

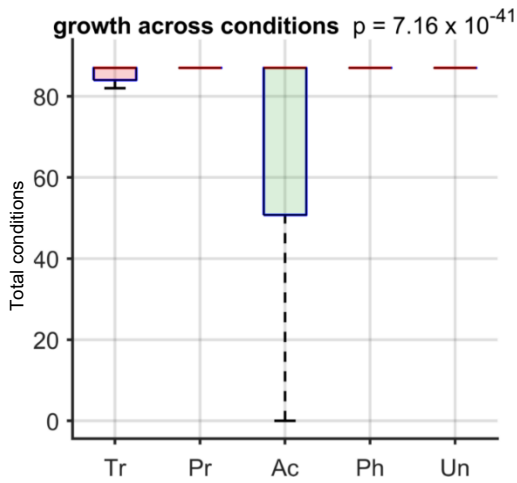
iScience, Volume 25

Supplemental information

**Metabolic signatures of regulation
by phosphorylation and acetylation**

Kirk Smith, Fangzhou Shen, Ho Joon Lee, and Sriram Chandrasekaran

1



2

3 Figure S1. Distribution of regulation based on gene essentiality across 87 different conditions, Related to Figure
4 1. These conditions comprise 56 different carbon sources including glucose, and 31 different nitrogen sources
5 including ammonium ions. The total number of conditions in which each gene deletion was viable was calculated.
6 This total number was then compared between targets of each regulatory mechanism. The box plots show that
7 acetylation preferentially regulates the genes that impact growth across the 87 conditions. The box plot whiskers
8 extend to the 99.3rd percentile of each distribution. The ANOVA p-value comparing the means is 7.1×10^{-41} .

9

10

11

12

13

14

15

16

17

18

19

20

21

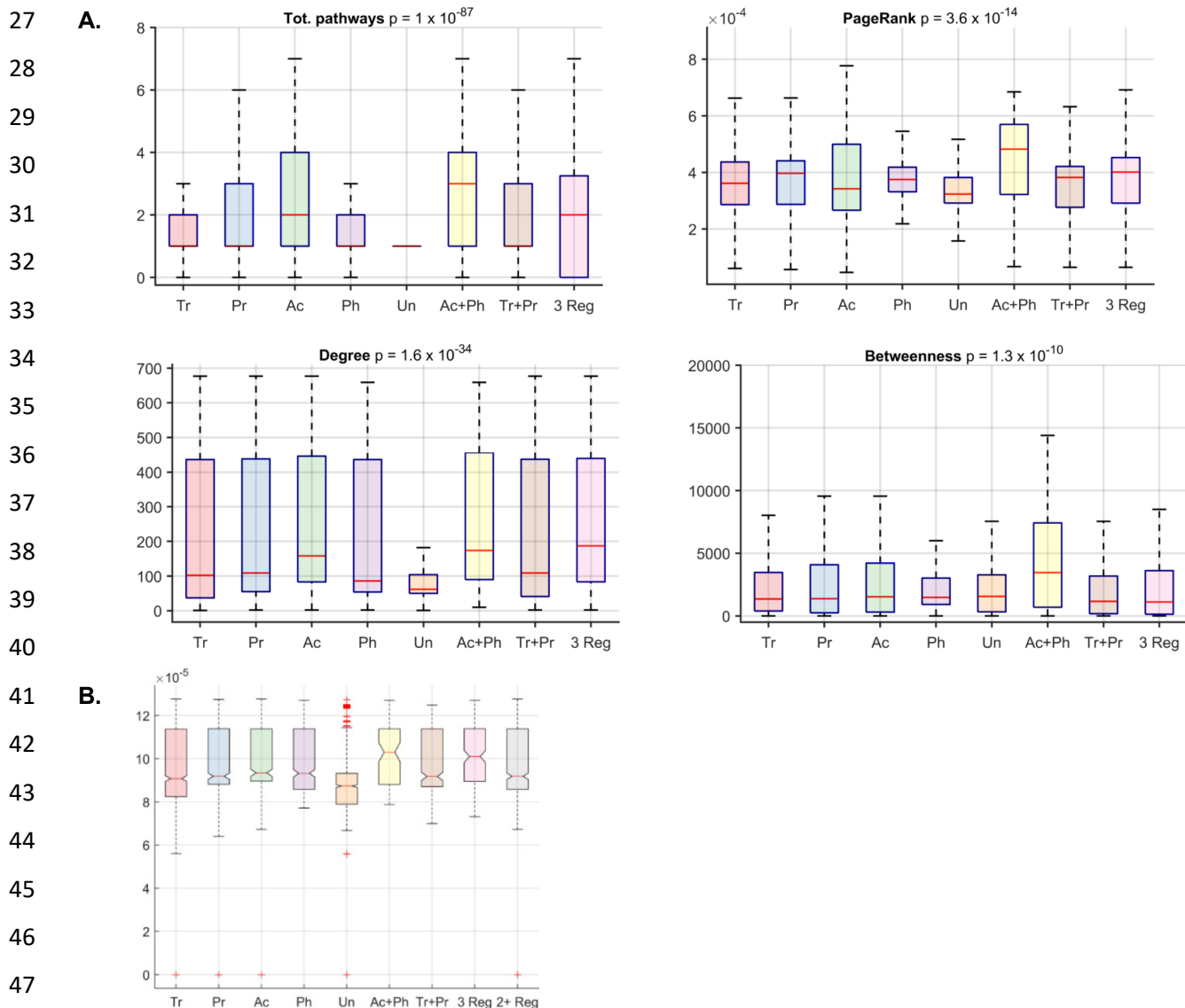
22

23

24

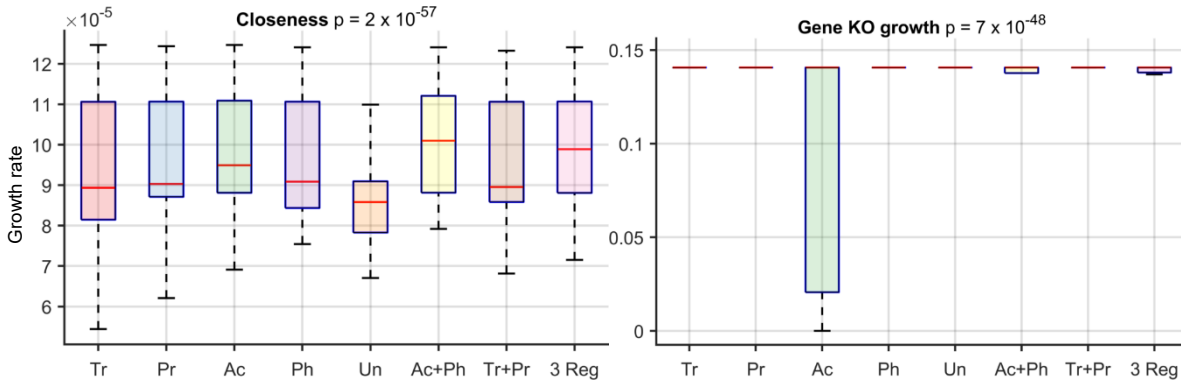
25

26

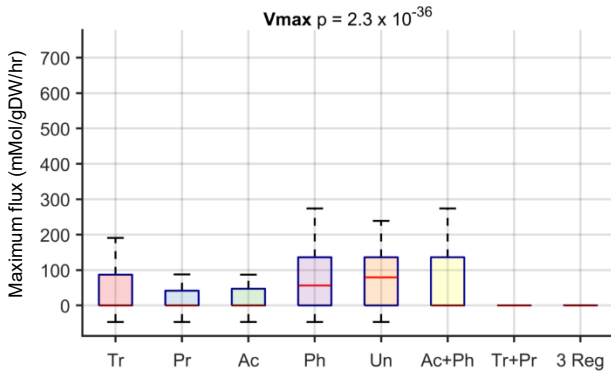


48 **Figure S2. Distribution of regulation based on topological properties of each reaction, Related to Figure 1. A.** Four
 49 different topological properties are shown in the box plots - the total number of annotated pathways each reaction
 50 participates (Tot. pathways), the number of times each reaction is traversed during a random walk between
 51 reactions in the network (PageRank), the total number of connected reactions (Degree) and the number of times
 52 each reaction appears on a shortest path between two reactions (Betweenness). These show that reactions that
 53 are regulated by any mechanism have a higher connectivity compared to those that are unregulated or regulated
 54 by unknown mechanisms. Furthermore, reactions regulated by both acetylation and phosphorylation had the
 55 highest connectivity across all metrics. The ANOVA p-value comparing the means is provided in the title.
 56 (Abbreviations: regulation by both transcription and post-transcription (Tr + Pr), both acetylation and
 57 phosphorylation (Ac + Ph), at least 3 regulators (3 Reg), and Unknown regulation (Un)). **B.** Demonstration of
 58 robustness of topological analysis. Highly connected metabolites (ATP ADP AMP NADH NAD) were removed
 59 from the yeast model prior to the calculation of topological parameters. The box plots compare the properties of
 60 enzymes regulated by transcription (Tr), post-transcription (Pr), acetylation (Ac), phosphorylation (Ph), both
 61 transcription and post-transcription (Tr + Pr), both acetylation and phosphorylation (Ac + Ph), or at least 3
 62 regulators (3 Reg). Reactions regulated by both acetylation and phosphorylation had the highest connectivity as
 63 measured by the Closeness. The ANOVA p-value comparing the means is $3e-46$ for closeness, $2e-29$ for degree
 64 (not shown) and $5e-15$ for pagerank (not shown).

65



66



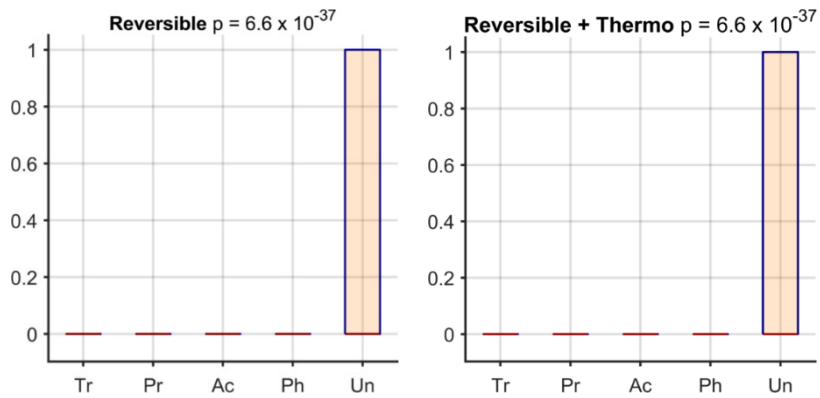
67

68 **Figure S3. Properties of reactions regulated by multiple mechanisms, Related to Figure 1.** The box plots compare
69 the properties of enzymes regulated by transcription, post-transcription, acetylation, phosphorylation with those
70 regulated by both transcription and post-transcription (Tr + Pr), both acetylation and phosphorylation (Ac + Ph), or
71 at least 3 regulators (3 Reg). This set of combinations among regulators was chosen as both acetylation and
72 phosphorylation are PTMs, and the transcriptome and proteome of yeast cells show significant correlation.
73 Reactions regulated by both acetylation and phosphorylation had the highest connectivity as measured by the
74 inverse sum of the distance from a reaction to all other reactions in the network (Closeness). Apart from
75 connectivity, reactions regulated by two different mechanisms did not share properties of reactions regulated by
76 each individual mechanism. For example, reactions regulated by acetylation and phosphorylation were not likely
77 to be essential or have high maximum flux. The ANOVA p-value comparing the means is provided in the title.

