

# UC Berkeley

## UC Berkeley Previously Published Works

### Title

Massively parallel cis-regulatory analysis in the mammalian central nervous system

### Permalink

<https://escholarship.org/uc/item/13d9s17j>

### Journal

Genome Research, 26(2)

### ISSN

1088-9051

### Authors

Shen, Susan Q  
Myers, Connie A  
Hughes, Andrew EO  
[et al.](#)

### Publication Date

2016-02-01

### DOI

10.1101/gr.193789.115

Peer reviewed

1 **Massively parallel *cis*-regulatory analysis in the mammalian central**  
2 **nervous system**

3  
4 Susan Q. Shen<sup>1</sup>, Connie A. Myers<sup>1</sup>, Andrew E. O. Hughes<sup>1</sup>, Leah C. Byrne<sup>2</sup>, John G. Flannery<sup>2</sup>,  
5 Joseph C. Corbo<sup>1\*</sup>

6  
7 <sup>1</sup>Department of Pathology and Immunology, Washington University in St. Louis, St. Louis, MO

8 <sup>2</sup>Helen Wills Neuroscience Institute, University of California, Berkeley, CA

9  
10 \* To whom correspondence should be addressed. Tel: +1 314 362 6254; Fax: +1 314 362 4096;  
11 Email: [jcorbo@wustl.edu](mailto:jcorbo@wustl.edu)

12 *Running title:* Cis-regulome analysis in the CNS

13 *Key words:* adeno-associated virus, cerebral cortex, *cis*-regulatory elements, DNase-seq,  
14 massively parallel reporter assay, retina

15

16 **ABSTRACT**

17

18 ***Cis*-regulatory elements (CREs, e.g., promoters and enhancers) regulate gene expression,**  
19 **and variants within CREs can modulate disease risk. Next-generation sequencing has**  
20 **enabled the rapid generation of genomic data that predict the locations of CREs, but a**  
21 **bottleneck lies in functionally interpreting these data. To address this issue, massively**  
22 **parallel reporter assays (MPRAs) have emerged, in which barcoded reporter libraries are**  
23 **introduced into cells and the resulting barcoded transcripts are quantified by next-**  
24 **generation sequencing. Thus far, MPRAs have been largely restricted to assaying short**  
25 **CREs in a limited repertoire of cultured cell types. Here, we present two advances that**  
26 **extend the biological relevance and applicability of MPRAs. First, we adapt exome**  
27 **capture technology to instead capture candidate CREs, thereby tiling across the targeted**  
28 **regions and markedly increasing the length of candidate CREs that can be readily**  
29 **assayed. Second, we package the library into adeno-associated virus (AAV), thereby**  
30 **allowing delivery of candidate CREs to target organs *in vivo*. As a proof-of-concept, we**  
31 **introduce a capture library of ~46,000 constructs, corresponding to ~3,500 DNase I**  
32 **hypersensitive (DHS) sites, into the mouse retina by *ex vivo* plasmid electroporation and**  
33 **into the mouse cerebral cortex by *in vivo* AAV injection. We demonstrate tissue-specific**  
34 ***cis*-regulatory activity of DHSs and provide examples of high-resolution truncation**  
35 **mutation analysis for multiplex parsing of CREs. Our approach should enable massively**  
36 **parallel functional analysis of a wide range of CREs in any organ or species that can be**  
37 **infected by AAV, such as non-human primates and human stem cell-derived organoids.**

38

## 39 INTRODUCTION

40

41 *Cis*-regulatory elements (CREs, e.g., promoters and enhancers) are DNA regions that  
42 regulate gene expression, and variants within CREs can contribute to phenotypic diversity,  
43 including disease susceptibility (Wray 2007; Albert and Kruglyak 2015). In the past several  
44 years, vast amounts of genomic data have been generated that predict the locations of  
45 hundreds of thousands of CREs in cell lines and primary tissues (Shen et al. 2012; The  
46 ENCODE Project Consortium 2012; Romanoski et al. 2015). As an avenue for the experimental  
47 validation of these predictions, massively parallel reporter assays (MPRAs, e.g., CRE-seq) have  
48 been developed, in which barcoded plasmid reporters are introduced into cells. Next-generation  
49 sequencing of the resulting barcoded transcripts provides a quantitative measure of CRE  
50 activity (Kwasnieski et al. 2012; Melnikov et al. 2012; Patwardhan et al. 2012; Arnold et al. 2013;  
51 White et al. 2013; Levo and Segal 2014; Shlyueva et al. 2014). Thus far, MPRAs have been  
52 largely restricted to assaying short CRE fragments (<150 bp) synthesized as oligonucleotide  
53 libraries on microarrays (Patwardhan et al. 2009; Baker 2011; White et al. 2013) and delivered  
54 into select mammalian cells accessible by transfection or electroporation. However, CREs are  
55 often hundreds of base pairs in length, and CRE activity depends crucially on the assayed cell  
56 type and its particular complement of transcription factors (TFs) (Davidson 2001). Therefore, we  
57 sought to expand the biological relevance and applicability of MPRAs by increasing the length of  
58 assayed CREs and by widening the repertoire of assayable cell types.

59 The retina and cerebral cortex are two parts of the central nervous system (CNS) with a  
60 shared forebrain origin, whose gene regulatory networks are topics of intense research interest  
61 (Swaroop et al. 2010; Wright et al. 2010; Bae et al. 2015; Nord et al. 2015). The genome-wide  
62 locations of putative CREs have been mapped in both tissues, using methods such as ChIP-seq  
63 and DNase-seq (Visel et al. 2009; Corbo et al. 2010; The ENCODE Project Consortium 2012;  
64 Wilken 2015). Compared to the cortex, the retina is more experimentally amenable to *cis*-  
65 regulatory analysis, in part because its cellular composition is more completely understood  
66 (Livesey and Cepko 2001; London et al. 2013). Electroporation can be used to efficiently deliver  
67 plasmid DNA into rod photoreceptors, which constitute the majority (~80%) of the cells in the  
68 retina (Jeon et al. 1998). We previously conducted CRE-seq by electroporating thousands of  
69 short CREs into the neonatal mouse retina *ex vivo* (Kwasnieski et al. 2012; White et al. 2013).  
70 Although hundreds of putative developmental forebrain enhancers have been assayed with one-  
71 at-a-time transgenic mouse reporter assays (Nord et al. 2013; Visel et al. 2013), never before  
72 has massively parallel *cis*-regulatory analysis been conducted in the mammalian CNS *in vivo*.

73 Here, we sought to overcome current technological hurdles by developing a 'capture-  
74 and-clone' approach for synthesizing CRE-seq libraries with a selectable range of fragment  
75 sizes for targeted *cis*-regulome analysis. As a built-in feature, our approach allows for truncation  
76 mutation analyses, which can identify regions within CREs that are critical for activity. We  
77 furthermore demonstrate the feasibility of conducting *in vivo* CRE-seq in the adult cerebral  
78 cortex by AAV-mediated delivery. Our approach provides a framework for the massively parallel  
79 functional analysis of CREs in a broad repertoire of organs and species *in vivo*.

80

81

## 82 RESULTS

83

### 84 Identification and characterization of candidate CRE regions

85

86 The genomic locations of CREs can be predicted by the patterns of phylogenetic  
87 conservation, the occurrence of transcription factor binding sites, and the presence of various  
88 chromatin features (Levo and Segal 2014; Shlyueva et al. 2014). DNase I hypersensitive (DHS)  
89 sites, which demarcate regions of open chromatin, are one of the most informative predictive

90 features of active CREs (Arvey et al. 2012; Natarajan et al. 2012; Kwasnieski et al. 2014).  
91 Moreover, DNase-seq data for a variety of primary mouse tissues are available as part of the  
92 Mouse ENCODE Project (Yue et al. 2014). To facilitate the direct comparison of a given CRE-  
93 seq library in retina and cerebral cortex, we generated a list of tissue-specific candidate CREs  
94 based on mouse DNase-seq data, corresponding to 1,000 DHS regions from adult retina and  
95 1,000 DHS regions from adult whole brain. Additionally, we included DHSs from two adult  
96 mouse non-neural tissues (1,000 DHSs from heart and 1,000 DHSs from liver) as controls  
97 (Supplemental Table S1). Together, this yielded 4,000 target DHS regions.

98 We first examined the genome-wide distributions of the 4,000 target DHS regions using  
99 GREAT and HOMER, two computational tools for annotating coding and non-coding regions  
100 (Heinz et al. 2010; McLean et al. 2010). The majority (75%) of the DHS regions were distal  
101 elements located more than 10 kb away from the nearest transcriptional start site (TSS)  
102 (Supplemental Fig. S1A). Almost all of the DHS regions fell within introns (46%) or intergenic  
103 regions (45%) (Supplemental Fig. S1B), similar to the genome-wide distribution of DHS regions  
104 in other cell types (Shu et al. 2011). A small number of DHSs (156/4,000 or 4%) were 'promoter-  
105 proximal', i.e., falling within -1 kb to +100 bp relative to the nearest TSS (Supplemental Fig.  
106 S1A). Among these, 77/156 (49%) were retinal DHSs, consistent with the previous observation  
107 that photoreceptor CREs often cluster around TSS's (Corbo et al. 2010).

108 Tissue-specific CREs are often enriched for the binding of TFs important for cell identity  
109 and function (Davidson 2001). Accordingly, we used HOMER (Heinz et al. 2010) to quantify  
110 enrichment of TF motifs in the target regions (Supplemental Table S2). For each set of tissue-  
111 specific target DHSs, we found strong enrichment of putative binding sites for TFs known to be  
112 important in that tissue. For example, among the top statistically significant enrichments for the  
113 retina, brain, heart, and liver DHSs were putative motifs for CRX (Chen et al. 1997; Freund et al.  
114 1997), ASCL1 (Kim et al. 2008), MEF2C (Edmondson et al. 1994), and ONECUT1 (also known  
115 as HNF6) (Clotman et al. 2005), respectively.

116 Since tissue-specific CREs are often associated with genes specifically expressed in the  
117 corresponding tissue (Natarajan et al. 2012; Heinz et al. 2015), we also examined the genes  
118 associated with the target DHSs based on the nearest TSS (Supplemental Table S1). Gene  
119 Ontology (GO) analysis (Carbon et al. 2009) revealed an enrichment for tissue-specific  
120 functions that corresponded to the tissue of DHS origin. For instance, among the top significant  
121 hits for the retina, brain, heart, and liver target DHSs were 'sensory perception of light stimulus',  
122 'nervous system development', 'cardiovascular system development', and 'organic substance  
123 metabolic process', respectively (Supplemental Table S3). Thus, the 4,000 target DHS regions  
124 were likely enriched for tissue-specific CREs.

