

¹ Deciphering the regulatory genome ² of *Escherichia coli*, one hundred ³ promoters at a time

⁴ William T. Ireland¹, Suzannah M. Beeler², Emanuel Flores-Bautista², Nathan M.
⁵ Belliveau² †, Michael J. Sweredoski³, Annie Moradian³, Justin B. Kinney⁴, Rob
⁶ Phillips^{1,2,5*}

*For correspondence:

phillips@pboc.caltech.edu (RP)

Present address: [†]Howard Hughes Medical Institute and Department of Biology, University of Washington, Seattle, WA 98195

⁷ Department of Physics, California Institute of Technology, Pasadena, CA 91125; ²Division of Biology and Biological Engineering, California Institute of Technology, Pasadena, CA 91125; ³Proteome Exploration Laboratory, Beckman Institute, California Institute of Technology, Pasadena, CA 91125; ⁴Simons Center for Quantitative Biology, Cold Spring Harbor Laboratory, Cold Spring Harbor, NY 11724; ⁵Department of Applied Physics, California Institute of Technology, Pasadena, CA 91125

¹⁴ Abstract

¹⁵ While advances in DNA sequencing have enabled widespread access to this groundbreaking technology, even in the case of arguably one of biology's best understood organisms, the bacterium *Escherichia coli*, for ≈ 65% of the promoters we remain completely ignorant of their regulation. Until ¹⁶ we have cracked this regulatory Rosetta Stone, efforts to read and write genomes will remain ¹⁷ haphazard. We introduce a new method for unraveling the regulatory genome of a broad array of ¹⁸ genes in *E. coli*. Specifically, we show that we can extract previously known regulatory information ¹⁹ and unmask the regulatory architectures for more than 80 uncharacterized promoters, and in ²⁰ many of those cases, which transcription factors mediate this regulation. The method introduced ²¹ here clears a path for fully characterizing the regulatory genome of model organisms, with the ²² potential of moving on to an array of other microbes of ecological and medical relevance.

²⁶ Introduction

²⁷ DNA sequencing is as important to biology as the telescope was (and is) to astronomy. We are ²⁸ now living in the age of genomics, where the ability to sequence DNA is cheap and easy. However, ²⁹ despite these incredible advances, how all of this genomic information is regulated and deployed ³⁰ remains largely enigmatic. Organisms must respond to their environments through regulation of ³¹ genes. Genomic methods often provide a "parts" list but tell us very little about how organisms ³² use those parts creatively and constructively in space and time. But we know that promoters ³³ apply all-important dynamic logical operations that control when and where genetic information ³⁴ is accessed. In this paper, we show how we can infer the logical and regulatory interactions that ³⁵ control bacterial decision making by tapping into the power of DNA sequencing as a biophysical ³⁶ tool. The method introduced here provides a framework for solving the problem of deciphering ³⁷ the regulatory genome by connecting perturbation and response, mapping information flow from ³⁸ base pair in the promoter to downstream gene expression, quite literally determining how much ³⁹ information each promoter base pair carries about the level of gene expression.

⁴⁰

41 The advent of RNA-Seq (*Mortazavi et al., 2008*) launched a new era in which sequencing could
42 be used as an experimental read-out of the biophysically interesting counts of mRNA, rather than
43 simply as a tool for collecting ever more complete organismal genomes. The slew of 'X'-Seq tech-
44 nologies that are available continues to expand at a dizzying pace, each serving their own creative
45 and insightful role: RNA-Seq, ChIP-Seq, Tn-Seq, SELEX, 5C, etc. (*Stuart and Satija, 2019*). In contrast
46 to whole genome screening sequencing approaches, such as Tn-Seq (*Goodall et al., 2018*) and
47 ChIP-Seq (*Gao et al., 2018*) which give a coarse-grained view of gene essentiality and regulation
48 respectively, another class of experiments known as massively-parallel reporter assays (MPRA)
49 has been used to study gene expression in a variety of contexts (*Sharon et al., 2012; Patwardhan
50 et al., 2012; Fulco et al., 2019*). One elegant study relevant to the bacterial case of interest here
51 by *Kosuri et al. (2013)* screened more than 10^4 combinations of promoter and ribosome binding
52 sites (RBS). Even more recently, they have utilized this system to search for regulation across the
53 genome (*Urtecho et al., 2019, 2020*), in a way we as being complementary to our own. While their
54 approach allows for a course-grained view of where regulation may be occurring, our approach
55 yields a base-pair-by-base-pair view of how exactly that regulation is being enacted.

56

57 One of the most exciting recent X-Seq tools based on massively-parallel reporter assays with
58 broad biophysical reach is the Sort-Seq approach developed by *Kinney et al. (2010)*. The great
59 utility of this technique is that it not only reveals where transcription factors are binding (which
60 many other techniques are capable of as well), but also that this sequencing-as-read-out approach
61 provides a base pair resolution mapping of how sequence controls the level of gene expression. The
62 results of such a massively-parallel reporter assay make it possible to build a biophysical model of
63 gene regulation to uncover how previously uncharacterized promoters are regulated. In particular,
64 high-resolution studies like those described here yield quantitative predictions about promoter
65 organization as described by energy matrices (*Kinney et al., 2010*) that allow us to unleash the tools
66 of statistical physics to describe the input-output properties of each of these promoters which can
67 be explored much further with in-depth experimental dissection like those done by *Razo-Mejia
68 et al. (2018)* and *Chure et al. (2019)*. In this sense, the Sort-Seq approach can provide a quantitative
69 framework to not only discover and quantitatively dissect regulatory interactions at the promoter
70 level, but provides an interpretable scheme to design genetic circuits with a desired expression
71 output (*Barnes et al., 2019*).

72

73 Earlier work from *Belliveau et al. (2018)* illustrated how Sort-Seq, in conjunction with mass
74 spectrometry to figure out what transcription factors are binding the putative binding sites, could
75 be used to dissect a handful of these promoters for which we were previously mired in regulatory
76 ignorance. However, a crucial drawback of the approach of *Belliveau et al. (2018)* is that while it is
77 high-throughput at the level of a single gene and the number of promoter variants it accesses, it was
78 unable to readily tackle multiple genes at once, still leaving much of the unannotated genome un-
79 touched. Given that even in one of biology's best understood organisms, the bacterium *Escherichia
80 coli*, for more than 65% of its genes, we remain completely ignorant of how those genes are regu-
81 lated (*Santos-Zavaleta et al., 2019; Belliveau et al., 2018*), a higher-throughput approach is needed.
82 If we hope to some day have a complete base pair resolution mapping of how genetic sequences re-
83 late to biological function, we must first be able to do so for the promoters of this "simple" organism.

84

85 What has been missing in uncovering the regulatory genome in organisms of all kinds is a large
86 scale method for inferring genomic logic and regulation. How is the problem solved? Here we
87 replace the low-throughput fluorescence-based Sort-Seq approach with a scalable RNA-Seq based
88 approach that makes it possible to attack multiple promoters at once, setting the stage for the
89 possibility of, to first approximation, uncovering the entirety of the regulatory genome. Accordingly,
90 we refer to this approach as Reg-Seq, which we unleash here on over one hundred promoters. The
91 Reg-Seq method has some key ideas that are critical for making the link between DNA sequence

92 and regulatory outcomes. The concept of the method is to perturb promoter regions by mutating
93 them and then using sequencing to read out both perturbation and gene expression. We generate a
94 broad diversity of promoter sequences for each promoter of interest experimentally and use mutual
95 information as a metric to measure information flow from that distribution of sequences to gene
96 expression. By doing this, Reg-Seq is able to collect causal information about candidate regulatory
97 sequences that is then complemented by mass spectrometry which allows us to find which transcrip-
98 tion factors mediate the action of those newly discovered candidate regulatory sequences. Hence,
99 Reg-Seq solves the causal problem of linking DNA sequence to regulatory logic and information flow.

100

101 To demonstrate the ability to scale up the power of the Sort-Seq method, we report here the re-
102 sults of applying our Reg-Seq approach to 113 *E. coli* genes, whose promoter architectures we were
103 able to elucidate in parallel. By taking the Sort-Seq approach from a gene-by-gene method to a more
104 whole-genome approach, we can begin to piece together not just how individual promoters are reg-
105 ulated, but also the nature of gene-gene interactions by revealing how certain transcription factors
106 serve to regulate multiple genes at once. This approach has the benefits of a high-throughput assay
107 while sacrificing little of the resolution afforded by the previous gene-by-gene basis, allowing us to
108 uncover a large swath of the *E. coli* regulome, in detailed resolution, in one set of experiments.

109

110 The organization of the remainder of the paper is as follows. In the Results section, we provide
111 a global view of the discoveries we made in our exploration of more than 100 promoters in *E. coli*
112 using Reg-Seq. These results are described in summary form in the paper itself with a full online
113 version of the results showing how different growth conditions elicit different regulatory responses.
114 This section also follows the overarching view of our results by examining several biological stories
115 that emerge from our data and serve as case studies in what has been revealed in our efforts to
116 uncover the regulatory genome. The Discussion section summarizes the method and the current
117 round of discoveries it has afforded with an eye to future applications of the method to further
118 elucidating the *E. coli* genome and opening up the quantitative dissection of other non-model
119 organisms. Lastly, in the Methods section and flushed out further in the Supplemental Information,
120 we describe our methodology and benchmark it against our own earlier Sort-Seq experiments to
121 show that using RNA-Seq as a readout of the expression of mutated promoters is equally reliable
122 as the fluorescence-based methods.

123

124 **Results**

125 **Selection of genes and methodology**

126 As shown in Figure 1, we have considered more than 100 genes from across the entire *E. coli* genome.
127 Our choices were based on a number of different factors including the desire to have a subset of
128 genes that serve as “gold standard” genes for which the hard work of generations of molecular
129 biologists have yielded deep insights into their regulation. By exploiting our method on these genes,
130 we were able to demonstrate that our Reg-Seq approach recovers not only what was already known
131 about binding sites and transcription factors for well-characterized promoters, but also to reveal
132 whether there are any important differences between the results of the methods presented here
133 and the previous generation experiments based on fluorescence and cell-sorting as a readout of
134 gene expression. These promoters of known regulatory architecture are complemented by an array
135 of previously uncharacterized genes that we selected in part using data from a recent proteomic
136 study using mass spectrometry that measured the copy number of different proteins over a wide
137 range (22 distinct) of growth conditions (*Schmidt et al., 2015*). We selected genes that exhibited a
138 wide variation in their copy number over the different growth conditions considered, reasoning
139 that differential expression across growth conditions implies that those genes are under regulatory
140 control.

141
142 As noted in the introduction, the original formulation of Reg-Seq was based on the use of fluores-
143 cence activated cell sorting as a way to uncover putative binding sites for previously uncharacterized
144 promoters (*Belliveau et al., 2018*). However, our previous method required a gene-by-gene ap-
145 proach to regulatory discovery. As a result, as shown in Figure 1 we have formulated a second
146 generation version that permits a high-throughput interrogation of the genome. Figure 1 provides
147 an overview of the method used in the work presented here. For each promoter we interrogate, we
148 generate a library of mutated variants and design each variant to express an mRNA with a unique
149 sequence barcode. By counting the frequency of each expressed barcode using RNA-Seq, we can as-
150 sess the differential expression from our promoter of interest based on the base-pair-by-base-pair
151 sequence of its promoter. Using the mutual information between mRNA counts and sequences, we
152 can develop an information footprint that reveals the importance of different bases in the promoter
153 region to the overall level of expression. We locate potential transcription factor binding regions by
154 looking for clusters of base pairs that have a significant effect on gene expression ($p < 0.01$). Further
155 details are found in the Methods section. Blue regions of the histogram correspond to hypothesized
156 activating sequences and red regions of the histogram correspond to hypothesized repressing
157 sequences. With the information footprint in hand, we can then determine energy matrices and
158 sequence logos. Given the putative binding sites, we can construct oligonucleotides that serve as
159 fishing hooks to fish out the transcription factors that bind to those putative binding sites. Given all
160 of this information, we can then formulate a schematized view of the newly discovered regulatory
161 architecture of the previously uncharacterized promoter. For the case schematized in Figure 1, the
162 experimental pipeline yields a complete picture of a simple repression architecture.
163

164 **Visual tools for data presentation**

165 Throughout our investigation of the more than 100 genes explored in this study, we repeatedly
166 rely on several key approaches to help make sense of the immense amount of data generated in
167 these experiments. As these different approaches to viewing the results will appear repeatedly
168 throughout the paper, here we familiarize the reader with five graphical representations referred to
169 respectively as information footprints, energy matrices, sequence logos, mass spectrometry enrich-
170 ment plots and regulatory cartoons, which taken all together provide a quantitative description of
171 previously uncharacterized promoters.

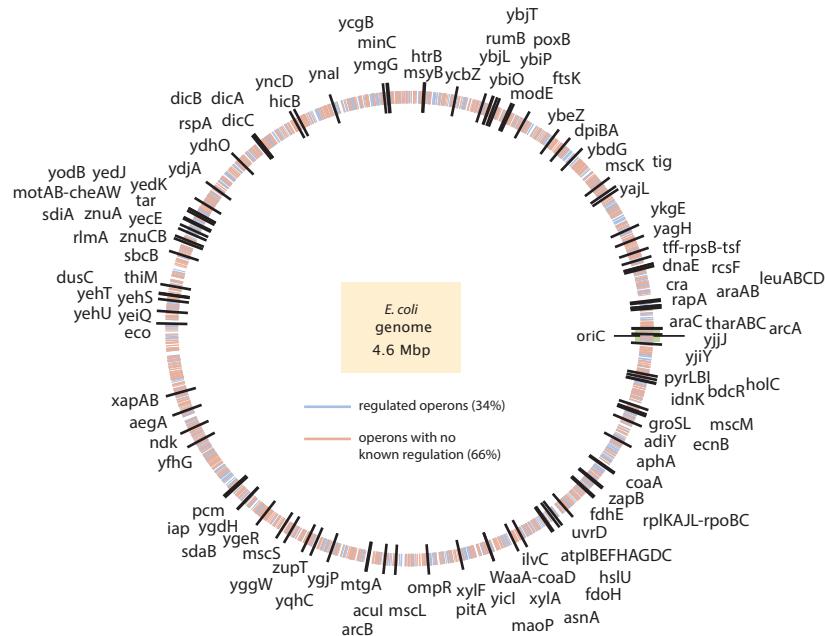
172
173 Information footprints: From our mutagenized libraries of promoter regions, we can build up a
174 base-pair-by-base-pair graphical understanding of how the DNA promoter sequence relates to gene
175 expression level in the form of the information footprint shown in the middle of Figure 1. The result
176 is nicely summarized in our information footprints, where each bar corresponds to an individual
177 base pair and the height of the bar represents how large of an effect mutations at this location
178 have on the gene expression level. Specifically, the value is the mutual information at base pair b
179 between mutation of a base pair at the site and expression level and is given by

$$I_b = \sum_{m,exp} p(m, exp) \log_2 \left(\frac{p(m, exp)}{p(m)p(exp)} \right) \quad (1)$$

180 where m represents whether or not the target base is mutated, and exp refers to reads from the
181 DNA library ($exp = 0$) or the mRNA transcripts ($exp = 1$). $p(m)$ in equation 1 refers to the probability
182 that a given sequencing read will be from a mutated base. $p(exp)$ is a normalizing factor that gives
183 the ratio of the number of DNA or mRNA sequencing counts to total number of counts.

