

ARTICLE OPEN



Quantitative prediction of conditional vulnerabilities in regulatory and metabolic networks using PRIME

Selva Rupa Christinal Immanuel 1 , Mario L. Arrieta-Ortiz¹, Rene A. Ruiz 1 , Min Pan¹, Adrian Lopez Garcia de Lomana^{1,5}, Eliza J. R. Peterson 1 and Nitin S. Baliga 1 . 1,2,3,4

The ability of *Mycobacterium tuberculosis* (Mtb) to adopt heterogeneous physiological states underlies its success in evading the immune system and tolerating antibiotic killing. Drug tolerant phenotypes are a major reason why the tuberculosis (TB) mortality rate is so high, with over 1.8 million deaths annually. To develop new TB therapeutics that better treat the infection (faster and more completely), a systems-level approach is needed to reveal the complexity of network-based adaptations of Mtb. Here, we report a new predictive model called PRIME (**P**henotype of **R**egulatory influences **I**ntegrated with **M**etabolism and **E**nvironment) to uncover environment-specific vulnerabilities within the regulatory and metabolic networks of Mtb. Through extensive performance evaluations using genome-wide fitness screens, we demonstrate that PRIME makes mechanistically accurate predictions of context-specific vulnerabilities within the integrated regulatory and metabolic networks of Mtb, accurately rank-ordering targets for potentiating treatment with frontline drugs.

npj Systems Biology and Applications (2021)7:43; https://doi.org/10.1038/s41540-021-00205-6

INTRODUCTION

Mycobacterium tuberculosis (Mtb) kills more people than any other microbe, and it has thus far resisted every attempt to bring the pandemic under control. Part of the pathogen's success is its ability to diversify itself phenotypically and survive both host and drug bactericidal action¹⁻³. Phenotypic heterogeneity (both stochastically and environmentally induced) seems to be an intrinsic characteristic of the pathogen and a major reason why standard chemotherapy of tuberculosis (TB) requires 6 months of treatment, and 5% of cases are not cured even then^{4,5}. To develop better interventions that account for pathogen heterogeneity, we need to identify the most important factors (e.g., transcriptional regulators) that create variation as well as the downstream effectors (e.g., regulatory target genes) that mediate drug tolerance.

Metabolic activity undoubtedly contributes to Mtb phenotypic heterogeneity and antibiotic tolerance. For example, changes in metabolism can affect the amount of drug target present⁶, the ability to generate toxic products⁷, and the efflux of antibiotics⁸. Mtb alters its growth and metabolism in response to stressful conditions through regulatory programs primarily encoded at the transcriptional level. Indeed, modeling host-related stresses in vitro produces large transcriptional changes in Mtb, particularly in metabolic pathways; consistently ~25% of differentially expressed genes are metabolic genes from (GSE116353)⁹, acidic pH (GSE165514), or nutrient limited (GSE165673) conditions. To develop effective antibiotic regimens, we need to understand at a systems-level and mechanistic-level how specific regulatory mechanisms conditionally activate and repress genes to redirect flux through metabolic networks to generate and support drug tolerant phenotypes. This mechanistic understanding will uncover new vulnerabilities in Mtb's regulatory and metabolic networks that can be rationally targeted in new drug regimens to achieve faster and complete clearance of the pathogen.

Previously, approaches to model the influence of transcriptional regulation on metabolism have used boolean logic (regulatory Flux Balance Analysis—rFBA)¹⁰, protein–DNA (P–D) interactions (Probabilistic Regulation Of Metabolism—PROM)^{11,12}, and regression-based regulatory influences (Integrated Deduced REgulation And Metabolism—IDREAM)¹³ to predict how transcriptional regulation of enzyme-coding genes modulates flux through their catalyzed reactions. Briefly, rFBA models the influence of transcriptional regulation on metabolism using boolean "on or off" states of metabolic genes, depending on the expression level of a transcription factor (TF) and its implicated role as a putative activator or repressor of that gene. The extensive manual curation required to develop rFBA and its inability to model TF activity as a continuous (i.e., not boolean) function greatly limits its application and accuracy. In contrast, PROM outperformed rFBA by using a probabilistic approach to model the regulation of a metabolic gene by a TF using a compendium of transcriptome profiles to calculate probabilities¹¹. However, PROM is limited in that it relies on a P-D interaction map for the regulatory network. P-D interactions are typically generated in a limited set of conditions by using an overexpressed TF as a bait to enrich and locate its genome-wide binding locations. P-D interactions are fraught with false positives (due to TF overexpression) and false negatives (due to lack of context for TF regulation across environmental conditions). Notwithstanding these caveats, PROM was useful in uncovering the mechanism by which pretomanid potentiates bedaquiline action on Mtb by disrupting a regulatory network that confers tolerance to the recently FDA-approved drug¹⁴. A third model, IDREAM addressed the shortcoming of using P-D interactions in PROM by constraining flux using TF regulatory influences from a predictive systems-scale Environment and Gene Regulatory Influence Network (EGRIN) model. An EGRIN model is inferred in two steps using (a) cMonkey, which identifies the specific context in which subsets of genes are co-regulated

¹Institute for Systems Biology, Seattle, WA, USA. ²Departments of Biology and Microbiology, University of Washington, Seattle, WA, USA. ³Molecular and Cellular Biology Program, University of Washington, Seattle, WA, USA. ⁴Lawrence Berkeley National Lab, Berkeley, CA, USA. ⁵Present address: Center for Systems Biology, University of Iceland, Reykjavik, Iceland. [©]email: eliza.peterson@isbscience.org; nitin.baliga@isbscience.org





(biclusters) by a conserved regulatory mechanism(s); and (b) Inferelator, which predicts TFs and environmental factors that causally influence the expression of genes within those biclusters^{15–17}. By integrating false discovery rates (FDR) for EGRIN-inferred regulatory influences, IDREAM achieved significantly better performance than rFBA and PROM in predicting synthetic lethal interactions between TFs and metabolic genes in yeast¹³. However, IDREAM does not incorporate quantitative environment-specific TF regulatory influences that are modeled by EGRIN, and is therefore also limited in accurately predicting environment-specific consequences of TF perturbations. For the reasons stated above, PROM, rFBA, and IDREAM are limited in their ability to predict environment-specific phenotypic consequences of perturbations to TFs.

Additionally, there are algorithms (e.g., OptORF¹⁸, EMILiO¹⁹, and BeReTa²⁰) that have the potential to predict the consequence of regulatory and metabolic network perturbations. They were originally designed to identify perturbations that maximize flux towards a desired metabolite and some of their features make them not well-suited for predicting systems-wide conditional outcomes of TF perturbation. For instance, OptORF¹⁸ and EMILiO¹⁹ use binary or fixed weights to model TF influences, which does not capture changes in relative strength of transcriptional regulation of metabolic genes across environments. By contrast, BeReTa²⁰ does take into account weighted, combinatorial influences of TFs, but the analysis is restricted to genes encoding reactions of specific pathways of interest to an industrial application. Thus, none of these algorithms were designed to predict systems level phenotypic consequences (e.g., fitness and growth rate) of perturbations to the transcriptional network.

Here, we report the development of Phenotype of Regulatory influences Integrated with Metabolism and Environment (PRIME), which incorporates environment-dependent combinatorial requlation of metabolic genes to mechanistically predict how individual TFs contribute to the phenotype of Mtb in any given environment. Through the use of comprehensive experimental validations, we demonstrate that PRIME significantly outperforms the previous methods in accurately predicting regulatory and metabolic genes that are conditionally required for growth on carbon sources that are specific for in vitro (glycerol) and in vivo (cholesterol) growth of Mtb. Further, PRIME has uncovered the interplay of regulatory and metabolic mechanisms that underlies Mtb's response to drug treatment. The accuracy of PRIME in predicting quantitative phenotypic effects of TF perturbations is demonstrated by high correlation between predicted and experimentally validated consequences of knocking out all metabolism-associated TFs (one-at-a-time) on isoniazid (INH) treatment-specific fitness of Mtb strains. Through this analysis, we have discovered new vulnerabilities in Mtb that can potentiate INH action, which are supported by experimental validation.