78

79

80



81

82 **Figure S4. Distribution of regulation based on reaction reversibility, Related to Figure 1.** Reversible reactions were
 83 highly likely to be not regulated by any of the four mechanisms. The left panel compares the distribution of
 84 regulation of reversible reactions based on the annotation from the Yeast 7 model (reversible reactions are set to
 85 1 and irreversible reactions are set to 0). The panel on the right uses an updated list based on thermodynamic
 86 analysis of the Yeast metabolic model by Martinez *et al* [49].

87

88

89

90

91

92

93

94

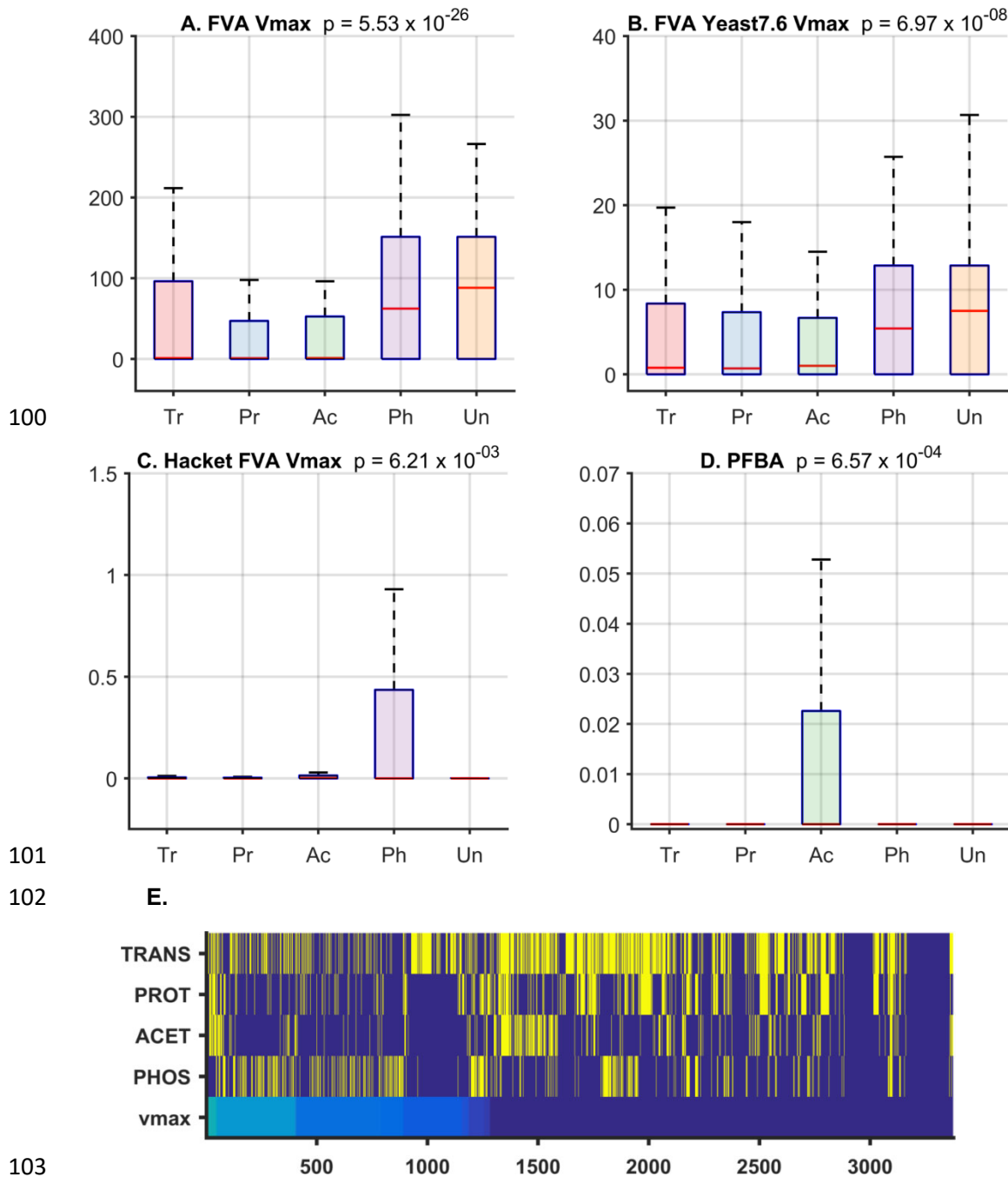
95

96

97

98

99



104 Figure S5. Distribution of regulation based on magnitude of maximum possible flux (mmol/gDW/hr) through each
 105 reaction, Related to Figure 1. The plots compare the distribution of regulation using flux calculated using various
 106 methods and models. The ANOVA p-value comparing the means is provided in the panel title of each plot. These
 107 results show that phosphorylated reactions are highly enriched among those reactions with high maximum flux. **A.**
 108 Maximum flux through each reaction was calculated using FVA using the Yeast 7 model without assuming that
 109 cells maximize their biomass (the default objective in FVA and FBA). The box plots compare the maximum flux
 110 value of reactions regulated by each mechanism. **B.** Maximum flux through each reaction was calculated using
 111 FVA without assuming that cells maximize their biomass using the Yeast 7.6 model (Yeast 7 model was used for
 112 all analyses). **C.** The flux through the model was first fit to the experimentally inferred flux data from Hackett *et*
 113 *a*[21]. The maximum flux through all reactions was then determined using FVA. **D.** The flux through each
 114 reaction was inferred from Parsimonious FBA (PFBA). Note that PFBA does not provide the maximum flux but the

115 flux value that minimizes the sum of flux through all reactions while maximizing the biomass objective. Hence it
116 does not reveal any futile cycles or redundancy in the network. **E.** The heatmap shows the distribution of
117 regulation based on magnitude of maximum possible flux (V_{max}) through of each reaction. Reactions are sorted
118 based on V_{max} inferred from FVA. The columns correspond to each reaction-gene pair. Those that are regulated
119 by each mechanism are shown in yellow, while those that are not regulated by a specific mechanism are in blue.

