

1 **Title**

2  
3 Using deep sequencing to characterize the biophysical mechanism of a transcriptional  
4 regulatory sequence

5  
6 **Authors**

7  
8 Justin B. Kinney<sup>1,2,†</sup>, Anand Murugan<sup>1</sup>, Curtis G. Callan Jr.<sup>1,3</sup>, Edward C. Cox<sup>4</sup>

9  
10 **Affiliations**

- 11  
12 1. Department of Physics, Princeton University, Princeton, NJ 08544, USA  
13 2. Lewis-Sigler Institute, Princeton University, Princeton, NJ 08544, USA  
14 3. Princeton Center for Theoretical Science, Princeton University, Princeton, NJ 08544  
15 USA  
16 4. Department of Molecular Biology, Princeton University, Princeton, NJ 08544, USA

17  
18 **Corresponding Author Information**

19  
20 Please email correspondence to [jkinney@cshl.edu](mailto:jkinney@cshl.edu).

21  
22 **Additional title page footnotes**

23  
24 † Current address: Cold Spring Harbor Laboratory, Cold Spring Harbor, NY, 11724, USA

25  
26 **Abstract**

27  
28 Cells use protein-DNA and protein-protein interactions to regulate transcription. A  
29 biophysical understanding of this process has, however, been limited by the lack of  
30 methods for quantitatively characterizing the interactions that occur at specific  
31 promoters and enhancers in living cells. Here we show how such biophysical  
32 information can be revealed by a simple experiment in which a library of partially  
33 mutated regulatory sequences are partitioned according to their *in vivo* transcriptional  
34 activities and then sequenced *en masse*. Computational analysis of the sequence data  
35 produced by this experiment can provide precise quantitative information about how the  
36 regulatory proteins at a specific arrangement of binding sites work together to regulate  
37 transcription. This ability to reliably extract precise information about regulatory  
38 biophysics in the face of experimental noise is made possible by a recently identified  
39 relationship between likelihood and mutual information (1). Applying our experimental  
40 and computational techniques to the *E. coli lac* promoter, we demonstrate the ability to  
41 identify regulatory protein binding sites *de novo*, determine the sequence-dependent  
42 binding energy of the proteins that bind these sites and, importantly, measure the *in vivo*  
43 interaction energy between RNA polymerase and a DNA-bound transcription factor. Our  
44 approach provides a generally applicable method for characterizing the biophysical  
45 basis of transcriptional regulation by a specified regulatory sequence. The principles of  
46 our method can also be applied to a wide range of other problems in molecular biology.

47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65  
66  
67  
68  
69  
70  
71  
72  
73  
74  
75  
76  
77  
78  
79  
80  
81  
82  
83  
84  
85  
86  
87  
88  
89  
90  
91  
92

"\body"

## **Introduction**

Cells regulate transcription primarily through the binding of proteins to DNA binding sites within cis-regulatory modules (CRMs). Understanding how CRMs use different arrangements of binding sites to encode regulatory programs remains a major challenge for molecular biology. High-throughput methods have spurred great progress in cataloging the genome-wide distribution of binding sites (2, 3), and many techniques exist for characterizing the sequence-specificity of individual regulatory proteins (4-7). However, determining how a specific CRM integrates information from multiple DNA-bound proteins still requires a laborious series of biochemical experiments that typically provide only qualitative information (reviewed in ref. 8).

The *E. coli lac* promoter (9,10) is one of the few CRMs whose function is well understood at the biophysical level (11,12). Kuhlman et al. (12) were the first to prove that a certain aspect of this system — the up-regulation of transcription by the protein CRP (13) — could be quantitatively explained by a simple energetic interaction between CRP and the  $\sigma^{70}$ -dependent RNA polymerase holoenzyme (henceforth RNAP). To do this, Kuhlman et al. measured transcriptional activity resulting from different *in vivo* concentrations of active CRP and showed that the resulting functional form of this activity was consistent with a simple thermodynamic model (14). By fitting the defining parameters of this model to their data, Kuhlman et al. were then able to measure the *in vivo* interaction energy between CRP and RNAP.

Despite its success, Kuhlman et al.'s approach is not feasible as a general method for studying CRMs. First, it requires quantitative control over the *in vivo* concentrations of all of the proteins that bind the CRM of interest. Secondly, CRMs typically contain multiple binding sites for each operative regulatory protein, making it difficult to determine the specific role of each site simply by varying the concentration of the regulator. This latter fact has created difficulty for similar studies of eukaryotic enhancers (15).

We hypothesized that measuring the activities of a large number of CRMs containing scattered point mutations could provide information similar to that produced by Kuhlman et al.'s method. This approach would be feasible for studying arbitrary CRMs, and would allow the effect of each individual binding site to be characterized. Because point mutations tend to preserve the spatial arrangement of binding sites, such measurements would allow one to interrogate the same protein-DNA complexes that allow the wild-type CRM to function. A similar approach had been tried by Schneider and Stormo in 1989 (16), but the recent advent of ultra-high-throughput sequencing, together with new techniques in machine learning (1) led us to believe that this approach could be much more powerful than had previously been realized.

In this paper we report the application of this mutagenesis-based approach to the *E. coli lac* promoter. Fluorescence-activated cell sorting (FACS; ref. 17) and 454

93 pyrosequencing (18) was used to characterize the activities of ~200,000 *lac* promoters  
94 mutagenized in a 75 bp region containing the CRP and RNAP binding sites (Fig. 1A).  
95 The resulting sequence data allowed us to identify these binding sites *de novo*,  
96 determine the sequence-dependent binding energy of both CRP and RNAP, and  
97 measure the *in vivo* interaction energy between these two proteins in their native DNA-  
98 bound configuration. We note that previous attempts to determine *in vivo* protein-protein  
99 interaction energies from sequence data (19,20) required unproven assumptions about  
100 how arbitrary arrangements of DNA-bound proteins interact; our approach does not.

101  
102 In this way, we demonstrate how deep sequencing can be used to measure protein-  
103 DNA and protein-protein interaction energies in living cells. This ability should be useful  
104 for addressing many different questions in molecular biology.

### 105 **Interrogating a CRM with flow cytometry and deep sequencing**

106  
107  
108 We performed six experiments on region [-75:-1] of the *E. coli lac* promoter (Figs. 1A  
109 and 1B). These experiments, summarized in Table 1, differed in the positions within  
110 region [-75:-1] that were mutagenized, the strain of *E. coli* used, and the physiological  
111 conditions under which *lac* promoter function was characterized.

