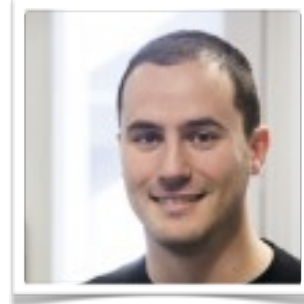


Structural Bioinformatics

Marc A. Marti-Renom

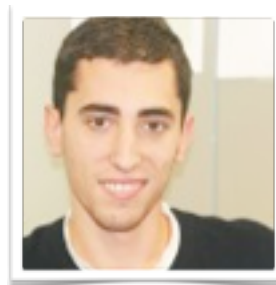
Genome Biology Group (CNAG)
Structural Genomics Group (CRG)

Structural Bioinformatics



David Dufour

- Estructura y biofísica de ácidos nucleicos y proteínas **25 febrero (DD)**
- Bases de datos de estructura de proteínas, ácidos nucleicos y pequeñas moléculas **11 marzo (DD)**
- Alineamiento y clasificación de estructura **25 marzo (DD)**
- Predicción de estructura tridimensional de ácidos nucleicos y proteínas **15 abril (DD)**



Francisco Martínez

- Docking de pequeñas moléculas en la superficie de estructura de proteínas **29 abril (FM)**
- Aplicaciones para el desarrollo de nuevos fármacos **13 mayo (FM)**

Outline...

COMPARATIVE MODELING

EXAMPLES

THE TROPICAL DISEASE INITIATIVE

Nomenclature

Homology: Sharing a common ancestor, may have similar or dissimilar functions

Similarity: Score that quantifies the degree of relationship between two sequences.

Identity: Fraction of identical aminoacids between two aligned sequences (case of similarity).

Target: Sequence corresponding to the protein to be modeled.

Template: 3D structure/s to be used during protein structure prediction.

Model: Predicted 3D structure of the target sequence.

Nomenclature

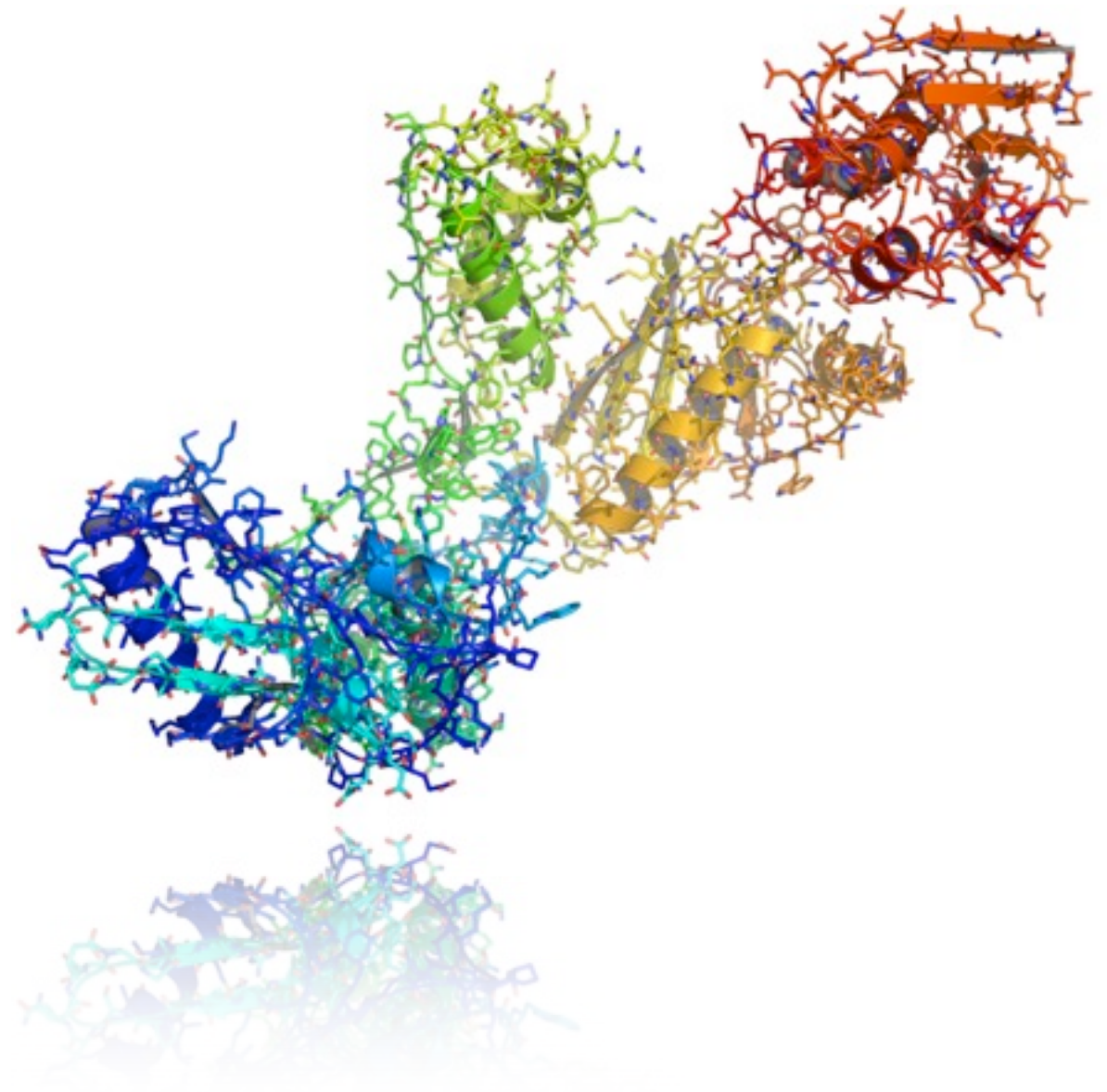
Fold: Three dimensional conformation of a protein sequence (usually at domain level).

Domain: Structurally globular part of a protein, which may independently fold.

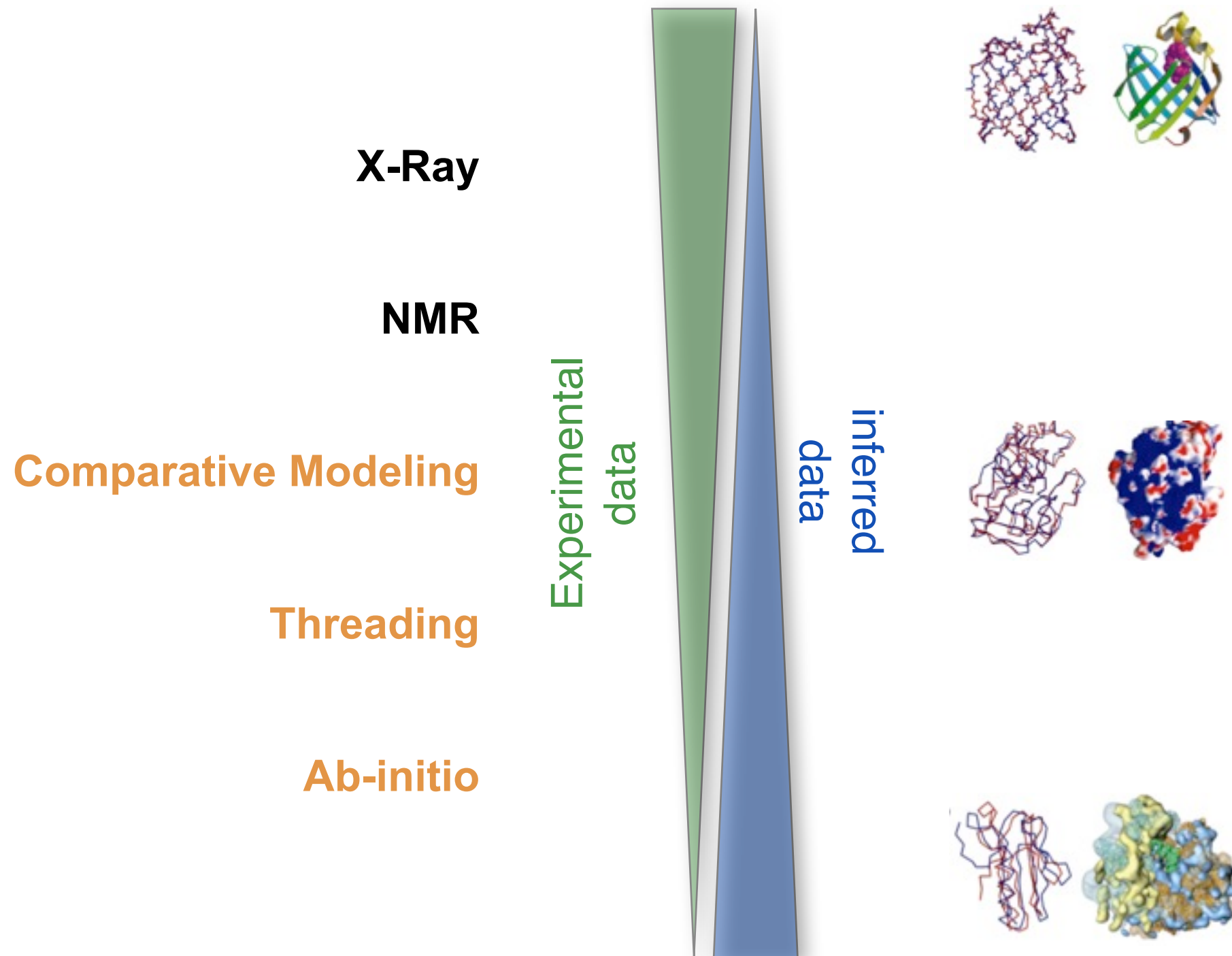
Secondary Structure: Regular sub-domain structures composed by alpha-helices, beta-sheets and coils (or loops).

Backbone: Protein structure skeleton composed by the carbon, nitrogen and oxygen atoms.

Side-Chain: Specific atoms identifying each of the 20 residues types.



protein prediction .vs. protein determination



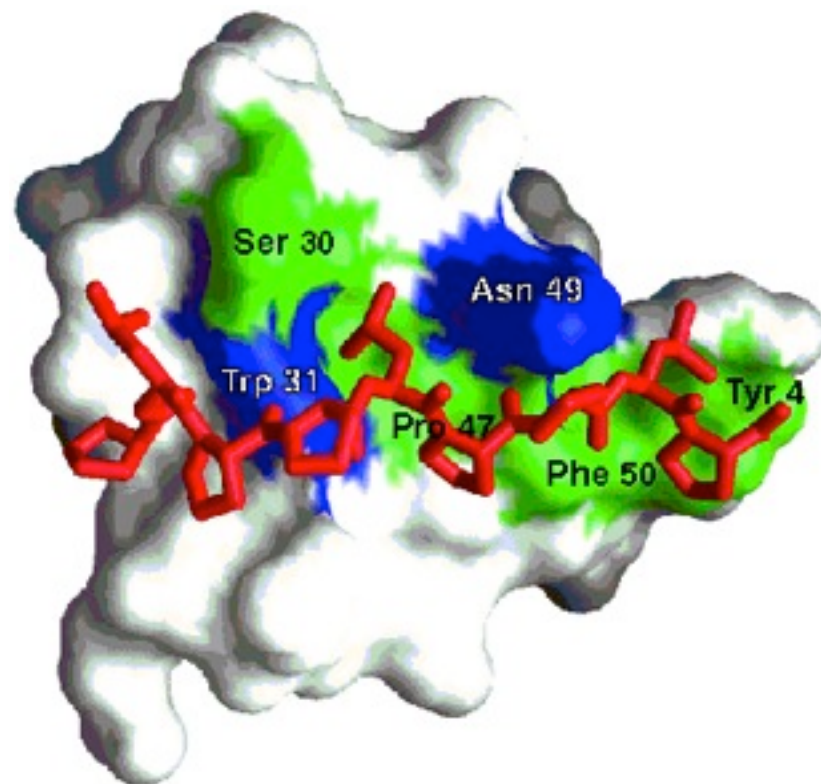
Why is it useful to know the **structure** of a protein, not only its sequence?

- ◆ The biochemical function (activity) of a protein is defined by its interactions with other molecules.
- ◆ The biological function is in large part a consequence of these interactions.
- ◆ The 3D structure is more informative than sequence because interactions are determined by residues that are close in space but are frequently distant in sequence.

YDL117W
(15-64)

10 20 30 40 50

KARYGWSGQTKGDLGFLEGDIMEVTRIAGSWFYGKLLRNKKCSGYFPHIE

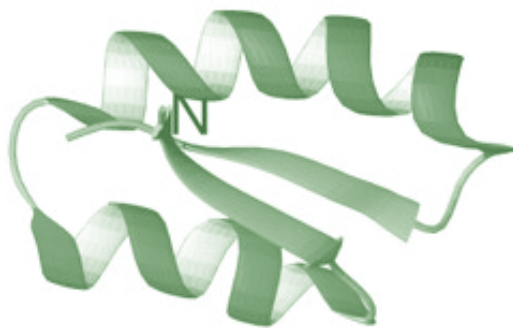
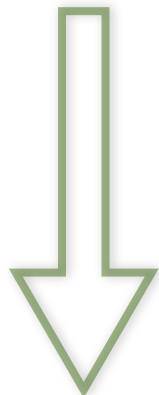


In addition, since evolution tends to conserve function and function depends more directly on structure than on sequence, **structure is more conserved in evolution than sequence.**

The net result is that **patterns in space are frequently more recognizable than patterns in sequence.**

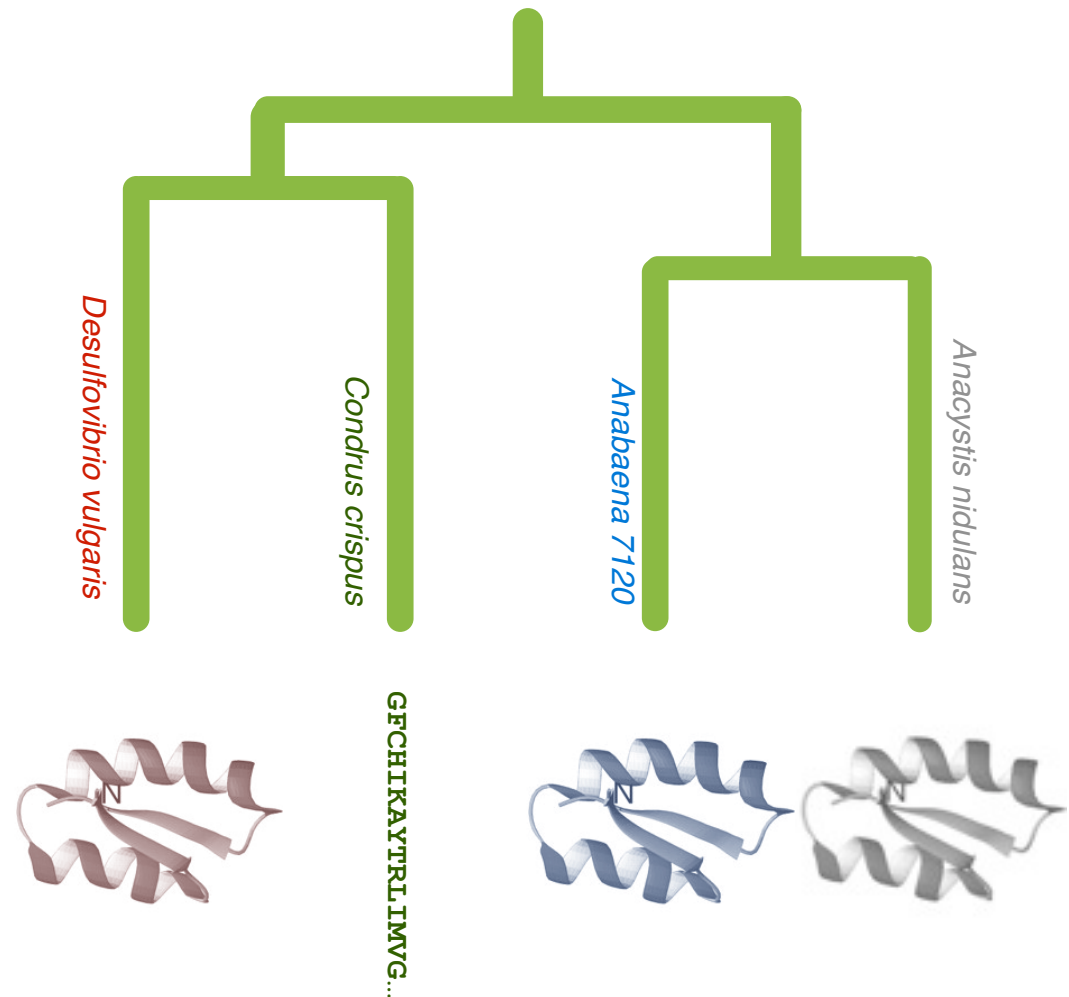
Principles of protein structure

GFCHIKAYTRLIMVG...



Folding (physics)

Ab initio prediction

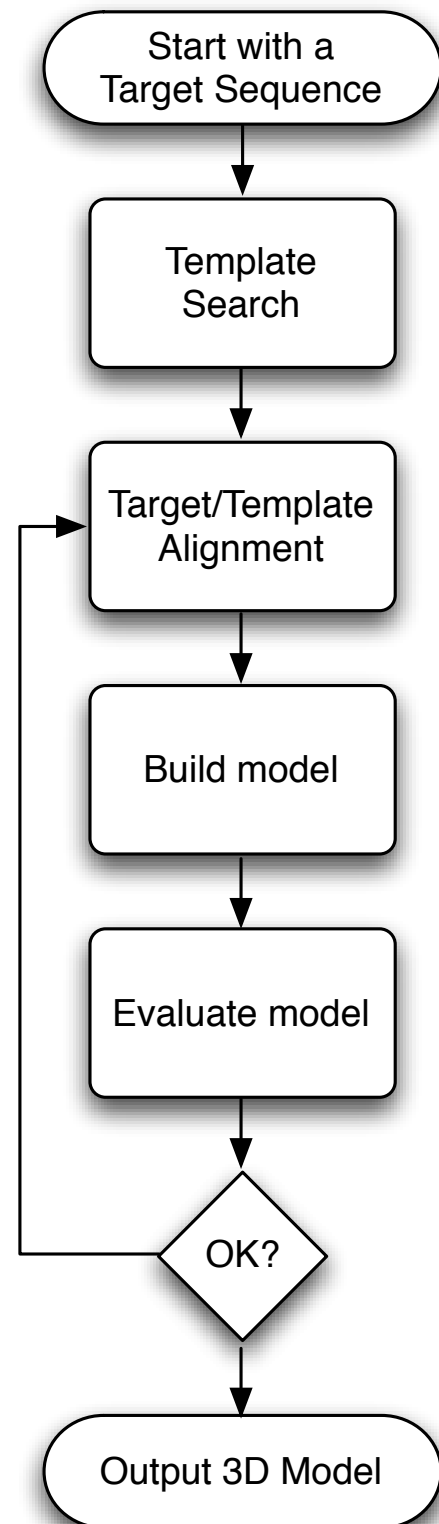


Evolution (rules)

Threading
Comparative Modeling

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

Comparative modeling by satisfaction of spatial restraints



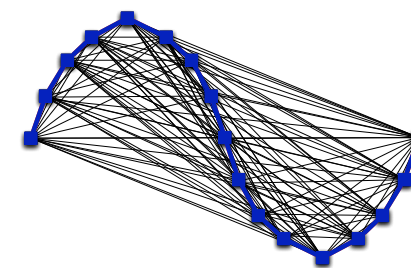
Given an alignment...

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

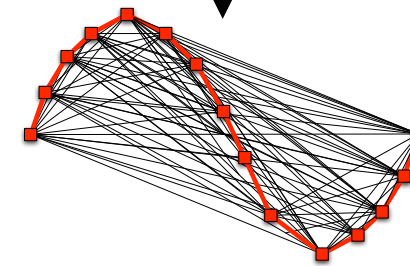
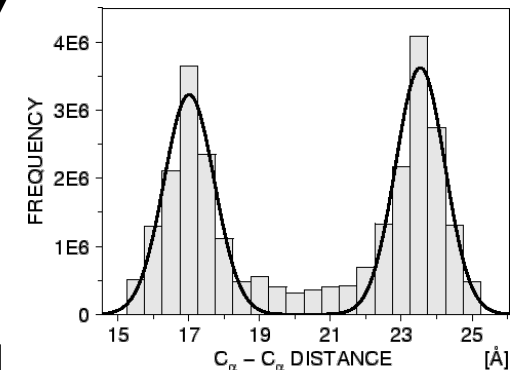
apply these features
as restraints on your
target sequence

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

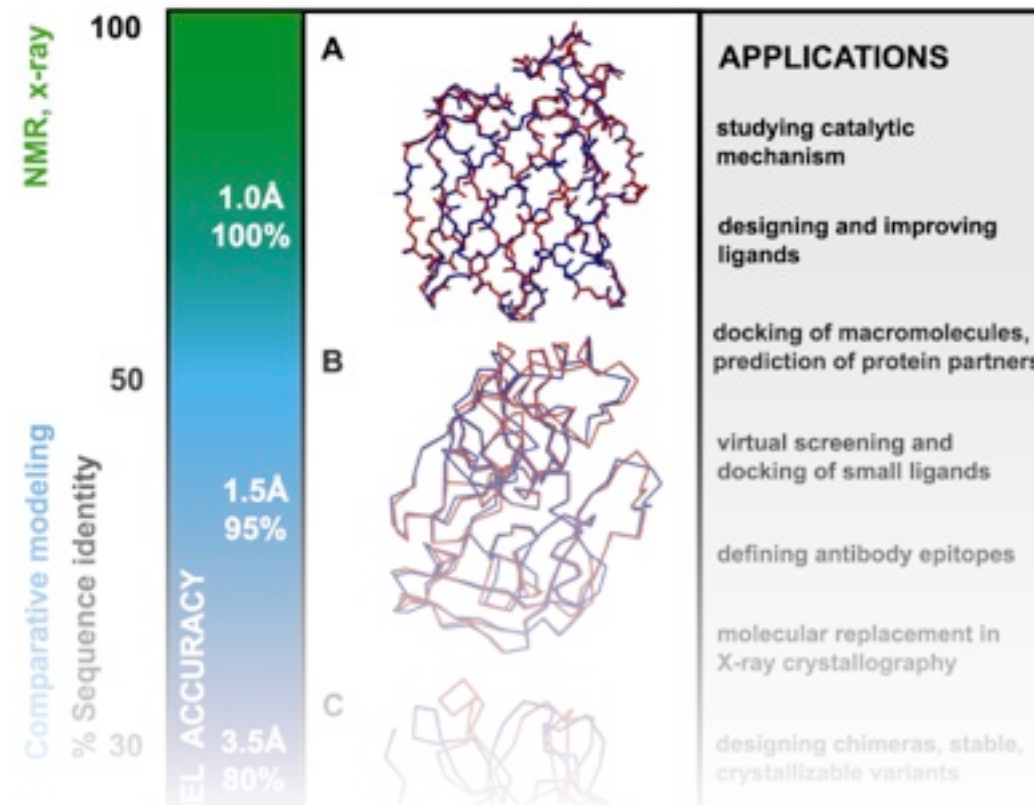
MSVIPKR--GNCEQTSE
ASILPKRLFGNCEQTSD



+



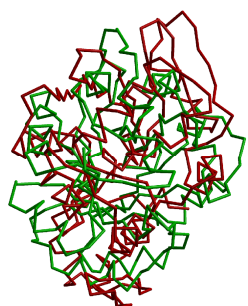
A. Šali & T. Blundell, *J. Mol. Biol.* 234, 779, 1993.
J.P. Overington & A. Šali, *Prot. Sci.* 3, 1582, 1994.
A. Fiser, R. Do & A. Šali, *Prot. Sci.*, 9, 1753, 2000.



