# 6th Permanent European School on Bioinformatics

**http://bioinfo.cipf.es/6ESB/**

**Preliminary** program for the BioSapiens 6th ESB

*26th April to 30th of April 2007*

- **Day 1 (April 26th) Databases**
  Trainees from the European Bioinformatics Institute (EBI).

- **Day 2 (April 27th) Analysis of microarray data**
  Trainees to be confirmed (Brazma's group)

- **Day 3 (April 28th) Proteins and protein families**
  Trainee Dr. Yaniv Lowenstein (Linial's group)

- **Day 4 (April 29th) Protein structure prediction**
  Trainees to be confirmed (Marti-Renom's group)

- **Day 5 (April 30th) Systems biology**
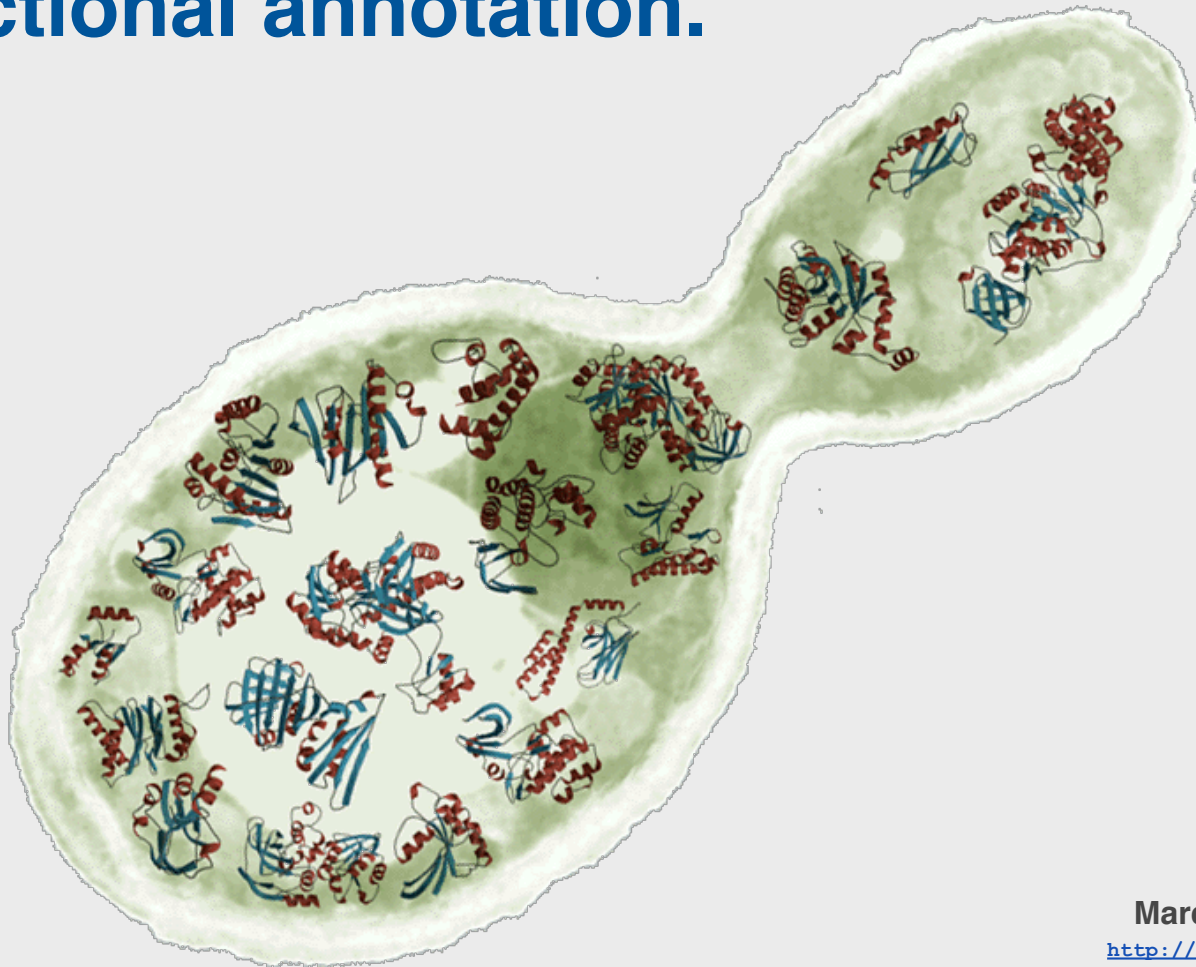  Trainee Dr. Ildefonso Cases (Valencia's groups)

Bioinformatics Department
Prince Felipe Resarch Center (CIPF), Valencia, Spain

# Comparative protein structure models for functional annotation.
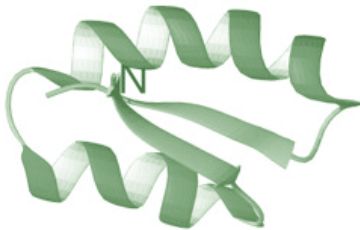
**Marc A. Marti-Renom**
http://bioinfo.cipf.es/sgu/

Structural Genomics Unit
Bioinformatics Department
Prince Felipe Resarch Center (CIPF), Valencia, Spain

PRINCIPE FELIPE
CENTRO DE INVESTIGACION

# Principles of protein structure

GFCHIKAYTRLIMVG...

Desulfovibrio vulgaris

Condrus crispus

GFCHIKAYTRLIMVG...

Anabaena 7120

Anacystis nidulans

Folding (physics)

*Ab initio* prediction

Evolution (rules)

Threading
Comparative Modeling

D. Baker & A. Sali. Science 294, 93, 2001.

3

# From domains to assemblies

domains

proteins

assemblies

~2.5 domains in a protein
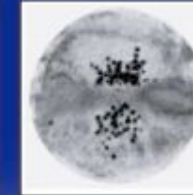a few domain partners per domain

