

1 **Quantitative prediction of conditional vulnerabilities in regulatory and**  
2 **metabolic networks of *Mycobacterium tuberculosis***

3

4 Selva Rupa Christinal Immanuel<sup>1</sup>, Mario L. Arrieta-Ortiz<sup>1</sup>, Rene A. Ruiz<sup>1</sup>, Min Pan<sup>1</sup>, Adrian Lopez Garcia  
5 de Lomana<sup>1</sup>, Eliza J. R. Peterson<sup>1,\*</sup>, Nitin S. Baliga<sup>1,2,3,4,\*</sup>

6

7 <sup>1</sup> Institute for Systems Biology, Seattle, WA, USA

8 <sup>2</sup> Departments of Biology and Microbiology, University of Washington, Seattle, WA, USA

9 <sup>3</sup> Molecular and Cellular Biology Program, University of Washington, Seattle, WA, USA

10 <sup>4</sup> Lawrence Berkeley National Lab, Berkeley, CA, USA

11 \* Corresponding author

12 Eliza J.R. Peterson: [eliza.peterson@isbscience.org](mailto:eliza.peterson@isbscience.org)

13 Nitin S. Baliga: [nitin.baliga@isbscience.org](mailto:nitin.baliga@isbscience.org) (lead contact)

14

15

16

17

18

19

20

21

22

23

24

25

26

27

28

29

30

31

32

33

34 Current affiliation for author: **ALGL** - Center for Systems Biology, University of Iceland, Reykjavik, Iceland.

35 **Abstract**

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

47

48

49

50

51

52

53

54

55

56

57

58

59

60

61

62

63

64

65

66

67

68

## 69 INTRODUCTION

70 *Mycobacterium tuberculosis* (Mtb) kills more people than any other microbe, and it has thus far  
71 resisted every attempt to bring the pandemic under control. Part of the pathogen's success is its ability  
72 to diversify itself phenotypically and survive both host and drug bactericidal action<sup>1-3</sup>. Phenotypic  
73 heterogeneity (both stochastically and environmentally induced) seems to be an intrinsic characteristic  
74 of the pathogen and a major reason why standard chemotherapy of tuberculosis (TB) requires 6 months  
75 of treatment, and 5% are not cured even then<sup>4,5</sup>. To develop better interventions that account for  
76 pathogen heterogeneity, we need to identify the most important factors (e.g., transcriptional regulators)  
77 that create variation as well as the downstream effectors (e.g., regulatory target genes) that mediate  
78 drug tolerance.

79 Metabolic activity undoubtedly contributes to Mtb phenotypic heterogeneity and antibiotic  
80 tolerance. For example, changes in metabolism can affect the amount of drug target present<sup>6</sup>, the ability  
81 to generate toxic products<sup>7</sup>, and the efflux of antibiotics<sup>8</sup>. Mtb alters its growth and metabolism in  
82 response to stressful conditions through regulatory programs primarily encoded at the transcriptional  
83 level. Indeed, modeling host-related stresses *in vitro* produces large transcriptional changes in Mtb,  
84 particularly in metabolic pathways; consistently ~25% of differentially expressed genes are metabolic  
85 genes from hypoxic (GSE116353)<sup>9</sup>, acidic pH (GSE165514), or nutrient limited (GSE165673) conditions.  
86 To develop effective antibiotic regimens, we need to understand at a systems- and mechanistic-level  
87 how specific regulatory mechanisms conditionally activate and repress genes to redirect flux through  
88 metabolic networks to generate and support drug tolerant phenotypes. This mechanistic understanding  
89 will uncover new vulnerabilities in Mtb's regulatory and metabolic networks that can be rationally targeted  
90 in new drug regimens to achieve faster and complete clearance of the pathogen.

91 Previously, approaches to model the influence of transcriptional regulation on metabolism have  
92 used boolean logic (Regulatory Flux Balance Analysis - rFBA)<sup>10</sup>, protein-DNA (P-D) interactions  
93 (Probabilistic Regulation of Metabolism - PROM)<sup>11,12</sup>, and regression-based regulatory influences  
94 (Integrated Deduced REgulation And Metabolism - IDREAM)<sup>13</sup> to predict how transcriptional regulation  
95 of enzyme-coding genes modulates flux through their catalyzed reactions. Briefly, rFBA models the  
96 influence of transcriptional regulation on metabolism using boolean "on or off" states of metabolic genes,  
97 depending on the expression level of the transcription factor (TF) and its implicated role as a putative  
98 activator or repressor of that gene. The extensive manual curation required to develop rFBA and its  
99 inability to model TF activity as a continuous (i.e., not boolean) function greatly limits its application and  
100 accuracy. In contrast, PROM outperformed<sup>11</sup> rFBA by using a probabilistic approach to model the  
101 regulation of a metabolic gene by a TF using a compendium of transcriptome profiles to calculate  
102 probabilities. However, PROM is limited in that it relies on a P-D interaction map for the regulatory

103 network. P-D interactions are typically generated in a limited set of conditions by using an overexpressed  
104 TF as a bait to enrich and locate its genome-wide binding locations. P-D interactions are fraught with  
105 false positives (due to TF overexpression) and false negatives (due to lack of context for TF regulation  
106 across environmental conditions). Notwithstanding these caveats, PROM was useful in uncovering the  
107 mechanism by which pretomanid potentiates bedaquiline action on Mtb by disrupting a regulatory  
108 network that confers tolerance to the recently approved FDA drug<sup>14</sup>. A third model, IDREAM addressed  
109 the shortcoming of using P-D interactions in PROM by constraining flux using TF regulatory influences  
110 from a predictive systems-scale environment and gene regulatory influence network (EGRIN) model.  
111 An EGRIN model is inferred in two steps using (a) cMonkey, which identifies the specific context in which  
112 subsets of genes are co-regulated (biclusters) by a conserved regulatory mechanism(s); and (b)  
113 Inferelator, which predicts TFs and environmental factors that causally influence the differential  
114 expression of genes within those biclusters<sup>15–17</sup>. By integrating confidence scores for EGRIN-inferred  
115 regulatory influences, IDREAM achieved significantly better performance than rFBA and PROM in  
116 predicting synthetic lethal interactions between TFs and metabolic genes in yeast<sup>13</sup>. However, IDREAM  
117 does not incorporate quantitative environment-specific TF regulatory influences that are modeled by  
118 EGRIN, and is therefore also limited in accurately predicting environment-specific consequences of TF  
119 perturbations. For the reasons stated above, PROM, rFBA, and IDREAM are limited in their ability to  
120 predict environment-specific phenotypic consequences of perturbations to TFs.

121 Additionally, there are algorithms (OptORF<sup>18</sup>, EMILiO<sup>19</sup> and BeReTa<sup>20</sup>) that have the potential to  
122 predict the consequence of regulatory and metabolic network perturbations. They were originally  
123 designed to identify perturbations that maximize flux towards a desired metabolite and some of their  
124 features make them not well-suited for predicting systems-wide conditional outcomes of TF perturbation.  
125 For instance, OptORF<sup>18</sup> and EMILiO<sup>19</sup> use binary or fixed weights to model TF influences, which does  
126 not capture changes in relative strength of transcriptional regulation of metabolic genes across  
127 environments. By contrast, BeReTa<sup>20</sup> does take into account weighted, combinatorial influences of TFs,  
128 but the analysis is restricted to genes encoding reactions of specific pathways of interest to an industrial  
129 application. Thus, none of these algorithms were designed to predict systems level phenotypic  
130 consequences (e.g., fitness and growth rate) of perturbations to the transcriptional network.

131 Here, we report the development of **Phenotype of Regulatory influences Integrated with**  
132 **Metabolism and Environment (PRIME)**, which incorporates environment-dependent combinatorial  
133 regulation of metabolic genes to mechanistically predict how individual TFs contribute to the phenotype  
134 of Mtb in any given environment. Through the use of comprehensive experimental validations, we  
135 demonstrate that PRIME significantly outperforms the previous methods in accurately predicting  
136 regulatory and metabolic genes that are conditionally required for growth on carbon sources that are

137 specific for *in vitro* (glycerol) and *in vivo* (cholesterol) growth of Mtb. Further, PRIME has uncovered the  
138 interplay of regulatory and metabolic mechanisms that underlies Mtb's response to drug treatment. The  
139 accuracy of PRIME in predicting quantitative phenotypic effects of TF perturbations is demonstrated by  
140 high correlation between predicted and experimentally validated consequences of knocking out all  
141 metabolism-associated TFs (one-at-a-time) on isoniazid (INH) treatment-specific fitness of Mtb strains.  
142 Through this analysis, we have discovered new vulnerabilities in Mtb that can potentiate INH action,  
143 which are supported by experimental validation.

144

## 145 RESULTS

### 146 CONDITION-SPECIFIC INTEGRATION OF REGULATION AND METABOLISM USING 147 PRIME

148 A causal and mechanistic model of the transcriptional regulatory network and its quantitative influence  
149 on metabolic flux is required to characterize how the 214<sup>21</sup> TFs encoded in the Mtb genome enable its  
150 physiological adaptations to disparate host relevant contexts including antibiotic treatment. We applied  
151 linear regression with TF activity (TFA) estimation using the Inferelator<sup>15,22</sup> to construct an EGRIN from  
152 a compendium of 664 transcriptomes for Mtb that represented transcriptional changes in 3,902 genes  
153 (potentially regulated by 142 TFs) across 77 environmental conditions including drug treatment, pH,  
154 oxygen and carbon source utilization (**Table S1**) (<http://www.colombos.net/>). Relative changes in the  
155 expression of every gene across all conditions were modeled as the sum of weighted influences of a  
156 minimal set of TFs. Altogether, 142 TFs were implicated in the regulation of 3,902 genes in the genome,  
157 acting through a combinatorial scheme represented by 4,820 regulatory influences, (see **Table S2** for  
158 details). EGRIN recapitulated 2,410 of the 4,546 TF- gene interactions in the Mtb P-D network with both  
159 physical binding (from ChIP-seq experiments) and functional evidence (from transcriptional  
160 profiling)<sup>21,23</sup>, and added weights ( $\beta$ ) to the influence of each TF on regulation of its target genes; here  
161 onwards we refer to this subset of 2,410 TF-gene interactions as the "EGRIN-PD Network" (**Table S2**).  
162 Thus, the Inferelator analysis added 2,410 novel TF regulatory influences that were not represented in  
163 the originally compiled P-D interaction network, accounting for 4,820 interactions in total, here onwards  
164 considered as the "EGRIN" network. Briefly, out of 4,820 interactions of EGRIN, 2410 interactions have  
165 P-D evidence (EGRIN P-D).

