

# 1 Prediction of Whole-Cell Transcriptional Response with 2 Machine Learning

**Mohammed Eslami**<sup>1,\*,+</sup>

**Amin Espah Borujeni**<sup>2,+</sup>

Hamid Doosthosseini<sup>2</sup>

**Matthew Vaughn**<sup>3</sup>

**Hamed Eramian**<sup>1</sup>

**Katie Clowers**<sup>4</sup>

D. Benjamin Gordon<sup>2</sup>

**Niall Gaffney**<sup>3</sup>

**Mark Weston**<sup>1</sup>

**Diveena Becker**<sup>4</sup>

Yuval Dorfan<sup>2</sup>

**John Fonner**<sup>3</sup>

**Joshua Urrutia**<sup>3</sup>

**Carolyn Corbet**<sup>4</sup>

**George Zheng**<sup>1</sup>

**Joe Stubbs**<sup>3</sup>

Alexander Cristofaro<sup>2,5</sup>

**Paul Maschhoff**<sup>4</sup>

**Jedediah Singer**<sup>6</sup>

Christopher A Voigt<sup>2</sup>

Enoch Yeung<sup>7,\*</sup>

3

<sup>1</sup>Netrias, LLC, Annapolis, MD 21409, USA, <sup>2</sup>Massachusetts Institute of Technology, Cambridge, MA 02139, USA, <sup>3</sup>Texas Advanced Computing Center, Austin, TX 78758, USA, <sup>4</sup>Ginkgo Bioworks, Inc., Boston, MA 02210, USA, <sup>5</sup>TScan Therapeutics, Inc., Waltham, MA 02451, USA, <sup>6</sup>Two Six Technologies, Arlington, VA 22203, USA, <sup>7</sup>University of California Santa Barbara, Santa Barbara, CA 93106, USA

\*To whom correspondence should be addressed. ([meslami@netrias.com](mailto:meslami@netrias.com), [eyeung@ucsb.edu](mailto:eyeung@ucsb.edu))

+Authors contributed equally.

4

## 5 Abstract

6 Applications in synthetic and systems biology can benefit from measuring whole-cell response to  
7 biochemical perturbations. Execution of experiments to cover all possible combinations of  
8 perturbations is infeasible. In this paper, we present the host response model (HRM), a machine  
9 learning approach that takes the cell response to single perturbations as the input and predicts the  
10 whole cell transcriptional response to the combination of inducers. We find that the HRM is able  
11 to qualitatively predict the directionality of dysregulation to a combination of inducers with an  
12 accuracy of >90% using data from single inducers. We further find that the use of known prior,  
13 known cell regulatory networks doubles the predictive performance of the HRM (an  $R^2$  from 0.3  
14 to 0.65). This tool will significantly reduce the number of high-throughput sequencing  
15 experiments that need to be run to characterize the transcriptional impact of the combination of  
16 perturbations on the host.

## 17 Introduction

18 Cells enact complex dynamics in response to environmental and biochemical perturbations. The  
19 perturbation can have a widespread effect so as to alter the dynamics of the whole cell through  
20 cascading effects that span through a cell's regulatory network. Combinations of the biochemical  
21 perturbations are thus not additive and can trigger complex responses such as heat shock<sup>1,2</sup>,  
22 osmotic shock<sup>3,4</sup>, or sudden shifts in nutrient availability<sup>5,6</sup>. Given the complexity of these  
23 responses to these perturbations, prior studies in perturbed whole cell response examine the role  
24 of a specific well-known biophysical perturbation, for which there is a natural intuition or

25 sensible alignment with known biophysical mechanisms<sup>7,8,9</sup>. These experiments are carefully  
26 performed, driven by biophysical knowledge and hypothesis-based modeling.

27

28 A natural extension of this research is to examine how the whole cell responds to inputs that are  
29 foreign to the natural workings of the cell. Further, suppose that a cell was presented with  
30 multiple inputs, each of which lent biophysical insight, but the goal was to predict how the cell  
31 responded combinatorially to these inputs. The scale of experiments required to measure  
32 response in such a combinatorially large condition space is infeasible in terms of cost, labor, and  
33 time. In such a setting, there is great value in developing discovery-based approaches for  
34 spotlighting biophysical mechanisms using data-driven algorithms<sup>10</sup>, such as nonlinear  
35 modeling<sup>11</sup> or machine learning<sup>12</sup>.

36

37 A domain that has grappled with modeling of combinatorially large condition spaces for  
38 prediction of specific cell response is that of drug combination/synergy prediction<sup>12,13</sup>. Machine  
39 and deep learning techniques are widely used to model pharmacodynamic and pharmacokinetic  
40 parameters of a drug and identify biomarkers of drug response given a large corpus of drug and  
41 response features<sup>14-17</sup>. The techniques used in these efforts require large training datasets that  
42 consist of specific cell responses to tens of thousands of drugs, a condition space that is often too  
43 large for high-throughput omics measurements, such as RNASeq, which provide insight into the  
44 whole-cell response. The ubiquity of high-throughput sequencing offers an opportunity to  
45 revolutionize the modeling of whole cell transcriptional response.

46

47 The ubiquity of high-throughput sequencing offers an opportunity to revolutionize the modeling  
48 of whole cell transcriptional response. A measure most often used to qualitatively and  
49 quantitatively assess a transcript's response is its dysregulation as compared to a control.  
50 Differential expression analysis (DEA) is a standard bioinformatics technique that measures  
51 response to perturbations as compared to a control condition<sup>18</sup>. DEA conducts custom  
52 normalization, dispersion modeling, and Bayesian optimization to account for biological and  
53 experimental variability in the data. It quantifies the transcriptional response to a perturbation in  
54 terms of *fold-change* and measures its statistical significance. Data-driven prediction using  
55 machine learning of transcription to date, however, has been limited to expression level  
56 predictions from sequences or images<sup>19-20</sup>. These techniques are prone to generalization errors  
57 that can arise from artifacts of normalization of counts data across experiments with  
58 combinatorically large condition spaces<sup>21</sup>.

59  
60 In this paper we present the host response model (HRM), a machine learning model that can  
61 predict whole-cell transcriptional response to a combination of biochemical perturbations using  
62 transcriptional response data from single perturbations. Biochemical perturbations, in this  
63 context, amounts to inducing a cell with a chemical. The HRM combines high-throughput  
64 sequencing with machine learning to infer links between experimental context, prior knowledge  
65 of cell regulatory networks, and the RNASeq data to predict differential expression of a gene.  
66 The HRM was tested in two organisms, *Escherichia coli* MG1655 (*E. coli*) and *Bacillus subtilis*  
67 Marburg 168 (*B. subtilis*). *E. coli* is a well-studied and characterized Gram-negative bacteria that  
68 served as a proof of concept in the development of the HRM. *B. subtilis* is a well characterized

69 and frequently used model organism for Gram-positive bacteria that was used to pressure test the  
 70 HRM. For conciseness, the figures for *E. coli* are provided in Supplementary information.

## 71 Results

### 72 Training and Validation of a Machine Learning Model

73 In this study, we train and test a model per organism with embeddings of prior known  
 74 transcriptional networks of the host cells to train three machine learning models for  
 75 combinatorial prediction of DEA from single inducers (**Figure 1A**). The best performing model  
 76 was selected using a validation dataset from two double-inducer conditions at two time points.  
 77 Experiments are then conducted with all remaining inducer combinations at two time points to  
 78 test the best performing model (in total 18 experimental conditions, 136 samples) (**Figure 1B**).

