

# 1 Compatibility logic of human enhancer and promoter sequences

2  
3 Drew T. Bergman<sup>1,2,\*</sup>, Thouis R. Jones<sup>1,\*</sup>, Vincent Liu<sup>3</sup>, Layla Siraj<sup>1,5</sup>, Helen Y. Kang<sup>3,4</sup>, Joseph  
4 Nasser<sup>1</sup>, Michael Kane<sup>1</sup>, Tung H. Nguyen<sup>1</sup>, Sharon R. Grossman<sup>1</sup>, Charles P. Fulco<sup>1,8</sup>, Eric S.  
5 Lander<sup>1,6,7,9</sup>, Jesse M. Engreitz<sup>1,3,4</sup>

6  
7 1. Broad Institute of MIT and Harvard, Cambridge, MA, USA.

8 2. Geisel School of Medicine at Dartmouth, Hanover, NH, USA.

9 3. Department of Genetics, Stanford University School of Medicine, Stanford, CA, USA.

10 4. BASE Initiative, Betty Irene Moore Children's Heart Center, Lucile Packard Children's Hospital, Stanford University  
11 School of Medicine, Stanford, CA, USA.

12 5. Biophysics Graduate Program, Harvard University, Cambridge, MA, USA.

13 6. Department of Biology, MIT, Cambridge, MA, USA.

14 7. Department of Systems Biology, Harvard Medical School, Boston, MA, USA.

15 8. Present address: Bristol Myers Squibb, Cambridge, MA, USA.

16 9. Currently on leave from the Broad Institute, MIT, and Harvard.

17 \* These authors contributed equally.

18

19

## 20 Abstract

21

22 Gene regulation in the human genome is controlled by distal enhancers that activate specific  
23 nearby promoters. One model for the specificity of enhancer-promoter regulation is that different  
24 promoters might have sequence-encoded preferences for distinct classes of enhancers, for  
25 example mediated by interacting sets of transcription factors or cofactors. This “biochemical  
26 compatibility” model has been supported by observations at individual human promoters and by  
27 genome-wide measurements in *Drosophila*. However, the degree to which human enhancers and  
28 promoters are intrinsically compatible or specific has not been systematically measured, and how  
29 their activities combine to control RNA expression remains unclear. To address these questions,  
30 we designed a high-throughput reporter assay called enhancer x promoter (ExP) STARR-seq and  
31 applied it to examine the combinatorial compatibilities of 1,000 enhancer and 1,000 promoter  
32 sequences in human K562 cells. We identify a simple logic for enhancer-promoter compatibility –  
33 virtually all enhancers activated all promoters by similar amounts, and intrinsic enhancer and  
34 promoter activities combine multiplicatively to determine RNA output ( $R^2=0.82$ ). In addition, two  
35 classes of enhancers and promoters showed subtle preferential effects. Promoters of  
36 housekeeping genes contained built-in activating sequences, corresponding to motifs for factors  
37 such as GABPA and YY1, that correlated with both stronger autonomous promoter activity and  
38 enhancer activity, and weaker responsiveness to distal enhancers. Promoters of context-specific  
39 genes lacked these motifs and showed stronger responsiveness to enhancers. Together, this  
40 systematic assessment of enhancer-promoter compatibility suggests a multiplicative model tuned  
41 by enhancer and promoter class to control gene transcription in the human genome.

## 42 Introduction

43  
44 The extent to which distal enhancers might activate specific types of promoters has been an  
45 outstanding question in human gene regulation. Since their initial discovery, enhancers have  
46 been defined in part based on their ability to activate multiple non-cognate promoter  
47 sequences<sup>1,2</sup>. High-throughput reporter assays have now confirmed that many enhancer  
48 sequences derived from the human genome have the capability to activate various human, viral,  
49 and synthetic promoters<sup>3-9</sup>.

50  
51 Yet, other observations have suggested that enhancers and promoters have some degree of  
52 intrinsic specificity. Early studies identified individual examples where particular enhancers or  
53 cofactors showed stronger activation with certain core promoters<sup>10-15</sup>. More recently, in  
54 *Drosophila*, studies using high-throughput reporter assays revealed that developmental and  
55 housekeeping gene promoters show >10-fold preferences for different classes of genomic  
56 enhancers<sup>16</sup>, have differing levels of sequence-encoded responsiveness to enhancer  
57 activation<sup>17</sup>, and respond differently to recruitment of various transcriptional cofactors<sup>18</sup>.  
58 Together, these studies have suggested a ‘biochemical compatibility’ model where different  
59 enhancers might have an intrinsic preference for activating different promoter sequences based  
60 on the transcription factors and cofactors they can recruit<sup>19,20</sup>.

61  
62 Despite these advances, the biochemical compatibility model has not been systematically tested  
63 for human enhancers and promoters. As such, it remains unclear whether compatibility classes  
64 of enhancers and promoters exist in the human genome, and, if so, how their activities combine  
65 and how such specificity is encoded.

## 66 High-throughput measurements of enhancer-promoter compatibility

67  
68 To investigate these questions, we developed an assay called enhancer x promoter (ExP)  
69 STARR-seq to test the ability of ~1,000 candidate enhancers to activate ~1,000 promoters. In  
70 this assay, we synthesize pools of enhancer and promoter sequences (here, 264-bp) and clone  
71 them in all pairwise combinations located ~340-bp apart in the revised human STARR-seq  
72 plasmid-based reporter vector (**Fig. 1a/S1a**)<sup>8</sup>. In STARR-seq assays, the enhancer sequence is  
73 transcribed and quantified using targeted RNA-seq to determine the level of expression of each  
74 plasmid<sup>4</sup>. For ExP STARR-seq, we introduce a unique 16-bp “plasmid barcode” adjacent to the  
75 enhancer sequence that allows us to determine which reporter transcripts are produced from  
76 which enhancer-promoter pairs. We transiently transfect this pool of plasmids into cells,  
77 measure the level of reporter transcripts produced, and calculate “STARR-seq expression” as  
78 the amount of RNA normalized to DNA input for each plasmid. This approach allows us to  
79 quantitatively measure the expression of hundreds of thousands of combinations of enhancer  
80 and promoter sequences, estimate the activities of individual enhancers and promoters, and test  
81 their compatibilities (see Methods).

82  
83 Hereafter, for clarity, we use the terms “enhancer sequences” and “promoter sequences” to  
84 refer to sequences cloned into the enhancer and promoter positions in the ExP STARR-seq  
85 assay, and “genomic enhancers” and “genomic promoters” to refer to the corresponding  
86 elements in the genome.

87  
88  
89 We applied ExP STARR-seq to examine the combinatorial activities of 1,000 enhancer and  
90 1,000 promoter sequences (**Table S1, Table S2**) in K562 erythroleukemia cells, which have  
91 been deeply profiled by the ENCODE Project<sup>21</sup> and where we have previously collected data

92 about which genomic enhancers regulate which genomic promoters using CRISPR interference  
93 (CRISPRi) screens<sup>22</sup>. Here, we selected promoter sequences to include (i) 65 genes studied in  
94 prior CRISPR screens; (ii) 735 additional genes sampled from across the genome to span a  
95 range of transcriptional activity (based on precision run-on sequencing (PRO-seq) data in K562  
96 cells); and (iii) 200 control sequences including random genomic control sequences that are not  
97 accessible by ATAC-seq, and dinucleotide shuffled sequences (**Fig. S1a**, see Methods). The  
98 promoter sequences were chosen to include approximately 20-bp downstream of the genomic  
99 transcription start site (as observed in capped analysis of gene expression (CAGE) data), and  
100 ~242-bp upstream (264 bp total, see Methods). In the enhancer position of ExP STARR-seq, we  
101 included (i) 131 accessible genomic elements we previously tested by CRISPRi; (ii) 669 other  
102 accessible genomic elements selected to span a range of quantitative H3K27ac and DNase-seq  
103 signals (centered on the summit of the DNase-seq peak); and (iii) 200 controls including random  
104 genomic control sequences and dinucleotide shuffled sequences (**Fig. S1a**, See Methods).

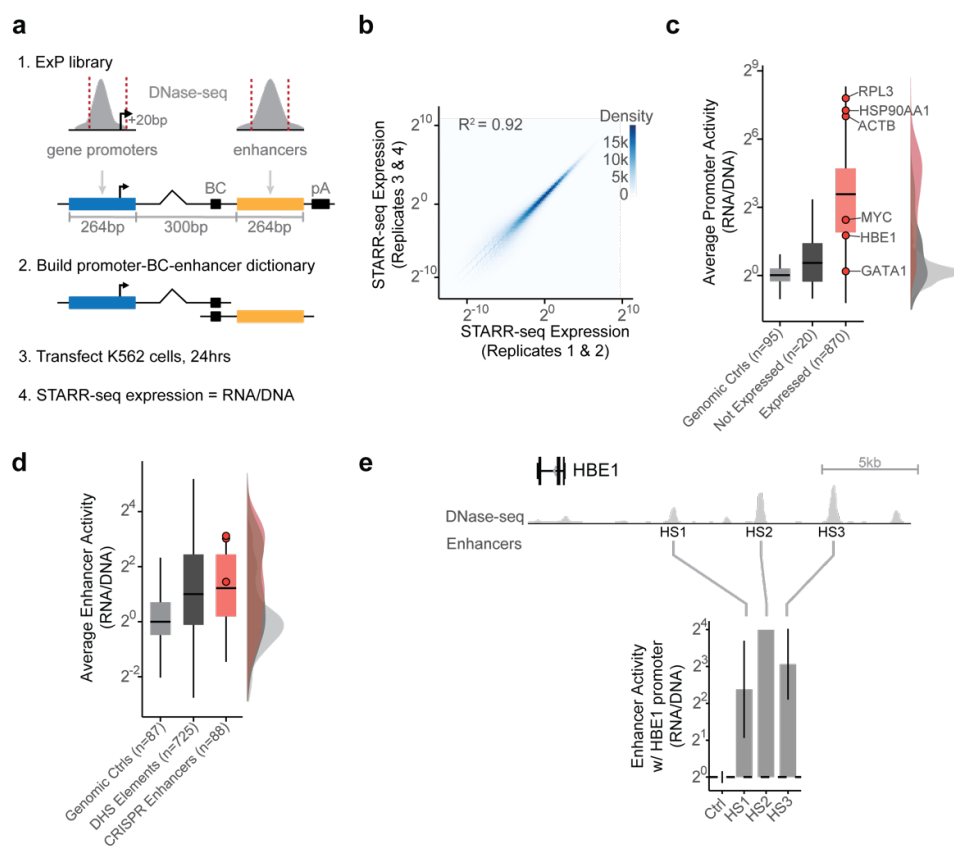
105  
106 We cloned these 1,000 enhancer and 1,000 promoter sequences in all pairwise combinations,  
107 transfected the plasmid pool into K562 cells in 4 biological replicates of 50 million cells each,  
108 and sequenced each STARR-seq RNA and input DNA library to a depth of at least 2.6 billion  
109 and 470 million reads, respectively. We focused our analysis on the 604,268 enhancer-promoter  
110 pairs where we obtained good coverage (see Methods). STARR-seq expression (RNA/DNA)  
111 varied over six orders of magnitude, and was highly reproducible, when comparing expression  
112 for individual plasmid barcodes between biological replicates ( $R^2 = 0.92$ , **Fig. 1b**), when  
113 comparing expression for an enhancer-promoter pair averaged across plasmid barcodes  
114 between biological replicates ( $R^2 = 0.92$ ), and when comparing expression for different plasmid  
115 barcodes for a given enhancer-promoter pair ( $R^2 = 0.62$ , **Fig. S1c-d**, see Methods).

116  
117 As expected, promoter sequences showed a very large (>1,500-fold) dynamic range of STARR-  
118 seq expression when paired with random genomic sequences in the enhancer position  
119 (“average promoter activity”). The strongest promoters in the dataset corresponded to  
120 housekeeping genes such as *RPL3*, *HSP90AA1*, and *ACTB*, and the weakest promoters  
121 included shuffled control sequences and non-expressed genes in K562 cells (**Fig. 1c**).  
122 Enhancer sequences also showed a wide (682-fold) range of STARR-seq expression in the  
123 dataset when averaged across promoters (“average enhancer activity”), and were on average 2-  
124 fold more active than random genomic control sequences (**Fig. 1d**). Enhancer and promoter  
125 activity from ExP STARR-seq were correlated with biochemical features of activity at the  
126 corresponding genomic elements, including with levels of chromatin accessibility, H3K27  
127 acetylation, and nascent gene and eRNA transcription (**Fig. S1e**).

128  
129 We also found that sequences derived from known genomic enhancers activated their cognate  
130 promoters in the ExP STARR-seq assay. For example, we included 3 enhancers in the beta-like  
131 globin locus control region (HS1-HS3) that are known to coordinate expression of hemoglobin  
132 subunits during erythrocyte development<sup>23,24</sup> and where CRISPRi perturbations in K562 cells  
133 reduce the expression of hemoglobin subunit epsilon 1 (*HBE1*) by 10-86%<sup>25,26</sup>. In ExP STARR-  
134 seq, each of these enhancers activated the *HBE1* promoter (by 5.21-15.9-fold versus random  
135 genomic controls, **Fig. 1e**). Similarly, an enhancer that we previously showed to regulate  
136 *GATA1* and *HDAC6* in the genome<sup>27</sup> led to 6.76 and 6.87-fold activation of the *GATA1* and  
137 *HDAC6* promoters in ExP STARR-seq, respectively (**Fig. S1f**).

138  
139 Taken together, these results show that ExP-STARR-seq produces quantitative and  
140 reproducible measurements of enhancer and promoter sequence activity over a large dynamic  
141 range.

142



143

144

145

### Fig. 1. Enhancer x Promoter STARR-seq

146 **a.** ExP STARR-seq method for measuring the activities of enhancer and promoter sequences and testing  
 147 their compatibilities. 264-bp sequences are selected and cloned in all pairwise combinations into the  
 148 promoter and enhancer positions of a plasmid vector, together with a plasmid barcode (BC). We build a  
 149 dictionary linking promoter-BC-enhancer triplets via sequencing (see **Fig. S1a**). We then transfect the  
 150 ExP STARR-seq plasmid pool into cells, where the promoter sequence on a given plasmid initiates  
 151 transcription of a polyadenylated RNA containing the plasmid barcode and enhancer. We sequence these  
 152 RNAs and calculate STARR-seq expression as the frequency of RNAs observed for each plasmid  
 153 normalized by the frequency of that plasmid in the input DNA plasmid pool.

154 **b.** Correlation of ExP STARR-seq expression between biological replicate experiments, calculated for  
 155 individual enhancer-promoter pairs with unique plasmid barcodes. Axes represent the average STARR-  
 156 seq expression (RNA/DNA) of two biological replicates. Density: number of enhancer-promoter plasmids.  
 157 **c.** Average promoter activity (STARR-seq expression when paired with random genomic controls in the  
 158 enhancer position) of promoter sequences derived from random genomic controls (set at 0), genes not  
 159 expressed in K562s, and all other gene promoters. Box is median and interquartile range, whiskers are  
 160 +/- 1.5 x IQR.

161 **d.** Average enhancer activity (STARR-seq expression of plasmids containing a given enhancer averaged  
 162 across all promoters) of enhancer sequences derived from random genomic controls, accessible  
 163 elements, and genomic enhancers validated in CRISPR experiments. Box and whiskers as in (c). Red  
 164 dots represent three enhancers near *HBE1* (see panel e).

165 **e.** Sequences derived from three genomic enhancers that regulate *HBE1* in the genome (HS1-HS3)  
 166 activate the *HBE1* promoter in ExP STARR-seq. Ctrl: Average of 44 random genomic control sequences  
 167 in the enhancer position that passed thresholds (see Methods). Error bars: 95% CI across plasmid  
 168 barcodes, n=110 (ctrl), 2 (HS1), 1 (HS2), 5 (HS3).

169

## 170 **Enhancer and promoter sequences are broadly compatible**

171  
172 We used this ExP STARR-seq dataset to test whether specific enhancers activate specific  
173 promoters. Surprisingly, virtually all active enhancer sequences activated all promoter  
174 sequences by similar amounts. For example, for a small subset of 5 enhancers and 5  
175 promoters, each with good coverage in the assay (median = 27 plasmid barcodes per pair),  
176 while the promoters spanned a 5.62-fold range of activities, the enhancers activated each  
177 promoter similarly (**Fig. 2a-b**). More generally, enhancers activated most promoters by similar  
178 amounts, with an average Spearman correlation across all pairs of promoters = 0.81 (**Fig. 2c,e,**  
179 **S2a**), and pairs of enhancers showed similar proportional activation of promoters, with an  
180 average Spearman = 0.72 (**Fig. 2d,f, S2b**). These observations indicate that, in this STARR-seq  
181 assay, there is broad compatibility between individual enhancer and promoter sequences — a  
182 striking difference from previous observations in *Drosophila*<sup>16,17</sup>.

183

184

## 185 **Enhancer and promoter activities combine approximately multiplicatively**

186  
187 This pattern of effects — where enhancers showed similar fold-activation across many  
188 promoters, and promoters showed similar levels of activation by many enhancers — suggested  
189 that intrinsic enhancer and promoter activities combine multiplicatively to produce the RNA  
190 output in STARR-seq. To quantify this, we correlated expression in the STARR-seq assay with  
191 intrinsic enhancer activity, intrinsic promoter activity, and the multiplicative product of intrinsic  
192 enhancer and promoter activities.

193

194 To do so, we fit the following Poisson count model:

195

$$196 \quad RNA \sim \text{Poisson}(k \times DNA \times P \times E),$$

197

198

199 where *RNA* is RNA reads counts per plasmid, *DNA* is DNA read counts per plasmid, *P* is the  
200 intrinsic promoter activity, *E* is intrinsic enhancer activity, and *k* is a free intercept term used to  
201 scale the activities of promoters, enhancers, and their pairings relative to the average of random  
202 genomic control sequences (see Methods). This multiplicative model assumes that there is no  
203 sequence or biochemical specificity between individual pairs of enhancers and promoters, and  
204 that differences in expression are solely due to differences in intrinsic enhancer and promoter  
205 activities. Hereafter, we define “intrinsic enhancer activity” and “intrinsic promoter activity” as the  
206 fits from this model, which yield similar estimates to the “average activities” calculated above  
207 (**Fig. S2c,d**) but better account for missing data and counting noise (see Methods). These  
208 estimates of activity were reproducible across replicate experiments and when comparing non-  
209 overlapping plasmid barcodes (**Fig. S2e,f**).

