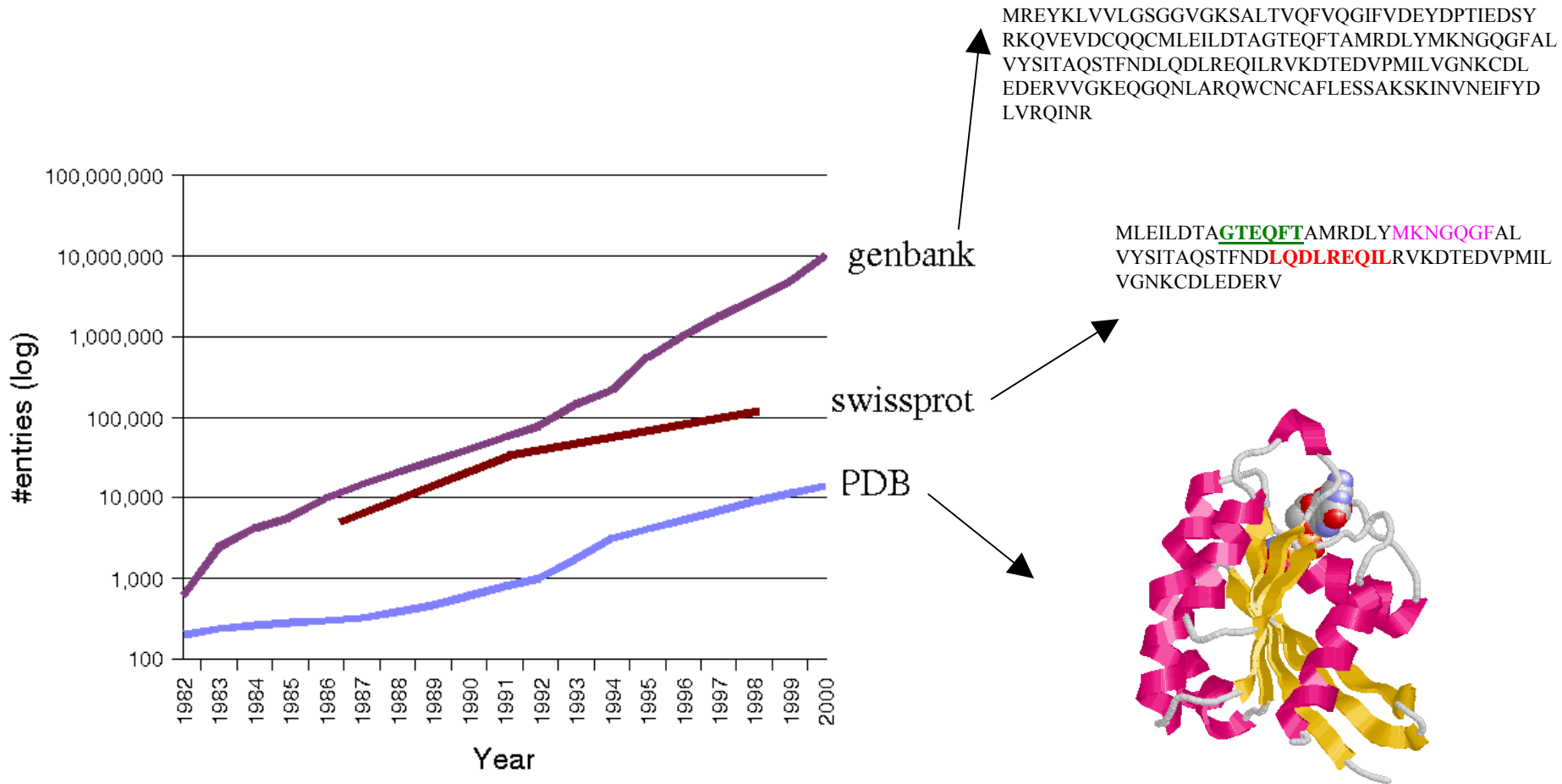Protein Sequence Analysis

# Exploiting Sequence Relationships for Predicting Protein Function

Florencio Pazos (CNB-CSIC)

