

Motifs, Profiles and Domains

Michael Tress

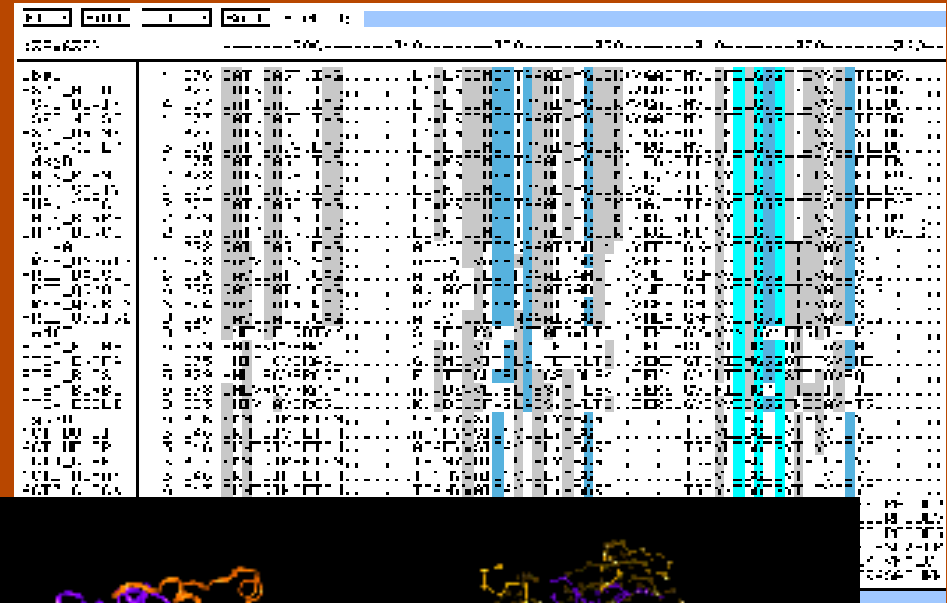
Protein Design Group

Centro Nacional de Biotecnología, CSIC

Comparing Two Proteins

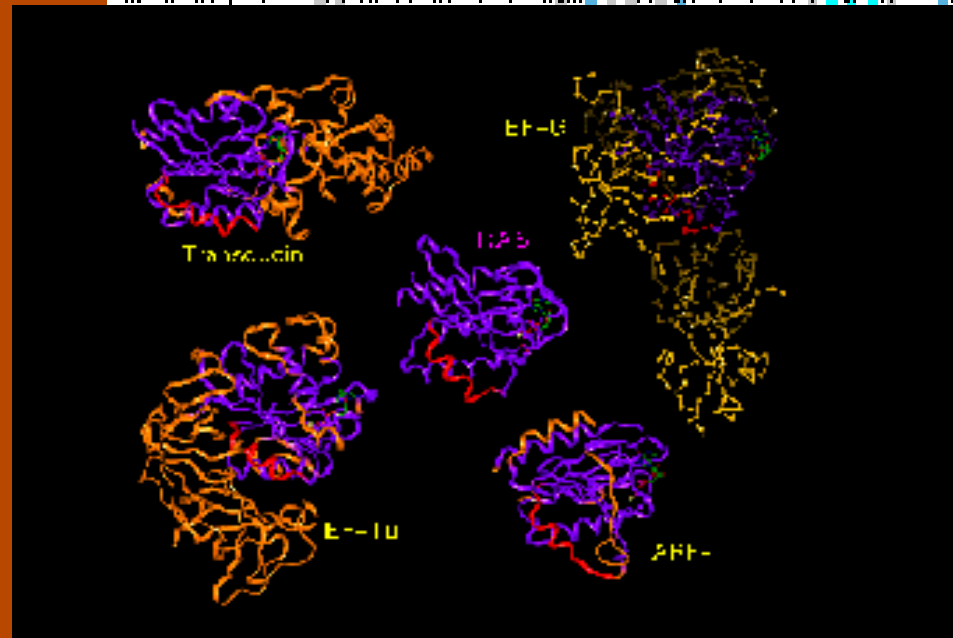
Sequence Alignment

Determining the pattern of evolution and identifying conserved regions that may be of structural or functional importance.



Structure Alignment

Structure is more conserved than sequence. If structural similarity exists it is still possible to infer common ancestry, even in the absence of sequence similarity



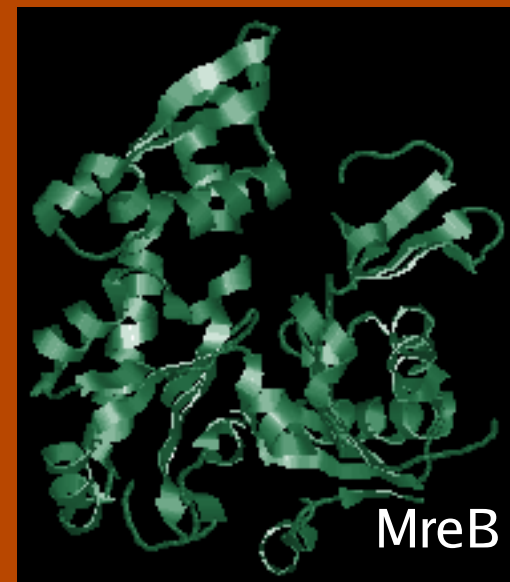
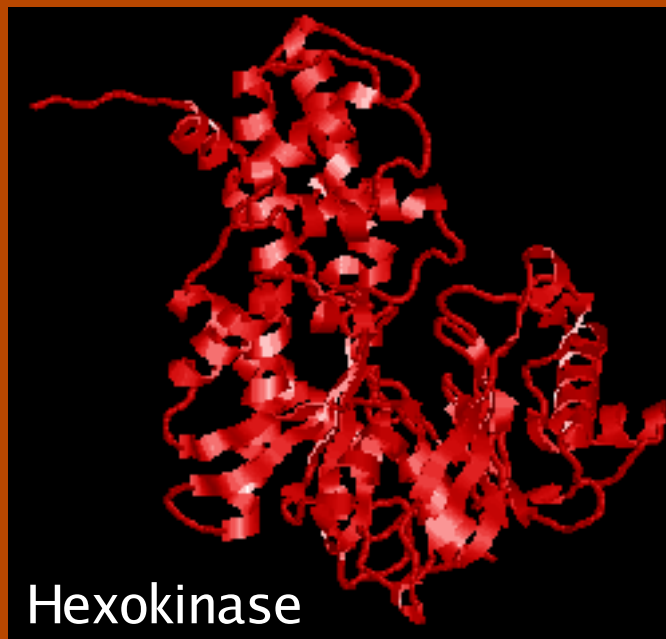
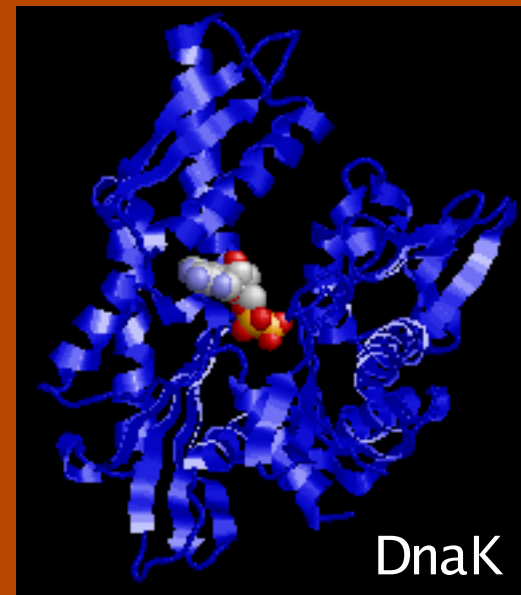
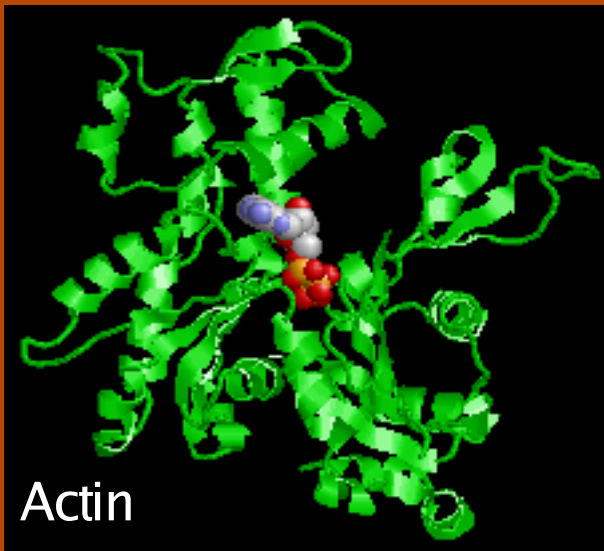
Sequence Analysis

Aims

- To identify similarity between sequences.
- To deduce common ancestry from sequence similarity.
- Sequence classification
- Sequence similarity based function prediction.
- Sequence similarity based structure prediction
- The prediction of functional and physical interaction between proteins

