

Another 100 genes

Tom Röschinger¹, Grace Solini¹, Kian Faizi¹, Rosalind Pan¹, Anika Nawar Choudhury², Lynn Yang³, Bob Jones³, Stephen Quake^{2, 3, 4}, and Rob Phillips^{1, 5, +}

¹*Division of Biology and Biological Engineering, California Institute of Technology, Pasadena, CA 91125, USA*

²*Chan Zuckerberg Biohub, San Francisco, CA 94158, USA*

³*Department of Bioengineering, Stanford University, Stanford, CA 94305, USA*

⁴*Department of Applied Physics, Stanford University, Stanford, CA 94305, USA*

⁵*Department of Physics, California Institute of Technology, Pasadena, CA 91125, USA*

⁺*Correspondence: phillips@pboc.caltech.edu*

1 Introduction

It has been more than sixty years since Jacob and Monod [1] shaped the way we think about transcriptional regulation in prokaryotes, yet, although about 10^{17} bases have been deposited in the SRA database[2], we have yet to obtain a full understanding of how all the genes of a single organism are regulated. Even in the case of one of biology’s best studied model organism, *Escherichia coli*, about two thirds of the genes lack any regulatory annotation (see S1.1). For other prokaryotic model organisms the numbers are similar (see S1.2, S1.3), while higher order model organisms such as *Saccharomyces cerevisiae* (see S1.4) and *C. elegans* (see S1.6) have close to no regulatory annotations, given the arguably more complex nature of gene regulation in eukaryotes. Understanding how genes are regulated is required to understand how an organism adapts its physiology on short time scales to environmental stresses, as well as evolutionary adaption on long time scales. In addition, gene regulation networks and their building blocks, such as transcription factor binding sites and RNA polymerase (RNAP) promoters, are key elements in the design of synthetic gene circuits [3–5].

With its ever increasing availability, Next Gen Sequencing (NGS) is primed to be the method of choice to discover transcription factor and RNAP binding sites. A vast array of methods exists that make it possible to identify binding sites of either specific proteins or for a broad spectrum of DNA binding factors. In methods like ChIP-Seq [6], proteins have to be cross linked to DNA, which often requires changing residues in the amino-acid sequence, such as for LacI in *E. coli* [7]. In addition, antibodies against the protein of interest have to be available, or the protein has to be modified to include a tag which can be targeted by antibodies. While the resolution of these methods is ever improving, it does not allow for a nucleotide resolution yet, making it difficult to identify changes in binding affinity caused by single mutations. Other methods such as ATAC-Seq [8, 9] and DNase-Seq [10] rely on open chromatin for binding site identification, and have almost exclusively been used for eukaryotic organisms. In general, identifying regulatory interactions from transcription factor occupancy alone can be misleading, since there can be high affinity binding sites in the genome, where there is no change in expression levels upon binding [11]. DAP-Seq [12] is a method similar to ChIP-Seq, however, instead of using immunoprecipitation to obtain DNA-TF pairs, purified and tagged TFs are incubated with fragmented genomic DNA. The method has been used to identify genome wide binding sites for TFs in *Clostridium thermocellum* [13] and *Riemerella anatipestifer* [14].

Another approach is to use RNA-Seq as readout for mutagenised promoter regions, where binding sites are identified as regions that, when mutated, lead to significant increase or decrease in expression of a repressor gene [15–17]. Here we studied the regulatory architecture of 104 genes, including energy matrices with nucleotide resolution that make it possible to build thermodynamic models to predict gene expression [17–20]. Additionally, we present major improvements to the method called Reg-Seq

[17], making further steps towards obtaining a method allowing to discover regulatory architectures genome wide. Reporter genes are chromosomally integrated into the *E. coli* genome, and reduced diversity in mRNA stability lead to more precise identification of binding sites. A vast array of growth conditions is used to show how certain binding sites can only be identified in a certain growth condition, such as (TR: name example). The identification of transcription factors was moved on from laborious mass spectrometry experiments, using *in vitro* binding assays as well as a library of transcription factor knockout strains. Finally, improved computational analysis increases the speed of data analysis and the accuracy of parameters that are used for thermodynamic models (TR: here I am thinking Rosalinds stuff).

2 Methods

2.1 Promoter sequence import

2.2 Reporter construct design

2.3 Barcode Mapping

2.4 Genome Integration

2.5 Growth Conditions

M9 minimal media was prepared without carbon source by adding to 700ml of water 200ml of 5x base salt solution (Difco M9 Minimal Salts, 5x: 239mM Disodium phosphate, 110mM Monopotassium phosphate, 42.8mM Sodium chloride, 93.5mM Ammonium chloride, autoclaved), 100µl of 1M Calcium chloride, 1ml of 1M Magnesium sulphate, 10ml of 100x Trace Elements (3mM Iron(III) chloride, 0.35mM Zinc chloride, 76.5µM Copper chloride dihydrate, 42µM Cobalt chloride hexahydrate, 160µM Boric Acid, 8.1µM Manganese chloride)[21] and 1ml of 1mg/ml thiamine. The total volume is filled up to 1L and the solution is filter sterilized (Fisherbrand 500ml Bottle Top Filter, 0.2µm aPES membrane). Carbon sources were prepared as aqueous solutions of 20% w/v.

2.6 Cultivation

The cultivation procedure was adapted from [22] with slight modifications. Frozen glycerol stocks of the library were grown overnight in 200ml of M9 minimal media with 0.5% Glucose in a 1L Erlenmeyer flask at 37C and shaken at 250rpm.

3 Results

3.1 Methodology

The methodology presented in this work is based on previous iterations of sequencing based approaches [17] and approaches based on fluorescence activated cell sorting, termed Sort-Seq [18, 19], with some significant modifications since the work by Ireland, et al. [17]. A detailed description of the methodology can be found in Methods and Supplementary Information. In short, for each promoter studied, we generated a library of mutant promoters with an average mutation rate of 0.1. Each mutant sequence is cloned upstream of a sfYFP reporter gene and tagged with a unique DNA barcode 3' end of the transcript, reducing potential biases in transcription rates induced by the random barcode (TR: would be great to have a citation for this). The reportes are integrated into the *E. coli* genome

using a modified ORBIT protocol [23], which allows for high-throughput, parallel genome integration of the entire library, see Methods and Supplementary Information for details. Using DNA and RNA sequencing, we measure the expression of the reporter gene for each unique set of mutations in a various set of growth conditions. We then use mutual information to identify bases whose identity is important for expression. The hypothesis is that mutations of bases within binding sites for transcription factors or RNAP have a much stronger effect on expression than mutations of bases that not bound by regulatory factors. (TR: add TF identification methods.) A summary of the method is shown in Figure 1.

3.2 Genes studied

104 genes were chosen for this study. 16 of these genes were chosen as so called “gold standards”. These genes have well annotated promoters and have been studied in detail in previous experiments [17, 19]. Including this set of genes allows us to compare the method presented in this work to previous iterations and verify the results, as well as find possible derivations or improvements. 18 genes were chosen that have been identified to have a high variation in protein copy number across a set of 22 growth conditions by Schmidt et al., 2016 [22]. These genes were chosen since a high variation in copy number suggests that there are regulatory proteins controlling the expression of the gene. Of these 18 genes, 9 had no function annotated at the time of this study. Another set of 13 genes was chosen from EcoCyc as part of the so-called y-ome[24], which is made up of genes not having any functional annotation. 18 genes were chosen for being part of toxin/anti-toxin systems. Two sets of genes were chosen from the work of Lamoureux et al.[25], where groups of genes which are controlled by the same transcription factor are identified, so called iModulons. We chose two newly identified groups, responding to the putative transcription factors YmfT and YgeV respectively. Finally, 6 genes were chosen for being part of gene regulatory networks with feed-forward-loop motifs. An entire list of genes can be found in (TR: some SI table).

3.3 Barcode Mapping

For each gene studied here, we designed a library of 1500 mutated promoter variants with an average mutation rate of 0.1 for each promoter that was identified with the gene in EcoCyc. If a gene did not have a promoter identified, we first looked for a possible TSS identified by Urtecho et al., 2018 [15], which then was taken as initial sequence for generating mutated variants. In some cases, no TSS could be identified and we used the model of LaFleur et al., 2022 [26], which predicts transcription start sites for σ^{70} promoters given a genomic sequence of *E. coli*, to find the site in the intergenic region leading up to the first coding region in the operon the gene was part of that had the highest predicted affinity for σ^{70} -factors. A total of 178619 sequences were ordered from Twist Biosciences (1500 mutated variants and wildtype sequence for 119 promoters). Random barcodes were added to the sequences by PCR and they were cloned into a plasmid vector for amplification and subsequent cloning steps. Details can be found in Supplementary information sections S3.1, S3.2, S4.1 and S5. 171591 (>96%) of sequences were identified in the plasmid library, and an additional 101 sequences were found, likely due to errors in the synthesis process of the oligonucleotide library. For seven genes, there were two TSSs annotated to the same position, likely by TSSs for different sigma factors. For five promoters (tnaCp, crpp1, lpp, yqaEp, dicCp) significantly less variants were recovered. Four of these (tnaCp, crpp1, lpp, yqaEp) have high AT content and thymine repeats in the beginning of the promoter sequence, possibly leading to increased error rates in sequencing (TR: look

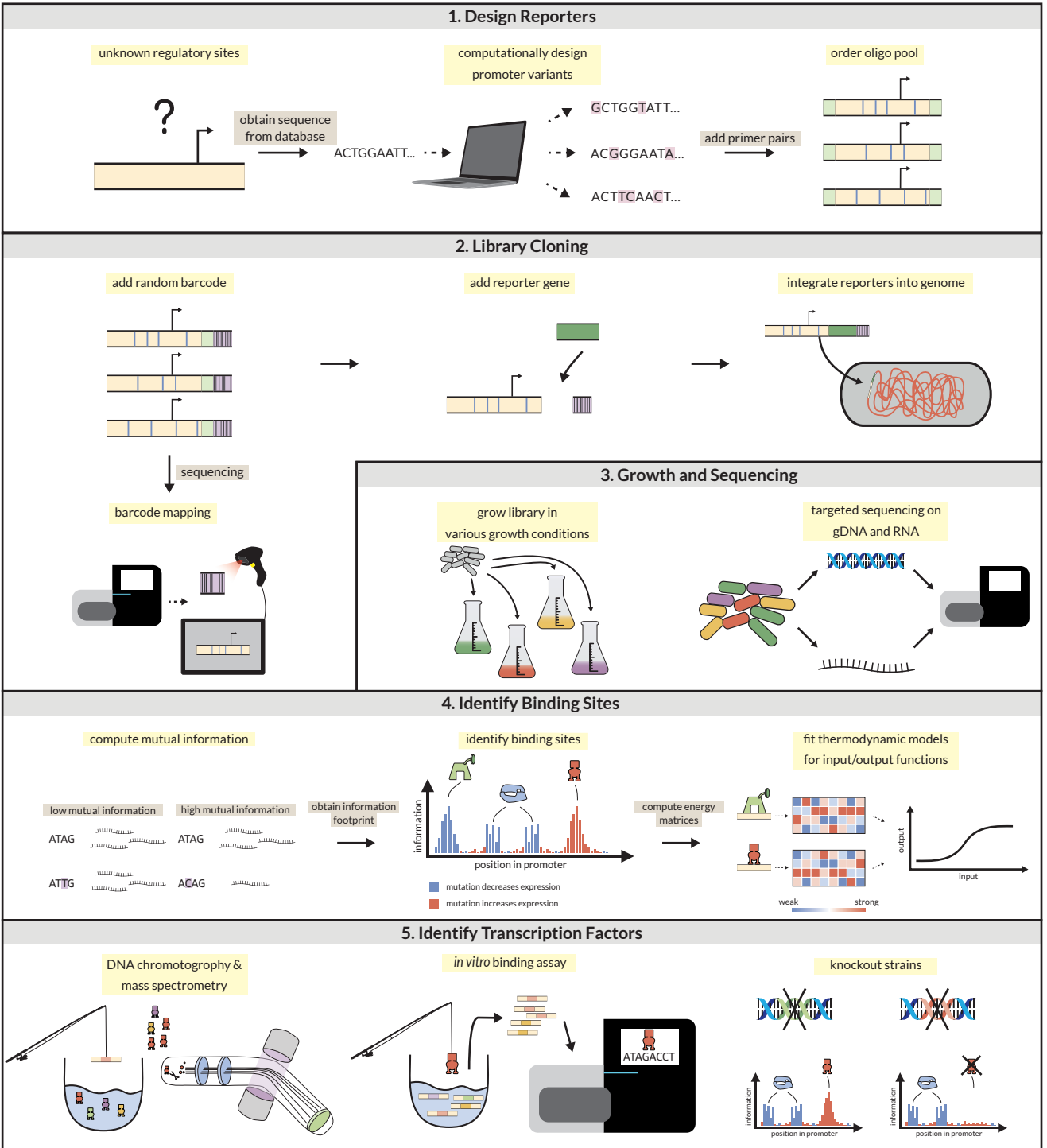


Figure 1. Method summary

this up). For dicCp, the recovered promoter variants have an increased mutation rate between bases -7 and 12 relative to the TSS, as shown in Fig. 2(B). This region is exactly the binding site identified for the transcription factor DicA in this promoter [27]. An increased mutation rate indicates that the binding affinity for the transcription factor DicA has to be reduced significantly, indicating that binding of DicA to its functional targets in the genome is essential for growth. CRISPR interference TF knockdown experiments have shown similar lethal effects when targeting DicA[28].

3.4 Expression measurements in 37 growth conditions

3.5 Transcription Factor identification

3.6 Growth Conditions

3.7 Gold Standard genes

3.8 Ethanol iModulon

YgeV has been predicted to be a regulator involved in purine catabolism, leading to the production of allantoin, which can be used as a sole nitrogen source [29]. There are 16 putative regulatory targets [25] for YgeV, including the *xdhABC* operon, which degrades xanthine to uric acid [29] in the purine catabolism pathway. *E. coli* can survive exposure to low ethanol concentrations up to 5%, which can even lead to increased DNA synthesis [30], but mostly leads to various stress responses such as an increased production of ROS. Growth in media supplemented with ethanol induced a change in gene expression for genes regulated by YgeV in a $\Delta baeR$ or $\Delta cpxR$ mutant strain. (TR: Is there a correlation between ethanol response and higher need for nitrogen?)

3.9 Oxidative stress response iModulon

The putative transcription factor YmfT regulates 14 out of 23 genes in the e14 prophage and is predicted to respond to oxidative stress [25]. Oxidative stress is caused by reactive oxygen species (ROS) such as H_2O_2 , which are highly reactive and damage DNA, the cell wall, proteins [31] etc., however, oxidized amino acids can also lead to conformational changes in transcription factors, such as OxyR and HypT, which induce DNA binding and subsequent regulation of genes involved in response to oxidative stress [31]. Hydrogen peroxide is produced endogenously in various pathways in *E. coli* and especially in high amounts when phenylethylamine is used as either carbon or nitrogen source [32], hence we used minimal media supplemented with 10 mM 2-phenylethylamine hydrochloride (PEA) as sole carbon source to induce stress responses to H_2O_2 and therefore oxidative stress. (TR: discuss findings of ymfT modulon and how it relates to oxyR, look at oxyR iModulons and possibly do another run including this gene.)

3.10 Antitoxin/Antibiotic genes

3.11 other y-ome genes

4 Discussion

- discuss how to scale to 1000 genes

5 Acknowledgements

- Bill Ireland (Discussion)

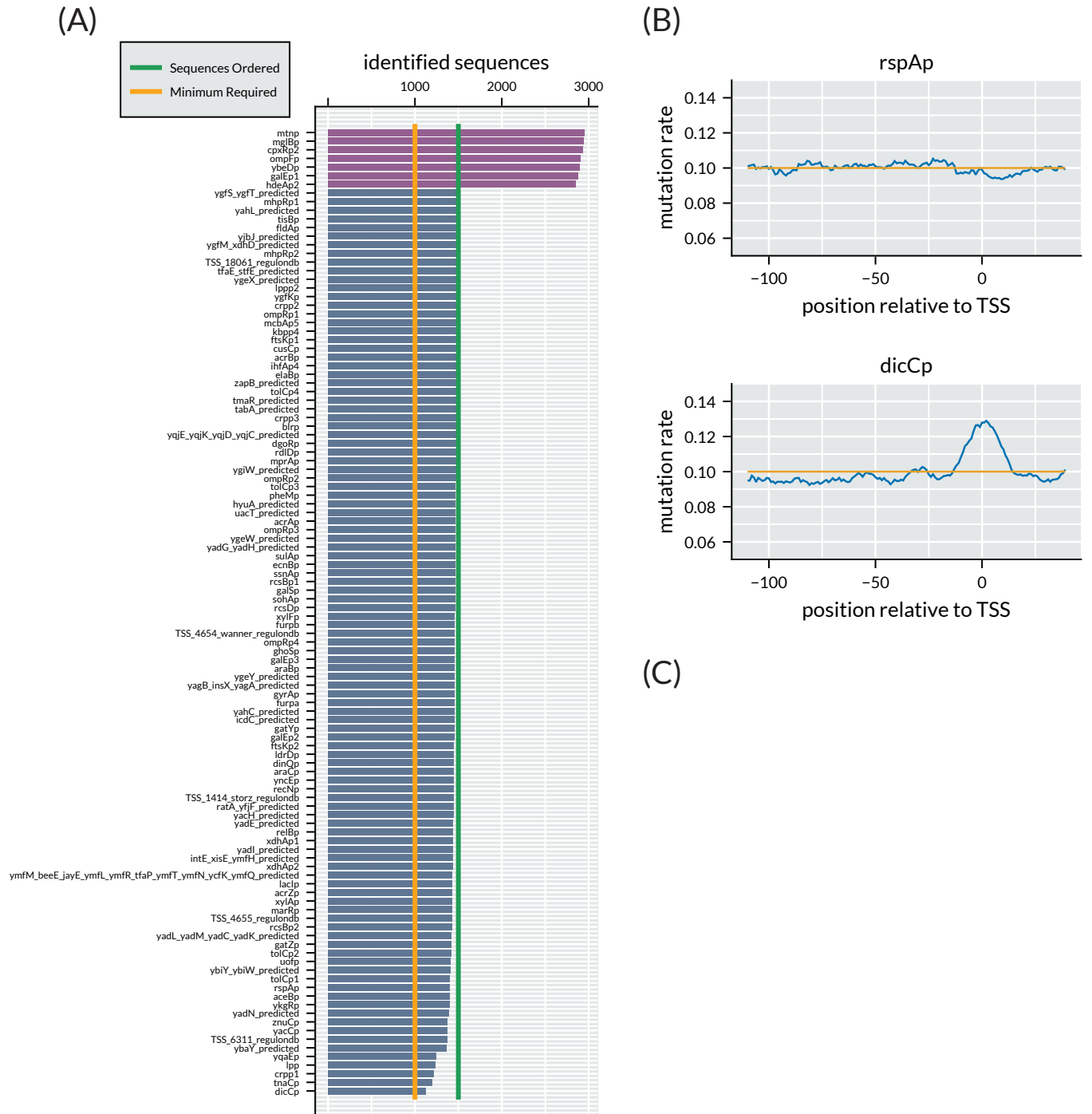


Figure 2. (A) Recovered promoter variants in barcode mapping. Average mutation rate per base across all promoter variants for galEp1 (B) and dicCp(C).

- Justin Kinney (Discussion)
- Stephan Grill (Discussion)
- Frank Jülicher (Discussion)
- Igor A. Antoshechkin (Sequencing). This work was supported by the Millard and Muriel Jacobs Genetics and Genomics Laboratory at California Institute of Technology.
- Victor Garcia-Ruiz (Cell Lysis)

6 To do list

- Complete Introduction
- Write SI about mutual information and determining binding sites
- expand paragraphs for genes chosen
- check experimental details in methods

References

- ¹F. Jacob and J. Monod, “Genetic regulatory mechanisms in the synthesis of proteins”, *Journal of molecular biology* **3**, 318–356 (1961).
- ²NIH, *Sra database growth*, <https://www.ncbi.nlm.nih.gov/sra/docs/sragrowth/>, Accessed: 01/11/2024, 2024.
- ³M. B. Elowitz and S. Leibler, “A synthetic oscillatory network of transcriptional regulators”, *Nature* **403**, 335–338 (2000).
- ⁴S. Mangan and U. Alon, “Structure and function of the feed-forward loop network motif”, *Proceedings of the National Academy of Sciences* **100**, 11980–11985 (2003).
- ⁵U. Alon, *An introduction to systems biology: design principles of biological circuits* (Chapman and Hall/CRC, 2006).
- ⁶H. S. Rhee and B. F. Pugh, “Chip-exo method for identifying genomic location of dna-binding proteins with near-single-nucleotide accuracy”, *Current protocols in molecular biology* **100**, 21–24 (2012).
- ⁷D. Rutkauskas, H. Zhan, K. S. Matthews, F. S. Pavone, and F. Vanzi, “Tetramer opening in laci-mediated dna looping”, *Proceedings of the National Academy of Sciences* **106**, 16627–16632 (2009).
- ⁸J. D. Buenrostro, B. Wu, H. Y. Chang, and W. J. Greenleaf, “Atac-seq: a method for assaying chromatin accessibility genome-wide”, *Current protocols in molecular biology* **109**, 21–29 (2015).
- ⁹Z. Li, M. H. Schulz, T. Look, M. Begemann, M. Zenke, and I. G. Costa, “Identification of transcription factor binding sites using atac-seq”, *Genome biology* **20**, 1–21 (2019).
- ¹⁰A. P. Boyle, S. Davis, H. P. Shulha, P. Meltzer, E. H. Margulies, Z. Weng, T. S. Furey, and G. E. Crawford, “High-resolution mapping and characterization of open chromatin across the genome”, *Cell* **132**, 311–322 (2008).
- ¹¹A. H. Yona, E. J. Alm, and J. Gore, “Random sequences rapidly evolve into de novo promoters”, *Nature communications* **9**, 1530 (2018).

