

Combining Cheminformatics Methods and Pathway Analysis to Identify Molecules with Whole-Cell Activity Against *Mycobacterium Tuberculosis*

Malabika Sarker · Carolyn Talcott · Peter Madrid · Sidharth Chopra · Barry A. Bunin · Gyanu Lamichhane · Joel S. Freundlich · Sean Ekins

Received: 31 August 2011 / Accepted: 16 March 2012 / Published online: 4 April 2012
© Springer Science+Business Media, LLC 2012

ABSTRACT

Purpose New strategies for developing inhibitors of *Mycobacterium tuberculosis* (Mtb) are required in order to identify the next generation of tuberculosis (TB) drugs. Our approach leverages the integration of intensive data mining and curation and computational approaches, including cheminformatics combined with bioinformatics, to suggest biological targets and their small molecule modulators.

Methods We now describe an approach that uses the TBCyc pathway and genome database, the Collaborative Drug Discovery database of molecules with activity against Mtb and their associated targets, a 3D pharmacophore approach and Bayesian models of TB activity in order to select pathways and metabolites and ultimately prioritize molecules that may be acting as substrate mimics and exhibit activity against TB.

Results In this study we combined the TB cheminformatics and pathways databases that enabled us to computationally search >80,000 vendor available molecules and ultimately test 23 compounds *in vitro* that resulted in two compounds (N-(2-furylmethyl)-N'-[(5-nitro-3-thienyl)carbonyl]thiourea and N'-[(5-nitro-3-thienyl)carbonyl]-N-(2-thienylmethyl)thiourea) proposed as mimics of D-fructose 1,6 bisphosphate, (MIC of 20 and 40 µg/ml, respectively).

Conclusion This is a simple yet novel approach that has the potential to identify inhibitors of bacterial growth as illustrated by compounds identified in this study that have activity against Mtb.

KEY WORDS Bioinformatics · cheminformatics · Collaborative drug discovery · *Mycobacterium tuberculosis* · pharmacophore

INTRODUCTION

Mycobacterium tuberculosis (Mtb), the causative agent of tuberculosis (TB), is estimated to maintain latent infection in approximately one-third of the world's population and kill 1.7–1.8 million people each year (1). The survival of Mtb relies on an array of cellular functions carried out by metabolites, enzymes, structural and regulatory proteins and RNAs. These essential functions can be targeted to kill or suppress the proliferation of Mtb. Soon after the genome sequence of the Mtb H37Rv strain was published (2), various laboratories focused on identifying genes essential for growth under *in vitro* and *in vivo* conditions (3). Classification of essential genes as targets is based on forward genetic approaches that consider a protein as a potential target if an essential gene encodes it (4). A target should be essential for growth and viability of the pathogen at least under the condition of host infection. During infection, Mtb appears to reside predominantly within the host lung

Electronic supplementary material The online version of this article (doi:10.1007/s11095-012-0741-5) contains supplementary material, which is available to authorized users.

M. Sarker · C. Talcott · P. Madrid · S. Chopra
SRI International
333 Ravenswood Avenue
Menlo Park, California 94025, USA

G. Lamichhane
Department of Medicine, Johns Hopkins School of Medicine
1550 Orleans St, Room 103
Baltimore, Maryland 21287, USA

B. A. Bunin · S. Ekins
Collaborative Drug Discovery
1633 Bayshore Highway, Suite 342
Burlingame, California 94010, USA

alveolar macrophages. Here the pathogen encounters and adapts to conditions that are considered to be unfavorable for growth such as a decrease in pH, depleted nutrition, hypoxia and reactive oxygen and nitrogen radicals (5). The genes of Mtb essential to perform such functions are not necessarily required under the *in vitro* growth conditions as the functions encoded by these genes are selectively required to survive and thrive in host imposed unfavorable conditions (6). Therefore, identifying the *in vivo* essential genes as potential targets is relevant for therapeutic intervention.

Another approach to select a target whose inhibition is of therapeutic value is to select metabolic pathways that are necessary for growth and proliferation of Mtb *in vivo* (7). This allows for a careful consideration of biological rationale and the metabolic role of the specific target within the context of a specific metabolic pathway. Functionality or reaction information about the target should be identified so that assays (both low- and high-throughput) can be built appropriately to mimic these *in vivo* conditions. The analysis of biosynthetic pathways helps determine alternative routes of synthesis of the essential proteins (7), highlighting areas of metabolism where degeneracy may make it difficult to deplete a given metabolite.

Discarding target enzymes from the pathogen which share a similarity with the host protein/s significantly lessens the probability of undesired host protein–drug interactions. This criterion, however, is not absolute. For example, successful antibiotics such as trimethoprim and quinolones display selectivity towards bacterial targets despite the existence of their human orthologs. Trimethoprim specifically inhibits bacterial dihydrofolate reductase despite 28% sequence identity with its human ortholog, and quinolones specifically inhibit bacterial gyrase A, despite 20% sequence similarity with human topoisomerase II (8). For selective targeting, substantial differences in the regions of the active site (presumably responsible for the difference in substrate specificity) have more significance than the overall 3D structures, which again are more critical than whole sequence similarity between orthologs (7). However, we offer that this

is a reasonable initial target filter criterion in order to limit the number of essential Mtb essential proteins that can be evaluated.

The rationale for exploring novel targets for TB is that the pipeline for therapeutics has not produced a new approved first line drug in over 40 years. Only a small fraction of TB proteins are known to be modulated by approved drugs and recent testing has targeted additional proteins; this has yet to result in a new drug (9,10). This also represents a pattern observed for other antibacterial targets, reflecting the difficulty of target-based high-throughput screening (11). In pharmaceutical companies, computational approaches are widely used to aid in drug discovery; these do not appear to have been as extensively applied for TB. For example, virtual screening of compound libraries is used as a complement to high-throughput screening *in vitro* for many diseases (12). A recent review pointed to some of the gaps in using such cheminformatics approaches in TB drug discovery (13). Alternative approaches include rational inhibitor design based on the substrate or product structure or on the reaction mechanism. The approach leverages the “chemical similarity principle” (14), which states that similar molecules likely have similar biological properties. Applied to small molecule metabolism, this principle has motivated the search for enzyme inhibitors chemically similar to their endogenous substrates. The approach has yielded many successes, including anti-metabolites such as trimethoprim, D-cycloserine, vancomycin, etc. Recently we have taken the mimic strategy utilizing 2D similarity and 3D pharmacophore searches of molecule databases using essential molecules as starting points (15) and have identified compounds with *in vitro* activity against TB. In this study, we have extended this work and taken an exhaustive approach to identifying essential targets that have to our knowledge not been interrogated for TB to identify small molecule inhibitors. We have then mined the known compounds with whole-cell activity and TB targets databases and used multiple cheminformatics tools to prioritize commercially available molecules for testing *in vitro*.

MATERIALS AND METHODS

Reagents and Molecules

All experimental compounds were purchased from Sigma-Aldrich, Maybridge or Asinex. Purities were required to be greater than 90% with a majority of compounds having a purity of greater than 95%. Compounds were all dissolved in dimethyl sulfoxide (Sigma Aldrich) at a stock concentration of 12.8 mg/ml immediately and then diluted for biological testing.

J. S. Freundlich
Departments of Pharmacology & Physiology and Medicine
Center for Emerging and Reemerging
Pathogens, UMDNJ—New Jersey Medical School
185 South Orange Avenue
Newark, New Jersey 07103, USA

S. Ekins (✉)
Collaborations in Chemistry
5616 Hilltop Needmore Road
Fuquay-Varina, North Carolina 27526, USA
e-mail: ekinssean@yahoo.com

Identification of Essential *In Vivo* Enzymes of *Mycobacterium Tuberculosis*

While there have been studies that evaluate the role of particular *M. tuberculosis* genes and define their potential as targets for new drugs (16) there have been none to our knowledge that take the following approach. Following intensive literature mining and manual curation, we extracted all the genes that are essential for Mtb growth *in vivo*. This involved

- i) the work of Sassetti and coworkers, who used a recombinant mycobacteriophage carrying a highly infectious transposon to develop a high-throughput technique called Transposon Site Hybridization (*TraSH*) and identified the Mtb genes required for growth both *in vitro* and *in vivo* in mice (17,18).
- ii) all published data by the Tuberculosis Animal Research and Gene Evaluation Taskforce (TARGET) in relation to the large collection of defined Mtb mutants (*Designer Arrays for Defined Mutant Analysis (DeADMan)*) that were used to identify the genes essential for growth in the lungs of mice (19), guinea pigs (20) and non-human primates (6).