RESULTS

Condition-specific integration of regulation and metabolism using PRIME

A causal and mechanistic model of the transcriptional regulatory network and its quantitative influence on metabolic flux is required to characterize how the 214²¹ TFs encoded in the Mtb genome enable its physiological adaptations to disparate host relevant contexts including antibiotic treatment. Using the Inferelator^{15,22,23}, which applies a Bayesian regression-based approach to estimate TF activity (TFA), we constructed an EGRIN network from a compendium of 664 transcriptomes for Mtb that represented transcriptional changes across 77 environmental conditions including drug treatment, pH, oxygen and carbon source utilization (Supplementary Data 1) (http://www.colombos.net/)²⁴. Relative changes in the expression of every gene across all

conditions were modeled as the sum of weighted influences for a minimal set of TFs. Altogether in the EGRIN network, 142 TFs were implicated in the regulation of 2905 genes acting through a combinatorial scheme represented by 4820 TF–gene interactions (see Supplementary Data 2 for details). The EGRIN network recapitulated 2410 of the 4546 TF–gene interactions from a P–D network of Mtb, which was derived through both physical binding (from ChIP-seq experiments) and functional evidence (from transcriptional profiling)^{21,25}. We refer to this subset of 2410 weighted TF influences as the "EGRIN-PD" network (Supplementary Data 2 and Supplementary Fig. 1).

We investigated the degree to which EGRIN and EGRIN-PD models captured the regulation of 1011 genes that encode enzymes implicated in catalyzing 1229 reactions in the iEK1011²⁶ model of the Mtb metabolic network. This analysis demonstrated that whereas EGRIN-PD modeled 1252 regulatory influences of 104 TFs on 605 genes associated with 409 metabolic reactions, EGRIN modeled 2568 regulatory influences of 129 TFs on 750 genes associated with 725 metabolic reactions. We leveraged the EGRIN and EGRIN-PD wiring diagrams and weights of regulatory influences inferred by the Inferelator to predict how change in the activity of a TF in a given environment manifests in altered flux through a metabolic reaction catalyzed by their regulated gene product. In order to integrate regulation with metabolism, we had to account for combinatorial regulation of metabolic genes, with each of 349 out of the 750 metabolic genes predicted to be putatively regulated by ≥2 TFs and 111 TFs were predicted to regulate ≥2 metabolic genes (Supplementary Fig. 2 and Supplementary Data 3), and association of ≥2 gene products to each of 313 reactions in Mtb.

PRIME calculates "quantitative influence" of each TF on a target gene by taking the product of the Inferelator-assigned regression weight (β) of that TF-gene interaction and the expression level of the TF. The absolute expression level of a TF is calculated as a scaled value of signal intensity (for microarray data) or read counts (for RNA-Seg) based on distribution of values across the transcriptome compendium (Fig. 1a; see "Methods" section). For a metabolic gene that is regulated by multiple TFs, PRIME calculates the relative contribution of each TF to the regulation of that gene in a given environment by dividing its quantitative influence with the sum of quantitative influences of all TFs that regulate that gene. In this scheme, a TF will have a large relative consequence on the expression of a metabolic gene in an environment in which the TF is active and in high abundance, and the influences of other TFs are minimal. But the relative contribution of the TF will be proportionally lower if other TFs are also actively regulating that gene in that environment. Thus, this approach accounts for regulation of a metabolic gene by multiple TFs, and it simultaneously corrects for environmentspecific changes in combinatorial regulatory schemes. For a TF that regulates multiple genes encoding enzymes or enzyme subunits for the same reaction, we consider the largest regulatory influence of that TF on any of those genes to predict its influence on flux through that reaction. Thus, together these advancements also account for complex combinatorial associations between regulation and metabolism to assign a single reaction influence factor (RIF, γ) to each TF-reaction association. The consequence of TF regulation (or knockout) on flux through a reaction is calculated by multiplying this final TF-induced relative inhibition of that reaction (i.e., RIF) to the maximum possible flux through that reaction. In this manner, by updating upper bounds of flux through all reactions catalyzed by regulated gene products of a specific TF, PRIME constrains the metabolic network to a new solution space, to enable the prediction of "environment-specific" growth consequences of perturbing a given TF (Fig. 1b and Supplementary Fig. 1).

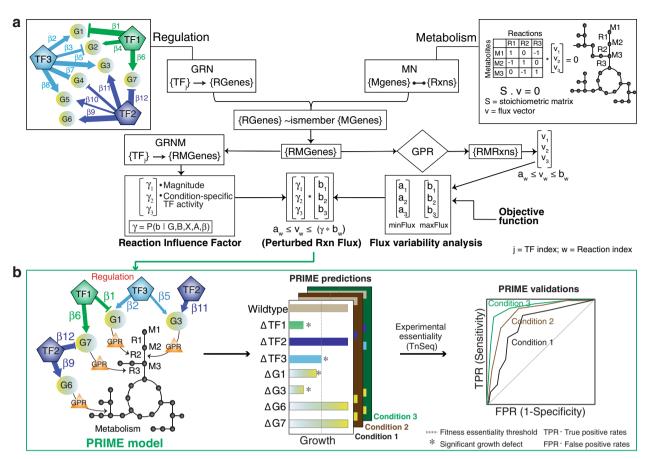


Fig. 1 Schematic for PRIME model development and performance assessment. a Schema for integration of gene regulation and metabolism. The gene regulatory network (GRN) models weighted regulatory influences of TFs on regulated genes (RGenes). A subset of the RGenes are enzyme-coding metabolic genes (Mgenes), whose functions are also modeled through gene-to-protein-to-reaction (GPR) mapping in a stoichiometric matrix representation of the metabolic network (MN). PRIME uses the integrated Gene Regulatory Network of Metabolism (GRNM) and a reaction flux influence estimator (ReFInE) to calculate the γ factor, which quantifies how the differential expression of multiple TFs and their weighted regulatory influences on a regulated metabolic gene (RMGene) manifest in altered flux (a_w: minimum flux; b_w: maximum flux) through the associated metabolic reaction (RMRxn) in a given environmental condition. b Illustration of condition-specific gene phenotype predictions and performance assessment. The example illustrates how PRIME predicts relative growth consequence of single gene knockouts in TFs (e.g., TF1, TF2, and TF3) and RMGenes (e.g., G1, G3, G6, and G7) in different contexts (e.g., Condition 1, 2, and 3). The vertical line in the barplot depicts a user-defined threshold in growth inhibition, below which a gene is deemed essential. Performance of PRIME is quantified using a receiver operating characteristic (ROC) curve based on accuracy of PRIME-predicted essential and non-essential genes in a given condition to experimentally determined phenotype consequences using transposon mutagenesis coupled with sequencing (TnSeq) in the same condition.

Performance assessment of PRIME

We investigated if the advancements in PRIME significantly improved its performance over PROM and IDREAM vis-à-vis accuracy of predicting environment-specific phenotypic consequences of knocking out TFs (Fig. 2a and Supplementary Fig. 3). To perform this comparison, we first updated the PROM model with the latest version of the Mtb P–D interaction map 12,21 and the current version of the metabolic network model iEK1011²⁶ (1011 genes encoding enzymes for 1229 reactions) that was used to construct PRIME. Using the methodology described in the original PROM paper^{11,12}, 2416 out of 2555 P-D interactions for 104 TFs were mapped to 605 genes assigned to 632 reactions in the iEK1011 metabolic network model. This represents a significant improvement in the overall coverage of TFs and metabolic genes in the PROM model (Table 1 and Fig. 2a). In parallel, we also developed the IDREAM model for Mtb by incorporating FDR values for 3643 regulatory influences for 142 TFs within the EGRIN model (FDR < 0.25) on a total of 641 genes associated with 639 reactions within iEK1011 (Table 1 and Fig. 2b). The slightly higher numbers of TFs and metabolic genes in IDREAM and PRIME (Fig. 2b) are because they use the Inferelator-derived regulatory network model, which has better coverage of genome-wide TF regulation across diverse environments, relative to the P–D interaction map generated in standard growth conditions that was used in PROM (Fig. 2c). However, the numbers of TF–gene interactions in PRIME and IDREAM were not identical. This is because IDREAM incorporates TF–gene interactions using FDRs derived using a modified elastic net²² and multiple Inferelator runs on different subsets of the transcriptome compendium. By contrast, PRIME uses a confidence score and the regression coefficient (β) deduced using Bayesian regression²² (Supplementary Fig. 4 and see "Methods" section). Hence, although the updated PROM and IDREAM models were not identical, they were similar to PRIME in terms of coverage of the total number of TFs and metabolic genes allowing comparative analysis of their performance (Table 1).

