# Predicting bacterial fitness in Mycobacterium tuberculosis with transcriptional regulatory network-informed interpretable machine learning

1  **Ethan Bustad[1], Edson Petry[2][‡], Oliver Gu[2][‡], Braden T. Griebel[1,3], Tige R. Rustad, David R.**
2  **Sherman[4], Jason H. Yang[2,5][\*], Shuyi Ma[1,3,6,7][\*]**

3  [1] Center for Global Infectious Disease Research, Seattle Children's Research Institute, Seattle WA,
4  USA

5  [2] Center for Emerging and Re-emerging Pathogens, Rutgers New Jersey Medical School, Newark NJ,
6  USA

7  [3] Department of Chemical Engineering, University of Washington, Seattle WA, USA

8  [4] Department of Microbiology, University of Washington, Seattle WA, USA

9  [5] Department of Microbiology, Biochemistry, & Molecular Genetics, Rutgers New Jersey Medical
10  School, Newark NJ, USA

11  [6] Department of Pediatrics, University of Washington, Seattle WA, USA

12  [7] Pathobiology Graduate Program, Department of Global Health, University of Washington, Seattle
13  WA, USA

14  [‡] These authors contributed equally.

15  **\* Correspondence:**
16  Corresponding Authors
17  Shuyi Ma, shuyi.ma@seattlechildrens.org
18  Jason H. Yang, jason.y@rutgers.edu

19

## Abstract

24  *Mycobacterium tuberculosis* (Mtb) is the causative agent of tuberculosis disease, the greatest source
25  of global mortality by a bacterial pathogen. Mtb adapts and responds to diverse stresses such as
26  antibiotics by inducing transcriptional stress-response regulatory programs. Understanding how and
27  when these mycobacterial regulatory programs are activated could enable novel treatment strategies
28  for potentiating the efficacy of new and existing drugs. Here we sought to define and analyze Mtb
29  regulatory programs that modulate bacterial fitness. We assembled a large Mtb RNA expression
30  compendium and applied these to infer a comprehensive Mtb transcriptional regulatory network and
31  compute condition-specific transcription factor activity profiles. We utilized transcriptomic and
32  functional genomics data to train an interpretable machine learning model that can predict Mtb
33  fitness from transcription factor activity profiles. We demonstrated that this transcription factor
34  activity-based model can successfully predict Mtb growth arrest and growth resumption under

35　hypoxia and reaeration using only RNA-seq expression data as a starting point. These integrative
36　network modeling and machine learning analyses thus enable the prediction of mycobacterial fitness
37　under different environmental and genetic contexts. We envision these models can potentially inform
38　the future design of prognostic assays and therapeutic intervention that can cripple Mtb growth and
39　survival to cure tuberculosis disease.

40

## 41　1. Introduction

42　　*Mycobacterium tuberculosis* (Mtb) remains a supremely successful pathogen, sickening 10.6
43　million people and killing over 1 million people worldwide each year [1]. An important factor for
44　Mtb's success is its ability to adapt to a broad range of host-associated and treatment-associated
45　stresses. The mechanisms underlying how Mtb dynamically regulates its growth and physiology in
46　response to stress response remains incompletely understood. Characterizing the gene regulatory
47　activities of transcription factors (TFs) under different environmental or stress conditions could help
48　inform interventions that modulate Mtb growth and survival to cure tuberculosis disease.

49　　Several groups have previously performed analyses to characterize Mtb's transcriptional
50　regulatory network (TRN) using experimental and computational approaches [2; 3; 4; 5; 6; 7; 8; 9].
51　These efforts have largely relied on two strategies: 1) detailed profiling of the molecular impact of
52　individual transcription factors (TFs) with recombinant induction and disruption strains, and/or 2)
53　statistically informed TRN inference using data from large transcriptome compendia.

54　　In principle, TRNs can be empirically assembled from measurements of TF-DNA binding
55　activities and gene expression profiles from conditions with known individual TF perturbations.
56　These data would enable the inference of direct regulatory interactions between TFs and their
57　putative target genes, which exhibit altered transcriptional expression in response to TF perturbations
58　and provide evidence of TF binding events proximal to a gene. To leverage this strategy, we
59　previously engineered a library of Mtb recombinant TF induction (TFI) strains [2; 6], from which we
60　profiled transcriptomes in 208 TFI strains by microarray analyses (GSE59086, [6; 10]) and detected
61　~16,000 ChIP-seq binding events for 154 TFs (~80% of all Mtb TFs) and 2,843 genes (~70% of all
62　Mtb genes) [3; 10]. These detailed ChIP-seq and transcriptional profiles have yielded important
63　insights into the regulatory programs active during Mtb broth culture. However, these experiments
64　possessed several technical limitations. For example, our microarray profiling efforts were unable to
65　measure changes in expression for 1,190 genes (~30% of Mtb genes) [6], and our ChIP-seq profiling
66　efforts were unable to detect TF binding associated with 1,040 genes (~26% of Mtb genes) [3].
67　Moreover, the existing profiles have focused specifically on regulatory behavior of the Mtb
68　laboratory strain H37Rv in log-phase growth in 7H9 media. Consequently, condition-specific
69　interactions relevant to other environments or Mtb strains were not captured. Thus, despite such
70　efforts, significant gaps remain in the ability to identify TF-gene regulatory interactions directly and
71　comprehensively by only experimental activities.

72　　Bioinformatic network inference approaches that utilize expression compendia comprising
73　transcriptome responses under diverse biological conditions are a useful complementary strategy to
74　recombinant strain profiling. These statistically informed approaches enable assessment of regulatory
75　interactions across the multitude of conditions present in a transcriptome compendium. However,
76　these computational network inference strategies are constrained by two limitations. First, large and
77　biologically diverse gene expression data are needed to fuel identification of high-confidence

78  statistical associations between TFs and putative target genes [11]. To meet this need, compendia of
79  expression data may be curated from public microarray [4; 10] or RNA-seq [7; 12; 13] data. Second,
80  statistical learning network inference algorithms differ in the assumptions made on the training data
81  and on the interpretation of TF-gene associations. These assumptions are often biologically
82  inaccurate. We previously performed such analyses and were able to only infer 598 clusters of
83  coregulated gene expression for 3,922 genes [4]. Others recently performed similar analyses and
84  inferred either 80 clusters for 3,906 genes [7] or 560 co-regulated gene modules for 3,912 genes [5].
85  These models have successfully revealed novel regulatory interactions impacting Mtb stress
86  adaptation, but none of these regulatory models may be precisely interpreted as TF regulatory
87  programs (as they only capture a fraction of Mtb's 214 TFs) and none can be used to directly
88  estimate TF activities (i.e., the extent of regulation that each TF exerts on its regulated target genes,
89  TFAs, [14]) under different experimental conditions. TRN inference efforts in other microbes,
90  including the DREAM5 challenge for *E. coli* and *S. aureus* [15], have found that robust TRNs may
91  be assembled by aggregating the regulatory relationships inferred by different statistical algorithms.
92  We hypothesized that implementing a similar "wisdom of crowds" approach to aggregate
93  complementary TRNs inferred via different statistical approaches would yield a more comprehensive
94  and higher quality Mtb TRN.

95  Here we assembled a biologically diverse and batch corrected Mtb RNA-seq gene expression
96  compendium. We integrated this RNA-seq compendium with the perturbative TFI microarray dataset
97  to infer a comprehensive Mtb transcriptional regulatory network that included all 214 TFs and all
98  4,027 genes present in our RNA-seq expression compendium. We used this TRN to estimate TFA
99  profiles corresponding to individual RNA expression profiles. We used the TFAs calculated from our
100  RNA-seq compendium to train an interpretable machine learning regression model that could predict
101  growth phenotypes previously measured in TF-induced strains [16]. We demonstrated that this
102  regression model can accurately predict Mtb fitness under stressful environmental conditions such as
103  hypoxia.

## 2. Methods

### 2.1 TFI microarray expression compendium assembly and normalization

106  Microarray expression data corresponding to TFI strains were downloaded from GEO
107  (GSE59086). Groups were assigned to each sample by the identity of each strain. The Rv2160A gene
108  fully encompasses the Rv2160c gene, so the Rv2160A and Rv2160c samples were combined into a
109  single Rv2160 TFI strain group. This resulted in 208 TFI strain groups. These 208 strain groups
110  included Rv0560, Rv3164c, and Rv3692 which were considered hypothetical TFs in TFI strain
111  construction [6], but later determined to not be true Mtb TFs [10]. However, for the purpose of the
112  analyses presented here, each of these 208 strains will be referred to as TFs. Smooth quantile
113  normalization [17] was performed using *PySNAIL* [18] using the assigned group definitions.

### 2.2 RNA-seq expression compendium assembly, quality control, and normalization

115  The NCBI Sequence Read Archive (SRA) was queried with "*Mycobacterium tuberculosis*" for
116  RNA expression samples containing raw FASTQ sequencing reads. 3,506 FASTQ sequencing reads
117  were downloaded and combined with FASTQ sequencing reads from 398 unpublished RNA-seq
118  profiles generated by our labs. We aligned these sequencing reads against the NC_000962.3 Mtb
119  H37Rv reference genome using Bowtie 2 [19]. Read counts were compiled using *featureCounts* [20].
120  Samples with fewer than 400,000 total gene counts and samples duplicated in our preliminary
121  compendium were excluded from further analysis. Sequencing counts between samples were

3

122 normalized by transcripts per kilobase million (TPM). Group definitions were manually added to
123 represent unique experimental conditions from each set of experiments; biological replicates for each
124 experimental condition were given the same group definitions. Smooth quantile normalization [17]
125 was performed using *PySNAIL* [18] using the assigned group definitions. Quality data, adapter and
126 quality trimming statistics, and alignment and counts metrics were compiled and assessed using
127 *MultiQC* [21].

### 2.3 UMAP visualization and cluster estimation

129 RNA expression compendia and TFAs were visualized by Uniform Manifold Approximation &
130 Projection (UMAP) [22]. Clusters were estimated by *DBSCAN* [23]. The ε hyperparameter was
131 optimized for each dataset by varying ε across 50 logarithmically distributed values from 0.1 to 10
132 and selecting the value of the elbow of the ε vs. Number of Outliers plot. This selection delivers the
133 minimum number of clusters that maximizes inclusion of samples without overfitting the data
134 (**Supplementary Figure S1**). UMAP and DBSCAN analyses were performed in Python using their
135 implementations in *umap-learn* and *scikit-learn* [24].