- ¹²A. Bartlett, R. C. O'Malley, S.-s. C. Huang, M. Galli, J. R. Nery, A. Gallavotti, and J. R. Ecker, "Mapping genome-wide transcription-factor binding sites using dap-seq", *Nature protocols* **12**, 1659–1672 (2017).
- ¹³S. D. Hebdon, A. T. Gerritsen, Y.-P. Chen, J. G. Marciano, and K. J. Chou, "Genome-wide transcription factor dna binding sites and gene regulatory networks in *clostridium thermocellum*", *Frontiers in Microbiology* **12**, 695517 (2021).
- ¹⁴Y. Zhang, Y. Wang, Y. Zhang, X. Jia, C. Li, Z. Zhou, S. Hu, and Z. Li, "Genome-wide analysis reveals that phop regulates pathogenicity in *riemerella anatipestifer*", *Microbiology Spectrum* **10**, e01883–22 (2022).
- ¹⁵G. Urtecho, A. D. Tripp, K. D. Insigne, H. Kim, and S. Kosuri, "Systematic dissection of sequence elements controlling $\sigma 70$ promoters using a genomically encoded multiplexed reporter assay in *escherichia coli*", *Biochemistry* **58**, 1539–1551 (2018).
- ¹⁶G. Urtecho, K. D. Insigne, A. D. Tripp, M. Brinck, N. B. Lubock, H. Kim, T. Chan, and S. Kosuri, "Genome-wide functional characterization of *escherichia coli* promoters and regulatory elements responsible for their function", *BioRxiv* (2020).
- ¹⁷W. T. Ireland et al., "Deciphering the regulatory genome of *escherichia coli*, one hundred promoters at a time", *Elife* **9**, e55308 (2020).
- ¹⁸J. B. Kinney, A. Murugan, C. G. Callan Jr, and E. C. Cox, "Using deep sequencing to characterize the biophysical mechanism of a transcriptional regulatory sequence", *Proceedings of the National Academy of Sciences* **107**, 9158–9163 (2010).
- ¹⁹N. M. Belliveau, S. L. Barnes, W. T. Ireland, D. L. Jones, M. J. Sweredoski, A. Moradian, S. Hess, J. B. Kinney, and R. Phillips, "Systematic approach for dissecting the molecular mechanisms of transcriptional regulation in bacteria", *Proceedings of the National Academy of Sciences* **115**, E4796–E4805 (2018).
- ²⁰S. L. Barnes, N. M. Belliveau, W. T. Ireland, J. B. Kinney, and R. Phillips, "Mapping dna sequence to transcription factor binding energy in vivo", *PLoS computational biology* **15**, e1006226 (2019).
- ²¹A. I. Flamholz et al., "Functional reconstitution of a bacterial co2 concentrating mechanism in *escherichia coli*", *Elife* **9**, e59882 (2020).
- ²²A. Schmidt et al., "The quantitative and condition-dependent *escherichia coli* proteome", *Nature biotechnology* **34**, 104–110 (2016).
- ²³S. H. Saunders and A. M. Ahmed, "Orbit for *e. coli*: kilobase-scale oligonucleotide recombineering at high throughput and high efficiency", *bioRxiv*, 2023–06 (2023).
- ²⁴S. Ghatak, Z. A. King, A. Sastry, and B. O. Palsson, "The y-ome defines the 35% of *escherichia coli* genes that lack experimental evidence of function", *Nucleic acids research* **47**, 2446–2454 (2019).
- ²⁵C. R. Lamoureux, K. T. Decker, A. V. Sastry, J. L. McConn, Y. Gao, and B. O. Palsson, "Precise 2.0-an expanded high-quality rna-seq compendium for *escherichia coli* k-12 reveals high-resolution transcriptional regulatory structure", *BioRxiv* (2021).
- ²⁶T. L. LaFleur, A. Hossain, and H. M. Salis, "Automated model-predictive design of synthetic promoters to control transcriptional profiles in bacteria", *Nature communications* **13**, 5159 (2022).
- ²⁷S. H. Yun, S. C. Ji, H. J. Jeon, X. Wang, S. W. Kim, G. Bak, Y. Lee, and H. M. Lim, "The *cnuk9e* hns complex antagonizes dna binding of *dica* and leads to temperature-dependent filamentous growth in *e. coli*", (2012).

- 239 ²⁸Y. Han, W. Li, A. Filko, J. Li, and F. Zhang, “Genome-wide promoter responses to crispr perturba-
240 tions of regulators reveal regulatory networks in escherichia coli”, *Nature communications* **14**, 5757
241 (2023).
- 242 ²⁹Y. Iwadate and J.-i. Kato, “Identification of a formate-dependent uric acid degradation pathway in
243 escherichia coli”, *Journal of bacteriology* **201**, e00573–18 (2019).
- 244 ³⁰T. Basu and R. Poddar, “Effect of ethanol on escherichia coli cells. enhancement of dna synthesis due
245 to ethanol treatment”, *Folia microbiologica* **39**, 3–6 (1994).
- 246 ³¹B. Ezraty, A. Gennaris, F. Barras, and J.-F. Collet, “Oxidative stress, protein damage and repair
247 in bacteria”, *Nature Reviews Microbiology* **15**, 385–396 (2017).
- 248 ³²S. Ravindra Kumar and J. A. Imlay, “How escherichia coli tolerates profuse hydrogen peroxide
249 formation by a catabolic pathway”, *Journal of bacteriology* **195**, 4569–4579 (2013).

Tom Röschinger¹, Grace Solini¹, Kian Faizi¹, Rosalind Pan¹, Anika Nawar Choudhury², Lynn Yang³, Bob Jones³, Stephen Quake^{2, 3, 4}, and Rob Phillips^{1, 5, +}

¹*Division of Biology and Biological Engineering, California Institute of Technology, Pasadena, CA 91125, USA*

²*Chan Zuckerberg Biohub, San Francisco, CA 94158, USA*

³*Department of Bioengineering, Stanford University, Stanford, CA 94305, USA*

⁴*Department of Applied Physics, Stanford University, Stanford, CA 94305, USA*

⁵*Department of Physics, California Institute of Technology, Pasadena, CA 91125, USA*

⁺*Correspondence: phillips@pboc.caltech.edu*

S1 Finding number of genes without regulatory annotation

S1.1 *E. coli* K12 MG1655

S1.2 *Bacillus Subtilis*

S1.3 *Pseudomonas Aeruginosa*

S1.4 *Saccharomyces cerevisiae*

S1.5 *Drosophila Melanogaster*

S1.6 *C. elegans*

S2 Reporter Sequence design

S3 Oligo Pool Design

S3.1 Identification of Transcription Start Sites

All oligo pools used in this work were manually designed. For each gene in our list we looked for promoters in Ecocyc [33] (accessed 12/08/2021) using the transcription start site if the promoter was found. If multiple promoters were identified, each promoter was included in the experiment. If no promoter was found, we looked for transcriptionally active sites in the data set from Urtecho et al, 2020[16]. In their work, every part of the genome was tested for transcription initiation in LB. If we could find a site that was identified as active close to the gene of interest, we chose this site as origin for computational promoter mutagenesis. If no transcription start site could be identified for a gene, the model from [26] was used to computationally predict a transcription start site in the intergenic region. The site predicted to be the most active within 500 bp upstream of the coding region was chosen as transcription start site since more than 99% of transcription start sites are within that region in *E. coli* K12 MG1655, see Fig. S3. Initially, we chose 119 promoters, however, 7 promoters (mglBp, hdeAp2, mtnp, ybeDp, cpxRp2, galEp1, ompFp), were duplicates of promoter annotated in Ecocyc. The duplicated promoters were treated as independent when mutated variants were created, leading to twice the number of variants.

S3.2 Computational Promoter Mutagenesis

Once a TSS is identified, the 160 bp region from 115 bp upstream of the TSS to 45 bp downstream is taken from the genome. It has been shown that most cis-regulation is happening within this window [34]. Based on the approach by [18], each promoter sequence is mutated randomly at a rate of 0.1

per position. 1500 mutated sequences are created per promoter, following the approach from [17], which creates sufficient mutational coverage across the window. The promoter oligonucleotides are flanked by restriction enzyme sites (*rs1* and *rs2* in Fig. S4) that are used in downstream cloning steps. The restriction sites are flanked by primer sites used to amplify the oligo pool. Primer sequences were chosen from a list of orthogonal primer pairs, designed to be optimal for cloning procedures [35]. Oligo pools were synthesized (TwistBioscience, San Francisco, CA, USA) and used for subsequent cloning steps.

S4 Library Cloning

S4.1 Cloning oligo pool into plasmid vector

The oligo pool was amplified using a 20bp forward primer (SC142) and a 40 bp reverse primer (SC143), which consists of 20bp primer binding site and 20bp overhang. PCR amplifications were run to minimal amplification (faint band on agarose gel) to minimize amplification bias. PCR products were cleaned and concentrated (DNA Clean & Concentrator-5, ZymoResearch) and used for a second amplification step. The 20 bp overhang from the first amplification was used as primer site for a reverse primer (SC172), which contains randomized 20 bp barcode, flanked by two restriction enzyme sites (*rs3* and *rs4* in Fig. S4). The forward primer is the same as in the first amplification step. PCR amplification is run again to minimal amplification to minimize amplification bias. PCR products are run on a 2% agarose TAE gel and subsequently extracted and purified (Zymoclean Gel DNA Recovery Kit, ZymoResearch). In the next step, restriction digest is performed on the outer restriction enzyme sites (*rs1* and *rs4* in Fig. S4). Unless noted otherwise, all restriction digests were run for 15 minutes at 37C. The plasmid vector was digested with different restriction enzymes which create compatible sticky ends. Most restriction enzyme sites are palindromes, so by choosing different enzymes with compatible ends, we avoid having palindromes flanking the plasmid inserts. This is important, since these sites are used for amplifications in the library preparation steps later in the protocol. (Maybe not needed to say). The oligo pool is combined with the plasmid vector using T7 DNA ligase (New England Biolabs, Ipswich, MA, USA) following the suppliers protocol. Ligation products were cleaned and concentrated (DNA Clean & Concentrator-5, ZymoResearch) and drop dialysis (MF-Millipore VSWP02500, MilliporeSigma, Burlington, MA, USA) was performed for 1h to improve sample purity. Electroporation using *E. coli* pir116 electrocompetent cells (Lucigen, Middleton, WI) was performed at 1.8kV in 1mm electroporation cuvettes, followed by 1h recovery at 37C and 250rpm in 1 ml LB-media (details here, the same for all following mentionings of LB). The entire cultures were plated on 150mm kanamycin (50µg/ml) + LB petri dishes and grown overnight. The following day, plates were scraped and the colonies resuspended. Freezer stocks were prepared using a 1:1 dilution of resuspended colonies and 50% glycerol. Cultures were inoculated with 5×10^8 cells in 200ml of LB + kanamycin (50µg/ml) and grown at 37C until saturation. Plasmid was extracted (ZymoPURE II Plasmid Maxiprep Kit, ZymoResearch) and used subsequent sequencing (see S5). The plasmid library is then used as template in a restriction digest using restriction enzymes *rs2* and *rs3*. The resulting product was cleaned and concentrated (NEB Monarch) and concentration measured on a Nanodrop. Similarly, the *riboJ::YFP* element was PCR amplified (primers SC191 and SC192), adding restriction sites as overhangs (see table S1). The PCR product was cleaned and concentrated (NEB Monarch) and digested with the respective restriction enzymes. The plasmid library is combined with the *RiboJ::sfYFP* element using 7 DNA ligase (New England Biolabs, Ipswich, MA, USA) following the suppliers protocol. Ligation products were cleaned and concentrated (NEB Monarch) and drop dialysis

Part	5' restriction site	3' restriction site
Plasmid Vector	XbaI	XhoI
RiboJ::YFP	ApaI	PtsI-HF
Oligo Pool	SpeI-HF	ApaI
Barcoding Primer	SbfI-HF	Sall-HF

Table S1. Restriction sites used. All enzymes were ordered from NEB (check which ones are high fidelity versions)

(MF-Millipore VSWP02500, MilliporeSigma, Burlington, MA, USA) was performed for 1h to improve sample purity. Electroporation using *E. coli* pir116 electrocompetent cells (Lucigen, Middleton, WI) was performed at 1.8kV in 1mm electroporation cuvettes, followed by 1h recovery at 37C and 250rpm in 1 ml LB-media. The entire cultures were plated on 150mm kanamycin (50 μ g/ml) + LB petri dishes and grown overnight. The following day, plates were scraped and the colonies resuspended. Freezer stocks were prepared using a 1:1 dilution of resuspended colonies and 50% glycerol. Cultures were inoculated with 5×10^8 cells in 200ml of LB + kanamycin (50 μ g/ml) and grown at 37C until saturation. Plasmid was extracted (ZymoPURE II Plasmid Maxiprep Kit, ZymoResearch) and used for subsequent genome integration.

S5 Barcode Mapping

The plasmid library is used for barcode mapping. Purified plasmid is PCR amplified using forward primer (SC185) outside the promoter region and a reverse primer outside the 20bp barcode (SC184). The PCR is run to minimal amplification (until a band is visible on an agarose gel), and the product is gel purified (NEB Monarch). The purified DNA was used as template for a second PCR using a primer (SC196) adding an Illumina P5 adapter to the promoter side, and a primer (SC199) adding an Illumina P7 adapter. The PCR is again run to minimal amplification and gel purified (NEB Monarch). The product was used for sequencing on a Illumina NextSeq P2 flow cell with pair end reads using primers SC185 for read 1, SC184 for read 2 and SC201 for the index read. Reads were filtered and merged using custom bash scripts, which are available in the Github repository. After processing, each promoter/barcode pair was identified in each read, and pairs with less than 3 total reads were discarded. An alignment algorithm was used to identify the identity of each sequenced promoter variant. This allowed to include additional promoter variants that were in the initial oligo pool due to synthesis errors in the production of the oligos. The barcode mapping was used in analysis of libraries grown in various growth conditions. The code used to perform processing of sequencing data can be found in the associated Github repository https://github.com/RPGroup-PBoC/1000_genes_ecoli/tree/main/code/processing/20220514_mapping. Processing is done with the help of various software modules [36–38]. Custom Julia code used for analysis and visualization of results can be found in the associated Github repository https://github.com/RPGroup-PBoC/1000_genes_ecoli/tree/main/code/processing/20220514_mapping

S5.1 Genome Integration

We used ORBIT to integrate the reporter libraries into the chromosome. A detailed description of the method and its efficiencies can be found in (Add scotts paper here). Wild type *E. coli* (K12 MG1655) are streaked on a LB plate and grown overnight at 37C. A single colony is picked and grown in 3ml of LB at 37C and shaken at 250rpm overnight. The overnight culture is diluted 1:1000 into fresh LB (e.g. 200ml) and grown at 37C and 250rpm until exponential phase (~ 0.4 OD 600nm). The cultures are then immediately put on ice and spun in a centrifuge at 5000g for 10min. Following the spin, the supernatant is discarded, and the cells are resuspended in deionized water at 4C at the same volume as the initial culture. The cells are spun again at 5000g for 10 min. This wash step is repeated 4 times with 10% glycerol. After the last wash, supernatant is discarded and cells are resuspended in the remaining liquid and distributed into 50 μ l aliquots. Aliquots are frozen on dry ice and kept at -80C until used for electroporation. For electroporation, aliquots are thawed on ice and 1mm electroporation cuvettes are pre-chilled on ice. 100ng of helper plasmid (link to helper plasmid file) is added to a 50 μ l cell aliquot and mixed by slowly pipetting up and down. The aliquot is then added to the electroporation cuvette and electroporation is performed at 1.8kV. The aliquot is recovered with 1ml of LB media prewarmed to 37C for an 1h prior to electroporation. The culture is recovered for 1h at 37C and shaken at 250rpm. After recovery, aliquots at various dilutions are plated on LB + gentamycin (check gent concentration). Plates are grown overnight and a single colony is picked to prepare frozen stocks as described above. To perform genome integration, the host strain carrying the helper plasmid is made electrocompetent (follow growing and washing steps described above), and the plasmid library is electroporated into the host strain. The cells are recovered in 3ml of prewarmed LB + 1% arabinose and shaken at 37C at 250rpm for 1h. The entire volume is plated on LB + kanamycin plates (TR: add concentration) and colonies are grown over night. The next day, colonies are scraped, resuspended in LB and diluted to optical density of 1 at 600nm. The helper plasmid used for genome integration causes growth deficits, hence, the library needs to get cured of the plasmid. Therefore, the library is inoculated with 0.5ml of culture at 1 OD in 200ml of LB, and grown until exponential phase at 37C shaken at 250 rpm. The helper plasmid carries the *sacB* gene, which is used for negative selection in the presence of sucrose. At exponential phase, the culture is plated on LB + 7.5% sucrose agarose plates. Plates are grown overnight, scraped and made into frozen stocks at an OD600 of 1. The frozen stocks are then ready for growth experiments.

S6 Growth Media and Culture Growth

Base Media Lysogeny Broth (LB) was prepared from powder (BD Difco, Tryptone 10g/L, Yeast extract 5g/L, Sodium Chloride 10g/L), and sterilized by autoclaving. M9 Minimal Media pre-mix without carbon source was prepared in the following way, similar to [22]: to 700ml of ultrapure water, 200 ml of 5 \times base salt solution (BD Difco, containing Disodium Phosphate (anhydrous) 33.9g/L, Monopotassium Phosphate 15g/L, Sodium Chloride 2.5g/L, Ammonium Chloride 5g/L, in H₂O, autoclaved), 10 ml of 100x trace elements (5g/L EDTA, 0.83 g/L FeCl₃-6H₂O, 84 mg/L ZnCl₂, 19 mg/L CuSO₄ - 5 H₂O, 10 mg/L CoCl₂ - 6H₂O in H₂O, 10 mg/L H₃BO₃, 1.6 mg/L MnCl₂ - 4H₂O, prepared as discribed in [21]), 1 ml 0.1 M CaCl₂ solution, in H₂O, autoclaved, 1 ml 1 M MgSO₄ solution, in H₂O, autoclaved and 1 ml of 1000 \times thiamine solution (1mg/ml in H₂O, filter sterilized) were added. The resulting solution was filled up to 1 l with water and filter sterilized. M9 minimal medium was complemented with carbon source by mixing appropriate amounts of carbon source free M9 minimal medium and carbon source stock solutions. Carbon source stock solutions were prepared as 20% solutions and filter sterilized.

Cultivation Overnight cultures were incubated from frozen stock in 200ml M9 Minimal Media with 0.5% Glucose and grown at 37C while shaken at 250rpm. Cultures were diluted 1:100 into the respective growth media (prewarmed to 37C, 200ml) and grown to exponential phase (OD600nm of 0.3). To ensure steady state growth, the cultures were diluted a second time 1:100 into the same growth media and grown again to exponential phase, ensuring at least 10 cell divisions in the growth media. At this step, cultures were either harvested or exposed to environmental stresses.

Specific Growth Conditions **Glucose:** 5ml of 20% Glucose solution added to 200ml of M9 Minimal Media pre-mix for a final concentration of 0.5%. **Xylose:** 5ml of 20% Xylose solution added to 200ml of M9 Minimal Media pre-mix for a final concentration of 0.5%. **Arabinose:** 5ml of 20% Arabinose solution added to 200ml of M9 Minimal Media pre-mix for a final concentration of 0.5%. **Galactose:** 2.3ml of 20% Xylose solution added to 200ml of M9 Minimal Media pre-mix for a final concentration of 0.5%. **Sodium Salicylate:** 1M Sodium Salicylate stock was prepared (TR: add vendor) and filter sterilized.

For each growth condition, cultures were inoculated from frozen stocks in 200ml of M9 Minimal Media with 0.5 % Glucose (details) and grown overnight at 37C shaken at 250rpm. In the morning, cultures were diluted 1:100 into the growth media of choice, which, unless noted otherwise is at a volume of 200ml. Cultures are grown at 37C until reaching exponential phase (OD600 of 0.4). Once the culture reaches exponential phase, 1ml aliquots are spun down at 5000g for 10min for subsequent gDNA extraction. Supernatant was discarded and pellets were frozen at -20C overnight. For RNA extraction, culture was diluted in RNAlater (Qiagen, add info) and 1ml aliquots are spun down at 5000g for 10min Supernatant was discarded and pellets were frozen at -20C overnight. (TR: Add details for special growth conditions)

S6.1 gDNA and RNA extractions

S7 Barcode Sequencing

S8 Promoter footprints

If a base of a binding site for a regulatory element in a promoter is mutated, the expression of the downstream gene is changed due to differences in binding affinity of the regulatory element (TR: cite Kinney, 2010 and Garcia, 2011; but maybe find some older/more original references). One can generate so called *footprints*, where the effect of a mutation in the promoter on expression levels can be quantified by various metrics. Here, we explore various ways to compute footprints and explain each method in detail. (TR: add the footprints from one real dataset to compare)

S8.1 Dataset

For a given promoter, there are $i = 1, \dots, n$ promoter variants, where each variant has m_i unique barcodes. Per barcode, there are c_{dna} counts from genomic DNA sequencing, as well as c_{rna} counts from RNAseq. DNA sequencing is performed to normalize the RNA sequencing data by the abundance of cells in the culture expressing the reporter from a specific promoter variant.

c_{dna}	c_{rna}	sequence
10	2	ACGTACGTAC
1	2	ACGTACGTTC
3	5	ACGTACGTTC
4	9	ACGTACGTTC
3	5	ACGTAAGAAC
3	6	ACGTAAGAAC
15	12	GCGTACGTAC
5	3	GCGTACGTAC
12	14	ACATACGTAC
2	3	ACATACGTAC
20	40	ACATACGTAC
5	3	ACGGATGTAC
5	1	ACGTACGTGA
10	1	ACGTACGTGA
2	10	ACGTCCATAC
2	10	ACGTCCATAC
4	13	ACGTCCGTAC
18	25	ACAAACGTAC
17	19	GCGTACGTAG
10	11	GCGTACGTAG
2	3	GGGTACGTAG

Table S2. Example dataset, arbitrarily generated. For each sequence, there are counts from RNA and DNA sequencing. Different counts for the same sequence come from unique barcodes, are therefore separate measurements. (TR: has to be updated to have the correct sequences for the figures below)

S8.2 Frequency Matrices

(TR: Not sure if I will actually write about it, just a different way of computing footprints I came up with based on comments by Frank Jülicher and Stephan Grill. Have try it on old data set.)