# Determining the structures of proteins and assemblies

Use structural information from any
        source: measurement, first principles, rules,
        resolution: low or high resolution
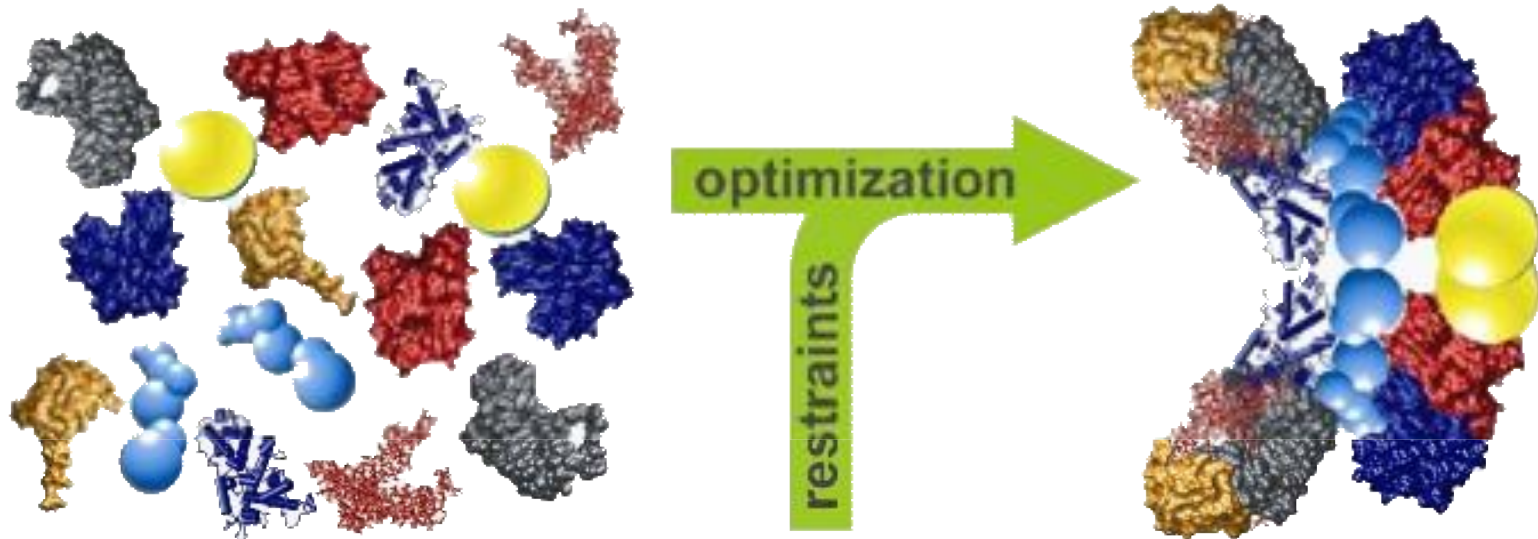to obtain the set of all models that are consistent with it.



| X-ray crystallography | NMR spectroscopy | 2D & single particle electron microscopy | electron tomography | immuno-electron microscopy | chemical cross-linking | affinity purification mass spectroscopy |
|---|---|---|---|---|---|---|
| subunit structure | subunit structure | | | | subunit structure | |
| subunit shape | subunit shape | subunit shape | subunit shape | | | |
| subunit-subunit contact | subunit-subunit contact | subunit-subunit contact | subunit-subunit contact | | subunit-subunit contact | subunit-subunit contact |
| subunit proximity | subunit proximity | subunit proximity | subunit proximity | subunit proximity | subunit proximity | subunit proximity |
| subunit stoichiometry | subunit stoichiometry | | | | | |
| assembly symmetry | assembly symmetry | assembly symmetry | assembly symmetry | assembly symmetry | | |
| assembly shape | assembly shape | assembly shape | assembly shape | | | |
| assembly structure | assembly structure | | | | | |

| FRET | site-directed mutagenesis | yeast two-hybrid system | gene/protein arrays | protein structure prediction | computational docking | bioinformatics |
|---|---|---|---|---|---|---|
| | | | | subunit structure | | |
| | | | | subunit shape | | |
| subunit-subunit contact | subunit-subunit contact | subunit-subunit contact | subunit-subunit contact | | subunit-subunit contact | subunit-subunit contact |
| subunit proximity | | subunit proximity | subunit proximity | | | |

# Modeling by satisfaction of spatial restraints
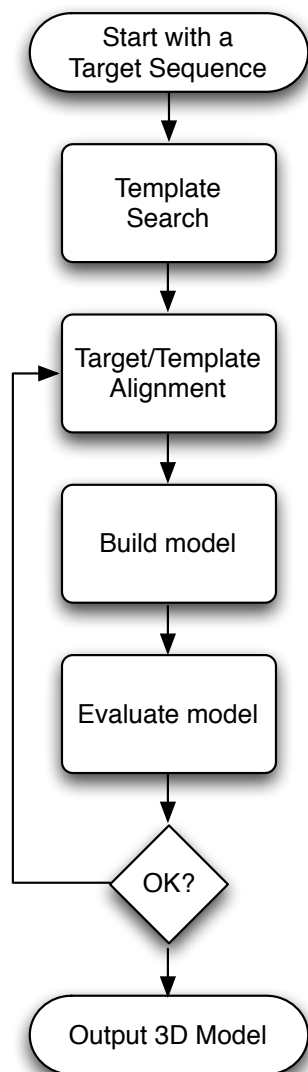
1) Representation of a system.
2) Scoring function (spatial restraints).
3) Optimization.

There is nothing but points and restraints on them.



optimization

restraints

*Alber et al. Structure, 13, 435 (2005)*

# Comparative modeling by satisfaction of spatial restraints
## MODELLER

Start with a Target Sequence

Template Search

Target/Template Alignment

Build model

Evaluate model
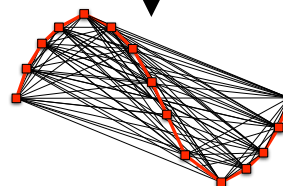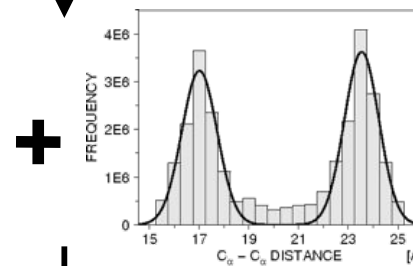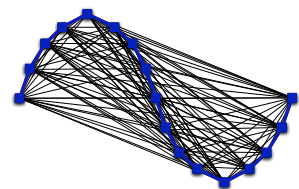
OK?

Output 3D Model

**Given an alignment...**

**extract spatial features from the template(s) and statistics from known structures**

**apply these features as restraints on your target sequence**

**optimize to find the best solution for the restraints to produce your 3D model**

MSVIPKR--GNCEQTSE
ASILPKRLFGNCEQTSD

+

FREQUENCY

4E6
3E6
2E6
1E6
0

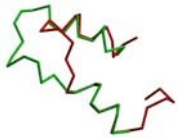15   17   19   21   23   25

$C_\alpha - C_\alpha$ DISTANCE        [Å]

A. Šali & T. Blundell. J. Mol. Biol. 234, 779, 1993.
J.P. Overington & A. Šali. Prot. Sci. 3, 1582, 1994.
A. Fiser, R. Do & A. Šali, Prot. Sci., 9, 1753, 2000.