### 2.4 Regulatory network inference methods

137 We implemented an ensemble of network inference methods by starting with a selection of
138 methods featured in the DREAM5 challenge [15]. These methods were selected based on diversity in
139 underlying statistical approach, predictive performance reported in the DREAM5 study, and the
140 availability of a working implementation. Our initial selection consisted of ARACNe [25; 26], CLR
141 [27], and GENIE3 [28]. We chose an ARACNe implementation that employs adaptive partitioning
142 for more efficient processing [25; 26]. We used an R implementation of CLR available on CRAN
143 from the *parmigene* package [29]. We used an R implementation of GENIE3 available on
144 BioConductor [30]. To supplement these methods, we incorporated two other more recent advances
145 in network inference approaches: cMonkey2 [31; 32] and iModulon [33]. We used a docker image
146 containing a Python implementation of cMonkey2, available at
147 https://hub.docker.com/r/weiju/cmonkey2. For iModulon, our desired output was different from the
148 output of this algorithm implemented by the original authors. We thus made a custom
149 implementation, borrowing heavily from https://github.com/SBRG/pymodulon and
150 https://github.com/SBRG/iModulonMiner, in Python. In addition, we also chose to implement a
151 regression strategy using Elastic Net regression, a more advanced technique than was used in
152 DREAM5. Elastic Net is a regularization method that takes advantage of the unique properties of
153 both the lasso (used extensively in DREAM5) and ridge regression [34]. Elastic Net performs better
154 than lasso or ridge regression when predictors may be correlated and under-determined [35]. We
155 modeled each gene individually on the expression of all the transcription factors, and used the
156 resulting coefficients to both select significant relationships and score those relationships; this
157 implementation was done in Python using *scikit-learn* [24]. Descriptions of each of these inference
158 methods are provided in **Supplementary Table 3**.

159 Each method was wrapped to produce a ranked list of putative TF regulator-target gene
160 relationships in order of the inferred strength of the regulatory relationship, from strongest to
161 weakest. Execution was done using docker images
162 (https://hub.docker.com/repositories/malabcgidr?search=network-inference). Auto-regulatory (self-
163 targeting) relationships were excluded. Method hyperparameters were chosen to match either original
164 publications or the DREAM5 challenge when possible. Execution for each method and optimization
165 of their corresponding hyperparameters was validated by testing against the evaluation scripts
166 provided in the supplemental material of [15; 32].

167   A network was generated for each combination of the two datasets (RNA-seq and TFI
168   microarray) and 6 inference methods, yielding 12 total constituent networks.

### 2.5  Inferred network truncation and aggregation

170   The constituent networks were large, as many of the network inference methods did not require a
171   cutoff threshold and did not perform multiple testing correction; the union of all inferred edges
172   constituted over 90% of the possible Mtb regulatory space (where 100% would be every TF
173   harboring a regulatory association with every Mtb gene). We therefore truncated each inferred
174   network to incorporate the unique perspective of each model without aggregating too many low-
175   confidence relationships. This was done by comparison with an independent validation set,
176   comprising a presumed unbiased sampling of the true population of regulatory relationships in Mtb.
177   This validation set was used to identify the extent of true positives in each network.

178   The validation data set was gleaned from Sanz et al., Material S1 [8]. The original list was
179   filtered for relationships whose supporting evidence included at least one high-confidence physical
180   methodology, namely values 4-9: LacZ-promoter fusion, GFP-promoter fusion, proteomic studies,
181   electrophoretic mobility shift assays (EMSA), one hybrid reporter system, and chip-on-chip. This
182   yielded a set of 433 high-confidence regulator-target relationships, including 51 regulators and 160
183   total target genes, that had little to no dependence on the transcriptional information used to build the
184   constituent networks.

185   A cutoff threshold was chosen for each network by binning the ranks of validation hits into 32
186   bins and truncating the network at the first bin where the number of hits fell below the expected level
187   of random overlap per bin. This level was calculated to equal the mean of a hypergeometric
188   distribution, with a population size equal to the total regulatory space of Mtb, a set of true positive
189   regulatory interactions identified by the Sanz validation set [8], and draws equal to the size of the
190   inferred network, taken without replacement. This shrunk each network to an average of about 10%
191   of its original size (3-28%) (**Supplementary Figure 2**). Three of the constituent networks displayed
192   insufficient enrichment against the validation dataset: ARACNe/TFI, cMonkey2/TFI, and
193   iModulon/TFI. Upon executing a Fisher's exact test to determine the chance of a random network
194   achieving the same enrichment, these three failed to pass a strict cutoff of 0.0001. They were thus
195   excluded from further aggregation.

196   The remaining truncated networks were then aggregated together, first into two combined
197   networks, one for each underlying input transcriptome dataset (RNA-seq compendium and TFI
198   microarray profile). Aggregation was performed by rank average as described in the DREAM5
199   challenge [15]. Repeating the enrichment analysis performed above, it was determined that the TFI
200   aggregate would benefit from additional truncation and was thus truncated using the same threshold
201   strategy described in the previous paragraph, whereas the RNA-seq network was already sufficiently
202   enriched. These two networks were then aggregated together again by rank average, yielding one
203   final aggregate network.

204   All these networks were validated against the Sanz et al. data set using the Matthews Correlation
205   Coefficient (MCC), as described previously [36; 37] (**Supplementary Figure 3**).

### 2.6 Principal Component Analysis

207   Principal component analysis (PCA) was performed on the inferred networks (after truncation),
208   the dataset-level aggregate networks, and the overall aggregate network, using the 16,792-

209 dimensional space represented by the ranks of edges shared across at least 3 of the inferred networks.
210 Any relevant edges not included in a given network were assigned a rank of 16,792, the size of the
211 space.

## 2.7 Regulatory directionality

213     The types of the regulatory connections (whether the TF up- or down-regulates the associated
214 gene) were explored using a combination of the regression models and measured TFI gene
215 expression values. Two elastic net models and two unpenalized linear models were used to infer
216 direction of regulation based on the sign of the regression coefficients, one of each for each dataset
217 (RNA-seq compendium and TFI microarray profile). We supplemented these regression associations
218 with the directionality of significant differential gene expression (i.e. upregulated vs. downregulated
219 expression) measured from the TFI microarray dataset. Linear models were fit in Python with the
220 *statsmodels* package. Coefficients with an FDR < 0.05 were selected as evidence. Elastic net models
221 with an $R^2$ of less than 0.8 were excluded; coefficients that were included by the remaining models
222 were selected as evidence. TFI differential expression from the microarray dataset was filtered using
223 an FDR < 0.05 and requiring at least 2-fold change in either direction. Elastic net models and TFI
224 differential expression were considered strong evidence, whereas the unpenalized linear models were
225 considered weak evidence. A flow chart depicting how the information from these models and
226 differential expression analyses were used to define up vs. down regulation is shown in
227 **Supplementary Figure 5**.

## 2.8 Comparing inferred networks against independent reference information

229     Additional orthogonal datasets were incorporated to corroborate the networks. All generated
230 networks were tested against a set of published ChIP-seq binding relationships gleaned from Minch,
231 et al. [3]. We took the intersection of their sets of statistically significant peaks (Supplementary Data
232 1 from [3]) and peaks in a canonical promoter region (Supplementary Data 3 from [3]) to yield 5,178
233 relationships, including 129 regulators and 2,271 total targets. The MCC was then calculated against
234 this data set for each network.

235     Gene ontology enrichment analysis was then performed to ascertain the extent to which TF
236 targeting could be used to gauge biological function within each group [38; 39]. For each TF, each
237 set of genes that our network identified as upregulated, downregulated, or regulated in both directions
238 by the regulator was analyzed for GO enrichment at an FDR < 0.05. All identified GO annotations
239 that had a child annotation also identified for a given TF were removed for the sake of simplicity
240 (**Supplementary Table 5C**). Results were filtered to regulators receiving at least 3 significant GO
241 enrichments for further manual inspection and analysis (**Supplementary Tables 5A, 5B**), and those
242 TFs with an annotated name and considered to have a testably specific functional role listed in the
243 Mycobrowser annotation [40] were juxtaposed for network validation (**Table 1**). GO analysis was
244 performed in Python using the *goatools* package [41]. Gene ontology data was taken from the 2024-
245 06-17 release of go-basic.obo from the Gene Ontology knowledgebase [42]
246 (https://purl.obolibrary.org/obo/go/releases/2024-06-17/go-basic.obo), and mappings to Mtb genes
247 were taken from the European Bioinformatics Institute GOA project, release 20240805
248 (https://ftp.ebi.ac.uk/pub/databases/GO/goa/proteomes/30.M_tuberculosis_ATCC_25618.goa).

## 2.9 Calculating transcription factor activity profiles from network component analysis

250     Transcription factor activities for each expression profile was computed using Robust Network
251 Component Analysis (ROBNCA) [43]. ROBNCA was implemented in Python, using code adapted

252  from
253  https://github.com/CovertLab/WholeCellEcoliRelease/tree/00cf7738cb8379c14d65ef632b2156bdf7c
254  23434/reconstruction/ecoli/scripts/nca [44].

## 2.10 Associating network activity with bacterial fitness

256  We built a model associating mycobacterial growth with TF activity, as inferred from measured
257  gene expression data. The GSE59086 microarray dataset was again used as a broad measure of TFI
258  conditions, with relative growth data for 194 matching TFI conditions added from Ma et al., 2021,
259  Table S1 as training data [16]. Expression levels in the form of log-2 fold-change were transformed
260  into putative TFAs using the control strengths calculated via NCA from the aggregate network and
261  RNA-seq compendium. A gradient boosted machine (GBM) model was trained to regress growth on
262  TFAs, using a grid search cross-validation scheme to optimize hyperparameters based on bounds
263  derived from [34], using the number of estimators to reward better performing models. The number
264  of estimators was then optimized with a simple grid search. The model was implemented in Python
265  using the *lightgbm* package [45; 46].

## 2.11 Hypoxia time-course experiment

267  Wildtype H37Rv (ATCC 27294) and H37Rv transformed with a control anhydrotetracycline
268  (ATc)-inducible expression vector (H37Rv::pEXCF-empty, which does not induce recombinant gene
269  expression) were cultured under in Middlebrook 7H9 with the oleic acid, bovine albumin, dextrose,
270  and catalase (OADC) supplement (Difco) and with 0.05% Tween 80 at 37°C. H37Rv::pEXCF-empty
271  was grown with the addition of 50 µg/ml hygromycin B to maintain the plasmid and induced with
272  100ng/mL ATc one day prior to onset of hypoxia. For hypoxia, strains were cultured in oxygen-
273  limited conditions (1% aerobic $O_2$ tension) for 7 days, followed by reaeration on day 7-12, initiated
274  by transferring cultures into continuously rolled bottles with 5:1 head space ratio using methods
275  described previously [2; 47; 48; 49]. Bacterial survival and growth were enumerated by plating for
276  colony forming units (CFU) on Middlebrook 7H10 solid media plates using standard microbiological
277  methods.

278  Transcriptomes were generated by RNA-seq from bacterial cultures sampled from the
279  aforementioned conditions using methods described previously [50]. Briefly, bacterial pellets
280  suspended in TRIzol were transferred to a tube containing Lysing Matrix B (QBiogene) and
281  vigorously shaken in a homogenizer. The mixture was centrifuged, and RNA was extracted from the
282  supernatant with chloroform, followed by RNA precipitation by isopropanol and high-salt solution
283  (0.8 M Na citrate, 1.2 M NaCl). Total RNA was purified using a RNeasy kit following the
284  manufacturer's recommendations (Qiagen). rRNA was depleted from samples using the RiboZero
285  rRNA removal (bacteria) magnetic kit (Illumina Inc., San Diego, CA). Illumina sequencing libraries
286  were prepared from the resulting samples using the NEBNext Ultra RNA Library Prep kit for
287  Illumina (New England Biolabs, Ipswich, MA) according to the manufacturer's instructions, and
288  using the AMPure XP reagent (Agencourt Bioscience Corporation, Beverly, MA) for size selection
289  and cleanup of adaptor-ligated DNA. We used the NEBNext Multiplex Oligos for Illumina (Dual
290  Index Primers Set 1) to barcode the libraries to enable sample multiplexing per sequencing run. The
291  prepared libraries were quantified using the Kapa quantitative PCR (qPCR) quantification kit and
292  sequenced at the University of Washington Northwest Genomics Center with the Illumina NextSeq
293  500 High Output v2 kit (Illumina Inc., San Diego, CA). The sequencing run generated an average of
294  75 million base-pair paired-end raw read counts per library. Read alignment and gene expression
295  estimation was carried out using a custom processing pipeline in R that harnesses the Bowtie 2

296  utilities [19; 51], which is publicly accessible at
297  https://github.com/robertdouglasmorrison/DuffyTools, and
298  https://github.com/robertdouglasmorrison/DuffyNGS.

