKELEQKLREYHTIIIANIEGFPADKLHDIRKKMRGM-AEIKVTKNTLFGIAAKNAG-----LDVS
RLA0_SULSO -----MKRLALALKQRKVASWVKELETELKNSNTILIGNLEGFPADKLHEIRKKLRGK-ATIKVTKNTLFGIAAKNAG-----IDIE
RLA0_AERPE MSVVS LVGQMYKREKPIPEWKTLMLELEELFSKHRVVFADLTGPTFVVRVRKKLWKK-YPMVAKKRIILRAMKAAGLE---LDDN
RLA0_PYRAE -MMLAIGKRRYVRTROYPARKVKIVSEATELLQKYPYVFLFDLHGLSSRILHEYRYRLRY-GVIKIIPKTLFKIAFTKVYGG---IPAE
RLA0_METAC -----MAEERHHTHEIPQWKKDEIENIKELIQSHKVFVGMVRIEGILATKIQKIRRDLDKDV-AVLKVSNTLTERALNQLG-----ETIP
RLA0_METMA -----MAEERHHTHEIPQWKKDEIENIKELIQSHKVFVGMVRIEGILATKIQKIRRDLDKDV-AVLKVSNTLTERALNQLG-----ESIP
RLA0_ARCFU -----MAAVRGS---PPEYKVRAVEEIKRMISPKPVVAIVSFRNVPAGQMQKIRREFRGGK-AEIKVVKNTLLERALDALG-----GDYL
RLA0_METKA MAVKAKGQPPSGYEPKVAEWKRREVKELKELMDEYENVGLVDLEGIPAPQLOEIRAKLRERDTIIRMSRNTLMRIALEEKLDER--PELE
RLA0_METTH -----MAHVAEWKKKEVQELHDLIKGYEVVGIANLADIPARQLOKMRQTLRDS-ALIRMSKKTLSLALAEKAGREL--ENVD
RLA0_METTL -----MITAESEHKIAPWKIEEVNKLKELLKNGQIVALVDMMEVPAQLOEIRDKIR-GTMTLKMSRNTLIERAIKEVAEETGNPEFA
RLA0_METVA -----MIDAKSEHKIAPWKIEEVNALKELLKSANVIALIDMMEVPAVQLOEIRDKIR-DQMTLKMSRNTLIKRAVEEVAEETGNPEFA