Collection of Metabolic Pathway and Reaction Information for the Essential Enzymes

Various TB-related databases (13) are available that cover diverse areas of TB research like genomes, pathway maps, phylogenetic trees, active compounds, large-scale screening data, resistance-associated mutations, targets, comparative analysis and gene expression data. In order to determine the biological role of the essential proteins of Mtb, we used TBCyc (<http://tbcyc.tdb.org/index.shtml>), an Mtb specific metabolic pathway database for our analysis. The TBCyc database was initially developed using SRI's Pathway Tools software that automatically generates a Pathway/Genome Database (PGDB) describing the genome and biochemical networks of the organism from the annotated genome sequence of Mtb (21,22). Automatic generation was followed by substantial additional curation. TBCyc provides a pathway-based visualization of the entire cellular biochemical network, called the cellular overview diagram, which supports interrogation and exploration of whole organism system-biology analyses. The cellular overview includes metabolic, transport, and signaling pathways, and other membrane and periplasmic proteins (see Fig. 1). The TBCyc metabolic pathways for the Mtb *in vivo* essential genes were extensively studied for analyzing the reactions, metabolites and other enzymes involved in the same pathway.

Comparison of Non-Human-Homologous Enzymes with Mtb *In Vivo* Essential Gene Set

Anishetty *et al.* (23) reported a thorough study on pairwise sequence comparison (BLASTp) between human and Mtb proteins. In this report, enzymes from the biochemical pathways of Mtb from the KEGG metabolic pathway database were compared with proteins from human with an e-value threshold cutoff of 0.005. Bacterial enzymes, which did not show similarity to any of the human proteins, below this threshold, were filtered out as potential drug targets. In total, they reported 185 proteins that were absent in humans. Sassetti *et al.* have also listed 49 essential Mtb proteins as unique to *Mycobacteria spp.* (18). In the current study we excluded putative essential Mtb proteins that are present in humans by comparing the list of the published non-human Mtb orthologs with the essential *in vivo* Mtb proteins that we extracted and curated from various studies.

Selection of Mtb Targets That Are Essential *In Vivo* But Not Homologous to Human Proteins and Not Known as TB Drug-Targets

Metabolic enzymes of Mtb that fulfill the criteria of being both essential *in vivo* and absent from humans were further analyzed to find out if they are already experimentally validated or *in silico* predicted targets of the known and FDA-approved TB drugs. This was achieved by searching the literature that had experimentally validated Mtb enzymes as a target for known TB drugs as well as reports predicting the *in silico* targets for the known TB drugs (24). The CDD TB database was also searched to find novel *in vivo* essential targets without screening hits.

In Silico Design of Small Molecule Inhibitors or Pharmacophores for Selected Enzyme Targets

The selection of the above-mentioned enzyme targets led to using their corresponding substrates as the starting point for pharmacophore models. Starting with each such structure, a 3D pharmacophore was developed using Accelrys Discovery Studio 2.5.5 (Accelrys, San Diego, CA) from 3D conformations of the metabolite. This identified key features, onto which was mapped a van der Waals surface for the molecule (15,25,26). The pharmacophore plus shape was then used to search 3D compound databases from well-known and widely used vendors including Maybridge ($N=57,181$ molecules), Asinex ($N=24,998$) and Sigma Aldrich (LOPAC $N=1200$) (for which up to 100 molecule conformations with the FAST conformer generation method with the maximum energy threshold of 20 kcal/mol, were created). The *in silico* hits were collated and uploaded in CDD, and Bayesian models for TB whole

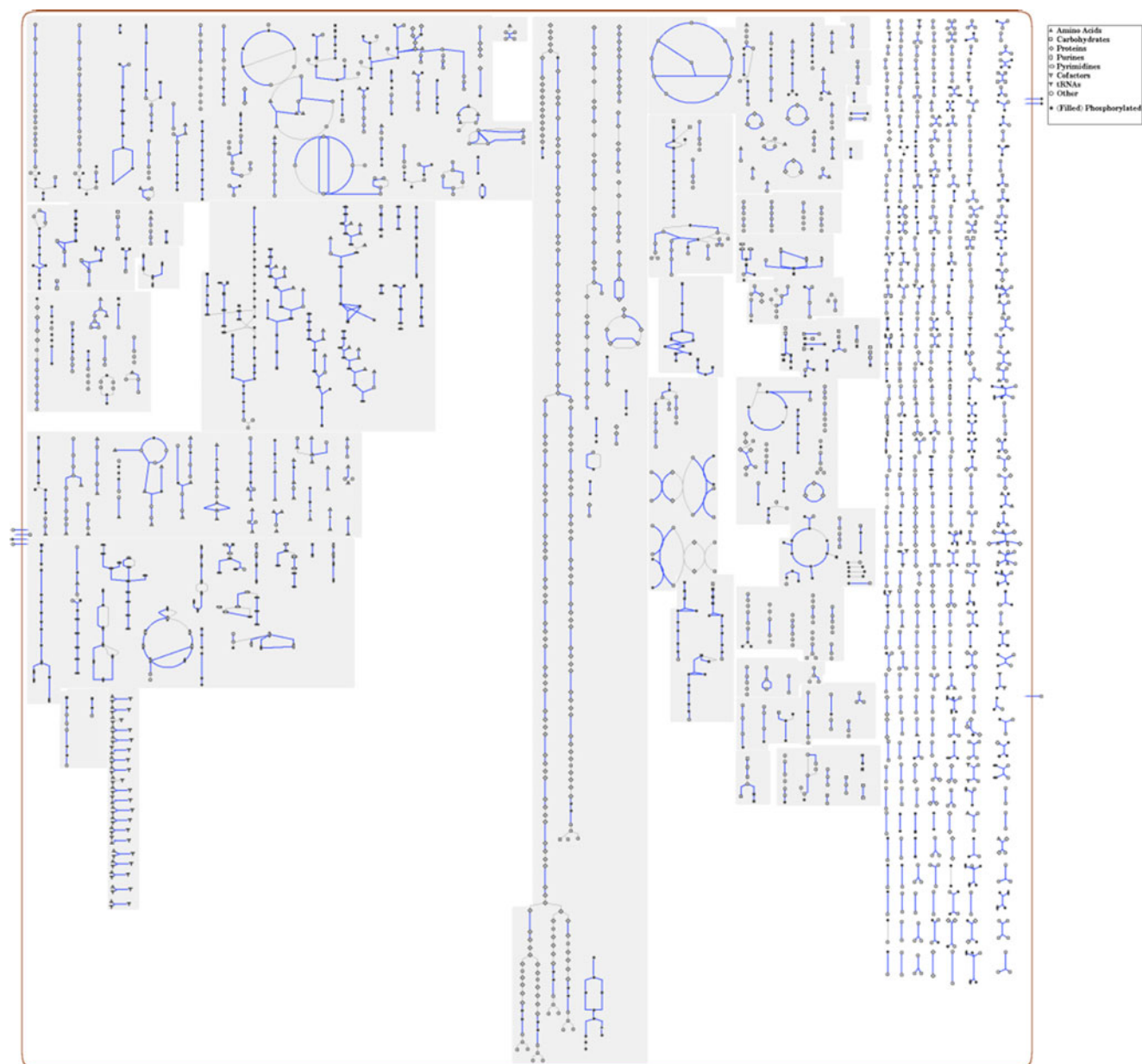


Fig. 1 The cellular overview diagram for *M. tuberculosis* H37Rv, from the TBCyc database (<http://tbcyc.tdb.org/index.shtml>).

cell activity (see discussion later) and SMARTS filters for reactivity (25,27,28) were run against the compounds and the data re-imported in CDD. Finally the compounds were filtered in CDD based on the Bayesian score and manual selection to retrieve compounds with ideal molecular properties for *in vitro* TB activity (25,27,28).

