# A "kernel" for the Tropical Disease Initiative
## *An open source approach to drug discovery*

**Marc A. Marti-Renom**

http://sgu.bioinfo.cipf.es

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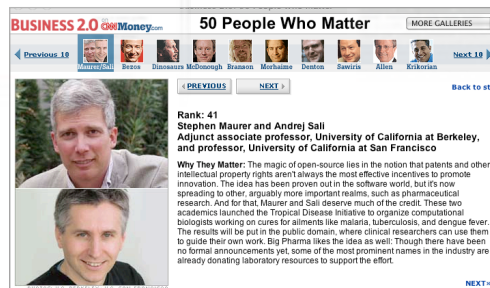Prince Felipe Resarch Center (CIPF), Valencia, Spain

# 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**

.TSL web site `http://TheSynapticLeap.org`
.Maurer and Sali 41th in "50 Who Matter"

**2008**

.TDI kernel `http://TropicalDisease.org/kernel`

**2009**

.TDI kernel papers published

# Initial feed-back...

**14 Mar 2005**

I think TDI is a unique and very interesting project. I would like so much to make something for it...

**So, where are we going? What's happening? What can we do?**

I still trust in open source drug discovery. :-))

Luca Brivio

I am interested in beginning rese[arch]
disease for underserved populati[on]
however, confused.
**If someone will tell me where te
begin on, I'd be greatful.**

Thank you kindly,
Adam Huber

**...any, the bottlenecks are?**
...eas and potential avenues to explore,
**...ction Plan!**

Regards,
Jacob Lester

**9 Mar 2005**
I'm a programmer, not a
the list active :)

**GNU started with RMS.
Linux started with Linu**
**You need someone gre**
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

**...stic that the rest**

Stephen Mark Maurer

3

# Initial feed-back...

**14 Mar 2005**

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

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

I still trust in open source drug discovery. :-))

Luca Brivio

**16 Feb 2005**

Hi,

**It would be interesting to know what, if any, the bottlenecks are?**
The Wiki site contains many interesting ideas and potential avenues to explore, but from what I can see it is **lacking an Action Plan!**

Regards,
Jacob Lester

I am interested in beginning rese... disease for underserved populati... however, confused.
**If someone will tell me where t... begin on, I'd be greatful.**

Thank you kindly,
Adam Huber

**9 Mar 2005**
I'm a programmer, not a ... the list active :)

**GNU started with RMS.**
**Linux started with Linu...**
**You need someone gre...**
**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

stic that the rest

Stephen Mark Maurer

4

# Initial feed-back...

**14 Mar 2005**

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

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

I still trust in open s

Luca Brivio

**16 Feb 2005**

Hi,

**bottlenecks are?**
potential avenues to explore,
**n!**

**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,
however, confused.
**If someone will tell me where to sign up and give me some research topics to
begin on,** **I'd be greatful.**

Thank you kindly,
Adam Huber

**9 Mar 2005**
I'm a programmer, not a
the list active :)

**GNU started with RMS.**
**Linux started with Linu**
**You need someone gre**
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

**stic that the rest**

Stephen Mark Maurer

5

# Initial feed-back...

**14 Mar 2005**

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

**So, where are we going? What's happening? Wha~~~~**

I still trust in open s~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~ **bottlenecks are?**
~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~potential avenues to explore,

Luca Brivio ~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~ **n!**

**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,

**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

**stic that the rest**

Stephen Mark Maurer

# Initial feed-back...

**14 Mar 2005**

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

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

I still trust in open s___

Luca Brivio

**16 Feb 2005**

Hi,

**bottlenecks are?**

potential avenues to explore,

**n!**

**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,

**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 nee___
papers or the website.

And you might consider merging into the bios___
critical mass for things to take off. Consider th___
together and make linux work.

Daniel Amelang

**19 Jan 2005**

**If we do the science well, I'm optimistic that the rest
of TDI will fall into place.**

Stephen Mark Maurer

7

# Open Source without a Kernel?



Linux distro timeline — Version 7.2 by NPU (nonplusx@gmail.com)

8

# Drug Discovery pipeline

# Drug Discovery pipeline



Cumulative cost
Success rate

shorter time...

Pre Lead | Year 1 | Year 2 | Year 3 | Year 4 | Year 5 | Year 6 | Year 7 | Year 8 | Year 9

**Target & Lead identification** → **Lead optimization** → **Preclinical** → **Phase I** → **Phase II** → **Phase III**

TDI

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

10

# 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**

I I

# Non-Profit organizations

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



**21 projects** in DNDi's portfolio, 2008

*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

Heart diseases

Rare diseases

DALY

Disease

13

# "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.

14

# Predicting binding sites in protein structure models.

# DBAli~v2.0~ database

## http://www.dbali.org



- ✓ **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**

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

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

Uses MAMMOTH for similarity detection

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

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

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

**http://www.dbali.org**



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

17

# Method

**DBAli tools**

Chain ID

AnnoLyze search

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

HTML output

**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

| | 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 |

18

# 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

20

# Sensitivity .vs. Precision

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

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

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

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

21

# Comparative docking



**Expansion**

co-crystalized protein/ligand

crystalized
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 |
|---|---|---|---|---|---|---|
| *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** |

24

# *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.*

25

# *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

26

# *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



*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).

# Acknowledgments

http://sgu.bioinfo.cipf.es
http://tropicaldisease.org - http://thesynapticleap.org

**COMPARATIVE MODELING**
**Andrej Sali**
**Narayanan Eswar**
**Ursula Pieper**
M. S. Madhusudhan
Min-Yi Shen
Ben Webb

**Tropical Disease Initiative**
**Stephen Maurer**
**Arti Rai**
**Andrej Sali**
**Ginger Taylor**
**Matthew Todd**
---
Bissan Al-Lazikani
James McKerrow
Brian Shoichet
David S. Roos

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

**FUNCTIONAL ANNOTATION**
Stefania Bosi
Anna Tramontano

كاوست
KAUST