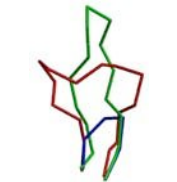# Comparative modeling by satisfaction of spatial restraints
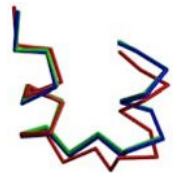## Types of errors and their impact
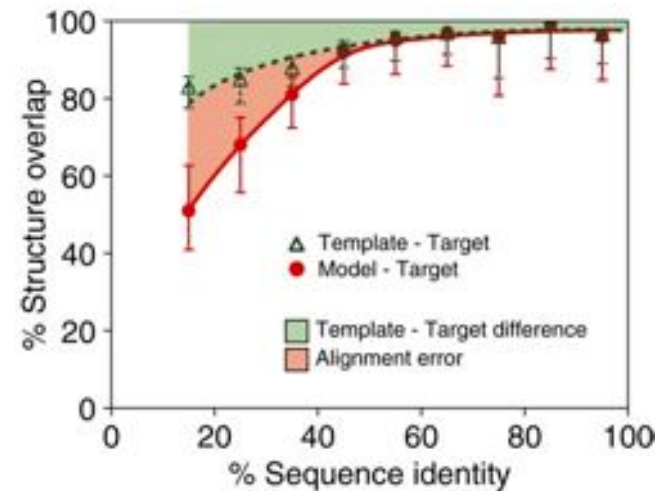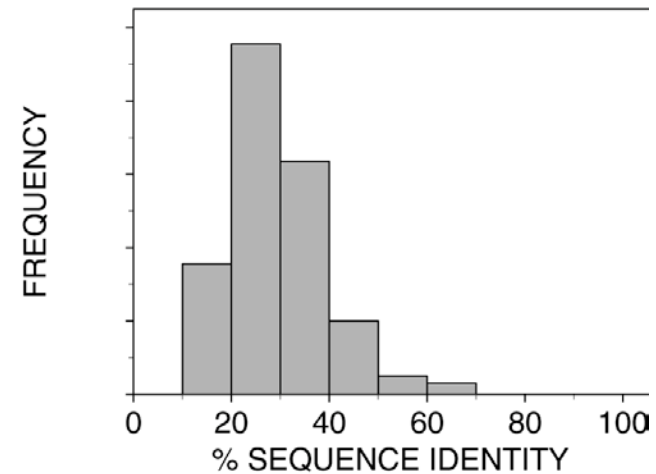
**Wrong fold**

**Miss alignments**

**Loop regions**

**Rigid body distortions**

**Side-chain packing**



*Marti-Renom et al. Ann Rev Biophys Biomol Struct (2000) 29, 291*

# ModBase Statistics

## Large-scale modeling of the TrEMBL-SWISSPROT databases

http://www.salilab.org/modbase/

| Sequences (total) | 1,930,692 |
|---|---|
| Sequences (modeled) | 1,084,784 |
| Models | 3,094,542 |



UCSF
University of California
San Francisco

*Pieper et al. NAR 34, D291 (2006)*

# Utility of protein structure models, despite errors

# For many protein structures function is *unknown*

| | Structural Genomics* | Traditional methods |
|---|---|---|
| **Annotated\*\*** | 654 | 28,342 |
| **Not Annotated** | 506 (43.6%) | 6,815 (19,4%) |
| **Total deposited** | 1,160 | 35,157 |

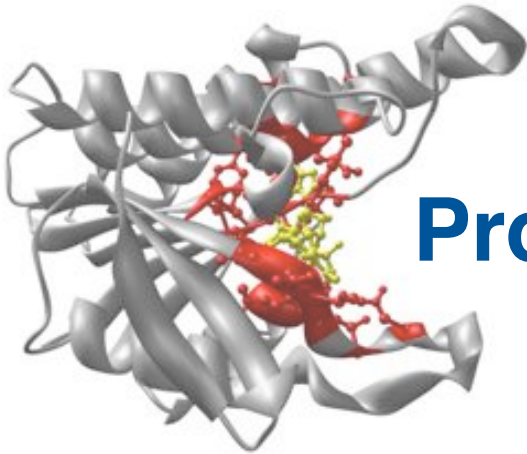*\* annotated as STRUCTURAL GENOMICS in the header of the PDB file*
*\*\*annotated with either CATH, SCOP, Pfam or GO terms in the MSD database*
*36,317 protein structures, as of August 8th, 2006*

# For **20%** protein structures function is *unknown*

| | Structural Genomics* | Traditional methods |
|---|---|---|
| **Annotated**** | 654 | 28,342 |
| **Not Annotated** | 506 (43.6%) | 6,815 (19,4%) |
| **Total deposited** | 1,160 | 35,157 |

*\* annotated as STRUCTURAL GENOMICS in the header of the PDB file*
*\*\*annotated with either CATH, SCOP, Pfam or GO terms in the MSD database*
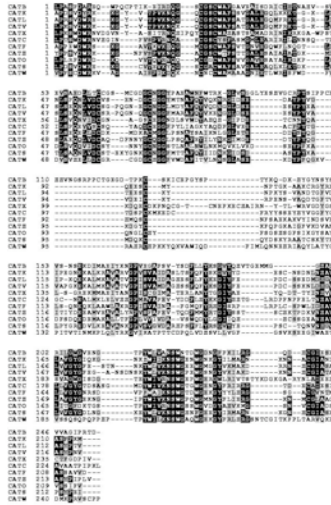*36,317 protein structures, as of August 8th, 2006*

# Protein function from structure
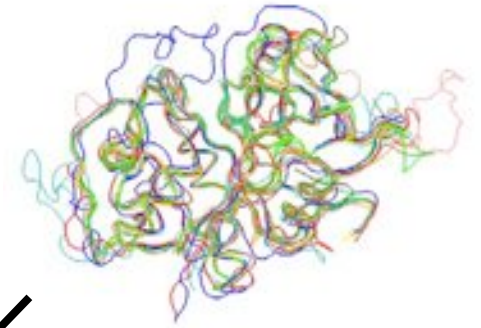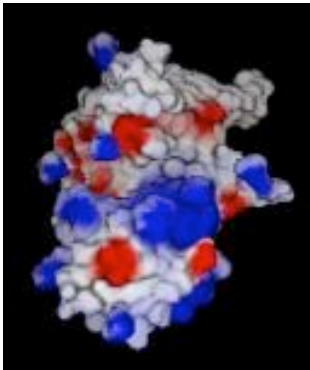## ab-initio *localization of binding sites*

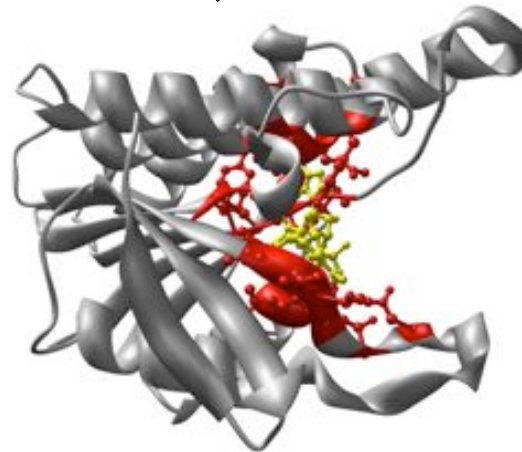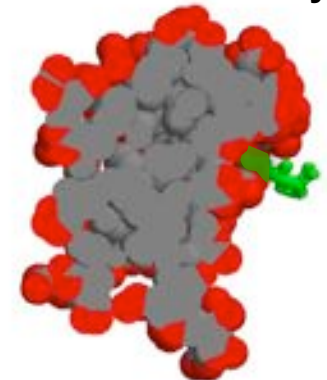# Representation