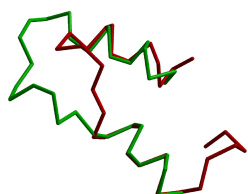
Accuracy and applicability of comparative models

Comparative modeling by satisfaction of spatial restraints

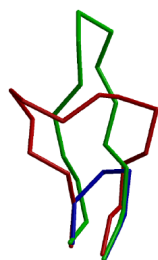
Types of errors and their impact



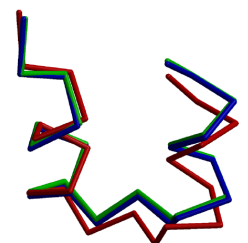
Wrong fold



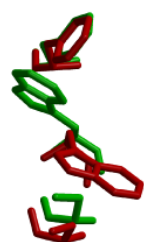
Miss alignments



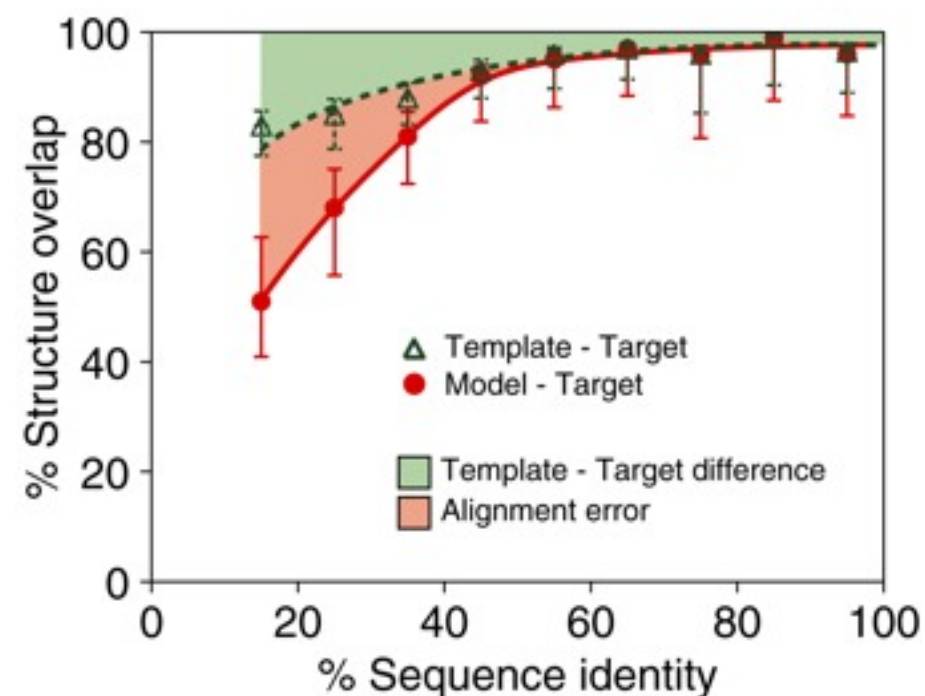
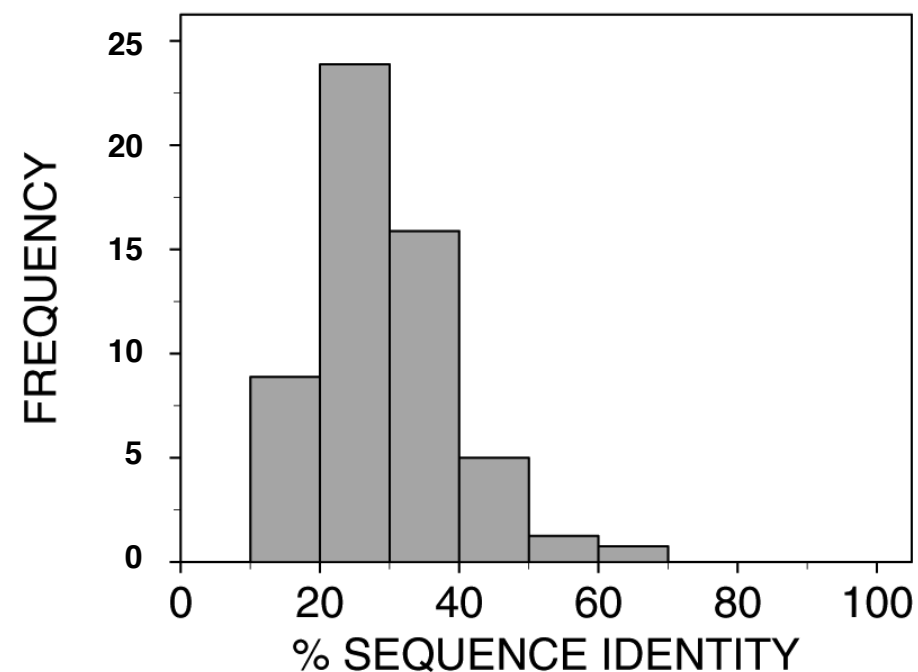
Loop regions



Rigid body distortions

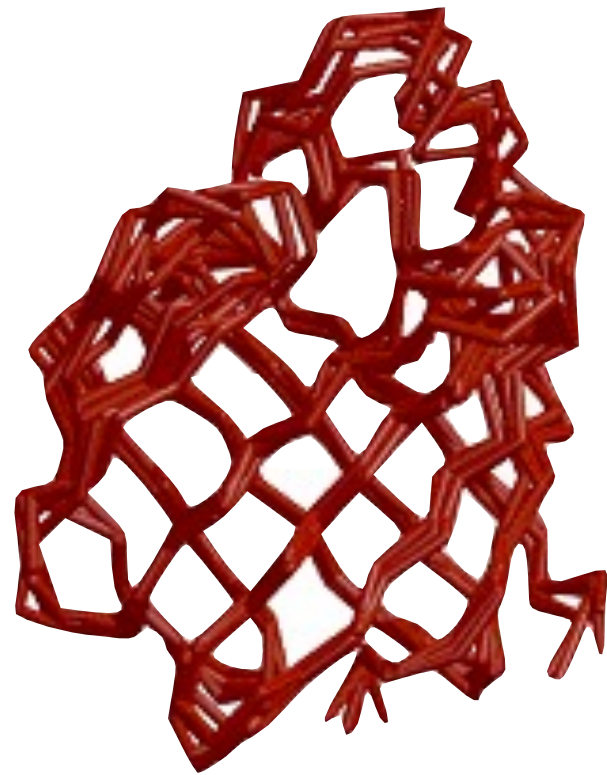


Side-chain packing



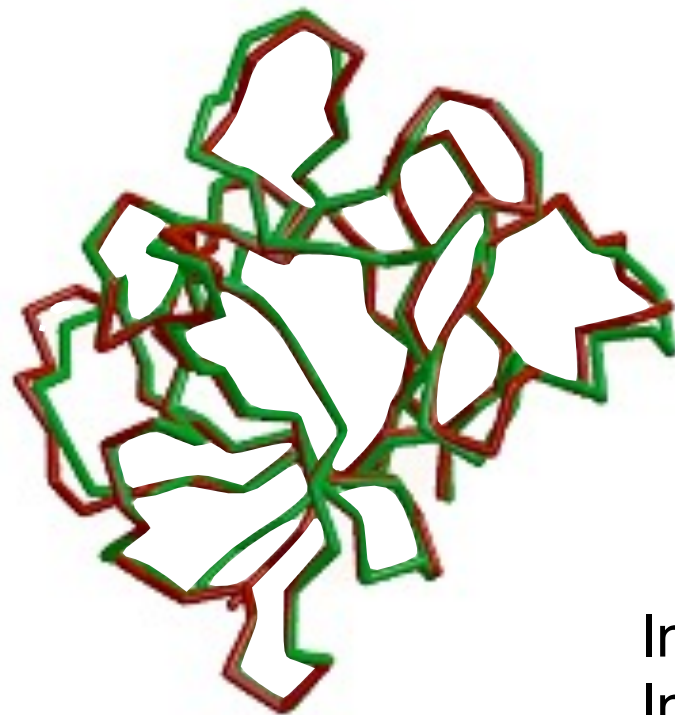
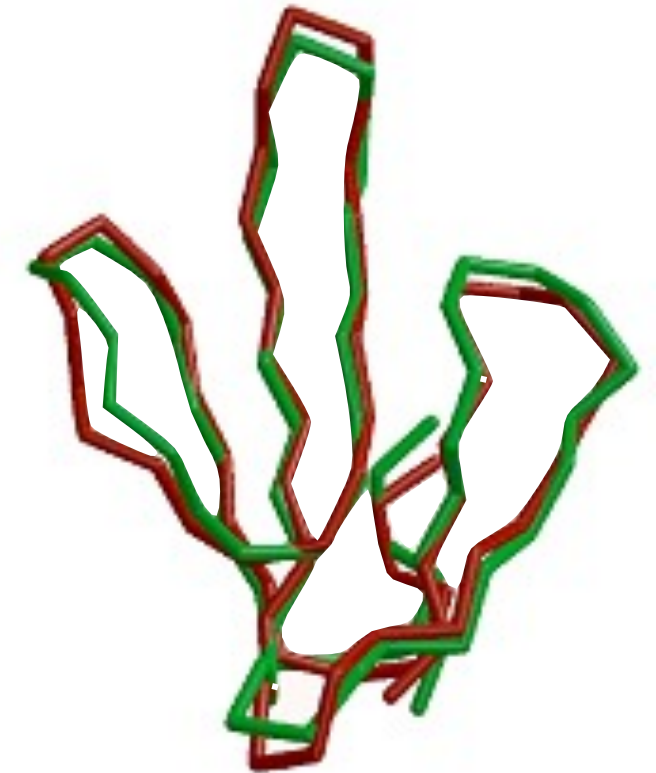
Marti-Renom et al. Ann Rev Biophys Biomol Struct (2000) 29, 291

“Biological” significance of modeling errors



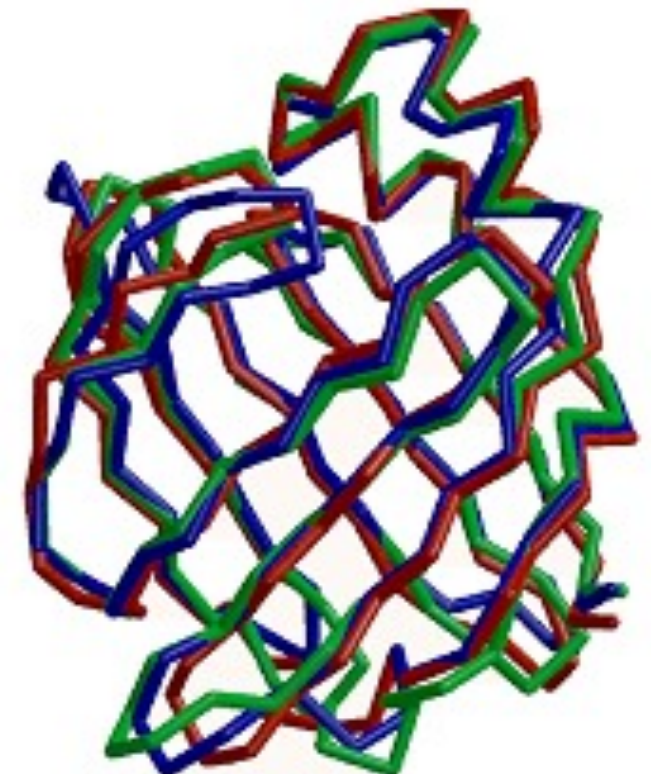
NMR
Ileal lipid-binding protein
1eal

NMR – X-RAY
Erabutoxin 3ebx
Erabutoxin 1era



X-RAY
Interleukin 1 β 41bi (2.9Å)
Interleukin 1 β 2mib (2.8Å)

CRABP II 1opbB
FABP 1ftpA
ALBP 1lib
40% seq. id.



Model Accuracy

HIGH ACCURACY

NM23 Seq id 77%

C α equiv 147/148
RMSD 0.41Å

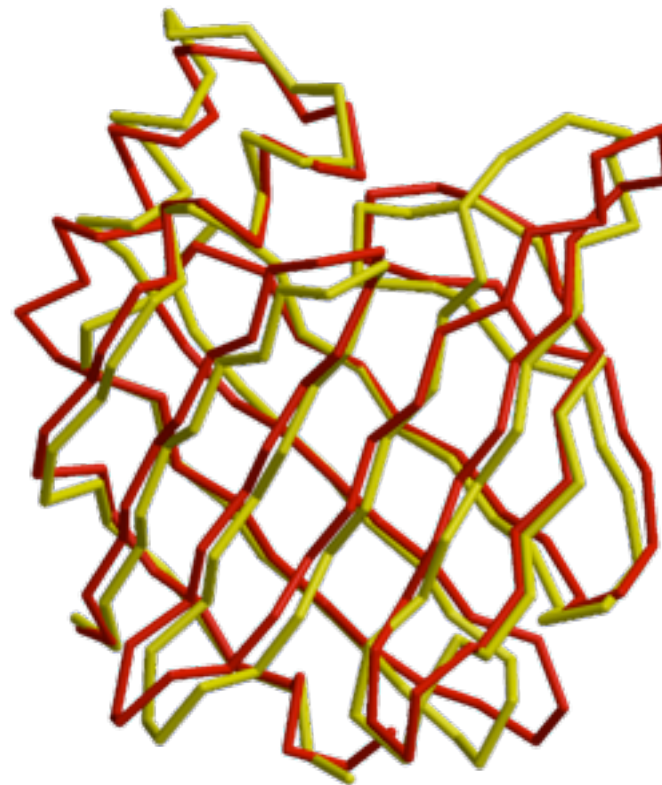


Sidechains
Core backbone
Loops

MEDIUM ACCURACY

CRABP Seq id 41%

C α equiv 122/137
RMSD 1.34Å



Sidechains
Core backbone
Loops
Alignment

LOW ACCURACY

EDN Seq id 33%

C α equiv 90/134
RMSD 1.17Å

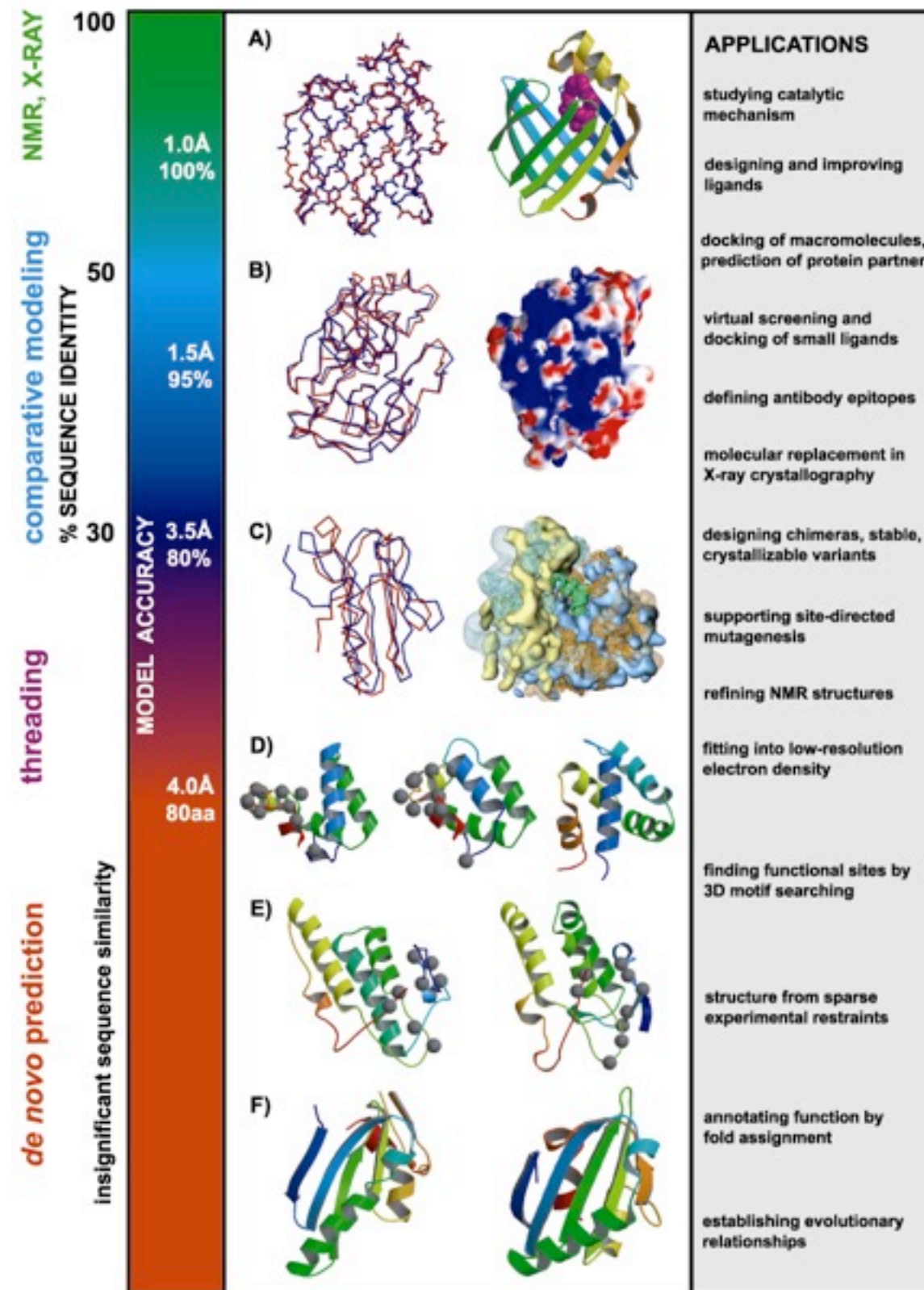


Sidechains
Core backbone
Loops
Alignment
Fold assignment

X-RAY / MODEL

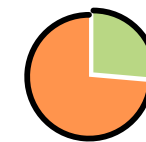
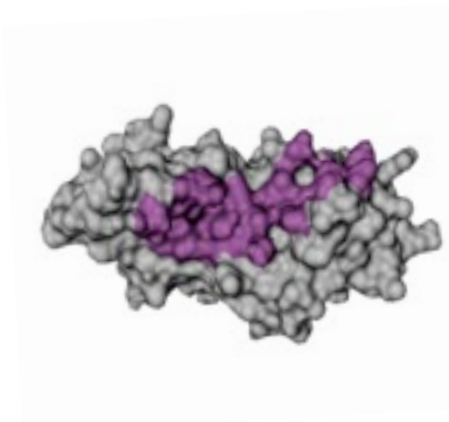
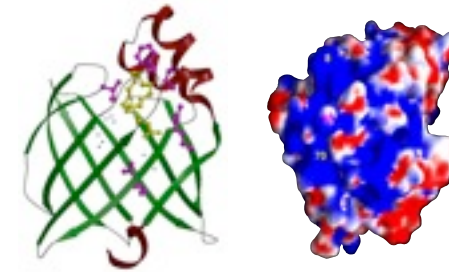
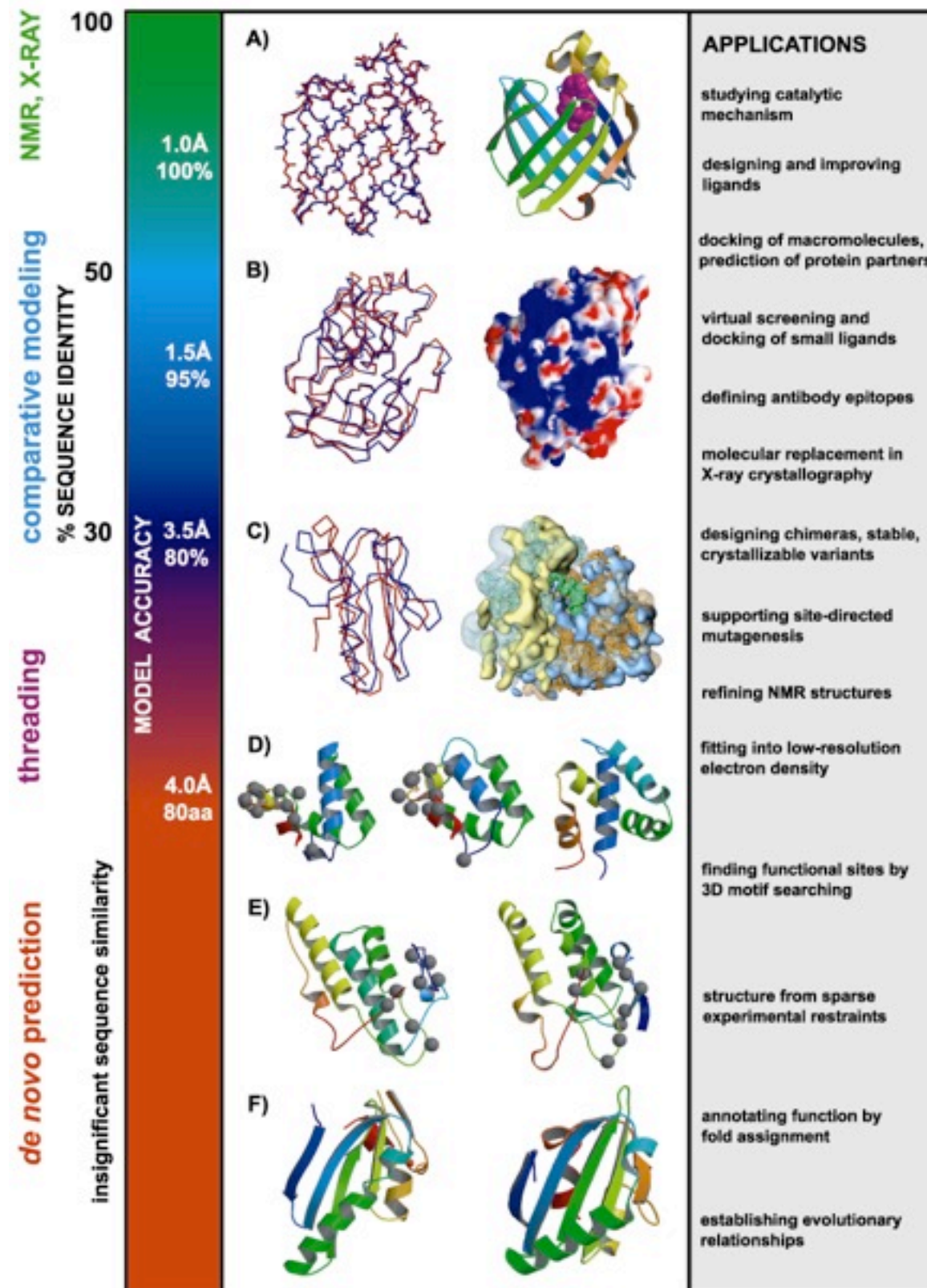
Marti-Renom et al. Annu.Rev.Biophys.Biomol.Struct. 29, 291-325, 2000.

