

1 **Title:** Genome-wide Functional Characterization of Escherichia coli Promoters and
2 Sequence Elements Encoding Their Regulation

3

4 **Authors:**

5 Guillaume Urtecho^{1,2*}, Kimberly D. Insigne^{3†}, Arielle D. Tripp⁴, Marcia S. Brinck⁵, Nathan B. Lubock⁶,
6 Christopher Acree⁷, Hwangbeom Kim⁶, Tracey Chan³, & Sriram Kosuri^{6,8,9}

7

8 **Affiliations:**

9 ¹ Molecular Biology Interdepartmental Doctoral Program, University of California, Los Angeles,
10 CA, 90095, USA

11 ² Department of Systems Biology, Columbia University, New York, NY, 10032, USA.

12 ³ Bioinformatics Interdepartmental Graduate Program, University of California, Los Angeles, CA 90095,
13 USA

14 ⁴ Department of Molecular, Cell, and Developmental Biology, University of California, Los
15 Angeles, CA, 90095, USA

16 ⁵ Department of Microbiology, Immunology, and Molecular Genetics, University of California, Los Angeles,
17 Los Angeles, California, USA

18 ⁶ Department of Chemistry and Biochemistry, University of California, Los Angeles, CA 90095, USA

19 ⁷ Department of Molecular Physiology and Biophysics, Vanderbilt University, Nashville, TN, 37232 USA

20 ⁸ Molecular Biology Institute, University of California, Los Angeles, CA 90095, USA

21 ⁹ UCLA-DOE Institute for Genomics and Proteomics, Quantitative and Computational Biology Institute, Eli
22 and Edythe Broad Center of Regenerative Medicine and Stem Cell Research, Jonsson Comprehensive
23 Cancer Center, University of California, Los Angeles, CA 90095, USA

24

25 *To whom correspondence should be addressed. Email: gu2144@cumc.columbia.edu

26

27 †The authors wish it to be known that, in their opinion, the first two authors should be regarded as joint
28 First Authors.

29

30 **Summary:**

31

32 Despite decades of intense genetic, biochemical, and evolutionary characterizations of bacterial
33 promoters, we lack the ability to identify or predict transcriptional activities of promoters using
34 primary sequence. Even in simple, well-characterized organisms such as *E. coli* there is little
35 agreement on the number, location, and strength of promoters. We use a genomically-encoded
36 massively parallel reporter assay to perform the first full characterization of autonomous
37 promoter activity across the *E. coli* genome. We measure promoter activity of >300,000
38 sequences spanning the entire genome and map 2,228 promoters active in rich media.
39 Surprisingly, 944 of these promoters were found within intragenic sequences and are associated
40 with conciliatory sequence adaptations by both the protein-coding regions and overlapping RNAP
41 binding sites. Furthermore, we perform a scanning mutagenesis of 2,057 promoters to uncover
42 sequence elements regulating promoter activity, revealing 3,317 novel regulatory elements.
43 Finally, we show that despite these large datasets and modern machine learning algorithms,
44 predicting endogenous promoter activity from primary sequence is still challenging.

45

46

47 Introduction

48

49 In 1961, François Jacob and Jacques Monod outlined the concept of the bacterial promoter
50 derived from an accumulation of genetic and biochemical studies of metabolic regulation in
51 *Escherichia coli*¹. Bacterial promoters have since become a foundation for understanding
52 molecular biology and gene regulation, with countless studies probing their genetic, evolutionary,
53 structural, thermodynamic and kinetic properties²⁻⁵. Several model promoters such as the *lac*, *trp*,
54 and phage promoters have been the subject of in-depth mechanistic studies for how RNA
55 polymerase (RNAP) recognizes promoter sequences, as well as the stepwise process to initiate
56 transcription⁶⁻⁸. In addition, many transcription factors have been described in similar detail,
57 revealing the processes through which these proteins modulate the behavior of RNAP and activity
58 of the promoter^{4,9-11}. The majority of the binding motifs for these transcription factors have been
59 studied at high resolution using modern methods¹²⁻¹⁵. In short, the myriad components that
60 define *E. coli* promoter function have been extensively cataloged and characterized, establishing
61 them as one of the most well-understood systems in molecular biology.

62

63 Despite this extensive knowledge, we still cannot answer many simple and fundamental
64 questions about *E. coli* promoters. For example, how many active promoters exist in *E. coli* at a
65 given growth condition? To what extent is promoter regulation responsible for protein level
66 remodeling during environmental changes? Given a sequence, can we predict whether a promoter
67 is encoded within it as well as its strength and/or regulation? Answers to these questions remain
68 difficult for many reasons. Although the consensus sequences for RNAP recognition motifs have
69 been known for decades, a simple search of the genome based on these motifs yields many false
70 positives. In fact, within a region, there are often sequences closer to the RNAP recognition motifs
71 than the actual functional promoter^{16,17}. Experimental efforts to identify promoters using 5' RNA-
72 Seq have found tens of thousands of putative transcription start sites (TSSs) that presumably
73 mark sites with functional promoter activity, however, there is little overlap between studies^{18,19}.
74 Furthermore, although many *E. coli* promoters have been verified with strong biochemical
75 evidence²⁰, identifying the cis-regulatory elements responsible for their activity is challenging. As
76 a consequence, roughly two-thirds of the 2,565 reported *E. coli* operons do not contain any
77 transcription factor binding site annotations^{20,21}. Finally, aside from a handful of thoroughly-
78 studied promoter sequences, we are still unable to quantitatively predict the activity or behavior
79 of promoters in the context of sequence perturbations such as moving, mutating, or removing
80 transcription factor binding sites.

81

82 There are several confounding factors which make it difficult to accurately gauge if a sequence
83 can confer promoter activity. First, recent work has shown that promoter activity varies depending
84 on location in the genome due to factors such as variance in chromosomal copy number²²⁻²⁴, the
85 distribution of transcription factors within a cell^{25,26}, and chromatin accessibility²⁷⁻²⁹ masking the
86 effects of cis-regulatory elements. Efforts to normalize these effects have utilized reporters on
87 high copy number plasmids that can saturate endogenous transcriptional machinery³⁰. Second,
88 inferring promoter strength from endogenous transcript production is problematic because these
89 transcripts often contain sequences that alter their processing and stability independent of the

90 promoter sequence^{31,32}. Third, multiple promoters within close proximity, whether co-directional
91 or opposing, can affect each other's strength and resulting transcription through mechanisms
92 such as RNAP collisions and antisense RNA³³⁻³⁶. Finally, not all sequences that initiate RNA
93 transcription are capable of producing mature and translatable RNA³⁷.

94
95 Here we investigated promoter regulation in *E. coli* using a massively-parallel reporter assay
96 (MPRA) designed to isolate promoter activity from other confounding factors influencing genetic
97 regulation³⁸. We measured promoter activity at 17,189 reported TSSs and found that a majority
98 are not autonomous promoter sequences capable of gene transcription. We then measured
99 promoter activity of 321,123 sheared genomic fragments spanning both strands of the *E. coli*
100 genome (8.5x coverage) to elucidate the promoter landscape in rich and minimal media. We then
101 systematically tiled these regions to precisely map promoter boundaries, revealing many regions
102 with multiple promoters in close proximity, as well as many antisense promoters within genes
103 that shape both codon usage and transcription levels. To characterize sequence motifs encoding
104 promoter activity, we performed systematic mutagenesis of 2,057 active promoters and identified
105 cis-regulatory elements affecting promoter activity. With this approach, we characterized the
106 regulatory effects of 568 transcription factor binding sites reported by RegulonDB as well as 2,583
107 novel sites, thereby providing functionally annotated profiles for promoters driving expression in
108 rich LB media for 1,158 of the 2,565 operons in *E. coli*. Lastly, we trained several machine learning
109 models on these datasets to better understand the features that may be used to classify *E. coli*
110 promoter sequences and quantitatively predict promoter function from sequence.

111

112 **Results**

113

114 **Functional characterization of 17,635 previously reported *E. coli* promoters reveals many are** 115 **transcriptionally inactive**

116

117 We first sought to validate promoters and TSSs identified by several genome-wide studies. We
118 assembled previously reported TSSs from three sources: the RegulonDB *E. coli* database²⁰ (8,486
119 TSSs), a directional RNA-Seq study by Wanner et. al¹⁸ (2,123 TSSs), and a RNA-Seq study by
120 Thomason et. al¹⁹ (14,868 TSSs). These three sources identify 23,798 unique TSSs active during
121 log-phase growth in rich media with little agreement regarding the location of TSSs between
122 studies and only 93 exact matches shared between all three (**Figure 1A**). Even when we collapsed
123 clusters of TSSs within 20 bp of each other to the most upstream TSS to minimize redundancy,
124 17,635 unique TSSs remained. These TSSs are likely some combination of true promoters and
125 false positives due to RNA processing, transcriptional noise, or experimental and computational
126 artifacts.

127

128 To see if these TSS regions could drive gene expression of a transcriptional reporter, we used a
129 genomic MPRA we developed³⁸ to quantitatively measure the autonomous promoter activity of
130 17,635 TSSs (**Figure 1B**). This system allows for single integration of large reporter libraries into
131 a defined locus. The promoter activity reporter is insulated by multiple transcriptional terminators
132 and the reporter transcript contains a RiboJ ribozyme sequence upstream of the RBS that

133 standardizes the transcript produced. For each TSS, we synthesized oligonucleotides spanning
134 120 bp upstream to 30 bp downstream of the TSS, which encodes the majority of promoter
135 activity driving expression at a given TSS³⁹. We included 96 well-characterized promoters from
136 the BioBricks registry⁴⁰ designed to span a wide range of expression levels to serve as positive
137 controls. We also included 500 negative controls that were selected 150 bp sequences from the
138 *E. coli* genome. Our criteria for selecting these sequences was that they are more than 200 bp
139 from any TSS reported in the aforementioned studies. We engineered these 18,222 unique
140 sequences to express a uniquely barcoded sfGFP transcript and subsequently integrated this
141 pooled library of reporter constructs into the *nth-ydgR* intergenic locus within the *E. coli*
142 chromosomal terminus using a recombination-mediated cassette exchange system⁴¹. We
143 determined promoter activity levels by performing targeted amplicon sequencing of the barcoded
144 sfGFP transcripts to quantify RNA-seq levels of each barcode normalized to their DNA-seq
145 abundances, and precisely measured expression for 97.5% (17,767/18,222) of TSSs in this library
146 (**Figure 1C**) with an average of 69.5 barcodes measured per library member (**Figure S1A**).
147 Expression measurements were consistent between replicates which were separately barcoded,
148 cloned, and quantified ($R^2 = 0.919$, $p < 2.2 \times 10^{-16}$) (**Figure S1B**). To call a TSS active we set a
149 threshold of at least greater than two standard deviations above the median of the negative
150 control distribution and normalized the data such that the threshold value was set to 1 (**Figure**
151 **1D**). Among the 17,635 original TSSs, we confidently quantified 17,189 (97.4%) and identified
152 2,670 exhibiting expression levels above our experimentally determined threshold (**Figure 1E**).
153 Notably, this number of active promoters is more consistent with the number of operons
154 identified using long-read sequencing to characterize full-length *E. coli* transcripts^{42,43}. Amongst
155 these 2,670 confirmed promoters, we recovered expression data for many well-known promoters
156 and three of the strongest corresponded to the 16S and 23S polycistronic operon, the most highly
157 expressed operon in the *E. coli* genome⁴⁴.

158
159 To confirm whether our set of negative controls were truly depleted of promoter activity, we tested
160 a set of 936 completely random 150 bp nucleotide sequences and compared the expression
161 levels to our negative controls (**Figure S2**). Despite overall low mean levels of expression
162 (Random sequences: 0.115, Negative controls: 0.036), 2.35% of random promoters drove
163 expression higher than our negative threshold whereas only 0.851% of negative controls
164 exceeded this threshold. A recent study found that 4/40 (10%) random 103 bp sequences
165 exhibited promoter activity⁴⁵ and suggests the frequency of promoter-like activity in overall
166 sequence space is seemingly very high. These results demonstrate that the negative controls
167 used in our assay are depleted in promoter activity, even compared to completely random
168 sequences, and implies that there is negative selection for spurious promoter activity across
169 certain regions of the *E. coli* genome.

170 171 **Chromosomal-position specific effects are consistent across diverse promoter sequences**

172
173 Several recent studies have shown that promoter expression levels can be highly variable
174 between genomic locations^{25,27,28}. However, these studies have primarily focused on individual
175 promoters in multiple locations, leaving uncertainty regarding whether these effects are

176 promoter-specific or represent a more widespread phenomenon impacting any promoter at a
177 given position. To study these chromosomal position effects across a wide range of promoters,
178 we integrated the entire TSS promoter library in both left and right chromosomal midreplichoes
179 and compared expression measurements between these positions and the *E. coli* chromosomal
180 terminus (**Figure S1C**). Promoter measurements remained highly consistent between locations,
181 although the two midreplichore positions exhibited slightly higher concordance with each other
182 ($r = 0.97, p < 2.2 \times 10^{-16}$), than either midreplichore to the terminus ($r = 0.95, p < 2.2 \times 10^{-16}$). Positive
183 control sequences, which do not contain regulatory elements in addition to the RNAP binding
184 sites, were highly correlated between all locations. We conclude that overall, diverse promoters
185 exhibit similar relative expression levels across genome-positions, although the absolute
186 expression may vary.

187

188 **Inactive TSS-associated promoters are enriched for -10 but not -35 $\sigma 70$ binding motifs**

189

190 A majority of *E. coli* promoters are regulated by the housekeeping sigma factor $\sigma 70^{46}$, and thus
191 we expected that active promoters would be enriched for the canonical $\sigma 70$ motifs. Promoters of
192 the $\sigma 70$ family are well known for containing two hexamer motifs, the -10 and -35 motifs, which
193 recruit RNAP and are named after their position relative to the TSS. We used a $\sigma 70$ position-weight
194 matrix (PWM)¹⁶ to analyze whether active TSS promoters were enriched for these motifs.
195 Although both active and inactive TSS-associated promoters were enriched for the canonical -10
196 motif compared to our negative controls (active: $p < 2.2 \times 10^{-16}$, inactive: $p = 6.2 \times 10^{-8}$), we found
197 the -35 scores of inactive promoters were generally no greater than negative controls ($p = 0.33$)
198 (**Figure 1F**). Conversely, active TSS-associated promoters contained significantly higher -35
199 scores than negative controls ($p = 1.4 \times 10^{-8}$) or inactive TSS-associated promoters ($p < 2.2 \times 10^{-16}$).
200 We conclude that inactive TSS-associated promoters lack -35 elements but may become
201 active in growth conditions where additional transcription factors mobilize and facilitate RNAP
202 positioning in the absence of a -35 motif.

203

204 **Genome-wide Identification of *E. coli* promoters**

205

206 Despite functionally screening 17,635 previously implicated TSS regions, we encountered
207 instances where essential operon promoters remained unidentified, suggesting that there were
208 still undiscovered promoters within the genome. For instance, despite screening several reported
209 TSS regions upstream of the essential *yrbA-murA* operon, none exhibited expression greater than
210 our activity threshold. To comprehensively detect all promoters, we cloned, barcoded, and
211 measured the transcriptional activity in LB of 321,123 sheared genomic fragments ranging
212 between 200 and 300 bp (median = 244 bp), providing $\sim 8.5x$ coverage per strand of the *E. coli*
213 genome (**Figure 2A, Figure S3A, Figure S3B**). We averaged the expression of fragments
214 overlapping each nucleotide position to achieve highly replicable values of strand-specific
215 promoter activity at single-nucleotide resolution (**Figure S3C**). This data may be viewed using our
216 online tool (<https://ecolipromoterdb.com>), revealing defined regions of promoter activity across
217 the entire *E. coli* genome (**Figure 2B**). We classify candidate promoter regions as contiguous
218 regions of at least 60 bp with promoter activity measurements higher than an empirically derived

219 threshold. This threshold was established to maximize recall of previously identified active TSSs
220 while minimizing the inclusion of inactive TSSs.

221
222 With the chosen threshold, we found 3,477 candidate promoter regions in LB that overlapped
223 2,293/2,670 (85.8%) active TSSs identified in LB, 3,193/14,493 (22.0%) inactive TSSs, and 47/482
224 (9.75%) negative controls. Active TSSs not overlapping a candidate promoter region generally
225 exhibited weak activity, which may indicate that greater sensitivity is achieved through testing of
226 oligo-array synthesized regions (**Figure S3D**). In many cases, candidate promoter regions
227 overlapped multiple TSS-associated promoters, both active and inactive, indicating the potential
228 for multiple promoters within individual regions (**Figure 2B**). Overall, we detected strong promoter
229 activity at active TSSs with little promoter activity at inactive TSS promoters, demonstrating
230 agreement between these independent methods for capturing genome-wide promoter activity
231 (**Figure 2C**).