## 126 **'Capture-and-clone' allows synthesis of targeted *cis*-regulome libraries**

127  
128 To overcome the length restrictions imposed by oligonucleotide array synthesis of CRE  
129 fragments (Cleary et al. 2004), we took advantage of DNA capture, a technique routinely used  
130 for exome sequencing. For exome capture, biotinylated RNA baits are designed to selectively  
131 hybridize with DNA fragments containing sequences of interest, i.e., exonic regions (Gnirke et al.  
132 2009). Here, we adapted this technology to target our CREs of interest (a subset of the putative  
133 '*cis*-regulome') instead of the exome. This approach offers important advantages. First, the input  
134 DNA pool can derive from any genomic DNA source. Hence, the *cis*-regulome of any single  
135 individual or groups of individuals can be assessed. Second, the input DNA pool can be size-  
136 selected for a range of fragment lengths, enabling inclusion of long CREs.

137 Using mouse (C57BL/6J) genomic DNA that was sheared by sonication and then size-  
138 selected to be ~400-500 bp (excluding adapter sequence), we captured with RNA baits tiling the  
139 central 300 bp (which is the median size of DHSs (Natarajan et al. 2012)) of the 4,000 target  
140 DHS regions. We amplified the captured fragments with primers containing restriction sites for

141 cloning into a barcoded vector library (Fig. 1A). Since the cloning was non-directional, both  
142 orientations were roughly equally represented, as expected (49% and 51% of fragments  
143 mapped to the plus and minus strands of the mm9 reference genome, respectively). Paired-end  
144 sequencing revealed a distribution of CRE fragment sizes with a median length of 464 bp (SD =  
145 72 bp) (Fig. 1B). Using two successive rounds of capture, we achieved a very high ‘on-target’  
146 rate: 98.5% of the captured fragments overlapped a target region. The median overlap for on-  
147 target fragments was 282 bp out of the 300 bp target, i.e., 94% of the target region length  
148 (Supplemental Fig. S2). Overall, 3,483 of the 4,000 (87%) targeted regions were represented,  
149 with a median coverage of 8 barcodes per represented region, for a total of 45,670 uniquely  
150 barcoded constructs (Fig. 1C).

151 The distribution of captured fragments across a representative chromosome is shown in  
152 Figure 2A. Notably, many loci exhibited a multiplicity of captured fragments corresponding to a  
153 single target region, resulting in a tiling of the DHS peak, as exemplified in Figure 2B-E. Hence,  
154 the ability to conduct CRE truncation mutation analysis at a given locus is a key built-in feature  
155 of our capture-and-clone approach.

156

### 157 **AAV packaging and delivery preserves CRE-seq library composition**

158

159 We next considered how to expand the repertoire of cell types accessible by CRE-seq.  
160 Whereas efficient plasmid delivery is limited to mitotic cells amenable to chemical transfection or  
161 electroporation (Mortimer et al. 1999; Karra and Dahm 2010), the ideal CRE-seq delivery  
162 vehicle would permit access to a variety of tissues, including post-mitotic tissues, and in a range  
163 of species. We reasoned that adeno-associated virus (AAV), a non-pathogenic virus commonly  
164 used for gene therapy studies, would be suitable for this purpose. AAV causes long-lasting  
165 infection in rodents and primates, and its tissue tropism ranges by serotype from promiscuous to  
166 cell-type selective (Mingozzi and High 2011). Moreover, unlike DNA delivered by lentivirus, the  
167 AAV-delivered DNA remains almost exclusively episomal, thereby permitting *cis*-regulatory  
168 analysis without the insertion site effects associated with integration into the host genome  
169 (McCarty et al. 2004).

170 After cloning in a TATA box-containing minimal promoter-green fluorescent protein (GFP)  
171 cassette (Fig. 1A), we transferred the library into a vector with inverted terminal repeats (ITRs),  
172 which are necessary for AAV packaging (Yan et al. 2005). This yielded the final plasmid library  
173 (Fig. 3A). To deliver the library into the retina, we conducted *ex vivo* electroporation of the  
174 plasmid library into the neonatal mouse retina, as in our past CRE-seq studies (Kwasnieski et al.  
175 2012; White et al. 2013). We generated three biological replicates, each consisting of multiple  
176 electroporated retinas.

177 To deliver the library into the cerebral cortex, we packaged the plasmid library into  
178 AAV9(2YF) and conducted *in vivo* stereotactic injections to infect adult primary motor cortex.  
179 AAV9 is a serotype that exhibits broad tissue tropism, and its tyrosine-mutated derivative  
180 AAV9(2YF) transduces neurons of the CNS with high efficiency and minimal host-mediated  
181 degradation of viral particles (Zhong et al. 2008; Zincarelli et al. 2008; Dalkara et al. 2012;  
182 Aschauer et al. 2013). We generated three biological replicates, each consisting of cerebral  
183 cortex tissue from a single injected mouse.

184 As evidence that AAV packaging and stereotactic injection did not adversely affect the  
185 composition of the library, we observed a strong correlation (Pearson  $r = 0.95$ ) between the  
186 relative abundance of individual barcoded constructs in the retina after delivery of the plasmid  
187 CRE-seq library and in the cerebral cortex after infection with the AAV-packaged CRE-seq  
188 library (Fig. 3B). Furthermore, 76% (34,824/45,670) of the on-target barcodes were ‘well-  
189 represented’ (i.e., had at least 10 raw DNA reads) in all six biological replicates (three replicates  
190 each for retina and cerebral cortex). These 34,824 barcodes covered 97% (3,375/3,483) of the  
191 targeted DHS regions that were represented in the initial post-capture library. These results

192 indicated good preservation of barcode abundance and diversity throughout the procedure, from  
193 the initial post-capture cloning to the delivery of the library.

194 We then examined the tissues histologically for evidence of library expression, as  
195 visualized by fluorescence microscopy. Upon examination of the electroporated retinas, we  
196 observed GFP-positive cells in the outer nuclear layer (ONL) of the retina, where the rod  
197 photoreceptor cell bodies reside (Fig. 3C). Moreover, the GFP-positive cells co-expressed the  
198 rod-specific *Rho*-CBR3-DsRed reporter (Corbo et al. 2010) (Supplemental Fig. S3A). These  
199 findings indicated that the GFP-positive cells were rod photoreceptors, which are the  
200 predominant cell type assayed by neonatal retinal electroporation.

201 Upon histological examination of the AAV-injected brains, we observed bilateral GFP-  
202 positive regions throughout all layers of the cerebral cortex (Fig. 3D), corresponding to GFP-  
203 expressing cells seen under higher magnification (Fig. 3E). Many of the GFP-positive cells were  
204 morphologically consistent with pyramidal neurons, with an apically oriented primary dendrite  
205 and an axon. Furthermore, GFP expression co-localized with RBFOX3 (also known as NeuN)  
206 (Mullen et al. 1992), a widely expressed marker of mature neurons (Supplemental Fig. S3B).  
207 Interestingly, there were bundles of GFP-positive axons crossing the midline in the corpus  
208 callosum (red arrow in Fig. 3D), indicating that interhemispheric projection neurons were among  
209 the cells that expressed the CRE-seq library.

210

### 211 **AAV-mediated CRE-seq demonstrates tissue-specific CRE activity of DHSs *in vivo***

212

213 Given the histological evidence for expression of the library in both tissues, we next  
214 quantified the *cis*-regulatory activity of individual constructs by next-generation sequencing. As  
215 quality control measures, we verified that the samples overall clustered by the assayed tissue  
216 type (retina vs. cerebral cortex). We also observed that the RNA read counts for individual  
217 barcodes were correlated among the three biological replicates for each tissue, although greater  
218 variability was observed among the cerebral cortex samples than the retinal samples  
219 (Supplemental Fig. S4 and Supplemental Table S4).

220 Since tissue-specific DHSs are believed to mediate tissue-specific *cis*-regulatory activity  
221 (Natarajan et al. 2012; Heinz et al. 2015), we first asked whether this was the case. For this  
222 analysis, we assigned the 'overall' *cis*-regulatory activity of a given DHS by averaging across  
223 corresponding barcoded constructs (as well as across biological replicates). Here, we included  
224 the ~3,000 DHSs with at least two barcoded constructs. When we examined the relationship  
225 between the DHS type (i.e., the tissue origin of the DHS) and CRE activity as assayed in the  
226 retina, we observed strong enrichment of retinal DHSs among highly expressed DHSs,  
227 especially among the top ~20% most highly expressed DHSs in the retina (Fig. 4A). Since  
228 averaging across barcoded constructs may not necessarily be the best metric of *cis*-regulatory  
229 activity for a given DHS, we also examined the expression of individual barcoded constructs.  
230 This again revealed the strong preference of the retina for expressing retinal DHSs (Fig. 4B).

231 Similarly, in the cerebral cortex, there was an enrichment of brain DHSs among highly  
232 expressed DHSs, especially among the top ~15% most highly expressed DHSs in the cortex  
233 (Fig. 4A). However, this enrichment was less pronounced than for retina: among the top 15%  
234 most highly expressed DHSs in the retina, 79% were retinal DHSs, while among the top 15%  
235 most highly expressed DHSs in the cerebral cortex, 42% were brain DHSs ( $p < 0.0001$ , Fisher's  
236 exact test). As seen from the individual barcoded constructs (Fig. 4B), there was a clear  
237 preference for brain DHSs among the most active constructs, but there was overall more  
238 promiscuous (less selective) activity of constructs in the cortex. The activity profile of non-brain  
239 DHSs in the cortex was right-shifted (increased) and overlapped to a greater extent with the  
240 activity profile of brain DHSs in the cortex, compared to the activity profile of non-retinal vs.  
241 retinal DHSs in the retina. Overall, these findings indicated that there was tissue-specific *cis*-

242 regulatory activity of DHSs in the retina and the cortex, with the retina exhibiting a stronger  
243 preference for retinal DHSs than the cortex exhibited for brain DHSs.

244

### 245 **Parameters that predict *cis*-regulatory activity**

246

247 We next asked whether certain parameters previously found to be associated with *cis*-  
248 regulatory activity were predictive of high activity in our assay. For each parameter examined in  
249 Figure 5, we considered the top 100 and top 200 most highly expressed DHSs for the tissue-  
250 appropriate DHS type (i.e., for the retina, we restricted our analysis to retinal DHSs, and for the  
251 cerebral cortex, we restricted our analysis to brain DHSs). Corresponding data for the liver and  
252 heart DHSs are provided in Supplemental Fig. S5. We first surveyed expression as a function of  
253 position relative to the center of the DHS target region, within a 1 kb window (Fig. 5A). While  
254 DNase-seq signals had a relatively narrow peak (~300 bp width) (Fig. 5B), *cis*-regulatory activity  
255 in both the retina and cortex had a much broader peak, plateauing in the central ~500 bp. The  
256 breadth of the *cis*-regulatory activity peaks likely reflects the longer length of the captured  
257 fragments (median length of 464 bp) and the large extent of overlap with the central 300 bp of  
258 the DHS regions (median overlap of 94%). Notably, we did not find a substantial relationship  
259 between the length of individual CRE fragments and CRE activity (Supplemental Fig. S6), or  
260 between distance from the nearest TSS and CRE activity (Supplemental Fig. S7).