112  
113 The full-wt experiment used a library of reporter constructs, derived from pUA66-lacZ  
114 (Fig. 1C; ref. 21), in which region [-75:-1] of the *lac* promoter was mutagenized at 12%  
115 per nucleotide, yielding  $9 \pm 3$  substitution mutations per sequence. Wild-type *E. coli*  
116 (strain MG1655) were transformed with this plasmid library, after which GFP expression  
117 was induced during exponential growth in minimal media supplemented with glucose,  
118 cAMP, and 1 mM isopropyl  $\beta$ -D-1-thiogalactopyranoside (IPTG). FACS (17) was used to  
119 sort induced cells into 10 different batches, each cell according to its measured  
120 fluorescence (Fig. 1D). PCR was then used to tag mutant CRMs according to the batch  
121 in which each CRM was found (Fig. 1E). 454 pyrosequencing (18) of the resulting PCR  
122 amplicons yielded a list of 51,835 mutant CRMs and corresponding batches. The batch  
123  $\mu$  associated with each CRM  $\sigma$  serves as a noisy and qualitative measurement of that  
124 CRM's *in vivo* transcriptional activity. Experiments *crp*-wt and *rnep*-wt were performed  
125 the same way, but using mutant CRMs in which only the CRP or RNAP binding site was  
126 mutagenized (Table 1).

127  
128 The full-500, full-150, and full-0 experiments were performed using the same plasmid  
129 library as in full-wt, but with transcriptional activity assayed in *E. coli* strain TK310 (12)  
130 grown in 500  $\mu$ M, 150  $\mu$ M, or 0  $\mu$ M cAMP, respectively. Cells were also sorted into 5  
131 batches instead of 10. Strain TK310 lacks adenylate cyclase (*cyaA*, needed for cAMP  
132 synthesis) and phosphodiesterase (*cpdA*, which degrades cAMP) and is therefore  
133 unable to control its intracellular cAMP levels (12). Growing TK310 cells in media  
134 supplemented with different concentrations of cAMP thus allowed us to control the  
135 active intracellular concentration of CRP (just as Kuhlman et al. did). Importantly,  
136 though, the mechanism of *lac* promoter function is the same in both MG1655 and  
137 TK310 cells.

138

139 In total we obtained 220,591 mutant CRMs, with each CRM  $\sigma$  assigned a noisy,  
140 qualitative measurement  $\mu$  of its transcriptional activity (Fig. 1B). These sequences  $\sigma$   
141 and measurements  $\mu$  comprise the only data used in the analysis that follows.

142

### 143 Information footprints reveal functional binding sites

144

145 Our first goal was to identify nucleotide positions that affect expression, thereby locating  
146 all functional binding sites within the probed region of the *lac* promoter. For this we used  
147 mutual information, a fundamental quantity from information theory that provides the  
148 most general measure of dependence between any two variables (22). For each  
149 nucleotide position  $i$ , we computed the number of sequences in each batch  $\mu$  having  
150 each of the four possible bases  $b_i$ . We then used this information to compute the mutual  
151 information  $I(b_i; \mu)$  between  $b_i$  and  $\mu$  (see SI Appendix Sec. 2e for details), thereby  
152 quantifying the effect of the base at position  $i$  on our measurements. Plotting the mutual  
153 information  $I(b_i; \mu)$  for each position  $i$  produced what we call an "information footprint".

154

155 Fig. 2A shows the information footprint produced by the full-wt experiment. The known  
156 binding sites of CRP and RNAP are clearly visible and each has the expected bipartite  
157 structure: CRP binds DNA as a homodimer at [-72:-51] with a 6 bp gap between its two  
158 DNA-binding domains, while RNAP binding to [-41:-1] results from the recognition of  
159 separate sequence elements centered roughly on positions -10 and -35 (9,13).

160

161 The information values displayed in Fig. 2A are small, ranging from  $\sim 0.05$  bits down to  
162 values indistinguishable from zero. This is not surprising, since each  $I(b_i; \mu)$  measures  
163 the effect on transcription of just one out of 75 positions. But because our data set is  
164 large ( $N = 51,856$ ), these information values are determined very precisely — typically  
165 to within  $\sim 4 \times 10^{-4}$  bits. Fig. 2A thus displays meaningful information values ranging over  
166 two orders of magnitude. This high level of sensitivity can reveal aspects of sequence  
167 function not detectable by other methods, e.g. (23). For instance, small but significant  
168 effects on expression were observed between the -10 and -35 elements and between  
169 the CRP and RNAP sites, regions not previously thought to influence transcription at the  
170 *lac* promoter. Indeed, only 10 positions ([-75, -70, -48, -42, -41, -23, -22, -21, -19, -3])  
171 out of 75 show an insignificant effect on expression (i.e.  $P > 0.05$ ).

172

173 Information footprints can further allow one to determine which of the identified sites are  
174 involved in the response to a specific biochemical signal or genetic perturbation. Fig. 2B  
175 shows the information footprint from the full-0 experiment, in which TK310 cells were  
176 induced in the absence of cAMP, thereby inactivating intracellular CRP. The lack of  
177 active CRP is reflected in the near-total loss of information at positions [-72:-51]:  
178 information values at all but three positions within this 22 bp site differ insignificantly  
179 from zero (i.e.  $P > 0.05$ ). An intermediate reduction in information occurs in the full-150  
180 footprint, while the full-500 footprint closely resembles that of full-wt (SI Appendix Fig.  
181 S3). We note that the small but significant information values at positions [-64, -57, -52]  
182 in the full-0 footprint might result from interactions between DNA and the  $\alpha$  subunits of  
183 RNAP (24), or from alternative RNAP binding sites (25).

184

185 An important caveat to this analysis is that the quantitative features of an information  
186 footprint ultimately depend on the details of one's experiment, including the level of  
187 mutagenesis used in the initial CRM library and the specific fluorescence gates used for  
188 sorting cells. So while qualitative differences between footprints from different  
189 experiments can be revealing, it is difficult to draw conclusions from more subtle  
190 quantitative differences, such as the different shapes of the RNAP footprint in Figs. 2A  
191 and 2B. But as we show in the next section, explicit biophysical models can be fit to  
192 data in a way that does not depend on such experimental details. Quantitative  
193 differences between models inferred from different experiments can, as a result, be  
194 revealing about underlying molecular mechanisms.