232

233 **Fine mapping of *E. coli* promoters within transcriptionally active regions**

234
235 Our survey of genomic fragments identified candidate regions of promoter activity that were well
236 above the expected size of typical promoters (**Figure S3E**)³⁹. To determine if these candidate
237 regions contained multiple promoters, we constructed a library of 48,379 150 bp oligos that tiled
238 the entire lengths of the 3,477 promoter regions identified in LB at 10 bp intervals (**Figure 2D**). For
239 candidate promoter regions under 150 bp, we synthesized a single oligo encoding the region
240 without including additional surrounding sequence context. We recovered highly replicable data
241 for 45,201(93.4%) of these variants with an average of 8 barcodes per variant (**Figure S3F, S3G**).
242 Using this approach, we could precisely pinpoint the boundaries of promoters by observing the
243 specific locations along the promoter region where tiled oligos exhibited changes in expression
244 levels. (**Figure 2E**). This analysis revealed that 1,889 of the previously identified promoter regions
245 contained one or more discrete promoters, including 278 regions containing multiple promoters
246 (**Figure 2F**). Notably, the number of promoters within a given region correlated with the size of
247 the candidate region (**Figure S3H**) but not necessarily the overall promoter activity of the region
248 (**Figure S3I**). In 1,465 candidate regions, no promoters were detected. These regions typically
249 measured under 150 bp in length, raising the possibility of being false positives or potentially
250 requiring additional transcription factors beyond the scope of the 150 bp regions assessed.
251 Altogether, this approach identified 2,228 distinct promoters active in LB. Furthermore, by
252 determining the overlap of all active oligos tiling a promoter, we were able to infer the minimal
253 sequence necessary for each promoter. When comparing the sizes of the minimal sequence
254 necessary for promoter activity, we observed an enrichment for sequences of approximately 40
255 bp, which is a typical size for $\sigma 70$ promoters⁴⁷⁻⁴⁹ (**Figure 2G**). We also observed an enrichment
256 for 150 bp minimal promoter regions, although these were generally weak indicating that our
257 resolution is limited when tiling weaker promoters. Overall, we were able to precisely map
258 boundaries for 2,228 promoters active in LB. Considering non-overlapping active promoters
259 identified during our TSS screen, we find 2,859 distinct promoters. Amongst these promoters, we
260 have identified promoters regulating 99 out of 100 randomly sampled essential genes including
261 the promoter for the essential *yrbA-murA* operon which was missed in the TSS screen

262 **(Supplementary Table 1)**. The missing promoter was for the *yjeE* gene, which exhibits an atypical
263 operon structure, wherein the first gene in the operon overlaps a gene encoded in the opposite
264 direction. Furthermore, we detected promoter activity in regions 100 bp upstream of 24 of 38
265 recently described small open reading frames (smORFs) identified by ribosome profiling⁵⁰,
266 indicating that these proteins may be transcriptionally-regulated independently of larger
267 neighboring genes (**Figure S4**).

268
269 **Intragenic promoters are widespread, often found in the antisense orientation, and alter**
270 **transcript levels and codon usage of the genes they are within.**

271
272 While promoters are commonly thought of as gene regulatory sequences upstream of transcribed
273 genes, they can also be found within genes and oriented to transcribe genes in the antisense
274 direction. We thus sought to explore these atypical promoters and their consequences on the *E.*
275 *coli* genome and transcriptome. Many studies have found pervasive antisense transcription in
276 prokaryotes⁵¹⁻⁵⁴, though there is controversy over the functional relevance and whether they are
277 just due to a noisy transcriptional apparatus⁵⁵. At the same time, it has been functionally shown
278 that antisense promoters can alter a sense gene's transcription, translation, and steady-state
279 message levels^{35,56}. Amongst the 2,228 promoters we precisely mapped, 1,131 were primarily
280 encoded within intergenic regions while 944 were found to fully or mostly overlap intragenic
281 regions (**Figure 3A, Figure S5A**). Notably, intragenic promoters exhibited a higher prevalence
282 within single-gene operons compared to individual genes within polycistronic operons ($p = 1.05$
283 $\times 10^{-9}$, $df = 1$, Chi-squared Test). Although intergenic promoters were predominantly positioned in
284 the sense orientation relative to the nearest downstream gene, 300 of the 944 intragenic
285 promoters were positioned antisense relative to the genes they overlapped. Interestingly,
286 intragenic promoter activity had greater correlation when comparing activity between growth
287 mediums, indicating that these regions may be primarily composed of constitutive promoter
288 elements (LB: $r = 0.648$, M9 minimal: $r = 0.787$, $p > 1 \times 10^{-16}$, Wilcoxon rank-sum test, **Figures S5B-**
289 **C**).

290
291 Given that we have determined the locations of the antisense promoters driving transcription, we
292 evaluated the genome-wide consequences of antisense promoters on the transcriptome. We
293 performed RNA-Seq on *E. coli* MG1655 grown in LB and compared the transcript coverage of all
294 genes with sense promoters, antisense promoters, and both sense and antisense promoters. We
295 found that overall, genes regulated by both sense and antisense promoters exhibited a two-fold
296 decrease in expression compared to strictly sense-regulated genes (**Figure 3B**). Notably, sense-
297 regulated genes exhibited similar promoter activity on average when compared to genes with
298 both sense and antisense promoters, indicating that the result cannot be attributed solely to
299 stronger promoters in sense-regulated genes. Genes with only antisense promoter activity
300 generally did not exhibit detectable sense transcription.

301
302 The significant overlap observed between protein-coding and promoter sequences is interesting
303 given the sequence specificity necessary to encode these distinct functions. Therefore we sought
304 to investigate how sequences navigate this constraint to accommodate diverse activities. After

305 comparing the amino acid composition within intragenic promoters, we observed a significant
306 enrichment of stop codons and a preference for amino acids encoded by codons with higher AT
307 nucleotide content (**Figure 3C**). Further inspection revealed specific codons that were
308 preferentially utilized within intragenic promoter regions (**Figure 3D**), with a notable bias observed
309 among arginine codons, showing a strong preference for AGA and AGG codons. The most
310 enriched codons within intragenic promoters were typically rare in the genome, which may
311 indicate a role of preferential codon usage in controlling promoter activity within genes. The
312 connection between rare codons and regulatory roles has been previously observed in the context
313 of N-terminal codon bias, where rare codons influenced expression levels through secondary
314 structure interactions⁵⁷. Moreover, the observed higher percentage of AT-content^{58,59} and rare
315 codons⁶⁰ may further support the notion that intragenic promoters are linked to horizontally-
316 acquired genes.

317
318 Next, we investigated how intragenic promoter sequences had adapted to conform to the
319 constraints of protein-coding sequence space. A peculiar feature of promoter sequences in *E.*
320 *coli*, is the presence of trinucleotides matching stop codons within the canonical -10 and -35 σ 70
321 motifs (-35: 5'-TTGACA, -10: 5'-TATAAT). Therefore, we hypothesized that the reuse of these
322 nucleotide patterns offers another mechanism by which the *E. coli* genome counteracts the
323 spurious evolution of intragenic promoters, thereby explaining their scarcity relative to the ease
324 by which they can evolve⁴⁵. We used a σ 70 PWM¹⁶ to identify the highest-scoring σ 70 motifs within
325 intragenic promoters and determined their relative coding frames. Interestingly, we observed a
326 lower frequency of -35 elements in +2 coding frames and the -35 motifs detected at +2 positions
327 exhibited significantly reduced resemblance to the canonical motif (**Figure S6A**). Similarly, -10
328 motifs were least frequently found in the +1 positions, although -10 motifs at this position did not
329 show lower overall scores (**Figure S6B**). The observed depletion of -35 motifs positioned in the
330 +2 reading frame and -10 motifs in the +1 reading frame is likely due to the fact that the canonical
331 sequences for these motifs would create stop codons within the protein if placed at these
332 positions. This suggests a simple, but effective preventative mechanism against the spurious
333 evolution of intragenic promoters that is inherent to their sequence motifs.

334

335 **The *E. coli* promoter landscape is dynamic in response to environmental conditions**

336
337 It is well understood that bacterial cells respond to environmental conditions through changes in
338 their transcriptional profiles⁶¹, however, it has not been shown how the global promoter landscape
339 changes to facilitate these cellular transitions. To explore this, we measured promoter activity of
340 our genomic fragment library in exponentially growing cells under glucose minimal media
341 conditions. Compared to LB, cells grown in glucose minimal media do not have access to
342 environmental amino acids and must synthesize these and other essential compounds on their
343 own⁶². We recovered replicable promoter activity measurements for 318,457 genomic fragments
344 in glucose minimal media, spanning the genome with 8.38x coverage (**Figure S7A, Figure S7B**).
345 We identify 3,321 candidate promoter regions in glucose minimal media with an average length
346 of 293 bp (**Figure S7C**). Although 2,466 of these regions overlapped with regions found in LB, we
347 found 960 only found in LB and 1,029 exclusive to M9 (**Figure 4A**). Many of these condition-

348 dependent promoter regions were weak compared to those identified in both conditions (**Figure**
349 **S7D**), nonetheless, each condition revealed distinct strongly activated regions unique to it. The
350 observed low activity of condition-unique promoters is similar to what has been observed in
351 synthetic inducible promoter systems, where tightly-regulated promoters often exhibit reduced
352 expression in induced conditions⁶³. To identify the most differentially-expressed promoters in
353 each condition, we extracted regions larger than 60 bp that exhibited greater than two-fold
354 difference in activity between conditions. With this criterion, we found 278 regions upregulated in
355 LB and 644 regions upregulated in glucose minimal media. In glucose minimal media, the greatest
356 increase in promoter activity occurred at *ryhB*, a Fur-regulated gene encoding a small RNA that
357 regulates iron-binding and iron-storing proteins when available iron is limited^{64,65} (**Figure S7E**). In
358 LB, the strongest activated region is positioned to drive expression of the *rbsDACBKR* operon,
359 which is essential for uptake and utilization of extracellular ribose⁶⁶ (**Figure S7E**).

360
361 For each condition, we matched activated intergenic and sense promoter regions with the nearest
362 downstream gene and found 159 genes poised for activation in LB and 392 genes poised for
363 activation in glucose minimal media. To see if promoter activation resulted in an increase in
364 expression of these genes, we compared RNA-Seq coverage of the genes with the top 100
365 strongest promoter activation in each condition (**Figure S8**). In each condition, promoter
366 activation resulted in a concomitant increase in RNA-Seq coverage (LB: $p = 1.1 \times 10^{-5}$, M9: $p = 1.9$
367 $\times 10^{-5}$, Wilcoxon rank-sum test). To see which cellular responses were being mobilized by
368 remodeling the promoter landscape, we used the RAST annotation engine^{67,68} to assign functional
369 categories to activated genes and identify enriched cellular processes. Genes downstream of
370 promoter regions activated in LB are predominantly associated with carbohydrate utilization
371 whereas genes downstream of promoters activated in glucose minimal media were associated
372 with amino acid utilization (**Figure 4B**). Overall, we find distinct condition-dependent activation of
373 promoter regions leading to changes in gene expression associated with carbohydrate utilization
374 in LB and amino acid utilization in glucose minimal media.

375
376 Next, we explored how these changes in the promoter landscape are mediated by transcriptional
377 machinery and evaluated the transcription factor binding site (TFBS) composition of promoter
378 regions activated in each condition. As opposed to traditional transcriptomebased
379 measurements which measure changes in downstream gene expression, this assay identifies
380 upstream regulatory regions that contribute to promoter activity in response to changing
381 conditions. By cross-referencing these activated promoter regions to TFBSs reported by
382 RegulonDB, we identified transcription factors facilitating these changes to the promoter
383 landscape (**Figure 4C**). Upon comparing TFBS content of these regions we found that binding
384 sites for several global transcriptional regulators⁶⁹, including IHF, Lrp, and Fis occurred at similar
385 frequencies between these conditions. Conversely, binding sites for Fur, another global
386 transcription factor, were enriched by roughly 20-fold within regions activated in glucose minimal
387 media compared to regions activated in LB. This transcription factor is essential for maintaining
388 iron homeostasis^{70,71}, and is a known regulator of *ryhB*, the most upregulated gene we found in
389 glucose minimal media. Binding sites for CRP were enriched by more than two-fold in regions
390 activated in LB compared to glucose minimal media. This transcription factor is activated in

391 glucose-limited conditions and so would likely not induce promoter activity in glucose minimal
392 media. Overall, we found 455 TFBSs within regions activated in LB and 637 annotations in regions
393 activated in glucose minimal media (**Figure 4D**). In addition to global regulators, we found many
394 TFBSs that appear exclusive to each condition targeting relatively few regulatory targets.
395 Interestingly, the combined contribution of non-global transcription factors activating 10 or fewer
396 sites were responsible for over a third of all activated promoter regions, underscoring the
397 significant involvement of local transcriptional regulators in driving the overall changes to the
398 transcriptome. Transcription factors MetJ, GadX, and GadW were exclusively found in regions
399 activated in glucose minimal media whereas FlhDC, GlpR, and CytR were the most enriched
400 amongst regions activated in LB.

401 402 **Mutational scanning of 2,057 *E. coli* promoters identifies regulatory elements controlling** 403 **transcription**

404
405 After globally identifying promoter regions in the bacterial genome, we sought to develop an
406 approach to identify sequence motifs regulating these promoters. Recent work has demonstrated
407 a high-resolution saturation mutagenesis approach to identify regulatory motifs within individual
408 uncharacterized promoters^{21,72}. Inspired by this work, we implemented a scanning mutagenesis
409 strategy to explore the sequence features that regulate active promoters. For 2,057 active TSS-
410 associated promoters identified in LB, we systematically scrambled individual 10 bp sequences
411 spanning the -120 to +30 positions at five bp intervals (**Figure 5A**). Using this approach, we would
412 expect that disrupting a repressor site would increase expression, whereas disrupting a RNAP or
413 activator site would decrease expression. These scrambled sequences were designed to
414 maximize distance from the original sequence while maintaining nucleotide content, ensuring
415 perturbation of any motifs at each position contributing to transcriptional regulation. In total, we
416 designed a library of 59,653 sequences consisting of 2,057 active TSS-associated promoters,
417 their scrambled variants, and the previously described set of negative and positive controls. We
418 measured promoter activity of this library as before and recovered replicable expression
419 measurements for 52,900/59,653 (89%) of this library in LB, with an average of seven barcodes
420 per variant (**Figure S9A, S9B**). Using this approach, we identify regions that either increased or
421 reduced expression across thousands of promoters in a single assay (**Figure 5B**). These
422 sequences were enriched at the -35 and -10 positions for regions that increased expression,
423 which is expected considering the majority of promoters are $\sigma 70$ dependent. However, many
424 sequences outside of these -10 and -35 regions were also found to contribute to regulation.

425
426 To validate our approach, we first examined the *lacZYA* promoter, a classic gene regulation model
427 whose sequence motifs are well characterized. This promoter is known to contain a variety of
428 regulatory motifs, including twin LacI repressor sites centered at +11 and -82⁷³, a CRP activator
429 site centered at -61⁷⁴, and a $\sigma 70$ RNAP binding site. Our analysis revealed distinct signals
430 corresponding to each of these sites, as well as quantitative measurements for their contribution
431 to expression (**Figure 5C**). Additionally, scanning mutagenesis of the previously characterized
432 *relBE* promoter achieved similar results, identifying a reported RelBE repressor site at the +1
433 position⁷⁵ as well as -10 and -35 $\sigma 70$ recognition motifs⁷⁵.

434
435 Considering that our approach effectively captures the effects of known binding sites, we
436 proceeded to investigate whether it could also identify regulatory sites within uncharacterized
437 promoters. Although we performed this scanning mutagenesis for 2,057 TSS-associated
438 promoters, here we highlight a few examples to demonstrate the utility of this method (**Figure**
439 **5D**). The cyclopropane fatty acyl phospholipid synthase gene, *cfa*, exhibits dynamic expression⁷⁶
440 and plays a crucial role in cell membrane integrity under acidic conditions⁷⁷. While there have
441 been several transcription factors implicated in regulation of *cfa*, the motifs responsible for its
442 direct regulation are still unknown. Our approach identified a candidate σ 70 promoter regulating
443 this gene with a -10 motif centered 34 nucleotides upstream of the reportedly associated TSS as
444 well as a -35 motif 57 bp upstream, implying that the reported TSS is likely not the primary site
445 for transcription initiation. Additionally, we identified two repressor sites—one located in the
446 spacer region and another upstream of the -35 motif. We also identified novel regulatory regions
447 for an uncharacterized promoter regulating *rpsL*, an essential gene and component of the 30S
448 ribosomal subunit. In this case, we identified a candidate σ 70 RNAP binding site with predictably
449 positioned -10 and -35 motifs, as well as an unknown repressor located over the transcription
450 start site. Notably, mutating the repressor site resulted in a threefold increase in promoter
451 expression. Although further experiments²¹ are necessary to identify the transcription factors
452 acting on these promoters, our results provide valuable insights by pinpointing the sequence
453 elements responsible for the regulation of these genes.