S8.3 Expression Shifts

Belliveau et al. (2018)[19] used so called *expression shifts* to compute footprints for mutagenized promoters. In their experiments, cells were sorted based on fluorescence, where the fluorescent reporter gene was expressed under the control of a mutagenized promoter variant, and subsequently sequenced. Therefore, each sequence had a bin associated with it, which is a read out for how strong the reporter is expressed relatively to the other promoter variants in the library. This approach can be adapted to our data set, where we first compute the average relative expression $\langle c \rangle_i$ for the i -th promoter variant across all of its unique barcodes,

$$\langle c \rangle_i = \frac{1}{m_i} \sum_{j=1}^{m_i} \frac{c_{\text{rna},j}}{c_{\text{dna},j}}. \quad (\text{S1})$$

Then, we determine how much relative expression is changed at each position if there is a mutation. If a base at position ℓ in promoter variant i is mutated, we denote that as $\sigma_{i,\ell} = 1$. Otherwise, if the base is wild type, we write $\sigma_{i,\ell} = 0$. Then, the change in relative expression due to mutation, the expression shift Δc_ℓ , at position ℓ is given by

$$\Delta c_\ell = \frac{1}{n} \sum_{i=1}^n \sigma_{i,\ell} \left(\langle c \rangle_i - \frac{1}{n} \sum_{k=1}^n \langle c \rangle_k \right). \quad (\text{S2})$$

The absolute value of expression shift can be hard to interpret, so indeed one can present it in terms of relative change to the mean expression, i.e., fold-change,

$$\delta c_\ell = \frac{\Delta c_\ell}{\langle c \rangle} = \frac{1}{n} \sum_{i=1}^n \sigma_{i,\ell} \left(\frac{\langle c \rangle_i}{\langle c \rangle} - 1 \right). \quad (\text{S3})$$

Figure ?? shows the expression shift footprint that is obtained for the test dataset. (TR: as well as the footprint for a real data set).

S8.4 Mutual Information

Mutual information is a measure of how much information is obtained about a random variable by measuring a different random variable. In the context of gene expression, this can be understood as the ability to predict changes in gene expression given a certain mutation on the promoter sequence. If there is no annotation, meaning it is unknown where RNAP or transcription factors bind, one can not make any predictions on the expression level of the downstream gene when observing a mutation in the promoter. In this case, there is low mutual information between sequence and expression level. On the other hand, if the promoter is annotated and one has binding energy matrices for all transcription factor binding sites and the RNAP binding site in hand, then one can precisely predict the change in gene expression given any point mutation based on thermodynamic models (TR: could cite a bunch of papers here), which is a case of high mutual information. Hence, by maximizing the mutual information between a model for the regulatory architecture and observed levels of gene expression, we can discover binding sites for transcription factors and subsequently, using equilibrium thermodynamic models and neural networks, compute binding energy matrices in real units of $k_B T$.

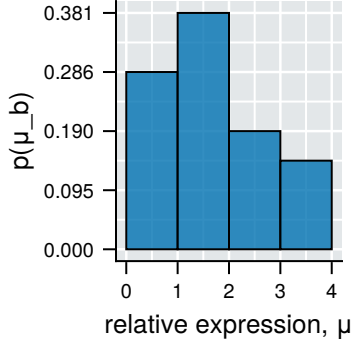


Figure S1. Possible binning of expression counts for example data set. (TR: Add footprint for real dataset.)

S8.4.1 Mutual Information based on Sequence Counts

The first way of computing mutual information at each position in the promoter is to take the base at each position as one random variable, and the expression of each sequence as other random variable. As measure for expression, we use RNA counts for each sequence normalized by DNA counts. In order to compute mutual information, we need to obtain a probability distribution $p_\ell(c, \mu)$, which gives the probability of finding a certain base c at position ℓ , and corresponding expression μ . One way of obtaining such a distribution is to find bins for the values of μ , denoted as μ_b , as shown in Figure S1. Then, mutual information is given by

$$I_\ell = \sum_{c=A,C,G,T} \sum_{\mu_b} p_\ell(c, \mu_b) \log_2 \left(\frac{p_\ell(c, \mu_b)}{p_\ell(c) p(\mu_b)} \right), \quad (\text{S4})$$

where $p_\ell(c)$ and $p(\mu_b)$ are the marginal distributions.

S8.4.2 Mutual Information based on Phenotype Matrices

A different way to utilize mutual information is to choose a phenotype as random variable instead of base identity. In this case, the phenotype Φ is a real number and is additive across the sequence, meaning that each position l with base c contributes $\Theta_{l:c}$ to the total phenotype. The contributions are independent, i.e., no epistasis effects are considered for this model. The phenotype Φ is then determined by the sum across all positions with a possible offset Θ_0 ,

$$\Phi = \Theta_0 + \sum_{l=1}^L \sum_c \Theta_{l:c} x_{l:c}, \quad (\text{S5})$$

where $x_{l:c}$ is a one-hot representation of the sequence with

$$x_{l:c} = \begin{cases} 1 & \text{if character } c \text{ occurs at position } l, \\ 0 & \text{otherwise} \end{cases} \quad (\text{S6})$$

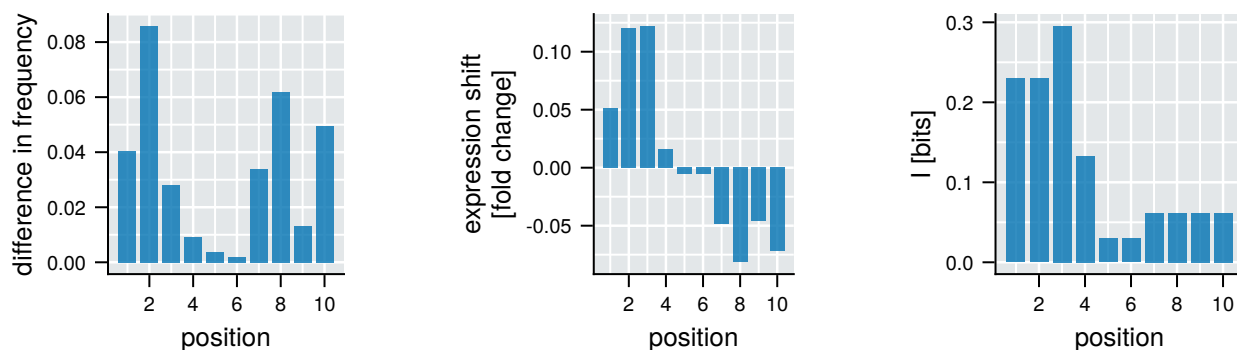


Figure S2. Different ways of computing footprints for test data set from table S2. Frequency Matrix left, Expression Shift middle, Mutual information right

where the notation is adapted from [39]. Without any knowledge of the regulatory architecture of the promoter, one can only make random guesses for the phenotype matrix. However, either using Metropolis-Hasting algorithms (TR: Reg-Seq, gotta decide how much to write about it) or Neural Networks (TR: MaveNN, will be included if we get good results with it), the phenotype matrix can be optimized in the sense its entries are more extreme where there are binding sites for regulatory elements in the sequence, since a mutation in that part of the sequence will have the strongest effect on gene expression. How extreme entries are can be quantified by using relative entropy, where the entries for each position on the sequence are first converted to a probability distribution using exponential weights, and then Kullback-Leiback-Divergence (KLD) between the resulting distribution and a uniform distribution is calculated. (TR: Expand by explaining how peaks are identified as binding sites.)

S8.4.3 Phenotype Matrices and Neural Networks, MaveNN

S8.4.4 Identifying Binding Energy Matrices

S9 Supplementary Files

- Plasmid Sequences with annotations + RiboJ::YFP
- pHelper sequence
- Primers
- list of restriction sites used in cloning
- Gene list
- Sequencing Data
- List of ordered sequences

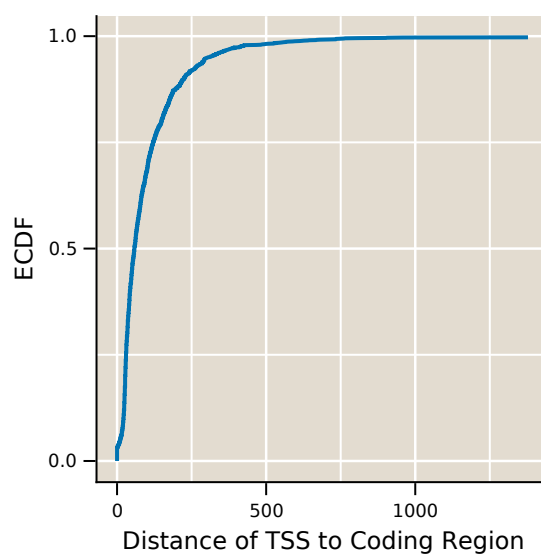


Figure S3. ECDF of distances of transcription start sites to the coding region for every operon in *E. coli* that has a transcription start site annotated in EcoCyc.

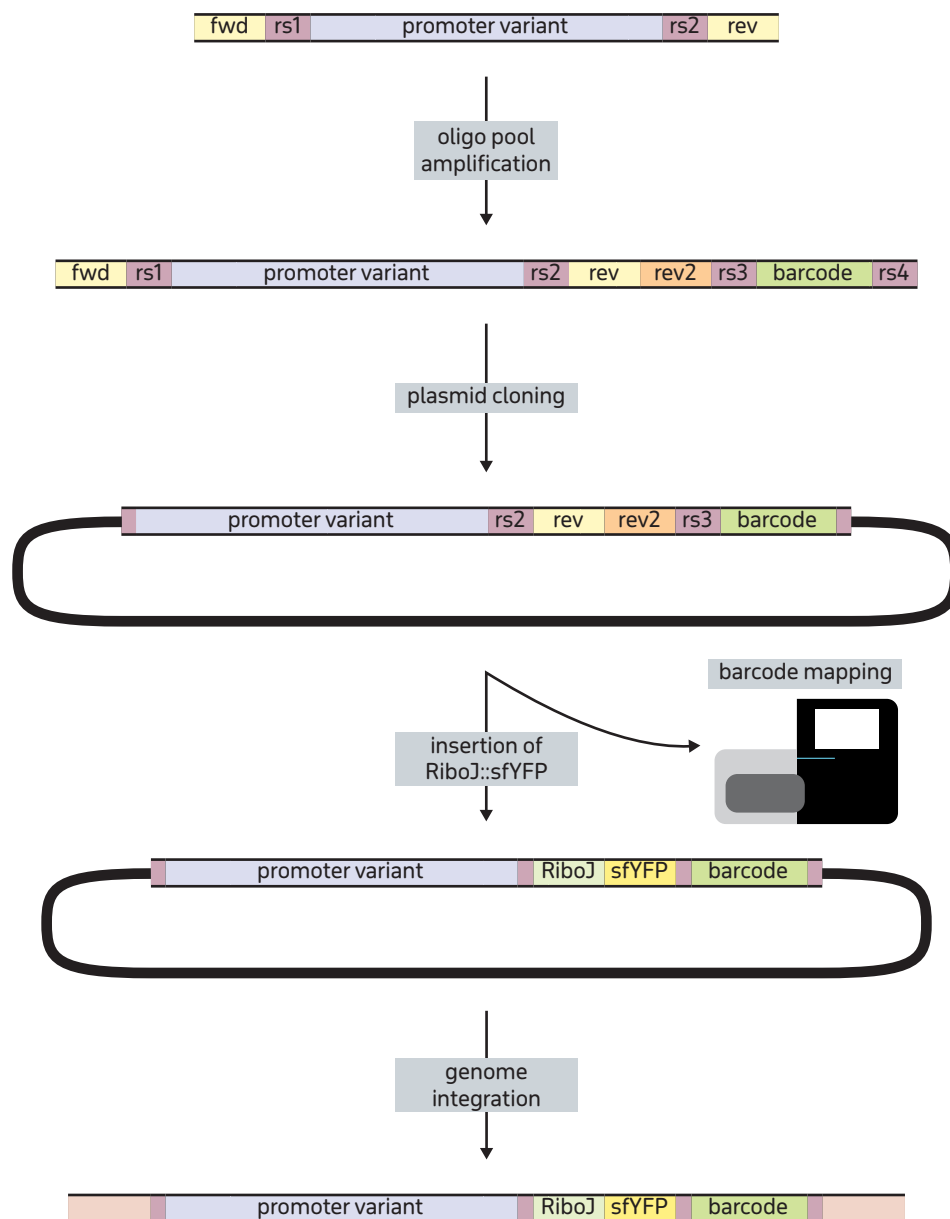


Figure S4. Placeholder figure for cloning scheme.

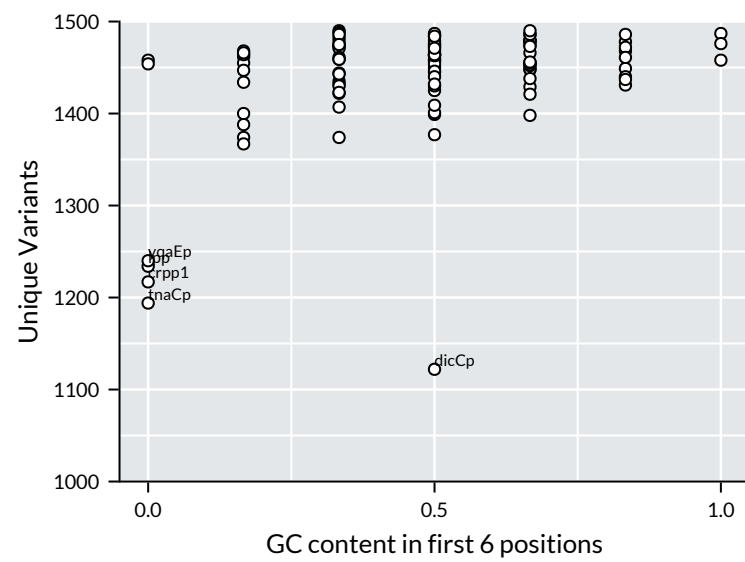


Figure S5

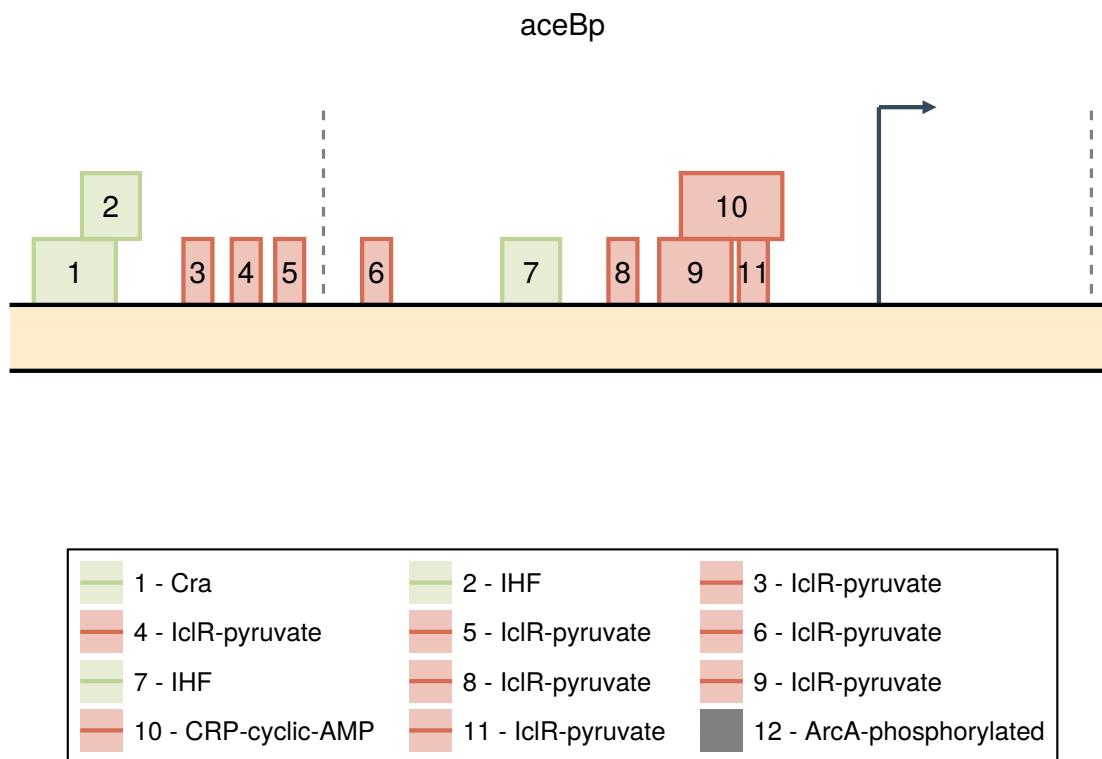


Figure S6

S10 Individual promoter architectures

S10.1 aceB

Malate synthase catalyzes the Claisen condensation of glyoxylate and acetyl-CoA, initially forming a malyl-CoA intermediate that is hydrolyzed to the products malate and CoA. [40][41]

