# Comparative docking for predicting molecular targets of known drugs
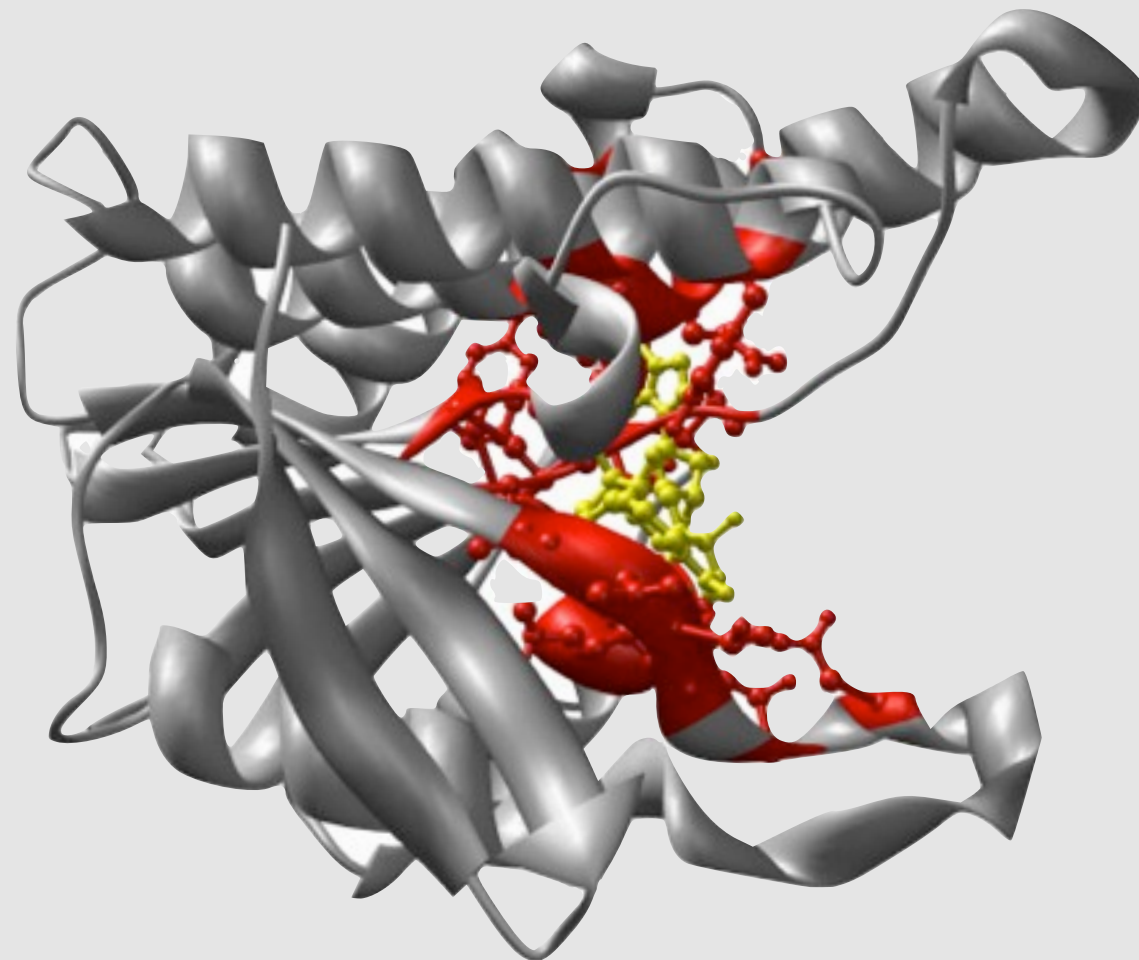## *A "kernel" for the Tropical Disease Initiative*
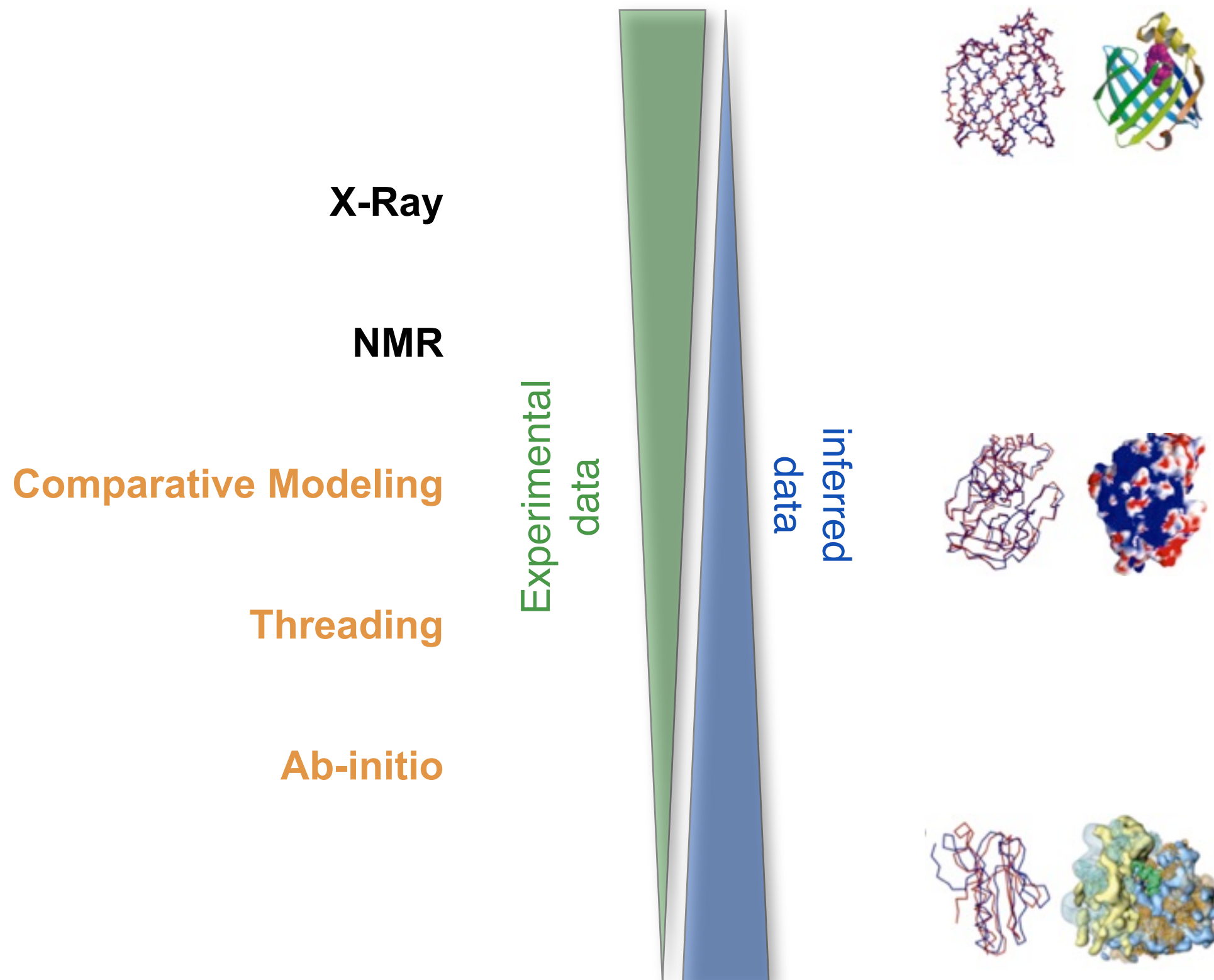
**Marc A. Marti-Renom**