*Florencio Pazos Cabaleiro*
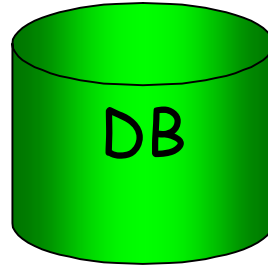*Protein Design Group (CNB-CSIC)*
*pazos@cnb.uam.es*

# Protein Sequences, Structures and Functions

MREYKLVVLGSGGVGKSALTVQFVQGIFVDEYDPTIEDSY
RKQVEVDCQQCMLEILDTAGTEQFTAMRDLYMKNGQGFAL
VYSITAQSTFNDLQDLREQILRVKDTEDVPMILVGNKCDL
EDERVVGKEQGQNLARQWCNCAFLESSAKSKINVNEIFYD
LVRQINR

MLEILDTA**GTEQFT**AMRDLY**MKNGQGF**AL
VYSITAQSTFND**LQDLREQIL**RVKDTEDVPMIL
VGNKCDLEDERV



genbank

swissprot

PDB

# General strategy

NewSequence

Similarity search (BLAST, FASTA, ...)

DB

Sequences with E-value BETTER than threshold

```
                                                              Score    E
Sequences producing significant alignments:                  (bits) Value

emb|CAB41881.1|    (Y18048) deoxynucleoside kinase [Drosophila m...    501  e-141
gb|AAF08104.1|AF105217_1    (AF105217) thymidine kinase 2 [Mus m...    190  1e-47
emb|CAC07190.1|    (AJ249341) mitochondrial thymidine kinase 2 [...    189  2e-47
dbj|BAB01587.1|    (AB046005) unnamed protein product [Macaca fa...    188  5e-47
emb|CAA71523.2|    (Y10498) thymidine kinase [Homo sapiens]           188  6e-47
sp|O00142|KITM_HUMAN    THYMIDINE KINASE 2, MITOCHONDRIAL >gi|19...    187  8e-47
gb|AAC51168.1|    (U80628) thymidine kinase 2 isoform B [Homo sa...    186  2e-46
ref|NP_031858.1|    deoxycytidine kinase >gi|1169273|sp|P43346|D...    129  2e-29
ref|NP_000779.1|    deoxycytidine kinase >gi|118447|sp|P27707|DC...    128  5e-29