195

### 196 Model fitting in the presence of uncharacterized noise

197

198 Next we used our data to infer quantitative models for how *in vivo* protein-DNA and  
199 protein-protein interactions modulate transcription. By 'model' we mean a mathematical  
200 function that assigns to each sequence  $\sigma$  a predicted value  $x$  for some quantity of  
201 interest, such as the DNA binding energy of a regulatory protein, or the rate of  
202 transcription resulting from the interactions between multiple proteins. To infer a given  
203 model, we first assumed a specific mathematical formula for the model, then determined  
204 the values of model parameters by matching the sequence-dependent predictions  $x$  to  
205 our sequence-specific measurements  $\mu$ .

206

207 Such model fitting could be done in the standard Bayesian way if we knew the  
208 experimental 'error model'  $p(\mu|x)$  — the probability of obtaining a measurement  $\mu$  given  
209 an underlying quantity of interest  $x$  (such as binding energy or transcription rate). But in  
210 our experiments it was virtually impossible to accurately determine  $p(\mu|x)$  *a priori*. Many  
211 difficult-to-characterize noise processes, including stochastic transcription, variations in  
212 cell size, and noise in FACS measurements, contribute to the quantitative form of  $p(\mu|x)$ .  
213 Also, in the case where  $x$  represents the DNA-binding energy of a protein, we do not  
214 know *a priori* how the binding of that protein affects transcription; indeed, this is  
215 something we want to learn from the data. Kinney et al. (1) overcame this problem by  
216 computing likelihood in the presence of an explicitly uncertain error model. They  
217 showed that regardless of what  $p(\mu|x)$  actually is, the likelihood of a candidate model is  
218 well approximated by

219

$$220 \quad \mathbf{[1]} \quad p(\text{data} | \text{model}) \approx \text{const} \times 2^{Nl(x;\mu)}$$

221

222 in the limit where the number  $N$  of independently measured sequences is sufficiently  
223 large. Here  $l(x;\mu)$  is the 'predictive information' of the model — the mutual information  
224 between model predictions  $x$  and measurements  $\mu$ . Although Kinney et al. focused on  
225 the analysis of microarray data, Eq. 1 is applicable to any data set consisting of a large  
226 number of sequences and corresponding measurements. Kinney et al.'s approach  
227 therefore provides a practical substitute for standard likelihood-based inference when  
228 the experimental error model is either difficult to characterize or is unknowable *a priori*  
229 (see SI Appendix Sec. 2b for more discussion).

230

231 In the analysis that follows, we used a custom parallel tempering Monte Carlo algorithm  
232 to sample model parameters according to the right-hand-side of Eq. 1 (see SI Appendix  
233 Sec. 2f). This allowed us to determine not just the best values for model parameters, i.e.  
234 which values maximize predictive information  $I(x;\mu)$ , but also the uncertainty in each  
235 parameter due to finite data. Because  $N > 10^4$  for all of the experiments described in  
236 this paper, even changes as small as  $10^{-4}$  bits in the value of  $I(x;\mu)$  led to substantial  
237 changes in model likelihood. The large amount of data produced by our experiments  
238 thus allowed us to determine precise quantitative values for model parameters even  
239 though our measurements were noisy and qualitative.

240

### 241 **In vivo interaction energies from sequence data**

242

243 Having located the binding sites of both CRP and RNAP *de novo* using information  
244 footprints, we sought an explicit model for each protein's sequence-dependent binding  
245 energy. For this we used 'energy matrix' models: each base within a protein's binding  
246 site was assumed to contribute additively to the overall binding energy. These simple  
247 models have been shown to accurately describe a number of transcription factors, e.g.  
248 (26,27), though there are known exceptions (28,29).

249

250 We fit an energy matrix for CRP to positions [-74:-49] using the full-wt data set. Energy  
251 matrix elements were sampled, according to Eq. 1, using the predictive information  
252  $I(\varepsilon_c;\mu)$  where  $\varepsilon_c$  is CRP's predicted binding energy. The resulting optimal matrix is shown  
253 in Fig. 3A. We similarly fit an energy matrix to positions [-41:-1] to model RNAP's  
254 binding energy  $\varepsilon_r$  (Fig. 3B). CRP and RNAP energy matrices were also inferred from our  
255 five other data sets (SI Appendix Fig. S4A). We note that at this stage of our analysis  
256 we were able to determine each matrix only up to an unknown multiplicative constant,  
257 not in physical units such as kcal/mol (see SI Appendix Sec. 2c).

258

259 Unlike information footprints, these energy matrices are meant to capture intrinsic  
260 properties of the regulatory proteins, properties that should not depend on specific  
261 ways cells were sorted or on the level of mutagenesis used in the CRM library. The  
262 optimal matrices inferred from our six experiments (4 matrices for CRP, 5 for RNAP) are  
263 nearly identical, supporting this interpretation: CRP matrix elements derived from  
264 different experiments correlate by  $> 95\%$  (SI Appendix Fig. S4B), while RNAP matrix  
265 elements exhibit  $> 92\%$  correlation (SI Appendix Fig. S4C). Furthermore, each of these  
266 matrix models performs better on every one of our data sets than do any of the models  
267 for either CRP or RNAP currently in the literature (with two minor exceptions; see SI  
268 Appendix Sec. 2g). We find this level of quantitative agreement between experiments  
269 remarkable, considering that our six experiments used different promoter libraries,  
270 different *E. coli* strains (MG1655 or TK310), different inducing conditions (500  $\mu\text{M}$ , 150  
271  $\mu\text{M}$ , or 0  $\mu\text{M}$  cAMP), and different fluorescence gates for sorting cells. This close  
272 agreement in the face of important experimental differences attests to both the  
273 usefulness and correctness of Eq. 1.

274

275 Our inferred CRP and RNAP energy matrices recapitulate much of what is known about  
276 the sequence specificities of these two proteins. The known consensus sequences —

277 TGTGA(N)<sub>6</sub>TCACA for CRP (4) and TTGACA(N)<sub>18</sub>TATAAT for RNAP (30) — exactly  
 278 match the lowest energy sequences predicted by nearly every one of our matrix models.  
 279 The one exception is the RNAP matrix fit to full-0 data, which predicts that TTGATA will  
 280 have slightly lower energy than TTGACA in the -35 region. We note that every one of  
 281 our RNAP matrices also predicts that having a 'G' at position -14 increases RNAP  
 282 binding strength. In the literature this 'G' is said to create an "extended -10 promoter",  
 283 and such promoters are known to have increased transcriptional activity. Our CRP  
 284 matrices are also in qualitative agreement with previous *in vitro* measurements (31),  
 285 though there are some quantitative discrepancies.