Chaperones (**dnak**), Proteins involved in bacterial cell wall production (**ftsA**, **mreB**), hexokinases (**hxx**), actin (**act**)

File	Feat	Column	Start	Picked:
(30x130)			300	310
			320	330
			340	350
			360	
1bc1	4	376	QAT<DAGT.IAG.....LWLRRIINEPTAFAAIAYGLDKKYGAERNV_IID_GGGTFDWSILTIEDG...	
HS7C_HMAN	4	377	QAT<DAGT.IAG.....LWLRRIINEPTAFAAIAYGLDKKYGAERNV_IID_GGGTFDWSILTIEDG...	
HS7C_BVIN	4	377	QAT<DAGT.IAG.....LWLRRIINEPTAFAAIAYGLDKKYGAERNV_IID_GGGTFDWSILTIEDG...	
HS7C_MUCD	4	377	QAT<DAGT.IAG.....LWLRRIINEPTAFAAIAYGLDKKYGAERNV_IID_GGGTFDWSILTIEDG...	
HS7D_DROMF	4	377	QAT<DAGT.IAG.....LWLRRIINEPTAFAAIAYGLDKKYGAERNV_IID_GGGTFDWSILTIEDG...	
HS7D_XENLA	5	378	QAT<DAGT.IAG.....LWLRRIINEPTAFAAIAYGLDKKYGAERNV_IID_GGGTFDWSILTIEDG...	
1dkg2	4	375	QAT<DAGR.IAG.....LVKRIINEPTAFAALAYGLDKK...TGNRTAVYD_GGGTFDWSILTIEDG...	
DNAK_PASMU	2	378	QAT<DAGR.IAG.....LVKRIINEPTAFAALAYGLDKK...TGNRTAVYD_GGGTFDWSILTIEDG...	
DNAK_SLLIY	1	377	QAT<DAGR.IAG.....LVKRIINEPTAFAALAYGLDKK...TGNRTAVYD_GGGTFDWSILTIEDG...	
DNAK_YTRCH	9	377	QAT<DAGR.IAG.....LVKRIINEPTAFAALAYGLDKK...TGNRTAVYD_GGGTFDWSILTIEDG...	
DNAK_B_RPS	2	379	QAT<DAGR.IAG.....LVKRIINEPTAFAALAYGLDKK...TGNRTAVYD_GGGTFDWSILTIEDG...	
DNAK_D_RDC	2	300	QAT<DAGR.IAG.....LVKRIINEPTAFAALAYGLDKK...TGNRTAVYD_GGGTFDWSILTIEDG...	
1jra2	4	372	RAT<DAGI.FAG.....ARKVFTTFPPAAATGANNH...YFFFSGNKVMATGGGTFEAWFEI...	
MREB_U27013	11	328	RAYVDAAK.SAG.....AREVYLVAEYFAAIGAGLP...VEEF_LGNMIVDILGGIIDIAYLEL...	
MRED_DACDU	0	325	RAYIDATR.QAG.....ARDAYPICEPTAFAIGANLP...YJCTGQKVVWIDGGTTEVAITEL3...	
MREB_Q2K3H5	6	325	RAYEDATK.QAG.....ARKVYTLCEPFAPAIGAGLP...YJCTGQKVVWIDGGTTEVAITEL3...	
MREB_U23066	6	324	RAYVDAAK.SAG.....AREVYLVAEYFAAIGAGLP...VEEF_LGNMIVDILGGIIDIAYLEL...	
MREB_Q211G6	9	326	RAYTFASS.QAG.....ARQVHTTFFPPAAATGAGLP...YHFAFGNVMATGGGTFEAWFEI...	
1e4ft	8	364	EMFYNFLQDTVK.....S.PFQLKSSLVSTHAGVLT...PEKDRGVVWVNGYNFTGLDFYKN...	
FTSA_CNTIR	5	379	HTIRKCYCHAGL.....V.VHCLVITPLA...TITLSD...GCKDFGTVIDMGGGQT...TAVM...D...	
FTSA_FNTFA	1	375	HTIRKCYCHAGL.....V.VHCLVITPLA...TITLSD...GCKDFGTVIDMGGGQT...TAVM...D...	
FTSA_BFC5U	5	379	HTIRKCYCHAGL.....V.VHCLVITPLA...TITLSD...GCKDFGTVIDMGGGQT...TAVM...D...	
FTSA_DIRDU	5	370	HTIRKCYCHAGL.....V.VHCLVITPLA...TITLSD...GCKDFGTVIDMGGGQT...TAVM...D...	
FTSA_EIO_I	8	363	HTIRKCYCHAGL.....V.VHCLVITPLA...TITLSD...GCKDFGTVIDMGGGQT...TAVM...D...	
1yag2	5	346	E<NTQIMFETFN.....TPAFYYSIQAVLELYASGRT.....TGIV_D3GDGYTHWVFEYF...	
ACT_FNTCT	5	346	E<NTQIMFETFN.....TPAFYYSIQAVLELYASGRT.....TGIV_D3GDGYTHWVFEYF...	
ACT_VELCR	5	346	E<NTQIMFETFN.....TPAFYYSIQAVLELYASGRT.....TGIV_D3GDGYTHWVFEYF...	
ACT4_CACIL	0	347	E<NTQIMFETFN.....TPAFYYSIQAVLELYASGRT.....TGIV_D3GDGYTHWVFEYF...	
ACTR_HMAN	5	346	E<NTQIMFETFN.....TPAFYYSIQAVLELYASGRT.....TGIV_D3GDGYTHWVFEYF...	
ACT5_C-1LK	6	347	E<NTQIMFETFN.....TPAFYYSIQAVLELYASGRT.....TGIV_D3GDGYTHWVFEYF...	
1qfa2	70	456	ADYVKLLH.KAKKKRGZYDANIYAVVNDVGTMTDGYD...DGHCEVG_IIGTG...THACYMELRIDLY	
HXK1_HMAN	78	456	ADYVKLLH.KAKKKRGZYDANIYAVVNDVGTMTDGYD...DGHCEVG_IIGTG...THACYMELRIDLY	
HXK1_BVIN	78	456	ADYVKLLH.KAKKKRGZYDANIYAVVNDVGTMTDGYD...DGHCEVG_IIGTG...THACYMELRIDLY	
HXK1_FCMA	68	443	HAYAEIIQ.TEIDKRF...YKCYAVVNDVGTMTDGYD...DGHCEVG_IIGTG...THACYMELRIDLY	
HXK2_DROME	128	505	KAVSLLQ.EADDRAGL<INTVAILNDVGTMTDGYD...DGHCEVG_IIGTG...THACYMELRIDLY	
HXK1_GIDL	95	405	CVVAALT.KAMLRKG.VDNVYALVNDVGTMTDGYD...DGHCEVG_IIGTG...THACYMELRIDLY	



FUNCTION PREDICTION PROTOCOL

Based on sequence similarity, structural analyses and information about interacting partners.

Protein
primary
sequence

Primary Database similarity search

- SwissProt / UniProt
- nr / SP+SPTreMBL
- COG / KOG
- PDB

*Orthologs / paralogs
MSA*

*Family assignment
Functional residues
Phylogenetic profile
Gene neighbourhood
Function prediction?*

*Function prediction
(cellular level)?*

Protein interactions
characterization

Secondary Database similarity search

- Prosite
- Pfam
- SMART
- PRINTS
- BLOCKS
- InterPro

*Protein motifs
Domain organization
Family assignment
Function prediction?*

Protein structure analyses

- SCOP / CATH classification
- Functional sites mapped on structure

Protein structure prediction

- 1D features
- 3D structure / fold prediction

Known / Predicted structure

*Function prediction
(molecular level)?*



Why Function Prediction is Important



The huge quantity of sequenced deposited in the protein databanks by genome sequencing projects will be dwarfed by the sequences from the environmental sequencing projects that are currently underway.

There is a growing **imbalance** between the number of sequences in databases and the information about structure and function regarding those sequences.

Functional and structural **predictions** could save a lot of time and effort and money, and they provide information about systems that are not experimentally treatable.

Why Pattern Recognition Methods Work. Mostly.

Sequence



The amino acid sequence of a protein **determines** its structure – in general similar sequences fold into similar structures.