http://sgu.bioinfo.cipf.es

Structural Genomics Unit
Bioinformatics & Genomics Department
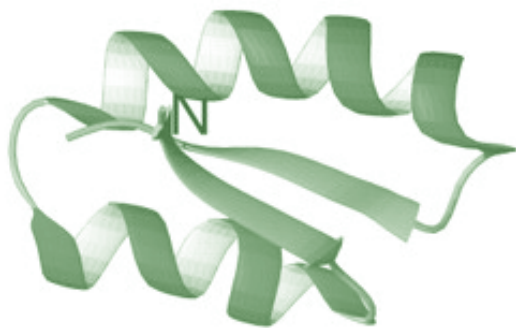Centro de Investigación Príncipe Felipe (CIPF), Valencia, Spain
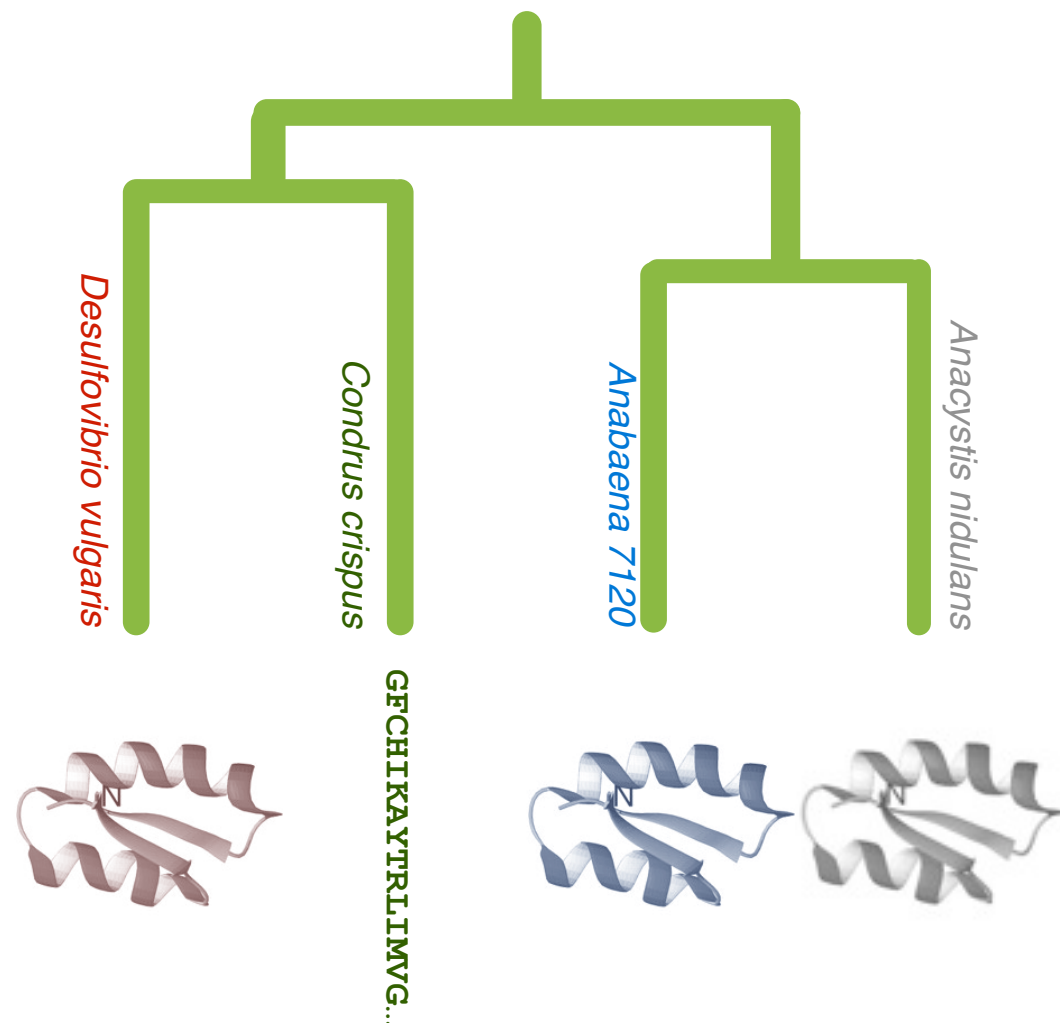
# protein prediction .vs. protein determination