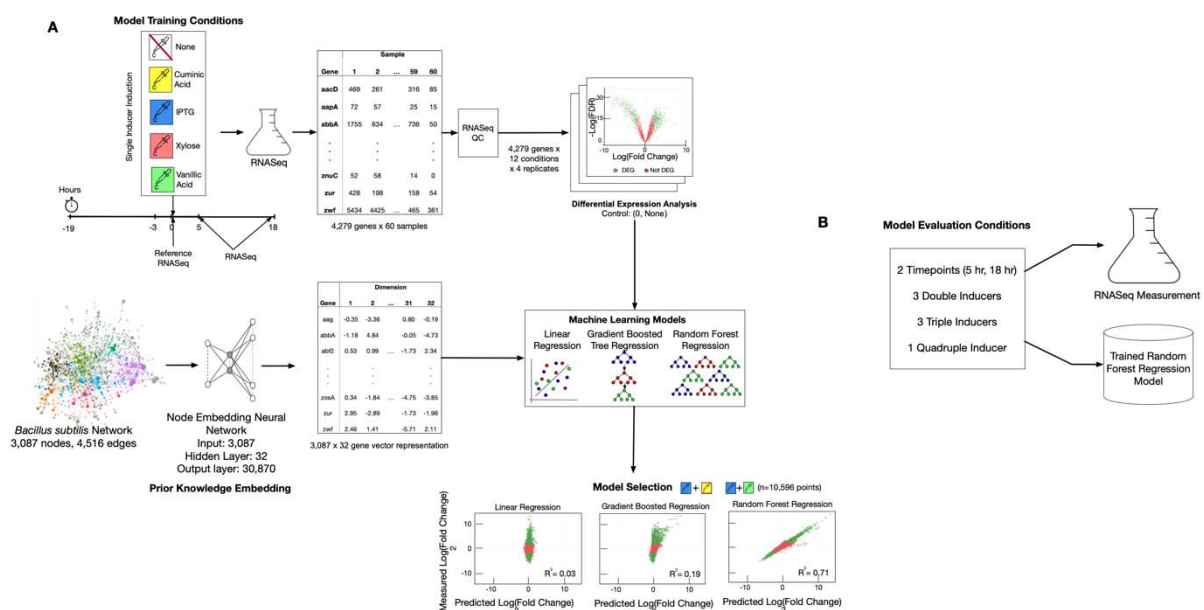


Figure 1: (A) Host response model experiment and knowledge integration for machine learning

training and validation. For *B. subtilis*, RNASeq data was generated for single inducer conditions at each time point and passed through a configuration based differential expression analysis pipeline. Three machine learning models were validated with two held out pairs of inducer combinations. (B) The best model for the HRM is selected and tested with experiments run from all remaining combinations of inducers.

79 The HRM is formulated as a transcriptional dysregulation model trained with differential  
80 expression data and prior knowledge of gene networks of the host. The full set of conditions and  
81 samples collected for the training, validation, and test sets can be found in Supplementary Table  
82 1 and a detailed description of the model can be found in the Methods section.

83

84 Training data for *E. coli* included used EColiNet<sup>22</sup> as the prior gene network and experimental  
85 data that consisted of two inducers, Isopropyl  $\beta$ -d-1-thiogalactopyranoside (IPTG) and arabinose  
86 at four time points (5, 6.5, 8, 18 hours). The test set was made up of the combination of the two  
87 inducers at all four time points. The non-induced, earliest time point (5 hours) was used as a  
88 control condition to measure host response using DEA for both training and testing data  
89 (Supplementary Figure 1). Training data for *B. subtilis* consisted of a transcriptional regulatory  
90 network<sup>23</sup> with experiments at two phases of growth (log and stationary), and four inducers:  
91 IPTG, cuminic acid (CA), vanillic acid (VA), and xylose. The data was passed through a QC and  
92 DEA pipeline which removed low-quality samples and genes (Supplementary Table 2).

93

94 A challenge faced by the HRM is that the number of differential expression comparisons can  
95 have many factors, and thus, many design formulas. A Python-based configurable toolkit, which  
96 we call *omics\_tools*, that parallelizes the execution of DEA for the large condition space was

97 developed to address this challenge. The tool aggregates the outputs from the parallelized runs  
98 and combines all the data into a single unified dataset where each row represents a gene, it's  
99 differential expression, statistical significance, and the condition for downstream machine  
100 learning. *Omicstools* uses edgeR<sup>24</sup> to conduct DEA across the set of design variables. A control  
101 condition of non-induced at the earliest time point was used to quantify the impact of induction  
102 and time. The same control condition was measured for all runs of the experiment that made up  
103 the training, validation and test set of data. The training corpus was formatted as follows: rows  
104 represented genes in each experimental condition, while columns consisted of features of the  
105 condition space, the node embedding features for the gene, and finally, the log fold change and  
106 associated statistical significance of the gene in the condition as compared to the control. The  
107 data can be found with the tutorial.

108  
109 We find, surprisingly, a subspace representation of the individual responses enables prediction of  
110 response to combination of inputs. Qualitative performance of the model was measured as the  
111 number of dysregulated genes whose direction (up/down) was predicted correctly. Quantitative  
112 performance was measured with an  $R^2$  metric comparing predicted versus actual fold-changes on  
113 a logarithmic scale.

## 114 HRM Predicts Transcriptional Response for E. coli

115 The first question to address was whether the set of differentially expressed genes had a large  
116 overlap between the train and test set. If so, then the task of machine learning would likely be  
117 trivial. We measured the Jaccard similarity between the set of genes of pairs of conditions to  
118 estimate the overlap (Supplementary Figure 2). The overlap between the conditions has a median

119 of 0.2 with a standard deviation of 0.25, indicating a significant difference between the train and  
120 test differentially expressed genes (DEGs).

121  
122 Three machine learning models were trained in two ways: using only genes present in the prior  
123 gene network versus using the whole transcriptome. To measure the impact of prior knowledge  
124 on the model, the best performing machine learning model with prior knowledge was selected  
125 and trained without prior knowledge to be used as a control. The criterion to label a gene as a  
126 DEG for *E. Coli* are genes with absolute  $\log_2(\text{Fold Change}) > 1.1$  and an FDR of  $< 0.01$ .  
127 Qualitatively and quantitatively it was clear that machine learning could accomplish the task, but  
128 the impact of prior knowledge was marginal for *E. coli* (Table 1).

129  
130 Table 1: *E. coli* results of qualitative and quantitative predictions for three different models using  
131 two training methods as compared to a control method for model selection.