Utility of protein structure models, despite errors

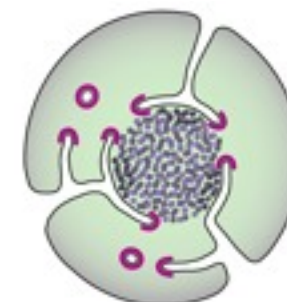


D. Baker & A. Sali. Science 294, 93, 2001.

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

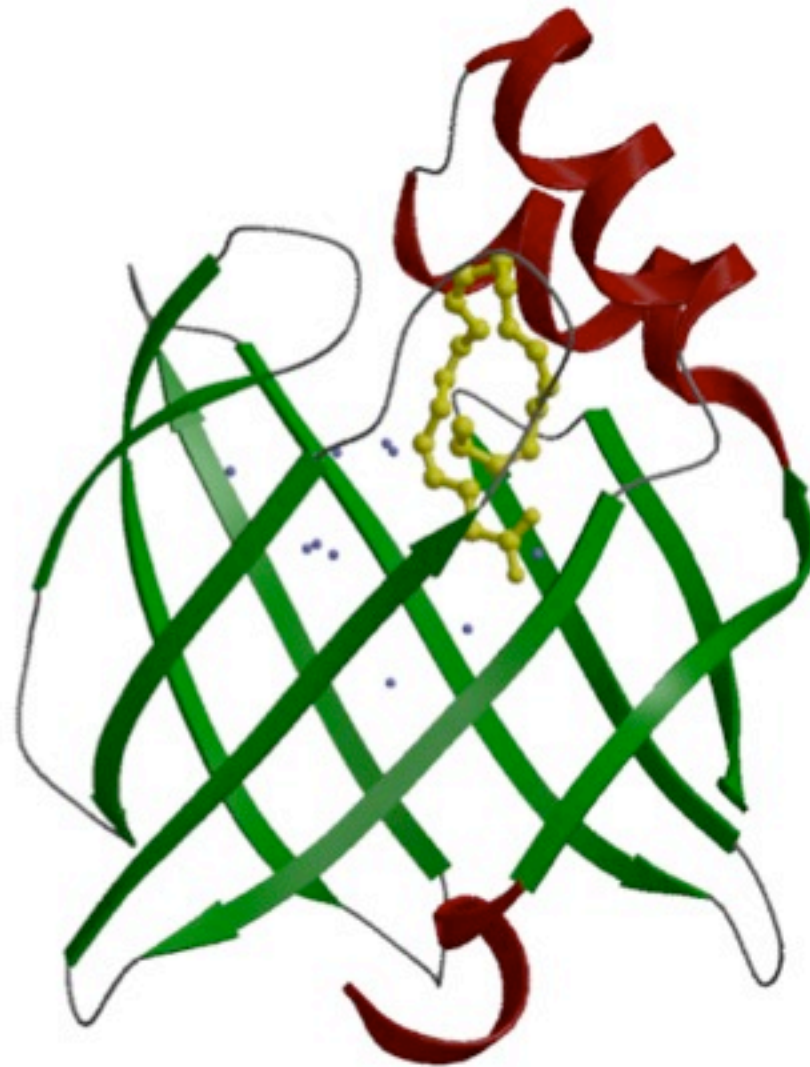
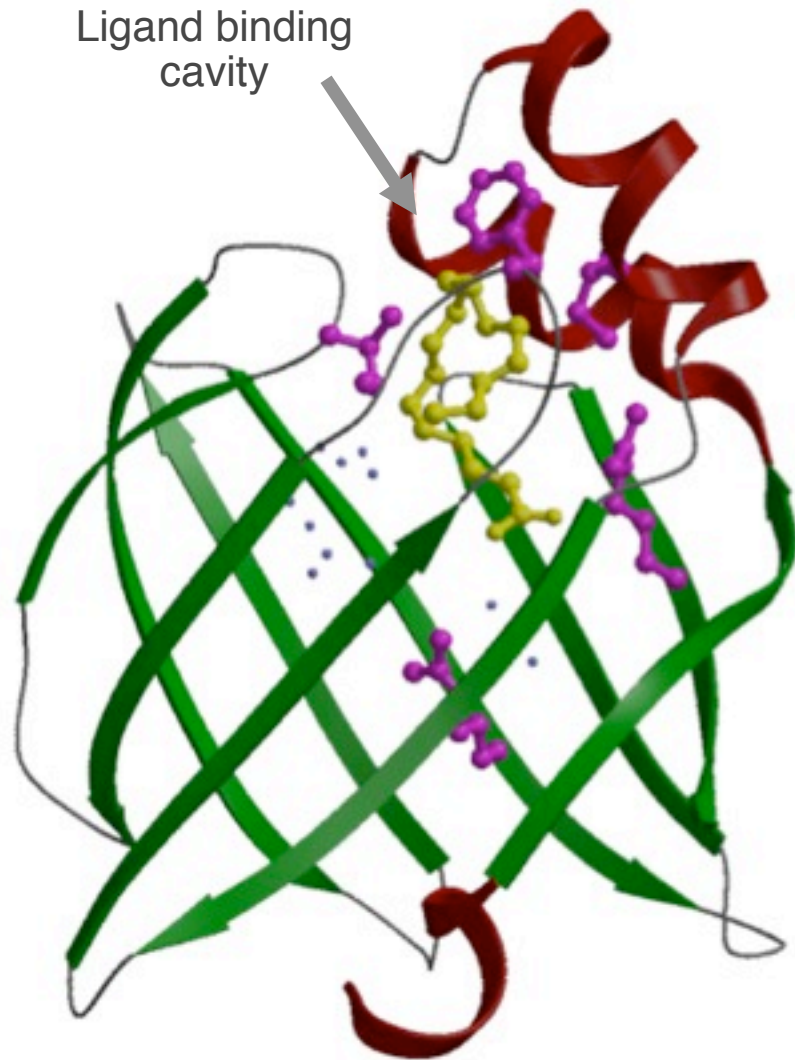
BLBP/oleic acid

Cavity is **not** filled

BLBP/docosahexaenoic acid

Cavity **is** filled

Ligand binding
cavity



1. BLBP binds fatty acids.

2. Build a 3D model.

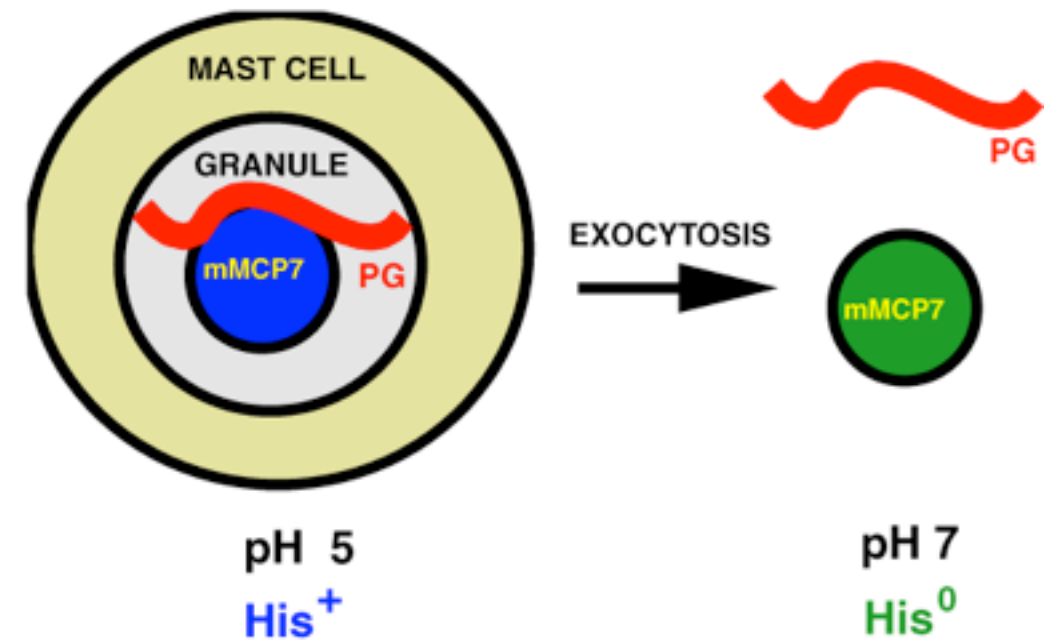
3. Find the fatty acid that fits most snugly into the ligand binding cavity.

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

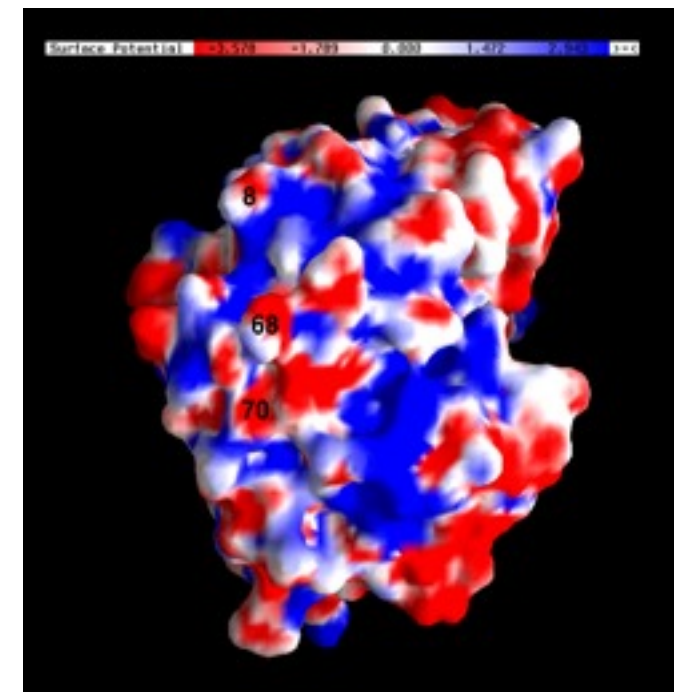
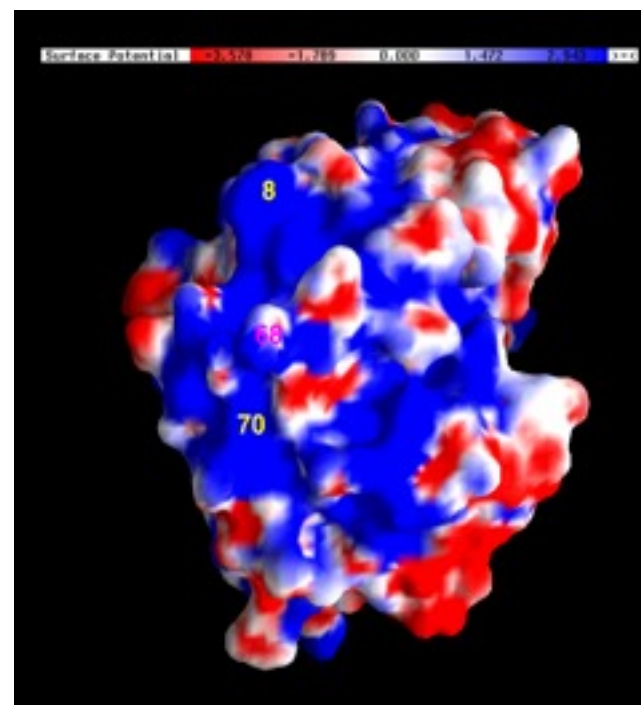
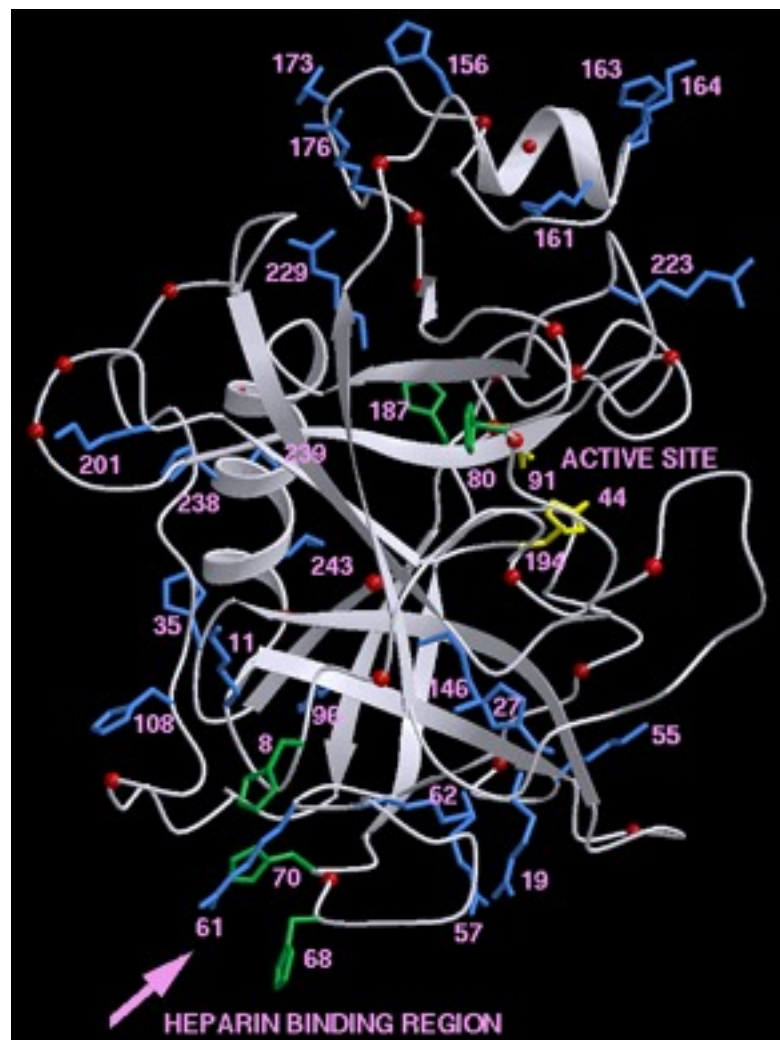
Do mast cell proteases bind proteoglycans? Where? When?

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

1. mMCPs bind negatively charged proteoglycans through electrostatic interactions
2. Comparative models used to find clusters of positively charged surface residues.
3. Tested by site-directed mutagenesis.



Huang *et al.* *J. Clin. Immunol.* **18**,169,1998.
Matsumoto *et al.* *J. Biol. Chem.* **270**,19524,1995.
Šali *et al.* *J. Biol. Chem.* **268**, 9023, 1993.



Structural analysis of missense mutations in human BRCA1 BRCT domains

Mirkovic et al. Structure-based assessment of missense mutations in human BRCA1: implications for breast and ovarian cancer predisposition. Cancer Res (2004) vol. 64 (11) pp. 3790-7

[CANCER RESEARCH 64, 3790-3797, June 1, 2004]

Structure-Based Assessment of Missense Mutations in Human BRCA1: Implications for Breast and Ovarian Cancer Predisposition

Nebojsa Mirkovic,¹ Marc A. Marti-Renom,² Barbara L. Weber,³ Andrej Sali,² and Alvaro N. A. Monteiro^{4,5}

¹Laboratory of Molecular Biophysics, Pels Family Center for Biochemistry and Structural Biology, Rockefeller University, New York, New York; ²Departments of Biopharmaceutical Sciences and Pharmaceutical Chemistry, and California Institute for Quantitative Biomedical Research, University of California at San Francisco, San Francisco, California; ³Abramson Family Cancer Research Institute, University of Pennsylvania, Philadelphia, Pennsylvania; ⁴Strang Cancer Prevention Center, New York, New York; and ⁵Department of Cell and Developmental Biology, Weill Medical College of Cornell University, New York, New York

ABSTRACT

The *BRCA1* gene from individuals at risk of breast and ovarian cancers can be screened for the presence of mutations. However, the cancer association of most alleles carrying missense mutations is unknown, thus creating significant problems for genetic counseling. To increase our ability to identify cancer-associated mutations in *BRCA1*, we set out to use the principles of protein three-dimensional structure as well as the correlation between the cancer-associated mutations and those that abolish transcriptional activation. Thirty-one of 37 missense mutations of known impact on the transcriptional activation function of BRCA1 are readily rationalized in structural terms. Loss-of-function mutations involve non-conservative changes in the core of the BRCA1 C-terminus (BRCT) fold or are localized in a groove that presumably forms a binding site involved in the transcriptional activation by BRCA1; mutations that do not abolish transcriptional activation are either conservative changes in the core or are on the surface outside of the putative binding site. Next, structure-based rules for predicting functional consequences of a given missense mutation were applied to 57 germ-line BRCA1 variants of unknown cancer association. Such a structure-based approach may be helpful in an integrated effort to identify mutations that predispose individuals to cancer.

INTRODUCTION

Many germ-line mutations in the human *BRCA1* gene are associated with inherited breast and ovarian cancers (1, 2). This information has allowed clinicians and genetic counselors to identify individuals at high risk for developing cancer. However, the disease association of over 350 missense mutations remains unclear, primarily because their relatively low frequency and ethnic specificity limit the usefulness of the population-based statistical approaches to identifying cancer-causing mutations. To address this problem, we use here the three-dimensional structure of the human BRCA1 BRCT domains to assess the transcriptional activation functions of BRCA1 mutants. Our study is made possible by the recently determined sequences (3–6) and three-dimensional structures of the BRCA1 homologs (7, 8). In addition, we benefited from prior studies that attempted to rationalize and predict functional effects of mutations in various proteins (9–12), including those of BRCA1 (13, 14).

BRCA1 is a nuclear protein that activates transcription and facilitates DNA damage repair (15, 16). The tandem BRCT domains at the

COOH-terminus of BRCA1 are involved in several of its functions, including modulation of the activity of several transcription factors (15), binding to the RNA polymerase II holoenzyme (17), and activating transcription of a reporter gene when fused to a heterologous DNA-binding domain (18, 19). Importantly, cancer-associated mutations in the BRCT domains, but not benign polymorphisms, inactivate transcriptional activation and binding to RNA polymerase II (18–21). These observations suggest that abolishing the transcriptional activation function of BRCA1 leads to tumor development and provides a genetic framework for characterization of BRCA1 BRCT variants.

MATERIALS AND METHODS

The multiple sequence alignment (MSA) of orthologous BRCA1 BRCT domains from seven species, including *Homo sapiens* (GenBank accession number U14680), *Pan troglodytes* (AF207822), *Mus musculus* (U68174), *Rattus norvegicus* (AF036760), *Gallus gallus* (AF355273), *Canis familiaris* (U50709), and *Xenopus laevis* (AF416868), was obtained by using program ClustalW (22) and contains only one gapped position (Supplementary Fig. 1). According to PSI-BLAST (23), the latter six sequences are the only sequences in the nonredundant protein sequence database at National Center for Biotechnology Information that have between 30% and 90% sequence identity to the human BRCA1 BRCT domains (residues 1649–1859).

The multiple structure-based alignment of the native structures of the BRCT-like domains was obtained by the SALIGN command in MODELLER (Supplementary Fig. 2). It included the experimentally determined structures of the two human BRCA1 BRCT domains (Protein Data Bank code 1JNX; Refs. 8, 24), rat BRCA1 BRCT domains (1L0B; Ref. 7), human p53-binding protein (1KZY; Ref. 7), human DNA-ligase III α (1IMO; Ref. 25), and human XRCC1 protein (1CDZ; Ref. 13). Structure variability was defined by the root-mean-square deviation among the superposed C α positions, as calculated by the COMPARE command of MODELLER. The purpose of these calculations was to gain insight into the variability of surface-exposed residues (left panel in Fig. 2). In conjunction with observed mutation clustering, these data may point to putative functional site(s) on the surface of BRCT repeats.

Comparative protein structure modeling by satisfaction of spatial restraints, implemented in the program MODELLER-6 (26), was used to produce a three-dimensional model for each of the 94 mutants. The crystallographic structure of the human wild-type BRCA1 BRCT domains was used as the template for modeling (8). The four residues missing in the crystallographic structure (1694 and 1817–1819) were modeled *de novo* (27). All of the models are available in the BRCA1 model set deposited in our ModBase database of comparative protein structure models (28).⁶

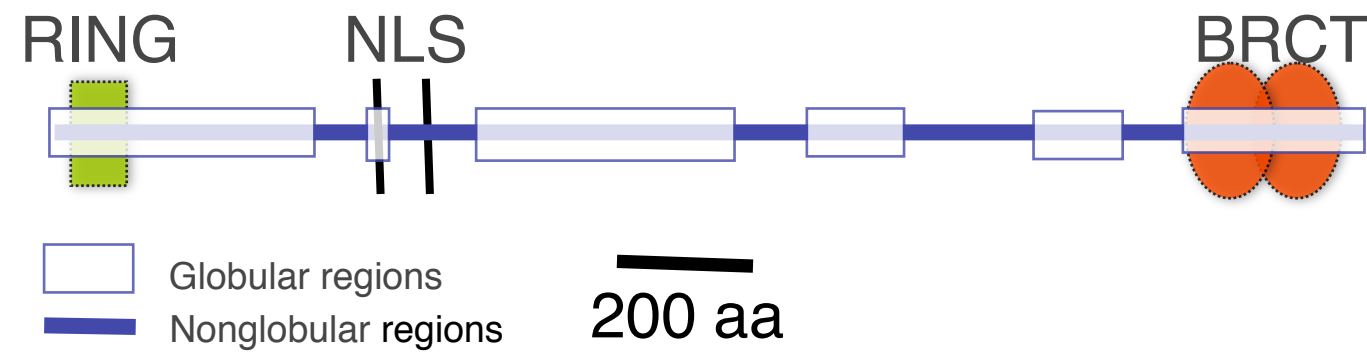
For the native structure of the human BRCT tandem repeat and each of the 94 mutant models, a number of sequence and structure features were calculated. These features were used in the classification tree in Fig. 3 (values for all 94 mutations are given in Supplementary Tables 1 and 2).

Buriedness. Accessible surface area of an amino acid residue was calculated by the program DSSP (29) and normalized by the maximum accessible surface area for the corresponding amino acid residue type. A residue was considered exposed if its accessible surface area was larger than 40Å² and if its relative accessible surface area was larger than 9% and buried otherwise. A mutation of a more exposed residue is less likely to change the structure and therefore its function.

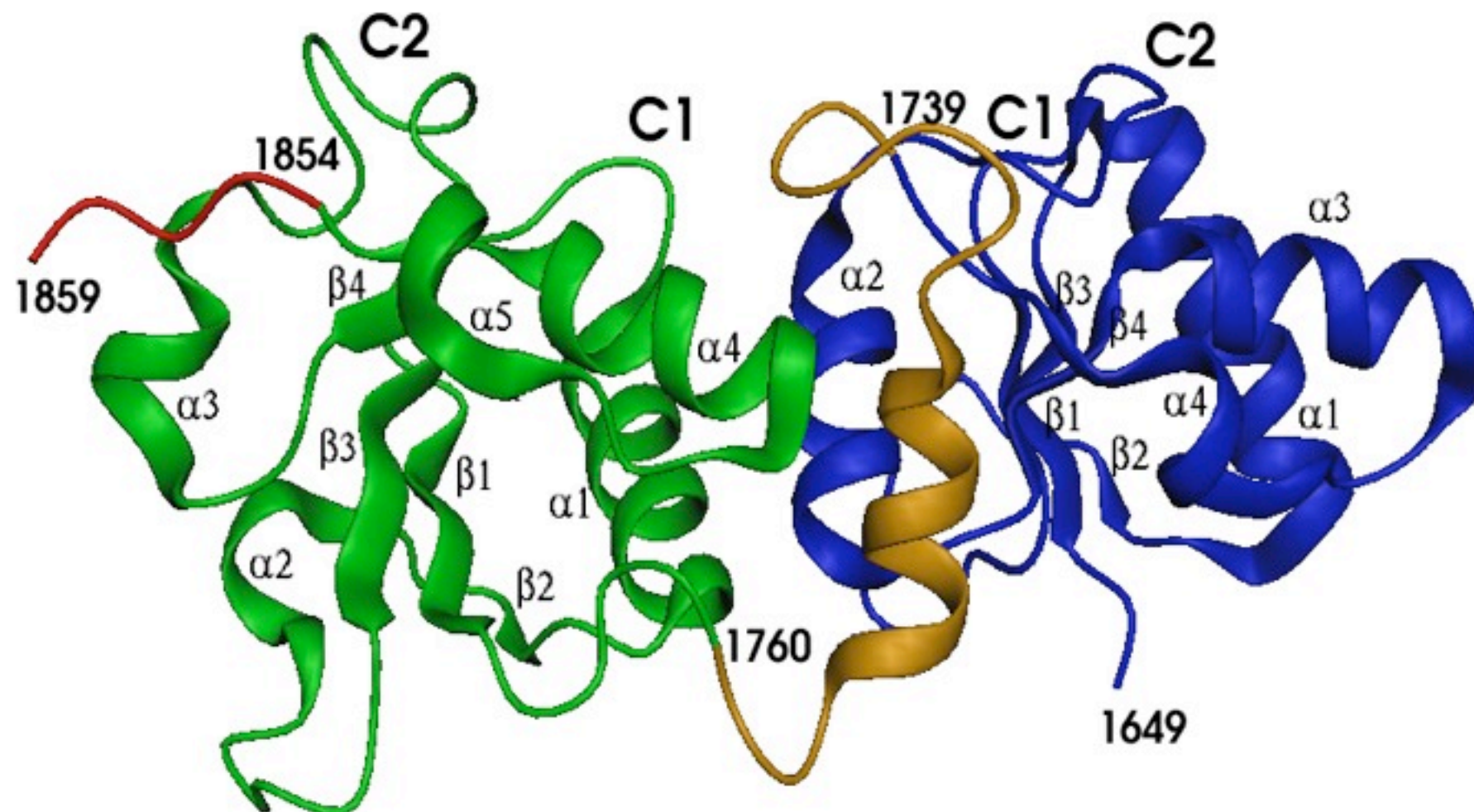
⁶ <http://salilab.org/modbase/>.