**X-Ray**

**NMR**

**Comparative Modeling**

**Threading**

**Ab-initio**

Experimental data

inferred data

# Principles of protein structure

GFCHIKAYTRLIMVG...
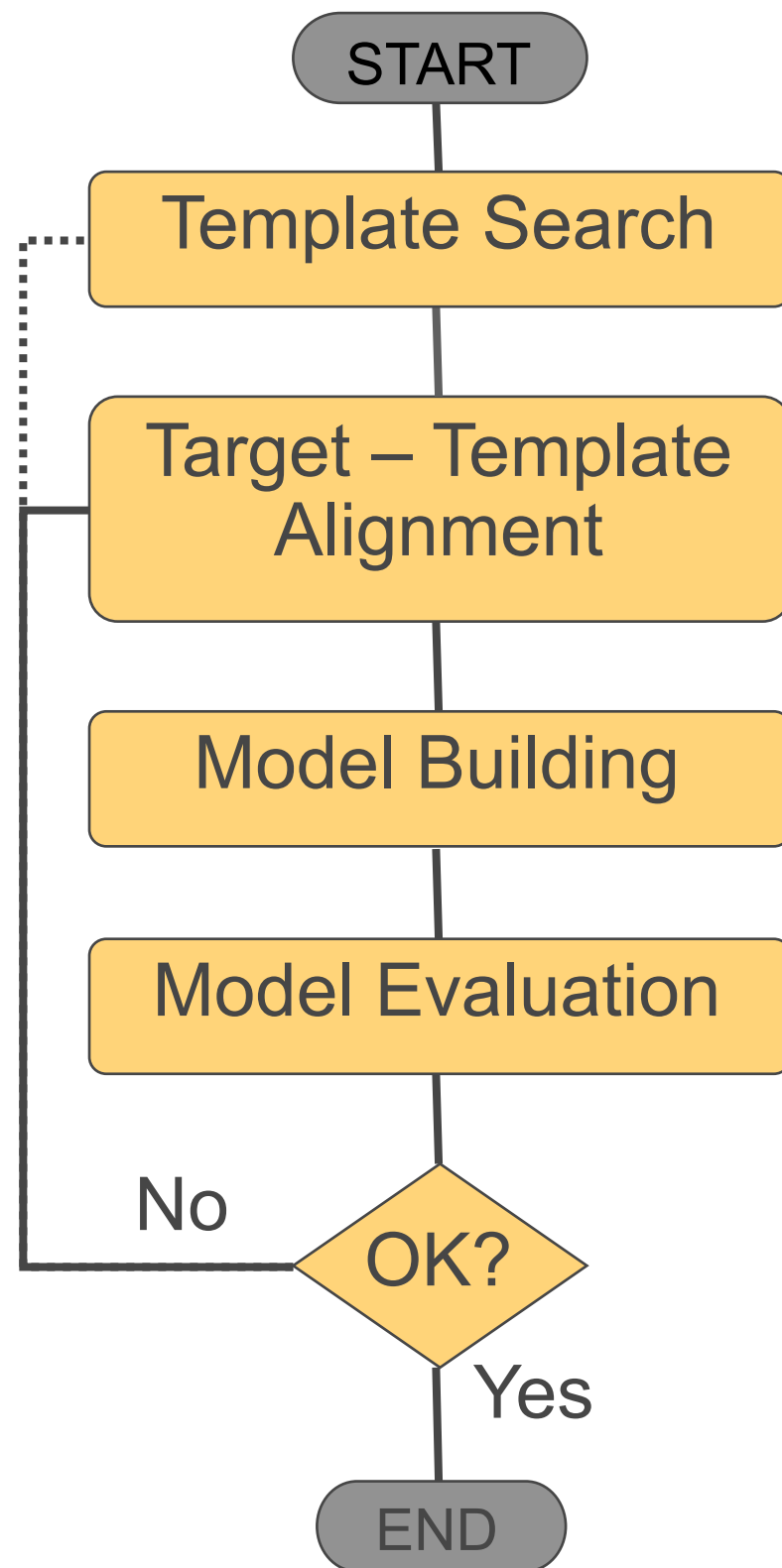


Folding (physics)