<b>Model Name</b>	<b>Prior Networks Used?</b>	<b>Training Method</b>	<b>Qualitative</b>	<b>Quantitative</b>
Gradient Boosted Regression	Yes	Genes only in network	58.39%	0.104
Gradient Boosted Regression	Yes	Whole Transcriptome	58.07%	0.103
Linear Regression	Yes	Genes only in network	51.85%	0.007
Linear Regression	Yes	Whole Transcriptome	51.73%	0.006
Random Forest Regression	Yes	Genes only in network	90.20%	0.887
Random Forest Regression	Yes	Whole Transcriptome	87.66%	0.846



Random Forest Regression	No	Control	89.59%	0.829
--------------------------	----	---------	--------	-------

132

### 133 Quantifying Prior Knowledge Impact for *B. subtilis*

134 Similar to the validation framework of *E. coli*, all single inducers and a subset of the double  
 135 inducer data was used to train/validate the model (Table 2). Prior knowledge has a profound  
 136 impact on the *B. subtilis* predictions showing that models trained and tested with genes only  
 137 present in the prior network achieves >90% accuracy. Most interestingly, a model that does not  
 138 use any prior knowledge of the host network achieved an  $R^2=0.306$ , while one that used prior  
 139 knowledge achieved  $R^2=0.708$ , a 2.5x increase in performance.

140

141 Table 2: *B. subtilis* results of qualitative and quantitative predictions for three different models  
 142 using two training methods as compared to a control method for model selection.

<b>Model Name</b>	<b>Prior Networks Used?</b>	<b>Training Method</b>	<b><i>Qualitative</i></b>	<b><i>Quantitative</i></b>
Gradient Boosted Regression	Yes	Genes only in network	51.20%	0.227
Gradient Boosted Regression	Yes	Whole Transcriptome	50.43%	0.194
Linear Regression	Yes	Genes only in network	47.24%	0.031

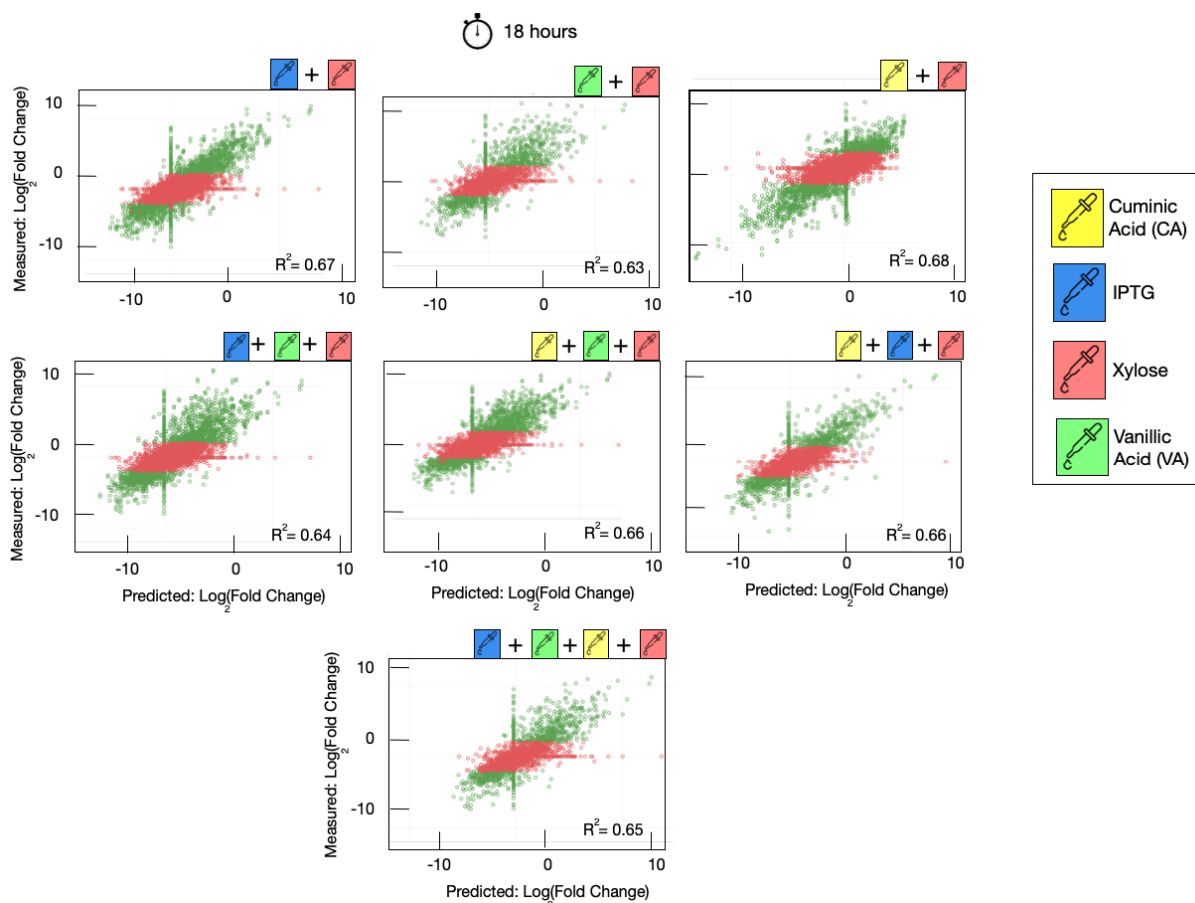
Linear Regression	Yes	Whole Transcriptome	47.27%	0.029
Random Forest Regression	Yes	Genes only in network	90.42%	0.917
Random Forest Regression	Yes	Whole Transcriptome	78.52%	0.708
Random Forest Regression	No	Control	53.04%	0.306

143

144 The discrepancy with *E. coli* can be explained by the heterogeneity of transcriptional response to  
145 the larger set of induction conditions. We computed a rank, or Spearman, correlation between the  
146 train and test induction conditions across all time points for both organisms. Specifically, this  
147 would be three comparisons for *E. coli* (double induced to none and single induced) and 39  
148 comparisons for *B. subtilis* (single induced to all combinations). A distribution of the statistic is  
149 shown in Supplementary Figure 3. The larger distribution observed in the conditions of *B.*  
150 *subtilis* versus *E. coli* explains the impact of prior knowledge.

## 151 Testing the HRM with All Inducer Combinations in *B. subtilis*

152 The best performing machine learning model, the random forest regressor, for the whole  
153 transcriptome was selected to test all remaining combinations. Specifically, predictions at  
154 remaining double, triple, and quadruple inducer conditions at the two time points were both  
155 qualitatively and quantitatively evaluated (Figure 2). Certain conditions could not be evaluated  
156 because there were less than two replicates that passed quality control (QC) (Supplementary  
157 Table 2 and Discussion for more details).