184
185 Furthermore, we color the bars based on whether mutations at this location on average lowered
186 gene expression (in blue, indicating an activating role) or increased gene expression (in red, indicat-
187 ing a repressing role). Within these footprints, we look for regions of contiguous base pairs which

(A) THE REGULATORY GENOME OF *ESCHERICHIA COLI*: PROMOTERS STUDIED



(B)

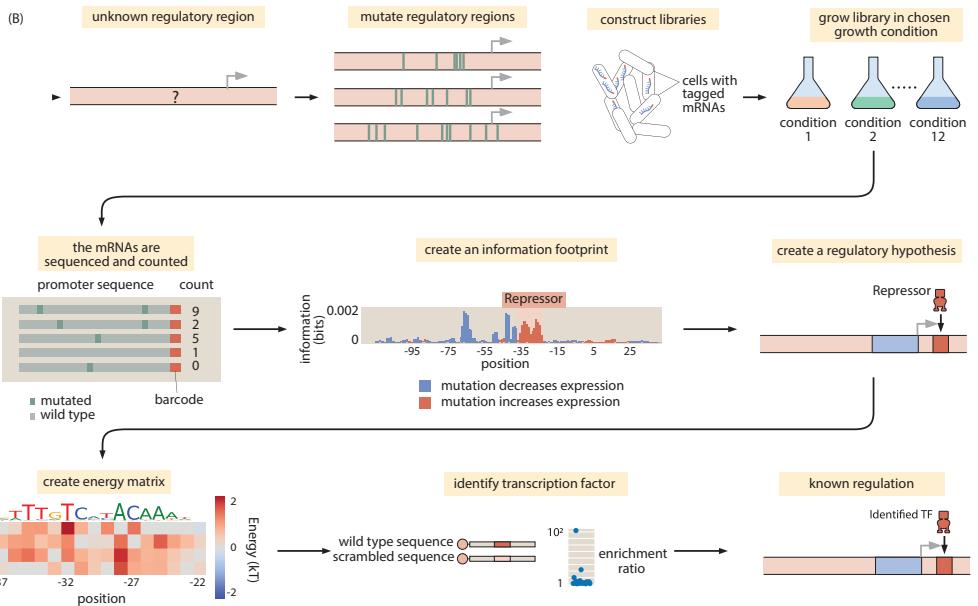


Figure 1. Dissection of the regulatory genome of *Escherichia coli*. (A) Illustration of current ignorance with respect to how genes are regulated in *E. coli*, with genes with previously annotated regulation (as reported on RegulonDB) denoted with blue ticks and genes with no previously annotated regulation denoted with red ticks. The 113 genes explored in this study are labeled in black, with (B) showing the Reg-Seq procedure by which we determine how a given promoter is regulated. This process is as follows: After constructing a promoter library driving expression of a randomized barcode (an average of 5 for each promoter), RNA-Seq is conducted to determine frequency of these mRNA barcodes across different growth conditions. By computing the mutual information between DNA sequence and mRNA barcode counts for each base pair in the promoter region, an "information footprint" is constructed yielding a regulatory hypothesis for the putative binding sites. Energy matrices, which describe the effect any given mutation has on DNA binding energy, and sequence logos are inferred for the putative transcription factor binding sites. Next, we identify which transcription factor preferentially binds to the putative binding site via DNA affinity chromatography followed by mass spectrometry. Finally, this procedure culminates in a coarse-grained cartoon-level view of our regulatory hypothesis for how this given promoter is regulated.

188 impact gene expression similarly (either increasing or decreasing), as these regions implicate the
189 influence of a transcription factor binding site. As can be seen throughout the paper (see Figure 3
190 for several examples of each of the main types of regulatory architectures) and the online resource,
191 we present such information footprints for every promoter we have considered, with one such
192 information footprint for every growth condition.

193
194 Energy matrices: Focusing on an individual putative transcription factor binding site as revealed
195 in the information footprint, we are interested in a more fine-grained, quantitative understanding
196 of how the underlying protein-DNA interaction is determined. An energy matrix displays this
197 information using a heat map format, where each column is a position in the putative binding site
198 and each row displays the effect on binding that results from mutating to that given nucleotide
199 (given as a change in the DNA-TF interaction energy upon mutation). These energy matrices are
200 scaled such that the wild type sequence is colored in white, mutations that improve binding are
201 shown in blue, and mutations that weaken binding are shown in red. These energy matrices encode
202 a full quantitative picture for how we expect sequence to relate to binding for a given transcription
203 factor, such that we can provide a prediction for the binding energy of every possible binding site
204 sequence as

$$\text{binding energy} = \sum_{i=1}^N \epsilon_i, \quad (2)$$

205 where the energy matrix is predicated on a linear binding model in which each base within the
206 binding site region contributes a specific value (ϵ_i for the i^{th} base in the sequence) to the total
207 binding energy. Energy matrices are either given in A.U. (arbitrary units), or if the gene has a simple
208 repression or activation architecture with a single RNAP site, are assigned kT energy units following
209 the procedure in *Kinney et al. (2010)* and validated on the *lac* operon in *Barnes et al. (2019)*.

210
211 Sequence logos: From an energy matrix, we can also represent a preferred transcription factor
212 binding site with the use of the letters corresponding to the four possible nucleotides, as is often
213 done with position weight matrices (*Schneider and Stephens, 1990*). In these sequence logos, the
214 size of the letters corresponds to how strong the preference is for that given nucleotide at that
215 given position, which can be directly computed from the energy matrix. This method of visualizing
216 the information contained within the energy matrix is more easily digested by human eyes and
217 allows for quick comparison among various binding sites.

218
219 Mass spectrometry enrichment plots: As the final piece of our experimental pipeline, we wish to
220 determine the identity of the transcription factor we suspect is binding to our putative binding site
221 that is represented in the energy matrix and sequence logo described above. While the details of
222 the DNA chromatin affinity chromatography and mass spectrometry can be found in the methods,
223 the results of these experiments are displayed in enrichment plots such as is shown in the bottom
224 panel of Figure 1(B). In these plots, the relative abundance of each protein binding to our site of
225 interest is quantified relative to a scrambled control sequence. The putative transcription factor is
226 the one we find to be highly enriched compared to all other DNA binding proteins ($p < 10^{-4}$).

227
228 Regulatory cartoons: The ultimate result of all these detailed base-pair-by-base-pair resolution
229 experiments yields a cartoon model of how we think the given promoter is being regulated. While
230 the cartoon serves as a convenient visual way to summarize our results, it's important to remember
231 that these cartoons are a shorthand representation of all the data in the four quantitative measures
232 described above and are in fact backed by quantitative predictions of how we expect the system to
233 behave. Throughout this paper we use consistent iconography to illustrate the regulatory architec-
234 ture of promoters, with activators and their binding sites in green, repressors in red, and RNAP in
235 blue.

236

237 **Specific results**

238 Figure 2 provides a summary of the discoveries made in the work done here using our next genera-
 239 tion Reg-Seq approach. Figure 2(A) provides a shorthand notation that conveniently characterizes
 240 the different kinds of regulatory motifs found in bacteria. In previous work, we have explored the
 241 entirety of what is known about the regulatory genome of *E. coli*, revealing that the most common
 242 motif is the (0,0) constitutive architecture, though we hypothesize that this is not a statement about
 243 the facts of the *E. coli* genome, but rather a reflection of our collective regulatory ignorance. The
 244 two most common regulatory architectures that emerged from our previous database survey are
 245 the (0,1) and (1,0) architectures, the simple repression motif and the simple activation motif, respec-
 246 tively. It is interesting to consider that the (0,1) architecture is in fact the repressor-operator model
 247 originally introduced in the early 1960s by Jacob and Monod as the concept of gene regulation
 248 emerged and to now see retrospectively the far-reaching importance of that architecture across the
 249 *E. coli* genome (*Jacob and Monod, 1961*).
 250

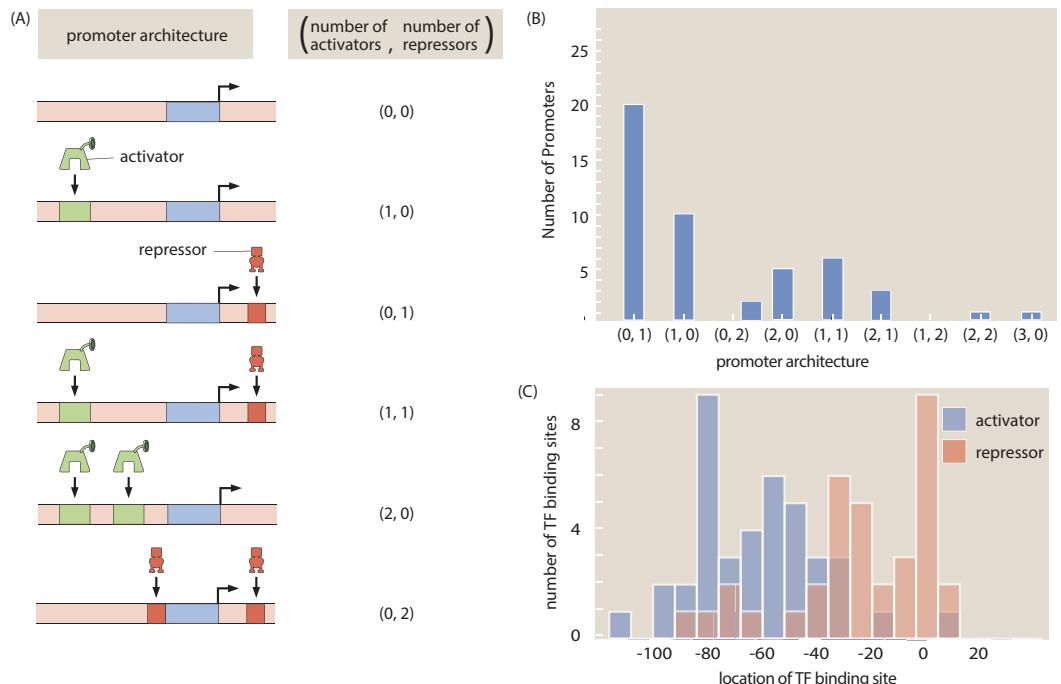


Figure 2. A summary of regulatory architectures discovered in this study. (A) The cartoons display a representative example of each type of architecture, along with the corresponding shorthand notation. (B) Counts of the different regulatory architectures discovered in this study as compared to known regulation found on RegulonDB. (C) Distribution of positions of binding sites discovered in this study for activators and repressors, respectively.

251 Figure 2(B) summarizes, for the 113 genes we considered, how many of them are simple activa-
 252 tion motifs, how many are simple repression motifs and so on. We observe that the most common
 253 motif to emerge from our work is the simple repression motif. Another relevant regulatory statistic
 254 is shown in Figure 2(C) where we see the distribution of binding site positions. Our own experience
 255 in the use of different quantitative modeling approaches to consider transcriptional regulation
 256 reveal that for now, we remain largely ignorant of how to account for transcription factor binding
 257 site position, and datasets like that presented here will begin to provide data that can help us un-
 258 cover how this parameter dictates gene expression. Indeed, with binding site positions and energy
 259 matrices in hand, we can systematically move these binding sites and explore the implications for
 260 the level of gene expression, providing a systematic tool to understand the binding-site position
 261 “knob”.

262

Figure 3 delves more deeply into the various regulatory architectures described in Figure 2 (B) by showing some of the promoters for which these architectures were discovered. In each of the cases shown in the figure, prior to the work presented here, these promoters had no regulatory information in knowledge bases like, Ecocyc (*Keseler et al., 2016*) and RegulonDB (*Santos-Zavaleta et al., 2019*). Now, using the sequencing methods explained above in conjunction with mass spectrometry we were able to identify candidate binding sites that were then used to construct oligonucleotides that were used to construct biotinylated DNA “bait” probes to pull down their corresponding putative transcription factors. Given this regulatory information, we are now poised to carry out a quantitative dissection of these promoters using theoretical tools from statistical physics to predict their input-output function and precision biophysical measurements to test those predictions.

274

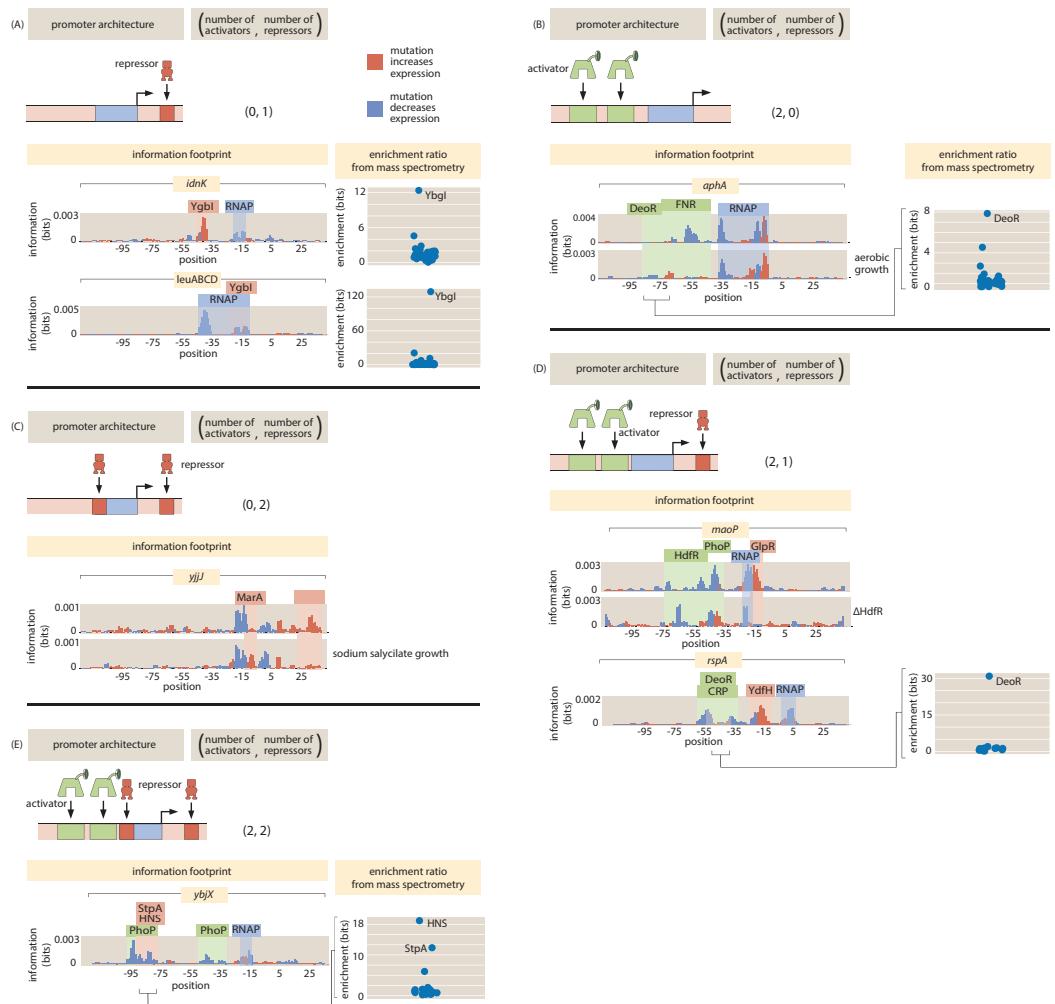


Figure 3. Newly discovered or updated regulatory architectures. Examples of information footprints, gene knockouts, and mass spectrometry data used to identify transcription factors for five genes. (A) Examples of simple repression architectures where the locations of the putative binding sites are highlighted in red and the identities of the bound transcription factors are revealed in the mass spectrometry data. (B) An example of a 2 activator architecture. (C) An example of a 2 repressor architecture. (D) An example of a 2 activator and 1 repressor architecture. (E) An example of a 2 activator and 2 repressor architecture.