454 455 **Global identification of 7,293 *E. coli* promoter regulatory motifs**

456
457 We expanded the scope of our analysis to systematically explore the regulatory motifs amongst
458 all 2,057 promoters tested. We used individual barcode measurements, across four replicates, to
459 find significant differences between the mean expression of the WT and mutated sequences
460 (Student's t-test with 1% FDR). Among the mutations that significantly altered expression, 1,885
461 increased expression whereas 5,408 decreased expression (**Figure 6A**). Mutated sites were
462 located throughout promoters and resulted in dramatic changes in expression, some over 100-
463 fold (**Figure S10A**). We observed markedly different distributions for the positions of sequences
464 that increased expression compared to those causing decreased expression (**Figure 6B**). Regions
465 that increased expression were enriched at the -10, -35, and -70 positions, which is consistent
466 with the σ 70 RNAP binding motif as well as the typical position of transcriptional activators
467 among class I bacterial promoters⁷⁸⁻⁸⁰. Regions that decrease expression were found to localize
468 to the TSS, spacer, and -35, which is consistent with known mechanisms of RNAP occlusion by
469 steric hindrance^{80,81}. Alternatively, repressive sites within the spacer could be negatively
470 influencing transcriptional initiation through transcription factor-independent mechanisms⁸².
471 Furthermore, we found that intergenic promoters contained more regions that altered promoter
472 activity when scrambled compared to intragenic promoters, implying that intragenic promoter
473 sequences contain more compact or fewer regulatory elements (**Figure S10B**).

474
475 Next, we cross-referenced these regulatory regions with the extensive collection of putative and
476 experimentally determined regulatory sites reported by RegulonDB⁸³. First, for all promoter

477 mutagenesis profiles, we merged adjacent regions found to influence promoter activity, resulting
478 in the identification of 1,414 regions that increase expression and 1,903 regions that decrease
479 expression. Sites were 20 bp on average (indicating they exhibited regulatory impacts across four
480 consecutive 10 bp scramble mutations spaced 5 bp apart) (**Figure S10C**) with effect sizes largely
481 independent of their lengths (**Figure S10D**). Of the 2,453 unique TFBSs reported by RegulonDB,
482 1,156 overlap with regulatory regions identified by our analysis and 49% (567/1,156) resulted in a
483 significant change in activity of the promoter. The effect we observed after disrupting these
484 reported TFBSs often did not agree with the annotated effect. Our scrambling results agreed with
485 the reported effect for 65% (185/253) of activators and 43% (196/450) of repressors (**Figure 6C**).
486 We presumed the lower concordance with repressors could be due to scrambling mutations
487 disrupting both a repressor and -35 or -10 element, resulting in a decrease in expression which
488 would appear to contradict a reported repressor site. Looking at the distribution of concordance
489 for merged scrambles by position relative to the TSS, we observed a higher proportion of
490 disagreement near the -35 and -10 elements, suggesting overlapping scrambles may be
491 disrupting crucial promoter elements in addition to reported repressor sites (**Figure S10E, S10F**).
492 This may be expected considering that many repressors operate by binding regions proximal to
493 the RNAP binding site. Regardless, we found several examples where the regulatory effects
494 predicted by RegulonDB were contradicted with strong evidence, which may indicate that the
495 effect of the reported annotation is incorrect or that these sites may support multiple transcription
496 factors (**Figure 6D**). Overall, we characterized regulatory sequences in promoters driving
497 expression of 1,158 of the 2,565⁸³ operons in *E. coli* as well as many other confirmed promoters.
498 Thus, we conclude that this approach is an efficient and effective method to rapidly characterize
499 regulatory motifs within thousands of experimentally verified promoter regions.

500

501 **Predicting promoter activity from sequence remains a challenge**

502

503 In this study we generated a powerful dataset linking 117,556 unique 150 bp sequences to a
504 quantitative measurement of *in vivo* promoter activity. Using this unique dataset, we evaluated
505 our ability to determine whether a promoter was active or inactive (classification) and the precise
506 level of activity (regression). We trained several machine learning models of varying complexity
507 for both classification and regression. As many sequences are highly similar due to library design
508 and close proximity of previously reported TSSs, we split the data into 75% for training (n = 87,164)
509 and 25% (n = 30,392) for testing according to genomic location, ensuring the two sets contain
510 sequences equidistant to the origin (see Methods). For classification, we determined a threshold
511 independently for each library based on the negative controls. Sequences are considered active
512 if their expression is greater than two standard deviations above the negative median value and
513 inactive if expression falls below this threshold.

514

515 We trained several different classifiers to predict whether a given sequence was active or inactive
516 (**Figure 7A**). All classifiers output the predicted probability for each class, rather than directly
517 predicting the class, allowing them to be compared using precision-recall curves. Further details
518 for all models are included in the methods. We trained a simple logistic regression based on four
519 biophysical features known to be associated with promoter strength: max -10 σ 70 motif position

520 weight matrix (PWM) score, max σ 70 -35 motif PWM score, paired -10 and -35 PWM score (PWMs
521 scanned jointly allowing for, 16, 17, or 18 gap between the -10 and -35), and percent GC content.
522 We trained this model only using variants from the TSS library, which contained the greatest
523 diversity, as the model was unable to converge when trained on the full dataset. For comparison,
524 we trained a gapped k-mer SVM (gkm-SVM) model with word-length 10 and 8 informative
525 columns ($L = 10$, $K = 8$) on the same training set, as this model is best suited for sample sizes
526 under 20,000 and observed decreased performance relative to the logistic regression (AUPRC =
527 0.43, AUPRC = 0.53, respectively). Furthermore, we created a feature set of all 3 to 6-mer
528 frequencies and trained a logistic regression, partial least squares discriminant analysis (PLS-
529 DA), and multi-layer perceptron (MLP). To observe the effects of reducing dimensionality, we
530 additionally trained on only 6-mer frequencies for the MLP and random forest. For the simpler
531 logistic regression and PLS-DA we performed an additional feature selection step based on the
532 performance of a random k-mer. All models performed similarly, with AUPRC ranging from 0.26
533 to 0.33.

534
535 There has been recent work predicting transcriptional regulatory activity from MPRA data using
536 convolutional neural networks (CNNs), which capture intricate sequence features without *a priori*
537 knowledge⁸⁴. Inspired by this work, we trained a CNN using the DragoNN toolkit which is built on
538 top of the keras python package⁸⁵. We performed hyperparameter tuning for a three-layer CNN
539 and achieved an AUPRC = 0.44. Next, we compared the CNN to other machine learning models
540 that require less hyperparameter tuning and are more interpretable. For comparison, we trained
541 a random forest on one-hot encoded DNA, which is not well suited to categorical features, and
542 achieved an AUPRC = 0.27. Furthermore, we trained this model using frequencies of 6-mers and
543 observed a slight increase in performance (AUPRC = 0.31). Overall, the CNN achieved the highest
544 AUPRC, but the logistic regression fit with biophysical features more accurately at higher levels
545 of recall. However, these two models may not be directly comparable, as the logistic regression
546 was trained on only the TSS library rather than the full dataset.

547
548 We separately trained all of the models described above, with the exception of gkm-SVM, for the
549 more difficult task of regression (**Figure 7B**). Additionally, we included a linear regression model
550 that fit to the four “mechanistic” features to predict log-transformed expression. We evaluated
551 each model using root mean squared error (RMSE) and R^2 between predicted and observed values
552 for promoter activity. Many models perform similarly to each other, with the CNN achieving the
553 highest R-squared and lowest RMSE (RMSE = 3.12, $R^2 = 0.31$, $p < 2.2 \times 10^{-16}$). We observe
554 improvement in the linear regression on log-transformed data compared to linear regression
555 without transformation, suggesting there are non-linear relationships that are presumably
556 captured by more complex models. Random forest on one-hot encoded DNA performs worse than
557 random forest on 6-mer frequencies, in line with the heuristic that random forests are not well
558 suited to categorical features. Overall, the CNN performs best in both classification and
559 regression, although simpler models have some predictive power and have the benefit of faster
560 training times.

561 Discussion

562 More than fifty years have passed since the first conceptions of what bacterial promoters were.
563 Today, *E. coli* promoters are arguably the most well-studied gene regulatory element and yet we
564 cannot seem to agree on basic questions of how many promoters exist, what elements define
565 their function, how constrained they are in sequence space, and how far are we from predicting
566 promoter activity from sequence. Systematic identification and characterization based on
567 transcriptional profiling is confounded by genomic location, RNA processing, stability, and
568 detection differences due to differences in sequences expressed.

569
570 Here we attempted to separate promoter activity from other mechanisms of gene regulation to
571 systematically identify promoter locations, strength, and internal structure genome-wide in rich
572 media conditions. We systematically probed previous predictions and combined them with more
573 unbiased approaches to better understand promoter architecture in *E. coli*. Overall, we found
574 2,859 ≤ 150 bp promoters during log-phase growth in LB, which is consistent with recent
575 estimations by RNA profiling using long-read sequencing technologies^{42,43} and in vitro
576 transcriptional assays⁸⁶. This included many promoters contained within genes, often in the
577 antisense direction, that had large effects on mRNA levels and constrained codon choice within
578 these genes. Despite the ability of our approach to interrogate promoter activity across the entire
579 genome, there are certainly many more condition-specific promoters that remain undiscovered.
580 Moreover, it is likely that we have not identified all active promoters even under the conditions
581 investigated in this study. It is essential to acknowledge that our approach to classifying
582 sequences as promoters is based on an empirically derived threshold. However, this is a
583 simplification as promoters that fall below the threshold could become active due to the influence
584 of other factors, such as message stability³¹ and genomic context^{25,27,28}. Taken together, these
585 measurements provide one of the richest datasets on autonomous promoter activity. Our data
586 suggests that all sequences have some propensity to be a promoter, and this propensity is further
587 modulated by other factors such as stability of the message produced or integration locus to
588 ultimately determine mRNA levels. Moreover, the frequency of promoter-like activity in overall
589 sequence space is seemingly very high. This view is consistent with the surprising ease by which
590 promoters evolve from random sequences^{45,87,88}.

591
592 Our scanning mutagenesis of active TSS-associated promoters identified 3,317 regions with no
593 corresponding TFBS annotation in RegulonDB, revealing that there is a great deal more we can
594 learn about how regulation is encoded in the *E. coli* genome. For regions that overlapped known
595 sites, an appreciable proportion disagreed with the reported effect. There could be several
596 explanations for this disagreement and the discovery of these missing annotations. First, it could
597 be that the predictions of TFBSs in RegulonDB are actually false positives due to promiscuous or
598 nonproductive binding events. This seems plausible considering a recent study of the global
599 transcription factor PhoB, which supports the notion that transcription factors engage in many
600 genomic binding events with apparent non regulatory functions⁸⁹. Second, some transcription
601 factors may possess condition-dependent behavior and the conditions tested in our study do not
602 capture the full scope of their regulatory program. Finally, it is plausible that a portion of the sites

603 we identify represent true functional sites that are missing from current annotation and should
604 be interesting targets for further dissection, such as identifying which transcription factors
605 operate at these motifs. Further studies using high resolution mutagenesis strategies²¹ will be an
606 effective approach to determining which sequences within promoters contribute to regulation
607 and further efforts to predict promoter sequence-function relationships.

608
609 To better understand how promoter activity is modulated by sequence, we trained a suite of
610 machine learning models to identify promoter sequences (classification) and predict the precise
611 level of activity (regression). These models varied in complexity, from simple linear regression
612 models based on a handful of known biological features to CNNs trained on raw sequence. Even
613 with the large training set and a wealth of mechanistic information, the performance of these
614 predictive models is limited. There are several possible explanations for why it remains a
615 challenge to classify or predict the activity of *E. coli* promoters. First, it is likely challenging to
616 develop a single generalizable model for all promoters as there are several families of sigma
617 factors with distinct motifs. Therefore, models that are sigma-factor specific may be more
618 tractable. Indeed, recent studies by us and others have leveraged large MPRA datasets
619 characterizing $\sigma 70$ promoters to develop a variety of statistical and biophysical models that
620 predict expression with surprisingly high accuracy^{38,90-92}. These findings suggest that overcoming
621 the challenges associated with promoter activity prediction is plausible with the appropriate
622 training sets and a reasonable scope of study. Second, although the range of our MPRA is quite
623 dynamic, accurate predictive models may require techniques with even greater quantitative
624 resolution, especially in the noise regime of the assay where most observations fall. Finally, we
625 might simply lack the basic models for how sequences define biological functions, such as
626 promoter activity, and thus we are looking in the wrong places for information. Recent efforts to
627 use much larger libraries of random DNA sequences to identify strong promoters may serve as a
628 better starting point to constrain computational models for how sequences affect function^{93,94}.

629
630 The experimental workflows demonstrated here enable the rapid and iterative exploration of how
631 sequence affects bacterial promoter function. The convergence of DNA synthesis technologies
632 with multiplexed assays for genetic function now allow an individual to routinely design, build and
633 test 10^4 - 10^5 designs on a monthly basis. Such empirical power has no equivalent in other physical
634 systems and has now reached the limits of human experimental design and planning. Thus,
635 understanding bacterial promoters might be one of the best problems to develop and test large-
636 scale design-of-experiment and active learning methodologies to build better predictors and
637 discriminate between different mechanistic models of function.

638
639
640
641
642
643
644
645

646 **Acknowledgements**

647

648 This work was supported by the National Science Foundation Graduate Research Fellowship
649 2015210106 and the HHMI Hanna H. Gray Postdoctoral Fellowship (GT15182) to G.U., National
650 Institutes of Health New Innovator Award DP2GM114829 to S.K., Searle Scholars Program (to
651 S.K.), U.S. Department of Energy (DE-FC02-02ER63421 to S.K.), UCLA, and Linda and Fred Wudl.
652 We thank the UCLA BSCRC high throughput sequencing core and Technology Center for
653 Genomics and Bioinformatics for technical assistance; Robert B. Phillips, Reid C. Johnson, and
654 Jeffery H. Miller for thoughtful feedback throughout this project; Matteo Pellegrini for
655 computational advice; Christina P. Burghard for advice on bioinformatics analysis; and all past
656 and present members of the Kosuri lab for technical feedback. Furthermore, we would also like
657 to thank David Gray and Lisa M. Golden for manuscript feedback. We would also like to thank
658 the invaluable resources of RegulonDB and EcoCyc as well as all contributors to these
659 collections. Lastly, we thank the UCLA Molecular Biology Interdepartmental Graduate Program
660 and UCLA Bioinformatics Interdepartmental Graduate Program.

661

662 **Author Contributions**

663

664 G.U., K.D.I., H.K., and S.K. designed the study. G.U., A.D.T., and M.B. developed and performed
665 experimental methods. N.B.L. developed the genomic fragmentation isolation method. G.U.,
666 K.D.I., and A.D.T. analyzed, and interpreted data. T.C. developed k-mer based multilayer
667 perceptron for promoter prediction. K.D.I. developed and implemented machine learning
668 approaches for promoter prediction. C.A. and G.U. created the interactive website for data
669 sharing. G.U., K.D.I., and S.K. wrote the manuscript.

670

671 **Declaration of Interests**

672

673 S.K. is cofounder and CEO and holds equity in Octant Inc. N.L. is an employee and holds equity
674 in Octant Inc. All other authors declare no competing interests.

675

676 **Declaration of Generative AI and AI-assisted technologies in the writing process**

677

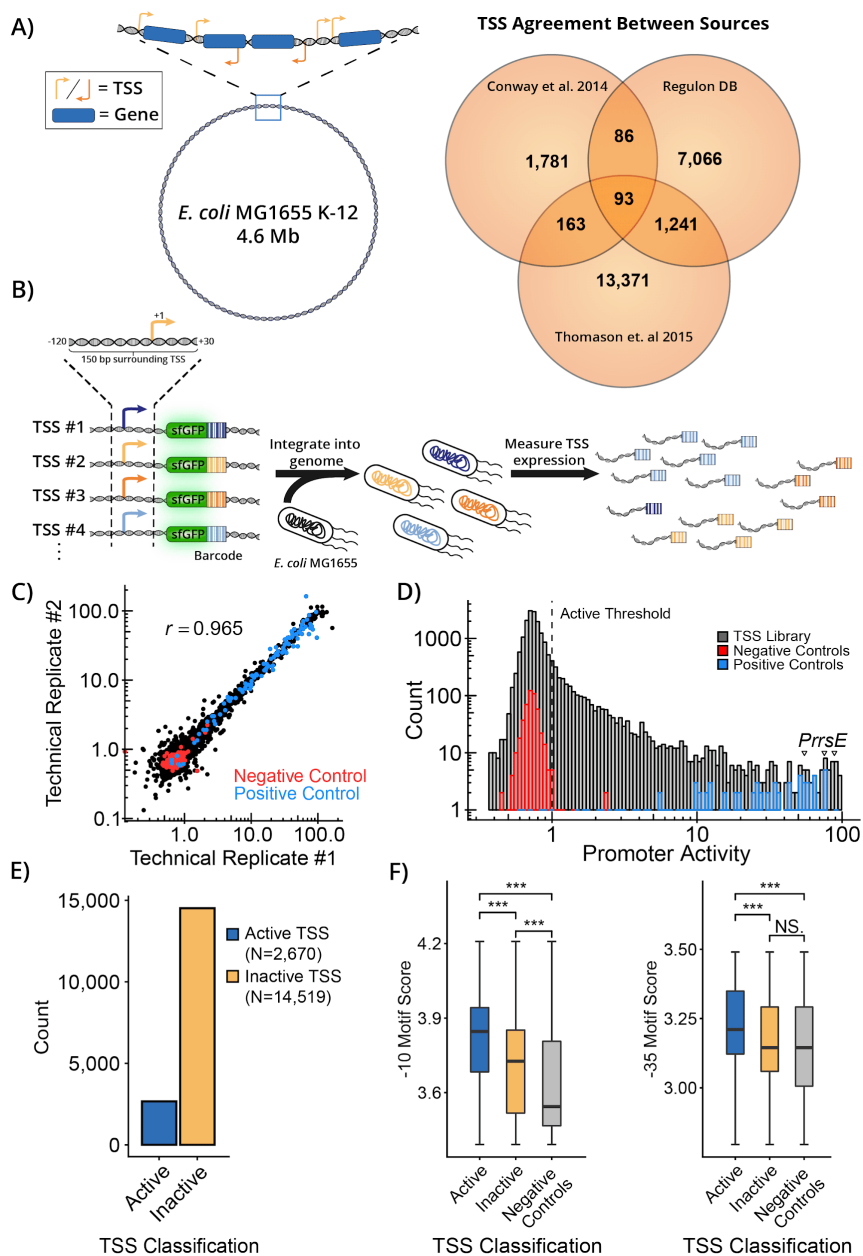
678 During the preparation of this work the author(s) used Chat GPT (GPT-3.5) in order to improve
679 language and clarity. After using this tool/service, the author(s) reviewed and edited the content
680 as needed and take(s) full responsibility for the content of the publication.