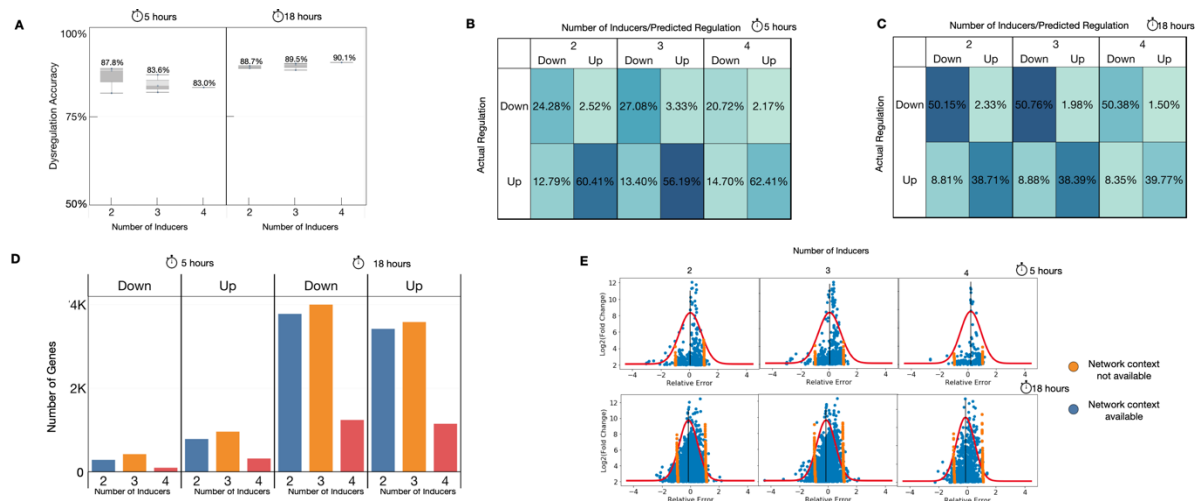


**Figure 2:** Predictions of transcriptional response --  $\log_2(\text{Fold change})$  -- for *B. subtilis* at 18 hours for over 2,000 genes. Conditions tested do not overlap with training and validation sets.

Missing conditions are due to sample quality control. Red points are genes that are not differentially expressed while green genes are ones that are differentially expressed.

158 The model always maintained its performance within statistical error as the number of inducer  
159 combinations increased (**Figure 3A**). Stationary phase predictions showed less variability and  
160 achieved >90% accuracy. The cells in log phase were of lower quality exhibiting lower OD and  
161 RNA integrity (Supplementary Figure 4) across all inducer conditions.

162



163

164 Figure 3: (A) Qualitative predictive performance of predictions on *B. subtilis* test set as it varies

165 by the number of inducers. The performance stays within statistical error as the number of

166 inducers increase. (B-C) Confusion matrices at 5 and 18 hours post-induction indicate that the

167 model predicts more down-regulated genes as up-regulated. (D) The number of up and down

168 regulated genes across inducers and time-points show 2x more up-regulated genes than down-

169 regulated genes at 5 hours which indicates a large class imbalance. The distribution is more

170 evenly balanced at 18 hours. (E) Quantitative error analysis shows most of the errors occur at

171  $|\log_2(\text{Fold change})| < 2.1$  or when the gene is not present in the network. As a matter of fact, the

172 gene's not present in the network place an upper bound on the performance of the model.

173

174 Confusion matrices to assess the qualitative predictions show that the model had a more difficult

175 time predicting up-regulated genes at both time points (**Figure 3B, 3C**). Namely, there are more

176 up-regulated genes predicted as down-regulated than there are down-regulated predicted as up-

177 regulated. We should note that the log phase of growth had double the number of up-regulated

178 genes as down-regulated ones (**Figure 3D**), commonly known as a class imbalance. This

179 provides an opportunity to get more up-regulated genes incorrect at that phase. When the class is

180 balanced, as in the stationary phase, the model's predictions improve but hit an upper bound  
181 which will be explained in the quantitative assessment.  
182 Quantitatively, the relative error is normally distributed with larger error at smaller  
183 dysregulations as those changes are harder to detect and predict. Genes with no network context  
184 available were mapped to a single point in the embedding space and so the model always  
185 predicted a constant value for all those genes (**Figure 3E**). The network is composed of 2,608  
186 genes from the total 4,266 genes in our reference strain. While this is >50% of the genome, genes  
187 with no network information make up only 25% of the DEGs. Even so, these genes make up the  
188 majority of qualitatively incorrect predictions and have greater than a single fold error (Table 3).  
189 This results in an upper bound on the performance of a model as the majority of the model's  
190 errors are due to a lack of network context for certain genes.

191

192 Table 3: Number of incorrect predictions with greater than a single fold change error are  
193 primarily made up of genes with no embedding information.

Number of Inducers in Test Conditions	Total number of predictions	Number of Incorrect Predictions due to absence of gene in network	Number of Incorrect Predictions due to presence of gene in network	Percent of Total Incorrect predictions due to absence of gene in network
2	3,765	1,016	134	88.4%
3	4,050	1,092	142	88.5%
4	1,240	313	41	88.4%

194

## 195 4 Discussion

196 In this paper, we present a machine learning model enriched with features from prior, known  
197 transcriptional networks to qualitatively and quantitatively predict dysregulation of genes to a  
198 combination of induction conditions at two phases of growth. We showed that the use of a prior  
199 host network adds useful information to a model for it to make predictions of gene dysregulation  
200 for unseen combinations of induction conditions.

201  
202 A natural next question one would have is how our predictions would impact analyses  
203 downstream of DEA, like enrichment analyses. These analyses help researchers gain mechanistic  
204 insight into gene lists generated from DEA <sup>30</sup>. The goal here is to see how different the  
205 mechanistic insights would be if a researcher uses the HRM's predictions from what they would  
206 have observed if they executed the experiments. To answer this question, we added annotations  
207 to the *B. subtilis* genes using SubtiWiki <sup>31</sup> to conduct enrichment analysis of predicted versus  
208 observed DEGs. No gene cluster file was publicly available to use standard enrichment tools. We  
209 created a gene cluster file using data from SubtiWiki (Supplementary Data 2) and used gene set  
210 enrichment analysis <sup>32,33</sup> to identify up and down regulated pathways for each condition (**Figure**  
211 **4A**). We evaluate the predictions by measuring the number of False Negatives (pathways that are  
212 down regulated, but we predict it is not down regulated), False Positives (pathways that are not  
213 down regulated, but we predict are down regulated), True Negatives (pathways that are not down  
214 regulated and we do not predict them to be down regulated), and, finally, True Positives  
215 (pathways that are downregulated and we predict them to be down regulated) (**Figure 4B**). The  
216 same assessment is made for each test condition on the upregulated pathways (**Figure 4C**). As  
217 one would expect, the inducers do not regulate many pathways and the model correctly identifies

218 most of those pathways (gray boxes). Every row is a level 2 category in SubtiWiki that is  
 219 composed of multiple pathways. The number in each box represents the number of pathways in  
 220 the category.  
 221



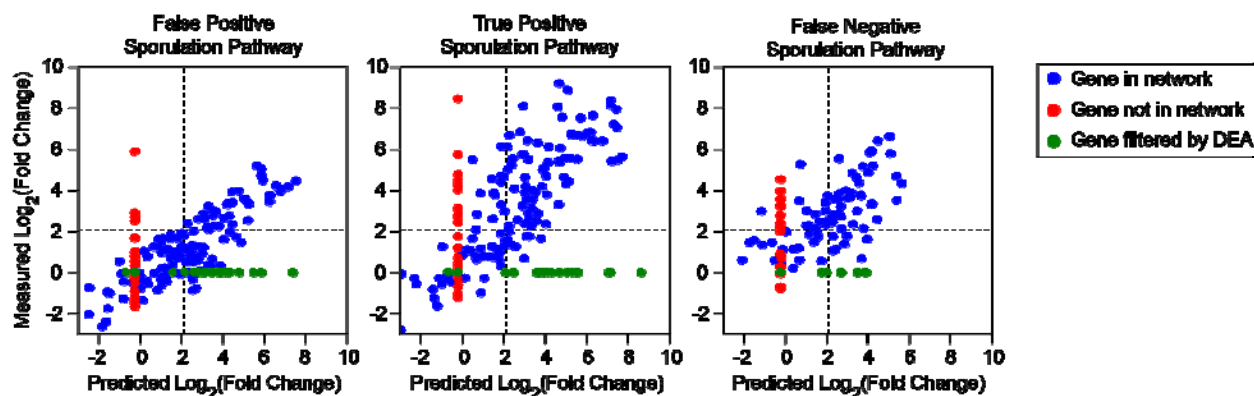
222  
 223 Figure 4: (A) Pathway analysis comparison of actual versus predicted set of DEGs using a gene  
 224 cluster file derived from SubtiWiki. (B-C) Level 2 categories of annotations in the database  
 225 composed of multiple down and up regulated pathways. The numbers inside each colored box  
 226 indicates the number of pathways that were FP, FN, TP, TN. As expected, the inducers do not  
 227 impact a majority of the pathways and the model accurately identifies those pathways.