Chaperonins

Structure



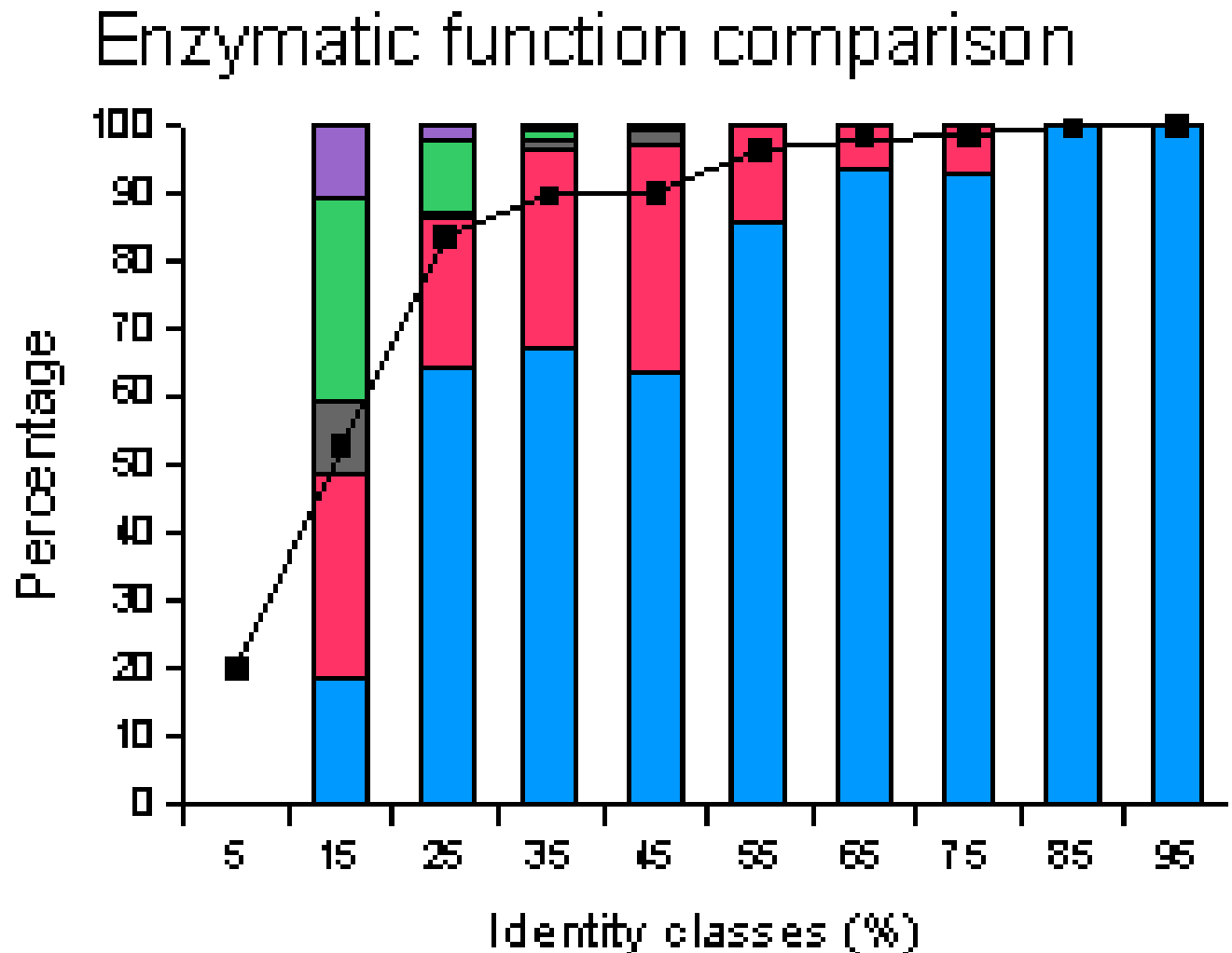
The structure of a protein **determines** its function. And in general similar protein structures perform similar functions.

Location

Function

Therefore, function prediction should be a simple matter of detecting similarity.

However, the relation between sequence similarity and function is far from simple.



Two Conclusions Drawn from Comparisons of Protein Sequences

First, protein sequences can be grouped in clusters or families based on sequence similarity.

Second, this similarity may be restricted to short stretches of sequence. These conserved regions of sequence are generally regions of functional importance.



Clusters

Clustering techniques aim to identify groups of sequences that are maximally similar between them, and minimally similar to sequences in other clusters.

The terms **superfamily**, **family** and **subfamily** are often used. All the members of a superfamily, family or subfamily, are evolutionarily related.

The members of a superfamily generally have a similar structure, even though it may not be possible to detect sequence similarity. The members of a subfamily have clear sequence similarity and the same function.



1. Patterns and Profiles

Constructed from multiple alignments of evolutionarily related sequences. They include consensus sequences, regular expressions, profiles (PSSMs) and hidden Markov models (HMMs).

2. Conserved sub-sequences

Domains are stretches of sequence that appear as **modules** within proteins. They are usually more obvious structurally. Certain domains can be found in a wide range of proteins. Domain shuffling is an accepted mechanism of protein evolution.

Motifs are short, conserved sub-sequences that usually correspond to active or functional sites.

Motifs

Proteins from the same family are often characterised by short conserved regions that are usually related with function, such as binding sites or active sites.

The structural and functional restrictions placed on these motifs means that they are usually conserved even over large evolutionary distances.

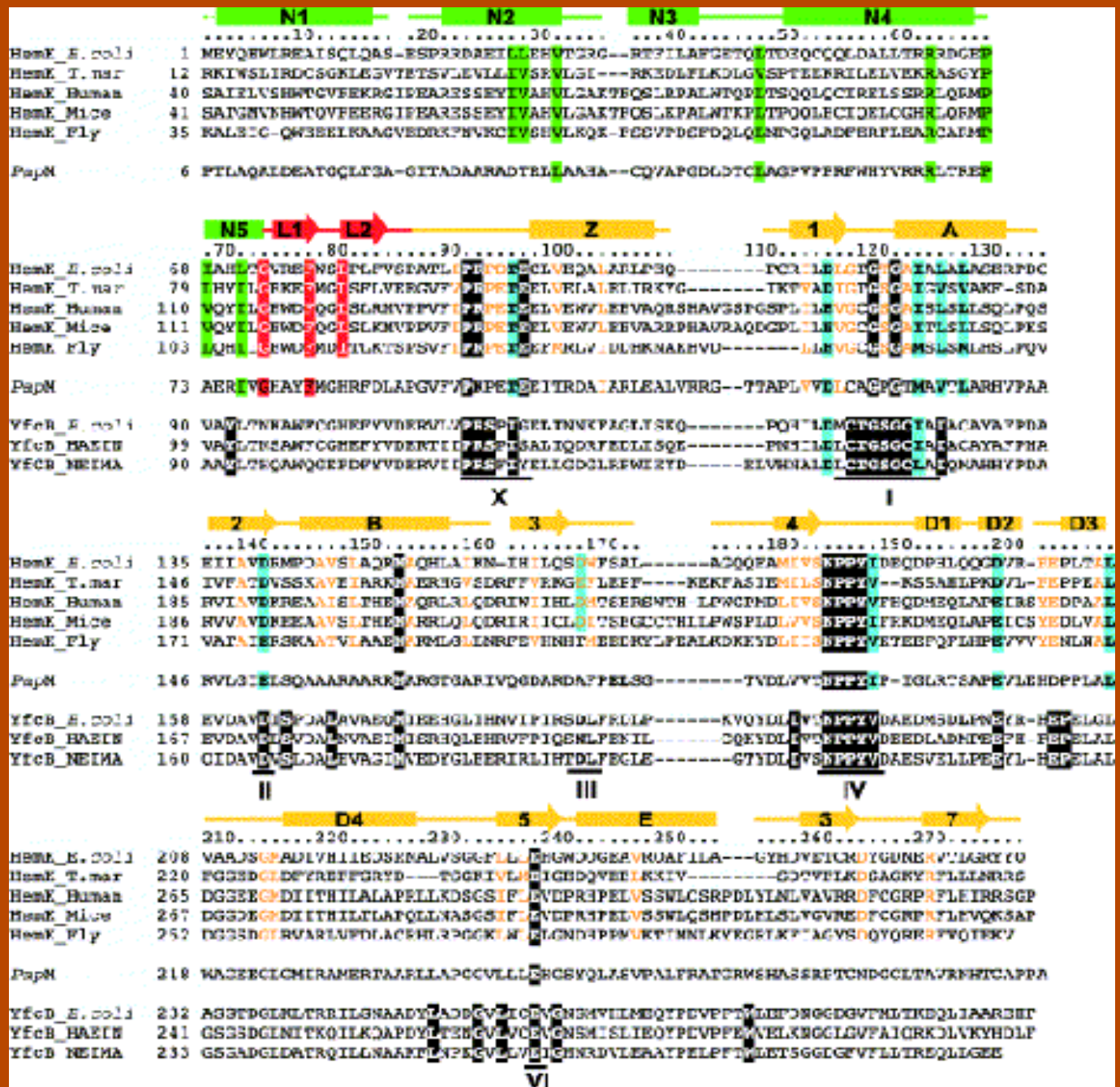
However, they are not usually detectable by sequence search techniques such as BLAST or FASTA.

