

1 **Information Content Differentiates Enhancers From Silencers in**
2 **Mouse Photoreceptors**

3 Ryan Z. Friedman¹, David M. Granas¹, Connie A. Myers², Joseph C. Corbo², Barak A. Cohen¹,
4 Michael A. White^{1*}

5

6 ¹Edison Family Center for Genome Sciences and Systems Biology, and Department of
7 Genetics, Washington University School of Medicine, St. Louis, MO, USA

8 ²Department of Pathology and Immunology, Washington University School of Medicine, St.
9 Louis, MO, USA

10 *Corresponding author (mawhite@wustl.edu)

11

12 ORCID IDs: 0000-0001-9013-8676 (RZF), 0000-0002-9323-7140 (JCC), 0000-0002-3350-2715
13 (BAC), 0000-0001-8511-6026 (MAW)

14

15 Abstract

16 Enhancers and silencers often depend on the same transcription factors (TFs) and are conflated
17 in genomic assays of TF binding or chromatin state. To identify sequence features that
18 distinguish enhancers and silencers, we assayed massively parallel reporter libraries of
19 genomic sequences targeted by the photoreceptor TF CRX in mouse retinas. Both enhancers
20 and silencers contain more TF motifs than inactive sequences, but relative to silencers,
21 enhancers contain motifs from a more diverse collection of TFs. We developed a measure of
22 information content that describes the number and diversity of motifs in a sequence and found
23 that, while both enhancers and silencers depend on CRX motifs, enhancers have higher
24 information content. The ability of information content to distinguish enhancers and silencers
25 targeted by the same TF illustrates how motif context determines the activity of *cis*-regulatory
26 sequences.

27 Introduction

28 Active *cis*-regulatory sequences in the genome are characterized by accessible chromatin and
29 specific histone modifications, which reflect the action of DNA-binding transcription factors (TFs)
30 that recognize specific sequence motifs and recruit chromatin modifying enzymes (Klemm et al.,
31 2019). These epigenetic hallmarks of active chromatin are routinely used to train machine
32 learning models that predict *cis*-regulatory sequences, based on the assumption that such
33 epigenetic marks are reliable predictors of genuine *cis*-regulatory sequences (Ernst & Kellis,
34 2012; Ghandi et al., 2014; Hoffman et al., 2012; Kelley et al., 2016; D. Lee et al., 2011; Sethi et
35 al., 2020; Zhou & Troyanskaya, 2015). However, results from functional assays show that many
36 predicted *cis*-regulatory sequences exhibit little or no *cis*-regulatory activity. Typically, 50% or
37 more of predicted *cis*-regulatory sequences fail to drive expression in Massively Parallel
38 Reporter Assays (MPRAs) (ENCODE Project Consortium et al., 2020; Kwasnieski et al., 2014),

39 indicating that an active chromatin state is not sufficient to reliably identify *cis*-regulatory
40 sequences.

41

42 Another challenge is that enhancers and silencers are difficult to distinguish by chromatin
43 accessibility or epigenetic state (Doni Jayavelu et al., 2020; Gisselbrecht et al., 2020; Pang &
44 Snyder, 2020; Petrykowska et al., 2008; Segert et al., 2021), and thus computational predictions
45 of *cis*-regulatory sequences often do not differentiate between enhancers and silencers.
46 Silencers are often enhancers in other cell types (Brand et al., 1987; Doni Jayavelu et al., 2020;
47 Gisselbrecht et al., 2020; Z. Huang et al., 2021; Jiang et al., 1993; Ngan et al., 2020; Pang &
48 Snyder, 2020), reside in open chromatin (Doni Jayavelu et al., 2020; D. Huang et al., 2019; Z.
49 Huang et al., 2021; Pang & Snyder, 2020), sometimes bear epigenetic marks of active
50 enhancers (Fan et al., 2016; Z. Huang et al., 2021), and can be bound by TFs that also act on
51 enhancers in the same cell type (Alexandre & Vincent, 2003; Grass et al., 2003; Z. Huang et al.,
52 2021; Iype et al., 2004; Jiang et al., 1993; Liu et al., 2014; Martínez-Montañés et al., 2013; Peng
53 et al., 2005; Rachmin et al., 2015; Rister et al., 2015; Stampfel et al., 2015; White et al., 2013).
54 As a result, enhancers and silencers share similar sequence features, and understanding how
55 they are distinguished in a particular cell type remains an important challenge (Segert et al.,
56 2021).

57

58 The TF cone-rod homeobox (CRX) controls selective gene expression in a number of different
59 photoreceptor and bipolar cell types in the retina (S. Chen et al., 1997; Freund et al., 1997;
60 Furukawa et al., 1997; Murphy et al., 2019). These cell types derive from the same progenitor
61 cell population (Koike et al., 2007; Wang et al., 2014), but they exhibit divergent, CRX-directed
62 transcriptional programs (Corbo et al., 2010; Hennig et al., 2008; Hughes et al., 2017; Murphy et

63 al., 2019). CRX cooperates with cell type-specific co-factors to selectively activate and repress
64 different genes in different cell types and is required for differentiation of rod and cone
65 photoreceptors (J. Chen et al., 2005; Hao et al., 2012; Hennig et al., 2008; Hsiau et al., 2007;
66 Irie et al., 2015; Kimura et al., 2000; Lerner et al., 2005; Mears et al., 2001; K. P. Mitton et al.,
67 2000; Murphy et al., 2019; Peng et al., 2005; Sanuki et al., 2010; Srinivas et al., 2006).
68 However, the sequence features that define CRX-targeted enhancers versus silencers in the
69 retina are largely unknown.

70

71 We previously found that a significant minority of CRX-bound sequences act as silencers in an
72 MPRA conducted in live mouse retinas (White et al., 2013), and that silencer activity requires
73 CRX (White et al., 2016). Here we extend our analysis by testing thousands of additional
74 candidate *cis*-regulatory sequences. We show that while regions of accessible chromatin and
75 CRX binding exhibit a range of *cis*-regulatory activity, enhancers and silencers contain more TF
76 motifs than inactive sequences, and that enhancers are distinguished from silencers by a higher
77 diversity of TF motifs. We capture the differences between these sequence classes with a new
78 metric, motif information content (Boltzmann entropy), that considers only the number and
79 diversity of TF motifs in a candidate *cis*-regulatory sequence. Our results suggest that CRX-
80 targeted enhancers are defined by a flexible regulatory grammar and demonstrate how
81 differences in motif information content encode functional differences between genomic loci with
82 similar chromatin states.

83

84

85 Results

86 We tested the activities of 4,844 putative CRX-targeted *cis*-regulatory sequences (CRX-targeted
87 sequences) by MPRA in live retinas. The MPRA libraries consist of 164 bp genomic sequences
88 centered on the best match to the CRX position weight matrix (PWM) (J. Lee et al., 2010)
89 whenever a CRX motif is present, and matched sequences in which all CRX motifs were
90 abolished by point mutation (Methods). The MPRA libraries include 3,299 CRX-bound
91 sequences identified by ChIP-seq in the adult retina (Corbo et al., 2010) and 1,545 sequences
92 that do not have measurable CRX binding in the adult retina but reside in accessible chromatin
93 in adult photoreceptors (Hughes et al., 2017) and have the H3K27ac enhancer mark in
94 postnatal day 14 (P14) retina (Ruzycki et al., 2018) (“ATAC-seq peaks”). We split the
95 sequences across two plasmid libraries, each of which contained the same 150 scrambled
96 sequences as internal controls (**Supplementary files 1 and 2**). We cloned sequences upstream
97 of the rod photoreceptor-specific *Rhodopsin* (*Rho*) promoter and a *DsRed* reporter gene,
98 electroporated libraries into explanted mouse retinas at P0 in triplicate, harvested the retinas at
99 P8, and then sequenced the RNA and input DNA plasmid pool. The data is highly reproducible
100 across replicates ($R^2 > 0.96$, **Figure 1—figure supplement 1**). After activity scores were
101 calculated and normalized to the basal *Rho* promoter, the two libraries were well calibrated and
102 merged together (two-sample Kolmogorov-Smirnov test $p = 0.09$, **Figure 1—figure**
103 **supplement 2, Supplementary file 3**, and Methods).

104

105 **Strong enhancers and silencers have high CRX motif content**

106 The *cis*-regulatory activities of CRX-targeted sequences vary widely (**Figure 1a**). We defined
107 enhancers and silencers as those sequences that have statistically significant activity that is at
108 least two-fold above or below the activity of the basal *Rho* promoter (Welch’s t-test, Benjamini-

109 Hochberg false discovery rate (FDR) $q < 0.05$, **Supplementary file 3**). We defined inactive
110 sequences as those whose activity is both within a two-fold change of basal activity and not
111 significantly different from the basal *Rho* promoter. We further stratified enhancers into strong
112 and weak enhancers based on whether or not they fell above the 95th percentile of scrambled
113 sequences. Using these criteria, 22% of CRX-targeted sequences are strong enhancers, 28%
114 are weak enhancers, 19% are inactive, and 17% are silencers; the remaining 13% were
115 considered ambiguous and removed from further analysis. To test whether these sequences
116 function as CRX-dependent enhancers and silencers in the genome, we examined genes
117 differentially expressed in *Crx*^{-/-} retina (Roger et al., 2014). Genes that are de-repressed are
118 more likely to be near silencers (Fisher's exact test $p = 0.001$, odds ratio = 2.1, $n = 206$) and
119 genes that are down-regulated are more likely to be near enhancers (Fisher's exact test $p =$
120 0.02 , odds ratio = 1.5, $n = 344$, Methods), suggesting that our reporter assay identified
121 sequences that act as enhancers and silencers in the genome. We sought to identify features
122 that would accurately classify these different classes of sequences.

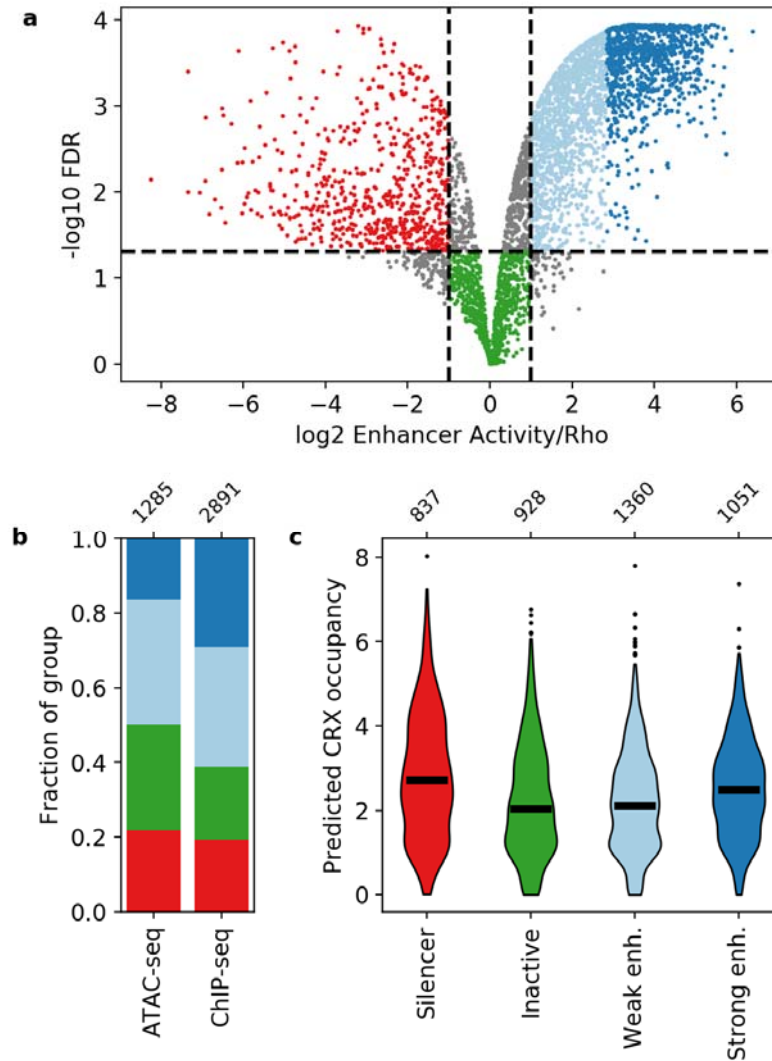
123

124 Neither CRX ChIP-seq binding status nor DNA accessibility as measured by ATAC-seq strongly
125 differentiates between these four classes (**Figure 1b**). Compared to CRX ChIP-seq peaks,
126 ATAC-seq peaks that lack CRX binding in the adult retina are slightly enriched for inactive
127 sequences (Fisher's exact test $p = 2 \times 10^{-7}$, odds ratio = 1.5) and slightly depleted for strong
128 enhancers (Fisher's exact test $p = 1 \times 10^{-21}$, odds ratio = 2.2). However, sequences with ChIP-
129 seq or ATAC-seq peaks span all four activity categories, consistent with prior reports that that
130 DNA accessibility and TF binding data are not sufficient to identify functional enhancers and
131 silencers (Doni Jayavelu et al., 2020; D. Huang et al., 2019; Z. Huang et al., 2021; Pang &
132 Snyder, 2020; White et al., 2013).

133

134 We examined whether the number and affinity of CRX motifs differentiate enhancers, silencers,
135 and inactive sequences by computing the predicted CRX occupancy (i.e. expected number of
136 bound molecules) for each sequence (White et al., 2013). Consistent with our previous work
137 (White et al., 2016), both strong enhancers and silencers have higher predicted CRX occupancy
138 than inactive sequences (Mann-Whitney U test, $p = 6 \times 10^{-10}$ and 6×10^{-17} respectively, **Figure**
139 **1c**), suggesting that total CRX motif content helps distinguish silencers and strong enhancers
140 from inactive sequences. However, predicted CRX occupancy does not distinguish strong
141 enhancers from silencers: a logistic regression classifier trained with five-fold cross-validation
142 only achieves an area under the receiver operating characteristic (AUROC) curve of
143 0.548 ± 0.023 and an area under the precision recall (AUPR) curve of 0.571 ± 0.020 (**Figure 2a**
144 and **Figure 2—figure supplement 1**). We thus sought to identify sequence features that
145 distinguish strong enhancers from silencers.

146



147

148 **Figure 1: Activity of putative *cis*-regulatory sequences with CRX motifs.** **a)** Volcano plot of
149 activity scores relative to the *Rho* promoter alone. Sequences are grouped as strong enhancers
150 (dark blue), weak enhancers (light blue), inactive (green), silencers (red), or ambiguous (grey).
151 Horizontal line, FDR $q = 0.05$. Vertical lines, 2-fold above and below *Rho*. **b)** Fraction of ChIP-
152 seq and ATAC-seq peaks that belong to each activity group. **c)** Predicted CRX occupancy of
153 each activity group. Horizontal lines, medians; enh., enhancer. Numbers at top of **(b and c)**
154 indicate n for groups.