261 Interestingly, higher DNase-seq scores were significantly associated with higher *cis*-  
262 regulatory activity in the retina but not in the cortex (Fig. 5B). A possible explanation is that the  
263 retinal DNase-seq data primarily reflect the chromatin state of rods, since they constitute the  
264 vast majority of cells in the retina (Jeon et al. 1998), and that the most strongly expressed DHSs  
265 are rod CREs. By comparison, the brain DNase-seq data reflect the chromatin state of a  
266 heterogeneous cell population, and the most strongly expressed DHSs in the cortex may be cell  
267 type-specific CREs highly active in only a subset of cells.

268 Next, we investigated GC content, which has been reported to be elevated within CREs.  
269 This elevation in GC content is thought to favor nucleosome occupancy in tissues where the  
270 CRE is not active, thereby repressing *cis*-regulatory activity in those tissues (Tillo and Hughes  
271 2009; Tillo et al. 2010; Fenouil et al. 2012; Wang et al. 2012; Hughes and Rando 2014). We  
272 previously published an enhancer study, in which short (84 bp) synthetic CREs were cloned  
273 upstream of a photoreceptor-specific proximal promoter. This study revealed a positive  
274 correlation between GC content and enhancer activity in the retina (White et al. 2013). Thus, we  
275 were surprised to find that here, the most active retinal DHSs in the retina had significantly lower  
276 GC content (Fig. 5C). However, a recent CRE-seq study using a minimal promoter also found  
277 lower GC content in highly active enhancers (Kwasnieski et al. 2014). Therefore, GC content  
278 appears to have distinct roles when the CRE acts as an autonomous element with a minimal  
279 promoter or as an enhancer with an active proximal promoter. Brain DHSs had a different  
280 pattern, with markedly elevated GC content centrally, and further increased GC content was  
281 seen among the most active brain DHSs in the cortex (Fig. 5C). The different effects of GC  
282 content in the two tissues may reflect AT-rich vs. GC-rich motifs of tissue-specific TFs, and/or  
283 the distinct preferences of tissue-specific TFs for AT-rich vs. GC-rich 'environments' surrounding  
284 the TF motif (Dror et al. 2015).

285 An ongoing debate in the field of genomics is the degree to which phylogenetic  
286 conservation at the DNA sequence level is an accurate predictor of functional CREs, given that  
287 there is rapid turnover of individual TF binding sites in the course of evolution (Dermitzakis and  
288 Clark 2002; Vierstra et al. 2014). We observed significantly higher vertebrate conservation (as  
289 measured by PhastCons scores (Siepel et al. 2005)) for the most strongly expressed retinal and  
290 brain DHSs in the retina and cortex, respectively. This elevated phylogenetic conservation  
291 occurred primarily within the central ~100 bp of DHSs (Fig. 5D). This distribution of phylogenetic



292 conservation is consistent with the previous observation that highly local (<100 bp) sequences  
293 confer substantial CRE activity (White et al. 2013).

294 We then considered TF motif content, which has been found to be predictive of *cis*-  
295 regulatory activity (Kwasnieski et al. 2014; Blatti et al. 2015). Here, we examined the  
296 enrichment of TF motifs among the DHSs with the highest or lowest activity in the retina and  
297 cortex, regardless of the type of DHS (Fig. 5E and Supplemental Table S5). In the retina, highly  
298 active DHSs were enriched for homeobox, E-box, nuclear receptor (NR), MADS-box, and  
299 CCAAT motifs, while in the cerebral cortex, highly active DHSs were enriched for MADS-box,  
300 zinc finger (ZF), and helix-turn-helix (HTH) motifs.

301 To assess the predictive power of these features (DNase-seq scores, GC content,  
302 PhastCons scores, and TF motifs), we created logistic regression models and visualized their  
303 performance with receiver operating characteristic (ROC) curves, with five-fold cross-validation  
304 to control for over-fitting (Supplemental Table S6). All constructs assayed in each tissue were  
305 classified as 'high' (top ~1% of ~36,000 constructs in retina, or top ~5% of ~39,000 constructs in  
306 cerebral cortex) vs. 'not high'. In the retina, DNase-seq was the single most predictive feature  
307 (AUC = 0.921), reflecting the strong tendency for highly active constructs to be retinal DHSs.  
308 Retinal CRX ChIP-seq peaks (Corbo et al. 2010) performed nearly as well (AUC = 0.892), likely  
309 reflecting the fact that CRX ChIP-seq peaks are essentially a subset of retinal DHSs (Wilken  
310 2015). Interestingly, a model based on 15 TF motifs also performed reasonably well (AUC =  
311 0.785). By comparison, in a prior CRE-seq study conducted in cell lines, a model using 50 TF  
312 motifs attained an AUC of 0.80 (Kwasnieski et al. 2014). The predictive values of GC content  
313 (AUC = 0.521) and PhastCons (AUC = 0.537) were weak. In the cerebral cortex, DNase-seq  
314 was likewise the single most predictive feature (AUC = 0.778). A model based on 13 TF motifs  
315 performed reasonably well (AUC = 0.734), while GC content (AUC = 0.608) and PhastCons  
316 (AUC = 0.659) had modest predictive power in the cortex. Notably, in both tissues, the  
317 combined model performed only slightly better than DNase-seq alone. Overall, these results  
318 reflect the degree of preference of the retina and cerebral cortex for expressing retinal DHSs  
319 and brain DHSs, respectively, while underscoring the importance of TF motifs in specifying CRE  
320 activity. Furthermore, these results underscore the power of open chromatin mapping  
321 techniques such as DNase-seq for identifying functional CREs.

322

### 323 **Tiling of captured fragments allows for truncation mutation analysis**

324

325 The potential for conducting truncation mutation analysis is an attractive and potentially  
326 powerful feature of the capture approach. We therefore sought to determine whether the results  
327 were comparable to those of a previously published 'traditional' one-at-a-time promoter analysis.  
328 NRL is a master regulator of rod photoreceptor development, required both for rod fate  
329 determination and maintenance (Mears et al. 2001; Swaroop et al. 2010). Past studies of the *Nrl*  
330 promoter region identified a 30 bp 'critical region' that is absolutely required for promoter activity.  
331 This critical region contains TF binding sites for CRX and RORB, both of which are required for  
332 *Nrl* expression (Kautzmann et al. 2011; Montana et al. 2011a). Since the *Nrl* promoter contained  
333 a retinal DHS that was targeted in our library, we compared the results of CRE-seq and a  
334 traditional promoter analysis that used fluorescence as a read-out of *cis*-regulatory activity  
335 (Montana et al. 2011a). Since promoters act directionally (Andersson et al. 2014; Duttke et al.  
336 2015), we compared CRE-seq constructs that were oriented in the same direction as the  
337 traditional promoter constructs. We found good agreement between the two assays overall (Fig.  
338 6A), despite differences in construct design (e.g., the CRE-seq constructs contained a minimal  
339 promoter, and the 3' ends of fragments varied). Importantly, both identified the same critical  
340 region within a block of phylogenetic conservation (Montana et al. 2011a). Thus, CRE-seq  
341 truncation analysis recapitulated the results of a traditional truncation mutation analysis.

342 Besides the *Nrl* promoter, we found additional instances of novel truncation mutation  
343 analyses afforded by the capture approach. As seen in Figure 6B, a retinal DHS in the intron of  
344 *Rbm20* showed strong activity in the retina and weak activity in the cortex. Intriguingly, our  
345 assay revealed a 12 bp critical region containing a predicted binding motif for CRX. This motif,  
346 'CTAATCCT' (on the negative strand) is a near-perfect match to the consensus motif,  
347 'CTAATCCC' (Lee et al. 2010).

348 Figure 6C depicts another truncation mutation analysis, this time for two brain DHSs  
349 (labeled '1' and '2') located <0.5 kb apart within an intron of *Bsn* (Bassoon). Bassoon is a  
350 presynaptic protein that is important for neurotransmitter release from glutamatergic (excitatory)  
351 neurons (Altrock et al. 2003). Both of these brain DHSs contained phylogenetically conserved  
352 regions, as observed by PhastCons (Siepel et al. 2005). Interestingly, while both had low *cis*-  
353 regulatory activity in the retina, DHS #1 had low activity in the cerebral cortex, whereas DHS #2  
354 had high activity in the cortex. Furthermore, given the extensive tiling of the region, the  
355 boundaries of activity could be determined at both the 5' and 3' ends of DHS #2.

356 Next, we present a brain DHS region with high *cis*-regulatory activity in the cerebral  
357 cortex (Fig. 6D). A critical region of ~150 bp in length was identified that overlapped a block of  
358 phylogenetic conservation. Incremental loss of bases in this region resulted in progressive  
359 decreases in *cis*-regulatory activity. Within this critical region, two TF motifs were identified: a  
360 consensus E-box motif (recognized by bHLH TFs) (Massari and Murre 2000), immediately next  
361 to a motif recognized by basic region leucine zipper (bZIP) proteins of the AP-1 family (Heinz et  
362 al. 2010). Like neural bHLH proteins, AP-1 family proteins are known to have important roles in  
363 regulating gene expression in the cerebral cortex (Raivich and Behrens 2006; Mongrain et al.  
364 2011).

365 Additional examples of truncation mutation analysis are presented in Supplemental  
366 Figure S8. Overall, we identified 46 retinal DHSs and 13 brain DHSs with examples of  
367 truncation mutation analysis, thus representing 4.6% and 1.3% of the 1000 retinal DHSs and  
368 1000 brain DHSs initially targeted in the library, respectively. We observed that for the loci with  
369 truncation mutation analyses, at least 8 barcoded constructs tiled across the DHS. For DHSs  
370 with at least 8 assayed barcodes, the fraction of loci with truncation mutation analyses was  
371 about 3-fold higher: 46/363 (12.7%) of retinal DHSs and 13/345 (3.8%) of brain DHSs.

372 Truncation mutation analyses rely on assaying long CRE fragments that tile across CRE  
373 regions. Previously, we conducted a CRE-seq enhancer study (White et al. 2013) in which short  
374 (84 bp) CREs (synthesized by oligonucleotide array) were assayed upstream of a rod  
375 photoreceptor-specific proximal promoter. These short CREs corresponded to retinal CRX  
376 ChIP-seq peaks, which are essentially a subset of retinal DHSs (Wilken 2015). Thus, we  
377 wondered whether, for a given CRE, our capture-and-clone approach identified active *cis*-  
378 regulatory sequences beyond the central region tested by the short CRE. Overall, there were  
379 176 CRE regions in the White et al. library that overlapped with assayed regions in the current  
380 library, all of which corresponded to retinal DHSs. Most (141/176 or 80%) regions were more  
381 active as short enhancers than as long autonomous elements (Supplemental Fig. S9A). This is  
382 not surprising, as it is known that some photoreceptor CREs exhibit strong activity as enhancers  
383 but minimal activity as autonomous elements (Corbo et al. 2010). Interestingly, in a minority  
384 (13/176 or 7%) of cases, the long autonomous elements exhibited substantially more activity,  
385 likely because they encompassed functional regions (e.g., critical regions and/or  
386 phylogenetically conserved regions) that were not found within the short CREs, as illustrated in  
387 Supplemental Figure S9B and S9C. Although the comparison of these two studies is limited by  
388 the differences in assay platforms and the small number of shared CREs, these results indicate  
389 that the capture-and-clone approach can provide additional *cis*-regulatory information beyond  
390 that of short CREs.