681

682 **Inclusion and Diversity**

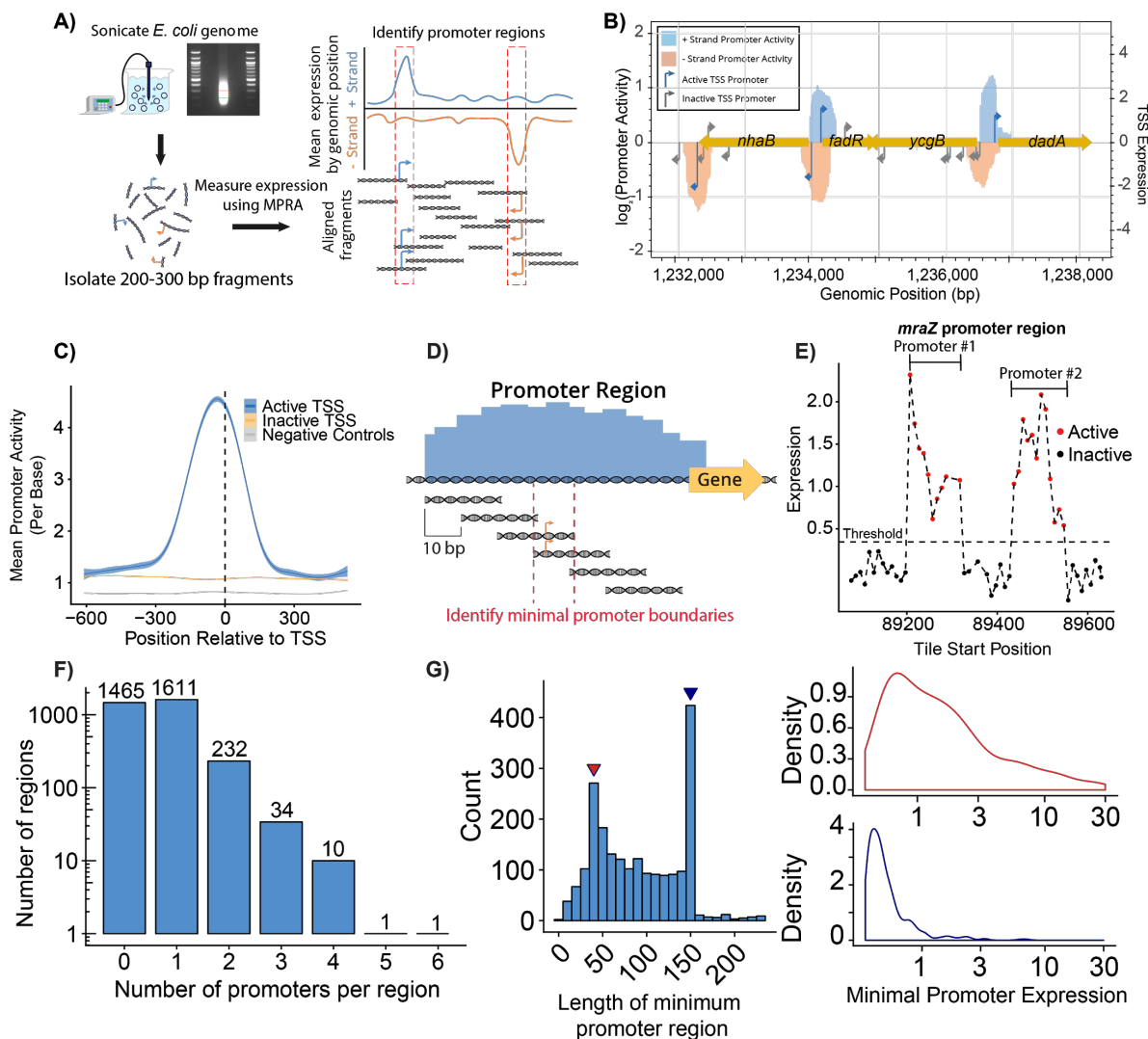
683

684 One or more of the authors of this paper self-identifies as an underrepresented ethnic minority
685 in their field of research or within their geographical location. One or more of the authors of this
686 paper self-identifies as a gender minority in their field of research. One or more of the authors of
687 this paper received support from a program designed to increase minority representation in
688 their field of research.



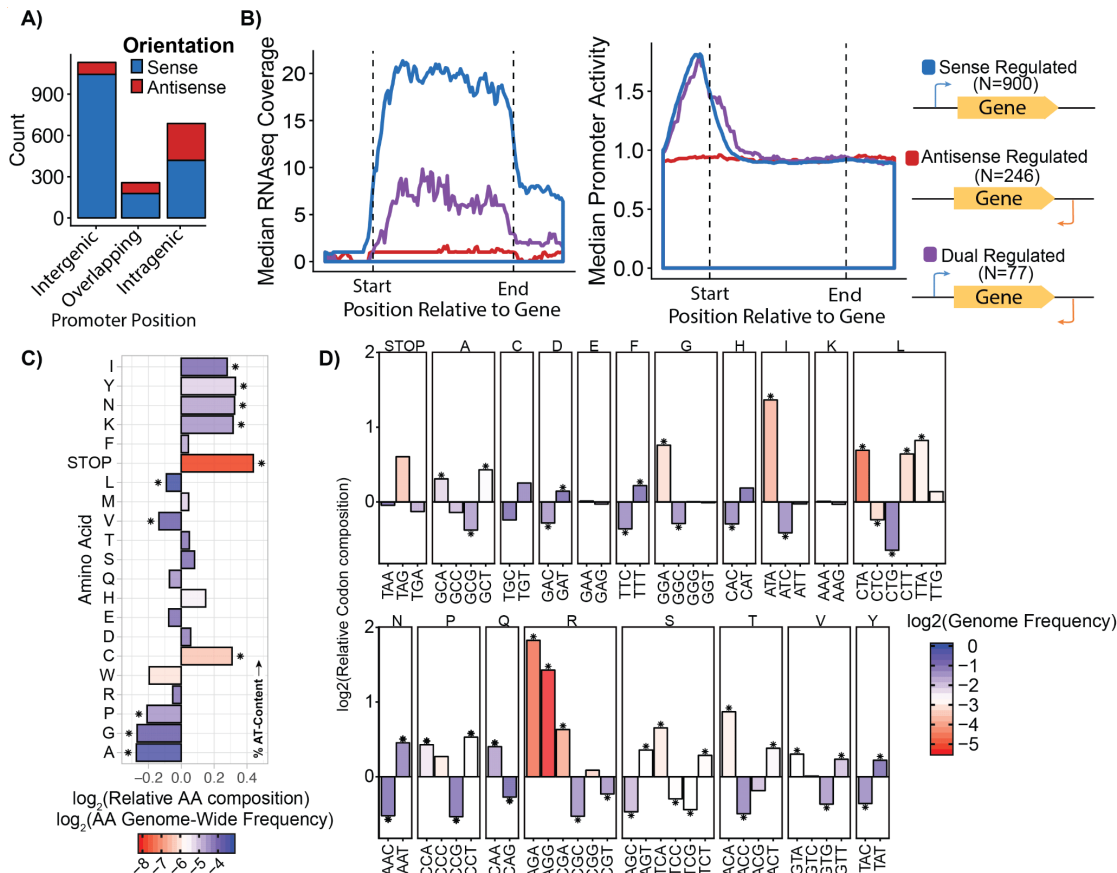
689
 690 **Figure 1) Functional characterization of 17,635 previously reported *E. coli* promoters.** **A)** Three
 691 sources of genome-wide promoter predictions show little agreement in the reported TSSs at the
 692 single-nucleotide level. **B)** We synthesized oligos overlapping the -120 to +30 bp context of
 693 17,635 reported TSSs and integrated construct into a fixed genomic landing pad. Measuring
 694 barcode expression using RNA-Seq captures quantitative measurements of transcriptional
 695 activity for individual TSSs. **C)** MPRA results are highly replicable across technical replicates ($r =$
 696 0.965 , $p < 2.2 \times 10^{-16}$). **D)** The TSS library measurements span over 100-fold with negative
 697 controls exhibiting low levels of expression and positive controls spanning the entire dynamic
 698 range. **E)** A majority of tested TSSs are inactive in LB. **F)** Active and inactive TSSs have
 699 significantly different mean PWM scores for -10 and -35 σ^{70} motifs (Wilcoxon rank-sum test,
 700 “***” = ≤ 0.001).

701

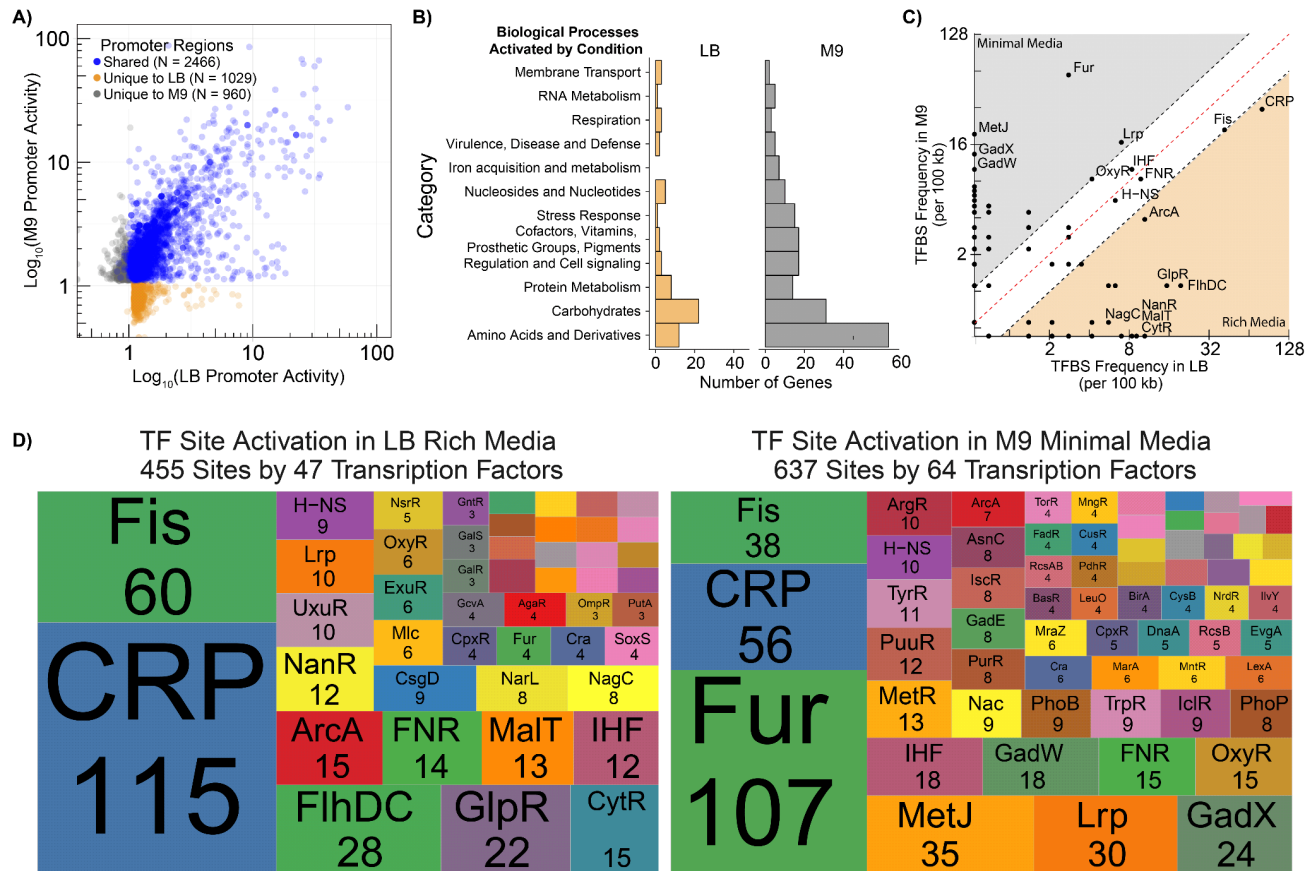


702
 703 **Figure 2) Genome-wide Identification of *E. coli* promoters.** **A)** 321,123 shared genomic
 704 fragments were screened using the same MPRA platform. The fragments were 200 to 300 bp in
 705 size giving an average 8.5x coverage across each strand of the *E. coli* genome. Promoter
 706 activity of each fragment was measured and averaged at each position to recover nucleotide-
 707 specific expression. **B)** We created a website to showcase the *E. coli* promoter landscape
 708 (<https://ecolipromoterdb.com/>). This section of the genome displayed in this figure contains
 709 five candidate promoter regions that appear within intergenic regions. **C)** Meta-analysis of mean
 710 promoter activity at experimentally validated active TSSs, inactive TSSs, and negative controls.
 711 **D)** Oligo tiling library identifies promoters within candidate promoter regions. We synthesized
 712 150 bp oligos tiling all promoter regions identified in rich media at 10 bp intervals. We then
 713 determine minimal promoter boundaries by identifying the overlap of transcriptionally active
 714 tiles. **E)** Oligo tile expression across the *mraZ* promoter shows two distinct promoters. Positions
 715 are defined according to the right-most genomic position of each 150 bp oligo. Dashed line
 716 indicates the threshold for active oligo tiles **F)** Distribution of the number of promoters per
 717 promoter region shows many regions contain multiple promoters. **G)** Left: Distribution of the

718 lengths of the minimal promoter boundaries shows enrichment for $\sigma 70$ -promoter sized regions
 719 (40 bp). Right: 40 bp minimal promoters (red) span a wide range of expression whereas 150 bp
 720 promoters are typically weak (blue).
 721

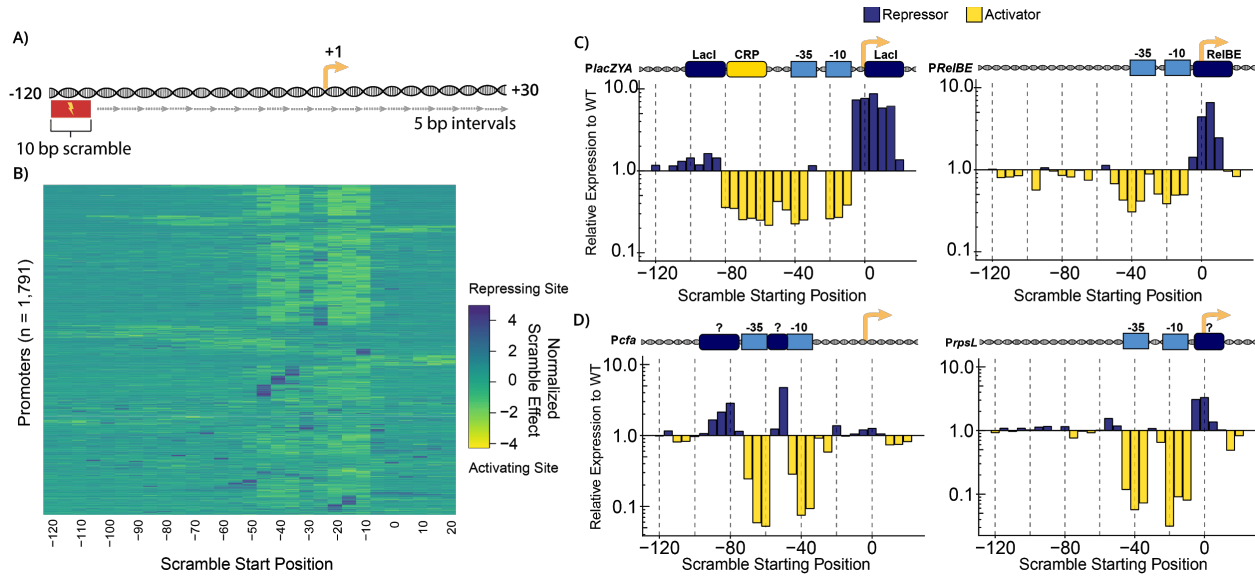


722 **Figure 3) Intragenic promoters are widespread, often found in the antisense orientation, and**
 723 **alter transcript levels and codon usage of the genes they are within. A)** Orientation and
 724 **positioning of identified promoters reveals many promoters are intragenic and antisense. B)**
 725 **Antisense promoters suppress gene expression genome-wide. Left: Meta-gene analysis of the**
 726 **median RNA-Seq coverage across all sense, antisense, and dual-regulated genes. Right: Meta-**
 727 **gene analysis of sense promoter activity at sense, antisense, and dual regulated genes. C)**
 728 **Intragenic promoters are enriched for specific amino acids relative to whole genome amino acid**
 729 **frequencies (Chi-squared test, “*” = $p < 0.05$). Amino acids are arranged by mean AT-content of**
 730 **all corresponding codons. D)** Specific, often rare, codons are enriched in intragenic promoters.
 731 Codon bias within intragenic promoters relative to whole genome. Bars are colored by the
 732 relative genome-wide usage compared to other synonymous codons (Chi-squared test, “*” = $p <$
 733 0.05).
 734