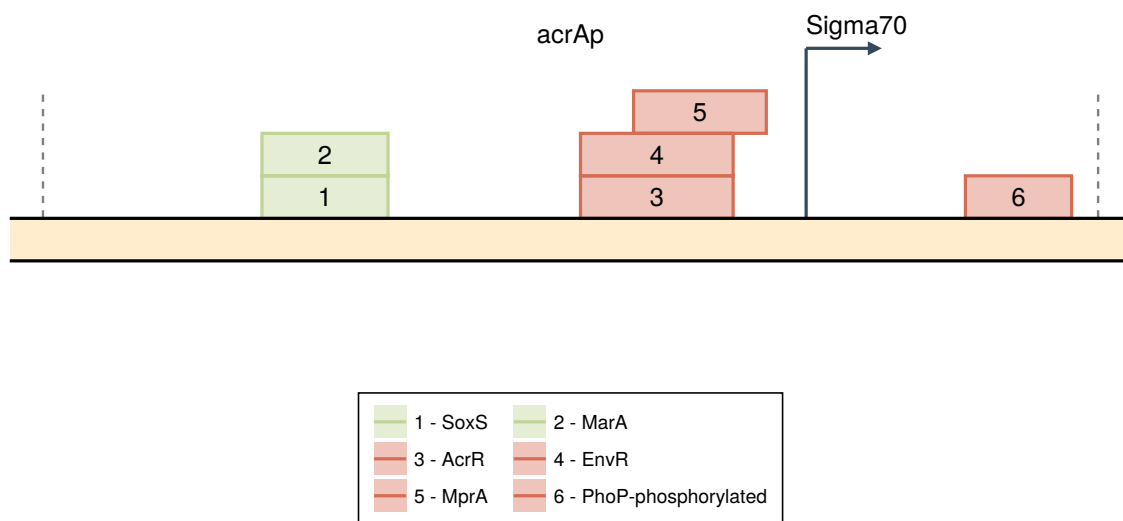


Figure S7

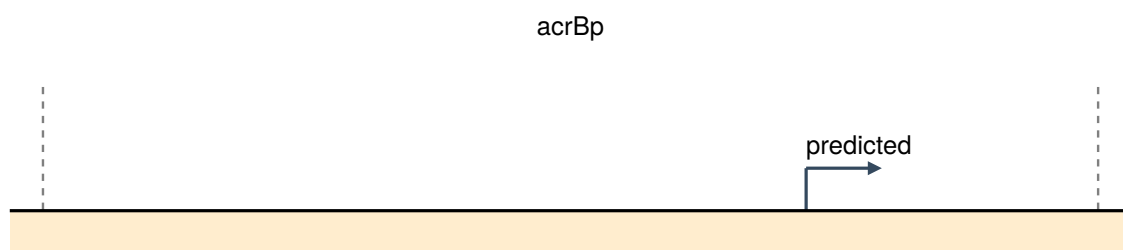


Figure S8

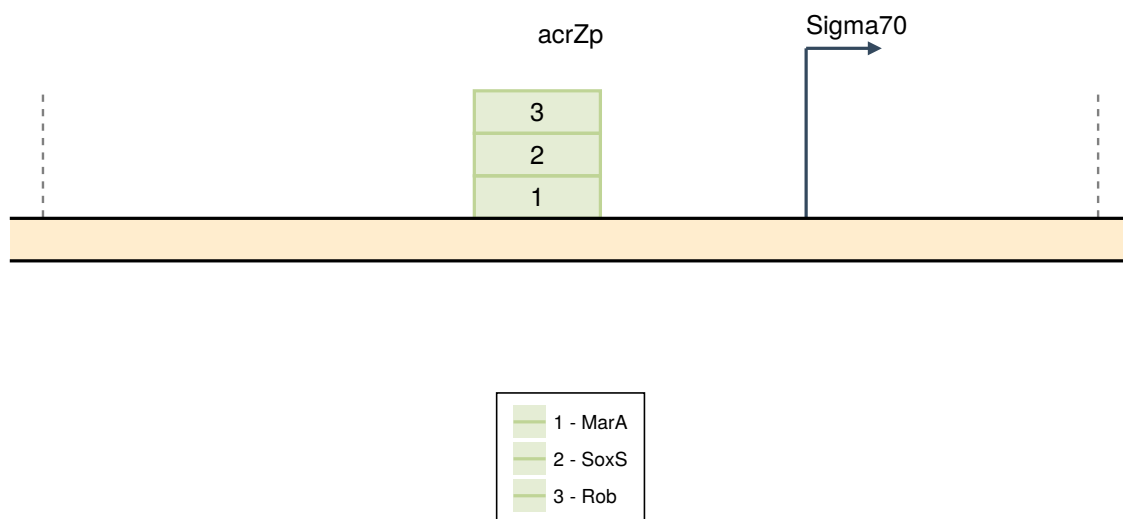


Figure S9

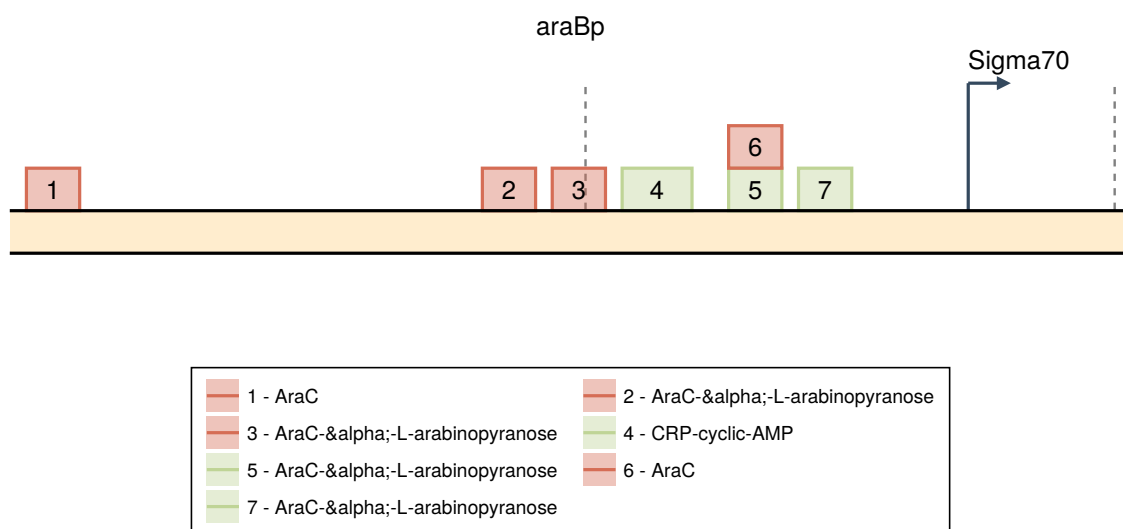


Figure S10

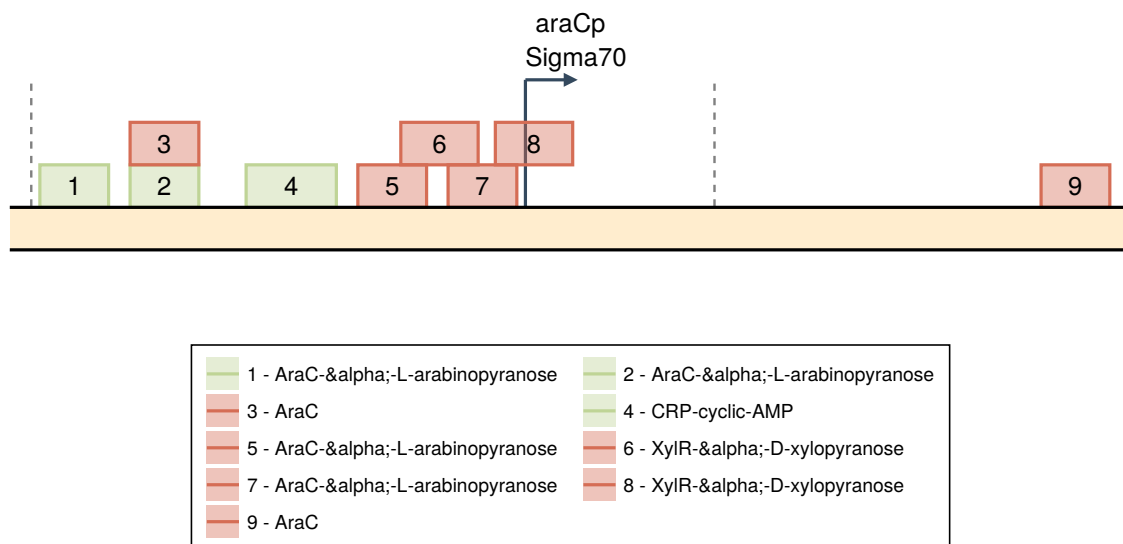


Figure S11

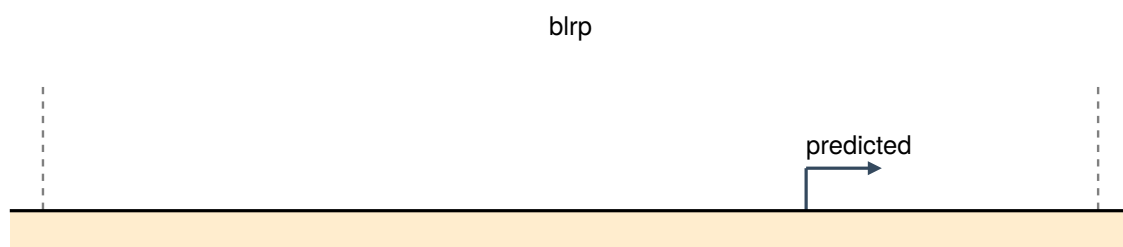


Figure S12

515	S10.2	acrA
516	S10.3	acrB
517	S10.4	acrZ
518	S10.5	araB
519	S10.6	araC
520	S10.7	blr
521	S10.8	cpxR

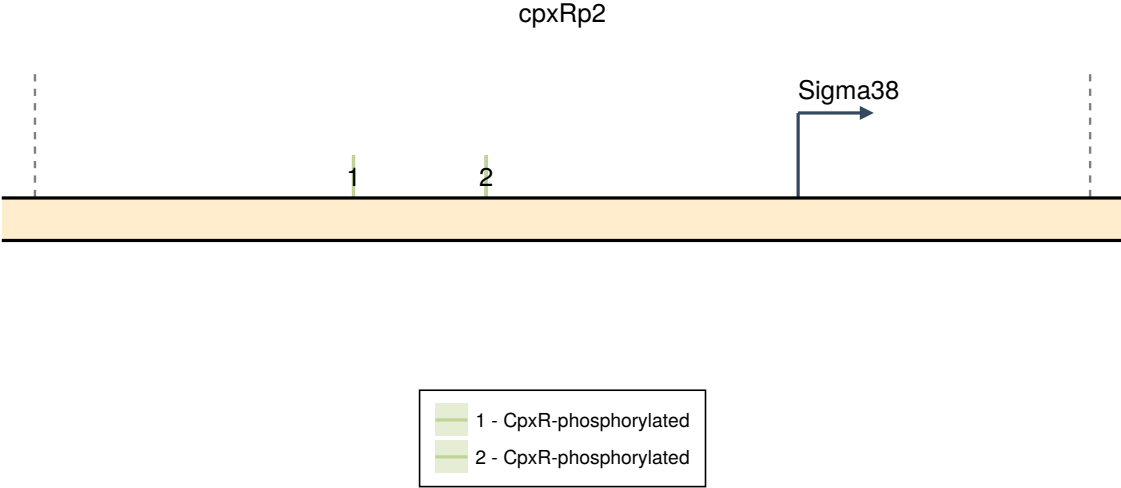


Figure S13

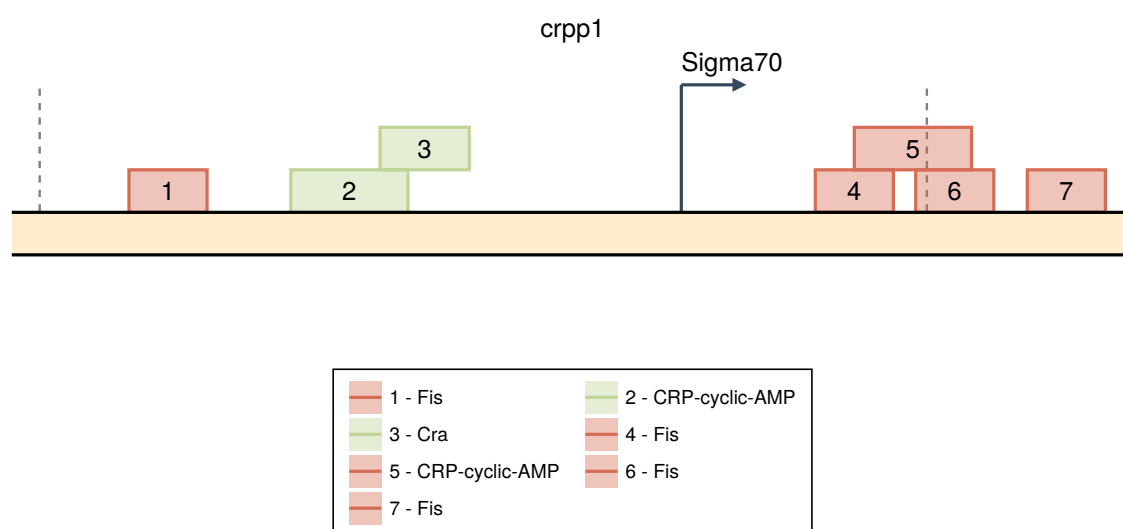


Figure S14

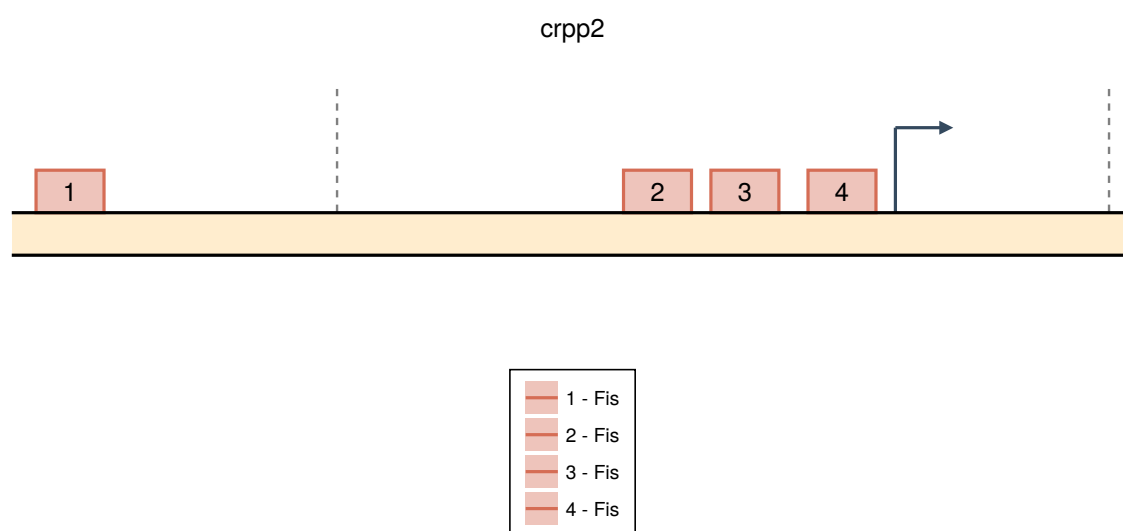


Figure S15

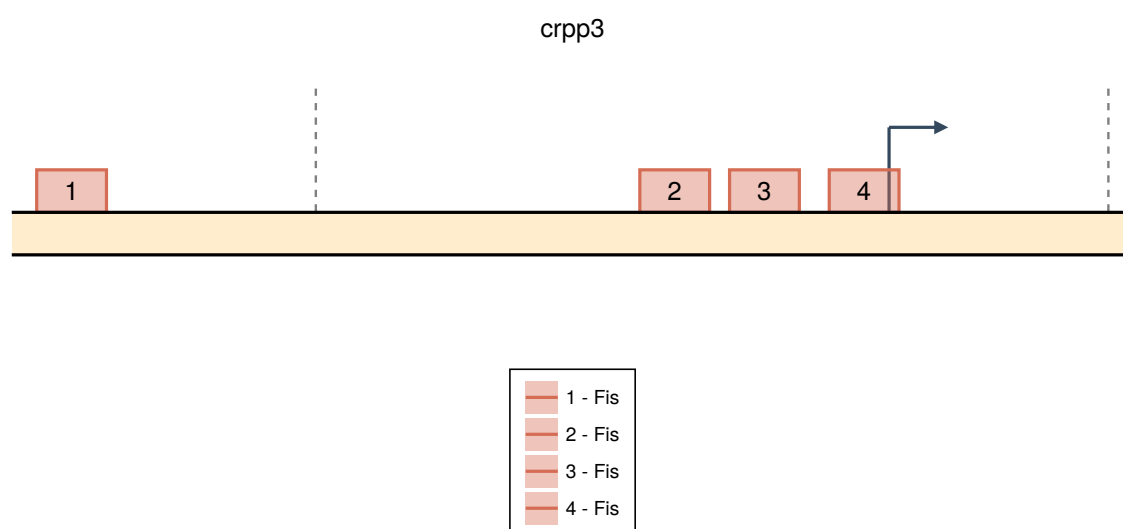


Figure S16

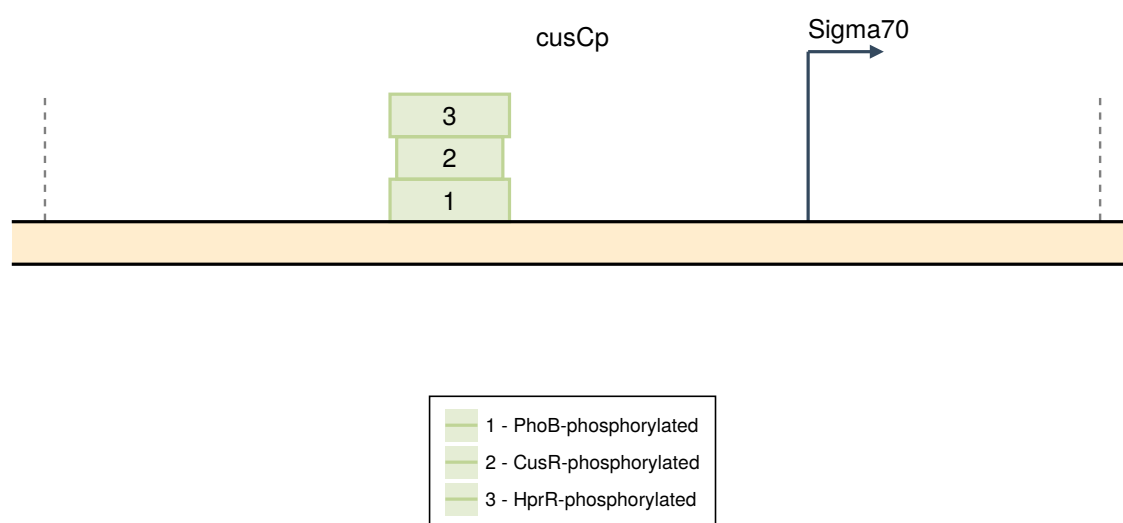


Figure S17

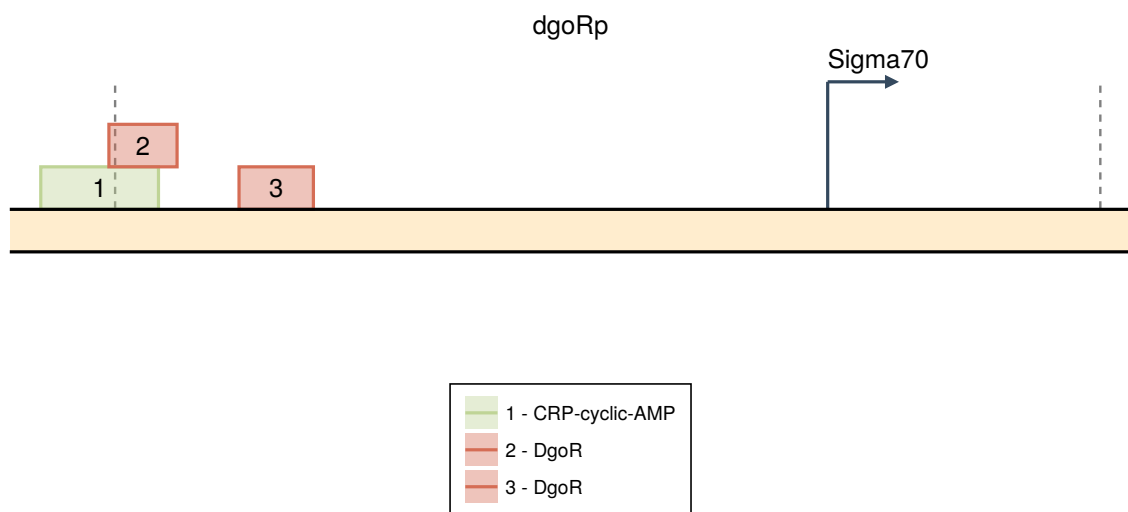


Figure S18

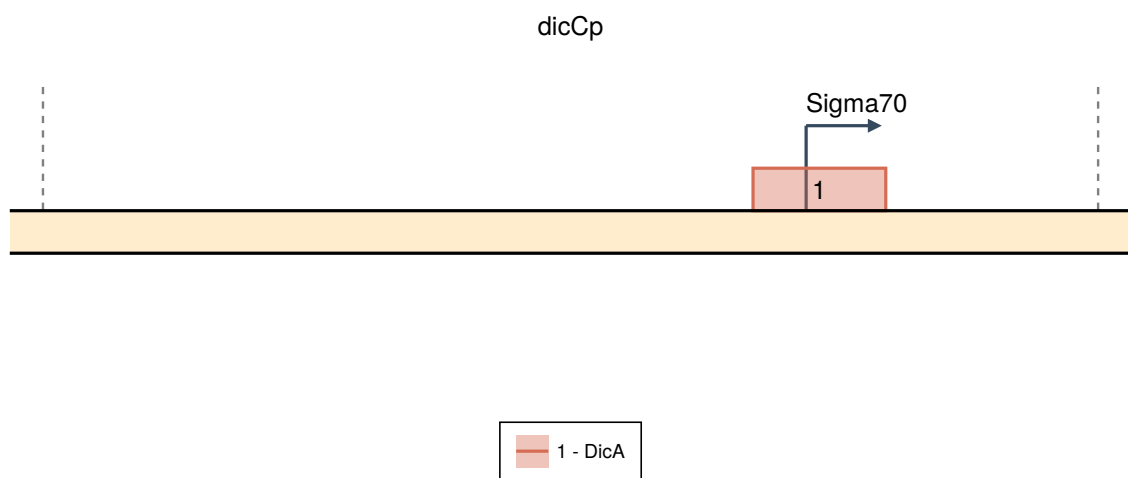


Figure S19

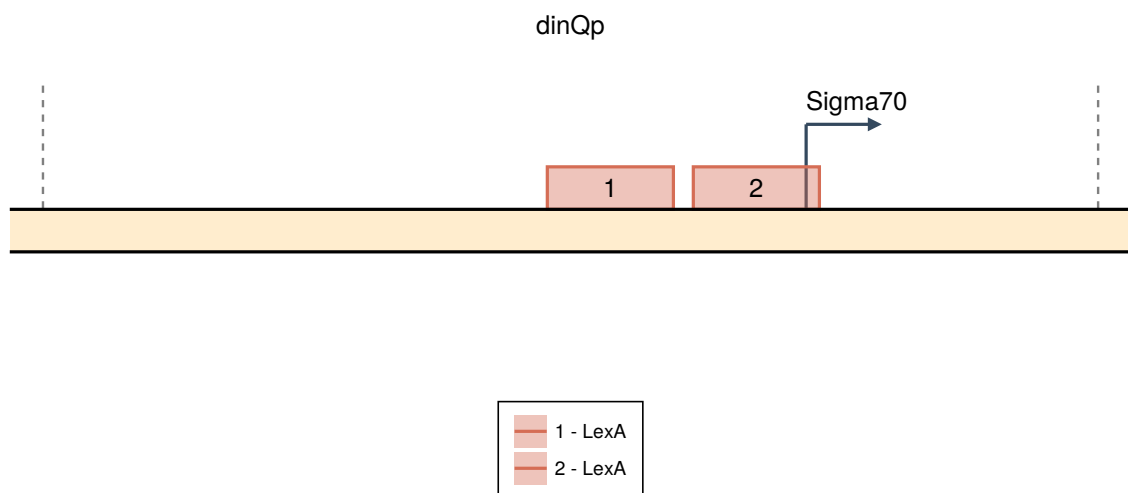


Figure S20

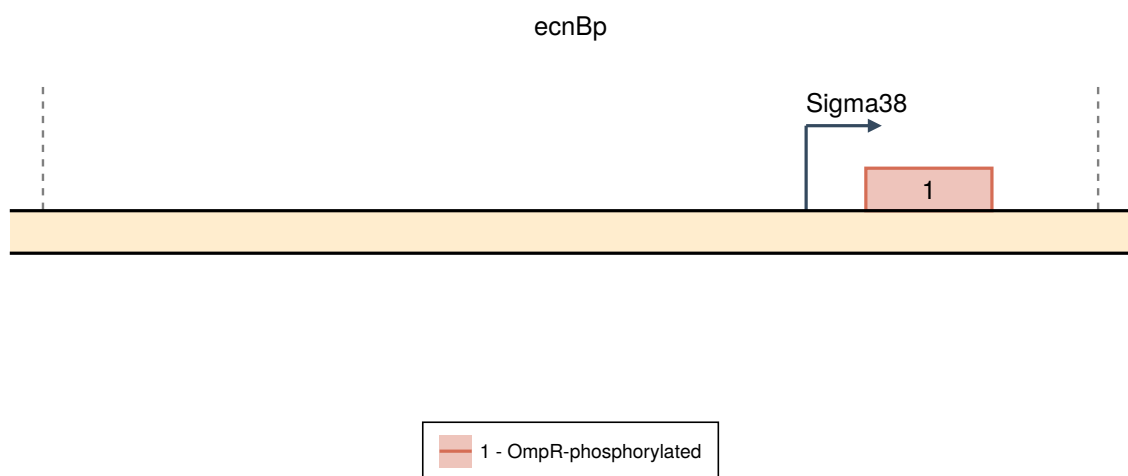


Figure S21

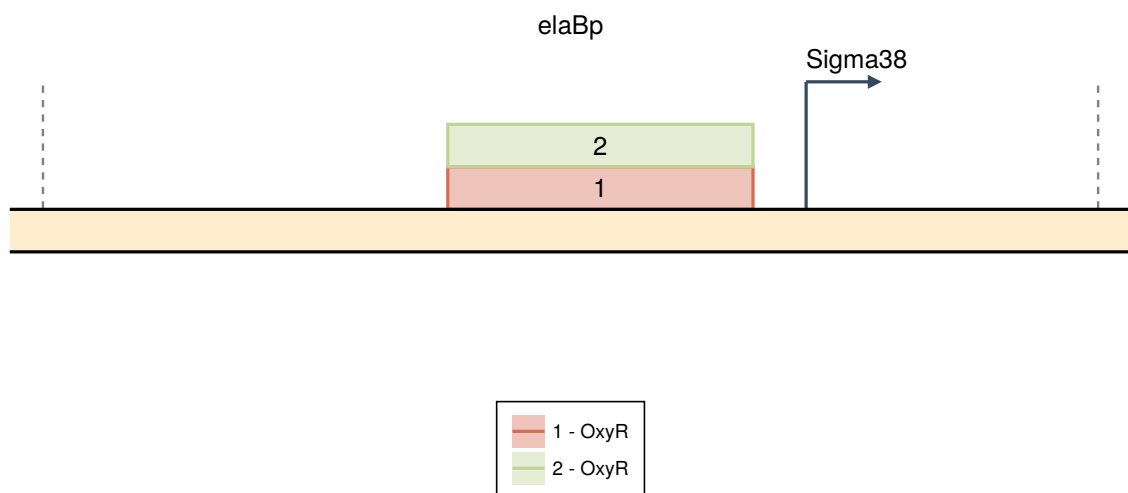


Figure S22

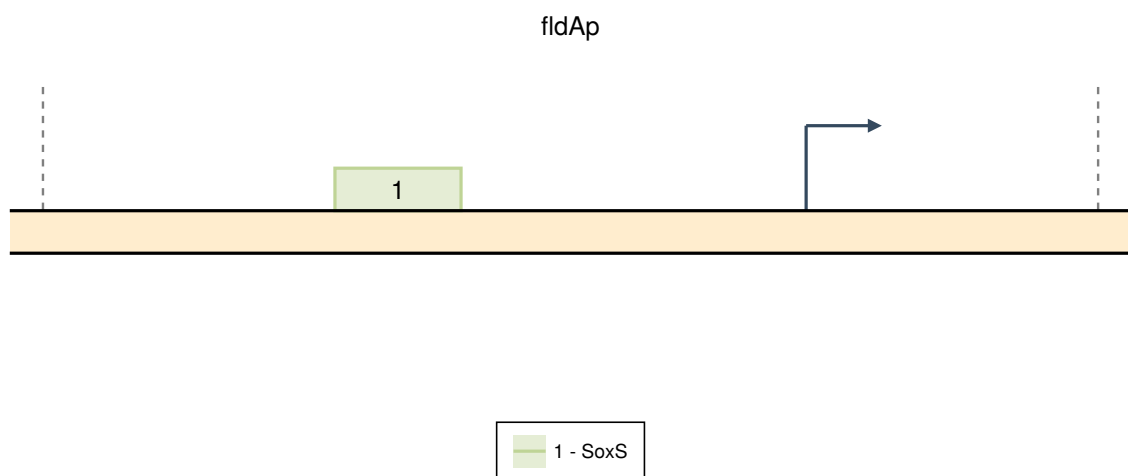


Figure S23

522	S10.9	crp
523	S10.9.1	crpp1
524	S10.9.2	crpp2
525	S10.9.3	crpp3
526	S10.10	cusC
527	S10.11	dgoR
528	S10.12	dicC
529	S10.13	dinQ
530	S10.14	ecnB
531	S10.15	elaB
532	S10.16	fldA
533	S10.17	ftsK
534	S10.17.1	ftsKp1