Human BRCA1 and its two BRCT domains



BRCA1 BRCT repeats, 1jnx



Williams, Green, Glover. *Nat.Struct.Biol.* 8, 838, 2001

CONFIDENTIAL



MYRIAD

BRCAAnalysis™

Comprehensive BRCA1-BRCA2 Gene Sequence Analysis Result

Niecee Singer, MS
Strang Cancer Prevention Center
428 E 72nd St
New York, NY 10021

SPECIMEN
Specimen Type: Blood
Draw Date: n/a
Accession Date: Oct 27, 2000
Report Date: Nov 17, 2000

PATIENT
Name:
Date of Birth: Feb 02, 1953
Patient ID:
Gender: Female
Accession #: 00019998
Requisition #: 56694

Physician: Fred Gilbert, MD

Test Result

| Gene Analyzed | Specific Genetic Variant |
|---------------|--------------------------|
| BRCA2 | H2116R |
| BRCA1 | None Detected |

Interpretation

GENETIC VARIANT OF UNCERTAIN SIGNIFICANCE

The BRCA2 variant H2116R results in the substitution of arginine for histidine at amino acid position 2116 of the BRCA2 protein. Variants of this type **may or may not** affect BRCA2 protein function. Therefore, the contribution of this variant to the relative risk of breast or ovarian cancer cannot be established solely from this analysis. The observation by Myriad Genetic Laboratories of this particular variant in an individual with a deleterious truncating mutation in BRCA2, however, reduces the likelihood that H2116R is itself deleterious.

Authorized Signature:

Brian E. Ward, Ph.D.
Laboratory Director

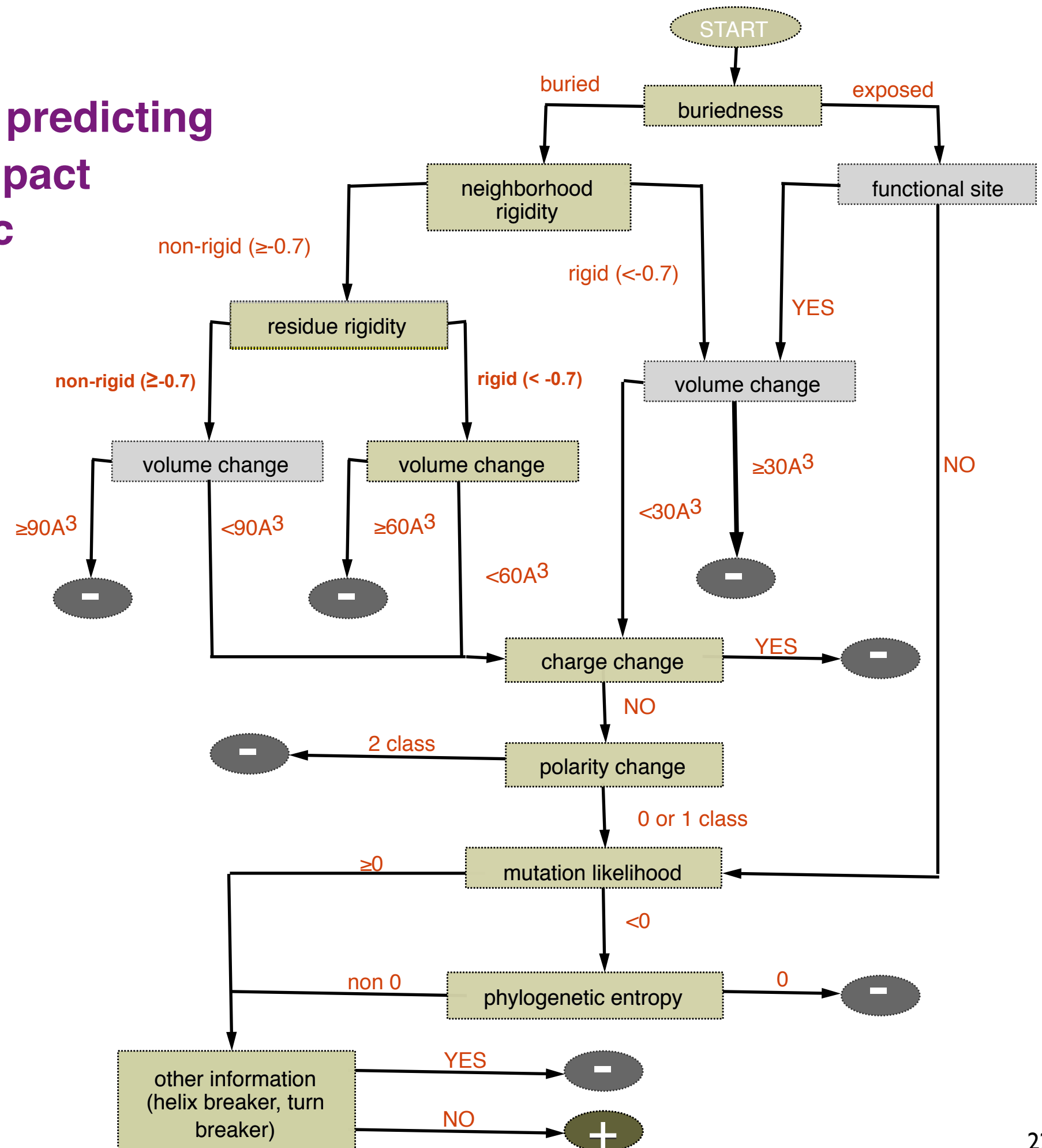

Thomas S. Frank, M.D.
Medical Director

These test results should only be used in conjunction with the patient's clinical history and any previous analysis of appropriate family members. It is strongly recommended that these results be communicated to the patient in a setting that includes appropriate counseling. The accompanying Technical Specifications summary describes the analysis, method, performance characteristics, nomenclature, and interpretive criteria of this test. This test may be considered investigational by some states. This test was developed and its performance characteristics determined by Myriad Genetic Laboratories. It has not been reviewed by the U.S. Food and Drug Administration. The FDA has determined that such clearance or approval is not necessary.

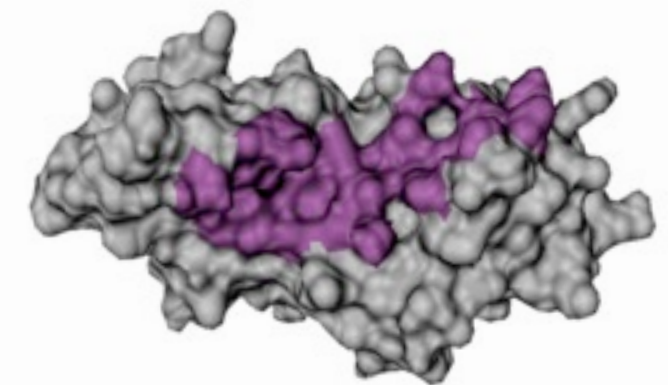
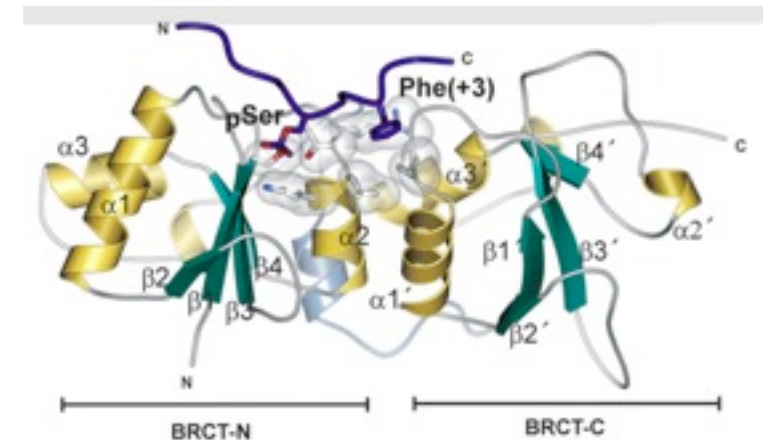
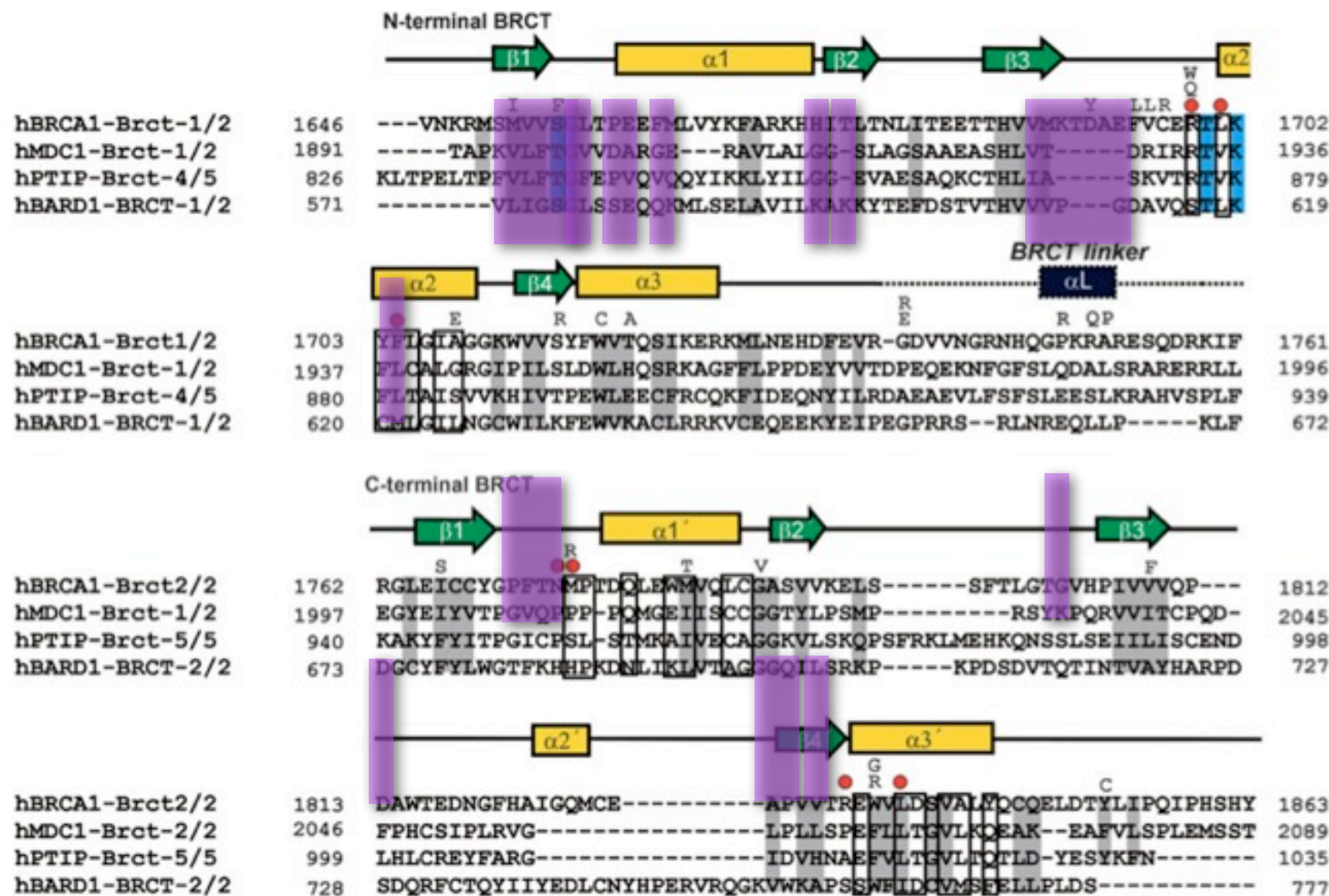
Missense mutations in BRCT domains by function

| | cancer associated | 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 RA1752P F1761I | 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 |

“Decision” tree for predicting functional impact of genetic variants



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

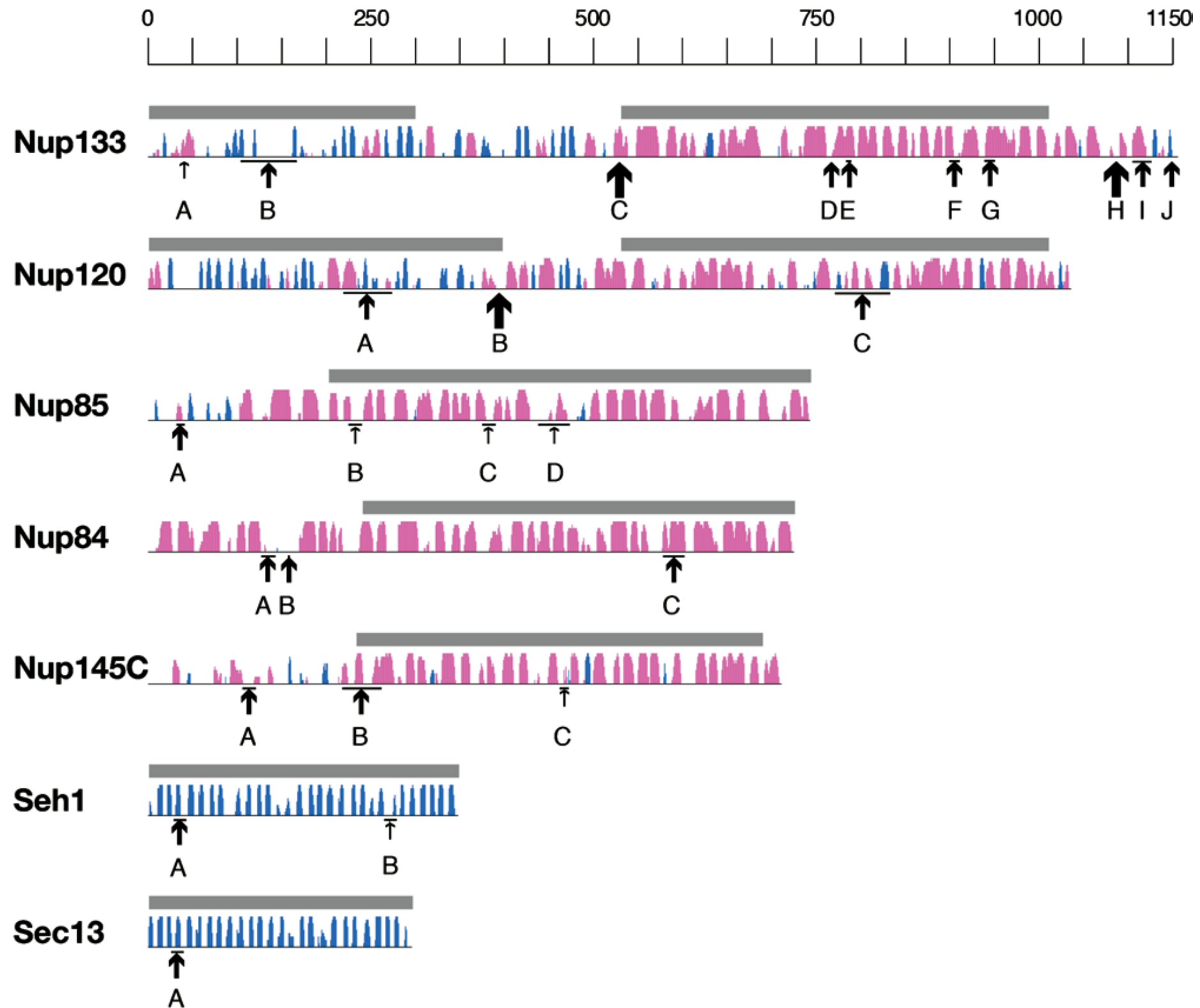
Mirkovic *et al.* 2004 Cancer Research. June 2004 64:3790

Common Evolutionary Origin of Coated Vesicles and Nuclear Pore Complexes

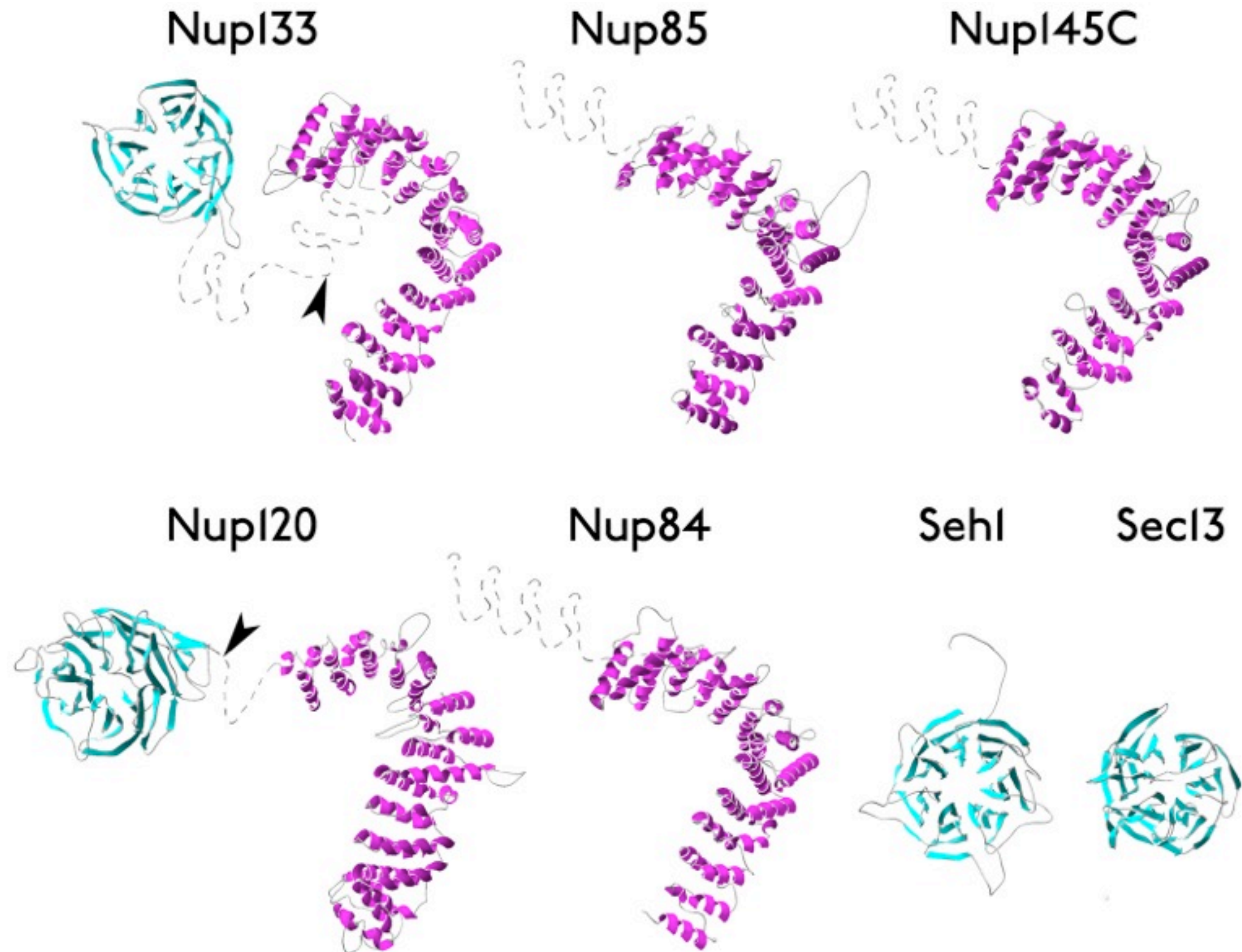
mGenThreader + SALIGN + MOULDER

D. Devos, S. Dokudovskaya, F. Alber, R. Williams, B.T. Chait, A. Sali, M.P. Rout.
Components of Coated Vesicles and Nuclear Pore Complexes Share a Common Molecular Architecture.
PLOS Biology **2(12)**:e380, 2004

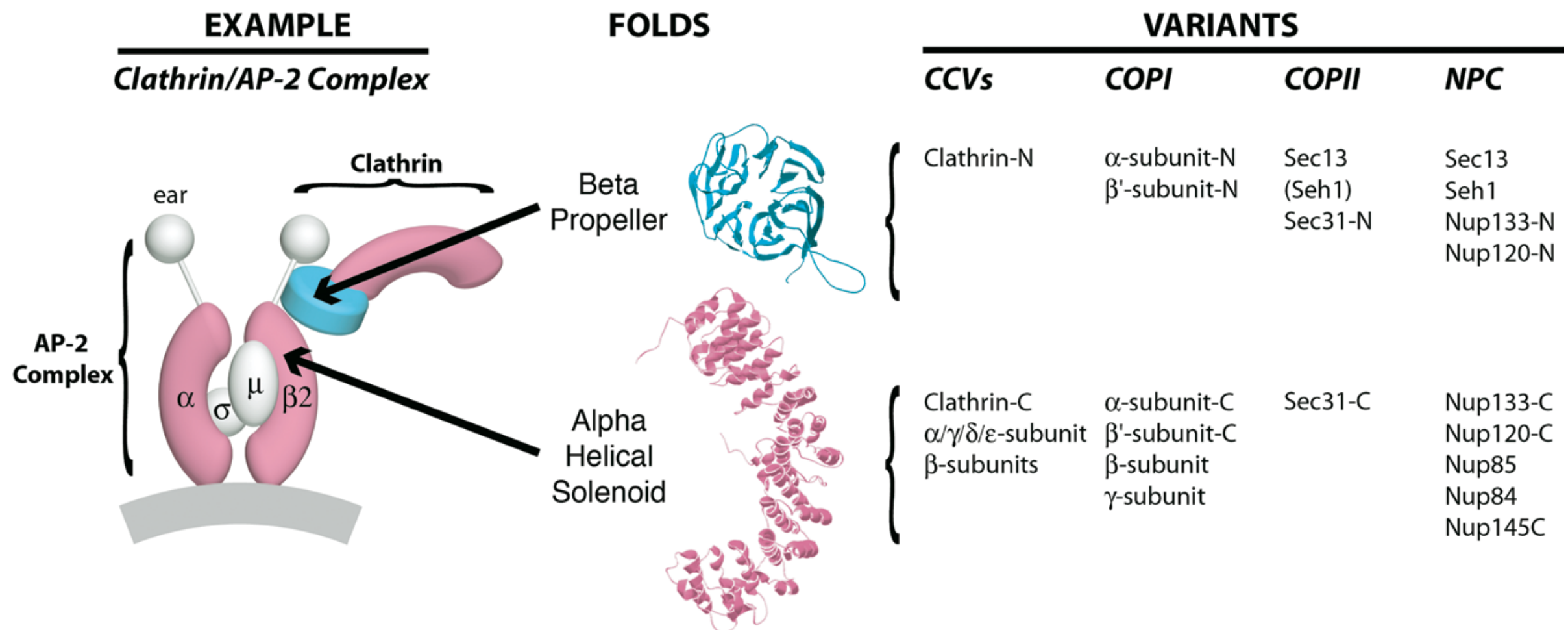
yNup84 complex proteins



All Nucleoporins in the Nup84 Complex are Predicted to Contain β -Propeller and/or α -Solenoid Folds



NPC and Coated Vesicles Share the β -Propeller and α -Solenoid Folds and Associate with Membranes