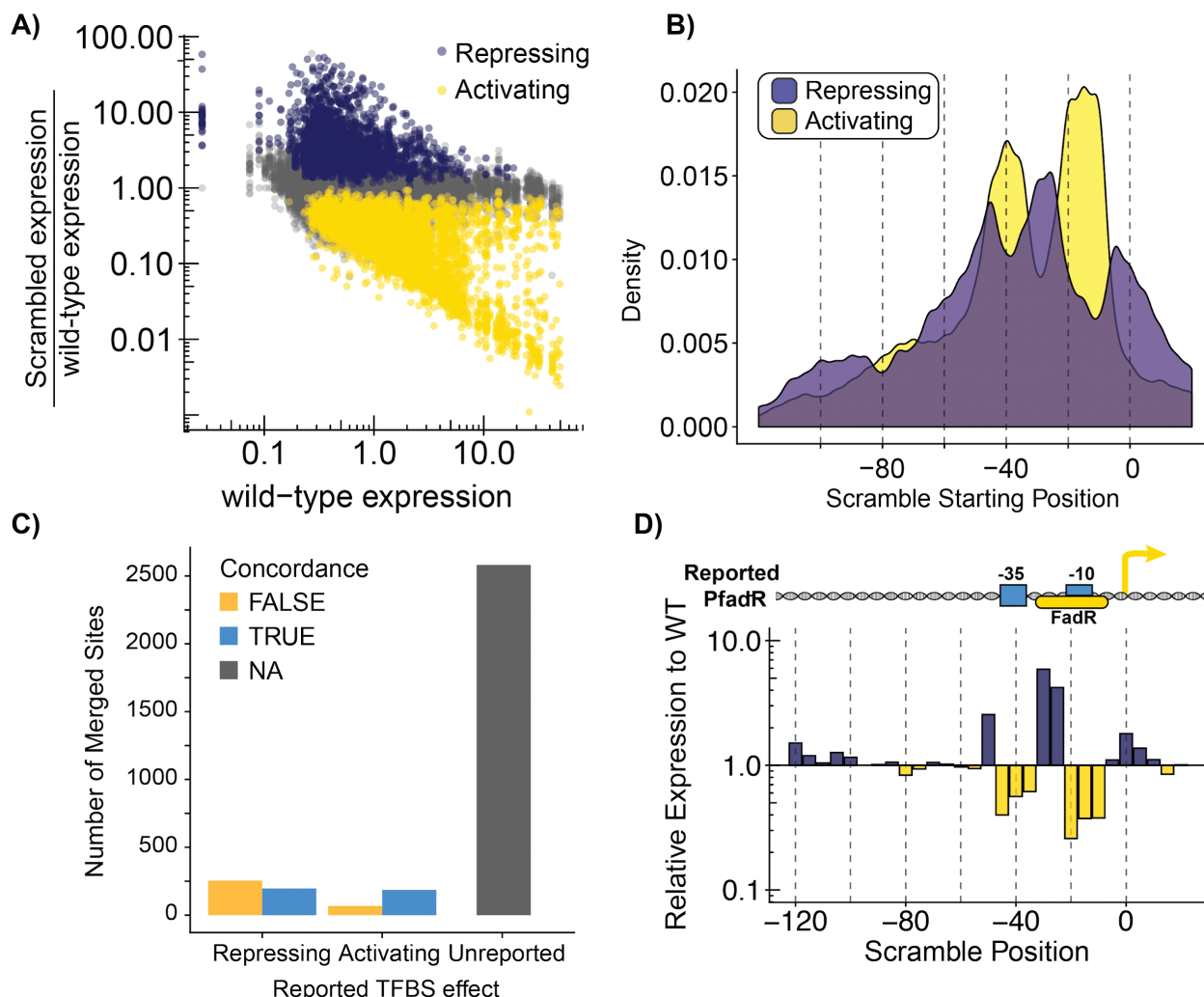
735
 736 **Figure 4) The *E. coli* promoter landscape dynamically responds to environmental conditions. A)**
 737 Shared and unique promoter regions are found between LB and glucose minimal media. **B)**
 738 Genes activated by promoters in glucose minimal media are enriched for amino acid-related
 739 genes according to RAST subsystem annotations. **C)** Occurrence of reported transcription
 740 factor binding sites in promoter regions activated in LB compared to glucose minimal media
 741 (M9). Black lines indicate 2-fold enrichment threshold. **D)** Number of binding sites per
 742 transcription factor within activated promoter regions. A median of four sites per transcription
 743 factors were activated in LB and a median of five sites in M9.

744
 745
 746



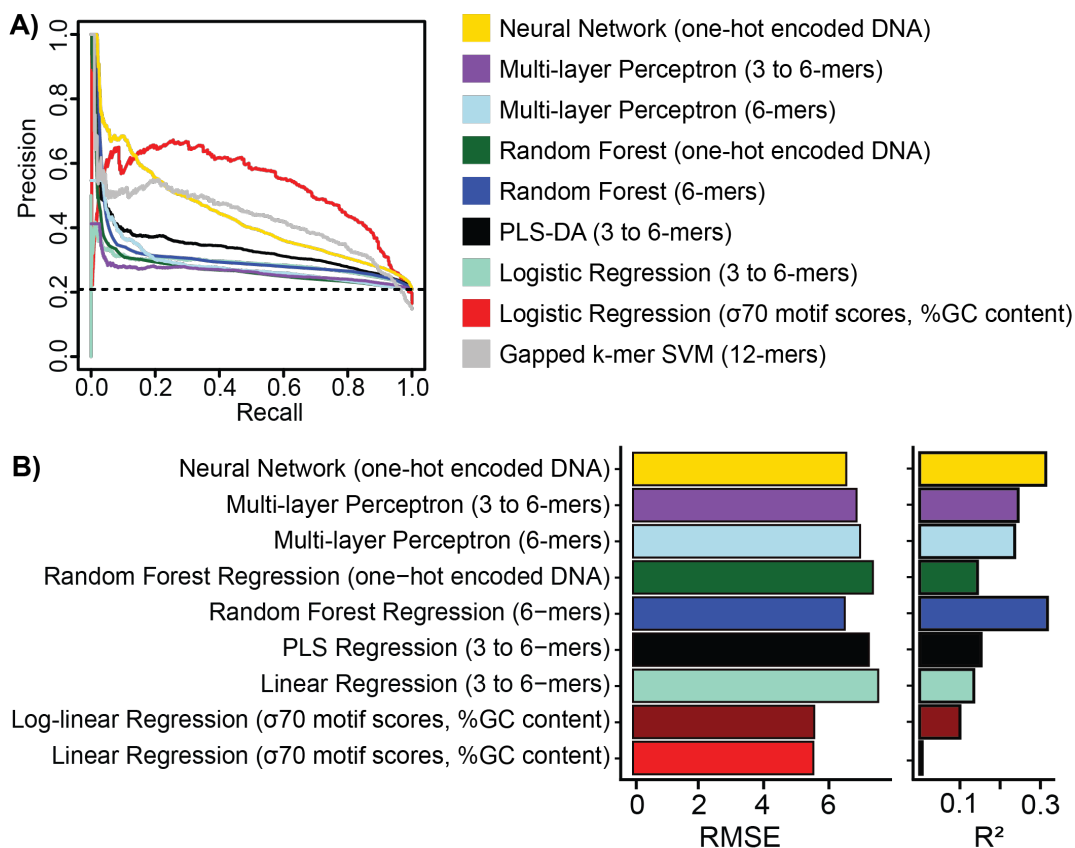
747
748 **Figure 5) Scanning mutagenesis of 2,057 TSS-associated promoters identifies known and**
749 **novel regulatory motifs.** **A)** Scanning mutagenesis of 2,057 *E. coli* promoters to identify
750 regulatory elements. For each promoter, 10 bp regions were mutated across the full length of
751 the promoter at 5 bp intervals. **B)** Mutating each position across *E. coli* promoters identifies
752 sequences that activate and repress promoter activity. Rows are rearranged using hierarchical
753 clustering and the intensities are normalized within each row. **C)** Scanning mutagenesis of the
754 well-characterized (Left) *lacZYA* and (Right) *relBE* promoters captures known regulatory
755 elements. **D)** Scanning mutagenesis of the newly characterized (Left) *cfa* and (Right) *rpsL*
756 promoters identifies regions encoding regulation within these promoters.

757
758
759
760
761
762
763



764
 765 **Figure 6) Global identification of 3,317 *E. coli* regulatory motifs by scanning mutagenesis.** **A)** We
 766 identified scrambled regulatory regions that significantly increase (N = 1,885) or decrease (N=5,408)
 767 expression when scrambled relative to the unscrambled promoter. Data are colored by whether the
 768 regulatory region activates or represses activity of the promoter. **B)** Activating promoter sequences are
 769 enriched at the -10, -35, and -80 positions whereas repressing sequences are enriched at +1, -20, and -50
 770 positions. **C)** Identified regulatory regions overlapping reported TFBS annotations shows mixed
 771 concordance with reported effects; 77.8% (2,583/3,317) of identified regulatory regions are unreported by
 772 RegulonDB. **D)** Scanning mutagenesis of the *FadR* promoter (bottom) identifies a repressing sequence
 773 near the -30 that has been reported to be activating (top).

774



775

776

777

778

779

780

781

782

783

784

785

786

787

788

789

790

791

792

793

794

795

796

797

798

Figure 7) **Various machine learning models for promoter activity classification and regression. A)** Performance of various models to classify promoter sequences. Convolutional neural networks performed best in the lower recall range, while logistic regression based on simple hand-crafted features performs better in the higher recall range. Dashed line represents the expected performance from random prediction using full library. **B)** Performance of regression models to predict a quantitative level of promoter activity. We evaluated performance using both root mean squared error (RMSE) and coefficient of determination (R^2) on the held-out test set. Similar to classification, convolutional neural networks performed the best with the lowest RMSE and highest R^2 .

799 **Methods**

800

801 **Strains**

802

803 All experiments were performed in the *E. coli* MG1655 background⁹⁵ which carries the following
804 genotype: F-, λ , *rph-1* (Yale Coli Genetic Stock Center no. 6300). For the genomically-integrated
805 MPRA, previously reported strains³⁸ with engineered landing pads in the right midreplicore
806 (*essQ-cspB* intergenic locus, Addgene no. 110244), chromosomal terminus (*nth-ydgR* intergenic
807 locus, Addgene no. 110245), and left midreplicore (*ybbD-ylbG* intergenic locus, Addgene no.
808 110243) were used. Briefly, these landing pads encode a fluorescent mCherry reporter as well as
809 chloramphenicol resistance, both of which are flanked by *loxP* sites for recombination-mediated
810 cassette exchange.

811

812 **TSS library design**

813

814 The TSS library incorporates all TSSs from the RegulonDB database⁸³ (Version 8.0,
815 genome version U00096.2) and those identified in two recent genome-wide TSS mapping
816 studies^{18,19}. Recent work provides evidence that most regulatory motifs fall within 100 bp
817 upstream of the TSS³⁹ and the initial transcribed region (+1 to +20) can also influence gene
818 expression. Thus, each TSS was synthesized embedded in its local sequence context -120 to +30
819 relative to the TSS, capturing a majority of the *cis*-regulatory elements. There were 23,798 unique
820 TSSs across all three sources, many of which were a few base pairs away from each other. We
821 minimized redundancy and collapsed together TSSs within 20 bp and selected the most upstream
822 TSS for our library, yielding 17,635 TSSs for the final synthesized library. Additionally, we included
823 500 negative controls from the *E. coli* genome that are assumed to have minimal regulatory
824 activity. These were selected from 150 bp regions that are more than 200 bp from a TSS (on either
825 strand), and many fall within coding regions. We included a set of 112 short synthetic positive
826 controls that were previously characterized^{40,96} and span a wide range of expression.

827

828 **TSS library barcoding and cloning**

829

830 The TSS library was synthesized by Twist Biosciences and delivered lyophilized as a 26 pmol
831 pool. The library was resuspended in 100 μ L of TE pH 8.0 and 1 μ L was amplified for 12 cycles
832 using GU72 and GU116 with NEB Q5 High-Fidelity 2x Master Mix (#M0492L). Unless otherwise
833 stated, all amplifications were performed using this polymerase mixture. This product was then
834 ran on a 2% TAE agarose gel and approximately 200 bp amplicons were extracted using a
835 Zymoclean Gel DNA Recovery Kit (#D4008). For barcoding, 1 ng of this eluate was amplified for
836 9 cycles using primers GU72 and GU73. Following cleaning using a Zymo Clean and Concentrator
837 Kit (#D40140), the library was digested using NEB's SbfI-HF and XhoI.

838

839 The plasmid backbone, pLibacceptorV2 (Addgene #106250) was digested using SbfI-HF and Sall-
840 HF with the addition of rSAP (NEB #M0371S). The digested library was ligated into
841 pLibacceptorV2 using T7 DNA Ligase (NEB #M0318S), cloned into 5-alpha Electrocompetent *E.*

842 *coli* (NEB #C2989K), and plated on LB + kanamycin (25 ug/mL) yielding approximately 2.3 million
843 colonies estimated by counting simultaneously plated dilutions. After allowing for 24 hours of
844 growth on plates, the library was scraped and resuspended in LB, and then 800 million cells
845 (based on OD₆₀₀) were inoculated in 450 mL LB + kanamycin (25 ug/mL) overnight. Unless stated
846 otherwise, all plasmids were isolated using a Qiagen Plasmid Plus Maxiprep Kit (#12963) and
847 concentrated using a Promega Wizard SV Gel and PCR Clean-up System (#A9281).

848
849 In order to clone the RiboJ::sfGFP reporter construct, the library was digested using NEB's Bsal-
850 HF and NheI-HF with the addition of rSAP. The reporter construct was digested using NEB's Bsal-
851 HF and NcoI-HF. Similarly to the previous cloning step, the reporter was cloned into the library
852 using T7 DNA Ligase, cloned into 5-alpha electrocompetent *E. coli*, and plated on LB + kanamycin
853 (25 ug/mL), yielding 6.8 million colonies. The completed plasmid library was isolated as stated
854 above.

855

856 **Isolation of genomic fragment library**

857

858 To isolate genomic fragments, 10 ug of *E. coli* MG1655 gDNA was sheared using a Covaris
859 focused ultra-sonicator. The settings used were as follows: Duty factor was set to 10%, Intensity
860 was set to 4, cycles/burst was set to 200, and time was 60 seconds. The sheared gDNA was ran
861 on a 3% TAE agarose gel and fragments between 200 and 300 bp were extracted using a
862 Zymoclean Gel DNA Recovery Kit and eluted in 18 uL water. All 18 uL of the extracted fragments
863 were end repaired using Enzymatics End Repair Mix (Part # Y9140-LC-L) following manufacturers
864 protocols, cleaned using 45 uL (1.8x volume) of Agencourt AMPure XP Beads (#A63880), and
865 eluted in 20 uL of water. The 20 uL eluate was A-tailed following the New England Biolabs
866 protocol:

867

868 Reaction:

869 20 uL End-repaired DNA
870 5 uL NEB Buffer 2 (10x)
871 0.5 uL dATP (10mM)
872 3 uL Klenow Fragment (3' -> 5' exo-) (Enzymatics #P7010-HC-L)
873 21.5 uL Nuclease-free water

874

875 The reaction was Incubated for 30 minutes at 37°C, then heat inactivated for 20 minutes at 75°C
876 before cleaning using 90 uL Agencourt AMPure XP beads and eluting in 20 uL water. Y-adapters
877 to facilitate fragment amplification and barcoding were ligated to the A-tailed fragments using
878 the following reaction mix:

879

880 Reaction:

881 20 uL A-tailed DNA
882 5 uL NEB T4 DNA Ligase Buffer (10x) (NEB #B0202S)
883 2 uL Y-adapter GU Y-Frag (25 uM)
884 1 uL NEB T4 DNA Ligase (NEB #M0202T)

885 22 uL Nuclease-free water

886

887 This reaction was incubated for 20 minutes at 25°C, heat inactivated for 20 minutes at 65°C, and
888 subsequently cleaned using 90 uL Agencourt AMPure XP beads and eluting in 12 uL nuclease-
889 free water.