166

167 We investigated the degree to which EGRIN and EGRIN-PD models captured the regulation of 1,011  
168 genes that encode enzymes implicated in catalyzing 1,229 reactions in the iEK1011<sup>24</sup> model of the *Mtb*  
169 metabolic network. This analysis demonstrated that whereas EGRIN-PD modeled 1,252 regulatory

170 influences of 104 TFs on 605 genes associated with 409 metabolic reactions, EGRIN modeled 2,568  
171 regulatory influences of 129 TFs on 750 genes associated with 725 metabolic reactions. We leveraged  
172 the EGRIN and EGRIN-PD wiring diagrams and weights of regulatory influences inferred by the  
173 Inferelator to predict how change in the activity of a TF in a given environment manifests in altered flux  
174 through a metabolic reaction catalyzed by their regulated gene product. In order to integrate regulation  
175 with metabolism, we had to account for combinatorial regulation of metabolic genes, with each of 349  
176 out of the 750 metabolic genes predicted to be putatively regulated by  $\geq 2$  TFs and 111 TFs predicted to  
177 regulate  $\geq 2$  metabolic genes (**Figure S1 and Table S3**), and association of  $\geq 2$  gene products to each of  
178 313 reactions in Mtb.

179

180 The quantitative influence of a TF on the regulation of a target gene in a given environment was  
181 calculated by multiplying the EGRIN-inferred regression weight ( $\beta$ ) of the TF influence with its absolute  
182 expression level in that environmental condition (i.e., a scaled value of signal intensity for microarray  
183 data or read counts for RNA-seq) based on distribution of values across the transcriptome compendium  
184 (**Figure 1A**; Methods). For a metabolic gene that is regulated by multiple TFs, we calculated the relative  
185 contribution of each TF to the regulation of that gene in a given environment by dividing its quantitative  
186 influence with the sum of quantitative influences of all TFs that regulate that gene. In this scheme, a TF  
187 will have a large relative consequence on the expression of a metabolic gene in an environment in which  
188 the TF is active and in high abundance, and the influences of other TFs are minimal. But the relative  
189 contribution of the TF will be proportionally lower if other TFs are also actively regulating that gene in  
190 that environment. Thus, this approach accounted for regulation of a metabolic gene by multiple TFs, and  
191 it simultaneously corrected for environment-specific changes in combinatorial regulatory schemes. For  
192 a TF that regulates multiple genes encoding enzymes or enzyme subunits for the same reaction, we  
193 considered the largest regulatory influence of that TF on any of those genes to predict its influence on  
194 flux through that reaction. Thus, together these advancements accounted for complex combinatorial  
195 associations between regulation and metabolism to assign a single relative influence factor ( $\gamma$ ) to each  
196 TF-reaction association. The consequence of TF regulation (or knockout) on flux through a reaction is  
197 calculated by multiplying the TF-induced relative inhibition of that reaction ( $1-\gamma$ ) to the maximum possible  
198 flux through that reaction. In this manner, by updating upper bounds of flux through all reactions  
199 catalyzed by regulated gene products of a specific TF, PRIME constrains the metabolic network to a  
200 new solution space, to enable the prediction of “environment-specific” growth consequences of  
201 perturbing a given TF which can be compared to conditional genome-wide fitness data for PRIME  
202 performance assessment (**Figure 1B**).

203

## 204 PERFORMANCE ASSESSMENT OF PRIME

205 In order to compare performance of PRIME to previously developed methods, we had to first update the  
206 PROM model with the latest version of the Mtb P-D interaction map<sup>12,21</sup> and the current version of the  
207 metabolic network model iEK1011<sup>24</sup> (1,011 genes encoding enzymes for 1,229 reactions) that was used  
208 to construct PRIME. Using the methodology described in the original PROM paper<sup>11,12</sup>, 2,416 out of  
209 2,555 P-D interactions for 104 TFs were mapped to 605 genes assigned to 632 reactions in the iEK1011  
210 metabolic network model. This represents a significant improvement in the overall coverage of TFs and  
211 metabolic genes in the PROM model (**Table 1, Figure 2A**). In parallel, we also developed the first  
212 IDREAM model for Mtb by incorporating confidence scores for 2,407 regulatory influences for 142 TFs  
213 within the EGRIN model (FDR <0.25) on a total of 641 genes associated with 639 reactions within  
214 iEK1011 (**Table 1, Figure 2B**). The slightly higher numbers of TFs and metabolic genes in IDREAM and  
215 PRIME (**Figure 2B**) are because they use the EGRIN model, which has better coverage of genome-  
216 wide TF regulation across diverse environments, relative to the P-D interaction map generated in  
217 standard growth conditions that was used in PROM (**Figure 2C**). In summary, the updated PROM and  
218 IDREAM models were similar to PRIME in terms of coverage of the total number of TFs and metabolic  
219 genes and suitable for comparing performance across the models. (**Table 1**).

220

221 We compared the performance of PRIME to PROM and IDREAM by assessing their accuracy (sensitivity  
222 and specificity) in predicting environment-specific growth inhibition upon TF deletion for Mtb cultured in  
223 minimal medium with glycerol or cholesterol as the carbon source. While Mtb is typically grown with  
224 glycerol during *in vitro* culture, the pathogen is capable of utilizing host-derived lipids, such as  
225 cholesterol, during infection. It is known that distinct metabolic genes and networks are associated with  
226 these two modes of growth. Accuracy of model predictions were evaluated using a leave-one-out cross  
227 validation (LOOCV) strategy<sup>25</sup> for comparison of model predictions to experimentally determined  
228 phenotypic consequences of transposon mutagenesis in genome-wide fitness screens (TnSeq) of Mtb  
229 cultured with glycerol or cholesterol<sup>26,27</sup>. Specifically, for each model we generated a set of receiver-  
230 operating characteristic (ROC) curves by plotting the true positive rate (i.e., proportion of model-  
231 predicted essential genes that were verified by experiment) and false positive rate (i.e., proportion of  
232 model-predicted essential genes that were experimentally determined to be non-essential) by leaving  
233 out one TF in each analysis. The distribution of area under the ROC curves (ROC-AUC) from the LOOCV  
234 analysis of model predictions of which TFs are essential for Mtb growth on cholesterol was used as a  
235 metric of performance. First, we evaluated predictions from PRIME using either EGRIN-PD or EGRIN,  
236 inferred using different Inferelator parameter settings as the source of regulatory influences, and  
237 concluded that the latter contributed to significantly better performance (**Figure S2**). Therefore, here

238 onwards all results reported for PRIME are based on regulatory influences from the EGRIN network.  
239 The LOOCV analysis demonstrated that the performance of PRIME was significantly better relative to  
240 PROM and IDREAM in both cholesterol and glycerol carbon sources (**Figure 3A, 3B** and **Figure S3**).  
241 In addition to providing a rigorous means for performance evaluation, the LOOCV<sup>25</sup> analysis also  
242 identified a clear division of TFs in terms of their ROC-AUC values for the PRIME model. Further analysis  
243 revealed that the top performing TFs (20 and 12 TFs for glycerol and cholesterol, respectively)  
244 contributed maximally (up to 65% of overall biomass accumulation) to the overall fitness of Mtb (**Table**  
245 **S4**). Out of 119 TFs with TnSeq data, the cholesterol fitness of 65% (77 TF KOs) were accurately  
246 predicted by PRIME, whereas IDREAM and PROM accurately predicted only 45% (53 TFs) and 30%  
247 (35 TFs), respectively (**Figure 3C**). Similarly, PRIME accurately predicted glycerol fitness for 92 out of  
248 119 TFs (77%), whereas IDREAM accurately predicted 55% (65 TFs) and PROM predicted 36% (43  
249 TFs) (**Figure 3D**). In general, PRIME, IDREAM and PROM predictions differed significantly (p-value  
250 <2.2e-16, t-test) both in the numbers and the context in which genes were called essential or non-  
251 essential.

252

253 Using PRIME, 22 and 7 TFs were accurately predicted (either essential or non-essential) for growth only  
254 with either glycerol or cholesterol, respectively, as determined by experimental fitness screening (**Figure**  
255 **3E**). Similarly, 51 and 25 metabolic genes were accurately predicted by PRIME for growth on either  
256 glycerol or cholesterol, respectively (**Figure 3F**). Among the PRIME predicted essential TFs, Rv2506,  
257 Rv3050c, Rv2760c, and Rv0348 are essential for growth on cholesterol, presumably because they  
258 conditionally regulate genes encoding enzymes or enzyme subunits catalyzing essential metabolic  
259 processes during cholesterol utilization (**Figure 3G**). For example, Rv2506 represses genes likely to be  
260 involved in branched-chain amino acid catabolism, which leads to the production of acetyl-coA and  
261 propionyl-coA<sup>28</sup>. Propionyl-coA is also an endpoint of cholesterol degradation and can be toxic to Mtb<sup>29</sup>.  
262 It is possible that Rv2506 repression of branched-chain amino acid metabolism genes prevents  
263 accumulation of toxic metabolic intermediates during growth on cholesterol. All in all, perturbation of  
264 cholesterol utilization in Mtb could induce metabolite intoxication<sup>29</sup>, unbalanced central metabolism<sup>30</sup> or  
265 lead to carbon starvation<sup>31</sup>. As such, TFs such as Rv2506, Rv3050c, Rv2760c and Rv0348 represent  
266 potential vulnerabilities in the cholesterol utilization pathways of Mtb that could be targeted by drugs.  
267 Notably, these TFs were also ascertained to be essential by the TnSeq screen performed with  
268 cholesterol as the carbon source<sup>26</sup> and are non-essential in glycerol (shown as inactive nodes in **Figure**  
269 **3H**). Other TFs (Rv1990c, Rv0023 and Rv0757) were predicted (and validated by TnSeq<sup>26</sup>) to be  
270 essential for growth with both carbon sources or only essential for growth on glycerol (e.g., Rv0238 and  
271 Rv1423).