391 Together, these examples illustrate that CRE-seq multiplex truncation mutation analysis  
392 can identify both known and novel critical regions. In some cases, the spatial resolution is high

393 enough to pinpoint candidate TF motifs required for activity. Thus, our assay has the ability not  
394 only to measure the overall activity of a candidate CRE, but also to demarcate the spatial  
395 boundaries of *cis*-regulatory activity.

396

### 397 **Traditional reporter assays confirm that critical bases identified by CRE-seq truncation** 398 **mutation analysis are required for activity**

399

400 To validate the ability of CRE-seq truncation mutation analysis to identify critical regions  
401 *de novo*, we utilized traditional reporter assays. We previously developed a quantitative  
402 fluorescence reporter assay in retinal explants that accurately measures CRE activity (Montana  
403 et al. 2011b; Kwasnieski et al. 2012). Thus, we selected three retinal DHS loci (including R64,  
404 which is the locus depicted in Figure 6B) with critical regions identified by CRE-seq truncation  
405 mutation analysis to test with the traditional approach (Fig. 7A). These critical regions contained  
406 bioinformatically predicted CRX sites, thus allowing us to test whether these CRX sites were  
407 required for *cis*-regulatory activity.

408 For each locus, we created a 'long' construct, a 'short' construct missing the critical  
409 region, and a 'mutant' construct identical to the 'long' construct except that a single point  
410 mutation was introduced in the predicted CRX site (Fig. 7A). The point mutation was an  
411 adenine-to-cytosine substitution at the fourth position of the CRX motif (thymine-to-guanine in  
412 the reverse orientation), which is predicted to inactivate the CRX site (Supplemental Table S7)  
413 (Lee et al. 2010; White et al. 2013). The constructs were directionally cloned upstream of the  
414 minimal promoter-GFP cassette in a non-AAV vector without barcodes in the 3' UTR, thus  
415 controlling for any effects of orientation, AAV vector sequence, or barcode sequence.

416 Each construct was individually electroporated into multiple retinas and quantified  
417 relative to a loading control, *Rho*-CBR3-DsRed (Fig. 7B). We observed that in each case, the  
418 long construct showed high activity, while the short construct showed extremely low activity.  
419 Notably, the mutant construct exhibited a low level of activity comparable to the activity of the  
420 short construct (Fig. 7C). Thus, for all three loci, we not only verified that the critical regions are  
421 required for activity, but also that these specific CRX sites are required. These experiments  
422 demonstrate that our approach identifies *bona fide* TF binding sites required for activity.

423

424

## 425 **DISCUSSION**

426

427 Here, we described an innovative 'capture-and-clone' approach for synthesizing CRE-  
428 seq libraries. We furthermore demonstrated the feasibility of using AAV-mediated CRE-seq to  
429 conduct massively parallel *cis*-regulatory analysis in the cerebral cortex *in vivo*. By comparing  
430 retina and cerebral cortex, we showed tissue-specific *cis*-regulatory activity of DHSs. By taking  
431 advantage of the truncation mutation analysis afforded by the tiling of captured fragments  
432 across targeted loci, we illustrated high-resolution, multiplex functional parsing of CREs.

433 Previously, high-throughput functional assays of CRE activity had been technologically  
434 limited with regards to the length of CREs that could be readily assayed (Levo and Segal 2014;  
435 Shlyueva et al. 2014). Our capture-and-clone approach provides a strategy for assaying  
436 candidate CREs with lengths of a desired range. Moreover, the capture approach can be used  
437 in conjunction with any existing MPRA-like approach, including those that already rely on DNA  
438 fragmentation (Dickel et al. 2014; Murtha et al. 2014). For example, STARR-seq (Arnold et al.  
439 2013) has been used to assess long DNA fragments obtained by whole-genome shotgun  
440 cloning of the *Drosophila* genome. However, the mouse and human genomes are ~25 times  
441 larger than the fly genome. Moreover, only ~5-10% of the mammalian genome is thought to be  
442 functionally constrained (Graur et al. 2013; Kellis et al. 2014; Rands et al. 2014). Therefore,

443 whole-genome shotgun cloning of mammalian genomes for *cis*-regulatory analysis is impractical.  
444 Instead, capture-and-clone permits targeted *cis*-regulome analysis.

445 We note that another group has recently coupled capture technology to STARR-seq  
446 (i.e., CapSTARR-seq) (Vanhille et al. 2015). Our approach differs from CapSTARR-seq in two  
447 key ways (Supplemental Table S8). First, we achieved higher on-target rates of capture (98.5%  
448 vs. 14%) due to a rigorous capture protocol to avoid non-specific pull-down of off-target DNA  
449 (Gnirke et al. 2009; Lee et al. 2009). Second, we conducted paired-end sequencing of the input  
450 library, whereas CapSTARR-seq mapped only one end of the fragments. Thus, we were able to  
451 harness the potential of capture-and-clone for truncation mutation analysis.

452 Capture-and-clone allows the testing of longer CREs, which presumably harbor more  
453 *cis*-regulatory information. However, there was essentially no correlation between fragment  
454 length and CRE activity. What accounts for this observation? One consideration is that the size  
455 range of assayed CRE fragments was relatively narrow. Another explanation, based on the  
456 truncation mutation analyses, is that some long fragments exhibited low activity due to the  
457 omission of critical regions. A third possibility is that some long CRE fragments included  
458 repressive sequences that decreased activity (Reynolds et al. 2013).

459 The capture-and-clone approach is particularly well suited for screening thousands of  
460 candidate CREs and identifying the most active CREs in a particular tissue of interest, thereby  
461 narrowing the list of CREs that may be relevant to a particular phenotype. For instance,  
462 genome-wide association studies (GWAS) and whole-genome sequencing studies have  
463 generated lists of thousands of disease-associated non-coding variants (Ward and Kellis 2012;  
464 Albert and Kruglyak 2015). To prioritize these lists and thereby accelerate the identification of  
465 causal variants, the locations of the candidate variants can be intersected with the locations of  
466 putative CREs. The *cis*-regulomes of unaffected and affected individuals can then be screened  
467 by capture-and-clone CRE-seq to identify CREs that exhibit the greatest differential activity  
468 between the unaffected and affected groups. Capture-and-clone is thus complementary to CRE-  
469 by-synthesis, which is better suited to precisely measuring the effects of specific variants (Levo  
470 and Segal 2014). Capture-and-clone can be used to assess a broad range of regions in any  
471 organism whose DNA and reference genome are available, although certain types of sequences  
472 are not amenable to targeted capture, namely repetitive regions (due to non-specific pull-down)  
473 and sequences with very high (>65%) or low (<25%) GC content (Mertes et al. 2011).

474 Prior to our study, the implementation of MPRAs in mammalian cells had been almost  
475 exclusively restricted to immortalized cell lines and cultured tissues (Shlyueva et al. 2014). The  
476 only mammalian tissue that had been assayed *in vivo* was the mouse liver, due to its ability to  
477 take up limited amounts of plasmid DNA via a hydrodynamic tail vein assay (Herweijer and  
478 Wolff 2007; Patwardhan et al. 2012). Here, we take a step forward by using AAV to conduct  
479 CRE-seq *in vivo* in the mammalian CNS.

480 One potential drawback of AAV is that packing constraints limit the size of the insert to  
481 less than 4.7 kb (Wu et al. 2010). Lentiviruses have greater carrying capacity (Kumar et al.  
482 2001), but their integration into the host genome poses the risk of integration site *cis*-regulatory  
483 effects (Clark et al. 1994). By contrast, AAV-mediated CRE-seq measures the *cis*-regulatory  
484 potential of elements independent of chromosomal context, thereby interrogating the function of  
485 the DNA sequences themselves. Interestingly, there is evidence that despite being episomal,  
486 the AAV vector is organized into nucleosomes (Penaud-Budloo et al. 2008). Another limitation  
487 of AAV is that the onset of expression is relatively slow, with maximal expression requiring up to  
488 several weeks (Day et al. 2014). This delay is due to the required conversion of the genome  
489 from single-stranded into double-stranded DNA. Recently, self-complementary AAV (scAAV)  
490 serotypes have been developed that exhibit more rapid transgene expression (McCarty 2008).  
491 As novel AAV serotypes for gene therapy continue to emerge (Wu et al. 2006; Daya and Berns  
492 2008), AAV-mediated CRE-seq will become increasingly powerful.

493 Why are some tissue-specific DHSs active and others inactive, even when assayed in  
494 the appropriate tissue? One reason is that DHSs demarcate not only active enhancers but also  
495 other types of regulatory elements (e.g., silencers and insulators) (Gross and Garrard 1988;  
496 Thurman et al. 2012). Here, we used a TATA-box containing minimal promoter to assay the  
497 autonomous *cis*-regulatory activity of the tested elements, rather than a tissue-specific proximal  
498 promoter to assay for enhancer/silencer activity (Butler and Kadonaga 2002). Only a minority  
499 (~10-20%) of mammalian promoters contain TATA boxes (Sandelin et al. 2007). Future use of  
500 tissue-specific proximal promoters may allow for more sensitive assays, especially as enhancer-  
501 promoter compatibility and TATA-box vs. DPE-containing promoters become better understood  
502 (Sandelin et al. 2007; van Arensbergen et al. 2014; Zabidi et al. 2015). Additionally, since some  
503 enhancers become active only in response to particular stimuli (Ostuni et al. 2013; Shlyueva et  
504 al. 2014), environmental perturbations may be necessary to unmask their *cis*-regulatory  
505 potential. Furthermore, the *cis*-regulatory landscape of a given tissue is dynamic across  
506 development, as illustrated by DNase-seq in the developing mouse retina and brain (Wilken  
507 2015). Future CRE-seq experiments at multiple developmental stages will help elucidate the  
508 temporal dynamics of CREs. Nonetheless, even with the TATA-box containing minimal  
509 promoter assayed in steady-state conditions, we demonstrated tissue-specific CRE activity.