210

211 Intrinsic promoter activity alone explained 49% of the variance in STARR-seq expression across  
212 all enhancer-promoter pairs (correlation with log<sub>2</sub> STARR-seq expression in pairs with at least 2  
213 plasmid barcodes, **Fig. 2g**), and intrinsic enhancer activity alone explained 28% of the variance  
214 (**Fig. 2h**). The multiplicative combination of intrinsic promoter and enhancer activities explained  
215 82% of the total variance (**Fig. 2i-k**).

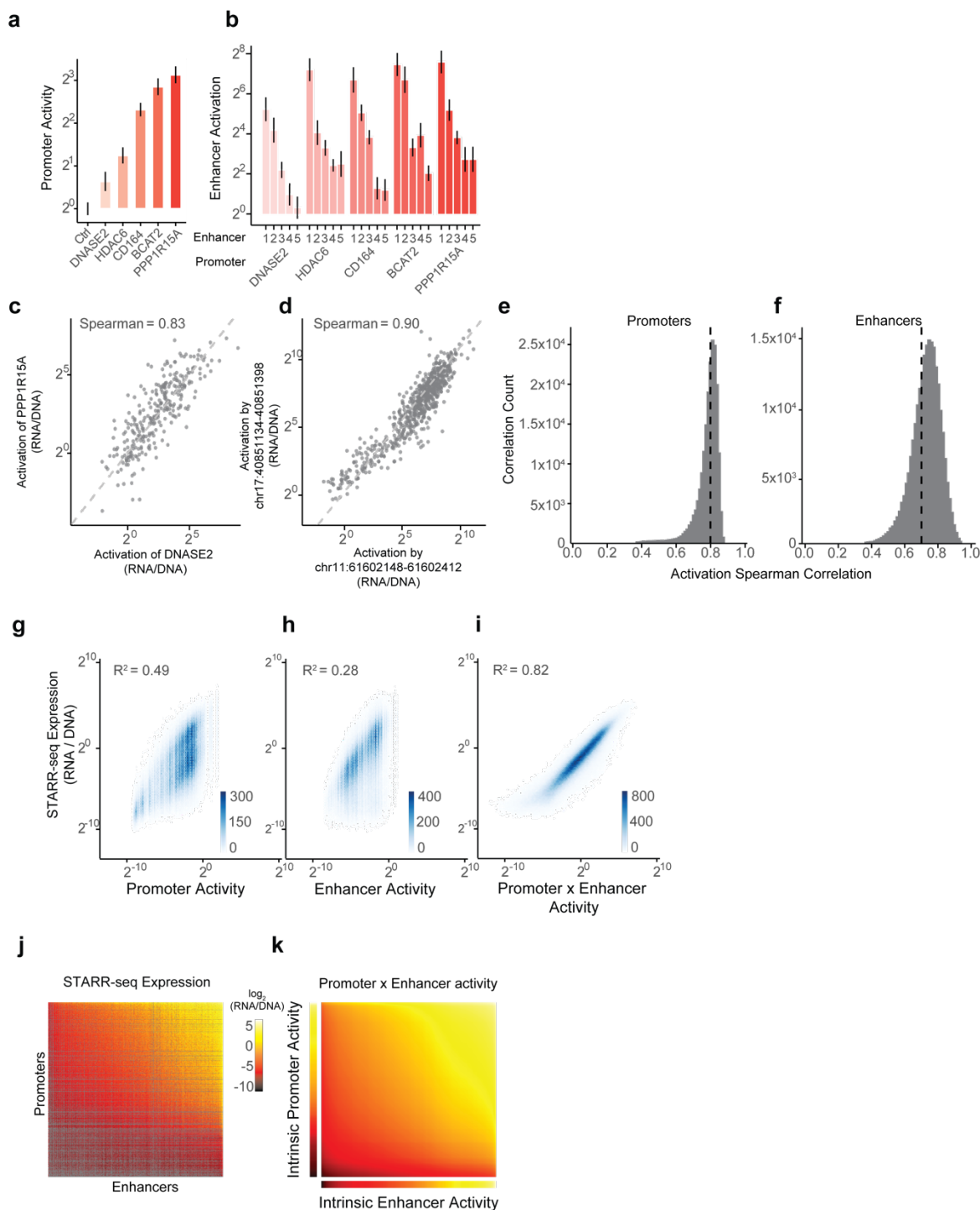
216

217 To confirm that this multiplicative relationship was not due to the specific design of our ExP  
218 STARR-seq assay, we cloned 7 enhancers from the *MYC* locus (1.0-2.2 kb) and 5 promoter  
219 sequences (138-908 bp, including the promoters of *MYC* and other nearby genes) in all

220 combinations into a different reporter plasmid in which the enhancer is located 1 kb upstream of  
221 the promoter, and measured the expression of these constructs using a luciferase reporter  
222 assay (**Fig. S2g, Table S3**). Again, despite a range of intrinsic promoter activities (**Fig. S2h**), all  
223 enhancer sequences activated all promoter sequences by a similar fold-change, and a  
224 multiplicative function of enhancer and promoter activities explained 78% of the total variance in  
225 the measurements (**Fig. S2i**).

226  
227 Thus, RNA expression in these reporter assays represents, to a first approximation, the  
228 multiplicative product of intrinsic enhancer activity and intrinsic promoter activity.





229  
230  
231  
232  
233  
234  
235  
236  
237

**Fig. 2. Enhancer and promoter activities combine multiplicatively**

**a.** Intrinsic promoter activity (expression versus random genomic controls in enhancer position) of five selected promoters. Error bars: 95% CI across plasmid barcodes (n=54-79).

**b.** Activation (expression versus random genomic controls in enhancer position) of 5 selected promoters by 5 selected enhancers (1 = chr11:61602148-61602412, 2 = chr19:49467061-49467325, 3 = chrX:48641342-48641606, 4 = chr19:12893216-12893480, 5 = chr17:40851134-40851398). Error bars: 95% CI across plasmid barcodes (n=12-56).

238 **c.** Correlation of enhancer activation for PPP1R15A and DNASE2 promoters. Each point is a shared  
239 enhancer sequence.  
240 **d.** Correlation of enhancer activation by chr17:40851134-40851398 and chr11:61602148-61602412  
241 enhancers. Each point is a shared promoter sequence.  
242 **e.** Distribution of pairwise correlations of enhancer activation between promoter sequences, as in **(c)**.  
243 black dotted line = mean Spearman correlation.  
244 **f.** Distribution of pairwise correlations of promoter activation between enhancer sequences, as in  
245 **(d)**. Black dotted line = mean Spearman correlation.  
246 **g-i.** Correlation of ExP STARR-seq expression with intrinsic promoter activity (**g**), intrinsic enhancer  
247 activity (**h**), and the product of intrinsic promoter and enhancer activities (**i**). Density color scale: number  
248 enhancer-promoter pairs.  
249 **j.** Heatmap of ExP STARR-seq expression across all pairs of promoter (vertical) and enhancer sequences  
250 (horizontal). Axes are sorted by intrinsic promoter and enhancer activities. Grey: missing data.  
251 **k.** Heatmap representing the multiplication of intrinsic promoter activity (vertical) with intrinsic enhancer  
252 activity (horizontal) from the Poisson model.  
253  
254  
255

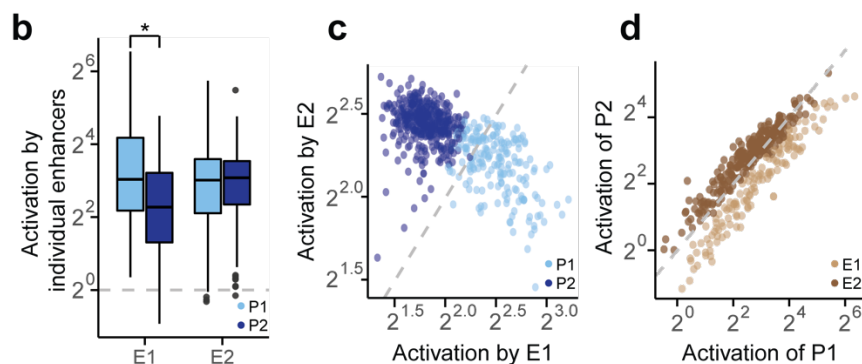
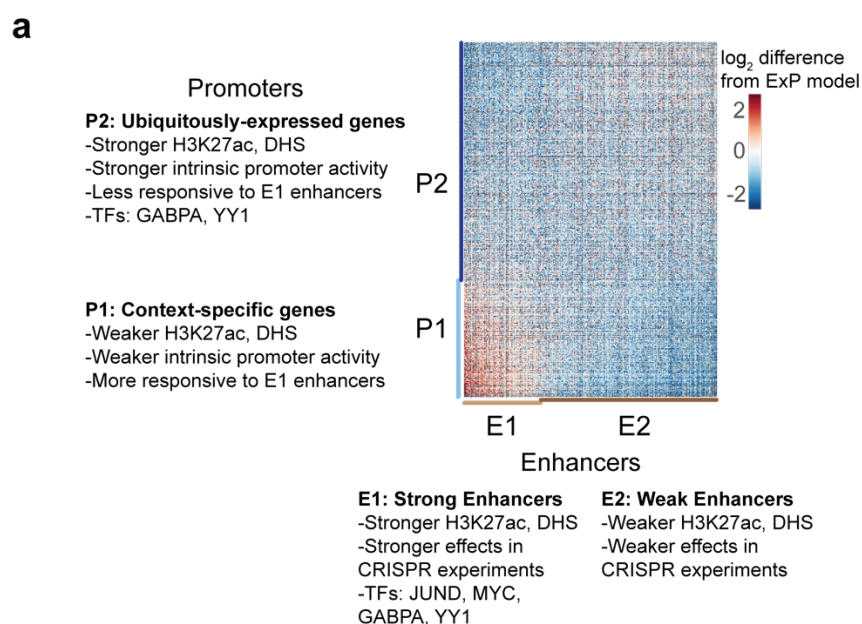
## 256 **Two functional classes of enhancer and promoter sequences**

257  
258 Although we did not observe a strong degree of specificity among enhancer and promoter  
259 sequences, we asked whether there might exist classes with more subtle, quantitative  
260 preferences. To do so, we calculated, for each enhancer-promoter pair, its deviation from the  
261 multiplicative enhancer x promoter model (observed STARR-seq expression versus the product  
262 of intrinsic enhancer activity and intrinsic promoter activity, see Methods).  
263

264 We identified two clusters of enhancer sequences (E1 and E2,  $n=126$  and  $290$  respectively) that  
265 showed differential effects with respect to two sets of promoter sequences (P1 and P2,  $n=192$   
266 and  $391$  respectively) (**Fig. 3a**). In particular, E1 enhancer sequences activated P1 promoters  
267 more strongly than P2 promoters (by 1.93-fold,  $P = 4.19e-08$ ,  $t$ -test), whereas E2 enhancer  
268 sequences activated promoters in both clusters approximately equally (1.05-fold stronger for P2  
269 versus P1,  $P = 0.424$ ,  $t$ -test; **Fig. 3b**). These sets of enhancers and promoters appeared to  
270 represent extremes of a graded scale: promoter responsiveness to E1 vs E2 enhancer  
271 sequences varied over a  $\sim 3$ -fold range (**Fig. 3c**, **Fig. S3d**, **Fig. S4b**), and enhancer activation of  
272 P1 vs P2 promoters varied over a  $\sim 2$ -fold range (**Fig. 3d**, **Fig. S3e**, **Fig. S4a**). Cluster  
273 assignments were highly stable to down-sampling of promoter and enhancer sequences (**Fig.**  
274 **S3g**, see Methods). Additional clusters (P0 and E0) contained sequences with very weak  
275 activity and/or missing data, and were excluded from further analysis (**Fig. S3a-c**).  
276

277 Together, these observations identify 2 classes of enhancer sequences and 2 classes of  
278 promoter sequences with subtle quantitative differences in compatibility. Accordingly, we next  
279 sought to characterize these classes of enhancer and promoter sequences and understand how  
280 such preferential effects might be encoded.  
281  
282  
283  
284





285  
286  
287  
288  
289  
290  
291  
292  
293  
294  
295  
296  
297  
298  
299  
300  
301  
302  
303  
304  
305

**Fig. 3. Compatibility classes of enhancers and promoters.**

**a.** Heatmap of deviations in enhancer-promoter STARR-seq expression from a multiplicative enhancer-promoter model (color scale: fold-difference between observed expression versus expression predicted by multiplicative model; gray: missing data). Vertical axis: promoter sequences grouped by class and sorted by responsiveness to E1 vs. E2 (see **b**); horizontal axis: enhancer sequences grouped by class and sorted by activation of P1 vs. P2 (see **c**).

**b.** Activation of P1 vs P2 promoters by E1 and E2 enhancer sequences (equivalently: Responsiveness to E1 vs E2 enhancer sequences). Boxes are median and interquartile range, whiskers are  $\pm 1.5 \times \text{IQR}$ . \* $P$ -value =  $4.2 \times 10^{-8}$ , two-sample  $t$ -test.

**c.** For each promoter, the average activation by (responsiveness to) E1 enhancer sequences ( $x$ -axis) versus the average activation by E2 enhancer sequences ( $y$ -axis). P1 promoters (light blue) are activated more strongly by E1 versus E2 enhancers.

**d.** For each enhancer, the average fold-activation when paired with P1 promoters ( $x$ -axis) versus P2 promoters ( $y$ -axis). E1 enhancers (light brown) more strongly activate P1 promoters.

## 306 **Classes of enhancer sequences correspond to strong and weak genomic enhancers**

307

308 To characterize the two classes of ExP STARR-seq enhancer sequences, we compared the  
309 classes with respect to biochemical features of their corresponding elements in the genome,  
310 sequence motifs, effects in CRISPR experiments, and other features.

311

312 E1 and E2 classes showed biochemical features of strong and weak genomic enhancers,  
313 respectively. The features most strongly associated with E1 versus E2 sequences in the  
314 genome included H3K27ac, DNase I hypersensitivity, AP-1 factor binding (JUN, ATF3), and  
315 other known activating transcription factors (**Fig. 3a**, **Fig. S5a-b**, **Table S4**). E2 sequences in  
316 the genome were also DNase accessible and sometimes bound these factors, but to a  
317 significantly lesser degree (**Fig. S7**). Consistent with these observations, E1 sequences had  
318 stronger effects on gene expression in CRISPR perturbation experiments, even when  
319 controlling for 3D contact with the target gene (**Fig. S5c**). While E1 sequences were more likely  
320 to be predicted to be enhancers in K562 cells (94% of E1 predicted to regulate a gene by the  
321 Activity-by-Contact (ABC) model, versus 49% of E2), both classes contained a large fraction of  
322 sequences predicted to be an enhancer in another cell type (90% of E1 and 70% of E2),  
323 suggesting that some E2 genomic elements may act as strong enhancers in other cell types.

324

325 These observations suggest that the differences in how these classes of enhancer sequences  
326 activate different promoters in ExP-STARR-seq could be related to their ability to recruit  
327 activating transcription factors (see below). We note that, despite these clear differences in  
328 genomic activity, the two classes of enhancer sequences showed, on average, similar levels of  
329 activity in the ExP-STARR-seq assay (**Fig. S3b**). This may reflect previous observations that the  
330 episomal STARR-seq assay often detects activity for sequences that do not appear to be active  
331 in their endogenous chromosomal context<sup>8,28</sup>.

332

333

## 334 **Classes of promoter sequences correspond to constitutive versus enhancer-responsive** 335 **genes**

336

337 The two classes of promoter sequences also showed striking differences in their functional  
338 annotations, intrinsic promoter activity, and responsiveness to enhancers in the genome.

339

340 We found that P2 promoter sequences were primarily derived from ubiquitously expressed  
341 genes (often called “housekeeping” genes), whereas P1 promoters corresponded to cell-type-  
342 or context-specific genes. For example, P2 promoters included beta actin (*ACTB*), all 37 tested  
343 ribosomal subunits (*e.g.*, *RPL13*, *RPS11*), components of the electron transport chain (*e.g.*,  
344 *NDUFA2*, *ATP5B*), and others (**Table S1**). In contrast, P1 promoters included erythroid-specific  
345 genes (*e.g.*, 3 hemoglobin genes, ferritin light chain (*FTL*)), context-specific transcription factors  
346 (*e.g.*, *KLF1*, *JUNB*, *REL*), and genes that are expressed in many cell types but at different  
347 levels, such as *MYC*. P1 and P2 promoters were associated with developmental and  
348 housekeeping gene ontology terms, respectively (**Fig. 4a**).

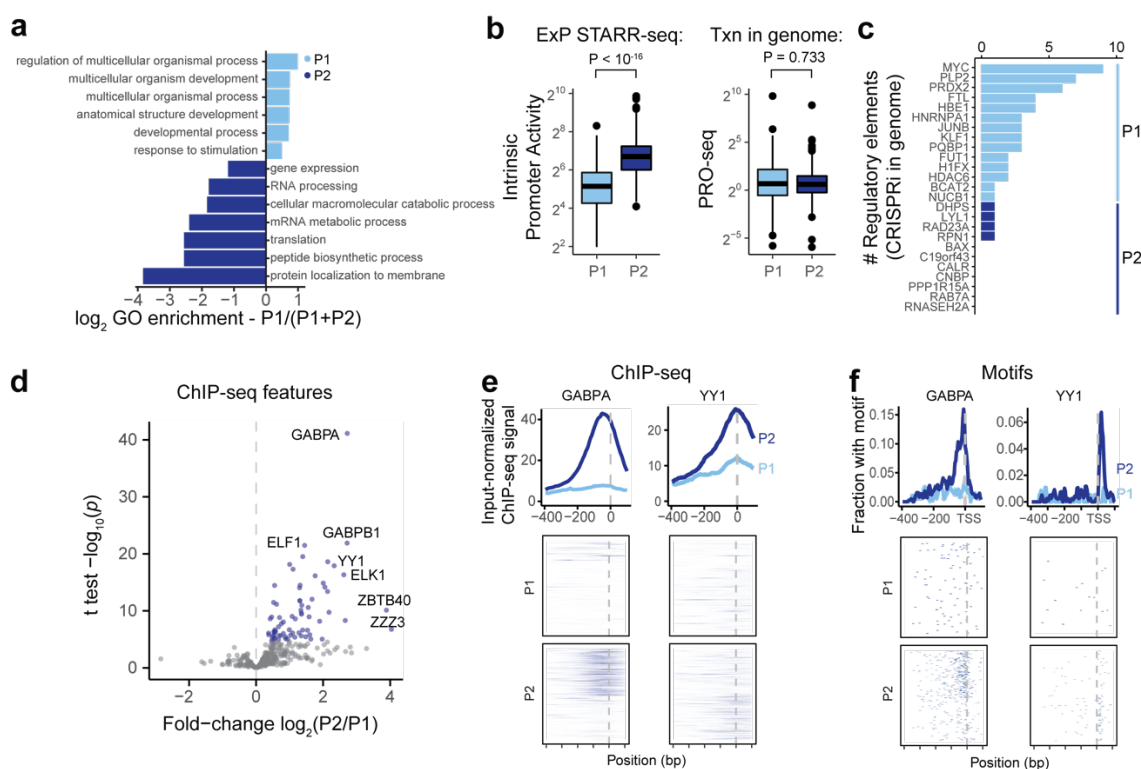
349

350 P1 promoters had on average 3.2-fold weaker intrinsic promoter activity than P2 promoters, as  
351 measured by ExP-STARR-seq ( $P < 10^{-16}$ , Mann-Whitney *U*-test; **Fig. 4b**), but showed similar  
352 levels of transcription in their native genomic locations, as measured by PRO-seq in the gene  
353 body ( $P = 0.733$ , Mann-Whitney *U*-test; **Fig. 4b**). This suggests that P1 promoters may be more  
354 dependent on genomic context for their level of transcription in the genome.