299  Gene expression data were transformed from log-2 fold-change to putative TFAs using the
300  control strengths calculated via NCA above and run through the GBM model to predict relative
301  fitness level of the Mtb culture as it progressed through the hypoxia time-course.

**2.12 False discovery rate correction**

303  False discovery rate correction was performed using the two-stage Benjamini-Krieger-Yekutieli
304  method [52].

## 3. Results

**3.1 Generation of a large and biologically diverse Mtb gene expression compendium for TRN inference**

309  Our previous attempts at TRN characterization utilized microarray expression profiles from
310  recombinant TFI strains as perturbative training data (GSE59086, [6]). However, while this dataset
311  enabled detailed characterization of transcriptional regulation of Mtb physiology during log-phase
312  broth culture, it possessed poor biological diversity. UMAP and DBSCAN analyses reveal that
313  expression profiles from these 698 microarray experiments and 208 TFI conditions only yielded 16
314  clusters of expression profiles (**Figure 1A**). This poor diversity likely arises from the original
315  experimental design for these data, in which each TFI strain was grown to log-phase in albumin-
316  dextrose-catalase (ADC)-supplemented 7H9 media before isolating RNA. UMAP and DBSCAN
317  analyses suggested that this TFI microarray dataset alone would be insufficient for predicting TFAs
318  corresponding to diverse experimental conditions. Moreover, microarray technologies have poor
319  sensitivity and dynamic range for quantifying gene expression [53]. We found that 101 genes in this
320  dataset did not possess expression measurements greater than 10 counts, indicating poor detection or
321  poor evidence for expression in these experiments (**Figure 1B**). In addition, the median absolute
322  deviation (MAD) was small ($< 1$) for nearly all genes, indicating the ability to detect gene expression
323  changes across conditions was limited. These analyses collectively motivated the need to assemble a
324  new RNA expression compendium.

325  We therefore collected samples from the NCBI Sequence Read Archive (SRA) and our own labs,
326  aligned, filtered, normalized, and batch corrected by smooth quantile normalization [17; 18] (see
327  **Methods** for details). Batch correction is an important pre-processing step for unifying data from
328  different sources that is frequently overlooked in Mtb RNA expression compendium analyses [4; 7;
329  12; 13]. After performing these pre-processing steps, our final compendium comprised 3,496 RNA-
330  seq samples from 1,288 experimental conditions (**Supplementary Table 1**). Expression counts for
331  the RNA-seq compendium can be queried at https://tfnetwork.streamlit.app/.

332  UMAP and DBSCAN analyses of the batch corrected RNA-seq expression compendium
333  validated its biological diversity (**Figure 1C**-**D, Supplementary Table 2**). We identified 142 unique
334  expression clusters. This RNA-seq transcriptome compendium exhibited significantly greater
335  dynamic range and variation in gene expression than in the TFI microarray dataset (**Figure 1D**). Of
336  note, genes with high variation (high MAD) were mostly well-characterized stress response genes

8

337    (e.g., Rv2031c *(hspX)*, Rv2626c *(hrp1)*, and Rv2623 *(TB31.7)*), with Rv2007c *(fdxA)* having higher
338    variation than the commonly studied Rv3133c *(devR)* stress response regulator. These are consistent
339    with expectation, as most stress response genes would be expected to only be induced in the presence
340    of their specific stressor.

## 3.2 Inferred transcriptional regulatory network interactions enrich for shared functional processes

343    Network inference studies in other bacteria have shown that combining regulatory interactions
344    from multiple different inference algorithms results in a TRN that outperform networks generated by
345    a single method [15]. To more comprehensively characterize Mtb regulatory interactions, we applied
346    a "wisdom of crowds" ensemble inference approach. We first applied a collection of regulatory
347    network inference tools to generate TRN models using individual methods (see **Methods**). These
348    tools were selected because they have been shown to be sensitive to distinct types of regulatory
349    relationships in other bacteria [15] or they have previously been successfully applied to infer
350    regulatory relationships in Mtb [4; 5; 7]. To further diversify the regulatory relationships inferred
351    from these approaches, we applied these tools to both our assembled RNA-seq compendium as well
352    as the TFI microarray dataset. Collectively, these inference activities yielded 12 networks that
353    describe 779,213 unique interactions between 214 regulators and 4,029 target genes. We truncated
354    these networks using a benchmark dataset of high confidence regulatory interactions with
355    biochemical evidence that was curated by Sanz et al. [8] (see **Methods**). We used this high
356    confidence regulatory interaction dataset to inform pruning of low-confidence regulatory
357    relationships inferred from each of the individual inference methods (**Supplementary Figure 2**),
358    yielding a shorter, more high-confidence network for each method. Principal component analysis of
359    these networks revealed substantial diversity in the regulatory interactions identified between the
360    different approaches applied to the two source datasets (**Figure 2B**).

361    We rank-aggregated the resulting 12 networks to consolidate regulatory relationships across the
362    individual inference methods. The resulting aggregate network has 68,226 regulatory interactions that
363    connect 214 transcriptional regulators with 4,027 target genes. Of these interactions, 37,236 are
364    associated with transcriptional activation across conditions, 15,820 interactions are associated with
365    transcriptional repression across conditions, 1,496 relationships are predicted to be either activating
366    or repressing, depending on the environmental condition, and 11,766 regulatory relationships have an
367    undetermined regulatory directionality (**Supplementary Table 4**). These interactions represent both
368    direct, biophysical regulatory events as well as indirect regulatory relationships mediated by
369    downstream regulators. These interactions also represent the union of regulatory relationships that are
370    active in at least a subset of all the different environmental conditions profiled in our assembled
371    source RNA-seq compendium and TF induction profiling datasets. Notably, not all these regulatory
372    relationships will be active under all environmental conditions. The distribution of regulatory
373    interactions per TF largely follows a power law distribution consistent with the scale free networks
374    found to represent transcriptional regulation in other bacteria (**Supplementary Figure 4**). We found
375    a deviation between the distribution of our aggregate network and the expected power law
376    distribution for regulators with relatively few target genes. This is likely due to the inclusion of
377    indirect regulatory relationships and relationships that are active under some but not all
378    environmental conditions. The networks can be viewed at https://tfnetwork.streamlit.app/, and the
379    TF-gene interactions are described in **Supplementary Table 4**.

380    To validate the connectivity of our aggregate network, we benchmarked it against experimentally
381    profiled TF binding data we previously profiled by ChIP-seq in the TFI strains under log-phase broth

382  culture [3]. To assemble a high-confidence regulatory association dataset, we included only
383  significant ChIP-seq peaks associated with TF binding in the promoter region of target genes. We
384  evaluated overlap between this high-confidence ChIP-seq regulatory interaction dataset and our
385  inferred regulatory networks with the Matthews correlation coefficient (MCC). We find that most of
386  the inferred networks that we generated had significant MCCs, and that the aggregate network
387  outperforms the majority of inferred networks using individual methods (**Supplementary Figure 3**),
388  whilst still retaining a large number of regulatory relationships (most of the better performing
389  individual inference networks have relatively few regulatory interactions).

390  We also assessed the extent to which the regulatory relationships captured by our aggregate
391  network preserved biological functional relationships between the regulating TFs and the target
392  genes. For TFs with clear literature characterization of its function, we found a high degree of
393  correspondence with the gene ontologies and annotated functions of its regulated target genes (**Table
394  1, Supplementary Table 5**). For example, Rv3574 (*kstR*) is a TF that has been linked to regulating
395  cholesterol metabolism [54], and the target genes associated with *kstR* in our aggregate network also
396  have gene ontology annotations linked to cholesterol metabolism (**Table 1**). Additionally, toxin-
397  antitoxin target genes were enriched for growth regulation, highlighting that the regulatory
398  relationships captured by the aggregate network include indirect regulatory relations. Collectively,
399  this suggests the significant ontology and functional annotation enrichments made for genes and TFs
400  that are currently poorly annotated represent testable hypotheses for function – this is one of the
401  major advances from the aggregate network.

### 3.3 Network component analyses reveal per-sample Mtb TF activities under different conditions

404  Understanding when TFs are actively exerting their regulatory influence on their target genes can
405  reveal mechanistic insights into bacterial physiology and stress response. Network component
406  analysis (NCA) is an efficient way of estimating these TFA profiles from expression data by using a
407  TRN to perform matrix decomposition [14]. Robust NCA (ROBNCA) is a variant of NCA that
408  improves the performance of NCA calculations on noisy data with outlier measurements [43]. We
409  applied ROBNCA to estimate TFAs corresponding to each sample in our TFI microarray and RNA-
410  seq compendium.

411  To first determine and validate the ROBNCA TFA estimation approach on our data, we
412  performed ROBNCA on the TFI microarray data using the aggregate network inferred only from the
413  TFI data, as well as on 10 randomized networks to be used as negative controls. We hypothesized
414  that if the estimated TFAs represent true TF activities, with high TFAs indicating strong net activator
415  activity and low TFAs indicating strong net repressor activity, then the percentile ranks of TFAs for
416  highly expressed TFs should be either very high or very low in their corresponding TFI strains. On
417  the other hand, if the ROBNCA-calculated TFAs were spurious, then the TFA percentile ranks
418  should be statistically indistinguishable from the TFA percentile ranks from randomized networks.

419  For each of the 208 TFI strains within the microarray expression dataset, we averaged the TFAs
420  for all TFs across their biological replicates. We rank ordered TFs by their activities for each TFI
421  strain, calculated the rank percentile activity of the induced TF for each TFI strain, and analyzed the
422  distribution of these percentiles (**Figure 3A**). For the TFI microarray network, 31 TFs were ranked in
423  the highest or lowest 15% of TFA ranks (greater than 1 standard deviation from the mean), implying
424  that these TFs were the dominant regulators active in their respective TFI strain profiling condition.
425  Interestingly, 91 TFs had TFAs in the middle 30% from 35-65%. These TFs were fairly uniformly
426  distributed suggesting their related transcriptional programs were likely cross-regulated by other TFs.

10

427 Importantly, this suggested that induction of TF expression alone may be insufficient for fully
428 inducing some transcriptional programs, thus supporting the use of TFAs over untransformed gene
429 expression for downstream analysis.