228

229 Each pathway consists of a set of genes. Pathway analysis uses a statistical test (like a Fisher's  
230 exact test) to identify dysregulated pathways by comparing the list of DEGs to each pathway's  
231 gene set. Thus, we should also check the precision of identifying the regulation of a pathway.  
232 This amounts to seeing how the predicted DEGs compare to the observed DEGs to identify if a  
233 pathway is up or down regulated. We picked a False Positive (FP), False Negative (FN), and  
234 True Positive (TP) sporulation pathway to assess the precision of the predictions (**Figure 5**).  
235 The FP was selected from the CA+xylose condition, while the FN and TP were selected from the  
236 IPTG+CA+VA+xylose condition, all at 18 hours. The genes that are in the bottom left quadrant  
237 of each plot are the ones that contribute to the pathway's down-regulation status in both  
238 predicted and actual settings. The genes in the top left quadrant contribute to the down-regulation  
239 status of the pathway for predicted models but not in the observations. Finally, the genes in the  
240 bottom right quadrant contribute to the pathway's down-regulation from the observations but not  
241 from the predicted model. The vertical set of orange genes are ones that are not present in the  
242 network and so the model predicts a constant for those expression values. The horizontal set of  
243 green genes are genes that did not pass QC and thus could not be validated. As indicated in the  
244 results section, these genes make up the majority of the genes that can be attributed to the  
245 enrichment errors. It is clear from this analysis that future work should consider jointly  
246 optimizing for differential expression as well as pathway inclusion.

247





248  
249 Figure 5: Predicted versus actual  $\text{Log}_2(\text{Fold change})$  of three pathways within the sporulation  
250 category that is a FP, TP, and FN. The genes in the top left quadrant contribute to the down-  
251 regulation status of the pathway for predicted models but not in the observations. Finally, the  
252 genes in the bottom right quadrant contribute to the pathway's down-regulation from the  
253 observations but not from the predicted model. The set of orange genes are ones that are not  
254 present in the network and so the model predicts a constant for those expression values. The  
255 horizontal set of green genes are genes that did not pass QC and thus could not be validated.

## 256 Methods

### 257 Two Stage Learning Model to Predict Differential Expression

258 The goal to predict transcriptional response in a combinatorically large condition space from  
259 single conditions makes an end-to-end learning model with many free parameters underspecified  
260 and prone to generalizability errors<sup>25</sup>. To address this issue, we instead used a two stage learning  
261 process:

- 262 1. **Node embedding of Prior Knowledge:** We applied the node2vec algorithm to derive  
263 vector features from the network that could be used in the downstream learning task<sup>26</sup>.

264 node2vec was selected because it is an unsupervised learning technique that balances depth  
265 and breadth first searches using a random walk to preserve both local and global  
266 connectivity structures of the genes in the network<sup>27-29</sup>. It can be parameterized by the  
267 length of the walks and the number of walks one takes from each node. A skip-gram model  
268 is then used to generate the vector embeddings in  $R^N$ , in our case N was 32. We chose N=32  
269 after an assessment of model predictions on the train/validation set sweeping N between 8,  
270 16, 32, and 64. The performance was not statistically different and so we chose a parameter  
271 that was large enough to provide the model with degrees of freedom to generalize but small  
272 enough to ensure the model does not overfit. We should note that the network need not be  
273 for the exact strain being used, but should have significant genomic overlap. In our case, the  
274 overlap of EColiNet vs the genes in the MG1655 strains was 3818/4111 genes, and for *B.*  
275 *subtilis* was 2608/4266, which is >50% for each organism. Genes that were not present in  
276 the network were mapped to the origin in  $R^{32}$ .

277 **2. Machine Learning Models:** We trained three machine learning models, a gradient boosted  
278 regressor, a linear regressor, and a random forest regressor, for their ability to predict the  
279 differential expression of a gene given the conditions of measurement and the derived  
280 network features. The models were trained on a regression task to minimize the error  
281 between predicted differential expression and the observed differential expression for a  
282 host's response to single inducers. The induction conditions were one hot encoded to enable  
283 the representation of multiple induction conditions. Since there were so few timepoints for  
284 *E. coli* and *B. subtilis*, it was not treated as a continuous variable and was also one hot  
285 encoded.

286

287 The output predictions were evaluated using an  $R^2$  metric comparing predicted to actual  
288 differential expression in  $\log_2(\text{Fold Change})$ . We should note that only genes that were  
289 differentially expressed were used to measure  $R^2$  as those are the ones most significant in  
290 differential expression analysis. For *E. coli*, we defined a gene to be differentially expressed if it  
291 had an absolute magnitude of  $\log_2(\text{Fold Change}) > 1.1$  and FDR of  $< 0.05$ , while for *B. subtilis*  
292 We defined a gene to be differentially expressed if it had an absolute magnitude of  $\log_2(\text{Fold}$   
293  $\text{Change}) > 2.1$  and FDR of  $< 0.05$ .

## 294 Sample Preparation and Processing

295 Wild type strains for *B. subtilis* (Bacillus Subtilis 168 Marburg) and *E. coli* (*E. coli* K-12  
296 MG1655) were cultured in M9 media consisting of 1X M9 media salts, 0.1mM  $\text{CaCl}_2$ , 1X Trace  
297 Salts, 1mM  $\text{MgSO}_4$ , 0.05mM  $\text{FeCl}_3$ /0.1mM  $\text{C}_6\text{H}_8\text{O}_7$ , 0.2% Casamino Acids, and 0.4%  
298 Glucose. The inducers used in this study were isopropyl  $\beta$ -D-1-thiogalactopyranoside (0.001  
299 M), arabinose (25mM), vanillic acid (0.001 M), cuminic acid (0.0001 M), and xylose (1%).

300

301 Glycerol stocks were inoculated into M9 media in shake flasks, and the culture was grown  
302 overnight for 18h at  $30^\circ\text{C}$  and 1000rpm. The following day, cultures were diluted to OD 0.1 in  
303 fresh M9 media and grown in 96-well plates under the same conditions for 3 h. For induction,  
304 cells were diluted a second time to OD 0.05 in the presence of inducers. Plates were incubated at  
305  $30^\circ\text{C}$  and 1000 rpm for 5 h and 18 h and cultured cells were harvested and fixed with either  
306 RNA protect (for *E. coli*) or methanol (*B. subtilis*).