**Sequence conservation**

**Surface geometry**

**Structure conservation**
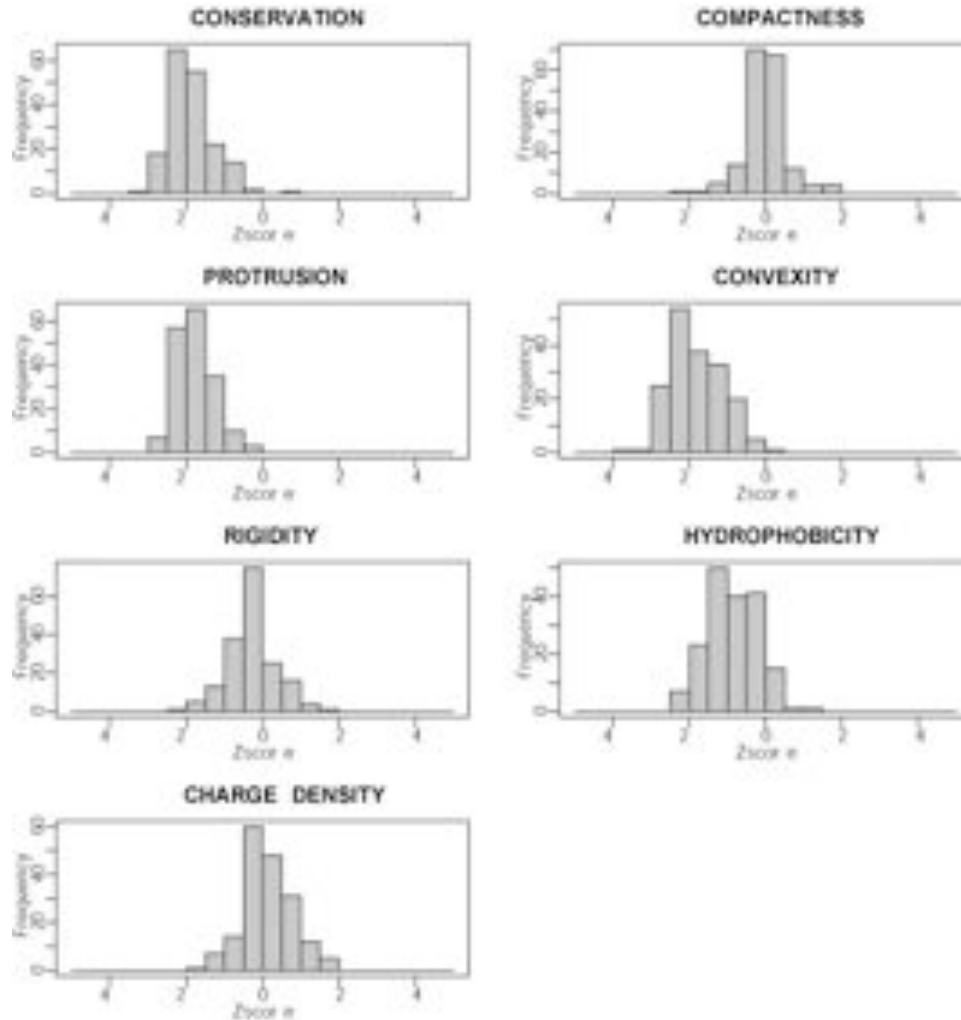
**Electrostatics**

**Solvent accessibility**
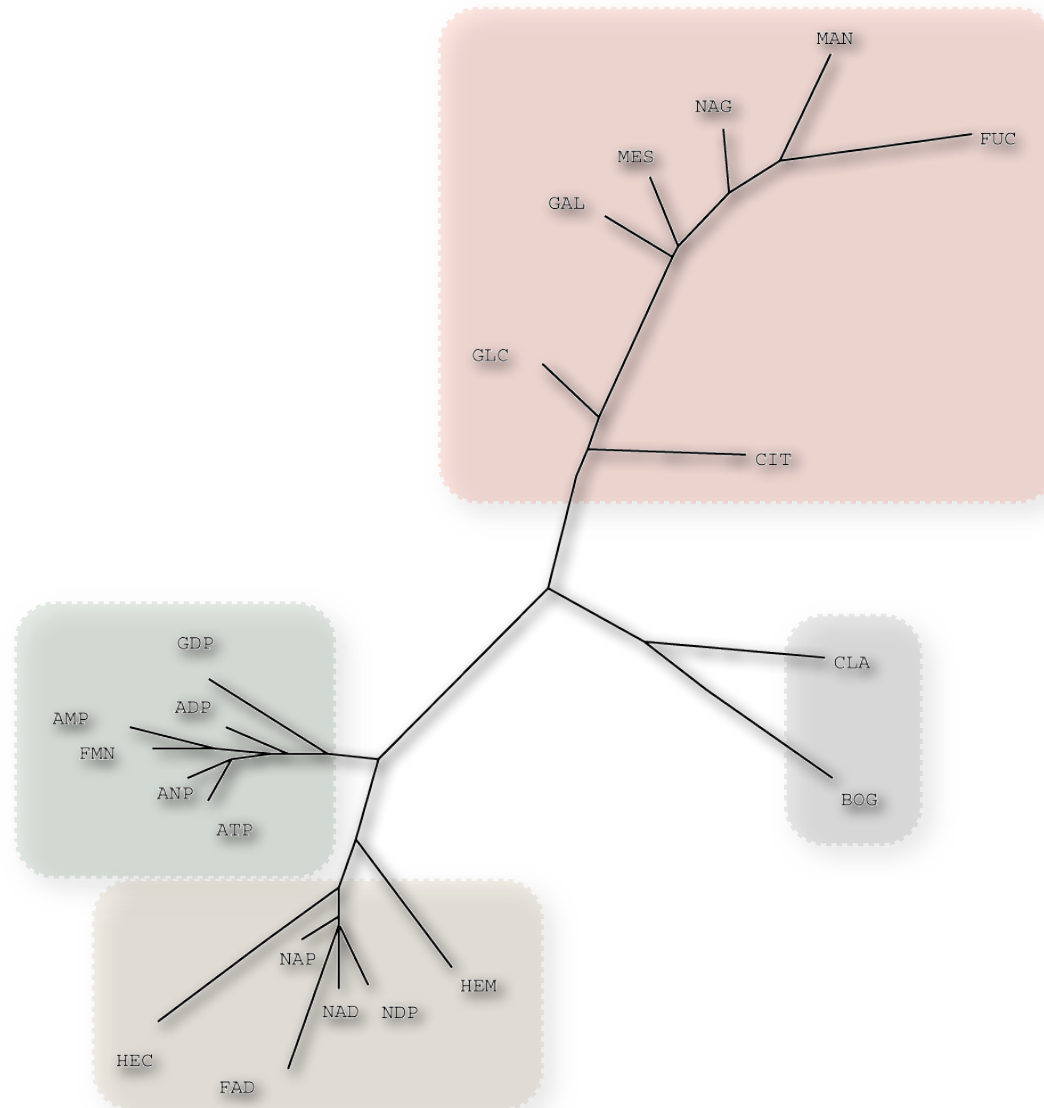
# Scoring

## NAD



$$w_k = \frac{1}{M} \sum_{\alpha=1}^{M} \tilde{f}_k^{(\alpha)}$$
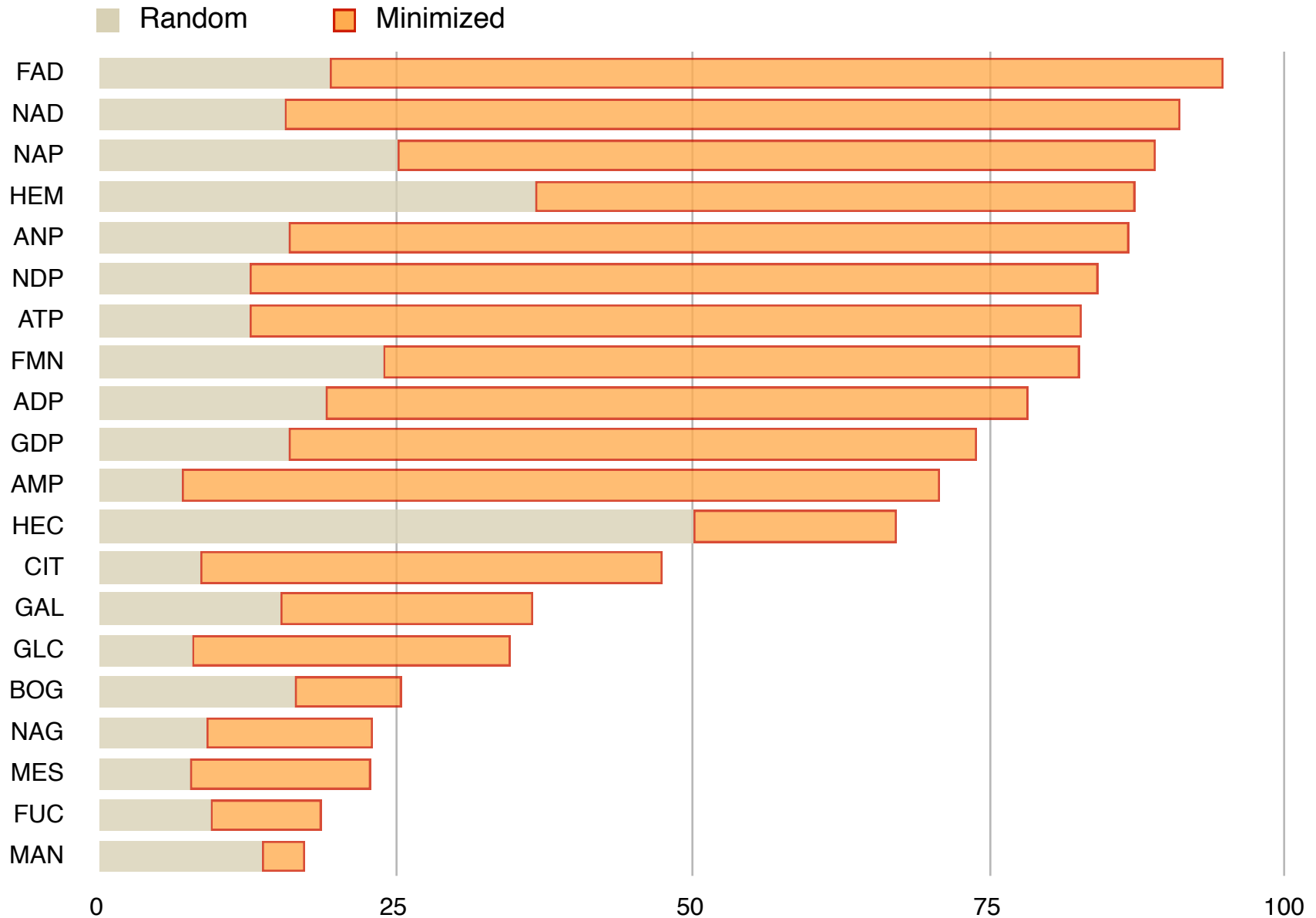
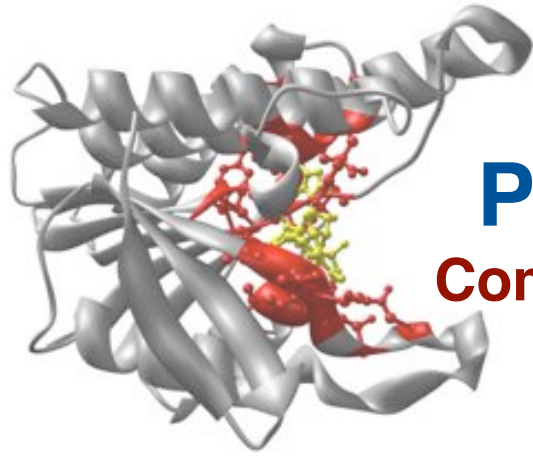*M* = number of proteins in training set

# Ligand fingerprints

| | Compactness | Conservation | Charge density | B-factor | Protrusion coefficient | Convexity score | Hydrophobicity |
|---|---|---|---|---|---|---|---|
| ADP | -1.266 | -2.009 | 0.447 | -0.414 | -1.521 | -1.388 | -0.118 |
| AMP | -1.62 | -1.962 | 0.341 | -0.381 | -1.909 | -1.944 | -0.518 |
| ANP | -1.007 | -2.227 | 0.176 | -0.392 | -1.706 | -1.595 | -0.14 |
| ATP | -1.122 | -2.156 | 0.228 | -0.274 | -1.845 | -1.768 | 0.038 |
| BOG | -2.067 | -0.012 | 0.552 | -0.465 | -0.356 | -0.49 | -0.781 |
| CIT | -2.948 | -1.58 | 0.563 | -0.527 | -0.922 | -0.838 | -0.113 |