Measurement of Antibacterial Activity Against *Mtb*

We used the resazurin (Alamar Blue) assay as the primary screen for activity against replicating *Mtb* (29). Each compound was tested over a range of concentrations to determine the MIC. The antimicrobial susceptibility test was performed

in a clear-bottomed, round well, 96-well microplate. Initial compounds were tested at 8 concentrations ranging between 40 and 0.31 $\mu\text{g/ml}$. After a growth medium containing $\sim 10^4$ bacteria was added to each well, the different dilutions of compounds were added. Controls included wells containing (1) the different concentrations of compounds only, to exclude autofluorescence in the presence of resazurin, (2) bacteria and growth medium, and (3) sterility control of the medium. Plates were incubated at 37°C for 5 days in an ambient incubator at which time 5 μl of 1% resazurin dye was added to each well. After 2 days of incubation, fluorescence was measured in a microplate fluorimeter with excitation at 530 nm and emission at 590 nm. The lowest drug concentration that inhibited

growth of $\geq 90\%$ of Mtb bacilli in the broth was considered the MIC value (30). Rifampicin and isoniazid were used as positive controls.

RESULTS

Identification of *In Vivo* Essential Enzymes of *Mycobacterium Tuberculosis*

We have collated for the first time all the genes that have so far been reported to be essential for Mtb growth *in vivo*. This gives us a non-redundant list of 314 genes. 194 genes are from mouse TraSH analysis, 31 genes are from a DeADMAN analysis that used mouse as the host, 18 genes are from an independent DeADMAN analysis that used guinea pig model and 108 genes are from a DeADMAN analysis that used non-human primate model of Mtb infection. There are overlaps between some of the studies. A Venn diagram (Fig. 2) shows the degree of intersection among the *in vivo*

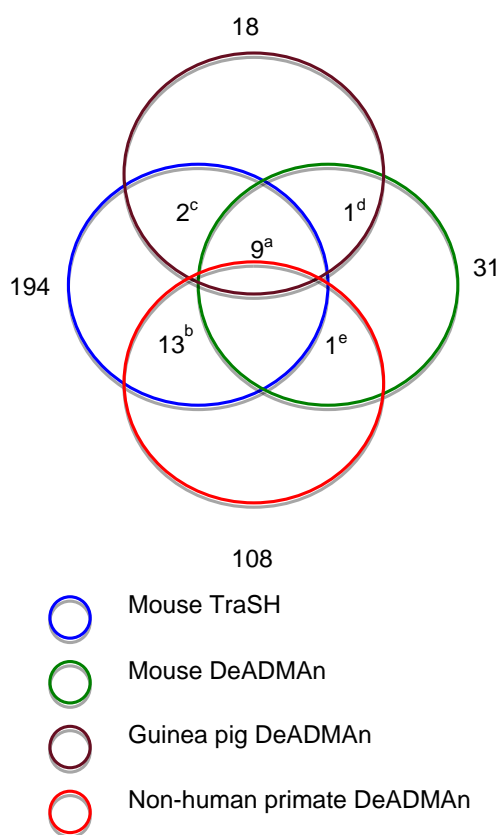


Fig. 2 A Venn diagram below shows the degree of association between the *in vivo* mutants of Mtb in different animal models. Genes are (a) *npr* (Rv0101), Rv0204c, *mkl* (Rv0655), *mmpL10* (Rv1183), *sugC* (Rv1238), *bioB* (Rv1589), Rv2224c, *mmpL7* (Rv2942), Rv3210c, (b) *mce1A* (Rv0169), *lprK* (Rv0173), Rv0687, *fadD21* (Rv1185c), Rv1371, *cobL* (Rv2072c), *dhrrA* (Rv2936), *lprN* (Rv3495c), Rv3683, Rv3871, *embC* (Rv3793), Rv2387, *fabG* (Rv3502c), (c) *fadE28* (Rv3544c), Rv3864, (d) Rv1798, (e) Rv0336.

mutants of Mtb in different models. It should be noted that functions encoded by many of the 314 genes are not yet known.

Collection of Metabolic Pathway and Reaction Information for the Essential Enzymes

We identified all the pathways that have one or more essential enzymes. TBCyc gives a total of 53 non-redundant pathways for the set of 314 *in vivo* essential genes. From this list of essential genes, *pcaA* (Rv0470c), *mmaA3* (Rv0643c), Rv1144, *fadA4* (Rv1323), *bioA* (Rv1568), *bioF1* (Rv1569), *bioB* (Rv1589), *argJ* (Rv1653), *pks12* (Rv2048c), *plsC* (Rv2483c), Rv2857c, *ddlA* (Rv2981c), *amiD* (Rv3375), *fabG* (Rv3502c), *fadA6* (Rv3556c), and *hycD* (Rv0084) belong to more than one TBCyc pathway. From the reactions catalyzed by the corresponding essential enzymes, substrates were identified. Their 2D structures, obtained from ChemSpider (www.chemspider.com, a free chemical structure database), were later used in our analysis for pharmacophore development.

Comparison of Enzymes with no Human Homologs with Mtb *In Vivo* Essential Gene Set

66 proteins were found to be both *in vivo* essential while having no human homologs. A list of 314 essential *in vivo* genes of Mtb along with 53 TBCyc pathways and 66 proteins with no human orthologs is provided as Supplementary Material 1 (“Essential-genes-*in vivo*-Mtb”) (Fig. 3a). These data are freely available in CDD (www.collaborativedrug.com). Each essential gene name is linked to the TB database, TBDB (<http://www.tbdb.org/>, Fig. 3b). All the pathways are linked to the TBCyc database for analysis and visualization of the pathways, reactions and metabolites. The PubMed abstracts can be accessed (via the PubMed identifiers) for essentiality and ortholog information. Where the 3D structures are available, the PDB (X-ray or NMR method) identifiers are given along with respective URLs for further details.

Selection of Targets That Are *In Vivo* Essential, Not Homologous to Human and Not Known as TB Drug-Targets

We produced a summary of published drugs for TB with known or predicted targets (Supplementary Material 2 TB drugs and literature compounds with targets, Fig. 3c) that has 14 known targets and 31 predicted targets for the already known 35 TB drugs. This dataset is also available in CDD along with a larger dataset of 666 literature compounds with antitubercular activity and their known targets, for which all the literature evidence is cited (Fig. 3d).

Only the new and unexplored enzymes were selected for further investigation. Supplementary Material 3 includes “Metabolites and their essential enzymes” (Fig. 3a). This table contains 12 such *in vivo* essential enzymes that are absent in human, have known reactions in TBCyc and are not targets of known TB drugs. The associated reactions, corresponding substrates and products (along with SMILES (31)) are annotated. This table was used for the cheminformatics analysis.

During this process, we identified several known drug targets including genes *embA* and *embC* (both encode enzymes that are essential *in vivo* and non-human orthologs) that are targeted by ethambutol (Supplementary Material 2 TB drugs and literature compounds with targets). Our findings (not used for the present analysis) also included several enzymes that are essential *in vitro* that had no human homologs and were already predicted targets for known drugs. These included MurD (mefloquine-predicted), KasA (cerulenin), RpoB (rifampin,

a.

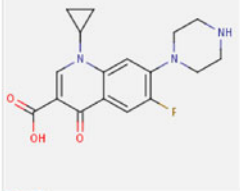
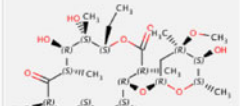
Molecule	TB target information			
	Gene Link	Essentiality	Essentiality Reference 1	Essentiality Reference 2
summary details <div>NO STRUCTURE</div> accD1 (Rv2502c) SRI TB Target Database	http://tbcyc.tdb.org/MT37RVV/NEW-IMAGE?type=GENE-IN-PWY&object=Rv2502C&detail-level=2	In Vivo Essential (0.255)	http://www.ncbi.nlm.nih.gov/pubmed/?term=14569030	
summary details <div>NO STRUCTURE</div> aceE (Rv2241) SRI TB Target Database	http://tbcyc.tdb.org/MT37RVV/NEW-IMAGE?type=GENE-IN-PWY&object=Rv2241&detail-level=2	In Vivo Essential (0.273)	http://www.ncbi.nlm.nih.gov/pubmed/?term=14569030	
summary details <div>NO STRUCTURE</div> amdD (Rv3375) SRI TB Target Database	http://tbcyc.tdb.org/MT37RVV/NEW-IMAGE?type=GENE-IN-PWY&object=Rv3375&detail-level=2	In Vivo Essential (0.39)	http://www.ncbi.nlm.nih.gov/pubmed/?term=14569030	
summary details <div>NO STRUCTURE</div> amt (Rv2920c)	http://tbcyc.tdb.org/MT37RVV/NEW-IMAGE?type=GENE-IN-PWY&object=Rv2920C&detail-level=2	In Vivo Essential	http://www.ncbi.nlm.nih.gov/pubmed/?term=20394526	

b.