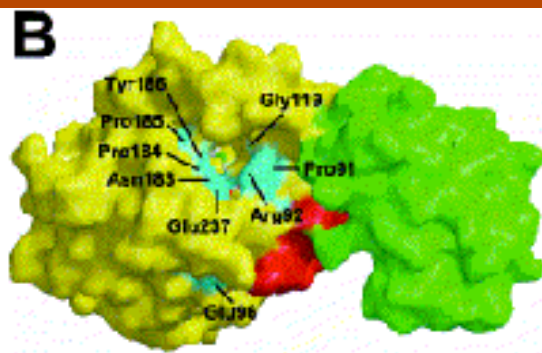
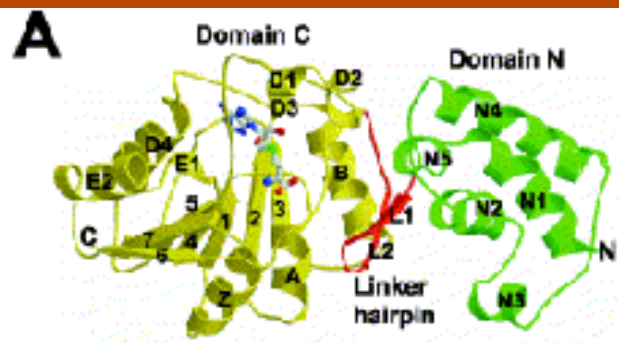
Instead there are motif databases, and tools for searching for motifs in sequences and for searching for sequences with a certain motif.

HemK

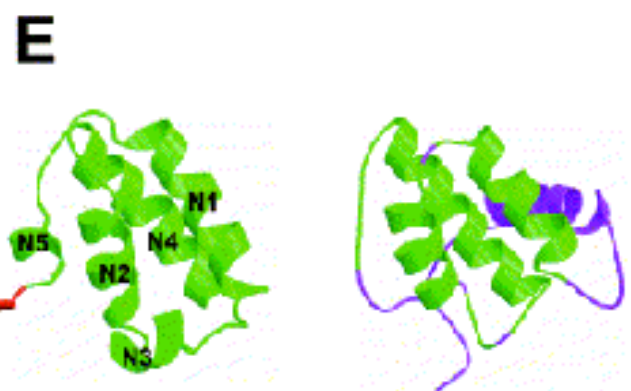
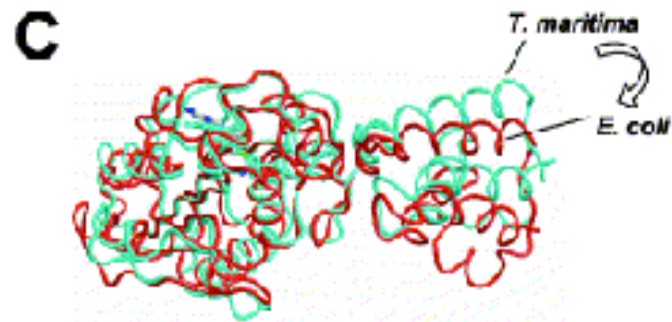
- an adenine methyltransferase?

The “NPPY” motif is characteristic of the DNA methylase family.



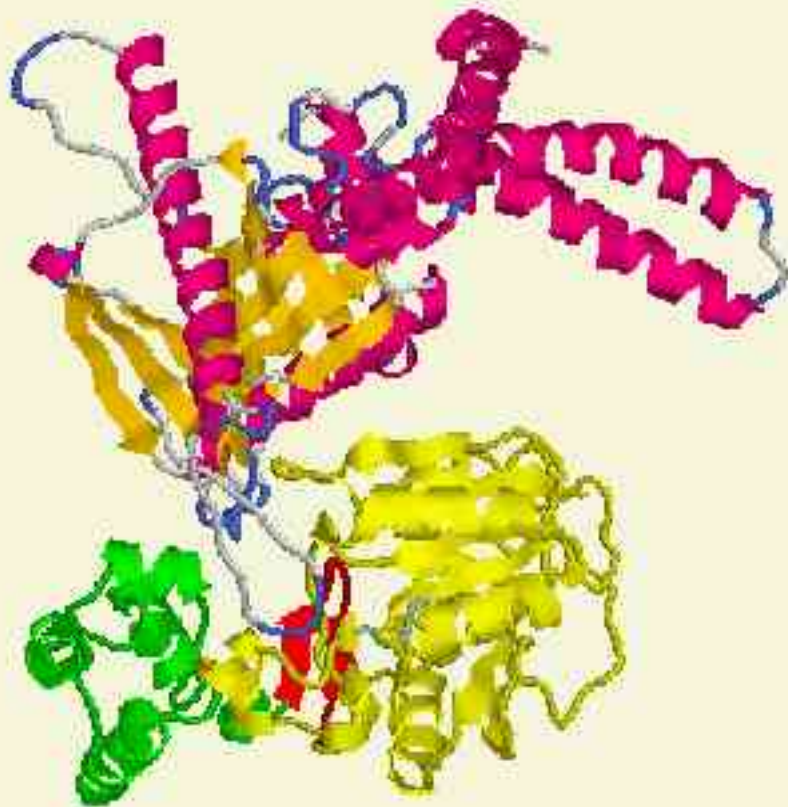


However HemK turns out to be the first known glutamine methyltransferase



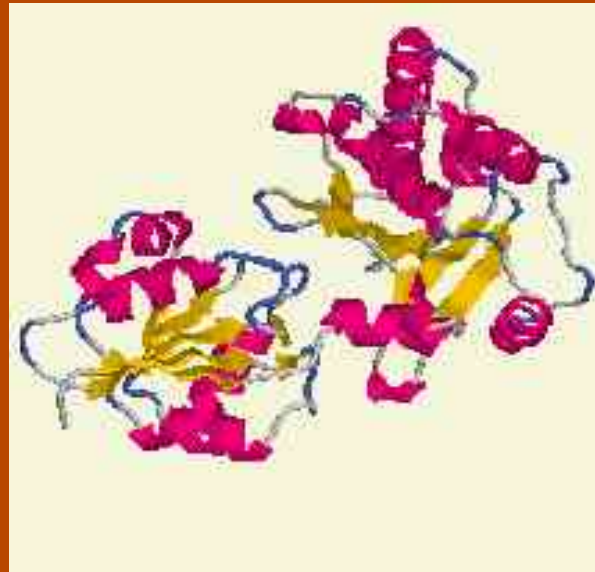
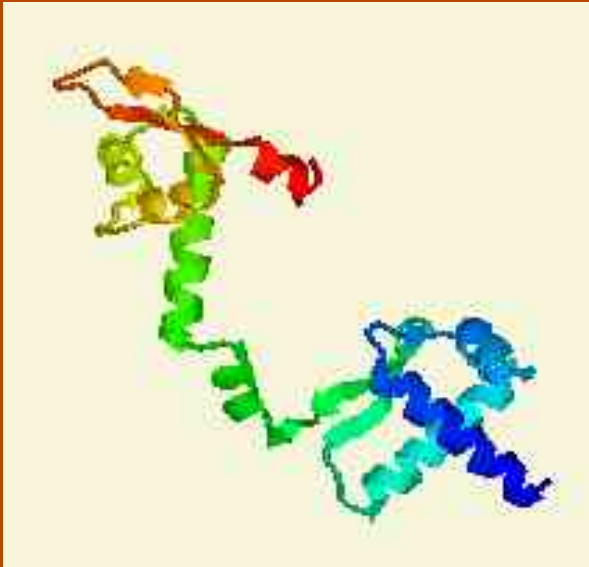
N-terminal domain of *E. coli* HemK

EPS15 homology domain



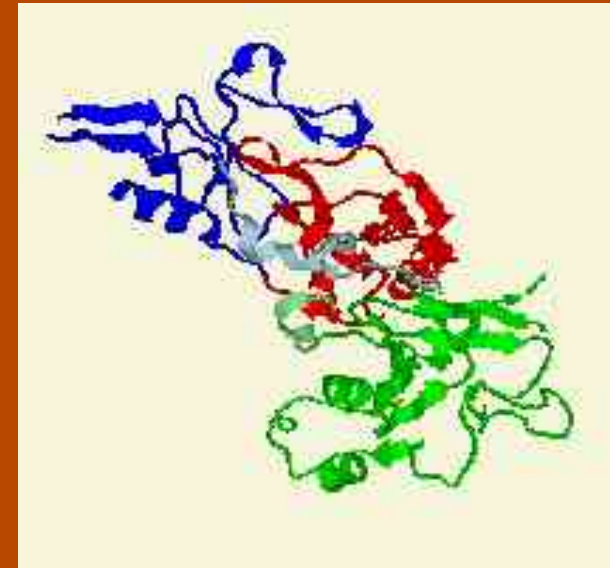
Domains

Domains are usually visible structurally, sometimes they are more obvious than others though ...