272

273 **PRIME RANK IDENTIFIES THE ESSENTIAL TRANSCRIPTIONAL FACTORS AND GENES FOR**  
274 **SURVIVAL DURING DRUG TREATMENT**

275 We used PRIME to investigate the regulatory and metabolic networks that drive physiological  
276 adjustments (e.g., cell wall modifications, shifts in metabolism and respiration) to enable the pathogen  
277 to survive and persist during drug treatment. To expose novel network vulnerabilities of Mtb in response  
278 to drug treatment, we generated transcriptome profiles of Mtb treated for 24 h with high- and low-doses  
279 of seven drugs (**Table S5**). The transcriptome profiles were analyzed using the PRIME model to identify  
280 the metabolic networks and their associated regulators that were essential for growth in the absence  
281 and presence of drug treatment. This analysis found clear distinction in TF essentiality between the  
282 untreated and drug-treated PRIME models and revealed that drug doses largely grouped together  
283 (**Figure 4A**). Interestingly, the TF essentiality profiles of rifampicin (a transcription inhibitor) were dose-  
284 dependent; the rifampicin profile at low-dose clustered separately, while the high-dose profile clustered  
285 with linezolid (a protein synthesis inhibitor). The resemblance to linezolid at high-dose suggests that a  
286 secondary effect of strong rifampicin-induced transcription inhibition also impacts translation.  
287 Furthermore, we observed that the TF essentiality profiles of isoniazid (inhibitor of cell wall synthesis)  
288 were quite distinct to the other six drugs. In fact, 58 TFs become conditionally essential in the presence  
289 of isoniazid because of their mechanistic role in regulating 569 metabolic reactions required for  
290 supporting growth during isoniazid treatment. This highlights the multitude of regulatory-metabolic  
291 networks associated with cell wall disruption in Mtb and the extreme vulnerability in cell wall metabolism.

292

293 Focusing on isoniazid (INH), we evaluated the accuracy of these predictions against experimentally-  
294 determined fitness values from a genome-wide TnSeq screen performed in the presence of a  
295 subinhibitory concentration of INH<sup>32</sup>. Notwithstanding the difference in dosage of drug treatment of the  
296 input transcriptome data used in the PRIME model (0.18 ug/mL, 1.8 ug/mL) and in the TnSeq fitness  
297 screen (27 ng/mL), the LOOCV analysis demonstrated high sensitivity and specificity of PRIME  
298 predictions of gene essentiality (max ROC AUC = 0.685), significantly outperforming PROM (max ROC  
299 AUC = 0.625) and IDREAM (ROC AUC = 0.6) (**Figure 4B**). We also used PRIME to rank order TFs  
300 based on their relative importance in supporting growth in the presence of INH, and compared these  
301 ranks to TnSeq determined importance of TFs. There was striking correlation (Spearman's rho = 0.695;  
302 p-value = 0.0001) in the rank ordering of TFs based on the predicted (PRIME) and observed (TnSeq)  
303 magnitude of growth inhibition of Mtb in the presence of INH upon knocking out each TF one-at-a-time  
304 (**Figure S4**). The correlation increased dramatically (Spearman's rho = 0.746, p-value = 0.0001) when  
305 only TFs implicated by EGRIN as regulators of essential metabolic reactions were considered in this

306 analysis, demonstrating the remarkable accuracy of PRIME in capturing how the differential regulation  
307 by TFs modulates flux through essential metabolic reactions to manifest at a phenotypic level (**Figure**  
308 **4C**). Notably, PRIME accurately predicted that knocking out the top 10 TFs one-at-a-time would result  
309 in at least 65% and up to 95% Mtb growth inhibition during INH treatment, but not in the absence of drug  
310 treatment, implicating these as conditional vulnerabilities for significantly potentiating INH treatment  
311 (**Table S6**).

312

313 To aid in the interpretation of PRIME predictions, we developed the PRIME pathway analysis (PPA) tool  
314 to uncover in a single-step the specific metabolic reaction(s) regulated by a TF that make it essential for  
315 growth in a given environmental condition. Given a TF, PPA identifies all reactions catalyzed by the  
316 genes it is predicted to regulate, rank orders the target genes based on the relative contribution of their  
317 gene product in driving flux towards biomass accumulation, and outputs a TF-metabolic gene-reaction  
318 map as a putative mechanism by which the TF is likely to be essential in a given environmental context.  
319 Using PPA, we identified the specific metabolic reactions that were mechanistically responsible for the  
320 conditional essentiality of 23 TFs validated by TnSeq data<sup>32</sup> to be essential in the presence of INH. For  
321 example, we discovered the mechanisms underlying the essentiality of Rv0827c, Rv1049, Rv1423,  
322 Rv1828 and Rv0472c for growth in the presence of INH (**Figure 4D**). Altogether, PPA uncovered that  
323 58 of the 142 TFs were conditionally essential for growth on INH because they conditionally regulate  
324 569 key reactions across 55 pathways, including 84 reactions within fatty acid metabolism and mycolic  
325 acid biosynthesis (target of INH). In so doing, PRIME has provided the most comprehensive systems  
326 level perspective into strategies to potentiate INH killing by targeting TFs that mediate Mtb's metabolic  
327 response to INH treatment.

328

## 329 **DISCUSSION**

330 We have demonstrated that by incorporating how TFs act contextually in combinatorial schemes to  
331 regulate gene expression, PRIME outperformed PROM and IDREAM in accurately predicting how  
332 transcriptional regulation redirects metabolic flux to manifest in environment-specific phenotypes of Mtb.  
333 The shortcoming of PROM can be attributed to its reliance on P-D interactions for regulatory network,  
334 which are plagued with false positive interactions (because overexpression of TFs can force non-  
335 functional binding across the genome) and false negative interactions because of lack of appropriate  
336 context (e.g., missing co-factors). Hence, a P-D interaction does not capture whether a TF is regulating  
337 a gene in a given condition, which is better modeled by regulatory influences inferred using regression  
338 analysis of transcript level changes in TFs and all genes across the genome. However, despite  
339 incorporating regulatory influences from the same EGRIN network, IDREAM performance was inferior

340 compared to PRIME, and in fact its performance in predicting gene essentiality in cholesterol and INH  
341 was worse than PROM. One explanation could be that relative to the number of P-D interactions used  
342 in PROM, IDREAM used nearly twice as many EGRIN-based regulatory influences that were inferred  
343 from a wide range of environmental contexts, without taking into account combinatorial regulatory  
344 schemes, weights of regulatory influences, or the absolute expression levels of TFs to prune regulatory  
345 edges that were not relevant for a given environmental context. Hence, reliance on a P-D interaction  
346 map, and even just the likelihood that a TF might regulate a gene based on regression analysis are both  
347 insufficient to capture the complex environment-dependent interplay of transcription and metabolism.  
348 Altogether, these comparative analyses have demonstrated that four key advancements in PRIME  
349 addressed the shortcomings of PROM and IDREAM: (i) PRIME took full advantage of EGRIN predictions  
350 to incorporate weights of TF regulatory influence on each gene; (ii) PRIME calibrated the relative  
351 influence of each TF on a given metabolic gene by accounting for all TFs that were also implicated in  
352 the regulation of that gene; (iii) PRIME accounted for regulation of multiple genes that encode enzymes  
353 for the same reaction by considering which gene(s) contributed maximally towards flux through that  
354 reaction in a given environmental context; and, finally (iv) PRIME considered the absolute expression  
355 level of each TF to evaluate the degree to which each regulatory influence was active in a given  
356 environment.

357

358 By demonstrating better accuracy in predicting environment-specific phenotypes of Mtb using EGRIN,  
359 PRIME overrides the need for a physical map of P-D interactions, which is difficult to generate for many  
360 organisms, across all environments of interest, and especially in some contexts, such as within infected  
361 tissue. In fact, the incompleteness of the P-D interaction map was demonstrated by the significant drop  
362 in the performance of PRIME upon excluding regulatory influences that were not supported by physical  
363 TF-gene interactions (i.e., EGRIN P-D). By contrast, EGRIN is inferred directly from a compendium of  
364 transcriptomes, which can be profiled across relevant environmental conditions with minimal  
365 manipulation (e.g., without overexpression of TFs) and even within infected cells using technologies like  
366 Path-seq<sup>33</sup>. As a consequence, EGRIN discovers a significantly larger number of novel regulatory  
367 mechanisms, including the combinatorial schemes and specific environmental contexts in which they  
368 are conditionally active. This explains why PRIME discovered mechanisms that become conditionally  
369 essential in the presence of INH, but also accurately predicted the relative importance of each TF for  
370 enhancing the potency of INH. Based on this observation, we posit that PRIME will be especially  
371 valuable to prioritize genes that represent novel context-dependent vulnerabilities that could be targeted  
372 to potentiate the action of any antibiotic and achieve faster clearance with a lower dosage. By enabling  
373 the *in-silico* discovery of vulnerabilities within the Mtb network, PRIME also overrides the need for large

374 scale transposon mutagenesis-based experiments (e.g., TnSeq, TraSH, HITS, etc), which are resource-  
375 intensive and difficult to perform across all conditions relevant to the lifecycle of Mtb. Instead, PRIME  
376 can be used to rank prioritize the strains and contexts in which to assay for an expected phenotype. This  
377 capability is particularly powerful considering the numerous mechanisms by which Mtb can be  
378 phenotypically different, with different antibiotic sensitivities. Additionally, there is growing evidence that  
379 upon gaining resistance to an antibiotic, the regulatory and metabolic networks within a pathogen are  
380 remodeled in order to reallocate resources for supporting the new phenotype<sup>34</sup>. Using PRIME, we can  
381 delineate novel vulnerabilities within these remodeled regulatory and metabolic networks to devise  
382 strategies for rationally disrupting the antibiotic resistance phenotype with a second drug.