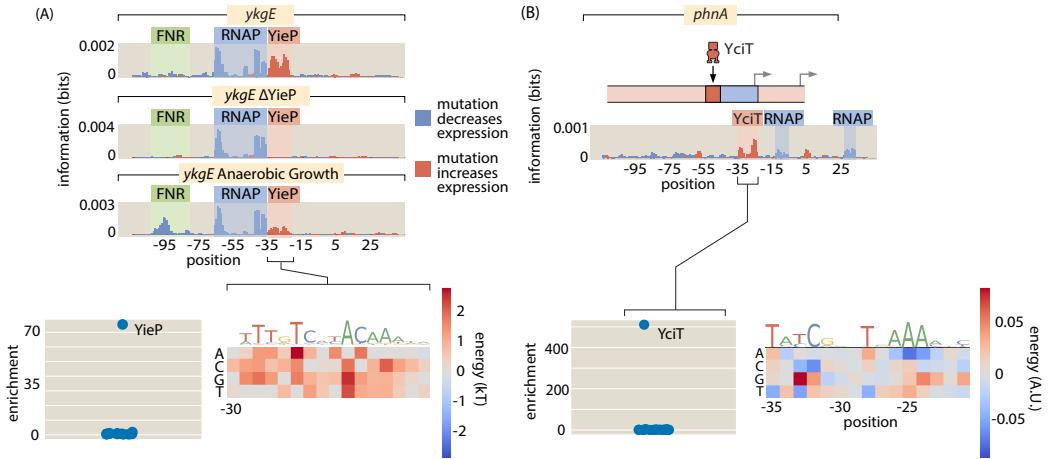


Figure 4. Examples of the power of Reg-Seq in the context of promoters with no previously known regulatory information. (A) From the information footprint of the *ykgE* promoter under different growth conditions, we can identify a repressor binding site downstream of the RNAP binding site. From the enrichment of proteins bound to the DNA sequence of the putative repressor as compared to a control sequence, we can identify YieP as the transcription factor bound to this site as it has a much higher enrichment ratio than any other protein. Lastly, the binding energy matrix for the repressor site along with corresponding sequence logo shows that the wild type sequence is the strongest possible binder and it displays an imperfect inverted repeat symmetry. (B) illustrates a comparable dissection for the *phnA* promoter.

275 A recent paper identified what they christened the y-ome in *E. coli* (Ghatak *et al.*, 2019), which
 276 they defined as the set of genes that don't have the annotation necessary to have an effect in a
 277 systems biology model. Their surprising result is the finding that roughly 35% of the genes in the *E.*
 278 *coli* genome are functionally unannotated. The situation is likely worse for other organisms, demon-
 279 strating the need for methods such as Reg-Seq to come to terms with these parts of genomes. For
 280 many of the genes in the y-ome, we remain similarly ignorant of how they are regulated. Figures 3
 281 and 4 provide several examples from the y-ome of genes and transcription factors for which little to
 282 nothing was previously known. As shown in Figure 4, our study has found the first examples that
 283 we are aware of in the entire *E. coli* genome of a binding site for YciT. These examples are intended
 284 to show the outcome of the methods developed here and to serve as an invitation to browse the
 285 online resource to see many examples of the regulation of y-ome genes.

286

287 The ability to find binding sites for both widely acting regulators and transcription factors which
 288 may have only a few sites in the whole genome allows us to get an in-depth and quantitative view
 289 of any given promoter. As indicated in Figures 4(A)-(B), in both cases, we were able to perform the
 290 relevant search and capture for the transcription factors that bind our putative binding sites. In
 291 both of these cases, we now hypothesize that these newly discovered binding site-transcription
 292 factor pairs exert their control through repression. The ability to extract the quantitative features of
 293 regulatory control through energy matrices means that we can take a nearly unstudied gene such
 294 as *ykgE*, which is regulated by an understudied transcription factor YieP, and quickly get to the point
 295 at which we can do quantitative modeling in the style that we and many others have performed on
 296 the *lac* operon (Kinney *et al.*, 2010; Bintu *et al.*, 2005; Barnes *et al.*, 2019; Garcia and Phillips, 2011).

297

298 One of the revealing case studies that demonstrates the broad reach of our approach for dis-
 299 covering regulatory architectures is offered by the insights we have gained into two widely acting
 300 regulators. In both cases, we have expanded the array of promoters that they are now known to
 301 regulate. Further, these two case studies illustrate that even for regulators with the potential to

302 change the entire landscape of *E. coli* gene regulation, there is a large gap in regulatory knowledge
303 and the approach advanced here has the power to discover new regulatory motifs. The newly
304 discovered binding sites in Figure 5(A) more than double the number of operons known to be
305 regulated by GlpR as reported in RegulonDB (*Santos-Zavaleta et al., 2019*). We found 5 newly
306 regulated operons in our data set, even though we were not specifically targeting GlpR regulation.
307 Although the numbers are too small to make good estimates, this does support the notion that
308 GlpR widely regulates and many of its sites would be found in a full search of the genome. The
309 regulatory roles revealed in Figure 5(A) also reinforce the evidence that GlpR is a repressor.

310

311 GlpR is known to bind the inducer glycerol-3-phosphate (G3P) (*Larsons et al., 1987*). When
312 G3P binds to GlpR, similarly to the way that the classic LacI repressor is induced by allolactose,
313 the probability of finding GlpR bound to DNA is reduced. In the three growth conditions in our
314 study containing glucose, of which a representative set of examples are shown in Figure 5(A), the
315 repressor binds strongly while in all other growth conditions the binding is greatly diminished but
316 not entirely abolished. As there is no previously known direct molecular interaction between GlpR
317 and glucose and the repression is reduced but not eliminated, the derepression in the absence of
318 glucose is likely an indirect effect. None of the genes found to be regulated by GlpR are also under
319 CRP regulation, and we cannot be sure what the exact cause of the indirect interaction is. However,
320 *gpsA* is activated by CRP (*Seoh and Tai, 1999*). Its product synthesizes the inducer of GlpR (G3P) and
321 so intracellular G3P levels will be lowered by the presence of glucose, which would repress those
322 genes under GlpR control.

323

324 Prior to this study, there were 4 operons known to be regulated by GlpR, each with between 4
325 and 8 GlpR binding sites (*Gama-Castro et al., 2016*). In these operons, the absence of glucose and
326 the partial induction of GlpR was not enough to prompt a notable change in gene expression (*Lin,*
327 *1976*). In these operons, GlpR acts as part of an AND gate, where high G3P concentration *and* an
328 absence of glucose is required for high gene expression. By way of contrast, we have discovered
329 operons whose regulation appears to be mediated by a single GlpR site per operon. With only a
330 single site, GlpR functions as a glucose sensor, as only the absence of glucose is needed to relieve
331 repression by GlpR.

332

333 The second widely acting regulator our study revealed, FNR, has 151 binding sites already
334 reported in RegulonDB and is well studied compared to most transcription factors (*Gama-Castro*
335 *et al., 2016*). However, the newly discovered FNR sites displayed in Figure 5(B) demonstrate that for
336 even well-understood transcription factors there is much still to be uncovered. Our information
337 footprints are in agreement with previous studies suggesting that FNR acts as an activator. In the
338 presence of O₂, dimeric FNR is converted to a monomeric form and its ability to bind DNA is greatly
339 reduced (*Myers et al., 2013*). Only in low oxygen conditions did we observe a binding signature
340 from FNR, and we show a representative example of the information footprint from one of 11
341 growth conditions with plentiful oxygen in Figure 5(B).

342

343 To end the summary of our results, we would like to further emphasize the value we gained from
344 having taken a hundred-gene approach rather than the gene-by-gene approach taken before. We
345 observe quantitatively how FNR affects the expression of *fdhE* both directly through transcription
346 factor binding (Figure 6(A)) and indirectly through increased expression of ArcA (Figure 6(B)). Also,
347 fully understanding even a single operon often requires investigating several regulatory regions as
348 we have in the case of *fdoGHI-fdhE* by investigating the main promoter for the operon as well as the
349 promoter upstream of *fdhE*. 36% of all multi-gene operons have at least one TSS which transcribes
350 only a subset of the genes in the operon (*Conway et al., 2014*). Regulation within an operon is even
351 more poorly studied than regulation in general. The main promoter for *fdoGHI-fdhE* is repressed
352 by a transcription factors. We see in Figure 6(B) that, as *fdhE* is regulated by ArcA and FNR, it

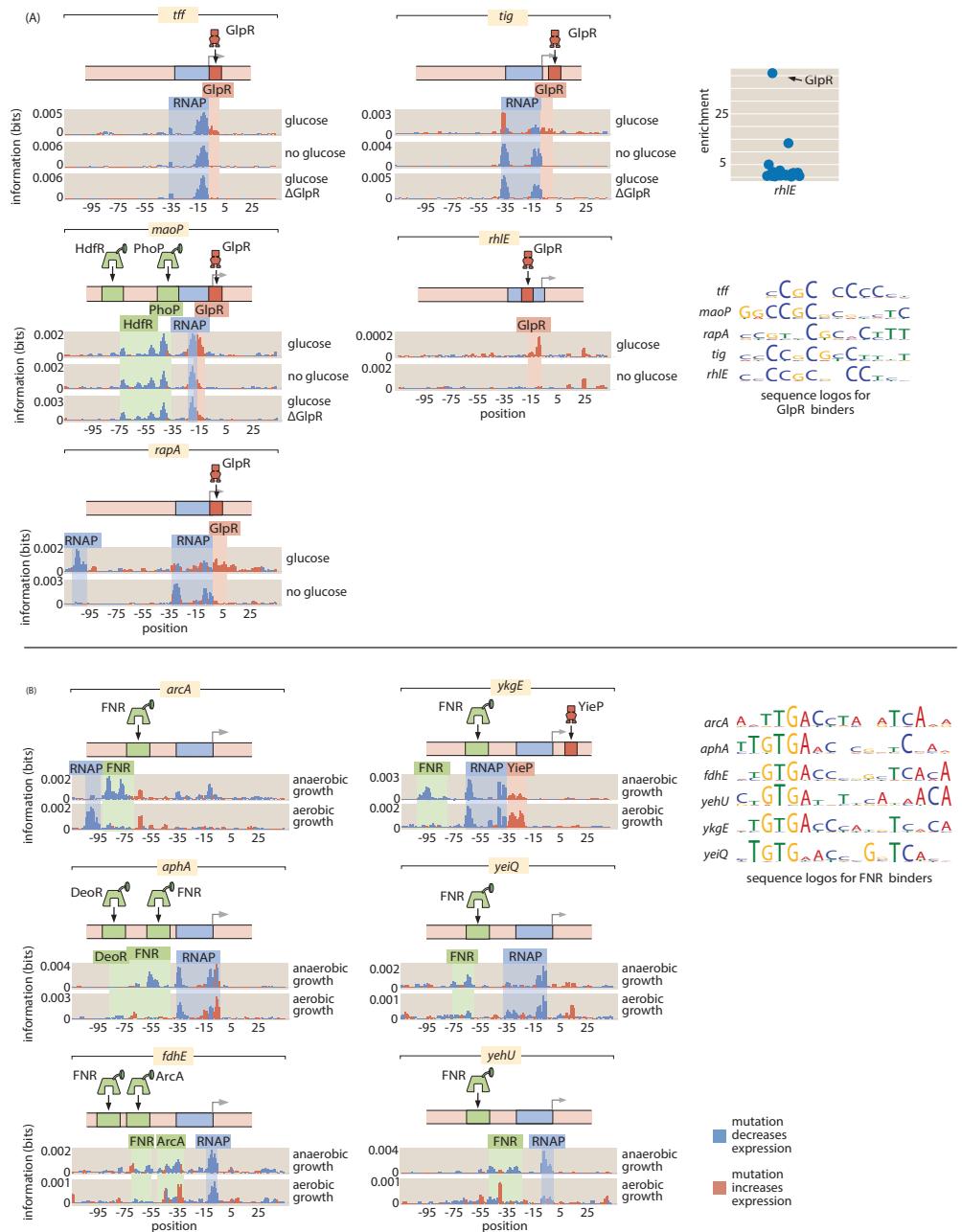


Figure 5. Reg-Seq analysis of broadly-acting transcription factors. By using Reg-Seq to dissect over 100 promoters at once, we can begin to understand not only how individual promoters are regulated, but also how individual transcription factors serve to regulate a suite of genes. (A) GlpR as a widely-acting regulator. Here we show the many promoters which we found to be regulated by GlpR, all of which were previously unknown. GlpR was demonstrated to bind to *rhlE* by mass spectrometry enrichment experiments as shown in the top right. Binding sites in the *tff*, *tig*, *maoP*, *rhlE*, and *rapA* have similar DNA binding preferences as seen in the sequence logos and each TF binding site binds strongly only in the presence of glucose. These similarities suggest that the same TF binds to each site. To test this hypothesis we knocked out GlpR and ran the Reg-Seq experiments for *tff*, *tig*, and *maoP*. We see that knocking out GlpR removes the binding signature of the TF. (B) FNR as a global regulator. FNR is known to be upregulated in anaerobic growth, and here we found it to regulate a suite of six genes. In growth conditions with prevalent oxygen the putative FNR sites are weakened, and the DNA binding preference of the six sites are shown to be similar from the sequence logos displayed on the right.