355

356 Genes corresponding to P1 promoters had more genomic regulatory elements in CRISPR  
 357 experiments. In data from previous studies, in which CRISPRi was used to perturb every  
 358 DNase-accessible element near selected promoters, the 14 genes corresponding to P1  
 359 promoters had an average of 3.6 (median: 3) distal enhancers in CRISPR experiments,  
 360 whereas the 11 genes corresponding to P2 promoters had only 0.36 (median: 0, **Fig. 4c**),  
 361 despite having similar numbers of nearby accessible elements (**Fig. S6a**). Distal enhancers for  
 362 P1 genes in the genome also had stronger effect sizes ( $P = 0.0071$ , *t*-test, **Fig. S6b**).  
 363

364 Together, these observations suggest that P1 promoter sequences correspond to context-  
 365 specific genes and depend more on distal enhancers for their transcriptional activation both in  
 366 ExP STARR-seq and in the genome, whereas P2 promoter sequences correspond to  
 367 constitutively expressed genes that are relatively less sensitive to distal enhancers in both  
 368 contexts.  
 369



370  
 371 **Fig. 4. Promoter classes correspond to enhancer-responsive versus constitutive genes**  
 372 **a.** Gene ontology  $\log_2$ -enrichment for P1 promoters using P1 and P2 promoters as a background set.  
 373 **b.** Intrinsic promoter activity for P1 vs P2 promoters (ExP STARR-seq) and genomic transcription level of  
 374 genes corresponding to P1 vs P2 promoters (PRO-seq reads per kilobase per million in gene bodies).  
 375 **c.** Number of activating genomic regulatory elements identified in comprehensive CRISPRi screens for  
 376 genes corresponding to P1 promoters ( $n=14$ ) and P2 promoters ( $n=11$ )<sup>22</sup>.  
 377 **d.** Volcano plot comparing ChIP-seq and other biochemical features for P2 versus P1 promoters (see  
 378 **Table S6**). X-axis: ratio of average signal at P2 versus P1 promoters. Blue points: features with  
 379 significantly higher signal at P2 promoters; no features have significantly higher signal at P1 promoters.  
 380 **e.** ChIP-seq signal for GABPA and YY1 in K562 cells at P1 and P2 promoters in the genome, aligned by  
 381 TSS (see Methods). Top: average ChIP signal (normalized to input) +/- 95% c.i. Bottom: signal at  
 382 individual genomic promoters.  
 383 **f.** Motif occurrences for GABPA and YY1 in P1 and P2 promoters, aligned by TSS.

## 384 TFs positioned at TSS distinguish constitutive from responsive promoters

385  
386 We next sought to identify sequence and chromatin features that distinguish P1 (“responsive”)  
387 from P2 (“constitutive”) promoters.

388  
389 We considered canonical core promoter motifs, which have been observed to differ between  
390 various subsets of promoters<sup>29–33</sup>, but did not find strong relationships. P1 and P2 promoter  
391 sequences had similar frequencies of the canonical ‘CA’ Initiator dinucleotide at the TSS (40.1%  
392 vs 35.3%, **Fig. S6c**), and corresponded to genes with similar patterns of dispersed versus  
393 focused TSSs in the genome (**Fig. S6d**). Consistent with previous studies comparing features of  
394 housekeeping versus other gene promoters<sup>29–33</sup>, P2 promoters had a slightly higher frequency  
395 of CpG dinucleotides (median 0.90 vs 0.81 normalized CpG content for P2 and P1 promoters,  
396 **Fig. S6e**), and P1 promoters had a 2-fold higher frequency of TATA box sequences upstream of  
397 the TSS (12.5% vs 6.1%), although only a small proportion of promoters contained this motif  
398 (**Fig. S6c**).

399  
400 Accordingly, we explored which other sequence features or TF binding measurements  
401 distinguished P2 constitutive from P1 responsive promoters. We examined 3,206 other features  
402 (including ChIP-seq measurements, TF motif predictions, and other features), and identified  
403 striking differences in the frequencies of certain transcription factor binding sites and motifs (**Fig.**  
404 **4d**, **Fig. S6f**, **Table S7**, see Methods). The most significantly enriched features included ChIP-  
405 seq signal for ETS family factors (GABPA, ELK1, ELF1), YY1, HCFC1, NR2C1, and C11orf30 /  
406 EMSY (**Fig. 4d**, **Fig. S7**). For example, two of the top factors (GABPA and YY1) together  
407 showed strong binding to a total of 64% of P2 promoters in the genome: 50% of P2 promoters  
408 showed strong GABPA binding (vs 8% of P1 promoters;  $P = 9.9 \times 10^{-22}$ , BH-corrected Fisher’s  
409 exact test), and 29% of P2 promoters showed strong YY1 binding (vs 5% of P1 promoters,  $P =$   
410  $9.4 \times 10^{-9}$ , BH-corrected Fisher’s exact test) (**Fig. 4e**). Notably, the sequence motifs for these  
411 factors showed positional preferences consistent with a function in regulating transcription  
412 initiation: the motif for GABPA was typically located 0-20 nucleotides upstream of the TSS  
413 (mode: –10), and for YY1 was often positioned at either +18 bp (both strands) or +2 bp  
414 (negative strand) from the TSS (**Fig. 4f**, **Fig. S6g**). Consistent with these factors playing a  
415 functional role, previous studies have found that adding GABPA or YY1 motifs to promoters  
416 increases gene expression in various reporter assays and cell types<sup>34–37</sup>.

417  
418 Together, these analyses suggest that P2 promoters can best be distinguished from P1  
419 promoters by the presence of certain transcription factors including GABPA and YY1, rather  
420 than canonical core promoter motifs.

421  
422  
423  
424

## 425 P2 constitutive promoters contain ‘built-in’ enhancer sequences

426  
427 We considered how transcription factors such as GABPA and YY1 might contribute to the  
428 reduced enhancer responsiveness of P2 versus P1 promoters. Interestingly, we noticed that  
429 these same factors showed strong binding in the genome not only at P2 promoters (**Fig. 4e,f**),  
430 but also at some E1 enhancers (**Fig. S5a**, **Fig. S7b**). For example, 3 of the genomic enhancers  
431 for *HBE1* (all classified as E1 in ExP STARR-seq) contained GABPA sequence motifs and  
432 showed strong GABPA binding by ChIP-seq, whereas the genomic promoter of *HBE1*  
433 (classified as P1) lacked these features (**Fig. 5a**).

434

435 These observations suggested that P2 promoters may have reduced responsiveness to E1  
436 enhancers because they contain some of the same motifs, potentially saturating some step in  
437 transcription. Accordingly, we explored the hypothesis that P2 promoters contain 'built-in' E1  
438 enhancer sequences that would increase promoter activity and decrease responsiveness to  
439 distal E1 enhancers.

440  
441 Consistent with this hypothesis, we found that (i) across all promoters, responsiveness to E1  
442 enhancers was inversely correlated with intrinsic promoter activity, in a way that appeared to  
443 saturate; (ii) P2 promoters had stronger enhancer activity than P1 promoters; and (iii) nearly all  
444 of the TF motifs enriched in P2 promoters were predictive of both promoter activity and  
445 enhancer activity:

446  
447 We first compared intrinsic promoter activity with responsiveness to E1 enhancers, and found  
448 that they were correlated both when considering all promoters in ExP STARR-seq (Pearson  $R =$   
449  $-0.62$ ,  $\log_2$  space; **Fig. 5b**) and when considering only P1 promoters ( $R = -0.51$ ). For example,  
450 comparing P1 promoters at opposite extremes, the *RAD23A* promoter (P2) had 11.8-fold higher  
451 intrinsic promoter activity compared to the *HBE1* promoter (P1), and was 2.1-fold less sensitive  
452 to E1 enhancers. As promoter activity increased, responsiveness to E1 enhancers decreased  
453 rapidly (for example, from ~9-fold average activation by E1 enhancers for the *SNAI3* P1  
454 promoter) and appeared to saturate at ~3-fold for most P2 promoter sequences (**Fig. S8a**).

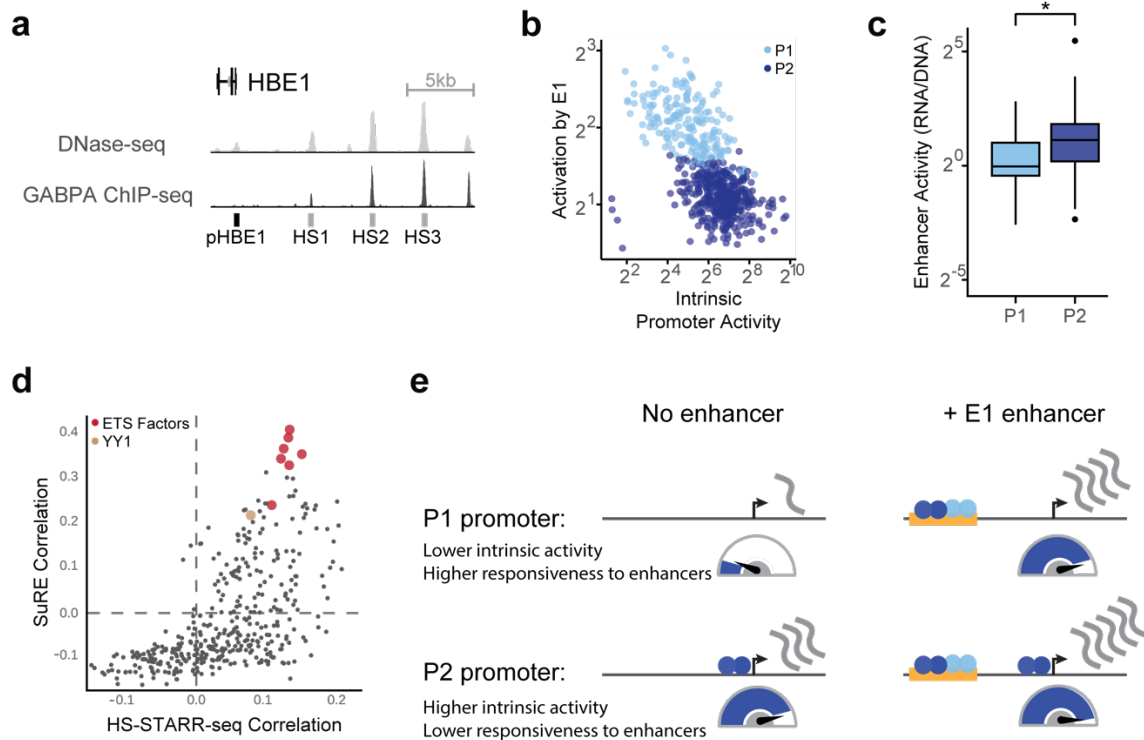
455  
456 We next tested whether P2 promoters had stronger intrinsic enhancer activity. To do so, we  
457 generated a second STARR-seq dataset in which we measured the enhancer activity of >8.9  
458 million sequences derived from DNase-accessible elements and promoters (by hybrid selection  
459 (HS)-STARR-seq, see Methods, **Fig. S8b-d**). In this dataset, many promoter elements tested in  
460 ExP STARR-seq (along with thousands of other accessible elements) were densely tiled (an  
461 average of ~11 fragments each covering at least 90% of the promoters tested in the ExP  
462 assay), allowing us to test the enhancer activity of entire P1 and P2 promoter sequences. P2  
463 promoters indeed showed ~2-fold higher intrinsic enhancer activity than P1 promoters in HS-  
464 STARR-seq ( $P = 1.14 \times 10^{-16}$ ,  $t$ -test, **Fig. 5c**), supporting a model where these promoters  
465 contain built-in enhancers.

466  
467 Finally, we examined whether the sequence motifs enriched in P2 promoters contribute to both  
468 enhancer activity and promoter activity. To do so, we examined data on enhancer activity from  
469 HS-STARR-seq along with another previous experiment that measured promoter activity for  
470 millions of random genomic fragments in K562 cells (SuRE<sup>38</sup>). 16 of the 17 motifs enriched in  
471 P2 promoters, including motifs for GABPA and YY1, were positively correlated with both  
472 enhancer activity and promoter activity (**Fig. 5d, Table S7**, see Methods).

473  
474 Together, these observations suggest a model for promoter sequence organization (**Fig. 5e**). P2  
475 promoters encode binding motifs for activating factors, including GABPA and YY1, that act as  
476 'built-in' enhancers for the promoter. This not only increases the autonomous activity of the  
477 promoter, but also reduces its responsiveness to distal enhancers. P1 promoters, in contrast,  
478 appear to exclude these activating factors, creating a sensitivity to distal enhancers.

479  
480





481  
482  
483  
484  
485  
486  
487  
488  
489  
490  
491  
492  
493  
494  
495  
496  
497  
498

**Fig. 5. P2 constitutive promoters contain built-in enhancer sequences**

**a.** DNase-seq and GABPA ChIP-seq binding at the HBE1 promoter (pHBE1) and HS1-HS3 enhancers.

**b.** Correlation between intrinsic promoter activity and responsiveness of promoters to E1 enhancers (average activation by E1 sequences, expressions vs. random genomic controls). Each point is one promoter.

**c.** Average enhancer activity in HS-STARR-seq (RNA/DNA) of P1 and P2 promoters.  $*P = 1.14 \times 10^{-16}$ , *t*-test.

**d.** For each of 400 sequence motifs that appeared in at least 5% of HS-STARR-seq fragments, correlation (Pearson *R*) of motif occurrence with intrinsic promoter activity (SuRE signal, y-axis) and with intrinsic enhancer activity (HS-STARR-seq signal among fragments not overlapping TSS, x-axis).

**e.** A model for enhancer-promoter compatibility. Enhancers multiplicatively scale the RNA output of promoters. P2 constitutive promoters contain built-in activating sequence motifs that both increase intrinsic promoter activity and reduce responsiveness to distal enhancers.



## 499 Discussion

500  
501 Since the discovery of the first enhancers forty years ago<sup>1,2</sup>, many enhancer and promoter  
502 sequences have been combined and found to be compatible<sup>3-9</sup>. At the same time, studies of  
503 individual natural or synthetic core promoters have been found to have some degree of  
504 specificity when combined with various transcriptional cofactors or enhancer sequences<sup>10-15</sup>.

505  
506 Here we develop and apply ExP STARR-seq to systematically quantify enhancer-promoter  
507 compatibility, and identify a simple rule for combining human enhancer and promoter activities.  
508 Enhancers are intrinsically compatible with many Pol II promoter sequences, and act  
509 multiplicatively to scale the RNA output of a promoter. As a result, independent control of  
510 intrinsic enhancer activity and intrinsic promoter activity can create significant variation in RNA  
511 expression: in our data, promoter activity and enhancer activity each vary over 3-4 orders of  
512 magnitude, and their multiplicative combination leads to >4-million-fold variation in STARR-seq  
513 expression. This finding of broad compatibility appears to be consistent with recent studies  
514 using reporters integrated into the genome, which found that human core promoters or  
515 enhancers are similarly scaled when they are inserted into different genomic loci<sup>39,40</sup>. This is  
516 also consistent with our previous finding that the effects of enhancers on nearby genes in the  
517 genome can be predicted with good accuracy using a model based only on genomic  
518 measurements of enhancer activity and distance-based 3D contacts, assuming no intrinsic  
519 enhancer-promoter specificity<sup>22</sup>. While there may be circumstances where promoters are  
520 responsive only to certain cofactors or enhancer sequences<sup>10-15</sup>, our observations indicate that  
521 biochemical specificity is not the dominant factor controlling the activity levels of human  
522 enhancers and promoters.

523  
524 Superimposed on this multiplicative function, we identify two classes of enhancers and  
525 promoters that show subtle preferences in activation. One class of promoters, corresponding  
526 largely to constitutively expressed (housekeeping) genes, is less responsive to distal enhancers  
527 both in ExP STARR-seq and in the genome, while the second class of promoters,  
528 corresponding to cell-type- or context-specific genes, is more responsive. Previous studies have  
529 identified numerous differences in sequence content and motifs between the promoters of  
530 housekeeping and context-specific genes<sup>29-33</sup>. We find that these promoters indeed show  
531 intrinsic differences in their levels of activity and responsiveness to enhancers. Interestingly, this  
532 pattern of promoter responsiveness also can be predicted by a simple logic: P2 promoters  
533 contain built-in activating sequences that increase both enhancer and promoter activity, which  
534 appears to reduce their responsiveness to distal enhancers. This model for human promoters  
535 appears to differ qualitatively from previous studies in *Drosophila*, which found that the  
536 promoters of housekeeping and developmentally regulated genes can both be highly  
537 responsive, but to distinct sets of enhancer sequences and cofactors<sup>16,18</sup>. We note that one  
538 methodological difference is that, whereas these previous studies focused on minimal core  
539 promoter sequences (100-138bp total), here we included more sequence context upstream of  
540 the TSS (264 bp total).