155

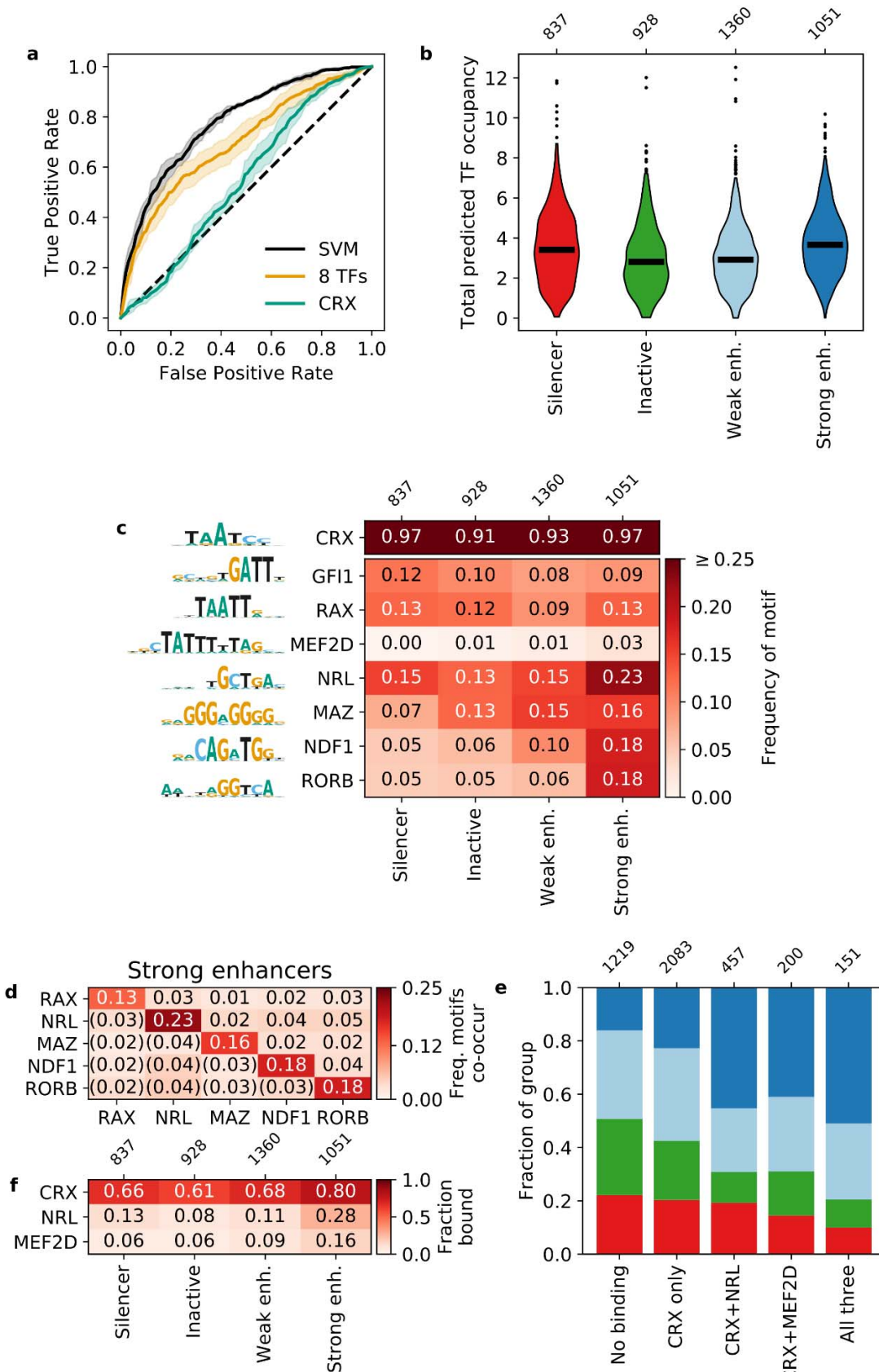
156 **Lineage-defining TF motifs differentiate strong enhancers from silencers**

157 We performed a *de novo* motif enrichment analysis to identify motifs that distinguish strong
158 enhancers from silencers and found several differentially enriched motifs matching known TFs.
159 For motifs that matched multiple TFs, we selected one representative TF for downstream
160 analysis, since TFs from the same family have PWMs that are too similar to meaningfully
161 distinguish between motifs for these TFs (**Figure 2—figure supplement 2**, Methods). Strong
162 enhancers are enriched for several motif families that include TFs that interact with CRX or are
163 important for photoreceptor development: NeuroD1/NDF1 (E-box-binding bHLH) (Morrow et al.,
164 1999), RORB (nuclear receptor) (Jia et al., 2009; Srinivas et al., 2006), MAZ or Sp4 (C2H2 zinc
165 finger) (Lerner et al., 2005), and NRL (bZIP) (Mears et al., 2001; K. P. Mitton et al., 2000). Sp4
166 physically interacts with CRX in the retina (Lerner et al., 2005), but we chose to represent the
167 zinc finger motif with MAZ because it has a higher quality score in the HOCOMOCO database
168 (Kulakovskiy et al., 2018). Silencers were enriched for a motif that resembles a partial K50
169 homeodomain motif but instead matches the zinc finger TF GF11, a member of the Snail
170 repressor family (Chiang & Ayyanathan, 2013) expressed in developing retinal ganglion cells
171 (Yang et al., 2003). Therefore, while strong enhancers and silencers are not distinguished by
172 their CRX motif content, strong enhancers are uniquely enriched for several lineage-defining
173 TFs.

174

175 To quantify how well these TF motifs differentiate strong enhancers from silencers, we trained
176 two different classification models with five-fold cross-validation. First, we trained a 6-mer
177 support vector machine (SVM) (Ghandi et al., 2014) and achieved an AUROC of 0.781 ± 0.013
178 and AUPR of 0.812 ± 0.020 (**Figure 2a** and **Figure 2—figure supplement 1**). The SVM
179 considers all 2,080 non-redundant 6-mers and provides an upper bound to the predictive power

180 of models that do not consider the exact arrangement or spacing of sequence features. We next
181 trained a logistic regression model on the predicted occupancy for eight lineage-defining TFs
182 (**Supplementary file 4**) and compared it to the upper bound established by the SVM. In this
183 model, we considered CRX, the five TFs identified in our motif enrichment analysis, and two
184 additional TFs enriched in photoreceptor ATAC-seq peaks (Hughes et al., 2017): RAX, a Q50
185 homeodomain TF that contrasts with CRX, a K50 homeodomain TF (Irie et al., 2015); and
186 MEF2D, a MADS box TF which co-binds with CRX (Andzelm et al., 2015). The logistic
187 regression model performs nearly as well as the SVM (AUROC 0.698 ± 0.036 , AUPR
188 0.745 ± 0.032 , **Figure 2a** and **Figure 2—figure supplement 1**) despite a 260-fold reduction from
189 2,080 to eight features. To determine whether the logistic regression model depends specifically
190 on the eight lineage-defining TFs, we established a null distribution by fitting 100 logistic
191 regression models with randomly selected TFs (Methods). Our logistic regression model
192 outperforms the null distribution (one-tailed Z-test for AUROC and AUPR, $p < 0.0008$, **Figure**
193 **2—figure supplement 3**), indicating that the performance of the model specifically requires the
194 eight lineage-defining TFs. To determine whether the SVM identified any additional motifs that
195 could be added to the logistic regression model, we generated *de novo* motifs using the SVM *k*-
196 mer scores and found no additional motifs predictive of strong enhancers. Finally, we found that
197 our two models perform similarly on an independent test set of CRX-targeted sequences (White
198 et al., 2013) (**Figure 2—figure supplement 3**). Since the logistic regression model performs
199 near the upper bound established by the SVM and depends specifically on the eight selected
200 motifs, we conclude that these motifs comprise nearly all of the sequence features captured by
201 the SVM that distinguish strong enhancers from silencers among CRX-targeted sequences.



203 **Figure 2: Strong enhancers contain a diverse array of motifs. a)** Receiver operating
204 characteristic for classifying strong enhancers from silencers. Solid black, 6-mer SVM; orange, 8
205 TF predicted occupancy logistic regression; aqua, predicted CRX occupancy logistic regression;
206 dashed black, chance; shaded area, 1 standard deviation based on five-fold cross-validation. **b**
207 **and c)** total predicted TF occupancy (**b**) and frequency of TF motifs (**c**) in each activity class. **d)**
208 Frequency of co-occurring TF motifs in strong enhancers. Lower triangle is expected co-
209 occurrence if motifs are independent. **e)** Frequency of activity classes, colored as in (**b**), for
210 sequences in CRX, NRL, and/or MEF2D ChIP-seq peaks. **f)** Frequency of TF ChIP-seq peaks
211 in activity classes. TFs in (**c**) are sorted by feature importance of the logistic regression model in
212 **(a)**.

213

214 **Strong enhancers are characterized by diverse total motif content**

215 To understand how these eight TF motifs differentiate strong enhancers from silencers, we first
216 calculated the total predicted occupancy of each sequence by all eight lineage-defining TFs and
217 compared the different activity classes. Strong enhancers and silencers both have higher total
218 predicted occupancies than inactive sequences, but total predicted occupancies do not
219 distinguish strong enhancers and silencers from each other (**Figure 2b, Supplementary file 5**).
220 Since strong enhancers are enriched for several motifs relative to silencers, this suggests that
221 strong enhancers are distinguished from silencers by the diversity of their motifs, rather than the
222 total number.

223

224 We considered two hypotheses for how the more diverse collection of motifs function in strong
225 enhancers: either strong enhancers depend on specific combinations of TF motifs (“TF identity
226 hypothesis”), or they instead must be co-occupied by multiple lineage-defining TFs, regardless

227 of TF identity (“TF diversity hypothesis”). To distinguish between these hypotheses, we
228 examined which specific motifs contribute to the total motif content of strong enhancers and
229 silencers. We considered motifs for a TF present in a sequence if the TF predicted occupancy
230 was above 0.5 molecules (**Supplementary file 4**), which generally corresponds to at least one
231 motif with a relative K_D above 3%. This threshold captures the effect of low affinity motifs that
232 are often biologically relevant (Crocker et al., 2015; Farley et al., 2015, 2016; Parker et al.,
233 2011). As expected, 97% of strong enhancers and silencers contain CRX motifs since the
234 sequences were selected based on CRX binding or significant matches to the CRX PWM within
235 open chromatin (**Figure 2c**). Compared to silencers, strong enhancers contain a broader
236 diversity of motifs for the eight lineage-defining TFs (**Figure 2c**). However, while strong
237 enhancers contain a broader range of motifs, no single motif occurs in a majority of strong
238 enhancers: NRL motifs are present in 23% of strong enhancers, NeuroD1 and RORB in 18%
239 each, and MAZ in 16%. Additionally, none of the motifs tend to co-occur as pairs in strong
240 enhancers: no specific pair occurred in more than 5% of sequences (**Figure 2d**). We also did
241 not observe a bias in the linear arrangement of motifs in strong enhancers (Methods). Similarly,
242 no single motif occurs in more than 15% of silencers (**Figure 2c**). These results suggest that
243 strong enhancers are defined by the diversity of their motifs, and not by specific motif
244 combinations or their linear arrangement.

245

246 The results above predict that strong enhancers are more likely to be bound by a diverse but
247 degenerate collection of TFs, compared with silencers or inactive sequences. We tested this
248 prediction by examining *in vivo* TF binding using published ChIP-seq data for NRL (Hao et al.,
249 2012) and MEF2D (Andzelm et al., 2015). Consistent with the prediction, sequences bound by
250 CRX and either NRL or MEF2D are approximately twice as likely to be strong enhancers
251 compared to sequences only bound by CRX (**Figure 2e**). Sequences bound by all three TFs are

252 the most likely to be strong or weak enhancers rather than silencers or inactive sequences.
253 However, most strong enhancers are not bound by either NRL or MEF2D (**Figure 2f**), indicating
254 that binding of these TFs is not required for strong enhancers. Our results support the TF
255 diversity hypothesis: CRX-targeted enhancers are co-occupied by multiple TFs, without a
256 requirement for specific combinations of lineage-defining TFs.

257

258 **Strong enhancers have higher motif information content than silencers**

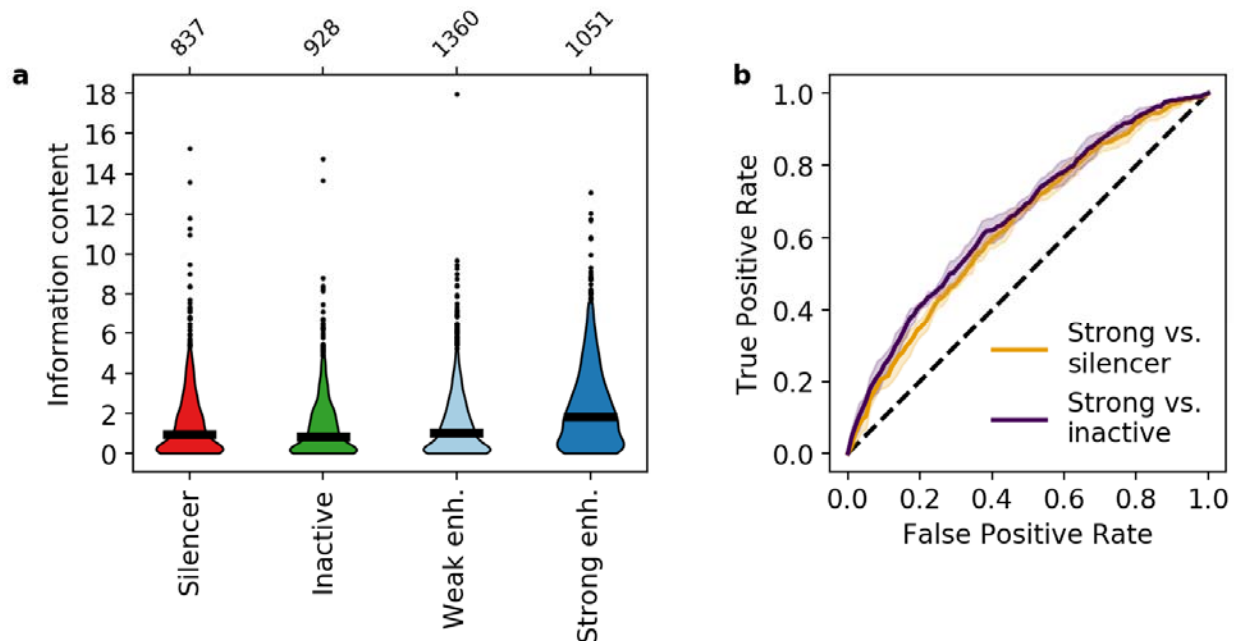
259 Our results indicate that both strong enhancers and silencers have a higher total motif content
260 than inactive sequences, while strong enhancers contain a more diverse collection of motifs
261 than silencers. To quantify these differences in the number and diversity of motifs, we computed
262 the information content of CRX-targeted sequences using Boltzmann entropy. The Boltzmann
263 entropy of a system is related to the number of ways the system's molecules can be arranged,
264 which increases with either the number or diversity of molecules (Phillips et al., 2012, Chapter
265 5). In our case, each TF is a different type of molecule and the number of each TF is
266 represented by its predicted occupancy for a *cis*-regulatory sequence. The number of molecular
267 arrangements is thus W , the number of distinguishable permutations that the TFs can be
268 ordered on the sequence, and the information content of a sequence is then $\log_2 W$ (Methods).

269

270 We found that on average, strong enhancers have higher information content than both
271 silencers and inactive sequences (Mann-Whitney U test, $p = 1 \times 10^{-23}$ and 7×10^{-34} respectively,
272 **Figure 3a, Supplementary file 5**), confirming that information content captures the effect of
273 both the number and diversity of motifs. Quantitatively, the average silencer and inactive
274 sequence contains 1.6 and 1.4 bits, respectively, which represents approximately three total
275 motifs for two TFs. Strong enhancers contain on average 2.4 bits, representing approximately

276 three total motifs for three TFs or four total motifs for two TFs. To compare the predictive value
277 of our information content metric to the model based on all eight motifs, we trained a logistic
278 regression model and found that information content classifies strong enhancers from silencers
279 with an AUROC of 0.634 ± 0.008 and an AUPR of 0.663 ± 0.014 (**Figure 3b** and **Figure 3—figure**
280 **supplement 1**). This is only slightly worse than the model trained on eight TF occupancies
281 despite an eight-fold reduction in the number of features, which is itself comparable to the SVM
282 with 2,080 features. The difference between the two logistic regression models suggests that
283 the specific identities of TF motifs make some contribution to the eight TF model, but that most
284 of the signal captured by the SVM can be described with a single metric that does not assign
285 weights to specific motifs. Information content also distinguishes strong enhancers from inactive
286 sequences (AUROC 0.658 ± 0.012 , AUPR 0.675 ± 0.019 , **Figure 3b** and **Figure 3—figure**
287 **supplement 1**). These results indicate that strong enhancers are characterized by higher
288 information content, which reflects both the total number and diversity of motifs.

289



290

291 **Figure 3: Information content classifies strong enhancers. a)** Information content for
292 different activity classes. **b)** Receiver operating characteristic of information content to classify
293 strong enhancers from silencers (orange) or inactive sequences (indigo).

