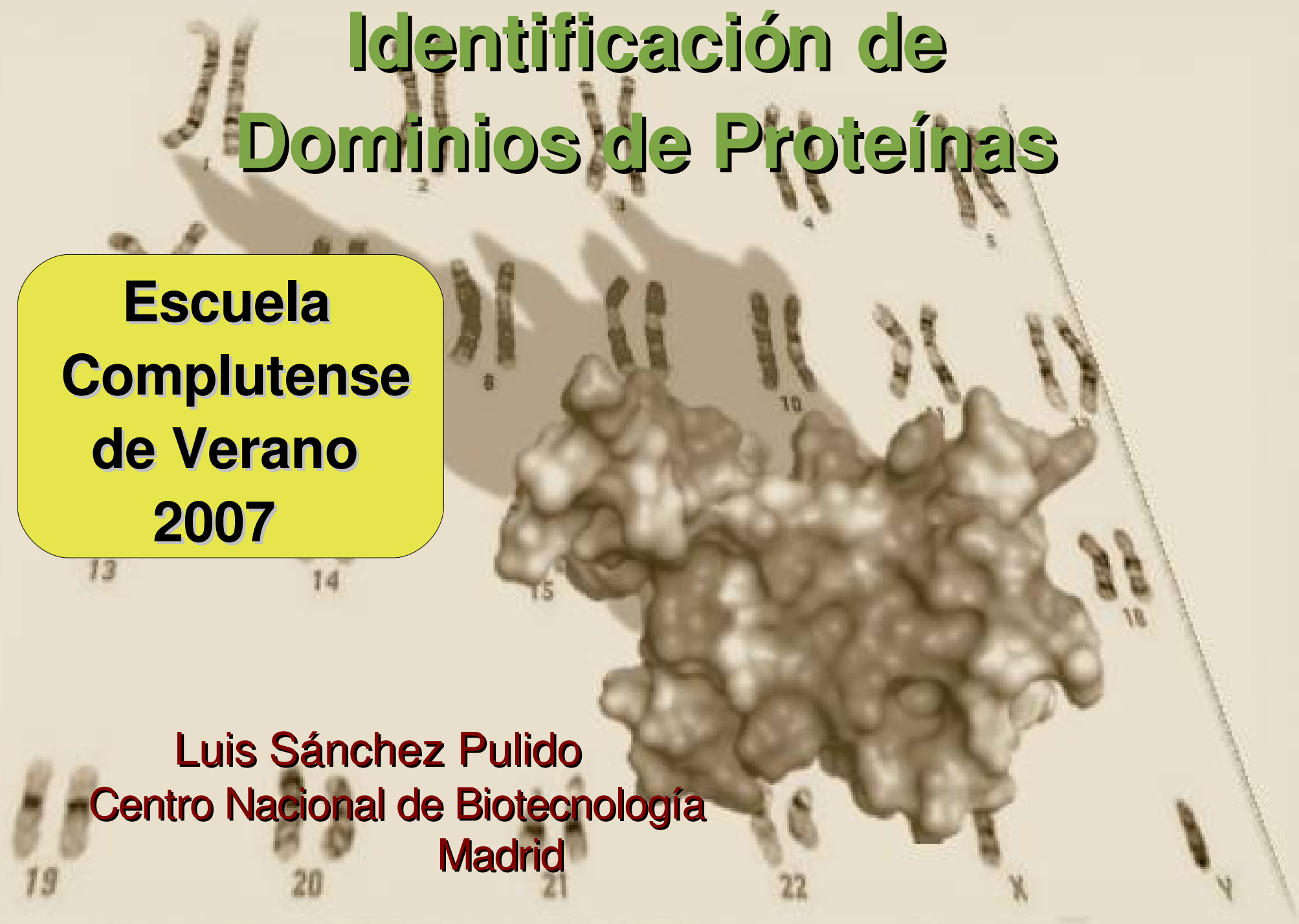




# Identificación de Dominios de Proteínas

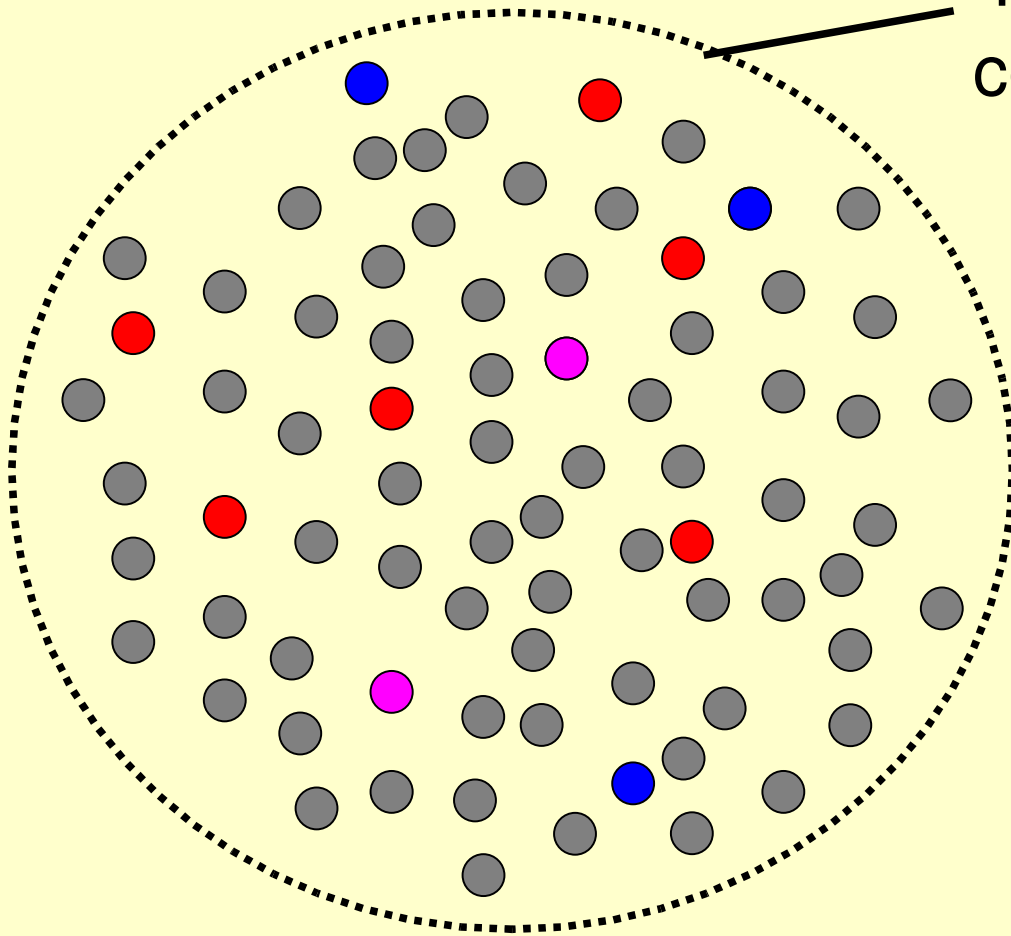
**Escuela  
Complutense  
de Verano  
2007**

**Luis Sánchez Pulido  
Centro Nacional de Biotecnología  
Madrid**

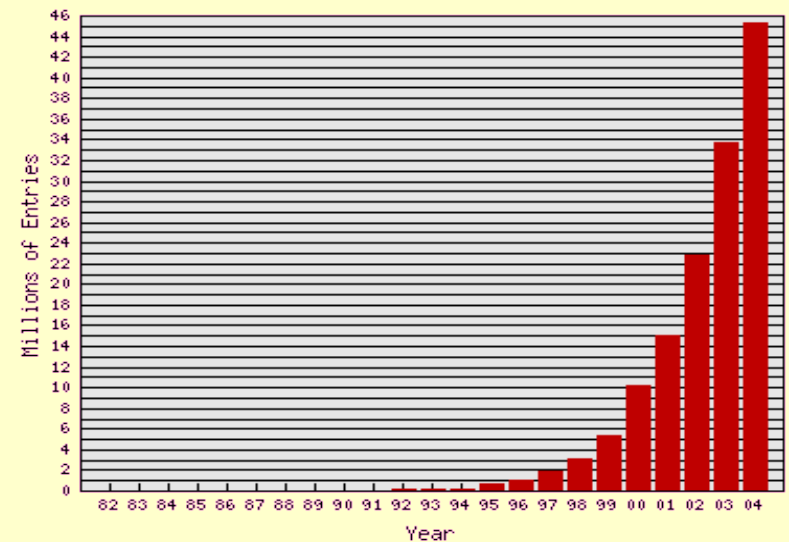


# ¿Por qué analizamos secuencias?

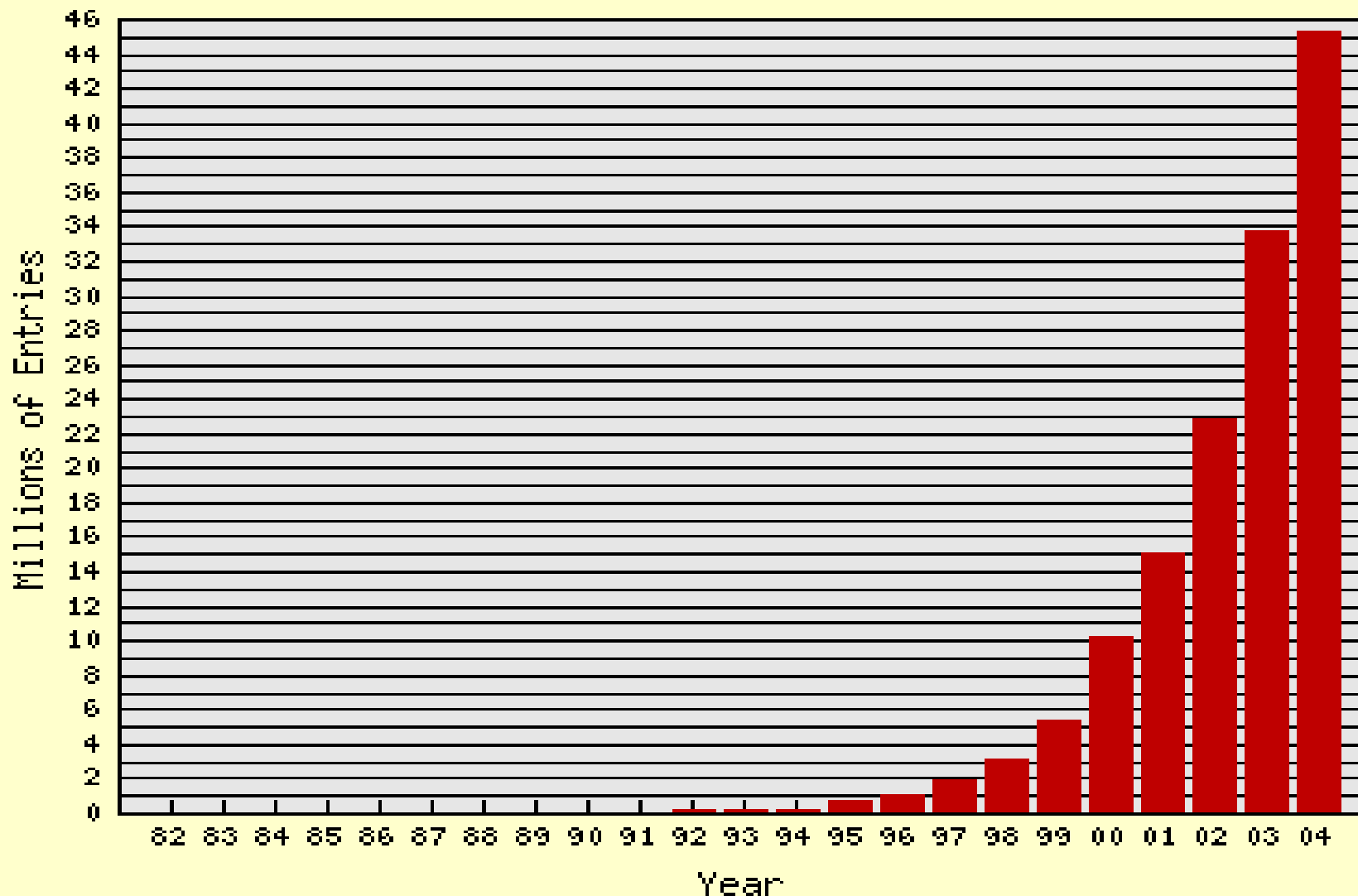
Proteínas de Secuencia conocida



- 3D
- Función
- Ambas
- ?????



# !!! ESTAMOS DESBORDADOS !!!



**Crecimiento de las bases de datos de secuencias.**

*Tomado de [www3.ebi.ac.uk/Services/DBStats/](http://www3.ebi.ac.uk/Services/DBStats/)*

***Gracias a la identificación de  
homología entre proteínas,  
podemos***

***TRANSFERIR INFORMACIÓN***

***Estructural y/o Funcional***

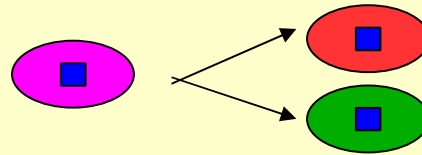
---

---

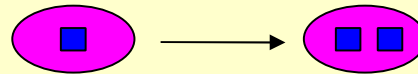
**Homólogos: par de proteínas con un ancestro común.**

**...y dependiendo del motivo de su divergencia:**

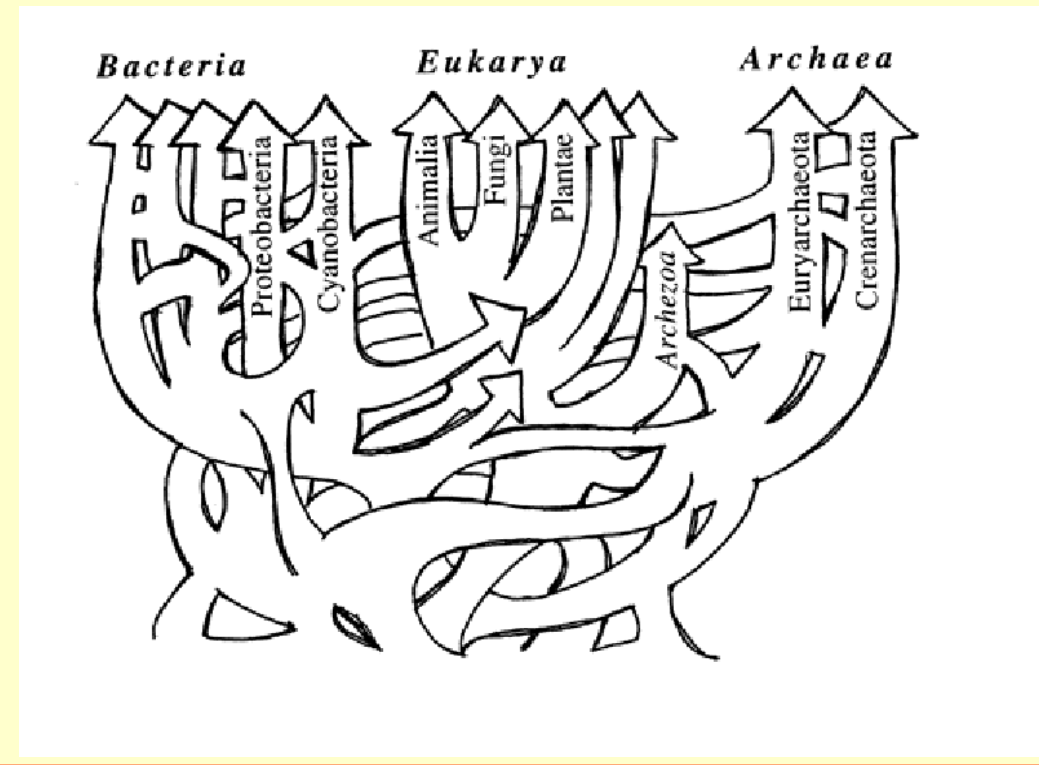
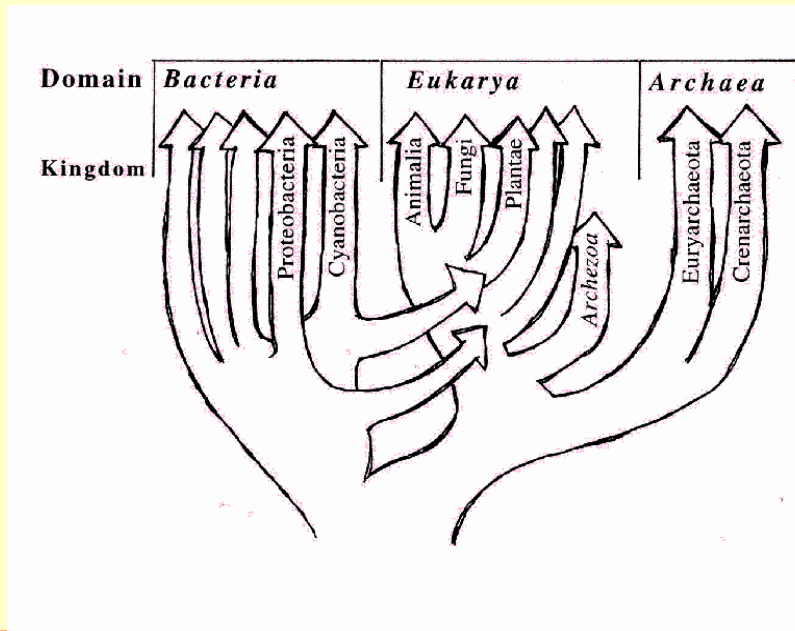
- **ortólogos** - especiación