383

384 PRIME will also be useful in biotechnology applications to further optimize the production of desired end  
385 products by rewiring the regulatory networks of metabolically engineered strains. Advancements in  
386 metabolic engineering have been effective in substantially increasing flux towards the production of a  
387 desired metabolite<sup>18–20,35</sup> but there is a limit to which metabolic engineering alone can improve the overall  
388 yield. It has been proposed that further enhancements in yield would require reprogramming of the  
389 regulatory network to control when genes of the engineered pathways are expressed, and to rationally  
390 up and down regulate competing metabolic pathways to maximize flux and resource allocation towards  
391 the desired objective. Hence, by using PRIME, metabolic engineering of high-yielding strain phenotypes  
392 can be identified. Although the capabilities of PRIME are elucidated extensively using Mtb as a model  
393 system in this study, we foresee the use and applications of PRIME in various organisms due to its  
394 scalability.

395

## 396 **METHODS**

### 397 **CONSTRUCTION OF EGRIN GENE REGULATORY NETWORK FOR *MYCOBACTERIUM*** 398 ***TUBERCULOSIS***

399 The Mtb EGRIN used in this study was constructed using the Inferelator algorithm<sup>15,22</sup> trained on a  
400 transcriptional compendium for Mtb with 3,902 genes across 664 experimental conditions (downloaded  
401 from the COLOMBOS database) and an experimentally supported signed Mtb P-D network (generated  
402 as previously described in<sup>33</sup>). The original transcriptional compendium contained a larger number of  
403 genes and conditions but was modified to remove genes and conditions with missing values. Briefly, we  
404 used the Inferelator to identify potential transcriptional regulators for the 3,902 Mtb genes in the  
405 expression compendium, as previously performed for other species<sup>22,36</sup>. The Inferelator first estimates  
406 the regulatory activities of each transcription factor activity (TFA) using the expression profile of TF

407 known targets (encoded in the signed P-D network). Then, the Inferelator uses a Bayesian Best Subset  
408 Regression to estimate the magnitude and sign (activation or repression) of potential interactions  
409 between TFs and genes. As before, we bootstrapped the expression data (20 times) to avoid regression  
410 overfitting. The Inferelator generates two scores for each TF-gene interaction, the corresponding  
411 regression coefficient (weight -  $\beta$ ) and a confidence score. The second score indicates the likelihood of  
412 the interaction. The final set of TF-gene interactions was defined with a 0.5 precision cutoff. This means  
413 that 50% of all interactions in the inferred network were already present in the signed P-D network used  
414 for training, while the other half corresponded to putative novel TF-gene interactions.

415

## 416 DEVELOPMENT OF PRIME

417 The PRIME algorithm has been developed by integrating weights ( $\beta$ ) from EGRIN with metabolic  
418 network (MN) models for phenotype prediction in a context-specific manner (wiring diagram in Fig. 1).  
419 PRIME requires 1) a MN in the format of constraints-based model<sup>37,38</sup> in systems biology markup  
420 language (SBML), an XML format as input, that are represented *in silico* in the form of a stoichiometric  
421 matrix, wherein every column corresponds to a reaction and every row corresponds to a metabolite.  
422 These constraints-based models were used to integrate the regulatory influences by updating the  
423 reaction flux, 2) a regulatory network containing TF and gene interactions (one array of regulators and  
424 one array of corresponding gene targets), 3) magnitude/weights ( $\beta$ ) of regulatory influences for each of  
425 the interactions (array of magnitudes) derived from Inferelator and 4) the gene expression data profiled  
426 under a specific condition (gene ids and their expression, provided as ratio to the control - in case of  
427 environment-specific predictions the ratio between initial  $t_0$  and final time point  $t_n$ ). The pipeline of  
428 PRIME initially links each metabolic gene in MN to its associated regulators considering the  
429 combinatorial effects, followed by applying the calculated relative influence factor. Specifically, we have  
430 introduced a new way to calculate the relative influence factor ( $\gamma$ ), a value that quantitatively constrains  
431 the reaction flux constraint space. The equations 1 to 5 consists of the details involved in each  
432 successive step within the algorithm.

433

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

$$435 \gamma_{i,w} = 1 - \left( \frac{\beta_{i,j} X'_j}{\sum_{j \in J} \beta_{i,j} X'_j} \right) \quad (Eq. 1)$$

436 where  $J$  is the subset of TFs that influence gene  $i$  and  $X'_j$  is the scaled expression of a TF  $j$  in a  
437 particular condition  $c$  of a coherent environmental context  $B$  as,

$$438 X'_j = \frac{X_{j,c} - \min X_j(B)}{\min X_j(B) - \max X_j(B)} \quad (Eq. 2)$$

439 Then, the regulatory influence that exerts the larger effect on reaction  $w$  across the set of metabolic  
440 genes  $i \in I$  of a given reaction has been identified as,

$$441 \quad g_w = \min \gamma_{i,w} \quad (\text{Eq. 3})$$

442 At this point, it is straightforward to incorporate calculated weights as new upper bounds,

$$443 \quad b^{PRIME} = b \circ g = (b)_w(g)_w \quad (\text{Eq. 4})$$

444 to the flux balance analysis (FBA)<sup>38</sup> formalism, assuming steady state metabolic concentrations, and  
445 defining the system mass balance as  $S \cdot v = 0$ , to maximize the objective function  $Z = c^T v$  such that  
446 fluxes are within the new boundary conditions,

$$447 \quad a \leq v \leq b^{PRIME} \quad (\text{Eq. 5})$$

448 The objectives in each prediction are defined during FBA optimization. The phenotype predictions  
449 mentioned in this study are the optimized biomass predicted by FBA. The complete PRIME algorithm  
450 package and details of the required input dataset is available for download from our GitHub Repository  
451 (<https://github.com/baliga-lab/PRIME>). All model simulations related to FBA were performed on  
452 MATLAB\_R2019a platform using the recent version of COBRA<sup>39</sup> (The COntstraint-Based Reconstruction  
453 and Analysis) toolbox. *In silico* gene essentiality predictions were performed using the COBRA toolbox  
454 'single-gene-deletion' function in MATLAB.

455

## 456 **INCORPORATING DRUG TREATMENT GENE EXPRESSION DATA ON METABOLIC MODEL**

457 The iEK1011<sup>24</sup> metabolic network (MN) model was used for all the predictions in this study. For drug-  
458 specific models, we applied the gene expression data from both drug-treated and untreated control  
459 experiments using the GIMME<sup>40</sup> algorithm on the iEK1011 MN model. This step was carried out to  
460 constrain the MN model to the specific condition being tested. We used GIMME because of the flexibility  
461 in defining objective function during implementation. The GIMME algorithm is implemented in the  
462 MATLAB\_R2019a platform, using the "GIMME.m function" in the COBRA Toolbox after processing the  
463 gene expression data through 'mapExpressionToReactions.m' function to convert the gene expression  
464 values as inputs to GIMME.

465

## 466 **PROM MODELS**

467 For developing PROM<sup>11,12</sup> models, we followed the PROM approach<sup>11</sup> to estimate the probability that a  
468 target gene is 'ON' or 'OFF' in the absence of the TF i.e., in the event of a TF knockout. This was  
469 calculated from a gene expression dataset as, Probability, P (Gene = 1|TF = 0) or P (TF = 1|Gene = 0).  
470 The gene expression threshold that delineated between the 'ON' and 'OFF' states was set as quantile  
471 (0.33) from the input expression data. These probabilities were then used to constrain the maximal fluxes  
472 of the reactions catalyzed by the gene products in the metabolic model as  $p \times V_{max}$ , where  $p$  is the

473 probability of the gene being on. The user defined “kappa” value was used as similar to earlier PROM  
474 models<sup>11</sup>. All PROM predictions and simulations were performed using PROM.m (MATLAB script) on  
475 the MATLAB\_R2019a platform. We used iEK1011 metabolic network model in XML format as input in  
476 the PROM. The P-D derived regulatory network was obtained from the study<sup>21</sup>, similar to the  
477 MTBPROMV2.0<sup>12</sup>.

478

## 479 **IDREAM MODELS**

480 For IDREAM<sup>13</sup> models, the GRN derived using EGRIN, was integrated with the PROM pipeline as it had  
481 been done previously for the yeast system<sup>13</sup>. We ran 200 iterations in EGRIN to calculate the confidence  
482 score for all predictions. For each gene, we estimated a false discovery rate (FDR) for each TF by  
483 counting the fraction of models that identified that factor as a regulator. Thus, if TF1 was predicted to  
484 regulate gene1 in 191 of 200 models, then the TF-gene interaction identified would have an FDR =  
485 0.045. We included only those interactions that passed an FDR cutoff of 0.25. We used EGRIN-derived  
486 GRN to integrate it with iEK1011 metabolic network model of Mtb using the PROM framework. The user  
487 defined “kappa” value was used as similar to earlier PROM models<sup>11</sup>. IDREAM does not rely on  
488 probabilities, hence the gene expression dataset was not used in IDREAM instead ‘prob\_prior’ in the  
489 PROM function was set based on the EGRIN FDR values for each TF-gene interaction. If the TF is an  
490 activator of a gene, we use the FDR value directly, if it is an inhibitor, we use 1-FDR value as ‘prob\_prior’.  
491 EGRIN network was derived using Inferelator in R (Inferelator.pkg.R) and PROM predictions and  
492 simulations were performed using PROM.m (MATLAB script) on the MATLAB\_R2019a platform as  
493 similar to PROM model development.