We compared the performance of PRIME, PROM and IDREAM by assessing their accuracy (sensitivity and specificity) in predicting environment-specific consequences of knocking out TFs on the growth of Mtb in minimal medium with glycerol or cholesterol as the carbon source. While Mtb is typically grown with glycerol during in vitro culture, the pathogen is capable of utilizing

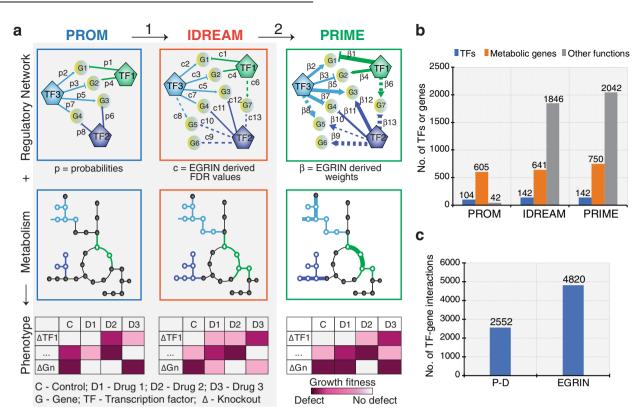


Fig. 2 PRIME model advancements. a Advancements in PRIME over previous methods (PROM and IDREAM) are indicated as (1) incorporation of regulatory influences from EGRIN (regression-based interactions are shown as dotted lines), which increases coverage of the regulatory network, (2) incorporation of the magnitude of regulatory influence of TFs on metabolic genes (β—shown as varying edge thickness) instead of probability (p) and FDR values (c) significantly improved the predictive accuracy of environment-specific gene essentiality. **b** Number of TFs and genes from PRIME, IDREAM, and PROM. **c** Number of TF-gene interactions identified using regression-based EGRIN and Protein–DNA (P–D) interactions from ChIP-seq data.

host-derived lipids, such as cholesterol, during infection. It is known that distinct metabolic genes and networks are associated with these two modes of growth. Accuracy of model predictions were evaluated using a leave-one-out cross validation (LOOCV) strategy²⁷ for comparison of model predictions to experimentally determined phenotypic consequences of transposon mutagenesis in genome-wide fitness screens (TnSeq) of Mtb cultured with glycerol or cholesterol^{28,29}. Specifically, for each model we generated a set of receiver-operating characteristic (ROC) curves by plotting the true positive rates (TPR) (i.e., proportion of modelpredicted essential genes that were verified by experiment) and false positive rates (FPR) (i.e., proportion of model-predicted essential genes that were experimentally determined to be nonessential) by leaving out one TF in each analysis. The distribution of area under the ROC curves (ROC AUC) from the LOOCV analysis of model predictions of which TFs are essential for Mtb growth was used as a metric of performance.

First, we investigated the influence of EGRIN-PD or EGRIN networks generated using different Inferelator 15,22,23 parameter settings (Zellner's g value) on the performance of PRIME. Briefly, in the Bayesian implementation of the Inferelator, a modified Zellner's g value is used to control the amount of variance in the model that is explained by prior information (i.e., a P–D network) 22 . In other words, a g value = 1 gives no preference to interactions supported by prior information over the novel TF–gene interactions. By increasing values of g>1, more preference is given to TF–gene interactions from the prior network, resulting in a higher recall of the EGRIN-PD network. We evaluated performance of PRIME using both EGRIN and EGRIN-PD regulatory networks generated using four g values. This analysis demonstrated that the performance of PRIME was

relatively robust to g values and significantly better with EGRIN over EGRIN-PD, irrespective of the g-value (Supplementary Fig. 5). The potential explanation for higher performance with EGRIN is that a regression-based approach uncovers a greater number of novel TF-gene interactions that are conditionally active in different environmental contexts, relative to a P-D interaction map that is typically generated in a limited number of environmental contexts using methodologies such as ChIP-seq. Therefore, here onwards all results reported for PRIME are based on regulatory influences from the EGRIN network.

The LOOCV analysis demonstrated that the performance of PRIME was significantly better relative to PROM and IDREAM in both cholesterol and glycerol carbon sources (Fig. 3a, b and Supplementary Fig. 6). In addition to providing a rigorous means for performance evaluation, the LOOCV analysis also identified a clear division of TFs in terms of their ROC-AUC values for the PRIME model. Further analysis revealed that the top performing TFs (20 and 12 TFs for glycerol and cholesterol, respectively) contributed maximally (up to 65% of overall biomass accumulation) to the overall fitness of Mtb (Supplementary Data 4). Out of 119 TFs with TnSeg data, the cholesterol-specific fitness consequences of knocking out 65% of all TFs (77 TF KOs) were accurately predicted by PRIME, whereas IDREAM and PROM accurately predicted only 45% (53 TFs) and 30% (35 TFs), respectively (Fig. 3c). Similarly, PRIME accurately predicted glycerol-specific fitness consequences for knocking out 92 out of 119 TFs (77%), whereas IDREAM accurately predicted 55% (65 TFs) and PROM predicted 36% (43 TFs) (Fig. 3d). In addition, we compared the performance of PRIME to IDREAM and IDREAMhybrid models, both generated with EGRIN FDR cutoffs of 0.05 and 0.25. PRIME outperformed both IDREAM and IDREAM-hybrid in

Table 1. Summary of PROM, IDREAM, and PRIME model features.						
Mtb model features	MTBPROM1.0 ¹¹	MTBPROM2.0 ¹²	PROMª	PROM ^a (EGRIN-PD)	IDREAM ^b	PRIME
Metabolic model	iNJ661	iSM810	iEK1011	iEK1011	iEK1011	iEK1011
Number of reactions	1025	938	1229	1229	1229	1229
Number of metabolic genes in the metabolic network	661	810 (759 genes in iEK1011)	1011	1011	1011	1011
Regulatory network	Balazsi 2008	Minch 2015	Minch 2015	EGRIN-PD ^c	EGRIN ^c by Elastic net (FDR < 0.25)	EGRIN ^c by Bayesian regression (Precision = 50%)
Number of transcription factors	30	104	104	131	142	142
Number of interactions	218	2555	2555	2410	3643	4820
Number of genes in the regulatory network (metabolic / total)	178/178	647/647	605/647	650/1509	641/2487	750/2905

Chandrasekaran et al. 11, Ma et al. 12.

predicting phenotypic consequences of TF KOs on Mtb growth in glycerol and cholesterol (Supplementary Fig. 7). Interestingly, PRIME EGRIN also outperformed PRIME EGRIN-PD, again reinforcing how a regression-based approach better captures meaningful environment-specific TF regulation relative to physical P-D interactions mapped using ChIP-seq in one condition with overexpressed TFs. Finally, we compared performance of PROM and PRIME in predicting gene essentiality using the same EGRIN-PD network, which excludes regression-inferred regulatory influences not supported by ChIP-seq. This comparison demonstrates that PRIME outperforms PROM even when both methods use the same EGRIN-PD network. Notably, the performance of PROM was similar irrespective of whether EGRIN-PD or the ChIP-seq derived network was used (Supplementary Fig. 8).

Using PRIME, 22 and 7 TFs were accurately predicted (either essential or non-essential) for growth only with either glycerol or cholesterol, respectively, as determined by experimental fitness screening (Fig. 3e). Similarly, 51 and 25 metabolic genes were accurately predicted by PRIME for growth on either glycerol or cholesterol, respectively (Fig. 3f and see Supplementary Data 4 for a complete list of PRIME-predicted and experimentally determined consequences of knocking out 142 TFs on Mtb growth on glycerol and cholesterol). Among the PRIME predicted essential TFs, Rv2506, Rv3050c, Rv2760c, and Rv0348 are essential for growth on cholesterol, presumably because they conditionally regulate genes encoding enzymes or enzyme subunits catalyzing essential metabolic processes during cholesterol utilization (Fig. 3g). For example, Rv2506 represses genes likely to be involved in branched-chain amino acid catabolism, which leads to the production of acetyl-coA and propionyl-coA³⁰. Propionyl-coA is also an endpoint of cholesterol degradation and can be toxic to Mtb³¹. It is possible that Rv2506 repression of branched-chain amino acid metabolism genes prevents accumulation of toxic metabolic intermediates during growth on cholesterol. All in all, perturbation of cholesterol utilization in Mtb could induce metabolite intoxication³¹, unbalanced central metabolism³² or lead to carbon starvation³³. As such, TFs such as Rv2506, Rv3050c, Rv2760c, and Rv0348 represent potential vulnerabilities in the cholesterol utilization pathways of Mtb that could be targeted by drugs. Notably, these TFs were also ascertained to be essential by the TnSeq screen performed with cholesterol as the carbon source²⁸ and are non-essential in glycerol (shown as inactive nodes in Fig. 3h). Other TFs (Rv1990c, Rv0023, and Rv0757) were predicted (and validated by TnSeq²⁸) to be essential for growth with both carbon sources or only essential for growth on glycerol (e.g., Rv0238 and Rv1423).