541  
542 A remaining challenge will be to link the sequences that control enhancer and promoter  
543 activities with effects on particular biochemical steps in transcription. In this regard, we find that  
544 GABPA and YY1 bind both to constitutive promoters and to distal enhancers, and are  
545 associated with increased enhancer activity, increased promoter activity, and reduced promoter  
546 responsiveness to distal enhancers. This suggests that distal enhancers may act, in part, on a  
547 particular rate-limiting step in transcription that can be saturated by inclusion of built-in activating  
548 sequences in a gene promoter. Indeed, a previous study found that adding GABPA and YY1  
549 motifs to several promoters led to an increase in RNA expression that saturates at 2 or 5 copies

550 of the motif, respectively.<sup>34</sup> Given the preferred positions of these motifs within 20 bp of the TSS  
551 — as well as previous findings that these proteins physically interact with general transcription  
552 factors<sup>41,42</sup> and/or influence transcriptional initiation and TSS selection<sup>36,43–45</sup> — such a rate-  
553 limiting step might involve assembly of the preinitiation complex. In addition to this step, our  
554 data are consistent with a model in which enhancers and promoters control additional steps in  
555 transcription that combine multiplicatively and do not saturate in the dynamic range of our  
556 assay. Examples of such processes that could combine multiplicatively include control of burst  
557 frequency and burst size<sup>46</sup>. Further work will be required to investigate these possibilities.

558  
559 Together, our findings support a simple logic for human enhancer-promoter compatibility, and  
560 will propel efforts to model gene expression, map the effects of human genetic variation, and  
561 design regulatory sequences for gene therapies.  
562

563 **Acknowledgements**

564 This work was supported by an NHGRI Genomic Innovator Award (R35HG011324 to J.M.E.);  
565 Gordon and Betty Moore and the BASE Research Initiative at the Lucile Packard Children's  
566 Hospital at Stanford University (J.M.E.); an NIH Pathway to Independence Award  
567 (K99HG009917 and R00HG009917 to J.M.E.); the Harvard Society of Fellows (J.M.E.); the  
568 Broad Institute (E.S.L.); an AQA Carolyn L. Kuckein Student Research Fellowship (D.T.B.); and  
569 by the National Institute of General Medical Sciences (T32GM007753, L.S.). We thank C.  
570 Vockley, V. Subramanian, and members of the Engreitz and Lander labs for discussions and  
571 technical assistance.

572

573 **Author Contributions**

574 D.T.B., C.P.F., T.R.J., and J.M.E. developed the ExP STARR-seq assay. D.T.B., M.K., and  
575 T.H.N. performed experiments. D.T.B., T.R.J., V.L., L.S., H.Y.K., J.N., S.R.G., and J.M.E.  
576 analyzed data. E.S.L. and J.M.E. supervised the work. All authors contributed to writing the  
577 manuscript.

578

579 **Competing Interests**

580 C.P.F. is now an employee of Bristol Myers Squibb. J.M.E. is a shareholder of Illumina, Inc. All  
581 other authors declare no competing interests.

582

583 **Data Availability**

584 Raw and processed data for ExP STARR-seq and HS STARR-seq can be found in NCBI GEO  
585 under accession number GSE184426. Luciferase data can be found in Supplementary Table  
586 S3.

587

588 **Tables**

589 **S1.** Gene promoters used in ExP STARR-seq

590 **S2.** Candidate enhancers used in ExP STARR-seq

591 **S3.** ExP-Luciferase elements and data

592 **S4.** Biochemical feature enrichment in E1 vs. E2 enhancers

593 **S5.** Transcription factor motif enrichment in E1 vs. E2 enhancers

594 **S6.** Biochemical feature enrichment in P1 vs. P2 promoters

595 **S7.** Transcription factor motif enrichment in P1 vs. P2 promoters

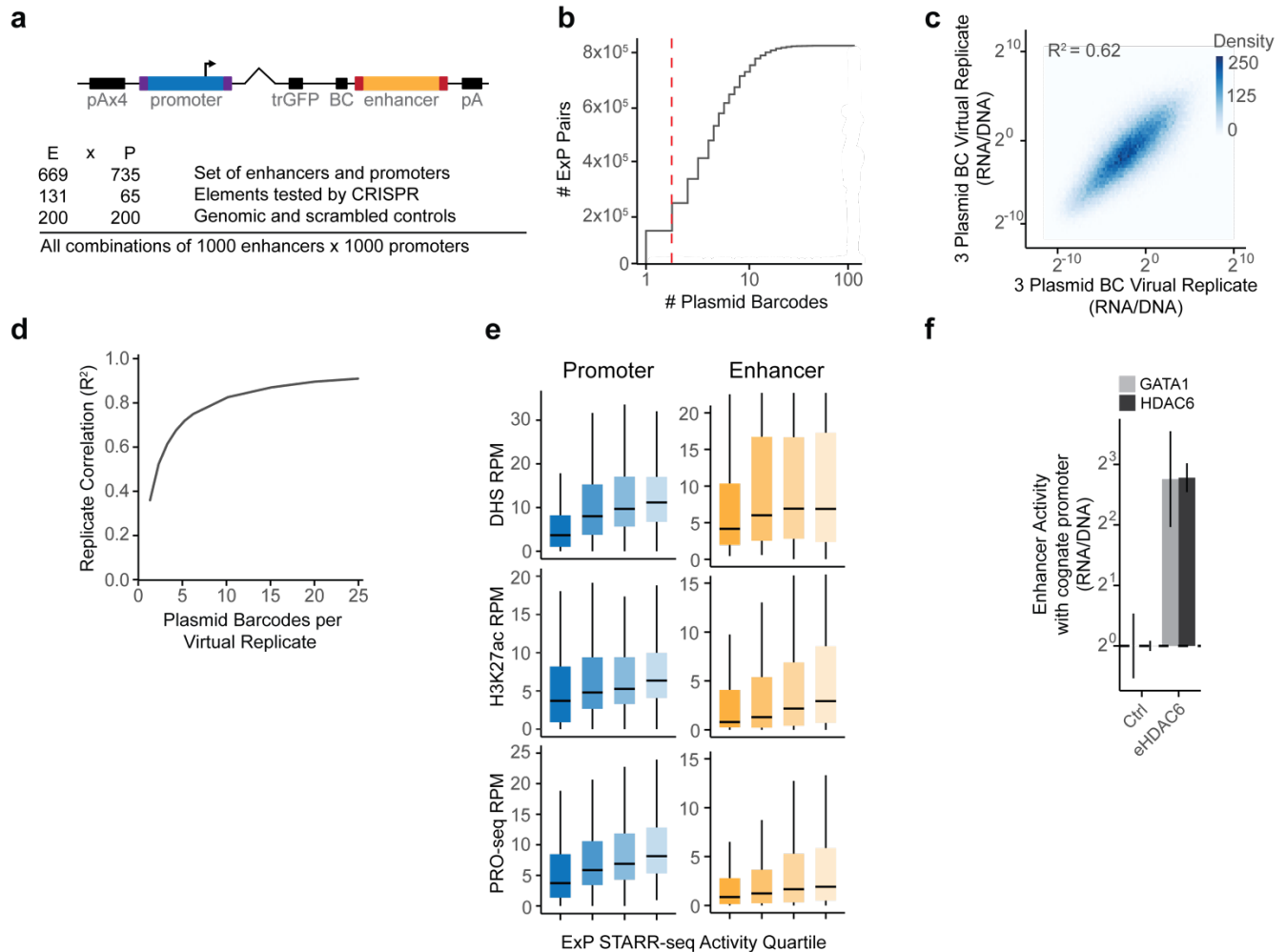
596 **S8.** Primer and oligo sequences

597 **S9.** ENCODE datasets used to annotate ExP enhancers and promoters

598 **S10.** Enhancer hybrid selection probe sequences

599 **S11.** Promoter hybrid selection probe sequences

## 600 Supplementary Figures

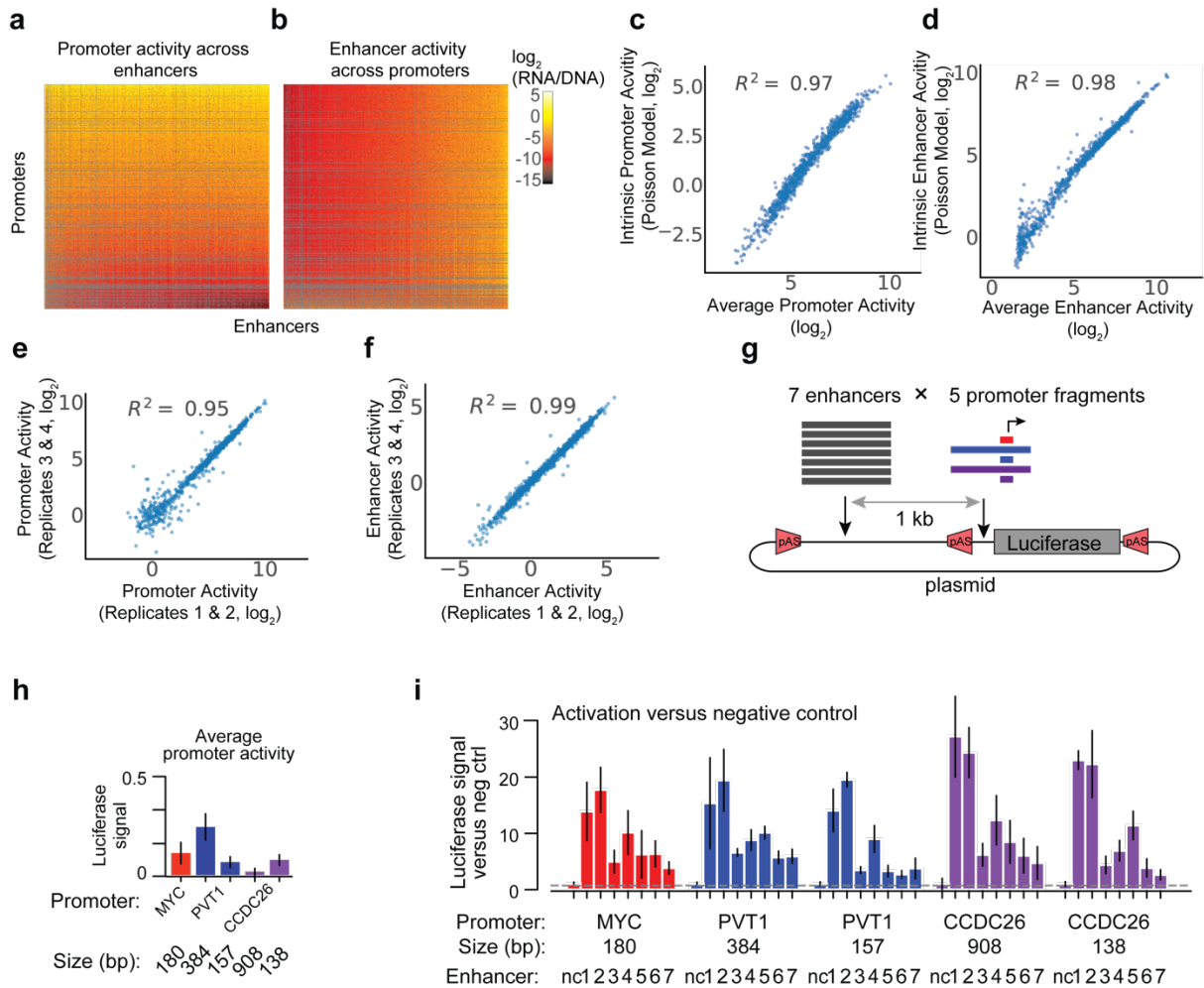


601  
602  
603

### Fig. S1. Design and reproducibility of ExP STARR-seq

604 **a.** ExP STARR-seq reporter construct (pA = polyadenylation signal; purple = promoter sequencing  
605 adaptors; angled = spliced sequence; trGFP = truncated GFP open reading frame; BC = 16bp N-mer  
606 plasmid barcode; red = enhancer sequencing adaptors) and 1000x1000 K562 library contents.  
607 **b.** Distribution of plasmid barcodes per enhancer-promoter pair, red dotted-line is threshold of two  
608 plasmid barcodes.  
609 **c.** Correlation between virtual replicates, formed by sampling two nonoverlapping groups of three plasmid  
610 barcodes from pairs with at least 6 barcodes, and averaging  $\log_2(\text{RNA/DNA})$  within groups.  
611 **d.** Correlation between virtual replicates as in (c) for increasing numbers of plasmid barcodes per pair in  
612 virtual replicates.  
613 **e.** DNase-seq, H3K27ac ChIP-seq, and PRO-seq (RPM) by increasing quartile of autonomous promoter  
614 activity and average enhancer activity in ExP STARR-seq. Box: median and interquartile range (IQR).  
615 Whiskers:  $\pm 1.5 \times \text{IQR}$ .  
616 **f.** Activation in ExP STARR-seq (expression versus genomic controls in distal position) of GATA1 and  
617 HDAC6 promoters by eHDAC6 (chrX:48641342-48641606). Ctrl = activity of promoters with random  
618 genomic controls in enhancer position. Error bars: 95% CI across plasmid barcodes. n = 7 (GATA1-ctrl),  
619 381 (HDAC6-ctrl), 4 (eHDAC6-GATA1), 37 (eHDAC6-HDAC6).  
620

621



622  
623  
624  
625

**Fig. S2. Comparison of methods of estimating enhancer and promoter activities and validation of multiplicative model using luciferase assays**

626 **a-b.** Heatmap of promoter activity (**a**, expression divided by intrinsic enhancer activity) or enhancer  
627 activity (**b**, expression divided by intrinsic promoter activity) across all pairs of promoter (vertical) and  
628 enhancer sequences (horizontal). Axes are sorted by intrinsic promoter and enhancer activities, as in Fig.  
629 2j. Grey: missing data.

630 **c-d.** Correlation between two estimates of promoter (**c**) and enhancer (**d**) activities. One method  
631 (“average activity”, x-axis) estimates activity calculated by averaging across elements, and the other  
632 method (“intrinsic activity”, y-axis) estimates activity by using coefficients estimated by a Poisson count  
633 model (see Methods).

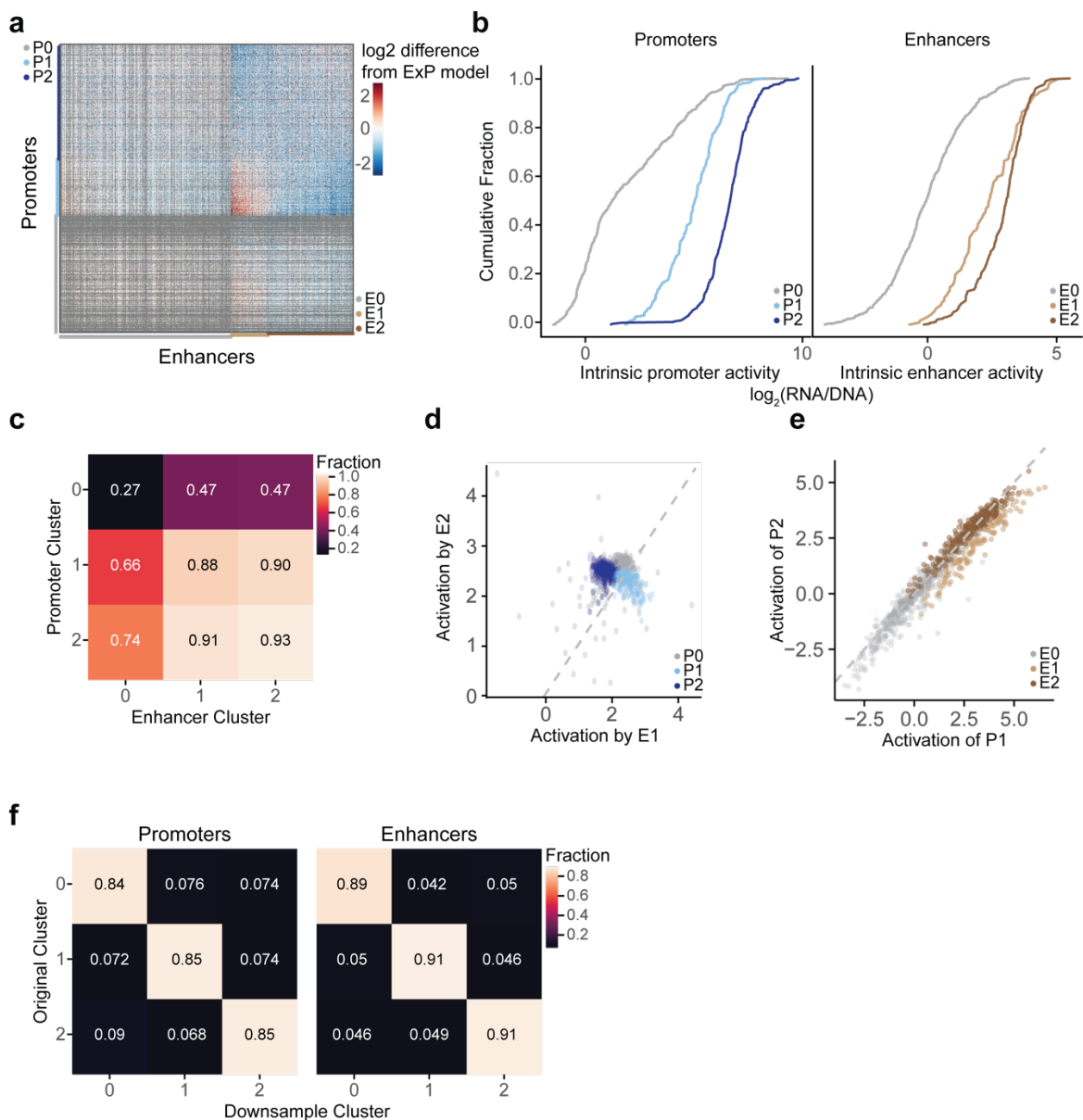
634 **e-f.** Correlation of intrinsic promoter (**e**) and enhancer (**f**) activity estimates from Poisson model using data  
635 from separate replicate experiments.

636 **g.** ExP luciferase reporter construct. Seven enhancer fragments, with flanking polyadenylation signals,  
637 were cloned upstream of five promoter fragments and measured via the dual luciferase assay.

638 **h.** Autonomous promoter activity of ExP luciferase (average luciferase signal of promoter with negative  
639 control) for 5 promoter sequences derived from 3 genes (*MYC*, *PVT1*, *CCDC26*). Error bars are 95% CI  
640 from three biological replicates.

641 **i.** Enhancer activation (luciferase signal versus negative control sequence in the enhancer position) of  
642 seven enhancers across five promoter fragments. Error bars are 95% CI from three biological replicates.