430 We performed similar calculations for each of the randomized networks (**Supplementary Figure**
431 **6**) and averaged the TFA rank percentiles for all TFs from each randomized network (**Figure 3B**).
432 We found that there were significantly fewer TFAs in the highest or lowest 15% of TFA ranks in
433 these randomized networks than the TFAs calculated from the TFI expression dataset (p = 1.66e-49,
434 z-test [55]). Similarly, there were significantly more TFAs in the middle 30% (p = 1.66e-49, z-test
435 [55]). These differences between the ROBNCA-calculated TFA percentile distributions between TFI
436 and randomized networks indicated that the TFAs estimated by ROBNCA were not spurious and
437 likely reported on true biological condition-specific activities.

438 We next applied ROBNCA to our RNA-seq compendium using the TRN inferred from the RNA-
439 seq compendium. UMAP and DBSCAN analyses revealed that the level of biological diversity of
440 ROBNCA-predicted TFAs was similar to the diversity within the expression compendium, with 112
441 clusters of TFAs across the 3,496 samples (versus 142 for untransformed expression; **Figure 3C**).
442 Amongst the TFs with the highest level of median activity were the essential nitric oxide-sensing
443 Rv3219 (*whiB1*), histone-like protein Rv2986c (*hupB*), and sigma factor Rv2703 (*sigA*) (**Figure 3D**).
444 Each of these would be expected to be constitutively active in live Mtb cells. Also consistent with
445 expectation, the well-characterized stress response regulators Rv3133c (*devR*), Rv1994c (*cmtR*),
446 Rv0827c (*kmtR*) and two-component system regulators Rv0602c (*tcrA*) and Rv0981 (*mprA*) were
447 amongst the TFs with the highest TFA MAD.

448 Interestingly, the distribution of TFAs appeared different from the distribution of TF expression
449 levels measured for each RNA-seq sample across the compendium (**Figure 3E**). We tested the
450 correlation of expression level vs. activity for each TF across the entire compendium and found that
451 expression and activity were only moderately correlated across the dataset (Pearson's r = 0.48 ± 0.16
452 median ± MAD) (**Figure 3F**). 31 TFs were strongly correlated (|Pearson's r| ≥ 0.7), 66 TFs were
453 moderately correlated (0.7 > |r| ≥ 0.5), and 61 TFs were weakly correlated (0.5 > |r| ≥ 0.3). Relatedly,
454 both median and MAD expression and activity were only weakly correlated across all TFs (median: r
455 = 0.43; MAD: r = 0.32). These analyses further support our observation that TF expression level is
456 not the sole determinant for TFAs for most TFs. Rather, expression and activity convey two distinct
457 but complementary insights into transcriptional regulation, highlighting the importance of accounting
458 for network interactions when investigating transcriptional regulation. In particular, we posit that TFs
459 with weak correlation between expression and activity may require allosteric or other post-
460 translational modification to trigger activation of transcriptional regulation. This hypothesis can be
461 tested in future studies.

462 **3.4 Transcription factor activity profiles can predict condition-specific bacterial fitness**

463 Because transcriptional regulation plays important roles in coordinating Mtb growth adaptations
464 under stress, we asked whether our regulatory network models could be used to predict fitness
465 consequences of TF regulatory activities. To test this hypothesis, we utilized gradient boosting
466 machine learning to construct an interpretable TFA regression model designed to predict the fitness
467 of each TFI strain during log-phase culture based on each strain's calculated TFA profiles. We
468 trained this model using the TFAs computed by ROBNCA from the RNA-seq compendium, paired
469 with TFI fitness measurements that we previously collected in a Transcriptional Regulator Induced

11

470 Phenotype (TRIP) screen [16]. This TFA–fitness regression model was able to explain 87% of the
471 observed variation of growth between the TFI strains in the TRIP screen (**Supplementary Figure 7**).

472      To determine if this TFA–fitness regression model could predict changes in Mtb fitness or growth
473 from new data that were not used to train the model (e.g., under differing experimental conditions),
474 we generated fitness predictions with our model using transcriptomes that we profiled from Mtb cells
475 undergoing hypoxia and reaeration stress. From the TFA profiles calculated for cells exposed to
476 hypoxia, the TFA–fitness regression model predicted a significant decrease in growth that persisted
477 for each of the timepoints profiled under hypoxia (**Figure 4A, Supplementary Figure 8**). From the
478 TFA profiles calculated for cells under reaeration, the model predicted a rebound in Mtb growth
479 comparable to growth levels experimentally measured during log-phase culture. The kinetics of the
480 shifts in growth predicted by the TFA–fitness regression model aligned well with the experimental
481 measurements of Mtb bacteriostasis in hypoxia, followed by growth during reaeration (**Figure 4A,**
482 **Supplementary Figure 8**). Importantly, the experimental growth data from the hypoxia-reaeration
483 time course aligned better with the predictions from the TFA regression model than from an
484 analogous regression model trained from TF expression data alone (**Supplementary Figure 10**).
485 These results further support our premise that TFAs more effectively capture condition-specific
486 transcriptional regulation than TF expression alone and implies that the activation and regulation of
487 transcriptional programs under hypoxia and reaeration may involve allosteric or other post-
488 transcriptional mechanisms.

489      Because the TFA–fitness regression model is openly interpretable, we examined which TFAs
490 most strongly predicted the fitness changes under hypoxia and reaeration. We found that our TFA–
491 fitness regression model predicts that growth restriction during hypoxia is primarily driven by the
492 activities of 7 TFs whose TFA profiles changed significantly during hypoxia (**Figure 4B**).
493 Importantly, each of these TFs have direct or indirect links to hypoxia in the literature
494 (**Supplementary Figure 9**, **Supplementary Table 7**), thus further validating these model predictions
495 and the use of TFAs as a lens into condition-specific stress response biology.

496

## 4. Discussion

498      Understanding the molecular drivers of phenotypic changes in an organism is a fundamental goal
499 of biological research. In this study, we applied machine learning approaches to construct an
500 interpretable TFA–fitness regression model that can utilize Mtb TRNs to predict experimentally
501 measured changes in Mtb growth state in diverse environmental conditions. Our models build upon
502 existing experimental profiling and network inference modeling efforts to characterize Mtb
503 transcriptional regulation by integrating the data and algorithms developed in these prior studies [2;
504 3; 4; 5; 6; 7; 14; 15; 43]. Moreover, by integrating Mtb fitness profiling data from TRIP, our models
505 have also enabled direct prediction of growth/survival phenotypic outcomes from condition-specific
506 gene expression data inputs.

507      Our "wisdom of crowds" approach for inferring transcriptional regulatory interactions yielded
508 significant enrichment of known regulatory relationships while also expanding the scope of
509 represented experimental conditions. Our resulting TRN is substantially larger than the networks
510 inferred by individual algorithms, while enriched for experimentally validated interactions. This
511 highlights the utility of ensemble inference algorithms, as has been previously shown for regulatory
512 network inference in other bacteria [15].

513      Importantly, our results demonstrate how network models can generate hypotheses on gene
514    function in at least two complementary ways. First, we show by gene ontology enrichment analysis
515    that there is significant correlation between the annotated function of a TF's target genes and the
516    condition-specific regulatory function of the TF. It is important to note that the regulatory
517    interactions identified by our aggregate TRN includes both direct regulatory interactions involving
518    physical interactions between a TF and its target gene as well as indirect associations mediated by
519    other factors. Both direct and indirect regulatory associations are important for coordinating changes
520    in bacterial physiology [56], so it is expected that both types of interactions share annotated
521    ontologies. Because ~25% of Mtb genes lack functional annotation [57], we think the regulatory
522    relationships identified in our TRN can aid basic microbiological efforts in investigating Mtb gene
523    function by generating hypotheses for the functions of these poorly characterized or unknown genes
524    (**Supplementary Table 5**).

525      Second, we show that TFA regression models can be trained to link condition-specific TFAs with
526    TF fitness in log-phase broth culture to predict Mtb fitness under stress. Notably, we show that our
527    TFA regression model was able to predict Mtb growth and bacteriostasis under hypoxia and
528    reaeration – environmental conditions not used in training the TFA regression model. Our results
529    biologically suggest that TFAs are a useful determinant of condition-specific changes in bacterial
530    growth, and that the estimated TFA is more predictive of growth phenotypes that TF expression
531    alone. This is consistent with expectation as Mtb uses transcriptional regulation to orchestrate
532    behavioral adaptations to varying environments, including in growth phenotypes. Our modeling also
533    enables inspection of which TFAs are driving the predicted bacterial fitness outcomes. This can
534    inform the generation of hypotheses on the mechanisms underlying how TFs and their corresponding
535    transcriptional programs are activated (e.g., via allosteric mechanisms and/or network interactions).
536    Our TRN and TFA–fitness models could potentially inform the identification of regulatory
537    mechanisms mediating Mtb response and adaptation to other clinically relevant stress conditions
538    where gene expression profiling data are available. The TFs and target genes highlighted by these
539    models may potentially represent future intervention targets aimed at modulating Mtb fitness in a
540    therapeutically beneficial way. In light of the growing crisis of antimicrobial resistance [58] and
541    multi- and extensively-drug-resistant tuberculosis [59], we think our approach will be important for
542    curing tuberculosis disease [60].

543      More broadly, our work here demonstrates how network models can be utilized for biologically
544    meaningful interpretable machine learning applications. A fundamental challenge in current machine
545    learning activities is the difficulty in understanding how a trained machine learning model makes
546    predictions [61; 62]. We previously demonstrated that machine learning regression models can be
547    used to elucidate metabolic mechanisms underlying antibiotic lethality in *E. coli* [63], as well as to
548    predict multidrug interaction outcomes in Mtb [50]. Our study here analogously extends this
549    approach by training a regression model on TFAs estimated from TRN analyses to predict changes in
550    Mtb growth state. The advantage of this strategy over other contemporary machine learning
551    approaches is the direct utilization of prior knowledge encompassed by biological network models,
552    which directly enable the generation of hypotheses for mechanisms linking network interactions to
553    cell phenotypes. These hypotheses can then be experimentally tested [50; 63] and used as the basis
554    for further mechanistic study [64] and investigation of translational potential.

555      Looking forward, we envision that this approach and our TFA regression model can be useful for
556    several facets of tuberculosis research. We demonstrated that our model can be used to predict
557    changes in Mtb growth state under environmental stress, which may inform the design of growth
558    state assays under conditions where standard microbiological tools are not feasible. There is

13

559 increasing appreciation that Mtb drug susceptibility is regulated by its environment [65; 66]. Our
560 TFA–fitness regression model can be used to elucidate the molecular mechanisms underlying these
561 phenotypes. Moreover, functional genetic datasets are becoming increasingly available using
562 different technologies [16; 67; 68; 69; 70; 71; 72]. These data can be applied to train next-generation
563 TFA–fitness regression models with improved predictive power. Finally, detailed characterizations of
564 Mtb clinical strains are now providing significant insights into the how mutations or other forms of
565 genomic diversity regulate drug susceptibility in human patients [72; 73; 74; 75]. We envision the
566 TRN and TFA-fitness regression framework established here can be extended not only to study the
567 mechanistic basis for differences between drug susceptibility amongst clinical isolates, but also to
568 anticipate the drug susceptibility of new clinical isolates as they become curated.