PRIME rank identifies the essential transcriptional factors and genes for survival during drug treatment

We used PRIME to investigate the regulatory and metabolic networks that drive physiological adjustments (e.g., cell wall modifications, shifts in metabolism and respiration) to enable the pathogen to survive and persist during drug treatment. To expose novel network vulnerabilities of Mtb in response to drug treatment, we generated transcriptome profiles of Mtb treated for 24 h with high-doses and low-doses of seven drugs (Supplementary Data 5). The transcriptome profiles were analyzed using the PRIME model to identify the metabolic networks and their associated regulators that were essential for growth in the absence and presence of drug treatment. This analysis found clear distinction in TF essentiality between the untreated and drugtreated PRIME models and revealed that drug doses largely group together (Fig. 4a). Interestingly, the TF essentiality profiles of rifampicin (a transcription inhibitor) were dose-dependent; the rifampicin profile at low-dose clustered separately, while the highdose profile clustered with linezolid (a protein synthesis inhibitor). The resemblance to linezolid at high-dose suggests that a secondary effect of strong rifampicin-induced transcription inhibition also impacts translation. Furthermore, we observed that the TF essentiality profiles of isoniazid (inhibitor of cell wall synthesis) were quite distinct to the other six drugs. In fact, 58 TFs become conditionally essential in the presence of isoniazid because of their mechanistic role in regulating 569 metabolic reactions required for supporting growth during isoniazid treatment. This highlights the multitude of regulatory-metabolic networks associated with cell wall disruption in Mtb and the extreme vulnerability in cell wall metabolism.

Focusing on isoniazid (INH), we evaluated the accuracy of these predictions against experimentally-determined fitness values from a genome-wide TnSeq screen performed in the presence of a subinhibitory concentration of INH 34 . Notwithstanding the difference in dosage of drug treatment of the input transcriptome data used in the PRIME model (0.18, 1.8 µg/mL) and in the TnSeq fitness screen (27 ng/mL), the LOOCV analysis demonstrated high sensitivity and specificity of PRIME predictions of gene essentiality (max ROC AUC = 0.685), significantly outperforming PROM (max ROC AUC = 0.625) and IDREAM (ROC AUC = 0.6) (Fig. 4b). We also used PRIME to rank order TFs based on their relative importance in

^aThe PROM model was updated in this study by incorporating the latest MN model for Mtb.

^bThe IDREAM model was constructed for Mtb in this study to evaluate performance relative to the other methods.

^cPlease see Supplementary Fig. 4 for the differences in the EGRIN models derived using Inferelator.

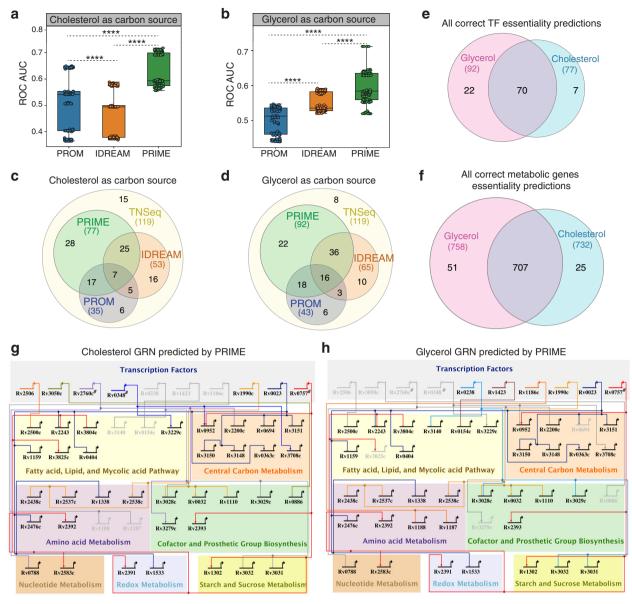
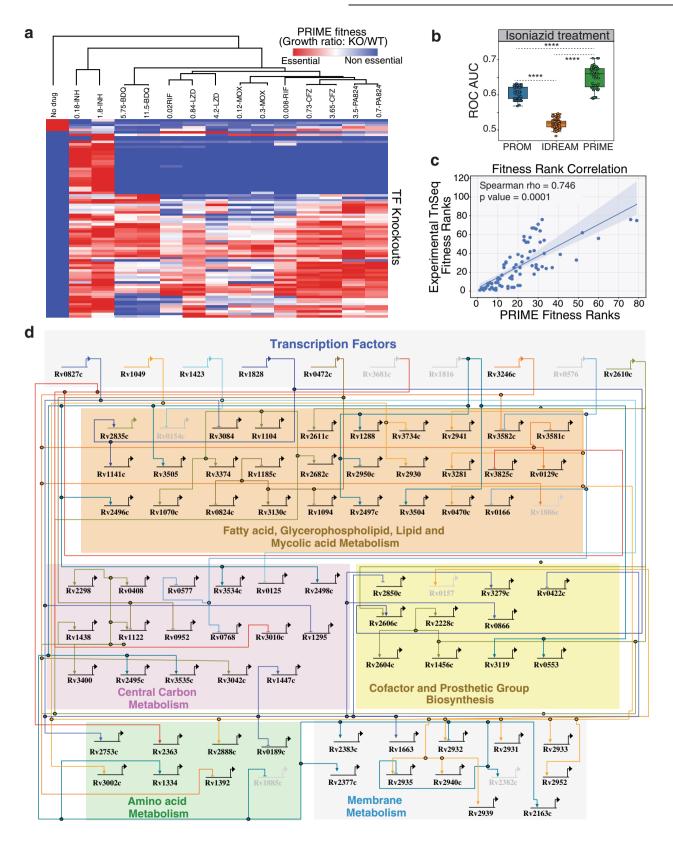


Fig. 3 Validation of PRIME predictions of conditional gene essentiality. Sensitivity and specificity of PRIME, PROM, and IDREAM predicted TF essentiality in **a**. cholesterol and **b** glycerol as determined by LOOCV analysis for the area under the receiver operating characteristic curve (ROC AUC). All the boxes in the boxplot indicate the upper and lower quartiles of the data and the middle line is the median with the whiskers extending to 1.5× interquartile range. Statistical significance was calculated as *p*-value with two sample *t*-test. ******p*-value < 0.0001. Comparison of all positive predictions (true positives and true negatives) for TF essentiality by PRIME, PROM, and IDREAM in **c** cholesterol and **d** glycerol. **e** The number of all correct PRIME predictions (true positives and true negatives) of TF knockouts across the two conditions (glycerol and cholesterol) that are validated by experimental TnSeq data. **f** The number of all correct PRIME predictions for deletion of all genes in the metabolic network across the two conditions that are validated by experimental TnSeq data. BioTapestry visualization showing a subset of the gene regulatory network of Mtb under growth in **g** cholesterol and **h** glycerol. TFs are grouped together in the top panel (represented by bent arrows), which extend to horizontal and vertical lines that connect to their regulatory gene targets. Highlighted TFs were predicted by the PRIME model to be essential and validated through TnSeq dataset in relevant conditions.

supporting growth in the presence of INH, and compared these ranks to TnSeq determined importance of TFs. There was striking correlation (Spearman's rho = 0.695; p-value = 0.0001) in the rank ordering of TFs based on the predicted (PRIME) and observed (TnSeq) magnitude of growth inhibition of Mtb in the presence of INH upon knocking out each TF one-at-a-time (Supplementary Fig. 9). The correlation increased (Spearman's rho = 0.746, p-value = 0.0001) when only TFs implicated by EGRIN as regulators of essential metabolic reactions were considered in this analysis, demonstrating the remarkable accuracy of PRIME in capturing how the differential regulation by TFs modulates flux through

essential metabolic reactions to manifest at a phenotypic level (Fig. 4c). Notably, PRIME accurately predicted that knocking out the top ten TFs one-at-a-time would result in at least 65% and up to 95% Mtb growth inhibition during INH treatment, but not in the absence of drug treatment, implicating these as conditional vulnerabilities for significantly potentiating INH treatment (Supplementary Data 6).