gb|AAF14342.1|U90524_1    (U90524) deoxyguanosine kinase [Mus mu...    126  2e-28
sp|P48769|DCK_RAT    DEOXYCYTIDINE KINASE (DCK) >gi|508570|gb|AA...      126  2e-28
ref|NP_038792.1|    deoxyguanosine kinase >gi|4877287|emb|CAB431...    126  3e-28
emb|CAB43122.1|    (AJ133750) deoxyguanosine kinase 2 [Mus muscu...    126  3e-28
sp|Q16854|DGK_HUMAN    DEOXYGUANOSINE KINASE PRECURSOR (DGUOK) >...    125  4e-28
ref|NP_039114.1|    ORF FPV151 Deoxycytidine kinase >gi|7271649|...    125  4e-28
pir||S15_15    deoxyguanosine kinase (EC 2.7.1.113) precursor - ...    123  2e-27
ref|NP_001020.1|    deoxyguanosine kinase >gi|1480198|emb|CAA660...    122  2e-27
ref|NP_039072.1|    ORF FPV059 Deoxycytidine kinase >gi|140631|s...     87  2e-16
ref|NP_041094.1|    thymidine kinase >gi|9626904|ref|NP_041174.1...     70  2e-11
pir||T03086    probable thymidine kinase (EC 2.7.1.21) - Chilo i...     67  2e-10
pir||F75535    deoxyguanosine kinase/deoxyadenosine kinase subun...     51  9e-06
ref|NP_048773.1|    contains ATP/GTP-binding site motif A;  simi...     51  1e-05
gb|AAG10455.1|AF279106_17    (AF279106) predicted deoxypurine ki...     49  5e-05
sp|Q59484|DGK2_LACAC    DEOXYGUANOSINE KINASE (DGUO KINASE) (DGK...     44  0.002
gb|AAC97156.1|    (U49397) unknown [Streptococcus pyogenes]           44  0.002
```
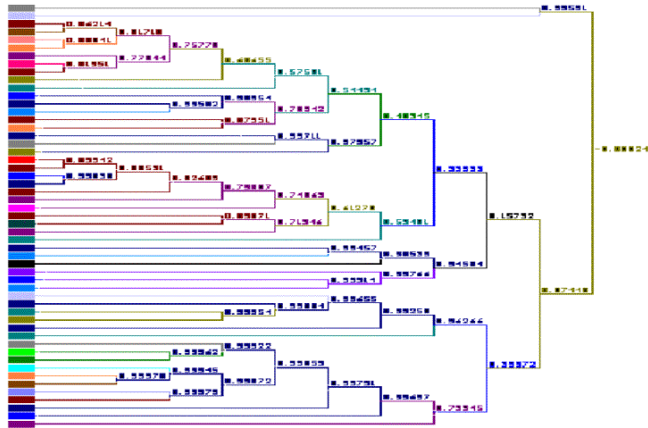
Score = 124 bits (309), Expect = 1e-27
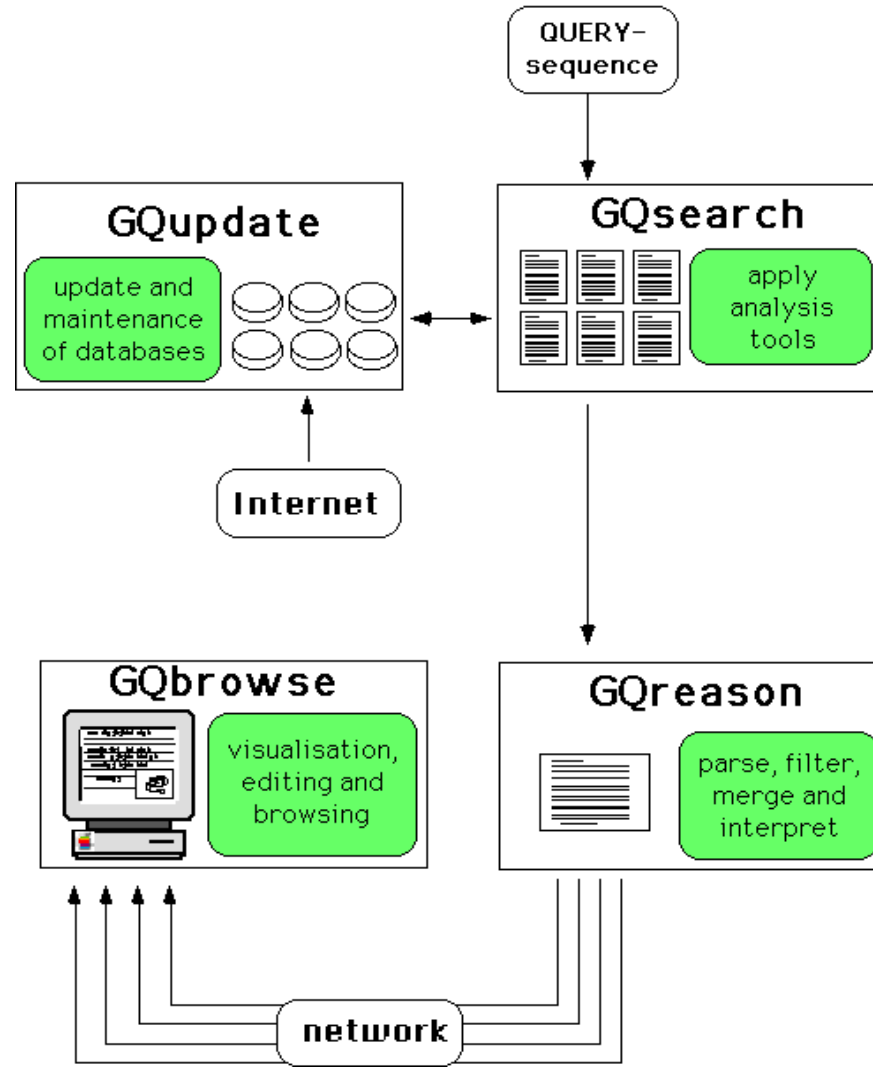Identities = 108/340 (31%), Positives = 156/340 (45%), Gaps = 62/340 (18%)