The screenshot displays the TB Database interface. At the top, it shows 36 structures (41 matches) and various action buttons like 'Show', 'Change display options', 'Plot results', 'Export results', and 'Add results to project'. The main table lists TB molecules and target information. The columns are: Molecule, Molecule name, Target gene, Target gene link, Drug for TB, Target gene pathway, and Target gene link. The first row shows CDD-151, a molecule with a chemical structure, linked to the target gene *glnK* (Rv0086). Arrows indicate links from the molecule name to the TB Database, from the target gene to PubMed, and from the target gene link to KEGG. Below the table, there are sections for 'Molecular Detection of Mutations Ass Resistance Compared with Conventional Mycobacterium tuberculosis' and 'DNA REPLICATION' with a diagram.

Fig. 3 Images of databases created in this project, which are freely available at www.collaborativedrug.com to illustrate the connection between molecular structure, gene link, pathway links and literature links. **(a)** *In vivo* essential targets database showing genelink, essentiality and essentiality references. **(b)** TB molecules and target information database connects molecule, gene, pathway and literature, showing links to TB database, PubMed and KEGG. **(c)** Drugs and targets database showing molecule structure, molecule name, target gene and gene link. **(d)** Literature compounds and targets database showing molecule structure, molecule name, target gene and gene link.

c.

TB molecules and target information				
Molecule	Molecule name	Target gene 1	Target gene 1 link	Drug for ta...to PubMed 1
CDD-151 summary details  CDD-151 SRI Group Vault	Ciprofloxacin	gyrA (Rv0006)	http://genome.tdb.org/annotation/genome/tdb/GeneDetails.html?sp=S7000000635248048	http://www.ncbi.nlm.nih.gov/pubterm=21300839
CDD-154 summary details  CDD-154 SRI Group Vault	Clarithromycin	dnaA (Rv0001) - Predicted	http://genome.tdb.org/annotation/genome/tdb/GeneDetails.html?sp=S7000000635248067	http://www.ncbi.nlm.nih.gov/pubterm=19301903

d.

 666 structures (690 matches) · Show structures · [Change display options](#) · [Plot results](#) · [Export results](#) · [Add results to project](#)

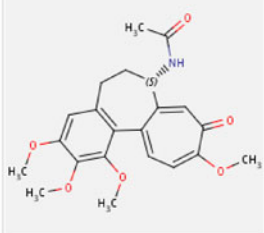

leads and targets				
Molecule	Target gene 1	Target gene 1 link	molecule fo...to PubMed 1	molecu
CDD-167 summary details  CDD-167 SRI Group Vault	tnk (Rv3247c)	http://genome.tdb.org/annotation/genome/tdb/GeneDetails.html?sp=S7000000635252378	http://www.ncbi.nlm.nih.gov/pubmed?term=15566289	
CDD-273 summary details  CDD-273 SRI Group Vault	tnk (Rv3247c)	http://genome.tdb.org/annotation/genome/tdb/GeneDetails.html?sp=S7000000635252378	http://www.ncbi.nlm.nih.gov/pubmed?term=15801836	

Fig. 3 (continued)

rifapentine, rifabutin), Alr (D-cycloserine–predicted), FolP1 (*p*-aminosalicylic acid–predicted) (Supplementary Material 2 TB drugs and literature compounds with targets).

In this study several enzymes, substrate metabolites, reactions and their pathways were selected based on the analysis described previously (Table I). The substrate metabolites of the essential enzymes were chosen as final targets for use with cheminformatics approaches. The cheminformatics methods included the construction of pharmacophores for individual substrates which provided a 3D shape and feature query for searching databases of compounds that could be purchased for testing.

In Silico Design of Metabolite Pharmacophores for Essential Enzyme Targets and Selection of Putative Substrates Mimics

842 molecules retrieved using the various pharmacophores based on substrate structures are suggested as potential mimics (Fig. 4). These molecules were run through the SMARTS filters (for chemical reactivity) and Bayesian models for whole-cell TB activity in Discovery Studio (28,32,33) and 234 were flagged as failing the SMARTS filters as they had features suggested as undesirable based on the default settings. All compounds were imported into CDD.

Table 1 Targets, Metabolites and Pathways Pursued in this Study

Essential gene	Pathway	Essential substrate/s
<i>bioB</i> (Rv1589)	Biotin biosynthesis	dethiobiotin
<i>thiE</i> (Rv0414c)	Thiamine biosynthesis	2-(4-methylthiazol-5-yl)ethyl phosphate and [(4-amino-2-methyl-pyrimidin-5-yl) methoxy-oxido-phosphoryl] phosphate
<i>cysE</i> (Rv2335)	Cysteine biosynthesis	L-serine and acetyl-CoA
<i>cobC</i> (Rv2231c)	No pathway assigned	L-threonine O-3-phosphate
<i>glpX</i> (Rv1099c)	glycolysis and gluconeogenesis	D-fructose 1,6-bisphosphate
<i>ppgK</i> (Rv2702)	Amino sugar and nucleotide sugar metabolism Gluconeogenesis	β-D-glucose
<i>arcA</i> (Rv1001)	arginine degradation V (arginine deiminase pathway)	L-arginine
<i>panD</i> (Rv3601c)	β-alanine biosynthesis IV	L-aspartate
<i>otsA</i> (Rv3490)	trehalose biosynthesis I	UDP-D-glucose and α-D-glucose 6-phosphate

The molecules were then sorted to focus on those passing SMARTS, molecular weight (MWT) 280–430 g/mol, logP 3–5, polar surface area PSA 50–100 Å², Bayesian score in the ‘single point model’ >0.3, Bayesian score in the ‘dose response model’ >1.37 and Bayesian score in the ‘Novartis model’ >1.11, signified predicted activity. These Bayesian score cutoff values and physicochemical parameter limits came from previous dataset analysis and model building to represent the boundary between active and inactive compounds against TB in whole cells (28,32,33). A set of 60 molecules was then sorted based on the Bayesian score dose response cut off (as this represents the highest quality dataset [compared to the single point model] using compounds with data from public datasets from Southern Research Institute (25)) and was exported to Excel before further filtering to manually exclude those already tested according to their presence in public databases in CDD. We also included 3 examples of compounds that had poor physicochemical properties (negative logP values, MWT<280) to further illustrate the importance of hydrophobicity on permeability and TB activity. We hypothesized that these would be inactive and/or would be unable to enter the cell. After sorting with the Bayesian model, 23 compounds for this study were imported into CDD, (Bayesian score dose response model range 1.6–11.8) including mimics of dethiobiotin (2), D-fructose 1,6-bisphosphate (17), UDP-glucose (3), L-serine (1) and L-arginine (1).