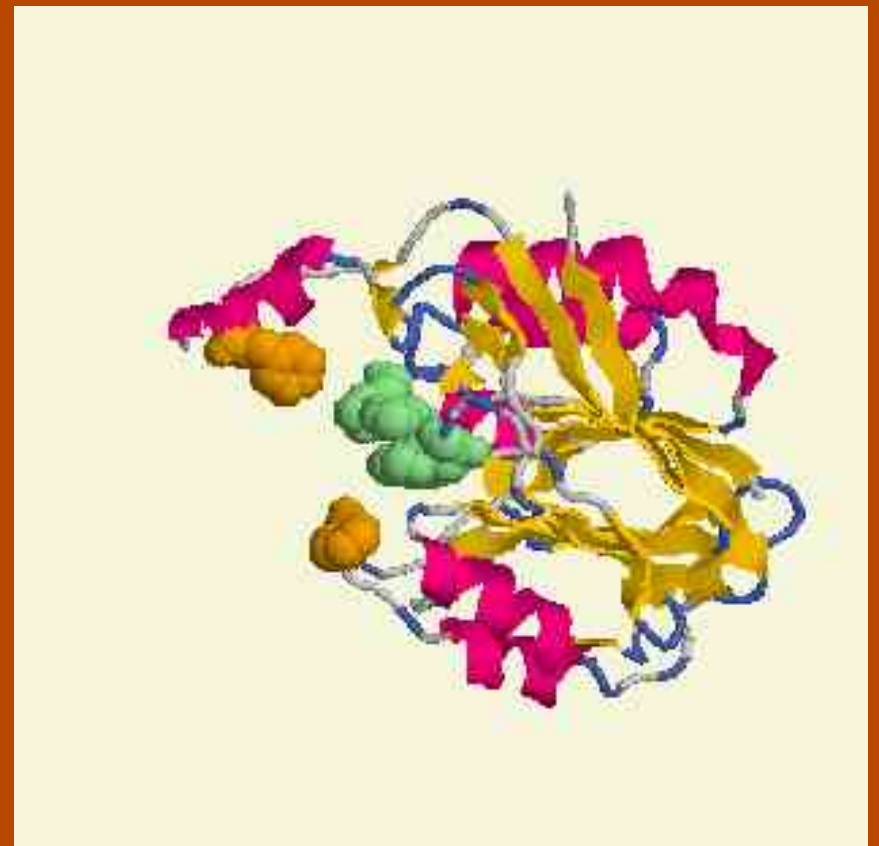
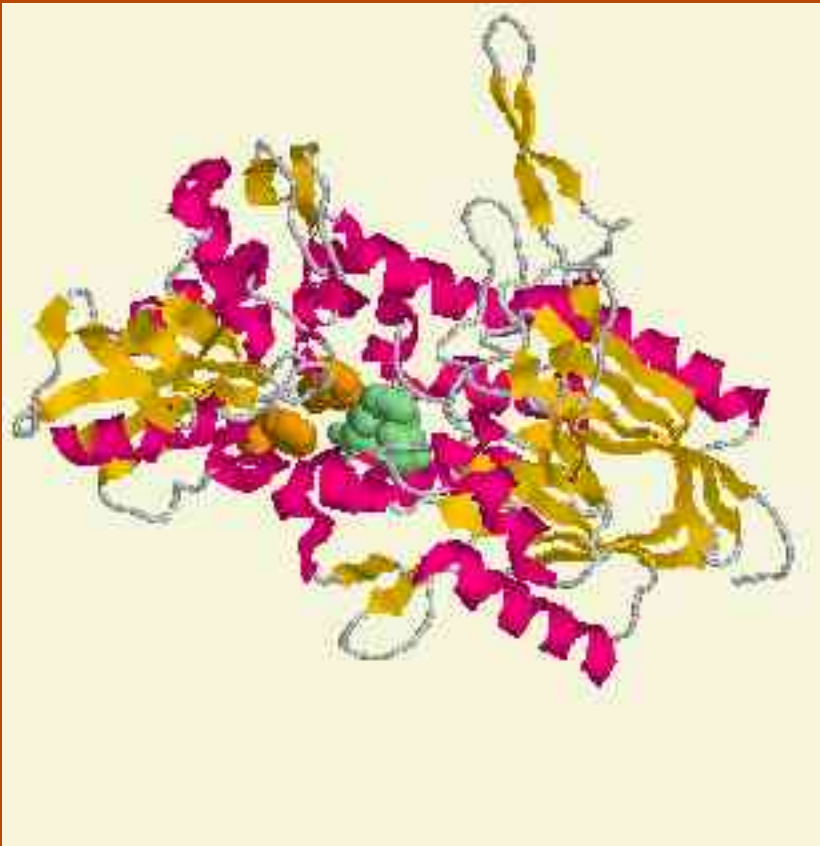
Diphtine synthase

Heat shock operon
repressor HrcA

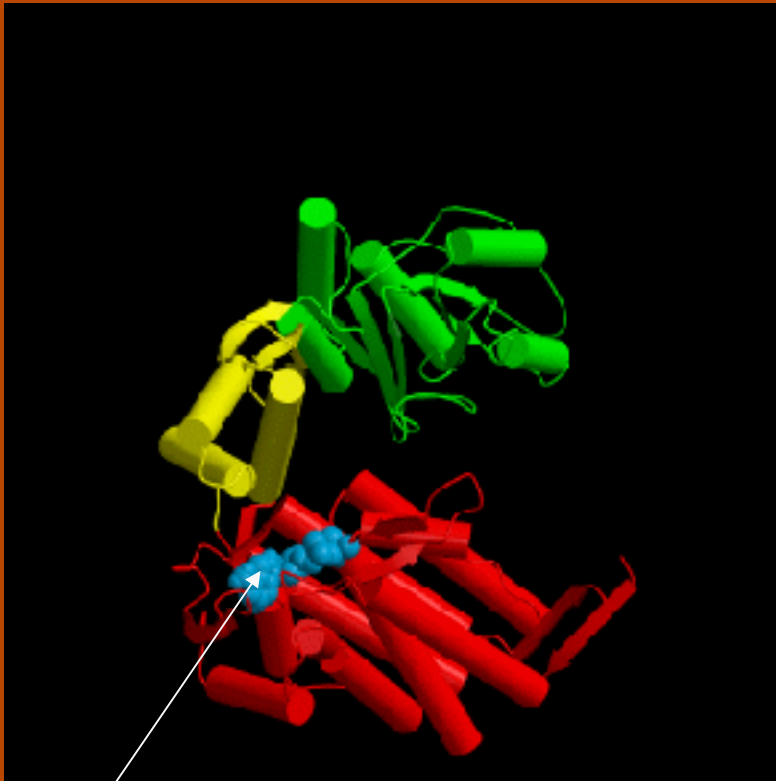


Domains

Domains often are associated with a single function, but sometimes both domains are involved – eq 1hpuD and T0200

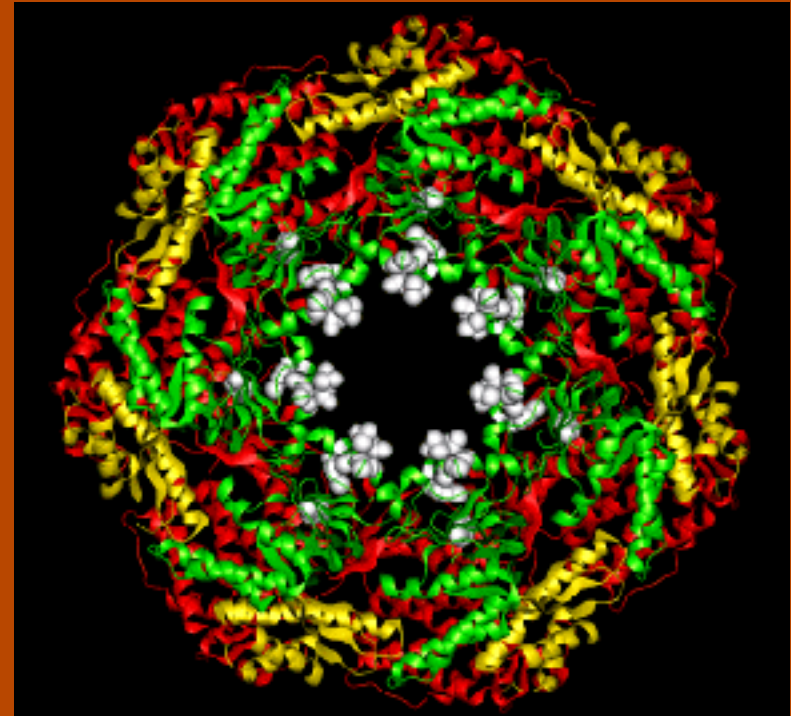


Domains and Chains



ATP

subunit



heptamer

Molecular chaperonin GroEL

(Dr Jianpeng Ma, Harvard Univ.)

Consensus Sequences

ALRDFATHDDF

SMTAEATHDSI

ECDQAATHEAS

80%

XXXXXATHXXX

50%

XXXXXATHDXX

Regular Expressions

ALRDFATHDDF

SMTAEATHDSI

ECDQAATHEAS

XXXXXATH[DE]xx

• N terminal: <, C-terminal: >

• Any amino acid: x

• Ambiguity: [A,B...] A or B...