RLA0_METJA -----METKVKAHVAPWKIEEVKTLKGLIKSKPVVAIVDMMDVPAPQLOEIRDKIR-DKVKLMSRNTLIIRALKEAAEELNNPKLA

```

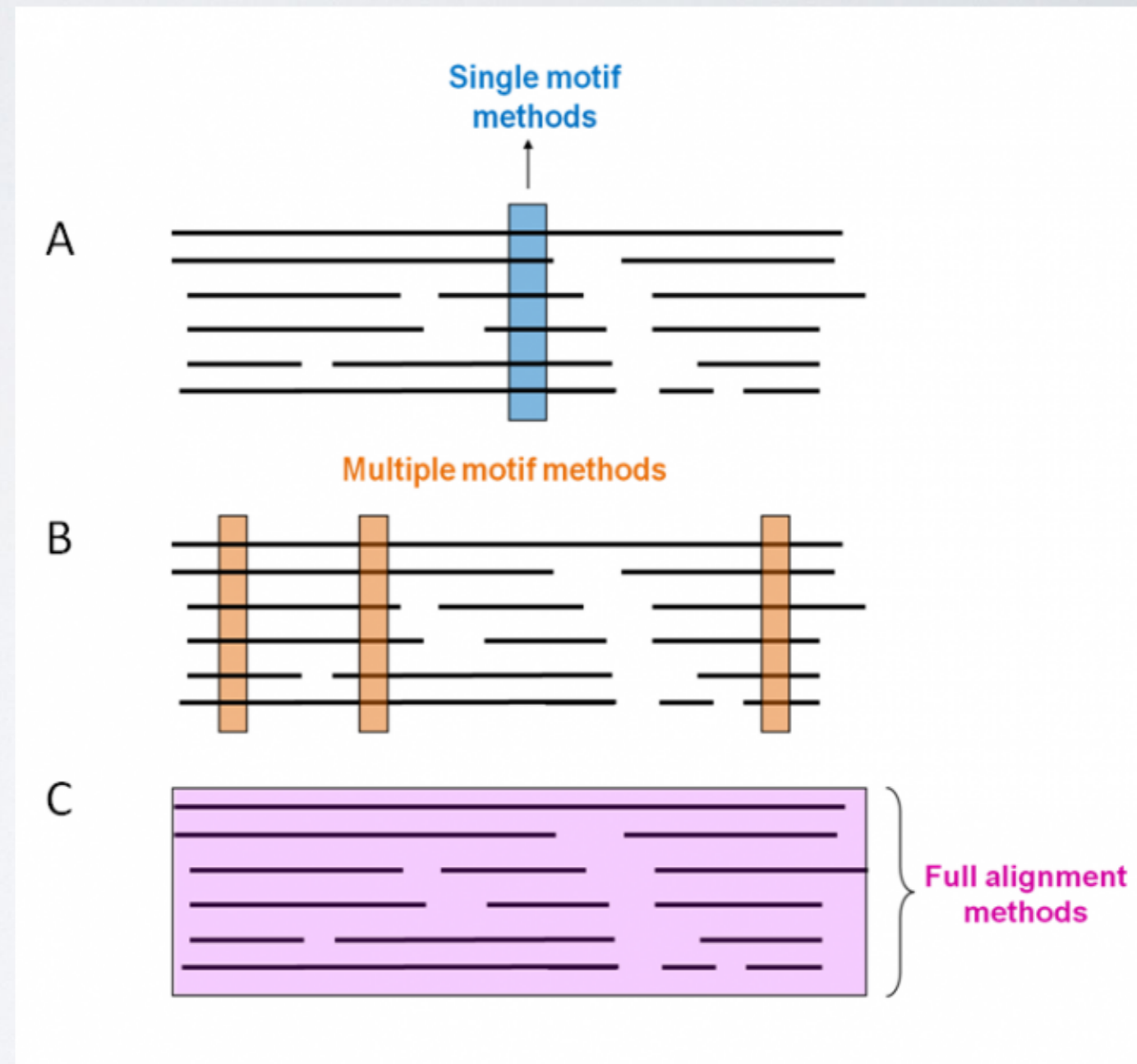


PROTEIN ALIGNMENT

Considering multiple sequences in the alignment so is more sensitive than BLAST. Only some residues are informative to classify the sequence
This is revealed through the multiple alignment.

HOW TO CLASSIFY

- Motif pattern
- Multiple Motifs - interspersed
 - Profile
 - Fingerprint
- Full alignment method with Hidden Markov Models



Sequence alignment



Motif



Extract pattern sequences

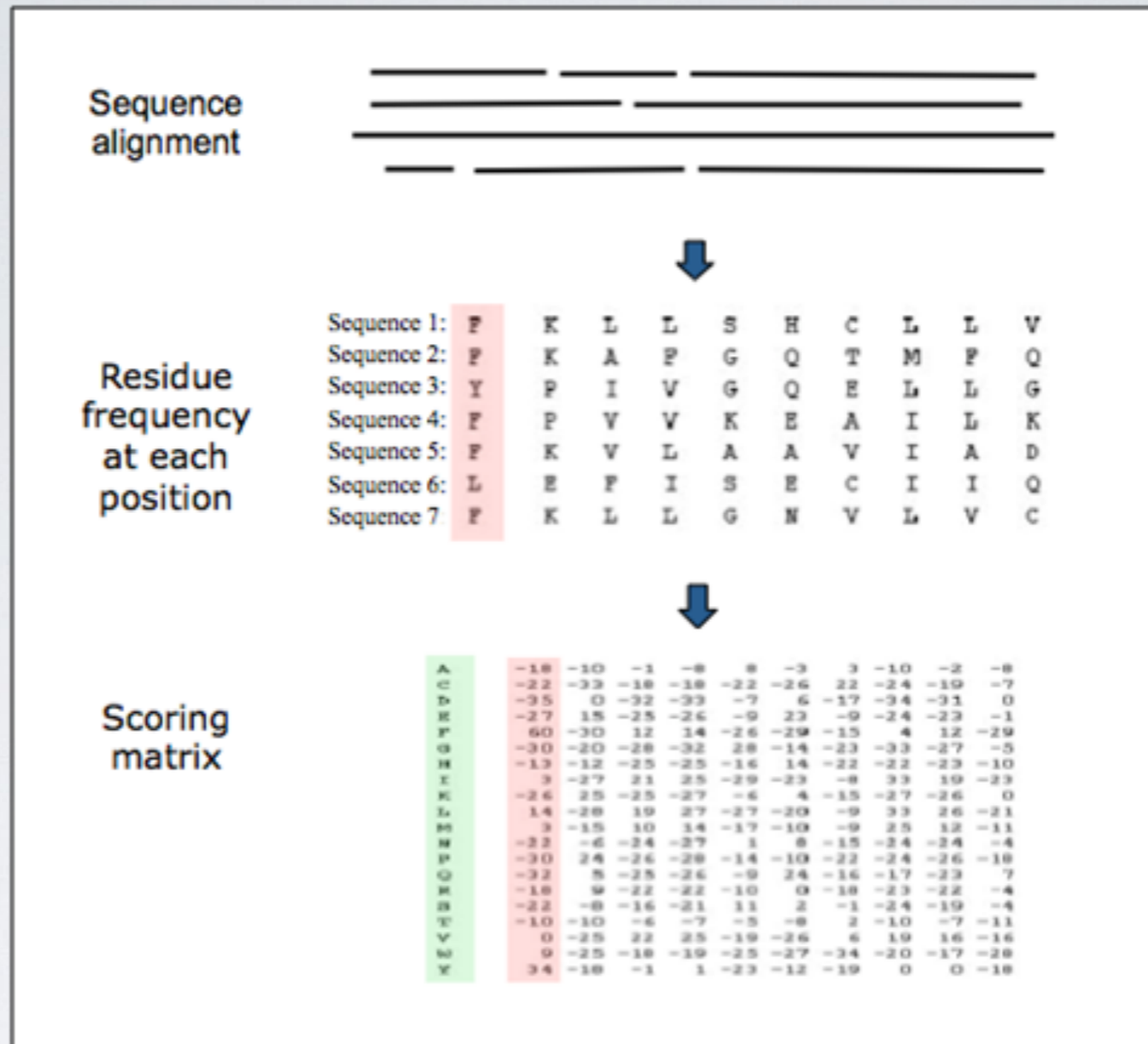


Build regular expression

[AC]-x-V-x(4)-{ED}

MOTIF DEFINED BY A PATTERN

Can write down the pattern with a series of letters and then logic called a Regular Expression



PROFILE

Created by converting a multiple alignment into a Position Specific Scoring Matrix - PSSM
 Amino acids at each position in the alignment are scored according to the frequency with which they occur

SCORING A PROFILE

	1	2	3
A	0.01	0.04	0.02