307

308 Total RNA was extracted using Magjet RNA extraction kit (Thermo) according to  
309 manufacturer's instructions. RNA quality was assessed using Tapestation (Agilent). KAPA RNA  
310 Hyperprep kit (Roche) was used for ribosomal RNA depletion and Illumina compatible library  
311 preparation. Prepared library was loaded on a Illumina sequencer to generate 150bp paired end  
312 reads.

313

314 Raw RNA-seq data was trimmed and quality filtered with trimmomatic (v0.36), reads were  
315 aligned with bwa (v0.7.17). After alignment with bwa, the resulting sam files were sorted by  
316 PICARD tools (v2.18.15) function SortSam, and then AddOrReplaceGroups is run on the sorted  
317 sam. Gene-level quantification of counts was performed using the featureCounts function of  
318 Rsubread (v1.34.4).

## 319 Samples and Transcript Quality Control for *B. subtilis*

320 A measure most often used to qualitatively and quantitatively assess a transcript's response is its  
321 dysregulation as compared to a control. Differential expression analysis (DEA) is a standard  
322 bioinformatics technique that measures this response to perturbations as compared to a control  
323 condition<sup>19</sup>. DEA conducts custom normalization, dispersion modeling, and Bayesian  
324 optimization to account for biological and experimental variability that is present in  
325 transcriptional counts data to quantify the transcriptional response to a perturbation and measure  
326 its statistical significance. While this method overcomes generalization errors that can arise from  
327 artifacts of normalization of counts data across experiments, it performs strict quality control  
328 (QC) at both the sample and gene level. These are listed below:

- 329 1. Sample QC = The significance tests to reject the null hypothesis of the differentially  
330 expressed genes require  $> 2$  samples per condition. If this criterion is not met then DEA  
331 cannot be conducted.
- 332 2. Gene QC = DEA tools fits the Cox-Reid profile-adjusted dispersion to a set of  
333 normalized expressions across conditions. Genes that do not fit this profile are labeled as  
334 outliers and removed from DEA.

335 Three Boolean metrics were used to measure sample quality:

- 336 1. Number of mapped reads  $\geq 500K$   
337 2. Count of all annotated genes  $\geq 500K$   
338 3. Replicate correlation of a condition  $> 0.9$

339

340 If any of these metrics did not pass, the sample would be flagged as a low quality sample and not  
341 used for downstream analysis. For the *B. subtilis* experiments, we also collected OD600  
342 measurements from a plate reader to correlate population with potential sample dropouts. We did  
343 not find a clear, discriminative correlation between this measurement and sample dropout for a  
344 condition, but the log phase measurements (timepoint 5.0) did have a lower OD on average and  
345 had twice as many samples that did not pass QC than stationary phase samples (timepoint 18.0).  
346 In conditions where only two replicates were available, differential expression analysis was not  
347 conducted and so those conditions could not be validated (Supplementary Table 2). All single  
348 inducer conditions for *B. subtilis* that passed QC were used to train the model. It should be noted,  
349 though, that if a sample passed QC, that did not mean all genes in that sample passed edgeR's  
350 outlier detection method. edgeR fits normalized counts to a Cox-Ried dispersion model with a  
351 Bayesian optimization algorithm. A gene is removed by edgeR if it does not fit this dispersion

352 model. While one can pass in a custom dispersion model per condition, we chose to use edgeR's  
353 default, as development of custom noise models across the condition space was out of scope of  
354 this effort (Supplementary Table 3).

## 355 Funding

356 Any opinions, findings and conclusions or recommendations expressed in this material are those of the  
357 author(s) and do not necessarily reflect the views of the Defense Advanced Research Projects Agency  
358 (DARPA), the Department of Defense, or the United States Government. This material is based upon  
359 work supported by the Defense Advanced Research Projects Agency (DARPA) and the Air Force  
360 Research Laboratory under Contract No. FA8750-17-C-0231 (and related contracts by SD2 Publication  
361 Consortium Members).

## 362 Data and Code Availability

363 The manuscript is accompanied with three code repositories that are fully documented with  
364 example python notebooks. The data for the publication is placed with the tutorials to ensure  
365 reproducibility of results.

- 366 1. A repository that includes the capability to train, validate, and test a machine learning  
367 model in a combinatorically large condition space. <https://github.com/sd2e/CDM>
- 368 2. A repository that includes a scaled, configurable differential expression analysis pipeline:  
369 [https://github.com/SD2E/omics\\_tools](https://github.com/SD2E/omics_tools).
- 370 3. A test-harness for machine learning models to make apples-to-apples comparisons of  
371 training and testing models: <https://github.com/SD2E/test-harness>

## 372 Author Contributions

373 M.E., A. E. B, H. E., and E. Y. led the analysis plan, analysis, and design of experiments for E.  
374 coli and B. subtilis. H.D. also contributed to B. subtilis experiments. D. B., C. B., P. M, K. C.,  
375 performed the experiments. J. U., M. W., M. V., G. Z., N. G., J. F., and J. S. designed and  
376 automated the data processing, quality control, and data/analysis collaboration infrastructure. M.  
377 E., G. Z, and A. C. developed, scaled, and automated the execution the differential expression  
378 analysis pipeline. M. E. and H. E. developed the machine learning test harness. M. E. and H. E.  
379 developed the Host Response Model Library. C. V., B.G., J.S., and Y.D are PIs that managed the  
380 program and provided technical guidance. M. E., A. E. B., and E. Y. prepared the manuscript.