494

## 495 **PERFORMANCE ASSESSMENT OF PRIME PREDICTIONS**

496 The predictive power of PRIME as a binary classifier (essential or non-essential) between the model  
497 predicted gene essentiality and experimentally defined gene essentiality (TnSeq) has been performed  
498 using receiver operating characteristic (ROC) curve. A gene was considered “essential” if its deletion  
499 reduced the biomass by >85%. By this analysis, the model classified each gene as “essential” or “non-  
500 essential”. We compared the gene essentiality predictions from Mtb grown under glycerol and  
501 cholesterol as carbon source with the available experimental TnSeq data<sup>26</sup> and deduced the confusion  
502 matrix to derive true positive rates (TPR) and false positive rates (FPR). We also took advantage of the  
503 follow-up study where Bayesian analysis was used to assign calls as essential and non-essential for the  
504 same TnSeq dataset<sup>27</sup>. We expanded the analysis of TnSeq data to classify essential and non-essential  
505 with a cutoff value of using cholesterol/glycerol ratio of 0.6 in order to assign calls for all the genes. This  
506 classification led to the elucidation of sensitivity and specificity of the model using ROC curve analysis.

507 Briefly, the gene expression data of Mtb profiled under growth on Glycerol (GSE52020) and Cholesterol  
508 (GSE13978) were used to generate condition-specific metabolic networks using GIMME. PRIME was  
509 applied on these models to predict gene and TF essentialities according to the condition tested. These  
510 predictions were then compared to the TnSeq data. A similar sensitivity and specificity analysis was  
511 performed while validating the performance of PRIME for INH-specific predictions using experimentally  
512 derived TnSeq data<sup>32</sup>. To construct the INH-specific metabolic models, we used INH-treated Mtb  
513 transcriptome sequencing (RNA-seq) data generated in this study (see below).

514

### 515 **PRIME PATHWAY ANALYSIS (PPA) PIPELINE**

516 The PRIME pathway analysis (PPA) pipeline was developed to derive the metabolic association of a  
517 specified TF in a simple process by accessing PRIME model genes and their interactions. The top  
518 ranked TFs and their associated metabolic genes are further linked to their metabolic processes using  
519 the PPA pipeline. PPA is provided as PRIMEanalysis.m (MATLAB script). All analyses related to PPA  
520 were performed in MATLAB\_R2019a platform. The illustration of PPA-derived essential gene regulatory-  
521 metabolic networks were deduced using BioTapestry tool (<http://www.biotapestry.org/>).

522

### 523 **DRUG TREATMENT CULTURING CONDITIONS**

524 Experiments were performed using *Mycobacterium tuberculosis* H37Rv grown with mild agitation at  
525 37°C in standard 7H9-rich media consisting of Middlebrook 7H9 broth supplemented with 10%  
526 Middlebrook ADC, 0.05% Tween-80, and 0.2% glycerol. Frozen 1 mL stocks of Mtb cells were added to  
527 7H9-rich medium and grown until the culture reached an OD<sub>600</sub> of ~0.4-0.8. The cells were then diluted  
528 to OD<sub>600</sub> of 0.05 and added to 7H9-rich medium containing drugs at the predetermined amounts.  
529 Samples, in biological triplicate, were collected at 24 h after drug treatment by centrifugation at high  
530 speed for 5 min, discarding supernatant and immediately flash freezing the cell pellet in liquid nitrogen.  
531 Cell pellets were stored at -80° C until RNA extraction was performed as previously described<sup>41</sup>.

532

### 533 **PROCESSING AND ANALYSIS OF RNA-SEQ DATA**

534 Sample collection and RNA-extraction was performed as described above. Total RNA samples were  
535 depleted of ribosomal RNA using the Ribo-Zero Bacteria rRNA Removal Kit (Illumina, San Diego, CA).  
536 Quality and purity of mRNA samples was determined with 2100 Bioanalyzer (Agilent, Santa Clara, CA).  
537 Samples were prepared with TrueSeq Stranded mRNA HT library preparation kit (Illumina, San Diego,  
538 CA). All samples were sequenced on the NextSeq sequencing instrument in a high output 150 v2 flow  
539 cell. Paired-end 75 bp reads were checked for technical artifacts using Illumina default quality filtering  
540 steps. Raw FASTQ read data were processed using the R package DuffyNGS<sup>42</sup>. Briefly, raw reads were



541 passed through a 2-stage alignment pipeline: (i) a pre-alignment stage to filter out unwanted transcripts,  
542 such as rRNA; and (ii) a main genomic alignment stage against the genome of interest. Reads were  
543 aligned to *M. tuberculosis* H37Rv (ASM19595v2) with Bowtie2<sup>43</sup>, using the command line option “very-  
544 sensitive.” BAM files from stage (ii) were converted into read depth wiggle tracks that recorded both  
545 uniquely mapped and multiply mapped reads to each of the forward and reverse strands of the  
546 genome(s) at single-nucleotide resolution. Gene transcript abundance was then measured by summing  
547 total reads landing inside annotated gene boundaries, expressed as both RPKM and raw read counts.  
548 We used the raw read counts as input for DESeq2<sup>44</sup> to obtain DESeq2 normalized counts. The RNA-  
549 seq data of Mtb response to drug exposure generated for this study are publicly available at the Gene  
550 Expression Omnibus under accession number GSE165673.

551

552

## 553 REFERENCES

- 554 1. Stanley, S. A. & Cox, J. S. Host–Pathogen Interactions During Mycobacterium tuberculosis  
555 infections. in *Springer International Publishing* **410**, 211–241 (2013).
- 556 2. Rienksma, R. A., Schaap, P. J., Martins dos Santos, V. A. P. & Suarez-Diez, M. Modeling the  
557 Metabolic State of Mycobacterium tuberculosis Upon Infection. *Front. Cell. Infect. Microbiol.* **8**,  
558 1–13 (2018).
- 559 3. Galagan, J. E. *et al.* The Mycobacterium tuberculosis regulatory network and hypoxia. *Nature*  
560 **499**, 178–183 (2013).
- 561 4. Chaulk, C. P. & Kazandjian, V. A. Directly observed therapy for treatment completion of  
562 pulmonary tuberculosis: Consensus statement of the public health tuberculosis guidelines  
563 panel. *J. Am. Med. Assoc.* **279**, 943–948 (1998).
- 564 5. Sarathy, J. P. *et al.* Extreme drug tolerance of mycobacterium tuberculosis in Caseum.  
565 *Antimicrob. Agents Chemother.* **62**, (2018).
- 566 6. Sarathy, J., Dartois, V., Dick, T. & Gengenbacher, M. Reduced drug uptake in phenotypically  
567 resistant nutrient-starved nonreplicating Mycobacterium tuberculosis. *Antimicrob. Agents*  
568 *Chemother.* **57**, 1648–1653 (2013).
- 569 7. de Steenwinkel, J. E. M. *et al.* Time-kill kinetics of anti-tuberculosis drugs, and emergence of  
570 resistance, in relation to metabolic activity of Mycobacterium tuberculosis. *J. Antimicrob.*  
571 *Chemother.* **65**, 2582–2589 (2010).
- 572 8. Rao, S. P. S., Alonso, S., Rand, L., Dick, T. & Pethe, K. The protonmotive force is required for  
573 maintaining ATP homeostasis and viability of hypoxic, nonreplicating Mycobacterium  
574 tuberculosis. *Proc. Natl. Acad. Sci. U. S. A.* **105**, 11945–11950 (2008).

- 575 9. Peterson, E. J. R. *et al.* Intricate Genetic Programs Controlling Dormancy in Mycobacterium  
576 tuberculosis. *Cell Rep.* **31**, (2020).
- 577 10. Covert, M. W., Schilling, C. H. & Palsson, B. Regulation of gene expression in flux balance  
578 models of metabolism. *J. Theor. Biol.* **213**, 73–88 (2001).
- 579 11. Chandrasekaran, S. & Price, N. D. Probabilistic integrative modeling of genome-scale metabolic  
580 and regulatory networks in Escherichia coli and Mycobacterium tuberculosis. *Proc. Natl. Acad.*  
581 *Sci.* **107**, 17845–17850 (2010).
- 582 12. Ma, S. *et al.* Integrated Modeling of Gene Regulatory and Metabolic Networks in Mycobacterium  
583 tuberculosis. *PLOS Comput. Biol.* **11**, e1004543 (2015).
- 584 13. Wang, Z. *et al.* Combining inferred regulatory and reconstructed metabolic networks enhances  
585 phenotype prediction in yeast. *PLOS Comput. Biol.* **13**, e1005489 (2017).
- 586 14. Peterson, E. J. R., Ma, S., Sherman, D. R. & Baliga, N. S. Network analysis identifies Rv0324  
587 and Rv0880 as regulators of bedaquiline tolerance in Mycobacterium tuberculosis. *Nat.*  
588 *Microbiol.* **1**, 16078 (2016).
- 589 15. Bonneau, R. *et al.* The Inferelator: an algorithm for learning parsimonious regulatory networks  
590 from systems-biology data sets de novo. *Genome Biol.* **7**, R36 (2006).
- 591 16. Brooks, A. N. *et al.* A system-level model for the microbial regulatory genome. *Mol. Syst. Biol.*  
592 **10**, 740 (2014).
- 593 17. Peterson, E. J. R. *et al.* A high-resolution network model for global gene regulation in  
594 Mycobacterium tuberculosis. *Nucleic Acids Res.* **42**, 11291–11303 (2014).
- 595 18. Kim, J. & Reed, J. L. OptORF: Optimal metabolic and regulatory perturbations for metabolic  
596 engineering of microbial strains. *BMC Syst. Biol.* **4**, 53 (2010).
- 597 19. Yang, L., Cluett, W. R. & Mahadevan, R. EMILiO: A fast algorithm for genome-scale strain  
598 design. *Metab. Eng.* **13**, 272–281 (2011).
- 599 20. Kim, M., Sun, G., Lee, D.-Y. & Kim, B.-G. BeReTa: a systematic method for identifying target  
600 transcriptional regulators to enhance microbial production of chemicals. *Bioinformatics* **33**, 87–  
601 94 (2017).
- 602 21. Minch, K. J. *et al.* The DNA-binding network of Mycobacterium tuberculosis. *Nat. Commun.* **6**,  
603 (2015).
- 604 22. Arrieta-Ortiz, M. L. *et al.* An experimentally supported model of the Bacillus subtilis global  
605 transcriptional regulatory network. *Mol. Syst. Biol.* **11**, 839 (2015).
- 606 23. Rustad, T. R. *et al.* Mapping and manipulating the Mycobacterium tuberculosis transcriptome  
607 using a transcription factor overexpression-derived regulatory network. *Genome Biol.* **15**, 502  
608 (2014).