To aid in the interpretation of PRIME predictions, we developed the PRIME pathway analysis (PPA) tool to uncover in a single-step the specific metabolic reaction(s) regulated by a TF that make it essential for growth in a given environmental condition. Given a



TF, PPA identifies all reactions catalyzed by the genes it is predicted to regulate, rank orders the target genes based on the relative contribution of their gene product in driving flux towards biomass accumulation, and outputs a TF-metabolic gene-reaction map as a putative mechanism by which the TF is likely to be essential in a given environmental context. Using PPA, we

identified the specific metabolic reactions that were mechanistically responsible for the conditional essentiality of 23 TFs validated by TnSeq data³⁴ to be essential in the presence of INH. For example, we discovered the mechanisms underlying the essentiality of Rv0827c, Rv1049, Rv1423, Rv1828, and Rv0472c for growth in the presence of INH (Fig. 4d). Altogether, PPA



Fig. 4 Drug-specific predictions of PRIME. a Heatmap of PRIME derived fitness for all TF knockouts in the presence of seven primary drugs and control at 24 h. The numbers indicate the concentration of drug used in μg/mL. INH isoniazid, BDQ bedaquiline, RIF rifampicin, LZD linezolid, MOX moxifloxacin, CFZ clofazamine, PA824 pretomanid. **b** Sensitivity and specificity of PRIME, PROM, and IDREAM predicted TF essentiality in the presence of INH as determined by LOOCV analysis for the area under the receiver operating characteristic curve (ROC AUC). Statistical significance was calculated as *p*-value with two-sample *t*-test. ****: *p*-value < 0.0001. All the boxes in the boxplot indicate the upper and lower quartiles of the data and the middle line is the median with the whiskers extending to 1.5× interquartile range. **c** Correlation of TnSeq experimental fitness ranking of TFs and PRIME derived fitness ranks. **d** BioTapestry visualization showing a subset of the gene regulatory network of Mtb with PRIME predictions during INH treatment. Some of the highlighted TFs were predicted as essential in the presence of INH (Rv0827c, Rv1049, and Rv0472c), while others were predicted essential in both the absence and presence of INH (Rv1423, Rv1828, Rv3246c, and Rv2610c). The lightened TFs were predicted essential in the untreated control but non-essential in the presence of INH (Rv3681c, Rv1816, and Rv0576). All of these PRIME predictions were validated by experimental fitness screening in relevant conditions.

uncovered that 58 of the 142 TFs were conditionally essential for growth on INH because they conditionally regulate 569 key reactions across 55 pathways, including 84 reactions within fatty acid metabolism and mycolic acid biosynthesis (target of INH). In so doing, PRIME has provided the most comprehensive systems level perspective into strategies to potentiate INH killing by targeting TFs that mediate Mtb's metabolic response to INH treatment.

DISCUSSION

We have demonstrated that by incorporating how TFs act contextually in combinatorial schemes to regulate gene expression, PRIME outperformed PROM, IDREAM and IDREAM hybrid in accurately predicting how transcriptional regulation redirects metabolic flux to manifest in environment-specific phenotypes of Mtb. The shortcoming of PROM can be attributed to its reliance on P-D interactions for regulatory network, which are plagued with false positive interactions (because overexpression of TFs can force non-functional binding across the genome) and false negative interactions because of lack of appropriate context (e.g., missing co-factors). Hence, a P-D interaction does not capture whether a TF is regulating a gene in a given condition, which is better modeled by regulatory influences inferred using regression analysis of transcript level changes in TFs and all genes across the genome. However, despite incorporating regulatory influences from EGRIN, IDREAM performance was inferior compared to PRIME, and in fact its performance in predicting gene essentiality in cholesterol and INH was worse than PROM. We also demonstrated that IDREAM failed to outperform PRIME even when FDR thresholds were relaxed (Supplementary Fig. 10). One explanation could be that relative to the number of P-D interactions used in PROM, IDREAM used nearly twice as many EGRIN-based regulatory influences that were inferred from a wide range of environmental contexts, without taking into account combinatorial regulatory schemes, weights of regulatory influences, or the absolute expression levels of TFs to prune regulatory edges that were not relevant for a given environmental context. Hence, reliance on a P-D interaction map, and even just the likelihood that a TF might regulate a gene based on regression analysis are both insufficient to capture the complex environmentdependent interplay of transcription and metabolism. Altogether, these comparative analyses have demonstrated that four key advancements in PRIME addressed the shortcomings of PROM and IDREAM: (i) PRIME took full advantage of EGRIN predictions to incorporate weights of TF regulatory influence on each gene; (ii) PRIME calibrated the relative influence of each TF on a given metabolic gene by accounting for all TFs that were also implicated in the regulation of that gene; (iii) PRIME accounted for regulation of multiple genes (that encode enzymes) for the same reaction by considering which gene(s) contributed maximally towards flux through that reaction in a given environmental context; and, finally (iv) PRIME considered the absolute expression level of each TF to evaluate the degree to which each regulatory influence was active in a given environment. We posit that PRIME performance could potentially improve even further upon uncovering and incorporating roles in metabolism for genes of unknown function, incorporating novel regulatory mechanisms (including post-transcriptional and post-translational regulatory mechanisms), and accounting for indirect influences of TFs on metabolism through the regulation of other TFs.

By demonstrating better accuracy in predicting environmentspecific phenotypes of Mtb using EGRIN, PRIME overrides the need for a condition-specific physical map of P-D interactions, which is difficult to generate for many organisms, across all environments of interest, and especially in some contexts, such as within infected tissue. In fact, the incompleteness of the P-D interaction map was demonstrated by the significant drop in the performance of PRIME upon excluding regulatory influences that were not supported by physical TF-gene interactions (i.e., EGRIN-PD). By contrast, EGRIN is inferred directly from a compendium of transcriptomes, which can be profiled across relevant environmental conditions with minimal manipulation (e.g., without overexpression of TFs) and even within infected cells using technologies like Path-seq³⁵. As a consequence, EGRIN discovers a significantly larger number of novel regulatory mechanisms, including the combinatorial schemes and specific environmental contexts in which they are conditionally active. This explains why PRIME discovered mechanisms that become conditionally essential in the presence of INH, but also accurately predicted the relative importance of each TF for enhancing the potency of INH. Based on this observation, we posit that PRIME will be especially valuable to prioritize genes that represent novel contextdependent vulnerabilities that could be targeted to potentiate the action of any antibiotic and achieve faster clearance with a lower dosage. By enabling the in-silico discovery of vulnerabilities within the Mtb network, PRIME also overrides the need for large scale transposon mutagenesis-based experiments (e.g., TnSeg, TraSH, HITS, etc), which are resource-intensive and difficult to perform across all conditions relevant to the lifecycle of Mtb. Instead, PRIME can be used to rank prioritize the strains and contexts in which to assay for an expected phenotype. This capability is particularly powerful considering the numerous mechanisms by which Mtb can be phenotypically different, with different antibiotic sensitivities. Additionally, there is growing evidence that upon gaining resistance to an antibiotic, the regulatory and metabolic networks within a pathogen are remodeled in order to reallocate resources for supporting the new phenotype³⁶. Using PRIME, we can delineate novel vulnerabilities within these remodeled regulatory and metabolic networks to devise strategies for rationally disrupting the antibiotic resistance phenotype with a second drug.

PRIME will also be useful in biotechnology applications to further optimize the production of desired end products by rewiring the regulatory networks of metabolically engineered strains. Advancements in metabolic engineering have been effective in substantially increasing flux towards the production of a desired metabolite^{18–20,37} but there is a limit to which metabolic engineering alone can improve the overall yield. It has been proposed that further enhancements in yield would require

reprogramming of the regulatory network to control when genes of the engineered pathways are expressed, and to rationally up and down regulate competing metabolic pathways to maximize flux and resource allocation towards the desired objective. Hence, by using PRIME, metabolic engineering of high-yielding strain phenotypes can be identified. Although the capabilities of PRIME are elucidated extensively using Mtb as a model system in this study, we foresee the use and applications of PRIME in various organisms due to its scalability.