120

121

122

123

124

125

126

127

128

129

130

131

132

133

134

135

136

137

138

139

140

141

142

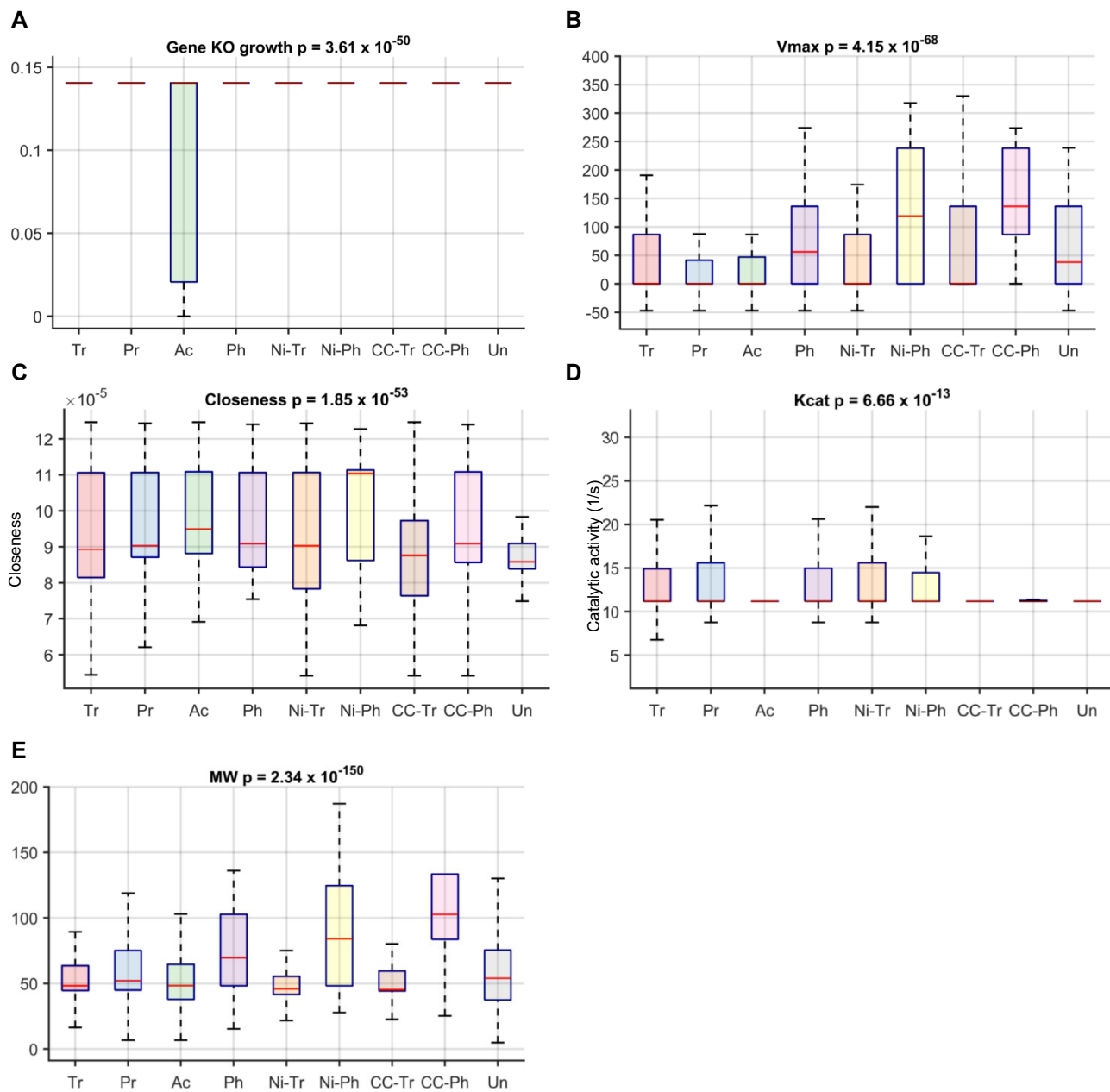
143

144

145

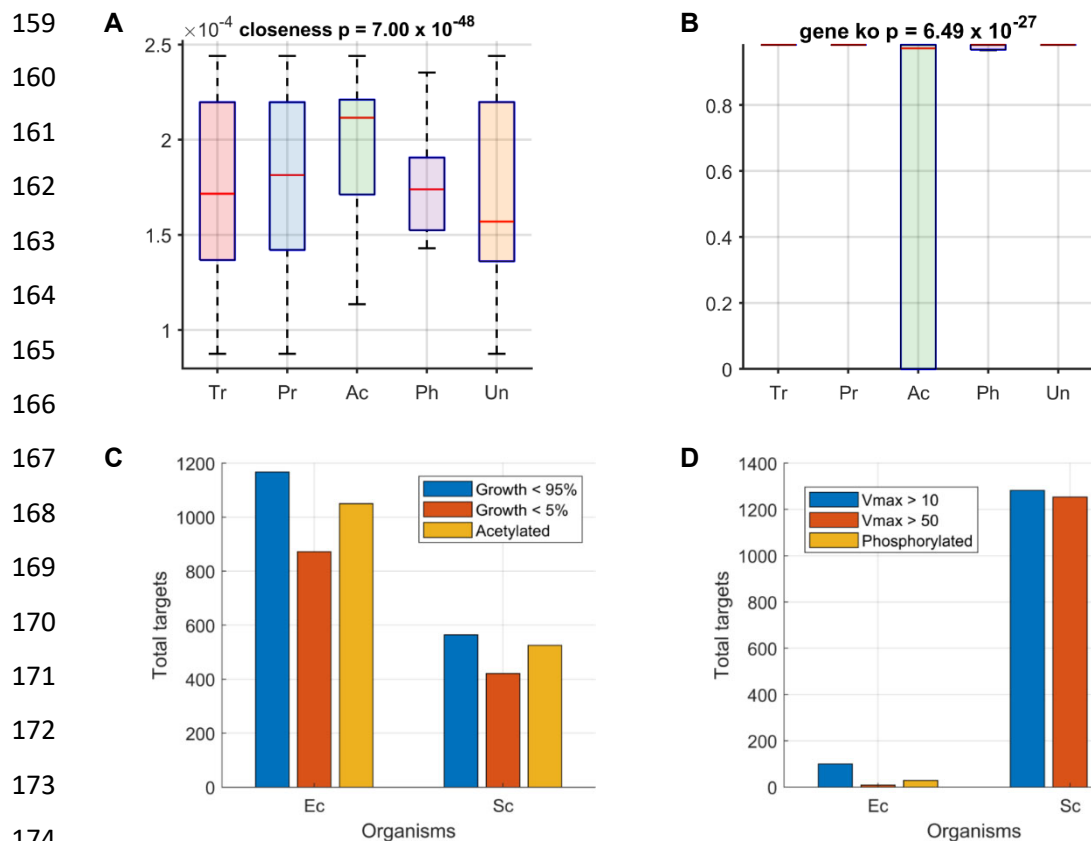
146

147



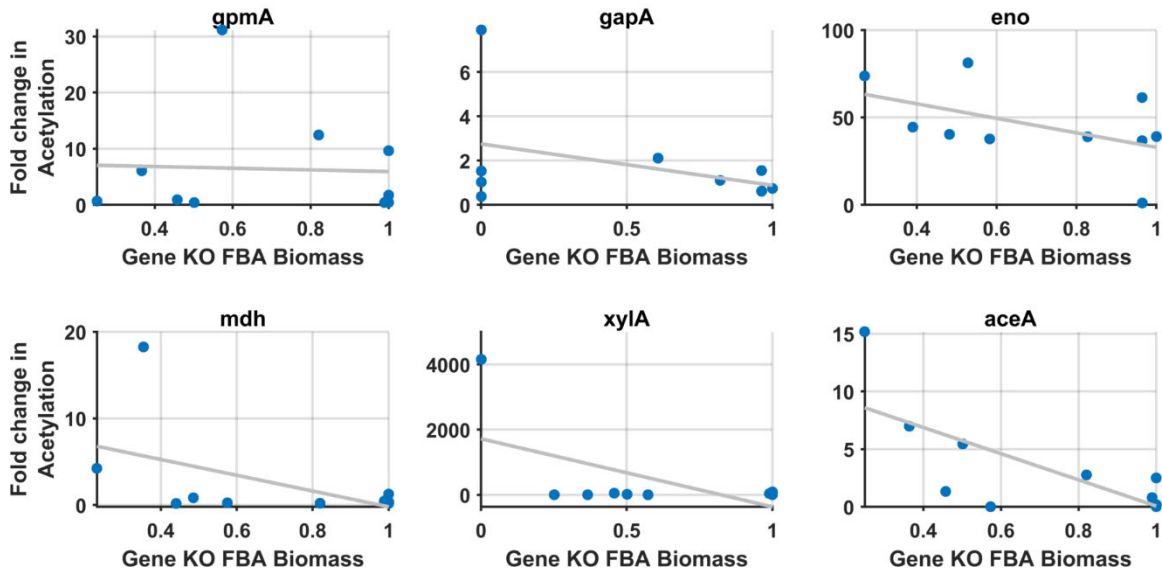
148 Figure S6. Comparison of the properties of enzymes in yeast regulated by each mechanism during the cell cycle
 149 (CC-Tr, CC-Ph) and nitrogen starvation (Ni-Tr, Ni-Ph), Related to Figure 1. Data from stationary phase conditions
 150 (transcription (Tr), post-transcription (Pr), acetylation (Ac), phosphorylation (Ph) or Unknown regulation (Un)) are
 151 shown for comparison. Similar to stationary phase, enzymes that impact growth when knocked out are likely to be
 152 acetylated (A), enzymes that catalyze reactions with high flux are likely to be regulated through phosphorylation in
 153 all three conditions (B), enzymes that are highly connected are likely to be regulated by one of the four
 154 mechanisms (C). No consistent difference across datasets was observed in regulation based on the enzyme
 155 catalytic activity (kcat) of the target enzyme (D) and enzymes regulated by phosphorylation on average tend to
 156 have high molecular weight (E). The Anova p-value comparing the differences in means is shown in the title.

157
 158



175 **Figure S7. Comparison of properties of enzymes in *E. coli* regulated by each mechanism, Related to Figure 2. A.**
 176 Similar to yeast, enzymes that are highly connected (i.e. high closeness) are likely to be regulated. **B.** Similar to
 177 our analysis in Figure 2A, which showed using the entire set of acetylated proteins the association between
 178 acetylation regulation and growth impacting enzymes, this figure shows that the subset of acetylated proteins
 179 regulated by the deacetylase *cobB* also show the same trend with reactions that impact growth when knocked out
 180 are highly likely to be acetylated and regulated by *cobB*. The Anova p-value comparing the differences in means
 181 is shown in the title. **C, D.** Comparison of total number of targets between species. Total number of regulation
 182 targets (i.e. gene-reactions) of PTMs in *E. coli* (Ec) and yeast (Sc) are compared with those that have high Vmax
 183 and are growth limiting in those species in the stationary phase condition.

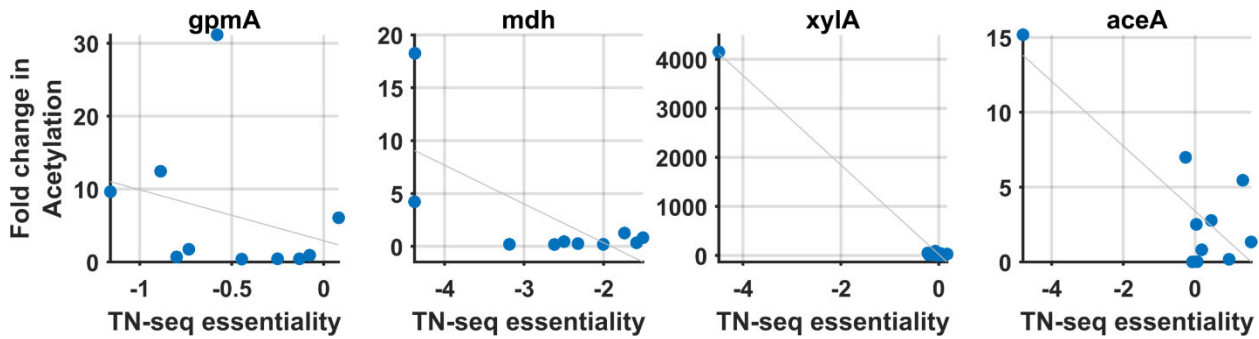
184
 185
 186



187

188

189 **Figure S8. Condition-specific essentiality is correlated with acetylation, Related to Figure 2.** The scatter plots
 190 show the association between the impact of a gene knockout on biomass from FBA with the acetylation levels of
 191 the corresponding protein in a given condition. On average, increased essentiality is associated with an increase
 192 in acetylation. All proteins with at least 2 fold change in acetylation between conditions and are part of the
 193 metabolic model are shown. The change in biomass relative to glucose is show in the x-axis. The correlations
 194 were observed even when the total absolute acetylation levels were considered instead of relative levels to
 195 proteins.

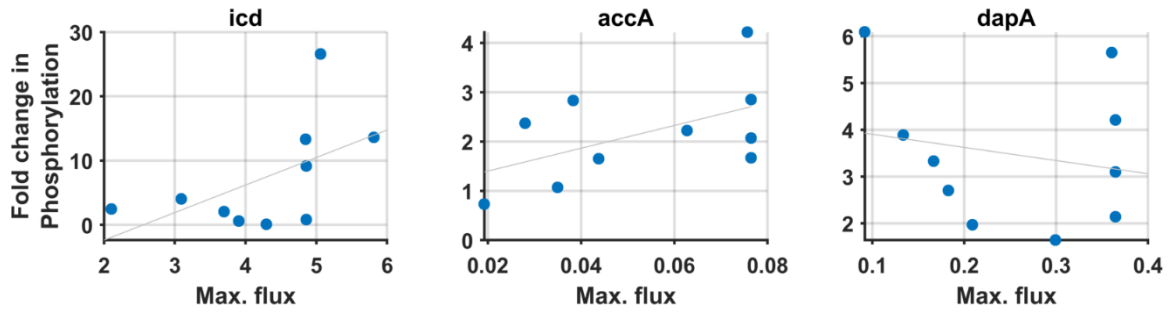


196

197

198 **Figure S9. Condition-specific essentiality from TN-seq is correlated with acetylation, Related to Figure 2.** The
 199 scatter plots show the association between the impact of a gene knockout on viability from Transposon
 200 mutagenesis screens with the acetylation levels of the corresponding protein in a given condition. All proteins in
 201 the metabolic model with available TN-seq data and acetylation data across conditions from Schmidt et al study
 202 are shown. Although FBA made false positive growth predictions for some enzymes such as XylA (Figure S8), our
 203 results were observed even with experimentally derived knockout screens, suggesting that this link between
 204 essentiality and acetylation is robust.

205



206

207

208 [Figure S10. Correlation between maximum flux and phosphorylation levels \(normalized to glucose\), Related to](#)
 209 [Figure 2.](#) All proteins that showed at least 2-fold change in phosphorylation levels between conditions are shown.
 210 This trend was observed with both the total phosphorylation levels and relative levels normalized to proteins.
 211 While in most cases a change in maximal flux or essentiality resulted in a change in regulation by PTMs (Figure
 212 2F), there were exceptions. For example, *dapA* did not show this trend suggesting that other factors likely
 213 influence regulation by PTMs in a combinatorial fashion.