353 will be upregulated in anaerobic conditions (*Compan and Touati, 1994*). The main TSS transcribes
 354 all four genes in the operon, while the secondary site shown in Figure 6(B) only transcribes *fdhE*,
 355 and therefore anaerobic conditions will change the stoichiometry of the proteins produced by the
 356 operon. At the higher throughput that we use in this experiment it becomes feasible to target
 357 multiple promoters within an operon as we have done with *fdoGHI-fdhE*. We can then determine
 358 under what conditions the typical operon assumption breaks down.
 359

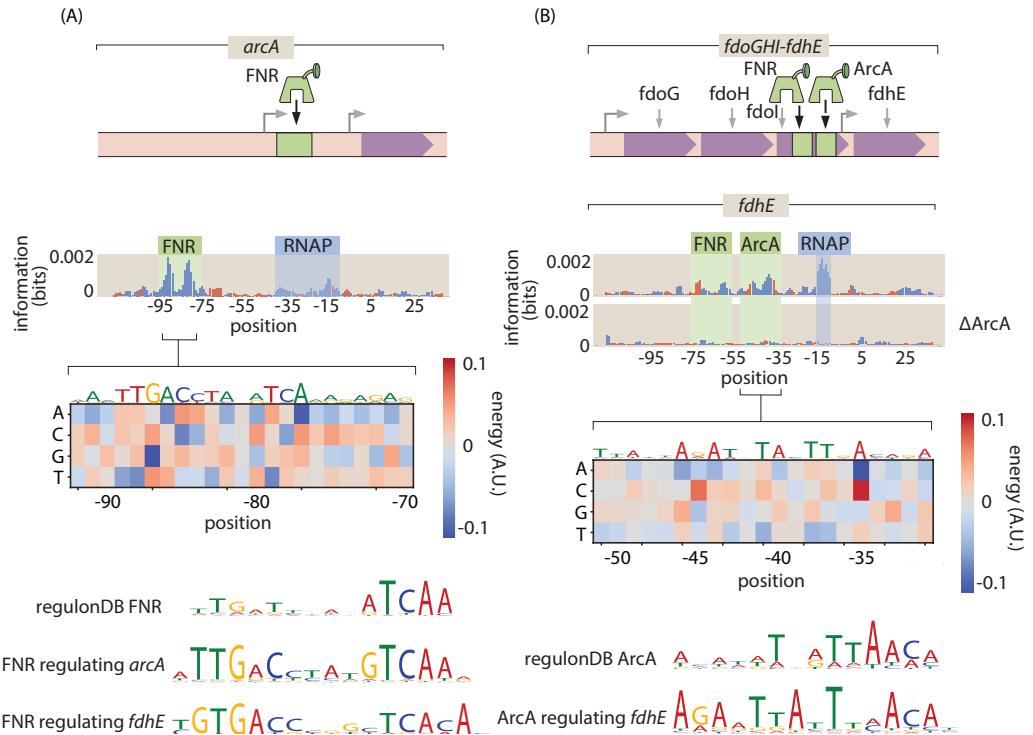


Figure 6. Inspection of an anaerobic respiration genetic circuit. By exploring over a 100 promoters at once, we can begin to piece together some genetic circuitry of the cell. A) Here we see not only how the *arcA* promoter is regulated, but also the role this transcription factor plays in the regulation of another promoter. B) Intra-operon regulation of *fdhE* by both FNR and ArcA. A TOMTOM search of the binding motif found that ArcA was the most likely candidate for the transcription factor. A knockout of ArcA demonstrates that the binding signature of the site, and its associated RNAP site, are no longer significant determinants of gene expression.

360 By examining the over 100 promoters considered here, grown under 12 growth conditions, we
 361 have a total of more than 1000 information footprints and data sets. In this age of big data, methods
 362 to explore and draw insights from that data are crucial. To that end, as introduced in Figure 7, we
 363 have developed an online resource (see www.rpgroup.caltech.edu/RNAseq_SortSeq/interactive_a)
 364 that makes it possible for anyone who is interested to view our data and draw their own biological
 365 conclusions. Information footprints for any combination of gene and growth condition are displayed
 366 via drop down menus. Each identified transcription factor binding site or transcription start site is
 367 marked, and energy matrices for all transcription factor binding sites are displayed. In addition,
 368 for each gene, we feature a simple cartoon-level schematic that captures our now current best
 369 understanding of the regulatory architecture and resulting mechanism.

370
 371 The interactive figure in question was invaluable in identifying transcription factors, such as
 372 GlpR, whose binding properties vary depending on growth condition. As sigma factor availability
 373 also varies greatly depending on growth condition, studying the interactive figure identified many of

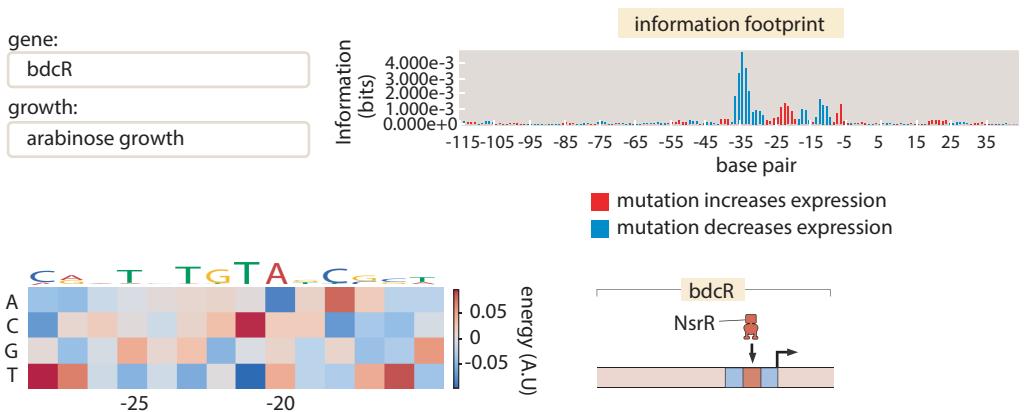


Figure 7. Representative view of the interactive figure that is available online. This interactive figure captures the entirety of our dataset. Each figure features a drop-down menu of genes and growth conditions. For each such gene and growth condition, there is a corresponding information footprint revealing putative binding sites, an energy matrix that shows the strength of binding of the relevant transcription factor to those binding sites and a cartoon that schematizes the newly-discovered regulatory architecture of that gene.

374 the secondary RNAP sites present. The interactive figure provides a valuable resource both to those
 375 who are interested in the regulation of a particular gene and those who wish to look for patterns in
 376 gene regulation across multiple genes or across different growth conditions.
 377

378 Discussion

379 The study of gene regulation is one of the centerpieces of modern biology. As a result, it is surprising
 380 that in the genome era, our ignorance of the regulatory landscape even of our best-understood
 381 model organisms remains so vast. Despite understanding the regulation of transcription initiation
 382 in bacterial promoters (*Browning and Busby, 2016*), and how to tune their expression, we lack
 383 an experimental framework to unravel understudied promoter architectures at scale. Until the
 384 work presented here, even in the case of one of biology's best understood organisms, the humble
 385 bacterium *E. coli*, over 65% of its genes have no known regulatory architecture (*Belliveau et al., 2018*).
 386 As such, in our view one of the grand challenges of the genome era is the need to uncover
 387 the regulatory landscape for each and every organism with a known genome sequence. Given the
 388 ability to read and write DNA sequence at will, we are convinced that to make that reading of DNA
 389 sequence truly informative about biological function and to give that writing the full power and
 390 poetry of what Crick christened "the two great polymer languages", we need a full accounting of
 391 how the genes of a given organism are regulated and how environmental signals communicate
 392 with the transcription factors that mediate that regulation – the so-called "allosterome" problem
 393 (*Lindsley and Rutter, 2006*). The work presented here provides a general methodology for attacking
 394 the former problem and even shows how by performing Reg-Seq in different growth conditions, we
 395 can make headway on the latter problem.
 396

397 The advent of cheap DNA sequencing offers the promise of beginning to achieve that grand
 398 challenge goal in the form of massively-parallel reporter assays reviewed in *Kinney and McCandlish*
 399 (*2019*). A particular implementation of such methods was christened Sort-Seq (*Kinney et al., 2010*)
 400 and was demonstrated in the context of well understood regulatory architectures. A next generation
 401 of the Sort-Seq method (*Belliveau et al., 2018*) established experiments through the use of affinity
 402 chromatography and mass spectrometry which made it possible to identify the transcription factors
 403 that bind the putative binding sites discovered in the Sort-Seq analysis. But there were critical

404 shortcomings in the method, not least of which was that it lacked the scalability to uncover the
405 regulatory genome on a genome-wide basis.

406
407 The work presented here builds on the foundations laid in the previous studies by invoking
408 RNA-Seq as a readout of the level of expression of the promoter mutant libraries needed to infer
409 information footprints and their corresponding energy matrices and sequence logos. The case
410 studies described in the main text of the paper provide evidence of the ability of the method to
411 deliver on the promise of beginning to systematically uncover the regulatory genome and the
412 extensive online resources hint at a way of systematically reporting those insights in a way that can
413 be used by the community at large to develop regulatory intuition for biological function and to
414 design novel regulatory architectures using energy matrices.

415
416 Several shortcomings remain in the approach introduced here. First, we are well aware of and
417 excited about regulatory action at a distance. Indeed, our laboratory has invested a significant
418 effort in exploring such long-distance regulatory action in the form of DNA looping in bacteria
419 and VDJ recombination in jawed vertebrates. It is well known that transcriptional control through
420 enhancers in eukaryotic regulation is central in contexts ranging from embryonic development to
421 hematopoiesis (*Melnikov et al., 2012*). The current incarnation of the methods described here has
422 focused on contiguous regions in the vicinity of the transcription start site. Clearly, to go further in
423 dissecting the entire regulatory genome, these methods will have to be extended to non-contiguous
424 regions of the genome.

425
426 The sequencing technology to use Reg-Seq to identify putative DNA binding sites across the
427 genome already exists, and this technique could be applied in parallel to recover putative binding
428 sites across the genome without losing resolution. However, a second key challenge faced by the
429 methods described here is that the mass spectrometry and the gene knockout confirmation aspects
430 of the experimental pipeline remain low throughput. To overcome this, we have begun to explore
431 next-generation experiments that will make it possible to accomplish transcription factor identifica-
432 tion at higher throughput. We are exploring multiplexed mass spectrometry measurements and
433 multiplexed Reg-Seq on libraries of gene knockouts as ways to break the identification bottleneck.

434
435 The findings from this study provide a foundation for systematically performing genome-wide
436 regulatory dissections. We have developed a method to pass from complete regulatory ignorance
437 to designable regulatory architectures in short order and we are hopeful that others will adopt
438 these methods with the ambition of uncovering the regulatory architectures that preside over their
439 organisms of interest.

440

441 **Methods**

442 **Library construction**

443 Promoter variants were synthesized on a microarray (TWIST Bioscience, San Francisco, CA). The
444 promoters were designed computationally to have a 10 % rate and to cover 160 bp. The mutations
445 were generated randomly, and the library was regenerated if the mutation rate did not fall within
446 0.5% of the target mutation rate. There are an average of 2200 unique promoter sequences per
447 gene. An average of 5 unique 20 base pair barcodes per variant promoter was used for the purpose
448 of counting transcripts. The barcode was inserted 110 base pairs from the 5' end of the mRNA,
449 containing 45 base pairs from the targeted regulatory region, 64 base pairs containing primer sites
450 used in the construction of the plasmid, and 11 base pairs containing a three frame stop codon.
451 All the sequences are listed in Supplementary Table 1. After the barcode there is a RBS and a GFP
452 coding region. Mutated promoters were PCR amplified and inserted by Gibson assembly into the

453 plasmid backbone of pJK14 (SC101 origin) (*Kinney et al., 2010*). Constructs were electroporated
454 into *E. coli* K-12 MG1655 (*Blattner, 1997*).
455

456 **RNA preparation and sequencing**

457 Cells were grown to an optical density of 0.3 and RNA was then stabilized using Qiagen RNA Protect
458 (Qiagen, Hilden, Germany). Lysis was preformed using lysozyme (Sigma Aldrich, Saint Louis, MO)
459 and RNA was isolated using the Qiagen RNA Mini Kit. Reverse transcription was preformed using
460 Superscript IV (Invitrogen, Carlsbad, CA) and a specific primer for the labeled mRNA. qPCR was
461 preformed to check the level of DNA contamination and the mRNA tags were PCR amplified and
462 Illumina sequenced. Within a single growth condition, all promoter variants for all regulatory regions
463 were tested in a single multiplexed RNA-seq experiment.
464

465 **Analysis of sequencing results**

466 To determine putative transcription factor binding sites, we first computed the effect of mutations
467 on gene expression at a base pair-by-base pair level using information footprints. To initially capture
468 all regions important for gene expression, regardless of whether or not they are transcription factors,
469 we identify regions where gene expression is changed significantly up or down by mutation ($p <$
470 0.01). We find a region of effect by averaging over 15 base pairs, as this is a reasonable size for a
471 transcription factor binding site. We include binding sites that are significant in any of the tested
472 growth conditions.

473 Many false positive will be secondary RNAP sites and we remove from consideration any sites
474 that resemble RNAP sites. We fit energy matrices to each of the possible binding sites and use
475 the preferred DNA sequence for binding to identify the RNAP sites. We use both visual inspection
476 (for example, does the site resemble the consensus TGNNTATAAT extended minus 10 for sigma 70
477 sites), and the TOMTOM tool to computationally compare the potential site to examples of sigma
478 70, sigma 38, and sigma 54 sites that we determined in this experiment. We discard any sites that
479 have a p-value of similarity with an RNAP site of less than 5×10^{-3} . If a single site contains an RNAP
480 site along with a transcription factor site we remove only those bases containing the probable RNAP
481 site. This leaves 83 identified transcription factor binding regions. We adjust the edges of binding
482 sites to only be wide enough to contain all base-pairs with a significant effect on gene expression (p
483 < 0.01).

484 For primary RNAP sites, we include a list of probable sigma factor identities as a Supplementary
485 Table. Sites are judged by similarity to consensus binding sites, with the most similar sigma factor
486 site reported based on TOMTOM comparison (*Gupta et al., 2007*). Those sites where there is signif-
487 icant uncertainty due to overlapping sites are omitted. Overlapping sites in general can pose issues
488 for this method. In many cases, looking at growth conditions where only one of the overlapping
489 sites is active is effective to establish site boundaries and energy matrices. For sites where no
490 adequate growth condition can be found, or when a TF overlaps with an RNAP site, the energy
491 matrix will not be reflective of the true DNA-protein interaction energies. If the TFs in overlapping
492 sites are composed of one activator and one repressor, then we use the point at which the effect
493 of mutation shifts from activator-like to repressor-like as a demarcation point between binding sites.
494

495 **DNA affinity chromatography and mass spectrometry**

496 Upon identifying a putative transcription factor binding site, we used DNA affinity chromatography,
497 as done in (*Belliveau et al., 2018*) to isolate and enrich for the transcription factor of interest. In
498 brief, we order biotinylation oligos of our binding site of interest (Integrated DNA Technologies,
499 Coralville, IA) along with a control, "scrambled" sequence, that we expect to have no specificity
500 for the given transcription factor. We then tether these oligos to magnetic streptavidin beads

501 (Dynabeads MyOne T1; ThermoFisher, Waltham, MA), which we then incubate overnight with whole
502 cell lysate grown in the presences of either heavy (with ^{15}N) or light (with ^{14}N) lysine for the experi-
503 mental and control sequences, respectively. The next day, proteins were recovered by digesting
504 the DNA with the PstI restriction enzyme (New England Biolabs, Ipswich, MA), whose cut site was
505 incorporated into all designed oligos.