294

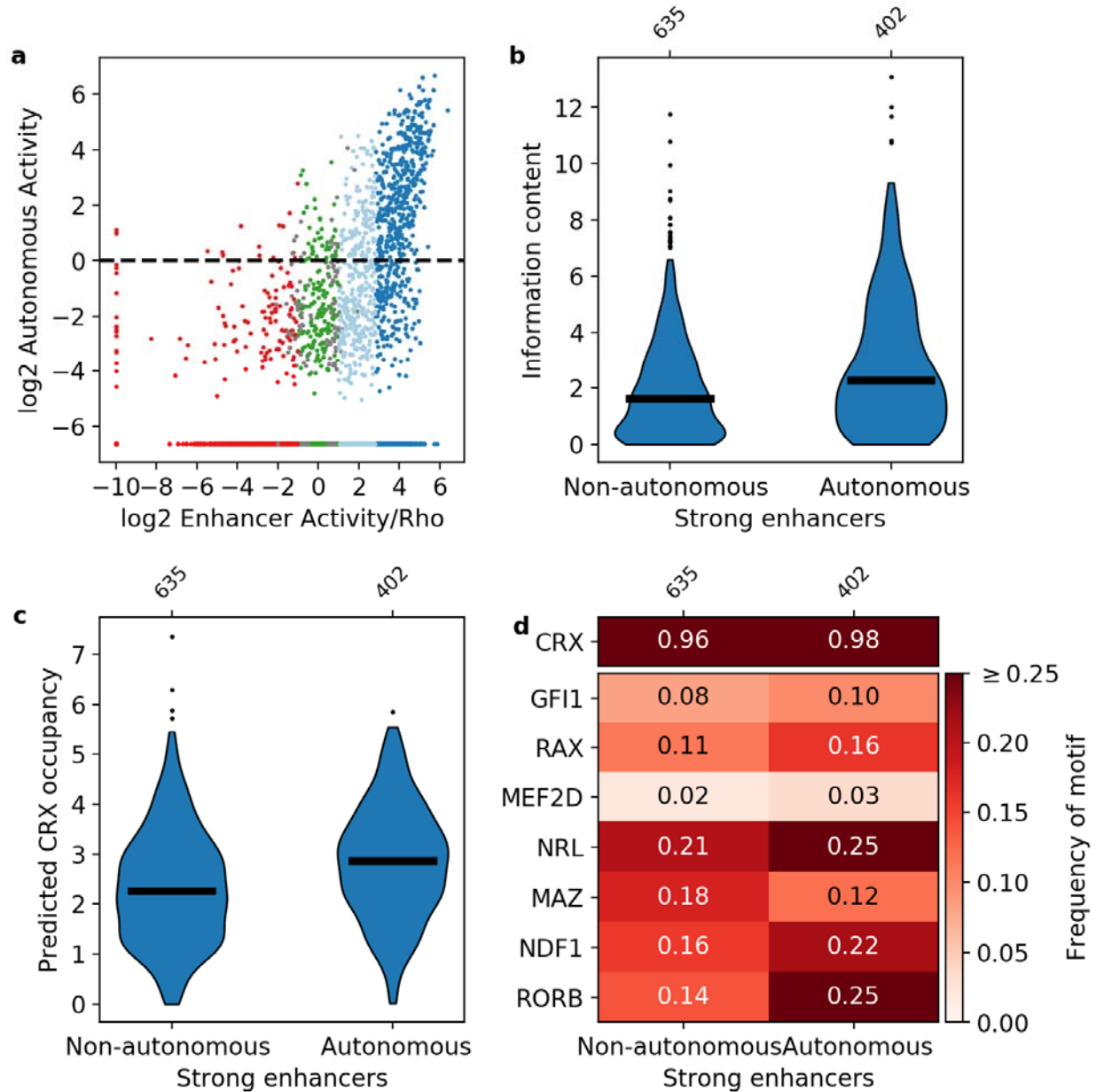
295 **Strong enhancers require high information content but not NRL motifs**

296 Our results show that except for CRX, none of the lineage-defining motifs occur in a majority of
297 strong enhancers. However, all sequences were tested in reporter constructs with the *Rho*
298 promoter, which contains an NRL motif and three CRX motifs (Corbo et al., 2010; Kwasnieski et
299 al., 2012). Since NRL is a key co-regulator with CRX in rod photoreceptors, we tested whether
300 strong enhancers generally require NRL, which would be inconsistent with our TF diversity
301 hypothesis. We removed the NRL motif by recloning our MPRA library without the basal *Rho*
302 promoter. If strong enhancers require an NRL motif for high activity, then only CRX-targeted
303 sequences with NRL motifs will drive reporter expression. If information content (i.e. total motif
304 content and diversity) is the primary determinant of strong enhancers, only CRX-targeted
305 sequences with sufficient motif diversity, measured by information content, will drive reporter
306 expression regardless of whether or not NRL motifs are present.

307

308 We replaced the *Rho* promoter with a minimal 23 bp polylinker sequence between our libraries
309 and *DsRed*, and repeated the MPRA (**Figure 1—figure supplement 1, Supplementary file 3**).
310 CRX-targeted sequences were designated as “autonomous” if they retained activity in the
311 absence of the *Rho* promoter ($\log_2(\text{RNA}/\text{DNA}) > 0$, Methods). We found that 90% of
312 autonomous sequences are from the enhancer class, while less than 3% of autonomous
313 sequences are from the silencer class (**Figure 4a**). This confirms that the distinction between
314 silencers and enhancers does not depend on the *Rho* promoter, which is consistent with our

315 previous finding that CRX-targeted silencers repress other promoters (Hughes et al., 2018;
316 White et al., 2016). However, while most autonomous sequences are enhancers, only 39% of
317 strong enhancers and 9% of weak enhancers act autonomously. Consistent with a role for
318 information content, autonomous strong enhancers have higher information content (Mann-
319 Whitney U test $p = 4 \times 10^{-8}$, **Figure 4b**) and higher predicted CRX occupancy (Mann-Whitney U
320 test $p = 9 \times 10^{-12}$, **Figure 4c**) than non-autonomous strong enhancers. We found no evidence
321 that specific lineage-defining motifs are required for autonomous activity, including NRL, which
322 is present in only 25% of autonomous strong enhancers (**Figure 4d**). Similarly, NRL ChIP-seq
323 binding (Hao et al., 2012) occurs more often among autonomous strong enhancers (41% vs.
324 19%, Fisher's exact test $p = 2 \times 10^{-14}$, odds ratio = 3.0), yet NRL binding still only accounts for a
325 minority of these sequences. We thus conclude that strong enhancers require high information
326 content, rather than any specific lineage-defining motifs.



327

328 **Figure 4: Sequence features of autonomous and non-autonomous strong enhancers. a)**

329 Activity of library in the presence (x-axis) or absence (y-axis) of the *Rho* promoter. Dark blue,

330 strong enhancers; light blue, weak enhancers; green, inactive; red, silencers; grey, ambiguous;

331 horizontal line, cutoff for autonomous activity. Points on the far left and/or very bottom are

332 sequences that were present in the plasmid pool but not detected in the RNA. **b-d)** Comparison

333 of autonomous and non-autonomous strong enhancers for information content **(b)**, predicted
334 CRX occupancy **(c)**, and frequency of TF motifs **(d)**.

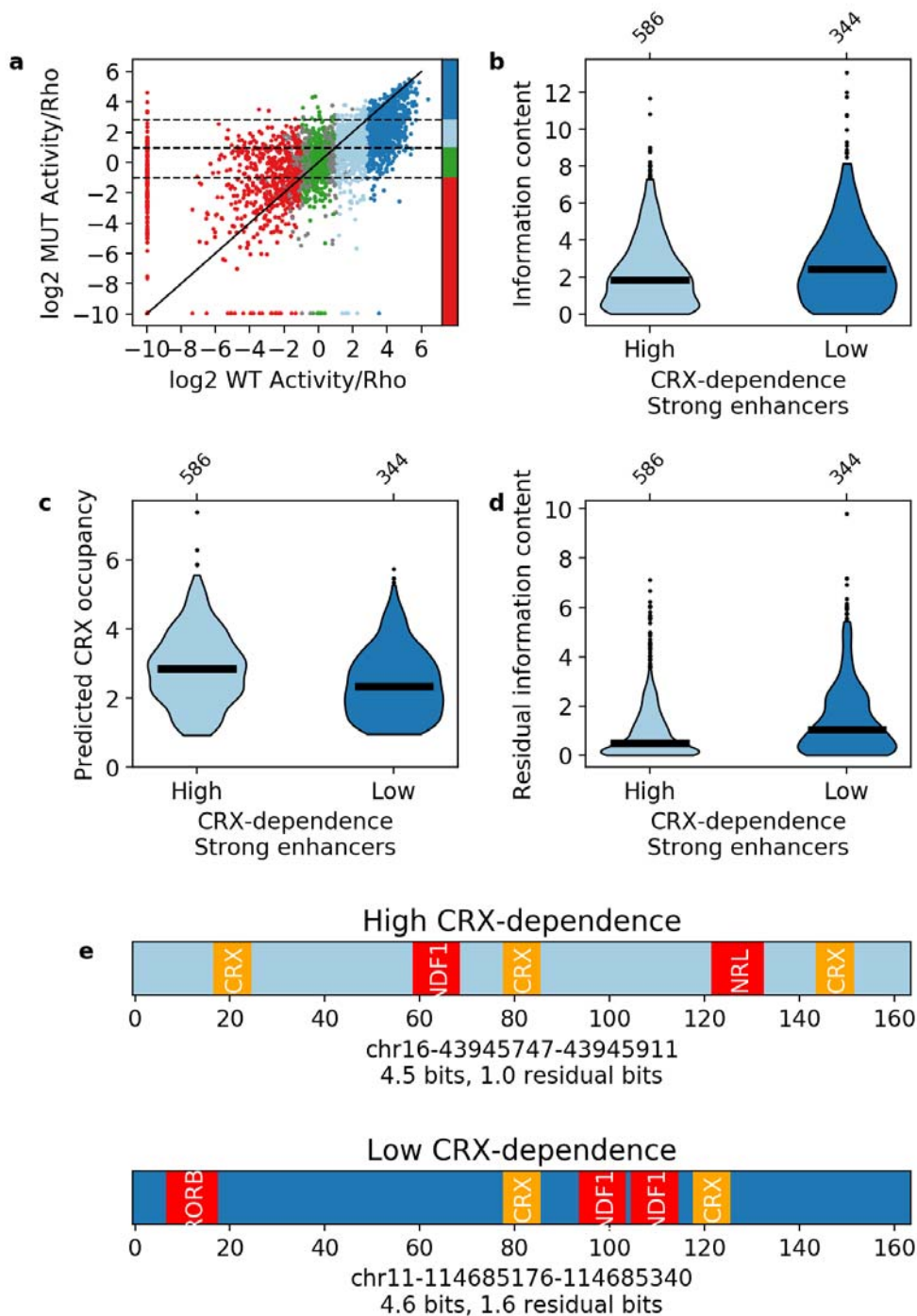
335

336 **TF motifs contribute independently to strong enhancers**

337 Our results indicate that information content distinguishes strong enhancers from silencers and
338 inactive sequences. Information content only takes into account the total number and diversity
339 of motifs in a sequence and not any potential interactions between them. The classification
340 success of information content thus suggests that each TF motif will contribute independently to
341 enhancer activity. We tested this prediction with CRX-targeted sequences where all CRX motifs
342 were abolished by point mutation (**Supplementary file 3**). Consistent with our previous work
343 (White et al., 2013), mutating CRX motifs causes the activities of both enhancers and silencers
344 to regress towards basal levels (Pearson's $r = 0.608$, **Figure 5a**), indicating that most enhancers
345 and silencers show some dependence on CRX. However, 40% of wild-type strong enhancers
346 show low CRX dependence and remain strong enhancers with their CRX motifs abolished.
347 Although strong enhancers with high and low CRX dependence have similar wild-type
348 information content (**Figure 5b**), strong enhancers with low CRX dependence have lower
349 predicted CRX occupancy than those with high CRX dependence (Mann-Whitney U test $p = 2 \times$
350 10^{-9} , **Figure 5c**), and also have higher “residual” information content (i.e. information content
351 without CRX motifs, Mann-Whitney U test $p = 1 \times 10^{-7}$, **Figure 5d**). Low CRX dependence
352 sequences have an average of 1.5 residual bits, which corresponds to three motifs for two TFs,
353 while high CRX dependence sequences have an average of 1.0 residual bits, which
354 corresponds to two motifs for two TFs (**Figure 5e**).

355

356 Strong enhancers with low and high CRX dependence have similar wild-type information
357 content and similar total predicted occupancy (**Figures 5b and e**). As a result, sequences with
358 more CRX motifs have fewer motifs for other TFs, suggesting that there is no evolutionary
359 pressure for enhancers to contain additional motifs beyond the minimum amount of information
360 content required to be active. To test this idea, we calculated the minimum number and diversity
361 of motifs necessary to specify a relatively unique location in the genome (Wunderlich & Mirny,
362 2009) and found that a 164 bp sequence only requires five motifs for three TFs (Methods).
363 These motif requirements can be achieved in two ways with similar information content that
364 differ only in the quantitative number of motifs for each TF. In other words, the number of motifs
365 for any particular TF is not important so long as there is sufficient information content. Taken
366 together, we conclude that each TF motif provides an independent contribution towards
367 specifying strong enhancers.



368

369 **Figure 5: Independence of TF motifs in strong enhancers.** a) Activity of sequences with and
 370 without CRX motifs. Points are colored by the activity group with CRX motifs intact: dark blue,
 371 strong enhancers; light blue, weak enhancers; green, inactive; red, silencers; grey, ambiguous;
 372 horizontal dotted lines and color bar represent the cutoffs for the same groups when CRX motifs

373 are mutated. Solid black line is the $y = x$ line. **b-d**) Comparison of strong enhancers with high
374 and low CRX dependence for information content (**b**), predicted CRX occupancy (**c**), and
375 residual information content (**d**). **(e)** Representative strong enhancers with high (top) or low
376 (bottom) CRX-dependence.

377

378

379 Discussion

380 Many regions in the genome are bound by TFs and bear the epigenetic hallmarks of active *cis*-
381 regulatory sequences, yet fail to exhibit *cis*-regulatory activity when tested directly. The
382 discrepancy between measured epigenomic state and *cis*-regulatory activity indicates that
383 enhancers and silencers consist of more than the minimal sequence features necessary to
384 recruit TFs and chromatin-modifying factors. Our results show that enhancers, silencers, and
385 inactive sequences in developing photoreceptors can be distinguished by their motif content,
386 even though they are indistinguishable by CRX binding or chromatin accessibility. We show that
387 both enhancers and silencers contain more TF motifs than inactive sequences, and that
388 enhancers also contain more diverse sets of motifs for lineage-defining TFs. These differences
389 are captured by our measure of information content. Information content, as a single metric,
390 identifies strong enhancers nearly as well as an unbiased set of 2,080 non-redundant 6-mers
391 used for an SVM, indicating that a simple measure of motif number and diversity can capture
392 the key sequence features that distinguish enhancers from other sequences that lie in open
393 chromatin.

394

395 The results of our information content classifier are consistent with the TF collective model of
396 enhancers (Junion et al., 2012; Spitz & Furlong, 2012): globally, active enhancers are specified
397 by the combinatorial action of lineage-defining TFs with little constraint on which motifs must co-
398 occur. We show that CRX-targeted enhancers are distinguished from inactive CRX-targeted
399 sequences by a larger, more diverse collection of TF motifs, and not any specific combination of
400 motifs. This indicates that enhancers are active because they have acquired the necessary
401 number of TF binding motifs, and not because they are defined by a strict regulatory grammar.
402 Sequences with fewer motifs may be bound by CRX and reside within open chromatin, but they
403 lack sufficient TF binding for activity. Such loose constraints would facilitate the *de novo*
404 emergence of tissue-specific enhancers and silencers over evolution and explain why critical
405 cell type-specific TF interactions, such as CRX and NRL in rod photoreceptors, occur at only a
406 minority of the active enhancers in that cell type (Hsiau et al., 2007; Hughes et al., 2018; White
407 et al., 2013).

408

409 Like enhancers, CRX-targeted silencers require higher motif content and are dependent on
410 CRX motifs, but they lack the TF diversity of enhancers. The lack of TF diversity in silencers
411 parallels the architecture of signal-responsive *cis*-regulatory sequences, which are silencers in
412 the absence of a signal and require multiple activators for induction (Barolo & Posakony, 2002).
413 Consistent with this, we previously showed using synthetic sequences that high occupancy of
414 CRX alone is sufficient to encode silencers while the addition of a single NRL motif converts
415 synthetic silencers to enhancers, and that genomic sequences with very high CRX motif content
416 repress a basal promoter that lacks NRL motifs (White et al., 2016). We found that
417 photoreceptor genes which are de-repressed upon loss of CRX are located near *cis*-regulatory
418 sequences with high CRX motif content, and that genes near regions that are bound only by
419 CRX are expressed at lower levels than genes near regions co-bound by CRX and NRL (White

420 et al 2016). In the current study, we find that silencers in our MPRA library are more likely to
421 occur near de-repressed photoreceptor genes, while strong enhancers are enriched near genes
422 that lose expression in *Crx*^{-/-} retina. These findings suggest that the low TF diversity and high
423 CRX motif content that characterize silencers in our MPRA library are also important for
424 silencing in the genome.