- **parálogos** – duplicación génica



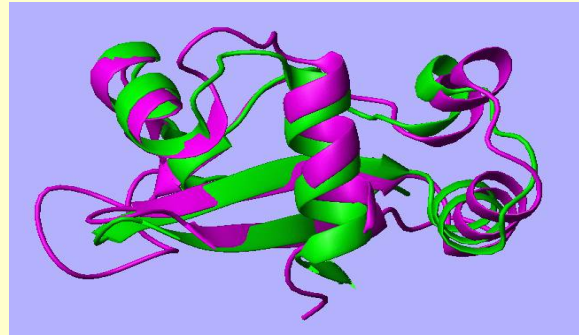
- **xenólogos** – transferencia horizontal



# ***TRANSFERIR INFORMACIÓN***

## **• Estructural**

*a partir de proteínas HOMÓLOGAS de estructura conocida por RayosX, RMN o ME*

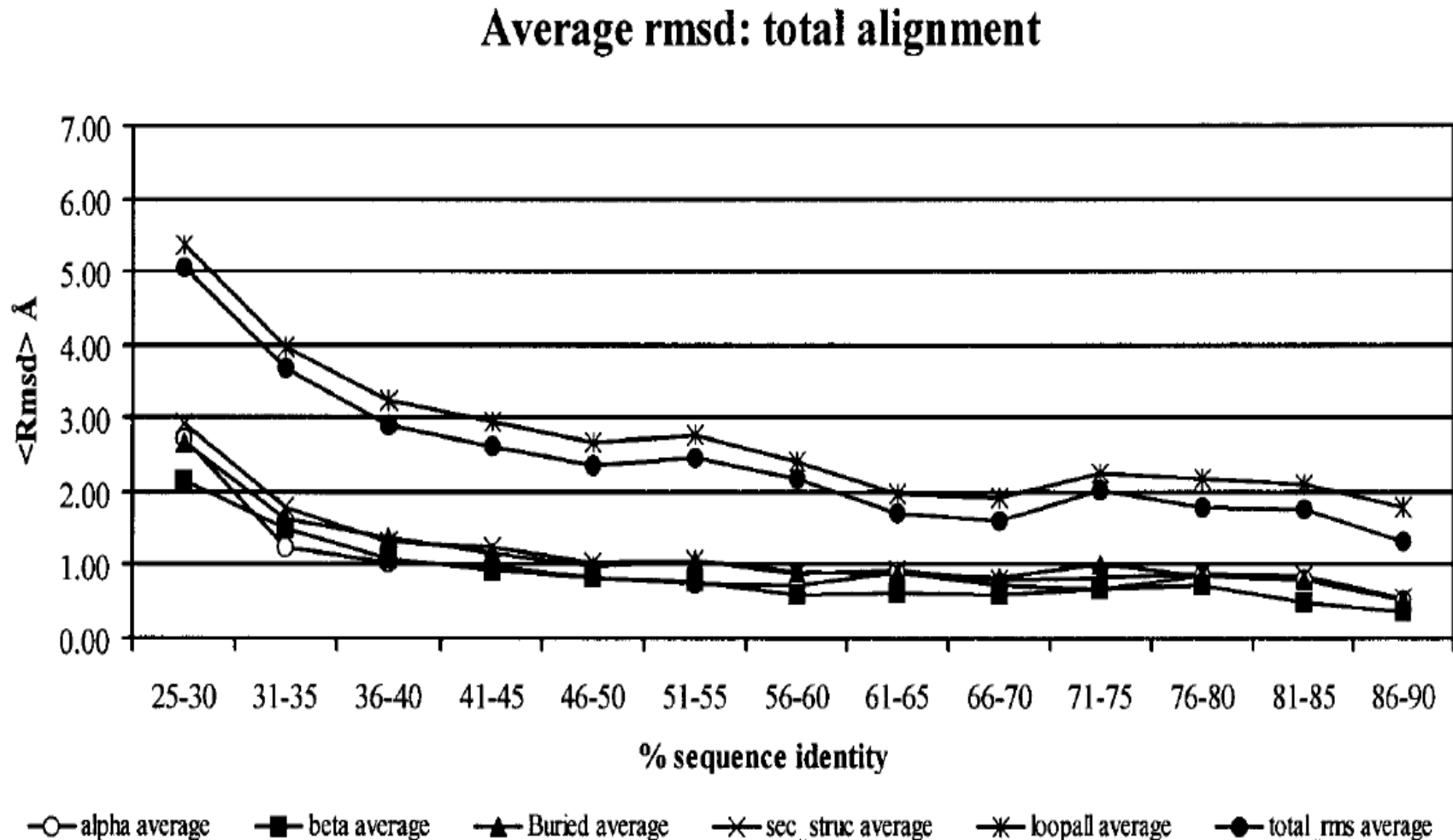


## **• Funcional**

*a partir de proteínas HOMÓLOGAS caracterizadas experimentalmente... y su contexto genómico y proteómico.*



# La estructura se conserva mejor que la secuencia!



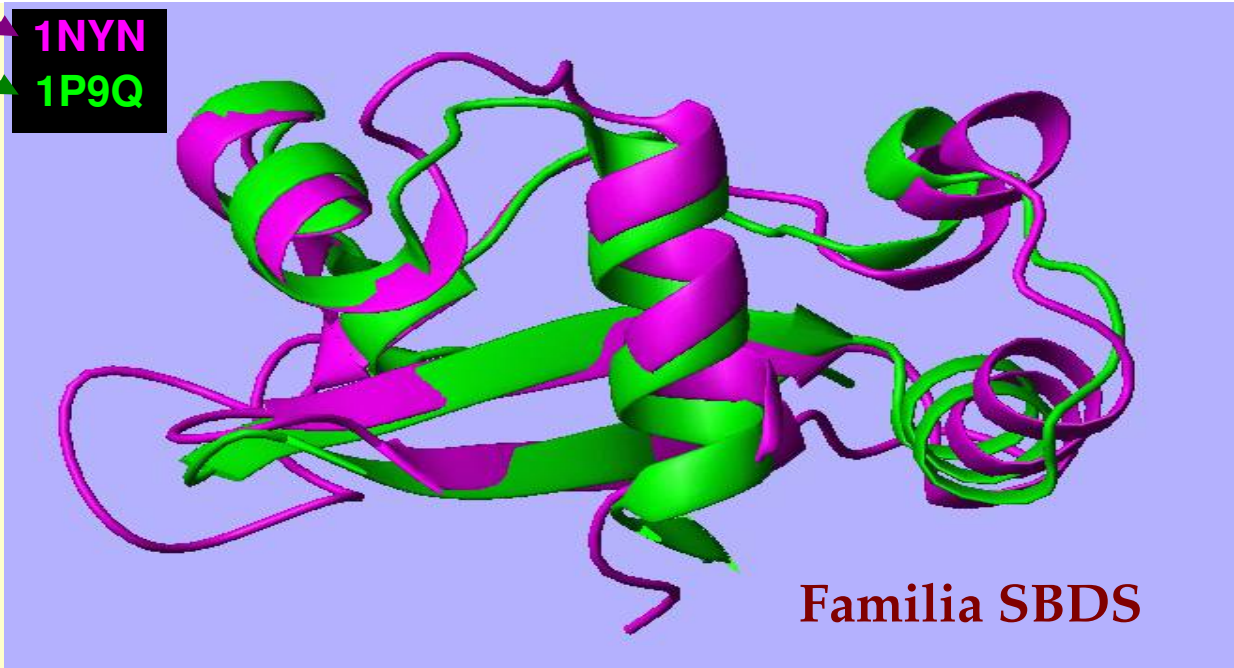
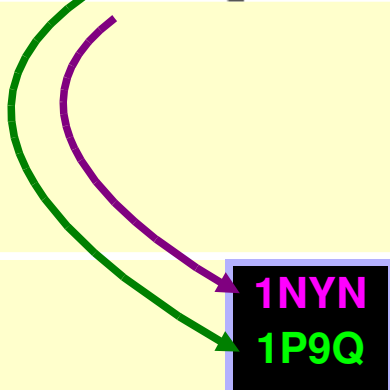
D'Alfonso G, Tramontano A, Lahm A.

Structural conservation in single-domain proteins: implications for homology modeling.

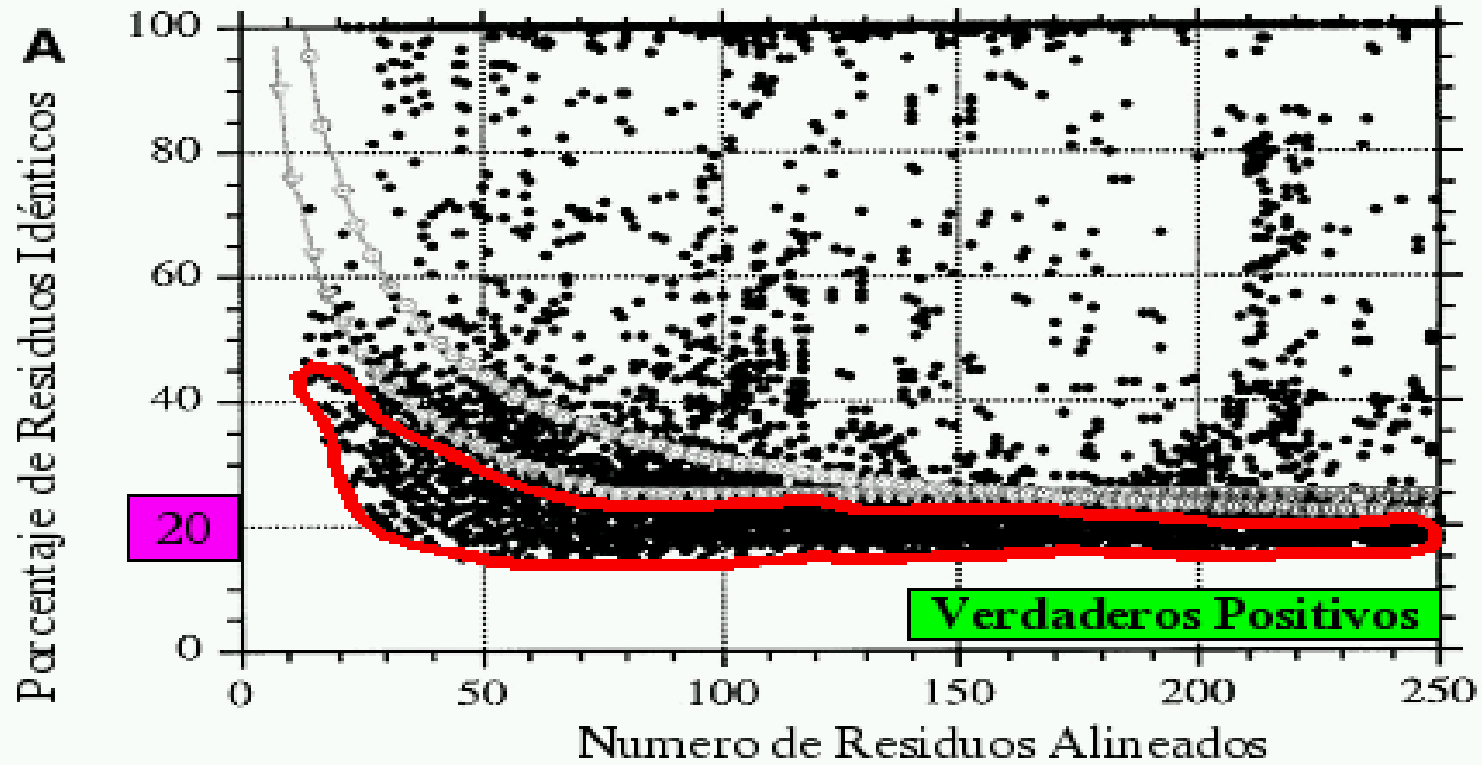
*J Struct Biol.* 134, 246-56. (2001)



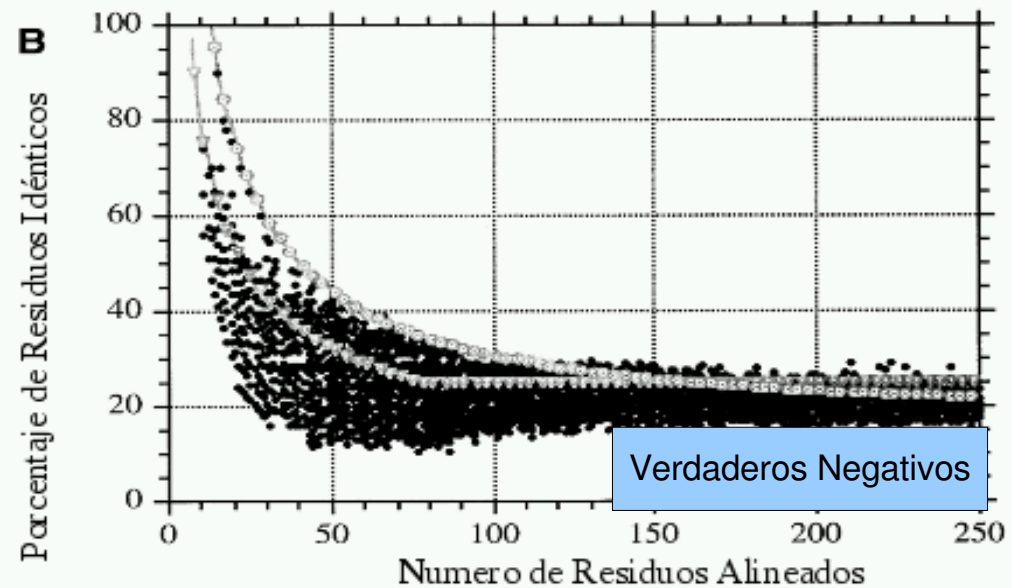
# A Remote Homology example:



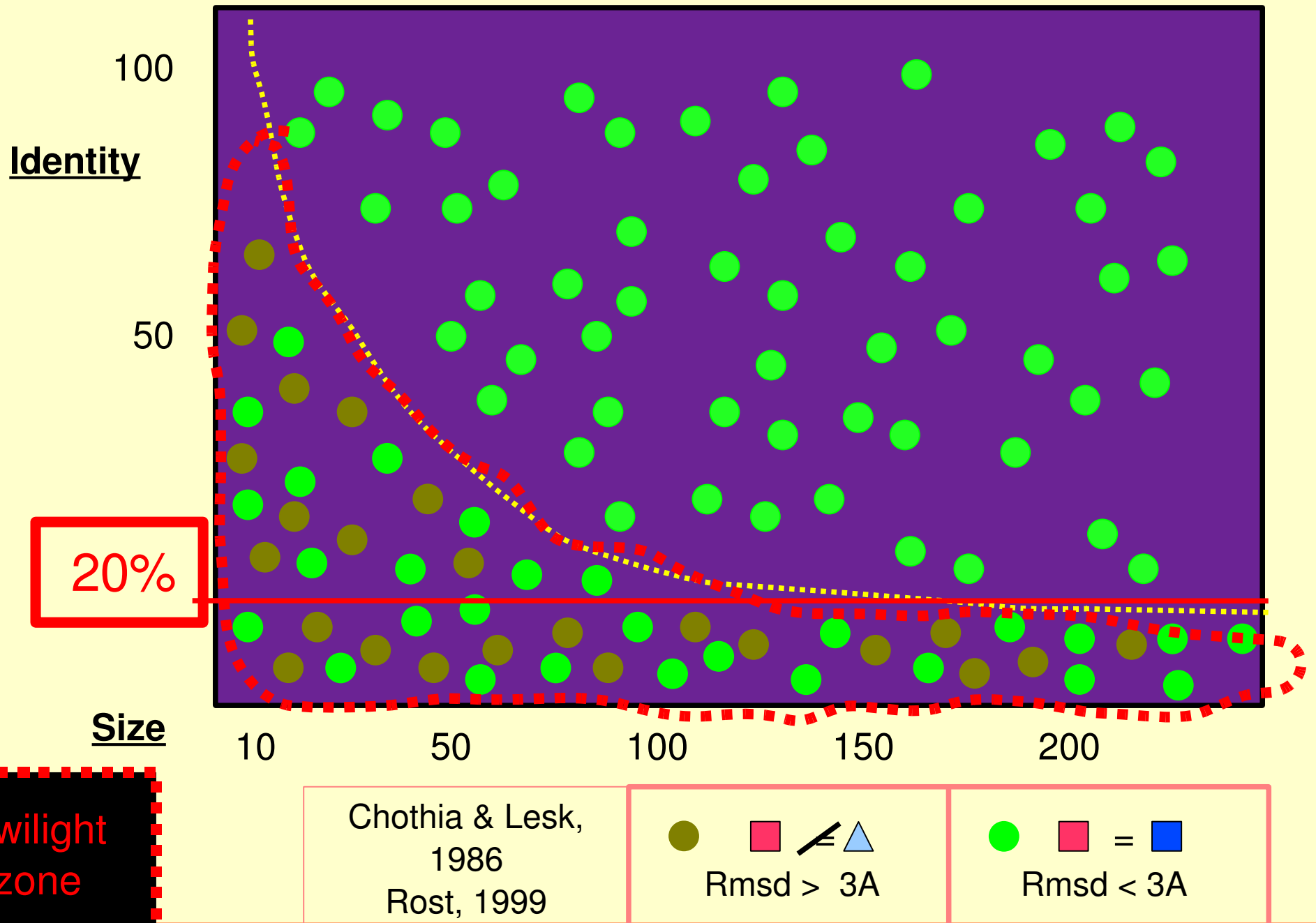
# Definiendo Homología Remota



Rost B. (1999)  
**Twilight zone of  
protein sequence alignments.**  
Protein Eng. 12:85-94.



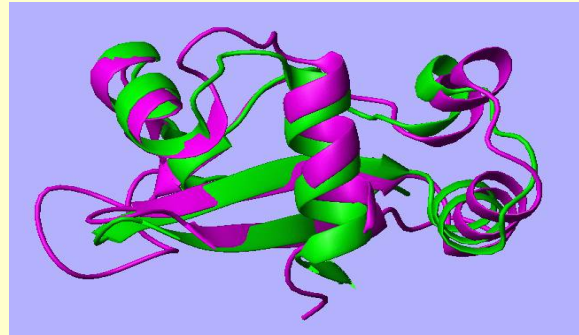
# Comparisons between pairs of sequences with known structure



# ***TRANSFERIR INFORMACIÓN***

## **• Estructural**

*a partir de proteínas HOMÓLOGAS de estructura conocida por RayosX, RMN o ME*



## **• Funcional**

*a partir de proteínas HOMÓLOGAS caracterizadas experimentalmente... y su contexto genómico y proteómico.*



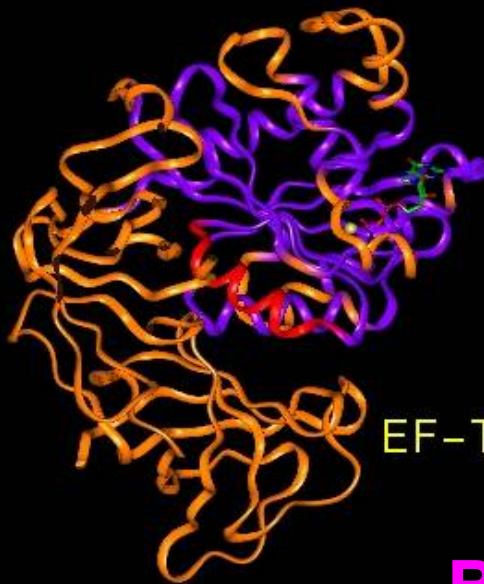
# ¿FUNCTION?



Transducin



RAS



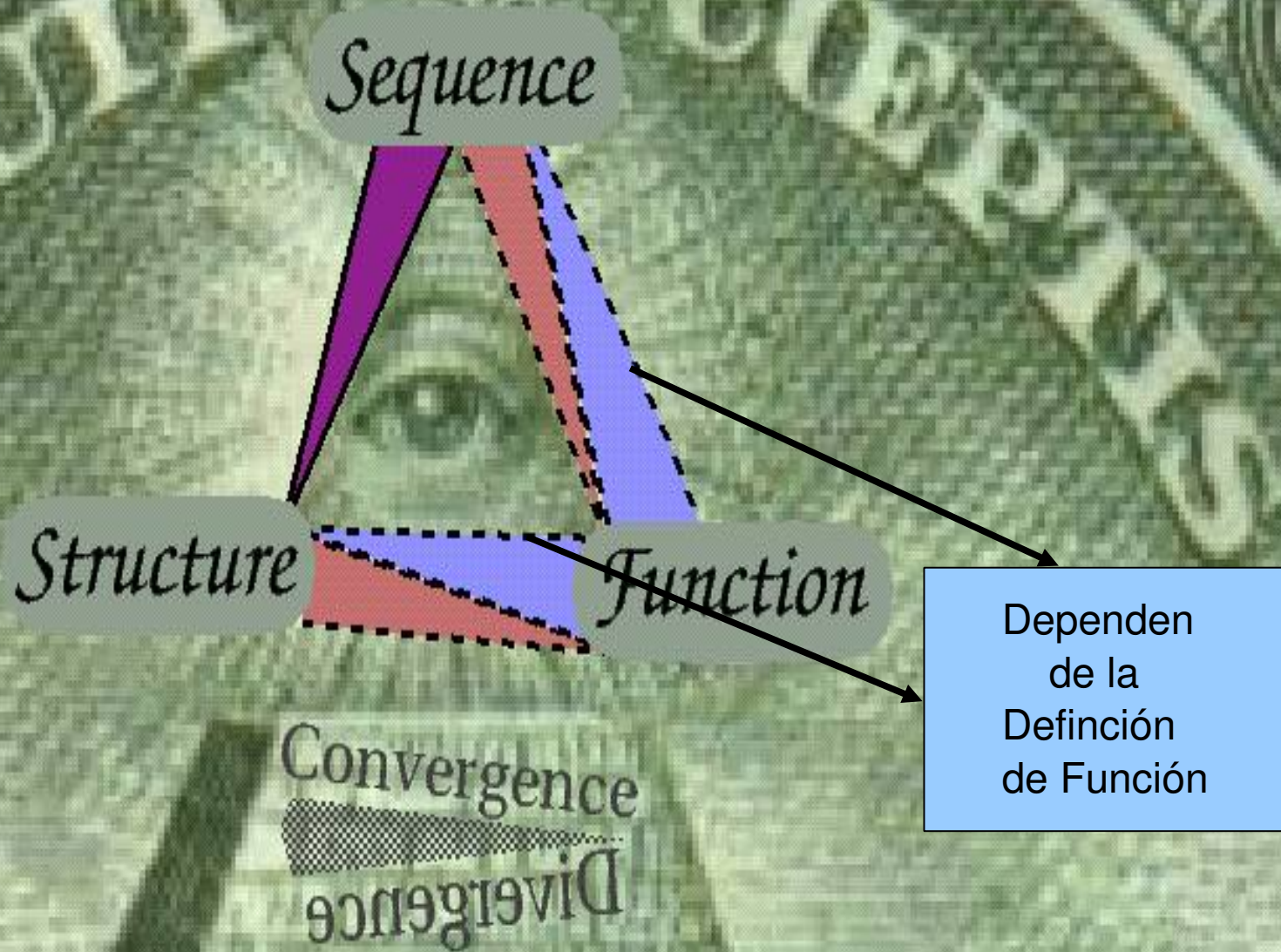
EF-Tu

They are homologous  
Proteins...

The Function could be  
very divergent

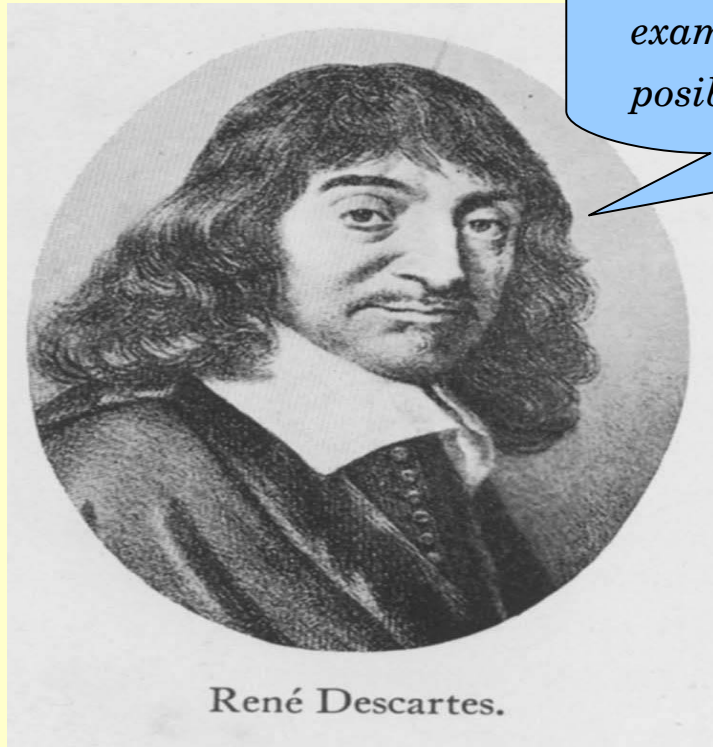
But... All of them bind GTP

# La Madre Naturaleza Consintió

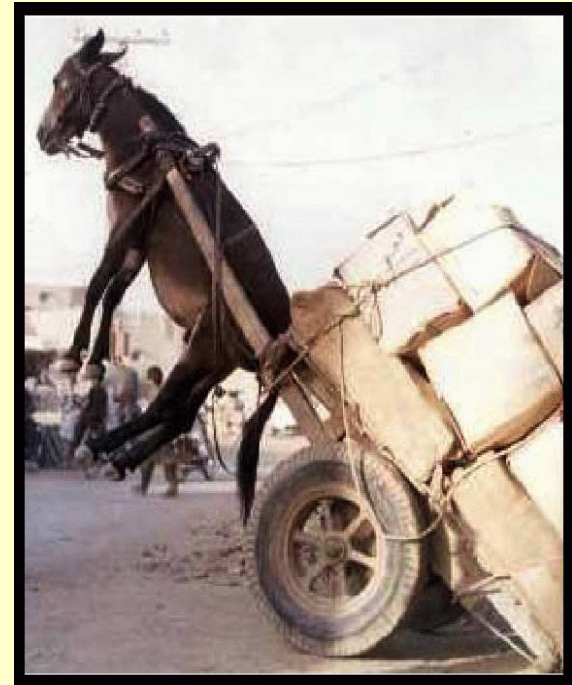


# Tarea Compleja la de transferir información estructural y/o funcional entre proteínas homólogas.

Qué Hacemos??

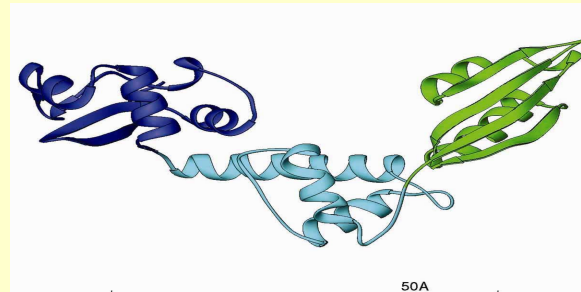


*Dividir cada una de las dificultades a examinar, en tantas partes como sea posible y necesario para resolverlas mejor*

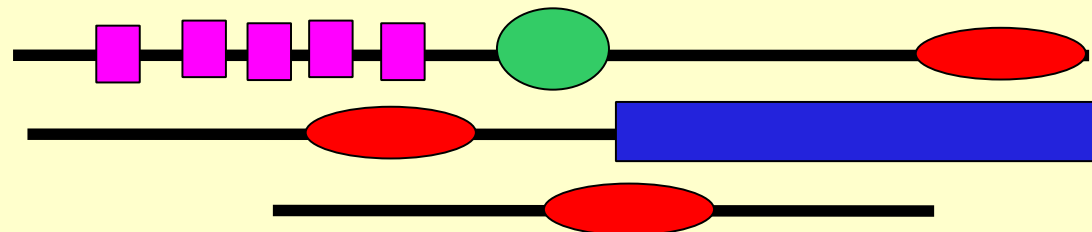


# DEFINICIÓN DE DOMINIO

Los dominios de proteínas han sido descritos, **desde un punto de vista estructural**, como unidades estructurales compactas y localmente independientes, caracterizadas usualmente por un núcleo hidrofóbico bien definido.

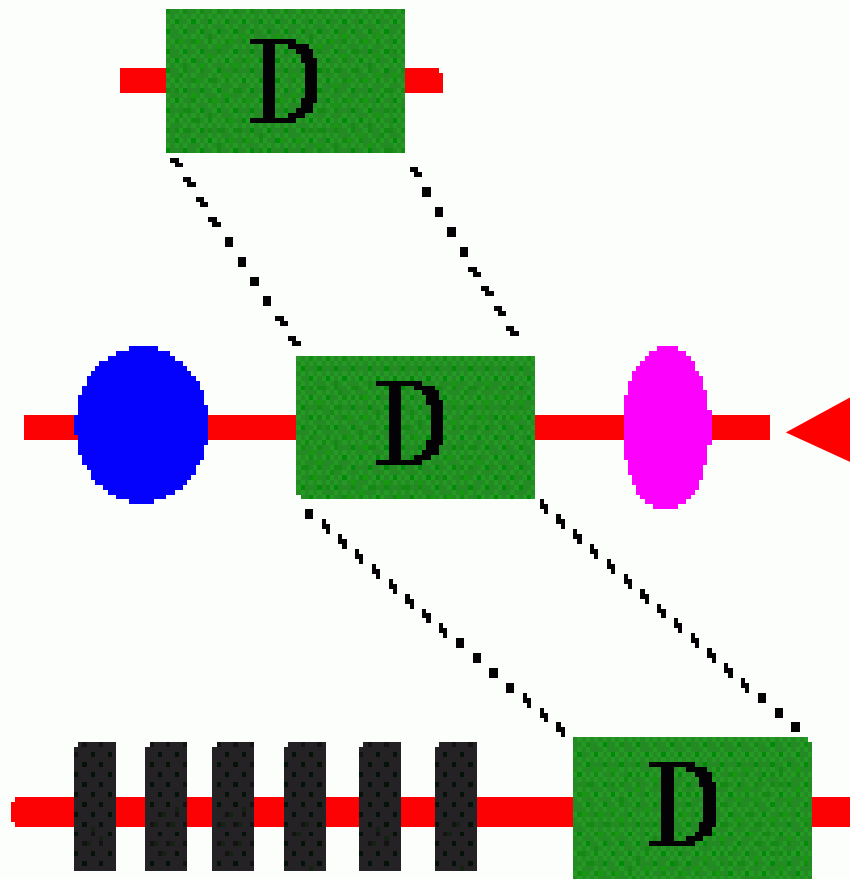


**Desde el punto de vista del análisis de secuencia**, los dominios se definen como regiones conservadas evolutivamente y adquieren mayor relevancia si son descritos como módulos móviles, es decir, presentes en diferentes familias de proteínas de arquitectura diversa.