Query: 1   ATTYNAVVSKSSSDGKTFKTIADAIASAPAGSTP-FVILIKNGVYNERLTITRN--NLHL 5
           + T NAVV+   S   FKT+A A+A+AP G T  ++I IK GVY E + +T+   N+
Sbjct: 269 SVTPNAVVAADGSGN--FKTVAAAVAAAPQGGTKRYIIRIKAGVYRENVEVTKKHKNII
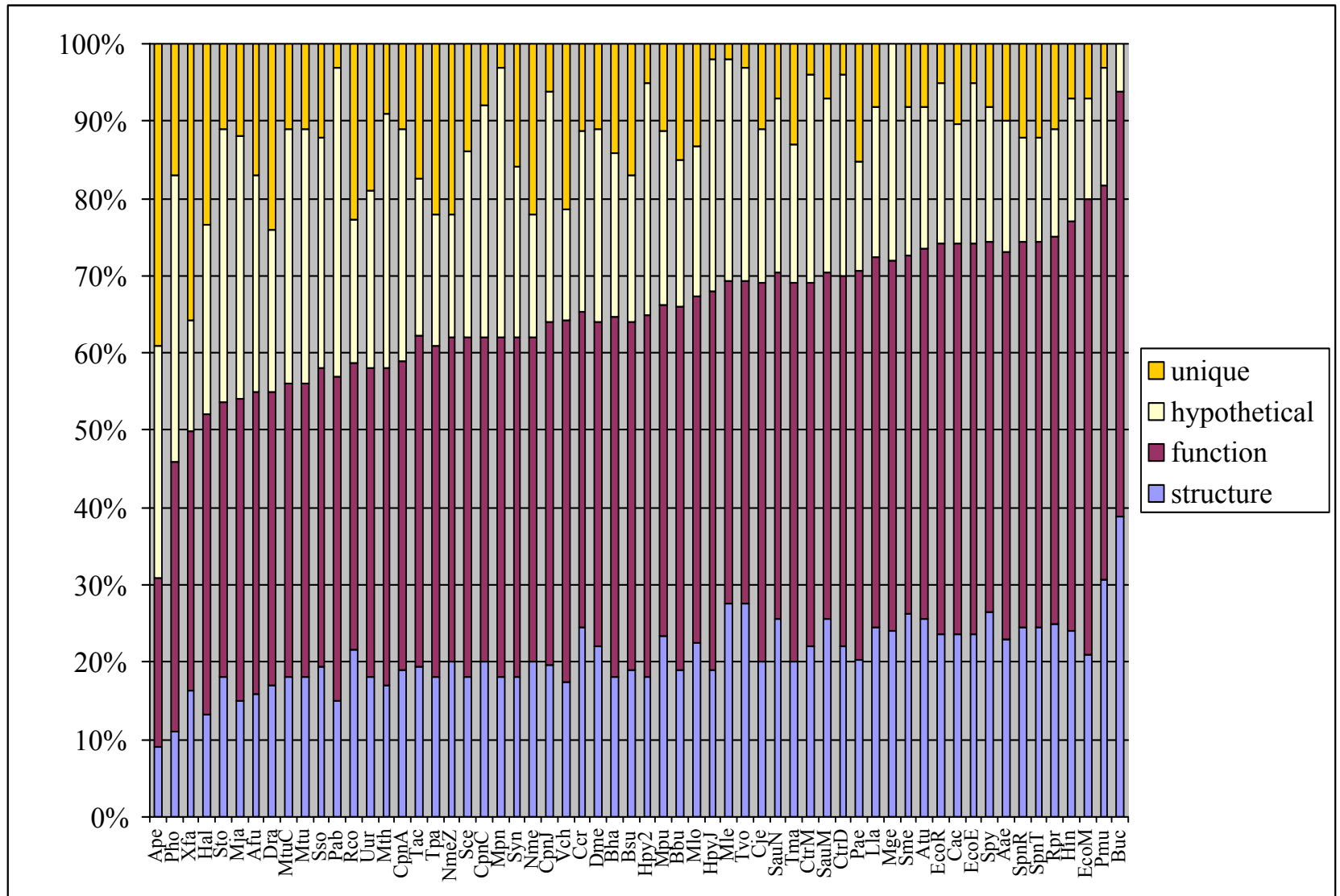
## Deduced homology