- 609 24. Kavvas, E. S. *et al.* Updated and standardized genome-scale reconstruction of *Mycobacterium*  
610 tuberculosis H37Rv, iEK1011, simulates flux states indicative of physiological conditions. *BMC*  
611 *Syst. Biol.* **12**, 25 (2018).
- 612 25. Webb, G. I. *et al.* Leave-One-Out Cross-Validation. in *Encyclopedia of Machine Learning* 600–  
613 601 (Springer US, 2011). doi:10.1007/978-0-387-30164-8\_469
- 614 26. Griffin, J. E. *et al.* High-Resolution Phenotypic Profiling Defines Genes Essential for  
615 Mycobacterial Growth and Cholesterol Catabolism. *PLoS Pathog.* **7**, e1002251 (2011).
- 616 27. DeJesus, M. A. *et al.* Bayesian analysis of gene essentiality based on sequencing of transposon  
617 insertion libraries. *Bioinformatics* **29**, 695–703 (2013).
- 618 28. Balhana, R. J. C. *et al.* *bkaR* is a TetR-type repressor that controls an operon associated with  
619 branched-chain keto-acid metabolism in *Mycobacteria*. *FEMS Microbiol. Lett.* **345**, 132–140  
620 (2013).
- 621 29. Lee, W., VanderVen, B. C., Fahey, R. J. & Russell, D. G. Intracellular *Mycobacterium*  
622 tuberculosis exploits host-derived fatty acids to limit metabolic stress. *J. Biol. Chem.* **288**, 6788–  
623 6800 (2013).
- 624 30. Martinot, A. J. *et al.* Mycobacterial Metabolic Syndrome: LprG and Rv1410 Regulate  
625 Triacylglyceride Levels, Growth Rate and Virulence in *Mycobacterium tuberculosis*. *PLoS*  
626 *Pathog.* **12**, (2016).
- 627 31. Pandey, A. K. & Sassetti, C. M. Mycobacterial persistence requires the utilization of host  
628 cholesterol. *Proc. Natl. Acad. Sci. U. S. A.* **105**, 4376–4380 (2008).
- 629 32. Xu, W. *et al.* Chemical Genetic Interaction Profiling Reveals Determinants of Intrinsic Antibiotic  
630 Resistance in *Mycobacterium tuberculosis*. *Antimicrob. Agents Chemother.* **61**, 1–15 (2017).
- 631 33. Peterson, E. J. *et al.* Path-seq identifies an essential mycolate remodeling program for  
632 mycobacterial host adaptation. *Mol. Syst. Biol.* **15**, e8584 (2019).
- 633 34. Arrieta-Ortiz, M. L. *et al.* Disrupting the ArcA regulatory network increases tetracycline  
634 susceptibility of Tet R *Escherichia coli* 2.3. *bioRxiv* 2020.08.31.275693 (2020).  
635 doi:10.1101/2020.08.31.275693
- 636 35. Burgard, A. P., Pharkya, P. & Maranas, C. D. Optknock: A bilevel programming framework for  
637 identifying gene knockout strategies for microbial strain optimization. *Biotechnol. Bioeng.* **84**,  
638 647–657 (2003).
- 639 36. Arrieta-Ortiz, M. L. *et al.* Inference of Bacterial Small RNA Regulatory Networks and Integration  
640 with Transcription Factor-Driven Regulatory Networks. *mSystems* **5**, (2020).
- 641 37. Varma, A. & Palsson, B. O. Metabolic Capabilities of *Escherichia coli*: I. Synthesis of  
642 Biosynthetic Precursors and Cofactors. *J. Theor. Biol.* **165**, 477–502 (1993).

- 643 38. Orth, J. D., Thiele, I. & Palsson, B. Ø. What is flux balance analysis? *Nat Biotechnol* **28**, 245–  
644 248 (2010).
- 645 39. Heirendt, L. *et al.* Creation and analysis of biochemical constraint-based models using the  
646 COBRA Toolbox v.3.0. *Nat. Protoc.* **14**, 639–702 (2019).
- 647 40. Becker, S. A. & Palsson, B. O. Context-Specific Metabolic Networks Are Consistent with  
648 Experiments. *PLoS Comput. Biol.* **4**, e1000082 (2008).
- 649 41. Peterson, E. J. R. *et al.* Intricate Genetic Programs Controlling Dormancy in *Mycobacterium*  
650 *tuberculosis*. *Cell Rep.* **31**, (2020).
- 651 42. Vignali, M. *et al.* NSR-seq transcriptional profiling enables identification of a gene signature of  
652 *Plasmodium falciparum* parasites infecting children. *J. Clin. Invest.* **121**, 1119–1129 (2011).
- 653 43. Langmead, B. & Salzberg, S. L. Fast gapped-read alignment with Bowtie 2. *Nat. Methods* **9**,  
654 357–359 (2012).
- 655 44. Love, M. I., Huber, W. & Anders, S. Moderated estimation of fold change and dispersion for  
656 RNA-seq data with DESeq2. *Genome Biol.* **15**, 550 (2014).

657

658 **Acknowledgements.** We thank members of the Baliga lab for critical discussions and feedback. This  
659 work is funded by Bill and Melinda Gates Foundation (INV-009322) and the National Institute of Allergy  
660 and Infectious Diseases of the National Institutes of Health (R01AI128215 and U19AI135976).

661

662 **Competing interest.** The authors declare no competing interest.

663

664 **Author Contributions.** **NSB** and **SRCI** conceptualized the study and designed the research. **SRCI**  
665 developed the PRIME algorithm, performed all the computational analyses related to PRIME and  
666 metabolic modelling, analyzed all data represented and designed all figures. **MLAO** performed  
667 computational analyses related to regulatory networks and contributed in the initial process of PRIME  
668 conceptualization. **RR** and **MP** performed the experiments related to drug treatment and transcriptomic  
669 profiling. **ALGL** contributed in PRIME method conceptualization in early stages. **EJRP** designed and  
670 analyzed transcriptome studies, and analyzed fitness data of PRIME analysis. **NSB** and **EJRP** provided  
671 overall supervision. **SRCI**, **EJRP**, and **NSB** wrote the manuscript. All authors read the manuscript and  
672 approved its content.