425

426 The contrast in motif diversity between enhancers and silencers that we observe could explain
427 how CRX achieves selective activation and repression of its target genes in multiple cell types
428 and across developmental time points (Murphy et al., 2019; Ruzycki et al., 2018). CRX itself is
429 required for silencing, and we previously showed that some silencers become active enhancers
430 in *Crx*^{-/-} retina (White et al., 2016). The mechanism of CRX-based silencing is unknown,
431 however CRX cooperates with other TFs that can sometimes act as repressors of cell type-
432 specific genes (J. Chen et al., 2005; Peng et al., 2005; Webber et al., 2008), while other
433 repressors can directly inhibit activation by CRX or its co-activators (Dorval et al., 2006;
434 Hlawatsch et al., 2013; Kenneth P. Mitton et al., 2003; Sanuki et al., 2010). In *Drosophila*
435 photoreceptors, selective silencing of opsin genes is controlled by cell-type specific expression
436 of a repressor, Dve, which acts on the same K50 homeodomain binding sites as a universally
437 expressed activator, Otd, a homolog of CRX (Rister et al., 2015). Other transcriptional activators
438 selectively act as repressors in the same cell type. GATA-1 represses the *GATA-2* promoter by
439 displacing CREB-binding protein (CBP), while at other genes GATA-1 binds CBP to activate
440 transcription (Grass et al., 2003). Selective repression by GATA-1 is also mediated by
441 chromatin occupancy levels and interaction with coregulators (Johnson et al., 2006), which is
442 consistent with our finding that sequence context enables a TF to both activate and repress
443 genes in the same cell type.

444

445 Given the central role of CRX in selectively regulating genes in multiple closely related cell types
446 (Murphy et al., 2019), we speculate that CRX-targeted silencers may contain sufficient
447 information to act as enhancers in other cell types in which a different set of co-activating TFs
448 are expressed. This hypothesis would be consistent with the finding that many silencers are
449 enhancers in other cell types (Doni Jayavelu et al., 2020; Gisselbrecht et al., 2020; Ngan et al.,
450 2020). Our work suggests that characterizing sequences by their motif information content
451 offers a way to identify these different classes of *cis*-regulatory sequences in the genome.

452

453

454 Materials and Methods

455 Key Resources Table

Reagent type (species) or resource	Designation	Source or reference	Identifiers	Additional information
strain, strain background (<i>Mus musculus</i> , male and female)	CD-1	Charles River	Strain code 022	
Recombinant DNA reagent	Library1	This paper		Listed in Supplementary file 1

Recombinant DNA reagent	Library2	This paper		Listed in Supplementary file 2
Recombinant DNA reagent	pJK01_Rho minprox-DsRed	Kwasnieski et al., 2012	AddGene plasmid # 173489	
Recombinant DNA reagent	pJK03_Rho_basal_DsRed	Kwasnieski et al., 2012	AddGene plasmid # 173490	
Sequence-based reagent	Primers	IDT		Listed in Supplementary file 6
Commercial assay or kit	Monarch PCR Cleanup Kit	New England Biolabs	T1030S	
Commercial assay or kit	Monarch DNA Gel Extraction Kit	New England Biolabs	T1020L	
Commercial assay or kit	TURBO DNA-free	Invitrogen	AM1907	
Commercial assay or kit	SuperScript III Reverse Transcriptase	Invitrogen	18080044	
Software, algorithm	Bedtools	https://bedtools.readthedocs.io/en/latest/	RRID: SCR_006646	
Software, algorithm	MEME Suite	https://meme-suite.org/	RRID: SCR_001783	
Software, algorithm	ShapeMF	https://github.com/h-samee/shape	DOI: 10.1016/j.cels.2018.12.001	

		-motif		
Software, algorithm	Numpy	https://numpy.org/	DOI: 0.1038/s41586-020-2649-2	
Software, algorithm	Scipy	https://www.scipy.org/	DOI:10.1038/s41592-019-0686-2	
Software, algorithm	Pandas	https://pandas.pydata.org/	DOI: 10.5281/zenodo.3509134	
Software, algorithm	Matplotlib	https://matplotlib.org/	DOI: 10.5281/zenodo.1482099	
Software, algorithm	Logomaker	https://github.com/jbkinney/logomaker	DOI: 10.1093/bioinformatics/btz921	

456

457 **Library Design**

458 CRX ChIP-seq peaks re-processed by Ruzycki *et al.* (Ruzycki et al., 2018) were intersected with
459 previously published CRX MPRA libraries (Hughes et al., 2018; White et al., 2013) and one
460 unpublished library to select sequences that had not been previously tested by MPRA. These
461 sequences were scanned for instances of CRX motifs using FIMO version 4.11.2 (Bailey et al.,
462 2009), a p-value cutoff of 2.3×10^{-3} (see below), and a CRX PWM derived from an
463 electrophoretic mobility shift assay (J. Lee et al., 2010). We centered 2622 sequences on the
464 highest scoring CRX motif. For 677 sequences without a CRX motif, we instead centered them
465 using the Gibbs sampler from ShapeMF (Github commit abe8421) (Samee et al., 2019) and a
466 motif size of 10.

467

468 For sequences unbound in CRX ChIP-seq but in open chromatin, we took ATAC-seq peaks
469 collected in 8 week FACS-purified rods, green cones, and *Nrl*^{-/-} blue cones (Hughes et al., 2017)
470 and removed sequences that overlapped with CRX ChIP-seq peaks. The remaining sequences
471 were scanned for instances of CRX motifs using FIMO with a p-value cutoff of 2.5×10^{-3} and the
472 CRX PWM. Sequences with a CRX motif were kept and the three ATAC-seq data sets were
473 merged together, intersected with H3K27ac and H3K4me3 ChIP-seq peaks collected in P14
474 retinas (Ruzycki et al., 2018), and centered on the highest scoring CRX motifs. We randomly
475 selected 1004 H3K27ac⁺H3K4me3⁻ sequences and 541 H3K27ac⁺H3K4me3⁺ to reflect the fact
476 that ~35% of CRX ChIP-seq peaks are H3K4me3⁺. After synthesis of our library, we discovered
477 11% of these sequences do not actually overlap H3K27ac ChIP-seq peaks (110/1004 of the
478 H3K4me3⁻ group and 60/541 of the H3K4me3⁺ group), but we still included them in the analysis
479 because they contain CRX motifs in ATAC-seq peaks.

480

481 All data was converted to mm10 coordinates using the UCSC liftOver tool (Haeussler et al.,
482 2019) and processed using Bedtools version 2.27.1 (Quinlan & Hall, 2010). All sequences in our
483 library design were adjusted to 164 bp and screened for instances of EcoRI, SpeI, SphI, and
484 NotI sites. In total, our library contains 4844 genomic sequences (2622 CRX ChIP-seq peaks
485 with motifs, 677 CRX ChIP-seq peaks without motifs, 1004 CRX⁻ATAC⁺H3K27ac⁺H3K4me3⁻
486 CRX motifs, and 541 CRX⁻ATAC⁺H3K27ac⁺H3K4me3⁺ CRX motifs), a variant of each sequence
487 with all CRX motifs mutated, 150 scrambled sequences, and a construct for cloning the basal
488 promoter alone.

489

490 For sequences centered on CRX motifs, all CRX motifs with a p-value of 2.5×10^{-3} or less were
491 mutated by changing the core TAAT to TACT (J. Lee et al., 2010) on the appropriate strand, as

492 described previously (Hughes et al., 2018; White et al., 2013). We then re-scanned sequences
493 and mutated any additional motifs inadvertently created.

494

495 To generate scrambled sequences, we randomly selected 150 CRX ChIP-seq peaks spanning
496 the entire range of GC content in the library. We then scrambled each sequence while
497 preserving dinucleotide content as previously described (White et al., 2013). We used FIMO to
498 confirm that none of the scrambled sequences contain CRX motifs.

499

500 We unintentionally used a FIMO p-value cutoff of 2.3×10^{-3} to identify CRX motifs in CRX ChIP-
501 seq peaks, rather than the slightly less stringent 2.5×10^{-3} cutoff used with ATAC-seq peaks or
502 mutating CRX motifs. Due to this anomaly, there may be sequences centered using ShapeMF
503 that should have been centered on a CRX motif, and these motifs would not have been mutated
504 because CRX motifs were not mutated in sequences centered using ShapeMF. However, any
505 intact CRX motifs would still be captured in the residual information content of the mutant
506 sequence.

507

508 **Plasmid Library Construction**

509 We generated two 15,000 libraries of 230 bp oligonucleotides (oligos) from Agilent Technologies
510 (Santa Clara, CA) through a limited licensing agreement. Our library was split across the two
511 oligo pools, ensuring that both the genomic and mutant forms of each sequence were placed in
512 the same oligo pool (**Supplementary files 1 and 2**). Both oligo pools contain all 150 scrambled
513 sequences as an internal control. All sequences were assigned three unique barcodes as
514 previously described (White et al., 2013). In each oligo pool, the basal promoter alone was

515 assigned 18 unique barcodes. Oligos were synthesized as follows: 5' priming sequence
516 (GTAGCGTCTGTCCGT)/EcoRI site/Library sequence/Spel site/C/SphI site/Barcode
517 sequence/NotI site/3' priming sequence (CAACTACTACTACAG). To clone the basal promoter
518 into barcoded oligos without any upstream *cis*-regulatory sequence, we placed the Spel site
519 next to the EcoRI site, which allowed us to place the promoter between the EcoRI site and the
520 3' barcode.

521

522 We cloned the synthesized oligos as previously described by our group (Kwasnieski et al.,
523 2012; White et al., 2013, 2016). Specifically, for each oligo pool we used 50 femtomoles of
524 template and 4 cycles of PCR in each of multiple 50 microliter reactions (New England Biolabs
525 (NEB), Ipswich, MA) (NEB Phusion) using primers MO563 and MO564 (**Supplementary file 6**),
526 2% DMSO, and an annealing temperature of 57C. PCR amplicons were purified from a 2%
527 agarose gel (NEB), digested with EcoRI-HF and NotI-HF (NEB), and then cloned into the EagI
528 and EcoRI sites of the plasmid pJK03 with multiple 20 microliter ligation reactions (NEB T4
529 ligase). The libraries were transformed into 5-alpha electrocompetent cells (NEB) and grown in
530 liquid culture. Next, 2 micrograms of each library was digested with SphI-HF and SpeI-HF (NEB)
531 and then treated with Antarctic phosphatase (NEB).

532

533 The *Rho* basal promoter and *DsRed* reporter gene was amplified from the plasmid pJK01 using
534 primers MO566 and MO567 (**Supplementary file 6**). The Polylinker and *DsRed* reporter gene
535 was amplified from the plasmid pJK03 using primers MO610 and MO567 (**Supplementary file**
536 **6**). The Polylinker is a short 23 bp multiple cloning site with no known core promoter motifs.
537 Inserts were purified from a 1% agarose gel (NEB), digested with NheI-HF and SphI-HF (NEB),

538 and cloned into the libraries using multiple 20 microliter ligations (NEB T4 ligase). The libraries
539 were transformed into 5-alpha electrocompetent cells (NEB) and grown in liquid culture.

540

541 **Retinal Explant Electroporation**

542 Animal procedures were performed in accordance with a Washington University in St. Louis

543 Institutional Animal Care and Use Committee approved vertebrate animals protocol.

544 Electroporation into retinal explants and RNA extraction was performed as described previously

545 (Hsiau et al., 2007; Hughes et al., 2018; Kwasnieski et al., 2012; White et al., 2013, 2016).

546 Briefly, retinas were isolated from P0 newborn CD-1 mice and electroporated in a solution with

547 30 micrograms library and 30 micrograms *Rho*-GFP. Electroporated retinas were cultured for

548 eight days, at which point they were harvested, washed three times with HBSS (ThermoFisher

549 Scientific/Gibco, Waltham, MA), and stored in TRIzol (ThermoFisher Scientific/Invitrogen,

550 Waltham, MA) at -80C. Five retinas were pooled for each biological replicate and three

551 replicates were performed for each library. RNA was extracted from TRIzol according to

552 manufacturer instructions and treated with TURBO DNase (Invitrogen). cDNA was prepared

553 using SuperScript RT III (Invitrogen) with oligo dT primers. Barcodes from both the cDNA and

554 the plasmid DNA pool were amplified for sequencing (described below). The resulting products

555 were mixed at equal concentration and sequenced on the Illumina NextSeq platform. We

556 obtained greater than 1300x coverage across all samples.

557

558 *Rho* libraries were amplified using primers MO574 and MO575 (**Supplementary file 6**) for 6

559 cycles at an annealing temperature of 66C followed by 18 cycles with no annealing step (NEB

560 Phusion) and then purified with the Monarch PCR kit (NEB). PCR amplicons were digested

561 using MfeI-HF and SphI-HF (NEB) and ligated to custom annealed adaptors with PE2 indexing

562 barcodes and phased P1 barcodes (**Supplementary file 6**). The final enrichment PCR used
563 primers MO588 and MO589 (**Supplementary file 6**) for 20 cycles at an annealing temperature
564 of 66C (NEB Phusion), followed by purification with the Monarch PCR kit. Polylinker libraries
565 were amplified using primers BC_CRX_Nested_F and BC_CRX_R (**Supplementary file 6**) for
566 30 cycles (NEB Q5) at an annealing temperature of 67C and then purified with the Monarch
567 PCR kit. Illumina adaptors were then added via two further rounds of PCR. First, P1 indexing
568 barcodes were added using forward primers P1_inner_A through P1_inner_D and reverse
569 primer P1_inner_nested_rev (**Supplementary file 6**) for 5 cycles at an annealing temperature
570 of 55C followed by 5 cycles with no annealing step (NEB Q5). PE2 indexing barcodes were then
571 added by amplifying 2 microliters of the previous reaction with forward primer P1_outer and
572 reverse primers PE2_outer_SIC69 and PE2_outer_SIC70 (**Supplementary file 6**) for 5 cycles
573 at an annealing temperature of 66C followed by 5 cycles with no annealing step (NEB Q5) and
574 then purified with the Monarch PCR kit.

575

576 **Data Processing**

577 All data processing, statistical analysis, and downstream analyses were performed in Python
578 version 3.6.5 using Numpy version 1.15.4 (Harris et al., 2020), Scipy version 1.1.0 (Virtanen et
579 al., 2020), and Pandas version 0.23.4 (McKinney, 2010); and visualized using Matplotlib version
580 3.0.2 (Hunter, 2007) and Logomaker version 0.8 (Tareen & Kinney, 2020). All statistical analysis
581 used two-sided tests unless noted otherwise.

582

583 Sequencing reads were filtered to ensure that the barcode sequence perfectly matched the
584 expected sequence (>93% reads in a sample for the *Rho* libraries, >86% reads for the
585 Polylinker libraries). For the *Rho* libraries, barcodes that had less than 10 raw counts in the DNA

586 sample were considered missing and removed from downstream analysis. Barcodes that had
587 less than 5 raw counts in any cDNA sample were considered present in the input plasmid pool
588 but below the detection limit and thus set to zero in all samples. Barcode counts were
589 normalized by reads per million (RPM) for each sample. Barcode expression was calculated by
590 dividing the cDNA RPM by the DNA RPM. Replicate-specific expression was calculated by
591 averaging the barcodes corresponding to each library sequence. After performing statistical
592 analysis (see below), expression levels were normalized by replicate-specific basal mean
593 expression and then averaged across biological replicates.

594

595 For the Polylinker assay, the expected lack of expression of many constructs required different
596 processing. Barcodes that had less than 50 raw counts in the DNA sample were removed from
597 downstream analysis. Barcodes were normalized by RPM for each replicate. Barcodes that had
598 less than 8 RPM in any cDNA sample were set to zero in all samples. cDNA RPM were then
599 divided by DNA RPM as above. Within each biological replicate, barcodes were averaged as
600 above but were not normalized to basal expression because there is no basal construct.
601 Expression values were then averaged across biological replicates. Due to the low expression
602 of scrambled sequences and the lack of a basal construct, we were unable to assess data
603 calibration with the same rigor as above.

604

605 **Assignment of Activity Classes**

606 Activity classes were assigned by comparing expression levels to basal promoter expression
607 levels across replicates. The null hypothesis is that the expression of a sequence is the same as
608 basal levels. Expression levels were approximately log-normally distributed, so we computed
609 the log-normal parameters for each sequence and then performed Welch's t-test. We corrected

610 for multiple hypotheses using the Benjamini-Hochberg False Discovery Rate (FDR) procedure.
611 We corrected for multiple hypotheses in each library separately to account for any potential
612 batch effects between libraries. The \log_2 expression was calculated after adding a pseudocount
613 of 1×10^{-3} to every sequence.

614

615 Sequences were classified as enhancers if they were 2-fold above basal and the q-value was
616 below 0.05. Silencers were similarly defined as 2-fold below basal and q-value less than 0.05.
617 Inactive sequences were defined as within a 2-fold change and q-value greater than or equal to
618 0.05. All other sequences were classified as ambiguous and removed from further analysis. We
619 used scrambled sequences to further stratify enhancers into strong and weak enhancers, using
620 the rationale that scrambled sequences give an empirical distribution for the activity of random
621 sequences. We defined strong enhancers as enhancers that are above the 95th percentile of
622 scrambled sequences and all other enhancers as weak enhancers.