| FAD | 0.505 | -2.108 | 0.366 | -0.702 | -1.735 | -1.725 | -0.75 |
| FMN | -1.132 | -1.98 | 0.382 | -0.387 | -1.803 | -1.886 | -0.695 |
| FUC | -3.43 | 0.016 | -0.295 | -0.123 | 0.002 | 0.132 | 0.459 |
| GAL | -3.186 | -0.538 | -0.234 | -0.068 | -0.906 | -0.987 | 0.298 |
| GDP | -1.061 | -1.471 | 0.409 | -0.81 | -1.472 | -1.423 | 0.182 |
| GLC | -2.813 | -1.247 | -0.207 | -0.399 | -1.247 | -1.337 | -0.089 |
| HEC | -0.172 | -0.912 | 0.286 | -0.325 | -1.153 | -1.27 | -1.282 |
| HEM | -0.651 | -1.571 | 0.683 | -0.51 | -1.797 | -1.937 | -1.47 |
| MAN | -3.72 | 0.131 | 0.105 | -0.52 | -0.605 | -0.509 | 0.405 |
| MES | -3.049 | -0.24 | -0.338 | -0.479 | -0.714 | -0.926 | 0.296 |
| NAD | -0.005 | -1.852 | 0.156 | -0.232 | -1.775 | -1.804 | -0.858 |
| NAG | -3.419 | -0.46 | -0.126 | -0.154 | -0.341 | -0.523 | -0.078 |
| NAP | -0.009 | -1.898 | 0.612 | -0.321 | -1.587 | -1.656 | -0.336 |
| NDP | 0.217 | -1.741 | 0.535 | -0.312 | -1.463 | -1.562 | -0.498 |