214

215

216

217

218

219

220

221

222

223

224

225

226

227

228

229

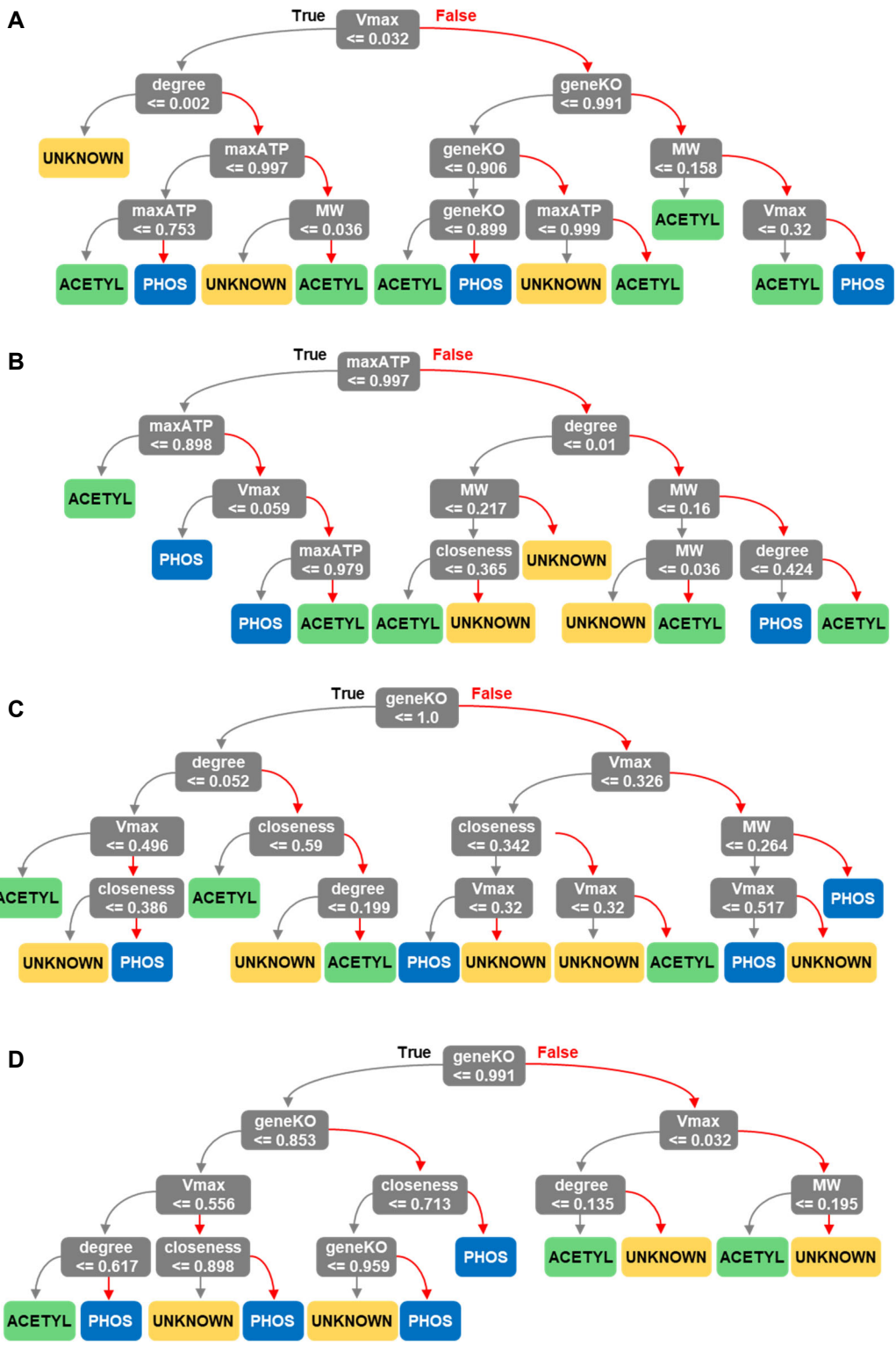
230

231

232

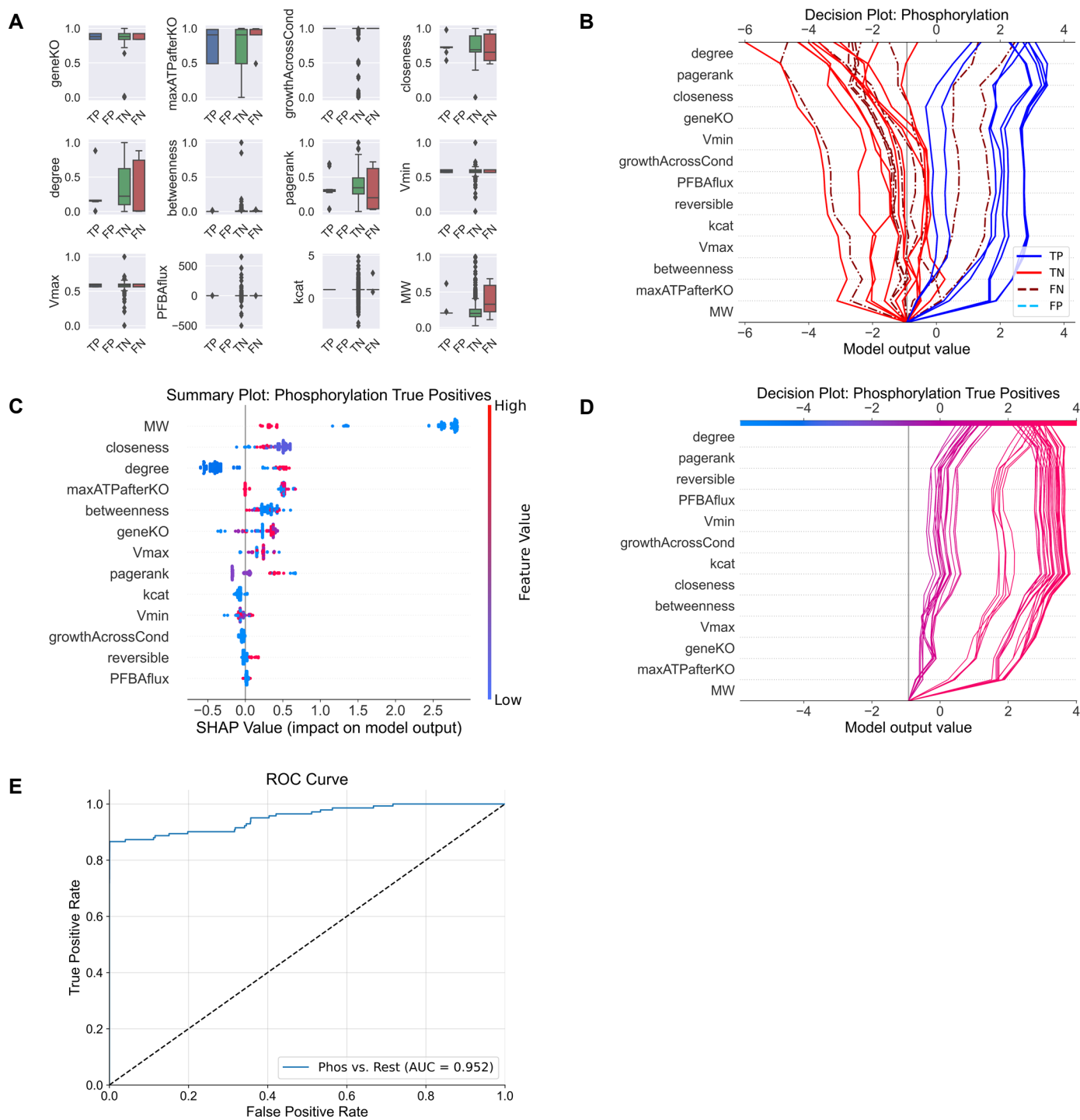
233

234
 235
 236
 237
 238
 239
 240
 241
 242
 243
 244
 245
 246
 247
 248
 249
 250
 251
 252
 253
 254
 255
 256
 257
 258
 259
 260
 261
 262
 263



264 Figure S11. Representative decision trees with maximum depth of 4, Related to Figures 3-5. Single decision tree
 265 models were trained for the multi-organism (A), *E. coli* (B), yeast (C), and mammalian (D) datasets. Only the top
 266 50% most important features, as identified in the Shapley analysis, were used to train the trees.

267
268
269
270
271
272
273
274
275
276
277
278
279
280
281
282
283
284
285
286
287
288
289
290
291
292
293



294 **Figure S12. Analysis of model predictions on the cell-cycle phosphorylation data, Related to Figures 3-5. A.**
295 Feature distributions for phosphorylated gene-reaction pairs are compared between true positive (TP), true
296 negative (TN) and false negative (FN) observations using boxplots. There were no false positives from this
297 validation test. **B** SHAP decision plot was created for 50 random observations to compare trends between the
298 classification groups. Values on the x-axis represent log odds of belonging to the phosphorylation class. **C** and **D.**
299 The phosphorylated gene-reaction pairs that were correctly classified (true positives) are displayed in a SHAP
300 summary plot (C) and decision plot (D). **E.** ROC curve for the model's phosphorylation predictions on the cell-
301 cycle data.

302
303
304
305
306
307
308
309
310
311
312
313
314
315
316
317
318
319
320
321
322
323
324
325
326
327
328
329
330
331
332
333
334
335

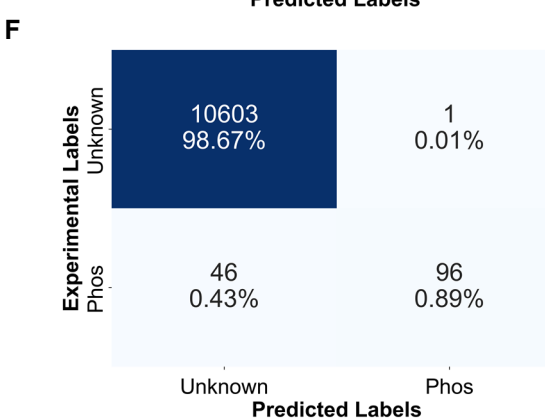
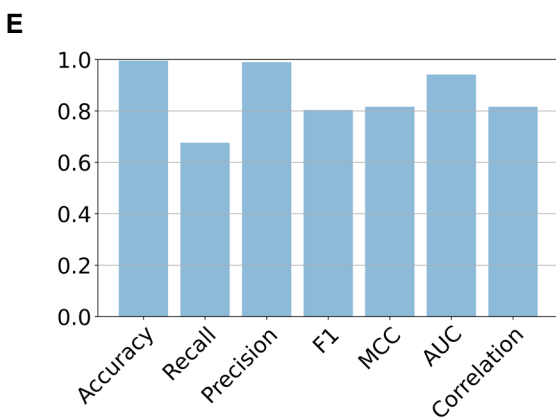
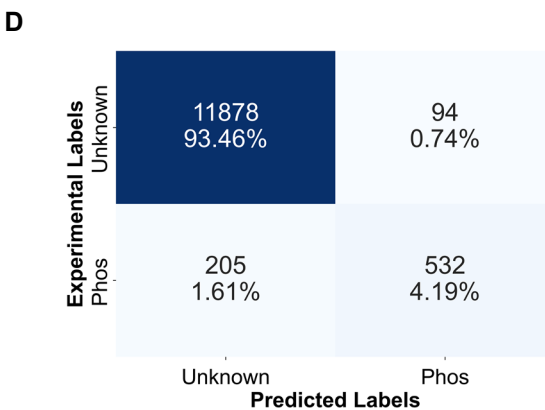
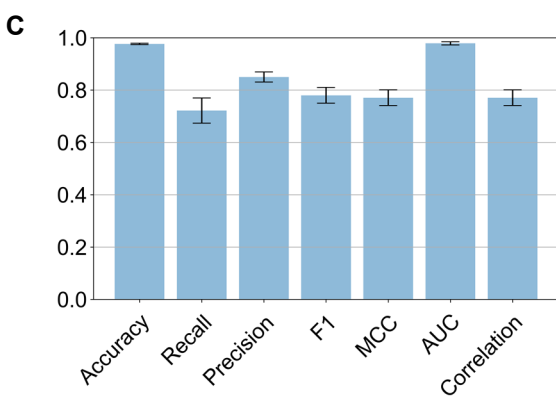
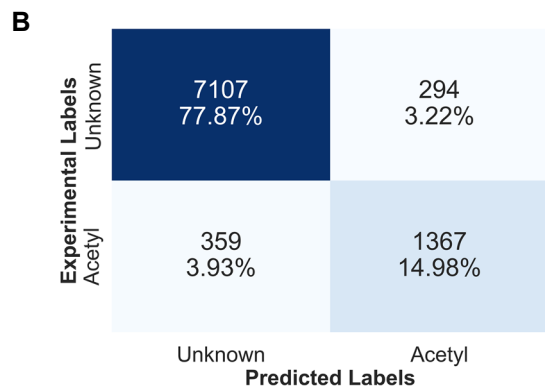
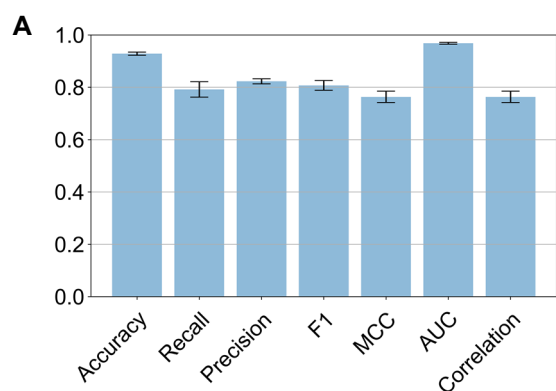
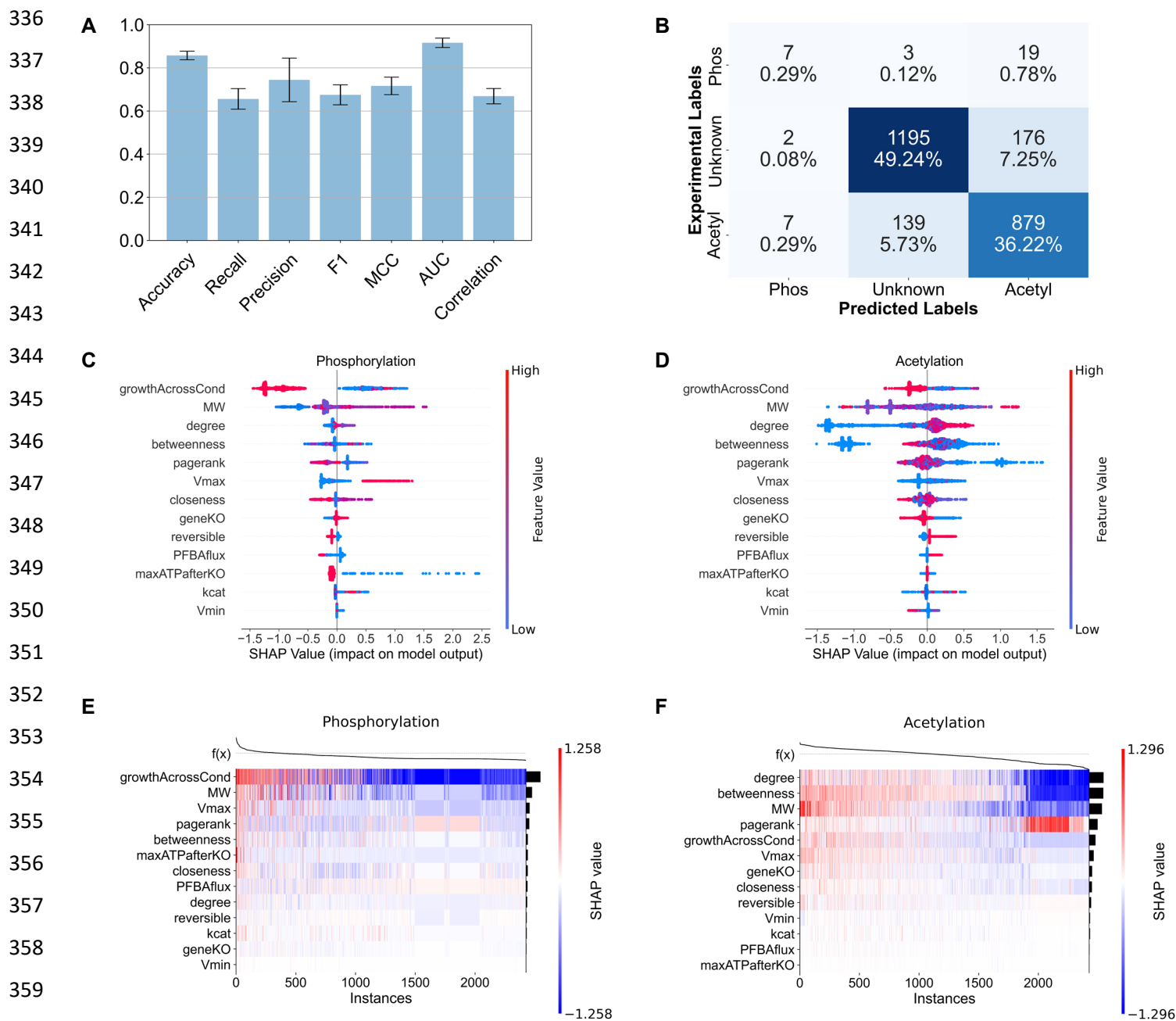