623

624 For the Polylinker assay, we did not have a basal construct as a reference point. Instead, we
625 defined a sequence to have autonomous activity if the average cDNA barcode counts were
626 higher than average DNA barcode counts, and non-autonomous otherwise. The \log_2 expression
627 was calculated after adding a pseudocount of 1×10^{-2} to every sequence.

628

629 **RNA-seq Analysis**

630 We obtained RNA-seq data from WT and Crx^{-/-} P21 retinas [REF Rogers 2014] processed into
631 a counts matrix for each gene by Ruzycki et al. [REF Ruzycki 2018]. Each sample was
632 normalized by read counts per million and replicates were averaged together. Genes with at

633 least a two-fold change between genotypes were considered differentially expressed. We
634 determined which differentially expressed genes are near a member of our library using
635 previously published associations between retinal ATAC-seq peaks and genes [REF Murphy
636 2019]. For de-repressed genes, we determined how often the nearest library member is a
637 silencer; for down-regulated genes, we determined how often the nearest library member is a
638 strong or weak enhancer.

639

640 **Motif Analysis**

641 We performed motif enrichment analysis using the MEME Suite version 5.0.4 (Bailey et al.,
642 2009). We searched for motifs that were enriched in one group of sequences relative to another
643 group using DREME-py3 with the parameters -mink 6 -maxk 12 -e 0.05 and compared the de
644 novo motifs to known motifs using TOMTOM on default parameters. We ran DREME using
645 strong enhancers as positives and silencers as negatives, and vice versa. For TOMTOM, we
646 used version 11 of the full mouse HOCOMOCO database (Kulakovskiy et al., 2018) with the
647 following additions from the JASPAR human database (Khan et al., 2018): NRL (accession
648 MA0842.1), RORB (accession MA1150.1), and RAX (accession MA0718.1). We added these
649 PWMs because they have known roles in the retina, but the mouse PWMs were not in the
650 HOCOMOCO database. We also used the CRX PWM that we used to design the library. Motifs
651 were selected for downstream analysis based on their matches to the de novo motifs, whether
652 the TF had a known role in retinal development, and the quality of the PWM. Because PWMs
653 from TFs of the same family were so similar, we used one TF for each DREME motif,
654 recognizing that these motifs may be bound by other TFs that recognize similar motifs. We did
655 not use any PWMs with a quality of “D”. We excluded DREME motifs without a match to the
656 database from further analysis; most of these resemble dinucleotides.

657

658 **Predicted Occupancy**

659 We computed predicted occupancy as previously described (White et al., 2013; Zhao et al.,
660 2009). Briefly, we normalized each letter probability matrix by the most probable letter at each
661 position. We took the negative log of this matrix and multiplied by 2.5, which corresponds to the
662 ideal gas constant times 300 Kelvin, to obtain an energy weight matrix. We used a chemical
663 potential μ of 9 for all TFs. At this value, the probability of a site being bound is at least 0.5 if the
664 relative K_D is at least 0.03 of the optimal binding site. We computed the predicted occupancy for
665 every site on the forward and reverse strands and summed them together to get a single value
666 for each TF.

667

668 To determine if there is a bias in the linear arrangement of motifs, we selected strong enhancers
669 with exactly one site occupied by CRX and exactly one site occupied by a second TF. We
670 counted the number of times the position of the second TF was 5' and 3' of the CRX site and
671 then performed a binomial test. We did not observe a statistically significant bias for any TF at
672 an FDR q-value cutoff of 0.05. We also performed this analysis for silencers with exactly one
673 site occupied by CRX and exactly one site occupied by NRL and did not observe a significant
674 difference in the 5' vs. 3' bias of strong enhancers vs. silencers (Fisher's exact test $p = 0.17$).

675

676 **Information Content**

677 To capture the effects of TF predicted occupancy and diversity in a single metric, we calculated
678 the motif information content using Boltzmann entropy. Boltzmann's equation states that the
679 entropy of a system S is related to the number of ways the molecules can be arranged

680 (microstates) W via the equation $S = k_B \log W$, where k_B is Boltzmann's constant (Phillips et al.,
681 2012, Chapter 5). The number of microstates is defined as $W = \frac{N!}{\prod_i N_i!}$ where N is the total
682 number of particles and N_i are the number of the i -th type of particles. In our case, the system is
683 the collection of predicted binding motifs for different TFs in a *cis*-regulatory sequence. We
684 assume each TF is a different type of molecule because the DNA binding domain of each TF
685 belongs to a different subfamily. The number of molecular arrangements W represents the
686 number of distinguishable ways that the TFs can be ordered on the sequence. Thus, N_i is the
687 predicted occupancy of the i -th TF and N is the total predicted occupancy of all TFs on the *cis*-
688 regulatory sequence. Because the predicted occupancies are continuous values, we exploit the
689 definition of the Gamma function, $\Gamma(N + 1) = N!$, to rewrite $W = \frac{\Gamma(N+1)}{\prod_i \Gamma(N_i+1)}$.

690

691 If we assume that each arrangement of motifs is equally likely, then we can write the probability
692 of arrangement $w = 1, \dots, W$ as $p_w = \frac{1}{W}$ and rewrite the entropy as $S = -\log\left(\frac{1}{W}\right) = -\log(p_w)$,
693 where we have dropped Boltzmann's constant since the connection between molecular
694 arrangements and temperature is not important. Because each arrangement is equally likely,
695 then $\frac{1}{W}$ is also the expected value of p_w and we can write the entropy as
696 $S = -E[\log(p_w)] = -\sum_w p_w \log(p_w)$, which is Shannon entropy. By definition, Shannon entropy
697 is also the expected value of the information content: $E[I] = -\sum_w p_w \log(p_w) = \sum_w p_w I(w)$
698 where the information content I of a particular state is $I(w) = -\log(p_w)$. Since we assumed
699 each arrangement is equally likely, then the expected value of the information content is also
700 the information content of each arrangement. Therefore, the information content of a *cis*-
701 regulatory sequence can be written as $I = -\log(p_w) = \log W$. We use log base 2 to express
702 the information content in bits.

703

704 With this metric, *cis*-regulatory sequences with higher predicted TF occupancies generally have
705 higher information content. Sequences with higher TF diversity have higher information content
706 than lower diversity sequences with the same predicted occupancy. Thus our metric captures
707 the effects of both TF diversity and total TF occupancy. For example, consider hypothetical TFs
708 A, B, and C. If motifs for only one TF are in a sequence, then W is always 1 and the information
709 content is always zero (regardless of total occupancy). The simplest case for non-zero
710 information content is one motif for A, one motif for B, and zero motifs for C (1-1-0). Then
711 $W = \frac{2!}{1!1!} = 2$ and $I = 1$ bit. If we increase predicted occupancy by adding a motif for A (2-1-0),
712 then $W = \frac{3!}{2!1!} = 3$ and $I = 1.6$ bits, which is approximately the information content of silencers
713 and inactive sequences. If we increase predicted occupancy again and add a second motif for B
714 (2-2-0), then $W = \frac{4!}{2!2!} = 6$ and $I = 2.6$ bits, which is approximately the information content of
715 strong enhancers. If instead of increasing predicted occupancy, we instead increase diversity by
716 replacing a motif for A with a motif for C (1-1-1), then $W = \frac{3!}{1!1!1!} = 6$ and once again $I = 2.6$ bits,
717 which is higher than the lower diversity case (2-1-0).

718

719 According to Wunderlich and Mirny (Wunderlich & Mirny, 2009), the probability of observing k
720 total motifs for m different TFs in a w bp window is $P(k) \sim \text{Poisson}(k; \lambda)$, where $\lambda = pmw$ and p
721 is the probability of finding a spurious motif in the genome. The expected number of windows
722 with k total motifs in a genome of length N is thus $E(k) = P(k) \cdot N$. In mammals, $N \approx 10^9$ and
723 Wunderlich and Mirny find that $p = 0.00025$ for multicellular eukaryotes. For $m = 3$ TFs and a
724 $w = 164$ bp window (which is the size of our sequences), $\lambda = 0.123$ and $E(5) = 1.6$, meaning
725 that five total motifs for three different TFs specifies an approximately unique 164 bp location in
726 a mammalian genome. Five total motifs for three different TFs can be achieved in two ways:

727 three motifs for A, one for B, and one for C (3-1-1), or two motifs for A, two for B, and one for C
728 (2-2-1). In the case of 3-1-1, $W = \frac{5!}{3!1!1!} = 20$ and $I = 4.3$ bits. In the case of 2-2-1,
729 $W = \frac{5!}{2!2!1!} = 30$ and $I = 4.9$ bits.

730

731 **Machine Learning**

732 The k -mer SVM was fit using gkmSVM (Ghandi et al., 2014). All other machine learning,
733 including cross-validation, logistic regression, and computing ROC and PR curves, was
734 performed using scikit-learn version 0.19.1 (Pedregosa et al., 2011). We wrote custom Python
735 wrappers for gkmSVM to allow for interfacing between the C++ binaries and the rest of our
736 workflow. We ran gkmSVM with the parameters -l 6 -k 6 -m 1. To estimate model performance,
737 all models were fit with stratified five-fold cross-validation after shuffling the order of sequences.
738 For the TF occupancy logistic regression model, we used L2 regularization. We selected the
739 regularization parameter C by performing grid search with five-fold cross-validation on the
740 values 10^{-4} , 10^{-3} , 10^{-2} , 10^{-1} , 1, 10^1 , 10^2 , 10^3 , 10^4 and selecting the value that maximized the F1
741 score. The optimal value of C was 0.01, which we used as the regularization strength when
742 assessing the performance of the model with other feature sets.

743

744 To assess the performance of the logistic regression model, we randomly sampled 8 PWMs
745 from the HOCOMOCO database and computed the predicted occupancy of each TF on each
746 sequence. We then fit a new logistic regression model with these features and repeated this
747 procedure 100 times to generate a background distribution of model performances.

748

749 To generate de novo motifs from the SVM, we generated all 6-mers and scored them against
750 the SVM. We then ran the `svmw_emalign.py` script from `gkmSVM` on the *k*-mer scores with the
751 parameters `-n 10 -f 2 -m 4` and a PWM length of 6, and then used TOMTOM to compare them
752 to the database from our motif analysis.

753

754 **Other Data Sources**

755 We used our previously published library (White et al., 2013) as an independent test set for our
756 machine learning models. We defined strong enhancers as ChIP-seq peaks that were above the
757 95th percentile of all scrambled sequences. There was no basal promoter construct in this
758 library, so instead we defined silencers as ChIP-seq peaks that were at least two-fold below the
759 \log_2 mean of all scrambled sequences.

760

761 Previously published ChIP-seq data for NRL (Hao et al., 2012) that was re-processed by
762 Hughes et al. (Hughes et al., 2017) and MEF2D (Andzelm et al., 2015) was used to annotate
763 sequences for *in vivo* TF binding. We converted peaks to mm10 coordinates using the UCSC
764 liftOver tool and then used Bedtools to intersect peaks with our library.

765 **Code and Data Availability**

766 The pJK01 and pJK03 plasmids have been deposited with AddGene. Raw sequencing data and
767 barcode counts have been uploaded to the NCBI GEO database under accession GSE165812.
768 All processed activity data, predicted occupancy, and information content values are available in
769 the supplementary material. All code for data processing, analysis, and visualization is available
770 on Github at <https://github.com/barakcohenlab/CRX-Information-Content>.

771 Acknowledgements

772 We thank Gary Stormo and members of the Cohen Lab for critically reading the manuscript and
773 helpful discussions; Philip A. Ruzycski, Andrew E.O. Hughes, and Timothy J. Cherry for providing
774 processed ChIP-seq and RNA-seq data; and Jessica Hoisington-Lopez and MariaLynn Crosby
775 from the DNA Sequencing Innovation Lab for assistance with high-throughput sequencing. RZF
776 was partially supported by a Kirschstein National Research Service Award from the National
777 Institutes for Health, F31HG011431. This work was supported by grants from the National
778 Institutes for Health, R01GM121755 (to MAW), R01EY027784 (to BAC), and R01 EY025196
779 and R01 EY03075 (to JCC).

780 Competing interests

781 The authors declare no competing interests.

782 Author Contributions

783 RZF: Conceptualization; Methodology; Software; Formal analysis; Investigation; Data curation;
784 Visualization; Funding acquisition; Writing - original draft

785 DMG: Investigation

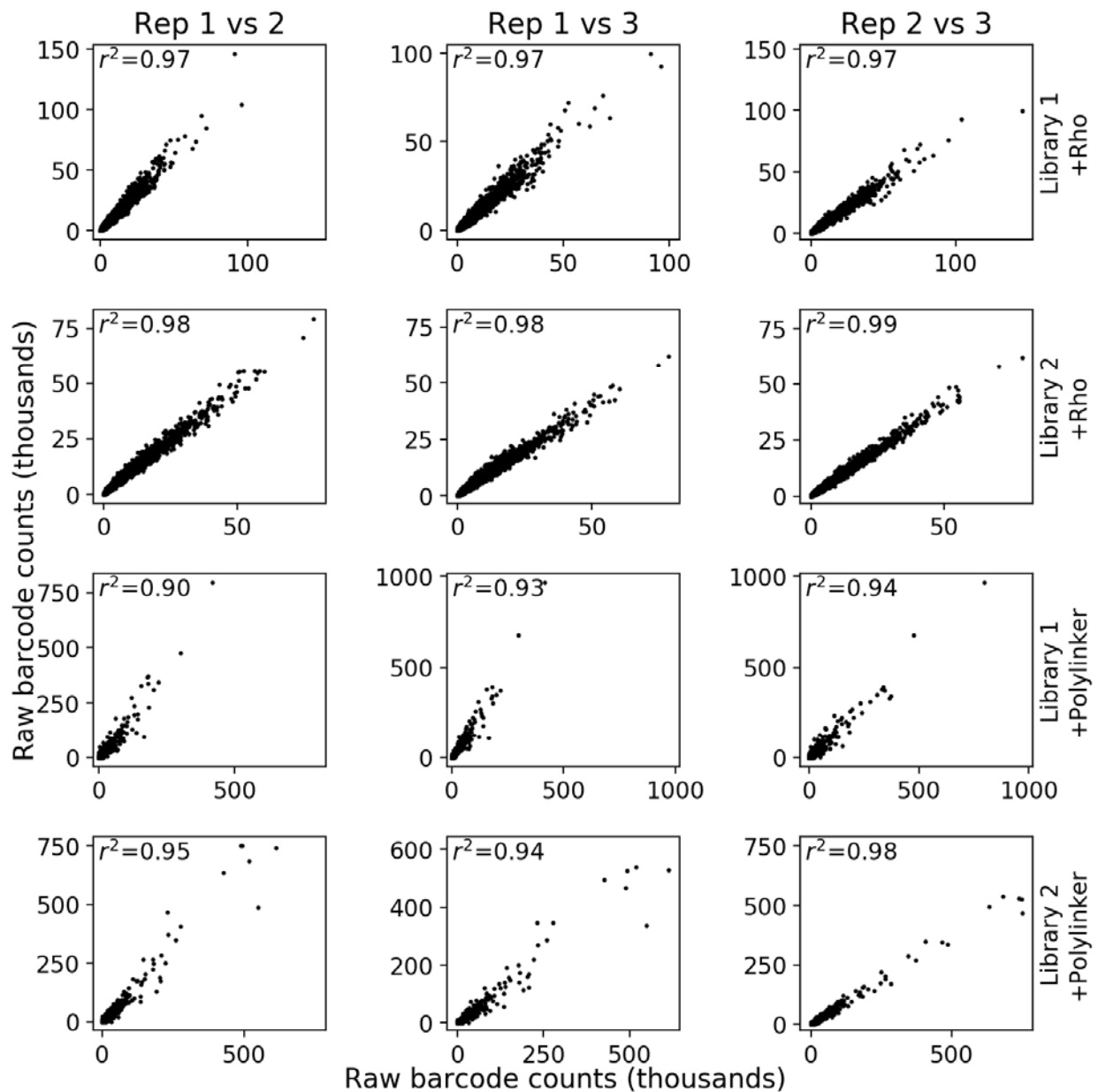
786 CAM: Investigation

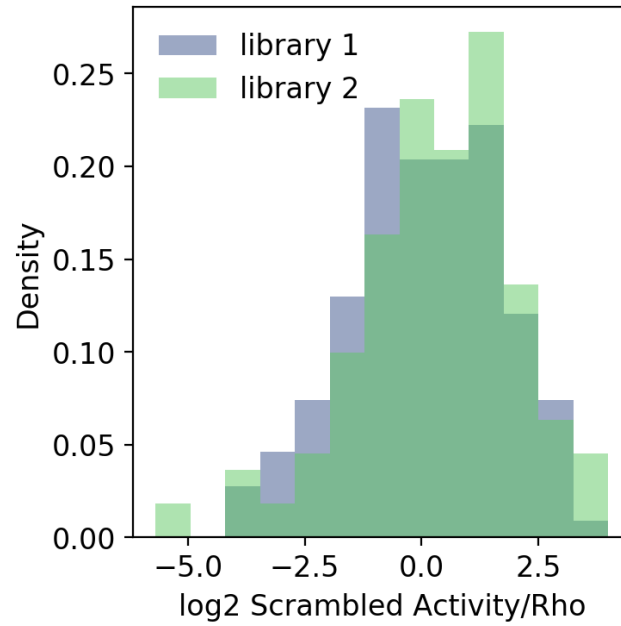
787 JCC: Supervision; Funding acquisition; Writing - original draft