15

# Ligand fingerprints
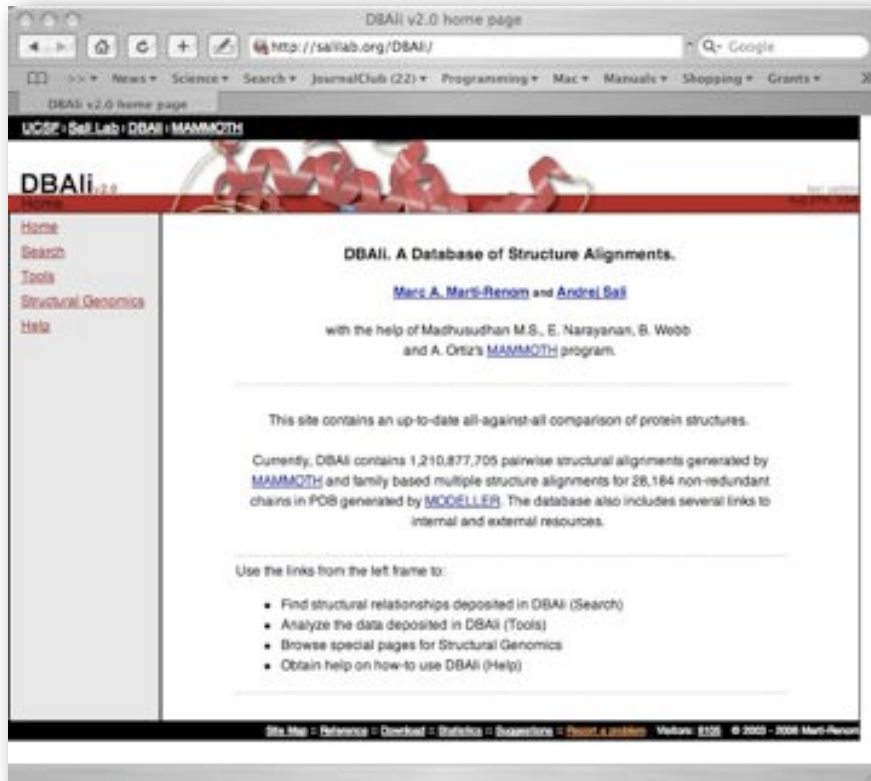
Prediction accuracy

# Protein function from structure
## Comparative annotation. AnnoLite and AnnoLyze.

# DBAli<sub>v2.0</sub> database

- ✓ **Fully-automatic**
- ✓ **Data is kept up-to-date with PDB releases**
- ✓ **Tools for "on the fly" classification of families.**
- ✓ **Easy to navigate**
- ✓ **Provides tools for structure analysis**

**Does not provide a stable classification similar to that of CATH or SCOP**

| Pairwise structure alignments | |
|---|---|
| Last update: | February 15th, 2007 |
| Number of chains: | 88,276 |
| Number of structure-structure comparisons: | 1,425,479,365 |
| Multiple structure alignments | |
| Last update: | January 23rd, 2007 |
| Number of representative chains: | 30,900 |
| Number of families: | 11,615 |

Uses MAMMOTH for similarity detection

- ✓ **VERY FAST!!!**
- ✓ **Good scoring system with significance**

*Ortiz AR, (2002) Protein Sci. 11 pp2606*

*Marti-Renom et al. 2001. Bioinformatics. 17, 746*

# DBAli v2.0 database

http://bioinfo.cipf.es/squ/services/DBAli/
http://www.salilab.org/DBAli/

# DBAli$_{v2.0}$ database

# AnnoLite

# Benchmark set

| | Number of chains |
|---|---|
| **Initial set*** | 50,223 |
| **FULL annotation**** | 10,997 |
| **Non-redundant set***** | 1,879 |

*data from BioMart  MSD.3 (release February 2005)
**annotated with CATH, SCOP, Pfam, EC, InterPro, and GO terms in the MSD database
**not two chains can be structurally aligned  within 2A, superimposing more than 60% of
their C  atoms and have a length difference inferior to 30aa

# Method



**DBAIi tools**

Chain ID

AnnoLite search

**Similar chains in DBAIi**

RMSD < 4A
% Seq Id *variable* (>15)
% Equivalent positions >75%
p-value >4

**BioMart protein annotation**

Annotations from MSD.msd database and descriptions from SCOP, CATH, InterPro, PFamA, ENZYME, and GO databases

Fischer´s 2x2 test for statistical significance

HTML output

# Scoring function

Fisher's 2x2 contingency test

| | Non-similar | Similar | Total |
|---|---|---|---|
| **Annotated** | a | b | a+b |
| **Not Annotated** | c | d | c+d |
| **Total** | a+c | b+d | n |

| 1b78A SCOP c.51.4.1 | Similar | Not similar | Total |
|---|---|---|---|
| **Annotated** | 4 | 2 | 6 |
| **Not Annotated** | 0 | 71,096 | 71,096 |
| **Total** | 4 | 71,098 | 71,102 |

$$p = \binom{a+b}{a}\binom{c+d}{c} \Big/ \binom{n}{a+c}$$

$$= \frac{(a+b)!(c+d)!(a+c)!(b+d)!}{n!\,a!\,b!\,c!\,d!}$$

$$p = 1.78e^{-19}$$

24

# Sensitivity .vs. Precision

| | Optimal cut-off | Sensitivity (%) Recall or TPR | Precision (%) |
|---|---|---|---|
| SCOP fold | 1e-6 | 92.7 | 88.4 |
| CATH fold | 1e-3 | 95.7 | 90.1 |
| InterPro | 1e-3 | 88.4 | 78.2 |
| PFam family | 1e-4 | 90.5 | 82.8 |
| EC number | 1e-4 | 93.3 | 79.7 |
| GO Molecular Function | 1e-1 | 84.3 | 80.9 |
| GO Biological Process | 1e-3 | 85.5 | 74.8 |
| GO Cellular Component | 1e-2 | 77.6 | 58.6 |

$$\text{Sensitivity} = \frac{TP}{TP + FN} \qquad \text{Precision} = \frac{TP}{TP + FP}$$

# AnnoLyze

# Benchmark

|  | Number of chains |
|---|---|
| **Initial set\*** | 78,167 |
| **LigBase\*\*** | 30,126 |
| **Non-redundant set\*\*\*** | 4,948 (8,846 ligands) |

*\*all PDB chains larger than 30 aminoacids in length (8th of August, 2006)*
*\*\*annotated with at least one ligand in the LigBase database*
*\*\*\*not two chains can be structurally aligned within 3A, superimposing more than 75% of their C atoms, result in a sequence alignment with more than 30% identity, and have a length difference inferior to 50aa*

|  | Number of chains |
|---|---|
| **Initial set\*** | 78,167 |
| **πBase\*\*** | 30,425 |
| **Non-redundant set\*\*\*** | 4,613 (11,641 partnerships) |

*\*all PDB chains larger than 30 aminoacids in length (8th of August, 2006)*
*\*\*annotated with at least one partner in the Base database*
*\*\*\*not two chains can be structurally aligned within 3A, superimposing more than 75% of their C atoms, result in a sequence alignment with more than 30% identity, and have a length difference inferior to 50aa*

# Method

# Scoring function

Ligands

Partners



Aloy *et al.* (2003) J.Mol.Biol. 332(5):989-98.