• or {A,B..} anything but A or B...

• Repeats: A(2,4) - A-A o A-A-A o A-A-A-A

[AC]-x-V-x(4)-{E,D}

[Ala or Cys]-any-Val-any-any-any-any-{any but Glu or Asp}

PSI-BLAST

PSI-BLAST (Position Specific Iterated BLAST) is accessible at the NCBI web server. Like BLAST, PSI-BLAST performs database searches with query sequences as input.

After the first BLAST search a multiple sequence alignment is constructed from significant local alignments found in the first pass.

From that multiple alignment a **profile** (or PSSM) is built .

Profiles

```

F K L L S H C L L V
F K A F G Q T M F Q
Y P I V G Q E L L G
F P V V K E A I L K
F K V L A A V I A D
L E F I S E C I I Q
F K L L G N V L V C

```

A	-18	-10	-1	-8	8	-3	3	-10	-2	-8
C	-22	-33	-18	-18	-22	-26	22	-24	-19	-7
D	-35	0	-32	-33	-7	6	-17	-34	-31	0
E	-27	15	-25	-26	-9	23	-9	-24	-23	-1
F	60	-30	12	14	-26	-29	-15	4	12	-29
G	-30	-20	-28	-32	28	-14	-23	-33	-27	-5
H	-13	-12	-25	-25	-16	14	-22	-22	-23	-10
I	3	-27	21	25	-29	-23	-8	33	19	-23
K	-26	25	-25	-27	-6	4	-15	-27	-26	0
L	14	-28	19	27	-27	-20	-9	33	26	-21
M	3	-15	10	14	-17	-10	-9	25	12	-11
N	-22	-6	-24	-27	1	8	-15	-24	-24	-4
P	-30	24	-26	-28	-14	-10	-22	-24	-26	-18
Q	-32	5	-25	-26	-9	24	-16	-17	-23	7
R	-18	9	-22	-22	-10	0	-18	-23	-22	-4
S	-22	-8	-16	-21	11	2	-1	-24	-19	-4
T	-10	-10	-6	-7	-5	-8	2	-10	-7	-11
V	0	-25	22	25	-19	-26	6	19	16	-16
W	9	-25	-18	-19	-25	-27	-34	-20	-17	-28
Y	34	-18	-1	1	-23	-12	-19	0	0	-18

Profiles are built based on the frequency of each amino acid at each position in the alignment. The relative frequency of each amino acid (or the no appearance) are weighted based on observed frequencies.

A is less probable than M because although it doesn't appear we know that M is similar to L, I, V y F.

PSI-BLAST II

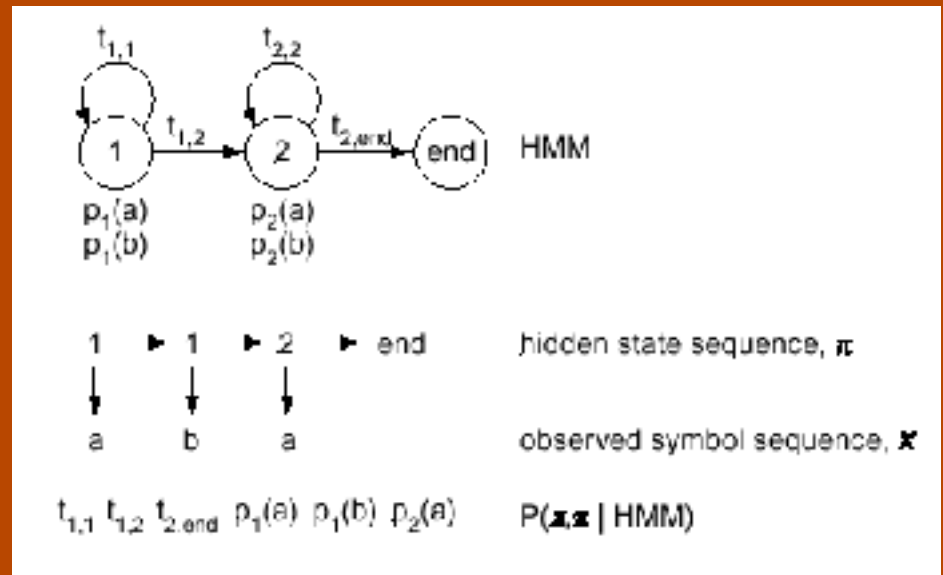
This profile is then used to search the database again, and any new significant hits are incorporated to the profile for further searches of the database.

The process iterates an arbitrary number of times or until no new sequences can be found in the database.

Hidden Markov Models I

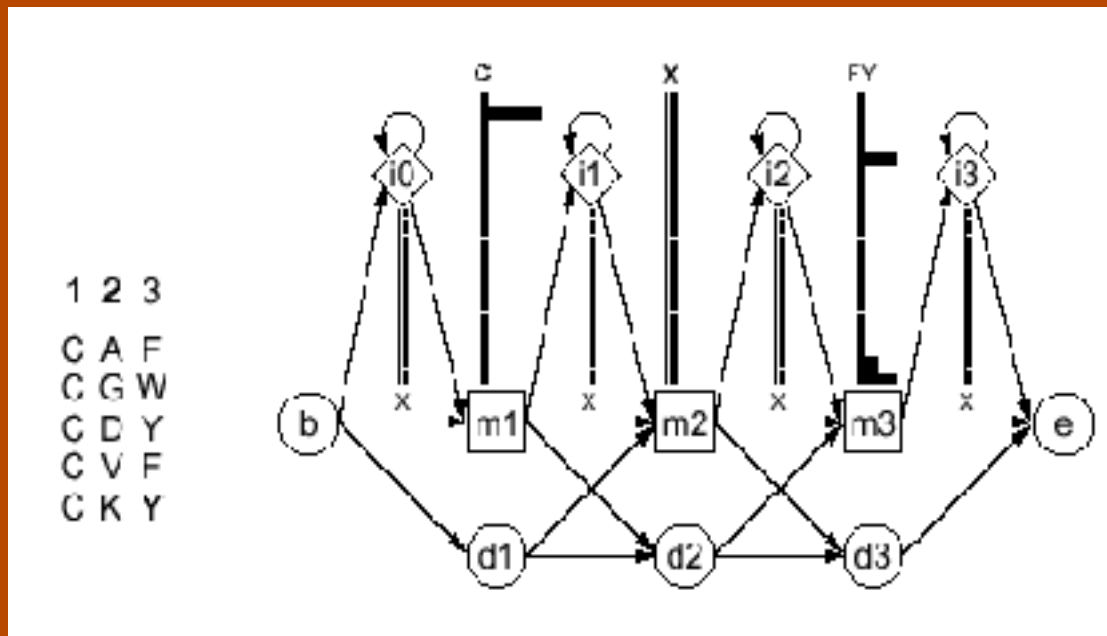
Hidden Markov models (HMMs) are statistical models of sequences.

They are similar to PSSMs, but the weights are calculated according to a probabilistic model that takes into account the previous positions.



Hidden Markov Models II

The model has three residues or states (m_1, m_2, m_3) with 20 probabilities to be a residue (bars), one state which is insertion (eg i_0) and one state which relates to deletion (eg d_1). The arrows represent the probability of the transition between each state.