Figure S13. Binary classification models for predicting acetylation and phosphorylation separately, Related to Figures 3-5. The pipeline for training the models was identical to process used for the multi-class model. **A, B** The 5-fold cross-validation results for the acetylation model. **C, D**. The 5-fold cross-validation results for the phosphorylation model. **E, F**. The phosphorylation model was used to predict the cell-cycle validation dataset, which includes the G1, S and G2 phases. Overall, these results show that the ternary classification model outperforms the binary classification models.



361 **Figure S14. Organism-specific ML models – E. coli, Related to Figures 3-5.** XGBoost model trained on the *E. coli*
 362 dataset. **A, B.** 5-fold cross-validation results. Bar graph shows the mean scores across the 5 folds with a 95%
 363 confidence interval. **C, D.** SHAP value summary plots for the phosphorylation and acetylation classes. **E, F.** SHAP
 364 value heatmaps for the phosphorylation and acetylation classes. Observations are clustered by the model output,
 365 $f(x)$.

366
367
368

369
370
371
372
373
374
375
376
377
378
379
380
381
382
383
384
385
386
387
388
389
390
391
392
393
394
395
396
397
398
399
400
401

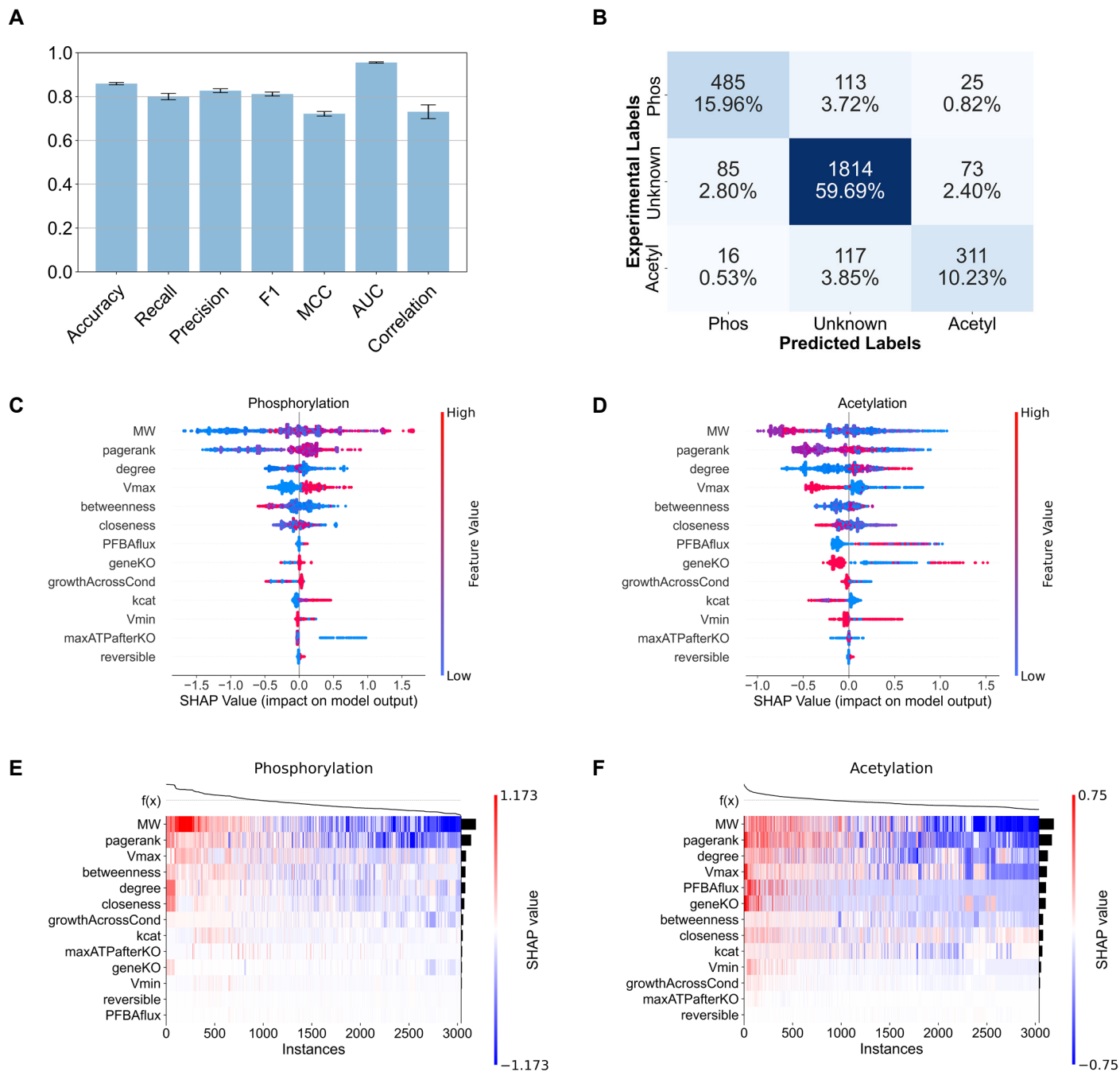
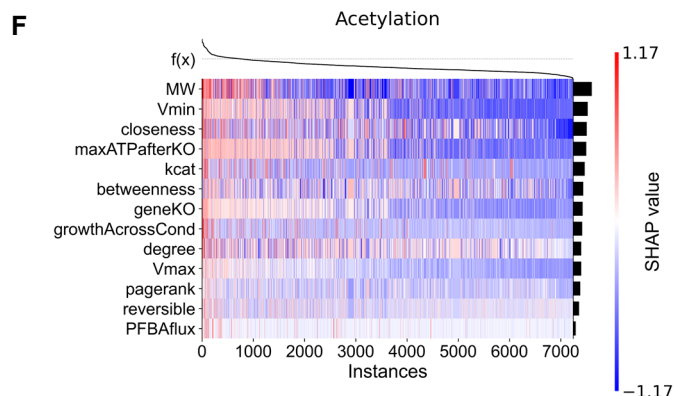
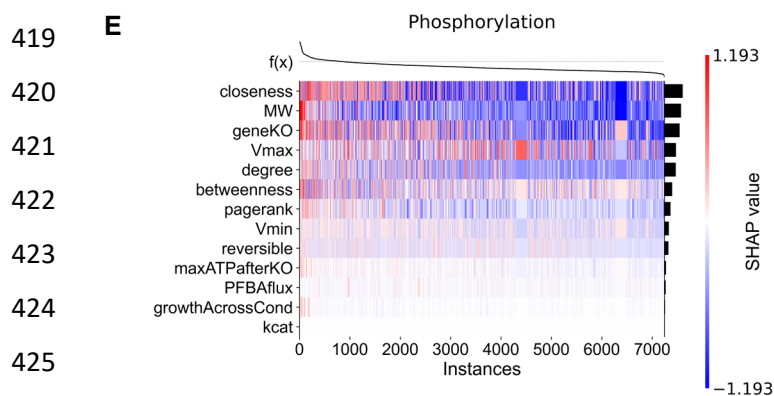
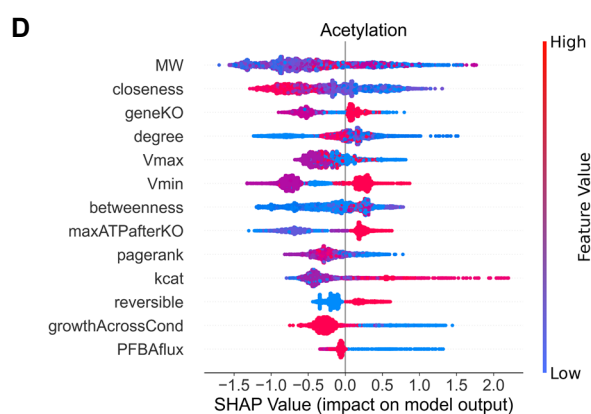
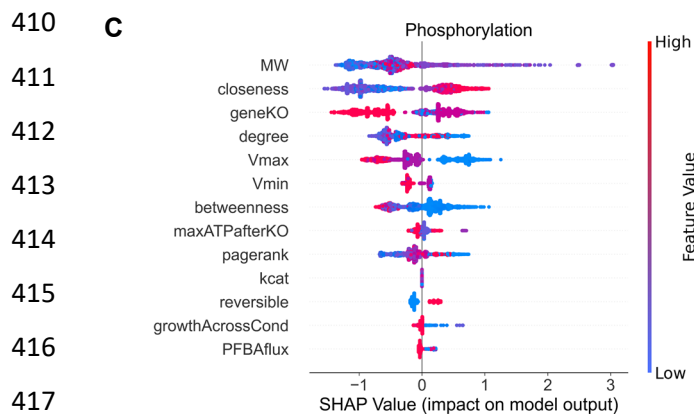
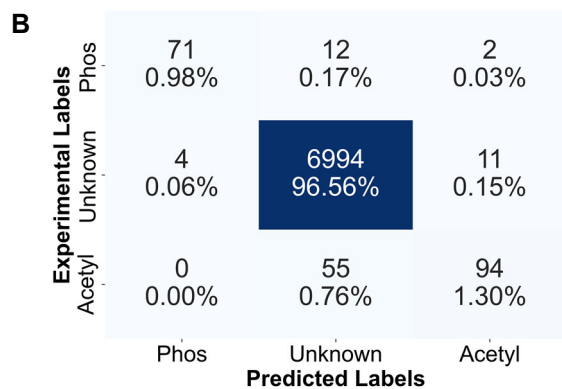
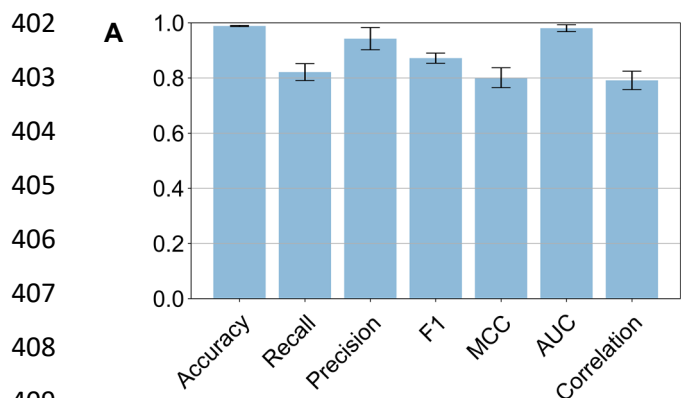
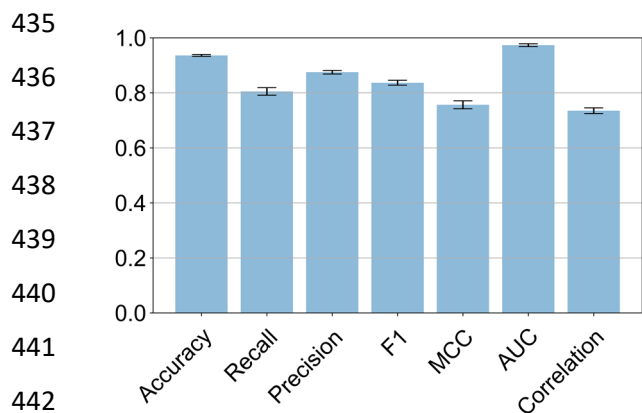