*Ab initio* prediction

Evolution (rules)

Threading
Comparative Modeling

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

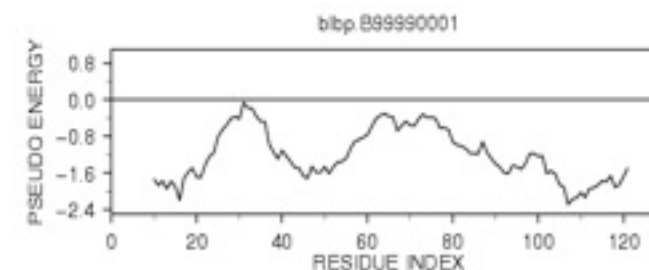# Steps in Comparative Protein Structure Modeling

START

Template Search

Target – Template Alignment

Model Building

Model Evaluation

OK?

No

Yes

END

TARGET

TEMPLATE

ASILPKRLFGNCEQTSDEGL
KIERTPLVPHISAQNVCLKID
DVPERLIPERASFQWMNDK

ASILPKRLFGNCEQTSDEGLKIERTPLVPHISAQNVCLKIDDVPERLIPE
MSVIPKRLYGNCEQTSEEAIRIEDSPIV---TADLVCLKIDEIPERLVGE
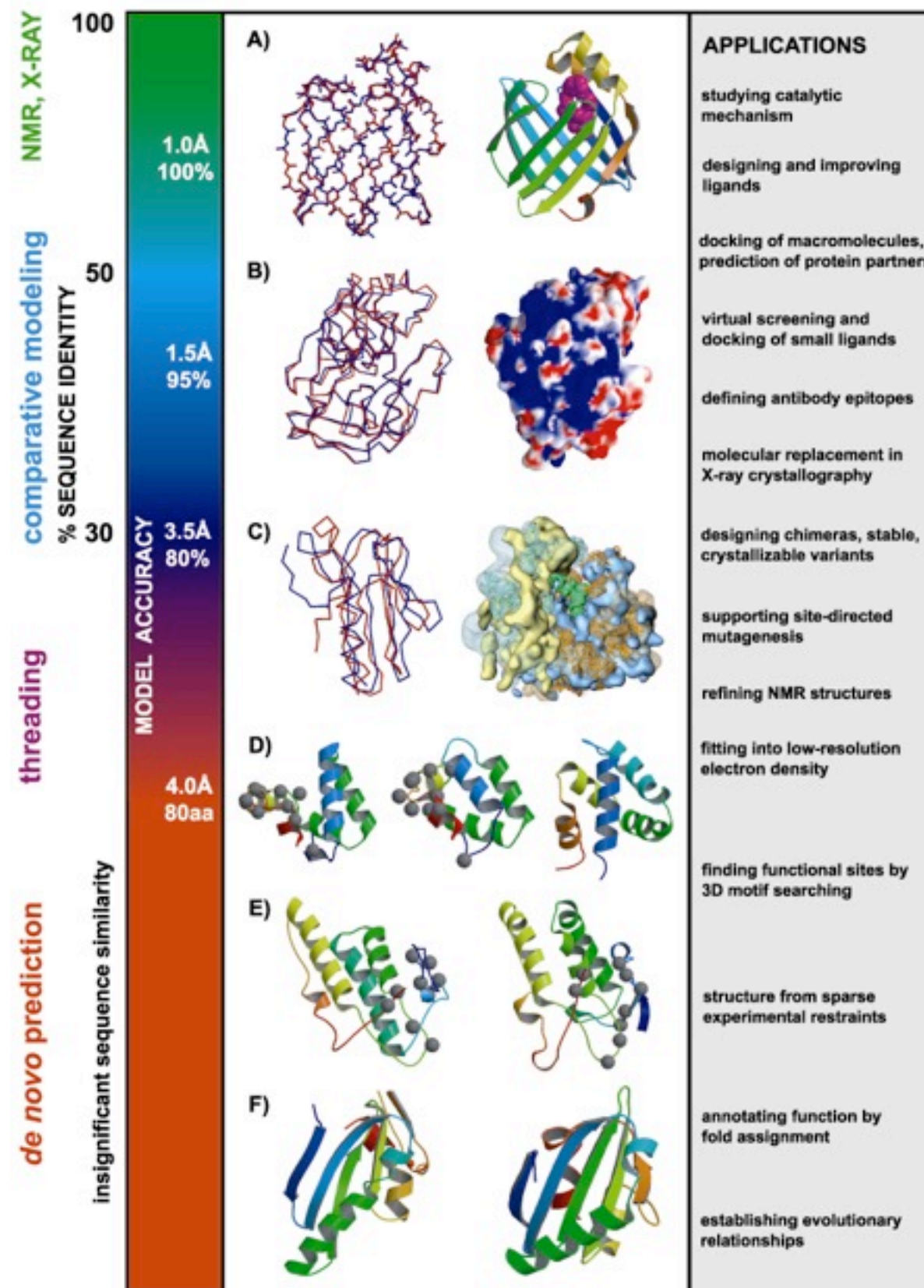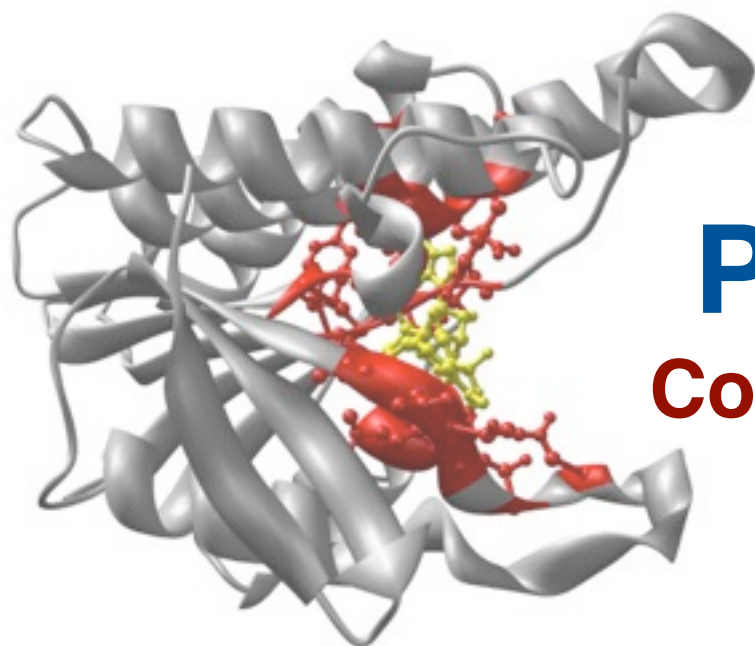
blbp.B99990001