METHODS

Construction of EGRIN gene regulatory network for Mycobacterium tuberculosis

The Mtb EGRIN used in this study was constructed using the Inferelator algorithm^{15,22,23} trained on a transcriptional compendium for Mtb with 3902 genes across 664 experimental conditions (downloaded from the COLOMBOS database) and an experimentally supported signed Mtb P-D network (generated as previously described in ref. 35). The original transcriptional compendium contained a larger number of genes and conditions but was modified to remove genes and conditions with missing values. Briefly, we used the Inferelator to identify potential transcriptional regulators for the 3902 Mtb genes in the expression compendium, as previously performed for other species^{23,38}. The Inferelator first estimates the regulatory activities of each transcription factor activity (TFA) using the expression profile of TF known targets (encoded in the signed P-D network). Then, the Inferelator uses a Bayesian Best Subset Regression to estimate the magnitude and sign (activation or repression) of potential interactions between TFs and genes. As before, we bootstrapped the expression data (20 times) to avoid regression overfitting. The Inferelator generates two scores for each TF-gene interaction, the corresponding regression coefficient (weight— β) and a confidence score. The second score indicates the likelihood of the interaction. The final set of TF-gene interactions was defined with a 0.5 precision cutoff. This means that 50% of all interactions in the inferred network were already present in the signed P-D network used for training, while the other half corresponded to putative novel TF-gene interactions.

Development of PRIME

The PRIME algorithm has been developed by integrating weights (β) from EGRIN with metabolic network (MN) models for phenotype prediction in a context-specific manner (wiring diagram in Fig. 1). PRIME requires 1) a MN in the format of constraints-based model^{39,40} in systems biology markup language (SBML), an XML format as input, that are represented in silico in the form of a stoichiometric matrix, wherein every column corresponds to a reaction and every row corresponds to a metabolite. These constraintsbased models are used to integrate the regulatory influences by updating the reaction flux, 2) a regulatory network containing TF and gene interactions (one array of regulators and one array of corresponding gene targets), 3) magnitude/weights (β) of regulatory influences for each of the interactions (array of magnitudes) derived from Inferelator and 4) the gene expression data profiled under a specific condition (gene ids and their expression, provided as ratio to the control—in case of environmentspecific predictions the ratio between initial t0 and final time point tn). The pipeline of PRIME initially links each metabolic gene in MN to its associated regulators considering the combinatorial effects, followed by applying the calculated reaction influence factor (RIF, y). Specifically, we have introduced a new way to calculate the RIF (y), a value that quantitatively constrains the reaction flux constraint space. The Eqs. 1-5 consists of the details involved in each successive step within the algorithm.

Given a TF j influencing a metabolic gene i of reaction w, we define RIF $(\gamma_{i,w})$ as,

$$\gamma_{i,w} = 1 - \left(\frac{\beta_{i,j}}{\sum_{j \in J} \beta_{i,j}} * X_j'\right),\tag{1}$$

where J is the subset of all TFs that influence gene i and X_j' is the scaled expression of a TF j in a particular condition c of a coherent environmental context B represented as

$$X_j' = \frac{X_{j,c} - \min X_j(B)}{\max X_j(B) - \min X_j(B)}.$$
 (2)

Then, in case of reactions where there are two or more gene products (I) involved in catalyzing the reaction w, the regulatory influence that exerts the larger effect on that reaction w across the set of metabolic genes $i \in I$ has been identified as,

$$g_{\rm w}=\min \gamma_{\rm i,w}$$
 (3)

Else, $g_w = \gamma_{i,w}$ when a single gene i is involved in the reaction w. At this point, it is straightforward to incorporate calculated weights as new upper bounds. The lower bound (a) and upper bound (b) used here were computed via flux variability analysis, which gives as output the flux span that can satisfy the constraints for a specific condition.

$$b^{\mathsf{PRIME}} = b \circ g = (b)_{\mathsf{w}}(g)_{\mathsf{w}} \tag{4}$$

This b^{PRIME} is the new upper bound calculated using PRIME to satisfy the flux balance analysis (FBA)⁴⁰ formalism, assuming steady state metabolic concentrations, and defining the system mass balance as S.v=0, to maximize the objective function $Z=c^{\mathsf{T}}v$ such that fluxes are within the new boundary conditions,

$$a \le \mathbf{v} \le b^{\mathsf{PRIME}}$$
 (5)

All variables and parameters used in the PRIME equations are listed in Supplementary Table 1. The objectives in each prediction are defined during FBA optimization. The phenotype predictions mentioned in this study are the optimized biomass predicted by FBA. The complete PRIME algorithm package and details of the required input dataset is available for download from our GitHub Repository (https://github.com/baliga-lab/PRIME). All model simulations related to FBA were performed on MATLAB_R2019a platform using the recent version of COBRA⁴¹ -The COnstraint-Based Reconstruction and Analysis toolbox. In silico gene essentiality predictions were performed using the COBRA toolbox "single-gene-deletion" function in MATLAB.

Incorporating drug treatment gene expression data on metabolic model

The iEK1011²⁶ model of Mtb was used for all the predictions in this study. For drug-specific models, we applied the gene expression data from both drug-treated and untreated control experiments using the GIMME⁴² algorithm on the iEK1011 MN model. This step was carried out to constrain the MN model to the specific condition being tested. We used GIMME because of the flexibility in defining objective function during implementation. The GIMME algorithm is implemented in the MATLAB_R2019a platform, using the "GIMME.m function" in the COBRA Toolbox after processing the gene expression data through "mapExpressionToReactions. m" function to convert the gene expression values as inputs to GIMME.

PROM models

For developing PROM models, we followed the PROM approach 11,12 to estimate the probability that a target gene is "ON" or "OFF" in the absence of the TF i.e., in the event of a TF knockout. This was calculated from a gene expression dataset as, Probability, P (Gene = 1|TF = 0) or P (TF = 1|Gene = 0). The gene expression threshold that delineated between the "ON" and "OFF" states was set as quantile (0.33) from the input expression data. These probabilities were then used to constrain the maximal fluxes of the reactions catalyzed by the gene products in the metabolic model as $p \times Vmax$, where p is the probability of the gene being on. The user defined "kappa" value was used as similar to earlier PROM models 11. All PROM predictions and simulations were performed using PROM.m (MATLAB script) on the MATLAB_R2019a platform. We used iEK1011 metabolic network model in XML format as input in the PROM. The P–D derived regulatory network was obtained from the study 21 , similar to the MTBPROMV2.0 12.

IDREAM and IDREAM-hybrid models

For IDREAM and IDREAM-hybrid models, the GRN derived using Inferelator, was integrated with the PROM pipeline as it had been done previously for the yeast system¹³. We ran 200 iterations in Inferelator to calculate the FDR for all predictions. For each gene, we estimated an FDR for each TF by counting the fraction of models that identified that factor as a regulator. Thus, if TF1 was predicted to regulate gene1 in 191 of 200 models, then te TF-gene interaction identified would have an FDR = 0.045. We included only those interactions that passed an FDR cutoff of 0.05 and 0.25. We used Inferelator-derived GRN to integrate it with iEK1011 metabolic network model of Mtb using the PROM framework. The user



defined "kappa" value was used as similar to earlier PROM models¹¹. IDREAM does not rely on probabilities, hence the gene expression dataset was not used in IDREAM instead 'prob_prior' in the PROM function was set based on the EGRIN FDR values for each TF-gene interaction. If the TF is an activator of a gene, we use the FDR value directly, if it is an inhibitor, we use 1-FDR value as "prob_prior". For IDREAM-hybrid models, the conditional probabilities calculated using gene expression dataset were used for the additional indirect interactions. EGRIN network was derived using Inferelator in R (Inferelator.pkg.R) and PROM predictions and simulations were performed using PROM.m (MATLAB script) on the MATLAB_R2019a platform as similar to PROM model development.

Performance assessment of PRIME predictions

The predictive power of PRIME as a binary classifier (essential or nonessential) between the model predicted gene essentiality and experimentally defined gene essentiality (TnSeq) has been performed using the ROC curve. A gene was considered "essential" if its deletion reduced the biomass by >85%. By this analysis, the model classified each gene as "essential" or "non-essential". We compared the gene essentiality predictions from Mtb grown under glycerol and cholesterol as carbon source with the available experimental TnSeq data²⁸ and deduced the confusion matrix to derive true positive rates (TPR) and false positive rates (FPR). We also took advantage of the follow-up study where Bayesian analysis was used to assign calls as essential and non-essential for the same TnSeq dataset²⁹. We expanded the analysis of TnSeq data to classify essential and non-essential with a cutoff value of using cholesterol/glycerol ratio of 0.6 in order to assign calls for all the genes. This classification led to the elucidation of sensitivity and specificity of the model using ROC curve analysis. Briefly, the gene expression data of Mtb profiled under growth on Glycerol (GSE52020) and Cholesterol (GSE13978) were used to generate condition-specific metabolic networks using GIMME. PRIME was applied on these models to predict gene and TF essentialities according to the condition tested. These predictions were then compared to the TnSeq data²⁸. A similar sensitivity and specificity analysis was performed while validating the performance of PRIME for INH-specific predictions using experimentally derived TnSeq data³⁴. To construct the INH-specific metabolic models, we used INH-treated Mtb transcriptome sequencing (RNA-Seq) data generated in this study (see below).