G	0.02	0.02	0.03
C	0.02	0.93	0.93
T	0.95	0.01	0.02

Probability Matrix

$$\log_2(0.01/0.25) = -4.6$$

	1	2	3
A	-4.6	-2.6	-3.6
G	-3.6	-3.6	-3.0
C	-3.6	1.8	1.8
T	1.9	-4.6	-3.6

Position Weight Matrix

. . . TCC . . .
 . . . TCG . . .
 . . . TCC . . .
 . . . TAC . . .
 . . . GCC . . .
 . . . TCC . . .

SCORING A PROFILE

	1	2	3
A	0.01	0.04	0.02
G	0.02	0.02	0.03
C	0.02	0.93	0.93
T	0.95	0.01	0.02

$$\log_2(0.01/0.25) = -4.6$$

Let's score the sequence
AGATCCTGCTCG

	1	2	3
A	-4.6	-2.6	-3.6
G	-3.6	-3.6	-3.0
C	-3.6	1.8	1.8
T	1.9	-4.6	-3.6

Position Weight Matrix

$$\underline{\text{AGATCCTGCTCG}} \quad \text{Score} = \begin{matrix} (A, 1) & (G, 2) & (A, 3) \\ -4.6 & + & -3.6 & + & -3.6 & = & -11.8 \end{matrix}$$

$$\underline{\text{AGATCCTGCTCG}} \quad \text{Score} = \begin{matrix} (G, 1) & (A, 2) & (T, 3) \\ -3.6 & + & -2.6 & + & -3.6 & = & -9.8 \end{matrix}$$

$$\underline{\text{AGATCCTGCTCG}} \quad \text{Score} = \begin{matrix} (T, 1) & (C, 2) & (C, 3) \\ 1.9 & + & 1.8 & + & 1.8 & = & 5.5 \end{matrix}$$

Score above 0 is a good score!