PSEUDO ENERGY

RESIDUE INDEX

*A. Šali, Curr. Opin. Biotech. 6, 437, 1995.*
*R. Sánchez & A. Šali, Curr. Opin. Str. Biol. 7, 206, 1997.*
*M.A. Marti-Renom et al. Ann. Rev. Biophys. Biomolec. Struct., 29, 291, 2000.*

# Utility of protein structure models, despite errors



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

# Protein function from structure
## Comparative binding site prediction by AnnoLyze.



*Marti-Renom et al. BMC Bioinformatics (2007)*

# 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*

7

# DBAIi$_{v2.0}$ database

## http://www.dbali.org

*Marti-Renom et al. Nucleic Acids Research (2007)*

DBAIi

Search

Tools

Special pages

Multiple
Pairwise
Get all similar
e-mail
DBAlit!
AnnoLite
AnnoLyze
ModClus from list
ModClus from chain
SALIGN
ModDom

Structural Genomics
Download
Statistics

Multiple alignment result

Pairwise alignment result

Table of structural similarities

Fast annotations result

Full annotations result

Multiple alignment result

Cluster results

Domain assignments

### DBAIi tools: mining the protein structure space

Marc A. Marti-Renom[1,*], Ursula Pieper[2], M. S. Madhusudhan[2], Andrea Rossi[2], Narayanan Eswar[2], Fred P. Davis[2], Fátima Al-Shahrour[3], Joaquín Dopazo[3] and Andrej Sali[2]

[1]Structural Genomics Unit, [2]Departments of Biopharmaceutical Sciences and Pharmaceutical Chemistry, and California Institute for Quantitative Biomedical Research, University of California at San Francisco, San Francisco, CA 94158-2330, USA and [3]Functional Genomics Unit, Bioinformatics Department, Centro de Investigación Príncipe Felipe (CIPF), Valencia, Spain

**ABSTRACT**

The DBAIi tools use a comprehensive set of structural alignments in the DBAIi database to leverage the structural information deposited in the Protein Data Bank (PDB). These tools include (i) the DBAlit program that allows users to input the 3D coordinates of a protein structure for comparison by MAMMOTH against all chains in the PDB; (ii) the AnnoLite and AnnoLyze programs that annotate a target structure based on its stored relationships to other structures; (iii) the ModClus program that clusters structures by sequence and structure similarities; (iv) the ModDom program that identifies domains as recurrent structural fragments and (v) an implementation of the COMPARER method in the SALIGN command in MODELLER that creates a multiple structure alignment for a set of related protein structures. Thus, the DBAIi tools, which are freely accessible *via* the World Wide Web at http://salilab.org/DBAIi/, allow users to mine the protein structure space by establishing relationships between protein structures and their functions.