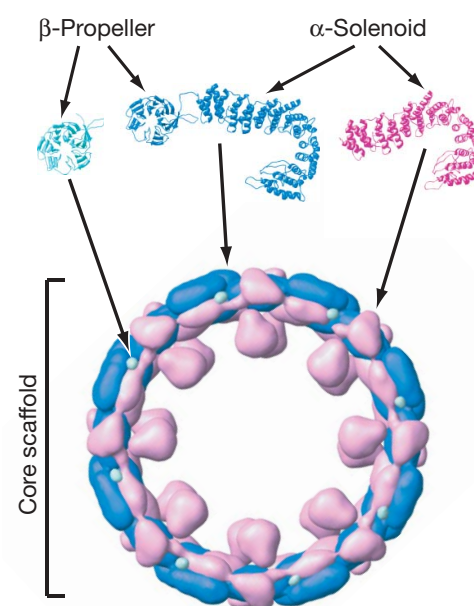
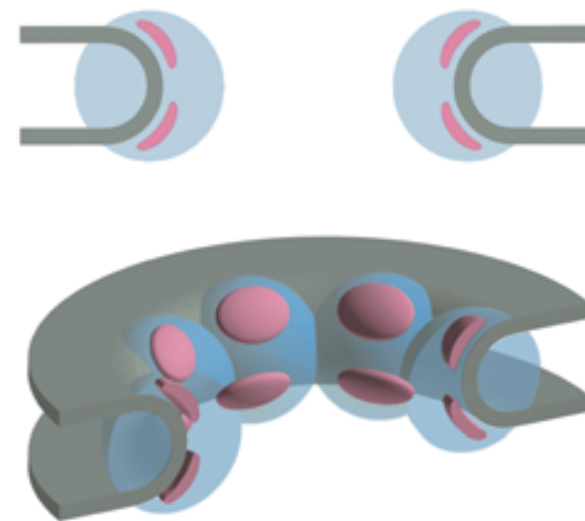


NPC and Coated Vesicles Both Associate with Membranes

Coated Vesicle

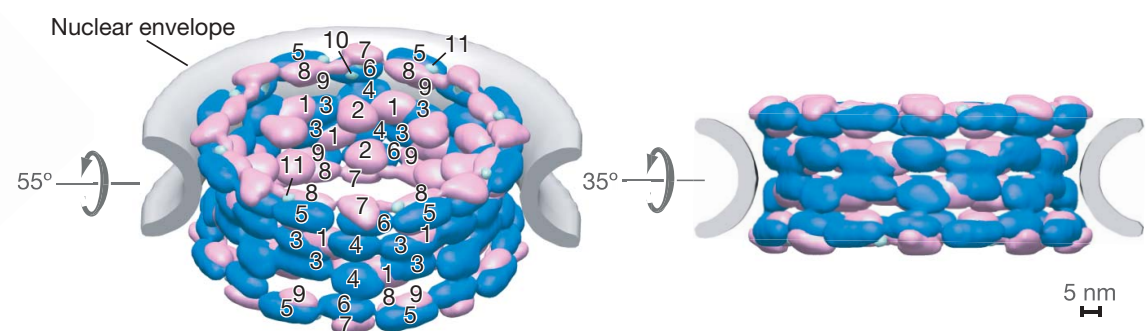


NPC model



Nup 84 complex

1 Nup192, 2 Nup188, 3 Nup170, 4 Nup157, 5 Nup133,
6 Nup120, 7 Nup85, 8 Nup84, 9 Nup145C, 10 Seh1, 11 Sec13

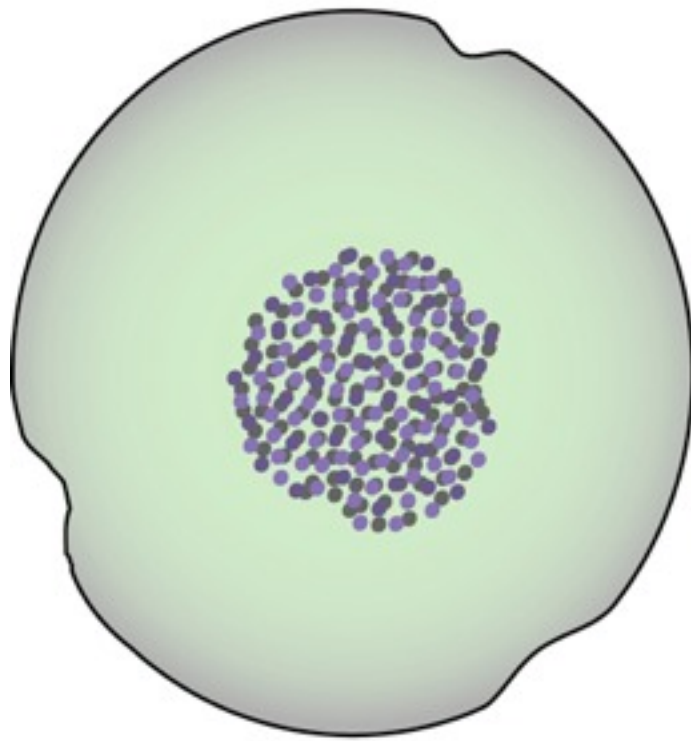


Alber et al. The molecular architecture of the nuclear pore complex. Nature (2007) vol. 450 (7170) pp. 695-701

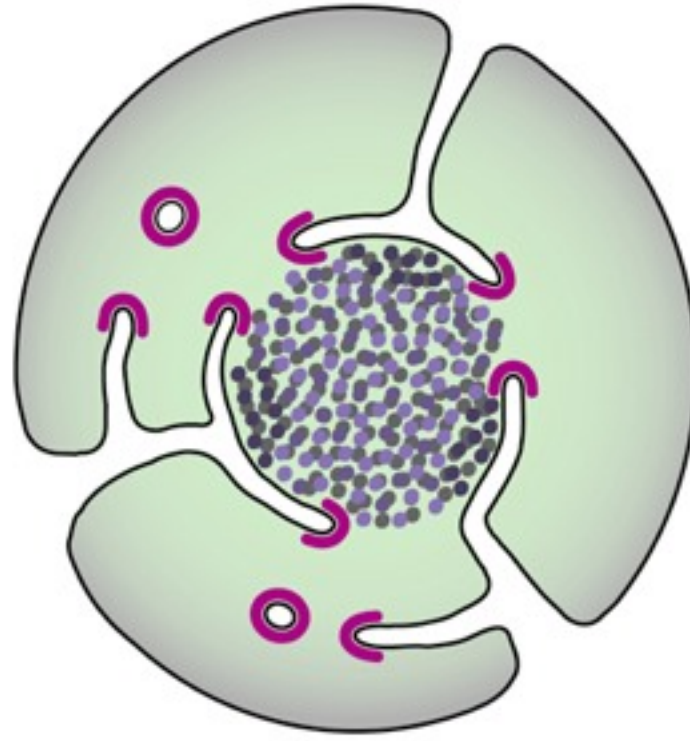
A Common Evolutionary Origin for Nuclear Pore Complexes and Coated Vesicles?

The proto-coatomer hypothesis

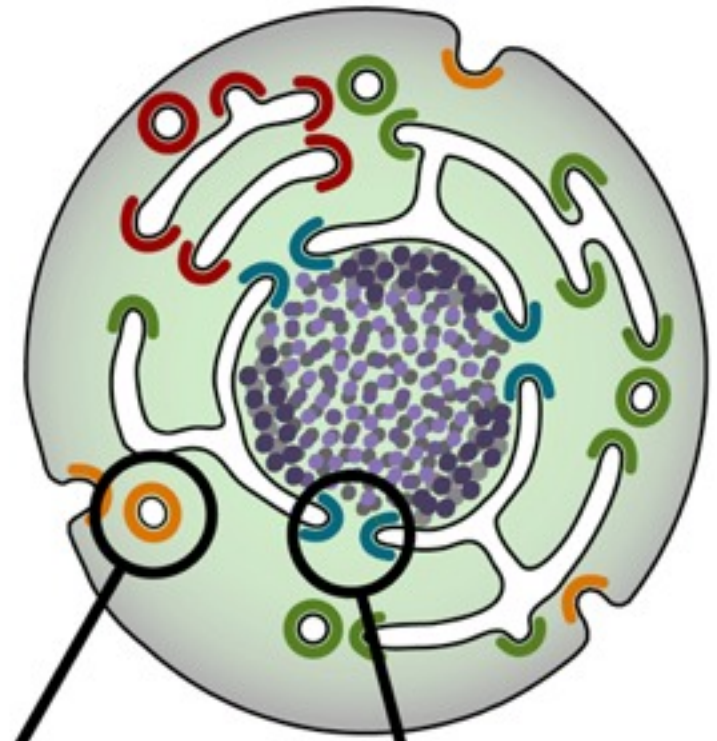
Prokaryote



Early Eukaryote

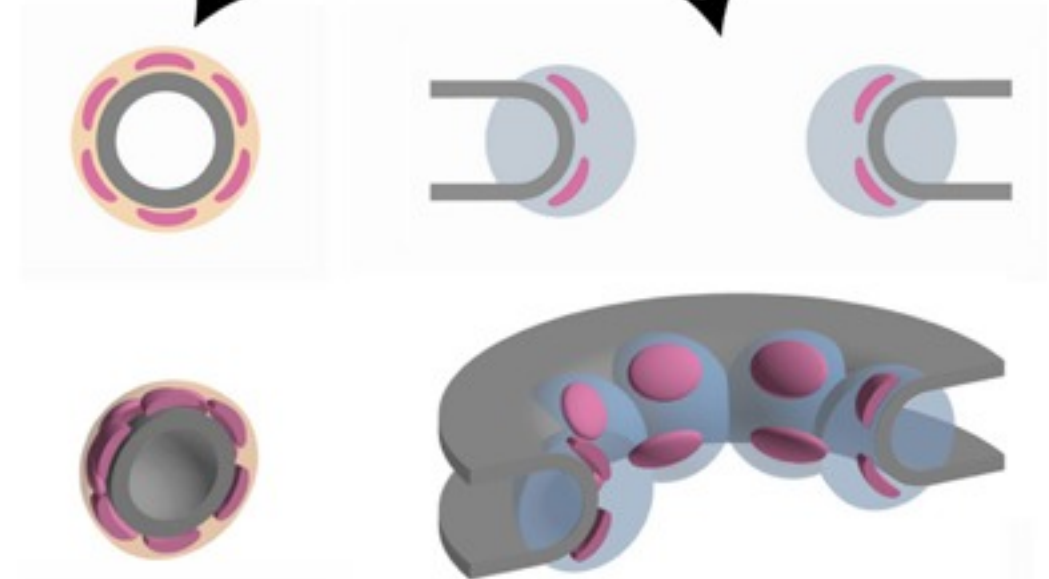


Modern Eukaryote



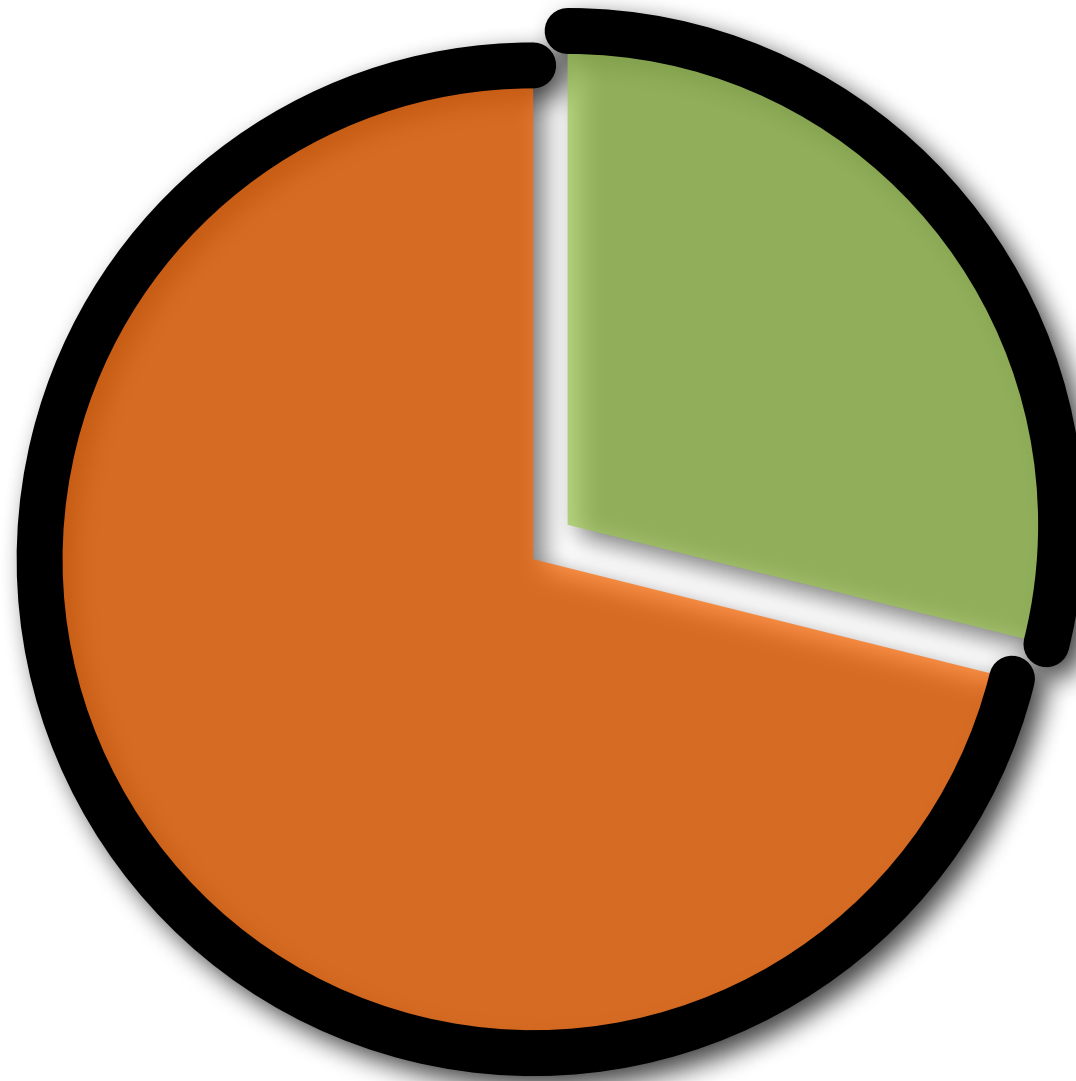
A simple coating module containing minimal copies of the two conserved folds evolved in proto-eukaryotes to bend membranes.

The progenitor of the NPC arose from a membrane-coating module that wrapped extensions of an early ER around the cell's chromatin.



Tropical Disease Initiative (TDI)

Predicting binding sites in protein structure models.



<http://www.tropicaldisease.org>



UCSF

Duke
UNIVERSITY

PRINCIPE FELIPE
CENTRO DE INVESTIGACION
CEREBRO DE MANIZABOY
E. SANTAFÉ DE BOGOTÁ

TDI *a story*



2004

- .Steve Maurer (Berkeley) and Arti Rai (Duke)
- .PLoS Medicine, Dec. 2004. Vol 1(3):e56

2005

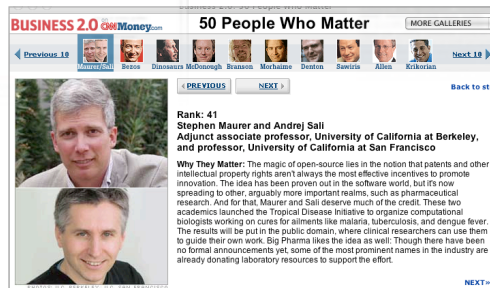
- .TDI web site <http://TropicalDisease.org>
- .Ginger Taylor and The Synaptic Leap

2006

- .Maurer and Sali 41th in “50 Who Matter”
- .TSL web site <http://TheSynapticLeap.org>

2008

- .TDI kernel <http://TropicalDisease.org/kernel>



Initial feed-back...

14 Mar 2005

I think TDI is a unique and very interesting project. I v
it...

So, where are we going? What's happening? What

I still trust in open s

Luca Brivio

16 Feb 2005

Hi,

10 Feb 2005

Hello,

My name is Adam Huber and I am a medical student at UNSW in Sydney Australia.
I am interested in beginning research focused on tropical and infectious
disease for underserved populations (A mission that seemingly matches TDI). I am,

bottlenecks are?

potential avenues to explore,

n!

9 Mar 2005

I'm a programmer, not a bioinformatician, but I stumbled across your site and thought I'd say something to keep
the list active :)

GNU started with RMS. He gave us programming/administration tools to play with.

Linux started with Linus. He released an operating system for us to play with.

**You need someone great in the field to release something for everyone to 'play with'. Then people start
sending patches...**

I know this is chicken-egg, but someone needs to point this out, since I haven't seen this brought up in the
papers or the website.

And you might consider merging into the bios.net effort mentioned already. Together, you just might reach the
critical mass for things to take off. Consider this like when people jumped off the HURD project to come
together and make linux work.

Daniel Amelang

Stephen Mark Maurer

stic that the rest

Linux distro timeline

Version 7.2 by NPU (nonplusx@gmail.com)

For the latest version, visit kde-files.org

Feel free to modify and spread. Mail me for updates, corrections and source flw/xcf files

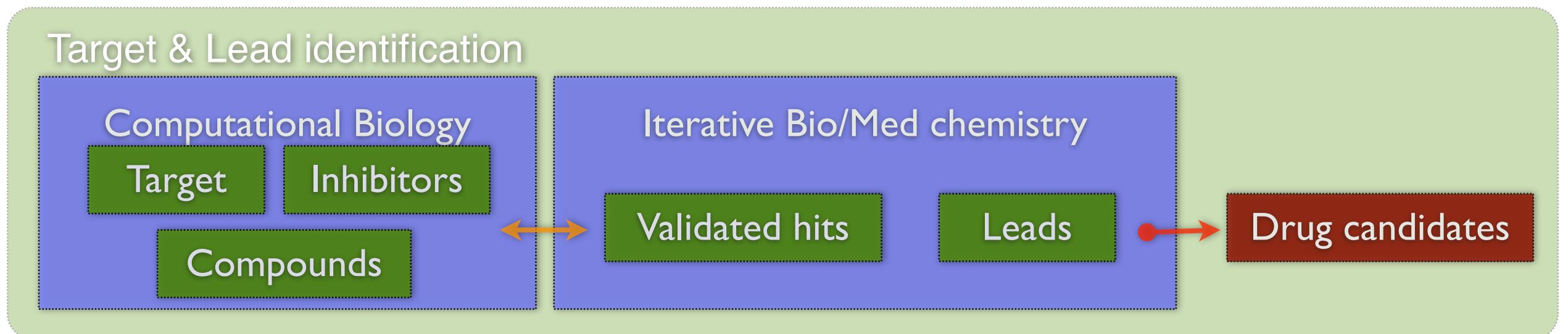
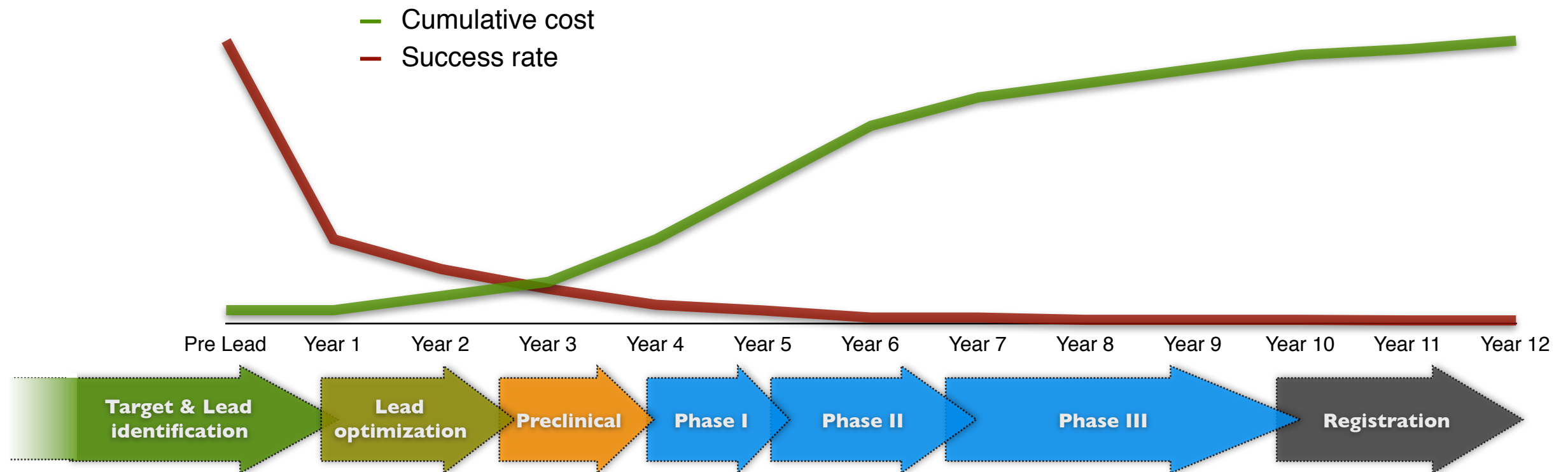
Based on "Línea del tiempo Distribuciones Linux" by A. Sandoval (microtecnologias.cl)

Additional info: distrowatch.com/wikipedia.org

The timeline shows the following distributions and their approximate start years:

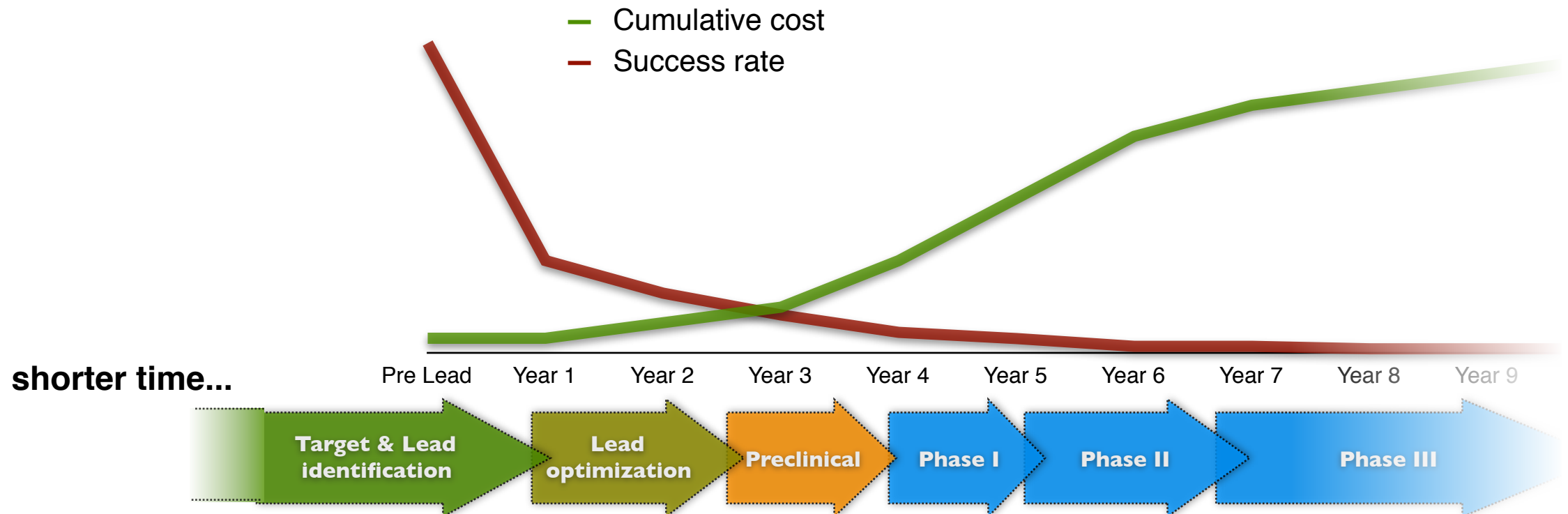
- 1991: GNU/Linux
- 1992: TAMU, MCC Interim, SLS, Slackware, S.u.S.E., Yggdrasil, LST, DLD / Delix, Red Hat
- 1993: Debian
- 1994: Caldera
- 1995: Conectiva
- 1996: Mandrake
- 1997: Turbolinux, Yellow Dog, Peanut, Red Flag
- 1998: Libranet, Storm, Astaro, Lindows, LinEx, Corel, Progeny, Xandros, Yoper, Sorcerer, Source Mage, Lunar, Vector, Beehive, Enoch, Stampede, CRUX, Rock Linux, Linux From Scratch, dyne:bolic, Ark, Redmond, Lytoris, PCLinuxOS, Mandriva, CentOS, Scientific, White Box, Specifix, rPath, Foresight, FoX, Ekaaty, BLAG, PLD, SELinux, EnGarde
- 1999: Ubuntu, Kubuntu, Edubuntu, Damn Small Linux, Symphony OS, KANOPPIX, Kanotix, Morphix, Minislack, Zenwalk, Frugalware, Sun JDS, RR4 / RR64, Sabayon, Kororaa, VidaLinux, Arch
- 2000: Ubuntu, Kubuntu, Edubuntu, Damn Small Linux, Symphony OS, KANOPPIX, Kanotix, Morphix, Minislack, Zenwalk, Frugalware, Sun JDS, RR4 / RR64, Sabayon, Kororaa, VidaLinux, Arch
- 2001: Ubuntu, Kubuntu, Edubuntu, Damn Small Linux, Symphony OS, KANOPPIX, Kanotix, Morphix, Minislack, Zenwalk, Frugalware, Sun JDS, RR4 / RR64, Sabayon, Kororaa, VidaLinux, Arch
- 2002: Ubuntu, Kubuntu, Edubuntu, Damn Small Linux, Symphony OS, KANOPPIX, Kanotix, Morphix, Minislack, Zenwalk, Frugalware, Sun JDS, RR4 / RR64, Sabayon, Kororaa, VidaLinux, Arch
- 2003: Ubuntu, Kubuntu, Edubuntu, Damn Small Linux, Symphony OS, KANOPPIX, Kanotix, Morphix, Minislack, Zenwalk, Frugalware, Sun JDS, RR4 / RR64, Sabayon, Kororaa, VidaLinux, Arch
- 2004: Ubuntu, Kubuntu, Edubuntu, Damn Small Linux, Symphony OS, KANOPPIX, Kanotix, Morphix, Minislack, Zenwalk, Frugalware, Sun JDS, RR4 / RR64, Sabayon, Kororaa, VidaLinux, Arch
- 2005: Ubuntu, Kubuntu, Edubuntu, Damn Small Linux, Symphony OS, KANOPPIX, Kanotix, Morphix, Minislack, Zenwalk, Frugalware, Sun JDS, RR4 / RR64, Sabayon, Kororaa, VidaLinux, Arch
- 2006: Ubuntu, Kubuntu, Edubuntu, Damn Small Linux, Symphony OS, KANOPPIX, Kanotix, Morphix, Minislack, Zenwalk, Frugalware, Sun JDS, RR4 / RR64, Sabayon, Kororaa, VidaLinux, Arch
- 2007: Ubuntu, Kubuntu, Edubuntu, Damn Small Linux, Symphony OS, KANOPPIX, Kanotix, Morphix, Minislack, Zenwalk, Frugalware, Sun JDS, RR4 / RR64, Sabayon, Kororaa, VidaLinux, Arch

Drug Discovery pipeline



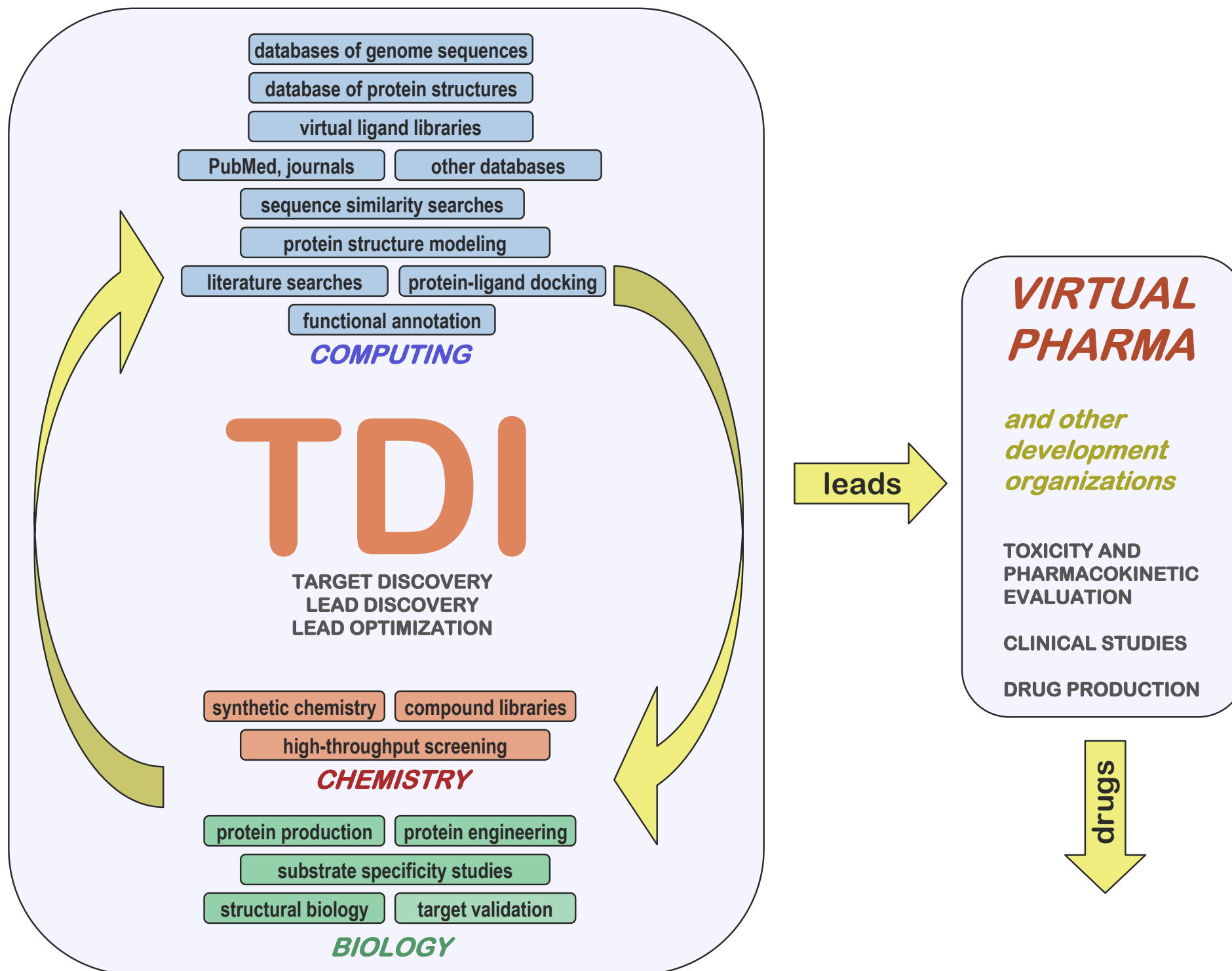
Adapted from: - Nwaka & Ridley. (2003) *Nature Reviews. Drug Discovery*. 2:919
 - Austin, Brady, Insel & collins. (2004) *Science*. 306:1138

Drug Discovery pipeline



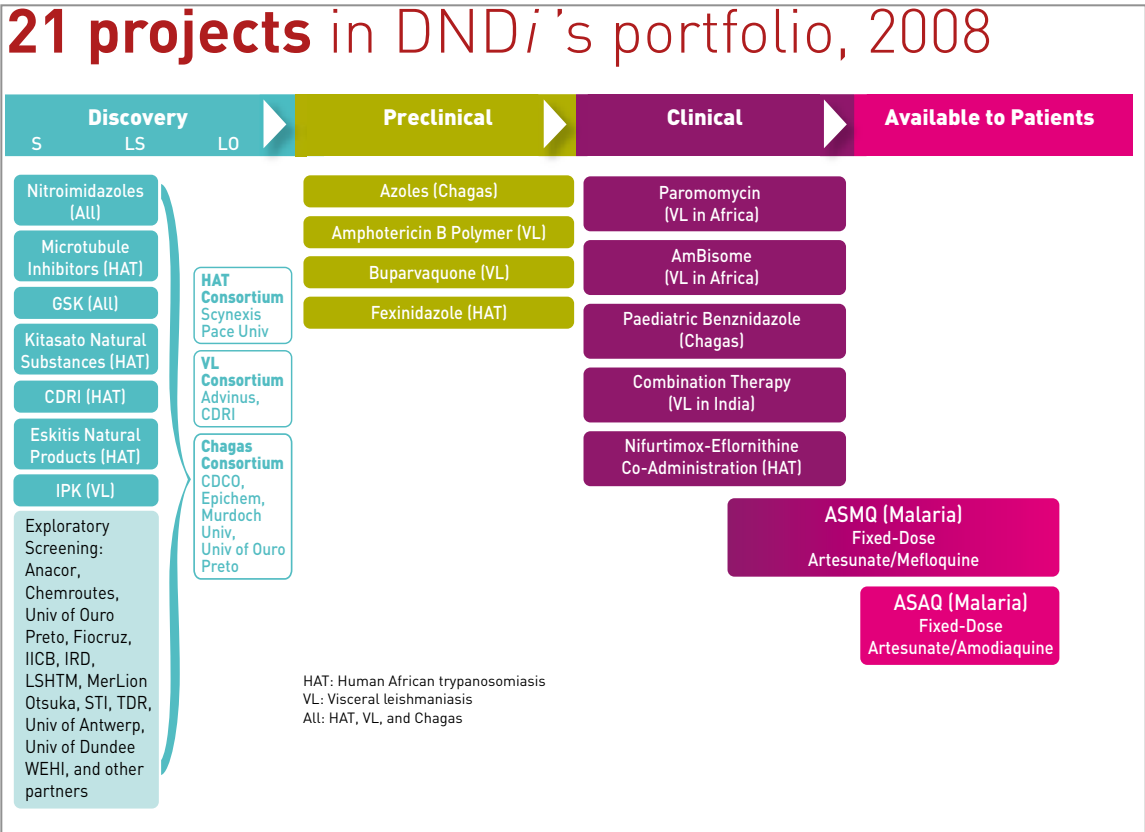
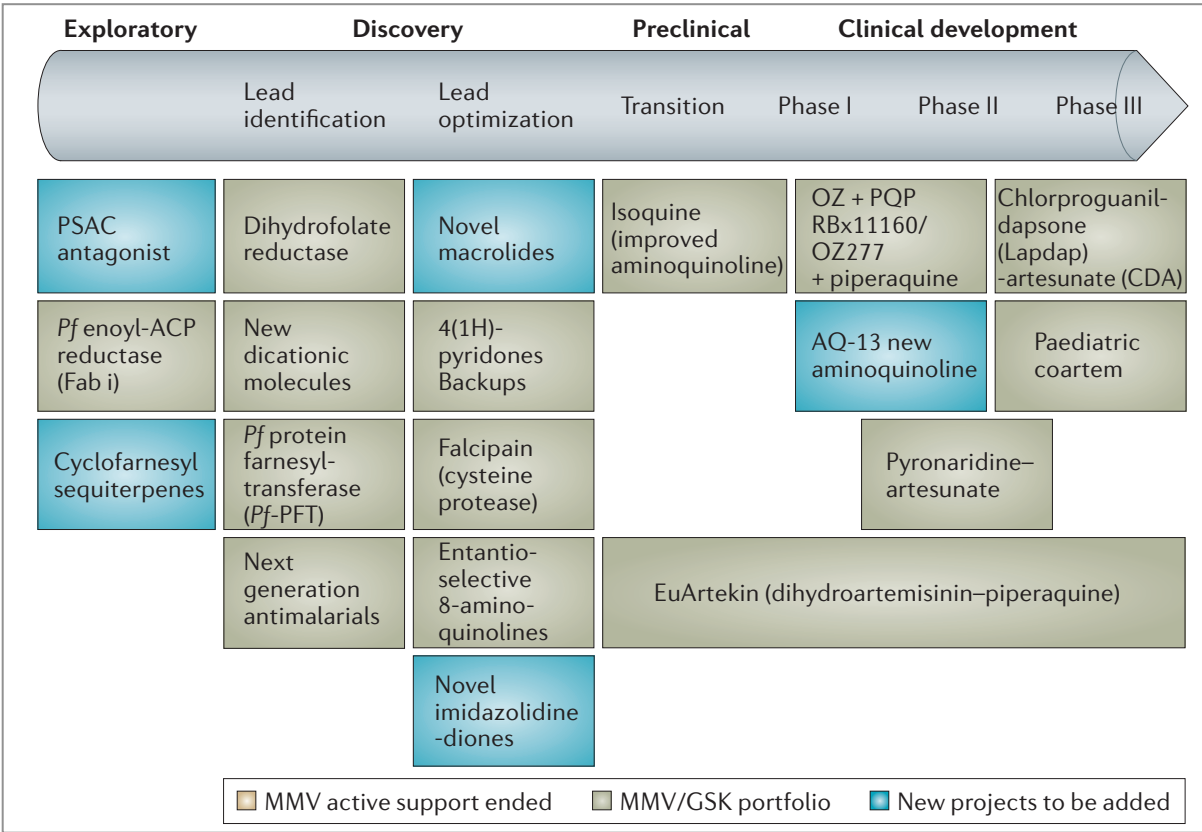
- + Completeness of genome projects (eg, Malaria)
- + New and more complete biological databases
- + New software and computers (cheaper and faster)
- + Internet == more people == less cost

TDI flowchart



Non-Profit organizations

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



Munos (2006) Nature Reviews. Drug Discovery.

Need is High in the Tail

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



Disease data taken from WHO, *World Health Report 2004*

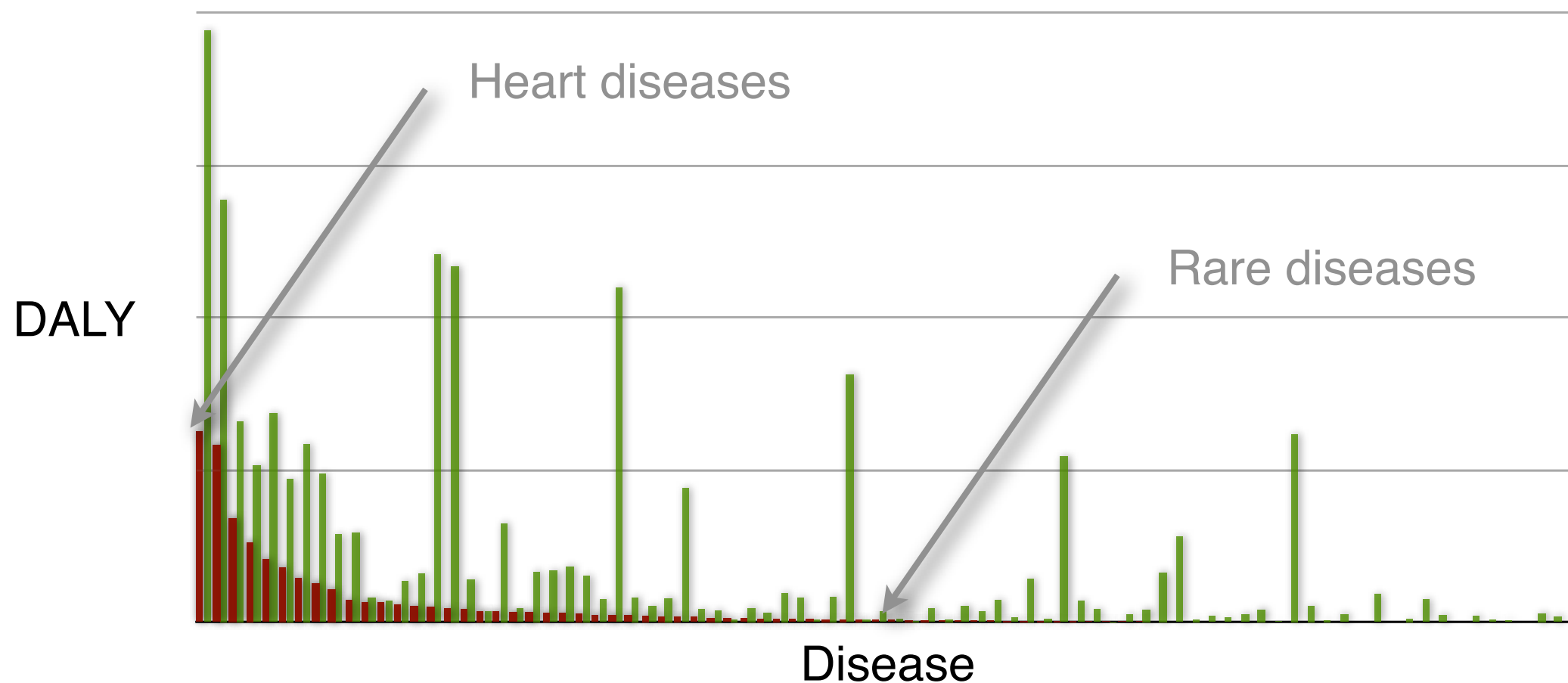
DALY - Disability adjusted life years

DALY is not a perfect measure of market size, but is certainly a good measure for importance.

DALYs for a disease are the sum of the years of life lost due to premature mortality (YLL) in the population and the years lost due to disability (YLD) for incident cases of the health condition. The DALY is a health gap measure that extends the concept of potential years of life lost due to premature death (PYLL) to include equivalent years of 'healthy' life lost in states of less than full health, broadly termed disability. One DALY represents the loss of one year of equivalent full health.

Need is High in the Tail

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



Disease data taken from WHO, *World Health Report 2004*

DALY - Disability adjusted life years

DALY is not a perfect measure of market size, but is certainly a good measure for importance.

DALYs for a disease are the sum of the years of life lost due to premature mortality (YLL) in the population and the years lost due to disability (YLD) for incident cases of the health condition. The DALY is a health gap measure that extends the concept of potential years of life lost due to premature death (PYLL) to include equivalent years of 'healthy' life lost in states of less than full health, broadly termed disability. One DALY represents the loss of one year of equivalent full health.

“Unprofitable” Diseases and Global DALY (in 1000’s)

| | |
|------------------------------|---------------|
| Malaria* | 46,486 |
| Tetanus | 7,074 |
| Lymphatic filariasis* | 5,777 |
| Syphilis | 4,200 |
| Trachoma | 2,329 |
| Leishmaniasis* | 2,090 |
| Ascariasis | 1,817 |
| Schistosomiasis* | 1,702 |
| Trypanosomiasis* | 1,525 |

| | |
|------------------------|------------|
| Trichuriasis | 1,006 |
| Japanese encephalitis | 709 |
| Chagas Disease* | 667 |
| Dengue* | 616 |
| Onchocerciasis* | 484 |
| Leprosy* | 199 |
| Diphtheria | 185 |
| Poliomyelitis | 151 |
| Hookworm disease | 59 |

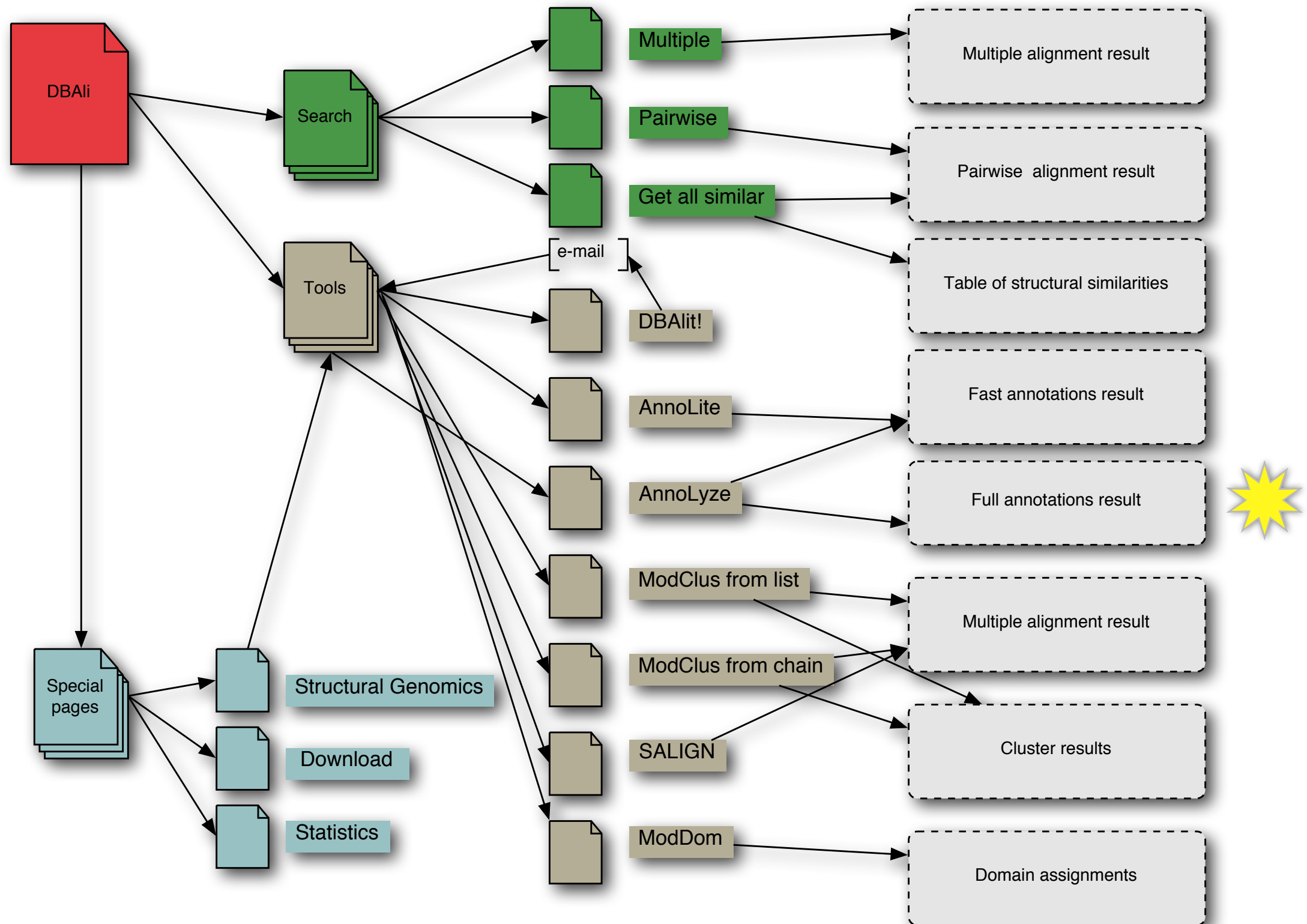
Disease data taken from WHO, *World Health Report 2004*

DALY - Disability adjusted life year in 1000’s.

* Officially listed in the WHO Tropical Disease Research [disease portfolio](#).