Figure S15. Organism-specific ML models – *S. cerevisiae*, Related to Figures 3-5. XGBoost model trained on the yeast dataset. **A**, **B**. 5-fold cross-validation results. Bar graph shows the mean scores across the 5 folds with a 95% confidence interval. **C**, **D**. SHAP value summary plots for the phosphorylation and acetylation classes. **E**, **F**. SHAP value heatmaps for the phosphorylation and acetylation classes.



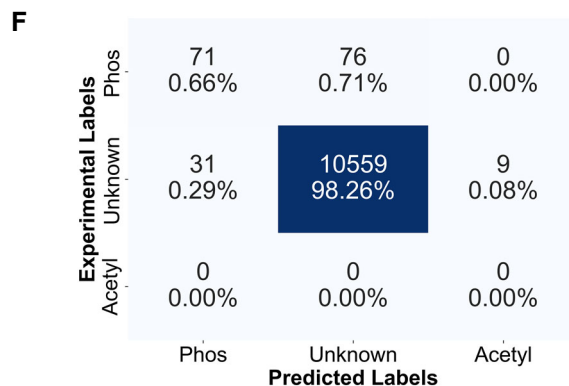
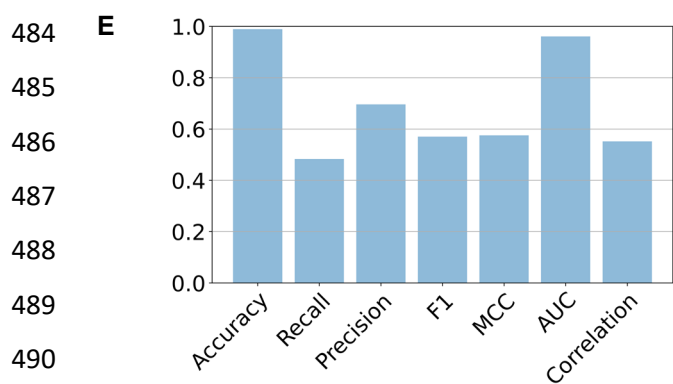
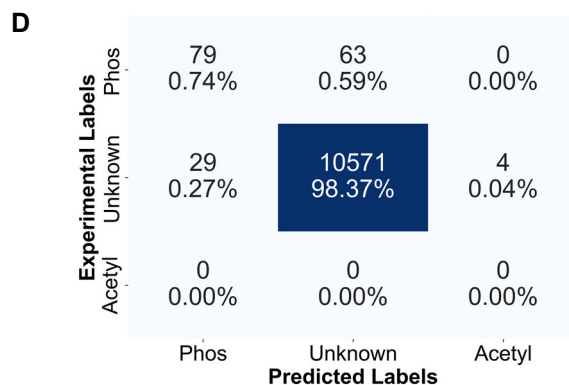
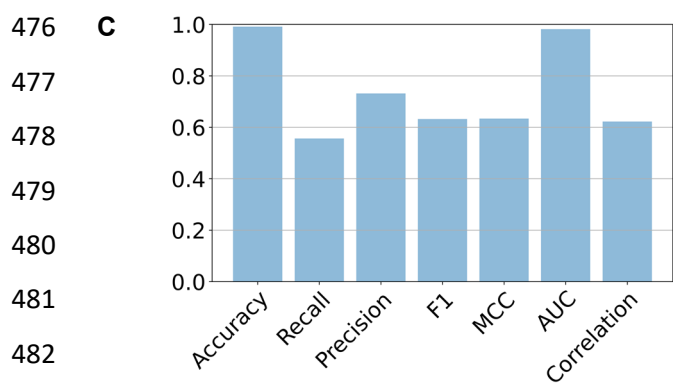
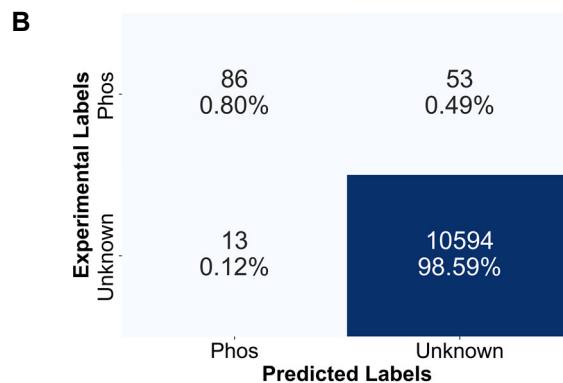
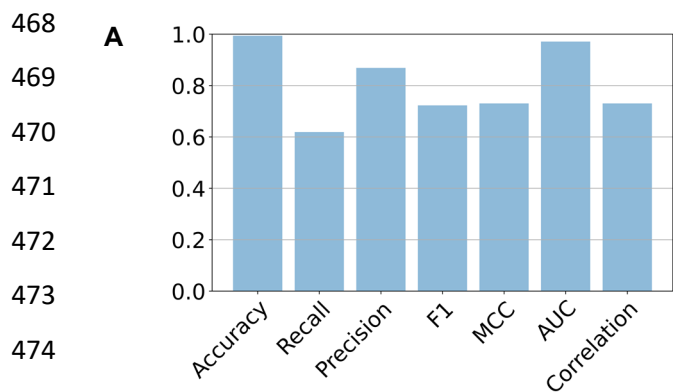
427 **Figure S16. Organism-specific ML models – mammalian cells, Related to Figures 3-5.** XGBoost model trained on
 428 the mammalian dataset. **A, B.** 5-fold cross-validation results. Bar graph shows the mean scores across the 5 folds
 429 with a 95% confidence interval. **C, D.** SHAP value summary plots for the phosphorylation and acetylation classes.
 430 **E, F.** SHAP value heatmaps for the phosphorylation and acetylation classes.



Experimental Labels	Predicted Labels		
	Phos	Unknown	Acetyl
Phos	530 3.54%	161 1.07%	46 0.31%
Unknown	66 0.44%	12330 82.29%	233 1.55%
Acetyl	18 0.12%	435 2.90%	1165 7.77%

Figure S17. Impact of including organism type in the ML model, Related to Figures 3-5. 5-fold cross-validation results for XGBoost model with organism-type included in the training data. Bar graph shows the mean scores across the 5 folds with a 95% confidence interval. The organism type was added as a categorical array where a 1 designated *E. coli*, 2 for yeast and 3 for human. The cross-validation results were extremely consistent with those from the primary model, suggesting that the model's decision-making is not influenced by organism type

443
444
445
446
447
448
449
450
451
452
453
454
455
456
457
458
459
460
461
462
463
464
465
466
467



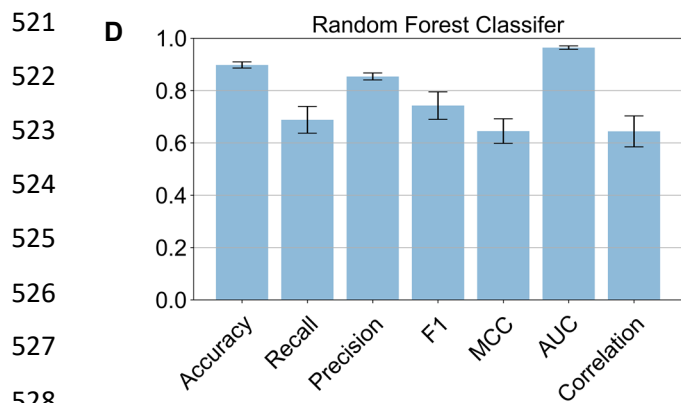
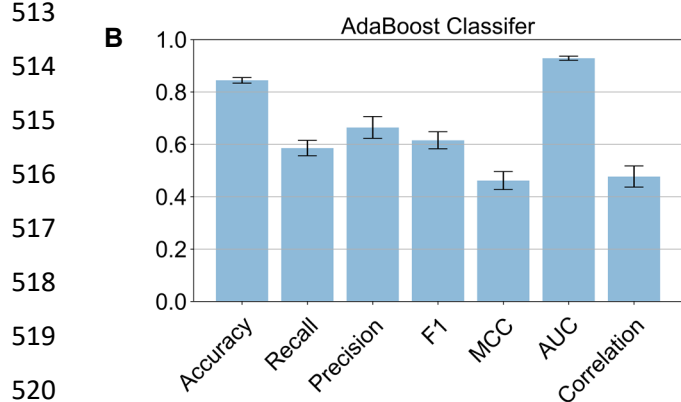
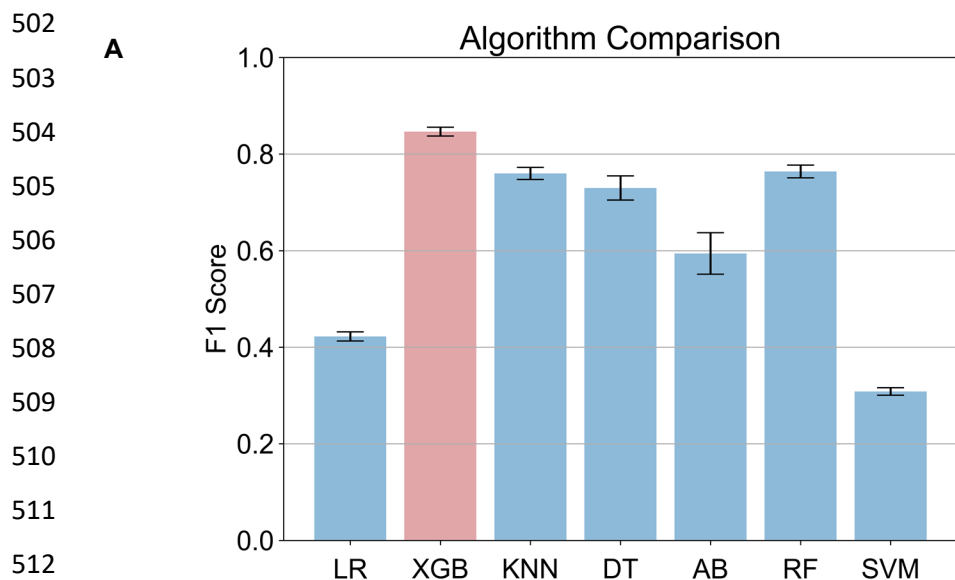
491 **Figure S18. Impact of training on different phases of the cell cycle, Related to Figures 3-5.** Models were trained by
492 replacing the G0 cell-cycle data from the training set with the feature matrix from the remaining phases: G1, S,
493 and G2. Each model was then used to predict the phosphorylated genes from the phases not featured in the
494 training. These results are shown here for the G1-model (**A, B**), S-model (**C, D**) and G2-model (**E, F**). All three
495 models, especially for S and G2, performed inferior to the primary CAROM-ML model in regard to this validation
496 test. These results suggest that S and G2 conditions have a distinct phosphorylation pattern from the remaining
497 conditions.