643  
644



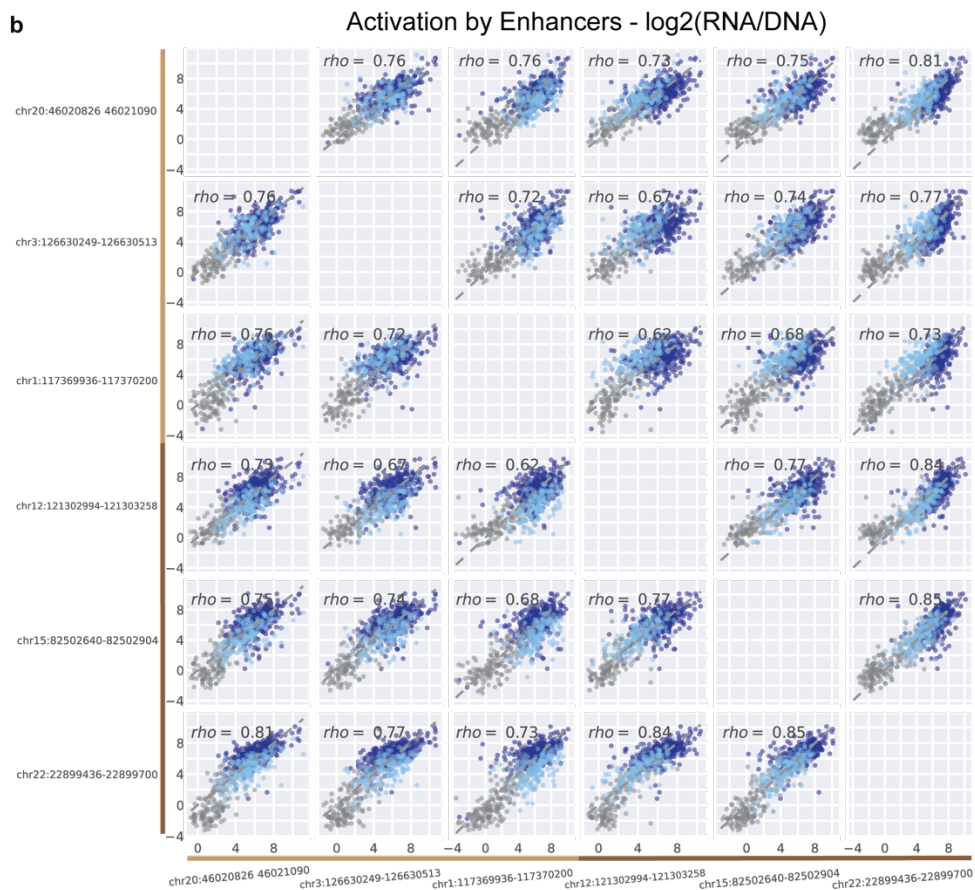
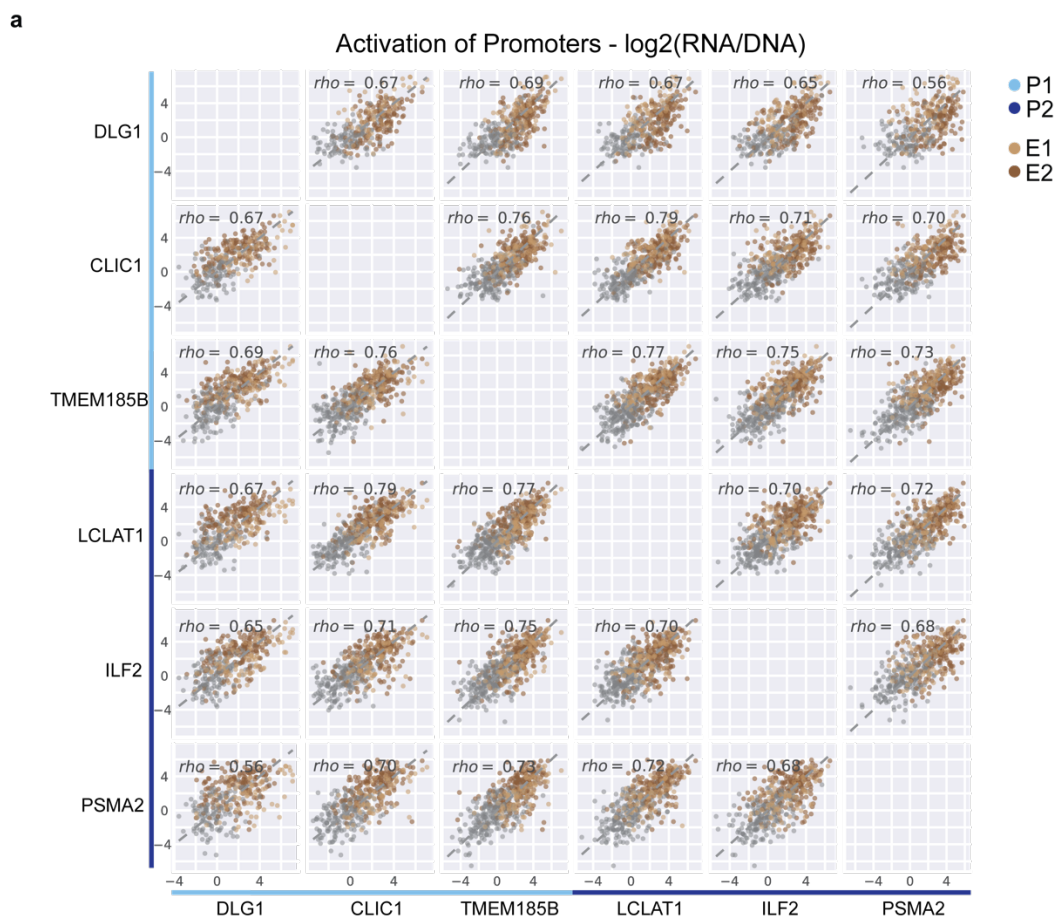
645  
646  
647

**Fig. S3. Enhancer and promoter cluster identification and reproducibility**

648 **a.** Heatmap of deviations in enhancer-promoter STARR-seq expression from a multiplicative enhancer-  
 649 promoter model (color scale: fold-difference between observed expression versus expression predicted  
 650 by multiplicative model; gray: missing data). Same as Fig 3a, except including clusters with weak  
 651 sequences and missing data (E0 and P0). Vertical axis: promoter sequences grouped by class and sorted  
 652 by responsiveness to E1 vs. E2; horizontal axis: enhancer sequences grouped by class and sorted by  
 653 activation of P1 vs. P2.  
 654 **b.** Distribution of intrinsic enhancer and promoter activity (expression versus genomic controls) by  
 655 cluster.  
 656 **c.** Fraction of enhancer-promoter pairs observed in ExP STARR-seq dataset ( $\geq 2$  plasmid barcodes)  
 657 by cluster.



- 658 **d.** Correlation of average promoter activation (expression versus genomic controls in enhancer position)  
659 by E2 versus E1 enhancer sequences. Each point is one promoter sequence. Same as Fig. 3c, except  
660 including P0 promoter sequences.
- 661 **e.** Correlation of average activation of P2 versus P1 promoters. Each point is one enhancer  
662 sequence. Same as Fig. 3d, except including E0 enhancer sequences.
- 663 **f.** Robustness of enhancer and promoter cluster assignments to downsampling of enhancer and promoter  
664 sequences. Clustering was repeated in 100 random downsamplings to 25% of promoter sequences and  
665 25% of enhancer sequences (6.25% of original matrix). Heatmap: Average fraction overlap between  
666 cluster assignments from the full and downsampled matrices.



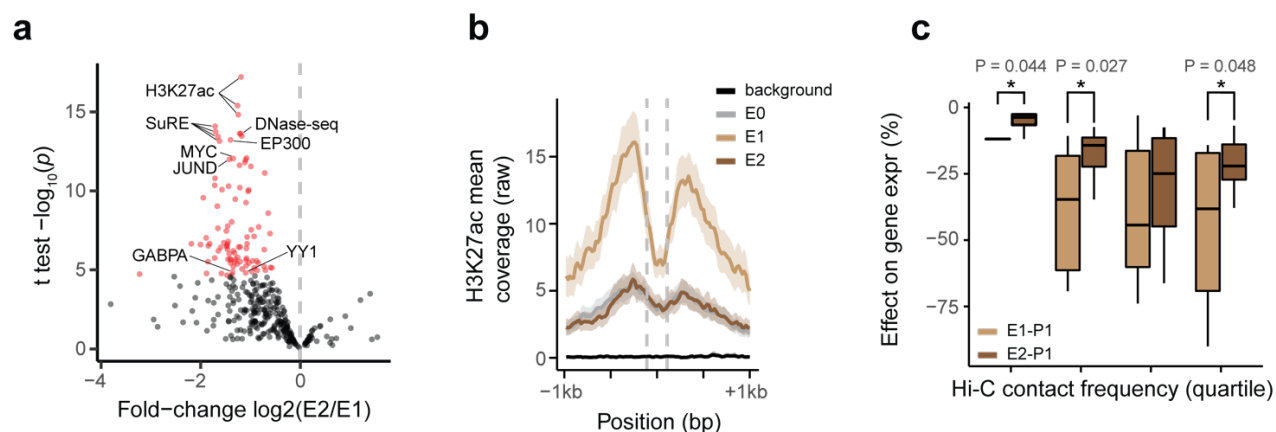
668 **Fig. S4. Classes of enhancer and promoter sequences show distinct patterns of**  
669 **activation and responsiveness.**

670 **a.** For 6 representative promoter sequences (3 P2 and 3 P1 sequences), the pairwise correlation of  
671 activation by enhancers (expression versus genomic controls in enhancer position, averaged across  
672 plasmid barcodes). Each point is one enhancer sequence.

673 **b.** For 6 representative enhancer sequences (3 E1 and 3 E2 sequences), the pairwise correlation of  
674 promoter activation (expression versus genomic controls in promoter position, averaged across plasmid  
675 barcodes). Each point is one promoter sequence.

676  
677

678  
679  
680

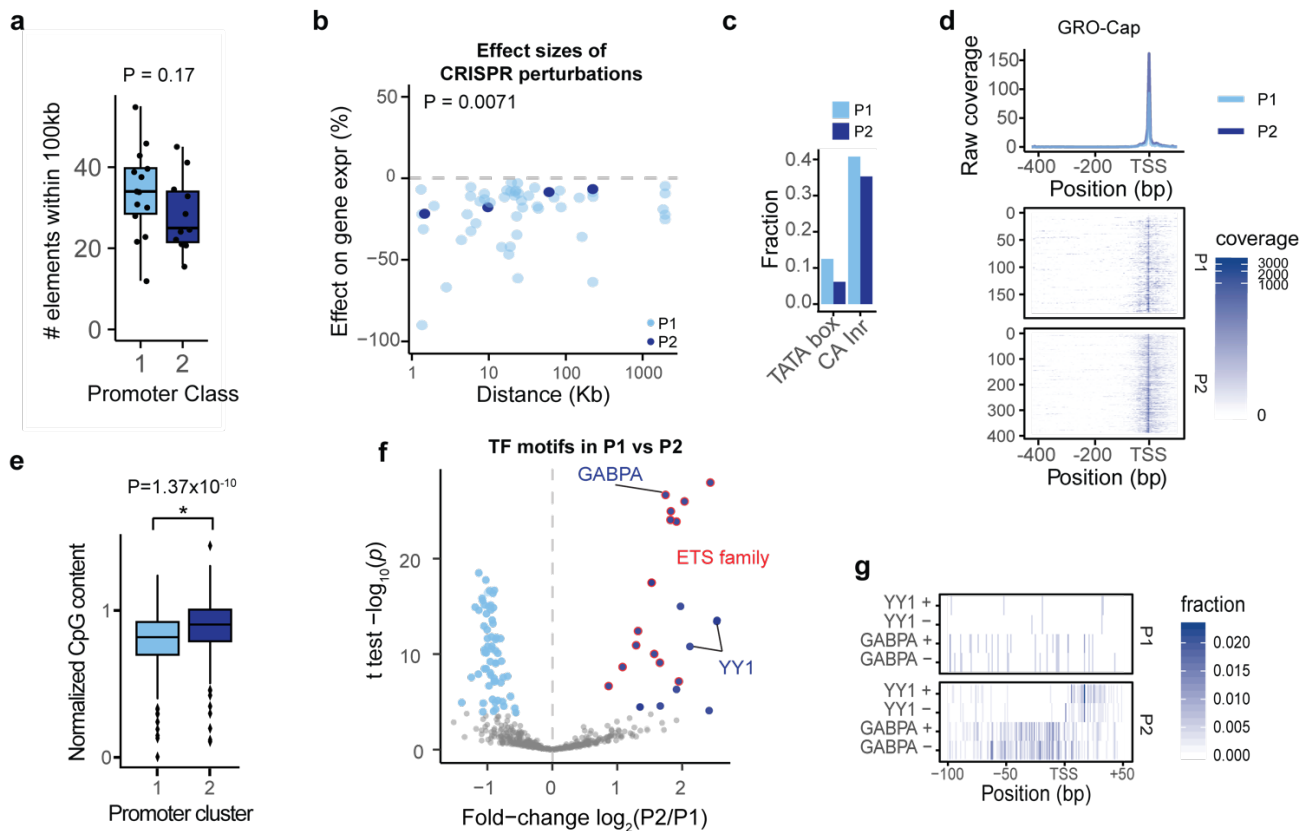


681  
682  
683

**Fig. S5. Classes of enhancer sequences correspond to strong and weak genomic enhancers**

684 **a.** Volcano plot comparing ChIP-seq and other genomic features for E2 versus E1 enhancer sequences  
685 (see **Table S4**). X-axis: ratio of average signal at P2 versus P1 promoters. Red dots: features with  
686 significantly higher signal at E1; no features have significantly higher signal at E2 enhancer sequences.  
687 **b.** Mean H3K27ac ChIP-seq coverage of genomic elements corresponding to E0, E1, E2, or genomic  
688 control enhancer sequences ( $\pm$  95% CI), aligned by DHS peak summit. Dotted lines mark bounds of the  
689 enhancer sequences used in ExP STARR-seq.  
690 **c.** % effect of genomic elements corresponding to E1 vs. E2 enhancer sequences on expression of genes  
691 corresponding to P1 promoters in CRISPRi screens, separated by quartiles of 3D contact frequency  
692 measured by Hi-C (0.39-11.9, 11.9-23.9, 23.9-58.3, 58.3-100). \* $P < 0.05$ , two-sample  $t$ -test.  
693

694  
695

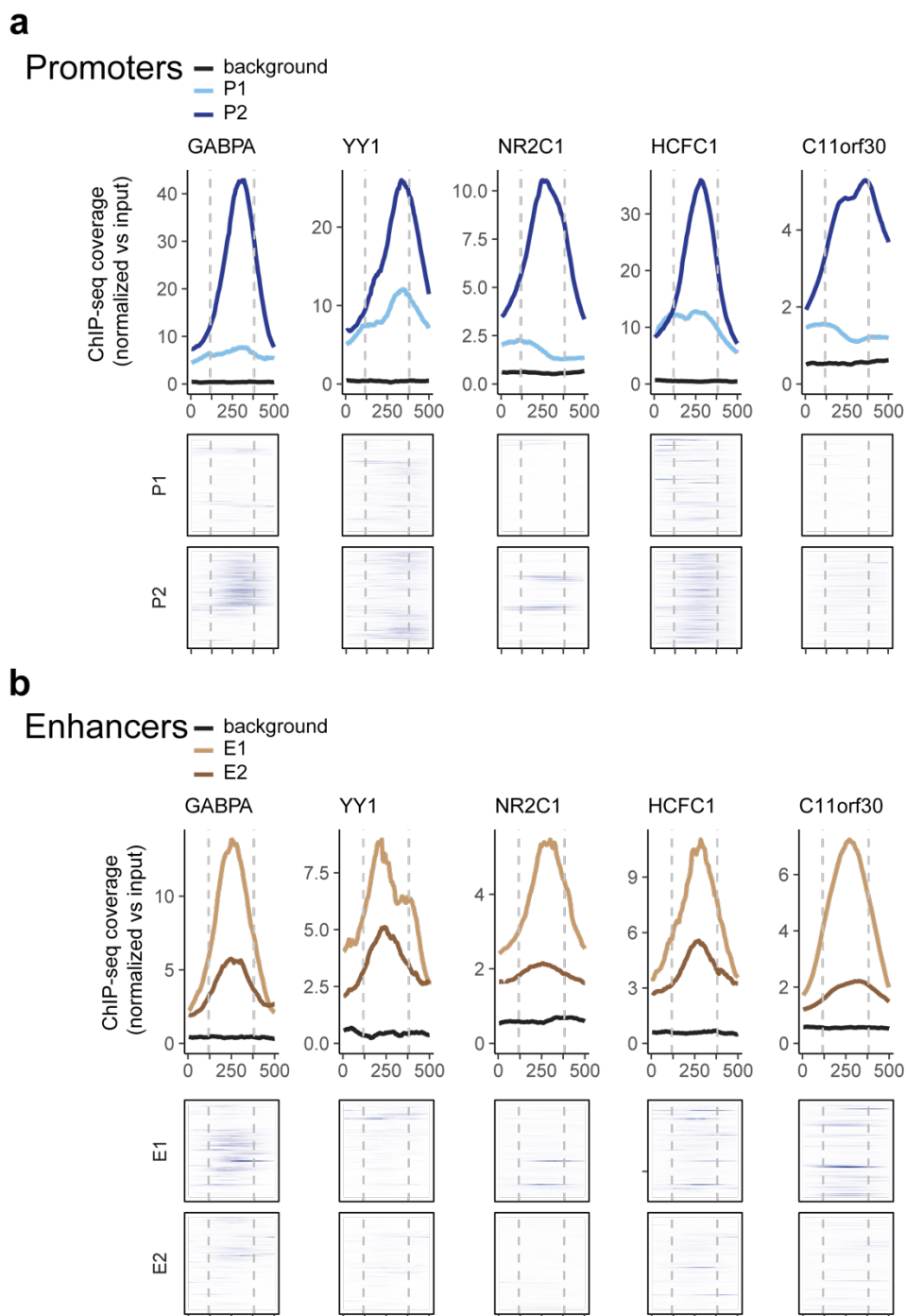


696  
697

### Fig. S6. Properties of promoter classes

- 698 **a.** Number of nearby accessible elements (within 100 Kb of the gene promoter, considering top 150,000  
699 DNase peaks in K562 cells as used in the ABC model<sup>22</sup>) for the 14 genes corresponding to P1 promoters  
700 and 11 genes corresponding to P2 promoters with comprehensive CRISPR tiling data.  $P = 0.17$ , Mann-  
701 Whitney U test.
- 702 **b.** % Effect of CRISPRi perturbations to genomic regulatory elements on genes corresponding to P1 vs.  
703 P2 promoters.  $P = 0.0071$ ,  $t$ -test.
- 704 **c.** Fraction of promoter sequences containing TATA or CA initiator core promoter motifs.
- 705 **d.** GRO-Cap coverage of genomic promoters aligned by TSS. Top: Mean coverage of genomic promoters  
706 corresponding to P1 vs. P2 classes. Bottom: Coverage across all individual promoters.
- 707 **e.** Normalized CpG-content of P1 and P2 promoter sequences, calculated as the ratio of observed to  
708 expected CpG =  $(\text{CpG fraction}) / ((\text{GC content})^2 / 2)$ .
- 709 **f.** Volcano plot comparing frequency of transcription factor motifs in P2 versus P1 promoter sequences  
710 (see **Table S7**). X-axis: ratio of average motif counts in P2 versus P1 promoter sequences. Light blue and  
711 dark blue dots: Motifs significantly more frequent in P1 or P2 promoter sequences, respectively. Red  
712 outline: significant motifs for ETS family transcription factors.
- 713 **g.** Fraction of P2 promoter sequences with YY1 and GABPA binding motifs by nucleotide position,  
714 aligned by TSS and separated by strand (see Methods).