510 Assaying autonomous activity and assaying enhancer activity are complementary  
511 approaches, as they appear to reflect different biological activities and properties of a given  
512 CRE. In the current study, we observed that GC content was associated with decreased  
513 autonomous CRE activity in the retina. Given the differences in the assays, this finding does not  
514 contradict our earlier retinal CRE-seq study (White et al. 2013), in which we observed a positive  
515 association between GC content and enhancer activity. In fact, the current result is consistent  
516 with a recent CRE-seq study in which GC content was associated with decreased autonomous  
517 activity of predicted enhancers in cell culture (Kwasnieski et al. 2014)..

518 In our study, the retina exhibited a stronger preference for retinal DHSs than the  
519 cerebral cortex exhibited for brain DHSs. Several explanations are possible. First, the cellular  
520 complexity of the brain is likely a major factor (Wurmbach et al. 2002). A recent DNase-seq  
521 study in the mouse brain observed that DHSs could be found around genes expressed in only a  
522 small percentage of neurons, such as cortical laminar-specific genes (Wilken 2015). Thus, a  
523 given 'brain DHS' may actually be a cell type-specific DHS that is active in a small population of  
524 cells. When averaged over the entire population of assayed cells, the cell type-specific activity  
525 of the DHS may be obscured. For tissues with highly heterogeneous cell populations such as  
526 the cerebral cortex, it should be possible to target specific subpopulations by combining AAV-  
527 mediated CRE-seq with fluorescence-activated cell sorting (FACS) of defined cell types (Okaty  
528 et al. 2011; Gisselbrecht et al. 2013; Dickel et al. 2014). Second, the minimal promoter used in  
529 this study contains a possible weak CRX site, whose affinity is predicted to be ~10% that of the  
530 CRX consensus motif (Chen and Zack 1996; Lee et al. 2010). Lastly, although DNA barcode  
531 representation was similar in the retina and cerebral cortex, the difference in delivery methods  
532 for the two tissues may have been a contributing factor.

533 In summary, we have developed a powerful and efficient strategy for constructing CRE-  
534 seq libraries that extends the size range of the CREs that can readily be assayed, using  
535 targeted *cis*-regulome capture. At the same time, we have demonstrated the feasibility of  
536 conducting CRE-seq *in vivo* in a mammalian tissue using AAV. As new assays for rapidly  
537 identifying the locations of putative cell type-specific CREs are developed, e.g., ATAC-seq  
538 (Buenrostro et al. 2013), our study sets the stage for the high-throughput functional screening of  
539 thousands of candidate CREs in a range of cell types and in a variety of model systems,  
540 including non-human primates and human induced pluripotent stem cell (iPSC)-derived  
541 organoids (Lancaster et al. 2013).

542  
543

544 **METHODS**

545

546 **Animals.** Mice were maintained on a 12-hour light/dark cycle at ~20-22 °C with free access to  
 547 food and water. Neonatal mice were euthanized by decapitation, and adult animals were  
 548 euthanized with CO<sub>2</sub> anesthesia followed by cervical dislocation, unless otherwise stated. All  
 549 experiments were conducted in accordance with the Guide for the Care and Use of Laboratory  
 550 Animals of the National Institutes of Health, and were approved by the Washington University in  
 551 St. Louis Institutional Animal Care and Use Committee.

552

553 **Reference genome.** The mouse reference genome used throughout was mm9.

554

555 **Identification of target tissue-specific DHS peaks.** We downloaded DHS data in narrowPeak  
 556 format from the Mouse ENCODE Project (Yue et al. 2014) for the following tissues (GEO  
 557 sample accessions are listed): whole brain age E14.5 (GSM1014197, replicate 1), whole brain  
 558 age E18.5 (GSM1014184, replicate 1), whole brain age 8 weeks (GSM1014151, replicate 1),  
 559 retina age P1 (GSM1014188), retina age P7 (GSM1014198), retina age 8 weeks  
 560 (GSM1014175), liver age E14.5 (GSM1014183, replicate 1), liver age 8 weeks (GSM1014195,  
 561 replicate 1), lung age 8 weeks (GSM1014194, replicate 1), kidney age 8 weeks (GSM1014193,  
 562 replicate 1), thymus age 8 weeks (GSM1014185, replicate 1), and heart age 8 weeks  
 563 (GSM1014166, replicate 1). We parsed these data using custom Perl scripts, tallying the  
 564 number of reads per 150 bp block across the mouse genome to give a DHS 'score'. We then  
 565 examined the top ~4,000 tissue-specific peaks each for brain age 8 weeks, retina age 8 weeks,  
 566 heart age 8 weeks, and liver age 8 weeks. For a peak to be identified as 'tissue-specific', it was  
 567 required to have a DHS score of >25 in the 8 week tissue of interest and <25 in samples derived  
 568 from other tissues (but the peak score for samples deriving from different developmental stages  
 569 of the same tissue type were not required to be <25). For instance, if the score for a retina age 8  
 570 weeks peak was >25 and the score for the corresponding retina age P7 peak was >25, but all  
 571 non-retinal peaks were <25, then that peak was called 'retina-specific'. After removing any  
 572 tissue-specific peaks that overlapped repetitive genomic sequences (~10% of peaks), we  
 573 selected the 1,000 peaks with the highest tissue-specific peak scores from each of adult brain,  
 574 retina, heart, and liver for inclusion as capture targets.

575

576 **Capture bait library design and synthesis.** For each of the 4,000 target regions, seven 80 bp  
 577 baits were designed to tile across the 300 bp region (sliding 37 bp at a time), for a total of 1.2  
 578 Mb and 28,000 baits. To check for potential off-target bait hybridization, bait candidates were  
 579 blasted against the mm9 genome, which was masked for the regions from which baits were  
 580 designed. By definition,  $T_m$  is the temperature at which 50% of the molecules are hybridized.  
 581 Bait candidates were accepted only if no BLAST hits (Altschul et al. 1990) with a predicted  $T_m >$   
 582 40.0 °C were found.

583

584 **GREAT analysis and Gene Ontology.** GREAT v2.0.2 analysis with mm9 as the reference  
 585 genome was implemented, using the 'single nearest gene' within 1000 kb as the algorithm for  
 586 associating genomic regions to genes, and using the whole genome as background and  
 587 excluding the 'include curated regulatory domains' option (McLean et al. 2010). The input to the  
 588 GREAT analysis was the list of 4,000 target DHS regions. Gene Ontology (GO) (Ashburner et al.  
 589 2000) enrichment analysis for 'biological process' in *Mus musculus* was implemented using  
 590 PANTHER (Mi et al. 2005) with AmiGO 2 v2.1.4 (Carbon et al. 2009). The input to the GO  
 591 analysis was the GREAT-generated list of genes associated with target DHSs ('region-to-gene'  
 592 associations).

593

594 **Restriction enzymes and PCR reagents.** Unless otherwise indicated, restriction enzymes  
595 were from New England Biolabs, and Phusion Hot Start Flex 2X Master Mix (New England  
596 Biolabs) was used for PCR. Primer sequences are listed in Supplemental Table S9.

597  
598 **Preparation of gDNA for capture.** Genomic DNA was purified from liver tissue of C57BL/6J  
599 mice and sonicated with Covaris E210 (duty 10%, intensity 4, cycles/burst 200, time 100 s). The  
600 freshly sonicated DNA was end repaired, 3' adenylated, ligated to commercial adapters, and  
601 enriched by PCR, using the TruSeq LT or TruSeq Nano Kit (Illumina) according to  
602 manufacturer's instructions (1 ug or 200 ng input gDNA, and 10 or 8 cycles of PCR,  
603 respectively). For final size selection and purification prior to capture, the samples were gel  
604 electrophoresed on 2% low melting point agarose and gel extracted with MinElute (Qiagen). To  
605 concentrate the samples in preparation for capture, the samples were speed vacuumed in  
606 LoBind tubes (Eppendorf).

607  
608 **Cis-regulome capture and preparation for cloning.** Capture was conducted in a similar  
609 manner as previously described (Gnirke et al. 2009). Two rounds of sequential capture were  
610 conducted to achieve high on-target rates (Lee et al. 2009). Briefly, for the first round of capture,  
611 a 9  $\mu$ L library mix was prepared, consisting of ~300 ng input (TruSeq LT or TruSeq Nano gDNA  
612 library), 2.5  $\mu$ g human Cot-1 DNA, 2.5  $\mu$ g salmon sperm DNA, and 0.6  $\mu$ L adapter blocking  
613 agent (MYcroarray). This solution was denatured at 95  $^{\circ}$ C for 5 min. Meanwhile, a 36.8  $\mu$ L  
614 hybridization mix was prepared, consisting of 5  $\mu$ L 20X SSPE (instead of the standard 20  $\mu$ L),  
615 0.8  $\mu$ L 0.5 M EDTA, 8  $\mu$ L 50X Denhardt's, 8  $\mu$ L 1% SDS, and 15  $\mu$ L RNase-free water. This  
616 solution was prewarmed at 65  $^{\circ}$ C for 3 min. A 6  $\mu$ L capture bait mix was prepared, consisting of  
617 50 ng (instead of the standard 500 ng) biotinylated baits and 1  $\mu$ L of SUPERase-In (Ambion).  
618 This solution was prewarmed at 65  $^{\circ}$ C for 2 min. Finally, 7  $\mu$ L of the library mix, 13  $\mu$ L of the  
619 hybridization mix, and all 6  $\mu$ L of the capture bait mix were incubated at 65  $^{\circ}$ C for ~24 hr. The  
620 reaction was then applied to Dynabeads MyOne Streptavidin C1 (Invitrogen) with washing and  
621 elution as described (Gnirke et al. 2009). Each capture reaction was purified with MinElute  
622 (Qiagen), with an elution volume of 30  $\mu$ L. Each eluate was speed vacuumed in a LoBind tube  
623 (Eppendorf) down to a volume of 3-4  $\mu$ L and used as the library 'input' for a single reaction in  
624 the second round of capture. The second round of capture was otherwise identical to the first.  
625 No PCR was conducted between the first and second rounds of capture. After the second round  
626 of capture, PCR was conducted using III\_NotI\_1XL and III\_NotI\_2XL primers (98  $^{\circ}$ C for 1 min,  
627 14-16 cycles: 98  $^{\circ}$ C for 10 sec, 58  $^{\circ}$ C for 30 sec, 72  $^{\circ}$ C for 1 min, followed by 72  $^{\circ}$ C for 5 min).  
628 The samples were PCR purified with MinElute (Qiagen), digested with NotI-high fidelity (HF) ,  
629 and gel extracted with MinElute (Qiagen). Two independent pools of capture products were  
630 generated, with each pool deriving from multiple capture reactions.

631  
632 **CRE-seq library construction.** To minimize the likelihood of cleaving captured fragments, the  
633 8-bp cutters NotI, FseI, and AseI were employed. To create the barcoded vector library for  
634 insertion of NotI-ended captured fragments, the *Rho* basal-DsRed construct (Hsiao et al. 2007)  
635 was modified with linkers on the 3' end of DsRed to replace a former NotI site with an EagI site  
636 and to add NsiI, FseI and AseI sites, and on the 5' end of the *Rho* basal promoter to add XbaI,  
637 NotI, and KpnI sites.