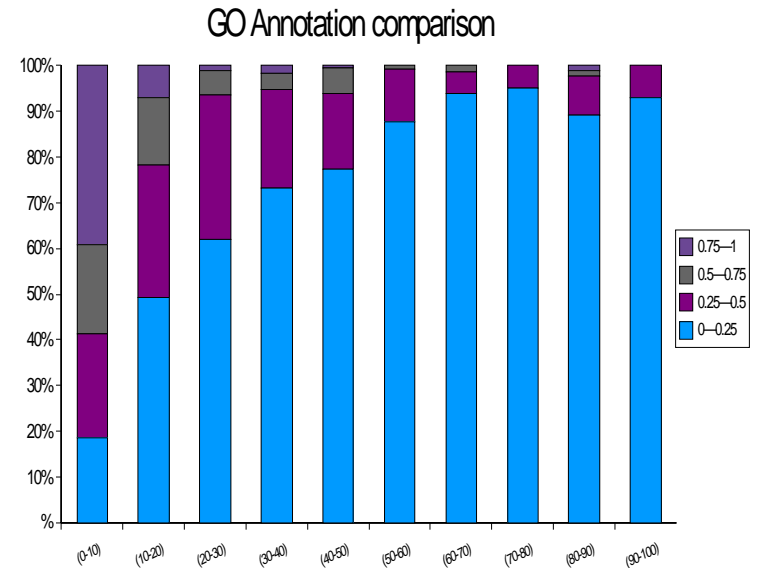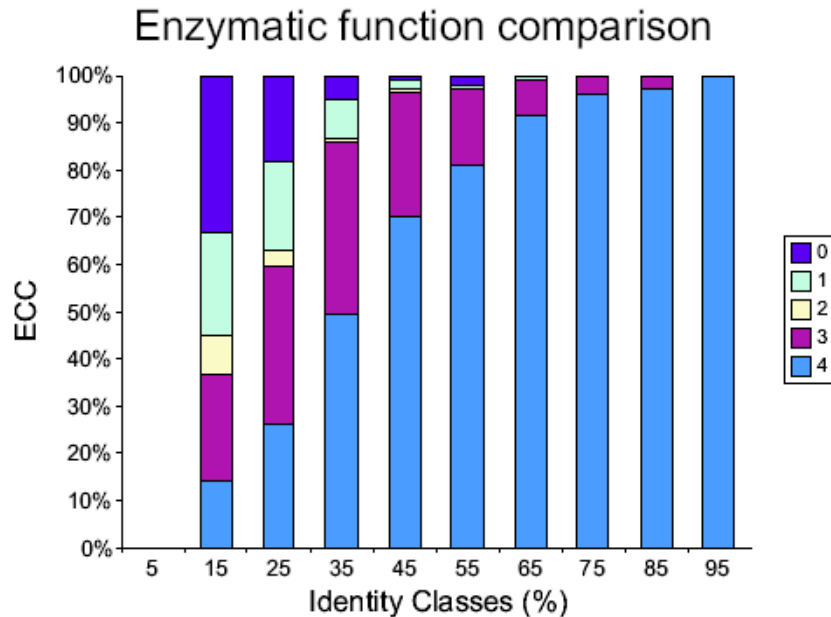# High throughput application - *GeneQuiz*