286  
 287 Next we sought a quantitative understanding of how the interaction of CRP with RNAP  
 288 affects transcription. Kuhlman et al. (12) previously showed that a simple biophysical  
 289 model, based on equilibrium statistical mechanics (reviewed in ref. 14) accounted well  
 290 for the effect of cAMP on *lacZ* expression in TK310 cells. We hypothesized that using  
 291 energy matrices to describe the binding energies of CRP and RNAP within Kuhlman et  
 292 al.'s model, then fitting all model parameters to our data *de novo*, would allow us to  
 293 recover Kuhlman et al.'s results, including their measurement of the interaction energy  
 294 between CRP and RNAP.

295  
 296 Following Kuhlman et al., we assumed that the rate of transcription  $\tau$  at the *lac* promoter  
 297 is proportional to the occupancy of RNAP at its binding site in thermal equilibrium. This  
 298 model is quantitatively expressed as

299  
 300 **[2]** 
$$\tau = \tau_{\max} \frac{C_r e^{-\varepsilon_r/RT} + C_c C_r e^{-(\varepsilon_c + \varepsilon_r + \varepsilon_i)/RT}}{1 + C_c e^{-\varepsilon_c/RT} + C_r e^{-\varepsilon_r/RT} + C_c C_r e^{-(\varepsilon_c + \varepsilon_r + \varepsilon_i)/RT}}$$

301  
 302 where RNAP occupancy is given by the sum of Boltzmann weights corresponding to  
 303 physical states in which RNAP is bound, divided by the sum of weights for all possible  
 304 states of the system. These Boltzmann weights depend on (i) the CRP and RNAP  
 305 binding energies  $\varepsilon_c$  and  $\varepsilon_r$ , which we express in kcal/mol and normalize to be zero at  
 306 each wild-type *lac* promoter site, (ii) the concentrations  $C_c$  and  $C_r$  of CRP and RNAP,  
 307 expressed in units of each wild-type site's dissociation constant and (iii) the CRP-RNAP  
 308 interaction energy  $\varepsilon_i$ , expressed in kcal/mol.  $\tau_{\max}$  is the transcription rate resulting from  
 309 full RNAP occupancy.  $R = 1.98 \times 10^{-3}$  kcal/mol °K is the gas constant and  $T = 310$  °K (37  
 310 °C) is the temperature at which cells were induced.

311  
 312 Using  $l(\tau, \mu)$  evaluated on full-wt data, we fit all of the parameters defining  $\tau$ , including  $\varepsilon_i$ ,  
 313  $C_c$ , and the elements of the energy matrices used to compute  $\varepsilon_c$  and  $\varepsilon_r$ . Doing so we  
 314 inferred a CRP-RNAP interaction energy  $\varepsilon_i = -3.26 \pm 0.41$  kcal/mol. This value is  
 315 consistent with Kuhlman et al.'s measurement of -3.4 kcal/mol (12) and represents the  
 316 first time (to our knowledge) that the analysis of sequence data has been shown to  
 317 successfully yield the *in vivo* interaction energy between two proteins. This procedure  
 318 also yielded an *in vivo* CRP concentration of  $C_c = [\text{CRP}]/K_d^{wt} = 10^{-1.2 \pm 0.2}$ . Fig. 3C shows  
 319 these values for  $\varepsilon_i$  and  $C_c$ , as well as the optimal energy matrices for  $\varepsilon_c$  and  $\varepsilon_r$  inferred  
 320 by fitting  $\tau$ . These matrices closely resemble those in Fig. 3A and 3B, but, unlike the

321 matrices we inferred by separately fitting  $\epsilon_c$  and  $\epsilon_r$ , their elements are determined  
322 explicitly in physical units of kcal/mol. We note that fitting  $\tau$  to full-wt data provided no  
323 information about the value of either  $\tau_{max}$  or  $C_r$  (see SI Appendix Sec. 2c).

324

### 325 **Testing biochemical mechanisms by fitting a single model to multiple data sets**

326

327 cAMP is known to alter *lac* promoter activity by affecting CRP's ability to bind DNA, not  
328 CRP's interaction with RNAP. Both of these possibilities, though, are consistent with the  
329 information footprints shown in Fig. 2. By contrast, the former hypothesis predicts that  
330 the CRP concentration  $C_c$  in our model for  $\tau$  should vary from experiment to experiment,  
331 while the latter predicts an experiment-dependent interaction energy  $\epsilon_i$ . To further test  
332 the validity of our approach, we fit a single model for  $\tau$  to all six of our data sets (see SI  
333 Appendix Sec. 2d). This multi-data-set model employed a single CRP energy matrix and  
334 a single RNAP energy matrix, but allowed for data-set-specific values for both  $\epsilon_i$  and  $C_c$ .

335

336 The experiment-specific values we inferred for  $\epsilon_i$  and  $C_c$  are shown in Figs. 4A and 4B.  
337 The  $\epsilon_i$  values determined for all six experiments are mutually consistent ( $P > 0.05$ ,  $\chi^2$   
338 test), whereas the six  $C_c$  values differ very significantly:  $C_c$  values for experiments full-0  
339 and full-150, which were performed in reduced concentrations of cAMP, were found to  
340 be much lower than the  $C_c$  values determined for the other four experiments. These  
341 results verify the well-known fact that cAMP affects the concentration of active CRP, not  
342 the strength of CRP's interaction with RNAP (9,13).

343

344 Repeating this inference using only one  $\epsilon_i$  value for all six experiments yielded the final  
345 inferred model of this paper (Fig. 4 and SI Appendix Fig. S5). The CRP-RNAP  
346 interaction energy inferred from all six data sets is  $\epsilon_i = -2.82 \pm 0.13$  kcal/mol. This differs  
347 from Kuhlman et al.'s measurement, but only by  $\sim 20\%$ . Also, the ratio of CRP  
348 concentrations inferred for full-150 and full-500 ( $C_c^{full-150}/C_c^{full-500} = 10^{-0.56 \pm 0.02}$ ) closely  
349 agrees with the ratio of cAMP (150  $\mu$ M/500  $\mu$ M =  $10^{-0.52}$ ) used in these two experiments.  
350 This agreement is consistent with Kuhlman et al.'s observation that, for TK310 cells, the  
351 concentration of active CRP is proportional to the exogenous concentration of cAMP  
352 (12). Indeed, even the expected trace amount of cAMP present in the full-0 induction  
353 media ( $\sim 50$  nM, a result of carry-over from the starter culture inoculum) is fully  
354 consistent with the much-reduced yet significantly nonzero value for  $C_c^{full-0}/C_c^{full-500} = 10^{-$   
355  $3.7 \pm 0.6$ .