Measurement of Antibacterial Activity Against Mtb

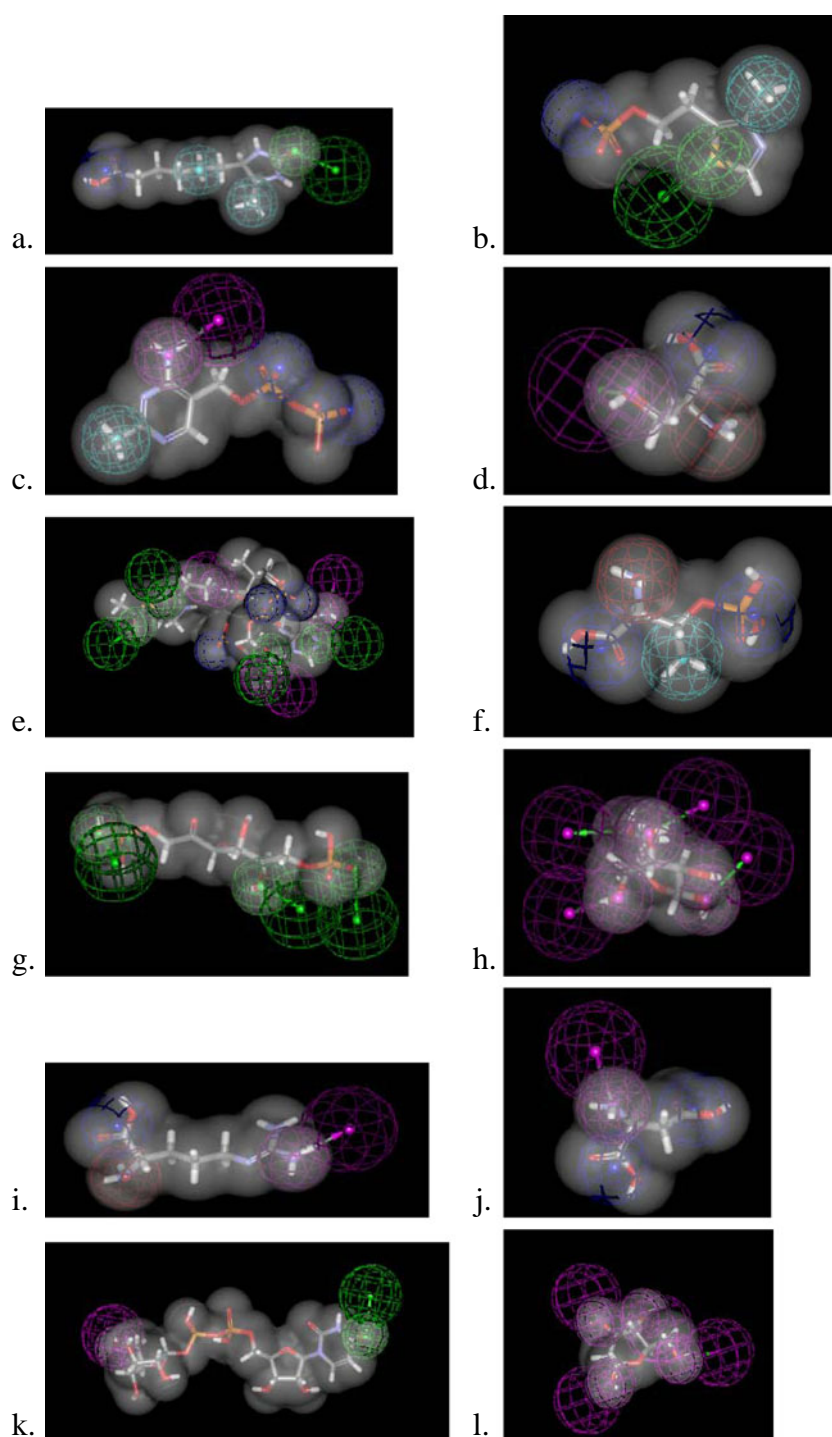
From the set of 23 compounds tested, two compounds showed moderate minimal inhibitory concentration (MIC) values against cultured Mtb. These are suggested to be mimics of D-fructose 1,6 -bisphosphate. N-(2-furylmethyl)-N'-[(5-nitro-3-thienyl)carbonyl]thiourea and N'-[(5-nitro-3-thienyl)carbonyl]-N''-(2-thienylmethyl)thiourea exhibited MIC values of 40 and

20 µg/ml, respectively (Fig. 5). The remaining compounds had MIC values >40 µg/ml (data not shown). Control MIC values for rifampicin and isoniazid were 0.0063 and 0.063 µg/ml, respectively, which are consistent with reported values in the literature as annotated in the CDD (TB efficacy data from the literature). All MIC data for compounds that showed activity were shared in the CDD database (Fig. 5c). It should be noted that as hypothesized the 3 compounds selected with poor logP and low MWT showed no activity against TB.

DISCUSSION

Relatively little attention has been paid to the integration of different types of biological, chemical and literature data for TB (13). Database integration is an important current trend in informatics-driven pharmaceutical discovery. Databases like TBCyc, SRI's BioCyc collection (34,35), and Pathway Logic models (36–39) are rich resources for biological networks and pathways. These knowledgebases provide systems level information for genomic, transcriptomic, proteomic and pathway context for proteins from more than 1100 organisms (prokaryotic and eukaryotic) including human. CDD, a widely used web-based drug discovery software platform, contains the CDD TB database, which incorporates biology, chemistry, molecular structure and physical property data for small molecules that are potentially valuable chemical tools, collated from the literature, patents and unpublished data obtained from the research network (25,28,40). Integration of target proteins and small molecule information through SRI databases, models, and analysis tools, and CDD TB database provide a synergistic computational environment for hypotheses testing, knowledge sharing, data archiving, data mining and drug discovery.

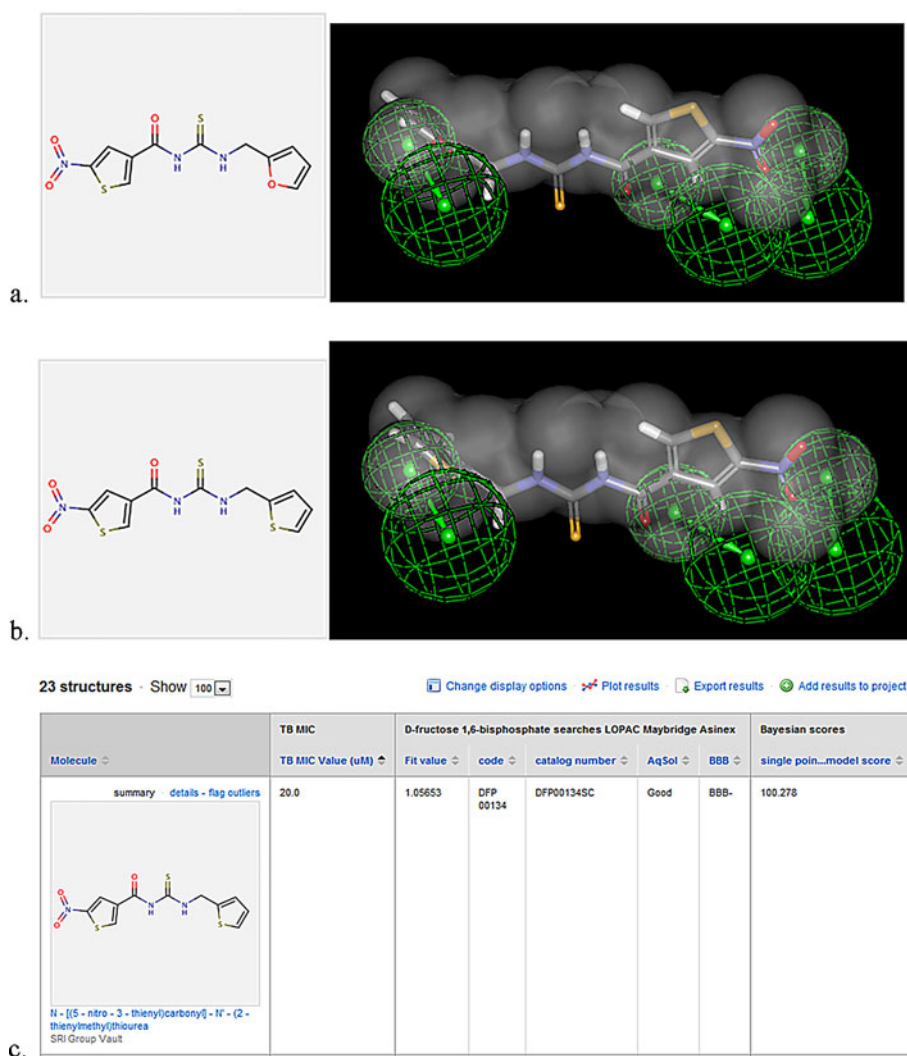
Fig. 4 *In vivo* essential metabolites and pharmacophores. (a) dethiobiotin, (b) 2-(4-methylthiazol-5-yl)ethyl phosphate, (c) [(4-amino-2-methyl-pyrimidin-5-yl)methoxy-oxido-phosphoryl] phosphate, (d) L-serine, (e) 2-[[[4-[[3-(2-acetylsulfanylethylamino)-3-oxo-propyl]amino]-3-hydroxy-2,2-dimethyl-4-oxo-butoxy]-oxido-phosphoryl]oxymethyl]-5-(6-aminopurin-9-yl)-4-hydroxy-tetrahydrofuran-3-yl] phosphate, (f) L-threonine O-3-phosphate, (g) D-fructose 1,6-bisphosphate, (h) β -D-glucose, (i) L-arginine, (j) L-aspartate, (k) UDP-D-glucose, l. α -D-glucose 6-phosphate. Starting with each such metabolite structure, a 3D pharmacophore was developed using Accelrys Discovery Studio 2.5.5 (Accelrys, San Diego, CA) from 3D conformations of the metabolite. This identified key features, onto which was mapped a van der Waals surface for the metabolite (15, 25, 26). Features represent: cyan = hydrophobic, green = hydrogen bond acceptor, purple = hydrogen bond donor, blue = negative ionizable and red = positive ionizable. The grey surface represents the van der Waals surface of the metabolite.