569

## Conflict of Interest

571 The authors declare that the research was conducted in the absence of any commercial or financial
572 relationships that could be construed as a potential conflict of interest.

## Author Contributions

574 E.B.: Formal Analysis, Investigation, Methodology, Software, Validation, Visualization, Writing –
575 original draft, Writing – review, editing; E.P.: Data curation, Formal analysis, Visualization, Writing
576 – review, editing; O.G.: Formal Analysis, Investigation, Methodology, Software, Validation,
577 Visualization, Writing – review, editing; B.T.G.: Data curation, Software, Visualization, Writing –
578 review, editing; T.R.R.: Investigation, Resources, Methodology, Writing – review, editing; D.R.S.:
579 Investigation, Resources, Methodology, Funding acquisition, Supervision, Writing – review, editing;
580 J.H.Y.: Conceptualization, Funding acquisition, Investigation, Methodology, Project administration,
581 Resources, Supervision, Visualization, Validation, Formal Analysis, Writing – original draft, Writing
582 – review, editing; S.M.: Conceptualization, Funding acquisition, Investigation, Methodology, Project
583 administration, Resources, Supervision, Visualization, Formal Analysis, Writing – original draft,
584 Writing – review, editing.

595

596

14

## Figures

**Figure 1**: A biologically diverse Mtb RNA expression compendium. (A) UMAP visualization of biological diversity in the TFI microarray data. TFI data were batch corrected by smooth quantile normalization before computing the UMAP. Density-based spatial clustering (DBSCAN) was performed on the UMAP to identify clusters of samples with similar gene expression. UMAP and DBSCAN analyses revealed 16 total expression clusters in the TFI dataset. (B) Median vs. median absolute deviation (MAD) plot of expression for each gene across the TFI dataset. Each point represents a gene. Median expression and MAD were calculated for each gene across the 698 samples. Colors reveal point density (yellow: high density, blue: low density). (C) UMAP visualization of samples from the normalized and batch corrected RNA-seq compendium determined by gene expression. UMAP and DBSCAN analyses reveal 142 clusters of samples with similar gene expression. (D) Median vs MAD plot of expression for each gene across the RNA-seq compendium.

**Figure 2**: Overview of aggregate network. **(A)** PCA was performed on each of the generated networks. The networks inferred from the RNASeq compendium (triangle symbols) cluster to the right, whereas the networks inferred from the recombinant TF induction transcriptomes (x symbols) fall to the left. The dataset-level aggregates each cluster loosely with the same-dataset constituent networks at the horizontal extremes, whereas the overall aggregate falls near the centroid of all networks. **(B)** Performance of each inferred and aggregate network, calculated against a set of TF–target gene relationships defined by a ChIP-Seq DNA-binding investigation of recombinant TFI strains [3], as measured by Matthews correlation coefficient (MCC). MCC quantifies the level of correlation between the two sets, with higher values indicating more correspondence. Blue bars depict the MCC for aggregate networks; the other colors depict the MCC for the individual inferred networks. Hatched bars indicate networks that were excluded from aggregation. The horizontal dashed line represents the 95th percentile MCC performance of 1000 randomly generated networks. Note that the excluded iModulon/TF induction network scores relatively highly by this metric, likely because of its size (~7k edges, versus an average of ~180k). See Methods for information about the exclusion criteria.

**Figure 3**: Compendium-wide transcription factor activities. (A) Distribution of the TFA rank percentiles for each induced TF in each strain from the TFI microarray dataset. ROBNCA was applied to the TFI microarray dataset using the network specifically inferred from the TFI dataset. For each sample, rank percentiles were computed for each TFA. TFAs were averaged across biological replicates for each TFI strain. Histogram depicts the percentile rank for TFAs corresponding to the over-expressed gene in each TFI strain. (B) Averaged distribution of TFA percentile ranks from ROBNCA using 10 randomized networks (**Supplementary Figure 6**). (C) UMAP visualization of samples from the normalized and batch corrected RNA-seq compendium as determined by TFA. UMAP and DBSCAN analyses reveal 112 clusters of samples with similar TFAs. (D) Median vs. MAD plot of activity for each TF across the RNA-seq compendium. (E) Median vs. MAD plot of expression for each TF across the RNA-seq compendium. (F) Distribution of Pearson's correlation coefficients between expression and activity for each TF across the RNA-seq compendium.

**Figure 4**: Machine learning model insights into Mtb growth through a hypoxic time-course. **(A)** *Top*: When Mtb grown for two days in log phase was subjected to hypoxic conditions (starting from day 0), the bacteria stopped growing for the duration of the imposed hypoxia, as indicated by the stable CFU between day 0 and day 7. When the culture was reintroduced to oxygen ("Reaeration", starting from day 7), the bacteria resumed growth, as indicated by significantly higher CFU after day 8.

15

642    *Bottom*: Our GBM model predicted a decrease in growth over the course of the period of hypoxia,
643    and an increase in growth again upon reaeration, based only on transcriptional data measured over
644    the course of the experiment. Each point represents an RNA-seq timepoint. **(B)** The GBM model can
645    be interrogated to determine the primary drivers of the phenotype it predicts; when comparing the
646    most impactful TFAs in hypoxic conditions (days 2-7) versus those in reliably reaerated conditions
647    (days 9-12), 7 TFs were predicted to be particularly influential to the reduced growth in hypoxia
648    versus reaeration, each contributing at least 5% of the total absolute impact predicted by the model.
649    Shown here is the mean TFA change for each of the impactful TFs across days 2-7; other TFAs show
650    no net activity change overall (see Methods for details on TFA change calculation).

651

652 **Tables**

653

654 Table 1. Network regulators: annotation versus gene set enrichment analysis of inferred regulon.

| Regulator | Name | Mycobrowser gene product and function information | Inferred Regulon GO Annots. (FDR <0.05) | |
|---|---|---|---|---|
| | | | # | Summary |
| Rv0353 | hspR | Probable MerR family heat shock protein transcriptional repressor. Involved in repression of heat shock proteins. Binds to three inverted repeats in the promoter region of the DnaK operon. Induced by heat shock. | 3 | heat response |
| Rv1657 | argR | Probable arginine repressor (AHRC). Regulates arginine biosynthesis genes. | 4 | cobalamin synthesis; UMP synthesis; C-N bond formation |
| Rv2215 | dlaT | Dihydrolipoamide acyltransferase, component of pyruvate dehydrogenase. Involved in TCA cycle; converts pyruvate to acetyl-CoA and $CO_2$. Also involved in defense against oxidative stress. | 51 | TCA cycle, respiration, downregulation of virulence factors |
| Rv2359 | zur | Probable zinc uptake regulation protein. Acts as a global negative controlling element, with $Zn^{2+}$ binds operator of repressed genes. | 8 | downregulating translation, iron import |
| Rv2374c | hrcA | Probable heat shock protein transcriptional repressor. Involved in repression of class I heat shock proteins. Prevents heat-shock induction of these operons. | 17 | transcription and translation |
| Rv2610c | pimA | Alpha-mannosyltransferase. Involved in the first mannosylation step in phosphatidylinositol mannoside biosynthesis (transfer of mannose residues onto PI). | 64 | amino acid and nucleobase synth., respiration, growth/proliferation |
| Rv2720 | lexA | Repressor. Represses genes involved in nucleotide excision repair and SOS response. Binds 14-bp palindromic sequence. | 10 | DNA binding, repair, cleavage |
| Rv3301c | phoY1 | Probable transcriptional regulatory protein PhoU-homolog 1. Involved in regulation of active transport of inorganic phosphate across the membrane. | 18 | ETC, oxidative phosphorylation |
| Rv3417c | groEL1 | 60 kDa chaperonin 1 (protein CPN60-1). Prevents misfolding, promotes refolding and proper assembly of unfolded polypeptides generated under stress conditions. | 15 | stress response |
| Rv3574 | kstR | Transcriptional regulatory protein (probably TetR-family). Involved in transcriptional mechanism. Predicted to control regulon involved in lipid metabolism. | 22 | cholesterol, lipid, and carbon metabolism |
| Rv0599c | vapB27 | Possible antitoxin. | 13 | growth regulation, toxin sequestration, RNase |
| Rv0608 | vapB28 | Possible antitoxin. | 12 | growth regulation, toxin sequestration, RNase |
| Rv0623 | vapB30 | Possible antitoxin. | 12 | growth regulation, toxin sequestration, RNase |

17

| **Rv1560** | vapB11 | Possible antitoxin. | 6 | growth regulation |
|---|---|---|---|---|
| **Rv1740** | vapB34 | Possible antitoxin. | 7 | growth regulation |
| **Rv1960c** | parD1 | Possible antitoxin. | 18 | growth regulation, toxin sequestration, RNase |
| **Rv2009** | vapB15 | Antitoxin. | 13 | growth regulation, RNase |
| **Rv2595** | vapB40 | Possible antitoxin. | 8 | growth regulation, toxin sequestration |
| **Rv2760c** | vapB42 | Possible antitoxin. | 4 | growth regulation, DNA repair |

655

656

657

**Supplementary Material**

1    **Supplementary Figure 1 UMAP.** Hyperparameter optimization was performed on UMAPs from the (A) TFI microarray compendium, (B) RNA-seq compendium, or (C) TFAs calculated from the RNA-seq compendium. ε was varied from 0.1 to 10 on a logarithmic scale and numbers of clusters (left), numbers of outliers (center), and maximum cluster size (right) were computed for each ε. ε was selected from the elbow of the outliers plot (ε = 0.281 for TFI data, 0.309 for RNA-seq compendium and estimated TFAs).

2    **Supplementary Figure 2 Inferred network validation.** Distribution of the ranks, in each network, of edges shared with the validation dataset from Sanz et al., 2011, [8] from each network. Each histogram is divided into 32 bins. Horizontal dashed lines represent the expected number of random matches between each network and the validation dataset. Truncation was performed on these networks at the first bin where the count dropped below the dashed line (see **Methods**). Panels with hashed backgrounds (B, F, and L) represent networks that were excluded from the aggregation due to insufficient enrichment.

3    **Supplementary Figure 3. Inferred network performance.** Performance of each inferred and aggregate network, calculated against a set of TF–target gene relationships identified by Sanz et al., 2011 [8] (see **Methods**), as measured by Matthews correlation coefficient (MCC). MCC quantifies the level of correlation between the two independent sets of relationships. Higher values indicate greater correlation. The blue bars depict the MCC for the dataset-level and overall aggregates. Other colors are used to depict the MCC for the individually inferred networks. Hatched bars indicate the networks that were excluded from aggregation. The horizontal dashed line represents the 95th percentile MCC performance of 1,000 randomly generated networks. See Methods for exclusion criteria.

4    **Supplementary Figure 4 TRN properties.** Out-degree distribution of TF-gene interactions (edges) from the overall aggregate network. This distribution significantly differs from a power law distribution on the left side of the plot, likely because the network includes indirect interactions. These will deflate counts of low-degree TFs (nodes) and inflate counts of higher-degree nodes.

5    **Supplementary Figure 5 Assignment of activating vs repressing regulatory interactions.** Flow chart depicting the logic used to assign directionality to regulatory relationships. Abbreviations used are defined in the legend in the bottom left.