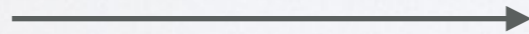
SCORING A PROFILE

	1	2	3
A	-4.6	-2.6	-3.6
G	-3.6	-3.6	-3.0
C	-3.6	1.8	1.8
T	1.9	-4.6	-3.6

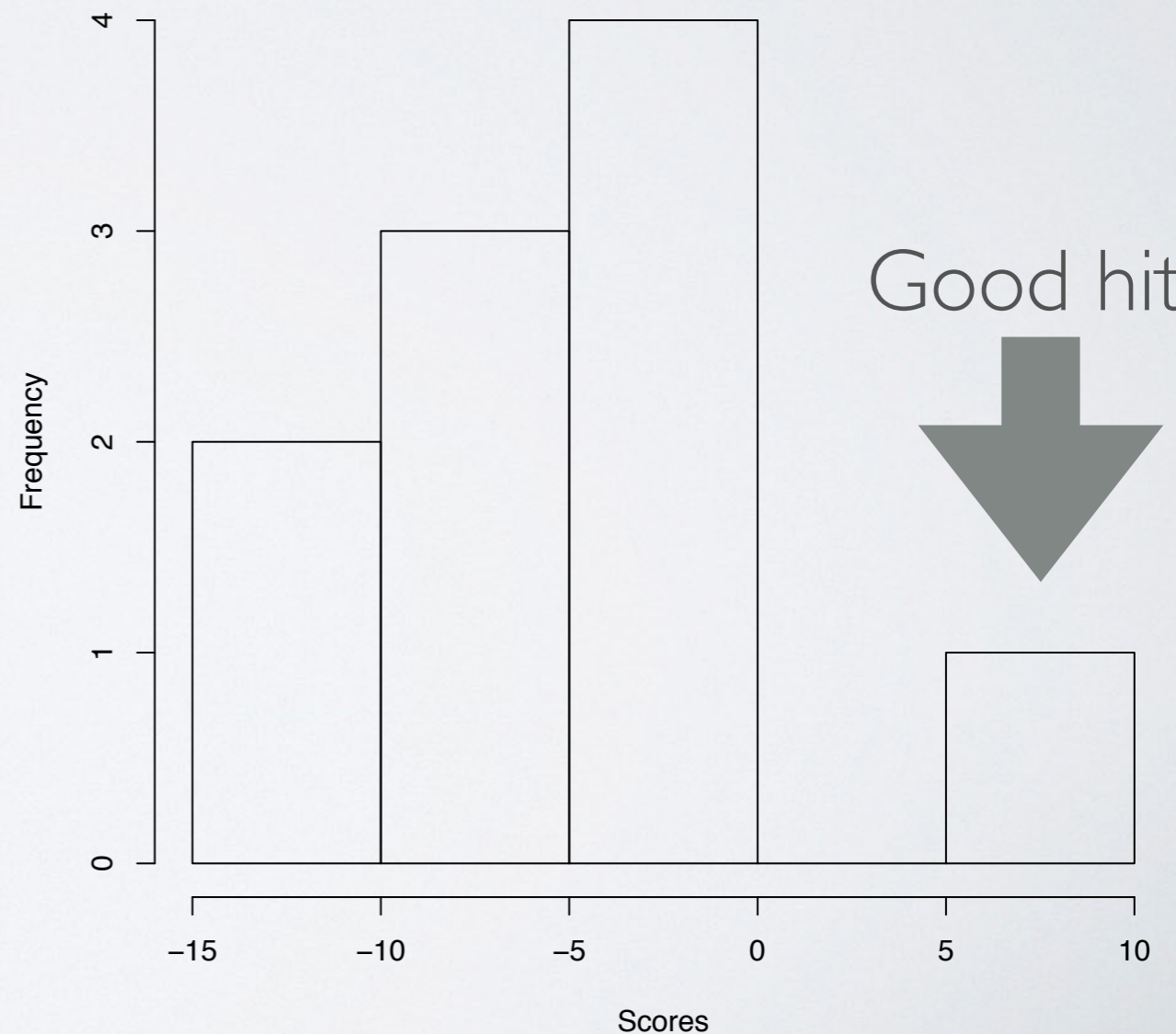
Consider the distribution of scores across the whole sequence to evaluate if there is a significance as well.

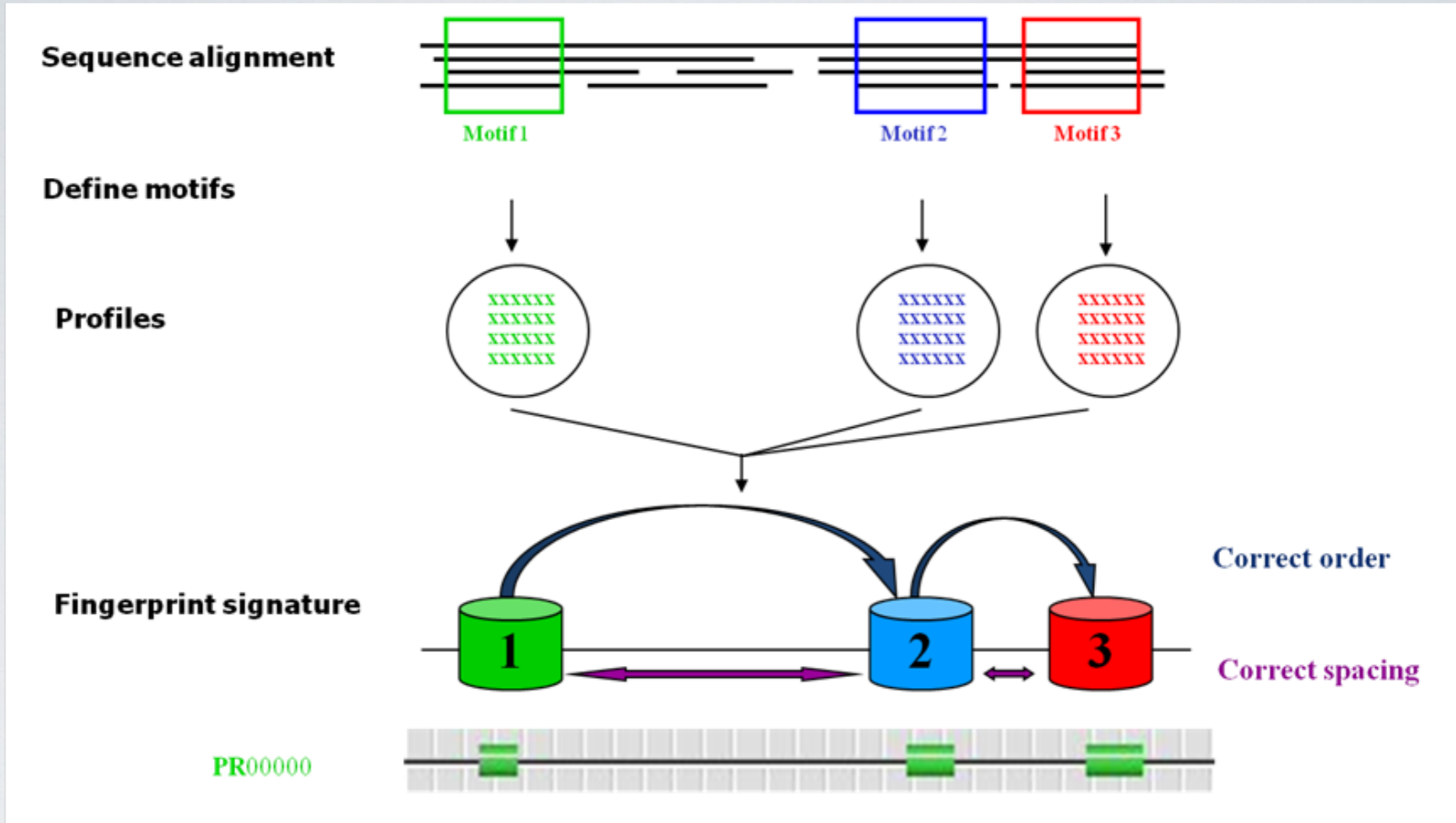
AGATCTTGCTCG

-11.8, -9.8, 5.5, etc



Distribution of scores





FINGERPRINT

Combination of motif or profile into a signature

CLCN1_HUMAN	F	P	L	V	L	I	L	F	S	A	L	F	C	H	L	I	S	P	Q	A	V	G	S	G	I	P	E	M	K	T	I	L	R	G	V	V	L	K	E	Y	L	T	M	K	A	F	V	A	K
CLCN1_RAT	F	P	L	I	L	I	L	F	S	A	L	F	C	Q	L	I	S	P	Q	A	V	G	S	G	I	P	E	M	K	T	I	L	R	G	V	V	L	K	E	Y	L	T	L	K	A	F	V	A	K