506
507 Protein samples were then prepared for mass spec by in-gel digestion using the Lys-C protease
508 (Wako Chemicals, Osaka, Japan) or in-solution digestion. Liquid chromatography coupled mass
509 spectrometry (LC-MS) was performed as previously described by (*Belliveau et al., 2018*), and are
510 further discussed in the SI appendix. SILAC labeling was performed by growing cells (Δ LysA) in
511 either heavy isotope form of lysine or its natural form.

512
513 It is also important to note that while we relied on the SILAC method to identify the TF identity
514 for each promoter, our approach doesn't require this specific technique. Specifically, our method
515 only requires a way to contrast between the copy number of proteins bound to a target promoter in
516 relation to a scrambled version of the promoter. In principle, one could use multiplexed proteomics
517 based on isobaric mass tags (*Pappireddi et al., 2019*) to characterize up to 10 promoters in parallel.
518 Isobaric tags are reagents used to covalently modify peptides by using the heavy-isotope distri-
519 bution in the tag to encode different conditions. The most widely adopted methods for isobaric
520 tagging are the isobaric tag for relative and absolute quantitation (iTRAQ) and the tandem mass tag
521 (TMT). This multiplexed approach involves the fragmentation of peptide ions by colliding with an
522 inert gas. The resulting ions are resolved in a second MS-MS scan (MS2).

523
524 Only a subset (20) of all identified transcription factor targets were analyzed by mass spec due to
525 limitations in scaling the technique to large numbers of targets. The transcription factors identified
526 by this method are enriched more than any other DNA binding protein ($p < 10^{-4}$).

528 **Construction of knockout strains**

529 Conducting DNA affinity chromatography followed by mass spectrometry on our putative binding
530 sites resulted in likely candidates for the transcription factors that are responsible for the infor-
531 mation contained at a given promoter region. To verify that a given transcription factor is, in fact,
532 acting to regulate a given promoter, we repeated the RNA sequencing experiments on strains with
533 the transcription factor of interest knocked out.

534
535 To construct the knockout strains, we ordered strains from the Keio collection (*Yamamoto et al.,*
536 *2009*) from the Coli Genetic Stock Center. These knockouts were then put in a MG1655 background
537 via phage P1 transduction and verified with Sanger sequencing. To remove the kanamycin resistance
538 that comes with the strains from the Keio collection, we transformed in the pCP20 plasmid, induced
539 FLP recombinase, and then selected for colonies that no longer grew on either kanamycin or
540 ampicillin. Finally, we transformed our desired promoter libraries into these constructed knockout
541 strains, allowing us to perform the RNA sequencing in the same context as the original experiments.

542

543 **Code and Data Availability**

544 All code used for processing data and plotting as well as the final processed data, plasmid sequences,
545 and primer sequences can be found on the GitHub repository (<https://github.com/RPGGroup->
546 PBoC/RNAseq_SortSeq) All raw sequencing data is available at the Sequence Read Archive (accession
547 no.SUB6723530). All inferred information footprints and energy matrices can be found on the
548 CalTech data repository doi:10.22002/D1.1331.

549 **Acknowledgments**

550 Experimentally, we would like to thank Jost Vielmetter and Nina Budaeva for providing access to
551 their Cell Disruptor. Brett Lomenick provided crucial help and advice with protein preparation. We
552 also thank Igor Antoshechkin for his help with sequencing at the Caltech Genomics Facility. We
553 also thank Stephanie Barnes, Griffin Chure, Hernan Garcia, Soichi Hirokawa, Heun Jin Lee, Nicholas
554 McCarty, Muir Morrison, Manuel Razo-Mejia and Gabe Salmon for useful discussion. Guillaume
555 Urtecho and Sri Kosuri have been instrumental in providing key advice and protocols at various
556 stages in the development of this work.

557

558 We would like to thank Guillaume Urtecho, Hernan Garcia, Steve Quake, Matt Thompson, Curt
559 Callan

560 Funding: We are deeply grateful for support from NIH Grants DP1 OD000217 (Director's Pioneer
561 Award) and 1R35 GM118043-01 (Maximizing Investigators Research Award) which made it possible
562 to undertake this multi-year project. N.M.B. was supported by an HHMI International Student
563 Research Fellowship. S.M.B was supported by the NIH training grant.,

564

565 **References**

- 566 Barnes SL, Belliveau NM, Ireland WT, Kinney JB, Phillips R. Mapping DNA sequence to transcription factor
567 binding energy *in vivo*. PLoS Computational Biology. 2019; 15(2):1–29. doi: [10.1371/journal.pcbi.1006226](https://doi.org/10.1371/journal.pcbi.1006226).
- 568 Belliveau NM, Barnes SL, Ireland WT, Jones DL, Sweredoski MJ, Moradian A, Hess S, Kinney JB, Phillips R.
569 Systematic approach for dissecting the molecular mechanisms of transcriptional regulation in bacteria.
570 Proceedings of the National Academy of Sciences of the United States of America. 2018; 115(21):E4796–E4805.
571 doi: [10.1073/pnas.1722055115](https://doi.org/10.1073/pnas.1722055115).
- 572 Bintu L, Buchler NE, Garcia HG, Gerland U, Hwa T, Kondev J, Phillips R. Transcriptional regulation by
573 the numbers: models. Current Opinion in Genetics & Development. 2005 Apr; 15(2):116–124. doi:
574 [10.1016/j.gde.2005.02.007](https://doi.org/10.1016/j.gde.2005.02.007).
- 575 Blattner FR. The Complete Genome Sequence of Escherichia coli K-12. Science. 1997 Sep; 277(5331):1453–1462.
576 doi: [10.1126/science.277.5331.1453](https://doi.org/10.1126/science.277.5331.1453).
- 577 Browning DF, Busby SJW. Local and global regulation of transcription initiation in bacteria. Nature Reviews
578 Microbiology. 2016; p. 638–650. doi: [10.1038/nrmicro.2016.103](https://doi.org/10.1038/nrmicro.2016.103).
- 579 Chure G, Razo-Mejia M, Belliveau NM, Einav T, Kaczmarek ZA, Barnes SL, Lewis M, Phillips R. Predictive shifts in
580 free energy couple mutations to their phenotypic consequences. Proceedings of the National Academy of
581 Sciences of the United States of America. 2019; 116(37):18275–18284. doi: [10.1073/pnas.1907869116](https://doi.org/10.1073/pnas.1907869116).
- 582 Compan I, Touati D. Anaerobic activation of arcA transcription in *Escherichia coli* : roles of Fnr and ArcA.
583 Molecular Microbiology. 1994 Mar; 11(5):955–964.
- 584 Conway T, Creecy JP, Maddox SM, Grissom JE, Conkle TL, Shadid TM, Teramoto J, San Miguel P, Shimada T,
585 Ishihama A, Mori H, Wanner BL. Unprecedented High-Resolution View of Bacterial Operon Architecture
586 Revealed by RNA Sequencing. mBio. 2014 Jul; 5(4):e01442–14. doi: [10.1128/mBio.01442-14](https://doi.org/10.1128/mBio.01442-14).
- 587 Cox J, Mann M. MaxQuant enables high peptide identification rates, individualized p.p.b.-range mass accuracies and proteome-wide protein quantification. Nature Biotechnology. 2008 Dec; 26(12):1367–1372. doi:
588 [10.1038/nbt.1511](https://doi.org/10.1038/nbt.1511).
- 589 Cox J, Matic I, Hilger M, Nagaraj N, Selbach M, Olsen JV, Mann M. A practical guide to the MaxQuant computational platform for SILAC-based quantitative proteomics. Nature Protocols. 2009 May; 4(5):698–705. doi:
590 [10.1038/nprot.2009.36](https://doi.org/10.1038/nprot.2009.36).
- 591 Fulco CP, Nasser J, Jones TR, Munson G, Bergman DT, Subramanian V, Grossman SR, Anyoha R, Doughty BR,
592 Patwardhan TA, Nguyen TH, Kane M, Perez EM, Durand NC, Lareau CA, Stamenova EK, Aiden EL, Lander
593 ES, Engreitz JM. Activity-by-contact model of enhancer-promoter regulation from thousands of CRISPR
594 perturbations. Nat Genet. 2019; 51(12):1664–1669.

- 597 **Gama-Castro S**, Salgado H, Santos-Zavaleta A, Ledezma-Tejeida D, Muñiz-Rascado L, García-Sotelo JS, Alquicira-
598 Hernández K, Martínez-Flores I, Pannier L, Castro-Mondragón JA, Medina-Rivera A, Solano-Lira H, Bonavides-
599 Martínez C, Pérez-Rueda E, Alquicira-Hernández S, Porrón-Sotelo L, López-Fuentes A, Hernández-Koutoucheva
600 A, Del Moral-Chavez V, Rinaldi F, et al. RegulonDB version 9.0: High-level integration of gene regula-
601 tion, coexpression, motif clustering and beyond. *Nucleic Acids Research*. 2016; 44(D1):D133–D143. doi:
602 10.1093/nar/gkv1156.
- 603 **Gao Y**, Yurkovich JT, Seo SW, Kabimoldayev I, Chen K, Sastry AV, Fang X, Mih N, Yang L, Eichner J, Cho Bk, Kim D,
604 Palsson BO. Systematic discovery of uncharacterized transcription factors in *Escherichia coli* K-12 MG1655.
605 *Nucleic Acids Research*. 2018; 46(20):10682–10696. doi: 10.1093/nar/gky752.
- 606 **Garcia HG**, Phillips R. Quantitative dissection of the simple repression input-output function. *Proceedings of the*
607 *National Academy of Sciences*. 2011 Jul; 108(29):12173–12178. doi: 10.1073/pnas.1015616108.
- 608 **Ghatak S**, King ZA, Sastry A, Palsson BO. The y-ome defines the 35% of *Escherichia coli* genes that lack
609 experimental evidence of function. *Nucleic Acids Research*. 2019; 47(5):2446–2454. doi: 10.1093/nar/gkz030.
- 610 **Goodall ECA**, Robinson A, Johnston IG, Jabbari S, Turner KA, Cunningham AF, Lund PA, Cole JA, Henderson IR.
611 The Essential Genome of *Escherichia coli* K-12. *mBio*. 2018; 9(1).
- 612 **Gupta S**, Stamatoyannopoulos JA, Bailey TL, Noble WS. Quantifying similarity between motifs. *Genome Biology*.
613 2007; 8(2). doi: 10.1186/gb-2007-8-2-r24.
- 614 **Jacob F**, Monod J. On the Regulation of Gene Activity. . 1961; p. 19.
- 615 **Keseler IM**, Mackie A, Santos-Zavaleta A, Billington R, Bonavides-Martínez C, Caspi R, Fulcher C, Gama-Castro
616 S, Kothari A, Krummenacker M, Latendresse M, Muñiz-Rascado L, Ong Q, Paley S, Peralta-Gil M, Subhraveti
617 P, Velázquez-Ramírez DA, Weaver D, Collado-Vides J, Paulsen I, et al. The EcoCyc database: reflecting new
618 knowledge about *Escherichia coli* K-12. *Nucleic Acids Research*. 2016; 45(D1):D543–D550. <https://doi.org/10.1093/nar/gkw1003>, doi: 10.1093/nar/gkw1003.
- 620 **Kinney JB**, McCandlish DM. Massively Parallel Assays and Quantitative Sequence–Function Relationships. *Annual*
621 *Review of Genomics and Human Genetics*. 2019; 20(1):99–127. doi: 10.1146/annurev-genom-083118-014845.
- 622 **Kinney JB**, Murugan A, Callan CG, Cox EC. Using deep sequencing to characterize the biophysical mechanism of
623 a transcriptional regulatory sequence. *Proceedings of the National Academy of Sciences of the United States*
624 *of America*. 2010; 107(20):9158–9163. doi: 10.1073/pnas.1004290107.
- 625 **Kosuri S**, Goodman DB, Cambray G, Mutualik VK, Gao Y. Composability of regulatory sequences controlling
626 transcription and translation in *Escherichia coli*. *Proceedings of the National Academy of Sciences of the*
627 *United States of America*. 2013; 110(34). doi: 10.1073/pnas.1301301110.
- 628 **Larsons TJ**, Ye S, Weissenborn DL, Hoffmann HJ. Purification and Characterization of the Repressor for the
629 sn-Glycerol 3-Phosphate Regulon of *Escherichia coli* K12. *Journal of Biological Chemistry*. 1987; 262(33):15869–
630 15874.
- 631 **Lin ECC**. Glycerol Dissimilation and its Regulation in Bacteria. . 1976; p. 44.
- 632 **Lindsley JE**, Rutter J. Whence cometh the allosterome? *Proceedings of the National Academy of Sciences of the*
633 *United States of America*. 2006; 103(28):10533–10535. doi: 10.1073/pnas.0604452103.
- 634 **Magoč T**, Salzberg SL. FLASH: fast length adjustment of short reads to improve genome assemblies. *Bioinfor-*
635 *matics*. 2011 Nov; 27(21):2957–2963. doi: 10.1093/bioinformatics/btr507.
- 636 **Melnikov A**, Murugan A, Zhang X, Tesileanu T, Wang L, Rogov P, Feizi S, Gnrke A, Jr CGC, Kinney JB, Kellis M,
637 Lander ES, Mikkelsen TS. Systematic dissection and optimization of inducible enhancers in human cells using
638 a massively parallel reporter assay. *Nature Biotechnology*. 2012; 30(3):271–277. doi: 10.1038/nbt.2137.
- 639 **Mortazavi A**, Williams BA, McCue K, Schaeffer L, Wold B. Mapping and quantifying mammalian transcriptomes
640 by RNA-Seq. *Nature Methods*. 2008; 5(7):621–628. doi: 10.1038/nmeth.1226.
- 641 **Myers KS**, Yan H, Ong IM, Chung D, Liang K, Tran F, Keleş S, Landick R, Kiley PJ. Genome-scale Analysis of
642 *Escherichia coli* FNR Reveals Complex Features of Transcription Factor Binding. *PLOS Genetics*. 2013 06;
643 9(6):1–24. <https://doi.org/10.1371/journal.pgen.1003565>, doi: 10.1371/journal.pgen.1003565.
- 644 **Pappireddi N**, Martin L, Wühr M. A Review on Quantitative Multiplexed Proteomics. *ChemBioChem*. 2019;
645 20(10):1210–1224. doi: 10.1002/cbic.201800650.