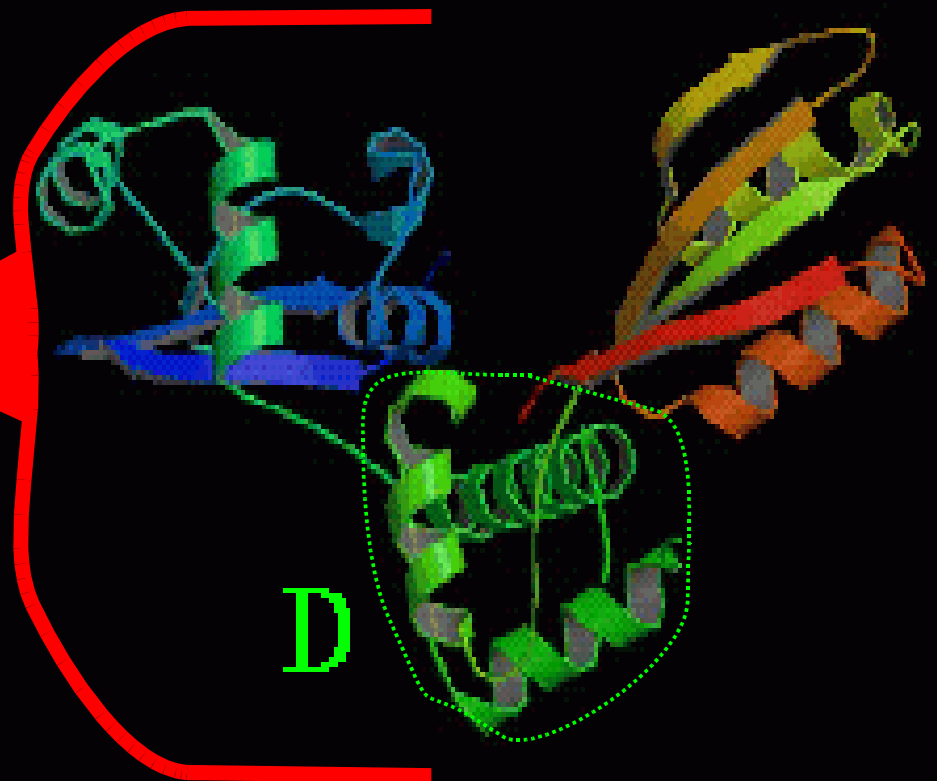




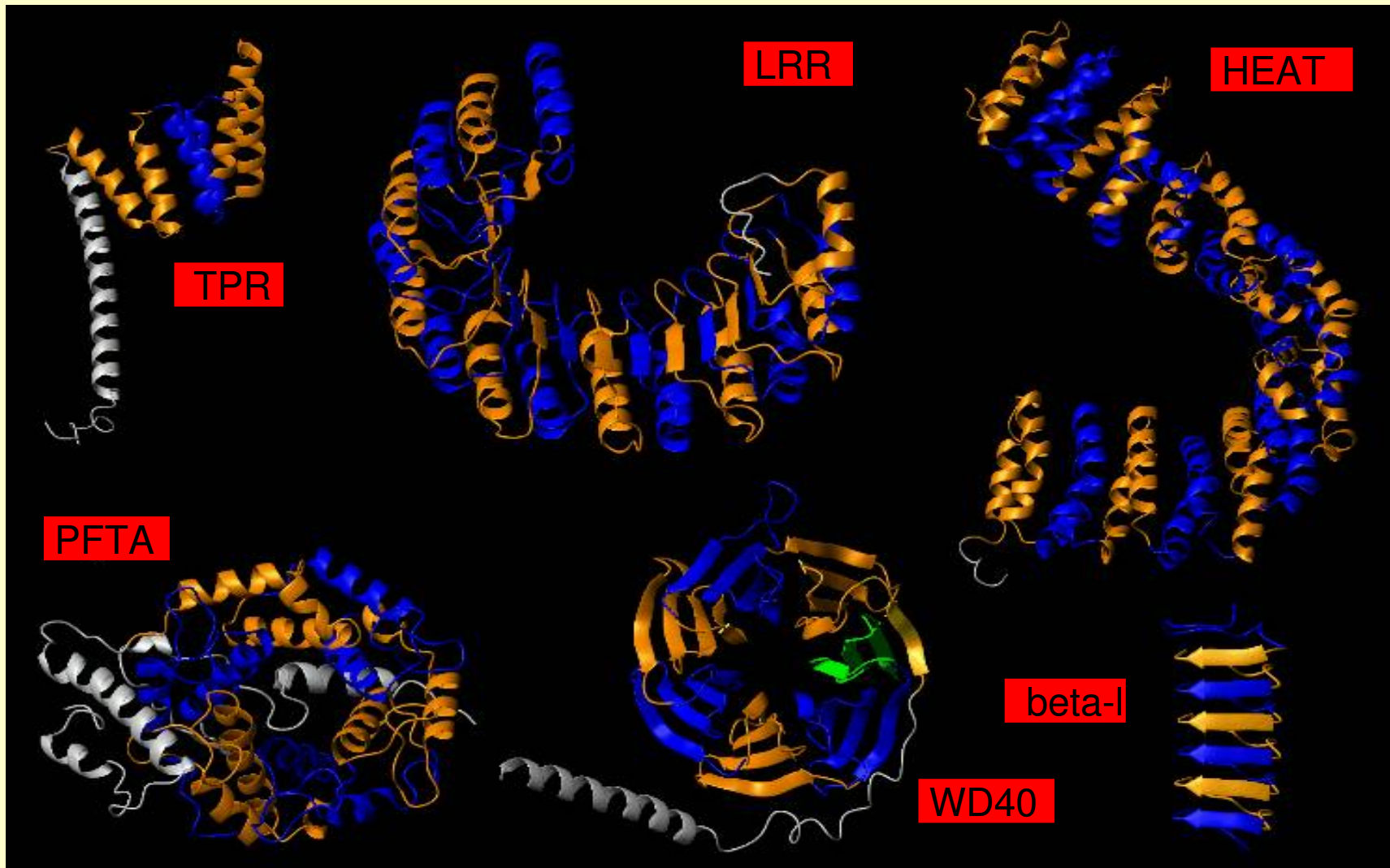
# Análisis de Secuencia



# Estructural



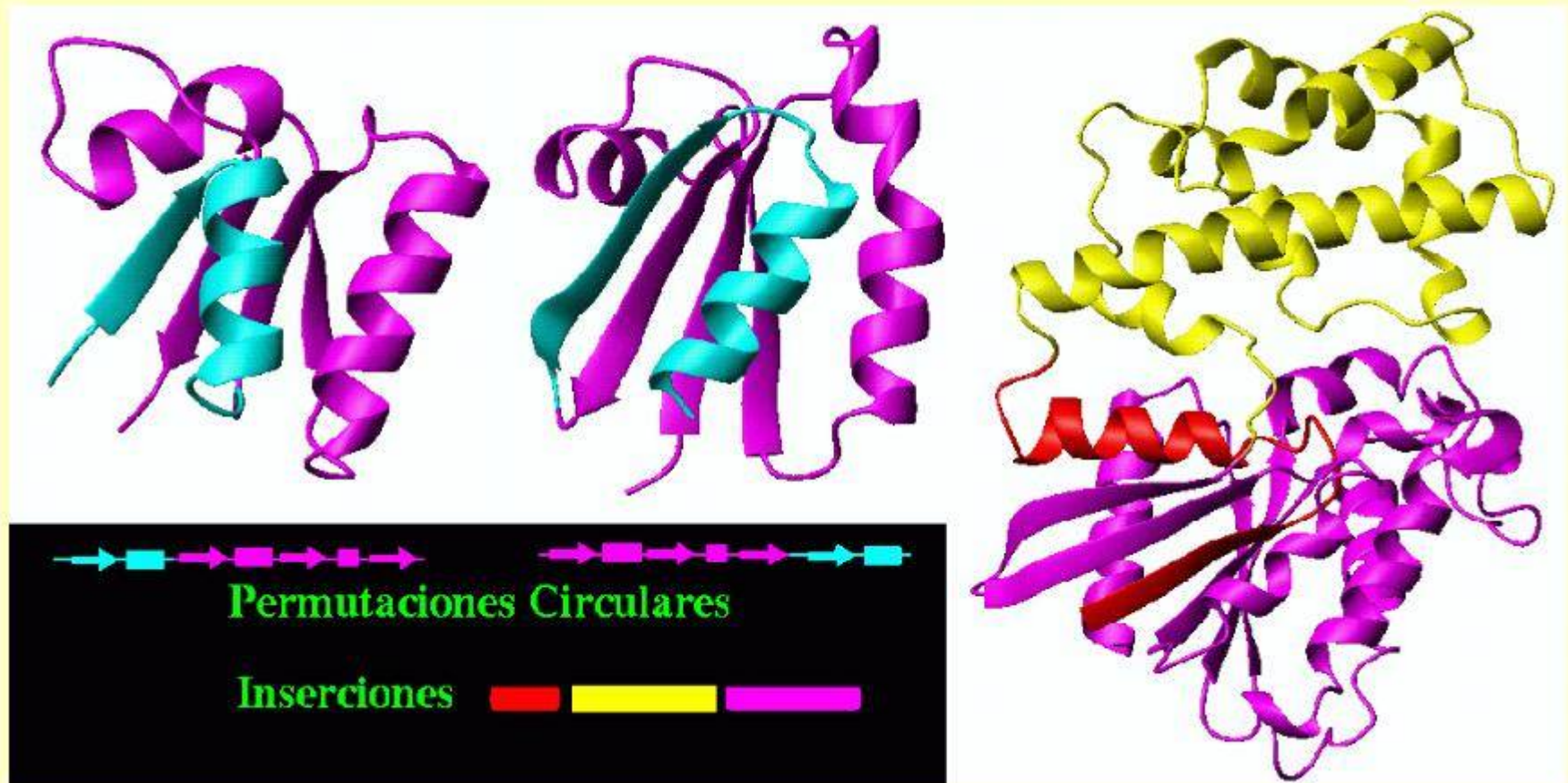
## REPEATS – In the limits of Domain Definition



Protein repeats. Short specialist review for the Encyclopedia of Genomics, Proteomics, and  
Cedida por: Perez-Iratxeta C, Andrade MA (2005) Bioinf. Ed. Wiley and Sons Ltd., UK.

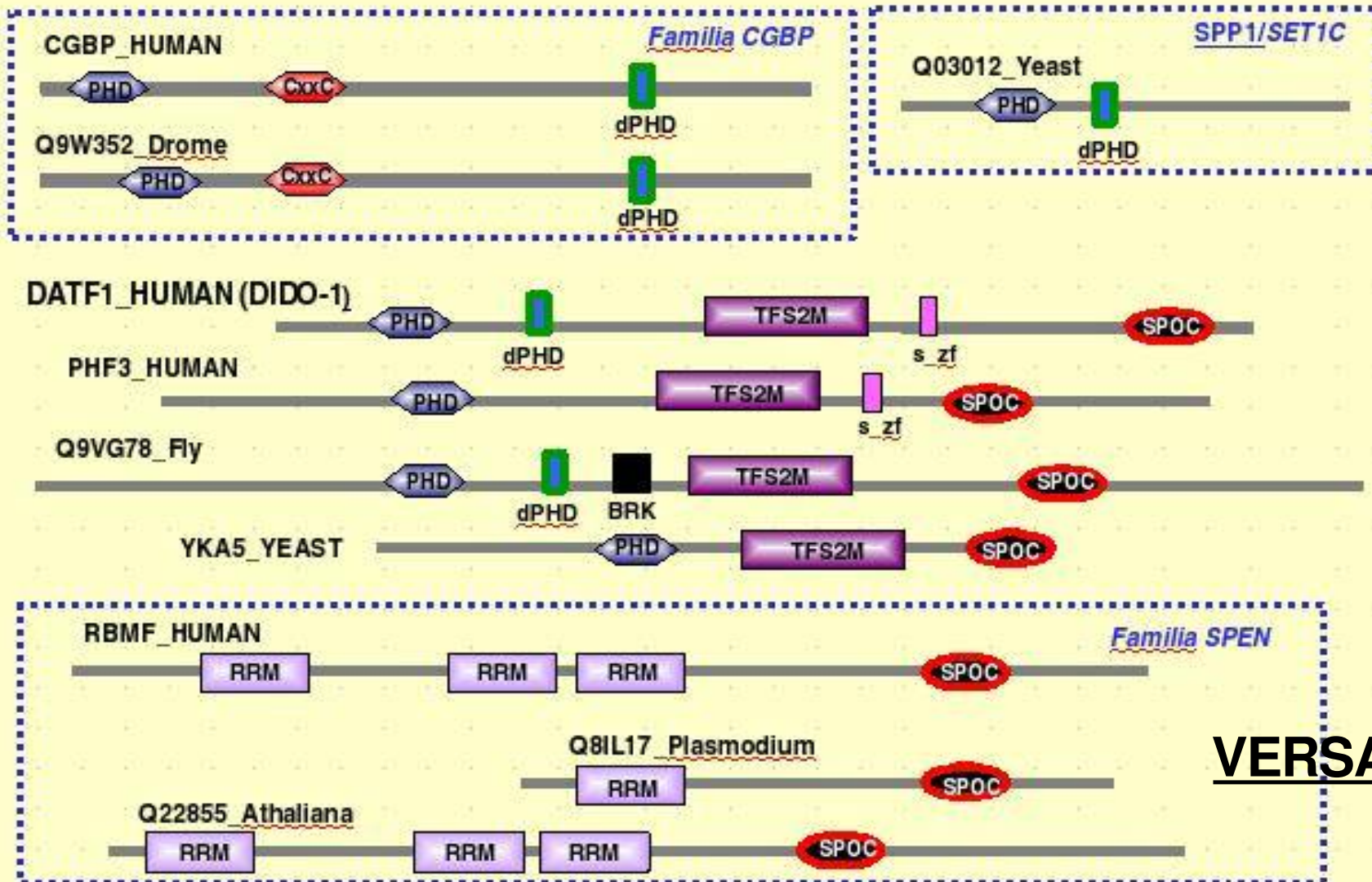
- **Protein irregularities that hinder sequence analysis**

- Low complexity regions
- Repeats, Trans-membrane and Coiled-coil regions (high mutation rates)
- and Fold irregularities, such as:  
Circular Permutations and Insertions