498

499

500

501



C

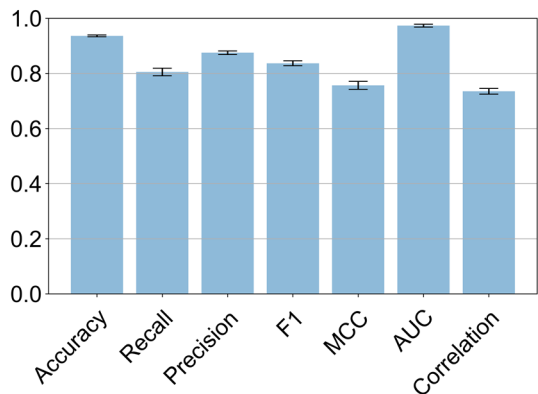
Experimental Labels	Predicted Labels		
	Phos	Unknown	Acetyl
Phos	220 1.73%	483 3.80%	34 0.27%
Unknown	228 1.79%	9660 76.01%	466 3.67%
Acetyl	19 0.15%	748 5.89%	851 6.70%

E

Experimental Labels	Predicted Labels		
	Phos	Unknown	Acetyl
Phos	341 2.68%	354 2.79%	42 0.33%
Unknown	54 0.42%	10047 79.05%	253 1.99%
Acetyl	3 0.02%	593 4.67%	1022 8.04%

529 **Figure S19. CAROM-ML model performance using various ML algorithms, Related to Figures 3-5.** 5-fold cross-
530 validation results were compared for various untuned algorithms, with F1 score used as the metric **(A)**. XGBoost,
531 colored in red, had the best performance and was therefore used for the main CAROM-ML model. AdaBoost **(B,**
532 **C)** and random forest **(D, E)** models were further tested by tuning their hyperparameters and performing 5-fold
533 cross-validation. For all bar graphs, the mean scores across the 5 folds are shown with a 95% confidence interval.

535
536
537
538
539
540
541
542
543
544
545
546
547
548
549
550
551
552
553
554
555
556
557
558
559
560
561
562
563
564
565
566
567
568
569



Experimental Labels	Predicted Labels		
	Phos	Unknown	Acetyl
Phos	530 3.54%	161 1.07%	46 0.31%
Unknown	66 0.44%	12330 82.29%	233 1.55%
Acetyl	18 0.12%	435 2.90%	1165 7.77%

Figure S20. Impact of retaining genes that do not have evidence for phosphorylation or acetylation Related to Figures 3-5. 5-fold cross-validation results for model trained on full set of genes is shown. Bar graph shows the mean scores across the 5 folds with a 95% confidence interval. For the main CAROM-ML model, online databases were used to compile a list of enzymes that have been found to be phosphorylated or acetylated in published studies. Non-annotated enzymes were removed from the training data. Here we show the results for the model which had these non-annotated enzymes included in the training data did not differ from the model with these genes removed during the model construction.

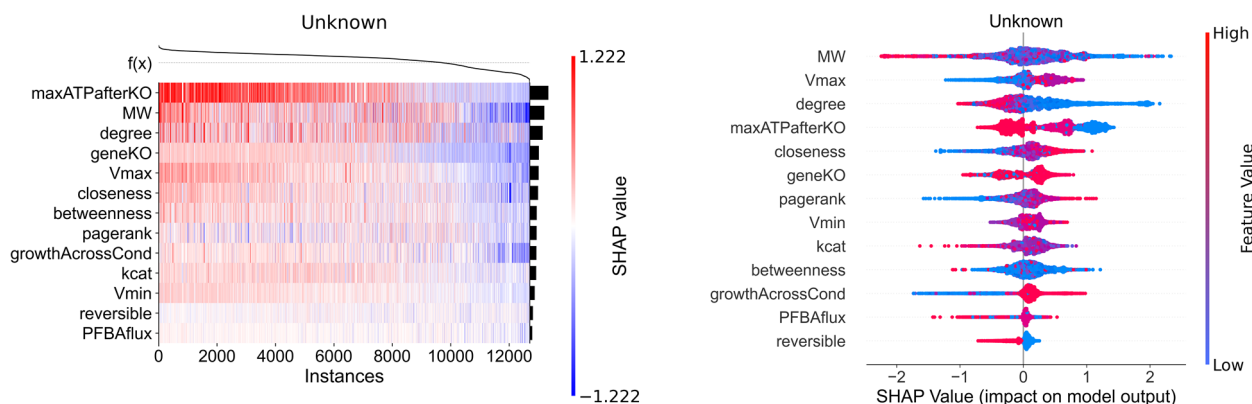


Figure S21. Interpretation of the CAROM-ML model using Shapley analysis: Unknown class, Related to Figures 3-5. Corresponding plots for the phosphorylation and acetylation classes are shown in Figure 4 in the main text.

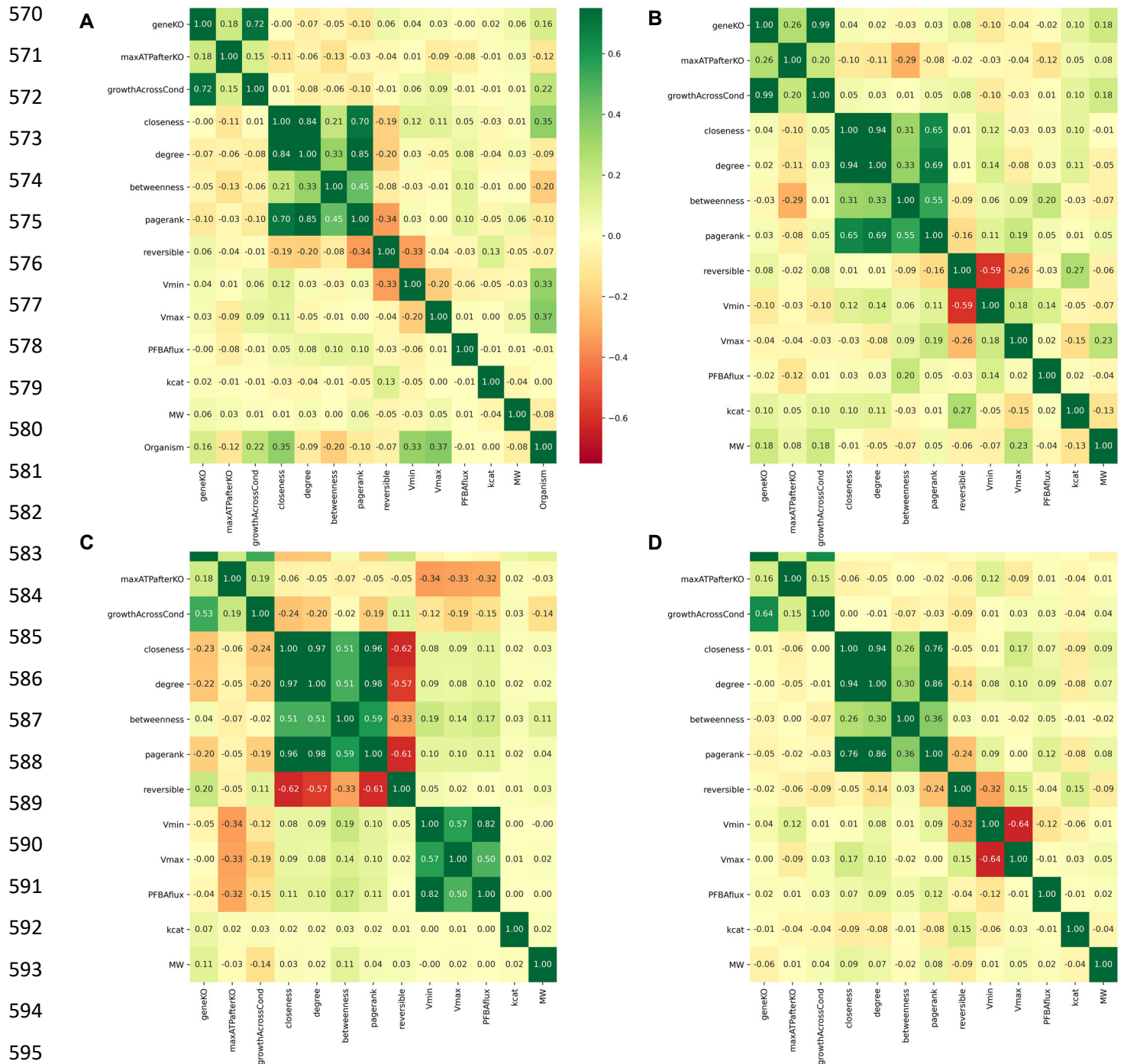


Figure S22. Correlation map of all model features, Related to Figures 3-5. Heatmap of Pearson's correlation between feature values for the following datasets: all organism types (A), yeast (B), *E. coli* (C), and human (D).