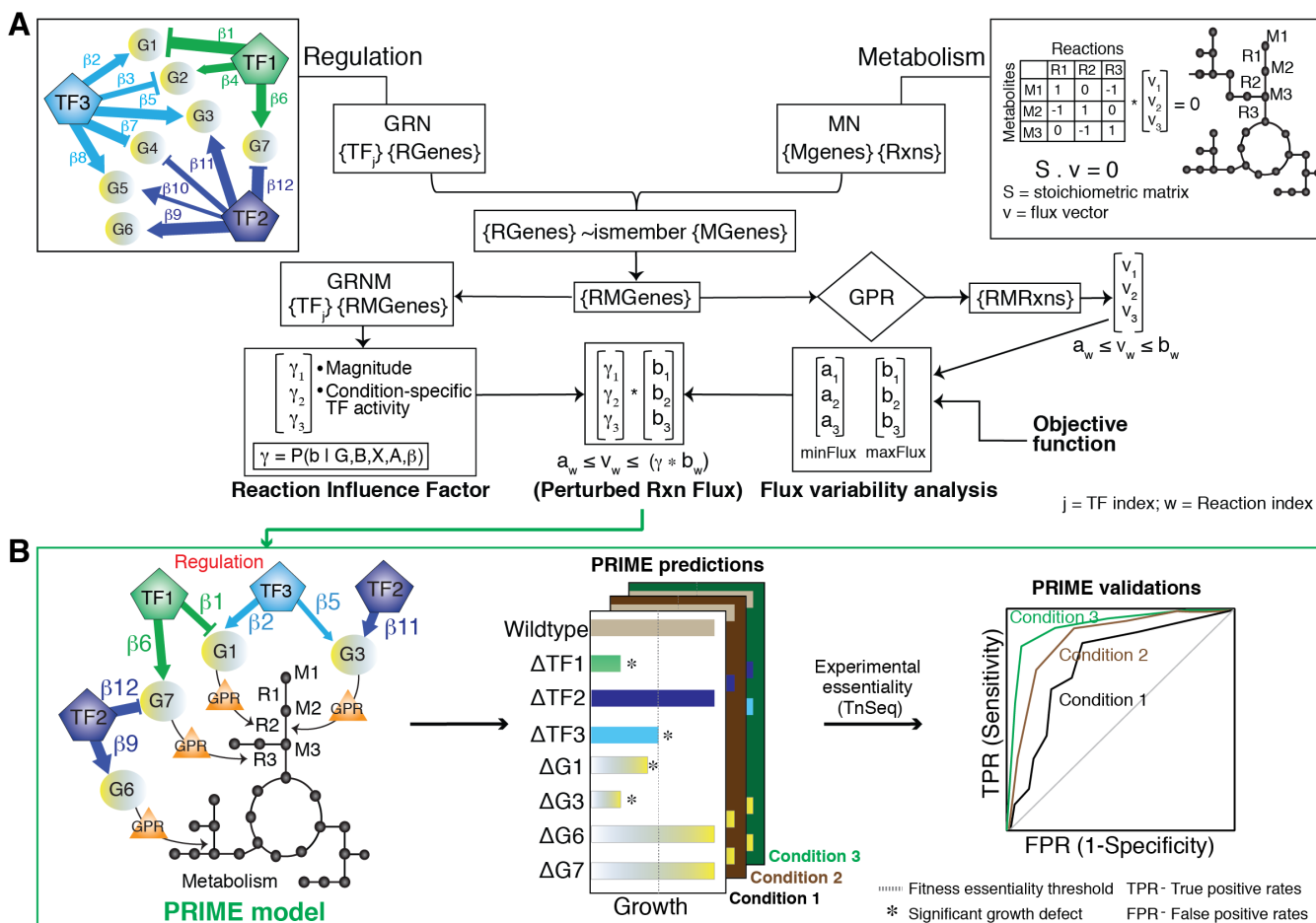
673

674 **Data availability.** Input files for PRIME used in this study are provided as **File S1**. All PRIME-generated  
675 data are provided as supplementary materials. PRIME code, with data and description for  
676 implementation, is available in GitHub repository: <https://github.com/baliga-lab/PRIME>. The RNA-seq

677 data generated for this study are available in the Gene Expression Omnibus under accession no.  
678 GSE165673.

679

680 **Figures and figure legends:**

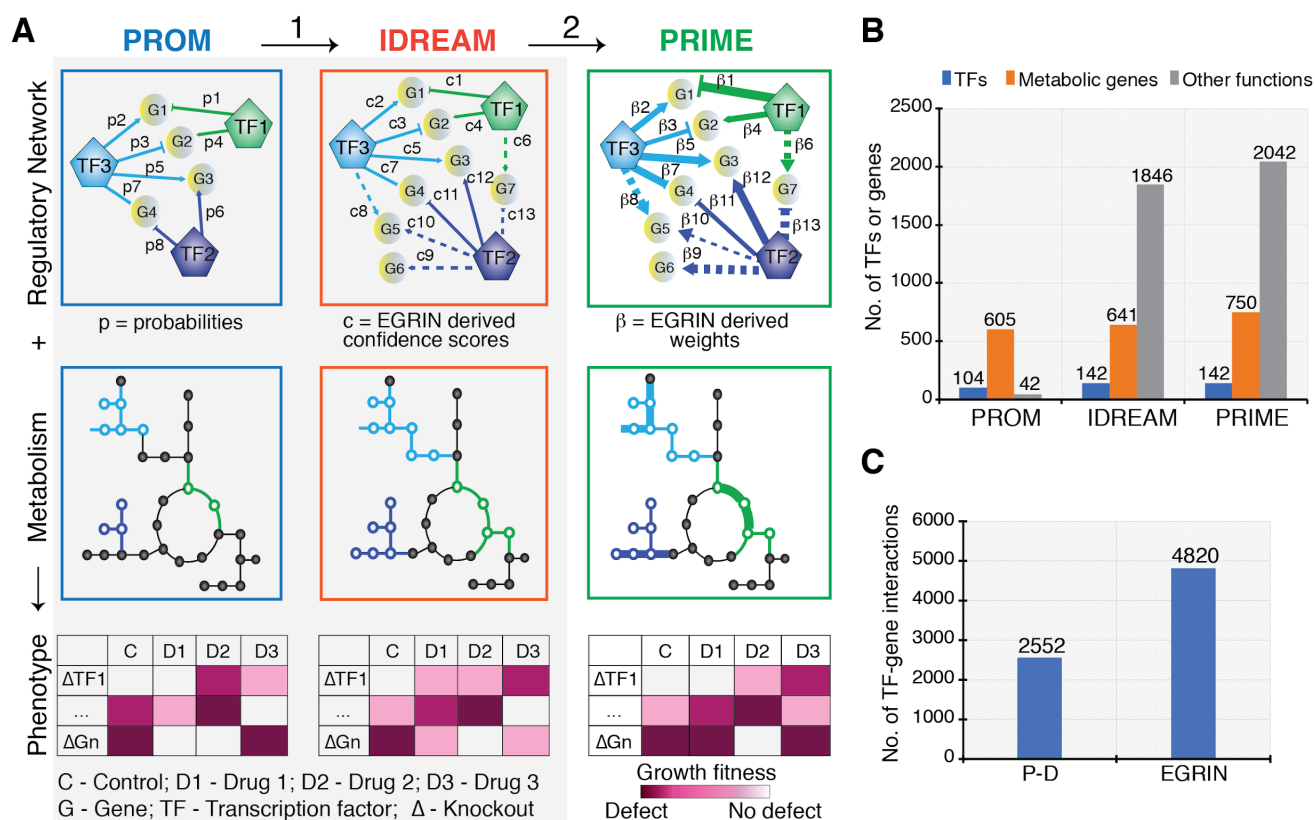


681

682 **Figure 1: Schematic for PRIME model development and performance assessment.** **A.** Schema for  
683 integration of gene regulation and metabolism. The gene regulatory network (**GRN**) models weighted  
684 regulatory influences of TFs on regulated genes (**RGenes**). A subset of the RGenes are enzyme-coding  
685 metabolic genes (**Mgenes**), whose functions are also modeled through gene-to-protein-to-reaction  
686 (GPR) mapping in a stoichiometric matrix representation of the metabolic network (**MN**). PRIME uses  
687 the integrated Gene Regulatory Network of Metabolism (**GRNM**) and a reaction flux influence estimator  
688 (**ReFINE**) to calculate the  $\gamma$  factor, which quantifies how the differential expression of multiple TFs and  
689 their weighted regulatory influences on a regulated metabolic gene (**RMGene**) manifests in altered flux  
690 (**a**: minimum flux; **b**: maximum flux) through the associated metabolic reaction (**RMRxn**) in a given  
691 environmental condition. **B.** Illustration of condition-specific gene phenotype predictions and  
692 performance assessment. The example illustrates how PRIME predicts relative growth consequence of

693 single gene knockouts in TFs (e.g., TF1, TF2 and TF3) and RMGenes (e.g., G1, G3, G6 and G7) in  
694 different contexts (e.g., Condition 1, 2, and 3). The vertical line in the barplot depicts a user-defined  
695 threshold in growth inhibition, below which a gene is deemed essential. Performance of PRIME is  
696 quantified using a Receiver Operating Characteristic (**ROC**) curve based on accuracy of PRIME-  
697 predicted essential and non-essential genes in a given condition to experimentally determined  
698 phenotype consequences using transposon mutagenesis coupled with sequencing (**TnSeq**) in the same  
699 condition.

700



701

702 **Figure 2: PRIME model advancements.** **A.** Advancements in PRIME over previous methods (PROM

703 and IDREAM) are indicated as (1) incorporation of regulatory influences from EGRIN (regression-based

704 interactions are shown as dotted lines), which increases coverage of the regulatory network, (2)

705 incorporation of the magnitude of regulatory influence of TFs on metabolic genes ( $\beta$  - shown as varying

706 edge thickness) instead of probability ( $p$ ) and confidence score ( $c$ ) significantly improved the predictive

707 accuracy of environment-specific gene essentiality. **B.** Number of TFs and genes from PRIME, IDREAM

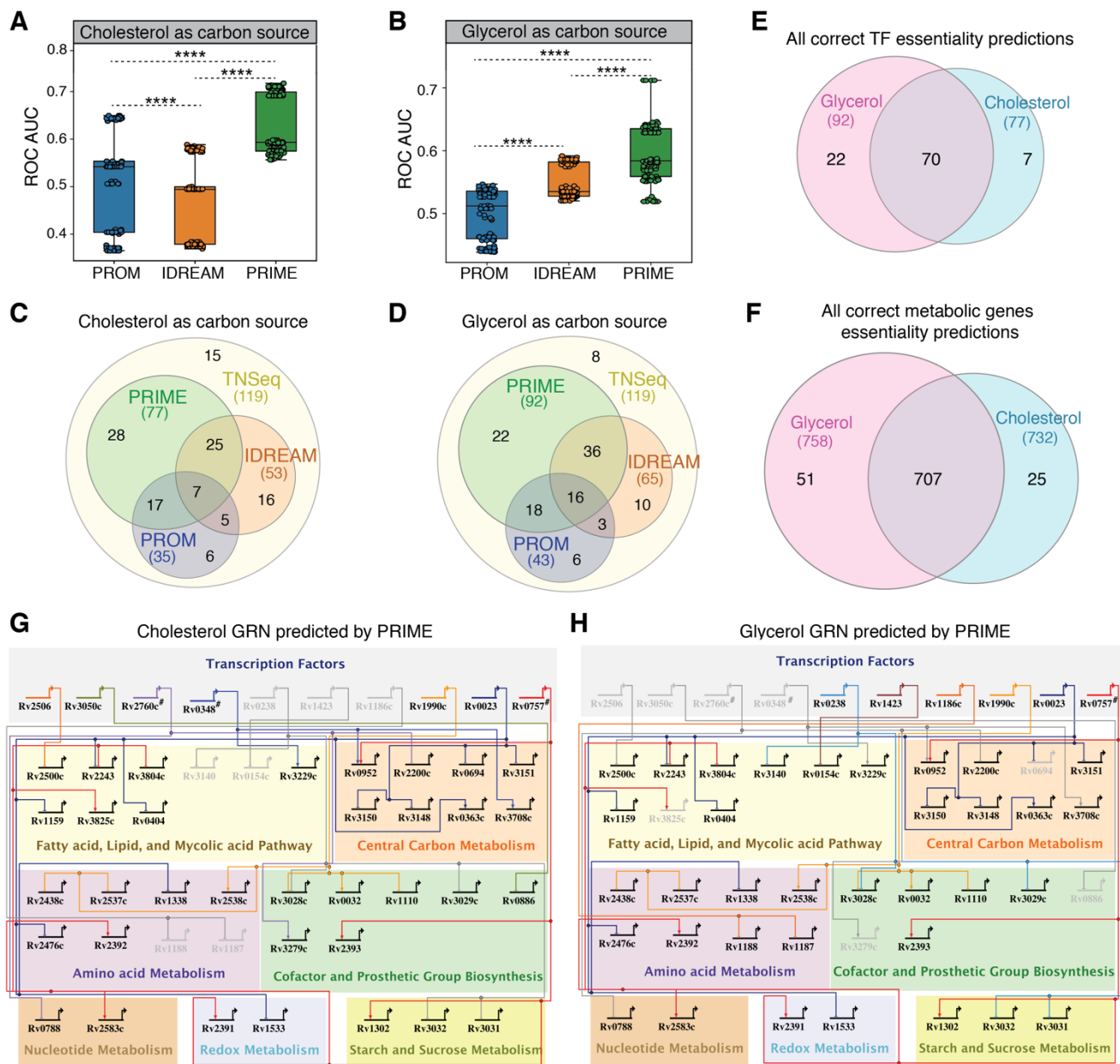
708 and PROM. **C.** Number of TF-gene interactions identified using regression-based EGRIN and Protein-