CLCN2_HUMAN	Y	P	V	V	L	I	T	F	S	A	G	F	T	Q	I	L	A	P	Q	A	V	G	S	G	I	P	E	M	K	T	I	L	R	G	V	V	L	K	E	Y	L	T	L	K	T	F	I	A	K
CLCN2_MOUSE	Y	P	V	V	L	I	T	F	S	A	G	F	T	Q	I	L	A	P	Q	A	V	G	S	G	I	P	E	M	K	T	I	L	R	G	V	V	L	K	E	Y	L	T	L	K	T	F	V	A	K
CLCN3_RAT	W	A	L	S	F	A	F	L	A	V	S	L	V	K	V	F	A	P	Y	A	C	G	S	G	I	P	E	I	K	T	I	L	S	G	F	I	I	R	G	Y	L	G	K	W	T	L	M	I	K
CLCN3_PONAB	W	A	L	S	F	A	F	L	A	V	S	L	V	K	V	F	A	P	Y	A	C	G	S	G	I	P	E	I	K	T	I	L	S	G	F	I	I	R	G	Y	L	G	K	W	T	L	M	I	K
CLCN3_RABIT	W	A	L	S	F	A	F	L	A	V	S	L	V	K	V	F	A	P	Y	A	C	G	S	G	I	P	E	I	K	T	I	L	S	G	F	I	I	R	G	Y	L	G	K	W	T	L	M	I	K