788 BAC: Conceptualization; Methodology; Supervision; Funding acquisition; Writing - original draft

789 MAW: Conceptualization; Methodology; Supervision; Funding acquisition; Writing - original draft

790 Figure Supplements





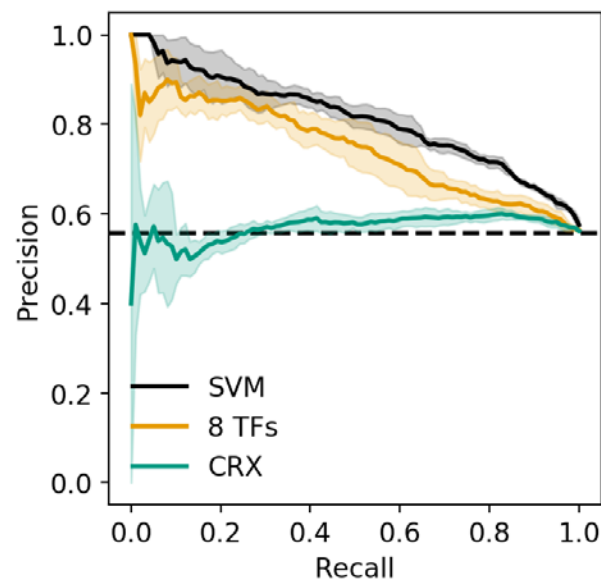
796

797 **Figure 1—figure supplement 2: Calibration of MPRA libraries with the *Rho* promoter.**

798 Probability density histogram of the same 150 scrambled sequences in two libraries after

799 normalizing to the basal *Rho* promoter.

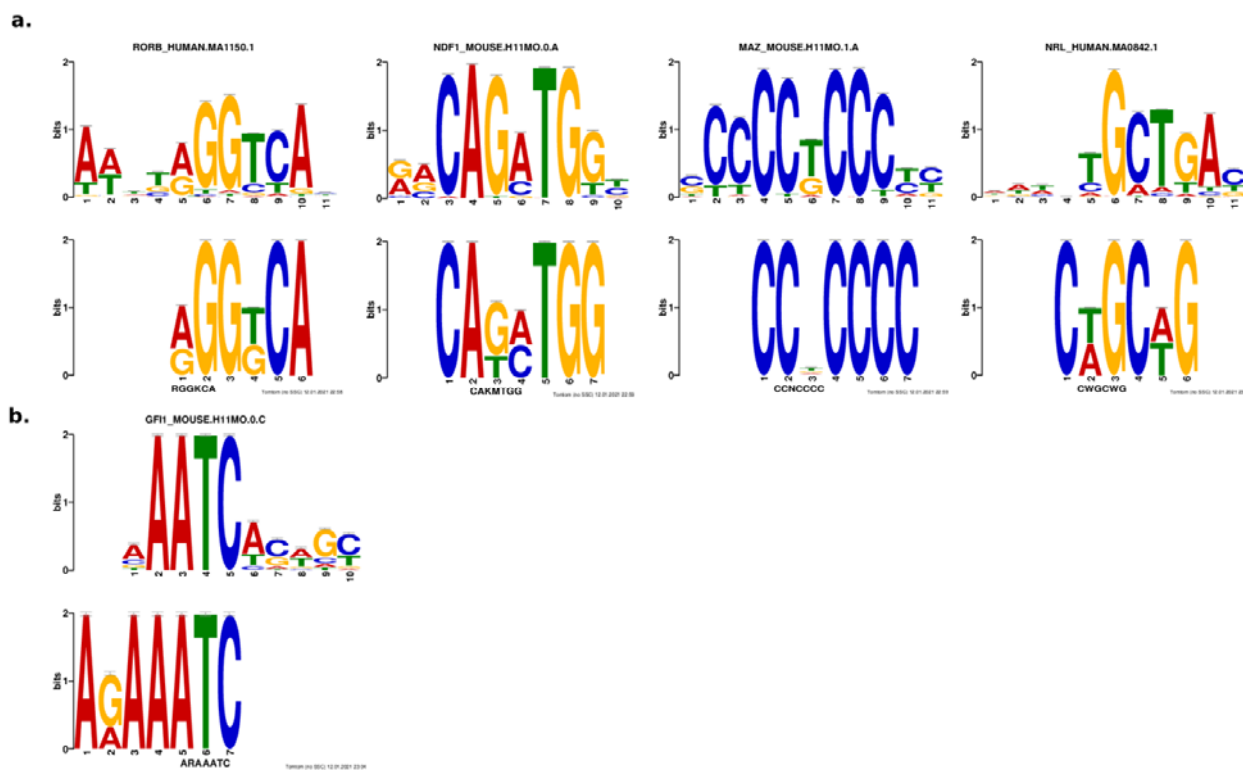
800



801

802 **Figure 2—figure supplement 1: Precision recall curve for strong enhancer vs. silencer**
803 **classifiers.** Solid black, 6-mer SVM; orange, 8 TF predicted occupancy logistic regression;
804 aqua, predicted CRX occupancy logistic regression; dashed black, chance; shaded area, 1
805 standard deviation based on five-fold cross-validation.

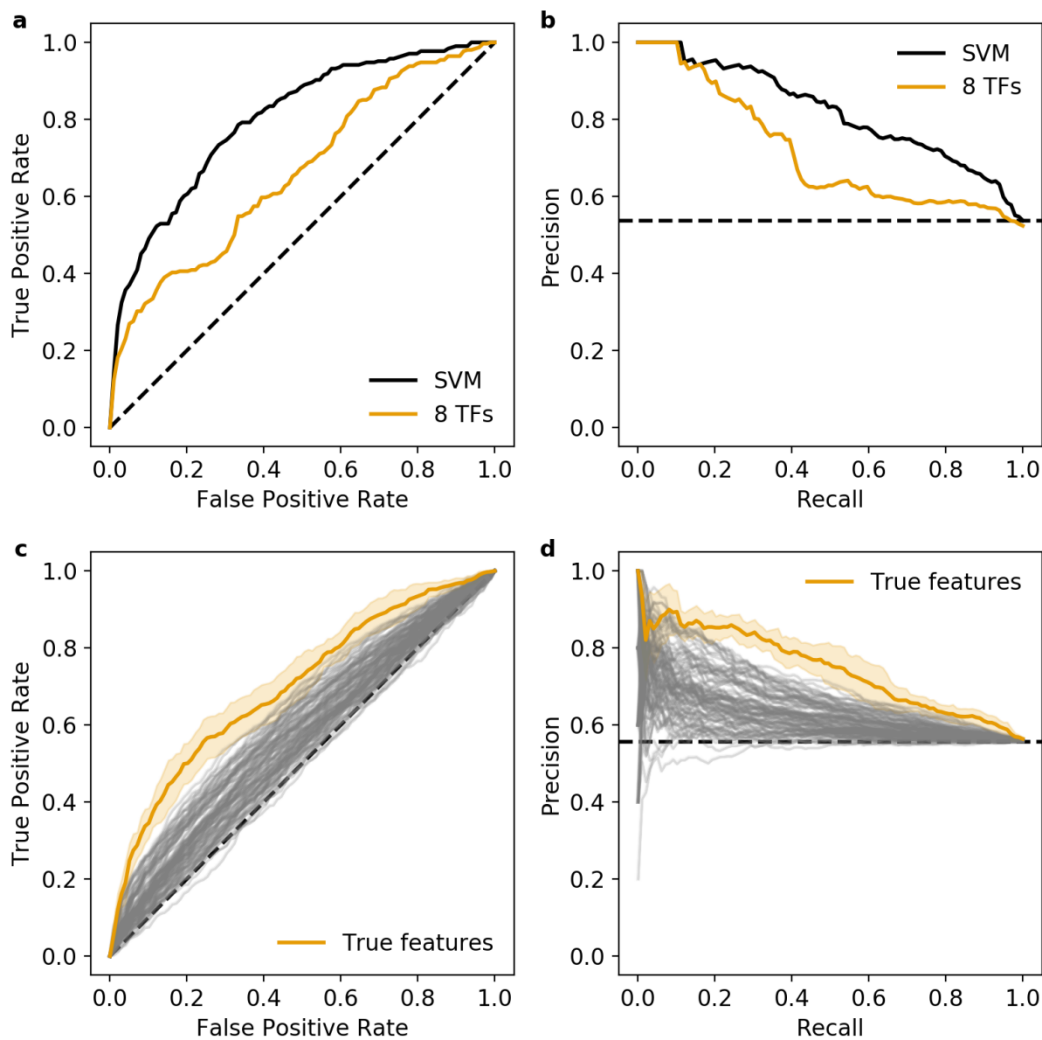
806



807

808 **Figure 2—figure supplement 2: Results from *de novo* motif analysis.** Motifs enriched in
809 strong enhancers **(a)** and silencers **(b)**. Bottom, *de novo* motif identified with DREME; top,
810 matched known motif identified with TOMTOM.

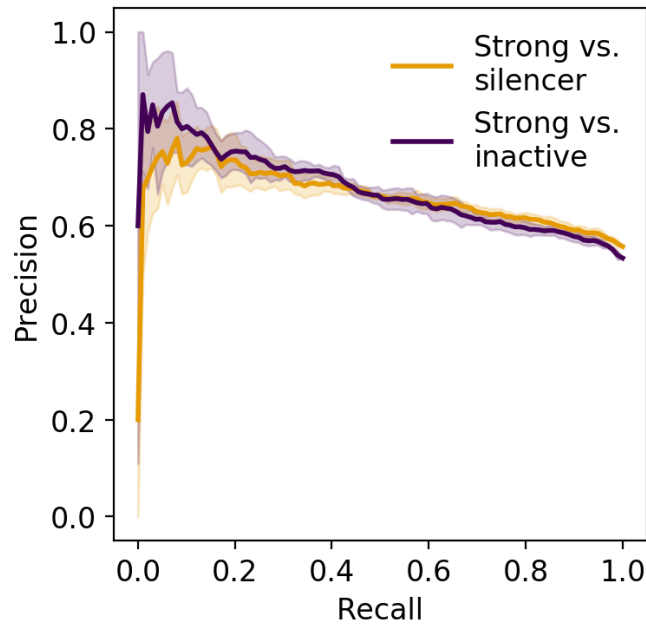
811



812

813 **Figure 2—figure supplement 3: Additional validation of the 8 TF predicted occupancy**
814 **logistic regression model. a and b)** Predictions of the 6-mer SVM (black) and 8 TF predicted
815 occupancy logistic regression model (orange) on an independent test set. **c and d)** Null
816 distribution of 100 logistic regression models trained using randomly selected motifs (grey)
817 compared to the true features (orange). Shaded area, 1 standard deviation based on five-fold
818 cross-validation. **a and c)** Receiver operating characteristic, **b and d)** precision recall curve.
819 Dashed black line represents chance in all panels.

820



821

822 **Figure 3—figure supplement 1: Precision recall curve of logistic regression classifier**
823 **using information content.** Orange, strong enhancer vs. silencer; indigo, strong enhancer vs.
824 inactive; shaded area, 1 standard deviation based on five-fold cross-validation.

825 Supplementary Files

826 **Supplementary file 1: FASTA file of all sequences in library 1.** Sequences were named
827 using the following nomenclature: “chrom-start-stop_ annotations_ variant”. “Chrom”, “start”, and
828 “stop” correspond to the mm10 genomic coordinates of the sequences in BED format.
829 “Annotations” is a four letter string where the first position indicates CRX binding status (ChIP-
830 seq peak or Unbound), the second position indicates CRX motif status (PWM hit, Shape motif,
831 or Both PWM and shape motif), the third position indicates ATAC-seq status (peak in Rods but
832 not cones, peak in Cones but not rods, peak in both rod and cone Photoreceptors, or peak in
833 None of the above), and the fourth position indicates histone ChIP-seq status (“Enhancer
834 marked” with H3K27Ac⁺H3K4me3⁺, “Promoter marked” with H3K27Ac⁺H3K4me3⁺, Q for
835 H3K27Ac⁻H3K4me3⁺, or Neither mark). “Variant” indicates whether the sequence is genomic

836 (“WT”), mutated CRX motifs (“MUT-allCrxCites”), scrambled shape motif (“MUT-shape”), or a
837 scrambled control (“scrambled”).

838

839 **Supplementary file 2: FASTA file of all sequences in library 2.** Sequences were named as
840 in **Supplementary file 1.**

841

842 **Supplementary file 3: Expression measurements and annotations of all sequences.**

843 Values are tab-delimited. Rows are named based on the sequence name from **Supplementary**
844 **files 1 and 2** without the “variant” information. Columns ending in “_WT” indicate the wild-type
845 sequence with the *Rho* promoter, “_MUT” as the CRX motif mutant sequence with the *Rho*
846 promoter, and “_POLY” as the wild-type sequence with the Polylinker. Sequences with the
847 scrambled shape motif were excluded from the “_MUT” columns. Columns are named as
848 follows: label, the sequence name from **Supplementary files 1 and 2** without the “variant”
849 information; expression, average activity of the sequence, NaN indicates sequence was missing
850 from the plasmid pool; expression_std, standard deviation of activity; expression_reps, number
851 of replicates in which the sequence was measured; expression_pvalue, p-value from Welch’s t-
852 test of log-normal data for the null hypothesis that the activity of the sequence with *Rho* is no
853 different than the *Rho* promoter alone; expression_qvalue, FDR-correction of the p-values;
854 library, which library contains the sequence; expression_log2, log2 average activity of the
855 sequence; group_name, activity classification of the sequence with the *Rho* promoter;
856 plot_color, hex code for visualization; variant, the “variant” portion of the sequence identifier;
857 wt_vs_mut_log2, log2 fold change between the wild-type and mutant version of the sequence,
858 NaN indicates the wild-type and/or mutant version was not measured; wt_vs_mut_pvalue, p-
859 value from Welch’s t-test for the null hypothesis that the wild-type and mutant sequences have

860 the same activity; wt_vs_mut_qvalue, FDR-correction of the p-values; autonomous_activity,
861 boolean value for if the wild-type sequence is autonomous with the Polylinker; crx_bound,
862 nrl_bound, and mef2d_bound, boolean values for if the sequence overlaps a ChIP-seq peak for
863 the corresponding TF; binding_group, string denoting each of the eight possible combinations of
864 CRX, NRL, and MEF2D binding.

865

866 **Supplementary file 4: Predicted occupancy scores for each TF and each sequence.**

867 Values are tab-delimited. Rows are named based on the sequence name from **Supplementary**
868 **files 1 and 2** including the “variant” information. Columns are the predicted occupancy scores
869 for the denoted TF.

870

871 **Supplementary file 5: Information content and related metrics for each sequence.** Values
872 are tab-delimited. Rows are named based on the sequence name from **Supplementary files 1**
873 **and 2**, including the “variant” information. Columns are named as follows: total_occupancy, total
874 predicted occupancy of all 8 TFs; diversity, number of TFs with predicted occupancy above 0.5;
875 entropy, information content (which is also entropy).

876

877 **Supplementary file 6: Primers used in this study.**

878 **References**

879 Alexandre, C., & Vincent, J.-P. (2003). Requirements for transcriptional repression and
880 activation by Engrailed in *Drosophila* embryos. *Development*, 130(4), 729–739.