638 To add 15-mer barcodes, two pools of 30 nmol oligos were synthesized with random 15  
639 bp sequences (Integrated DNA Technologies) as BC\_F and BC\_R. The two pools were  
640 annealed and ligated into the AseI and NsiI sites of the vector. After transformation of 5-alpha  
641 chemically competent *E. coli* (New England Biolabs) and overnight growth in liquid culture, a  
642 total of  $\sim 9.5 \times 10^6$  colonies were harvested (as estimated from plating a small aliquot) and  
643 purified with the PureLink HiPure Plasmid Maxiprep Kit (Invitrogen). The barcoded vector library

644 was then digested with EagI-HF and dephosphorylated with alkaline phosphatase (Roche). The  
645 captured fragments were digested with NotI-HF and cloned into the EagI site of the vector  
646 library with 5-alpha chemically competent *E. coli* (New England Biolabs). A total of ~80,000  
647 colonies were scraped from LB/ampicillin agar plates, grown for ~2 hours in liquid LB/ampicillin  
648 culture, and purified with the PureLink HiPure Plasmid Maxiprep Kit (Invitrogen).

649 After paired-end sequencing to determine the CRE-barcode correspondence (described  
650 below), the minimal promoter-eGFP cassette was cloned into the FseI and AseI sites. The  
651 minimal promoter is the previously described 'Rho basal' minimal promoter, which contains a  
652 TATA box ('CATAA'), and which by itself does not have detectable activity in electroporated  
653 retina (Hsiau et al. 2007). The minimal promoter-eGFP cassette was created by replacing  
654 DsRed with eGFP (Zhang et al. 1996) in the Rho basal-DsRed construct (Hsiau et al. 2007).  
655 After transformation with 5-alpha chemically competent *E. coli* (New England Biolabs) and  
656 overnight growth in liquid culture, a total of  $\sim 2.7 \times 10^6$  colonies were harvested (as estimated by  
657 plating a small aliquot) and purified with the PureLink HiPure Plasmid Maxiprep Kit (Invitrogen).

658 The AAV-ITR vector was prepared by digesting the pAAV2.1-RHO-eGFP vector (Allocca  
659 et al. 2007) with NheI and XhoI, and replacing the RHO-eGFP cassette with a linker containing  
660 an EagI site. To transfer the library into the AAV-ITR vector, the entire CRE-minimal promoter-  
661 eGFP-polyA cassette was subjected to PCR using 5' Tak and NotI\_polyA\_R1 primers (98 °C for  
662 1 min, 10 cycles: 98 °C for 10 sec, 64 °C for 30 sec, 72 °C for 1 min 30 sec, followed by 72 °C  
663 for 5 min). The PCR product was digested with NotI-HF (New England Biolabs) and cloned into  
664 the EagI site of the AAV-ITR vector. After transformation of 5-alpha chemically competent *E. coli*  
665 (New England Biolabs) and overnight growth in liquid culture, a total of  $\sim 2.5 \times 10^6$  colonies (as  
666 estimated by plating a small aliquot) were harvested and purified with the PureLink HiPure  
667 Plasmid Maxiprep Kit (Invitrogen). ITR integrity was verified by restriction digest. Note that the  
668 final NotI digestion removes any captured fragments initially cloned in as NotI multimers, leaving  
669 only the 3'-most captured fragment.

670  
671 **Paired-end sequencing for CRE-barcode correspondence.** Prior to insertion of the promoter-  
672 reporter cassette, the library was prepared for paired-end sequencing as follows. PCR  
673 amplification was conducted using primers LibPCR\_F and LibPCR\_R (98 °C for 1 min, 8 cycles:  
674 98 °C for 10 sec, 64 °C for 30 sec, 72 °C for 1 min, followed by 72 °C for 5 min). The product  
675 was digested with NotI-HF and SacII, gel purified with MinElute (Qiagen), and ligated to P1\_NotI  
676 and PE2\_SacII adapters with T4 DNA ligase (New England Biolabs), using an equimolar mix of  
677 P1\_NotI indexed adapters to facilitate nucleotide balance. The ligation products were PCR  
678 amplified to enrich for molecules that had both P1 and PE2 adapters, using primers JKP4F and  
679 JKP4R (98 °C for 1 min, 14 cycles: 98 °C for 10 sec, 65 °C for 30 sec, 72 °C for 1 min, followed  
680 by 72 °C for 5 min). The final product was gel-extracted on 2% low melting point agarose and  
681 verified on an Agilent Bioanalyzer. Two lanes of MiSeq 2x250 bp sequencing were run at a  
682 loading concentration of 1.6-2 pM and 12-15% spiked-in Phi-X DNA (Illumina).

683  
684 **Analysis of paired-end sequencing for CRE-barcode correspondence.** Barcodes and  
685 captured fragment sequences were extracted based on flanking bases. Captured fragment  
686 sequences were aligned as paired reads to mm9 using Bowtie 2 v2.1.0 (Langmead and  
687 Salzberg 2012) with an allowed maximum insert size of 1000 bp ('-X 1000' setting). SAM files  
688 were converted to BAM files using SAMtools v0.1.19 (Li et al. 2009) and then to BED files using  
689 BEDTools v2.22.1 (Quinlan and Hall 2010). Only paired reads that mapped concordantly were  
690 used. Fragments were examined for overlap with the 4,000 target DHS regions (which were  
691 each 300 bp). If a fragment overlapped two adjacent target regions, it was assigned to the target  
692 region with the most bases of overlap. Barcodes were required to be 14-16 bp in length.  
693 Barcodes with multiple CRE fragment associations, and PCR-duplicate CRE fragments



694 associated with multiple barcodes (~1.6% of fragments), were filtered. A list of 'on-target' CRE  
695 correspondences for 45,670 barcoded constructs (minimum 10 reads) resulted. To determine  
696 the 'off-target' rate, the number of barcoded constructs that did not overlap a target DHS was  
697 found to be 712. Hence, ~98.5% of fragments were on-target.

698  
699 **Retinal explant electroporation and culture for CRE-seq.** Electroporation and explant culture  
700 of mouse retinas were performed as described previously (Montana et al. 2011b). In brief,  
701 retinas were dissected from newborn (P0) CD-1 mouse pups and coelectroporated with  
702 0.5 µg/µL AAV-ITR plasmid CRE-seq library and 0.5 µg/µL *Rho*-CBR3-DsRed, a rod-specific  
703 construct for visualizing electroporation efficiency (Corbo et al. 2010). Retinas were grown in  
704 explant culture and harvested 8 days later. Five retinas were pooled for each CRE-seq  
705 biological replicate.

706  
707 **Viral production.** Recombinant AAV9(2YF) was produced and purified as previously described  
708 (Grieger et al. 2006). To summarize, HEK293 cells at ~80% confluency were cotransfected with  
709 the AAV-ITR plasmid CRE-seq library, p-Helper plasmid, and AAV9(2YF) rep/cap plasmid  
710 (Dalkara et al. 2012). Cells were harvested 72 hours after transfection, and the virus was  
711 purified by iodixanol gradient ultracentrifugation, followed by buffer exchange. The viral titer, as  
712 determined by dot blot or quantitative PCR, ranged from  $5 \times 10^{12}$  to  $1 \times 10^{14}$  vg/mL (Zolotukhin  
713 et al. 2002; Aurnhammer et al. 2012).

714  
715 **Stereotactic cortical injection.** Stereotactic cortical injections were performed in a manner  
716 similar to that described (Cetin et al. 2006). Briefly, female CD-1 mice (age 4-6 weeks) were  
717 anesthetized with isoflurane. Each mouse received bilateral injections. For each injection, a  
718 small craniotomy was performed and 1 µL of AAV9(2YF) CRE-seq library was delivered into the  
719 primary motor cortex (stereotactic coordinates: dorsal/ventral axis 0.52 mm, anterior/posterior  
720 axis 1 mm, medial/lateral axis 1.5 mm). Animals were harvested 4-5 weeks after injection. The  
721 brain was sliced coronally and a fluorescent dissecting scope (Leica MZ16 F) was used to  
722 visualize GFP-positive regions, which were isolated by microdissection. Each CRE-seq  
723 biological replicate consisted of GFP-positive cortical tissue from a single animal.

724  
725 **Isolation of RNA and DNA and preparation for sequencing.** Tissues were rapidly harvested  
726 and rinsed in cold sterile HBSS with calcium and magnesium (Gibco) and stored at -80 °C in  
727 TRIzol (Invitrogen). Samples were homogenized in TRIzol, and RNA and DNA were isolated  
728 according to the manufacturer's instructions. RNA samples were treated with TURBO DNase  
729 (Ambion) to remove potential DNA contamination. RNA and DNA were prepared for sequencing  
730 essentially as previously described (Kwasnieski et al. 2012). RNA was reverse-transcribed with  
731 SuperScript III (Invitrogen) using oligo-dT primers. The resulting first-strand cDNA was treated  
732 with RNaseH. Both the cDNA and DNA samples were subjected to PCR to amplify the barcode  
733 sequence in the 3' UTR of GFP using the forward primer SSP1F and the reverse primer JKP3R  
734 (98 °C for 1 min, 22 cycles for DNA or 26 cycles for cDNA: 98 °C for 10 s, 60 °C for 30 s, 72 °C  
735 for 30 s, followed by 72 °C for 5 min). This resulted in PCR products flanked by *EagI* and *EcoRI*  
736 restriction enzyme sites. The products were purified with PureLink PCR Purification Kit  
737 (Invitrogen) and digested with *EagI*-HF and *EcoRI*. After digestion, the samples were gel  
738 purified with Qiagen Gel Extraction Kit and ligated to P1\_ *EagI* and PE2\_ *EcoRI* adapters using  
739 T4 DNA ligase (New England Biolabs). To enrich for molecules that had both P1 and PE2  
740 adapters, the ligation products were PCR amplified with primers JKP4F and JKP4R (98 °C for 1  
741 min, 20 cycles: 98 °C for 30 sec, 65 °C for 30 sec, 72 °C for 30 sec, followed by 72 °C for 5 min).  
742 The final product was gel purified from 2% low melting point agarose and verified on an Agilent  
743 Bioanalyzer.