- Amino acids relatively well conserved across all chloride channel protein family members
- Amino acids uniquely conserved in chloride channel protein 3 subfamily members

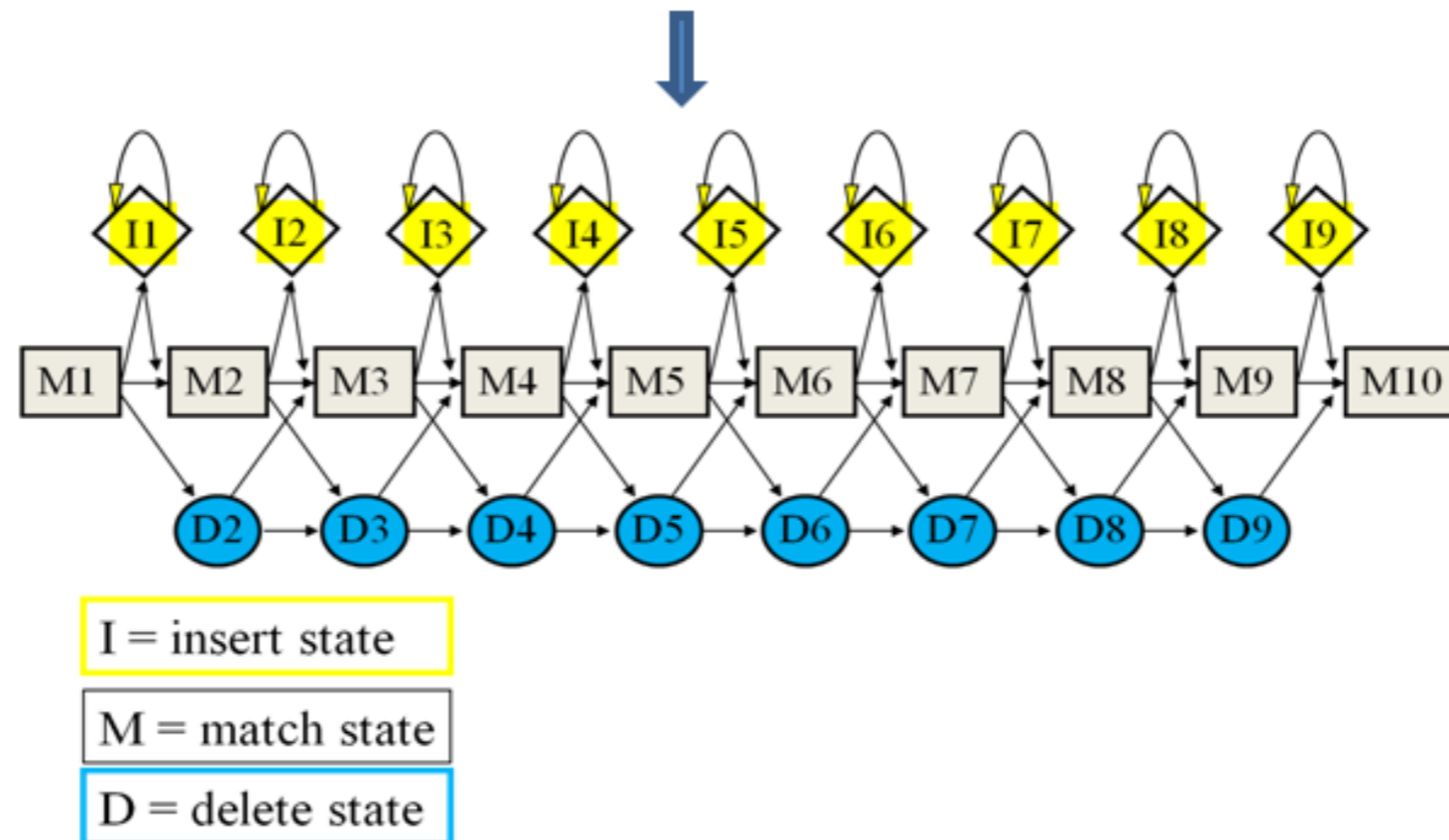
WHY ARE FINGERPRINTS USEFUL?

Can capture and model small differences between sub-families can capture the individual differences.

In this example of a chloride channel protein family identified by blue box a subset can be further classified into channel 3 - subfamily

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



HIDDEN MARKOV MODELS

Can model the alignment by capturing insertion and deletions and probabilistically score sequence similarity.