890

891 **Barcoding and cloning of genomic fragment library**

892

893 To barcode the genomic fragments, 1 uL of the processed fragments was amplified for 13 cycles
894 using GU72 and GU116. This product was then cleaned using a Zymo Clean and Concentrator Kit
895 and eluted in 12 uL nuclease-free water. For barcoding, 1 ng of this eluate was amplified for 10
896 cycles using primers GU72 and GU73. Following cleaning using a Zymo Clean and Concentrator
897 Kit (#D40140), the library was digested using NEB's SbfI-HF and XhoI.

898

899 This library was cloned following the same protocols as the TSS library. The transformation of
900 the barcoded library yielded approximately 3.3 million colonies and the transformation after
901 addition of the RiboJ::sfGFP yielded approximately 1.25 million colonies.

902

903 **Genomic promoter tiling library design**

904

905 We used a custom peak caller on the single-nucleotide resolution strand-specific expression
906 pileup generated from our genomic fragment library to define "peaks" of promoter activity. Our
907 peak calling method is simple and conservative, as we wanted to tile the most active regions and
908 keep the library size reasonable. We defined a peak as a continuous region with expression above
909 an empirically determined threshold. We considered a continuous range of thresholds and for
910 each evaluated the percentage of active TSSs, from our previous library, contained in a peak and
911 determined an expression level of 1.1 was sufficient and captured 90% of active TSSs (data not
912 shown). We required that each peak be at least 60 bp, and merged adjacent peaks that were within
913 40 bp, yielding 1753 and 1724 peaks for the minus and plus strands, respectively. We tiled each
914 peak by synthesizing 150 bp windows across the region, with no overlap between adjacent tiles,
915 yielding 48,491 peak tiles. Additionally, we included 1000 randomly generated 150 bp sequences
916 to test what fraction of random sequence can drive expression. We included the same set of
917 positive and negative controls as described in the TSS library design.

918

919 **Genomic promoter tiling library barcoding and cloning**

920

921 The active TSS mutagenesis library was synthesized by Agilent and delivered lyophilized as a 10
922 pmol pool. The library was resuspended in 100 uL of TE pH 8.0 and 1 uL was amplified for 10
923 cycles using GU120 and GU121. This product was then cleaned using a Zymo Clean and
924 Concentrator Kit and eluted in 12 uL nuclease-free water. For barcoding, 1 ng of this eluate was
925 amplified for 8 cycles using primers GU120 and GU122. Following cleaning using a Zymo Clean
926 and Concentrator Kit (#D40140), the library was digested using NEB's SbfI-HF and XhoI.

927

928 This library was cloned following the same protocols as the TSS library. The transformation of
929 the barcoded library yielded approximately 1.5 million colonies and the transformation after
930 addition of the RiboJ::sfGFP yielded approximately 5.2 million colonies.

931

932 **Active TSS mutagenesis design**

933

934 We systematically mutagenized all active TSSs from our initial TSS library to design a follow-up
935 library. We used 500 negative controls to classify the TSS library into active and inactive TSSs.
936 We set the active threshold at two standard deviations above the median expression for the
937 negative controls, resulting in 2,670 active TSSs. We mutagenized the active sequence by
938 scrambling 10 bp windows, sliding across the 150 bp at 5 bp intervals, resulting in 5 bp of overlap
939 between adjacent scrambles. We scrambled the sequence using the existing 10 bp to preserve
940 nucleotide content and selected the scramble that was most dissimilar to the original sequence
941 out of 100 scrambling attempts. Our final library included 59,653 scrambled sequences and 2,057
942 unscrambled sequences. We also included the same set of negative and positive controls as
943 described above for the TSS library, for a total library size of 62,322.

944

945 **Active TSS mutagenesis library barcoding**

946

947 The active TSS mutagenesis library was synthesized by Agilent and delivered lyophilized as a 10
948 pmol pool. The library was resuspended in 100 uL of TE pH 8.0 and 1 uL was amplified for 12
949 cycles using GU123 and GU124. This product was then cleaned using a Zymo Clean and
950 Concentrator Kit and eluted in 12 uL nuclease-free water. For barcoding, 1 ng of this eluate was
951 amplified for 10 cycles using primers GU123 and GU125. Following cleaning using a Zymo Clean
952 and Concentrator Kit (#D40140), the library was digested using NEB's SbfI-HF and XhoI.

953

954 This library was cloned following the same protocols as the TSS library. The transformation of
955 the barcoded library yielded approximately 3.7 million colonies and the transformation after
956 addition of the RiboJ::sfGFP yielded approximately 5.2 million colonies.

957

958 **Library Barcode mapping**

959

960 We used PCR to individually barcode each library sequence to quantitatively measure expression
961 in our MPRA. Prior to genome integration, DNA-sequencing was performed to computationally
962 map barcodes to sequences. A custom barcode mapper developed by Nathan Lubock⁹⁷ was
963 used to collapse reads into a barcode-sequence map. We used two filtering steps for barcode
964 quality. First, we required a minimum number of reads for every barcode, assuming reads that
965 appear once or twice correspond to sequencing errors. Second, BBMap⁹⁸ was used to align the
966 reads associated with a given barcode, and discarded barcodes that map to sequences that are
967 too dissimilar to one another. A Levenshtein distance of 30 was used to discard barcodes that
968 map to two very distinct sequences, while still allowing for a small number of sequence errors.

969

970

971 **Library integration into specific genomic loci**

972

973 Library integration was performed as previously described³⁸.

974

975 The isolated plasmid library was digested with Sall-HF and NheI-HF to eliminate incompletely
976 cloned plasmid before transformation into electrocompetent MG1655 with a landing pad
977 engineered in the *nth-ydgR* locus and plating on LB + kanamycin (25 ug/mL). Colonies were
978 resuspended in LB and 800 million cells were inoculated into 250 mL LB + kanamycin (25 ug/mL)
979 and grown overnight. Several 2 mL frozen aliquots were made of this overnight culture.

980

981 The library was integrated into the *nth-ydgR* locus as follows. A frozen aliquot of MG1655 with a
982 landing pad engineered in the reverse orientation at the *nth-ydgR* intergenic locus was
983 transformed with the library and grown overnight in 200 mL LB + kanamycin (25 ug/mL).
984 Following overnight growth, 400 million cells of this culture were seeded into 250 mL LB +
985 kanamycin (25 ug/mL) + 0.2% arabinose (g/mL) and grown for 24 hours. After integration of the
986 library, the plasmid backbone was removed through heat-curing. From the 24 hour induced
987 culture, 800 million cells were inoculated into 80 mL of LB + kanamycin (25 ug/mL) and grown at
988 42 °C for approximately 1.5 hours before reaching an OD₆₀₀ = 0.3. Upon reaching exponential
989 growth, 200 million cells from this culture library were plated and grown for 16 hours at 42 °C.
990 Heat-cured plates were scraped and resuspended in LB and 400 million cells were inoculated into
991 200 mL LB + kanamycin (25 ug/mL). This culture, consisting of our integrated and heat-cured
992 library, was grown overnight at 37 °C and several frozen 2 mL aliquots were made.

993

994 To test the TSS library in the *essQ-cspB* and *ybbD-ylbG* midreplichore regions, the same protocol
995 was followed using strains engineered with landing pads in these intergenic regions.

996

997 **Library growth and harvest for expression measurements**

998

999 To measure expression of all promoter libraries, libraries were grown and harvested as previously
1000 described³⁸ with minor changes to culture conditions.

1001

1002 For each library and biological replicates, a 2 mL frozen aliquot of the library was inoculated in
1003 200 mL LB (BD#244620) with 25 ug/mL of kanamycin and grown at 30 °C overnight. The
1004 overnight cultures were used to seed new cultures at OD₆₀₀ = .0005 and grown for approximately
1005 5.5 hours at 30 °C until reaching an OD₆₀₀ between = 0.5 and 0.55. The genomic fragment library
1006 was also grown in Minimal Media (Fisher Scientific #DF0485-17) with 0.2% glucose (g/mL) and
1007 25 ug/mL of kanamycin for 10 hours at 30 °C until reaching an OD₆₀₀ between = 0.5 and 0.55.
1008 Cultures were rapidly cooled to 0 °C in an ice slurry for two minutes. Three 50 mL aliquots were
1009 pelleted at 4 °C by centrifugation at 13,000xg for two minutes and the supernatants were poured
1010 out before snap-freezing the pellets in liquid nitrogen. Three 5 mL aliquots of each library were
1011 harvested using the same approach to be processed for genomic DNA extractions.

1012

1013

1014 **RNA and DNA sequencing library preparation**

1015

1016 RNA was extracted from 50 mL library pellets using a Qiagen RNEasy Midi kit (#75142) and 45
1017 ug of each extract was concentrated using a Qiagen Minelute Cleanup Kit (#74204). Barcoded
1018 cDNA was generated from 25 ug of each concentrated RNA extract using Thermo Fisher
1019 SuperScript IV (#18090010) primed with GU101. The manufacturer's protocol was followed aside
1020 from extending the reaction time to 1 hour at 52 °C. The cDNA reaction was cleaned using a Zymo
1021 Research DNA Clean and Concentrator kit (#D40140) before amplification. Barcoded cDNA was
1022 amplified via PCR for 13 cycles using primers GU59 and GU102. This reaction was cleaned using
1023 a Zymo Research DNA Clean and Concentrator Kit and 1 uL of this reaction was used in a second
1024 PCR for indexing and addition of flow cell adapters. The second PCR was for 8 cycles and utilized
1025 primers GU102 with either GU61, GU62, GU63, or GU64 (which add separate 6 bp indices).

1026

1027 gDNA was extracted from 5 mL cell library pellets using a Qiagen Genra Puregene kit (#158567).
1028 Barcoded DNA was amplified from 1 ug of gDNA via PCR for 12-15 cycles using primers GU59
1029 and GU60. The reaction was subsequently cleaned using a Zymo Research DNA Clean and
1030 Concentrator kit. To add sequencing adapters and indices to the library, 1 ng of this reaction was
1031 subject to a second PCR for 8 cycles using primers GU70 with either GU63, GU64, GU65, or GU66
1032 (which add separate 6 bp indices). RNA and DNA sequencing libraries were cleaned using a Zymo
1033 Research Clean and Concentrator Kit (#D40140) before quantification using an Agilent
1034 TapeStation.

1035

1036 For each library, eight separate sequencing libraries were prepared: Four sequencing libraries for
1037 each RNA/DNA with two biological replicates and two technical replicates of each biological
1038 replicate. Biological replicates originated from separately grown and harvested glycerol stocks of
1039 each library. For each biological replicate, two RNA/gDNA extractions and sequencing library
1040 preparations (technical replicates) were performed in parallel. Libraries were submitted to the
1041 Broad Stem Cell Research Center at UCLA for sequencing on a HiSeq2500 or to the UCLA
1042 Translational Pathology Core Laboratory for sequencing on a NextSeq500. Raw sequencing data
1043 and promoter expression measurements are available on NCBI's Gene Expression Omnibus
1044 (Accession no. GSE144621).

1045

1046 **RNA-Seq of MG1655 in M9 minimal Media and Rich LB media**

1047

1048 To compare the promoter landscape to local transcriptional levels, RNA-Seq was performed on
1049 MG1655 grown in M9 minimal media (BD Difco #248510) supplemented with 0.2% glucose, 2 mM
1050 magnesium sulfate, and 0.1 mM calcium chloride. Similarly, RNA-Seq was performed for MG1655
1051 grown in LB (BD#244620). Cells growth and RNA preps were prepared as previously described
1052 (see methods section titled: library growth and harvest for expression measurements). Samples
1053 were prepared using an Illumina TruSeq® Stranded mRNA Library Prep (#20020594) following
1054 manufacturers protocols to achieve strand-specific coverage. We note that no rRNA depletion
1055 was performed to preserve the fully intact transcriptional landscape. Samples were submitted to
1056 the UCLA TCGB sequencing core and sequenced on a Hiseq 4000.

1057 **Standardizing Promoter Expression Quantification and Activity Thresholding**

1058

1059 We processed the TSS, scramble, and peak tiling libraries using the same computational pipeline
1060 to facilitate comparisons between libraries. First, we use a set of 96 short synthetic positive
1061 controls, designed to span a range of activity^{40,96}, to fit a robust linear regression (rlm function
1062 from MASS package) with the TSS library as the reference. Each library is standardized
1063 independently to the TSS library using the set of positive controls present in both libraries. Next,
1064 for each library we independently determined the level of background noise based on the median
1065 of 500 negative controls and subtracted this background from the newly fitted measurements.
1066 These steps standardize our data so we can train jointly across all datasets.

1067

1068 **-10 Motif and -35 Motif characterization**

1069

1070 A position weight matrix from bTSSfinder was used to identify and score the best match to the -
1071 10 and -35 motifs within active tss-associated promoters, inactive tss-associated promoters, and
1072 a set of 500 negative controls. Best scores were reported regardless of position within the
1073 sequence. For all pairwise comparisons of active tss-associated promoters, inactive tss-
1074 associated promoters, and the negative controls, the distributions of motif scores were compared
1075 and a student's t-test was performed to determine significance.

1076

1077 **Genomic fragment processing, alignment and promoter landscape quantification**

1078

1079 To calculate fragment expression, we used measurements from DNA-seq and RNA-seq and
1080 excluded fragments with low expression (< 0.1) or high variance (5-fold difference in relative
1081 expression between biological replicates). To identify the coordinates of genomic fragments
1082 assayed using the MPRA, fragment sequences were aligned using bowtie2⁹⁹ (version 2.3.4.3).
1083 To determine nucleotide-resolution calculations for promoter activity, we utilize the script,
1084 frag_expression_pileup.py. This script outputs WIG files in a strand-specific manner detailing the
1085 median expression of fragments overlapping each nucleotide position.

1086

1087 **Identification of minimal promoter regions**

1088

1089 To identify minimal sequences necessary for promoter activity, contiguous stretches of
1090 candidate promoter region peak tiles were grouped and the minimal shared overlapping region
1091 was identified. Peak tiles above the expression threshold were identified and grouped together if
1092 they shared an overlap of at least 110 bp of their 150 bp total length. The minimal region
1093 necessary for promoter activity was found by determining the overlap of the outermost
1094 sequences within a contiguous stretch of tiles.

1095

1096 **Determining promoter-gene associations**

1097

1098 To assign genomic promoter peaks to their regulated genes, peaks were first assigned specific
1099 nucleotide positions by identifying the maximum activity score within a peak. Promoter peaks

1100 were considered intragenic if their maximum scoring nucleotide overlapped with a gene
1101 coordinate. For peaks whose maximum scoring nucleotides were within intergenic regions,
1102 regulated genes were assigned by identifying the first downstream gene within 500 bp. Once gene
1103 associations were identified, promoter peaks were labeled sense or antisense depending on
1104 whether the regulated gene shared strand orientation with the promoter peak

1105

1106 **RNA-Seq alignment and genome transcript coverage**

1107

1108 RNA-Seq analysis was performed using the script *RNAseq_LB_processing.sh* or
1109 *RNAseq_M9_processing.sh*. This script trims reads using the trimmomatic software (ver.
1110 0.36+dfsg-3) and aligned to the MG1655 reference genome (U00096.2) using Hisat2¹⁰⁰ (ver.
1111 2.1.0-1). Genome nucleotide-resolution coverage was determined using Samtools depth (ver. 1.7-
1112 1). Meta-analysis across gene groups (as in figure 3B), was performed using Deeptools¹⁰¹ (ver.
1113 2.5.6). Gene expression coverage (as in figure 4B) was calculated using custom script
1114 *wig_over_bed.py*, which calculates the total (.wig) coverage across reported *E. coli* genes. In all
1115 cases, default parameters were used with the exception of allowing for strand-specific
1116 quantifications.

1117

1118 **Amino acid and codon bias within intragenic promoters**

1119

1120 Amino and codon usage was characterized within intragenic promoters and compared to all *E.*
1121 *coli* coding regions. To identify intragenic promoters, minimal regions necessary for promoter
1122 activity were identified by cross referencing genomic coordinates to reported genes. Reported
1123 gene coordinates were acquired from RegulonDB Version 8.0⁸³. Once intragenic promoters were
1124 identified, nucleotide triplets were extracted while conserving the reading frame of the
1125 overlapping gene. Similarly, nucleotide triplets were extracted from all reported *E. coli* coding
1126 regions after filtering out sequences which did not have nucleotide lengths of a multiple of three.
1127 For these extracted sequences, codon frequencies were normalized to their relative abundance
1128 compared to other codons encoding the same amino acid. Amino acid frequencies were
1129 normalized to the total number of amino acids within each group. Significantly enriched or
1130 depleted codons were identified by performing a chi-squared test within each amino acid group
1131 and adjusting the p-value using FDR. Significantly enriched or depleted amino acids were
1132 identified by performing a chi-squared test for each amino acid relative to the total pool of amino
1133 acids and adjusting the p-value using FDR.