# Sensitivity .vs. Precision

| | Optimal cut-off | Sensitivity (%) Recall or TPR | Precision (%) |
|---|---|---|---|
| **Ligands** | 30% | 71.9 | 13.7 |
| **Partners** | 40% | 72.9 | 55.7 |

$$\text{Sensitivity} = \frac{TP}{TP + FN} \qquad \text{Precision} = \frac{TP}{TP + FP}$$

# Example (2azwA)
## *Structural Genomics Unknown Function*

Molecule: MutT/nudix family protein

# Can we use models to infer function?



*T. cruzi*

# What is the physiological ligand of Brain Lipid-Binding Protein?

Predicting features of a model that are not present in the template

BLBP/oleic acid

Cavity is not filled

Ligand binding cavity

BLBP/docosahexaenoic acid

Cavity is filled

1. BLBP binds fatty acids.

2. Build a 3D model.

3. Find the fatty acid that fits most snuggly into the ligand binding cavity.

L. Xu, R. Sánchez, A. Šali, N. Heintz, J. Biol. Chem. 271, 24711, 1996.

# Structural analysis of missense mutations in human BRCA1 BRCT domains

Nebojsa Mirkovic, Marc A. Marti-Renom, Barbara L. Weber, Andrej Sali and Alvaro N.A. Monteiro

**Cancer Research (June 2004). 64:3790-97**



Cannot measure the functional impact of every possible SNP at all positions in each protein! Thus, prediction based on general principles of protein structure is needed.

# Missense mutations in BRCT domains by function

|  | cancer associate | not cancer associated | ? |
|---|---|---|---|
| **no transcription activation** | C1697R R1699W A1708E S1715R P1749R M1775R | | M1652K L1657P E1660G H1686Q R1699Q K1702E Y1703HF 1704S    L1705PS 1715NS1 722FF17 34LG173 8EG1743 RA1752 PF1761I    F1761S M1775E M1775K L1780P I1807S V1833E A1843T |
| **transcription activation** | | M1652I A1669S | V1665M D1692N G1706A D1733G M1775V P1806A |
| **?** | | | M1652T V1653M L1664P T1685A T1685I M1689R D1692Y F1695L V1696L R1699L G1706E W1718C    W1718S T1720A W1730S F1734S E1735K V1736A G1738R D1739E D1739G D1739Y V1741G H1746N    R1751P R1751Q R1758G L1764P I1766S P1771L T1773S P1776S D1778N D1778G D1778H M1783T    C1787S G1788D G1788V G1803A V1804D V1808A V1809A V1809F V1810G Q1811R P1812S N1819S    A1823T V1833M W1837R W1837G S1841N A1843P T1852S P1856T P1859R |



35

# Putative binding site on BRCA1



Putative binding site predicted in 2003
and accepted for publication on March 2004.

Williams *et al.* 2004 Nature Structure Biology. **June 2004 11**:519

Mirkovic *et al.* 2004 Cancer Research. **June 2004 64**:3790

# *S. cerevisiae* ribosome



Fitting of comparative models into 15Å cryo-electron density map.

43 proteins could be modeled on 20-56% seq.id. to a known structure.

The modeled fraction of the proteins ranges from 34-99%.

C. Spahn, R. Beckmann, N. Eswar, P. Penczek, A. Sali, G. Blobel, J. Frank. Cell 107, 361-372, 2001.

37

# Modeling & cryoEM



Topf etal. JMB, 357, 1655 (2006)

# The Nucleopore complex Cell evolution (?)



Prokaryote · Early Eukaryote · Modern Eukaryote

Nup84 · Nup85 · Nup145C · Nup120 · Nup133

# Tropical Disease Initiative (TDI)
*Predicting binding sites in protein structure models.*



**http://www.tropicaldisease.org**

# Need is High in the Tail

■ DALY Burden Per Disease in Developed Countries

■ DALY Burden Per Disease in Developing Countries

Heart diseases

Rare diseases

DALY

Disease

# Need is High in the Tail



DALY Burden Per Disease in Developed Countries
DALY Burden Per Disease in Developing Countries

Heart diseases

Rare diseases

DALY

Disease

# "Unprofitable" Diseases and Global DALY (in 1000's)

| | | | | |
|---|---|---|---|---|
| **Malaria*** | **46,486** | | Trichuriasis | 1,006 |
| Tetanus | 7,074 | | Japanese encephalitis | 709 |
| **Lymphatic filariasis*** | **5,777** | | **Chagas Disease*** | **667** |
| Syphilis | 4,200 | | **Dengue*** | **616** |
| Trachoma | 2,329 | | **Onchocerciasis*** | **484** |
| **Leishmaniasis*** | **2,090** | | **Leprosy*** | **199** |
| Ascariasis | 1,817 | | Diphtheria | 185 |
| **Schistosomiasis*** | **1,702** | | Poliomyelitise | 151 |
| **Trypanosomiasis*** | **1,525** | | Hookworm disease | 59 |

Disease data taken from WHO, *World Health Report 2004*
DALY - Disability adjusted life year in 1000's.
*  Officially listed in the WHO Tropical Disease Research disease portfolio.

# TDI flowchart

*Sali, Rai, Maurer. PLoS Medicine (2004)*
*Kepler, et al. Australian Journal of Chemistry (2006)*

# Modeling Genomes

*data from models generated by ModPipe (Eswar, Pieper & Sali)*

*A good model has MPQS of 1.1 or higher*

44

# Comparative docking



**1. Expansion**

co-crystalized protein/ligand

**2. Inheritance**

model

crystalized protein

template

# Ligand "expanded" space

**from 6,859 templates used in "good" models**

| Expansion cut-off | Templates | Expanded | Unique |
|---|---|---|---|
| 30% | 4,639 | 64,800 | 3,178 |
| 50% | 4,242 | 37,945 | 3,030 |
| 70% | 3,323 | 20,603 | 2,786 |

# Ligand "inherited" space

second cut-offs

**Using a 70% "expansion" cut-off**

| Inheritance cut-offs | Models | Inherited | Unique |
|---|---|---|---|
| **90% / 70%** | 5,181 | 23,286 | 1,137 |
| **90% / 80%** | 4,383 | 17,842 | 1,027 |
| **90% / 90%** | 3,462 | 11,803 | 827 |

# Distribution of models with inherited ligands

## from 3,882 "good" models
## using a 90% / 90% "inherited" cut-offs



| | | |
|---|---|---|
| ■ | C.hominis | 183 |
| ■ | C.parvum | 219 |
| ■ | L.major | 488 |
| ■ | M.leprae | 286 |
| ■ | M.tuberculosis | 404 |
| ■ | P.falciparum | 271 |
| ■ | P.vivax | 267 |
| ■ | T.brucei | 440 |
| ■ | T.cruzi | 730 |
| ■ | T.gondii | 174 |