- 646 **Patwardhan RP**, Hiatt JB, Witten DM, Kim MJ, Smith RP, May D, Lee C, Andrie JM, Lee SI, Cooper GM, Ahituv N,
647 Pennacchio LA, Shendure J. Massively parallel functional dissection of mammalian enhancers *in vivo*. *Nat Biotechnol.* 2012; 30(3):265–70.
- 649 **Razo-Mejia M**, Barnes SL, Belliveau NM, Chure G, Einav T, Lewis M, Phillips R. Tuning Transcriptional Regulation
650 through Signaling: A Predictive Theory of Allosteric Induction. *Cell Systems*. 2018; 6(4):456–469.e10. doi:
651 [10.1016/j.cels.2018.02.004](https://doi.org/10.1016/j.cels.2018.02.004).
- 652 **Santos-Zavaleta A**, Salgado H, Gama-castro S, Laura G, Ledezma-tejeida D, Mishael S, Garc S, Alquicira-hern K,
653 Jos L, Pe P, Ishida-guti C, Vel DA, Moral-ch D, Galagan J, Collado-vides J. RegulonDB v 10 . 5 : tackling challenges
654 to unify classic and high throughput knowledge of gene regulation in *E. coli* K-12. *Nucleic Acids Research*.
655 2019; 47(November 2018):212–220. doi: 10.1093/nar/gky1077.
- 656 **Schmidt A**, Kochanowski K, Vedelaar S, Ahrne E, Volkmer B, Callipo L, Knoops K, Bauer M, Aebersold R, Heine-
657 mann M. The quantitative and condition-dependent *Escherichia coli* proteome. *Nature biotechnology*. 2015;
658 34(1):104–110. doi: [10.1038/nbt.3418](https://doi.org/10.1038/nbt.3418).
- 659 **Schneider TD**, Stephens RM. Sequence logos: a new way to display consensus sequences. *Nucleic Acids Research*. 1990; 18(20):6097–6100. doi: [10.1093/nar/18.20.6097](https://doi.org/10.1093/nar/18.20.6097).
- 660 **Seoh HK**, Tai PC. Catabolic repression of secB expression is positively controlled by cyclic AMP (cAMP) receptor
661 protein-cAMP complexes at the transcriptional level. *Journal of Bacteriology*. 1999 Mar; 181(6):1892–1899.
- 663 **Sharon E**, Kalma Y, Sharp A, Raveh-Sadka T, Levo M, Zeevi D, Keren L, Yakhini Z, Weinberger A, Segal E. Infer-
664 ring gene regulatory logic from high-throughput measurements of thousands of systematically designed
665 promoters. *Nat Biotechnol*. 2012; 30(6):521–30.
- 666 **Stuart T**, Satija R. Integrative single-cell analysis. *Nature Reviews Genetics*. 2019; 20(May):257–272. doi:
667 [10.1038/s41576-019-0093-7](https://doi.org/10.1038/s41576-019-0093-7).
- 668 **Urtecho G**, Tripp AD, Insigne KD, Kim H, Kosuri S. Systematic Dissection of Sequence Elements Controlling
669 sigma70 Promoters Using a Genomically Encoded Multiplexed Reporter Assay in *Escherichia coli*. *Biochemistry*.
670 2019; 58(11):1539–1551.
- 671 **Urtecho G**, Insigne K, Tripp AD, Brinck M, Lubock NB, Kim H, Chan T, Kosuri S. Genome-wide Functional Charac-
672 terization of *Escherichia col* Promoters and Regulatory Elements Responsible for their Function. *bioRxiv*. 2020;
673 <https://www.biorxiv.org/content/early/2020/01/06/2020.01.04.894907>, doi: [10.1101/2020.01.04.894907](https://doi.org/10.1101/2020.01.04.894907).
- 674 **Yamamoto N**, Nakahigashi K, Nakamichi T, Yoshino M, Takai Y, Touda Y, Furubayashi A, Kinjyo S, Dose H,
675 Hasegawa M, Datsenko KA, Nakayashiki T, Tomita M, Wanner BL, Mori H. Update on the Keio collec-
676 tion of *Escherichia coli* single-gene deletion mutants. *Molecular systems biology*. 2009; 5:335–335. doi:
677 [10.1038/msb.2009.92](https://doi.org/10.1038/msb.2009.92).

679

Comparison of results for known binding sites

680

681

682

683

684

685

686

687

688

689

690

691

692

693

694

695

696

697

698

699

700

701

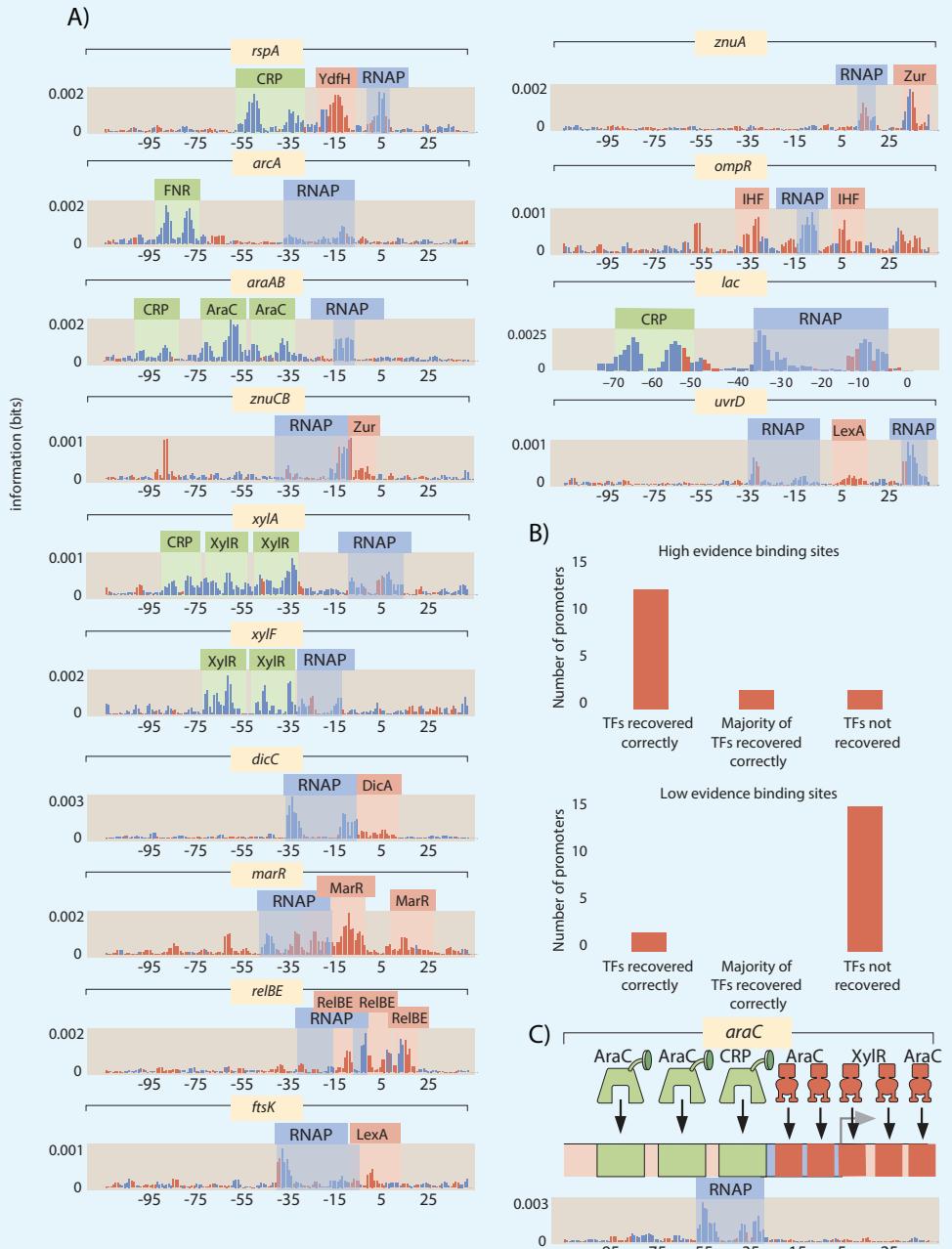
702

703

704

The work presented here is effectively a third-generation of the use of Sort-Seq methods for the discovery of regulatory architecture. The primary difference between the present work and previous generations is the use of RNA-Seq rather than fluorescence and cell sorting as a readout of the level of expression of our promoter libraries. As such, there are many important questions to be asked about the comparison between the earlier methods and this work. We attack that question in several ways. First, as shown in Appendix 1 Figure 2, we have performed a head-to-head comparison of the two approaches to be described further in this section. Second, as shown in the next section, our list of candidate promoters included roughly 20% for which the last 50 years of molecular biology has resulted in a knowledge (sometimes only partial) of their regulatory architecture. In these cases, we examined the extent to which our methods recover the known features of regulatory control about those promoters.

We have tested over 20 genes for which there is already some substantial regulatory knowledge reported in the literature. The successes and failures of this test are detailed in Appendix 1 Figure 1. For those promoters which have strong evidence of a binding site, as determined by RegulonDB, we recover all relevant transcription factor binding sites for 12 out of 16 cases, the majority of relevant binding sites for 2 out of 16 cases, and miss all or most of the regulation for 2 promoters. We identify a total of 22 previously known high evidence binding sites.



705
706
707
708
709
710

Appendix 1 Figure 1. An analysis of binding sites in the literature. (A) Information footprints for known and properly recovered binding sites. (B) A summary of how well the SortSeq results conform to literature results. The sites that are low evidence in the literature are determined by RegulonDB. (C) The information footprint and known binding sites for the *araC* promoter. Despite all the binding sites present, the only binding signature that appears is for a RNAP.

These results showcase that our method largely agrees with the established literature but also highlights several areas in which our method is prone to missing regulatory elements. One failure mode is caused by the presence of strong secondary binding sites. For example, in the *araC* promoter, as shown in Appendix 1 Figure 1(C), the only binding signatures that appear in the information footprint are from a secondary RNAP site. The secondary site seems to be expressed constitutively, and in the cases where the primary start site is even

718 partially repressed, the secondary start site will dominate transcription and obscure the
719 many binding sites that are in this promoter.

720
721 If there are large numbers of regulatory elements, the data will often only show the few
722 most important elements. If we look at the *mar* promoter in Appendix 1 Figure 1 C), we can
723 only see the signature of the two MarR sites even though CpxR, Fis, and CRP are all known to
724 bind to the promoter. MarR is a strong enough repressor that mutating any of the other
725 transcription factor sites is unlikely to meaningfully change gene expression unless the MarR
726 site is also mutated. This illustrates that the regulatory architectures discovered in this study
727 represent a lower bound on what exists in each promoter.
728

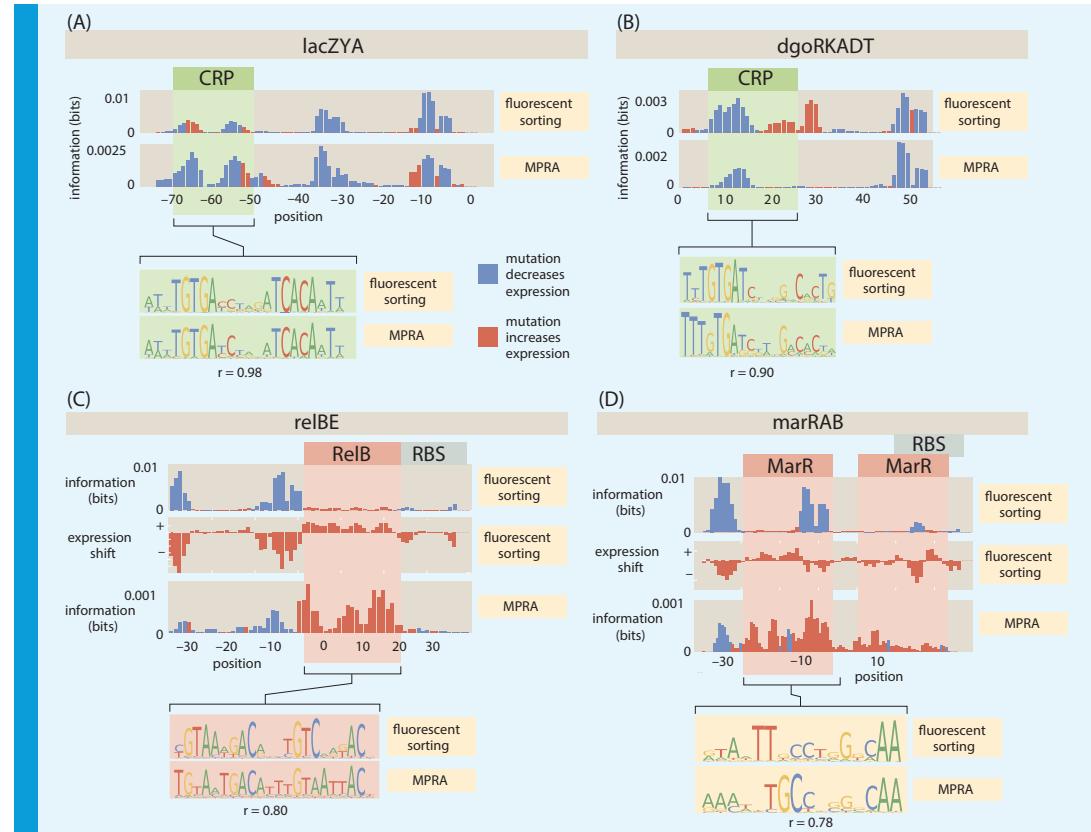
729 Finally, for some genes such as *dicA* there was no known TSS prior to the experiment.
730 Although there is a small regulatory region between *dicA* and its neighboring gene, this does
731 not insure that we will include the strongest RNAP sites. Better mapping of transcription
732 start sites could improve our method.
733

734 We next consider low evidence binding sites. Other research determined the locations of
735 the low evidence sites through gene expression analysis and sequence comparison to con-
736 sensus sequences. For 5 promoters in our list, the binding sites location itself is not known,
737 only that the TF in question regulates the gene. For these promoters we recover the known
738 regulation in only 2 out of 15 cases. Comparison to consensus sequences can be unreliable
739 and generate false positives when the entirety of the *E. Coli* genome is considered. Gene
740 expression analysis alone has difficulty ruling out indirect effects of a given transcription
741 factor on gene expression and regulation determined by this method may occur outside of
742 the 160 bp mutation window we consider. As our results recover high evidence sites well,
743 the poor recovery of sites based on sequence gazing and gene expression analysis most
744 likely indicates that these methods are unreliable for determining binding locations.
745