356

357 We note, however, that there are also puzzling quantitative oddities in our results. Our  
358 CRP concentration ratios are in good agreement with expectations, but the absolute  
359 values we inferred for  $C_c$  (full-500:  $10^{-1.13 \pm 0.06}$ ; full-150:  $10^{-1.70 \pm 0.07}$ ) are about 20-fold  
360 lower than the corresponding values claimed by Kuhlman et al. (SI Appendix Sec. 2g).  
361 Also, in the six- $\epsilon_i$  fit, the average  $\epsilon_i$  is  $-2.34 \pm 0.09$  kcal/mol, substantially less than the  
362  $-2.82 \pm 0.13$  kcal/mol we inferred using a common  $\epsilon_i$  for all six experiments. Finally, all of  
363 our inferred CRP energy matrices are asymmetric and predict that CRP binds in the  
364 energetically unfavorable orientation at the wild-type CRP site. We think these strange  
365 results probably result from our models for  $\epsilon_c$ ,  $\epsilon_r$ , and  $\tau$  being too simplistic. The issue of  
366 model selection is complicated, however, and is beyond the scope of this paper.



367

## 368 **Model validation using predictive information**

369

370 Predictive information can be used to determine how well a proposed thermodynamic  
371 model integrates binding energies into a single transcriptional output. The predictive  
372 information  $I(\tau, \mu)$  quantifies how well the biophysical model in Eq. 2 accounts for our  
373 measurements  $\mu$ . However, we can also directly compute the predictive information  
374  $I(\epsilon_c, \epsilon_r; \mu)$  of the pair of binding energies,  $\epsilon_c$  and  $\epsilon_r$ , without any model for  $\tau$  (SI Appendix  
375 Sec. 2e). The mere fact that  $\tau$  is a function of  $\epsilon_c$  and  $\epsilon_r$  (in Eq. 2) means that  $I(\tau, \mu) \leq$   
376  $I(\epsilon_c, \epsilon_r; \mu)$ , i.e. the predictive information of  $\tau$  is bounded above by the predictive  
377 information of the pair of energies  $\epsilon_c$  and  $\epsilon_r$ . This is a direct consequence of the Data  
378 Processing Inequality, a basic result in information theory (22). In our case, equality  
379 between  $I(\tau, \mu)$  and  $I(\epsilon_c, \epsilon_r; \mu)$  can be achieved only if  $\tau$  preserves all of the  
380 transcriptionally relevant information encoded in the predicted values for  $\epsilon_c$  and  $\epsilon_r$ . We  
381 emphasize that there is no *a priori* guarantee that any quantity  $\tau$  can do this, let alone a  
382 quantity derived from a simple biophysical model.

383

384 Remarkably we find (on full-wt data) that  $I(\tau, \mu) = 0.732 \pm 0.007$  bits, which is identical  
385 within error bars to  $I(\epsilon_c, \epsilon_r; \mu) = 0.732 \pm 0.006$  bits. We believe this equality is an  
386 important validation of the specific thermodynamic formula (Eq. 2) used to represent  $\tau$ .  
387 This agreement also argues that our energy matrix models for  $\epsilon_c$  and  $\epsilon_r$  provide a valid  
388 representation of physical binding energy; if they did not, their use in the Boltzmann  
389 exponents in Eq. 2 would be unlikely to yield sensible results. We note that a simpler  
390 model in which transcription depends only linearly on  $\epsilon_c$  and  $\epsilon_r$  — which would be  
391 appropriate if CRP and RNAP bound DNA only as a complex — achieves only  $I(\epsilon_c + \epsilon_r; \mu)$   
392  $= 0.647 \pm 0.005$  bits when all parameters are fit *ab initio* to full-wt data. This is  
393 significantly less than  $I(\epsilon_c, \epsilon_r; \mu)$ . Our specific thermodynamic model for  $\tau$ , with its  
394 physically motivated functional form, therefore provides a marked improvement over the  
395 more naive linear model.

396

## 397 **Discussion**

398

399 The approach we present here can be applied to a wide variety of CRMs in a number of  
400 different organisms. No prior knowledge of a CRM's sequence architecture is needed.  
401 All that is required is that (i) the CRM of interest function on a reporter construct and (ii)  
402 a large library of reporter constructs be introduced into cells so that each cell receives a  
403 single mutant CRM. After the activity-based-partitioning and sequencing of mutant  
404 CRMs, information footprints can be used to identify all functionally relevant positions  
405 within the probed sequence. Footprints from experiments performed in growth  
406 conditions or genetic backgrounds that are known to affect expression can further help  
407 one identify which 'clumps' of informative positions correspond to discrete binding sites,  
408 as well as which of these sites are involved in transducing specific intracellular signals.  
409 Eq. 1 then allows one to infer mathematical models describing the sequence-dependent  
410 binding energy of each site's cognate protein. The same fitting procedure can also be  
411 used to build biophysical models of the *in vivo* interactions between multiple DNA-bound

412 proteins. Such inference requires no quantitative model of experimental noise, thus  
413 allowing experiments performed in very different ways to produce nearly identical  
414 results.

415  
416 GFP reporter plasmids have been constructed for almost all *E. coli* promoters (21).  
417 Starting from these plasmids and using our protocols, it should be possible to  
418 biophysically characterize the vast majority of promoters in *E. coli*. Similar experiments  
419 can likely be performed for most CRMs in yeast. While we demonstrated our technique  
420 on a CRM containing only two protein-binding sites, the great sensitivity of this  
421 approach should allow the simultaneous effects of many DNA-binding proteins  
422 (including nucleosomes) to be discerned from a single experiment. Ultimately, our  
423 method should be useful for characterizing the detailed functional architecture of CRMs  
424 active in a wide variety of culturable cells, including stem cells and cancer cell lines. The  
425 biophysical characterization of CRMs in living animals, however, will likely require  
426 significant changes to the experimental approach described here.

427  
428 The underlying principles of our approach should also be useful for studying many  
429 different systems in which sequence-dependent interactions (e.g. protein-DNA, protein-  
430 RNA, protein-protein or protein-ligand interactions) play a central role in establishing  
431 some activity of interest. Experiments like ours can be performed both *in vivo* and *in*  
432 *vitro*; one simply needs to partition a large number of mutant sequences according to  
433 each one's activity. Our analysis method should then allow quantitative models of  
434 sequence-dependent function to be rigorously inferred from the resulting sequence  
435 data, regardless how the activity-dependent partitioning of sequences is accomplished.

### 436 437 **Experimental Procedures**

438  
439 See SI Appendix Sec. 1 for an expanded explanation of our experimental procedures.