# The role of domains in protein evolution

## Shuffling, Accretion and Supra-Domains



**VERSATILITY !!**

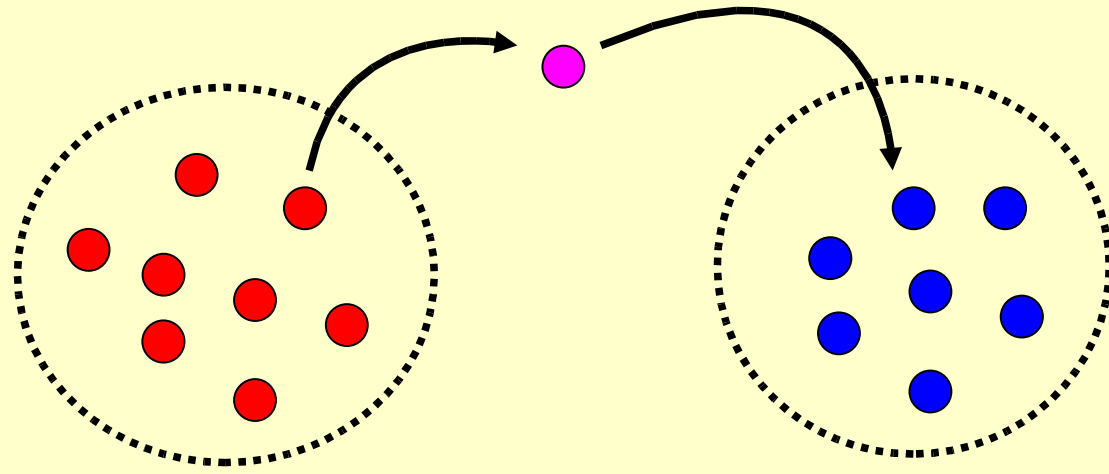
# METHODS ON DOMAIN ORIENTED SEQUENCE ANALYSIS



# Detection of homologous protein sequences

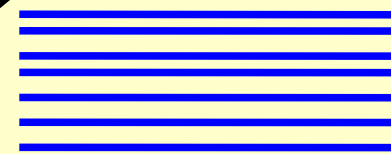
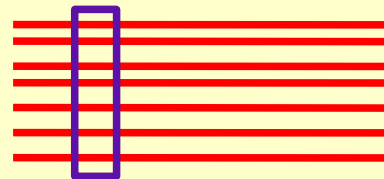
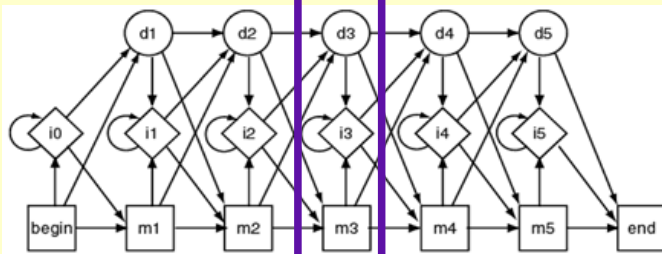
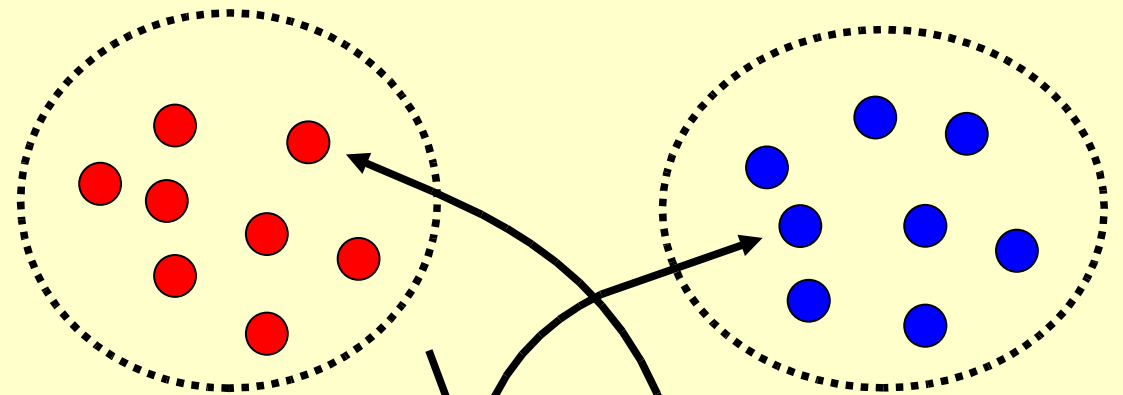
Two strategies

ISS (Blast, FASTA)



Profiles

(PSSMs: PsiBlast, HMMs)



¡Recíproco!


 Search Pfam

- Protein name or sequence
- Keyword
- Domain query
- DNA sequence
- Taxonomy query

Pfam is a large collection of multiple sequence alignments and hidden Markov models covering many common protein domains and families. For any family in Pfam you can:

- Look at multiple alignments
- View protein domain architectures
- Examine species distribution
- Follow links to other databases
- View known protein structures

For more information on Pfam, on using this site, or on the changes between Pfam releases 19.0 and 20.0, click [here](#).

Pfam can be used to view the domain organisation of proteins. A typical example is shown below. Notice that a single protein can belong to several Pfam families.



74% of protein sequences have at least one match to Pfam. This number is called the sequence coverage and is shown in the pie chart on the right.

Pfam is a database of two parts, the first is the curated part of Pfam containing over 8296 protein families. To give Pfam a more comprehensive coverage of known proteins we automatically generate a supplement called Pfam-B. This contains a large number of small families taken from the [PRODOM](#) database that do not overlap with Pfam-A. Although of lower quality Pfam-B families can be useful when no Pfam-A families are found.

### Version 20.0

May 2006, 8296 families



### Web feed

You can use the RSS feed to keep updated about Pfam releases

[XML](#) [RSS](#)

### Enter your keyword(s) here

 Go Example

### Enter a SWISS-PROT 48.1 or TrEMBL 31.1 name or accession number

 Go Example



**SMART MODE:**  
**NORMAL**  
**GENOMIC**

- Simple
- Modular
- Architecture
- Research
- Tool

Schultz et al. (1998) *Proc. Natl. Acad. Sci. USA* **95**, 5857-5864  
 Letunic et al. (2006) *Nucleic Acids Res* **34**, D257-D260

[HOME](#)
[SETUP](#)
[FAQ](#)
[ABOUT](#)
[GLOSSARY](#)
[WHAT'S NEW](#)
[FEEDBACK](#)

### Sequence analysis

You may use either a Uniprot/Ensembl sequence identifier (ID) / accession number (ACC) or the protein sequence itself to request the SMART service.

Sequence ID or ACC

Sequence

### Architecture analysis

You can search for proteins with combinations of specific domains in different species or taxonomic ranges. You can input the domains directly into "Domain selection" box, or use "GO terms query" to get a list of domains. See [What's New](#) for more info.

Domain selection

Example: **TyrKc AND SH3 AND NOT SH2**

GO terms query

Example: **membrane AND signal transduction**

Taxonomic selection

Select a taxonomic range via the selection box or type it into the text box below:

Examples: **Dictyostelium**





[Remove menu]



- InterPro home
- Text Search
- InterProScan
- Databases
- Documentation
  - ▶ Tutorial
  - ▶ Project Outlines
  - ▶ Collaborators
  - ▶ Example Entry
  - ▶ Dataflow Scheme
  - ▶ Release Notes
  - ▶ User Manual

## InterPro Home

InterPro is a database of protein families, domains and functional sites in which identifiable features found in known proteins can be applied to unknown protein sequences.

Further information on InterPro can be found in the [documentation](#) - see links on the left hand side.

For information, comments and/or suggestions on the InterPro database, please contact us at [EBI Support](#).

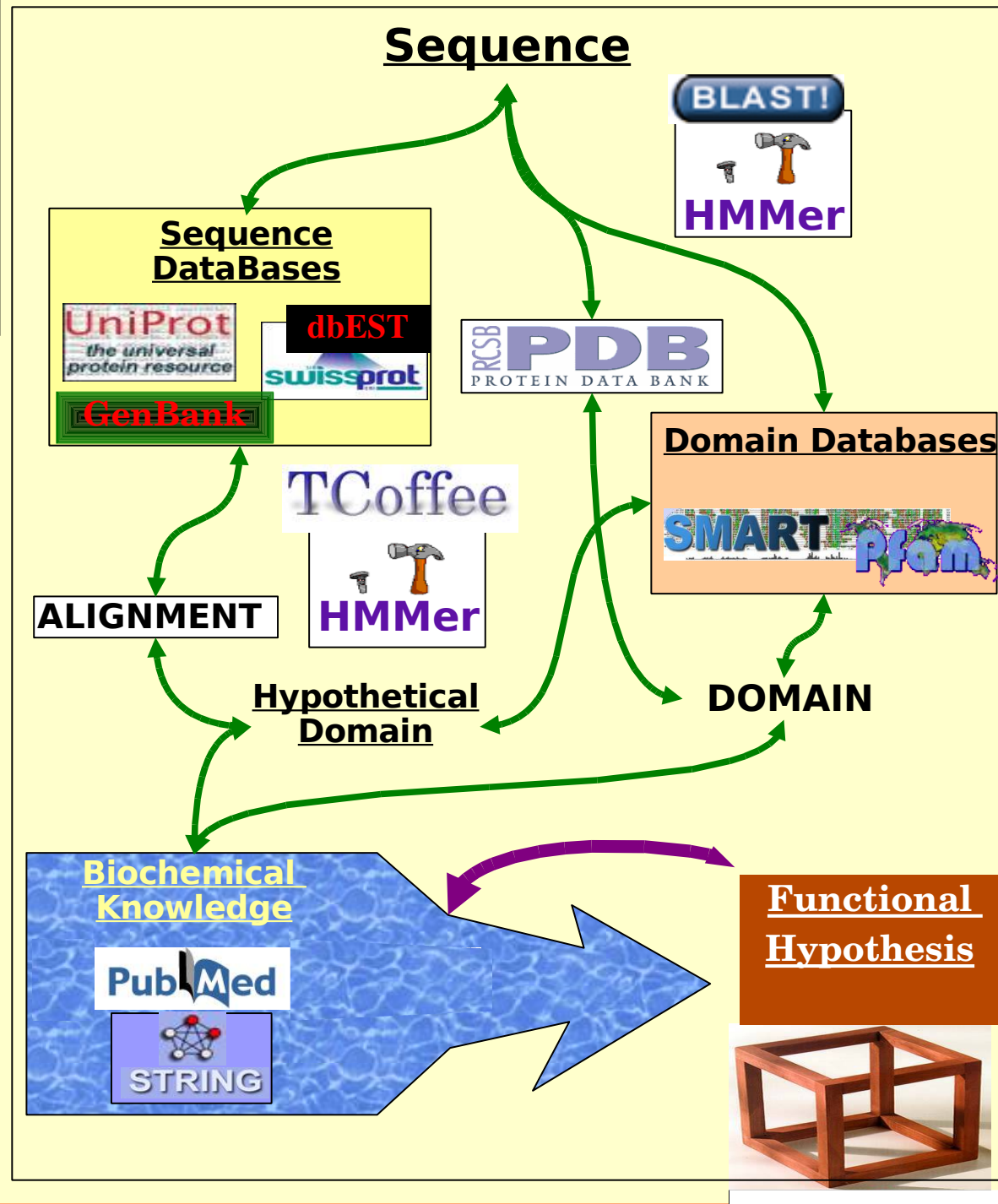
Search

Search - [help](#) - example: [kinase](#)

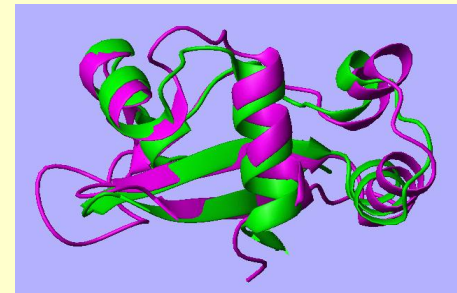
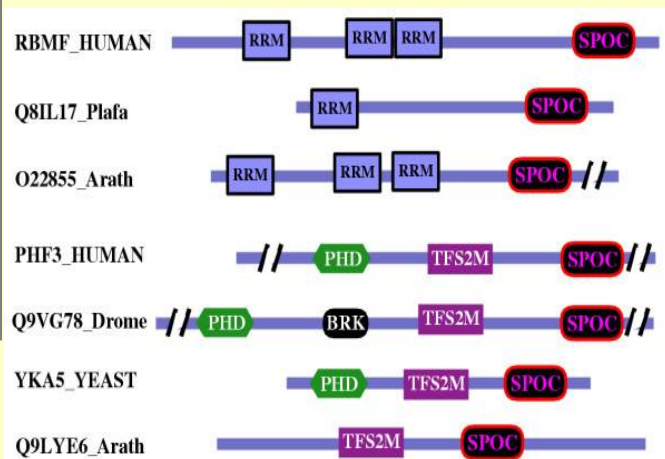
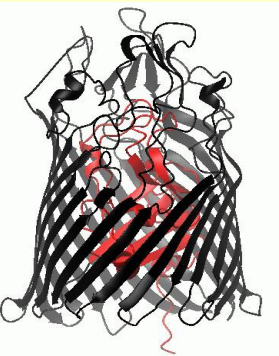
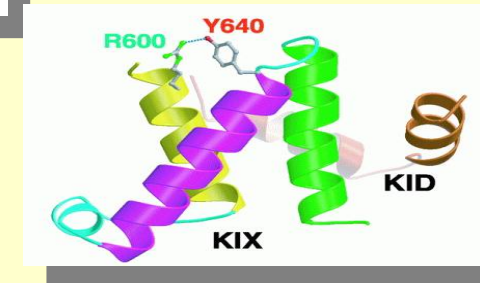
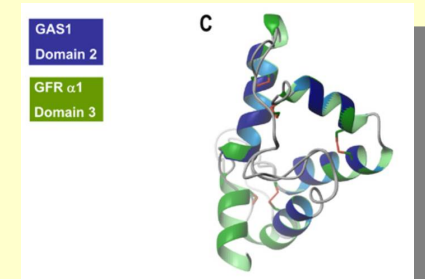
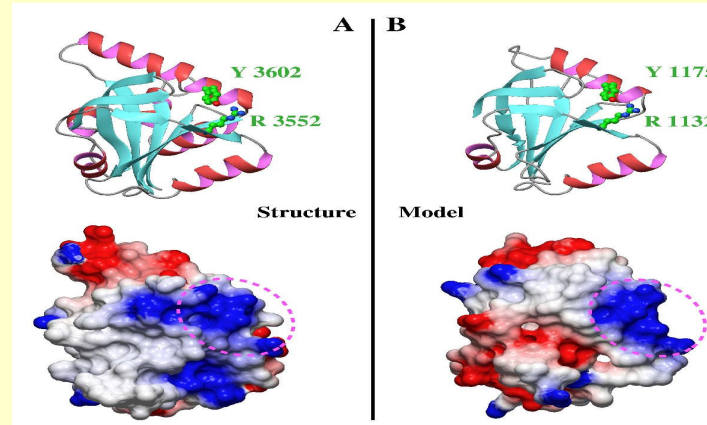
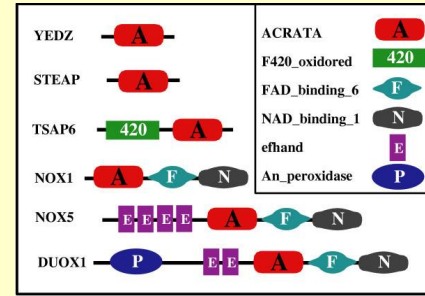
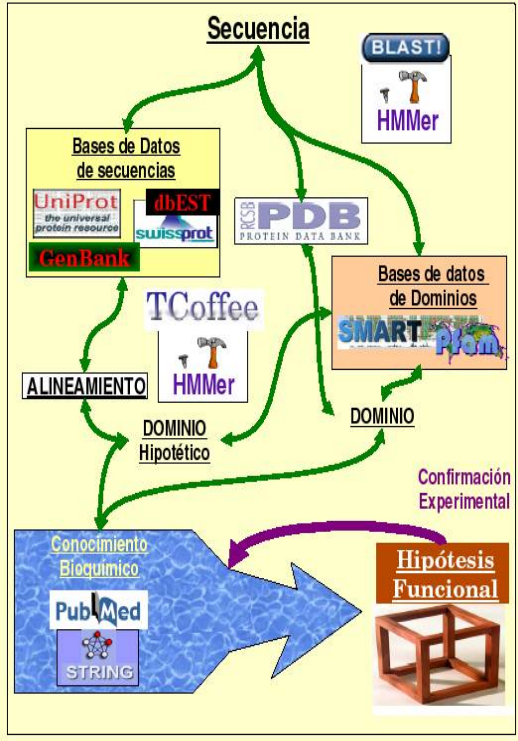
Search Entries



# Domain Oriented Sequence Analysis Flow-Chart



# REAL-LIFE EXAMPLES



# SPOC: A widely distributed domain associated with cancer, apoptosis and transcription.

Sanchez-Pulido L, Rojas AM, Van Wely K, Martinez-A C, Valencia A.