715  
716



717  
718  
719

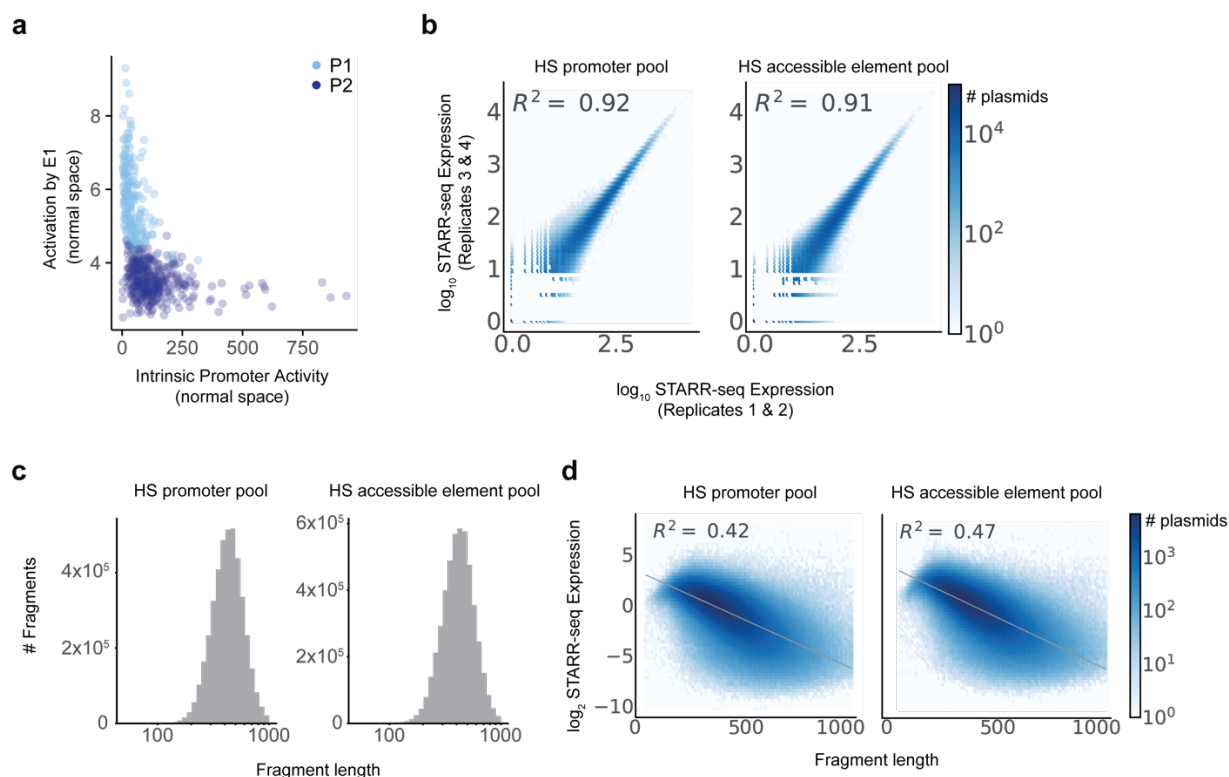
**Fig. S7. Transcription factors enriched at P2 promoters are also enriched at E1 enhancers**

720 **a.** ChIP-seq signal for 5 transcription factors in K562 cells at P1 and P2 promoters in the genome, aligned  
721 by boundaries of the 264-bp ExP STARR-seq promoter sequence (see Methods). Top: average ChIP-seq  
722 signal normalized to input. Bottom: signal at individual genomic promoters. Black line: average for random  
723 genomic control sequences.

724 **b.** ChIP-seq signal at E1 and E2 enhancers in the genome. Black line: average for random genomic  
725 control sequences.



726



727

728

729

**Fig. S8. Responsiveness to E1 enhancers versus intrinsic promoter activity, and metrics for hybrid-selection STARR-seq experiments**

730

731

732

733

734

735

736

737

738

739

740

741

**a.** Correlation between intrinsic promoter activity and responsiveness of promoters to E1 enhancers (average activation by E1 sequences, expressions vs. random genomic controls). Each point is one promoter. Same as Fig. 5b, but in normal scale instead of  $\log_2$  scale.

**b.** Correlation of HS-STARR-seq expression between biological replicate experiments for promoter and accessible element pools, calculated for individual elements with unique plasmid barcodes. Axes represent the average STARR-seq expression (RNA/DNA,  $\log_{10}$  scale) of two biological replicates. Density: number of plasmids.

**c.** Fragment length distribution in HS-STARR-seq in promoter and accessible element pools, of fragments with at least 25 DNA counts.

**d.** STARR-seq expression (y-axis) and fragment length (x-axis) relationship in HS-STARR-seq. Density: number of plasmids.

## 742 **Methods**

743

### 744 **Genome build**

745

746 All analyses and coordinates are reported using human genome reference hg19.

747

### 748 **Design of ExP STARR-seq**

749

750 We designed ExP STARR-seq to systematically measure the intrinsic, sequence-encoded  
751 compatibility or specificity of many pairs of human enhancer and promoter sequences. The key  
752 design features we considered when developing this assay were the ability to measure the  
753 activity of individual enhancer-promoter sequence combinations, to precisely quantify the  
754 expression of each enhancer-promoter pair, and to test hundreds of thousands of combinations  
755 in order to identify patterns of compatibility or specificity across a large number of human  
756 sequences.

757

758 Accordingly, we designed a new variant of the STARR-seq high-throughput plasmid reporter  
759 assay called enhancer x promoter (ExP) STARR-seq. In both STARR-seq and ExP STARR-  
760 seq, enhancer sequences are cloned downstream of a promoter, transfected into cells, and  
761 transcribed to produce a reporter mRNA transcript, which is then sequenced to quantify the  
762 relative expression levels of plasmids containing different enhancer sequences<sup>4</sup>. In ExP  
763 STARR-seq, we modified the cloning and RNA sequencing strategy to enable testing different  
764 enhancer sequences in combination with different promoter sequences.

765

766 To clone combinations of enhancer and promoter sequences into a reporter plasmid, we  
767 synthesized 264-bp enhancer and promoter sequences in an oligo pool format, PCR amplified  
768 enhancer and promoter sequences separately, and inserted them into the hSTARR-seq\_SCP1  
769 vector\_blocking 4 vector<sup>8</sup> in the promoter position (replacing the original SCP1 promoter) or  
770 enhancer position in a single pooled cloning step using Gibson assembly to generate all  
771 pairwise combinations of chosen enhancer and promoter sequences (**Fig. 1a, Fig. S1a**). We  
772 chose this specific STARR-seq vector with 4 polyA sequences upstream of the promoter  
773 position because it was specifically designed in order to avoid spurious transcription initiation  
774 from the origin of replication<sup>8</sup>, which would interfere with the STARR-seq signal from the cloned  
775 enhancer-promoter pairs. This STARR-seq vector also includes 5' and 3' splice sites upstream  
776 of the enhancer that allows using a PCR primer targeting the splice junction to specifically  
777 amplify cDNA derived from the reporter mRNA while avoiding amplifying the plasmid DNA  
778 sequence.

779

780 To quantify the reporter mRNA transcripts and determine which enhancer-promoter pair they  
781 correspond to, we further adapted the cloning and RNA sequencing design. In the standard  
782 STARR-seq assay, the reporter mRNA contains the enhancer sequence but not the full  
783 promoter sequence, and therefore cannot determine from which promoter a given reporter  
784 mRNA is derived. Accordingly, in ExP STARR-seq we introduced a random 16-mer sequence  
785 located just upstream of the enhancer sequence that we use as a “plasmid barcode” to identify  
786 which enhancer reporter mRNAs are derived from which enhancer-promoter pairs (**Fig. 1a**).  
787 After cloning the plasmid pool, we map which plasmid barcodes correspond to which promoters  
788 by applying Illumina high-throughput sequencing to a PCR amplicon containing the promoter  
789 sequence and plasmid barcode. From this, we build a dictionary to look up, for a given reporter  
790 mRNA containing a plasmid barcode and enhancer sequence, which enhancer-promoter-  
791 plasmid barcode construct that mRNA is derived from.

792

793 Finally, we selected the number of constructed tested (~1 million pairs of enhancer and  
794 promoter sequences cloned, with an average of 6.3 plasmid barcodes per pair) and sequencing  
795 depth (>1 billion reads per replicate) to enable highly precise measurements of expression for  
796 each enhancer-promoter pair. We obtained high reproducibility of enhancer-promoter  
797 expression levels between biological replicates ( $R^2 = 0.92$ ), allowing us to develop quantitative  
798 models of how enhancer and promoter activities combine.

799

800 Altogether, this approach enables precisely quantifying the expression levels of thousands of  
801 combinations of enhancer and promoter sequences.

802

803

804

### 805 **Selection of enhancer and promoter sequences for ExP STARR-seq**

806

807 To explore the compatibility of human enhancers and promoters, we selected 1000 promoter  
808 and 1000 enhancer sequences, including sequences from the human genome spanning a range  
809 of expression or activity levels, and dinucleotide shuffled controls. Based on available lengths of  
810 oligonucleotide pool synthesis, each sequence was 264bp.

811

812 Promoters: We selected the 1000 promoter sequences to include:

- 813 • 65 genes whose enhancers have previously been studied in CRISPR experiments in  
814 K562 cells<sup>22</sup>
- 815 • 715 genes sampled to span a range of potential promoter activities, including the 200  
816 most highly expressed genes in K562 cells, based on CAGE signal at their TSS<sup>21</sup> and a random  
817 sample of 515 other expressed genes (>1 TPM in RNA-seq data<sup>27</sup>).
- 818 • 20 genes that are not expressed or lowly expressed in K562 cells (<1 TPM), and that are  
819 expressed in both GM12878 and HCT-116 cells (in the top 70% of genes by TPM based on  
820 RNA-seq<sup>21</sup>).
- 821 • 100 random genomic sequences as negative controls (+ strand)
- 822 • 100 dinucleotide shuffles of these random genomic sequences

823

824 For the selected genes, we synthesized a 264-bp sequence including approximately 244 bp  
825 upstream and 20 bp downstream of the TSS. Here, we defined the TSS as the center of the 10-  
826 bp window with the most CAGE 5' read counts within 1 Kbp of a RefSeq TSS. For lowly  
827 expressed genes (which lack clear CAGE signal), we used the hg19 RefSeq-annotated TSS.  
828 For genes studied in Fulco *et al.* 2019, we adjusted the assigned 10bp TSS window by manual  
829 examination of the CAGE if necessary.

830

831 Enhancers: We selected the 1000 enhancer sequences to include:

- 832 • 131 elements previously studied with CRISPR<sup>22</sup>, including (i) all distal elements (i.e., >1  
833 Kb from an annotated TSS) with significant effects in previous CRISPRi tiling screens (activating  
834 or repressive), (ii) all distal elements predicted by the Activity-by-Contact model to regulate one  
835 of the tested genes in K562 cells<sup>22</sup>, and (iii) two promoter elements for PVT1 that also act as  
836 enhancers for MYC<sup>22</sup>. We selected 264-bp regions centered on the overlapping DHS narrow  
837 peak. For the small number of CRISPR elements that did not overlap a narrow peak, we tiled  
838 the corresponding element with 264-bp windows overlapping by 50 bp.
- 839 • 200 DNase peaks with the strongest predicted enhancer activity, and 351 other DNase  
840 peaks sampled evenly across the range of predicted enhancer activity. Here, we considered all  
841 distal DHS peaks in K562 cells (DHS narrow peaks<sup>22</sup>) and calculated predicted enhancer  
842 activity as the geometric mean of DNase I hypersensitivity and H3K27ac ChIP-seq read counts

843 in K562 cells in the ~500-bp candidate enhancer regions used by the ABC model in Fulco et al.  
844 2019<sup>22</sup>. Some candidate ABC elements in this set span more than one DHS peak, in which case  
845 we divided the predicted enhancer activity equally among each overlapping peak. We  
846 downloaded introns from the UCSC Genome Browser 'refGene' track (version 2017-06-24), and  
847 removed any peaks overlapping an annotated splice donor or acceptor site. We then selected  
848 264-bp regions centered on the remaining DHS narrow peaks.

- 849 • 100 random genomic sequences as negative controls
- 850 • 100 dinucleotide shuffles of these random genomic sequences

851

852 All enhancer sequences were taken from the hg19 reference in the + strand direction.

853

854

## 855 **Library Cloning**

856

857 We ordered 264bp sequences in an oligo array format from Twist Bioscience with separate  
858 pairs of 18bp adaptors (total length = 300bp) for enhancers (5' = GCTAACTTCTACCCATGC, 3'  
859 = GCAAGTTAAGTAGGTCGT) and promoters (5' = TCATGTGGGACATCAAGC, 3' =  
860 GCATAGTGAGTCCACCTT). We then PCR amplified enhancers and promoters separately  
861 from the same array using Q5 high-fidelity DNA polymerase (NEB M0492). We amplified  
862 enhancers in four 50uL PCR reactions (98°C for 30 seconds; 15 cycles of 98°C for 15 seconds,  
863 61°C for 15 seconds, and 72°C for 20 seconds) using primers (forward:

864 TAGATTGATCTAGAGCATGCANNNNNNNNNNNNNNNNNNGAGTACTGGTATGTTTCAGCTAACT  
865 TCTACCCATGC, reverse:

866 TCGAAGCGGCCGCGCCGAATTCGTCATTCCATGGCATCTCACGACCTACTTAACTTGC)

867 which add an additional 17bp on either side, a 16bp N-mer plasmid barcode upstream, and  
868 homology arms for Gibson assembly on either side of the enhancer sequence. We amplified

869 promoters in four 50uL PCR reactions (98°C for 30 seconds; 4 cycles of 98°C for 15 seconds,  
870 61°C for 15 seconds, 72°C for 20 seconds; 11 cycles of 98°C for 15 seconds and 72°C for 20

871 seconds) using primers (forward: CTCTGGCCTAACTGGCCGGTACGAGTGAGCTCTCGTTCA

872 TCATGTGGGACATCAAGC, reverse:

873 CCCAGTGCCTCACGACCGGGCCTGGTAGCAAGCTTAGATAAGGTGGACTCACTATGC)

874 which add an additional 17bp and homology arms for Gibson assembly on either side of the  
875 promoter sequence. We purified the PCR products using 0.8X volume of AMPure XP beads

876 (Beckman Coulter, A63881) and pooled the reactions together while keeping enhancers and  
877 promoters separate.

878

879 We digested the human STARR-seq screening vector (hSTARR-seq\_SCP1 vector\_blocking 4,  
880 Addgene #99319) with both Thermo SgrDI and BshTI (AgeI) (replaced with enhancer

881 sequence), then NEB KpnI and Apal (replaced with promoter sequence), with purification using  
882 0.8X volume AMPure XP after each digestion. We then recombined 500ng of this digestion

883 (including ~4.4kb of backbone vector and 250bp of filler sequence including a spliced region  
884 and truncated GFP ORF) with 150ng of both the purified enhancer and promoter products using

885 Gibson assembly (NEB, E2611) for 1 hour at 50°C in a 40uL reaction and purified the reaction  
886 using 1X volume AMPure XP with 3 total ethanol washes.

887

888 We electroporated the assembled libraries into Lucigen Endura Electrocompetent cells (60242)  
889 using 0.1cm cuvettes (BioRad) using the Gene Pulser Xcell Microbial System (BioRad) (10 uF,  
890 600 Ohms, 1800 Volts) following the manufacturer's recommendations. We expanded the  
891 transformations for 12 hours in LB with carbenicillin while also estimating the number of  
892 transformed colonies by plating a serial dilution of transformation mixture as previously  
893 described<sup>47</sup>. We midiprep the expanded transformations with ZymoPURE II Plasmid  
894 Midiprep (D4200).

895  
896

### 897 **Building the Barcode-Promoter Dictionary**

898

899 We introduced a unique 16-bp "plasmid barcode" adjacent to the enhancer sequence to allow  
900 us to determine from which promoter each transcript originated, which, together with the self-  
901 transcribed enhancer, allow us to map each transcript to a promoter-enhancer pair.

902

903 To build the map from 16-bp plasmid barcodes to promoters we PCR-amplified a fragment  
904 containing both the promoter and plasmid barcode from the plasmid library (98°C for 1 minute  
905 and 16 cycles of 98°C for 10 seconds, 66°C for 15 seconds, and 72°C for 25 seconds,  
906 ExP\_P1\_fwd\_I2: AATGATACGGCGACCACCGAGATCTACAC[index-  
907 2]GGGAGGTATTGGACAGGC, ExP\_P3\_rev:  
908 CAAGCAGAAGACGGCATAACGAGATGCATGGGTAGAAGTTAGCTGAAC) and sequenced the  
909 promoter position with paired-end reads (using custom sequencing primers  
910 ExP\_P1\_fwd\_seq\_R1: GAGTGAGCTCTCGTTCATCATGTGGGACATCAAGC,  
911 ExP\_P2\_rev\_seq\_R2: TGGTAGCAAGCTTAGATAAGGTGGACTCACTATGC) and the plasmid  
912 barcode with an index read (using custom sequencing primer ExP\_fwd\_BC\_seq:  
913 GTCCCAATTCTTGTGAATTAGATTGATCTAGAGCATGCA). We mapped these sequences to  
914 a specially constructed index of the promoter sequences using bowtie2 (X: -q --met-stderr --  
915 maxins 2000 -p 4 --no-mixed --dovetail --fast). We dropped any BC-promoter pairs with  
916 singleton reads, then removed ambiguous pairings (more than one promoter for the same BC),  
917 and finally thresholded pairs with at least 5 reads to build the Barcode-Promoter dictionary.