440  
441 **Strains:** *E. coli* strains MG1655 (wild-type) and TK310 ( $\Delta cyaA \Delta cpdA \Delta lacY$ ) were  
442 kindly provided by Thomas Kuhlman. Except where noted, TK310 cells were maintained  
443 in media supplemented with 500  $\mu$ M cAMP in order to prevent *crp*-mediated  
444 suppression of  $\Delta cyaA$  (32).

445  
446 **Library construction:** Mutant *lac* promoter libraries were synthesized by IDT using  
447 defined mixtures of nucleoside phosphoramidites. Plasmid libraries consisted of  $\sim 2 \times 10^6$   
448 independently cloned plasmids in which region [-75:-1] was exactly replaced without the  
449 introduction of artificial restriction sites (SI Appendix Sec. 1a)

450  
451 **Sorting:** Cells were grown in exponential phase for  $\geq 10$  generations, diluted into buffer,  
452 and stored on ice for 0-24 hr prior to sorting. A BD Biosciences FACSVantage SE with  
453 DiVa was then used to sort 100,000 cells into each batch based solely on GFP  
454 fluorescence. Plating revealed  $\sim 70,000$  viable cells per batch.

455  
456 **Amplicon generation and sequencing:** Minipreped plasmid from each FACS batch  
457 was used as template for amplicon-generating PCR. Two control sets of amplicons

458 were also generated from pUA66-lacZ plasmid. The 47 resulting amplicon libraries (45  
459 FACS batches + 2 wild-type controls) were collated and sequenced by Roche using the  
460 Genome Sequencer FLX platform. This yielded 448,416 sequences, 308,309 of which  
461 passed our quality filters. Unfortunately, an analysis of these sequences indicated a  
462 large post-sort reduction in sequence diversity (SI Appendix Sec. 1f). To guarantee that  
463 each sequence was independently sorted, we discarded all but one copy of each  
464 sequence in each batch, leaving a total of 220,591 sequences across our six  
465 experiments.

466  
467 **Sequence data and analysis results:** Our 454 sequence data is available on the NCBI  
468 website under accession number SRA012345. Additional information, including  
469 processed sequence reads and inferred model parameters, is available at  
470 <http://www.princeton.edu/~ccallan/sortseq09/>.

471

## 472 **Acknowledgements**

473

474 We are deeply grateful to Christina DeCoste, who assisted with all flow cytometry  
475 instrumentation. Amy Caudy, Thomas Kuhlman, and Paul Wiggins provided critical  
476 advice on various experimental issues. We further benefited from discussions with  
477 William Bialek, Thomas Gregor, Stanislav Shvartsman, Antoinette Sutto, Gasper Tkačik,  
478 and Michael Zhang. The work of CGC, JBK and AM was supported in part by National  
479 Science Foundation grant PHY-0650617. The work of JBK and EC was supported in  
480 part by National Institute of Health grants GM078591 and GM071508. The work of JBK  
481 was supported in part by the Simons Foundation. The work of CGC was supported in  
482 part by US Department of Energy grant DE-FG02-91ER40671. The authors declare no  
483 conflicts of interest.

484

## 485 **References**

486

487 1. Kinney JB, Tkačik G, Callan C (2007) Precise physical models of protein-DNA  
488 interaction from high-throughput data. *Proc Natl Acad Sci USA* 104:501-506.

489

490 2. Ren B, et al. (2000) Genome-wide location and function of DNA binding proteins.  
491 *Science* 290:2306-2309.

492

493 3. Johnson D, Mortazavi A, Myers R, Wold B (2007) Genome-wide mapping of *in vivo*  
494 protein-DNA interactions. *Science* 316:1497-1502.

495

496 4. Berg O, von Hippel P (1988) Selection of DNA binding sites by regulatory proteins. II.  
497 The binding specificity of cyclic AMP receptor protein to recognition sites. *J Mol Biol*  
498 200:709-723.

499

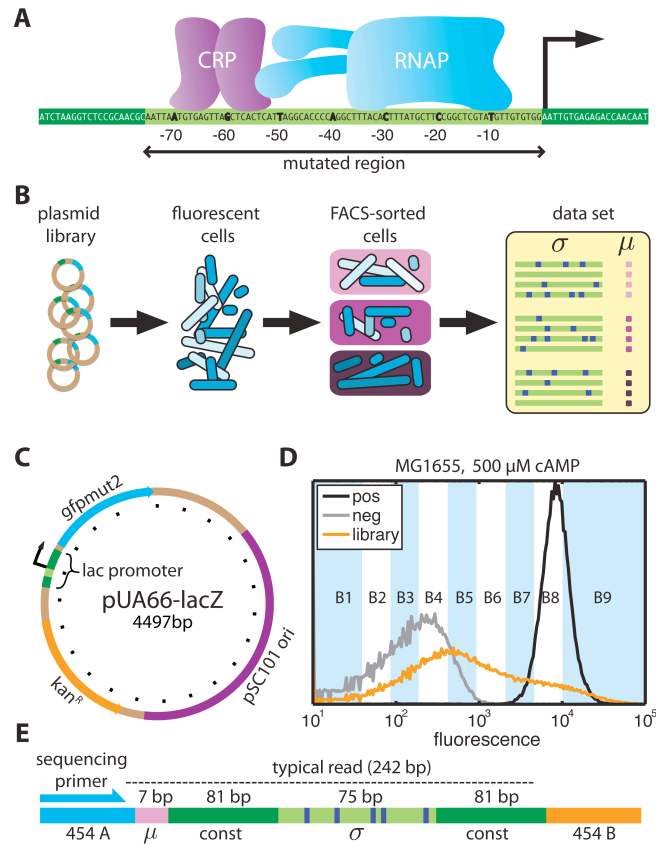
500 5. Tuerk C, Gold L (1990) Systematic evolution of ligands by exponential enrichment:  
501 RNA ligands to bacteriophage T4 DNA polymerase. *Science* 249:505-510.