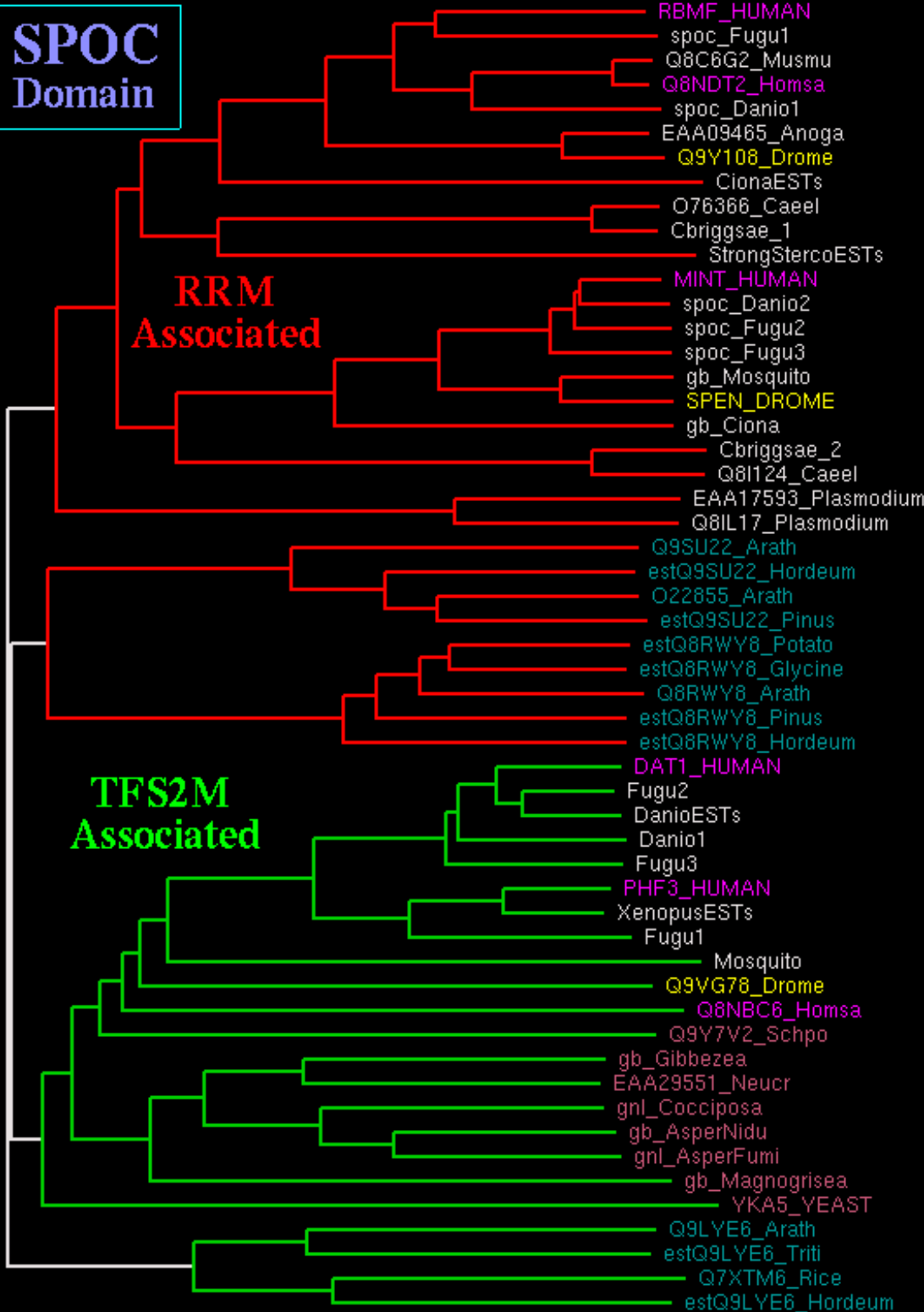
## CNB-CSIC

DATF1\_HUMAN (DIDO-1)

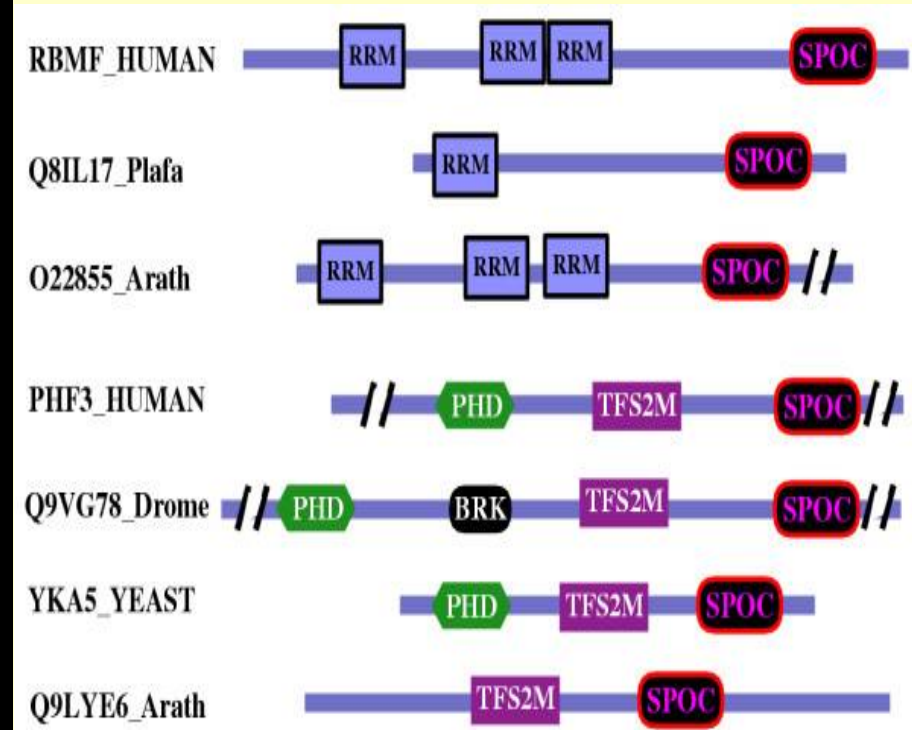


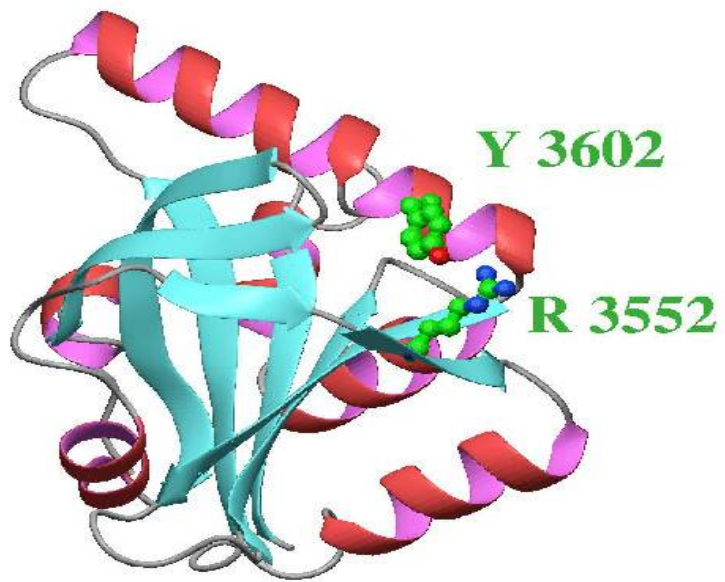
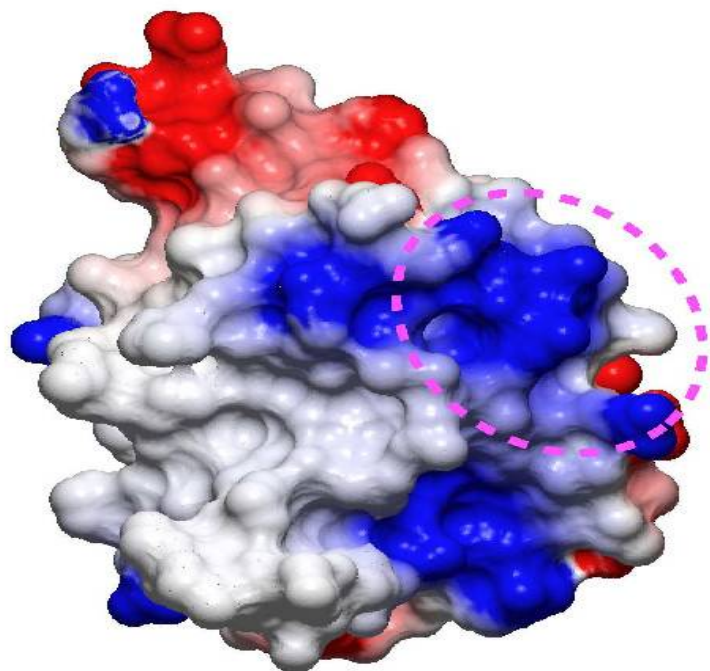
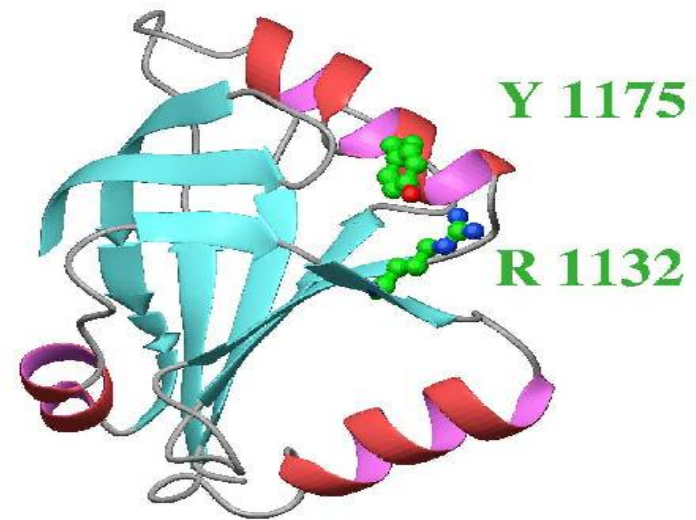
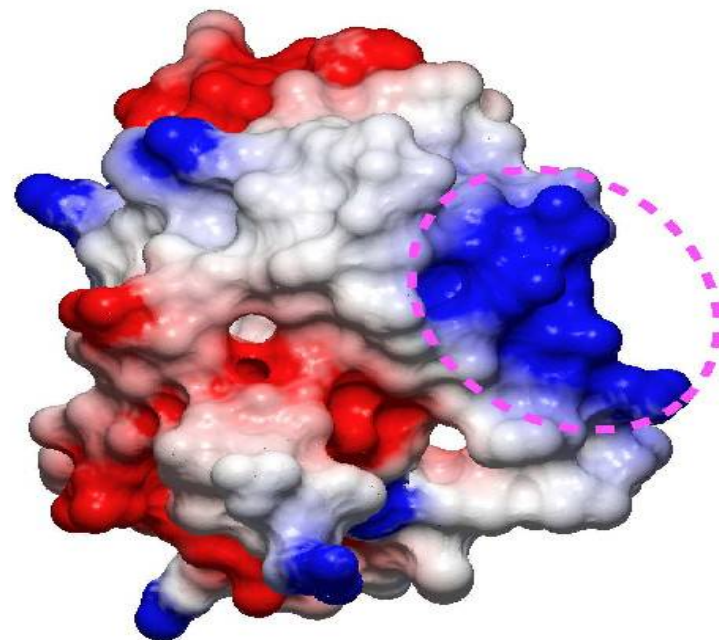
DAT1_HUMAN	1093	WKGFINMQSVAKFVTKAYPVSG	....CFDYLEDLPD	.....	TIHIGGRIAPKTWWDYVGKLGKSSVSK	ELCLRFPHPATE	.....	EEEVAYISLYSYFSSRGRFGVWANNR	.....	HVKDLYLIP	1199	
PHD2DPred		..SEEE	HHHHEEEEEEEE	....HHH	EEEEEE	HHHHHHHHHH	.....	EEEEEE	.....	HHHHHHHHHH	SEEE	EEEEEE
estDio_Fish		WKGFINMHSVAKFVTKAYLVSG	....SFENIKEDLPD	.....	TIHIGGRILPHTWWDYVGKLGKTSLSK	ELSLIRFPHPATE	.....	EEEVAYVSLFSYFSSRKRFGVWANGK	.....	RIKDLYLIP		
PHF3_HUMAN	1209	WKGFINMPSVAKFVTKAYPVSG	....SPEYLTEDLPD	.....	SIQVGGRIISPQTWWDYVEKIKASGTK	EICVVRFPVTE	.....	EDQISYTLFFAYFSSRKRKYGWAANNM	.....	QVKDLYLIP	1315	
estDio_Frog		WQGFNMPVAKFLIKAYPVSG	....SLEHLAEDLPE	.....	SIQVGGRIISPQTWWDYVDKIKASGTK	ETCLVRFPVTE	.....	EDQISYTLFFSYPSSRKRKYGWAANNM	.....	QVKDLYLIP		
Q9VG78_Drome	1621	WSGTLKMIDLADFEIVMYPVQG	....NCHQLGNLMP	.....	QMDVIGRITRVNWDYVEKIKASGTK	EVVIVNIFPASP	.....	SETYKFDLFFEYLDSPQRLGVLGVEDS	.....	QIRDFYIFP	1727	
unf_Aspnidu		WHGRVVMNPVAEFSSPAKHVAGADLS	GRIPWNDLIPS	.....	TLIDGRIKIQSAGEYLQGLRPSQST	DVSVVAISSPDSS	.....	KDKSNFDKLFDPYQGRBRYGVMGKHPL	.....	AVRDTYLIP		
Q9Y7V2_Schpo	484	WTGKVKMATVSEFHANALNLFEDV	SASHLFEILSA	.....	TALIEGRISVSSVLQYFHALRKTPSK	EIIAVLFPVTE	.....	QNSQGFIDLYDFVKRNRYSVGLHKS	.....	SVKDAYIIP	591	
YKA5_YEAST	442	YPLGLEFTGYLNYIGASQKLR	DIFKEAIGDG	.....	KLYVEGRLLPTTAAAPYLKEISCSR	AILVYQLFSPNDS	.....	ESKTFADVVDSLENKGRIAGIKPKTR	.....	YEKDFYIVP	547	
Q9LYE6_Arath	542	WDGILQLSMSSVVPVAGIPKSG	EKAETSEWPA	.....	MVEVKGRVRLSGFGKFIQELPKSRT	RALMVMYLAYKDGISES	.....	QRGSLIEVIDSYVA.DQRVGYAEP	.....	SGVELYLCP	648	
est_Triti		WEGAIQLTSLSLTNVVAIFKSG	EKPSGKEWSS	.....	LIEIKGRVKLSAPQDFLEQLPKSRS	RAIMVTELCWKEGSSBS	.....	GRQLSQTIDSYIA.DBRVGLAEP	.....	DGLELYLCP		

# SPOC Domain



## Domain located in, at least, Two Architectures



**A****Structure****B****Model**

Gas1 is related to the GFR $\alpha$  family and regulates Ret signaling

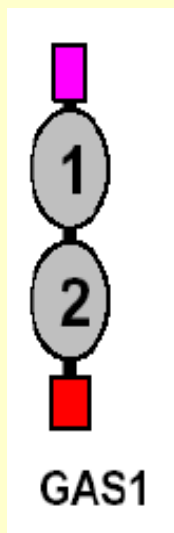
Cabrera J.R., Sánchez-Pulido L., Rojas A.M., Valencia A.,

Mañes S., Naranjo J.R. & Mellstrom B. (2005)

CNB – CSIC

PROTEÍNA INICIAL: GAS1 (Growth Arrest Specific 1)

FUNCIÓN: Regulación de procesos apoptóticos.