918

919

### 920 **Cell Culture**

921

922 We maintained cells at a density between 100,000 and 1,000,000 cells per ml in RPMI-1640  
923 (Thermo Fisher Scientific) with 10% heat-inactivated FBS, 2 mM L-glutamine and 100 units per  
924 ml streptomycin and 100 mg ml<sup>-1</sup> penicillin by diluting cells 1:8 in fresh medium every 3 days.  
925 Cell lines were regularly tested for mycoplasma.

926

927

### 928 **Library Transfection**

929

930 We nucleofected 10 million K562 cells with 15µg of the ExP plasmid library in 100µL cuvettes  
931 with the Lonza 4D-Nucleofector using settings and protocols specified by the manufacturer for  
932 K562 cells (T-016). We pooled 5 nucleofections together during recovery to form 50 million cell  
933 biological transfection replicates and generated 4 replicates for a total of 200 million total cells.  
934 After 24 hours, we harvested the cells in Qiagen buffer RLT (79216) and proceeded with  
935 STARR-seq library preparation.

936

937

938



## 939 **STARR-seq Library Preparation**

940  
941 We proceeded with STARR-seq library preparation using an adapted protocol from Arnold  
942 2013<sup>4</sup>. We split the 50 million-cell transfection replicates in half and extracted total RNA using 3  
943 Qiagen RNeasy mini columns (74134), performing the on-column DNase step. We isolated  
944 polyA+ mRNA using the Qiagen Oligotex mRNA kit for the 1000 x 1000 ExP dataset (note this  
945 kit has been discontinued, we now use the Poly(A)Purist MAG kit from Thermo Fisher Scientific,  
946 AM1922). Following mRNA elution, we treated with TURBO DNase (Thermo Fisher Scientific,  
947 AM2238) in 100uL reactions at 37°C for 30 minutes, then added an additional 2uL of TURBO  
948 DNase and incubated at 37°C for 15 minutes. We purified the RNA following DNA digestion with  
949 Zymo RNA Clean & Concentrator 5 (R1013). We reverse transcribed the polyA+ mRNA using  
950 Thermo SuperScriptIV using the STARR\_RT primer (CAAACATCAATGTATCTTATCATG) in  
951 20uL reactions according to manufacturer's recommendations. We included 1uL of ribonuclease  
952 inhibitor RNaseOUT (Invitrogen, 10777019). Following reverse transcription, we added 1uL of  
953 RNaseH (Thermo Fisher Scientific, EN0201) and incubated at 37°C for 20 minutes. We purified  
954 the cDNA with 1.8X volume of AMPure XP beads. We next selectively amplified the reporter  
955 transcript using intron-spanning junction primers with Q5 polymerase in 50uL reactions (98°C  
956 for 45 seconds and 15 cycles of 98°C for 15 seconds, 65°C for 30 seconds, and 72°C for 70  
957 seconds, jPCR\_fwd: TCGTGAGGCACTGGGCAG\*G\*T\*G\*T\*C, jPCR\_rev:  
958 CTTATCATGTCTGCTCGA\*A\*G\*C, \* = phosphorothioate bonds). Following purifications with  
959 0.8X volume of AMPure XP beads, we performed a test final sequencing-ready PCR with a  
960 dilution of the junction PCR product to determine the optimal cycle number, then proceeded with  
961 the final PCR using Q5 polymerase in 50uL reactions (98°C for 45 seconds and ~9 cycles of  
962 98°C for 10 seconds, 65°C for 30 seconds, and 72°C for 30 seconds, ExP\_GFP\_fwd\_I2:  
963 AATGATACGGCGACCACCGAGATCTACAC[index-2]GGCTTAAGCATGGCTAGCAAAG,  
964 ExP\_P4\_rev: CAAGCAGAAGACGGCATACGAGATTCATTCCATGGCATCTCACG. We purified  
965 the final libraries with 2 rounds of 0.8X volume of SPRISelect (Beckman Coulter, B23318).

966  
967

## 968 **Alignment and counting of STARR-seq data**

969  
970 To characterize activity in the STARR-seq assay, we define "STARR-seq expression" for a  
971 given plasmid (corresponding to a particular promoter, enhancer, and plasmid barcode) as the  
972 expression of the reporter RNA transcript normalized to the abundance of that plasmid in the  
973 input DNA pool.

974  
975 To quantify STARR-seq expression, we sequenced the library of RNA transcripts produced from  
976 replicate transfections (described above) along with the DNA input with paired-end reads (using  
977 custom sequencing primers ExP\_P3\_fwd\_seq\_R1:  
978 GAGTACTGGTATGTTTCAGCTAACTTCTACCCATGC, ExP\_P4\_rev\_seq\_R2:  
979 TCATTCCATGGCATCTCACGACCTACTTAATTGC) and the plasmid barcode with an index  
980 read (using custom sequencing primer ExP\_fwd\_BC\_seq:  
981 GTCCCAATTCTTGTGGAATTAGATTGATCTAGAGCATGCA). We aligned reads for both the  
982 RNA and DNA libraries to the designed enhancer sequences using bowtie2 (bowtie2 options: -q  
983 --met 30 --met-stderr --maxins 2000 -p 16 --no-discordant --no-mixed --fast).

984  
985 We counted reads separately from PCR replicates derived from each biological replicate of 50M  
986 transfected cells, and scaled each of the PCR replicates within a biological replicate such that  
987 they had the same total normalized counts, equal to the maximum across all PCR replicates.  
988 We combined counts into per-biological replicate counts for further processing. We used the



989 BC-promoter dictionary to identify the promoter associated with each transcript. We used the  
990 same mapping and BC-promoter assignment process for DNA.

991

992 For subsequent analysis, we discarded plasmids that had fewer than 25 DNA reads or fewer  
993 than 1 RNA transcript reads from further processing.

994

995

### 996 **Computing technical reproducibility and influence of plasmid barcode sequences**

997

998 To assess the technical reproducibility of ExP-STARR-seq, we first compared STARR-seq  
999 expression between biological replicate experiments. Specifically, we first combined data from  
1000 biological replicates 1 & 2 and 3 & 4. Next, we correlated  $\log_2(\text{RNA/DNA})$  for these groups  
1001 before (**Fig. 1b**) and after (**Fig. S1e**) averaging across plasmid barcodes corresponding to the  
1002 same enhancer-promoter pair.

1003

1004 We next assessed the variation between plasmids with the same enhancer and promoter  
1005 sequences but different random 16-bp plasmid barcodes, because these 16 nucleotides of  
1006 random sequence might contain transcription factor motifs or other sequences that affect  
1007 STARR-seq expression. To do so, we combined data from all biological replicate experiments  
1008 and created two “virtual replicates” for each enhancer-promoter pair by splitting the  
1009 corresponding plasmid barcodes into two groups. For example, an enhancer-promoter pair with  
1010 6 plasmid barcodes was split into 2 virtual replicates each with 3 barcodes). We averaged  $\log_2$   
1011 STARR-seq expression within enhancer-promoter pairs (across different barcodes) and  
1012 correlated these virtual replicates. We compared versions of this analysis for increasing  
1013 thresholds on the minimum number of barcodes in each virtual replicate (**Fig. S1c,d**).

1014

1015

### 1016 **Estimating enhancer and promoter activities — naïve averaging approach**

1017

1018 We sought to compare the intrinsic activities of different enhancer and promoter sequences in  
1019 ExP STARR-seq — that is, the contribution of a given enhancer or promoter sequence to  
1020 STARR-seq expression, relative to other sequences. We estimated enhancer activity and  
1021 promoter activity in two ways: by a simple averaging method, and by fitting a multiplicative  
1022 Poisson count model (see next section).

1023

1024 As a first approach to estimate promoter activity, we calculated, for each promoter sequence,  
1025 the average  $\log_2$  STARR-seq expression when that promoter is paired with random genomic  
1026 sequences in the enhancer position (**Fig. 1c**). This quantity represents the “basal” or  
1027 “autonomous” expression level of the promoter, in the absence of a strong activating sequence  
1028 in the enhancer position.

1029

1030 As a first approach to estimate enhancer activity, we calculated, for each enhancer sequence,  
1031 the average  $\log_2$  STARR-seq expression of all pairs including that enhancer sequence (**Fig. 1d**).

1032

1033 As noted above, we fit this model on the set of plasmids with at least 25 DNA reads, and at least  
1034 1 RNA read. In addition, to reduce noise in our promoter and enhancer activity estimates, we  
1035 required at least two separate plasmid barcodes per promoter-enhancer pair. These filters  
1036 resulted in 604,268 promoter-enhancer pairs across 4,512,907 total unique plasmids (~ 7.5  
1037 plasmids per pair) that were used to estimate promoter and enhancer activity.

1038

1039 In practice, this averaging method of calculating enhancer and promoter activity was inaccurate  
1040 and biased, for several reasons. First, the averaging method does not consider the variance  
1041 introduced by sampling & counting noise in sequencing, which is significant because many  
1042 promoter-enhancer pairs have low RNA read counts. Second, the averaging method does not  
1043 account for differences introduced due to missing data. In the 1000 enhancer x 1000 promoter  
1044 data matrix, many entries are missing either due to low RNA counts (resulting from counting and  
1045 sampling noise, or low expression) or due to low DNA counts (resulting from variation  
1046 introduced in cloning the plasmid library). As a result of these factors, the averaging method  
1047 produces biased (inflated) estimates of activity for weaker enhancer and promoter sequences  
1048 because the expression of plasmids containing these sequences is more likely to drop below  
1049 the threshold of detection given our sequencing depth (**Fig. S2c-d**).

1050  
1051 Because this model explained the data well, we used this same model to estimate intrinsic  
1052 enhancer and promoter activity.

1053  
1054

### 1055 **Estimating intrinsic enhancer and promoter activities — multiplicative model**

1056  
1057 We fit a count-based Poisson model to address the limitations of using a simple averaging  
1058 approach to estimate intrinsic enhancer and promoter activities (see previous section), and to  
1059 quantify the extent to which the ExP STARR-seq data can be explained by a simple  
1060 multiplicative function of intrinsic enhancer and promoter activities. In this multiplicative model,  
1061 all enhancers are assumed to activate all promoters by the same fold-change, without  
1062 enhancer-promoter interaction terms.

1063  
1064 Specifically, we estimate enhancer and promoter activities from ExP STARR-seq data by fitting  
1065 the observed RNA read counts to a multiplicative function of observed DNA input read counts,  
1066 intrinsic enhancer activity, and intrinsic promoter activity:

$$1067 \quad \quad \quad RNA \sim \text{Poisson}(k \times DNA \times P \times E),$$

1069  
1070 In this formula,  $P$  is the intrinsic promoter activity of promoter sequence  $p$ ,  $E$  is intrinsic  
1071 enhancer activity of enhancer sequence  $e$ , and  $k$  is a global scaling/intercept term that accounts  
1072 for factors that control the relative counts of DNA and RNA such as sequencing depth.

1073  
1074 We fit these parameters using block coordinate descent on the negative log-likelihood of the  
1075 distribution above, initially fixing  $k=0$ , then alternatively optimizing (i) promoter activities while  
1076 holding enhancer activities constant, and (ii) enhancer activities while holding promoter activities  
1077 constant.

1078  
1079 We then re-normalized enhancer activities and promoter activities by the mean activity of  
1080 random genomic sequences, and adjusted the scaling factor  $k$  accordingly.

1081  
1082 In practice, this model produces similar estimates to simply taking the mean value of an  
1083 enhancer sequence across all promoters, and vice versa, but accounts for missing data points  
1084 in the 1000x1000 matrix, and provides a more robust estimate for very weak enhancers or  
1085 promoters, which produce relatively little RNA and are therefore difficult to measure in this  
1086 STARR-seq experiment except when paired with a strong element in the other group (**Fig. S2c-**  
1087 **d**).

1088  
1089

## 1090 **Computing and clustering residuals from the multiplicative model:**

1091

1092 We explored whether enhancer-promoter compatibility could explain variation in STARR-seq  
1093 expression beyond that explained by the multiplicative model. To do so, we looked for shared  
1094 behaviors between groups of promoters and enhancers by clustering them according to their  
1095 residual error from the Poisson model described above.

1096

1097 For each enhancer-promoter pair, we used the Poisson model above to compute predicted RNA  
1098 given the input DNA counts and estimates of intrinsic enhancer and promoter activities. We  
1099 then compute a transformed residual as

1100

$$\log_2(\text{predicted RNA} + \text{pseudocount}) - \log_2(\text{observed RNA} + \text{pseudocount}),$$

1102

1103 where pseudocount = 10 to stabilize variance of the estimates across the range of values for  
1104 RNA<sup>48</sup>. We filtered to all enhancer-promoter pairs with at least two barcodes, and calculated the  
1105 mean of the residuals across barcodes to form a (sparse) 1000x1000 matrix of residuals  
1106 indexed by promoter and enhancers.

1107

1108 We clustered this matrix independently along rows and columns (treating missing pairs as  
1109 having a residual of 0) using K-means with 3 clusters, labeling the clusters as 0, 1, and 2 such  
1110 that they had increasing mean activity estimates in the Poisson model. One cluster each of  
1111 enhancers and promoters (E0 and P0) contained sequences that were missing many data  
1112 points due to their weaker activity leading to dropout due to low RNA expression. The sparsity of  
1113 data for the E0 and P0 clusters prevented accurate characterization of compatibility, and so we  
1114 excluded these clusters from subsequent analysis.

1115

1116

## 1117 **Assessing reproducibility of the clusters:**

1118

1119 We evaluated whether the clustering we observed in the residuals was a general trend of the  
1120 data, or an artifact of a few promoters or enhancers. To test this possibility, we randomly  
1121 downsampled the residual matrix to 25% of promoters and 25% of enhancers (6.25% of the  
1122 total data) 100 times, and clustered the subsets. We found that the original (full-data) cluster  
1123 identity of a promoter or enhancer predicted the downsampled cluster with greater than 80%  
1124 accuracy (**Fig. S3g**).

1125

1126

## 1127 **Estimating enhancer activity with specific promoter classes, and promoter 1128 responsiveness to specific enhancer classes:**

1129

1130 We evaluated whether certain promoters were more responsive when paired with different  
1131 enhancer classes, and whether certain enhancers had more activity when paired with promoters  
1132 from different classes (**Fig. 3c,d**).

1133

1134 To explore differences in enhancer activity when paired with different promoter classes, we fit  
1135 the Poisson model (described above) separately to two different subsets of the data: (i) all  
1136 enhancer sequences paired with P1 or genomic background promoter sequences (yielding an  
1137 estimate of the activity of an enhancer sequence on a P1 promoter), and (ii) all enhancer  
1138 sequences paired with P2 or genomic background promoter sequences (yielding an estimate of  
1139 the activity of an enhancer sequence on a P2 promoter).

1140

1141 Similarly, to estimate promoter responsiveness to either E1 or E2 enhancers, we fit the Poisson  
1142 model to the subsets: (iii) all promoters paired with E1 or genomic background enhancer  
1143 sequences (yielding an estimate of the responsiveness of a promoter sequence to E1  
1144 enhancers), and (iv) all promoters paired with E2 or genomic background enhancer sequences  
1145 (yielding an estimate of the responsiveness of a promoter sequence to E2 enhancers).

1146  
1147 We used the genomic background promoter sequences to set a common baseline.

### 1148 1149 1150 **Annotating enhancer and promoter sequences with genomic features and sequence** 1151 **motifs**

1152  
1153 To annotate enhancer and promoter sequences with features of transcription factor (TF) binding  
1154 of the corresponding genomic elements, we downloaded list of Human TF ChIP-seq  
1155 narrowpeak files from the ENCODE Project<sup>21</sup>, and annotated each enhancer or promoter  
1156 sequence with the maximum signalValue column for any overlapping peak (or 0 signal, for no  
1157 overlap). We then compared the fold-change in signal between classes of sequences (**Fig. 4d**,  
1158 **Fig. S5a, Table S9**).

1159  
1160 To annotate enhancer and promoter sequences with transcription factor motifs, we used FIMO<sup>49</sup>  
1161 (default parameters, including  $p$ -value threshold of  $10^{-6}$ ) to identify matches for HOCOMOCO  
1162 v11 CORE motifs<sup>50</sup>. We then compared the fold-change in motif counts between classes of  
1163 sequences (**Fig. S6f, Table S5, Table S7**).

1164  
1165 For comparing features between E1 and E2 enhancers, we compared motif, ChIP-seq, and  
1166 other features between the E1 and E2 enhancer sequences that overlapped the summit of a  
1167 DNase peak.

1168  
1169 For analyzing the proportion of P2 promoters bound by various factors, we defined "strongly  
1170 bound" as having ChIP-seq signal greater than 20% of maximum ChIP-seq signal among P1  
1171 and P2 promoters.

### 1172 1173 1174 **Comparison of CRISPR-derived regulatory elements for P1 vs P2 promoters**

1175  
1176 To compare the number and effect sizes of genomic regulatory elements for P1 and P2  
1177 promoters, we analyzed CRISPRi tiling screens from previous studies that perturbed all DNase  
1178 accessible sites around selected genes<sup>22,26,27</sup>. We counted the number of activating distal  
1179 regulatory elements — *i.e.*, distal, non-promoter DNase accessible sites whose perturbation led  
1180 to a significant reduction in gene expression (**Fig. 4c**). We also compared the effect sizes on  
1181 gene expression for these same activating distal regulatory elements (**Fig. S5c, Fig. S6b**).