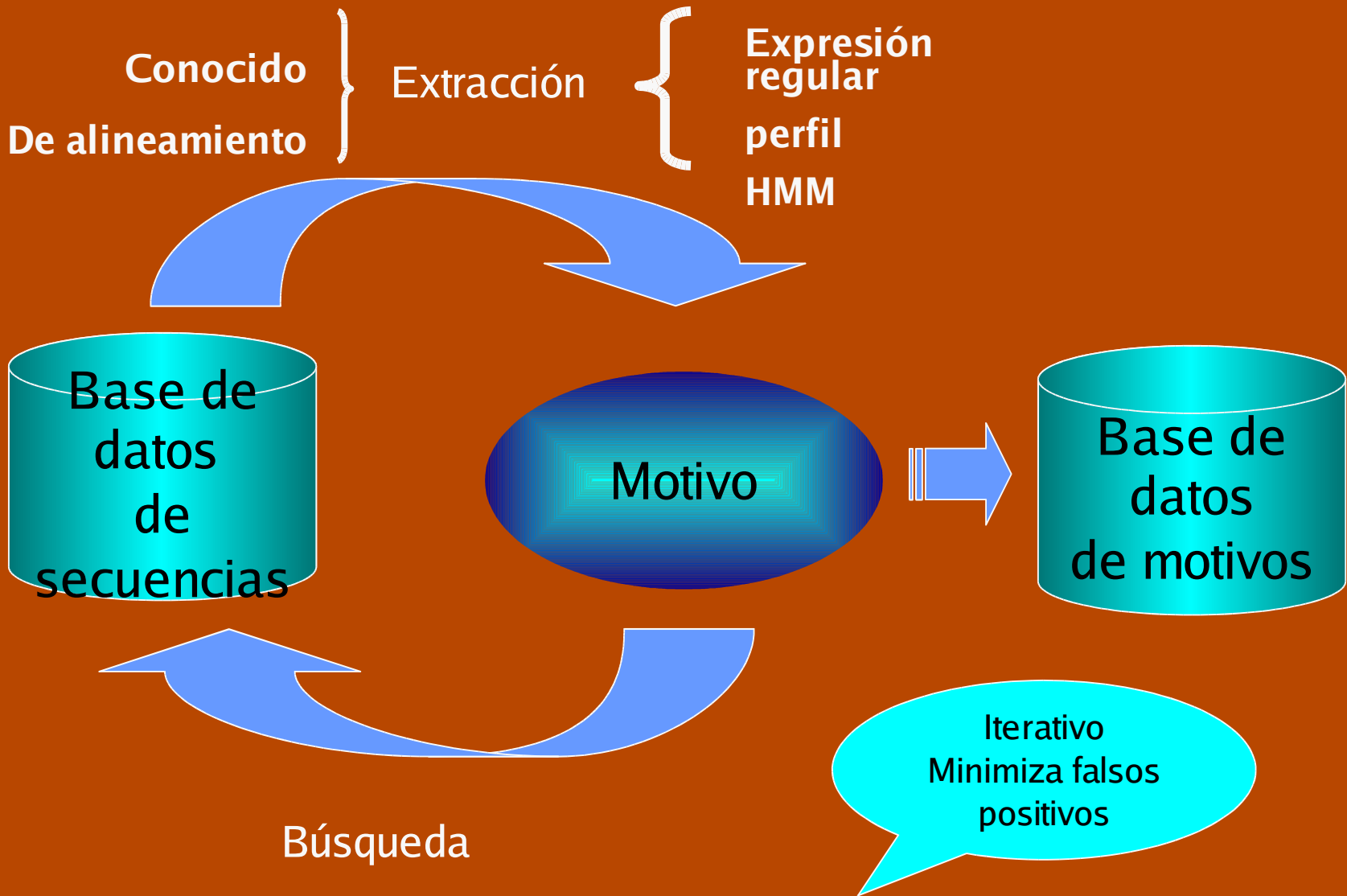
Motif Construction

Known motifs (published or found in known families)

Eg **PROSITE**

Empirical motifs: obtained directly from alignments. In some cases their significance is unknown.

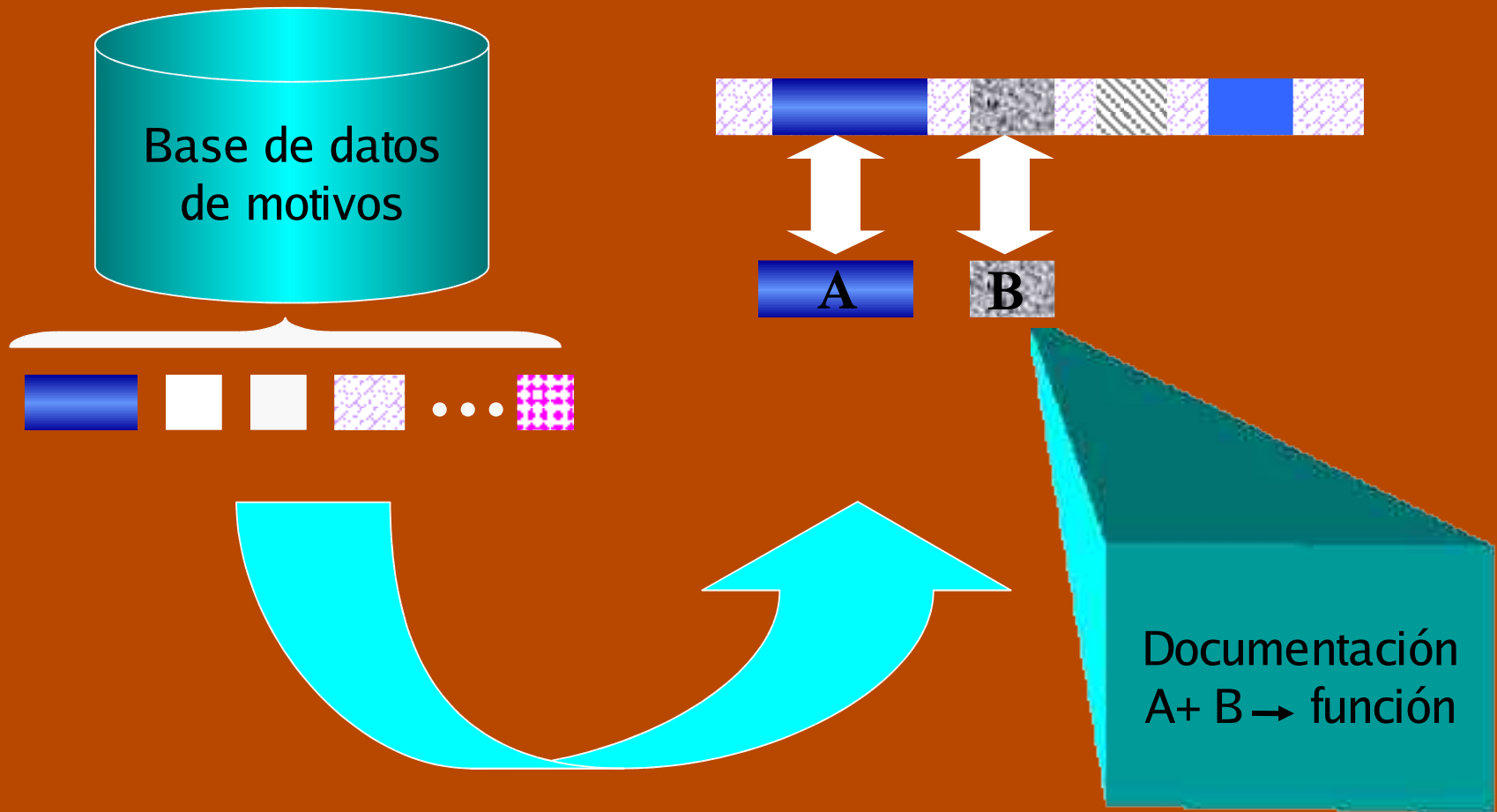
Eg **Pfam**.



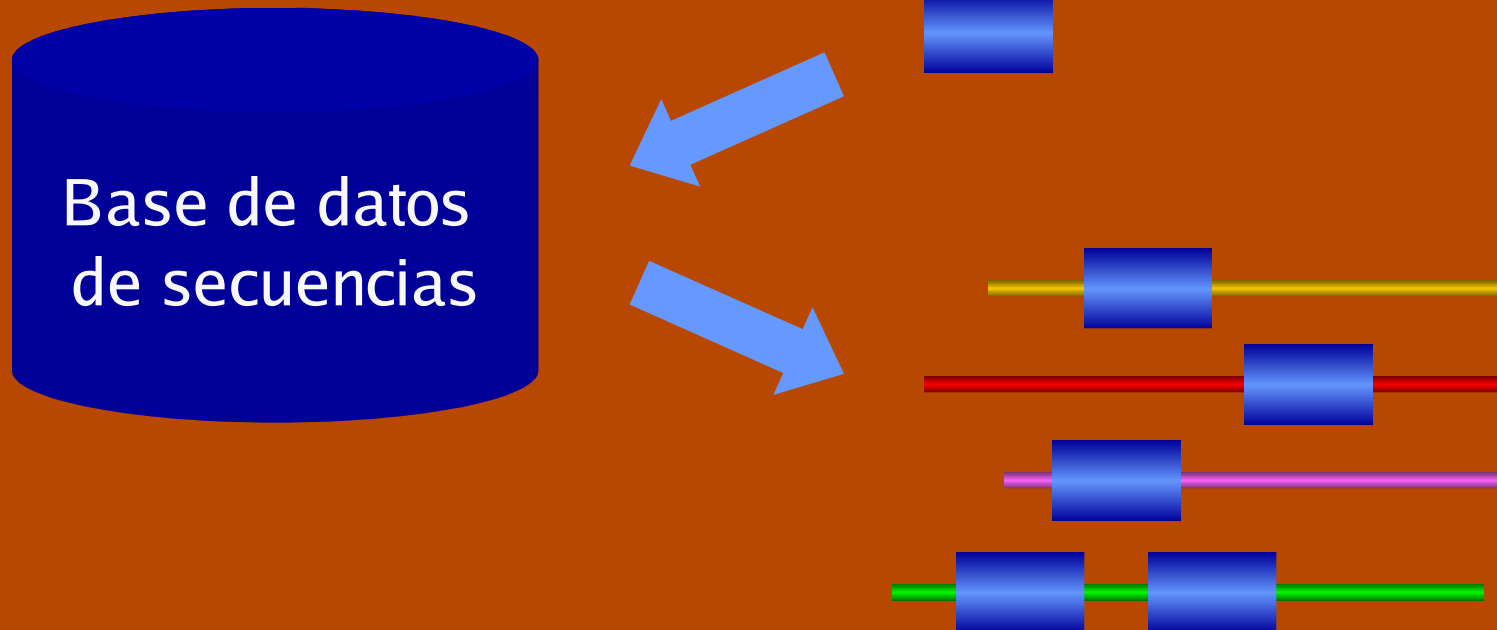
Using Databases and Tools Based on Motifs and Domains

1. **Function Prediction:** Identifying motifs in a query sequence
2. **Finding remote similarities:** Searching all sequences in a database with a certain motif.