6    **Supplementary Figure 6 TFA rank percentiles for randomized networks**. TRNs were randomized 10 times. For each random network, ROBNCA was used to compute TFAs for the TFI dataset. Rank percentiles were assigned to each TFA for each TFI microarray profile and averaged across replicates for each TFI strain. Plotted are TFA rank percentile distributions for all over-expressed TFs corresponding to their respective TFI strain from each randomized network.

7    **Supplementary Figure 7. TFA-fitness regression model performance.** (A) Fitness values predicted by the gradient boosted machine (GBM) model versus the experimentally measured values supplied to the model upon training. The line of best fit depicts the relationship between predicted and measured values. The slope of this line is slightly less than 1, indicating that the regression model modestly underestimates relative fitness changes. The model achieved a

706    coefficient of determination ($R^2$) of 0.87 against its training set, indicating that the model can
707    explain 87% of the variation in fitness from the TRIP screen. (B) Residuals of the model
708    predictions versus measured values form a roughly normal distribution, indicating a lack of
709    bias and overall reliable predictive ability.
710

711  8    **Supplementary Figure 8 TFA hypoxia prediction.** Our TFA-fitness regression model
712    predicted a decrease in growth over the course of the period of hypoxia, and an increase in
713    growth again upon reaeration, based only on transcriptional data measured over the course of
714    the experiment (each point represents an RNA-seq timepoint), in both the empty plasmid strain
715    (blue) and wild-type H37Rv (orange).
716

717    **Supplementary Figure 9 Hypoxia-responsive TFAs.** The TFA-fitness regression model can
718    be interrogated to determine drivers of hypoxia by comparing the most impactful TFAs under
719    hypoxia (days 2-7) versus reaeration (days 9-12). 7 TFs were most important for predicting
720    reduced growth under hypoxia versus reaeration. Each contributes at least 5% to total model
721    predictions. Depicted is the mean change in TFA for each of the impactful TFs across days 2-7
722    (orange) versus days 9-12 (cyan). Other TFAs show negligible changes in activity across
723    hypoxia or (see Methods for details on calculations for changes in TFA).
724

725  9    **Supplementary Figure 10 TF expression hypoxia prediction.** Hypoxia and reaeration fitness
726    changes predicted by a GBM model trained using only TF expression data instead of TFAs.
727

728  10    **Supplementary Table 1 Expression data from the TFI microarray dataset.** Batch
729    correction group assignments for each sample in the TFI microarray dataset. Smooth quantile
730    normalized and microarray expression for all genes and all samples in the TFI microarray
731    dataset. Median and MAD expression for each gene. Group assignments were used by the
732    PySNAIL smooth quantile normalization algorithm for batch correction [18].
733

734  11    **Supplementary Table 2 Expression data from the RNA-seq expression compendium.**
735    Batch correction group assignments for each sample in the RNA-seq compendium. Group
736    assignments were used by the PySNAIL smooth quantile normalization algorithm for batch
737    correction [18]. Median and MAD expression for each gene.
738

739  12    **Supplementary Table 3 Network inference methods.** Description of transcriptional
740    regulatory network inference methods.
741

742  13    **Supplementary Table 4. Aggregate network directionality of regulation.** Summary of the
743    assignments of activating (up) vs. repressing (down) regulatory interactions for all TF-gene
744    regulatory interactions in the aggregate transcriptional regulatory network (TRN).
745

746  14    **Supplementary Table 5 TF Gene Ontology assignments.** GO enrichment for each
747    transcriptional program regulated by each TF inferred by our aggregate TRN. (A) Annotated
748    functions and a summary of GO enrichments found for targets from selected TFs. All TFs with
749    at least 3 significant GO enrichment terms and a non-locus gene name in Mycobrowser [40].
750    45 TFs meet these criteria. These data validate the accuracy of our network, as one would
751    expect an accurate regulatory network to have target sets significantly enriched for the known
752    functions of each TF. (B) Remaining TFs with at least 3 significant GO enrichments assigned
753    by our analysis but without an annotated gene name (36 additional TFs). These data represent

754 predictions for potentially novel TF functions. (C) All GO enrichments identified by our
755 analysis were corrected for FDR with a cutoff of 0.05.
756
757 15 **Supplementary Table 6 Transcription Factor Activities.** Median and MAD expression and
758 activity for each TF in the RNA-seq compendium. Pearson correlation coefficient between TF
759 expression and TFA for each TF across all samples in the RNA-seq compendium.
760
761 16 **Supplementary Table 7** Overview of the top 7 most important TFAs for predicting fitness
762 under hypoxia as identified by our TFA regression model, validated by published evidence for
763 mechanistic activation under hypoxia [3; 6; 76; 77; 78; 79; 80; 81; 82; 83; 84; 85; 86].

764

## Data Availability Statement

The transcriptome datasets analyzed for this study can be found in the supplemental material and at https://tfnetwork.streamlit.app. The code and software implementations associated with this study can be found at https://github.com/Ma-Lab-Seattle-Childrens-CGIDR/Mtb-TFA-fitness-regression and https://hub.docker.com/repositories/malabcgidr.

## References

[1] WHO, Global tuberculosis report 2023, 2023.

[2] J.E. Galagan, K. Minch, M. Peterson, A. Lyubetskaya, E. Azizi, L. Sweet, A. Gomes, T. Rustad, G. Dolganov, I. Glotova, T. Abeel, C. Mahwinney, A.D. Kennedy, R. Allard, W. Brabant, A. Krueger, S. Jaini, B. Honda, W.H. Yu, M.J. Hickey, J. Zucker, C. Garay, B. Weiner, P. Sisk, C. Stolte, J.K. Winkler, Y. Van de Peer, P. Iazzetti, D. Camacho, J. Dreyfuss, Y. Liu, A. Dorhoi, H.J. Mollenkopf, P. Drogaris, J. Lamontagne, Y. Zhou, J. Piquenot, S.T. Park, S. Raman, S.H. Kaufmann, R.P. Mohney, D. Chelsky, D.B. Moody, D.R. Sherman, and G.K. Schoolnik, The Mycobacterium tuberculosis regulatory network and hypoxia. Nature 499 (2013) 178-83.

[3] K.J. Minch, T.R. Rustad, E.J. Peterson, J. Winkler, D.J. Reiss, S. Ma, M. Hickey, W. Brabant, B. Morrison, S. Turkarslan, C. Mawhinney, J.E. Galagan, N.D. Price, N.S. Baliga, and D.R. Sherman, The DNA-binding network of Mycobacterium tuberculosis. Nat Commun 6 (2015) 5829.

[4] E.J. Peterson, D.J. Reiss, S. Turkarslan, K.J. Minch, T. Rustad, C.L. Plaisier, W.J. Longabaugh, D.R. Sherman, and N.S. Baliga, A high-resolution network model for global gene regulation in Mycobacterium tuberculosis. Nucleic Acids Res 42 (2014) 11291-303.

[5] E.J.R. Peterson, A.N. Brooks, D.J. Reiss, A. Kaur, J. Do, M. Pan, W.J. Wu, R. Morrison, V. Srinivas, W. Carter, M.L. Arrieta-Ortiz, R.A. Ruiz, A. Bhatt, and N.S. Baliga, MtrA modulates Mycobacterium tuberculosis cell division in host microenvironments to mediate intrinsic resistance and drug tolerance. Cell Rep 42 (2023) 112875.

[6] T.R. Rustad, K.J. Minch, S. Ma, J.K. Winkler, S. Hobbs, M. Hickey, W. Brabant, S. Turkarslan, N.D. Price, N.S. Baliga, and D.R. Sherman, Mapping and manipulating the Mycobacterium tuberculosis transcriptome using a transcription factor overexpression-derived regulatory network. Genome Biol 15 (2014) 502.

[7] R. Yoo, K. Rychel, S. Poudel, T. Al-Bulushi, Y. Yuan, S. Chauhan, C. Lamoureux, B.O. Palsson, and A. Sastry, Machine Learning of All Mycobacterium tuberculosis H37Rv RNA-seq Data Reveals a Structured Interplay between Metabolism, Stress Response, and Infection. mSphere 7 (2022) e0003322.

[8] J. Sanz, J. Navarro, A. Arbues, C. Martin, P.C. Marijuan, and Y. Moreno, The transcriptional regulatory network of Mycobacterium tuberculosis. PLoS One 6 (2011) e22178.

[9] G. Balazsi, A.P. Heath, L. Shi, and M.L. Gennaro, The temporal response of the Mycobacterium tuberculosis gene regulatory network during growth arrest. Mol Syst Biol 4 (2008) 225.

[10] S. Turkarslan, E.J.R. Peterson, T.R. Rustad, K.J. Minch, D.J. Reiss, R. Morrison, S. Ma, N.D. Price, D.R. Sherman, and N.S. Baliga, A comprehensive map of genome-wide gene regulation in Mycobacterium tuberculosis. Scientific Data 2 (2015).

[11] J.M. Escorcia-Rodriguez, E. Gaytan-Nunez, E.M. Hernandez-Benitez, A. Zorro-Aranda, M.A. Tello-Palencia, and J.A. Freyre-Gonzalez, Improving gene regulatory network inference and assessment: The importance of using network structure. Front Genet 14 (2023) 1143382.

[12] H. Poonawala, Y. Zhang, S. Kuchibhotla, A.G. Green, D.M. Cirillo, F. Di Marco, A. Spitlaeri, P. Miotto, and M.R. Farhat, Transcriptomic responses to antibiotic exposure in Mycobacterium tuberculosis. Antimicrob Agents Chemother 68 (2024) e0118523.

[13] C. Bei, J. Zhu, P.H. Culviner, M. Gan, E.J. Rubin, S.M. Fortune, Q. Gao, and Q. Liu, Genetically encoded transcriptional plasticity underlies stress adaptation in Mycobacterium tuberculosis. Nat Commun 15 (2024) 3088.

[14] J.C. Liao, R. Boscolo, Y.L. Yang, L.M. Tran, C. Sabatti, and V.P. Roychowdhury, Network component analysis: reconstruction of regulatory signals in biological systems. Proc Natl Acad Sci U S A 100 (2003) 15522-7.

[15] D. Marbach, J.C. Costello, R. Kuffner, N.M. Vega, R.J. Prill, D.M. Camacho, K.R. Allison, M. Kellis, J.J. Collins, and G. Stolovitzky, Wisdom of crowds for robust gene network inference. Nat Methods 9 (2012) 796-804.

[16] S. Ma, R. Morrison, S.J. Hobbs, V. Soni, J. Farrow-Johnson, A. Frando, N. Fleck, C. Grundner, K.Y. Rhee, T.R. Rustad, and D.R. Sherman, Transcriptional regulator-induced phenotype screen reveals drug potentiators in Mycobacterium tuberculosis. Nat Microbiol 6 (2021) 44-50.

[17] S.C. Hicks, K. Okrah, J.N. Paulson, J. Quackenbush, R.A. Irizarry, and H.C. Bravo, Smooth quantile normalization. Biostatistics 19 (2018) 185-198.

[18] P.H. Hsieh, C.M. Lopes-Ramos, M. Zucknick, G.K. Sandve, K. Glass, and M.L. Kuijjer, Adjustment of spurious correlations in co-expression measurements from RNA-Sequencing data. Bioinformatics 39 (2023).