502

- 503 6. Meng X, Brodsky MH, Wolfe SA (2005) A bacterial one-hybrid system for determining  
504 the DNA-binding specificity of transcription factors. *Nature Biotechnol* 23:988-994.  
505
- 506 7. Berger M, et al. (2006) Compact, universal DNA microarrays to comprehensively  
507 determine transcription-factor binding site specificities. *Nat Biotechnol* 24:1429-1435.  
508
- 509 8. Carey MF, Peterson CL, Smale ST (2009) *Transcriptional regulation in eukaryotes:  
510 concepts, strategies, and techniques* (Cold Spring Harbor, NY: Cold Spring Harbor  
511 Laboratory Press).  
512
- 513 9. Ptashne M, Gann A (2002) *Genes and signals* (Cold Spring Harbor, NY: Cold Spring  
514 Harbor Laboratory).  
515
- 516 10. Müller-Hill B (1996) *The lac operon: a short history of a genetic paradigm* (Berlin:  
517 Walter de Gruyter).  
518
- 519 11. Vilar JMG, Leibler S (2003) DNA looping and physical constraints on transcription  
520 regulation. *J Mol Biol* 331:981-989.  
521
- 522 12. Kuhlman T, Zhang Z, Saier MH, Hwa T (2007) Combinatorial transcriptional control  
523 of the lactose operon of Escherichia coli. *Proc Natl Acad Sci USA* 104:6043-6048.  
524
- 525 13. Busby S, Ebricht RH (1999) Transcription activation by catabolite activator protein  
526 (CAP). *J Mol Biol* 293:199-213.  
527
- 528 14. Bintu L, et al. (2005) Transcriptional regulation by the numbers: models. *Curr Opin  
529 Genet Dev* 15:116-124.  
530
- 531 15. Fakhouri WD, et al. (2010) Deciphering a transcriptional regulatory code: modeling  
532 short-range repression in the Drosophila embryo. *Mol Syst Biol* 6:341.  
533
- 534 16. Schneider T, Stormo G (1989) Excess information at bacteriophage T7 genomic  
535 promoters detected by a random cloning technique. *Nucleic Acids Res* 17:659-674.  
536
- 537 17. Herzenberg L, Sweet R, Herzenberg L (1976) Fluorescence-activated cell sorting.  
538 *Sci Am* 234:108-117.  
539
- 540 18. Margulies M, et al. (2005) Genome sequencing in microfabricated high-density  
541 picolitre reactors. *Nature* 437:376-380.  
542
- 543 19. Segal E, Raveh-Sadka T, Schroeder M, Unnerstall U, Gaul U (2008) Predicting  
544 expression patterns from regulatory sequence in Drosophila segmentation. *Nature*  
545 451:535-540.  
546
- 547 20. Gertz J, Siggia ED, Cohen BA (2009) Analysis of combinatorial cis-regulation in  
548 synthetic and genomic promoters. *Nature* 457:215-218.

549  
550 21. Zaslaver A, et al. (2006) A comprehensive library of fluorescent transcriptional  
551 reporters for Escherichia coli. *Nat Methods* 3:623-628.  
552  
553 22. Cover TM, Thomas JA (1991) *Elements of information theory* (New York, NY:  
554 Wiley).  
555  
556 23. Patwardhan RP, et al. (2009) High-resolution analysis of DNA regulatory elements  
557 by synthetic saturation mutagenesis. *Nature Biotechnol* 27:1173-1175.  
558  
559 24. Ross W, et al. (1993) A third recognition element in bacterial promoters: DNA  
560 binding by the alpha subunit of RNA polymerase. *Science* 262:1407-1413.  
561  
562 25. Reznikoff WS (1992) The lactose operon-controlling elements: a complex paradigm.  
563 *Mol Microbiol* 6:2419-2422.  
564  
565 26. Takeda Y, Sarai A, Rivera VM (1989) Analysis of the sequence-specific interactions  
566 between Cro repressor and operator DNA by systematic base substitution experiments.  
567 *Proc Natl Acad Sci USA* 86:439-443.  
568  
569 27. Sarai A, Takeda Y (1989) Lambda repressor recognizes the approximately 2-fold  
570 symmetric half-operator sequences asymmetrically. *Proc Natl Acad Sci USA* 86:6513-  
571 6517.  
572  
573 28. Man TK, Stormo GD (2001) Non-independence of Mnt repressor-operator  
574 interaction determined by a new quantitative multiple fluorescence relative affinity  
575 (QuMFRA) assay. *Nucleic Acids Res* 29:2471-2478.  
576  
577 29. Maerkl SJ, Quake SR (2007) A systems approach to measuring the binding energy  
578 landscapes of transcription factors. *Science* 315:233-237.  
579  
580 30. Harley CB, Reynolds RP (1987) Analysis of E. coli promoter sequences. *Nucl Acids*  
581 *Res* 15:2343-2361.  
582  
583 31. Gunasekera A, Ebright YW, Ebright RH (1992) DNA sequence determinants for  
584 binding of the Escherichia coli catabolite gene activator protein. *J Biol Chem* 267:14713-  
585 14720.  
586  
587 32. Harman JG, Dobrogosz WJ (1983) Mechanism of CRP-mediated cya suppression in  
588 Escherichia coli. *J Bacteriol* 153:191-199.  
589  
590

590 **Main text figures and tables**  
 591



592 **Figure 1.** Overview of the experiments A) Our experiments used *lac* promoters mutagenized in region [-  
 593 75:-1], which binds both CRP and RNAP. B) Plasmids containing mutant *lac* promoters driving GFP  
 594 expression were transformed into *E. coli*. Induced cells were then partitioned using FACS. Deep  
 595 sequencing of the mutant promoters in each FACS batch yielded a long list of sequences  $\sigma$  with  
 596 corresponding measurements  $\mu$ . C) Plasmid pUA66-lacZ (21), a very-low-copy-number plasmid on which  
 597 the wild-type *lac* promoter drives the expression of GFP; tick mark spacing is 200 bp. D) Fluorescence  
 598 distributions of MG1655 cells containing the full-wt plasmid library (orange), the pUA66-lacZ plasmid  
 599 (black), or a negative control plasmid pJK10 (SI Appendix Fig. S1) in which region [-75:-1] of the *lac*  
 600 promoter was deleted (gray). In the full-wt experiment, batches B1-B9 received cells from the indicated  
 601 fluorescence ranges, while batch B0 received cells randomly sampled from the initial library. E) Each  
 602 PCR amplicon contained a 7 bp DNA barcode indicating the batch  $\mu$  in which each sequence  $\sigma$  was  
 603 found. 454 pyrosequencing (18) yielded reads of about 242 bp covering the indicated regions.  
 604  
 605  
 606  
 607

607

**Table 1. Summary of our six experiments**

<b>Data Set</b>	<b>Mut. Region</b>	<b>Mut. Rate</b>	<b>Strain</b>	<b>cAMP (<math>\mu</math>M)</b>	<b>No. <math>\mu</math></b>	<b>No. Reads</b>
full-wt	[-75:-1]	12%	MG1655	500	10	51,835
crp-wt	[-74:-49]	24%	MG1655	500	10	46,986
rnap-wt	[-39:-4]	15%	MG1655	500	10	45,461
full-500	[-75:-1]	12%	TK310	500	5	23,431
full-150	[-75:-1]	12%	TK310	150	5	24,334
full-0	[-75:-1]	12%	TK310	0	5	28,544