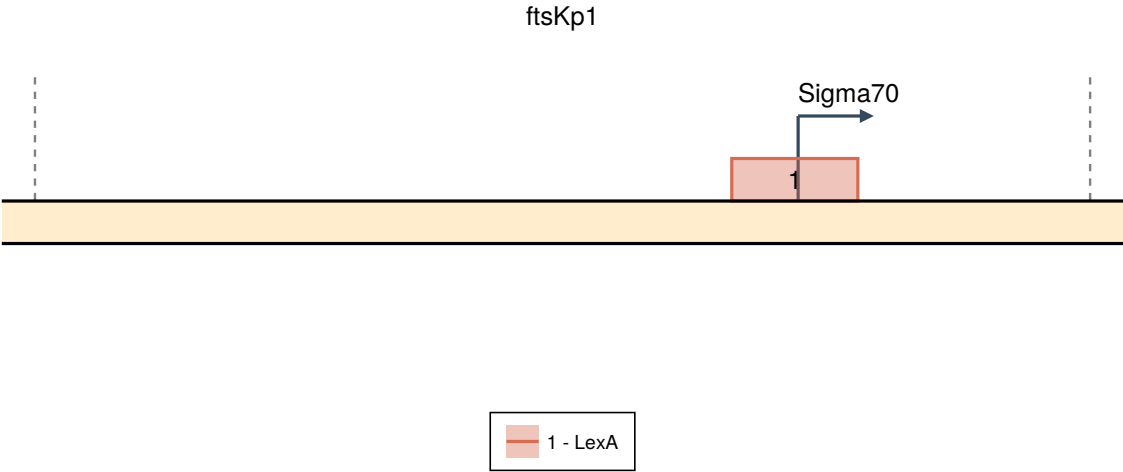


Figure S24

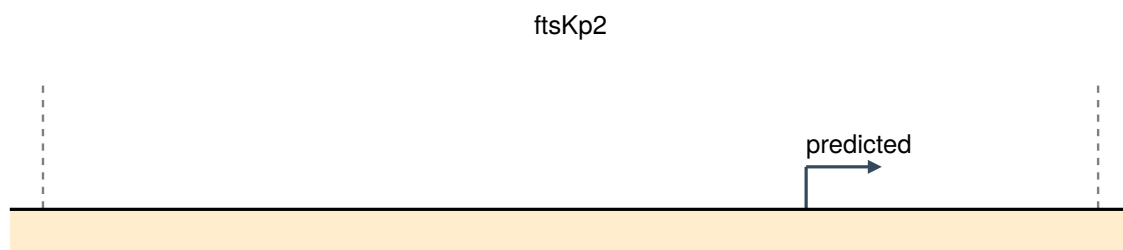


Figure S25

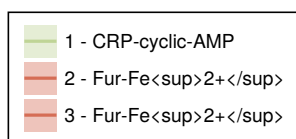


Figure S26

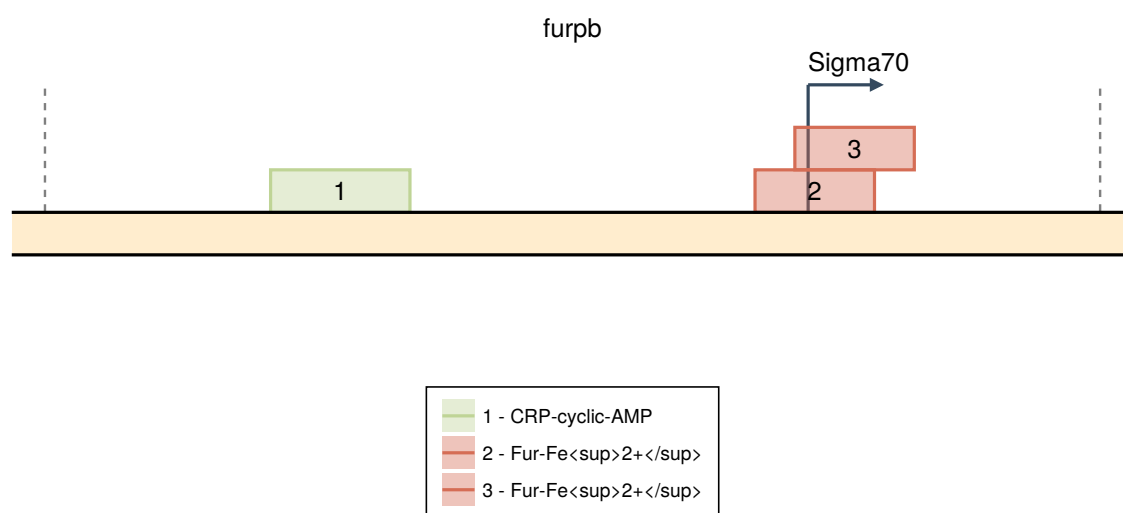


Figure S27

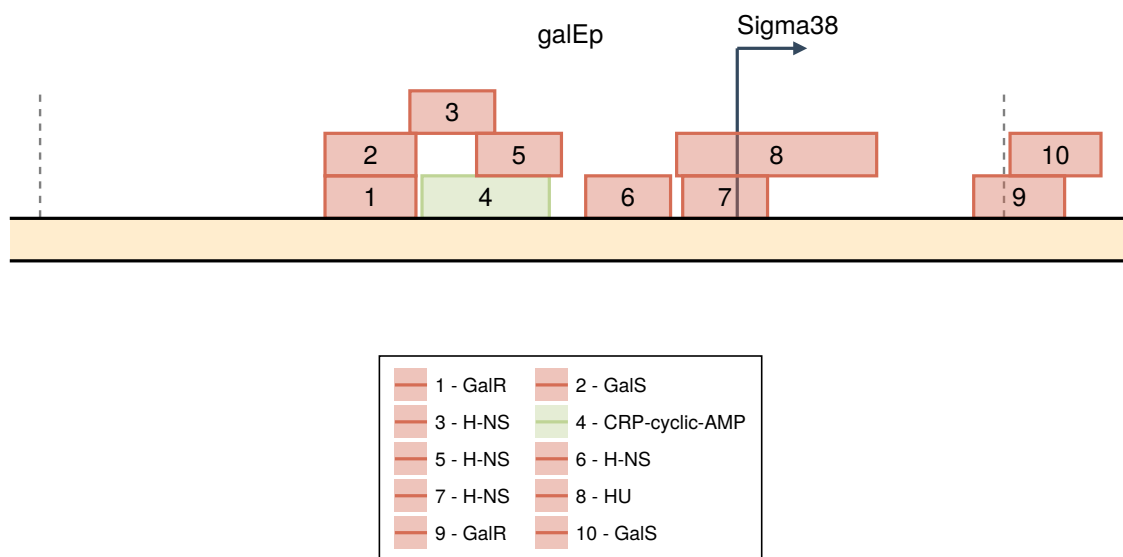


Figure S28

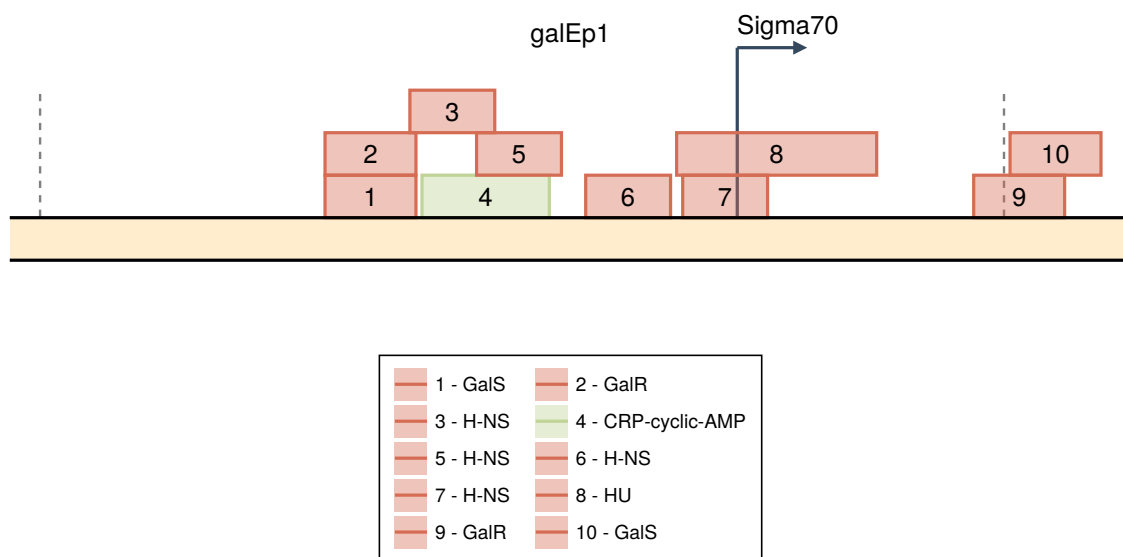


Figure S29

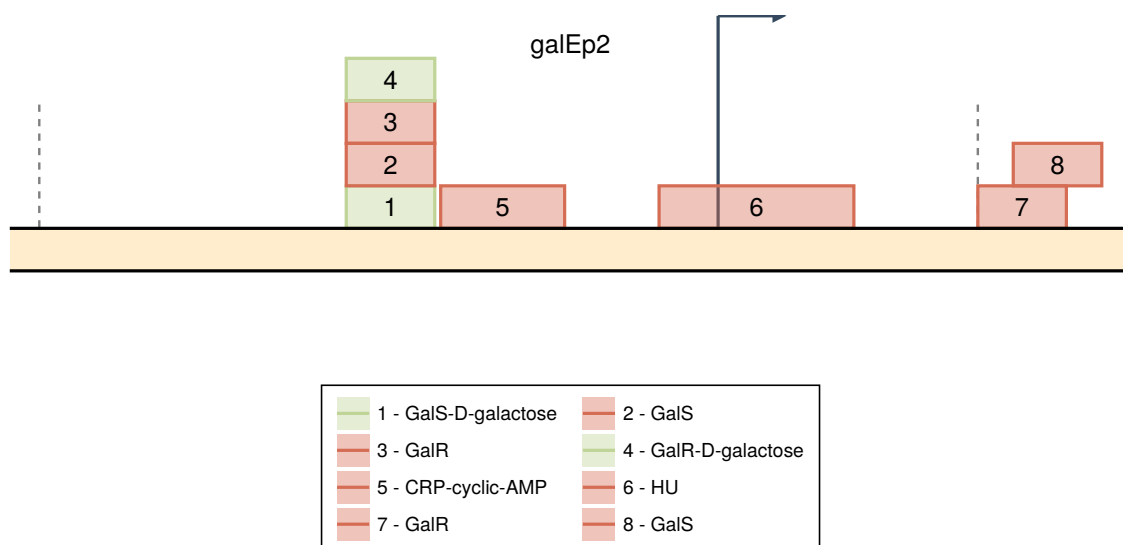


Figure S30



Figure S31

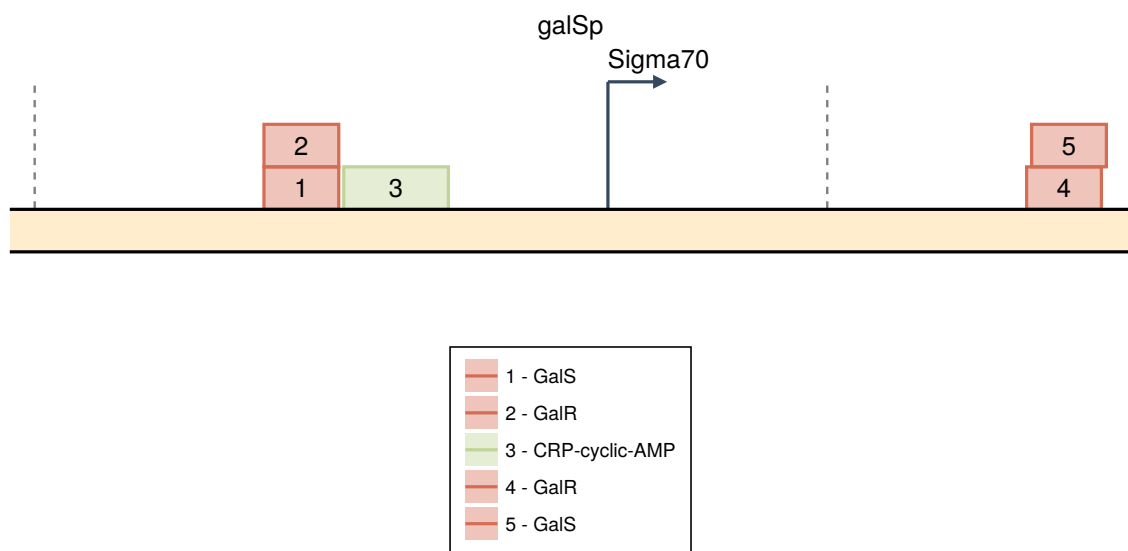


Figure S32

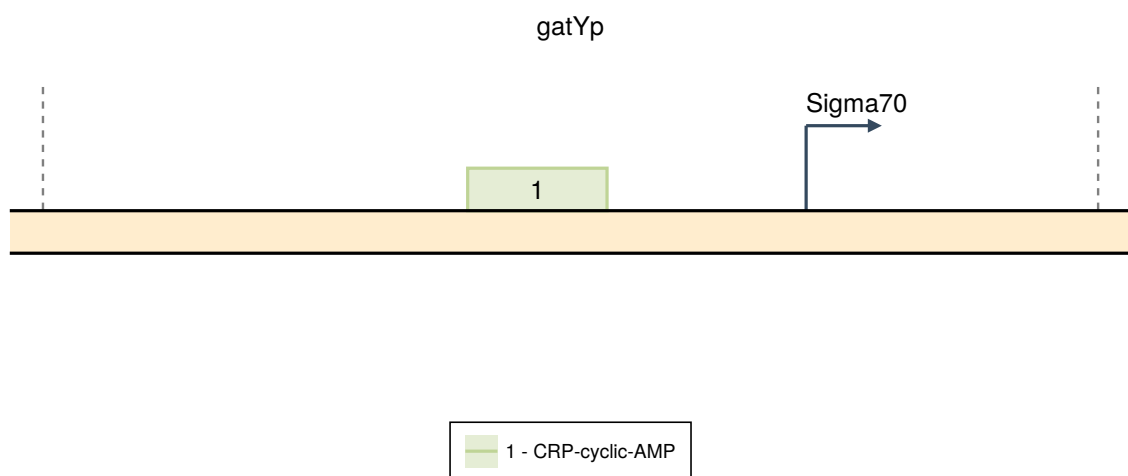


Figure S33



Figure S34

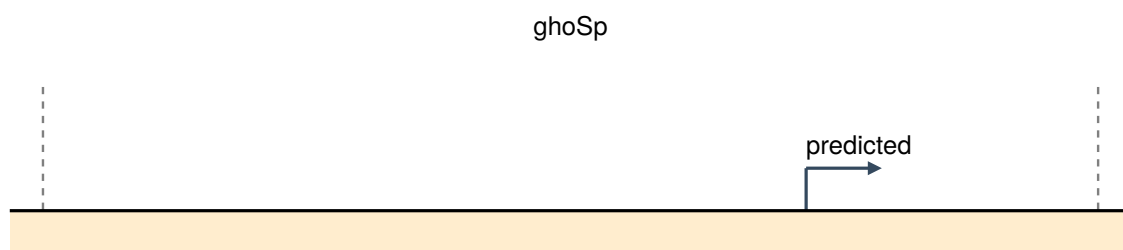


Figure S35

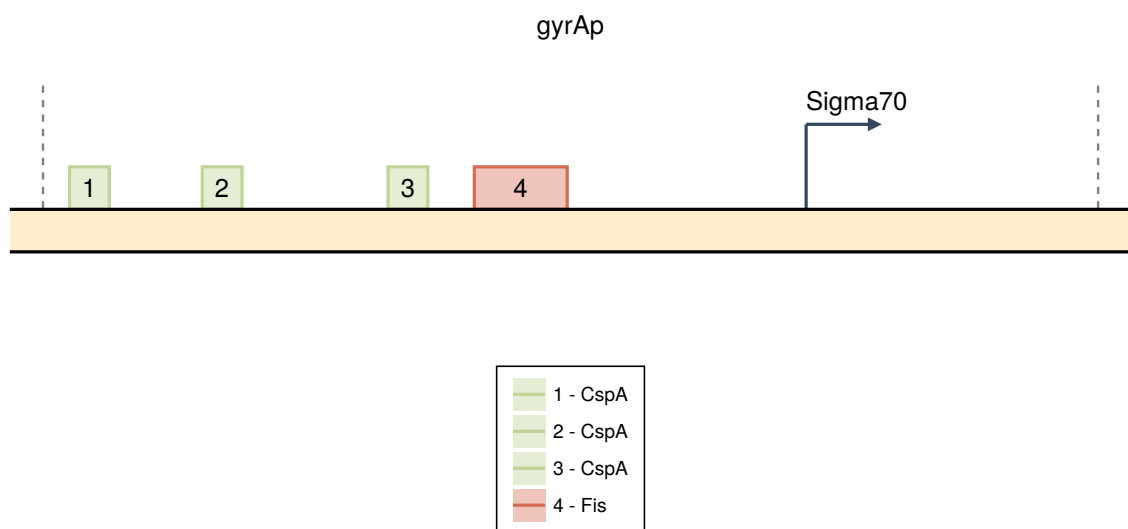


Figure S36

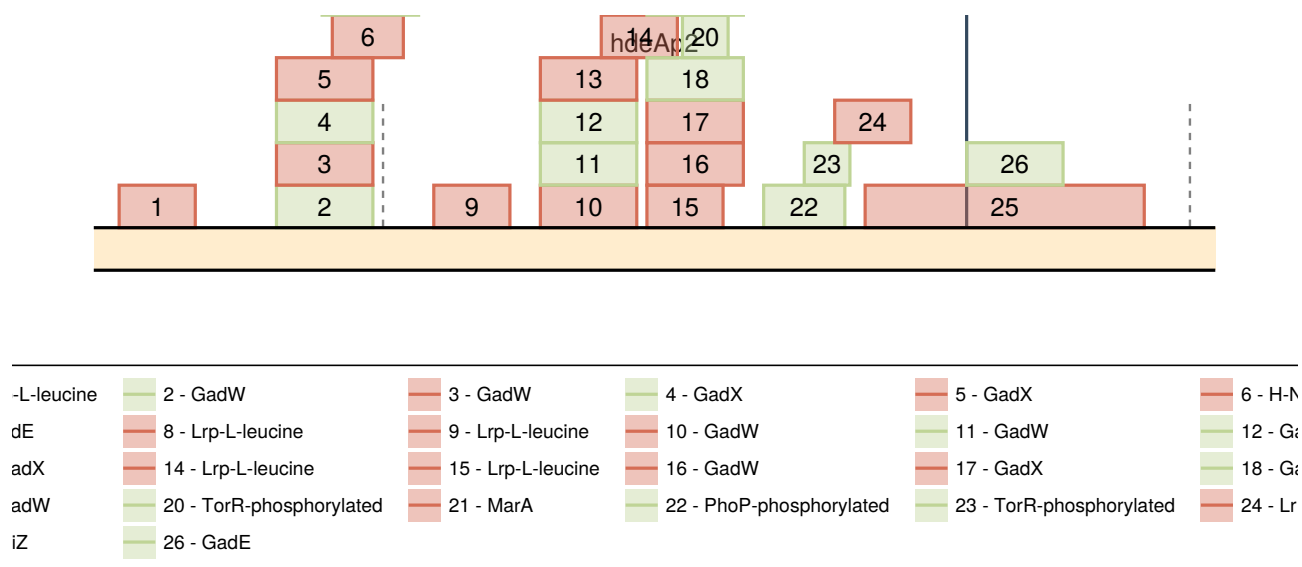


Figure S37

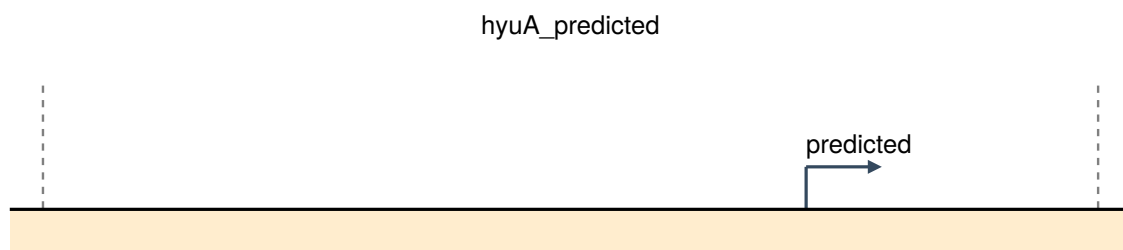


Figure S38

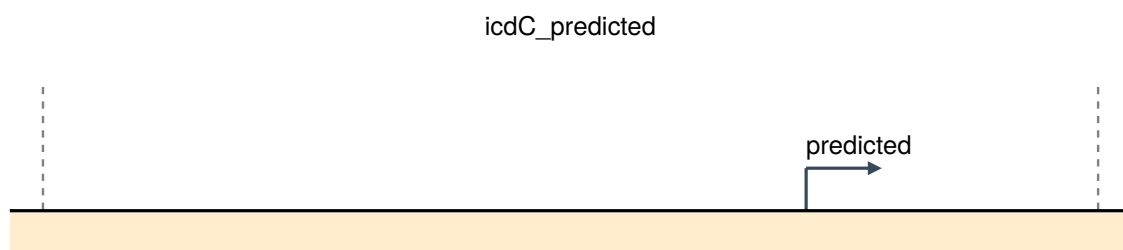


Figure S39

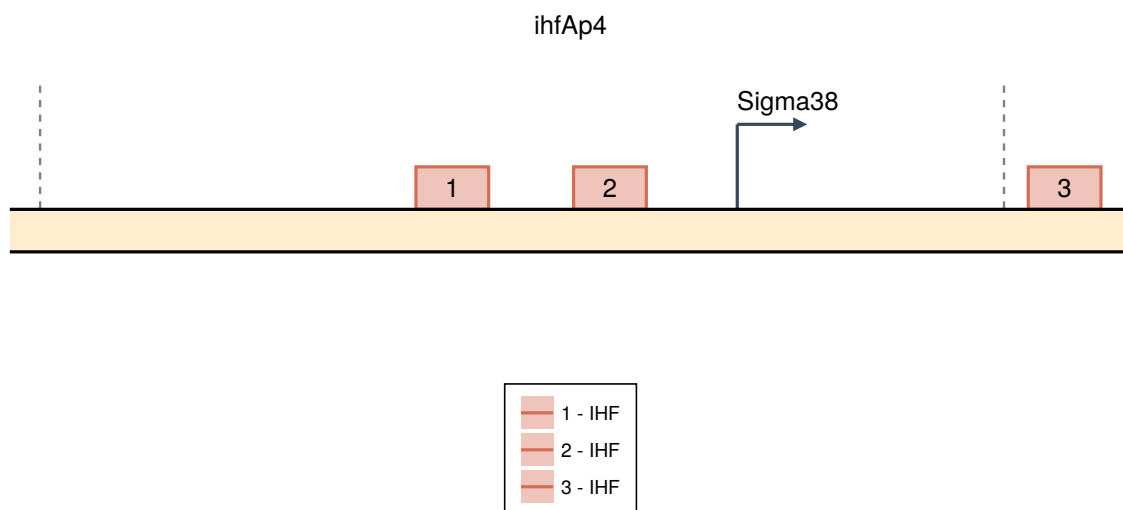


Figure S40

535 S10.17.2 ftsKp2

536 S10.18 fur

537 S10.18.1 furpa

538 S10.18.2 furpb

539 S10.19 galE

540 S10.19.1 galEp

541 S10.19.2 galEp1

542 S10.19.3 galEp2

543 S10.19.4 galEp3

544 S10.20 galS

545 S10.21 gatY

546 S10.22 gatZp

547 S10.23 ghoS

548 S10.24 gyrA

549 S10.25 hdeA

550 S10.26 hyuA

551 S10.27 icdC

552 S10.28 ihfA

553 S10.29 kbp

554 S10.30 lacI

555 S10.31 ldrD

556 S10.32 lpp

557 S10.33 lppp

558 S10.33.1 lppp2

559 S10.34 marR

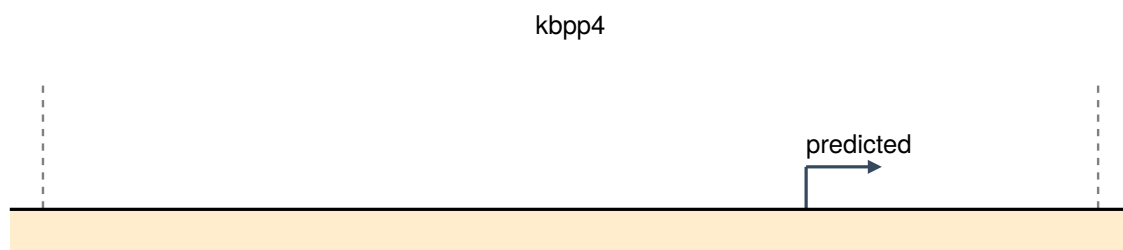


Figure S41

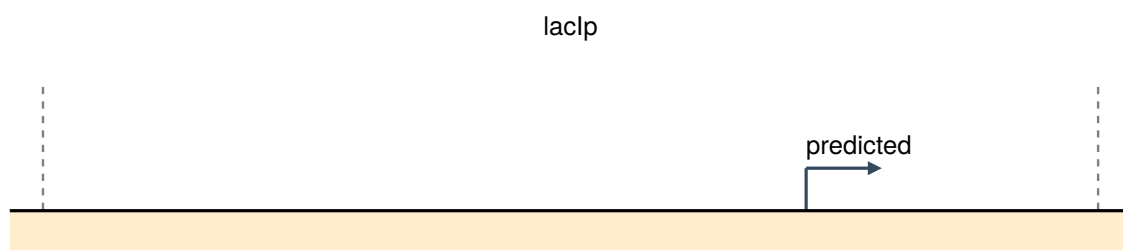


Figure S42

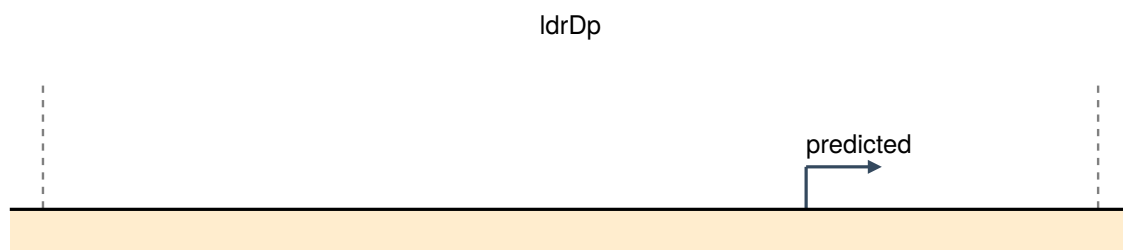


Figure S43

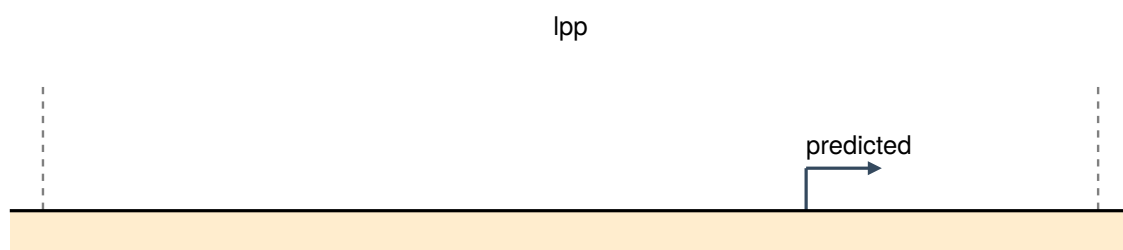


Figure S44

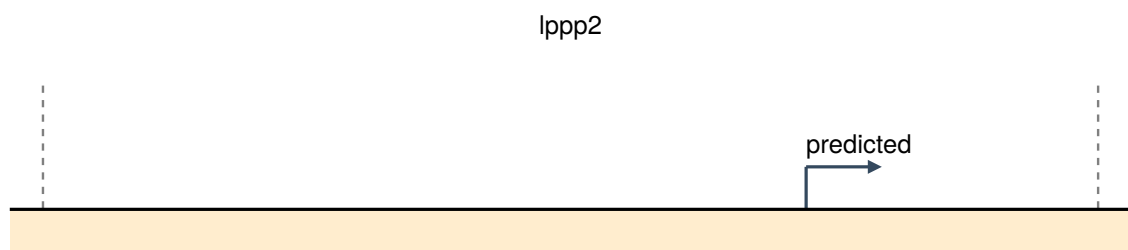


Figure S45

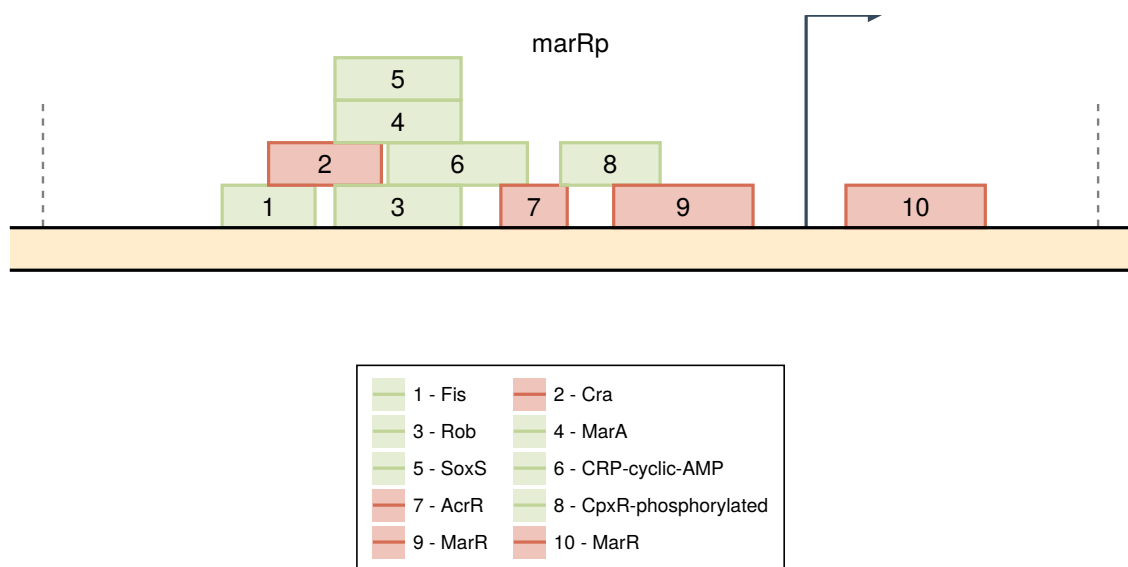


Figure S46

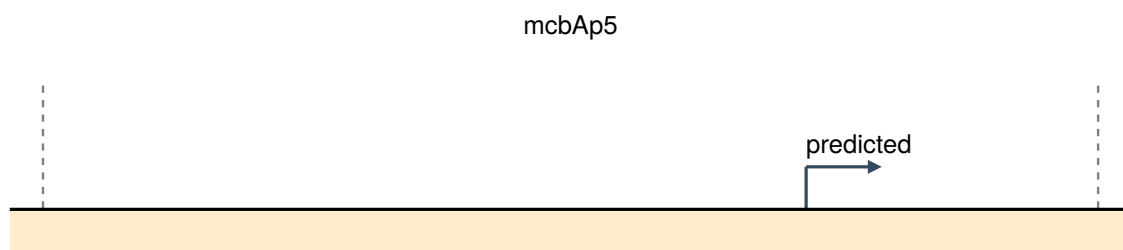


Figure S47

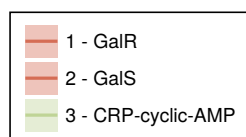
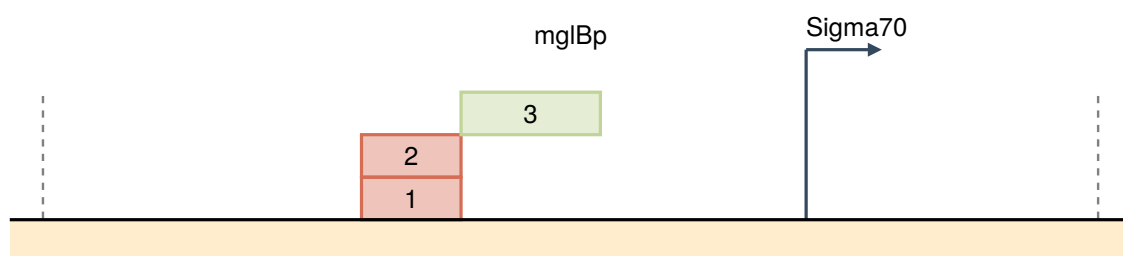