## 381 Competing Interests

382 The authors report no competing interests.

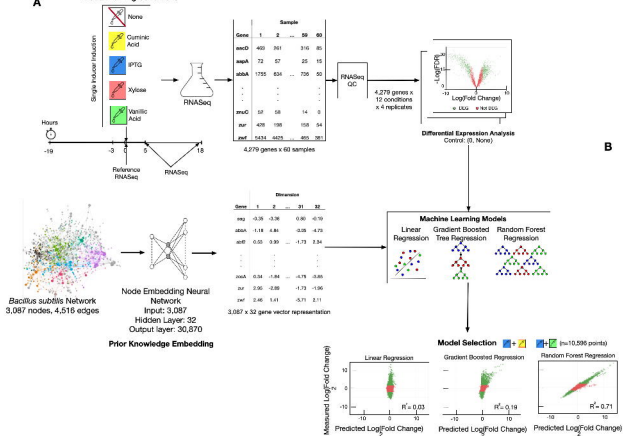
383

## 384 References

- 385  
386 1. Helmann, J. D. *et al.* Global transcriptional response of *Bacillus subtilis* to heat shock. *J.*  
387 *Bacteriol.* **183**, 7318–7328 (2001).  
388 2. Gao, H. *et al.* Global transcriptome analysis of the heat shock response of *Shewanella*  
389 *oneidensis*. *J. Bacteriol.* **186**, 7796–7803 (2004).  
390 3. Hengge-Aronis, R. Back to log phase: sigma S as a global regulator in the osmotic control  
391 of gene expression in *Escherichia coli*. *Mol. Microbiol.* **21**, 887–893 (1996).  
392 4. Soufi, B. *et al.* Global analysis of the yeast osmotic stress response by quantitative  
393 proteomics. *Mol. Biosyst.* **5**, 1337–1346 (2009).  
394 5. Erickson, D. W. *et al.* A global resource allocation strategy governs growth transition  
395 kinetics of *Escherichia coli*. *Nature* **551**, 119–123 (2017).  
396 6. Alexander, H., Rouco, M., Haley, S. T. & Dyhrman, S. T. Transcriptional response of  
397 *Emiliana huxleyi* under changing nutrient environments in the North Pacific Subtropical  
398 Gyre. *Environ. Microbiol.* **22**, 1847–1860 (2020).  
399 7. Dunn, T. M., Hahn, S., Ogden, S. & Schleif, R. F. An operator at -280 base pairs that is  
400 required for repression of *araBAD* operon promoter: addition of DNA helical turns between  
401 the operator and promoter cyclically hinders repression. *Proc. Natl. Acad. Sci. USA* **81**,  
402 5017–5020 (1984).  
403 8. Harmer, T., Wu, M. & Schleif, R. The role of rigidity in DNA looping-unlooping by AraC.  
404 *Proc. Natl. Acad. Sci. USA* **98**, 427–431 (2001).  
405 9. Martin, K., Huo, L. & Schleif, R. F. The DNA loop model for *ara* repression: AraC protein  
406 occupies the proposed loop sites in vivo and repression-negative mutations lie in these  
407 same sites. *Proc. Natl. Acad. Sci. USA* **83**, 3654–3658 (1986).  
408 10. Brunton, S. L., Proctor, J. L. & Kutz, J. N. Discovering governing equations from data by  
409 sparse identification of nonlinear dynamical systems. *Proc. Natl. Acad. Sci. USA* **113**,  
410 3932–3937 (2016).  
411 11. Champion, K., Lusch, B., Kutz, J. N. & Brunton, S. L. Data-driven discovery of coordinates  
412 and governing equations. *Proc. Natl. Acad. Sci. USA* **116**, 22445–22451 (2019).  
413 12. Adam, G. *et al.* Machine learning approaches to drug response prediction: challenges and  
414 recent progress. *NPJ Precis. Oncol.* **4**, 19 (2020).  
415 13. Fitzgerald, J. B., Schoeberl, B., Nielsen, U. B. & Sorger, P. K. Systems biology and  
416 combination therapy in the quest for clinical efficacy. *Nat. Chem. Biol.* **2**, 458–466 (2006).  
417 14. Kuru, H. I., Tastan, O. & Cicek, A. E. Matchmaker: A deep learning framework for drug  
418 synergy prediction. *BioRxiv* (2020). doi:10.1101/2020.05.24.113241  
419 15. Li, J., Tong, X.-Y., Zhu, L.-D. & Zhang, H.-Y. A machine learning method for drug  
420 combination prediction. *Front. Genet.* **11**, 1000 (2020).  
421 16. Chen, G., Tsoi, A., Xu, H. & Zheng, W. J. Predict effective drug combination by deep belief  
422 network and ontology fingerprints. *J. Biomed. Inform.* **85**, 149–154 (2018).  
423 17. Xue, Y., Ding, M. Q. & Lu, X. Learning to encode cellular responses to systematic  
424 perturbations with deep generative models. *NPJ Syst. Biol. Appl.* **6**, 35 (2020).  
425 18. Costa-Silva, J., Domingues, D. & Lopes, F. M. RNA-Seq differential expression analysis:  
426 An extended review and a software tool. *PLoS One* **12**, e0190152 (2017).  
427 19. Schmauch, B. *et al.* A deep learning model to predict RNA-Seq expression of tumours from  
428 whole slide images. *Nat. Commun.* **11**, 3877 (2020).  
429 20. Schmidt, F., Kern, F. & Schulz, M. H. Integrative prediction of gene expression with  
430 chromatin accessibility and conformation data. *Epigenetics Chromatin* **13**, 4 (2020).  
431 21. Abbas-Aghababazadeh, F., Li, Q. & Fridley, B. L. Comparison of normalization approaches  
432 for gene expression studies completed with high-throughput sequencing. *PLoS One* **13**,



- 433 e0206312 (2018).  
434 22. Kim, H., Shim, J. E., Shin, J. & Lee, I. EcoliNet: a database of cofunctional gene network for  
435 Escherichia coli. *Database (Oxford)* **2015**, (2015).  
436 23. Arrieta-Ortiz, M. L. *et al.* An experimentally supported model of the Bacillus subtilis global  
437 transcriptional regulatory network. *Mol. Syst. Biol.* **11**, 839 (2015).  
438 24. Robinson, M. D., McCarthy, D. J. & Smyth, G. K. edgeR: a Bioconductor package for  
439 differential expression analysis of digital gene expression data. *Bioinformatics* **26**, 139–140  
440 (2010).  
441 25. D'Amour, A. *et al.* Underspecification Presents Challenges for Credibility in Modern  
442 Machine Learning. *arXiv* (2020).  
443 26. Grover, A. & Leskovec, J. node2vec: Scalable Feature Learning for Networks. *KDD* **2016**,  
444 855–864 (2016).  
445 27. Kim, M., Baek, S. H. & Song, M. Relation extraction for biological pathway construction  
446 using node2vec. *BMC Bioinformatics* **19**, 206 (2018).  
447 28. Ata, S. K. *et al.* Integrating node embeddings and biological annotations for genes to  
448 predict disease-gene associations. *BMC Syst. Biol.* **12**, 138 (2018).  
449 29. Nelson, W. *et al.* To embed or not: network embedding as a paradigm in computational  
450 biology. *Front. Genet.* **10**, 381 (2019).  
451 30. Reimand, J. *et al.* Pathway enrichment analysis and visualization of omics data using  
452 g:Profiler, GSEA, Cytoscape and EnrichmentMap. *Nat. Protoc.* **14**, 482–517 (2019).  
453 31. Zhu, B. & Stülke, J. SubtiWiki in 2018: from genes and proteins to functional network  
454 annotation of the model organism Bacillus subtilis. *Nucleic Acids Res.* **46**, D743–D748  
455 (2018).  
456 32. Chen, E. Y. *et al.* Enrichr: interactive and collaborative HTML5 gene list enrichment  
457 analysis tool. *BMC Bioinformatics* **14**, 128 (2013).  
458 33. Kuleshov, M. V. *et al.* Enrichr: a comprehensive gene set enrichment analysis web server  
459 2016 update. *Nucleic Acids Res.* **44**, W90-7 (2016).  
460

**A****Model Training Conditions****B****Model Evaluation Conditions**

2 Timepoints (5 hr, 18 hr)

3 Double Inducers

3 Triple Inducers

1 Quadruple Inducer

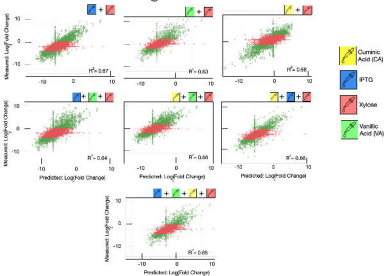


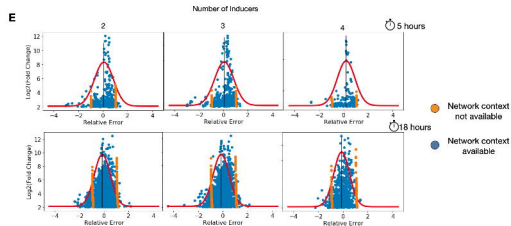
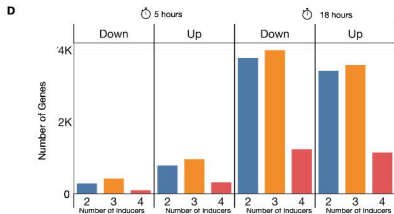
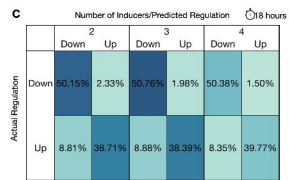
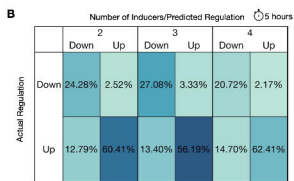
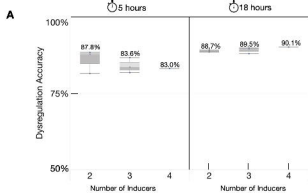
RNASeq Measurement