PRIME Pathway Analysis (PPA) pipeline

The PRIME pathway analysis (PPA) pipeline was developed to derive the metabolic association of a specified TF in a simple process by accessing PRIME model genes and their interactions. The top ranked TFs and their associated metabolic genes are further linked to their metabolic processes using the PPA pipeline. PPA is provided as PRIMEanalysis.m (MATLAB script). All analyses related to PPA were performed in MATLAB_R2019a platform. The illustration of PPA-derived essential gene regulatory-metabolic networks were deduced using BioTapestry tool (http://www.biotapestry.org/).

Drug treatment culturing conditions

Experiments were performed using *Mycobacterium tuberculosis* H37Rv grown with mild agitation at 37 °C in standard 7H9-rich media consisting of Middlebrook 7H9 broth supplemented with 10% Middlebrook ADC, 0.05% Tween-80, and 0.2% glycerol. Frozen 1 mL stocks of Mtb cells were added to 7H9-rich medium and grown until the culture reached an OD₆₀₀ of $\sim\!0.4$ –0.8. The cells were then diluted to OD₆₀₀ of 0.05 and added to 7H9-rich medium containing drugs at the predetermined amounts. Samples, in biological triplicate, were collected at 24 h after drug treatment biological triplicate, were collected at 24 h after drug treatment and immediately flash freezing the cell pellet in liquid nitrogen. Cell pellets were stored at -80 °C until RNA extraction was performed.

RNA extraction

Cell pellets stored at $-80\,^{\circ}\text{C}$ were resuspended in $600\,\mu\text{L}$ of fresh lysozyme solution in Tris-EDTA buffer pH 8.0 (5 mg per mL). The resuspended cells were transferred to a tube containing Lysing Matrix B (MP Biomedicals, Santa Ana, CA) and incubated at 37 °C for 30 min. Following incubation, $60\,\mu\text{L}$ (1/10th volume of lysate volume) of 10% SDS was added and thrubes were vigorously shaken at maximum speed for 30 s in a FastPrep 120 homogenizer (MP Biomedicals) three times. Tubes were centrifuged for 1 min (maximum speed), then $66\,\mu\text{L}$ of 3 M sodium acetate pH 5.2 added

and mixed well. Acid phenol (pH 4.2) was added at 726 μ L and tubes were inverted to mix well (~60 times). Samples were incubated at 65 °C for 5 min, inverting tubes to mix samples every 30 s. Then, centrifuged at 14,000 rpm for 5 min and the upper aqueous phase was transferred to a new tube. 3 M sodium acetate (pH 5.2) was added at 1/10th volume along with 3x volumes of 100% ethanol. Sample was mixed well and incubated at -20 °C for 1 hr or overnight. Following incubation, samples were centrifuged at 14,000 rpm for 30 min at 4 °C, ethanol was discarded and 500 μ L of 70% ethanol was added. Samples were centrifuged again at 14,000 rpm for 10 min at 4 °C, supernatant discarded, and any residual ethanol removed using pipet. Pellet was allowed to air dry, resuspended in 30–40 μ L of RNase free water and quantified by Nanodrop (Thermo Scientific). This was followed by in solution genomic DNA digestion using RQ1 DNase (Promega) following manufacturer's recommendation.

Processing and analysis of RNA-Seq data

Sample collection and RNA-extraction was performed as described above. Total RNA samples were depleted of ribosomal RNA using the Ribo-Zero Bacteria rRNA Removal Kit (Illumina, San Diego, CA). Quality and purity of mRNA samples was determined with 2100 Bioanalyzer (Agilent, Santa Clara, CA). Samples were prepared with TrueSeq Stranded mRNA HT library preparation kit (Illumina, San Diego, CA). All samples were sequenced on the NextSeq sequencing instrument in a high output 150 v2 flow cell. Pairedend 75 bp reads were checked for technical artifacts using Illumina default quality filtering steps. Raw FASTQ read data were processed using the R package DuffyNGS⁴³. Briefly, raw reads were passed through a 2-stage alignment pipeline: (i) a pre-alignment stage to filter out unwanted transcripts, such as rRNA; and (ii) a main genomic alignment stage against the genome of interest. Reads were aligned to M. tuberculosis H37Rv (ASM19595v2) with Bowtie2⁴⁴, using the command line option "verysensitive." BAM files from stage (ii) were converted into read depth wiggle tracks that recorded both uniquely mapped and multiply mapped reads to each of the forward and reverse strands of the genome(s) at singlenucleotide resolution. Gene transcript abundance was then measured by summing total reads landing inside annotated gene boundaries, expressed as both RPKM and raw read counts. We used the raw read counts as input for DESeq2⁴⁵ to obtain DESeq2 normalized counts. The RNA-Seq data of Mtb response to drug exposure generated for this study are publicly available at the Gene Expression Omnibus under accession number GSE165673.

DATA AVAILABILITY

Input files for PRIME used in this study are provided as Supplementary Data 7. All PRIME-generated data are provided as supplementary materials. The RNA-Seq data generated for this study are available in the Gene Expression Omnibus under accession no. GSE165673.

CODE AVAILABILITY

PRIME code, with data and description for implementation, is available in GitHub repository: https://github.com/baliga-lab/PRIME.

Received: 5 April 2021; Accepted: 2 November 2021; Published online: 06 December 2021

REFERENCES

- Stanley, S. A. & Cox, J. S. Host-pathogen interactions during Mycobacterium tuberculosis infections. Curr. Top. Microbiol. Immunol. 410, 211–241 (2013).
- Rienksma, R. A., Schaap, P. J., Martins dos Santos, V. A. P. & Suarez-Diez, M. Modeling the metabolic state of *Mycobacterium tuberculosis* upon infection. *Front. Cell. Infect. Microbiol.* 8, 1–13 (2018).
- Galagan, J. E. et al. The Mycobacterium tuberculosis regulatory network and hypoxia. Nature 499, 178–183 (2013).
- Chaulk, C. P. & Kazandjian, V. A. Directly observed therapy for treatment completion of pulmonary tuberculosis: Consensus statement of the public health tuberculosis guidelines panel. *J. Am. Med. Assoc.* 279, 943–948 (1998).
- Sarathy, J. P. et al. Extreme drug tolerance of Mycobacterium tuberculosis in Caseum. Antimicrob. Agents Chemother. 62, e02266–17 (2018).
- Sarathy, J., Dartois, V., Dick, T. & Gengenbacher, M. Reduced drug uptake in phenotypically resistant nutrient-starved nonreplicating *Mycobacterium tuberculosis*. Antimicrob. Agents Chemother. 57, 1648–1653 (2013).