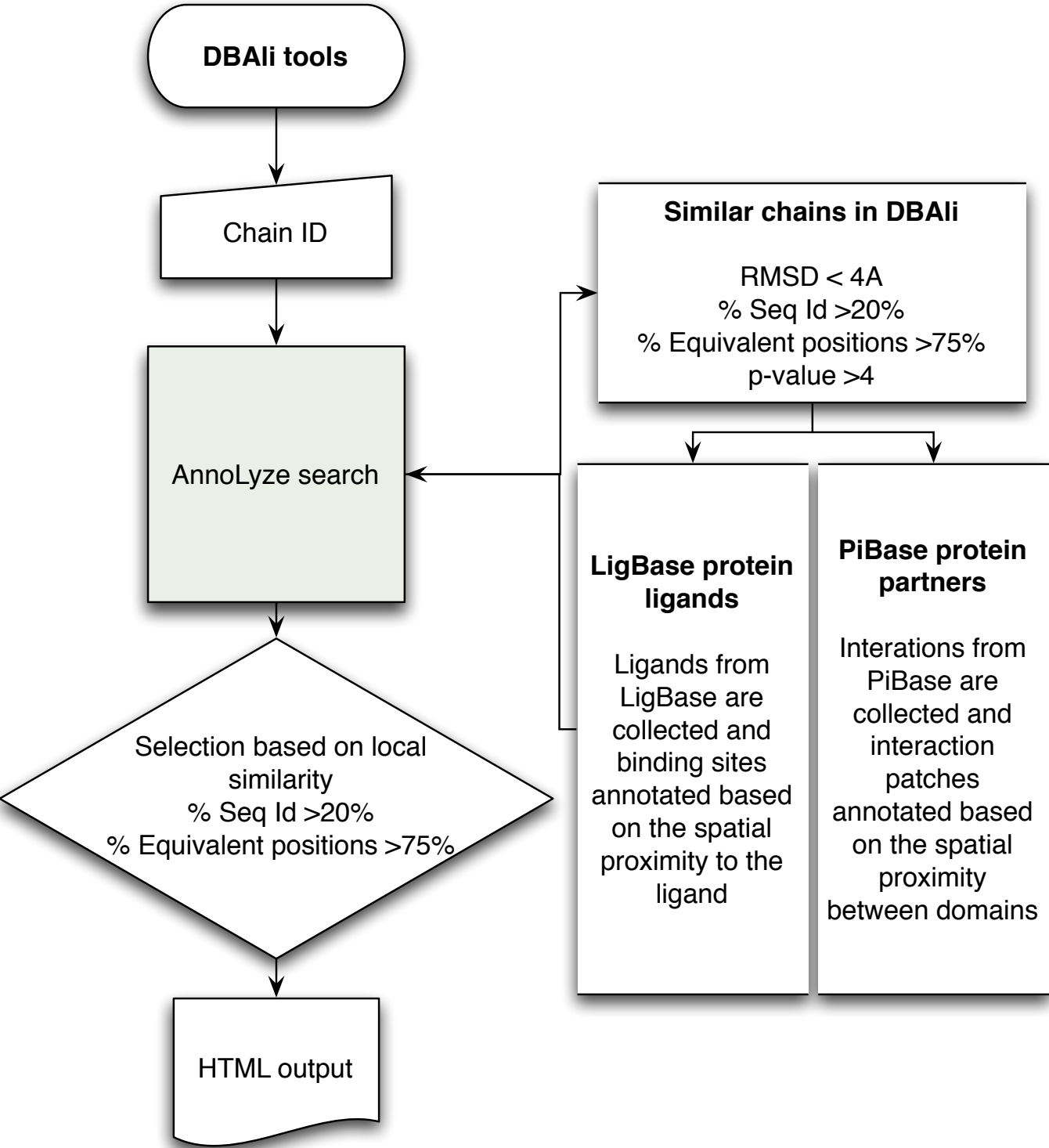
DBAli_{v2.0} database

<http://www.dbali.org>



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

Method



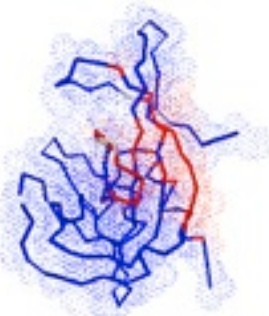
Inherited ligands: 4

| 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 |
| BOG | 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 |



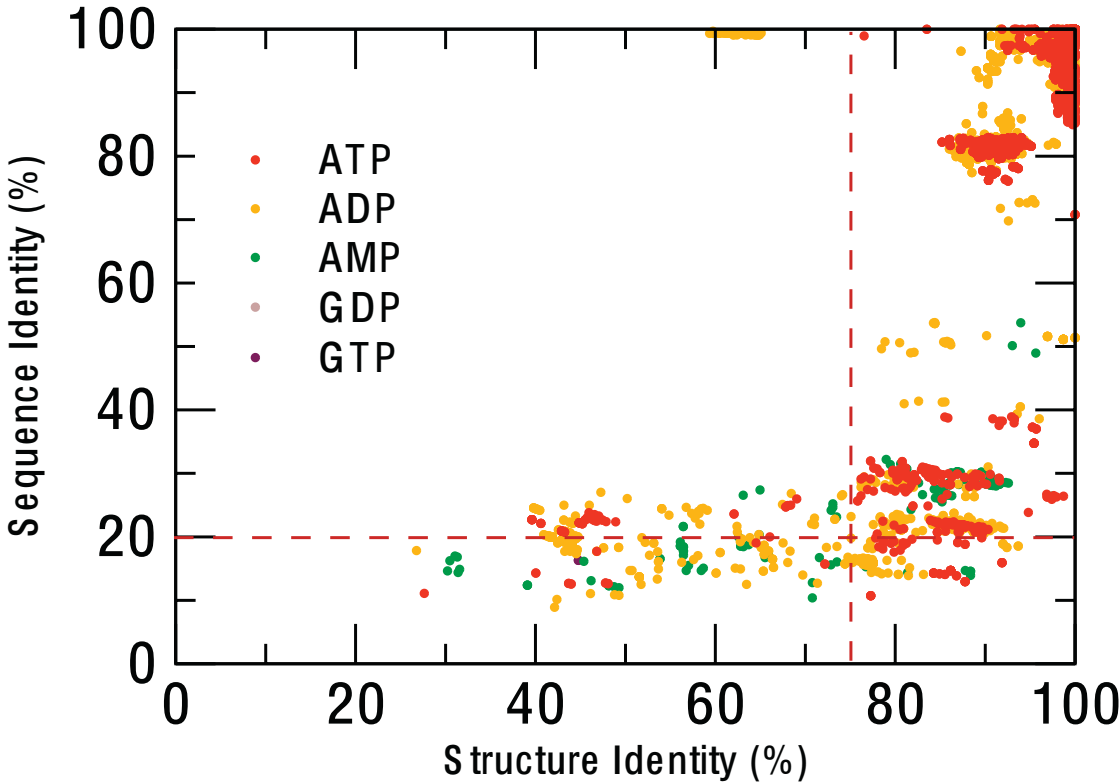
Inherited partners: 1

| 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 |

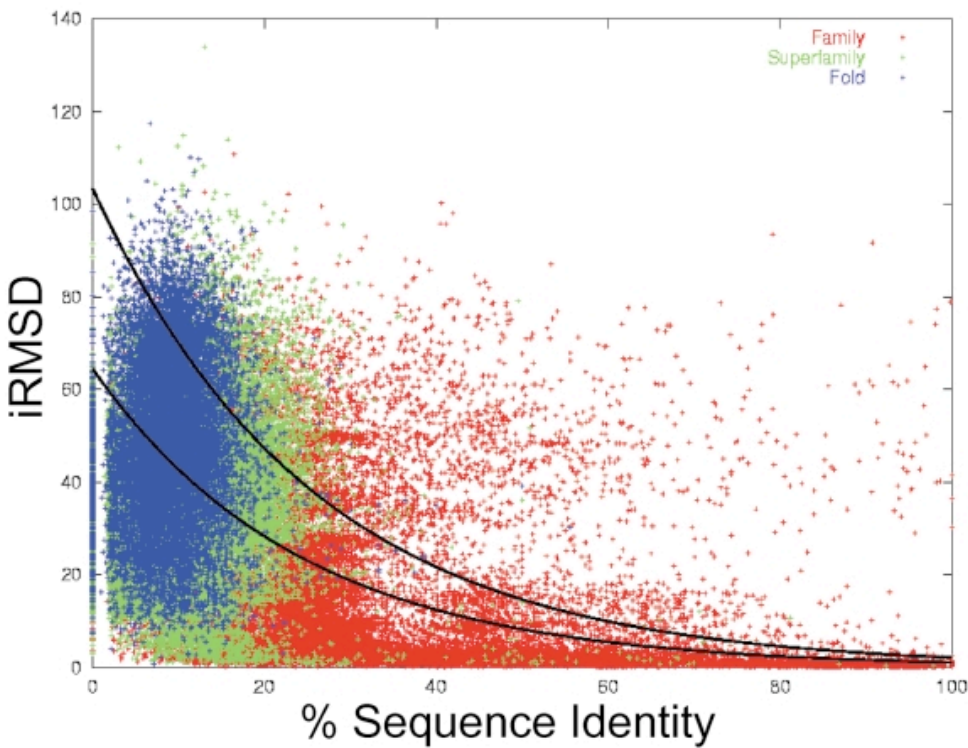


Scoring function

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 3Å, 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*

| | 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 πBase database*

****not two chains can be structurally aligned within 3Å, 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 (%) Recall or TPR | Precision (%) |
|---------|-----------------|----------------------------------|---------------|
| Ligands | 30% | 71.9 | 13.7 |

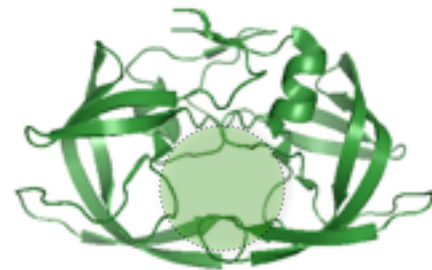
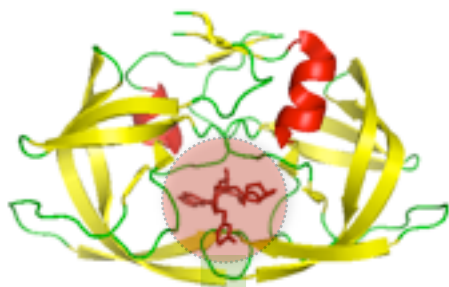
$$\text{Sensitivity} = \frac{TP}{TP + FN} \quad \text{Precision} = \frac{TP}{TP + FP}$$

~90-95% of residues correctly predicted

Comparative docking

Expansion

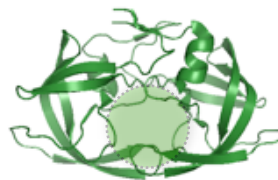
co-crystallized protein/ligand



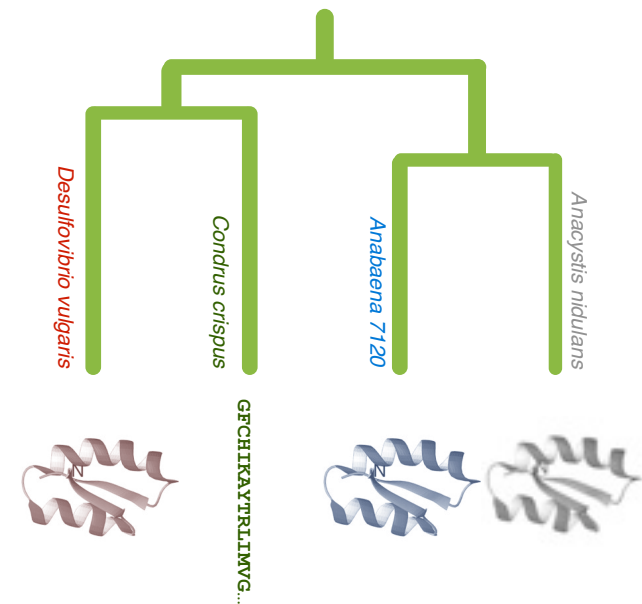
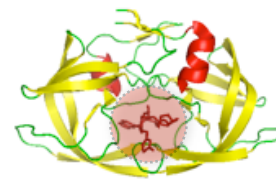
crystallized protein

2. Inheritance

model



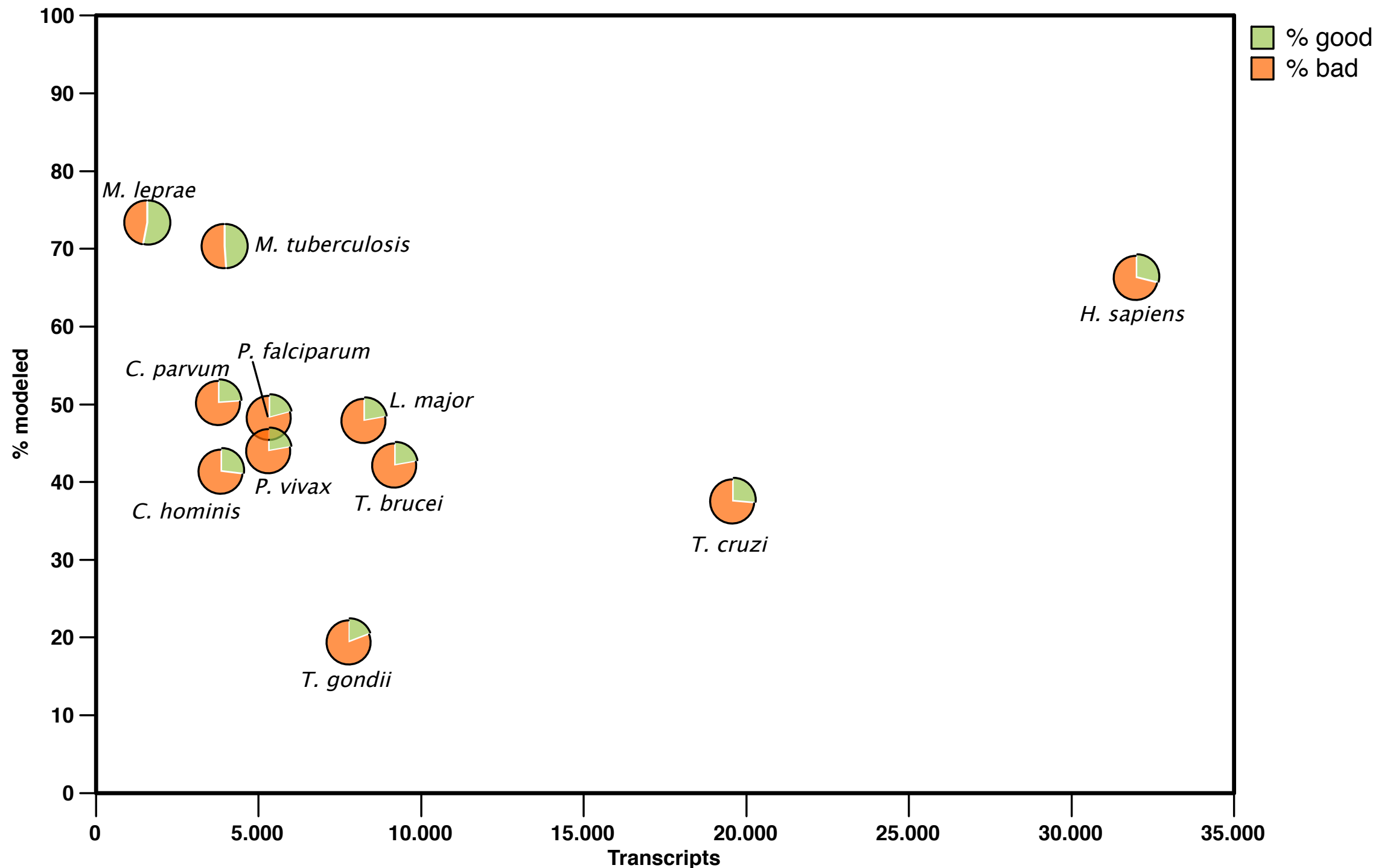
template



1. Modeling

Modeling Genomes

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



A good model has MPQS of 1.0 or higher

Summary table

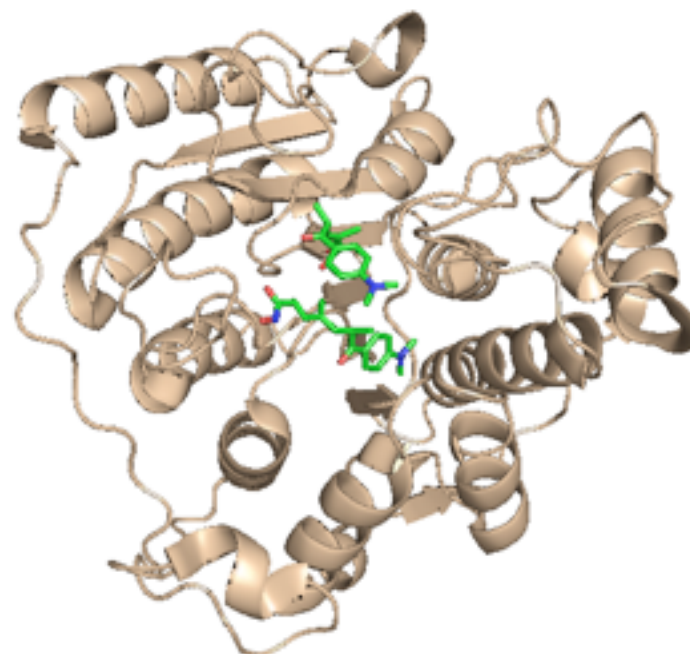
models with inherited ligands

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

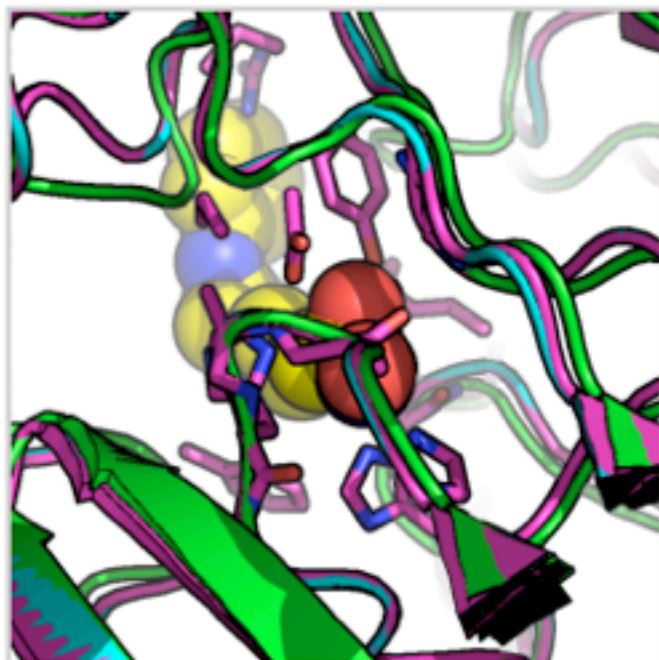
| | Transcripts | Modeled targets | Selected models | Inherited ligands | Similar to a drug | Drugs |
|------------------------|---------------|-----------------|-----------------|-------------------|-------------------|------------|
| <i>C. hominis</i> | 3,886 | 1,614 | 666 | 197 | 20 | 13 |
| <i>C. parvum</i> | 3,806 | 1,918 | 742 | 232 | 24 | 13 |
| <i>L. major</i> | 8,274 | 3,975 | 1,409 | 478 | 43 | 20 |
| <i>M. leprae</i> | 1,605 | 1,178 | 893 | 310 | 25 | 6 |
| <i>M. tuberculosis</i> | 3,991 | 2,808 | 1,608 | 365 | 30 | 10 |
| <i>P. falciparum</i> | 5,363 | 2,599 | 818 | 284 | 28 | 13 |
| <i>P. vivax</i> | 5,342 | 2,359 | 822 | 268 | 24 | 13 |
| <i>T. brucei</i> | 7,793 | 1,530 | 300 | 138 | 13 | 6 |
| <i>T. cruzi</i> | 19,607 | 7,390 | 3,070 | 769 | 51 | 28 |
| <i>T. gondii</i> | 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 | IO | Template | IO | Model | IO | 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

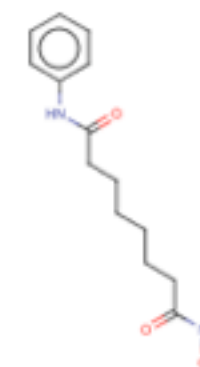
Small Molecule; Approved; Investigational

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.



L. major Histone deacetylase 2 + Vorinostat

Literature

Proc. Natl. Acad. Sci. USA
Vol. 93, pp. 13143–13147, November 1996
Medical Sciences

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

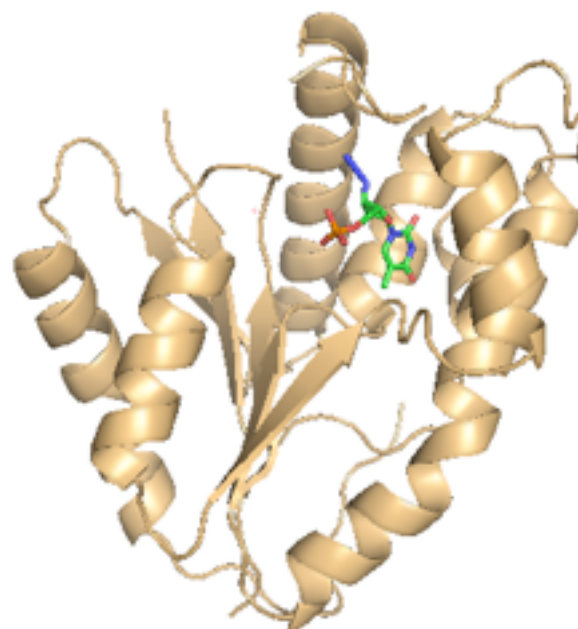
ANTIMICROBIAL AGENTS AND CHEMOTHERAPY, Apr. 2004, p. 1435–1436
0066-4804/04/\$08.00+0 DOI: 10.1128/AAC.48.4.1435–1436.2004
Copyright © 2004, American Society for Microbiology. All Rights Reserved.