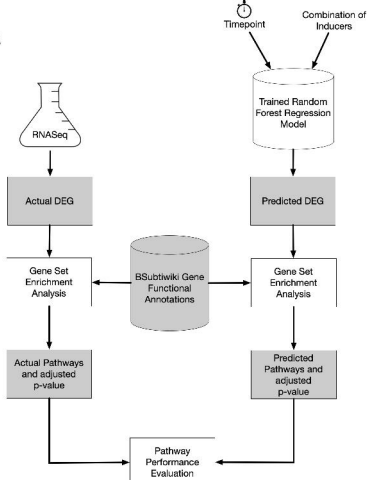


Trained Random Forest Regression Model

18 hours





**A****B**

	Up Regulated Pathways						18 Hours					
(Post-) Exponential lifestyles	17	17	17	17	17	17	17	17	17	17	17	
Additional metabolic pathways	36	1	35	36	35	36	1	35	36			
Amino acid/ nitrogen metabolism	29	1	27	1	29	28	1	29	28	1		
Carbon metabolism	33	33	33	33	33	33	33	33	33			
Cell envelope and cell division	21	21	21	21	21	21	21	21	21			
Coping with stress	1	20	1	20	1	20	1	20	1	20		
Detoxification reactions	1	1	1	1	1	1	1	1	1	1		
Elp-dependent proteins	1	1	1	1	1	1	1	1	1	1		
Electron transport and ATP synthesis	9	9	9	9	9	9	9	9	9	9		
Essential genes	1	1	1	1	1	1	1	1	1	1		
Genetics	16	16	16	16	16	16	16	16	16			
GTP-binding proteins	1	1	1	1	1	1	1	1	1	1		
Homeostasis	13	13	13	13	13	13	13	13	13			
Lifestyle/ miscellaneous	1	1	1	1	1	1	1	1	1	1		
Lipid metabolism	7	1	7	1	7	1	7	1	7	1		
Membrane proteins	1	1	1	1	1	1	1	1	1	1		
Mobile genetic elements	2	2	2	2	2	2	2	2	2	2		
NcRNA	1	1	1	1	1	1	1	1	1	1		
Nucleotide metabolism	9	9	9	9	9	9	9	9	9	9		
Phosphoproteins	9	9	9	9	9	9	9	9	9	9		
Poorly characterized/ putative enzymes	1	1	1	1	1	1	1	1	1	1		
Prophages	6	6	6	6	6	6	6	6	6	6		
Protein synthesis, mods, degradation	27	27	27	27	27	27	27	27	27	27		
Proteins of unknown function	1	1	1	1	1	1	1	1	1	1		
Pseudogenes	1	1	1	1	1	1	1	1	1	1		
Quasi-essential genes	1	1	1	1	1	1	1	1	1	1		
Regulation of gene expression	28	28	28	28	28	28	28	28	28	28		
RNA synthesis and degradation	12	12	12	12	12	12	12	12	12	12		
Secreted proteins	1	1	1	1	1	1	1	1	1	1		
Short peptides	1	1	1	1	1	1	1	1	1	1		
Sporulation	2	19	2	19	2	19	2	19	2	19	20	
Targets of second messengers	4	4	4	4	4	4	4	4	4	4		
Transporters	1	41	1	40	1	41	1	40	1	41	41	
Universally conserved proteins	1	1	1	1	1	1	1	1	1	1		

■ False Negative  
■ False Positive  
■ True Negative  
■ True Positive

Inducers	CA	IPTG	VA	Xylose
CA	-	+	+	+
IPTG	+	-	+	+
VA	+	+	-	+
Xylose	+	+	+	-

**C**

	Down Regulated Pathways												18 Hours					
(Post-) Exponential lifestyles	2	14	2	15	2	14	1	16	1	15	1	15	2	14	1	15		
Additional metabolic pathways	2	34	1	35	2	34	1	35	2	34	1	35	2	34	1	35		
Amino acid/ nitrogen metabolism	1	28	1	29	1	28	29	29	29	29	29	29	29	29	29	29		
Carbon metabolism	33	33	33	33	33	33	33	33	33	33	33	33	33	33	33	33		
Cell envelope and cell division	21	21	21	21	21	21	21	21	21	21	21	21	21	21	21	21		
Coping with stress	2	21	2	21	2	21	2	21	2	21	2	21	2	21	2	21		
Detoxification reactions	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1		
Elp-dependent proteins	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1		
Electron transport and ATP synthesis	9	9	9	9	9	9	9	9	9	9	9	9	9	9	9	9		
Essential genes	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1		
Genetics	1	15	1	15	1	15	16	16	16	16	16	16	16	16	16	16		
GTP-binding proteins	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1		
Homeostasis	1	12	1	12	1	12	13	13	13	13	13	13	13	13	13	13		
Lifestyle/ miscellaneous	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1		
Lipid metabolism	9	9	9	9	9	9	9	9	9	9	9	9	9	9	9	9		
Membrane proteins	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1		
Mobile genetic elements	2	2	2	2	2	2	2	2	2	2	2	2	2	2	2	2		
NcRNA	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1		
Nucleotide metabolism	7	7	7	7	7	7	7	7	7	7	7	7	7	7	7	7		
Phosphoproteins	9	9	9	9	9	9	9	9	9	9	9	9	9	9	9	9		
Poorly characterized/ putative enzymes	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1		
Prophages	6	6	6	6	6	6	6	6	6	6	6	6	6	6	6	6		
Protein synthesis, mods, degradation	27	27	27	27	27	27	27	27	27	27	27	27	27	27	27	27		
Proteins of unknown function	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1		
Pseudogenes	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1		
Quasi-essential genes	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1		
Regulation of gene expression	28	28	28	28	28	28	28	28	28	28	28	28	28	28	28	28		
RNA synthesis and degradation	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12		
Secreted proteins	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1		
Short peptides	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1		
Sporulation	2	22	2	22	2	22	22	22	22	22	22	22	22	22	22	22		
Targets of second messengers	4	4	4	4	4	4	4	4	4	4	4	4	4	4	4	4		
Transporters	1	41	1	41	1	41	42	42	42	42	42	42	42	42	42	42		
Universally conserved proteins	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1		

■ False Negative  
■ False Positive  
■ True Negative  
■ True Positive

Inducers	CA	IPTG	VA	Xylose
CA	-	+	+	+
IPTG	+	-	+	+
VA	+	+	-	+
Xylose	+	+	+	-