npj

- de Steenwinkel, J. E. M. et al. Time-kill kinetics of anti-tuberculosis drugs, and emergence of resistance, in relation to metabolic activity of *Mycobacterium* tuberculosis. J. Antimicrob. Chemother. 65, 2582–2589 (2010).
- Rao, S. P. S., Alonso, S., Rand, L., Dick, T. & Pethe, K. The protonmotive force is required for maintaining ATP homeostasis and viability of hypoxic, nonreplicating Mycobacterium tuberculosis. Proc. Natl Acad. Sci. USA 105, 11945–11950 (2008).
- 9. Peterson, E. J. R. et al. Intricate genetic programs controlling dormancy in *Mycobacterium tuberculosis. Cell Rep.* **31**, 107577 (2020).
- Covert, M. W., Schilling, C. H. & Palsson, B. Regulation of gene expression in flux balance models of metabolism. J. Theor. Biol. 213, 73–88 (2001).
- Chandrasekaran, S. & Price, N. D. Probabilistic integrative modeling of genomescale metabolic and regulatory networks in *Escherichia coli* and *Mycobacterium* tuberculosis. Proc. Natl Acad. Sci. 107, 17845–17850 (2010).
- 12. Ma, S. et al. Integrated modeling of gene regulatory and metabolic networks in *Mycobacterium tuberculosis*. *PLoS Comput. Biol.* **11**, e1004543 (2015).
- Wang, Z. et al. Combining inferred regulatory and reconstructed metabolic networks enhances phenotype prediction in yeast. PLOS Comput. Biol. 13, e1005489 (2017).
- Peterson, E. J. R., Ma, S., Sherman, D. R. & Baliga, N. S. Network analysis identifies Rv0324 and Rv0880 as regulators of bedaquiline tolerance in *Mycobacterium* tuberculosis. Nat. Microbiol. 1, 16078 (2016).
- Bonneau, R. et al. The Inferelator: an algorithm for learning parsimonious regulatory networks from systems-biology data sets de novo. Genome Biol. 7, R36 (2006).
- Brooks, A. N. et al. A system-level model for the microbial regulatory genome. Mol. Syst. Biol. 10, 740 (2014).
- Peterson, E. J. R. et al. A high-resolution network model for global gene regulation in Mycobacterium tuberculosis. Nucleic Acids Res. 42, 11291–11303 (2014).
- Kim, J. & Reed, J. L. OptORF: optimal metabolic and regulatory perturbations for metabolic engineering of microbial strains. BMC Syst. Biol. 4, 53 (2010).
- Yang, L., Cluett, W. R. & Mahadevan, R. EMILiO: a fast algorithm for genome-scale strain design. *Metab. Eng.* 13, 272–281 (2011).
- Kim, M., Sun, G., Lee, D.-Y. & Kim, B.-G. BeReTa: a systematic method for identifying target transcriptional regulators to enhance microbial production of chemicals. *Bioinformatics* 33, 87–94 (2017).
- Minch, K. J. et al. The DNA-binding network of Mycobacterium tuberculosis. Nat. Commun. 6, 1–10 (2015).
- Greenfield, A., Hafemeister, C. & Bonneau, R. Robust data-driven incorporation of prior knowledge into the inference of dynamic regulatory networks. *Bioinformatics* 29, 1060–1067 (2013).
- Arrieta-Ortiz, M. L. et al. An experimentally supported model of the Bacillus subtilis global transcriptional regulatory network. Mol. Syst. Biol. 11, 839 (2015).
- Moretto, M. et al. COLOMBOS v3.0: leveraging gene expression compendia for cross-species analyses. Nucleic Acids Res. 44, D620–D623 (2016).
- Rustad, T. R. et al. Mapping and manipulating the Mycobacterium tuberculosis transcriptome using a transcription factor overexpression-derived regulatory network. Genome Biol. 15, 502 (2014).
- Kavvas, E. S. et al. Updated and standardized genome-scale reconstruction of *Mycobacterium tuberculosis* H37Rv, iEK1011, simulates flux states indicative of physiological conditions. *BMC Syst. Biol.* 12, 25 (2018).
- 27. Webb, G. I. et al. Encyclopedia of Machine Learning 600–601 (Springer, 2011).
- Griffin, J. E. et al. High-resolution phenotypic profiling defines genes essential for mycobacterial growth and cholesterol catabolism. PLoS Pathog. 7, e1002251 (2011).
- DeJesus, M. A. et al. Bayesian analysis of gene essentiality based on sequencing of transposon insertion libraries. *Bioinformatics* 29, 695–703 (2013).
- Balhana, R. J. C. et al. bkaR is a TetR-type repressor that controls an operon associated with branched-chain keto-acid metabolism in Mycobacteria. FEMS Microbiol. Lett. 345, 132–140 (2013).
- Lee, W., VanderVen, B. C., Fahey, R. J. & Russell, D. G. Intracellular Mycobacterium tuberculosis exploits host-derived fatty acids to limit metabolic stress. J. Biol. Chem. 288, 6788–6800 (2013).
- Martinot, A. J. et al. Mycobacterial metabolic syndrome: LprG and Rv1410 regulate triacylglyceride levels, growth rate and virulence in *Mycobacterium tuber*culosis. PLoS Pathog. 12, e1005351 (2016).
- Pandey, A. K. & Sassetti, C. M. Mycobacterial persistence requires the utilization of host cholesterol. *Proc. Natl Acad. Sci. USA* 105, 4376–4380 (2008).
- Xu, W. et al. Chemical genetic interaction profiling reveals determinants of intrinsic antibiotic resistance in *Mycobacterium tuberculosis*. *Antimicrob. Agents Chemother.* 61, 1–15 (2017).
- Peterson, E. J. et al. Path-seq identifies an essential mycolate remodeling program for mycobacterial host adaptation. Mol. Syst. Biol. 15, e8584 (2019).
- Arrieta-Ortiz, M. L. et al. Disrupting the ArcA regulatory network increases tetracycline susceptibility of Tet R *Escherichia coli*. Preprint at https://www.biorxiv.org/content/10.1101/2020.08.31.275693v1 (2020).
- Burgard, A. P., Pharkya, P. & Maranas, C. D. Optknock: A bilevel programming framework for identifying gene knockout strategies for microbial strain optimization. *Biotechnol. Bioeng.* 84, 647–657 (2003).

- Arrieta-Ortiz, M. L. et al. Inference of bacterial small RNA regulatory networks and integration with transcription factor-driven regulatory networks. mSystems 5, e00057–20 (2020).
- Varma, A. & Palsson, B. O. Metabolic capabilities of *Escherichia coli*: I. Synthesis of biosynthetic precursors and cofactors. *J. Theor. Biol.* 165, 477–502 (1993).
- Orth, J. D., Thiele, I. & Palsson, B. Ø. What is flux balance analysis? *Nat. Biotechnol.* 28, 245–248 (2010).
- 41. Heirendt, L. et al. Creation and analysis of biochemical constraint-based models using the COBRA Toolbox v.3.0. *Nat. Protoc.* **14**, 639–702 (2019).
- Becker, S. A. & Palsson, B. O. Context-specific metabolic networks are consistent with experiments. *PLoS Comput. Biol.* 4, e1000082 (2008).
- 43. Vignali, M. et al. NSR-seq transcriptional profiling enables identification of a gene signature of Plasmodium falciparum parasites infecting children. *J. Clin. Invest.* **121**. 1119–1129 (2011).
- 44. Langmead, B. & Salzberg, S. L. Fast gapped-read alignment with Bowtie 2. *Nat. Methods* **9**, 357–359 (2012).
- Love, M. I., Huber, W. & Anders, S. Moderated estimation of fold change and dispersion for RNA-seq data with DESeq2. Genome Biol. 15, 550 (2014).

ACKNOWLEDGEMENTS

We thank members of the Baliga lab for critical discussions and feedback. We thank Vivek Srinivas for his suggestions related to PRIME code. This work is funded by Bill and Melinda Gates Foundation (INV-009322), the National Institute of Allergy and Infectious Diseases of the National Institutes of Health (R01Al141953, R01Al128215 and U19Al135976) and the National Science Foundation (NSF IIBR RoL 2042948).

AUTHOR CONTRIBUTIONS

N.S.B. and S.R.C.I. conceptualized the study and designed the research. S.R.C.I. developed the PRIME algorithm, performed all the computational analyses related to PRIME and metabolic modeling, analyzed all data represented and designed all figures. M.L.A.O. performed computational analyses related to regulatory networks and contributed to PRIME conceptualization and performance assessment. R.R. and M.P. performed the experiments related to drug treatment and transcriptomic profiling. A.L.G.L. contributed in PRIME method conceptualization in early stages. E.J. R.P. designed and analyzed transcriptome studies, and analyzed fitness data of PRIME analysis. N.S.B. and E.J.R.P. provided overall supervision. S.R.C.I., E.J.R.P., and N.S.B. wrote the manuscript. All authors read the manuscript and approved its content.

COMPETING INTERESTS

The authors declare no competing interests.

ADDITIONAL INFORMATION

Supplementary information The online version contains supplementary material available at https://doi.org/10.1038/s41540-021-00205-6.

Correspondence and requests for materials should be addressed to Eliza J. R. Peterson or Nitin S. Baliga.

Reprints and permission information is available at http://www.nature.com/reprints

Publisher's note Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.



Open Access This article is licensed under a Creative Commons Attribution 4.0 International License, which permits use, sharing,

adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons license, and indicate if changes were made. The images or other third party material in this article are included in the article's Creative Commons license, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons license and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this license, visit http://creativecommons.org/licenses/by/4.0/.

© The Author(s) 2021