48

# Summary table

models with inherited ligands

**from 16,284 good models, 295 inherited a ligand/substance with at least one compound already approved by FDA and ready to be used from ZINC**

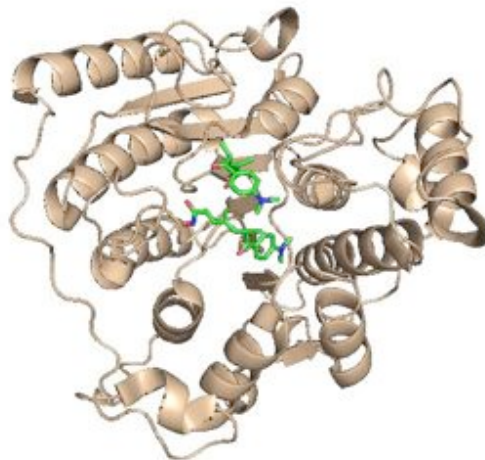|  | Transcripts | Good | Ligands | Lipinski | Lipinski+ZINC | FDA+ZINC |
|---|---|---|---|---|---|---|
| *C. hominis* | 3,886 | 886 | 183 | 131 | 28 | 12 (10) |
| *C. parvum* | 3,806 | 949 | 219 | 145 | 30 | 12 (10) |
| *L. major* | 8,274 | 1,845 | 488 | 334 | 84 | 44 (34) |
| *M. leprae* | 1,605 | 1,321 | 286 | 189 | 39 | 29 (25) |
| *M. tuberculosis* | 3,991 | 2,887 | 404 | 285 | 71 | 44 (37) |
| *P. falciparum* | 5,363 | 1,057 | 271 | 191 | 48 | 20 (16) |
| *P. vivax* | 5,342 | 1,042 | 267 | 177 | 37 | 18 (15) |
| *T. brucei* | 921 | 1,795 | 440 | 309 | 94 | 46 (36) |
| *T. cruzi* | 19,607 | 3,915 | 730 | 493 | 127 | 62 (52) |
| *T. gondii* | 7,793 | 587 | 174 | 124 | 28 | 8 (7) |
| **TOTAL** | **60,588** | **16,284** | **3,462** | **2,378** | **586** | **295 (242)** |

# Example of inheritance (expansion)

*LmjF21.0680 from* L. major *"Histone deacetylase 2" (model 1)*
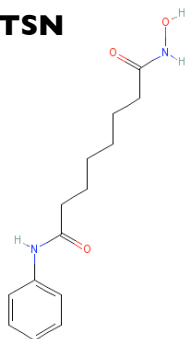
**Template 1t64A a human HDAC8 protein.**



|  | Origen | Formula | Name | Cov. | Seq, Id. (%) |
|---|---|---|---|---|---|
| **ZN** | X-ray | $Zn^{2+}$ | Zinc ion | -- | -- |
| **NA** | X-ray | $Na^+$ | Sodium ion | -- | -- |
| **CA** | X-ray | $Ca^{2+}$ | Calcium ion | -- | -- |
| **TSN** | X-ray | $C_{17} H_{22} N_2 O_3$ | Trichostatin A | -- | -- |
| **SHH** | Expanded | $C_{14} H_{20} N_2 O_3$ | Octadenioic acid hudroxyamide phenylamide | 100.00 | 83.8 |

# Example of inheritance (inheritance)

*LmjF21.0680 from L. major "Histone deacetylase 2" (model 1)*

| | Formula | Name | Cov. | Seq, Id. (%) | Residues |
|---|---|---|---|---|---|
| **TSN** | $C_{17} H_{22} N_2 O_3$ | Trichostatin A | 100.00 | 90.9 | 90 131 132 140 141 167 169 256 263 293 295 |
| **SHH** | $C_{14} H_{20} N_2 O_3$ | Octadenioic acid hudroxyamide phenylamide | 100.00 | 90.9 | |

**TSN**



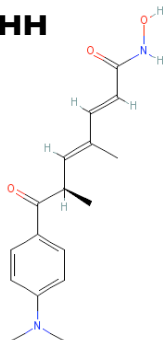### suberoylanilide hydroxamic acid

**Pharmacological Action:**
Anti-Inflammatory Agents, Non-Steroidal
Antineoplastic Agents
Enzyme Inhibitors
Anticarcinogenic Agents

Inhibits histone deacetylase I and 3

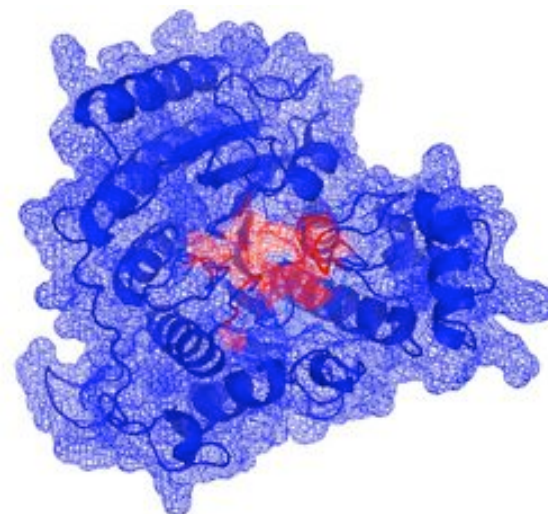| | LmjF21.0680.1.pdb |
|---|---|
| **Template** | 1t64A |
| **Seq. Id (%)** | **38.00** |
| **MPQS** | 1.47 |



**SHH**



### trichostatin A

**Pharmacological Action:**
Antibiotics, Antifungal
Enzyme Inhibitors
Protein Synthesis Inhibitors

chelates zinc ion in the active site of histone deacetylases, resulting in preventing histone unpacking so DNA is less available for transcription

# Example of inheritance (CDD-Roos-literature)

*LmjF21.0680 from* L. major *"Histone deacetylase 2" (model 1)*

## Apicidin: A novel antiprotozoal agent that inhibits parasite histone deacetylase

(cyclic tetrapeptide/Apicomplexa/antiparasitic/malaria/coccidiosis)

SANDRA J. DARKIN-RATTRAY*[†], ANNE M. GURNETT*, ROBERT W. MYERS*, PAULA M. DULSKI*, TAMI M. CRUMLEY*, JOHN J. ALLOCCO*, CHRISTINE CANNOVA*, PETER T. MEINKE[‡], STEVEN L. COLLETTI[‡], MARIA A. BEDNAREK[‡], SHEO B. SINGH[§], MICHAEL A. GOETZ[§], ANNE W. DOMBROWSKI[§], JON D. POLISHOOK[§], AND DENNIS M. SCHMATZ*

Departments of *Parasite Biochemistry and Cell Biology, [‡]Medicinal Chemistry, and [§]Natural Products Drug Discovery, Merck Research Laboratories, P.O. Box 2000, Rahway, NJ 07065

## Antimalarial and Antileishmanial Activities of Aroyl-Pyrrolyl-Hydroxyamides, a New Class of Histone Deacetylase Inhibitors

# "take home" message

# Acknowledgments