The development of the CDD database has been described previously with applications for collaborative malaria (40) and TB research (25,28). The literature data on Mtb drug discovery has been curated and over ~20 Mtb specific datasets are hosted, representing well over 300,000 compounds derived from patents, literature and high throughput screening (HTS) data. CDD have recently made several large HTS datasets of compounds

for TB and malaria available publically (41). We have also undertaken a manual evaluation of these and other datasets using a simple descriptor analysis as well as readily available substructure alerts or “filters” (28,32,33). By creating a very large collaborative database CDD TB, we have been able to compare inactive and active molecules against Mtb and show which molecular properties are important for activity in whole cells (25,27,28). We have

Fig. 5 Two suggested mimics of D-fructose 1,6 bisphosphate (see also Fig. 4g for comparison) (a) DFP000133SC and (b) DFP000134SC with MIC values of 40 and 20 $\mu\text{g/ml}$, respectively. These molecules are also shown mapped to the pharmacophore and shape based on D-fructose 1,6-bisphosphate. (c) Image of data stored and securely collaboratively shared in CDD, showing molecule structure, MIC, pharmacophore and Bayesian model predictions etc. Green pharmacophore features represent hydrogen bond acceptors. The grey surface represents the van der Waals surface of the metabolite.



previously performed multiple computational analyses that provided strong preliminary evidence for the value of the TB machine learning (Bayesian) models used in this study for prioritizing the compounds (25,27,28). We have observed from 4 to over 10 fold enrichment factors. These results also showed that computational models generated with whole-cell screening data from one laboratory rank ordered compounds screened and identified as Mtb hits by independent laboratories, according to different assays (27). In total these analyses present strong evidence that such models can be used for prioritizing compounds herein.

Preliminary experiments showed two compounds (N-(2-furymethyl)-N'-[(5-nitro-3-thienyl)carbonyl]thiourea and N-[(5-nitro-3-thienyl)carbonyl]-N'-(2-thienylmethyl)thiourea) which inhibit the growth of Mtb, and may represent a starting point for further optimization. These two compounds were suggested as mimics of D-fructose 1,6-bisphosphate, exhibiting FitValues of 0.79 and 1.05, respectively, for the 3D-pharmacophore of the

metabolite. Intriguingly, these FitValues ranked them 470 and 377, respectively out of 608 compounds that were scored from a total of >80,000 molecules in the Maybridge, Asinex and LOPAC databases. Future work could evaluate some of the compounds scored with higher FitValues but which may have scored poorly with our other filters. Also, the two acylthioureas ((N-(2-furymethyl)-N'-[(5-nitro-3-thienyl)carbonyl]thiourea and N-[(5-nitro-3-thienyl)carbonyl]-N'-(2-thienylmethyl)thiourea)) exhibit Tanimoto similarities of 0.28 and 0.24, respectively, in comparison with D-fructose 1,6-bisphosphate when using MDL public key fingerprints (in Accelrys Discovery Studio). This implies that the pharmacophore method can identify compounds that are not similar in 2D to the starting molecule used.

It is important to note that the pharmacophore model of D-fructose 1,6-bisphosphate was created with the phosphates treated as hydrogen-bond acceptors. We have previously demonstrated that a “relaxed” pharmacophore

model can be useful in treating negative charges as solely hydrogen-bond acceptors (Fig. 4g and 5a, b) in the case of a metabolite with two negatively-charged groups at physiologic pH. This relaxation avoids the return of compounds with two formal negative charges as putative substrate mimics, which could be severely limited in their ability to cross the waxy Mtb cell wall, in the absence of active transport.

It is noteworthy that both putative mimics are of the acylthiourea chemotype, solely differing by the conservative replacement of a furan with a thiophene. This chemical type has been identified amongst hits in whole-cell phenotypic screens, looking for growth inhibition of cultured Mtb, without mention of a specific biological target. The published SRI screen of an approximately 100,000-member commercial diversity library disclosed this hit class *versus* H37Rv (42). Visual inspection of this dataset utilizing CDD (TAACF CB2 set) demonstrated a wide range of acylthiourea hits (>50% inhibition at 10 µg/mL compound), with alkyl, aryl, and heteroaryl substituents at the termini. Similar observations were made with the Southern Research Institute screen of approximately 215,000 compounds from the MLSCN SMR library (43) using CDD (MLSMR). This suggests the privileged nature of this chemotype and/or its ability to serve as a prodrug through activation of the thione moiety, in analogy to the thiourea isoxyl (44). Kachhadia and colleagues previously reported the synthesis and biological testing of

a series of acylthioureas, intriguingly containing a substituted benzothiophene attached via its 2-position to the acyl moiety. The eleven analogs, tested at a concentration of 6.25 µg/mL, inhibited the growth of H37Rv by 10–69% (45).

The two acylthioureas in this work were suggested as mimics of D-fructose 1,6-bisphosphate, a substrate of the enzyme fructose-1,6-bisphosphatase II (FBPase II; EC 3.1.3.11). This enzyme is encoded by the gene *glpX* (Rv1099c) of Mtb, which is a key enzyme of gluconeogenesis. FBPase II catalyzes the hydrolysis of fructose 1,6-bisphosphate to form fructose 6-phosphate and orthophosphate. This reaction is the reverse of that catalyzed by phosphofructokinase in glycolysis, and the catalytic product, fructose 6-phosphate, which is an important precursor in various biosynthetic pathways, is used to generate important structural components of the cell wall and glycolipids in mycobacteria. In all organisms, gluconeogenesis is an important metabolic pathway that allows the cells to synthesize glucose from non-carbohydrate precursors, such as organic acids, amino acids, and glycerol. Until recently, five different classes of FBPases have been identified based on their amino acid sequences (FBPases I to V). Eukaryotes possess only the FBPase I-type enzyme, but all five types exist in various prokaryotes. The *Mtb* FBPase II constitutes the only known FBPase in Mtb and has no human homologue. The *glpX* transposon mutant was predicted to be attenuated in TraSH experiments (17,18),