881 <https://doi.org/10.1242/dev.00286>

- 882 Andzelm, M. M., Cherry, T. J., Harmin, D. A., Boeke, A. C., Lee, C., Hemberg, M., Pawlyk, B.,
883 Malik, A. N., Flavell, S. W., Sandberg, M. A., Raviola, E., & Greenberg, M. E. (2015).
884 MEF2D drives photoreceptor development through a genome-wide competition for tissue-
885 specific enhancers. *Neuron*, *86*(1), 247–263. <https://doi.org/10.1016/j.neuron.2015.02.038>
- 886 Bailey, T. L., Boden, M., Buske, F. A., Frith, M., Grant, C. E., Clementi, L., Ren, J., Li, W. W., &
887 Noble, W. S. (2009). MEME SUITE: tools for motif discovery and searching. *Nucleic Acids*
888 *Research*, *37*(Web Server issue), W202–W208. <https://doi.org/10.1093/nar/gkp335>
- 889 Barolo, S., & Posakony, J. W. (2002). Three habits of highly effective signaling pathways:
890 Principles of transcriptional control by developmental cell signaling. *Genes and*
891 *Development*, *16*(10), 1167–1181. <https://doi.org/10.1101/gad.976502>
- 892 Brand, A. H., Micklem, G., & Nasmyth, K. (1987). A yeast silencer contains sequences that can
893 promote autonomous plasmid replication and transcriptional activation. *Cell*, *51*(5), 709–
894 719. [https://doi.org/10.1016/0092-8674\(87\)90094-8](https://doi.org/10.1016/0092-8674(87)90094-8)
- 895 Chen, J., Rattner, A., & Nathans, J. (2005). The rod photoreceptor-specific nuclear receptor
896 Nr2e3 represses transcription of multiple cone-specific genes. *The Journal of*
897 *Neuroscience: The Official Journal of the Society for Neuroscience*, *25*(1), 118–129.
898 <https://doi.org/10.1523/JNEUROSCI.3571-04.2005>
- 899 Chen, S., Wang, Q.-L., Nie, Z., Sun, H., Lennon, G., Copeland, N. G., Gilbert, D. J., Jenkins, N.
900 A., & Zack, D. J. (1997). Crx, a Novel Otx-like Paired-Homeodomain Protein, Binds to and
901 Transactivates Photoreceptor Cell-Specific Genes. *Neuron*, *19*(5), 1017–1030.
902 [https://doi.org/10.1016/S0896-6273\(00\)80394-3](https://doi.org/10.1016/S0896-6273(00)80394-3)
- 903 Chiang, C., & Ayyanathan, K. (2013). Snail/Gfi-1 (SNAG) family zinc finger proteins in
904 transcription regulation, chromatin dynamics, cell signaling, development, and disease.
905 *Cytokine & Growth Factor Reviews*, *24*(2), 123–131.
906 <https://doi.org/10.1016/j.cytogfr.2012.09.002>
- 907 Corbo, J. C., Lawrence, K. A., Karlstetter, M., Myers, C. A., Abdelaziz, M., Dirkes, W., Weigelt,

- 908 K., Seifert, M., Benes, V., Fritsche, L. G., Weber, B. H. F., & Langmann, T. (2010). CRX
909 CHIP-seq reveals the cis-regulatory architecture of mouse photoreceptors. *Genome*
910 *Research*, 20(11), 1512–1525. <https://doi.org/10.1101/gr.109405.110>
- 911 Crocker, J., Abe, N., Rinaldi, L., McGregor, A. P., Frankel, N., Wang, S., Alsawadi, A., Valenti,
912 P., Plaza, S., Payre, F., Mann, R. S., & Stern, D. L. (2015). Low Affinity Binding Site
913 Clusters Confer Hox Specificity and Regulatory Robustness. *Cell*, 160(1-2), 191–203.
914 <https://doi.org/10.1016/J.CELL.2014.11.041>
- 915 Doni Jayavelu, N., Jajodia, A., Mishra, A., & Hawkins, R. D. (2020). Candidate silencer
916 elements for the human and mouse genomes. *Nature Communications*, 11(1), 1061.
917 <https://doi.org/10.1038/s41467-020-14853-5>
- 918 Dorval, K. M., Bobechko, B. P., Fujieda, H., Chen, S., Zack, D. J., & Bremner, R. (2006). CHX10
919 targets a subset of photoreceptor genes. *The Journal of Biological Chemistry*, 281(2), 744–
920 751. <https://doi.org/10.1074/jbc.M509470200>
- 921 ENCODE Project Consortium, Moore, J. E., Purcaro, M. J., Pratt, H. E., Epstein, C. B., Shores,
922 N., Adrian, J., Kawli, T., Davis, C. A., Dobin, A., Kaul, R., Halow, J., Van Nostrand, E. L.,
923 Freese, P., Gorkin, D. U., Shen, Y., He, Y., Mackiewicz, M., Pauli-Behn, F., ... Weng, Z.
924 (2020). Expanded encyclopaedias of DNA elements in the human and mouse genomes.
925 *Nature*, 583(7818), 699–710. <https://doi.org/10.1038/s41586-020-2493-4>
- 926 Ernst, J., & Kellis, M. (2012). ChromHMM: automating chromatin-state discovery and
927 characterization. *Nature Methods*, 9(3), 215–216. <https://doi.org/10.1038/nmeth.1906>
- 928 Fan, R., Toubal, A., Goñi, S., Drareni, K., Huang, Z., Alzaid, F., Ballaire, R., Ancel, P., Liang, N.,
929 Damdimopoulos, A., Hainault, I., Soprani, A., Aron-Wisnewsky, J., Foufelle, F., Lawrence,
930 T., Gautier, J.-F., Venteclef, N., & Treuter, E. (2016). Loss of the co-repressor GPS2
931 sensitizes macrophage activation upon metabolic stress induced by obesity and type 2
932 diabetes. *Nature Medicine*, 22(7), 780–791. <https://doi.org/10.1038/nm.4114>
- 933 Farley, E. K., Olson, K. M., Zhang, W., Brandt, A. J., Rokhsar, D. S., & Levine, M. S. (2015).

- 934 Suboptimization of developmental enhancers. *Science*, 350(6258), 325–328.
935 <https://doi.org/10.1126/science.aac6948>
- 936 Farley, E. K., Olson, K. M., Zhang, W., Rokhsar, D. S., & Levine, M. S. (2016). Syntax
937 compensates for poor binding sites to encode tissue specificity of developmental
938 enhancers. *Proceedings of the National Academy of Sciences of the United States of*
939 *America*, 113(23), 6508–6513. <https://doi.org/10.1073/pnas.1605085113>
- 940 Freund, C. L., Gregory-Evans, C. Y., Furukawa, T., Papaioannou, M., Looser, J., Ploder, L.,
941 Bellingham, J., Ng, D., Herbrick, J. A. S., Duncan, A., Scherer, S. W., Tsui, L. C., Loutradis-
942 Anagnostou, A., Jacobson, S. G., Cepko, C. L., Bhattacharya, S. S., & McInnes, R. R.
943 (1997). Cone-rod dystrophy due to mutations in a novel photoreceptor-specific homeobox
944 gene (CRX) essential for maintenance of the photoreceptor. *Cell*, 91(4), 543–553.
945 [https://doi.org/10.1016/S0092-8674\(00\)80440-7](https://doi.org/10.1016/S0092-8674(00)80440-7)
- 946 Furukawa, T., Morrow, E. M., & Cepko, C. L. (1997). Crx, a Novel otx-like Homeobox Gene,
947 Shows Photoreceptor-Specific Expression and Regulates Photoreceptor Differentiation.
948 *Cell*, 91(4), 531–541. [https://doi.org/10.1016/S0092-8674\(00\)80439-0](https://doi.org/10.1016/S0092-8674(00)80439-0)
- 949 Ghandi, M., Lee, D., Mohammad-Noori, M., & Beer, M. A. (2014). Enhanced Regulatory
950 Sequence Prediction Using Gapped k-mer Features. *PLoS Computational Biology*, 10(7),
951 e1003711. <https://doi.org/10.1371/journal.pcbi.1003711>
- 952 Gisselbrecht, S. S., Palagi, A., Kurland, J. V., Rogers, J. M., Ozadam, H., Zhan, Y., Dekker, J.,
953 & Bulyk, M. L. (2020). Transcriptional Silencers in Drosophila Serve a Dual Role as
954 Transcriptional Enhancers in Alternate Cellular Contexts. *Molecular Cell*, 77(2), 324–337.
955 <https://doi.org/10.1016/j.molcel.2019.10.004>
- 956 Grass, J. A., Boyer, M. E., Pal, S., Wu, J., Weiss, M. J., & Bresnick, E. H. (2003). GATA-1-
957 dependent transcriptional repression of GATA-2 via disruption of positive autoregulation
958 and domain-wide chromatin remodeling. *Proceedings of the National Academy of Sciences*
959 *of the United States of America*, 100(15), 8811–8816.

- 960 <https://doi.org/10.1073/pnas.1432147100>
- 961 Haeussler, M., Zweig, A. S., Tyner, C., Speir, M. L., Rosenbloom, K. R., Raney, B. J., Lee, C.
962 M., Lee, B. T., Hinrichs, A. S., Gonzalez, J. N., Gibson, D., Diekhans, M., Clawson, H.,
963 Casper, J., Barber, G. P., Haussler, D., Kuhn, R. M., & Kent, W. J. (2019). The UCSC
964 Genome Browser database: 2019 update. *Nucleic Acids Research*, *47*(D1), D853–D858.
965 <https://doi.org/10.1093/nar/gky1095>
- 966 Hao, H., Kim, D. S., Klocke, B., Johnson, K. R., Cui, K., Gotoh, N., Zang, C., Gregorski, J.,
967 Gieser, L., Peng, W., Fann, Y., Seifert, M., Zhao, K., & Swaroop, A. (2012). Transcriptional
968 regulation of rod photoreceptor homeostasis revealed by in vivo NRL targetome analysis.
969 *PLoS Genetics*, *8*(4), e1002649. <https://doi.org/10.1371/journal.pgen.1002649>
- 970 Harris, C. R., Millman, K. J., van der Walt, S. J., Gommers, R., Virtanen, P., Cournapeau, D.,
971 Wieser, E., Taylor, J., Berg, S., Smith, N. J., Kern, R., Picus, M., Hoyer, S., van Kerkwijk,
972 M. H., Brett, M., Haldane, A., Del Río, J. F., Wiebe, M., Peterson, P., ... Oliphant, T. E.
973 (2020). Array programming with NumPy. *Nature*, *585*(7825), 357–362.
974 <https://doi.org/10.1038/s41586-020-2649-2>
- 975 Hennig, A. K., Peng, G.-H., & Chen, S. (2008). Regulation of photoreceptor gene expression by
976 Crx-associated transcription factor network. *Brain Research*, *1192*, 114–133.
977 <https://doi.org/10.1016/J.BRAINRES.2007.06.036>
- 978 Hlawatsch, J., Karlstetter, M., Aslanidis, A., Lückoff, A., Walczak, Y., Plank, M., Böck, J., &
979 Langmann, T. (2013). Sterile alpha motif containing 7 (samd7) is a novel crx-regulated
980 transcriptional repressor in the retina. *PloS One*, *8*(4), e60633.
981 <https://doi.org/10.1371/journal.pone.0060633>
- 982 Hoffman, M. M., Buske, O. J., Wang, J., Weng, Z., Bilmes, J. A., & Noble, W. S. (2012).
983 Unsupervised pattern discovery in human chromatin structure through genomic
984 segmentation. *Nature Methods*, *9*(5), 473–476. <https://doi.org/10.1038/nmeth.1937>
- 985 Hsiao, T. H.-C., Diaconu, C., Myers, C. A., Lee, J., Cepko, C. L., & Corbo, J. C. (2007). The Cis-

- 986 regulatory logic of the mammalian photoreceptor transcriptional network. *PloS One*, 2(7),
987 e643. <https://doi.org/10.1371/journal.pone.0000643>
- 988 Huang, D., Petrykowska, H. M., Miller, B. F., Elnitski, L., & Ovcharenko, I. (2019). Identification
989 of human silencers by correlating cross-tissue epigenetic profiles and gene expression.
990 *Genome Research*, 29(4), 657–667. <https://doi.org/10.1101/gr.247007.118>
- 991 Huang, Z., Liang, N., Goñi, S., Damdimopoulos, A., Wang, C., Ballaire, R., Jager, J., Niskanen,
992 H., Han, H., Jakobsson, T., Bracken, A. P., Aouadi, M., Venteclef, N., Kaikkonen, M. U.,
993 Fan, R., & Treuter, E. (2021). The corepressors GPS2 and SMRT control enhancer and
994 silencer remodeling via eRNA transcription during inflammatory activation of macrophages.
995 *Molecular Cell*, 81(5), 953–968.e9. <https://doi.org/10.1016/j.molcel.2020.12.040>
- 996 Hughes, A. E. O., Enright, J. M., Myers, C. A., Shen, S. Q., & Corbo, J. C. (2017). Cell Type-
997 Specific Epigenomic Analysis Reveals a Uniquely Closed Chromatin Architecture in Mouse
998 Rod Photoreceptors. *Scientific Reports*, 7(January), 43184.
999 <https://doi.org/10.1038/srep43184>
- 1000 Hughes, A. E. O., Myers, C. A., & Corbo, J. C. (2018). A massively parallel reporter assay
1001 reveals context-dependent activity of homeodomain binding sites in vivo. *Genome*
1002 *Research*, 28(10), 1520–1531. <https://doi.org/10.1101/gr.231886.117>
- 1003 Hunter, J. D. (2007). Matplotlib: A 2D Graphics Environment. *Computing in Science &*
1004 *Engineering*, 9(3), 90–95. <https://doi.org/10.1109/MCSE.2007.55>
- 1005 Irie, S., Sanuki, R., Muranishi, Y., Kato, K., Chaya, T., & Furukawa, T. (2015). Rax
1006 Homeoprotein Regulates Photoreceptor Cell Maturation and Survival in Association with
1007 Crx in the Postnatal Mouse Retina. *Molecular and Cellular Biology*, 35(15), 2583–2596.
1008 <https://doi.org/10.1128/MCB.00048-15>
- 1009 Iype, T., Taylor, D. G., Ziesmann, S. M., Garmey, J. C., Watada, H., & Mirmira, R. G. (2004).
1010 The transcriptional repressor Nkx6.1 also functions as a deoxyribonucleic acid context-
1011 dependent transcriptional activator during pancreatic beta-cell differentiation: evidence for

- 1012 feedback activation of the nkx6.1 gene by Nkx6.1. *Molecular Endocrinology*, 18(6), 1363–
1013 1375. <https://doi.org/10.1210/me.2004-0006>
- 1014 Jia, L., Oh, E. C. T., Ng, L., Srinivas, M., Brooks, M., Swaroop, A., & Forrest, D. (2009).
1015 Retinoid-related orphan nuclear receptor RORbeta is an early-acting factor in rod
1016 photoreceptor development. *Proceedings of the National Academy of Sciences of the*
1017 *United States of America*, 106(41), 17534–17539. <https://doi.org/10.1073/pnas.0902425106>
- 1018 Jiang, J., Cai, H., Zhou, Q., & Levine, M. (1993). Conversion of a dorsal-dependent silencer into
1019 an enhancer: evidence for dorsal corepressors. *The EMBO Journal*, 12(8), 3201–3209.
1020 <https://doi.org/10.1002/j.1460-2075.1993.tb05989.x>
- 1021 Johnson, K. D., Kim, S.-I., & Bresnick, E. H. (2006). Differential sensitivities of transcription
1022 factor target genes underlie cell type-specific gene expression profiles. *Proceedings of the*
1023 *National Academy of Sciences of the United States of America*, 103(43), 15939–15944.
1024 <https://doi.org/10.1073/pnas.0604041103>
- 1025 Junion, G., Spivakov, M., Girardot, C., Braun, M., Gustafson, E. H., Birney, E., & Furlong, E. E.
1026 M. (2012). A transcription factor collective defines cardiac cell fate and reflects lineage
1027 history. *Cell*, 148(3), 473–486. <https://doi.org/10.1016/j.cell.2012.01.030>
- 1028 Kelley, D. R., Snoek, J., & Rinn, J. L. (2016). Basset: Learning the regulatory code of the
1029 accessible genome with deep convolutional neural networks. *Genome Research*, 26(7),
1030 990–999. <https://doi.org/10.1101/gr.200535.115>
- 1031 Khan, A., Fornes, O., Stigliani, A., Gheorghe, M., Castro-Mondragon, J. A., van der Lee, R.,
1032 Bessy, A., Chèneby, J., Kulkarni, S. R., Tan, G., Baranasic, D., Arenillas, D. J., Sandelin,
1033 A., Vandepoele, K., Lenhard, B., Ballester, B., Wasserman, W. W., Parcy, F., & Mathelier,
1034 A. (2018). JASPAR 2018: update of the open-access database of transcription factor
1035 binding profiles and its web framework. *Nucleic Acids Research*, 46(D1), D1284.
1036 <https://doi.org/10.1093/nar/gkx1188>
- 1037 Kimura, A., Singh, D., Wawrousek, E. F., Kikuchi, M., Nakamura, M., & Shinohara, T. (2000).