602

603

604

605

606

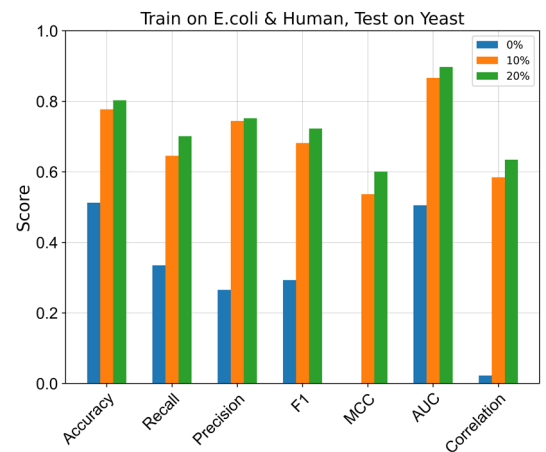
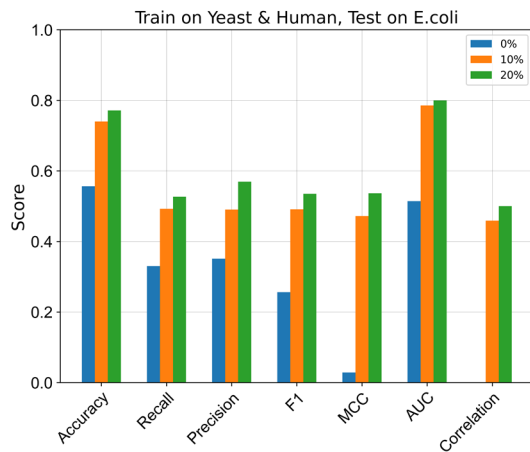
607

608

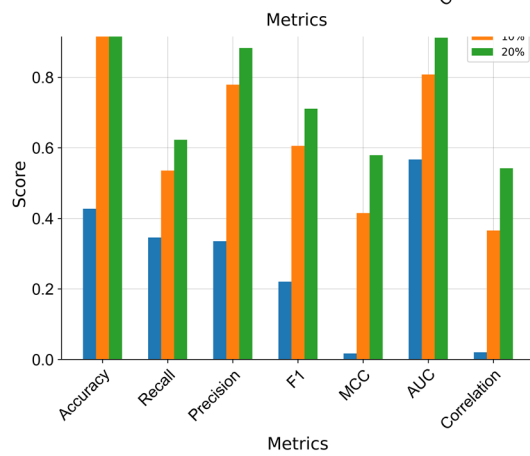
609

610

611



612



619

620 **Figure S23. Predicting on unseen organisms, Related to Figures 3-5.** XGBoost models were trained on the data
 621 from two organisms and used to make predictions on the third (e.g. train on E. coli and yeast, test on
 622 mammalian). Data from the test organism was moved to the training data in increments of 0%, 10% and 20%.
 623 Model performance improved significantly after including a small number of samples from the test organism in the
 624 training dataset.

625

626

627

628

629

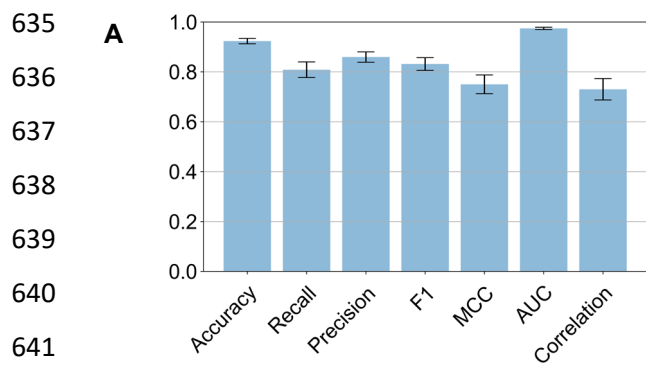
630

631

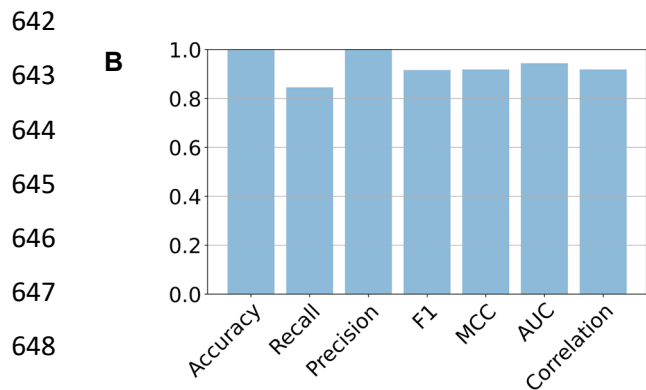
632

633

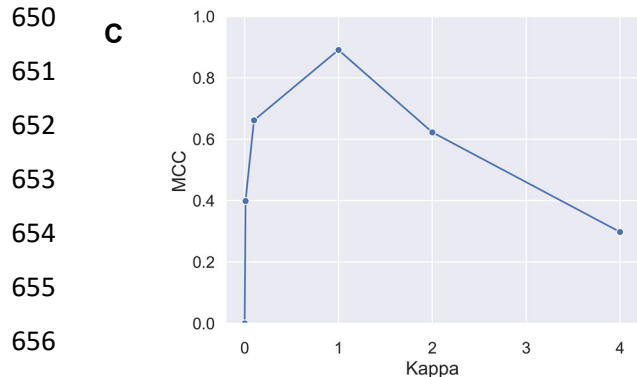
634



Experimental Labels	Predicted Labels		
	Phos	Unknown	Acetyl
Phos	522 4.11%	159 1.25%	56 0.44%
Unknown	81 0.64%	9995 78.65%	278 2.19%
Acetyl	15 0.12%	386 3.04%	1217 9.58%



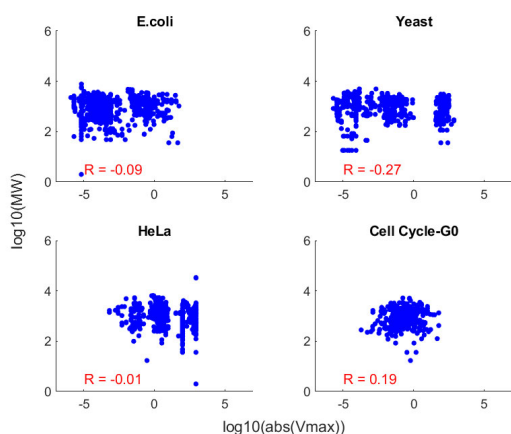
Experimental Labels	Predicted Labels	
	Phos	Unknown
Phos	120 1.12%	22 0.20%
Unknown	0 0.00%	10604 98.68%



657 **Figure S24. Impact of adjusting flux-related feature parameters on the ML results, Related to Figures 3-5.** For the
 658 ML analysis, the Vmax and Vmin features were constrained to magnitudes below 100 in order to reduce the effect of
 659 unconstrained reactions and the variability across organism types. Here we show that the CAROM-ML model is
 660 robust to increasing the threshold to the 900 mmol/gDW/hr value used for the ANOVA testing. A supplementary
 661 model was trained on the *E. coli*, yeast, HeLa and G0 phase data after adjusting this threshold. **A.** Results from
 662 training the model using 5-fold cross-validation. Bar graph shows the mean scores across the 5 folds with a 95%
 663 confidence interval. **B.** The model was used to predict on the cell cycle validation dataset, which includes the G1,
 664 S and G2 phases. **C.** The cell cycle metabolic models were generated using dynamic flux analysis (DFA) with a
 665 default value of 1 for kappa, the optimization weight that is applied to the metabolomics data relative to the
 666 biomass objective. Changes to kappa therefore affect the flux- and growth-related features. The ML model's
 667 performance on the cell cycle dataset was fairly robust as kappa was incrementally changed from 1e-3 to 4,
 668 however the default value of 1 provided the best results. Setting it 0 or very low values results in the model not
 669 learning any differences between the cell cycle phases as expected. At very high values, the DFA model overfits
 670 to the metabolomics and is affected by noise in the measurement. The default value of 1 provides a good trade off
 671 in separating signal from noise.

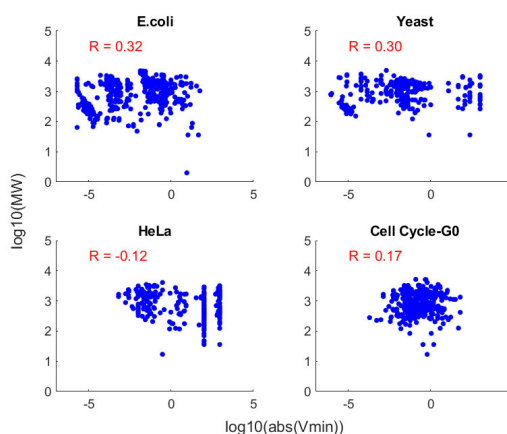
672

A



673

B



C

Vmax/min Adjusted for MW

Acetyl	102 0.95%	41 0.38%	0 0.00%
Phos	40 0.37%	10478 97.51%	0 0.00%
Unknown	0 0.00%	85 0.79%	0 0.00%
	Acetyl	Phos	Unknown

Accuracy=0.985
Precision=0.569
Recall=0.570
F1 Score=0.569
MCC=0.562

674

675 **Figure S25. Relationship between flux and MW per reaction, Related to Figures 3-5.** Here we address whether
 676 the flux features generated from flux variability analysis, V_{max} and V_{min} , are strongly correlated with the
 677 molecular weight (MW) of the metabolites present in the corresponding reactions. We did not find a significant
 678 correlation between MW of the metabolites and the predicted fluxes. For this analysis, “MW” represents the sum
 679 of MW for all metabolites in a given reaction. The plots show the relationships between **(A)** V_{max} vs. MW and **(B)**
 680 V_{min} vs. MW for each organism on log scales. The Spearman’s correlation shown on each plot suggest there is
 681 not a consistent relationship between flux and the MW present in a given reaction. While a negative correlation is
 682 expected, in some cases we see a positive correlation. This suggests that there is not a strong relationship
 683 between the two. **C.** A separate XGBoost model was trained after adjusting the V_{max}/min features for the MW per
 684 reaction by multiplying the fluxes with the MW. The model’s performance slightly worsened compared to the main
 685 CAROM-ML model.

686

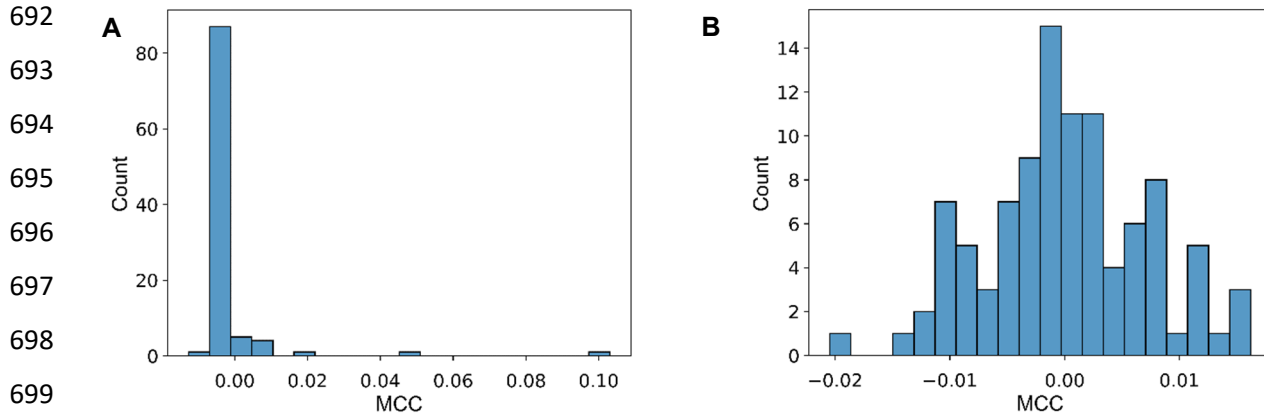
687

688

689

690

691



700 [Figure S26. Random permutation models as benchmark for the CAROM-ML model, Related to Figures 3-5.](#) We
 701 generated 100 random permutations of the class labels for the CAROM-ML training dataset. The models
 702 generated with these permutations achieved scores close to $MCC=0$, as expected for a random model. **A.** For
 703 each permutation, an XGBoost model was trained using the CAROM-ML feature dataset and the shuffled class
 704 labels, then used to predict on the cell cycle G1/S/G2 dataset. **B.** For each permutation, a subset of the shuffled
 705 class labels from the training dataset was used to guess the G1/S/G2 class labels.

706

707

708

709