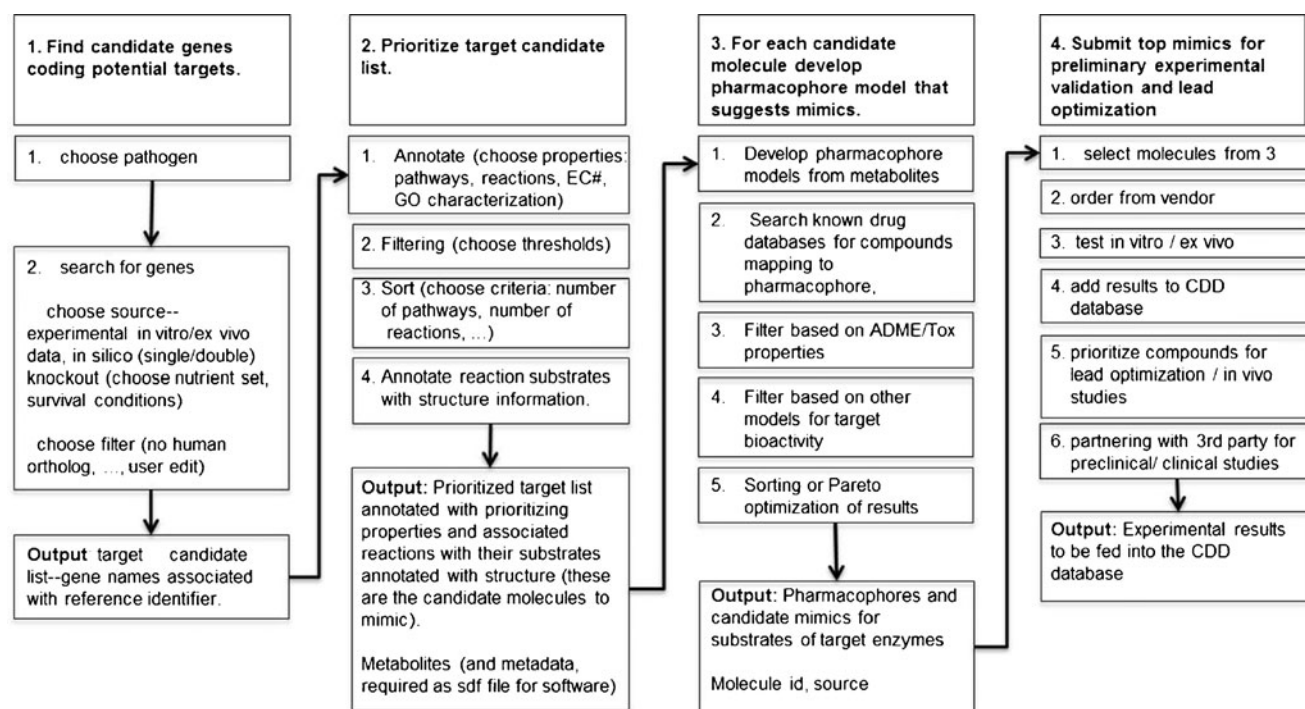


Fig. 6 Proposed generalized workflow for molecule discovery.

indicating a probable role of this enzyme in mycobacterial pathogenesis (46). In addition, FBPase II is an essential enzyme for Mtb *in vivo* and has not yet been targeted by any approved TB drugs. All the evidence collected in this study suggested it as a potential target for the mimic approach. Further experimental validation of the two postulated mimics of D-fructose 1,6-bisphosphate will be ultimately needed to confirm this.

The two Mtb growth inhibitors disclosed in this work were found via a multi-tiered, integrative informatics workflow that consists of a sequence of four main tasks as shown in the Fig. 6. Each task takes data produced from the previous task and produces data as input for the following task. Central to the translation from drug target to putative small molecule inhibitor is a strategy that may be viewed as intermediate between high-throughput screening and rational structure-based drug design. Intriguingly, it is possible that an approved drug might be found as a metabolite mimic and through repurposing could represent a novel antitubercular agent with little if any need for optimization prior to clinical trials (47). To date, an exhaustive screening of known drugs has not been performed by NIAID TAACF or others (48). Efforts to date have screened only a fraction of the known drugs, although thorough *in silico* screening is feasible using cheminformatics methods, such as those discussed in this work. In the current study, substrate mimicry afforded 2 hits, representing a 10% hit rate (if the three compounds intentionally selected to have suboptimal properties are excluded), that is higher than high throughput screening hit rates (frequently <1%) (49, 50). Such an approach may be a more efficient way to screen the vast array of known drugs or commercially available compounds for activity against Mtb.

ACKNOWLEDGMENTS & DISCLOSURES

S.E. kindly acknowledges CDD colleagues for developing the CDD TB database as well as the many TB research collaborators. M.S. and C.T. acknowledge the Biocyc group and TBDB for access to tools and data. J.S.F. acknowledges generous start-up funding from UMDNJ-New Jersey Medical School. The CDD TB database was made possible with funding from the Bill and Melinda Gates Foundation (Grant#49852 “Collaborative drug discovery for TB through a novel database of SAR data optimized to promote data archiving and sharing”). The project described was supported by Award Number R41AI088893 from the National Institute of Allergy And Infectious Diseases. The content is solely the responsibility of the authors and does not necessarily represent the official views of the National Institute Of Allergy And Infectious Diseases or the National Institutes of Health.

S.E. is a consultant for Collaborative Drug Discovery.

REFERENCES

- Balganesh TS, Alzari PM, Cole ST. Rising standards for tuberculosis drug development. *Trends Pharmacol Sci*. 2008;29:576–81.
- Cole ST. Learning from the genome sequence of *Mycobacterium tuberculosis* H37Rv. *FEBS Lett*. 1999;452:7–10.
- Weinand JR, Rubin EJ. The many roads to essential genes. *Tuberculosis* (Edinburgh, Scotland). 2008;88 Suppl 1:S19–24.
- Camacho LR, Ensergueix D, Perez E, Gicquel B, Guilhot C. Identification of a virulence gene cluster of *Mycobacterium tuberculosis* by signature-tagged transposon mutagenesis. *Mol Microbiol*. 1999;34:257–67.
- Wayneand LG, Hayes LG. An *in vitro* model for sequential study of shutdown of *Mycobacterium tuberculosis* through two stages of nonreplicating persistence. *Infect Immun*. 1996;64:2062–9.
- Dutta NK, Mehra S, Didier PJ, Roy CJ, Doyle LA, Alvarez X, Ratterree M, Be NA, Lamichhane G, Jain SK, Lacey MR, Lackner AA, Kaushal D. Genetic requirements for the survival of tubercle bacilli in primates. *J Infect Dis*. 2010;201:1743–52.
- Ostermanand AL, Begley TP. A subsystems-based approach to the identification of drug targets in bacterial pathogens. *Prog Drug Res*. 2007;64(131):133–70.
- Moir DT, Shaw KJ, Hare RS, Vovis GF. Genomics and antimicrobial drug discovery. *Antimicrob Agents Chemother*. 1999;43:439–46.
- Sacchettini JC, Rubin EJ, Freundlich JS. Drugs *versus* bugs: in pursuit of the persistent predator *Mycobacterium tuberculosis*. *Nat Rev Microbiol*. 2008;6:41–52.
- Ballel L, Field RA, Duncan K, Young RJ. New small-molecule synthetic antimycobacterials. *Antimicrob Agents Chemother*. 2005;49:2153–63.
- Payne DA, Gwynn MN, Holmes DJ, Pompliano DL. Drugs for bad bugs: confronting the challenges of antibacterial discovery. *Nat Rev Drug Disc*. 2007;6:29–40.
- Schneider G. Virtual screening: an endless staircase? *Nat Rev Drug Discov*. 2010;9:273–6.
- Ekins S, Freundlich JS, Choi I, Sarker M, Talcott C. Computational databases, pathway and cheminformatics tools for tuberculosis drug discovery. *Trends Microbiol*. 2011;19:65–74.
- Adams JC, Keiser MJ, Basuino L, Chambers HF, Lee DS, Wiest OG, Babbitt PC. A mapping of drug space from the viewpoint of small molecule metabolism. *PLoS Comput Biol*. 2009;5:e1000474.
- Lamichhane G, Freundlich JS, Ekins S, Wickramaratne N, Nolan S, Bishai WR. Essential metabolites of *M. tuberculosis* and their mimics. *Mbio*. 2011;2:e00301–00310.
- McAdam RA, Quan S, Smith DA, Bardarov S, Betts JC, Cook FC, Hooker EU, Lewis AP, Woollard P, Everett MJ, Lukey PT, Bancroft GJ, Jacobs Jr WR, Duncan K. Characterization of a mycobacterium tuberculosis H37Rv transposon library reveals insertions in 351 ORFs and mutants with altered virulence. *Microbiology*. 2002;148:2975–86.
- Sasseti CM, Boyd DH, Rubin EJ. Genes required for mycobacterial growth defined by high density mutagenesis. *Mol Microbiol*. 2003;48:77–84.
- Sassetiand CM, Rubin EJ. Genetic requirements for mycobacterial survival during infection. *Proceedings of the National Academy of Sciences of the United States of America*. 2003;100:12989–94.
- Lamichhane G, Tyagi S, Bishai WR. Designer arrays for defined mutant analysis to detect genes essential for survival of *Mycobacterium tuberculosis* in mouse lungs. *Infect Immun*. 2005;73:2533–40.
- Jain SK, Hernandez-Abanto SM, Cheng QJ, Singh P, Ly LH, Klinkenberg LG, Morrison NE, Converse PJ, Nuernberger E, Grosset J, McMurray DN, Karakousis PC, Lamichhane G, Bishai WR. Accelerated detection of *Mycobacterium tuberculosis* genes