**INTRODUCTION**

The number of known protein structures deposited in the Protein Data Bank (PDB) has grown exponentially over the years (1). This trend is expected to continue, partly due to the structural genomics efforts (2,3). Currently, there are ~41 000 protein structures deposited in the PDB, containing ~88 000 protein chains. These protein structures constitute a structural space that can be mined t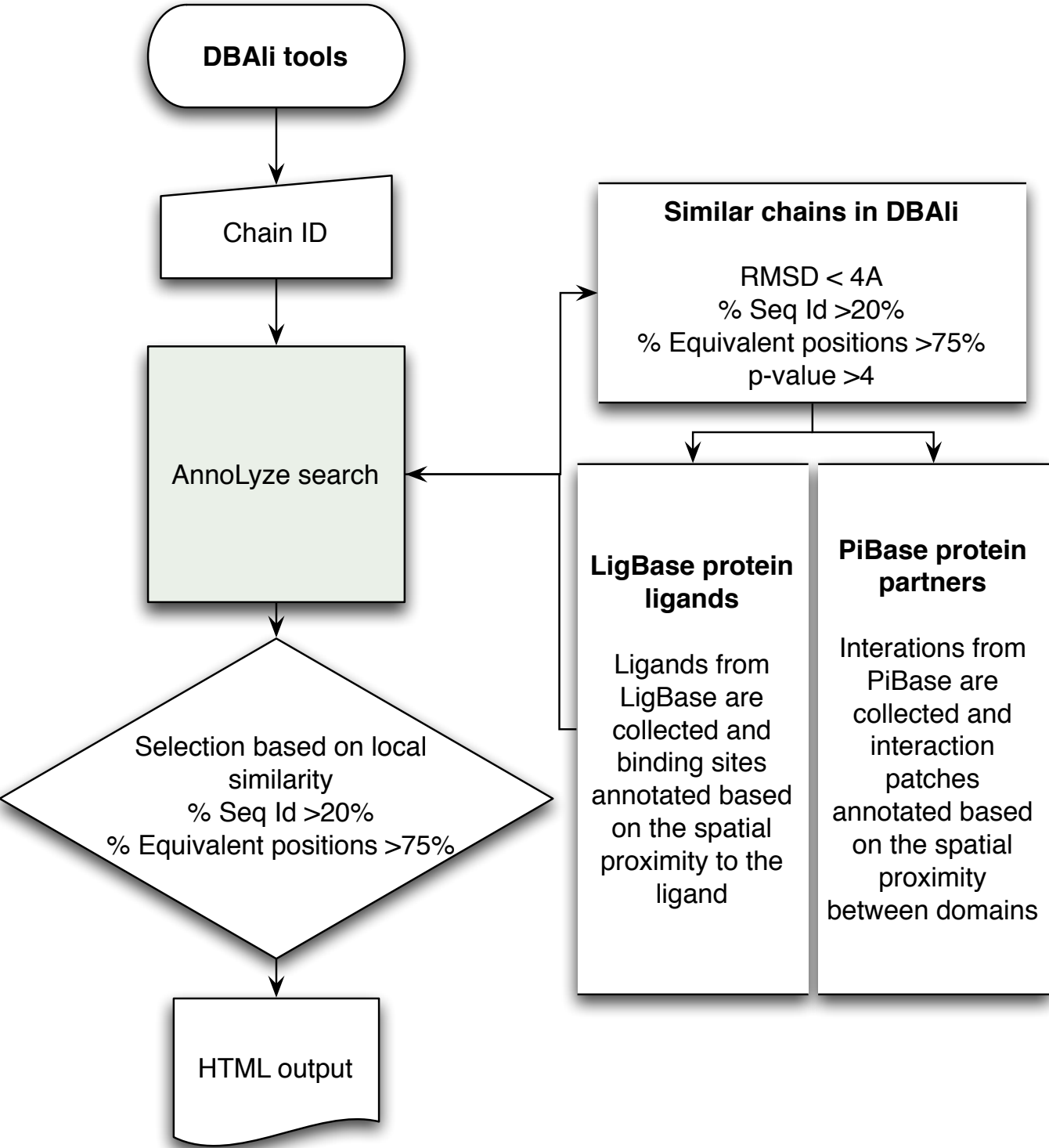o facilitate the understanding, assignment and modification of protein function. Previously developed databases for the classification of protein structure domains, such as SCOP [http://scop.mrc-lmb.cam.ac.uk/scop/ (4)] or CATH [http://www.cathdb.info (5)], and servers for functional annotation of protein structures, such as ProFunc [http://www.ebi.ac.uk/thornton-srv/databases/ProFunc/ (6,7)], ProKnow [http://www.doe-mbi.ucla.edu/Services/ProKnow (8)] and Phunctioner [http://www.sbg.bio.ic.ac.uk (9)], provide an effective way of describing and annotating the protein structure space. However, none of these servers combine a comprehensive database of protein structural alignments with tools for automatically annotating protein structures.

Here, we describe five tools that aid in the analysis of the data stored in DBAIi, our comprehensive relational database of pairwise and multiple structural alignments (10). These tools include (i) the DBAlit program that allows users to input their structure for comparison by MAMMOTH (11) against all chains in the PDB; (ii) the AnnoLite and AnnoLyze programs that annotate a target structure based on its stored relationships to other structures; (iii) the ModClus program that clusters structures by sequence and structure similarities; (iv) the ModDom program that identifies recurrent fragments, including domains, from structure; and (v) an implementation of the COMPARER method (12) in the SALIGN command in MODELLER that creates a multiple structure alignment for a set of related protein structures. The DBAIi tools allow users to establish relationships between protein structures and their fragments in a flexible and dynamic manner.