[19] B. Langmead, and S.L. Salzberg, Fast gapped-read alignment with Bowtie 2. Nat Methods 9 (2012) 357-9.

[20] Y. Liao, G.K. Smyth, and W. Shi, featureCounts: an efficient general purpose program for assigning sequence reads to genomic features. Bioinformatics 30 (2014) 923-30.

[21] P. Ewels, M. Magnusson, S. Lundin, and M. Kaller, MultiQC: summarize analysis results for multiple tools and samples in a single report. Bioinformatics 32 (2016) 3047-8.

[22] L. McInnes, J. Healy, and J. Melville, UMAP: Uniform Manifold Approximation and Projection for Dimension Reduction. arXiv 1802.03426 (2020).

[23] M. Ester, H.-P. Kriegel, J. Sander, and X. Xu, A Density-Based Algorithm for Discovering Clusters in Large Spatial Databases with Noise, Knowledge Discovery and Data Mining, 1996.

[24] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, and E. Duchesnay, Scikit-learn: Machine Learning in Python. J Mach Learn Res 12 (2011) 2825-2830.

847  [25] A.A. Margolin, I. Nemenman, K. Basso, C. Wiggins, G. Stolovitzky, R. Dalla Favera, and A.
848       Califano, ARACNE: an algorithm for the reconstruction of gene regulatory networks in a
849       mammalian cellular context. BMC Bioinformatics 7 Suppl 1 (2006) S7.

850  [26] A. Lachmann, F.M. Giorgi, G. Lopez, and A. Califano, ARACNe-AP: gene network reverse
851       engineering through adaptive partitioning inference of mutual information. Bioinformatics 32
852       (2016) 2233-5.

853  [27] J.J. Faith, B. Hayete, J.T. Thaden, I. Mogno, J. Wierzbowski, G. Cottarel, S. Kasif, J.J. Collins,
854       and T.S. Gardner, Large-scale mapping and validation of Escherichia coli transcriptional
855       regulation from a compendium of expression profiles. PLoS Biol 5 (2007) e8.

856  [28] V.A. Huynh-Thu, A. Irrthum, L. Wehenkel, and P. Geurts, Inferring regulatory networks from
857       expression data using tree-based methods. PLoS One 5 (2010).

858  [29] G. Sales, and C. Romualdi, parmigene--a parallel R package for mutual information estimation
859       and gene network reconstruction. Bioinformatics 27 (2011) 1876-7.

860  [30] S. Aibar, C.B. Gonzalez-Blas, T. Moerman, V.A. Huynh-Thu, H. Imrichova, G. Hulselmans, F.
861       Rambow, J.C. Marine, P. Geurts, J. Aerts, J. van den Oord, Z.K. Atak, J. Wouters, and S.
862       Aerts, SCENIC: single-cell regulatory network inference and clustering. Nat Methods 14
863       (2017) 1083-1086.

864  [31] D.J. Reiss, N.S. Baliga, and R. Bonneau, Integrated biclustering of heterogeneous genome-wide
865       datasets for the inference of global regulatory networks. BMC Bioinformatics 7 (2006) 280.

866  [32] D.J. Reiss, C.L. Plaisier, W.J. Wu, and N.S. Baliga, cMonkey2: Automated, systematic,
867       integrated detection of co-regulated gene modules for any organism. Nucleic Acids Res 43
868       (2015) e87.

869  [33] A.V. Sastry, Y. Gao, R. Szubin, Y. Hefner, S. Xu, D. Kim, K.S. Choudhary, L. Yang, Z.A.
870       King, and B.O. Palsson, The Escherichia coli transcriptome mostly consists of independently
871       regulated modules. Nat Commun 10 (2019) 5536.

872  [34] T. Hastie, R. Tibshirani, and J. Friedman, The Elements of Statistical Learning: Data Mining,
873       Inference and Prediction, Second Edition, Springer, 2008.

874  [35] H. Zou, and T. Hastie, Regularization and Variable Selection Via the Elastic Net. Journal of the
875       Royal Statistical Society Series B: Statistical Methodology 67 (2005) 301-320.

876  [36] D. Chicco, Ten quick tips for machine learning in computational biology. BioData Min 10
877       (2017) 35.

878  [37] B.W. Matthews, Comparison of the predicted and observed secondary structure of T4 phage
879       lysozyme. Biochimica et biophysica acta 405 (1975) 442-51.

880  [38] M. Ashburner, C.A. Ball, J.A. Blake, D. Botstein, H. Butler, J.M. Cherry, A.P. Davis, K.
881       Dolinski, S.S. Dwight, J.T. Eppig, M.A. Harris, D.P. Hill, L. Issel-Tarver, A. Kasarskis, S.
882       Lewis, J.C. Matese, J.E. Richardson, M. Ringwald, G.M. Rubin, and G. Sherlock, Gene
883       ontology: tool for the unification of biology. The Gene Ontology Consortium. Nat Genet 25
884       (2000) 25-9.

885  [39] C. Gene Ontology, S.A. Aleksander, J. Balhoff, S. Carbon, J.M. Cherry, H.J. Drabkin, D. Ebert,
886       M. Feuermann, P. Gaudet, N.L. Harris, D.P. Hill, R. Lee, H. Mi, S. Moxon, C.J. Mungall, A.
887       Muruganugan, T. Mushayahama, P.W. Sternberg, P.D. Thomas, K. Van Auken, J. Ramsey,
888       D.A. Siegele, R.L. Chisholm, P. Fey, M.C. Aspromonte, M.V. Nugnes, F. Quaglia, S.
889       Tosatto, M. Giglio, S. Nadendla, G. Antonazzo, H. Attrill, G. Dos Santos, S. Marygold, V.

Strelets, C.J. Tabone, J. Thurmond, P. Zhou, S.H. Ahmed, P. Asanitthong, D. Luna Buitrago, M.N. Erdol, M.C. Gage, M. Ali Kadhum, K.Y.C. Li, M. Long, A. Michalak, A. Pesala, A. Pritazahra, S.C.C. Saverimuttu, R. Su, K.E. Thurlow, R.C. Lovering, C. Logie, S. Oliferenko, J. Blake, K. Christie, L. Corbani, M.E. Dolan, H.J. Drabkin, D.P. Hill, L. Ni, D. Sitnikov, C. Smith, A. Cuzick, J. Seager, L. Cooper, J. Elser, P. Jaiswal, P. Gupta, P. Jaiswal, S. Naithani, M. Lera-Ramirez, K. Rutherford, V. Wood, J.L. De Pons, M.R. Dwinell, G.T. Hayman, M.L. Kaldunski, A.E. Kwitek, S.J.F. Laulederkind, M.A. Tutaj, M. Vedi, S.J. Wang, P. D'Eustachio, L. Aimo, K. Axelsen, A. Bridge, N. Hyka-Nouspikel, A. Morgat, S.A. Aleksander, J.M. Cherry, S.R. Engel, K. Karra, S.R. Miyasato, R.S. Nash, M.S. Skrzypek, S. Weng, E.D. Wong, E. Bakker, et al., The Gene Ontology knowledgebase in 2023. Genetics 224 (2023).

[40] A. Kapopoulou, J.M. Lew, and S.T. Cole, The MycoBrowser portal: a comprehensive and manually annotated resource for mycobacterial genomes. Tuberculosis (Edinb) 91 (2011) 8-13.

[41] D.V. Klopfenstein, L. Zhang, B.S. Pedersen, F. Ramirez, A. Warwick Vesztrocy, A. Naldi, C.J. Mungall, J.M. Yunes, O. Botvinnik, M. Weigel, W. Dampier, C. Dessimoz, P. Flick, and H. Tang, GOATOOLS: A Python library for Gene Ontology analyses. Sci Rep 8 (2018) 10872.

[42] S. Carbon, and C. Mungall, Gene Ontology Data Archive (2024-06-17) [Data set]. Zenodo (2024).

[43] A. Noor, A. Ahmad, E. Serpedin, M. Nounou, and H. Nounou, ROBNCA: robust network component analysis for recovering transcription factor activities. Bioinformatics 29 (2013) 2410-8.

[44] T.A. Ahn-Horst, L.S. Mille, G. Sun, J.H. Morrison, and M.W. Covert, An expanded whole-cell model of E. coli links cellular physiology with mechanisms of growth rate control. NPJ Syst Biol Appl 8 (2022) 30.

[45] Y. Shi, G. Ke, Z. Chen, S. Zheng, and T.-Y. Liu, Quantized Training of Gradient Boosting Decision Trees. in: S. Koyejo, S. Mohamed, A. Agarwal, D. Belgrave, K. Cho, and A. Oh, (Eds.), 2022, pp. 18822--18833.

[46] G. Ke, Q. Meng, T. Finley, T. Wang, W. Chen, W. Ma, Q. Ye, and T.-Y. Liu, LightGBM: A Highly Efficient Gradient Boosting Decision Tree. in: I. Guyon, U.V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, (Eds.), 2017.

[47] A.M. Sherrid, T.R. Rustad, G.A. Cangelosi, and D.R. Sherman, Characterization of a Clp protease gene regulator and the reaeration response in Mycobacterium tuberculosis. PLoS One 5 (2010) e11622.

[48] D.R. Sherman, M. Voskuil, D. Schnappinger, R. Liao, M.I. Harrell, and G.K. Schoolnik, Regulation of the Mycobacterium tuberculosis hypoxic response gene encoding alpha -crystallin. Proc Natl Acad Sci U S A 98 (2001) 7534-9.

[49] Y. Yuan, D.D. Crane, R.M. Simpson, Y.Q. Zhu, M.J. Hickey, D.R. Sherman, and C.E. Barry, 3rd, The 16-kDa alpha-crystallin (Acr) protein of Mycobacterium tuberculosis is required for growth in macrophages. Proc Natl Acad Sci U S A 95 (1998) 9578-83.

[50] S. Ma, S. Jaipalli, J. Larkins-Ford, J. Lohmiller, B.B. Aldridge, D.R. Sherman, and S. Chandrasekaran, Transcriptomic Signatures Predict Regulators of Drug Synergy and Clinical Regimen Efficacy against Tuberculosis. mBio 10 (2019).

933 [51] H. Li, B. Handsaker, A. Wysoker, T. Fennell, J. Ruan, N. Homer, G. Marth, G. Abecasis, R.
934      Durbin, and S. Genome Project Data Processing, The Sequence Alignment/Map format and
935      SAMtools. Bioinformatics 25 (2009) 2078-9.

936 [52] Y. Benjamini, A.M. Krieger, and D. Yekutieli, Adaptive linear step-up procedures that control
937      the false discovery rate. Biometrika 93 (2006) 491-507.

938 [53] Z. Wang, M. Gerstein, and M. Snyder, RNA-Seq: a revolutionary tool for transcriptomics. Nat
939      Rev Genet 10 (2009) 57-63.

940 [54] S.L. Kendall, P. Burgess, R. Balhana, M. Withers, A. Ten Bokum, J.S. Lott, C. Gao, I. Uhia-
941      Castro, and N.G. Stoker, Cholesterol utilization in mycobacteria is controlled by two TetR-
942      type transcriptional regulators: kstR and kstR2. Microbiology (Reading) 156 (2010) 1362-
943      1371.