Vol. 48, No. 4

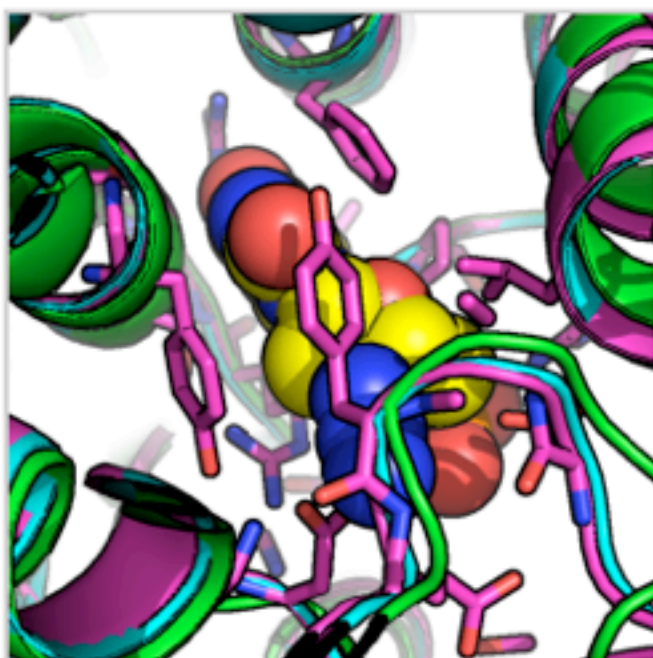
Antimalarial and Antileishmanial Activities of Aroyl-Pyrrolyl-Hydroxyamides, a New Class of Histone Deacetylase Inhibitors

P. falciparum thymidylate kinase + zidovudine

Template 3tmkA a yeast thymidylate kinase.



| PDB | IO | Template | IO | Model | IO | 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

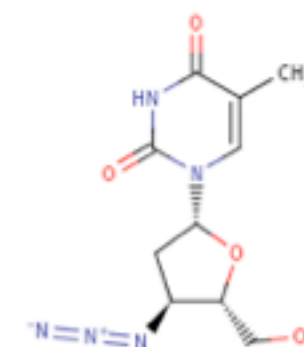
Small Molecule; Approved

Drug categories:

Anti-HIV Agents
Antimetabolites
Nucleoside and Nucleotide Reverse Transcriptase Inhibitors

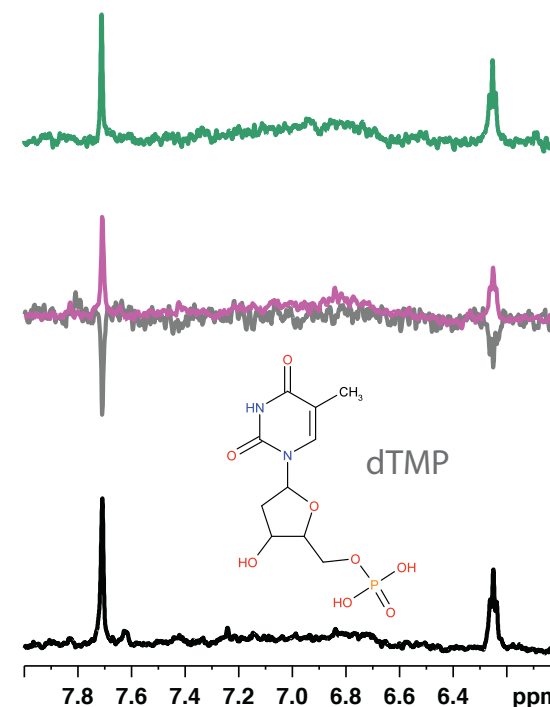
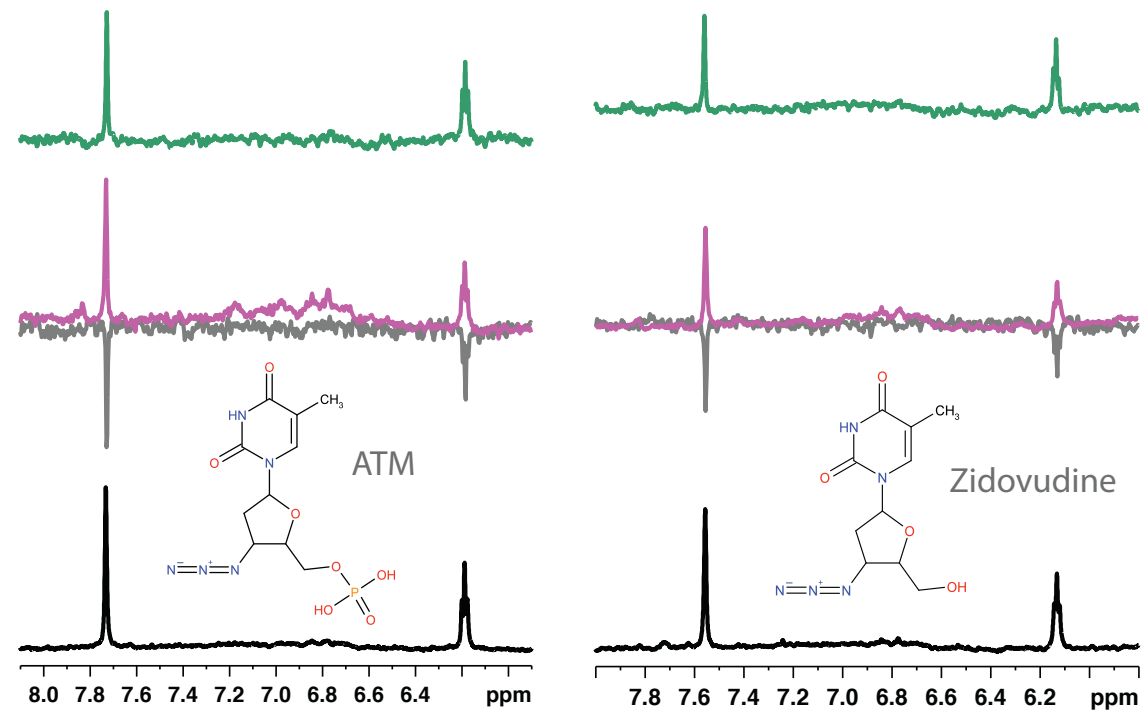
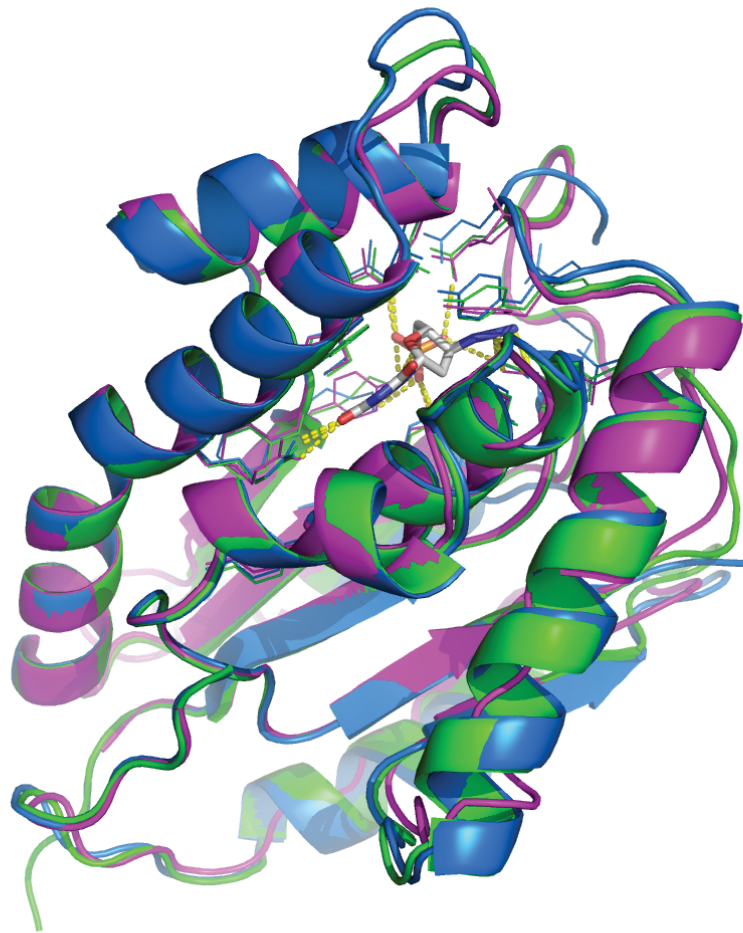
Drug Indication:

For the treatment of human immunovirus (HIV) infections.



P. falciparum thymidilate kinase + zidovudine

NMR Water-LOGSY and STD experiments



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

TDI's kernel

<http://tropicaldisease.org/kernel>

TDI Kernel database » Q9GU59

http://tropicaldisease.org/kernel/q9gu59/ RSS Inquisitor

the **T**ropical **D**isease **I**nitiative *an open source drug discovery project*

You are browsing version 1.0 (2008/05/01) of the TDI Kernel.

Posted on 05.07.08 to Target. Grab the feed. No comments yet. Add your thoughts or trackback from your own site. Edit this entry.

Putative histone deacetylase. predicted to bind 1 ligands [SHH]

UniPort id: **Q9GU59** [*C. parvum*]

Target keywords: ; Anticarcinogenic Agents; Antineoplastic Agents; Transcription; Chromatin regulator; Anti-Inflammatory Agents, Non-Steroidal; Enzyme Inhibitors; Q9GU59; Transcription regulation.; Nucleus

Do you consider this target suitable for drug discovery: ★★★★★ (No Ratings Yet)

Binding site prediction to approved drugs (need help reading this page?):

| PDB | IC ₅₀ | Template | ms | Model | | Ligand | Exact | SupStr | SubStr | Similar |
|-----------------------|------------------|-----------------------|------------|---------------------------------|--------------|---------------------|-------------------------|-------------------------|-------------------------|-------------------------|
| 1c3sA | 83.33/90.00 | 1t64A | 37.09/1.47 | cgd6_1380.1.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.

Shown ligand [SHH](#)

OCTANEDIOICACIDHYDROXYAMIDEIPHENYLAMIDE

expanded from [1c3sA](#) to template [1t64A](#) used for building a 3D model of [cgd6_1380.1.pdb](#). Download the coordinates [data/Q9GU59/Q9GU59.SHH.952.pdb](#)

Highest rated target:
★ A7UD81 (5 out of 5)

2008 : Open Access.
Powered by WordPress.
Theme by Upstart Blogger.

TDI's kernel

<http://tropicaldisease.org/kernel>

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

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

CORRESPONDENCE

A kernel for the Tropical Disease Initiative

To the Editor:

Identifying proteins that are good drug targets and finding drug leads that bind to them is generally a challenging problem. It is particularly difficult for neglected tropical diseases, such as malaria and tuberculosis, where research resources are relatively scarce¹. Fortunately, several developments improve our ability to deal with drug discovery for neglected diseases: first, the sequencing of many complete genomes of organisms that cause tropical diseases; second, the determination of a large number of protein structures; third, the creation of compound libraries, including already-

approved drugs; and fourth, the availability of improved bioinformatics analysis, including methods for comparative protein structure modeling, binding site identification, virtual ligand screening and drug design. Therefore, we are now in a position to increase the odds of identifying high-quality drug targets and drug leads for neglected tropical diseases. Here we encourage a collaboration among scientists to engage in drug discovery for tropical diseases by providing a 'kernel' for the Tropical Disease Initiative (TDI, <http://www.tropicaldisease.org/>)². As the Linux kernel did for open source code development, we suggest that the TDI kernel may help overcome a major stumbling block, in this case, for open source drug discovery: the absence of a critical mass of preexisting work that volunteers can build on incrementally. This kernel complements several other initiatives on neglected tropical diseases^{3–5}, including collaborative web portals (e.g., <http://www.thesynapticleap.org/>), public-

private partnerships (e.g., <http://www.mmv.org/>) and private foundations (e.g., <http://www.gatesfoundation.org/>); for an updated list of initiatives, see the TDI website above.

The TDI kernel was derived with our software pipeline^{6,7} for predicting structures of protein sequences by comparative modeling, localizing small-molecule binding sites on the surfaces of the models and predicting ligands that bind to them. Specifically, the pipeline linked 297 proteins from ten pathogen genomes with already approved drugs that were developed for treating other diseases (Table 1). Such links, if proven experimentally, may significantly increase the efficiency of target identification, target validation, lead discovery, lead optimization and clinical trials. Two of the kernel targets were tested for their binding to a known drug by NMR spectroscopy, validating one of our predictions (Fig. 1 and Supplementary Data online). It is difficult to assess the accuracy of our computational predictions based on this limited experimental testing. Thus, we encourage other investigators to donate their expertise and facilities to test additional predictions. We hope the testing will occur within the

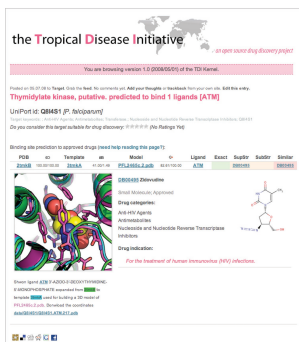


Figure 1 TDI kernel snapshot of the web page for the *Plasmodium falciparum* thymidylate kinase target (<http://tropicaldisease.org/kernel/q8i4s1/>). Our computational pipeline predicted that thymidylate kinase from *P. falciparum* binds ATM (3'-azido-3'-deoxythymidine-5'-monophosphate), a supra-structure of the zidovudine drug approved for the treatment of HIV infection. The binding of this ligand to a site on the kinase was experimentally validated by one-dimensional Water-LOGSY⁹ and saturation transfer difference¹⁰ NMR experiments.

open source context, where results are made available with limited or no restrictions.

A freely downloadable version of the TDI kernel is available in accordance with the Science Commons protocol for implementing open access data (<http://sciencecommons.org/projects/publishing/open-access-data-protocol/>), which prescribes standard academic attribution and facilitates tracking of work but imposes no other restrictions. We do not seek intellectual property rights in the actual discoveries based on the TDI kernel, in the hope of reinvigorating drug discovery for neglected tropical diseases⁸. By minimizing restrictions on the data, including viral terms that would be inherited by all derivative works, we hope to attract as many eyeballs as we possibly can to use and improve the kernel. Although many of the drugs in the kernel are proprietary under diverse types of rights, we believe that the existence of public domain pairs of targets and compounds will reduce the royalties that patent owners can charge and sponsors must pay. This should decrease the large sums of money governments and

| Table 1 TDI kernel genomes | | | | |
|-----------------------------------|--------------------------|------------------------------|----------------------|--------------------|
| Organism ^a | Transcripts ^b | Modeled targets ^c | Similar ^d | Exact ^e |
| <i>Cryptosporidium hominis</i> | 3,886 | 666 | 20 | 13 |
| <i>Cryptosporidium parvum</i> | 3,806 | 742 | 24 | 13 |
| <i>Leishmania major</i> | 8,274 | 1,409 | 43 | 20 |
| <i>Mycobacterium leprae</i> | 1,605 | 893 | 25 | 6 |
| <i>Mycobacterium tuberculosis</i> | 3,991 | 1,608 | 30 | 10 |
| <i>Plasmodium falciparum</i> | 5,363 | 818 | 28 | 13 |
| <i>Plasmodium vivax</i> | 5,342 | 822 | 24 | 13 |
| <i>Toxoplasma gondii</i> | 7,793 | 300 | 13 | 6 |
| <i>Trypanosoma cruzi</i> | 19,607 | 3,070 | 51 | 28 |
| <i>Trypanosoma brucei</i> | 9,210 | 1,386 | 39 | 21 |
| Total | 68,877 | 11,714 | 297 | 143 |

^aOrganisms in bold are included in the World Health Organization (Geneva) Tropical Disease portfolio. ^bNumber of transcripts in each genome. ^cNumber of targets with at least one domain accurately modeled (that is, MODPIPE quality score of at least 1.0). ^dNumber of modeled targets with at least one predicted binding site for a molecule with a Tanimoto score¹¹ of at least 0.39 to a drug in DrugBank¹². ^eNumber of modeled targets with at least one predicted binding site for a molecule in DrugBank.

320

VOLUME 27 NUMBER 4 APRIL 2009 NATURE BIOTECHNOLOGY

320

VOLUME 27 NUMBER 4 APRIL 2009 NATURE BIOTECHNOLOGY

| Organism | Transcripts | Modeled targets | Similar | Exact |
|-----------------------------------|-------------|-----------------|---------|-------|
| <i>Cryptosporidium hominis</i> | 3,886 | 666 | 20 | 13 |
| <i>Cryptosporidium parvum</i> | 3,806 | 742 | 24 | 13 |
| <i>Leishmania major</i> | 8,274 | 1,409 | 43 | 20 |
| <i>Mycobacterium leprae</i> | 1,605 | 893 | 25 | 6 |
| <i>Mycobacterium tuberculosis</i> | 3,991 | 1,608 | 30 | 10 |
| <i>Plasmodium falciparum</i> | 5,363 | 818 | 28 | 13 |
| <i>Plasmodium vivax</i> | 5,342 | 822 | 24 | 13 |
| <i>Toxoplasma gondii</i> | 7,793 | 300 | 13 | 6 |
| <i>Trypanosoma cruzi</i> | 19,607 | 3,070 | 51 | 28 |
| <i>Trypanosoma brucei</i> | 9,210 | 1,386 | 39 | 21 |
| Total | 68,877 | 11,714 | 297 | 143 |

the kernel. The kernel is available in accordance with the Science Commons protocol for implementing open access data (<http://sciencecommons.org/projects/publishing/open-access-data-protocol/>), which prescribes standard academic attribution and facilitates tracking of work but imposes no other restrictions. We do not seek intellectual property rights in the actual discoveries based on the TDI kernel, in the hope of reinvigorating drug discovery for neglected tropical diseases⁸. By minimizing restrictions on the data, including viral terms that would be inherited by all derivative works, we hope to attract as many eyeballs as we possibly can to use and improve the kernel. Although many of the drugs in the kernel are proprietary under diverse types of rights, we believe that the existence of public domain pairs of targets and compounds will reduce the royalties that patent owners can charge and sponsors must pay. This should decrease the large sums of money governments and

OPEN ACCESS Freely available online



A Kernel for Open Source Drug Discovery in Tropical Diseases

Leticia Orti^{1,2}, Rodrigo J. Carbajo², Ursula Pieper³, Narayanan Eswar^{3a}, Stephen M. Maurer⁴, Arti K. Rai⁵, Ginger Taylor⁶, Matthew H. Todd⁷, Antonio Pineda-Lucena², Andrej Sali^{3a}, Marc A. Marti-Renom^{1*}

1 Structural Genomics Unit, Bioinformatics and Genomics Department, Centro de Investigación Principe Felipe, Valencia, Spain, **2** Structural Biology Laboratory, Medicinal Chemistry Department, Centro de Investigación Principe Felipe, Valencia, Spain, **3** Department of Bioengineering and Therapeutic Sciences, Department of Pharmaceutical Chemistry, and California Institute for Quantitative Biosciences, University of California San Francisco, San Francisco, California, United States of America, **4** Gould School of Law, University of Southern California, Los Angeles, California, United States of America, **5** School of Law, Duke University, Durham, North Carolina, United States of America, **6** The Synaptic Leap, San Ramon, California, United States of America, **7** School of Chemistry, University of Sydney, Sydney, New South Wales, Australia

Abstract

Background: Conventional patent-based drug development incentives work badly for the developing world, where commercial markets are usually small to non-existent. For this reason, the past decade has seen extensive experimentation with alternative R&D institutions ranging from private-public partnerships to development prizes. Despite extensive discussion, however, one of the most promising avenues—open source drug discovery—has remained elusive. We argue that the stumbling block has been the absence of a critical mass of preexisting work that volunteers can improve through a series of granular contributions. Historically, open source software collaborations have almost never succeeded without such “kernels”.

Methodology/Principal Findings: Here, we use a computational pipeline for: (i) comparative structure modeling of target proteins, (ii) predicting the localization of ligand binding sites on their surfaces, and (iii) assessing the similarity of the predicted ligands to known drugs. Our kernel currently contains 143 and 297 protein targets from ten pathogen genomes that are predicted to bind a known drug or a molecule similar to a known drug, respectively. The kernel provides a source of potential drug targets and drug candidates around which an online open source community can nucleate. Using NMR spectroscopy, we have experimentally tested our predictions for two of these targets, confirming one and invalidating the other.

Conclusions/Significance: The TDI kernel, which is being offered under the Creative Commons attribution share-alike license for free and unrestricted use, can be accessed on the World Wide Web at <http://www.tropicaldisease.org>. We hope that the kernel will facilitate collaborative efforts towards the discovery of new drugs against parasites that cause tropical diseases.

Citation: Orti L, Carbajo RJ, Pieper U, Eswar N, Maurer SM, et al. (2009) A Kernel for Open Source Drug Discovery in Tropical Diseases. *PLoS Negl Trop Dis* 3(4): e418. doi:10.1371/journal.pntd.0000418

Editor: Timothy G. Geary, McGill University, Canada

Received: December 29, 2008; **Accepted:** March 23, 2009; **Published:** April 21, 2009

Copyright: © 2009 Orti et al. This is an open-access article distributed under the terms of the Creative Commons Attribution License, which permits unrestricted use, distribution, and reproduction in any medium, provided the original author and source are credited.

Funding: MAM-R acknowledges the support from a Spanish Ministerio de Educación y Ciencia grant (BIO2007/66670). AS acknowledges the support from the Sandler Family Supporting Foundation and the National Institutes of Health (R01 GM54762, U54 GM074945, P01 AI035707, and P01 GM71790). AP-L acknowledges the support from a Spanish Ministerio de Ciencia e Innovación grant (SAF2008-01845). RJC acknowledges the support from the Ramon y Cajal Program of the Spanish Ministerio de Educación y Ciencia. We are also grateful for computer hardware gifts to AS from Ron Conway, Mike Homer, Intel, IBM, Hewlett-Packard, and NetApp. The funders had no role in study design, data collection and analysis, decision to publish, or preparation of the manuscript.

Competing Interests: The authors have declared that no competing interests exist.

* E-mail: sali@salilab.org (AS); mmarti@cipf.es (MAM-R)

† Current address: DuPont Knowledge Center, Hyderabad, India

Introduction

There is a lack of high-quality protein drug targets and drug leads for neglected diseases [1,2]. Fortunately, many genomes of organisms that cause tropical diseases have already been sequenced and published. Therefore, we are now in a position to leverage this information by identifying potential protein targets for drug discovery. Atomic-resolution structures can facilitate this task. In the absence of an experimentally determined structure, comparative modeling can provide useful models for sequences that are detectably related to known protein structures [3,4]. Approximately half of known protein sequences contain domains that can be currently predicted by comparative modeling [5,6]. This coverage

will increase as the number of experimentally determined structures grows and modeling software improves. A protein model can facilitate at least four important tasks in the early stages of drug discovery [7]: prioritizing protein targets for drug discovery [8], identifying binding sites for small molecules [9,10], suggesting drug leads [11,12], and optimizing these leads [13–15].

Here, we address the first three tasks by assembling our computer programs into a software pipeline that automatically and on large-scale predicts protein structures, their ligand binding sites, and known drugs that interact with them. As a proof of principle, we applied the pipeline to the genomes of ten organisms that cause tropical diseases (“target genomes”). We also experimentally tested two predicted drug-target interactions using Nuclear Magnetic

www.plosntds.org

1

April 2009 | Volume 3 | Issue 4 | e418

www.plosntds.org

APRIL 2009 | VOLUME 3 | ISSUE 4 | e418

the kernel. The kernel is available in accordance with the Science Commons protocol for implementing open access data (<http://sciencecommons.org/projects/publishing/open-access-data-protocol/>), which prescribes standard academic attribution and facilitates tracking of work but imposes no other restrictions. We do not seek intellectual property rights in the actual discoveries based on the TDI kernel, in the hope of reinvigorating drug discovery for neglected tropical diseases⁸. By minimizing restrictions on the data, including viral terms that would be inherited by all derivative works, we hope to attract as many eyeballs as we possibly can to use and improve the kernel. Although many of the drugs in the kernel are proprietary under diverse types of rights, we believe that the existence of public domain pairs of targets and compounds will reduce the royalties that patent owners can charge and sponsors must pay. This should decrease the large sums of money governments and

the kernel. The kernel is available in accordance with the Science Commons protocol for implementing open access data (<http://sciencecommons.org/projects/publishing/open-access-data-protocol/>), which prescribes standard academic attribution and facilitates tracking of work but imposes no other restrictions. We do not seek intellectual property rights in the actual discoveries based on the TDI kernel, in the hope of reinvigorating drug discovery for neglected tropical diseases⁸. By minimizing restrictions on the data, including viral terms that would be inherited by all derivative works, we hope to attract as many eyeballs as we possibly can to use and improve the kernel. Although many of the drugs in the kernel are proprietary under diverse types of rights, we believe that the existence of public domain pairs of targets and compounds will reduce the royalties that patent owners can charge and sponsors must pay. This should decrease the large sums of money governments and

Acknowledgments

<http://marciuslab.org>

<http://cnag.cat> · <http://crg.cat>

<http://integrativemodeling.org>

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)

FUNDING

CNAG

MINECO

Era-Net Pathogenomics

HFSP

MODEL ASSESSMENT

Francisco Melo (CU)

Alejandro Panjkovich (CU)

NMR

Antonio Pineda-Lucena

Leticia Ortí

Rodrigo J. Carbajo

MAMMOTH

Angel R. Ortiz

3D Genomes

George Church (Harvard)

Job Dekker (UMASS)

Jeane Lawrence (UMASS)

Lucy Shapiro (Stanford)

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)

GeMoA

LLuís Ballell (GSK)

Brigitte Gicquel (IP)

Olivier Neyrolles (IPBS)

Marc A. Marti-Renom (CNAG)

Matthias Wilmanns (EMBL)

CNAG available job positions

Bioinformatics

- Postdoctoral Researcher for **Blueprint** Project
- Postdoctoral Data Analyst **RD Connect** Project
- PhD Student on **high-performance algorithms** for processing genomic data
- Senior Postdoctoral Researcher in **Functional Bioinformatics**

Laboratory

NO OPENINGS

More information at <http://www.cnag.cat/jobs>