• Scharf, M., Schneider, R., Casari, G., Bork, P., Valencia, A., Ouzounis, C. and Sander, C. (1994) GeneQuiz: a workbench for sequence analysis. *Proc Int Conf Intell Syst Mol Biol.*, **2**, 348-353.
• Hoersch, S., Leroy, C., Brown, N.P., Andrade, M.A. and Sander, C. (2000) The GeneQuiz web server: protein functional analysis through the Web. *Trends Biochem Sci.*, **25**, 33-35.
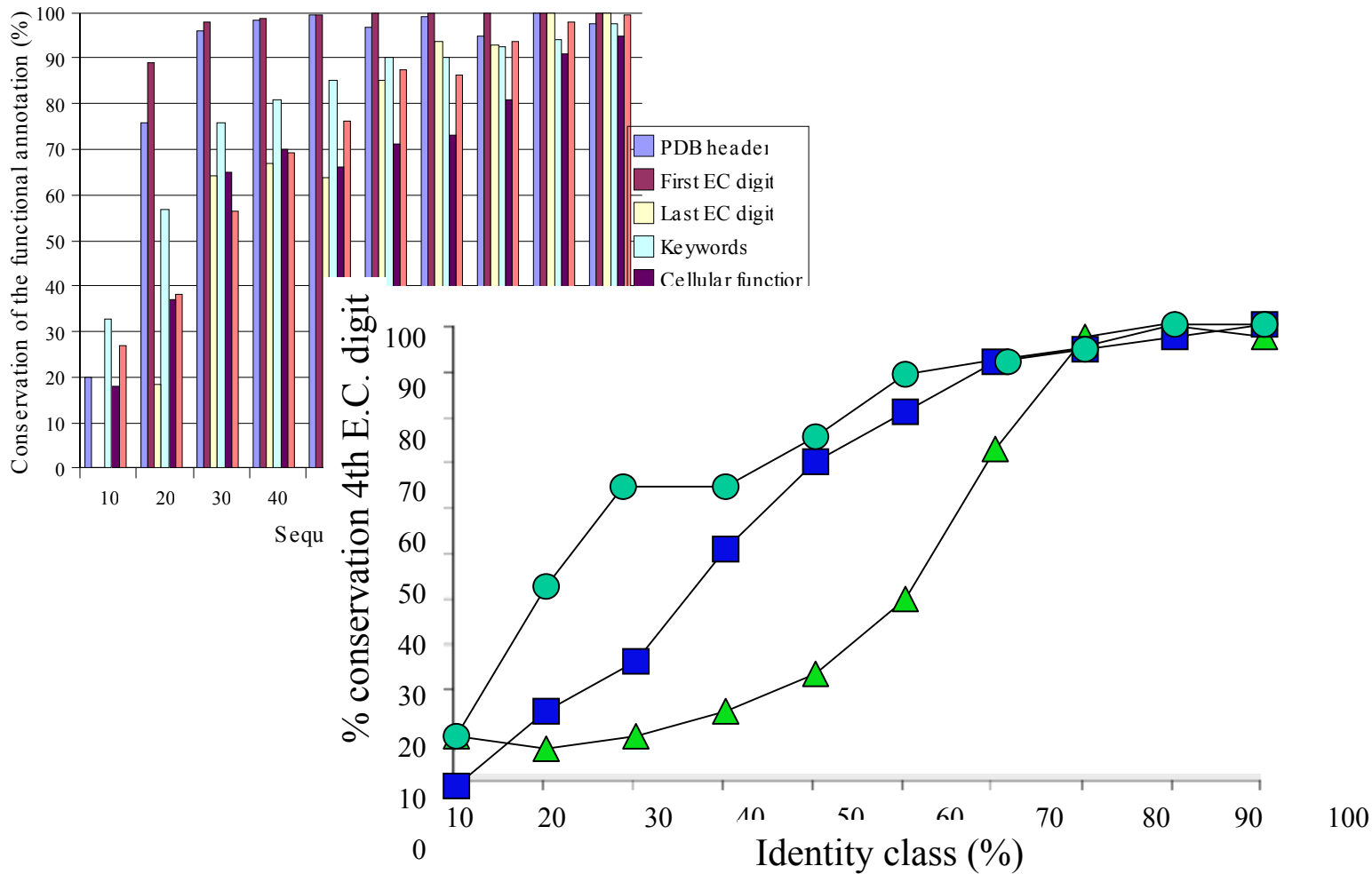
# High throughput application - *GeneQuiz*