744  
745 **Illumina sequencing for CRE-seq barcode abundance.** For each tissue, the three cDNA  
746 samples and three corresponding DNA samples were multiplexed and run on a single lane of  
747 Illumina HiSeq 2000 (1x50 bp) at a loading concentration of 8 pM with 10% spiked-in Phi-X DNA.  
748

749 **CRE-seq data analysis.** Samples were demultiplexed and the barcode was extracted based on  
750 flanking sequences. Reads were tabulated to obtain the raw RNA and DNA counts for each  
751 barcode. Only barcodes with at least 10 raw DNA reads in all 3 biological replicates of a tissue  
752 were included (36,005 barcodes for retina and 38,826 barcodes for cerebral cortex). For each  
753 barcode, the RNA count was normalized to the total RNA counts in the sample, and the DNA  
754 count was normalized to the total DNA counts in the sample. The normalized expression was  
755 the ratio of the normalized RNA count to the normalized DNA count. A pseudocount of 0.001  
756 was added to the normalized expression, and the  $\log_2$  was taken. The average of the  $\log_2$   
757 values across biological replicates was the 'mean expression ( $\log_2$  units)'.  
758

759 **Histology.** Retinal explants were rinsed twice with PBS and fixed in 4% paraformaldehyde/PBS  
760 for 30-60 min at room temperature, equilibrated in 30% sucrose/PBS, and embedded in Tissue-  
761 Tek O.C.T. (Sakura). Retinal cryosections (12-14  $\mu\text{m}$ ) were prepared and stored at -20 °C until  
762 imaging. For stereotactically injected brains, animals were deeply anesthetized with  
763 ketamine/xylazine and then transcardially perfused with heparin/PBS followed by 4%  
764 paraformaldehyde/PBS. Animals were decapitated and the brains were dissected in PBS and  
765 post-fixed in 4% paraformaldehyde/PBS at 4 °C for at least a day. Vibratome sections (200  $\mu\text{m}$ )  
766 were prepared from agarose-embedded brain slices and then optically cleared with  
767 glycerol/PBS (Selever et al. 2011). Brain slices were treated with sodium borohydride to  
768 minimize autofluorescence (Clancy and Cauller 1998). For anti-RBFOX3 (also known as anti-  
769 NeuN) staining of free-floating vibratome sections, the sections were blocked with 4% normal  
770 donkey serum (NDS)/0.25% Triton X-100/PBS for at least 1 hr at room temperature with gentle  
771 agitation, incubated with rabbit anti-RBFOX3 antibody (ABN78; EMD Millipore) (1:50, diluted in  
772 4% NDS/0.1% Triton X-100/PBS) overnight at 4 °C with gentle agitation, washed with 0.1%  
773 Triton X-100/PBS, incubated with Alexa Fluor 555 donkey anti-rabbit (A-31572; Molecular  
774 Probes) (1:800, diluted in 4% NDS/0.1% Triton X-100/PBS) for 1 hr at room temperature with  
775 gentle agitation, and washed with 0.1% Triton X-100/PBS. All brain slices were stored in PBS at  
776 4 °C until imaging. For imaging, tissue was mounted with Vectashield (Vectorlabs) and  
777 coverslipped. Confocal imaging was conducted with a laser confocal microscope (Zeiss LSM  
778 700) and ZEN 2009 software (Zeiss). Flat-mount imaging of an untreated brain slice (Fig. 3D)  
779 was conducted with an inverted fluorescent microscope (Nikon Eclipse TE300) and MetaMorph  
780 software (Molecular Devices). Images were processed with Adobe Photoshop.  
781

782 **Cluster analysis of biological replicates.** Hierarchical clustering and principal component  
783 analysis (PCA) were used to assess the underlying structure of CRE expression across retina  
784 and brain replicates. For hierarchical clustering, the sample distance was defined as one minus  
785 the Pearson correlation coefficient (calculated across the normalized expression of the ~35,000  
786 barcodes with at least 10 DNA reads in all six samples), and clustering was implemented using  
787 average linkage. PCA was performed via singular value decomposition on scaled,  
788 centered expression data (i.e., zero-centered values with unit variance).  
789

790 **Analysis of TF motif enrichment in low vs. high-expressing DHSs.** To compare the motif  
791 content of low- and high-expressing constructs (Fig. 5E), a list of brain and retina TF motifs  
792 were obtained as follows. DNase-seq reads for adult brain (GSM1014151, replicate 1) and adult  
793 retina (GSM1014175) were downloaded and aligned to mm9 with Bowtie 2 v2.2.3 (Langmead

794 and Salzberg 2012). DNase-seq peaks were then called using MACS2 v2.1.0 (Zhang et al.  
795 2008). For *de novo* motif discovery, peaks were first partitioned by HOMER v4.7 annotations  
796 ('promoter,' 'intronic,' and 'intergenic') (Heinz et al. 2010). *De novo* motif discovery was then  
797 performed independently for each of these classes of peaks from brain and retina, with the final  
798 motif list consisting of all motifs identified at a threshold of  $p < 1 \times 10^{-50}$ . To compare similar  
799 numbers of DHSs in the 'high' and 'low' categories, individual barcoded constructs were ranked  
800 by average expression in each tissue. The highest-expressing constructs that constituted 100  
801 distinct DHS target regions (regardless of DHS tissue origin) were classified as 'high' in that  
802 tissue, and the lowest-expressing constructs that constituted 100 distinct DHS target regions  
803 (regardless of DHS tissue origin) were classified as 'low' in that tissue (DNA read count was  
804 used to break ties). Finally, overlapping intervals were merged, and the resulting regions were  
805 scored for motif enrichment (binomial test, via HOMER) relative to a background of ~50,000  
806 random mm9 sequences matched for size and dinucleotide content.

807  
808 **Receiver operating characteristic (ROC) curves.** To quantify the extent to which sequence  
809 features and epigenomic data could predict expression (Fig. 5F), we implemented multiple  
810 logistic regression as a means of classifying whether or not individual constructs were among  
811 those with the highest expression (similar to the approach described by (Kwasniewski et al.  
812 2014)). Briefly, all assayed constructs (~36,000 constructs for retina and ~39,000 constructs for  
813 cerebral cortex) were partitioned by expression into 'high' and 'not high' expression groups.  
814 'High' was defined here as mean expression across replicates ( $\log_2$  units) of  $>-2$  for constructs  
815 assayed in retina (~95<sup>th</sup> percentile), and  $>2$  for constructs assayed in the cerebral cortex (~99<sup>th</sup>  
816 percentile) (see Figure 4B). Our model included terms for GC content (averaged across the  
817 CRE fragment), phylogenetic conservation (30-way vertebrate PhastCons, averaged across the  
818 CRE fragment) (Siepel et al. 2005), brain or retina DNase-seq data ( $\log_2((\text{read depth}+1)/\text{CRE}$   
819  $\text{size}))$ , retina CRX ChIP-seq data ( $\log_2((1/2)*(\text{read depth of two WT CRX ChIP-seq replicates} +$   
820  $1)/\text{CRE size}))$  (Corbo et al. 2010), and individual TF motifs (the number of each motif in each  
821 CRE fragment, as identified by HOMER). CRX ChIP-seq data were only included in the retina  
822 model, and distinct TFs were considered for retina and cerebral cortex models. TF motifs for  
823 each tissue were identified as described above (17 motifs for retina, and 13 motifs for cerebral  
824 cortex; see Supplemental Table S5). Two retinal motifs (YY1 and ZBTB33) were omitted from  
825 the model, as they were observed fewer than 100 times across the ~36,000 constructs, and  
826 hence 15 motifs were in the retina TF motif model. The performance (AUC) of models was  
827 quantified using the ROCR package in R (Sing et al. 2005). Five-fold cross-validation was used  
828 to control for over-fitting.

829  
830 **Expression scores for browser screenshots.** For Figure 6A, the scales for the heat maps are  
831 indicated. Elsewhere, heat maps were generated according to the default grayscale on the  
832 UCSC Genome Browser (Karolchik et al. 2014), using custom bed tracks that were generated  
833 as follows. For each biological replicate, a bed track was created using the useScore=1 attribute  
834 for intensity shading of individual barcoded constructs using a 'bed score'. The 'bed score' was  
835 obtained by adding 10 to the  $\log_2$  expression and multiplying by 75. For each tissue, an  
836 'average signal' bedGraph track was created by segmenting the tiled regions and averaging the  
837 bed scores across replicates and barcodes. A segment was required to be encompassed by at  
838 least 2 barcoded constructs to be included in the 'average signal' track. The windowing function  
839 was set to 'mean'. A smoothing window function (10 pixels) was applied to the average signal  
840 tracks, which were displayed on the following scales: 0 to 1400 for retina, and 300 to 1200 for  
841 cortex.

842  
843 **Synthesis of individual constructs for validation.** The R28 constructs were cloned as  
844 EcoRV/KpnI fragments. To create the long and short R28 constructs, the R28\_L/R28\_R and

845 R28\_S/R28\_R primer pairs were used, respectively. To create the mutant R28 construct,  
846 R28\_MT was ordered as a double-stranded gene block (Integrated DNA Technologies). The  
847 R62 constructs were cloned as EcoRI/XbaI fragments. To create the long and short R62  
848 constructs, the R62\_L/R62\_R and R62\_S/R62\_R primer pairs were used, respectively. To  
849 create the mutant R62 construct, R62\_MT was ordered as a double-stranded gene block  
850 (Integrated DNA Technologies). The R64 constructs were cloned as EcoRV/KpnI fragments. To  
851 create the long, short, and mutant R64 constructs, the R64\_L/R64\_R, R64\_S/R64\_R, and  
852 R64\_MT/R64\_R primer pairs were used, respectively. For the PCR reactions, C57BL/6J gDNA  
853 was the template. The CREs were digested and cloned upstream of the minimal promoter-  
854 eGFP cassette in the *Rho* basal-eGFP vector, which was created from *Rho* basal-DsRed (Hsiau  
855 et al. 2007) by replacing DsRed with eGFP at XmaI and NotI sites. Test constructs were  
856 confirmed with Sanger sequencing that encompassed the entire CRE.

857  
858 **Validation of individual constructs by fluorescent reporter assays.** Electroporation, explant  
859 culture, and quantification of fluorescence were performed essentially as previously described  
860 (Montana et al. 2011b). In brief, as for CRE-seq, retinas were dissected from newborn (P0) CD-  
861 1 mouse pups. Here, they were coelectroporated with 0.5 µg/uL of the test construct and 0.5  
862 µg/uL *Rho*-CBR3-DsRed (Corbo et al. 2010). Retinas were cultured for 8 days, fixed, and then  
863 whole mounted for quantitative imaging of fluorescent intensity (GFP intensity normalized to  
864 DsRed intensity), using a monochromatic camera (Hamamatsu ORCA-AG) and MetaMorph  
865 software (Molecular Devices). For each retina, five regions were quantified in ImageJ and  
866 averaged. SEM was calculated based on normalized fluorescence measurements across  
867 retinas (n = 10-12 retinas per test construct). Representative whole mount images using a color  
868 camera (Olympus DP70) were also taken.

869  
870 **Comparison with CapSTARR-seq.** The raw sequence data for the CapSTARR-seq (Vanhille  
871 et al. 2015) input library (GEO accession number GSM1463994) were downloaded and mapped  
872 to mm9 with Bowtie 2 v2.1.0 (Langmead and Salzberg 2012).