746 **Comparison between Reg-Seq by RNA-Seq and fluorescent sorting**

747 As the basis for comparing the results of the fluorescence-based Sort-Seq approach with our
748 RNA-Seq-based approach, we use information footprints, expression shifts and sequence
749 logos as our metrics. Appendix 1 Figure 2 shows examples of this comparison for four
750 distinct genes of interest. Appendix 1 Figure 2(A) shows the results of the two methods for
751 the *lacZYA* promoter with special reference to the CRP binding site. Both the information
752 footprint and the sequence logo identify the same binding site.
753

754 Appendix 1 Figure 2(B) provides a similar analysis for the *dgoRKADT* promoter where
755 once again the information footprints and the sequence logos from the two methods are in
756 reasonable accord. Appendix 1 Figure 2(C) provides a quantitative dissection of the *relBE*
757 promoter which is repressed by RelBE. Here we use both information footprints and expres-
758 sion shifts as a way to quantify the significance of mutations to different binding sites across
759 the promoter. Finally, Appendix 1 Figure 2(D) shows a comparison of the two methods for
760 the *marRAB* promoter. The two approaches both identify a MarR binding site.
761



762
763 **Appendix 1 Figure 2.** A summary of four direct comparisons of measurements using fluorescence and
764 sorting and using RNA-Seq. (A) CRP binds upstream of RNAP in the *lacZYA* promoter. Despite the
765 different measurement techniques the two inferred energy matrices, for the CRP binding site have a
766 Pearson correlation coefficient $r = 0.98$. (B) The *dgoRKADT* promoter is activated by CRP in the presence
767 of galactonate. The FACS measurements were taken in the JK10 strain in the presence of 500mM cAMP.
768 In both cases, a type II activator binding site can be identified based on the signals in the information
769 footprint and expression shift plot in the area indicated in orange. Additionally the quantitative
770 agreement between the CRP binding preference matrices are strong, with $r = 0.9$. (C) The *relBE* promoter
771 is repressed by RelBE. The inferred matrices between the two measurement methods have $r = 0.8$. (D)
772 The *marRAB* promoter is repressed by MarR. The features we can observe in the information footprint
773 reflect this under measurement with both FACS or RNAseq. The inferred binding matrices shown have
774 $r = 0.78$. The right most MarR site overlaps with a ribosome binding site. The overlap has a stronger
775 obscuring effect on the sequence specificity of the FACS measurement, which measures protein levels
776 directly, than it does on the output of the RNAseq measurement.
777

778 We note that the first aim of our methods is regulatory discovery. We would like to be
779 able to determine how previously uncharacterized promoters are regulated and ultimately,
780 this is a question of binding-site and transcription factor identification. For that task, we
781 do not require perfect correspondence between the two methods. With regulatory sites
782 identified, our next objective is the determination of energy matrices that will allow us to
783 turn binding site strength into a tunable knob that can nearly continuously tune the strength
784 of transcription factor binding, thus altering gene expression in predictable ways as already
785 shown in our earlier work (**Barnes et al., 2019**). The r-values between energy matrices range
786 from 0.78 to 0.96, indicating reasonable to very good agreement. Reg-Seq appears to be, if
787 anything, more accurate than previous methods as it has higher relative information content
788 in known areas of transcription factor binding and also does not have repressor-like bases
789 on CRP sites as in Appendix 1 Figure 2 (A) and (B).

790

792
793 **Extended details of analysis methods**794 **Information footprints**

795 We use information footprints as a tool for hypothesis generation to identify regions which
 796 may contain transcription factor binding sites. In general, a mutation within a transcription
 797 factor site is likely to severely weaken that site. We look for groups of positions where
 798 mutation away from wild type has a large effect on gene expression. Our data sets consist
 799 of nucleotide sequences, the number of times we sequenced the construct in the plasmid
 800 library, and the number of times we sequenced its corresponding mRNA. A simplified data
 801 set on a 4 nucleotide sequence then might look like

Sequence	Library Sequencing Counts	mRNA Counts
ACTA	5	23
ATTA	5	3
CCTG	11	11
TAGA	12	3
GTGC	2	0
CACA	8	7
AGGC	7	3

802 One possible calculation is to take all sequences which have base b at position i and
 803 determine the number of mRNAs produced per read in the sequencing library. By comparing
 804 the values for different bases we could determine how large of an effect mutation has on
 805 gene expression. However, in this paper we will use mutual information to quantify the
 806 effect of mutation, as [Kinney et al. \(2010\)](#) demonstrated could be done successfully. In
 807 Table 1 the frequency of the different nucleotides in the library at position 2 is 40% A, 32% C,
 808 14% G and 14% T. Cytosine is enriched in the mRNA transcripts over the original library, as it
 809 now composes 68% of all mRNA sequencing reads while A, G, and T only compose only 20%,
 810 6%, and 6% respectively. Large enrichment of some bases over others occurs when base
 811 identity is important for gene expression. We can quantify how important using the mutual
 812 information between base identity and gene expression level. Mutual information is given at
 813 position i by
 814
 815

$$I_b = \sum_{m,exp} p(m,exp) \log_2 \left(\frac{p(m,exp)}{p(m)p(exp)} \right) \quad (3)$$

816 $p(b)$ in equation 3 refers to the probability that a given sequencing read will be from a
 817 mutated base. $p(exp)$ is a normalizing factor that gives the ratio of the number of DNA or
 818 mRNA sequencing counts to total number of counts.

819 The mutual information quantifies how much a piece of knowledge reduces the entropy
 820 of a distribution. At a position where base identity matters little for expression level, there
 821 would be little difference in the frequency distributions for the library and mRNA transcripts.
 822 The entropy of the distribution would decrease only by a small amount when considering
 823 the two types of sequencing reads separately.

824 We are interested in quantifying the degree to which mutation away from a wild type
 825 sequence affects expression. Although there are obviously 4 possible nucleotides, we can
 826 classify each base as either wild-type or mutated so that b in equation 3 represents only

828
829
830
831
832

these two possibilities.

833
834
835
836
837
838
839
840
841
842

If mutations at each position are not fully independent, then the information value calculated in equation 1 will also encode the effect of mutation at correlated positions. If having a mutation at position 1 is highly favorable for gene expression and is also correlated with having a mutation at position 2, mutations at position 2 will also be enriched amongst the mRNA transcripts. Position 2 will appear to have high mutual information even if it has minimal effect on gene expression. Due to the DNA synthesis process used in library construction, mutation in one position can make mutation at other positions more likely by up to 10 percent. This is enough to cloud the signature of most transcription factors in an information footprint calculated using equation 1.

843
845
846
847
848

We need to determine values for $p_i(m|exp)$ when mutations are independent, and to do this we need to fit these quantities from our data. We assert that

$$\langle mRNA \rangle \propto e^{-\beta E_{eff}} \quad (4)$$

849
850
851
852
853
854
855
856
857

is a reasonable approximation to make. $\langle mRNA \rangle$ is the average number of mRNAs produced by that sequence for every cell containing the construct and E_{eff} is an effective energy for the sequence that can be determined by summing contributions from each position in the sequence. There are many possible underlying regulatory architectures, but to demonstrate that our approach is reasonable let us first consider the simple case where there is only a RNAP site in the studied region. We can write down an expression for average gene expression per cell as

858
859

$$\langle mRNA \rangle \propto p_{bound} \propto \frac{\frac{p}{N_{NS}} e^{-\beta E_p}}{1 + \frac{p}{N_{NS}} e^{-\beta E_p}} \quad (5)$$

860
861
862
863
864
865

Where p_{bound} is the probability that the RNAP is bound to DNA and is known to be proportional to gene expression in *E. coli* (**Garcia and Phillips, 2011**), E_p is the energy of RNAP binding, N_{NS} is the number of nonspecific DNA binding sites, and p is the number of RNAP. If RNAP binds weakly then $\frac{p}{N_{NS}} e^{-\beta E_p} \ll 1$. We can simplify equation 5 to

866
867

$$\langle mRNA \rangle \propto e^{-\beta E_p}. \quad (6)$$

868
869
870

If we assume that the energy of RNAP binding will be a sum of contributions from each of the positions within its binding site then we can calculate the difference in gene expression between having a mutated base at position i and having a wild type base as

871
872
873
874

$$\frac{\langle mRNA_{WT_i} \rangle}{\langle mRNA_{Mut_i} \rangle} = \frac{e^{-\beta E_{p_{WT_i}}}}{e^{-\beta E_{p_{Mut_i}}}} \quad (7)$$

$$\frac{\langle mRNA_{WT_i} \rangle}{\langle mRNA_{Mut_i} \rangle} = e^{-\beta(E_{p_{WT_i}} - E_{p_{Mut_i}})}. \quad (8)$$

875
876
878

In this example we are only considering single mutation in the sequence so we can further simplify the equation to

879
880

$$\frac{\langle mRNA_{WT_i} \rangle}{\langle mRNA_{Mut_i} \rangle} = e^{-\beta \Delta E_p_i}. \quad (9)$$

881
882

We can now calculate the base probabilities in the expressed sequences. If the probability of finding a wild type base at position i in the DNA library is $p_i(m = WT|exp = 0)$ then

$$p_i(m = WT|exp = 1) = \frac{p_i(m = WT|exp = 0) \frac{\langle mRNA_{WT_i} \rangle}{\langle mRNA_{Mut_i} \rangle}}{p_i(m = Mut|exp = 0) + p_i(m = WT|exp = 0) \frac{\langle mRNA_{WT_i} \rangle}{\langle mRNA_{Mut} \rangle}} \quad (10)$$

$$p_i(m = WT|exp = 1) = \frac{p_i(m = WT|exp = 0) e^{-\beta \Delta E_{P_i}}}{p_i(m = Mut|exp = 0) + p_i(m = WT|exp = 0) e^{-\beta \Delta E_{P_i}}} \quad (11)$$

Under certain conditions, we can also infer a value for $p_i(m|exp = 1)$ using a linear model when there are any number of activator or repressor binding sites. We will demonstrate this in the case of a single activator and a single repressor, although a similar analysis can be done when there are greater numbers of transcription factors. We will define $P = \frac{p}{N_{NS}} e^{-\beta E_P}$. We will also define $A = \frac{a}{N_{NS}} e^{-\beta E_A}$ where a is the number of activators, and E_A is the binding energy of the activator. We will finally define $R = \frac{r}{N_{NS}} e^{-\beta E_R}$ where r is the number of repressors and E_R is the binding energy of the repressor. We can write

$$\langle mRNA \rangle \propto p_{bound} \propto \frac{P + PAe^{-\beta \epsilon_{AP}}}{1 + A + P + R + PAe^{-\beta \epsilon_{AP}}} \quad (12)$$

If activators and RNAP bind weakly but interact strongly, and repressors bind very strongly, then we can simplify equation 12. In this case $A \ll 1$, $P \ll 1$, $PAe^{-\epsilon_{AP}} \gg P$, and $R \gg 1$. We can then rewrite equation 12 as

$$\langle mRNA \rangle \propto \frac{PAe^{-\beta \epsilon_{AP}}}{R} \quad (13)$$

$$\langle mRNA \rangle \propto e^{-\beta(-E_P - E_A + E_R)} \quad (14)$$

As we typically assume that RNAP binding energy, activator binding energy, and repressor binding can all be represented as sums of contributions from their constituent bases, the combination of the energies can be written as a total effective energy E_{eff} which is a sum of contributions from all positions within the binding sites.

We fit the parameters for each base using a Markov Chain Monte Carlo Method. Two MCMC runs are conducted using randomly generated initial conditions. We require both chains to reach the same distribution to prove the convergence of the chains. We do not wish for mutation rate to affect the information values so we set the $p(WT) = p(Mut) = 0.5$ in the information calculation. The information values are smoothed by averaging with neighboring values.

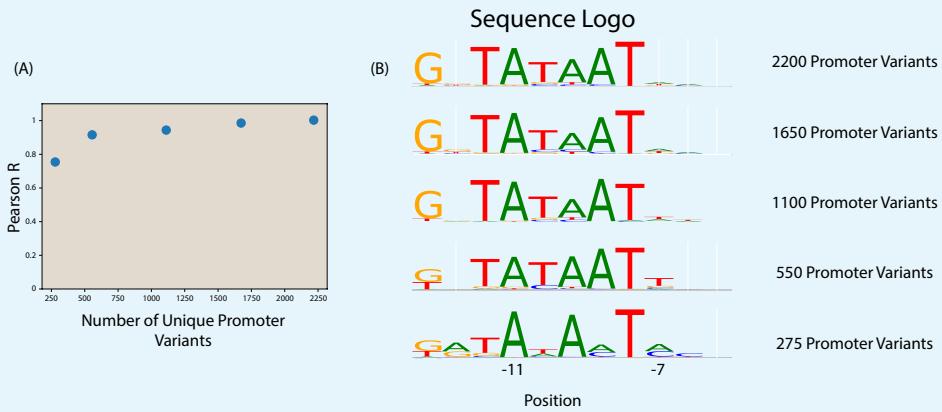
Analysis of mass spectrometry results

Mass spectrometry results were processed using MaxQuant ([Cox and Mann, 2008](#)) ([Cox et al., 2009](#)). Spectra were searched against the UniProt *E. coli* K-12 database as well as a contaminant database (256 sequences). LysC was specified as the digestion enzyme. Proteins were considered if they were known to be transcription factors, or were predicted to bind DNA (using gene ontology term GO:0003677, for DNA-binding in BioCyc). The reported binding TFs were enriched above all other other DNA binding proteins with $p < 10^{-4}$. The p-value was calculated with a two sample t-test using the python function `scipy.stats.ttest_ind_from_stats`. The uncertainty in the background DNA binding protein

921
 922
 923
 924
 925
 926 ratios was calculated using all the enrichment ratios of proteins not determined to bind. The
 927 uncertainty in the enrichment ratio of the binding proteins was calculated from the spread in
 928 their enrichment ratios across the many (≈ 10) mass spec runs where the protein in question
 929 was not an active TF. Any dataset where the protein in question is not identified is excluded
 930 from the analysis.

931 **Uncertainty due to number of independent sequences**
 932

933 1400 promoter variants were ordered from TWIST Bioscience for each promoter studied.
 934 Due to errors in synthesis, the final number of variants received was an average of 2200
 935 per promoter. To test whether the number of promoter variants is a significant source of
 936 uncertainty in the experiment we computationally reduced the number of promoter variants
 937 used in the analysis of the *zapAB*-10 RNAP region. Each sub-sampling was performed 3
 938 times. The results, as displayed in Appendix 2 Figure 1, show that there is only a small effect
 939 on the resulting sequence logo until the library has been reduced to approximately 500
 940 promoter variants.