Figure S48

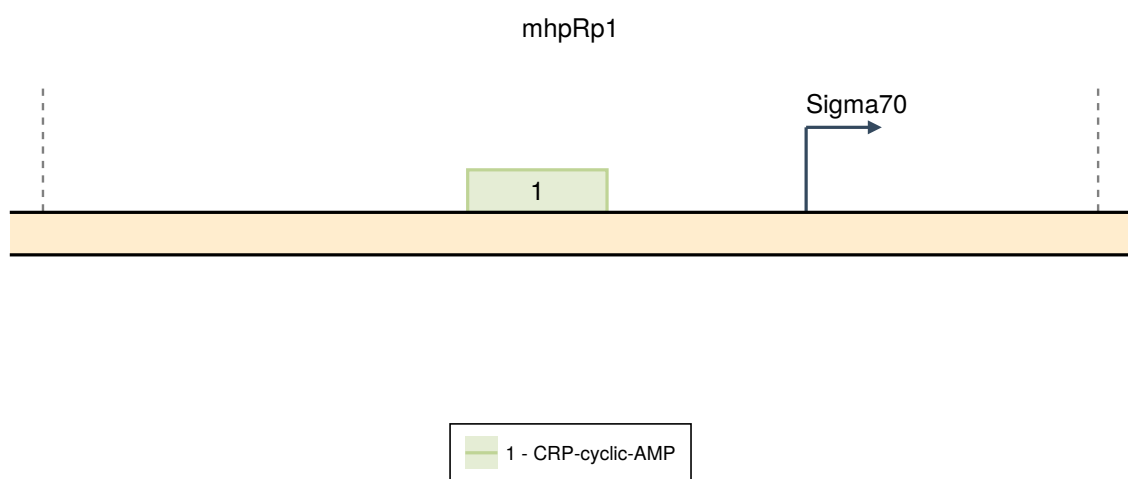


Figure S49

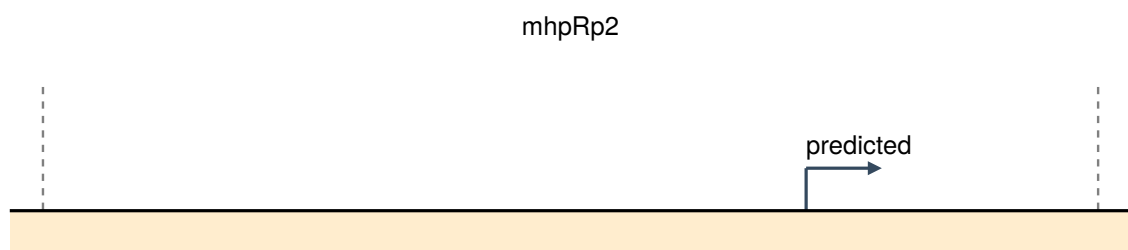


Figure S50

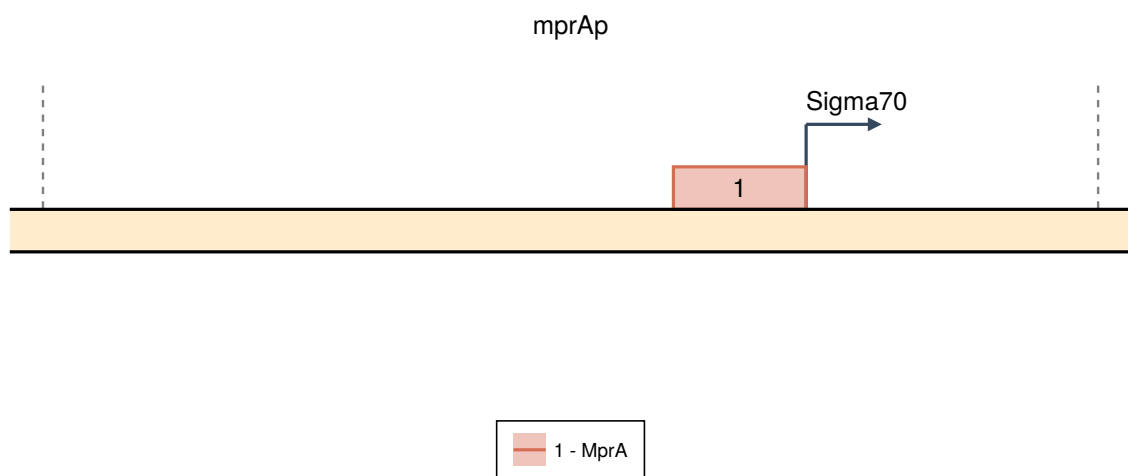


Figure S51

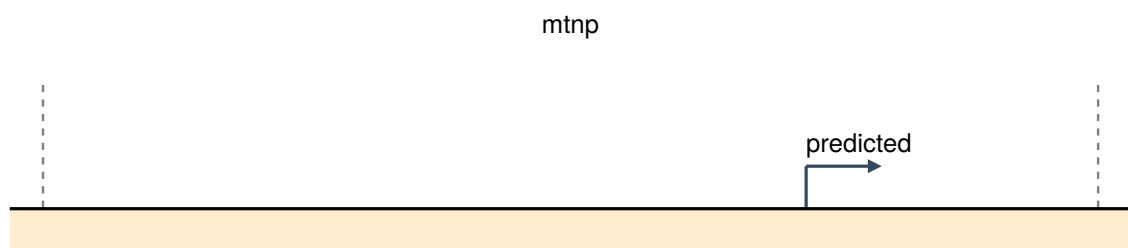


Figure S52

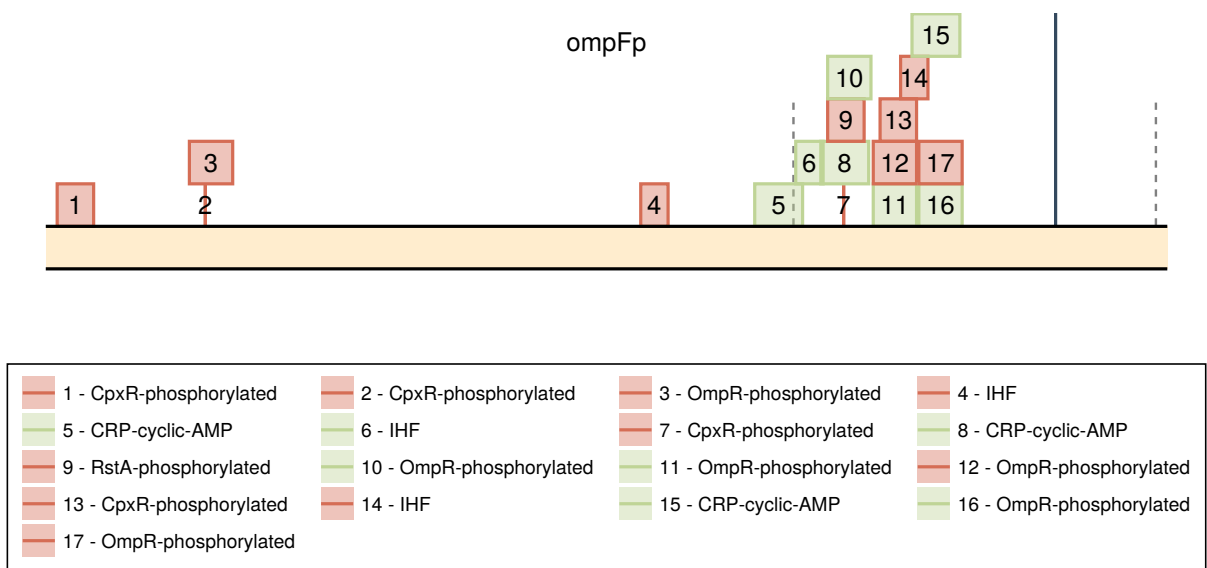


Figure S53

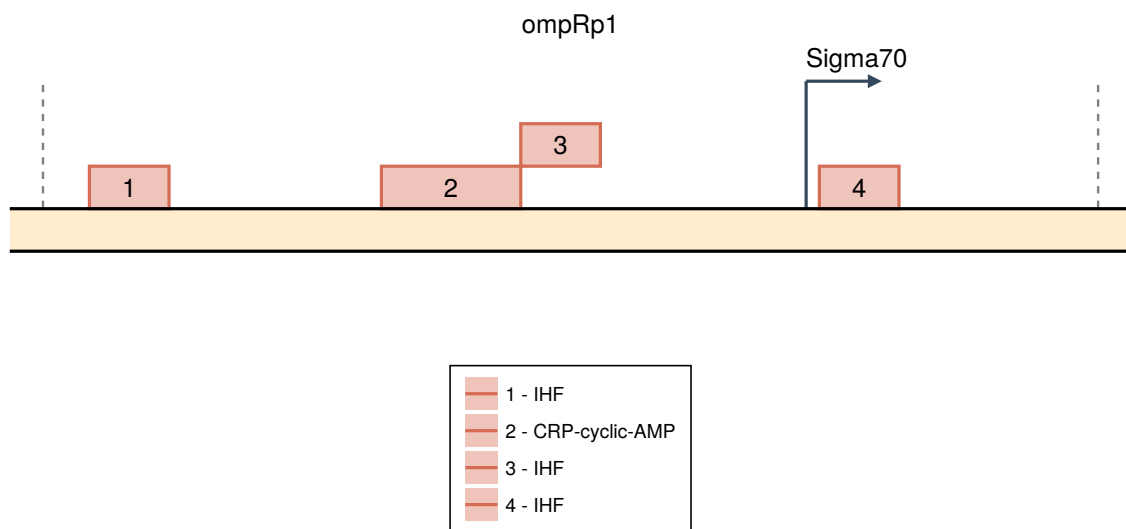


Figure S54

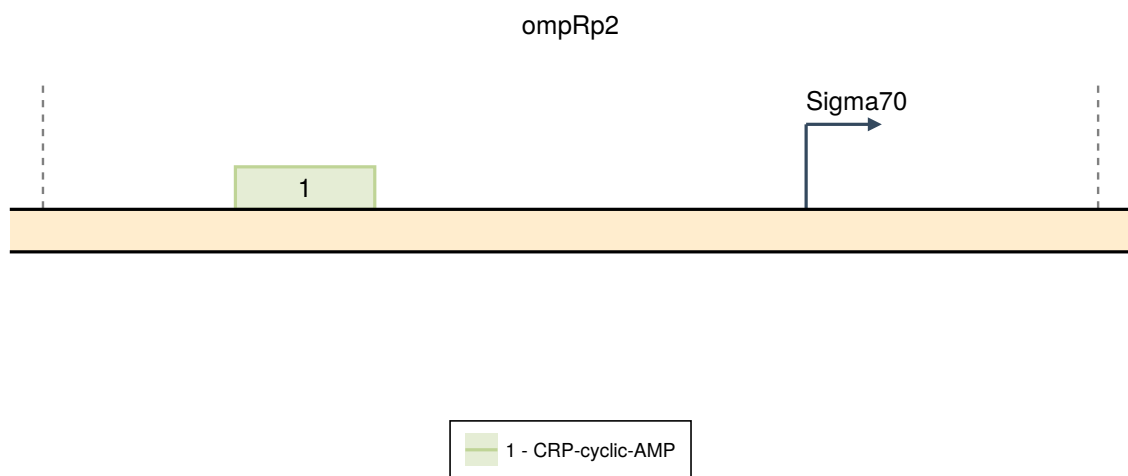


Figure S55

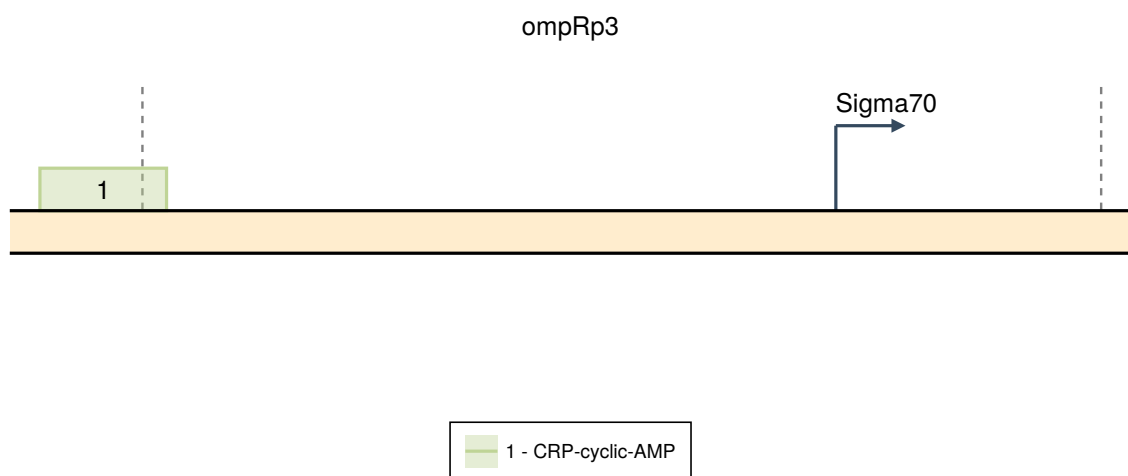


Figure S56

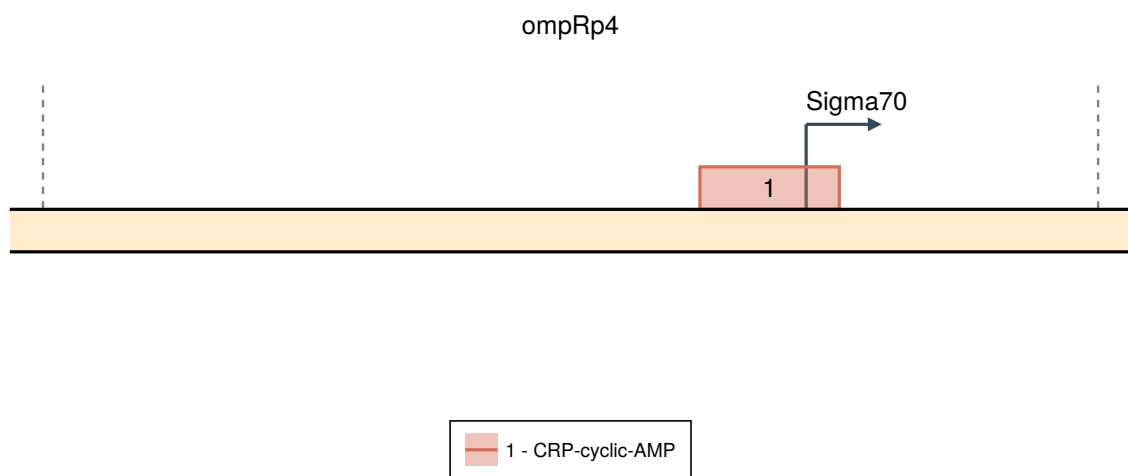


Figure S57

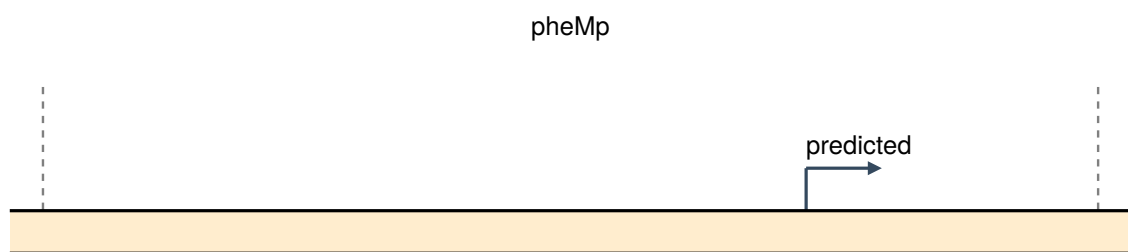


Figure S58

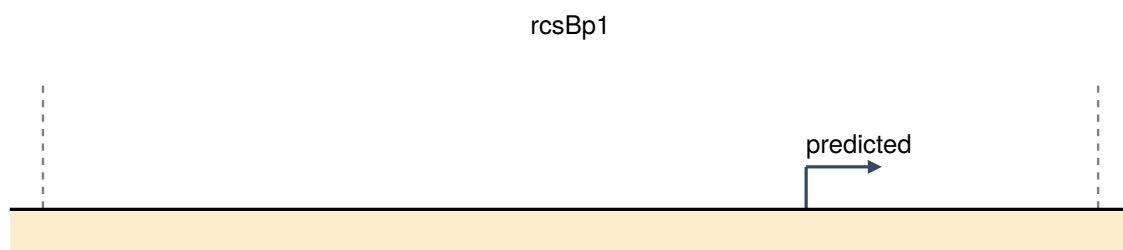


Figure S59

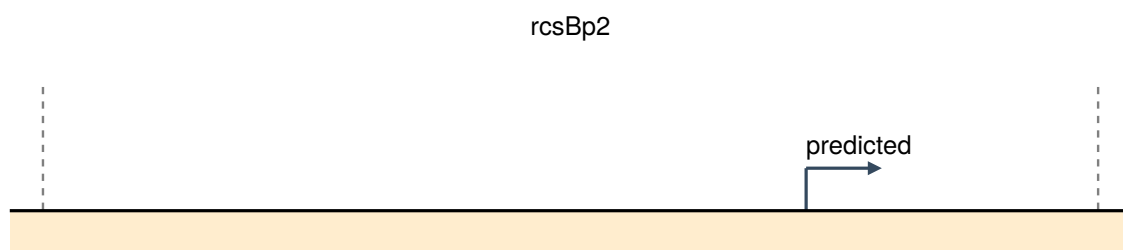


Figure S60

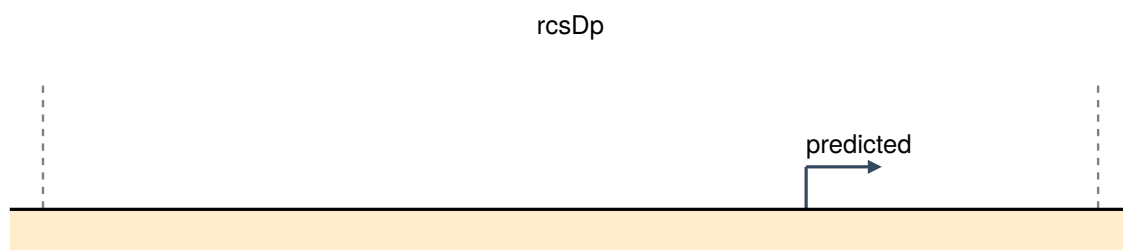


Figure S61

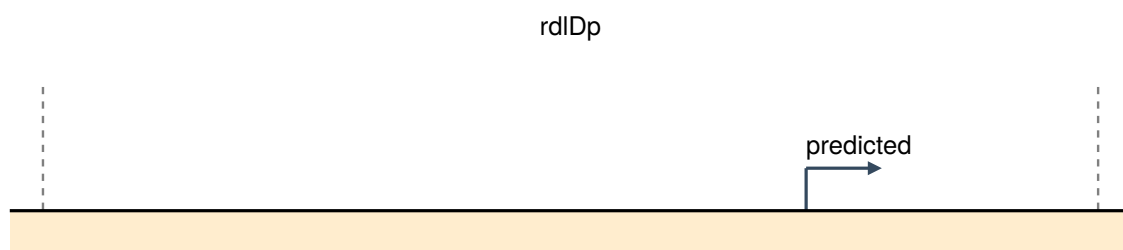


Figure S62

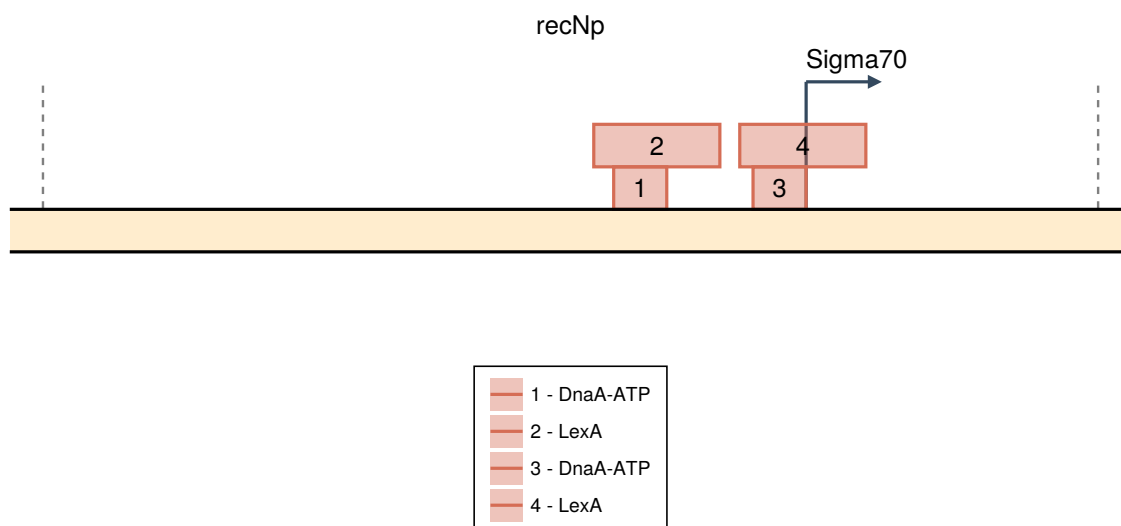


Figure S63

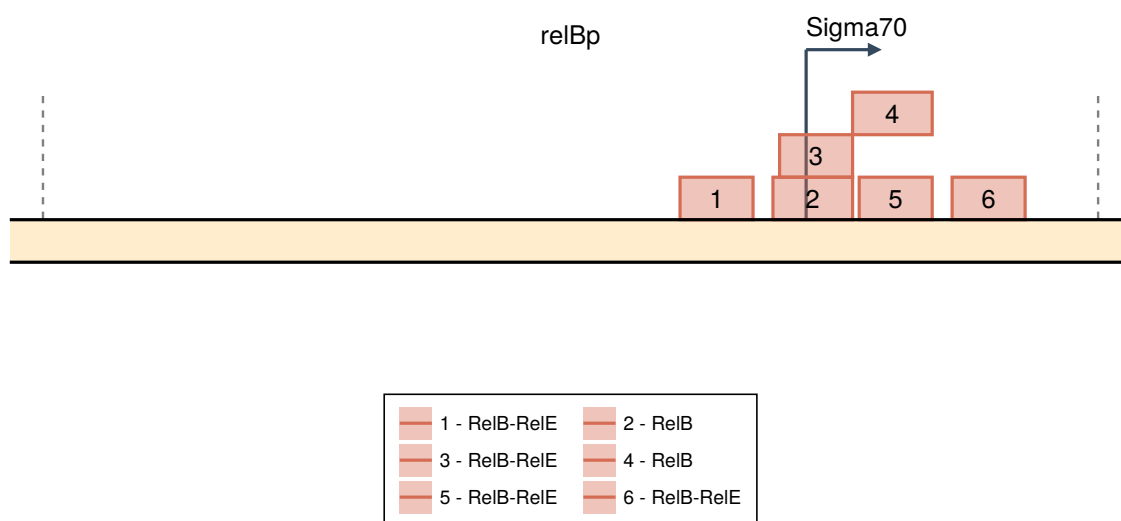


Figure S64

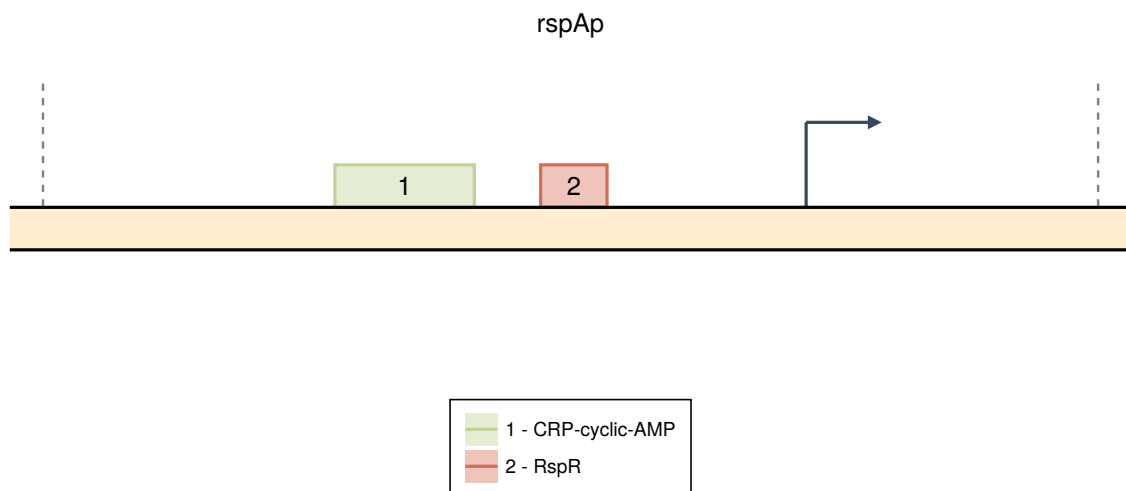


Figure S65

584 as substrates [42]. *rspA* has annotated binding sites for

585 **S10.49** *sohA*

586 **S10.50** *ssnA*

587 **S10.51** *sulA*

588 **S10.52** *tabA*

589 **S10.53** *tisB*

590 **S10.54** *tmaR*

591 **S10.55** *tnaC*

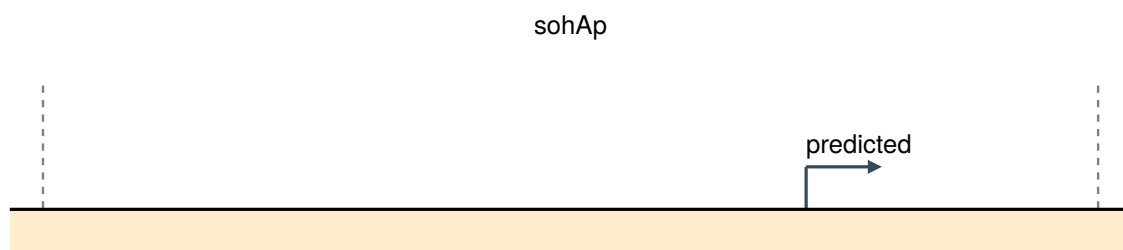


Figure S66

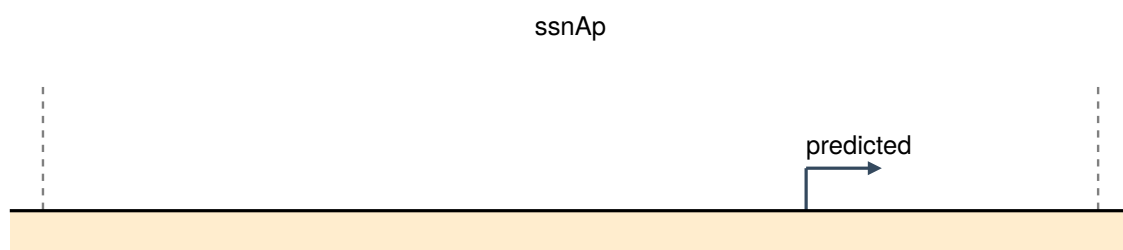


Figure S67

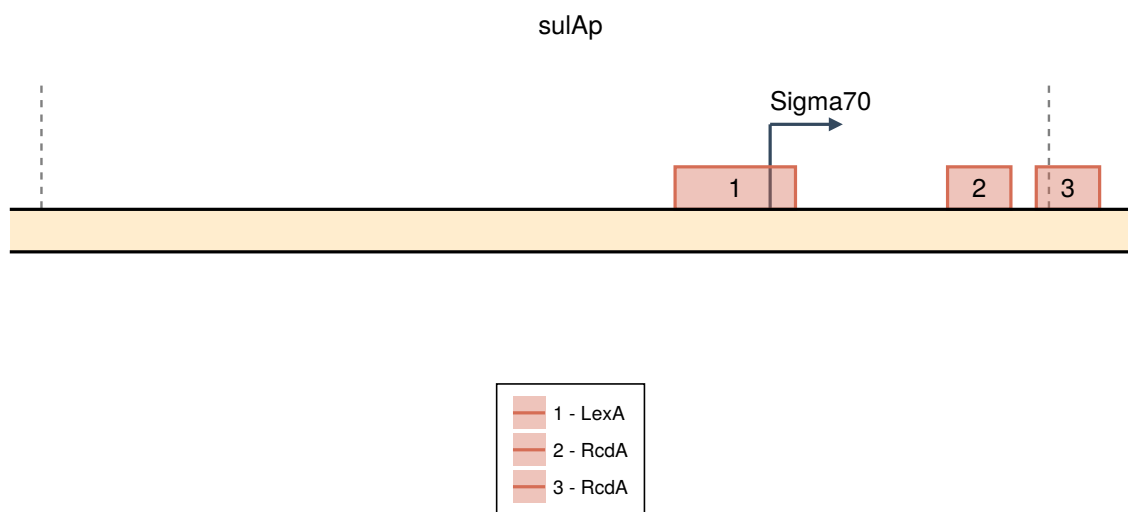


Figure S68

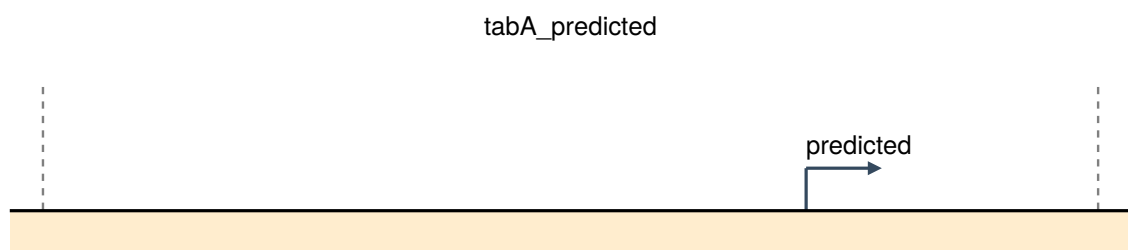


Figure S69

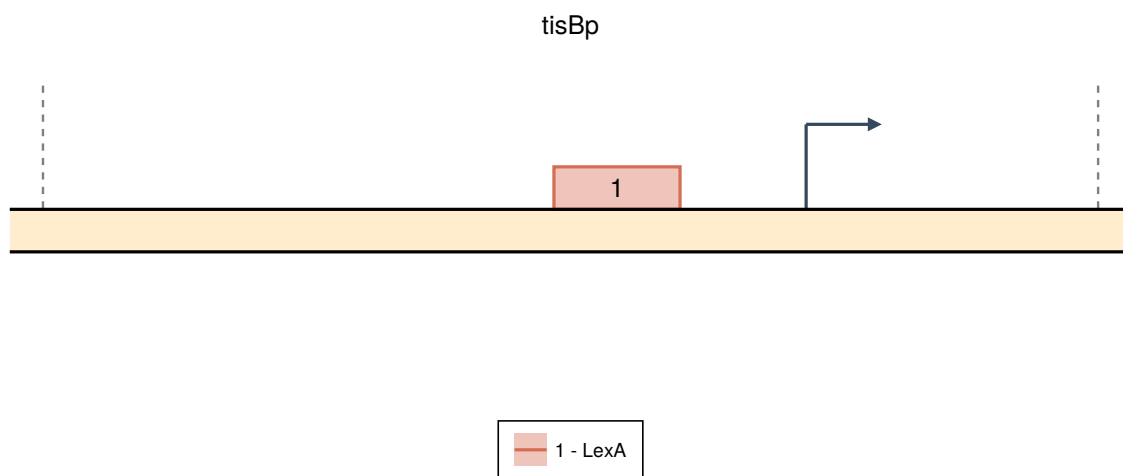


Figure S70

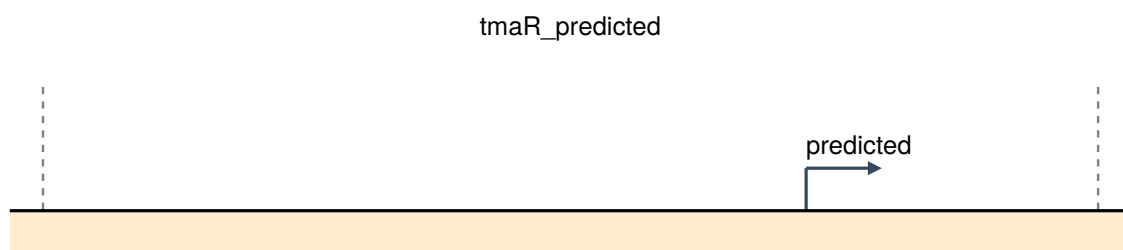


Figure S71

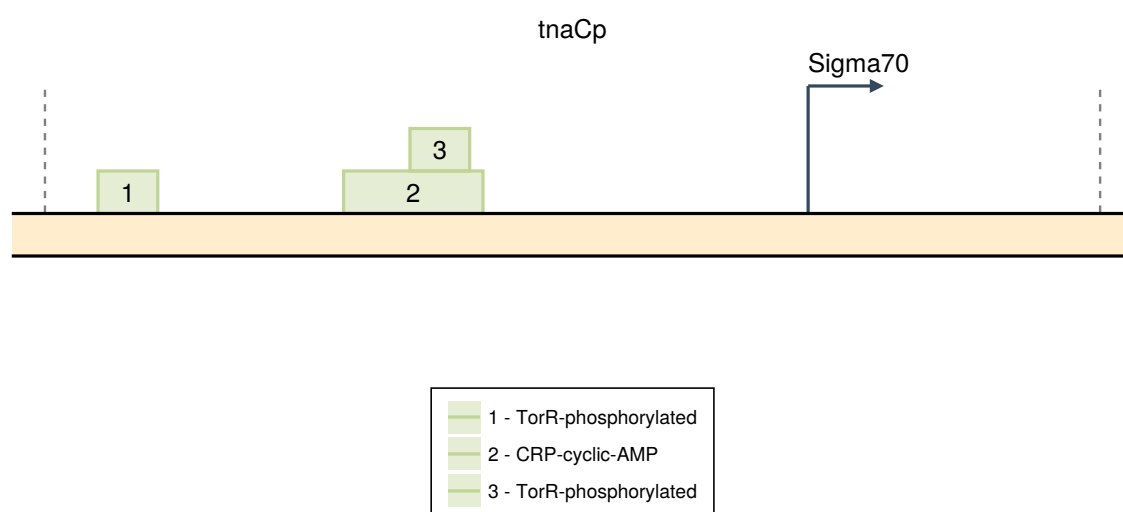