608

609

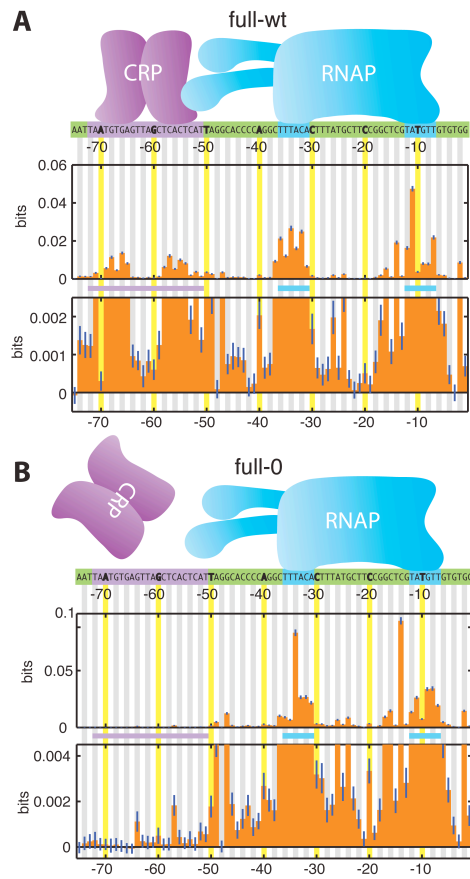
610

611

612

613

**Table 1.** Shown for each of our six experiments are the mutagenized region of the *lac* promoter, the per-position substitution rate, the *E. coli* strain used, the cAMP concentration used for induction, the number of batches into which cells were sorted, and the final number of filtered, non-redundant reads.

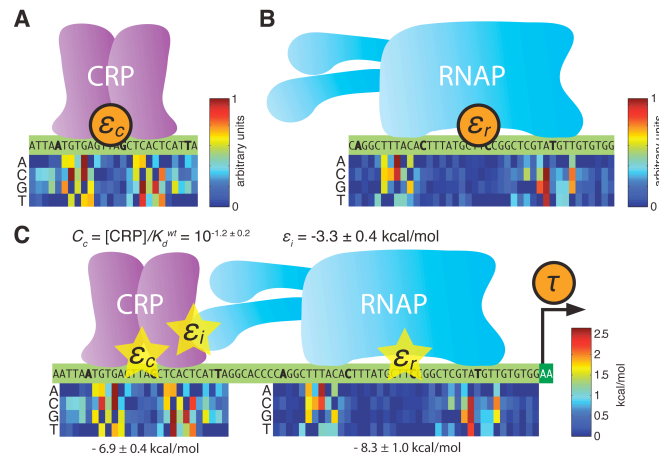


614  
 615  
 616  
 617  
 618  
 619  
 620  
 621

**Figure 2.** Information footprints. A) Footprint from full-wt data, aligned with known protein-DNA contact positions (highlighted). The lower plot is a 20X magnification of the upper plot. Blue lines indicate uncertainties due to finite sample effects (SI Appendix Sec. 2e). B) Footprint of the full-0 data set, in which intracellular CRP was inactive. SI Appendix Fig. S3 shows information footprints from all six experiments.



621



622

623

624

625

626

627

628

629

630

631

632

633

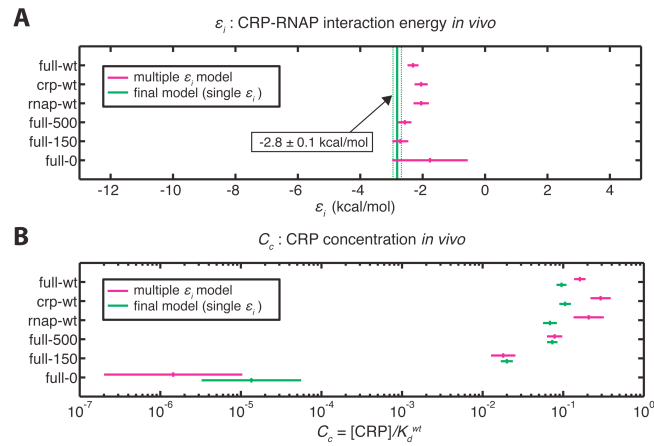
634

635

636

**Figure 3.** Models fit to full-wt data. A) The CRP energy matrix fit to [-74:-49] by maximizing  $l(\epsilon_c; \mu)$  on full-wt data. B) The RNAP energy matrix fit to [-41:-1] by maximizing  $l(\epsilon_r; \mu)$  on full-wt data. In panels A and B, each matrix column lists the energy contributions of the four possible bases at the aligned position within the site. Matrix elements range from 0 to 1 (in arbitrary units) with the lowest element in each column set to zero by convention. SI Appendix Fig. S4 shows the CRP and RNAP matrices derived from all six of our data sets. C) The thermodynamic model for  $\tau$  inferred using  $l(\tau; \mu)$  in Eq. 1. Optimal CRP and RNAP energy matrices are shown with elements expressed in kcal/mol (1 kcal/mol =  $1.62 k_b T$  at  $T = 310$  °K). It is useful to define each wild-type *lac* promoter site as having zero energy. We therefore add an energy shift, shown below each matrix, when computing  $\epsilon_c$  and  $\epsilon_r$ . Doing this means that  $C_c$  is the intracellular CRP concentration in units of the dissociation constant of the wild-type (zero energy) site. Values quoted for  $\epsilon_i$  and  $C_c$  are mean  $\pm$  RMSD values determined from the parameter ensembles sampled using parallel tempering Monte Carlo.

636



637  
638  
639  
640  
641  
642  
643  
644

**Figure 4.** Parameters fit to all six data sets. A) CRP-RNAP interaction energies  $\epsilon_i$  (mean  $\pm$  RMSD) inferred by fitting  $\tau$  to all six data sets, using either data-set-specific values for  $\epsilon_i$  (magenta) or a single  $\epsilon_i$  for all six data sets (green). B) CRP concentrations  $C_c$  inferred for these same multi-data-set models. SI Appendix Fig. S5 shows full ensemble distributions for  $\epsilon_i$  and the six  $C_c$  parameters of the final model, together with mean and RMSD values for all the CRP and RNAP matrix elements.