944 [55] R. Sprinthall, Basic Statistical Analysis, Pearson Education, 2011.

945 [56] U. Alon, Network motifs: theory and experimental approaches. Nat Rev Genet 8 (2007) 450-61.

946 [57] S.J. Modlin, A. Elghraoui, D. Gunasekaran, A.M. Zlotnicki, N.A. Dillon, N. Dhillon, N. Kuo, C.
947      Robinhold, C.K. Chan, A.D. Baughn, and F. Valafar, Structure-Aware Mycobacterium
948      tuberculosis Functional Annotation Uncloaks Resistance, Metabolic, and Virulence Genes.
949      mSystems 6 (2021) e0067321.

950 [58] G.B.D.A.R. Collaborators, Global burden of bacterial antimicrobial resistance 1990-2021: a
951      systematic analysis with forecasts to 2050. Lancet (2024).

952 [59] M. Farhat, H. Cox, M. Ghanem, C.M. Denkinger, C. Rodrigues, M.S. Abd El Aziz, H. Enkh-
953      Amgalan, D. Vambe, C. Ugarte-Gil, J. Furin, and M. Pai, Drug-resistant tuberculosis: a
954      persistent global health concern. Nat Rev Microbiol 22 (2024) 617-635.

955 [60] M.N. Anahtar, J.H. Yang, and S. Kanjilal, Applications of Machine Learning to the Problem of
956      Antimicrobial Resistance: an Emerging Model for Translational Research. Journal of clinical
957      microbiology 59 (2021) e0126020.

958 [61] S. Lobentanzer, P. Rodriguez-Mier, S. Bauer, and J. Saez-Rodriguez, Molecular causality in the
959      advent of foundation models. Mol Syst Biol 20 (2024) 848-858.

960 [62] V. Chen, M. Yang, W. Cui, J.S. Kim, A. Talwalkar, and J. Ma, Applying interpretable machine
961      learning in computational biology-pitfalls, recommendations and opportunities for new
962      developments. Nat Methods 21 (2024) 1454-1461.

963 [63] J.H. Yang, S.N. Wright, M. Hamblin, D. McCloskey, M.A. Alcantar, L. Schrübbers, A.J.
964      Lopatkin, S. Satish, A. Nili, B.O. Palsson, G.C. Walker, and J.J. Collins, A White-Box
965      Machine Learning Approach for Revealing Antibiotic Mechanisms of Action. Cell 177
966      (2019) 1649-1661.e9.

967 [64] A.J. Lopatkin, and J.H. Yang, Digital Insights Into Nucleotide Metabolism and Antibiotic
968      Treatment Failure. Front Digit Health 3 (2021).

969 [65] J. Larkins-Ford, Y.N. Degefu, N. Van, A. Sokolov, and B.B. Aldridge, Design principles to
970      assemble drug combinations for effective tuberculosis therapy using interpretable pairwise
971      drug response measurements. Cell Rep Med 3 (2022) 100737.

972 [66] J. Larkins-Ford, T. Greenstein, N. Van, Y.N. Degefu, M.C. Olson, A. Sokolov, and B.B.
973      Aldridge, Systematic measurement of combination-drug landscapes to predict in vivo
974      treatment outcomes for tuberculosis. Cell Syst 12 (2021) 1046-1063 e7.

[67] M.A. DeJesus, E.R. Gerrick, W. Xu, S.W. Park, J.E. Long, C.C. Boutte, E.J. Rubin, D. Schnappinger, S. Ehrt, S.M. Fortune, C.M. Sassetti, and T.R. Ioerger, Comprehensive Essentiality Analysis of the Mycobacterium tuberculosis Genome via Saturating Transposon Mutagenesis. MBio 8 (2017).

[68] B. Bosch, M.A. DeJesus, N.C. Poulton, W. Zhang, C.A. Engelhart, A. Zaveri, S. Lavalette, N. Ruecker, C. Trujillo, J.B. Wallach, S. Li, S. Ehrt, B.T. Chait, D. Schnappinger, and J.M. Rock, Genome-wide gene expression tuning reveals diverse vulnerabilities of M. tuberculosis. Cell 184 (2021) 4579-4592 e24.

[69] S. Li, N.C. Poulton, J.S. Chang, Z.A. Azadian, M.A. DeJesus, N. Ruecker, M.D. Zimmerman, K.A. Eckartt, B. Bosch, C.A. Engelhart, D.F. Sullivan, M. Gengenbacher, V.A. Dartois, D. Schnappinger, and J.M. Rock, CRISPRi chemical genetics and comparative genomics identify genes mediating drug potency in Mycobacterium tuberculosis. Nat Microbiol 7 (2022) 766-779.

[70] W. Xu, M.A. DeJesus, N. Rucker, C.A. Engelhart, M.G. Wright, C. Healy, K. Lin, R. Wang, S.W. Park, T.R. Ioerger, D. Schnappinger, and S. Ehrt, Chemical Genetic Interaction Profiling Reveals Determinants of Intrinsic Antibiotic Resistance in Mycobacterium tuberculosis. Antimicrob Agents Chemother 61 (2017).

[71] P.O. Oluoch, E.-I. Koh, M.K. Proulx, C.J. Reames, K.G. Papavinasasundaram, K.C. Murphy, M.D. Zimmerman, V. Dartois, and C.M. Sassetti, Chemical genetic interactions elucidate pathways controlling tuberculosis antibiotic efficacy during infection. bioRxiv (2024) 2024.09.04.609063.

[72] A.F. Carey, J.M. Rock, I.V. Krieger, M.R. Chase, M. Fernandez-Suarez, S. Gagneux, J.C. Sacchettini, T.R. Ioerger, and S.M. Fortune, TnSeq of Mycobacterium tuberculosis clinical isolates reveals strain-specific antibiotic liabilities. PLoS Pathog 14 (2018) e1006939.

[73] N.D. Hicks, J. Yang, X. Zhang, B. Zhao, Y.H. Grad, L. Liu, X. Ou, Z. Chang, H. Xia, Y. Zhou, S. Wang, J. Dong, L. Sun, Y. Zhu, Y. Zhao, Q. Jin, and S.M. Fortune, Clinically prevalent mutations in Mycobacterium tuberculosis alter propionate metabolism and mediate multidrug tolerance. Nat Microbiol 3 (2018) 1032-1042.

[74] S. Stanley, C.N. Spaulding, Q. Liu, M.R. Chase, D.T.M. Ha, P.V.K. Thai, N.H. Lan, D.D.A. Thu, N.L. Quang, J. Brown, N.D. Hicks, X. Wang, M. Marin, N.C. Howard, A.J. Vickers, W.M. Karpinski, M.C. Chao, M.R. Farhat, M. Caws, S.J. Dunstan, N.T.T. Thuong, and S.M. Fortune, Identification of bacterial determinants of tuberculosis infection and treatment outcomes: a phenogenomic analysis of clinical strains. Lancet Microbe 5 (2024) e570-e580.

[75] C.C. The, A data compendium associating the genomes of 12,289 Mycobacterium tuberculosis isolates with quantitative resistance phenotypes to 13 antibiotics. PLoS Biol 20 (2022) e3001721.

[76] J.E. Cronan, The Escherichia coli FadR transcription factor: Too much of a good thing? Mol Microbiol 115 (2021) 1080-1085.

[77] M.A. Forrellad, M.V. Bianco, F.C. Blanco, J. Nunez, L.I. Klepp, C.L. Vazquez, L. Santangelo Mde, R.V. Rocha, M. Soria, P. Golby, M.G. Gutierrez, and F. Bigi, Study of the in vivo role of Mce2R, the transcriptional regulator of mce2 operon in Mycobacterium tuberculosis. BMC microbiology 13 (2013) 200.

[78] V. Gopinath, S. Raghunandanan, R.L. Gomez, L. Jose, A. Surendran, R. Ramachandran, A.R. Pushparajan, S. Mundayoor, A. Jaleel, and R.A. Kumar, Profiling the Proteome of

27

1019         Mycobacterium tuberculosis during Dormancy and Reactivation. Mol Cell Proteomics 14
1020         (2015) 2160-76.

1021   [79] C. Larsson, B. Luna, N.C. Ammerman, M. Maiga, N. Agarwal, and W.R. Bishai, Gene
1022         expression of Mycobacterium tuberculosis putative transcription factors whiB1-7 in redox
1023         environments. PLoS One 7 (2012) e37516.

1024   [80] J. Li, X. Wang, W. Gong, C. Niu, and M. Zhang, Crystallization and preliminary X-ray analysis
1025         of Rv1674c from Mycobacterium tuberculosis. Acta Crystallogr F Struct Biol Commun 71
1026         (2015) 354-7.

1027   [81] R. Manganelli, L. Cioetto-Mazzabo, G. Segafreddo, F. Boldrin, D. Sorze, M. Conflitti, A.
1028         Serafini, and R. Provvedi, SigE: A master regulator of Mycobacterium tuberculosis. Front
1029         Microbiol 14 (2023) 1075143.

1030   [82] S. Mehra, and D. Kaushal, Functional genomics reveals extended roles of the Mycobacterium
1031         tuberculosis stress response factor sigmaH. J Bacteriol 191 (2009) 3965-80.

1032   [83] B. Ramos, S.V. Gordon, and M.V. Cunha, Revisiting the expression signature of pks15/1
1033         unveils regulatory patterns controlling phenolphtiocerol and phenolglycolipid production in
1034         pathogenic mycobacteria. PLoS One 15 (2020) e0229700.

1035   [84] S. Yousuf, R.K. Angara, A. Roy, S.K. Gupta, R. Misra, and A. Ranjan, Mce2R/Rv0586 of
1036         Mycobacterium tuberculosis is the functional homologue of FadR(E. coli). Microbiology
1037         (Reading) 164 (2018) 1133-1145.

1038   [85] T.R. Rustad, M.I. Harrell, R. Liao, and D.R. Sherman, The enduring hypoxic response of
1039         Mycobacterium tuberculosis. PLoS One 3 (2008) e1502.

1040   [86] K.R. Nicholson, R.M. Cronin, A.R. Menon, M.K. Jennisch, D.M. Tobin, and P.A. Champion,
1041         The EspN transcription factor is an infection-dependent regulator of the ESX-1 system in M.
1042         marinum. bioRxiv (2023).

1043

1044

**A** TFI Microarray Data

**B** TFI Microarray Data

**C** RNA-seq Compendium

**D** RNA-seq Compendium

**A**

PC 1 (41% of variance) vs PC 2 (15% of variance)

**B**

MCC (Test)

Legend:
- Aggregate
- cMonkey2
- ARACNe
- Elasticnet
- CLR
- GENIE3
- iModulon
- RNASeq (▲)
- TFI (✕)
- Both (●)

**A** TFI Network

**B** Average from Random Networks

**C** Transcription Factor Activities

112 Clusters

**D** Transcription Factor Activities

*Rv1985c*

*sigA*

*tcrA*

*mprA*

*devR*

*kmtR*

*whiB1*

*hupB*

**E** Transcription Factor RNA-seq Expression

*whiB6*

*devR*

*kmtR*

*csoR*

*sigB*

*whiB1*

*hupB*

**F**

A — Hypoxia timecourse

B — Most Impactful TFAs