- essential for bacterial survival in guinea pigs, compared with mice. *J Infect Dis*. 2007;195:1634–42.
21. Reddy TB, Riley R, Wymore F, Montgomery P, DeCaprio D, Engels R, Gellesch M, Hubble J, Jen D, Jin H, Koehrsen M, Larson L, Mao M, Nitzberg M, Sisk P, Stolte C, Weiner B, White J, Zachariah ZK, Sherlock G, Galagan JE, Ball CA, Schoolnik GK. TB database: an integrated platform for tuberculosis research. *Nucleic Acids Res*. 2009;37:D499–508.
 22. Galagan JE, Sisk P, Stolte C, Weiner B, Koehrsen M, Wymore F, Reddy TB, Zucker JD, Engels R, Gellesch M, Hubble J, Jin H, Larson L, Mao M, Nitzberg M, White J, Zachariah ZK, Sherlock G, Ball CA, Schoolnik GK. TB database 2010: overview and update. *Tuberculosis (Edinburgh, Scotland)*. 2010;90:225–35.
 23. Anishetty S, Pulimi M, Pennathur G. Potential drug targets in *Mycobacterium tuberculosis* through metabolic pathway analysis. *Comput Biol Chem*. 2005;29:368–78.
 24. Prathipati P, Ma NL, Manjunatha UH, Bender A. Fishing the target of antitubercular compounds: *in silico* target deconvolution model development and validation. *J Proteome Res*. 2009;8:2788–98.
 25. Ekins S, Bradford J, Dole K, Spektor A, Gregory K, Blondeau D, Hohman M, Bunin B. A collaborative database and computational models for tuberculosis drug discovery. *Mol BioSystems*. 2010;6:840–51.
 26. Zheng X, Ekins S, Rauffman J-P, Polli JE. Computational models for drug inhibition of the human apical sodium-dependent bile acid transporter. *Mol Pharm*. 2009;6:1591–603.
 27. Ekinsand S, Freundlich JS. Validating new tuberculosis computational models with public whole cell screening aerobic activity datasets. *Pharmaceut Res*. 2011;28:1859–69.
 28. Ekins S, Kaneko T, Lipinski CA, Bradford J, Dole K, Spektor A, Gregory K, Blondeau D, Ernst S, Yang J, Goncharoff N, Hohman M, Bunin B. Analysis and hit filtering of a very large library of compounds screened against *Mycobacterium tuberculosis*. *Mol BioSyst*. 2010;6:2316–24.
 29. Palomino JC, Martin A, Camacho M, Guerra H, Swings J, Portaels F. Resazurin microtiter assay plate: simple and inexpensive method for detection of drug resistance in *Mycobacterium tuberculosis*. *Antimicrob Agents Chemother*. 2002;46:2720–2.
 30. Collinsand L, Franzblau SG. Microplate alamar blue assay *versus* BACTEC 460 system for high-throughput screening of compounds against *Mycobacterium tuberculosis* and *Mycobacterium avium*. *Antimicrob Agents Chemother*. 1997;41:1004–9.
 31. Weininger D. SMILES 1. Introduction and encoding rules. *J Chem Inform Comput Sci*. 1988;28:31.
 32. Ekinsand S, Williams AJ. Meta-analysis of molecular property patterns and filtering of public datasets of antimalarial “hits” and drugs. *MedChemComm*. 2010;1:325–30.
 33. Ekinsand S, Williams AJ. When pharmaceutical companies publish large datasets: an abundance of riches or fool’s gold? *Drug Disc Today*. 2010;15:812–5.
 34. <http://biocyc.org>.
 35. Karp PD. Pathway databases: a case study in computational symbolic theories. *Science*. 2001;293:2040–4.
 36. <http://pl.csl.sri.com>.
 37. Tiwari A, Talcott C, Knapp M, Lincoln P, Laderoute K. Analyzing pathways using SAT-based approaches. In: Ania H, Horimoto K, Kutsia T, editors. *Algebraic biology*, vol. 4545. 2007. p. 155–69.
 38. Talcott C, Eker S, Knapp M, Lincoln P, Laderoute K. Pathway logic modeling of protein functional domains in signal transduction. *Pac Symp Biocomput* 2004;568–580.
 39. Talcott C. Symbolic modeling of signal transduction in pathway logic. In: Perrone LF, Wieland FP, Liu J, Lawson BG, Nicol DM, Fujimoto RM, editors. 2006 winter simulation conference. 2006. p. 1656–65.
 40. Hohman M, Gregory K, Chibale K, Smith PJ, Ekins S, Bunin B. Novel web-based tools combining chemistry informatics, biology and social networks for drug discovery. *Drug Disc Today*. 2009;14:261–70.
 41. Gamo F-J, Sanz LM, Vidal J, de Cozar C, Alvarez E, Lavandera J-L, Vanderwall DE, Green DVS, Kumar V, Hasan S, Brown JR, Peishoff CE, Cardon LR, Garcia-Bustos JF. Thousands of chemical starting points for antimalarial lead identification. *Nature*. 2010;465:305–10.
 42. Ananthan S, Faaleolea ER, Goldman RC, Hobrath JV, Kwong CD, Laughon BE, Maddry JA, Mehta A, Rasmussen L, Reynolds RC, Secrist 3rd JA, Shindo N, Showe DN, Sosa MI, Suling WJ, White EL. High-throughput screening for inhibitors of *Mycobacterium tuberculosis* H37Rv. *Tuberculosis (Edinburgh, Scotland)*. 2009;89:334–53.
 43. Maddry JA, Ananthan S, Goldman RC, Hobrath JV, Kwong CD, Maddox C, Rasmussen L, Reynolds RC, Secrist 3rd JA, Sosa MI, White EL, Zhang W. Antituberculosis activity of the molecular libraries screening center network library. *Tuberculosis (Edinburgh, Scotland)*. 2009;89:354–63.
 44. Kordulakova J, Janin YL, Liav A, Barilone N, Dos Vultos T, Raugier J, Brennan PJ, Gicquel B, Jackson M. Isoxyl activation is required for bacteriostatic activity against *Mycobacterium tuberculosis*. *Antimicrob Agents Chemother*. 2007;51:3824–9.
 45. Kachhadia VV, Patel MR, Joshi HS. Heterocyclic systems containing S/N regioselective nucleophilic competition: facile synthesis, antitubercular and antimicrobial activity of thiohydantoin and iminothiazolidinones containing the benzo[b]thiophene moiety. *J Serb Chem Soc*. 2005;70:153–61.
 46. Gutka HJ, Rukseree K, Wheeler PR, Franzblau SG, Movahedzadeh F. *glpX* gene of *mycobacterium tuberculosis*: heterologous expression, purification, and enzymatic characterization of the encoded fructose 1,6-bisphosphatase II. *Appl Biochem Biotechnol*. 2011;164:1376–89.
 47. Ekins S, Williams AJ, Krasowski MD, Freundlich JS. *In silico* repositioning of approved drugs for rare and neglected diseases. *Drug Disc Today*. 2011;16:298–310.
 48. Loughheed KE, Taylor DL, Osborne SA, Bryans JS, Buxton RS. New anti-tuberculosis agents amongst known drugs. *Tuberculosis (Edinburgh, Scotland)*. 2009;89:364–70.
 49. Polgar T, Baki A, Szendrei GI, Keseru GM. Comparative virtual and experimental high-throughput screening for glycogen synthase kinase-3 β inhibitors. *J Med Chem*. 2005;48:7946–59.
 50. Doman TN, McGovern SL, Witherbee BJ, Kasten TP, Kurumbail R, Stallings WC, Connolly DT, Shoichet BK. Molecular docking and highthroughput screening for novel inhibitors of protein tyrosine phosphatase-1B. *J Med Chem*. 2002;45:2213–21.