1134

1135 **Comparison of condition-dependent promoter and gene activation between rich and minimal 1136 media**

1137

1138 To identify condition specific promoters, coordinates of candidate promoter regions identified in
1139 both M9 and LB conditions were compared to identify overlaps. Coordinates of promoter peaks
1140 were cross compared between conditions using the bedtools intersect tool (bedtools v2.27.1)
1141 and considered unique to a particular condition if they had no overlap between conditions. To
1142 identify regions that were activated between conditions, we compared the relative promoter

1143 activity between conditions at all positions in the genome and identified stretches greater than
1144 60 bp that exhibited over 2-fold difference in activity. Regions were called using custom script
1145 run_differential_wig.sh available on the Github repository. To identify genes being expressed by
1146 differentially active regions, intergenic differentially active regions and matched these to the
1147 nearest downstream gene within 500 bp.

1148

1149 **Identification of SEED subsystem annotations enriched in differentially activated genes**

1150

1151 To identify genetic functions associated with condition-dependent genes, the *E. coli* MG1655 K-
1152 12 genome (Genbank: U00096.2) was annotated using the SEED and RAST webserver^{67,68}. Genes
1153 within 500 bp downstream of promoter regions activated by condition were identified and
1154 associated with activation in LB or minimal media. For each media condition, genes were grouped
1155 by functional categories and the number of genes for each category was tallied.

1156

1157 **Identification of condition dependent TFBSs**

1158

1159 The TFBS content of promoter peaks unique to each condition was evaluated by cross-
1160 referencing with TFBSs reported by RegulonDB⁸³ (Release 8.8). Genomic regions activated in
1161 each condition were assigned TFBSs based on overlapping genomic coordinates using the
1162 bedtools intersect tool (bedtools v2.27.1) with default parameters and ignoring strand
1163 assignments. Incidents of each TFBS overlap were quantified between conditions and normalized
1164 to incidents per 100,000 bp of promoter peak sequence.

1165

1166

1167 **Identification of statistically significant scrambling promoter variants**

1168

1169 We identified scrambling promoter variants that significantly altered expression compared to the
1170 wild-type (WT) variant in the script scramble_ttest.Rmd. We considered each scramble and
1171 barcode combination as an independent observation, rather than summarizing expression as an
1172 average across all barcodes. A two-sample two-sided Student's t-test (t.test) was performed to
1173 test for a significant difference in mean expression levels between barcodes for a scrambled
1174 variant and barcodes for the corresponding WT variant. We performed multiple testing correction
1175 and identified 1,885 scrambles that increase expression and 5,408 that decrease expression
1176 relative to the WT variant, at a false discovery rate of 1%.

1177

1178 Next, bedtools merge was used to merge overlapping adjacent scramble variants to produce
1179 "merged" scrambles. These merged sites correspond to a continuous scrambled region that
1180 induced significant changes in expression. We identified 1,414 merged scrambles that increased
1181 expression and 1,903 merged scrambles that decreased expression, and scrambles were merged
1182 separately based on effect.

1183

1184

1185

1186 **Comparison of identified regulatory regions to RegulonDB annotations**

1187

1188 We compared our identified merged scramble sites to existing RegulonDB annotations. We used
1189 bedtools intersect and required that 10% of the TFBS overlapped with a merged scramble site to
1190 count as an overlap. Next, we assessed whether the expression effect seen in our MPRA agreed
1191 with the direction of effect of the TFBS as indicated in RegulonDB. A merged scramble site was
1192 marked as “concordant” if any of the component scrambles agreed with existing annotation, and
1193 not concordant otherwise.

1194

1195 **Machine learning models**

1196

1197 We implemented several machine learning models, independently trained for both classification
1198 and regression. All reproducible code is provided in the Github
1199 (https://github.com/KosuriLab/ecoli_promoter_mpra.git) and we will briefly describe each model
1200 and the appropriate parameters or implementation details.

1201

1202 *Data processing*

1203

1204 We standardized all datasets as detailed above in “Universal Promoter Expression Quantification
1205 and Activity Thresholding”. Next, we split our data, using custom scripts, into 75%/25% for
1206 training/testing based on genomic location, ensuring the splits are equidistant from the origin, to
1207 avoid overfitting (`define_genome_splits.py`). Briefly, we split the genome into eight chunks, with
1208 the first and last chunk adjacent to the origin of replication. We designated the second and
1209 seventh chunk as the test set and remaining chunks as training set. This splitting maintains
1210 roughly the same distance from the origin between the training and test sets to avoid any potential
1211 effects of genome location. Many of our library designs include high overlap between adjacent
1212 positions in the genome. Splitting by genome location mitigates inflated performance due to
1213 highly similar sequences present in both train and test sets. Across the three libraries (TSS, peak
1214 tiling, scramble) there are 87,164 training samples and 30,392 test samples.

1215

1216 We trained models for both regression and classification. Our data was skewed toward negative
1217 examples, with many samples near our determined threshold. For classification, we created a
1218 buffer around the threshold and only include sequences with expression ≤ 0.75 as negatives and
1219 ≥ 1.25 as positives and labeled sequences as active or inactive. Our training set was reduced to
1220 53,326 samples and testing set to 18,567 samples.

1221

1222 We used the classification models to predict probabilities, instead of the class, to derive
1223 precision-recall curves.

1224

1225 *Simple model with promoter features*

1226

1227 For the models in this section we created features only for the TSS library because it is closest to
1228 endogenous sequence and is a smaller dataset. The training and test sets were split by genomic
1229 location, as described above, with 13,118 training samples and 4549 testing samples.

1230

1231 We created a simple model which incorporates four features related to promoter function. We
1232 calculated the maximum position weight matrix (PWM) score using motifs from bTSSfinder¹⁰²
1233 for both the -10 and -35 core promoter motifs. We scanned the -10 and -35 PWM individually and
1234 took the max score at any position using scoring functions from the Bioconductor package
1235 Biostrings¹⁰³. Next, we scanned the sequence with -10 and -35 PWM jointly, allowing either 16, 17,
1236 or 18bp spacing in between the PWMs, reflecting common spacer lengths between core motifs.
1237 We assigned the “paired” max score as the max score at any position in the sequence across the
1238 three length options. Finally, we calculated the GC content (percentage) as this has been shown
1239 to be negatively correlated with promoter strength¹⁰⁴. We constructed models in R with these four
1240 features and fit 1) a linear regression (lm), 2) a linear regression on the log-transformed
1241 expression values (lm), and 3) a logistic regression (glm, family = ‘binomial’, type = ‘response’).

1242

1243 We trained the gapped k-mer SVM (gkm-SVM¹⁰⁵) model on only the TSS dataset because the
1244 model is suited for training sets < 20,000. The training and test sets were split by genome position
1245 as described above. We specified a word length = 10 with 8 informative columns (L = 10, K = 8).

1246

1247 *K-mer frequencies and simple models (linear regression, logistic regression, partial least squares*
1248 *regression, partial least squares discriminant analysis)*

1249

1250 All of the models described in the remaining sections were trained using all three combined
1251 datasets, as described above.

1252

1253 We created a feature set based on k-mer frequencies, with k-mers ranging in length from 3 to 6-
1254 mers. We generated feature sets and trained models in python. For simpler models we performed
1255 an additional feature selection step using custom scripts (kmer_feature_generator.py).

1256

1257 We trained four models:

- 1258 ● linear regression (statsmodel.api.OLS)
- 1259 ● logistic regression (sklearn.linear_model.LogisticRegression())
- 1260 ● partial least squares regression (sklearn.cross_decomposition.PLSRegression())
- 1261 ● partial least squares discriminant analysis
1262 (sklearn.cross_decomposition.PLSRegression() on binary dependent variable)

1263

1264 For each k-mer, we computed the frequency in a set of random genomic sequences, the same
1265 length and size of the training set. We include a k-mer if the absolute correlation with expression
1266 is greater than the “random” k-mer frequency, resulting in 4800/5440 filtered k-mers. We chose
1267 partial least squares regression because it projects the input features onto a new space and is
1268 better equipped to handle a large number of features with high collinearity.

1269 *Random forest regression and classification*

1270

1271 Next, we trained a random forest, for both regression
1272 (`sklearn.ensemble.RandomForestRegressor()`) and classification
1273 (`sklearn.ensemble.RandomForestClassifier()`). We train on one-hot encoded DNA as a
1274 comparison to the neural network model, although random forest is not well suited to categorical
1275 input features. To compensate for this, we trained the random forest using frequencies of all 6-
1276 mers and observed improved performance.

1277

1278 *Multi-layer perceptron and neural networks*

1279

1280 We trained a multi-layer perceptron for both regression (`sklearn.neural_network.MLPRegressor()`)
1281 and classification (`sklearn.neural_network.MLPClassifier()`). MLPs are a class of feedforward
1282 artificial networks and are “vanilla” neural networks consisting of an input layer, hidden layer, and
1283 output layer. We used two different feature sets: frequency of all 3- to 6-mers and frequency of
1284 only 6-mers. Feature sets were standardized with `sklearn.preprocessing.StandardScaler()` to
1285 remove mean and scale to unit variance. We trained all four models with the following
1286 parameters: `alpha = 0.005`, `hidden_layer_sizes=(800, 30)`, `solver = 'lbfgs'`, `random_state=1`,
1287 `max_iter=10000`, `early_stopping=True`, `learning_rate='adaptive'`, `tol=1e-8`.

1288

1289 We trained a convolutional neural network (CNN) on one-hot encoded DNA sequence for both
1290 regression and classification. We performed hyperparameter tuning and training using ⁸⁴, a toolkit
1291 for working with CNNs built on keras. We performed a random hyperparameter search for a three-
1292 layer CNN for 100 combinations and the optimal parameters are listed below.

1293

1294 Regression:

- 1295 ● Dropout: 0.1340735187802852
- 1296 ● Pooling width: 16
- 1297 ● Convolutional filter width (for each layer): 16, 17, 18
- 1298 ● Number of filters (for each layer): 19, 39, 54

1299

1300 Classification:

- 1301 ● Dropout: 0.45541334972592196
- 1302 ● Pooling width: 7
- 1303 ● Convolutional filter width (for each layer): 8, 29, 29
- 1304 ● Number of filters (for each layer): 99, 87, 60

1305

1306

1307

1308

1309

1310

1311

1312
1313
1314
1315

- 1316 1. Jacob, F. & Monod, J. Genetic regulatory mechanisms in the synthesis of proteins. *J. Mol.*
1317 *Biol.* **3**, 318–356 (1961).
- 1318 2. Feklístov, A., Sharon, B. D., Darst, S. A. & Gross, C. A. Bacterial sigma factors: a historical,
1319 structural, and genomic perspective. *Annu. Rev. Microbiol.* **68**, 357–376 (2014).
- 1320 3. Haugen, S. P., Ross, W. & Gourse, R. L. Advances in bacterial promoter recognition and its
1321 control by factors that do not bind DNA. *Nat. Rev. Microbiol.* **6**, 507–519 (2008).
- 1322 4. Lee, D. J., Minchin, S. D. & Busby, S. J. W. Activating transcription in bacteria. *Annu. Rev.*
1323 *Microbiol.* **66**, 125–152 (2012).
- 1324 5. Browning, D. F. & Busby, S. J. W. Local and global regulation of transcription initiation in
1325 bacteria. *Nat. Rev. Microbiol.* **14**, 638–650 (2016).
- 1326 6. Johnson, X. B. & Hinton, D. M. Escherichia coli RNA polymerase recognition of a σ 70-
1327 dependent promoter requiring a -35 DNA element and an extended-10 TGN motif. *J.*
1328 *Bacteriol.* **188**, 8352–8359 (2006).
- 1329 7. Hook-Barnard, I. G. & Hinton, D. M. Transcription initiation by mix and match elements:
1330 flexibility for polymerase binding to bacterial promoters. *Gene Regul. Syst. Bio.* **1**, 275–293
1331 (2007).
- 1332 8. Murakami, K. S., Masuda, S., Campbell, E. A., Muzzin, O. & Darst, S. A. Structural basis of
1333 transcription initiation: an RNA polymerase holoenzyme-DNA complex. *Science* **296**, 1285–
1334 1290 (2002).
- 1335 9. Liu, B., Hong, C., Huang, R. K., Yu, Z. & Steitz, T. A. Structural basis of bacterial transcription
1336 activation. *Science* **358**, 947–951 (2017).
- 1337 10. Newberry, K. J. & Brennan, R. G. The structural mechanism for transcription activation by

- 1338 MerR family member multidrug transporter activation, N terminus. *J. Biol. Chem.* **279**,
1339 20356–20362 (2004).
- 1340 11. Lawson, C. L. *et al.* Catabolite activator protein: DNA binding and transcription activation.
1341 *Curr. Opin. Struct. Biol.* **14**, 10–20 (2004).
- 1342 12. Kinney, J. B., Murugan, A., Callan, C. G. & Cox, E. C. Using deep sequencing to characterize
1343 the biophysical mechanism of a transcriptional regulatory sequence. *Proceedings of the*
1344 *National Academy of Sciences* vol. 107 9158–9163 Preprint at
1345 <https://doi.org/10.1073/pnas.1004290107> (2010).
- 1346 13. Ishihama, A., Shimada, T. & Yamazaki, Y. Transcription profile of Escherichia coli: genomic
1347 SELEX search for regulatory targets of transcription factors. *Nucleic Acids Res.* **44**, 2058–
1348 2074 (2016).
- 1349 14. Peano, C. *et al.* Characterization of the Escherichia coli σ (S) core regulon by Chromatin
1350 Immunoprecipitation-sequencing (ChIP-seq) analysis. *Sci. Rep.* **5**, 10469 (2015).
- 1351 15. Bonocora, R. P., Smith, C., Lapierre, P. & Wade, J. T. Genome-Scale Mapping of Escherichia
1352 coli σ 54 Reveals Widespread, Conserved Intragenic Binding. *PLoS Genet.* **11**, e1005552
1353 (2015).
- 1354 16. Huerta, A. M. & Collado-Vides, J. Sigma70 promoters in Escherichia coli: specific
1355 transcription in dense regions of overlapping promoter-like signals. *J. Mol. Biol.* **333**, 261–
1356 278 (2003).
- 1357 17. Rhodius, V. A., Mutalik, V. K. & Gross, C. A. Predicting the strength of UP-elements and full-
1358 length E. coli σ E promoters. *Nucleic Acids Res.* **40**, 2907–2924 (2012).
- 1359 18. Conway, T. *et al.* Unprecedented High-Resolution View of Bacterial Operon Architecture
1360 Revealed by RNA Sequencing. *mBio* vol. 5 Preprint at [https://doi.org/10.1128/mbio.01442-](https://doi.org/10.1128/mbio.01442-14)
1361 14 (2014).
- 1362 19. Thomason, M. K. *et al.* Global transcriptional start site mapping using differential RNA

- 1363 sequencing reveals novel antisense RNAs in Escherichia coli. *J. Bacteriol.* **197**, 18–28
1364 (2015).
- 1365 20. Gama-Castro, S. *et al.* RegulonDB (version 6.0): gene regulation model of Escherichia coli K-
1366 12 beyond transcription, active (experimental) annotated promoters and Textpresso
1367 navigation. *Nucleic Acids Res.* **36**, D120–D124 (2008).
- 1368 21. Belliveau, N. M. *et al.* Systematic approach for dissecting the molecular mechanisms of
1369 transcriptional regulation in bacteria. *Proc. Natl. Acad. Sci. U. S. A.* **115**, E4796–E4805
1370 (2018).
- 1371 22. Block, D. H. S., Hussein, R., Liang, L. W. & Lim, H. N. Regulatory consequences of gene
1372 translocation in bacteria. *Nucleic Acids Res.* **40**, 8979–8992 (2012).
- 1373 23. Sousa, C., de Lorenzo, V. & Cebolla, A. Modulation of gene expression through
1374 chromosomal positioning in Escherichia coli. *Microbiology* **143 (Pt 6)**, 2071–2078 (1997).
- 1375 24. Yin, J. *et al.* Effects of chromosomal gene copy number and locations on
1376 polyhydroxyalkanoate synthesis by Escherichia coli and Halomonas sp. *Appl. Microbiol.*
1377 *Biotechnol.* **99**, 5523–5534 (2015).
- 1378 25. Kuhlman, T. E. & Cox, E. C. Gene location and DNA density determine transcription factor
1379 distributions in Escherichia coli. *Mol. Syst. Biol.* **8**, 610 (2012).
- 1380 26. Sendy, B., Lee, D. J., Busby, S. J. W. & Bryant, J. A. RNA polymerase supply and flux through
1381 the lac operon in Escherichia coli. *Philos. Trans. R. Soc. Lond. B Biol. Sci.* **371**, (2016).
- 1382 27. Scholz, S. A. *et al.* High-Resolution Mapping of the Escherichia coli Chromosome Reveals
1383 Positions of High and Low Transcription. *Cell Syst* **8**, 212–225.e9 (2019).
- 1384 28. Bryant, J. A., Sellars, L. E., Busby, S. J. W. & Lee, D. J. Chromosome position effects on gene
1385 expression in Escherichia coli K-12. *Nucleic Acids Res.* **42**, 11383–11392 (2014).
- 1386 29. Brambilla, E. & Sclavi, B. Gene regulation by H-NS as a function of growth conditions
1387 depends on chromosomal position in Escherichia coli. *G3* **5**, 605–614 (2015).