Figure S72

592 **S10.56** tolC

593 **S10.56.1** tolCp1

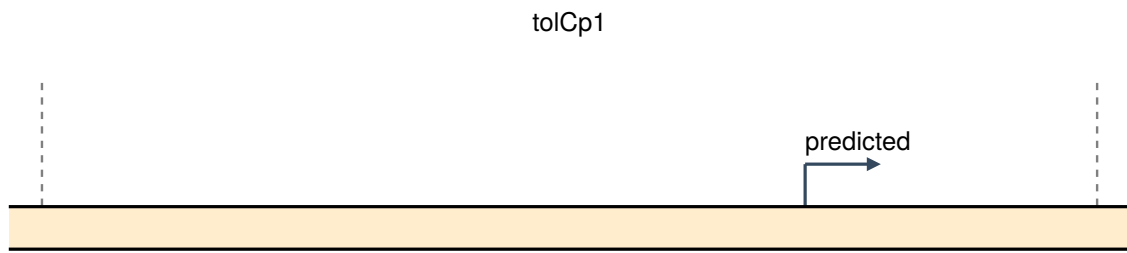


Figure S73

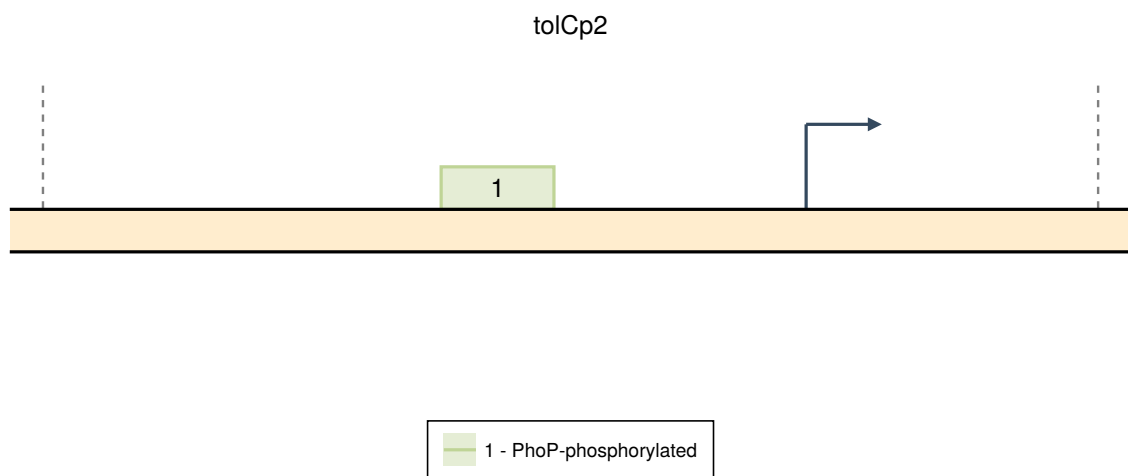


Figure S74

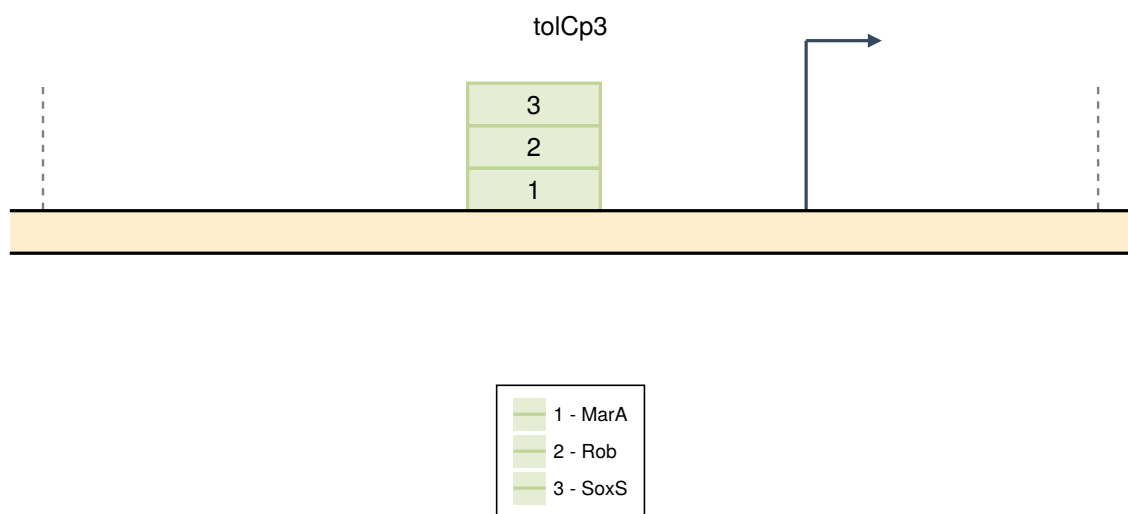


Figure S75

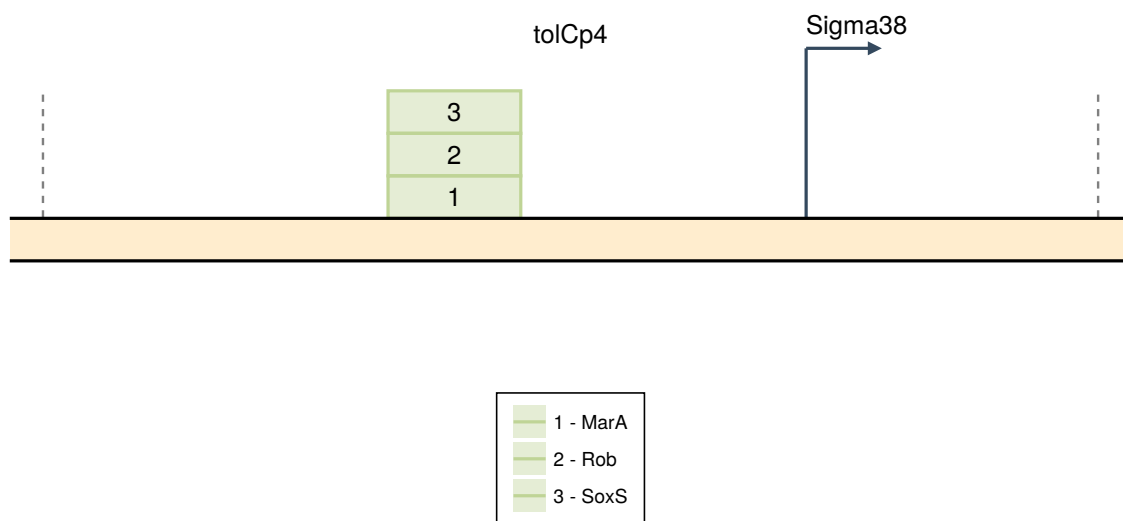


Figure S76

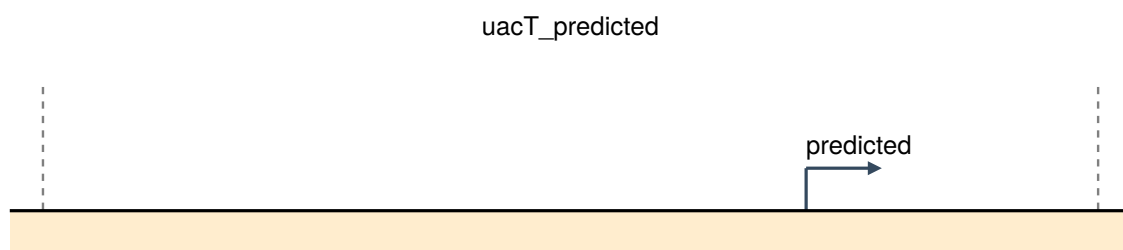


Figure S77

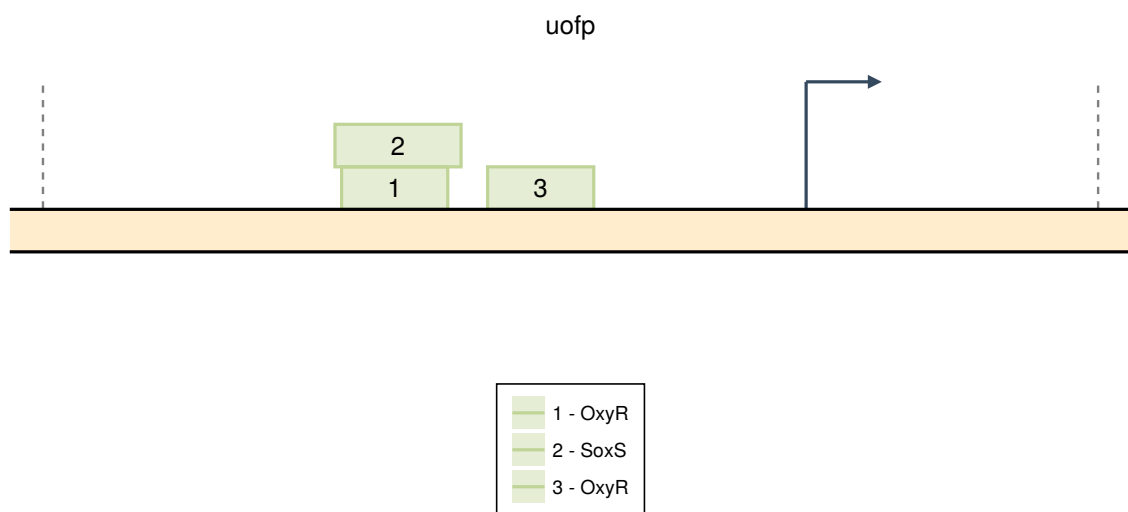


Figure S78



Figure S79



Figure S80

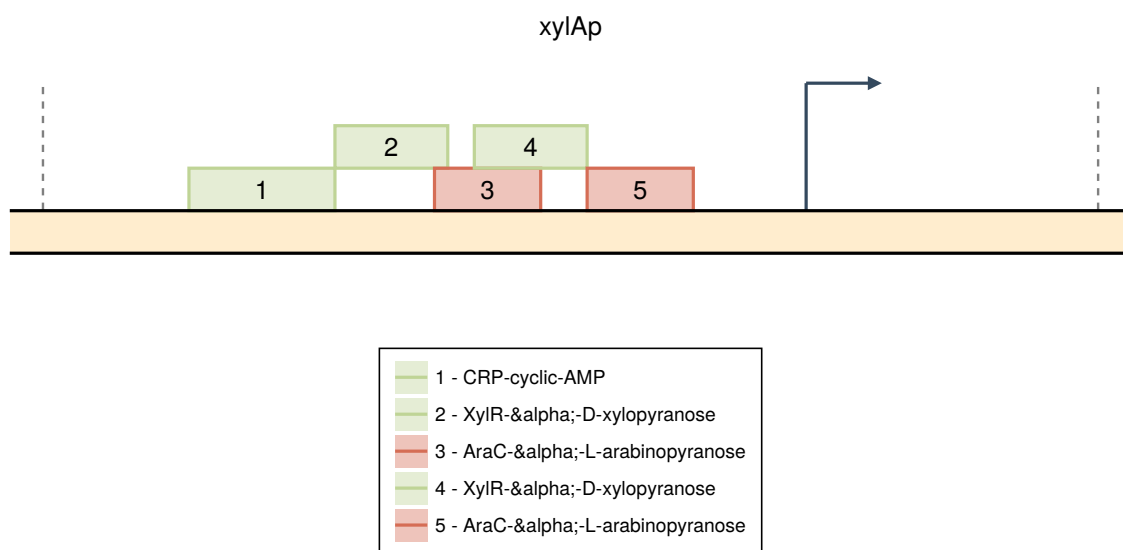


Figure S81

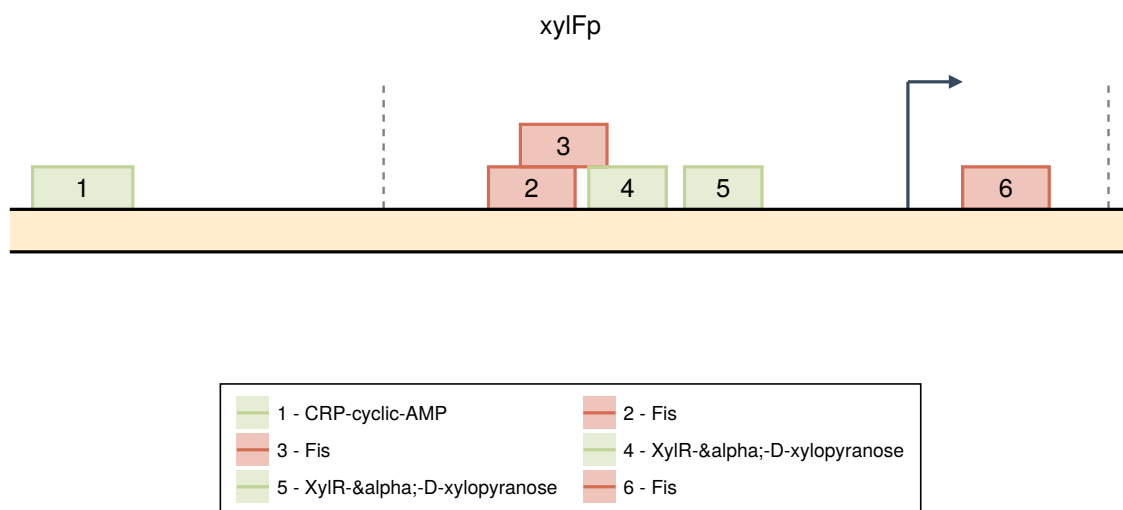


Figure S82

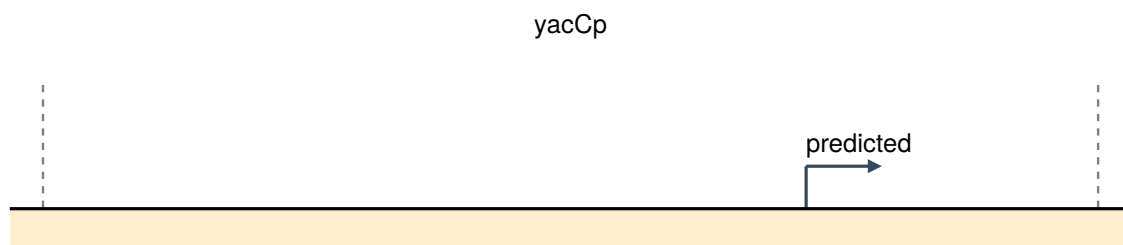


Figure S83

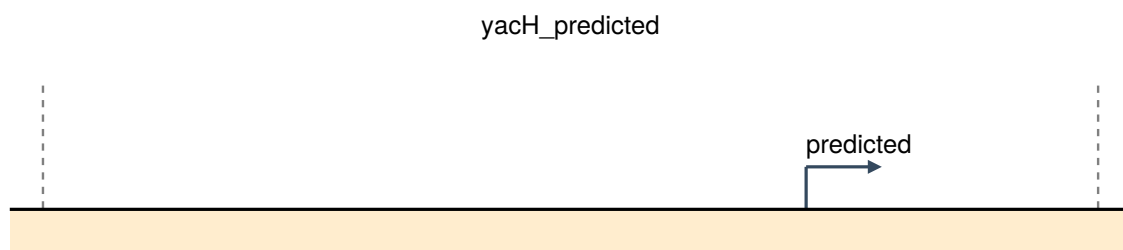


Figure S84

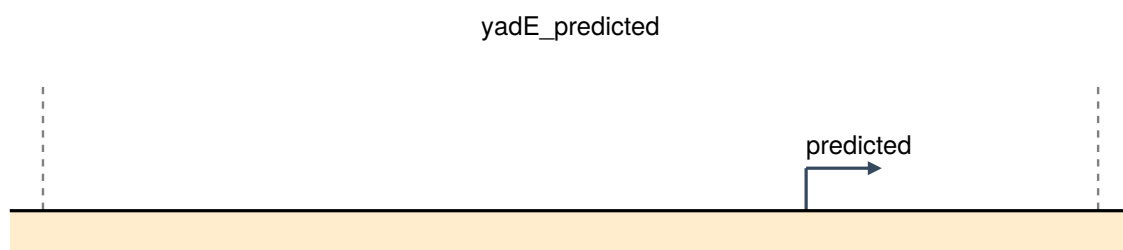


Figure S85

yadG_yadH_predicted

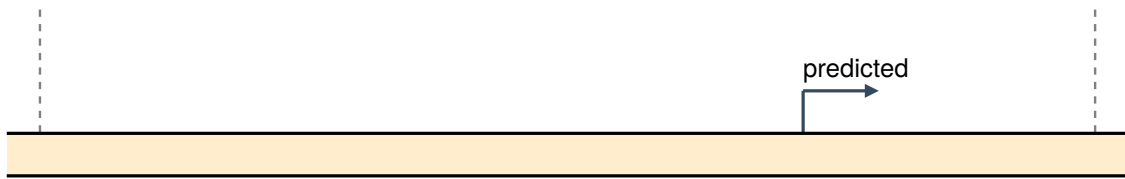


Figure S86

yadI_predicted



Figure S87

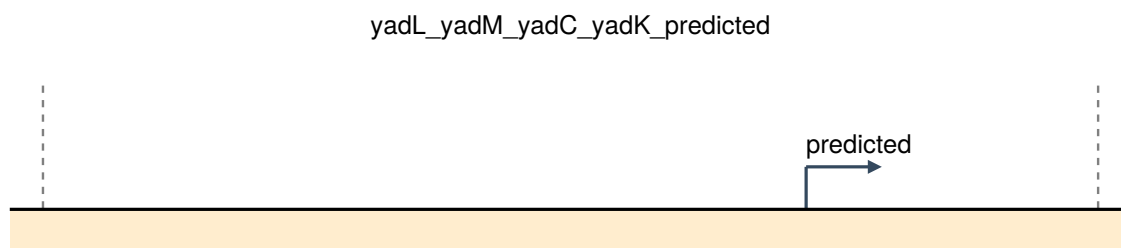


Figure S88

594	S10.56.2	tolCp2
595	S10.56.3	tolCp3
596	S10.56.4	tolCp4
597	S10.57	uacT
598	S10.58	uof
599	S10.59	xdhA
600	S10.59.1	xdhAp1
601	S10.59.2	xdhAp2
602	S10.60	xylA
603	S10.61	xylF
604	S10.62	yacC
605	S10.63	yacH
606	S10.64	yadE
607	S10.65	yadG-yadH
608	S10.66	yadI
609	S10.67	yadL-yadM-yadC-yadK
610	S10.68	yadN

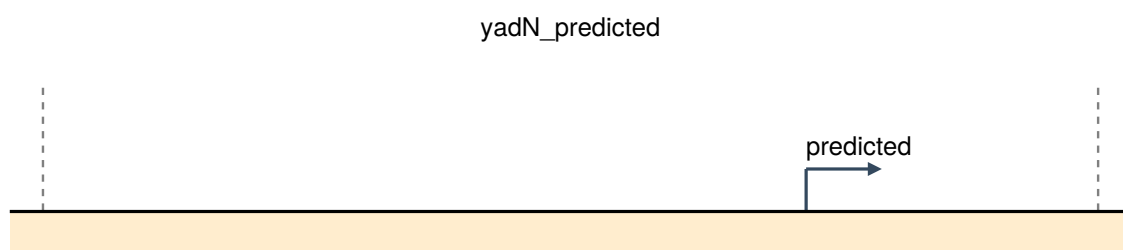


Figure S89

Supplemental References

- ³³I. M. Keseler et al., “Ecocyc: a comprehensive database of escherichia coli biology”, *Nucleic acids research* **39**, D583–D590 (2010).
- ³⁴M. Rydenfelt, R. S. Cox III, H. Garcia, and R. Phillips, “Statistical mechanical model of coupled transcription from multiple promoters due to transcription factor titration”, *Physical Review E* **89**, 012702 (2014).