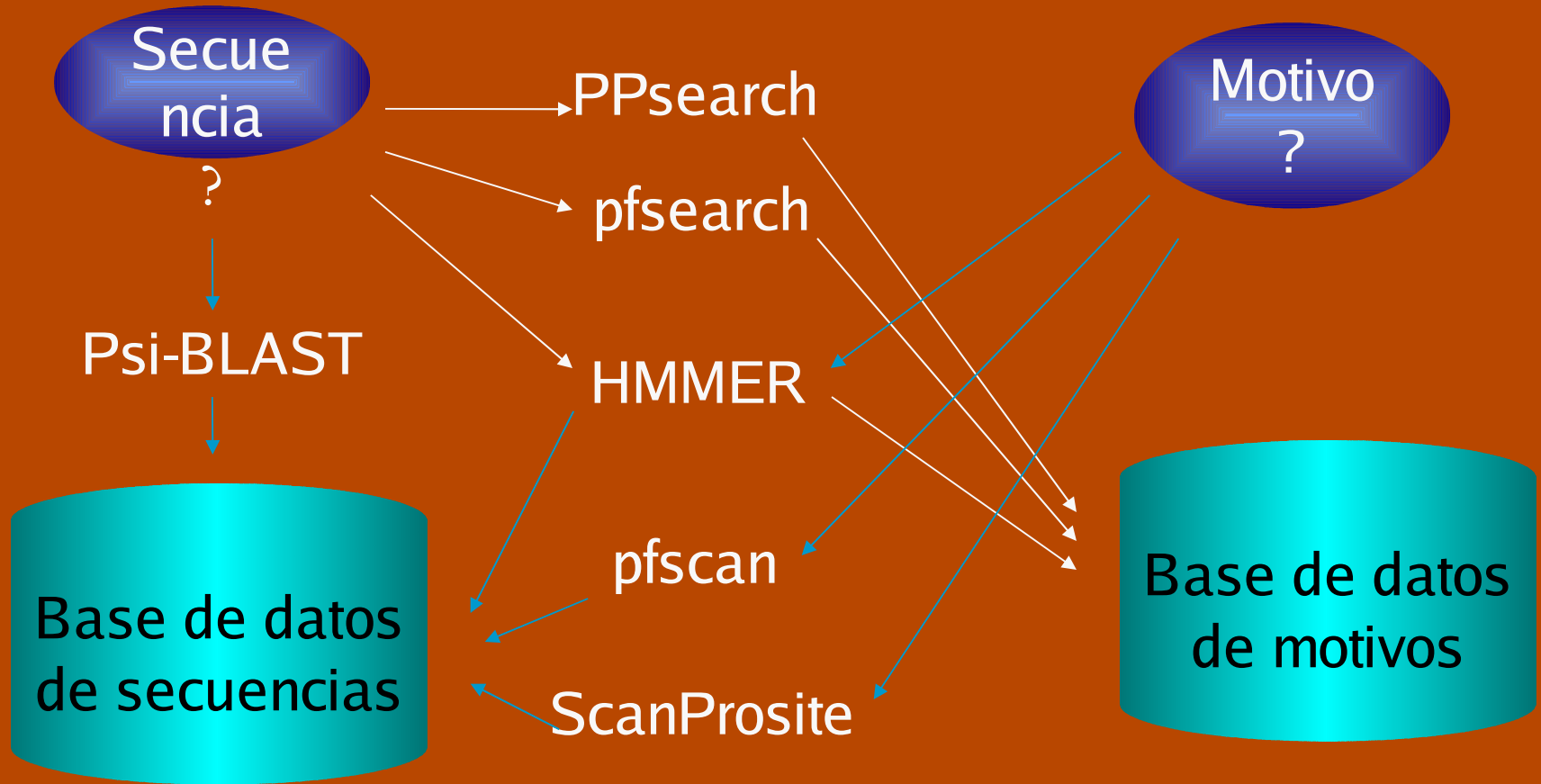
Function Prediction



Searching for Remote Similarity with Motifs



Tools



Databases

PROSITE Regular Expressions and profiles based on known motifs from SwissProt.

BLOCKS Profiles, based on PROSITE.

PRINTS Profiles, based on known motifs

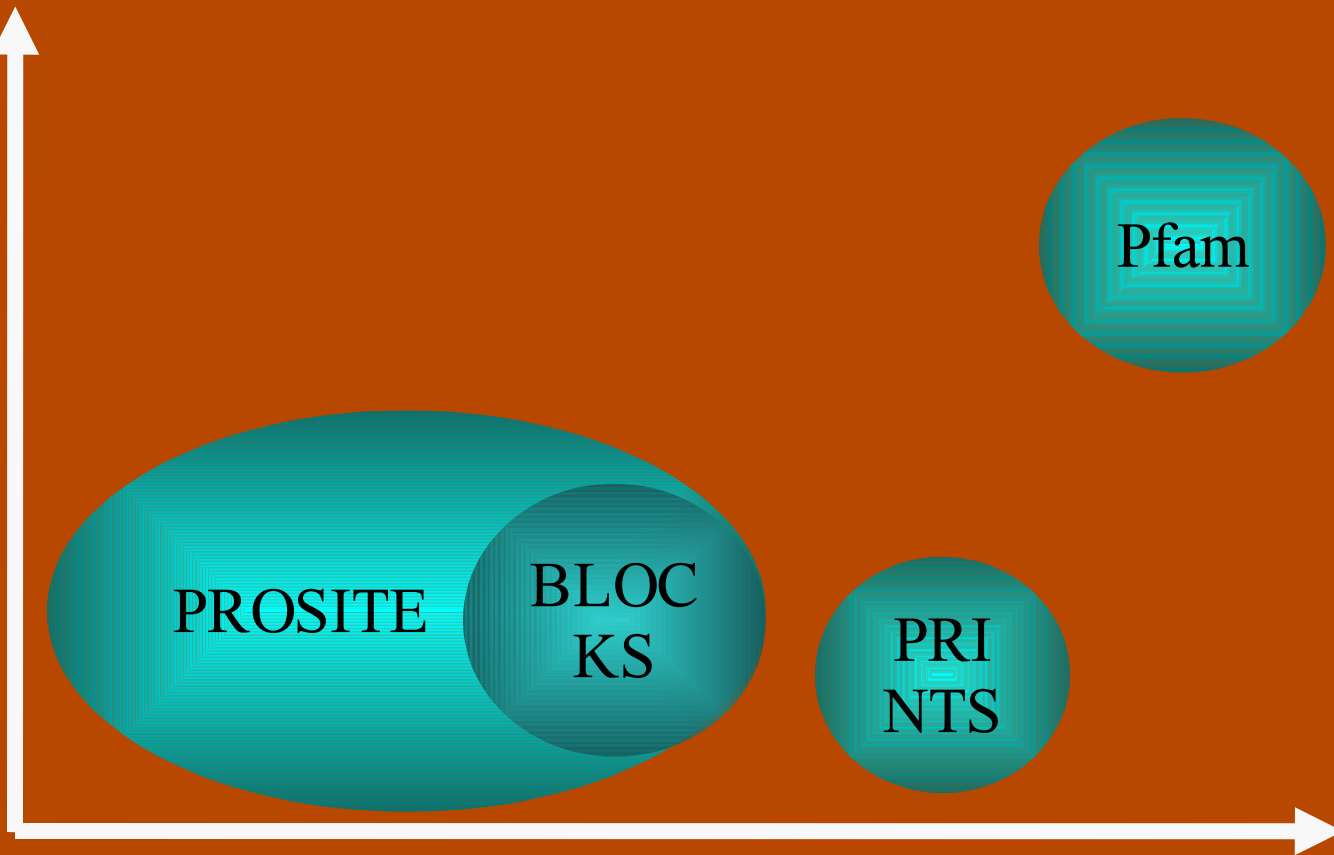
Pfam Profiles from HMM. Motifs generated automatically using SwissProt + SP-TrEMBL

SMART Domains identified with PSI-BLAST and HMMs.

InterPro Hierarchical database of domains and motifs that integrates information from Prosite, PRINTS, Pfam, Prodom, SMART and TIGRFAMs.

Información

SP-TRMBL
SwissProt



Expresiones
regulares

Perfiles
simples múltiples

HMMs

Precisión

This presentation contains material from:

Manuel J Gómez, CAB

Oswaldo Trelles, UMA

Joaquín Dopazo, CNIO

Paulino Gómez Puertas, CBM