941 **Appendix 2 Figure 1.** A comparison of RNAP -10 site sequence logos. (A) This figure shows the Pearson
 942 correlation coefficient between the energy matrix models inferred from the full dataset (2200 unique
 943 promoter variants) and that from a computationally restricted dataset. (B) Sequence logos of the RNAP
 944 -10 region from each sub-sampled dataset.
 945

946 **TOMTOM motif comparison**
 947

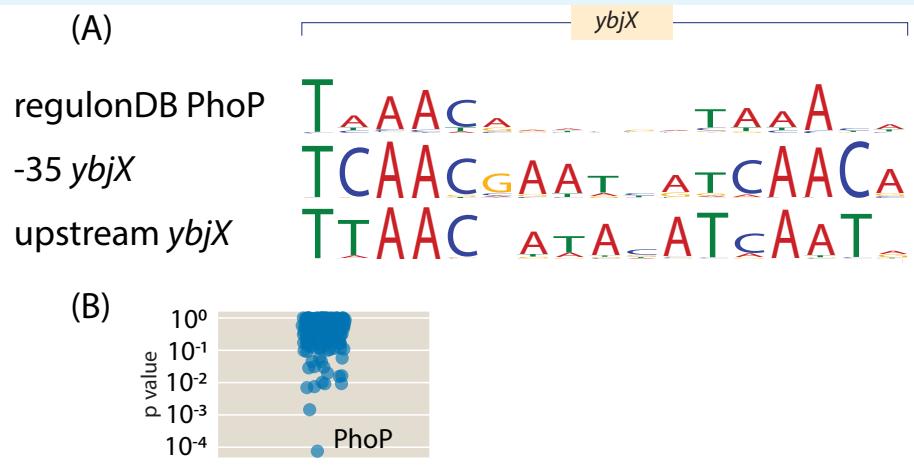
948 We used an alternative approach to discover the TF identity regulating a given promoter,
 949 using a motif comparison tool. TOMTOM (*Gupta et al., 2007*) is a tool that uses a statisti-
 950 cal method to infer if a putative motif resembles to any previously discovered motif in a
 951 database. Of interest, it accounts for all possible offsets between the motifs. Moreover,
 952 it uses a suite of metrics to compare between motifs such as Kullback-Leibler divergence,
 953 Pearson correlation, euclidean distance, among others.

954 We performed comparisons of the motifs generated from our energy matrices to those
 955 generated from all known transcription factor binding sites in RegulonDB. Appendix 2
 956 Figure 2 shows a result of TOMTOM, where we compared the motif derived from the -35
 957 region of the *ybjX* promoter and found a good match with the motif of PhoP from RegulonDB.
 958

959 The information derived from this approach was then used to guide some of the TF knock-
 960 out experiments, in order to validate its interaction with a target promoter characterized
 961 by the loss of the information footprint. Furthermore, we also used TOMTOM to search for

960
961
962
963
964
965

similarities between our own database of motifs, in order to generate regulatory hypotheses in tandem. This was particularly useful when looking at the group of GlpR binding sites found in this experiment.



966
967
968
969
970

Appendix 2 Figure 2. Motif comparison using TOMTOM. Searching our energy motifs against the RegulonDB database using TOMTOM allowed us to guide our TF knockout experiments. Here we show the sequence logos of the PhoP transcription factor from RegulonDB (top) and the one generated from the *ybjX* promoter energy matrix. E-value = 0.01 using Euclidean distance as a similarity matrix.

973
974 **Extended details of experimental design**975 **Choosing target genes**

976 Genes in this experiment were chosen to cover several different categories. 29 genes had
 977 some information on their regulation already known to validate our method under a number
 978 of conditions. 37 were chosen because the work of *Schmidt et al. (2015)* demonstrated
 979 that gene expression changed significantly under different growth conditions. A handful of
 980 genes such as *minC*, *maoP*, or *fdhE* were chosen because we found either their physiological
 981 significance interesting, as in the case of the cell division gene *minC* or that we found the gene
 982 regulatory question interesting, such for the intra-operon regulation demonstrated by *fdhE*.
 983 The remainder of the genes were chosen because they had no regulatory information, often
 984 had minimal information about the function of the gene, and had an annotated transcription
 start site (TSS) in RegulonDB.

985 **Choosing transcription start sites**

986 A known limitation of the experiment is that the mutational window is limited to 160 bp. As
 987 such, it is important to correctly target the mutation window to the location around the most
 988 active TSS. To do this we first prioritized those TSS which have been extensively experimen-
 989 tally validated and catalogued in RegulonDB. Secondly we selected those sites which had
 990 evidence of active transcription from RACE experiments and were listed in RegulonDB. If the
 991 intergenic region was small enough, we covered the entire region with our mutation window.
 992 If none of these options were available, we used computationally predicted start sites.

993 **Sequencing**

994 All sequencing was carried out by either the Millard and Muriel Jacobs Genetics and Ge-
 995 nomics Laboratory at Caltech (HiSeq 2500) on a 100 bp single read flow cell or using the
 996 sequencing services from NGX Bio on a 250 bp or 150 base paired end flow cell. The total
 997 library was first sequenced by PCR amplifying the region containing the variant promoters
 998 as well as the corresponding barcodes. This allowed us to uniquely associate each random
 999 20 bp tag with a promoter variant. Any tag which was associated with a promoter variant
 1000 with insertions or deletions was removed from further analysis. Similarly, any tag that was
 1001 associated with multiple promoter variants was also removed from the analysis. The paired
 1002 end reads from this sequencing step were then assembled using the FLASH tool *Magoč and*
 1003 *Salzberg (2011)*. Any sequence with PHRED score less than 20 was removed using the FastX
 1004 toolkit. Additionally, when sequencing the initial library, sequences which only appear in the
 1005 dataset once were not included in further analysis in order to remove possible sequencing
 1006 errors.

1007 For all the MPRA experiments, only the region containing the random 20 bp tag was
 1008 sequenced, since the tag can be matched to a specific promoter variant using the initial
 1009 library sequencing run described above. For a given growth condition, each promoter yielded
 1010 50,000 to 500,000 usable sequencing reads. Under some growth conditions, genes were not
 1011 analyzed further if they did not have at least 50,000 reads.

1012 To determine which base pair regions were statistically significant a 99 % confidence
 1013 interval was constructed using the MCMC inference to determine the uncertainty.

1016
1017
1018
1019
1020
1021
1022
1023
1024
1025
1026
1027
1028
1029
1030
1031
1032
1033
1034
1035
1036
1037
1038
1039
1040

Growth Conditions

The growth conditions studied in this experiment include differing carbon sources such as growth in M9 with 0.5% Glucose, M9 with acetate (0.5%), M9 with arabinose (0.5%), M9 with Xylose (0.5%) and arabinose (0.5%), M9 with succinate (0.5%), M9 with fumarate (0.5%), M9 with Trehalose (0.5%), and LB. In each case cell harvesting was done at an OD of 0.3. These growth conditions were chosen so as to span a wide range of growth rates, as well as to illuminate any carbon source specific regulators.

We also used several stress conditions such as heat shock, where cells were grown in M9 and were subjected to a heat shock of 42 degrees for 5 minutes before harvesting RNA. We grew in low oxygen conditions. Cells were grown in LB in a container with minimal oxygen, although some will be present as no anaerobic chamber was used. This level of oxygen stress was still sufficient to activate FNR binding, and so activated the anaerobic metabolism. We also grew cells in M9 with Glucose and 5mM sodium selenite.

Growth with zinc was performed at a concentration of 5mM ZnCl₂ and growth with iron was performed by first growing cells to an OD of 0.3 and then adding FeCl₂ to a concentration of 5mM and harvesting RNA after 10 minutes. Growth without cAMP was accomplished by the use of the JK10 strain which does not maintain its cAMP levels.

All knockout experiments were performed in M9 with Glucose except for the knockouts for *arcA*, *hdfR*, and *phoP* which were grown in LB.

Additional discussion of results

Binding sites regulating divergent operons



Appendix 3 Figure 1. A figure displaying two cases in which we see transcription factor binding sites that we have found to regulate both of the two divergently transcribed genes.

In addition to discovering new binding sites, we have discovered additional functions of known binding sites. In particular, in the case of *bdcR*, the repressor for the divergently transcribed gene *bdcA*, is also shown to repress *bdcR* in Appendix 3 Figure 2 (A). Similarly in Appendix 3 Figure 2 (B) IlvY is shown to repress *ilvC* in the absence of inducer. Divergently transcribed operons that share regulatory regions are plentiful in *E. coli*, and although there are already many known examples of transcription factor binding sites regulating several different operons, there are almost certainly many examples of this type of transcription that have yet to be discovered.

Multi-purpose binding sites allow for more genes to be regulated with fewer binding sites. However, they can also serve to sharpen the promoter's response to environmental cues. In

1054

1055

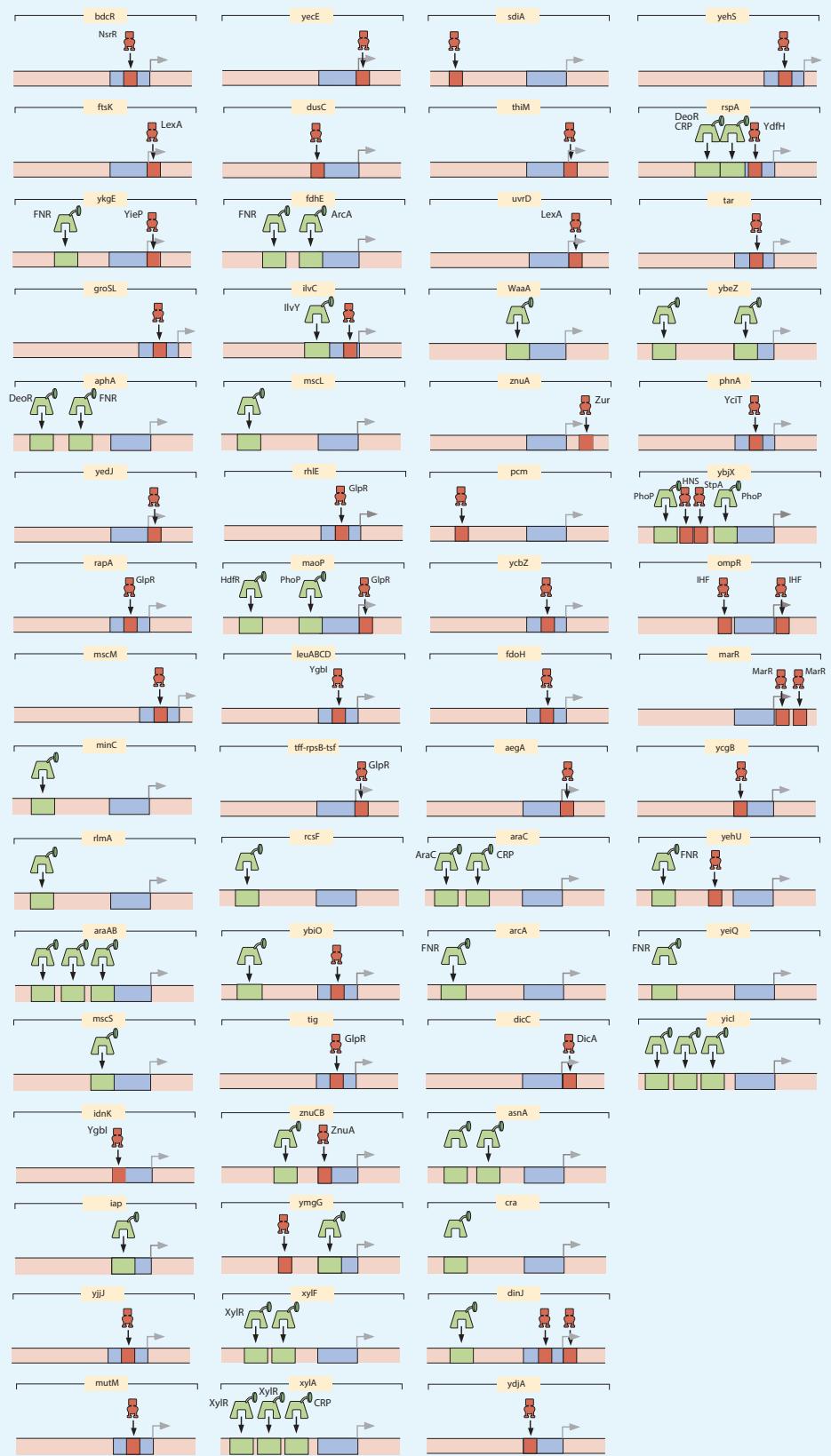
1056 the case of *ilvC*, IlvY is known to activate *ilvC* in the presence of inducer. However, we now
1057 see that it also represses the promoter in the absence of that inducer. The production of
1058 *ilvC* is known to increase by approximately a factor of 100 in the presence of inducer. The
1059 magnitude of the change is attributed to the cooperative binding of two IlvY binding sites,
1060 but the lowered expression of the promoter due to IlvY repression in the absence of inducer
1061 is also a factor.

1062



Further results.

Regulatory cartoons



1068

Appendix 4 Figure 1. All regulatory cartoons for genes studied in this experiment.

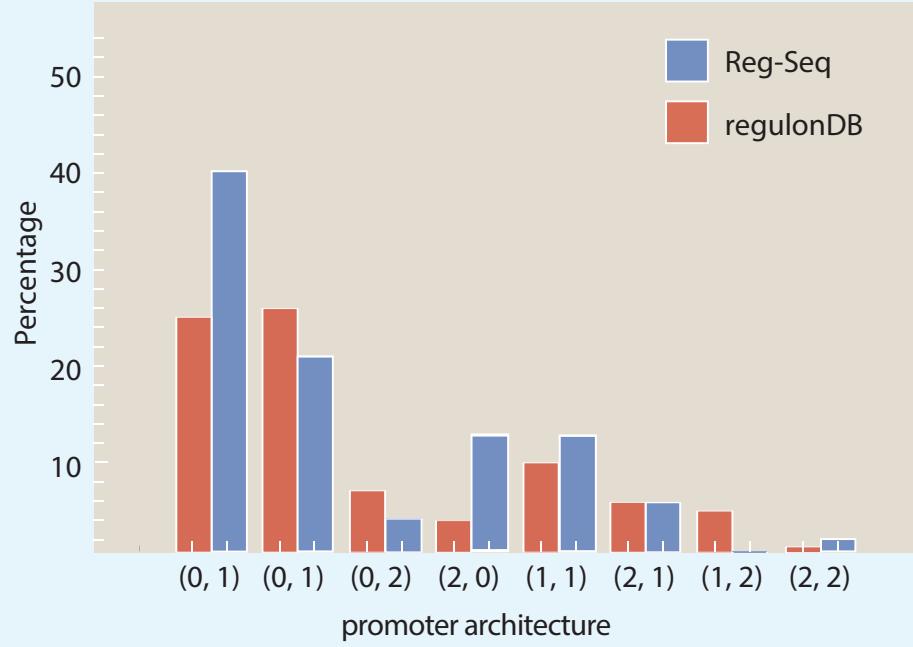
1069

Comparison of results to regulonDB

1070

1071

1073

**Appendix 4 Figure 2.** A comparison of the types of architectures found in regulonDB to the architectures with newly discovered binding sites found in the Reg-Seq experiment.