The DBAIi database is briefly introduced first. Next, we describe each of the five tools that make use of the structural alignments deposited in DBAIi. Finally, we discuss the use of the DBAIi tools to analyze a structure determined by the New York Structural Genomics Research Consortium (NYSGXRC).

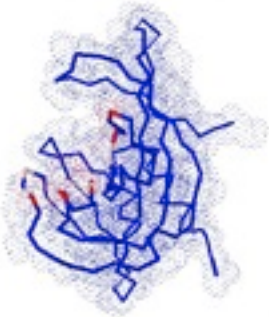*To whom correspondence should be addressed. Tel: +34 96 3289680; Fax: +34 96 3289701; Email: mmarti@cipf.es

Monday, July 12, 2010

# Method

**DBAli tools**

Chain ID

AnnoLyze search

**Similar chains in DBAli**

RMSD < 4A
% Seq Id >20%
% Equivalent positions >75%
p-value >4

**LigBase protein ligands**

Ligands from LigBase are collected and binding sites annotated based on the spatial proximity to the ligand

**PiBase protein partners**

Interations from PiBase are collected and interaction patches annotated based on the spatial proximity between domains

Selection based on local similarity
% Seq Id >20%
% Equivalent positions >75%

HTML output

Inherited ligands: 4

| Ligand | Av. binding site seq. id. | Av. residue conservation | Residues in predicted binding site (size proportional to the local conservation) |
|--------|------|------|------|
| MO2 | 59.03 | 0.185 | 48 49 52 62 63 66 67 113 116 |
| CRY | 20.00 | 0.111 | 23 29 31 37 44 48 49 83 85 94 96 103 121 |
| 8OG | 20.00 | 0.111 | 19 20 21 48 49 51 96 98 136 |
| ACY | 15.87 | 0.163 | 23 29 31 37 44 45 81 83 85 94 96 98 103 121 135 |

Inherited partners:1

| Partner | Av. binding site seq. id. | Av. residue conservation | Residues in predicted binding site (size proportional to the local conservation) |
|--------|------|------|------|
| d.113.1.1 | 23.68 | 0.948 | 19 20 50 51 52 53 54 55 56 57 58 77 78 79 80 81 82 83 84 85 93 95 97 99 134 135 138 142 145 |

9

# Scoring function

Ligands

Partners



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

# 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 Cα atoms, result in a sequence alignment with more than 30% identity, and have a length difference inferior to 50aa

# Sensitivity .vs. Precision

| | Optimal cut-off | Sensitivity (%)<br>Recall or TPR | Precision (%) |
|---|---|---|---|
| **Ligands** | 30% | 71.9 | 13.7 |

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

*Marti-Renom et al. BMC Bioinformatics (2007) Volume 8. Suppl S4*

**~90-95% of residues correctly predicted**

Monday, July 12, 2010

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

Molecule: MutT/nudix family protein

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



[http://www.tropicaldisease.org](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

Monday, July 12, 2010

# TDI flowchart



databases of genome sequences

database of protein structures

virtual ligand libraries

PubMed, journals    other databases

sequence similarity searches

protein structure modeling

literature searches    protein-ligand docking

functional annotation

**COMPUTING**

# TDI

**TARGET DISCOVERY**
**LEAD DISCOVERY**
**LEAD OPTIMIZATION**

synthetic chemistry    compound libraries

high-throughput screening

**CHEMISTRY**

protein production    protein engineering

substrate specificity studies

structural biology    target validation

**BIOLOGY**

leads

**VIRTUAL PHARMA**

*and other development organizations*

**TOXICITY AND PHARMACOKINETIC EVALUATION**

**CLINICAL STUDIES**

**DRUG PRODUCTION**

drugs

16

# Non-Profit organizations

## *Open-Source + Out-Source = low cost business model*



*Munos (2006) Nature Reviews. Drug Discovery.*

# "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.

18

# Comparative docking



**Expansion**

co-crystalized protein/ligand

crystalized
protein

template

**2. Inheritance**

model

**1. Modeling**

# Modeling Genomes

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

*A good model has MPQS of 1.0 or higher*

Monday, July 12, 2010

# Summary table

## models with inherited ligands

**29,271 targets with good models, 297 inherited a ligand/substance similar to a known drug in DrugBank**

| | Transcripts | Modeled targets | Selected models | Inherited ligands | Similar to a drug | Drugs |
|---|---|---|---|---|---|---|
| *C. hominis* | 3,886 | 1,614 | 666 | 197 | 20 | 13 |
| *C. parvum* | 3,806 | 1,918 | 742 | 232 | 24 | 13 |
| *L. major* | 8,274 | 3,975 | 1,409 | 478 | 43 | 20 |
| *M. leprae* | 1,605 | 1,178 | 893 | 310 | 25 | 6 |
| *M. tuberculosis* | 3,991 | 2,808 | 1,608 | 365 | 30 | 10 |
| *P. falciparum* | 5,363 | 2,599 | 818 | 284 | 28 | 13 |
| *P. vivax* | 5,342 | 2,359 | 822 | 268 | 24 | 13 |
| *T. brucei* | 7,793 | 1,530 | 300 | 138 | 13 | 6 |
| *T. cruzi* | 19,607 | 7,390 | 3,070 | 769 | 51 | 28 |
| *T. gondii* | 9,210 | 3,900 | 1,386 | 458 | 39 | 21 |
| **TOTAL** | **68,877** | **29,271** | **11,714** | **3,499** | **297** | **143** |