### 1182 1183 1184 **Luciferase assays**

1185  
1186 We tested the ability of each of 7 large *MYC* enhancer fragments to activate the promoters of 3  
1187 genes in the *MYC* locus — *MYC*, *PVT1*, and *CCDC26* — using a classic plasmid luciferase-  
1188 based enhancer assay. The 7 *MYC* enhancers were defined as the 1.0-2.2 kb sequences  
1189 identified in our previous *MYC* proliferation-based CRISPRi screen<sup>27</sup>, and a 1 kb bacterial  
1190 plasmid sequence was used as a negative control sequence. We cloned promoter fragments  
1191 into plasmids in combination with each of these sequences. The promoter fragments

1192 corresponded to the dominant transcription start site of each gene in K562 cells (as determined  
1193 by CAGE). For each of *PVT1* and *CCDC26* — which do not appear to be regulated by most of  
1194 the 7 *MYC* enhancers in the genome — we cloned two promoter fragments of different lengths  
1195 to determine if nearby sequences might encode biochemical specificity. We designed an  
1196 insertion site ~1 kb upstream of the promoter in the plasmid for inserting each enhancer  
1197 sequence (**Fig. S2g**), and we flanked this region with polyadenylation signals in either direction  
1198 to avoid measuring luciferase activity driven from transcripts initiating from the enhancer  
1199 elements themselves. Luciferase assays using the Dual-Luciferase Reporter Assay (Promega)  
1200 were performed as previously described<sup>27</sup> in biological triplicate. For each experiment, we  
1201 calculated the fold-change in luciferase signal (Firefly / Renilla) for enhancer versus negative  
1202 control (**Fig. S2i**).

1203  
1204

### 1205 **Assessing the cell-type specificity of E1 and E2 enhancers**

1206

1207 We tested whether E1 and E2 enhancer sequences from ExP STARR-seq overlapped elements  
1208 predicted to act as enhancers by the ABC model in K562 cells or in 128 other cell types and  
1209 tissues. To do so, we intersected the E1 and E2 enhancer sequences with the ~200-bp regions  
1210 predicted by the ABC model to act as enhancers for at least 1 nearby expressed gene, as  
1211 previously defined<sup>51</sup>. The ABC enhancer-gene predictions from this previous study<sup>51</sup> are  
1212 available at <https://www.engreitzlab.org/resources/>.

1213

1214

### 1215 **Aligning promoters by transcription start site**

1216

1217 For each 264-bp promoter sequence, we defined the primary transcription start site (TSS) as  
1218 the nucleotide with the highest stranded 5' signal in GRO-Cap data in K562 cells  
1219 (GSM1480321)<sup>52</sup>. This primary TSS position was used for plotting genomic signals relative to  
1220 TSS and in analyses of motif positioning (e.g., for GABPA and YY1).

1221

1222

### 1223 **Analysis of motif position relative to TSS**

1224

1225 We used FIMO<sup>49</sup> to scan for HOCOMOCO motifs in promoters including for GABPA  
1226 (GABPA\_HUMAN.H11MO.0.A), YY1 (YY1\_HUMAN.H11MO.0.A), and the TATA box  
1227 (TBP\_HUMAN.H11MO.0.A). We reported positional preferences as the distance between the  
1228 primary transcription start site from GRO-cap (see above) and the center of the motif. For  
1229 example, GABPA, the most common position was -10 relative to the TSS (*i.e.* with the second  
1230 'G' in the core 'GGAA' motif located at position -10).

1231

1232

### 1233 **Hybrid selection STARR-seq (HS-STARR-seq) to measure enhancer activity for millions 1234 of genomic sequences**

1235

1236 We conducted two STARR-seq experiments to measure the enhancer activity of millions of long  
1237 genomic sequences tiling across human enhancer and promoter sequences. To generate these  
1238 tiling sequences, we used a hybrid selection strategy, similar to previous approaches<sup>53</sup>.  
1239 Specifically, we purified genomic DNA from K562 cells, tagged DNA using Tn5 and gel size  
1240 selection to a size range of approximately 300-700 bp (**Fig. S8c**), and conducted hybrid  
1241 selection using RNA probes as previously described<sup>54</sup> targeting either (i) all gene promoters  
1242 ("HS promoter pool") or (ii) all accessible elements ("HS accessible element pool") in K562 cells



1243 (see **Table S10** and **Table S11** for probe sequences). We amplified these sequences using  
1244 primers including a UMI (CapStarrFa\_N10 primer:  
1245 tagatTGAtCTAGAGCATGCACCGGCAAGCAGAAGACGGCATACGAGATNNNNNNNNNNATG  
1246 TCTCGTGGGCTCGGAGATGT and CapStarrR primer:  
1247 CGAAGCGGCCGGCCGAATTCGTTCGATCGTTCGGCAGCGTCAGATGTG) and cloned these  
1248 selected sequences into the hSTARR-seq\_ori vector<sup>8</sup>, which uses the bacterial origin of  
1249 replication (ORI) sequence as the promoter for the reporter transcript, by Gibson assembly. In  
1250 the final HS promoter and accessible element Pools, 9% and 12% of fragments mapped to their  
1251 intended targets, respectively, and each element was tiled by a median of 20 and 55  
1252 sequences. We conducted the rest of the STARR-seq experiment as described above,  
1253 transfecting 50 million cells per replicate for each of 4 replicates.

1254  
1255 We sequenced the input DNA libraries to a depth of 880 million and 810 million reads (promoter  
1256 and accessible element pools, respectively), and the RNA libraries to a depth of 1.1 billion reads  
1257 (both pools). We aligned reads to the hg19 genome using bowtie2 (options: -q --met-stderr --  
1258 maxins 1000 -p 4 --no-discordant --no-mixed). We discarded fragments with fewer than 25  
1259 aligned DNA reads. Biological replicates were highly correlated ( $R^2 = 0.92$  and  $0.91$  for  
1260 promoter and accessible element pools) (**Fig. S8b**).

1261  
1262 We analyzed this data by computing a  $\log_2$  activity per fragment equal to the  $\log_2(\text{RNA} / \text{DNA})$ .  
1263 and correcting for a fragment-length bias. We noted that STARR-seq expression was highly  
1264 inversely correlated with the length of the enhancer sequence, even among random genomic  
1265 fragments that did not overlap putative regulatory elements, which could result from biases in  
1266 library preparation and sequencing. To adjust for this, we fit a linear regression (separately for  
1267 the two pools) and subtracted this regression from the  $\log_2(\text{RNA} / \text{DNA})$  activity to give a bias-  
1268 corrected activity. We then correlated motifs with bias-corrected activity. To estimate enhancer  
1269 activity of promoters from the ExP, we found HS-STARR-seq fragments that overlapped at least  
1270 90% of an ExP promoter and averaged their activity scores.  
1271



1272 **References**

- 1273 1. Banerji, J., Rusconi, S. & Schaffner, W. Expression of a beta-globin gene is enhanced by  
1274 remote SV40 DNA sequences. *Cell* **27**, 299–308 (1981).
- 1275 2. Banerji, J., Olson, L. & Schaffner, W. A lymphocyte-specific cellular enhancer is located  
1276 downstream of the joining region in immunoglobulin heavy chain genes. *Cell* **33**, 729–740  
1277 (1983).
- 1278 3. Melnikov, A. *et al.* Systematic dissection and optimization of inducible enhancers in human  
1279 cells using a massively parallel reporter assay. *Nat. Biotechnol.* **30**, 271–277 (2012).
- 1280 4. Arnold, C. D. *et al.* Genome-wide quantitative enhancer activity maps identified by STARR-  
1281 seq. *Science* **339**, 1074–1077 (2013).
- 1282 5. Kermekchiev, M., Pettersson, M., Matthias, P. & Schaffner, W. Every enhancer works with  
1283 every promoter for all the combinations tested: could new regulatory pathways evolve by  
1284 enhancer shuffling? *Gene Expr.* **1**, 71–81 (1991).
- 1285 6. Tewhey, R. *et al.* Direct Identification of Hundreds of Expression-Modulating Variants using  
1286 a Multiplexed Reporter Assay. *Cell* **172**, 1132–1134 (2018).
- 1287 7. Klein, J. C. *et al.* A systematic evaluation of the design and context dependencies of  
1288 massively parallel reporter assays. *Nat. Methods* **17**, 1083–1091 (2020).
- 1289 8. Muerdter, F. *et al.* Resolving systematic errors in widely used enhancer activity assays in  
1290 human cells. *Nat. Methods* **15**, 141–149 (2018).
- 1291 9. Nguyen, T. A. *et al.* High-throughput functional comparison of promoter and enhancer  
1292 activities. *Genome Res.* **26**, 1023–1033 (2016).
- 1293 10. Emami, K. H., Navarre, W. W. & Smale, S. T. Core promoter specificities of the Sp1 and  
1294 VP16 transcriptional activation domains. *Mol. Cell. Biol.* **15**, 5906–5916 (1995).
- 1295 11. Ohtsuki, S., Levine, M. & Cai, H. N. Different core promoters possess distinct regulatory  
1296 activities in the *Drosophila* embryo. *Genes Dev.* **12**, 547–556 (1998).
- 1297 12. Emami, K. H., Jain, A. & Smale, S. T. Mechanism of synergy between TATA and initiator:  
1298 synergistic binding of TFIID following a putative TFIIA-induced isomerization. *Genes Dev.*  
1299 **11**, 3007–3019 (1997).
- 1300 13. Butler, J. E. F. Enhancer-promoter specificity mediated by DPE or TATA core promoter  
1301 motifs. *Genes & Development* vol. 15 2515–2519 (2001).
- 1302 14. Yean, D. & Gralla, J. Transcription reinitiation rate: a special role for the TATA box.  
1303 *Molecular and Cellular Biology* vol. 17 3809–3816 (1997).
- 1304 15. Wefald, F. C., Devlin, B. H. & Williams, R. S. Functional heterogeneity of mammalian  
1305 TATA-box sequences revealed by interaction with a cell-specific enhancer. *Nature* **344**,  
1306 260–262 (1990).
- 1307 16. Zabidi, M. A., Arnold, C. D., Schernhuber, K. & Pagani, M. Enhancer–core-promoter  
1308 specificity separates developmental and housekeeping gene regulation. *Nature* (2015).
- 1309 17. Arnold, C. D. *et al.* Genome-wide assessment of sequence-intrinsic enhancer  
1310 responsiveness at single-base-pair resolution. *Nat. Biotechnol.* **35**, 136–144 (2017).
- 1311 18. Haberle, V. *et al.* Transcriptional cofactors display specificity for distinct types of core  
1312 promoters. *Nature* **570**, 122–126 (2019).
- 1313 19. van Arensbergen, J., van Steensel, B. & Bussemaker, H. J. In search of the determinants of

- 1314 enhancer–promoter interaction specificity. *Trends in Cell Biology* vol. 24 695–702 (2014).
- 1315 20. Li, X. & Noll, M. Compatibility between enhancers and promoters determines the  
1316 transcriptional specificity of gooseberry and gooseberry neuro in the *Drosophila* embryo.  
1317 *EMBO J.* **13**, 400–406 (1994).
- 1318 21. ENCODE Project Consortium. An integrated encyclopedia of DNA elements in the human  
1319 genome. *Nature* **489**, 57–74 (2012).
- 1320 22. Fulco, C. P. *et al.* Activity-by-contact model of enhancer-promoter regulation from  
1321 thousands of CRISPR perturbations. *Nat. Genet.* **51**, 1664–1669 (2019).
- 1322 23. Wall, L., deBoer, E. & Grosveld, F. The human beta-globin gene 3' enhancer contains  
1323 multiple binding sites for an erythroid-specific protein. *Genes Dev.* **2**, 1089–1100 (1988).
- 1324 24. Tuan, D. Y., Solomon, W. B., London, I. M. & Lee, D. P. An erythroid-specific,  
1325 developmental-stage-independent enhancer far upstream of the human 'beta-like globin'  
1326 genes. *Proc. Natl. Acad. Sci. U. S. A.* **86**, 2554–2558 (1989).
- 1327 25. Thakore, P. I. *et al.* Highly specific epigenome editing by CRISPR-Cas9 repressors for  
1328 silencing of distal regulatory elements. *Nat. Methods* **12**, 1143–1149 (2015).
- 1329 26. Klann, T. S. *et al.* CRISPR-Cas9 epigenome editing enables high-throughput screening for  
1330 functional regulatory elements in the human genome. *Nat. Biotechnol.* **35**, 561–568 (2017).
- 1331 27. Fulco, C. P. *et al.* Systematic mapping of functional enhancer-promoter connections with  
1332 CRISPR interference. *Science* **354**, 769–773 (2016).
- 1333 28. Liu, Y. *et al.* Functional assessment of human enhancer activities using whole-genome  
1334 STARR-sequencing. *Genome Biol.* **18**, 219 (2017).
- 1335 29. Haberle, V. & Stark, A. Eukaryotic core promoters and the functional basis of transcription  
1336 initiation. *Nat. Rev. Mol. Cell Biol.* **19**, 621–637 (2018).
- 1337 30. Lenhard, B., Sandelin, A. & Carninci, P. Metazoan promoters: emerging characteristics and  
1338 insights into transcriptional regulation. *Nat. Rev. Genet.* **13**, 233–245 (2012).
- 1339 31. Fan, K., Moore, J. E., Zhang, X.-O. & Weng, Z. Genetic and epigenetic features of  
1340 promoters with ubiquitous chromatin accessibility support ubiquitous transcription of cell-  
1341 essential genes. *Nucleic Acids Res.* **49**, 5705–5725 (2021).
- 1342 32. Xi, H. *et al.* Identification and characterization of cell type-specific and ubiquitous chromatin  
1343 regulatory structures in the human genome. *PLoS Genet.* **3**, e136 (2007).
- 1344 33. Landolin, J. M. *et al.* Sequence features that drive human promoter function and tissue  
1345 specificity. *Genome Res.* **20**, 890–898 (2010).
- 1346 34. Weingarten-Gabbay, S. *et al.* Systematic interrogation of human promoters. *Genome Res.*  
1347 **29**, 171–183 (2019).
- 1348 35. Sahu, B., Hartonen, T., Pihlajamaa, P., Wei, B. & Dave, K. Sequence determinants of  
1349 human gene regulatory elements. *bioRxiv* (2021).
- 1350 36. Yu, M. *et al.* GA-binding protein-dependent transcription initiator elements. Effect of helical  
1351 spacing between polyomavirus enhancer a factor 3(PEA3)/Ets-binding sites on initiator  
1352 activity. *J. Biol. Chem.* **272**, 29060–29067 (1997).
- 1353 37. Curina, A. *et al.* High constitutive activity of a broad panel of housekeeping and tissue-  
1354 specific cis-regulatory elements depends on a subset of ETS proteins. *Genes Dev.* **31**,  
1355 399–412 (2017).

- 1356 38. van Arensbergen, J. *et al.* Genome-wide mapping of autonomous promoter activity in  
1357 human cells. *Nat. Biotechnol.* **35**, 145–153 (2017).
- 1358 39. Hong, C. K. Y. & Cohen, B. A. Genomic environments scale the activities of diverse core  
1359 promoters. doi:10.1101/2021.03.08.434469.
- 1360 40. Maricque, B. B., Chaudhari, H. G. & Cohen, B. A. A massively parallel reporter assay  
1361 dissects the influence of chromatin structure on cis-regulatory activity. *Nat. Biotechnol.*  
1362 (2018) doi:10.1038/nbt.4285.
- 1363 41. Chiang, C. M. & Roeder, R. G. Cloning of an intrinsic human TFIID subunit that interacts  
1364 with multiple transcriptional activators. *Science* **267**, 531–536 (1995).
- 1365 42. Austen, M., Lüscher, B. & Lüscher-Firzlaff, J. M. Characterization of the transcriptional  
1366 regulator YY1. The bipartite transactivation domain is independent of interaction with the  
1367 TATA box-binding protein, transcription factor IIB, TAFII55, or cAMP-responsive element-  
1368 binding protein (CPB)-binding protein. *J. Biol. Chem.* **272**, 1709–1717 (1997).
- 1369 43. Sucharov, C., Basu, A., Carter, R. S. & Avadhani, N. G. A novel transcriptional initiator  
1370 activity of the GABP factor binding ets sequence repeat from the murine cytochrome c  
1371 oxidase Vb gene. *Gene Expr.* **5**, 93–111 (1995).
- 1372 44. Carter, R. S. & Avadhani, N. G. Cooperative binding of GA-binding protein transcription  
1373 factors to duplicated transcription initiation region repeats of the cytochrome c oxidase  
1374 subunit IV gene. *J. Biol. Chem.* **269**, 4381–4387 (1994).
- 1375 45. Usheva, A. & Shenk, T. YY1 transcriptional initiator: protein interactions and association  
1376 with a DNA site containing unpaired strands. *Proc. Natl. Acad. Sci. U. S. A.* **93**, 13571–  
1377 13576 (1996).
- 1378 46. Larsson, A. J. M. *et al.* Genomic encoding of transcriptional burst kinetics. *Nature* **565**,  
1379 251–254 (2019).
- 1380 47. Wang, T., Lander, E. S. & Sabatini, D. M. Large-Scale Single Guide RNA Library  
1381 Construction and Use for CRISPR-Cas9-Based Genetic Screens. *Cold Spring Harb.*  
1382 *Protoc.* **2016**, db.top086892 (2016).
- 1383 48. Anscombe, F. J. THE TRANSFORMATION OF POISSON, BINOMIAL AND NEGATIVE-  
1384 BINOMIAL DATA. *Biometrika* vol. 35 246–254 (1948).
- 1385 49. Grant, C. E., Bailey, T. L. & Noble, W. S. FIMO: scanning for occurrences of a given motif.  
1386 *Bioinformatics* **27**, 1017–1018 (2011).
- 1387 50. Kulakovskiy, I. V. *et al.* HOCOMOCO: towards a complete collection of transcription factor  
1388 binding models for human and mouse via large-scale ChIP-Seq analysis. *Nucleic Acids*  
1389 *Res.* **46**, D252–D259 (2018).
- 1390 51. Nasser, J. *et al.* Genome-wide enhancer maps link risk variants to disease genes. *Nature*  
1391 **593**, 238–243 (2021).
- 1392 52. Core, L. J. *et al.* Analysis of nascent RNA identifies a unified architecture of initiation  
1393 regions at mammalian promoters and enhancers. *Nat. Genet.* **46**, 1311–1320 (2014).
- 1394 53. Vanhille, L. *et al.* High-throughput and quantitative assessment of enhancer activity in  
1395 mammals by CapStarr-seq. *Nat. Commun.* **6**, 6905 (2015).
- 1396 54. Engreitz, J. M. *et al.* Local regulation of gene expression by lncRNA promoters,  
1397 transcription and splicing. *Nature* **539**, 452–455 (2016).