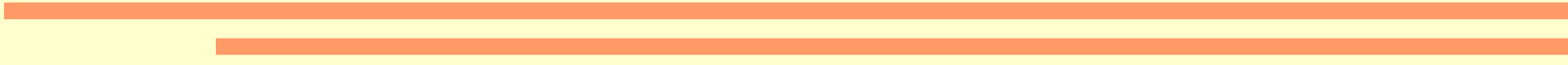
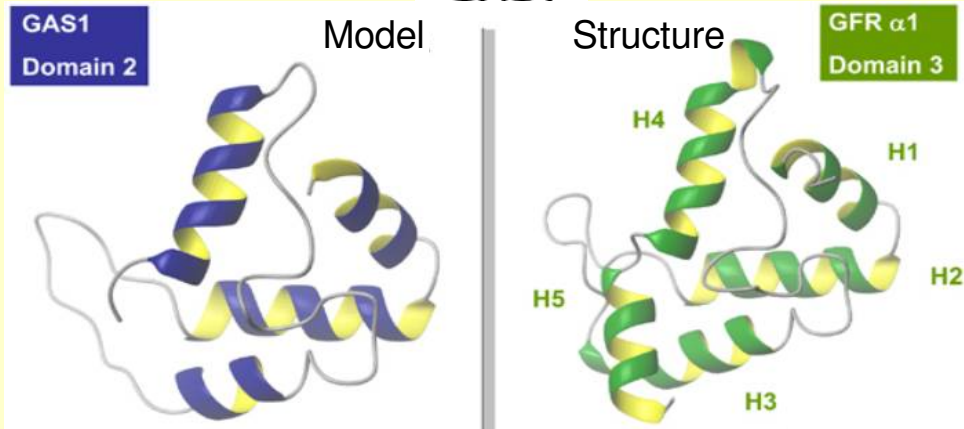
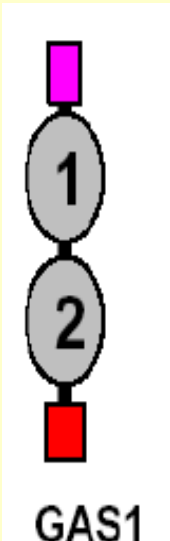


En un primer abordaje:

- Péptido señal
- Duplicación Interna
- GPI-Anclaje

# GAS1

GAS1_HUMAN_1	48	CWQALLCCQGE.PECSYAYNOYAEACAPVLAQHGGGDAPGAAAAAFPASAASFSSRWRCP...	SHCISALIQLNHTRRGF.....	ALEDCDCA.....	QDENCKSTKRAIE..FC	146	
GAS1_HUMAN_2	166	CTEARRRCDRD.SRCNLALSRYLTYCGKV.....	FNGLRCT...DECRVIEDMLAMPKVA.....	LLNDCVCD...	GLERPICESVKENMAR..LC	243	
GAS1_MOUSE_1	47	CWQALLCCQGE.PDCSYAYSQYAEACAPVLAQRGGADAPG..PAGAFPASAASSPRWRCP...	SHCISALIQLNHTRRGF.....	ALEDCDCA.....	QDEHCRSTKRAIE..FC	143	
GAS1_MOUSE_2	162	CTEARRRCDRD.SRCNLALSRYLAYCGKL.....	FNGLRCT...DECRAVIEDMLAVPKAA.....	LLNDCVCD...	GLERPICESVKENMAR..LC	239	
estGas1_Frog_1	25	CWQAMMRQCEE.AECSYAYRQYVDACSSVLPRPGGEA.....	ASSSSSSSSSSRRRCF...	SHCISALIQLNHTRWGF.....	ALEDCDCA.....	MDETCRATKRAIE..FC	116
estGas1_Frog_2	138	CMEARNICEGD.WRCGMSLSRYLTKCGRL.....	FDGLRCT...DECKEVIEDMMRVPKAM.....	LLSECECD...	GHERPICESIKENMAR..LC	215	
Gas1_Fish_1	31	CWKAILKCHGD.PDCHYAYDQYLYACASVI.....	SGEHQKCP...SHCISLIQLNRTQSGF.....	ALEDCDCA.....	LDPVCRSAKQAIIE..FC	107	
Gas1_Fish_2	116	CTEARLEGEAD.PSSAMKDYLFHCRKL.....	FGGERCT...EECRRVIADMRSIPKAQ.....	QLDTCVCD...	GAERNICEYIKASMKT..FC	193	