# *L. major* Histone deacetylase 2 + Vorinostat

## Template 1t64A a human HDAC8 protein.



| PDB | ⎐ | Template | ⛢ | Model | ↦ | Ligand | Exact | SupStr | SubStr | Similar |
|---|---|---|---|---|---|---|---|---|---|---|
| 1c3sA | 83.33/80.00 | 1t64A | 36.00/1.47 | LmjF21.0680.1.pdb | 90.91/100.00 | SHH | DB02546 | DB02546 | DB02546 | DB02546 |



**DB02546 Vorinostat**

Small Molecule; Approved; Investigational

**Drug categories:**

Anti-Inflammatory Agents, Non-Steroidal

Anticarcinogenic Agents

Antineoplastic Agents

Enzyme Inhibitors

**Drug indication:**

*For the treatment of cutaneous manifestations in patients with cutaneous T-cell lymphoma who have progressive, persistent or recurrent disease on or following two systemic therapies.*



22

# *L. major* Histone deacetylase 2 + Vorinostat

## *Literature*

## 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

# *P. falciparum* tymidylate kinase + zidovudine

## Template 3tmkA a yeast tymidylate kinase.



| PDB | ◖◗ | Template | ▥ | Model | ↪ | Ligand | Exact | SupStr | SubStr | Similar |
|---|---|---|---|---|---|---|---|---|---|---|
| 2tmkB | 100.00/100.00 | 3tmkA | 41.00/1.49 | PFL2465c.2.pdb | 82.61/100.00 | ATM | | DB00495 | | DB00495 |



### DB00495 Zidovudine

Small Molecule; Approved

**Drug categories:**

Anti-HIV Agents

Antimetabolites

Nucleoside and Nucleotide Reverse Transcriptase
Inhibitors

**Drug indication:**

*For the treatment of human immunovirus (HIV) infections.*

# *P. falciparum* thymidylate kinase + zidovudine

NMR *Water-LOGSY* and *STD* experiments

ATM

Zidovudine

dTMP

*Leticia Ortí, Rodrigo J. Carbajo, and Antonio Pineda-Lucena*

# TDI's kernel

## http://tropicaldisease.org/kernel

# TDI's kernel

## http://tropicaldisease.org/kernel

L. Orti *et al.*, *Nat Biotechnol* **27**, 320 (2009).

L. Orti *et al.*, *PLoS Negl Trop Dis* **3**, e418 (2009).





27

# Acknowledgments

## COMPARATIVE MODELING
**Andrej Sali**
M. S. Madhusudhan
**Narayanan Eswar**
Min-Yi Shen
**Ursula Pieper**
Ben Webb
Maya Topf (Birbeck College)

## MODEL ASSESSMENT
David Eramian
Min-Yi Shen
Damien Devos

## FUNCTIONAL ANNOTATION
Andrea Rossi (Rinat-Pfizer)
Fred Davis (Janelia Fram)

## MODEL ASSESSMENT
Francisco Melo (CU)
Alejandro Panjkovich (CU)

## NMR
**Antonio Pineda-Lucena**
**Leticia Ortí**
**Rodrigo J. Carbajo**

## MAMMOTH
**Angel R. Ortiz**

## FUNCTIONAL ANNOTATION
**Fatima Al-Shahrour**
**Joaquin Dopazo**

## BIOLOGY
Jeff Friedman (RU)
James Hudsped (RU)
Partho Ghosh (UCSD)
Alvaro Monteiro (Cornell U)
Stephen Krilis (St.George H)

**Tropical Disease Initiative**
**Stephen Maurer (UC Berkeley)**
**Arti Rai (Duke U)**
**Andrej Sali (UCSF)**
**Ginger Taylor (TSL)**
**Matthew Todd (U Sydney)**

**CCPR Functional Proteomics**
Patsy Babbitt (UCSF)
Fred Cohen (UCSF)
Ken Dill (UCSF)
Tom Ferrin (UCSF)
John Irwin (UCSF)
Matt Jacobson (UCSF)
Tack Kuntz (UCSF)
Andrej Sali (UCSF)
Brian Shoichet (UCSF)
Chris Voigt (UCSF)

**EVA**
Burkhard Rost (Columbia U)
Alfonso Valencia (CNB/UAM)

**CAMP**
Xavier Aviles (UAB)
Hans-Peter Nester (SANOFI)
Ernst Meinjohanns (ARPIDA)
Boris Turk (IJS)
Markus Gruetter (UE)
Matthias Wilmanns (EMBL)
Wolfram Bode (MPG)