DATABASES OF PROTEIN DOMAINS

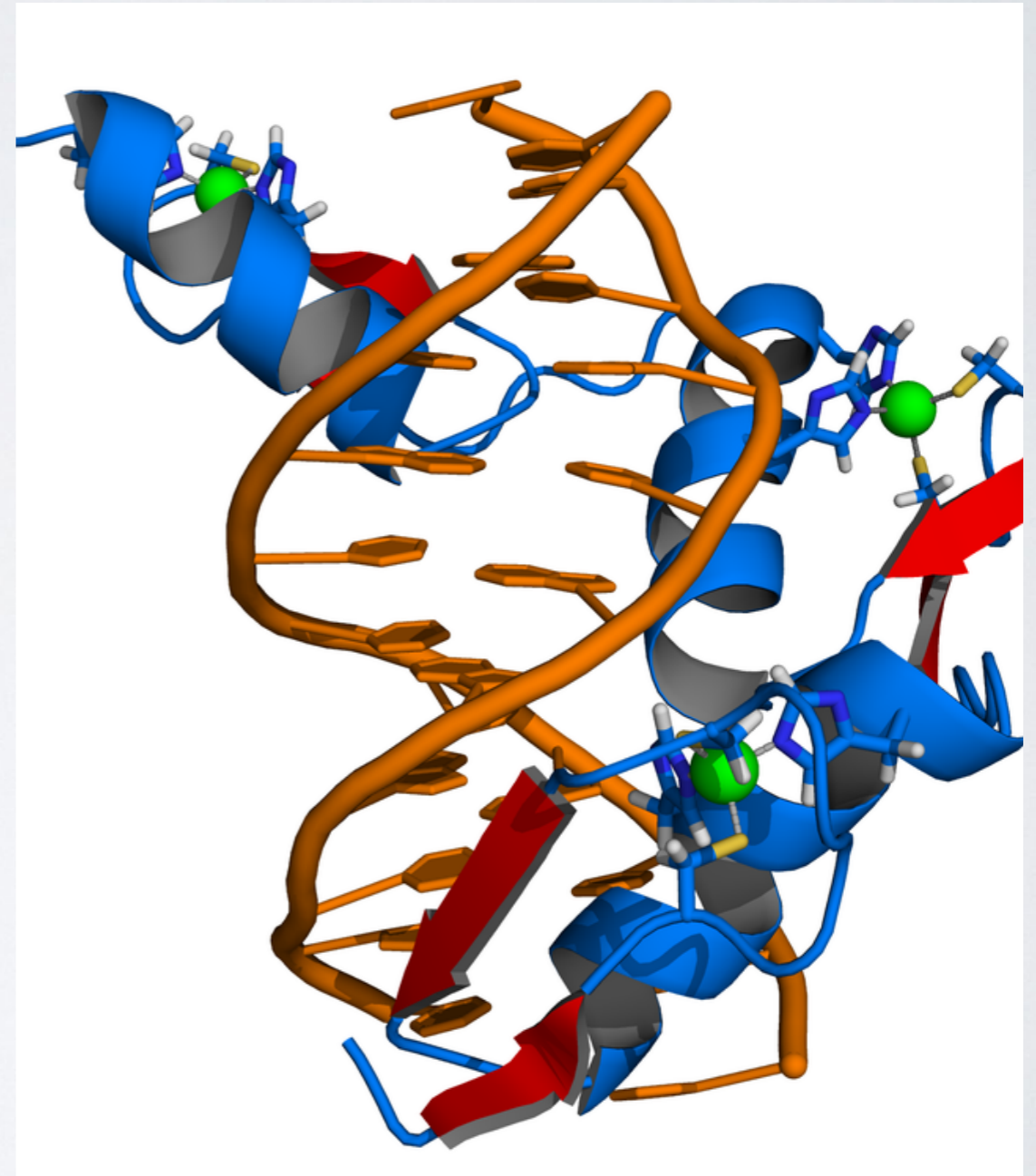
- Pfam - Protein <https://pfam.xfam.org/>; Panther - <http://www.pantherdb.org/>; SMART - <http://smart.embl-heidelberg.de/>
 - databases of HMMs of domains
- Prosite - <https://prosite.expasy.org/> - motifs
- Interpro - <https://www.ebi.ac.uk/interpro/>
 - HMMs + Profiles + fingerprints

PROTEINS TO FUNCTION

- Together these domains and classifications can provide ways to link an unknown sequence to function
- If domains that makeup the protein are known can make guess about the protein function even if a homolog does not have a known function in other species
- Domains can be shared among many types of proteins
- Shuffling of domains and motifs can provide new function

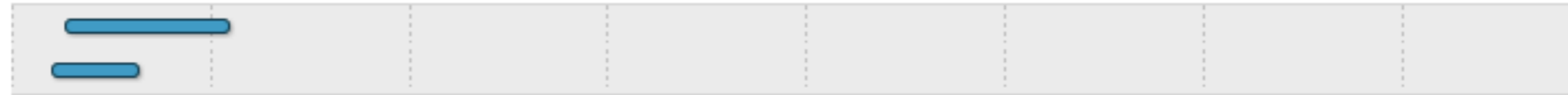
SOME PROTEIN DOMAINS

- Zinc-finger - https://en.wikipedia.org/wiki/Zinc_finger
- Typically bind DNA, protein, RNA
- Often part of transcription factors