709 DNA (P-D) interactions from ChIP-seq data.

710

711

712



713

714

**Figure 3: Validation of PRIME predictions of conditional gene essentiality.** Sensitivity and

715

specificity of PRIME, PROM, and IDREAM predicted TF essentiality in **A.** cholesterol and **B.** glycerol as

716

as determined by LOOCV analysis for the area under the receiver operating characteristic curve (ROC

717

AUC). Statistical significance was calculated as  $p$ -value with two sample t-test. \*\*\*\*:  $p$ -value < 0.0001.

718

Comparison of all positive predictions (true positives and true negatives) for TF essentiality by PRIME,

719

PROM, and IDREAM in **C.** cholesterol and **D.** glycerol. **E.** The number of all correct PRIME predictions

720

(true positives and true negatives) of TF knockouts across the two conditions (glycerol and cholesterol)

721

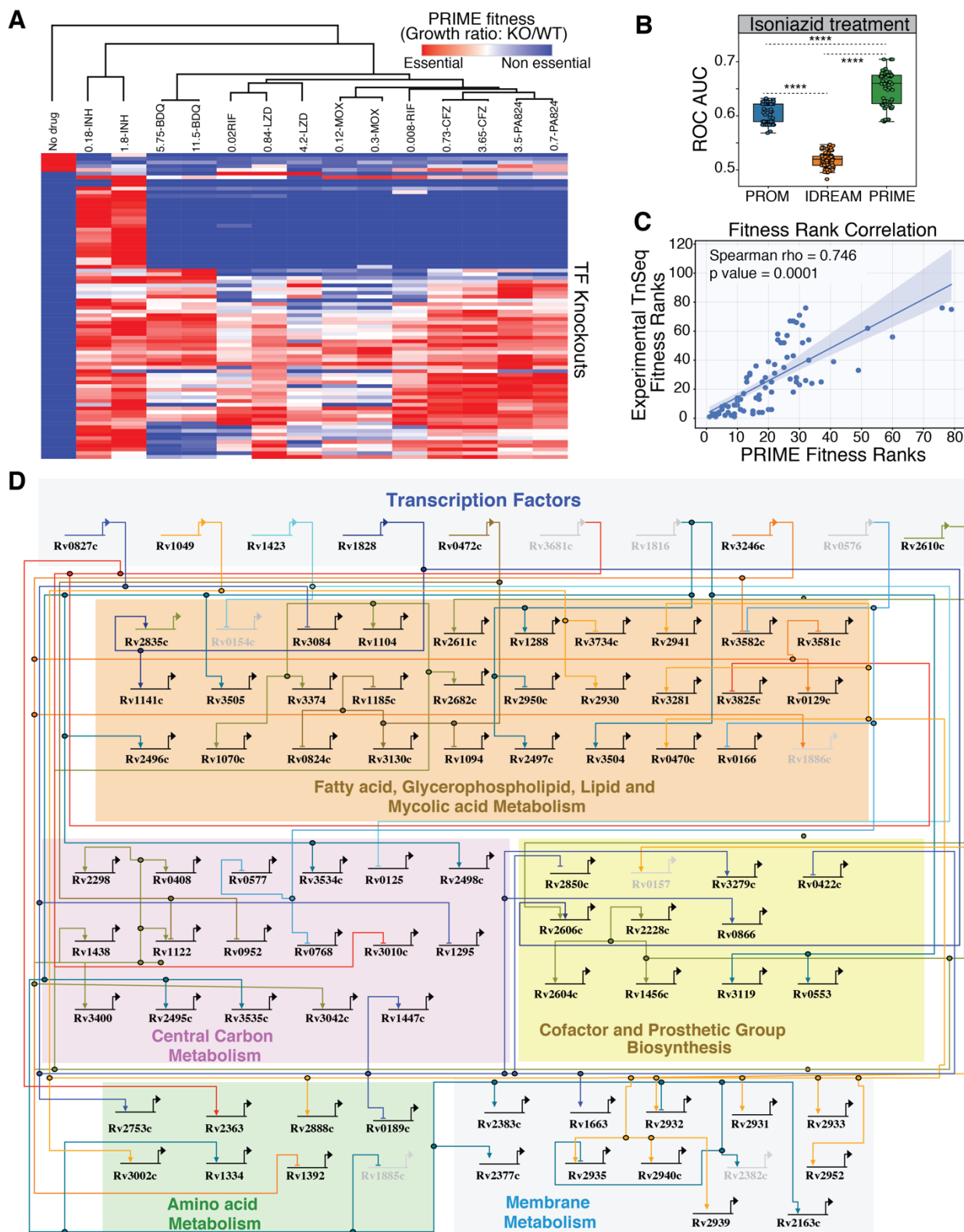
that are validated by experimental TnSeq data. **F.** The number of all correct PRIME predictions for

722

deletion of all genes in the metabolic network across the two conditions that are validated by



723 experimental TnSeq data. BioTapestry visualization showing a subset of the gene regulatory network of  
724 Mtb under growth in **G.** cholesterol and **H.** glycerol. TFs are grouped together in the top panel  
725 (represented by bent arrows), which extend to horizontal and vertical lines that connect to their  
726 regulatory gene targets. Highlighted TFs were predicted by the PRIME model to be essential and  
727 validated through TnSeq dataset in relevant conditions.  
728



729

730

731

**Figure 4. Drug-specific predictions of PRIME.** A. Heatmap of PRIME derived fitness for all TF knockouts in the presence of 7 primary drugs and control at 24 h. The numbers indicate the

732 concentration of drug used in  $\mu\text{g/mL}$ . INH: isoniazid, BDQ: bedaquiline, RIF: rifampicin, LZD: linezolid,  
733 MOX: moxifloxacin, CFZ: clofazamine, PA824: pretomanid. **B.** Sensitivity and specificity of PRIME,  
734 PROM, and IDREAM predicted TF essentiality in the presence of INH as determined by LOOCV analysis  
735 for the area under the receiver operating characteristic curve (ROC AUC). Statistical significance was  
736 calculated as  $p$ -value with two-sample t-test. \*\*\*\*:  $p$ -value < 0.0001. **C.** Correlation of TnSeq  
737 experimental fitness ranking of TFs and PRIME derived fitness ranks. **D.** BioTapestry visualization  
738 showing a subset of the gene regulatory network of Mtb with PRIME predictions during INH treatment.  
739 Some of the highlighted TFs were predicted as essential in the presence of INH (Rv0827c, Rv1049 and  
740 Rv0472c), while others were predicted essential in both the absence and presence of INH (Rv1423,  
741 Rv1828, Rv3246c, and Rv2610c). The lightened TFs were predicted essential in the untreated control  
742 but non-essential in the presence of INH (Rv3681c, Rv1816, and Rv0576). All of these PRIME  
743 predictions were validated by experimental fitness screening in relevant conditions.

744

745

746 **Tables:**

747

748 **Table 1.** Summary of PROM, IDREAM, and PRIME model features

Mtb Model Features	Chandrasekaran, 2010	Ma, 2015	Present Study		
	MTBPROM1.0	MTBPROM2.0	PROM*	IDREAM*	PRIME
Metabolic model	iNJ661	iSM810	iEK1011	iEK1011	iEK1011
Number of reactions	1025	938	1229	1229	1229
Number of metabolic genes in the metabolic network	661	810 (759 genes in iEK1011)	1011	1011	1011
Regulatory network	Balazsi 2008	Minch 2015	Minch 2015	EGRIN (FDR<0.25)	EGRIN (Precision=50%)
Number of transcription factors	30	104	104	142 <sup>#</sup>	142 <sup>#</sup>
Number of interactions	218	2555	2555	3643 <sup>#</sup>	4820 <sup>#</sup>
Number of genes in the regulatory network ( <b>metabolic / total</b> )	178 / 178	647 / 647	605 / 647	641 / 2487	750 / 2905

749 \*The PROM model was updated in this study by incorporating the latest metabolic network (MN) model for Mtb;  
750 the IDREAM model was constructed in this study to evaluate performance relative to the other methods

751 <sup>#</sup>PRIME uses the same EGRIN network as IDREAM, but incorporates the weights of regulation of each metabolic  
752 enzyme to update the constraint on reaction fluxes through the MN.

753