# How reliable is the similarity-based Functional transfer?



Enzymatic function comparison

GO Annotation comparison

• Devos, D. and Valencia, A. (2000) Practical limits of function prediction. *Proteins*, **41**, 98-107.

• Valencia, A. (2005) Automatic annotation of protein function. *Curr Opin Struct Biol*, **15**, 267-274.

# How reliable is the similarity-based Functional transfer?

Valencia, A. (2005) Automatic annotation of protein function. *Curr Opin Struct Biol*, **15**, 267-274.

# So… Which degree of error can we expect?



| | PDB header | First EC Digit | Last EC Digit | Key- words | Funct. class | Binding site |
|---|---|---|---|---|---|---|
| **Mg** | 4 | 2 | 37 | 23 | 35 | 40 |
| **Hi** | 4 | 2 | 31 | 20 | 33 | 34 |
| **Mj** | 4 | 2 | 32 | 20 | 33 | 34 |

Devos, D. and Valencia, A. (2001) Intrinsic errors in genome annotation. *Trends Genet.*, **17**, 429-431.

# Environmental Genomics (Sargasso Sea)

Differentia l sequence

Sequences Found by

SS position in the sequence

| | Symptoms | Consequences |
|---|---|---|
| Sequence Comparisons | 1. 40% higher isoleucine, asparagine and lysine content<br><br>2. Sequences shorter and more fragments<br><br>3. Little overlap at 90% identity between current databases and Sargasso Sea | Less homologues found when searching Sargasso Sea resource with BLAST |
| Multiple Alignments and Families | 1. The distribution of sequences found by PSIBLAST differs between the Sargasso Sea and current databases<br><br>2. PSIBLAST profiles drift more<br><br>3. Profiles lose evolutionary information and decrease in quality | Worse annotation of function<br><br>Slightly worse deifintion of functional regions<br><br>Sequences ÑlostÓ from the profile in extreme cases |
| 3D Models | 1. More sequences found for alignments<br><br>2. Alignments vary in quality | Potentially worse 3D models |

Sargasso see sequences speaks trouble for biologists!

Difference in alignment correctness

*M. Tress*

# More complex strategies
## *FunCut*



ISS

**A. Homology Search**

Query Protein

**BLAST**

Sequence Space

**ISS**

**Rounds**

| 1 | 2 | 3 | 4 | 5 |
|---|---|---|---|---|

**Problem Protein Annotation (GO, EC, keywords, ...)**

**Neighbor Groups Annotation**

• Abascal, F. and Valencia, A. (2002) Clustering of proximal sequence space for the identification of protein families. *Bioinformatics.*, **18**, 908-921.
• Abascal, F. and Valencia, A. (2003) Automatic annotation of protein function based on family identification. *Proteins*, **53**, 683-692.

# Integrating annotation services

# Integrating annotation services



- **Utilises SOAP interfaces to simultaneously access:**

  ➢ **ProFunc (EBI, Hinxton, UK)**

  ➢ **CATH (UCL, London, UK)**

  ➢ **FUNcut (CNB, Madrid, Spain)**

  ➢ **STRING (EMBL, Heidelberg, Germany)**