- 1388 30. Brewster, R. C. *et al.* The transcription factor titration effect dictates level of gene
1389 expression. *Cell* **156**, 1312–1323 (2014).
- 1390 31. Chen, H., Shiroguchi, K., Ge, H. & Xie, X. S. Genome-wide study of mRNA degradation and
1391 transcript elongation in Escherichia coli. *Molecular Systems Biology* vol. 11 808–808
1392 Preprint at <https://doi.org/10.15252/msb.20159000> (2015).
- 1393 32. Esquerre, T. *et al.* Dual role of transcription and transcript stability in the regulation of gene
1394 expression in Escherichia coli cells cultured on glucose at different growth rates. *Nucleic
1395 Acids Res.* **42**, 2460–2472 (2013).
- 1396 33. Shearwin, K., Callen, B. & Egan, J. Transcriptional interference – a crash course. *Trends in
1397 Genetics* vol. 21 339–345 Preprint at <https://doi.org/10.1016/j.tig.2005.04.009> (2005).
- 1398 34. Callen, B. P., Shearwin, K. E. & Egan, J. B. Transcriptional interference between convergent
1399 promoters caused by elongation over the promoter. *Mol. Cell* **14**, 647–656 (2004).
- 1400 35. Brophy, J. A. N. & Voigt, C. A. Antisense transcription as a tool to tune gene expression.
1401 *Mol. Syst. Biol.* **12**, 854 (2016).
- 1402 36. Warman, E. A. *et al.* Widespread divergent transcription from bacterial and archaeal
1403 promoters is a consequence of DNA-sequence symmetry. *Nat Microbiol* **6**, 746–756 (2021).
- 1404 37. Yus, E. *et al.* Transcription start site associated RNAs in bacteria. *Mol. Syst. Biol.* **8**, 585
1405 (2012).
- 1406 38. Urtecho, G., Tripp, A. D., Insigne, K. D., Kim, H. & Kosuri, S. Systematic Dissection of
1407 Sequence Elements Controlling σ_{70} Promoters Using a Genomically Encoded Multiplexed
1408 Reporter Assay in Escherichia coli. *Biochemistry* **58**, 1539–1551 (2019).
- 1409 39. Garcia, H. G. *et al.* Operator sequence alters gene expression independently of transcription
1410 factor occupancy in bacteria. *Cell Rep.* **2**, 150–161 (2012).
- 1411 40. Mutalik, V. K. *et al.* Precise and reliable gene expression via standard transcription and
1412 translation initiation elements. *Nat. Methods* **10**, 354–360 (2013).

- 1413 41. Enyeart, P. J. *et al.* Generalized bacterial genome editing using mobile group II introns and
1414 Cre-lox. *Mol. Syst. Biol.* **9**, 685 (2013).
- 1415 42. Yan, B., Boitano, M., Clark, T. A. & Ettwiller, L. SMRT-Cappable-seq reveals complex operon
1416 variants in bacteria. *Nat. Commun.* **9**, 3676 (2018).
- 1417 43. Grünberger, F., Ferreira-Cerca, S. & Grohmann, D. Nanopore sequencing of RNA and cDNA
1418 molecules in. *RNA* **28**, 400–417 (2022).
- 1419 44. Schneider, D. A., Ross, W. & Gourse, R. L. Control of rRNA expression in Escherichia coli.
1420 *Curr. Opin. Microbiol.* **6**, 151–156 (2003).
- 1421 45. Yona, A. H., Alm, E. J. & Gore, J. Random sequences rapidly evolve into de novo promoters.
1422 *Nat. Commun.* **9**, 1530 (2018).
- 1423 46. Cho, B.-K., Kim, D., Knight, E. M., Zengler, K. & Palsson, B. O. Genome-scale reconstruction
1424 of the sigma factor network in Escherichia coli: topology and functional states. *BMC Biol.*
1425 **12**, 4 (2014).
- 1426 47. Krummel, B. & Chamberlin, M. J. RNA chain initiation by Escherichia coli RNA polymerase.
1427 Structural transitions of the enzyme in early ternary complexes. *Biochemistry* **28**, 7829–
1428 7842 (1989).
- 1429 48. Hawley, D. K. & McClure, W. R. Compilation and analysis of Escherichia coli promoter DNA
1430 sequences. *Nucleic Acids Res.* **11**, 2237–2255 (1983).
- 1431 49. He, W., Jia, C., Duan, Y. & Zou, Q. 70ProPred: a predictor for discovering sigma70 promoters
1432 based on combining multiple features. *BMC Syst. Biol.* **12**, 44 (2018).
- 1433 50. Weaver, J., Mohammad, F., Buskirk, A. R. & Storz, G. Identifying Small Proteins by Ribosome
1434 Profiling with Stalled Initiation Complexes. *MBio* **10**, (2019).
- 1435 51. Dornenburg, J. E., Devita, A. M., Palumbo, M. J. & Wade, J. T. Widespread antisense
1436 transcription in Escherichia coli. *MBio* **1**, (2010).
- 1437 52. Georg, J. & Hess, W. R. cis-antisense RNA, another level of gene regulation in bacteria.

- 1438 *Microbiol. Mol. Biol. Rev.* **75**, 286–300 (2011).
- 1439 53. Güell, M. *et al.* Transcriptome complexity in a genome-reduced bacterium. *Science* **326**,
1440 1268–1271 (2009).
- 1441 54. Ju, X., Li, D. & Liu, S. Full-length RNA profiling reveals pervasive bidirectional transcription
1442 terminators in bacteria. *Nat Microbiol* **4**, 1907–1918 (2019).
- 1443 55. Lloréns-Rico, V. *et al.* Bacterial antisense RNAs are mainly the product of transcriptional
1444 noise. *Sci Adv* **2**, e1501363 (2016).
- 1445 56. Brantl, S. Regulatory mechanisms employed by cis-encoded antisense RNAs. *Curr. Opin.*
1446 *Microbiol.* **10**, 102–109 (2007).
- 1447 57. Goodman, D. B., Church, G. M. & Kosuri, S. Causes and effects of N-terminal codon bias in
1448 bacterial genes. *Science* **342**, 475–479 (2013).
- 1449 58. Lamberte, L. E. *et al.* Horizontally acquired AT-rich genes in *Escherichia coli* cause toxicity
1450 by sequestering RNA polymerase. *Nat Microbiol* **2**, 16249 (2017).
- 1451 59. Singh, S. S. *et al.* Widespread suppression of intragenic transcription initiation by H-NS.
1452 *Genes Dev.* **28**, 214–219 (2014).
- 1453 60. Bolotin, E. & Hershberg, R. Horizontally Acquired Genes Are Often Shared between Closely
1454 Related Bacterial Species. *Front. Microbiol.* **8**, 1536 (2017).
- 1455 61. Constantinidou, C. *et al.* A reassessment of the FNR regulon and transcriptomic analysis of
1456 the effects of nitrate, nitrite, NarXL, and NarQP as *Escherichia coli* K12 adapts from aerobic
1457 to anaerobic growth. *J. Biol. Chem.* **281**, 4802–4815 (2006).
- 1458 62. Tao, H., Bausch, C., Richmond, C., Blattner, F. R. & Conway, T. Functional Genomics:
1459 Expression Analysis of *Escherichia coli* Growing on Minimal and Rich Media. *J. Bacteriol.*
1460 **181**, 6425–6440 (1999).
- 1461 63. Yu, T. C. *et al.* Multiplexed characterization of rationally designed promoter architectures
1462 deconstructs combinatorial logic for IPTG-inducible systems. *Nat. Commun.* **12**, 325

- 1463 (2021).
- 1464 64. Massé, E. & Gottesman, S. A small RNA regulates the expression of genes involved in iron
1465 metabolism in *Escherichia coli*. *Proc. Natl. Acad. Sci. U. S. A.* **99**, 4620–4625 (2002).
- 1466 65. Massé, E., Vanderpool, C. K. & Gottesman, S. Effect of RyhB small RNA on global iron use in
1467 *Escherichia coli*. *J. Bacteriol.* **187**, 6962–6971 (2005).
- 1468 66. Barroga, C. F., Zhang, H., Wajih, N., Bouyer, J. H. & Hermodson, M. A. The proteins encoded
1469 by the rbs operon of *Escherichia coli*: I. Overproduction, purification, characterization, and
1470 functional analysis of RbsA. *Protein Science* vol. 5 1093–1099 Preprint at
1471 <https://doi.org/10.1002/pro.5560050611> (2008).
- 1472 67. Aziz, R. K. *et al.* The RAST Server: rapid annotations using subsystems technology. *BMC*
1473 *Genomics* **9**, 75 (2008).
- 1474 68. Overbeek, R. *et al.* The SEED and the Rapid Annotation of microbial genomes using
1475 Subsystems Technology (RAST). *Nucleic Acids Res.* **42**, D206–14 (2014).
- 1476 69. Fang, X. *et al.* Global transcriptional regulatory network for *Escherichia coli* robustly
1477 connects gene expression to transcription factor activities. *Proc. Natl. Acad. Sci. U. S. A.*
1478 **114**, 10286–10291 (2017).
- 1479 70. Braun, V. Iron uptake by *Escherichia coli*. *Front. Biosci.* **8**, s1409–21 (2003).
- 1480 71. Lee, J.-W. & Helmann, J. D. Functional specialization within the Fur family of
1481 metalloregulators. *Biometals* **20**, 485–499 (2007).
- 1482 72. Ireland, W. T. *et al.* Deciphering the regulatory genome of *Escherichia coli*, one hundred
1483 promoters at a time. (2020) doi:10.7554/eLife.55308.
- 1484 73. Flashner, Y. & Gralla, J. D. Dual mechanism of repression at a distance in the lac operon.
1485 *Proc. Natl. Acad. Sci. U. S. A.* **85**, 8968–8972 (1988).
- 1486 74. Czarniecki, D., Noel, R. J., Jr & Reznikoff, W. S. The -45 region of the *Escherichia coli* lac
1487 promoter: CAP-dependent and CAP-independent transcription. *J. Bacteriol.* **179**, 423–429

- 1488 (1997).
- 1489 75. Li, G.-Y., Zhang, Y., Inouye, M. & Ikura, M. Structural mechanism of transcriptional
1490 autorepression of the Escherichia coli RelB/RelE antitoxin/toxin module. *J. Mol. Biol.* **380**,
1491 107–119 (2008).
- 1492 76. Grogan, D. W. & Cronan, J. E., Jr. Cloning and manipulation of the Escherichia coli
1493 cyclopropane fatty acid synthase gene: physiological aspects of enzyme overproduction. *J.*
1494 *Bacteriol.* **158**, 286–295 (1984).
- 1495 77. Chang, Y. Y. & Cronan, J. E., Jr. Membrane cyclopropane fatty acid content is a major factor
1496 in acid resistance of Escherichia coli. *Mol. Microbiol.* **33**, 249–259 (1999).
- 1497 78. Ebright, R. H. Transcription activation at Class I CAP-dependent promoters. *Mol. Microbiol.*
1498 **8**, 797–802 (1993).
- 1499 79. Williams, S. M., Savery, N. J., Busby, S. J. & Wing, H. J. Transcription activation at class I
1500 FNR-dependent promoters: identification of the activating surface of FNR and the
1501 corresponding contact site in the C-terminal domain of the RNA polymerase alpha subunit.
1502 *Nucleic Acids Res.* **25**, 4028–4034 (1997).
- 1503 80. Browning, D. F. & Busby, S. J. The regulation of bacterial transcription initiation. *Nat. Rev.*
1504 *Microbiol.* **2**, 57–65 (2004).
- 1505 81. Rojo, F. Repression of transcription initiation in bacteria. *J. Bacteriol.* **181**, 2987–2991
1506 (1999).
- 1507 82. Liu, M., Tolstorukov, M., Zhurkin, V., Garges, S. & Adhya, S. A mutant spacer sequence
1508 between -35 and -10 elements makes the Plac promoter hyperactive and cAMP receptor
1509 protein-independent. *Proc. Natl. Acad. Sci. U. S. A.* **101**, 6911–6916 (2004).
- 1510 83. Salgado, H. *et al.* RegulonDB v8.0: omics data sets, evolutionary conservation, regulatory
1511 phrases, cross-validated gold standards and more. *Nucleic Acids Research* vol. 41 D203–
1512 D213 Preprint at <https://doi.org/10.1093/nar/gks1201> (2013).

- 1513 84. Paggi, J. *et al.* Predicting Transcriptional Regulatory Activities with Deep Convolutional
1514 Networks. *bioRxiv* 099879 (2017) doi:10.1101/099879.
- 1515 85. Movva, R. *et al.* Deciphering regulatory DNA sequences and noncoding genetic variants
1516 using neural network models of massively parallel reporter assays. *PLoS One* **14**, e0218073
1517 (2019).
- 1518 86. Schmidt, P., Brandt, D., Busche, T. & Kalinowski, J. Characterization of Bacterial
1519 Transcriptional Regulatory Networks in through Genome-Wide In Vitro Run-Off
1520 Transcription/RNA-seq (ROSE). *Microorganisms* **11**, (2023).
- 1521 87. Horwitz, M. S. & Loeb, L. A. Promoters selected from random DNA sequences. *Proceedings*
1522 *of the National Academy of Sciences* vol. 83 7405–7409 Preprint at
1523 <https://doi.org/10.1073/pnas.83.19.7405> (1986).
- 1524 88. Wolf, L., Silander, O. K. & van Nimwegen, E. Expression noise facilitates the evolution of
1525 gene regulation. *Elife* **4**, (2015).
- 1526 89. Fitzgerald, D. M., Stringer, A. M., Smith, C., Lapierre, P. & Wade, J. T. Genome-Wide Mapping
1527 of the Escherichia coli PhoB Regulon Reveals Many Transcriptionally Inert, Intragenic
1528 Binding Sites. *MBio* **14**, e0253522 (2023).
- 1529 90. LaFleur, T. L., Hossain, A. & Salis, H. M. Automated model-predictive design of synthetic
1530 promoters to control transcriptional profiles in bacteria. *Nat. Commun.* **13**, 5159 (2022).
- 1531 91. Einav, T. & Phillips, R. How the avidity of polymerase binding to the–35/–10 promoter sites
1532 affects gene expression. *Proceedings of the National Academy of Sciences* 201905615
1533 (2019).
- 1534 92. Selvarajoo, K. *Computational Biology and Machine Learning for Metabolic Engineering and*
1535 *Synthetic Biology*. (Springer Nature, 2022).
- 1536 93. de Boer, C. G. *et al.* Deciphering eukaryotic gene-regulatory logic with 100 million random
1537 promoters. *Nat. Biotechnol.* (2019) doi:10.1038/s41587-019-0315-8.

- 1538 94. Vaishnav, E. D. *et al.* The evolution, evolvability and engineering of gene regulatory DNA.
1539 *Nature* **603**, 455–463 (2022).
- 1540 95. Blattner, F. R. *et al.* The complete genome sequence of Escherichia coli K-12. *Science* **277**,
1541 1453–1462 (1997).
- 1542 96. Kosuri, S. *et al.* Composability of regulatory sequences controlling transcription and
1543 translation in Escherichia coli. *Proc. Natl. Acad. Sci. U. S. A.* **110**, 14024–14029 (2013).
- 1544 97. Jones, E. M. *et al.* Structural and Functional Characterization of G Protein-Coupled
1545 Receptors with Deep Mutational Scanning. *bioRxiv* 623108 (2019) doi:10.1101/623108.
- 1546 98. Bushnell, B. BMap short read aligner. (2016).
- 1547 99. Langmead, B. & Salzberg, S. L. Fast gapped-read alignment with Bowtie 2. *Nat. Methods* **9**,
1548 357–359 (2012).
- 1549 100. Kim, D., Langmead, B. & Salzberg, S. L. HISAT: a fast spliced aligner with low memory
1550 requirements. *Nat. Methods* **12**, 357–360 (2015).
- 1551 101. Ramírez, F. *et al.* deepTools2: a next generation web server for deep-sequencing data
1552 analysis. *Nucleic Acids Res.* **44**, W160–5 (2016).
- 1553 102. Shahmuradov, I. A., Mohamad Razali, R., Bougouffa, S., Radovanovic, A. & Bajic, V. B.
1554 bTSSfinder: a novel tool for the prediction of promoters in cyanobacteria and Escherichia
1555 coli. *Bioinformatics* **33**, 334–340 (2017).
- 1556 103. Pagès, H., Aboyoun, P., Gentleman, R. & DebRoy, S. Biostrings: Efficient manipulation of
1557 biological strings. *R package version 2*, (2017).
- 1558 104. Johns, N. I. *et al.* Metagenomic mining of regulatory elements enables programmable
1559 species-selective gene expression. *Nat. Methods* **15**, 323–329 (2018).
- 1560 105. Ghandi, M. *et al.* gkmSVM: an R package for gapped-kmer SVM. *Bioinformatics* **32**, 2205–
1561 2207 (2016).

1562

1563

1564