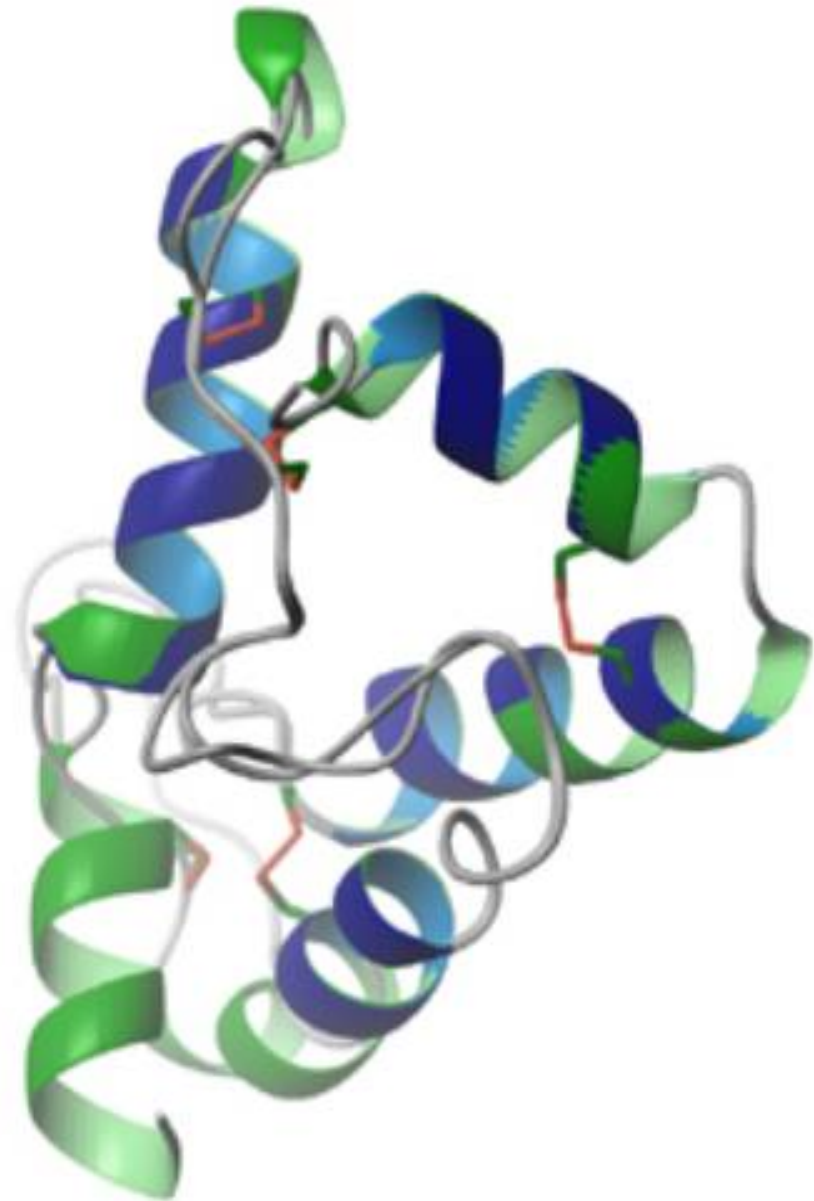


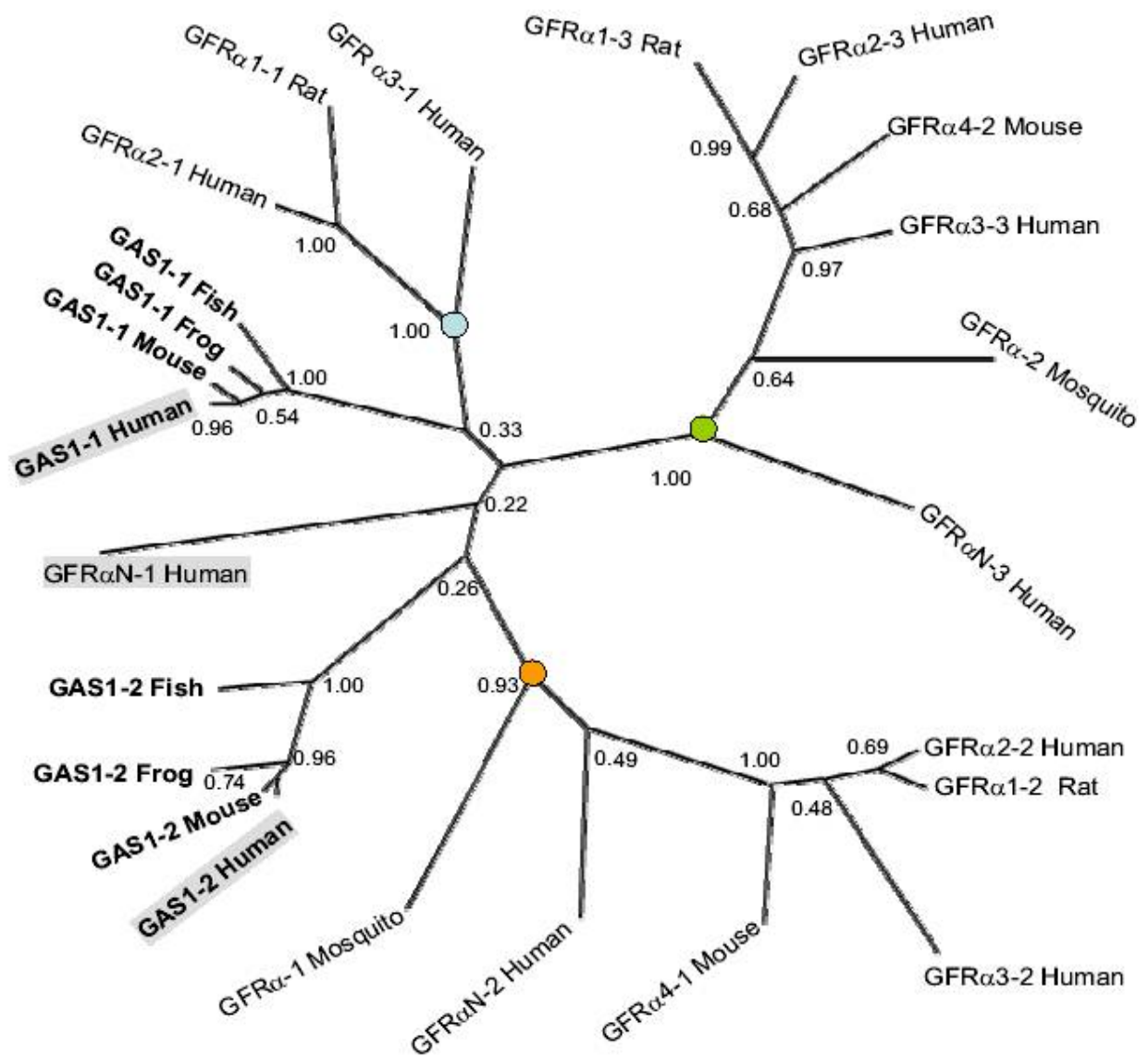
GAS1

Domain 2

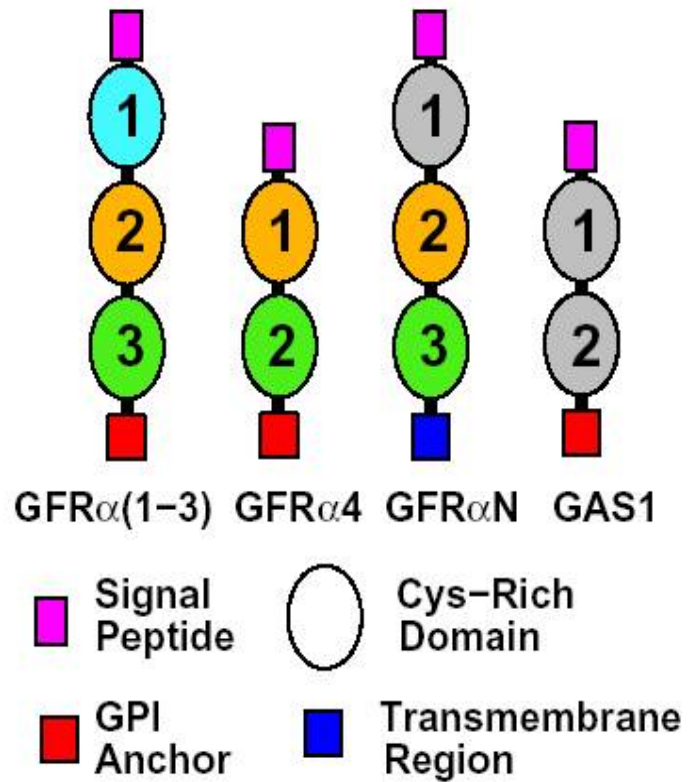
GFR  $\alpha$ 1

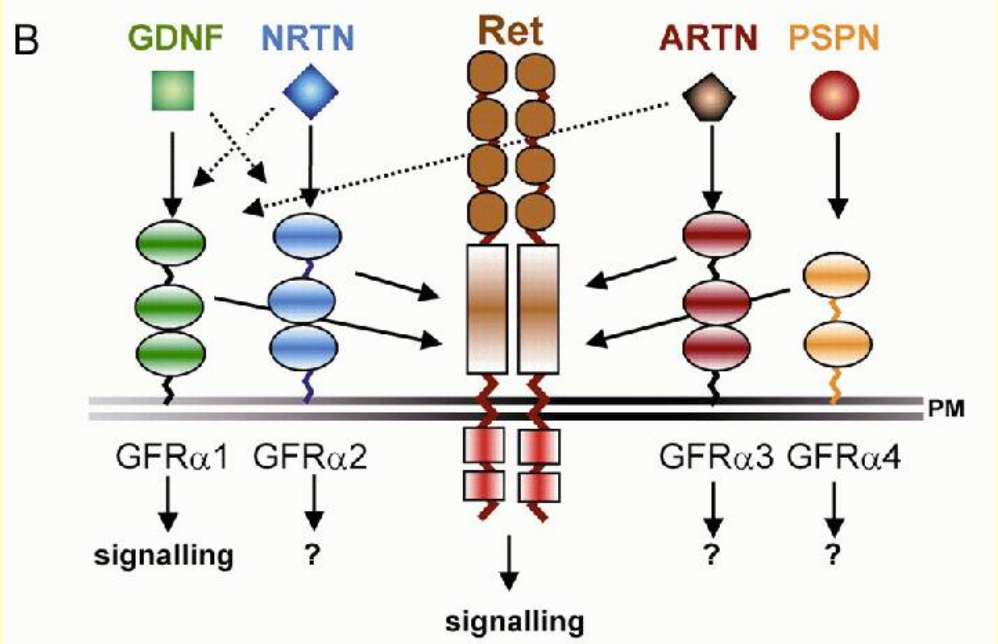
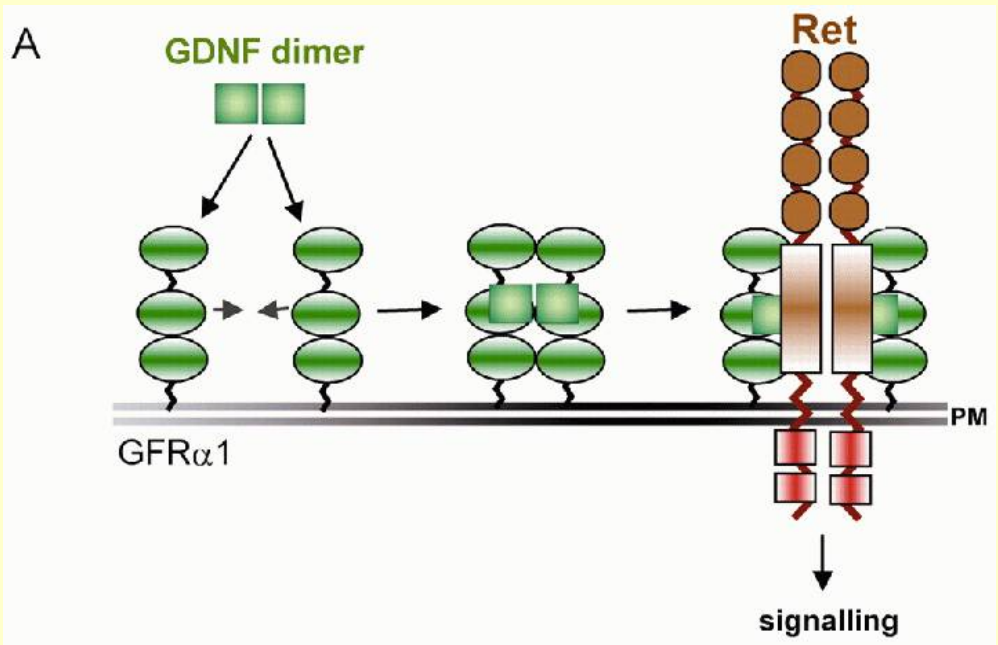
Domain 3





• Arquitectura similar





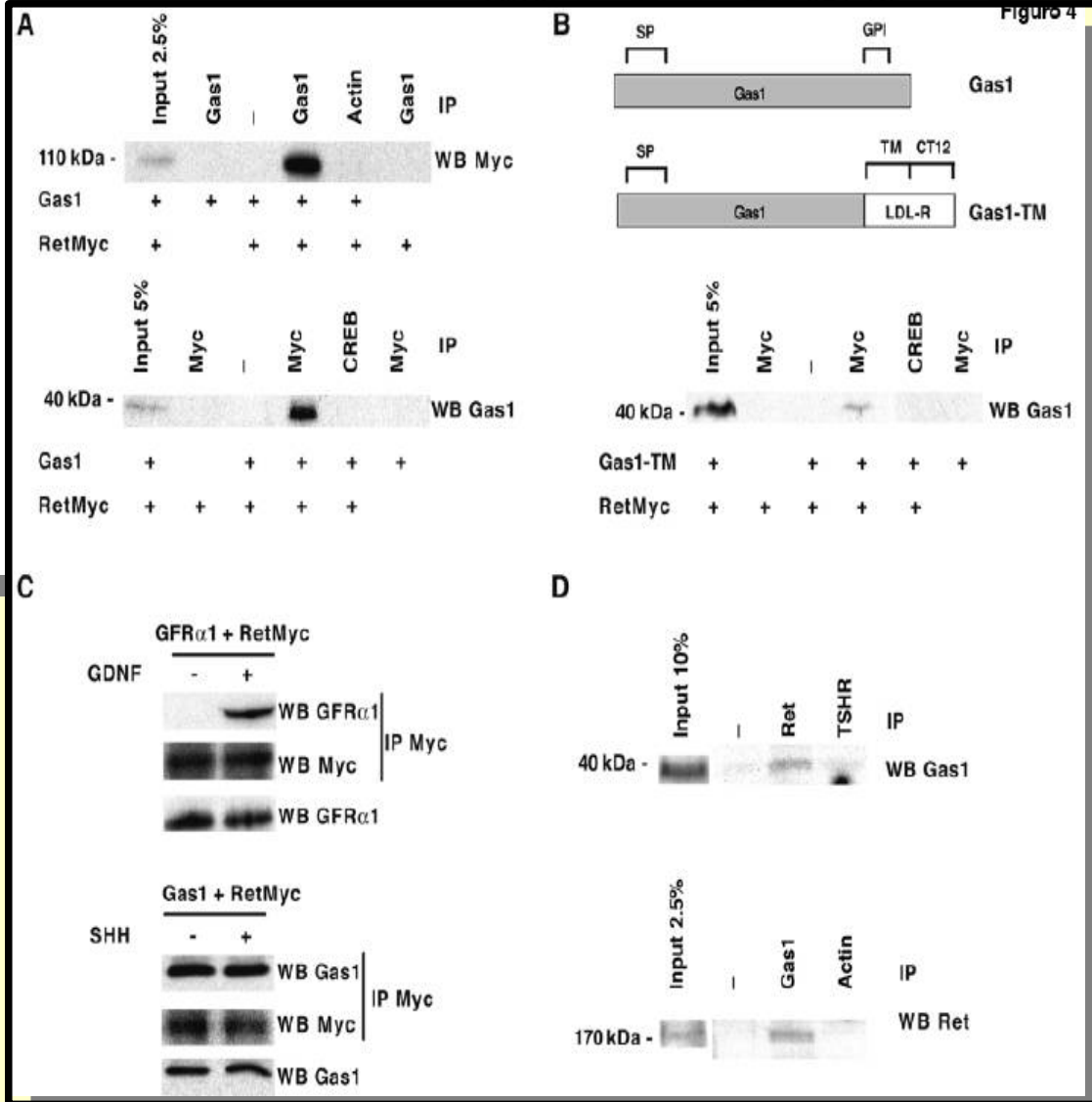
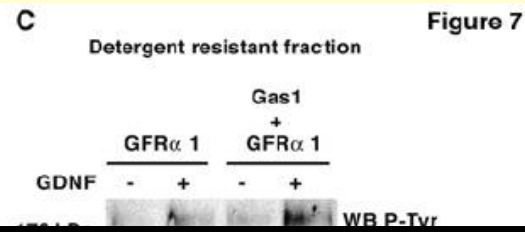
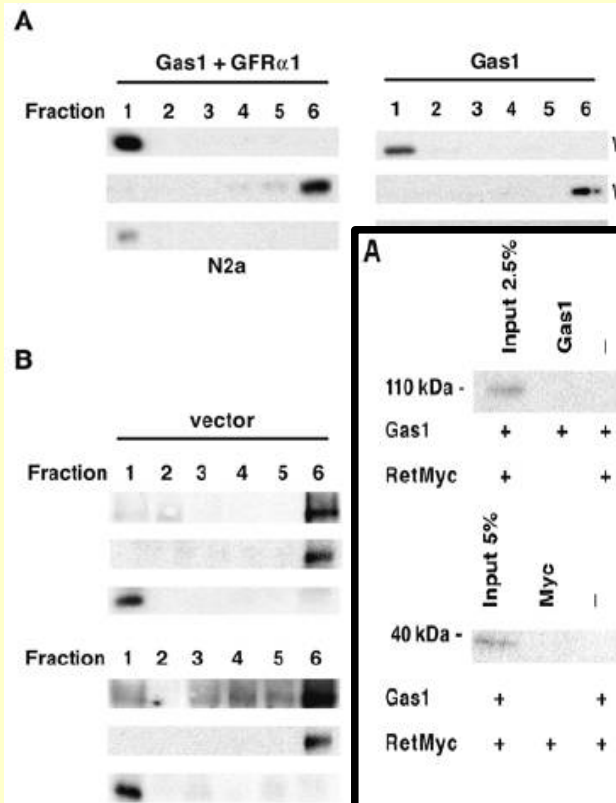
JR Cabrera  
CNB-CSIC

# BIBLIOGRAFÍA

Diferentes factores tróficos o ligandos:

- GDNF Glial cell Derived Neurotrophic Factor
- NRTN Neurturin
- ARTN Artemin
- PSPN Persephin

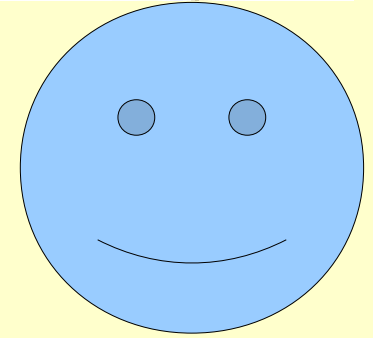
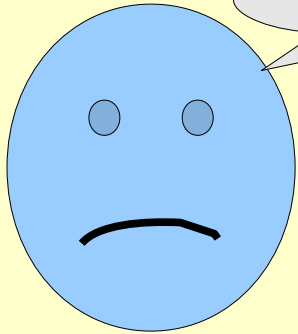
La Familia **GFR $\alpha$**  converge en la transducción de la señal en su **interacción** con la quinasa Ret -----> **¿Y GAS1...?**



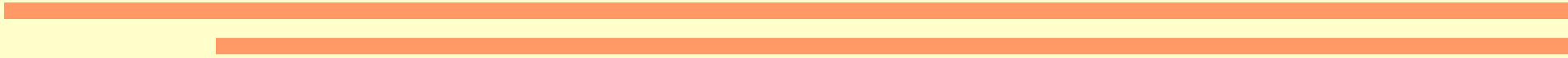
Computational predictions supported by experimental analysis.



DTRGHYFASSTNDR  
???????????????

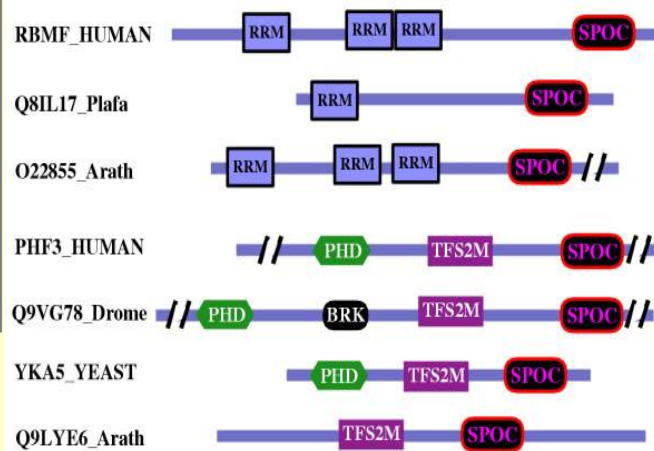
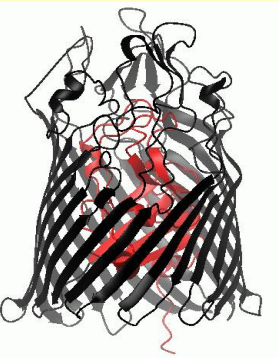
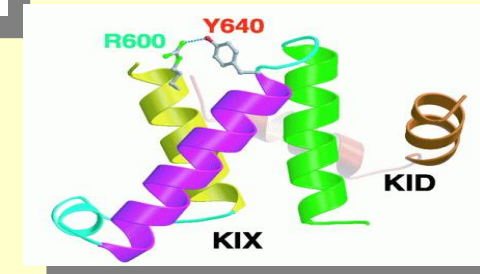
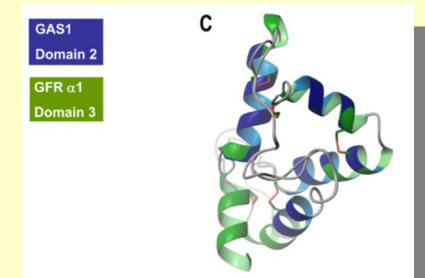
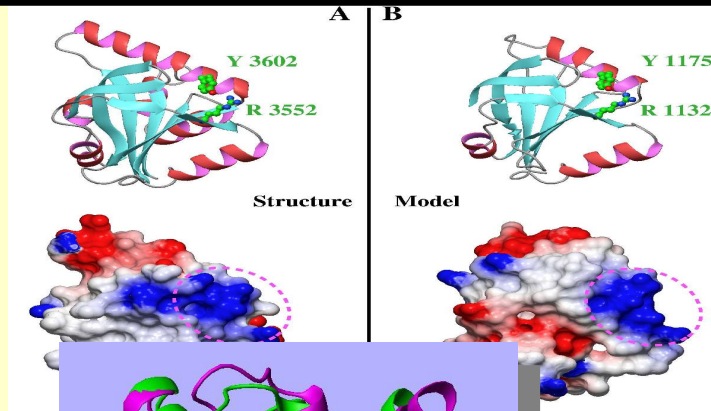
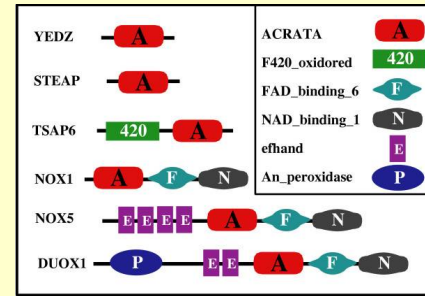
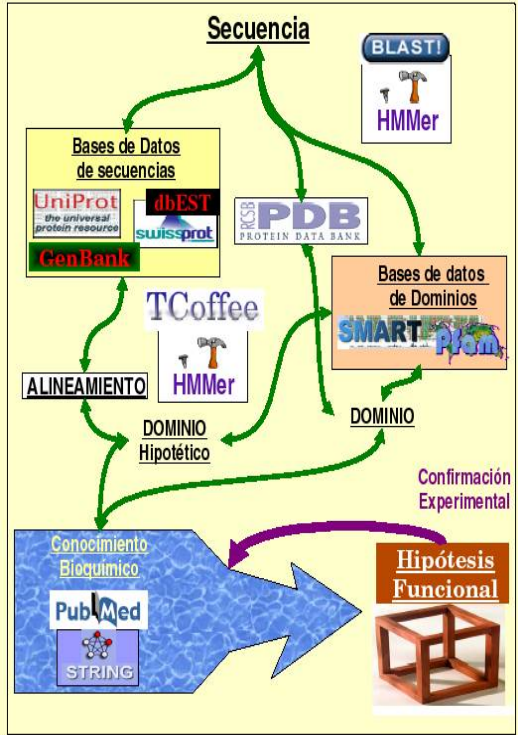


*Jorge Ruben,  
Santos  
&  
Ana*



# CONCLUSIONE

**¡Ogni Proteina `e un Mondo!**



*"As a general guide to functional annotation, it should be kept in mind that current methods for genome analysis, even the most powerful and sophisticated of them, facilitate, but do not supplant the work of a human expert."*

*Eugene Koonin.*

# Basic References:

Zuckerlandl E, y Pauling L. (1965)

**Evolutionary divergence and convergence in proteins.**

*Evolving Genes and Proteins,*  
*Academic Press, New York, 97-166.*

Bork P, Gibson TJ. (1996)

**Applying motif and profile searches.**

*Methods Enzymol. 266:162-184.*

Iyer LM, Aravind L, Bork P, Hofmann K, Mushegian AR, Zhulin IB, Koonin EV. (2001)

**Quoderat demonstrandum? The mystery of experimental validation of apparently erroneous computational analyses of protein sequences.**

*Genome Biol. 2:RESEARCH0051.*

*Questions:*

*sanchez@cnb.uam.es*

---

---