- ³⁵S. K. Subramanian, W. P. Russ, and R. Ranganathan, “A set of experimentally validated, mutually orthogonal primers for combinatorially specifying genetic components”, *Synthetic Biology* **3**, ysx008 (2018).
- ³⁶B. Bushnell, *Bbmap: a fast, accurate, splice-aware aligner*, tech. rep. (Lawrence Berkeley National Lab.(LBNL), Berkeley, CA (United States), 2014).
- ³⁷S. Chen, Y. Zhou, Y. Chen, and J. Gu, “Fastp: an ultra-fast all-in-one fastq preprocessor”, *Bioinformatics* **34**, i884–i890 (2018).
- ³⁸O. Tange, *Gnu parallel 20230322 ('arrest warrant')*, GNU Parallel is a general parallelizer to run multiple serial command line programs in parallel without changing them., Mar. 2023.
- ³⁹A. Tareen, M. Kooshkbaghi, A. Posfai, W. T. Ireland, D. M. McCandlish, and J. B. Kinney, “Mave-nn: learning genotype-phenotype maps from multiplex assays of variant effect”, *Genome biology* **23**, 1–27 (2022).
- ⁴⁰K. Yamamoto and A. Ishihama, “Two different modes of transcription repression of the escherichia coli acetate operon by iclr”, *Molecular microbiology* **47**, 183–194 (2003).
- ⁴¹H. Huang, P. Zhou, J. Xie, et al., “Molecular mechanisms underlying the function diversity of transcriptional factor iclr family”, *Cellular Signalling* **24**, 1270–1275 (2012).
- ⁴²J. A. Gerlt, P. C. Babbitt, and I. Rayment, “Divergent evolution in the enolase superfamily: the interplay of mechanism and specificity”, *Archives of biochemistry and biophysics* **433**, 59–70 (2005).