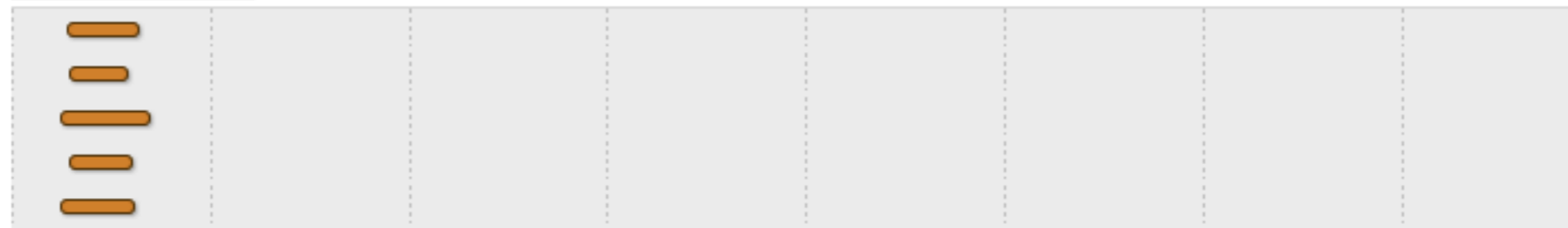
Detailed signature matches

H IPR036864 Zn(2)-C6 fungal-type DNA-binding domain superfamily



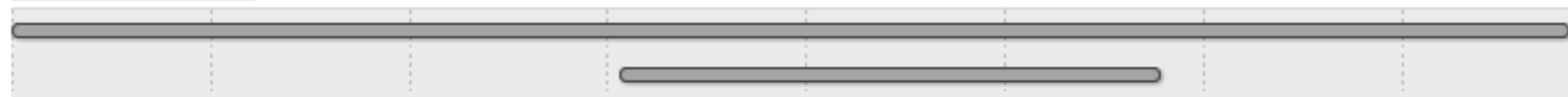
- ▶ G3DSA:4.10.24...
- ▶ SSF57701 (Zn2/Cys6 ...)

D IPR001138 Zn(2)-C6 fungal-type DNA-binding domain



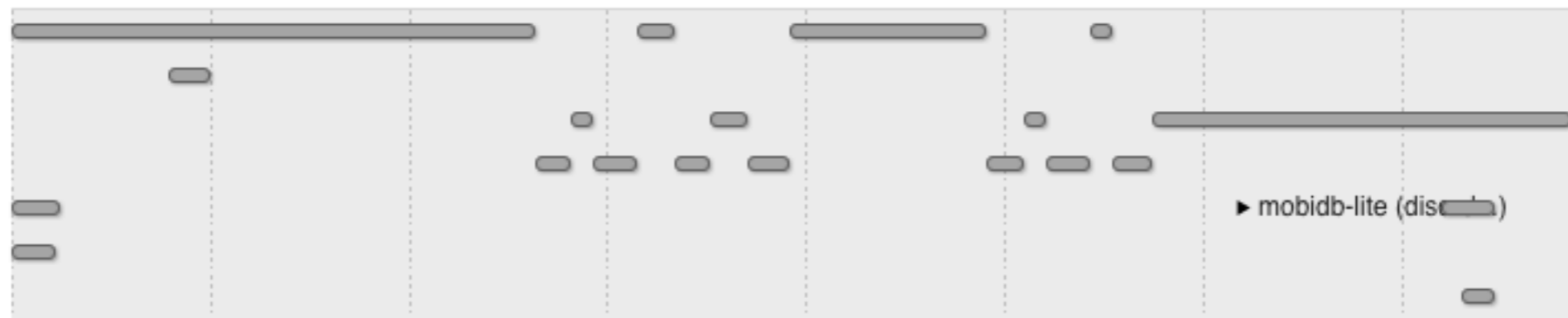
- ▶ PF00172 (Zn_clus)
- ▶ PS00463 (ZN2_CY6_FU...)
- ▶ SM00066 (gal4_2)
- ▶ PS50048 (ZN2_CY6_FU...)
- ▶ cd00067 (GAL4)

i no IPR Unintegrated signatures



- ▶ PTHR31405 (FAMILY N...)
- ▶ cd12148 (fungal_TF_MHR)

Other features



- ▶ CYTOPLASMIC_D... (C...)
- ▶ Coil
- ▶ NON_CYTOPLASM... (N...)
- ▶ TRANSMEMBRANE (Tran...)
- ▶ mobidb-lite (dis...)
- ▶ mobidb-lite (Polyam...)
- ▶ mobidb-lite (Polar)

Residue annotation

TAKE HOME POINTS

- Proteins can be classified by their similarity
- Parts of proteins can further be found to be conserved and functional
- Motifs, Profiles, Fingerprints, HMMs
- Domains discovery can be used to assign function