- 1038 Both PCE-1/RX and OTX/CRX interactions are necessary for photoreceptor-specific gene
1039 expression. *The Journal of Biological Chemistry*, 275(2), 1152–1160.
1040 <https://doi.org/10.1074/jbc.275.2.1152>
- 1041 Klemm, S. L., Shipony, Z., & Greenleaf, W. J. (2019). Chromatin accessibility and the regulatory
1042 epigenome. *Nature Reviews. Genetics*, 20(4), 207–220. [https://doi.org/10.1038/s41576-](https://doi.org/10.1038/s41576-018-0089-8)
1043 018-0089-8
- 1044 Koike, C., Nishida, A., Ueno, S., Saito, H., Sanuki, R., Sato, S., Furukawa, A., Aizawa, S.,
1045 Matsuo, I., Suzuki, N., Kondo, M., & Furukawa, T. (2007). Functional roles of Otx2
1046 transcription factor in postnatal mouse retinal development. *Molecular and Cellular Biology*,
1047 27(23), 8318–8329. <https://doi.org/10.1128/MCB.01209-07>
- 1048 Kulakovskiy, I. V., Vorontsov, I. E., Yevshin, I. S., Sharipov, R. N., Fedorova, A. D., Rumynskiy,
1049 E. I., Medvedeva, Y. A., Magana-Mora, A., Bajic, V. B., Papatsenko, D. A., Kolpakov, F. A.,
1050 & Makeev, V. J. (2018). HOCOMOCO: towards a complete collection of transcription factor
1051 binding models for human and mouse via large-scale ChIP-Seq analysis. *Nucleic Acids*
1052 *Research*, 46(D1), D252–D259. <https://doi.org/10.1093/nar/gkx1106>
- 1053 Kwasnieski, J. C., Fiore, C., Chaudhari, H. G., & Cohen, B. A. (2014). High-throughput
1054 functional testing of ENCODE segmentation predictions. *Genome Research*, 24(10), 1595–
1055 1602. <https://doi.org/10.1101/gr.173518.114>
- 1056 Kwasnieski, J. C., Mogno, I., Myers, C. A., Corbo, J. C., & Cohen, B. A. (2012). Complex effects
1057 of nucleotide variants in a mammalian cis-regulatory element. *Proceedings of the National*
1058 *Academy of Sciences of the United States of America*, 109(47), 19498–19503.
1059 <https://doi.org/10.1073/pnas.1210678109>
- 1060 Lee, D., Karchin, R., & Beer, M. A. (2011). Discriminative prediction of mammalian enhancers
1061 from DNA sequence. *Genome Research*, 21(12), 2167–2180.
1062 <https://doi.org/10.1101/gr.121905.111>
- 1063 Lee, J., Myers, C. A., Williams, N., Abdelaziz, M., & Corbo, J. C. (2010). Quantitative fine-tuning

1064 of photoreceptor cis-regulatory elements through affinity modulation of transcription factor
1065 binding sites. *Gene Therapy*, 17(11), 1390–1399. <https://doi.org/10.1038/gt.2010.77>

1066 Lerner, L. E., Peng, G. H., Gribanova, Y. E., Chen, S., & Farber, D. B. (2005). Sp4 is expressed
1067 in retinal neurons, activates transcription of photoreceptor-specific genes, and synergizes
1068 with Crx. *The Journal of Biological Chemistry*, 280(21), 20642–20650.
1069 <https://doi.org/10.1074/jbc.M500957200>

1070 Liu, Y.-R., Laghari, Z. A., Novoa, C. A., Hughes, J., Webster, J. R. M., Goodwin, P. E.,
1071 Wheatley, S. P., & Scotting, P. J. (2014). Sox2 acts as a transcriptional repressor in neural
1072 stem cells. *BMC Neuroscience*, 15, 95. <https://doi.org/10.1186/1471-2202-15-95>

1073 Martínez-Montañés, F., Rienzo, A., Poveda-Huertes, D., Pascual-Ahuir, A., & Proft, M. (2013).
1074 Activator and repressor functions of the Mot3 transcription factor in the osmostress
1075 response of *Saccharomyces cerevisiae*. *Eukaryotic Cell*, 12(5), 636–647.
1076 <https://doi.org/10.1128/EC.00037-13>

1077 McKinney, W. (2010). Data structures for statistical computing in python. *Proceedings of the 9th*
1078 *Python in Science Conference*, 445, 51–56.
1079 <http://conference.scipy.org/proceedings/scipy2010/pdfs/mckinney.pdf>

1080 Mears, A. J., Kondo, M., Swain, P. K., Takada, Y., Bush, R. A., Saunders, T. L., Sieving, P. A.,
1081 & Swaroop, A. (2001). Nrl is required for rod photoreceptor development. *Nature Genetics*,
1082 29(4), 447–452. <https://doi.org/10.1038/ng774>

1083 Mitton, K. P., Swain, P. K., Chen, S., Xu, S., Zack, D. J., & Swaroop, A. (2000). The leucine
1084 zipper of NRL interacts with the CRX homeodomain. A possible mechanism of
1085 transcriptional synergy in rhodopsin regulation. *The Journal of Biological Chemistry*,
1086 275(38), 29794–29799. <https://doi.org/10.1074/jbc.M003658200>

1087 Mitton, K. P., Swain, P. K., Khanna, H., Dowd, M., Apel, I. J., & Swaroop, A. (2003). Interaction
1088 of retinal bZIP transcription factor NRL with Flt3-interacting zinc-finger protein Fiz1:
1089 possible role of Fiz1 as a transcriptional repressor. *Human Molecular Genetics*, 12(4), 365–

- 1090 373. <https://doi.org/10.1093/hmg/ddg035>
- 1091 Morrow, E. M., Furukawa, T., Lee, J. E., & Cepko, C. L. (1999). NeuroD regulates multiple
1092 functions in the developing neural retina in rodent. *Development*, *126*(1), 23–36.
1093 <https://www.ncbi.nlm.nih.gov/pubmed/9834183>
- 1094 Murphy, D. P., Hughes, A. E., Lawrence, K. A., Myers, C. A., & Corbo, J. C. (2019). Cis-
1095 regulatory basis of sister cell type divergence in the vertebrate retina. *eLife*, *8*, e48216.
1096 <https://doi.org/10.7554/eLife.48216>
- 1097 Ngan, C. Y., Wong, C. H., Tjong, H., Wang, W., Goldfeder, R. L., Choi, C., He, H., Gong, L., Lin,
1098 J., Urban, B., Chow, J., Li, M., Lim, J., Philip, V., Murray, S. A., Wang, H., & Wei, C.-L.
1099 (2020). Chromatin interaction analyses elucidate the roles of PRC2-bound silencers in
1100 mouse development. *Nature Genetics*, *52*(3), 264–272. [https://doi.org/10.1038/s41588-020-](https://doi.org/10.1038/s41588-020-020-0581-x)
1101 [0581-x](https://doi.org/10.1038/s41588-020-0581-x)
- 1102 Pang, B., & Snyder, M. P. (2020). Systematic identification of silencers in human cells. *Nature*
1103 *Genetics*, *52*(3), 254–263. <https://doi.org/10.1038/s41588-020-0578-5>
- 1104 Parker, D. S., White, M. A., Ramos, A. I., Cohen, B. A., & Barolo, S. (2011). The cis-regulatory
1105 logic of Hedgehog gradient responses: key roles for gli binding affinity, competition, and
1106 cooperativity. *Science Signaling*, *4*(176), ra38. <https://doi.org/10.1126/scisignal.2002077>
- 1107 Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., Blondel, M.,
1108 Prettenhofer, P., Weiss, R., Dubourg, V., & Others. (2011). Scikit-learn: Machine learning in
1109 Python. *The Journal of Machine Learning Research*, *12*, 2825–2830.
1110 <http://www.jmlr.org/papers/volume12/pedregosa11a/pedregosa11a.pdf>
- 1111 Peng, G. H., Ahmad, O., Ahmad, F., Liu, J., & Chen, S. (2005). The photoreceptor-specific
1112 nuclear receptor Nr2e3 interacts with Crx and exerts opposing effects on the transcription
1113 of rod versus cone genes. *Human Molecular Genetics*, *14*(6), 747–764.
1114 <https://doi.org/10.1093/hmg/ddi070>
- 1115 Petrykowska, H. M., Vockley, C. M., & Elnitski, L. (2008). Detection and characterization of

- 1116 silencers and enhancer-blockers in the greater CFTR locus. *Genome Research*, 18(8),
1117 1238–1246. <https://doi.org/10.1101/gr.073817.107>
- 1118 Phillips, R., Kondev, J., Theriot, J., & Garcia, H. (2012). *Physical Biology of the Cell*. Garland
1119 Science. <https://play.google.com/store/books/details?id=t2SzDwAAQBAJ>
- 1120 Quinlan, A. R., & Hall, I. M. (2010). BEDTools: a flexible suite of utilities for comparing genomic
1121 features. *Bioinformatics*, 26(6), 841–842. <https://doi.org/10.1093/bioinformatics/btq033>
- 1122 Rachmin, I., Amsalem, E., Golomb, E., Beerli, R., Gilon, D., Fang, P., Nechushtan, H., Kay, G.,
1123 Guo, M., Yiqing, P. L., Foo, R. S.-Y., Fisher, D. E., Razin, E., & Tshori, S. (2015). FHL2
1124 switches MITF from activator to repressor of Erbin expression during cardiac hypertrophy.
1125 *International Journal of Cardiology*, 195, 85–94. <https://doi.org/10.1016/j.ijcard.2015.05.108>
- 1126 Rister, J., Razzaq, A., Boodram, P., Desai, N., Tsanis, C., Chen, H., Jukam, D., & Desplan, C.
1127 (2015). Single-base pair differences in a shared motif determine differential Rhodopsin
1128 expression. *Science*, 350(6265), 1258–1261. <https://doi.org/10.1126/science.aab3417>
- 1129 Roger, J. E., Hiriyanna, A., Gotoh, N., Hao, H., Cheng, D. F., Ratnapriya, R., Kautzmann, M.-A.
1130 I., Chang, B., & Swaroop, A. (2014). OTX2 loss causes rod differentiation defect in CRX-
1131 associated congenital blindness. *The Journal of Clinical Investigation*, 124(2), 631–643.
1132 <https://doi.org/10.1172/JCI72722>
- 1133 Ruzycski, P. A., Zhang, X., & Chen, S. (2018). CRX directs photoreceptor differentiation by
1134 accelerating chromatin remodeling at specific target sites. *Epigenetics & Chromatin*, 11(1),
1135 42. <https://doi.org/10.1186/s13072-018-0212-2>
- 1136 Samee, M. A. H., Bruneau, B. G., & Pollard, K. S. (2019). A De Novo Shape Motif Discovery
1137 Algorithm Reveals Preferences of Transcription Factors for DNA Shape Beyond Sequence
1138 Motifs. *Cell Systems*, 8(1), 27–42.e6. <https://doi.org/10.1016/j.cels.2018.12.001>
- 1139 Sanuki, R., Omori, Y., Koike, C., Sato, S., & Furukawa, T. (2010). Panky, a novel photoreceptor-
1140 specific ankyrin repeat protein, is a transcriptional cofactor that suppresses CRX-regulated
1141 photoreceptor genes. *FEBS Letters*, 584(4), 753–758.

- 1142 <https://doi.org/10.1016/j.febslet.2009.12.030>
- 1143 Segert, J. A., Gisselbrecht, S. S., & Bulyk, M. L. (2021). Transcriptional Silencers: Driving Gene
1144 Expression with the Brakes On. *Trends in Genetics: TIG*.
1145 <https://doi.org/10.1016/j.tig.2021.02.002>
- 1146 Sethi, A., Gu, M., Gumusgoz, E., Chan, L., Yan, K.-K., Rozowsky, J., Barozzi, I., Afzal, V.,
1147 Akiyama, J. A., Plajzer-Frick, I., Yan, C., Novak, C. S., Kato, M., Garvin, T. H., Pham, Q.,
1148 Harrington, A., Mannion, B. J., Lee, E. A., Fukuda-Yuzawa, Y., ... Gerstein, M. (2020).
1149 Supervised enhancer prediction with epigenetic pattern recognition and targeted validation.
1150 *Nature Methods*, 17(8), 807–814. <https://doi.org/10.1038/s41592-020-0907-8>
- 1151 Spitz, F., & Furlong, E. E. M. (2012). Transcription factors: from enhancer binding to
1152 developmental control. *Nature Reviews. Genetics*, 13(9), 613–626.
1153 <https://doi.org/10.1038/nrg3207>
- 1154 Srinivas, M., Ng, L., Liu, H., Jia, L., & Forrest, D. (2006). Activation of the blue opsin gene in
1155 cone photoreceptor development by retinoid-related orphan receptor beta. *Molecular*
1156 *Endocrinology*, 20(8), 1728–1741. <https://doi.org/10.1210/me.2005-0505>
- 1157 Stampfel, G., Kazmar, T., Frank, O., Wienerroither, S., Reiter, F., & Stark, A. (2015).
1158 Transcriptional regulators form diverse groups with context-dependent regulatory functions.
1159 *Nature*, 528(7580), 147–151. <https://doi.org/10.1038/nature15545>
- 1160 Tareen, A., & Kinney, J. B. (2020). Logomaker: beautiful sequence logos in Python.
1161 *Bioinformatics*, 36(7), 2272–2274. <https://doi.org/10.1093/bioinformatics/btz921>
- 1162 Virtanen, P., Gommers, R., Oliphant, T. E., Haberland, M., Reddy, T., Cournapeau, D.,
1163 Burovski, E., Peterson, P., Weckesser, W., Bright, J., van der Walt, S. J., Brett, M., Wilson,
1164 J., Millman, K. J., Mayorov, N., Nelson, A. R. J., Jones, E., Kern, R., Larson, E., ... SciPy
1165 1.0 Contributors. (2020). SciPy 1.0: fundamental algorithms for scientific computing in
1166 Python. *Nature Methods*, 17(3), 261–272. <https://doi.org/10.1038/s41592-019-0686-2>
- 1167 Wang, S., Sengel, C., Emerson, M. M., & Cepko, C. L. (2014). A gene regulatory network

1168 controls the binary fate decision of rod and bipolar cells in the vertebrate retina.
1169 *Developmental Cell*, 30(5), 513–527. <https://doi.org/10.1016/j.devcel.2014.07.018>

1170 Webber, A. L., Hodor, P., Thut, C. J., Vogt, T. F., Zhang, T., Holder, D. J., & Petrukhin, K.
1171 (2008). Dual role of Nr2e3 in photoreceptor development and maintenance. *Experimental*
1172 *Eye Research*, 87(1), 35–48. <https://doi.org/10.1016/j.exer.2008.04.006>

1173 White, M. A., Kwasnieski, J. C., Myers, C. A., Shen, S. Q., Corbo, J. C., & Cohen, B. A. (2016).
1174 A Simple Grammar Defines Activating and Repressing cis-Regulatory Elements in
1175 Photoreceptors. *Cell Reports*, 17(5), 1247–1254.
1176 <https://doi.org/10.1016/j.celrep.2016.09.066>

1177 White, M. A., Myers, C. A., Corbo, J. C., & Cohen, B. A. (2013). Massively parallel in vivo
1178 enhancer assay reveals that highly local features determine the cis-regulatory function of
1179 ChIP-seq peaks. *Proceedings of the National Academy of Sciences of the United States of*
1180 *America*, 110(29), 11952–11957. <https://doi.org/10.1073/pnas.1307449110>

1181 Wunderlich, Z., & Mirny, L. A. (2009). Different gene regulation strategies revealed by analysis
1182 of binding motifs. *Trends in Genetics: TIG*, 25(10), 434–440.
1183 <https://doi.org/10.1016/j.tig.2009.08.003>

1184 Yang, Z., Ding, K., Pan, L., Deng, M., & Gan, L. (2003). Math5 determines the competence
1185 state of retinal ganglion cell progenitors. *Developmental Biology*, 264(1), 240–254.
1186 <https://doi.org/10.1016/j.ydbio.2003.08.005>

1187 Zhao, Y., Granas, D., & Stormo, G. D. (2009). Inferring binding energies from selected binding
1188 sites. *PLoS Computational Biology*, 5(12), e1000590.
1189 <https://doi.org/10.1371/journal.pcbi.1000590>

1190 Zhou, J., & Troyanskaya, O. G. (2015). Predicting effects of noncoding variants with deep
1191 learning-based sequence model. *Nature Methods*, 12(10), 931–934.
1192 <https://doi.org/10.1038/nmeth.3547>

1193