## 873 874 **DATA ACCESS**

875  
876 The sequence data from this study have been submitted to the NCBI Gene Expression  
877 Omnibus (GEO; <http://www.ncbi.nlm.nih.gov/geo/>) under accession number GSE68247. Custom  
878 tracks for the UCSC Genome Browser (Karolchik et al. 2014) are provided in Supplemental  
879 Table S10.

## 880 881 **ACKNOWLEDGMENTS**

882  
883 The authors would like to thank Karen Lawrence and Jennifer Enright for contributing to the  
884 design of the barcoded vector library and sequencing adapters, Jean-Marie Rouillard of  
885 MYcroarray for capture advice, Ronald Perez of the Animal Surgery Core at the Hope Center for  
886 Neurological Disorders for stereotactic cortical injections, Mingjie Li of the Viral Vectors Core at  
887 the Hope Center for Neurological Disorders for assistance with viral production, and the  
888 Genome Technology Access Center in the Department of Genetics at Washington University  
889 School of Medicine for sequencing services. We would also like to thank Michael A. White for  
890 helpful discussion and Shuyi Ma for critical reading of the manuscript. This work was supported  
891 by the Foundation Fighting Blindness (J.G.F.), Simons Foundation Autism Research Initiative  
892 (grant number 275579 to J.C.C.) and the National Institutes of Health (HG006790 and  
893 EY018826 to J.C.C., EY022975 to J.G.F., EY024958 to J.C.C. and J.G.F., and 5T32EY013360  
894 to S.Q.S.).

895

896 **DISCLOSURE DECLARATION**

897

898 The authors have no disclosures to declare.

899 **FIGURE LEGENDS**

900

901 **Figure 1. ‘Capture-and-clone’ allows synthesis of CRE-seq libraries with long CREs.** (A)  
 902 Schematic of the capture-and-clone approach. Size-selected, adapter-ligated genomic DNA was  
 903 hybridized to biotinylated RNA baits that tiled across candidate CRE regions of interest.  
 904 Captured fragments were cloned into a barcoded vector library with unique 15-mer barcodes.  
 905 Paired-end sequencing revealed the CRE-barcode correspondence. A minimal promoter-GFP  
 906 reporter cassette was subsequently cloned into the library. (B) Histogram showing the  
 907 distribution of the lengths of captured fragments that were cloned into the barcoded vector  
 908 library, based on paired-end sequencing. The median length was 464 bp. (C) Histogram  
 909 showing the distribution of target coverage, i.e., the number of captured fragments that  
 910 overlapped a 300 bp target region. Of the 4,000 targeted regions, 3,483 regions were  
 911 represented by at least one construct. The median coverage among represented regions was 8.  
 912 Not shown in graph: 517 non-represented regions and 114 target regions with a coverage  
 913 of >50.

914

915 **Figure 2. Tiling of captured fragments across target regions.** Capture baits were designed  
 916 based on adult (8 week old C57BL/6J) DNase-seq data from Mouse ENCODE (Yue et al. 2014).  
 917 Paired-end sequencing revealed the locations of individual barcoded, captured-and-cloned  
 918 fragments. The UCSC Genome Browser (mm9) (Karolchik et al. 2014) screenshots depict: (A)  
 919 Captured fragments for an entire representative chromosome (chr7). ‘Off-target’ fragments, i.e.,  
 920 those that did not overlap a 300 bp target bait region, are also shown. Examples of captured  
 921 fragments: (B) around a retina-specific locus, *Rho* (rhodopsin), (C) in an intron of a brain-  
 922 specific locus, *Grin2a* (glutamate receptor, ionotropic, NMDA2a [epsilon 1]), (D) in the 5’  
 923 UTR/promoter region of a heart-specific locus, *Tnni3* (troponin I, cardiac 3), and (E) downstream  
 924 of a liver-specific locus, *Alb* (albumin). Note that some DNase-seq peaks visible in the  
 925 screenshots were not included as targets for capture. PhastCons depict 30-way vertebrate  
 926 phylogenetic conservation (Siepel et al. 2005).

927

928 **Figure 3. Delivery of capture CRE-seq library into mouse retina *ex vivo* and cerebral**  
 929 **cortex *in vivo*.** (A) Schematic of the CRE-seq library delivery approach. The plasmid library can  
 930 be directly electroporated into the retina *ex vivo*. Alternatively, the library can be packaged into  
 931 AAV and delivered via stereotactic injection into the cerebral cortex *in vivo*. (B) Scatterplot  
 932 comparing the relative abundance of ~45,000 individual barcoded constructs in the plasmid  
 933 library delivered into the retina, and in the AAV-packaged library delivered into cortex, as  
 934 measured by barcode DNA reads summed across the three biological replicates for each tissue  
 935 and then normalized to the total number of barcode DNA reads. Each data point represents a  
 936 unique barcoded construct. DNA reads were well-correlated (Pearson  $r = 0.95$ ), indicating  
 937 fidelity of barcode representation after AAV packaging and delivery. Off-target constructs and  
 938 constructs with 0 reads in all samples were excluded. Not shown: 4 points falling outside the  
 939 depicted plot range (included in the calculation of Pearson  $r$ ). Red line, linear regression. (C)  
 940 Confocal image of a retina that was electroporated with the plasmid library and cryosectioned  
 941 after 8 days in culture. ONL, outer nuclear layer. INL, inner nuclear layer. (D) Flat-mount image  
 942 of a coronal slice from a brain injected with the AAV-packaged library bilaterally into the primary  
 943 motor cortex and harvested ~4 weeks later. (D’) Schematic corresponding to the flat-mount  
 944 image. Note the bilateral GFP-positive regions in the cortex, as well as bundles of GFP-positive  
 945 axons in the corpus callosum (red arrow). (E) Confocal image of a cortical region infected with  
 946 the AAV-packaged library.

947

948 **Figure 4. Tissue-specific *cis*-regulatory activity of DHSs.** (A) Frequency distribution of DHSs  
 949 ranked by *cis*-regulatory activity (bin size: 5 percentile) as measured in the retina (top) or

950 cerebral cortex (bottom). In the retina, ~15% DHSs had undetectable activity and hence were  
 951 binned together. Averages were taken across biological replicates and barcodes for a given  
 952 target DHS. Only DHSs with at least 2 barcoded constructs were included in this analysis  
 953 (~3,000 DHSs). Frequencies were normalized to the total number of DHSs in each category. To  
 954 test for enrichment, chi-squared test was performed (one-tailed): \*\*\* $p < 10^{-4}$ , \*\* $p < 0.01$ , \* $p < 0.05$ .  
 955 (B) Scatterplot showing the expression of individual barcoded constructs as assayed in the  
 956 cerebral cortex (x-axis) vs. retina (y-axis). Each dot represents an individual construct. For each  
 957 construct, the average measurement across the three biological replicates for each tissue was  
 958 taken. The ~35,000 barcodes that were well-represented (at least 10 DNA reads) in all six  
 959 samples were included in the analysis. Gray, blue, red, and orange dots denote constructs with  
 960 CRE fragments that overlap retina, brain, heart, and liver DHSs, respectively. The dotted gray  
 961 box encompasses constructs that are strongly active in the retina, and the dotted blue box  
 962 encompasses constructs that are strongly active in the cortex.  
 963

964 **Figure 5. Parameters that predict CRE activity.** (A) to (D) Retinal DHSs as assayed in the  
 965 retina (left) and brain DHSs as assayed in the cerebral cortex (right). Each panel shows a 1 kb  
 966 centered window. Only DHSs with at least 2 barcodes were included in this analysis, i.e., 710  
 967 retinal DHSs in retina (black lines, left) and 696 brain DHSs in cortex (black lines, right). The top  
 968 100 (red lines, left) and top 200 (orange lines, left) retinal DHSs expressed in the retina and the  
 969 top 100 (red lines, right) and top 200 (orange lines, right) brain DHSs expressed in the cortex  
 970 are shown. To compare the top 100 DHSs vs. the rest of the DHSs in each group, two-tailed  
 971 student's t-test was calculated for the means within the 1 kb window, except for PhastCons  
 972 scores, which was calculated within the central 100 bp. \*\*\* $p < 0.001$ , \*\* $p < 0.01$ , \* $p < 0.05$ , N.S., not  
 973 significant. (A) *Cis*-regulatory activity, as measured by mean expression in  $\log_2$  units. For each  
 974 assayed DHS, at each base position across the 1 kb window, the expression values of the  
 975 individual barcoded constructs whose CREs overlapped the position were averaged across  
 976 biological replicates. (B) DNase-seq score (Yue et al. 2014). (C) GC content, calculated in 50 bp  
 977 windows, sliding 25 bp at a time. The fractions denote the proportion of DHSs that were  
 978 promoter-proximal (i.e., located within -1 kb to +100 bp relative to the nearest TSS) based on  
 979 GREAT annotations (McLean et al. 2010). (D) Phylogenetic conservation as measured by 30-  
 980 way vertebrate PhastCons (Siepel et al. 2005). (E) Enrichment for TF motifs among low vs.  
 981 high-expressing DHSs in each tissue, without restriction on the type of DHS (see Methods).  
 982 Only significant motifs are shown ( $p < 0.05$  in at least one category). For motifs enriched in both  
 983 tissues, the logo from the tissue with the more significant enrichment is shown. Abbreviations:  
 984 HD, homeodomain; NR, nuclear receptor; ZF, zinc finger; HTH, helix-turn-helix. (F) Receiver  
 985 operator characteristic (ROC) curves show the performance of logistic regression models for GC  
 986 content, PhastCons, TF motifs, retina or brain DNase-seq, or a combined model. A model  
 987 based on CRX ChIP-seq (Corbo et al. 2010) was included for the retina only. The area under  
 988 the curve (AUC) for each model is indicated. For cross-validation results, see Supplemental  
 989 Table S6.  
 990

991 **Figure 6. Truncation mutation analysis by CRE-seq.** (A) Example of a truncation mutation  
 992 analysis at the *Nrl* promoter via a traditional one-at-a-time reporter assay (Montana et al. 2011b)  
 993 vs. capture-and-clone CRE-seq. For the traditional reporter constructs, the 3' end extends  
 994 beyond the window depicted in the figure. For the CRE-seq data, only barcoded constructs in  
 995 the same orientation as the *Nrl* promoter are shown. The yellow highlighted region corresponds  
 996 to a known critical region with CRX and RORB motifs (Andre et al. 1998; Montana et al. 2011b).  
 997 The minus strand of DNA is displayed. In (A) and (B), the CRX motif (from HOMER (Heinz et al.  
 998 2010)) is based on CRX ChIP-seq data (Corbo et al. 2010). The reverse orientation of the CRX  
 999 motif is displayed. Additional examples of CRE-seq truncation mutation analysis: (B) Retinal  
 1000 DHS with retina-specific expression. The critical region identified by CRE-seq (pink) contains a

1001 putative CRX motif. (C) Two adjacent brain DHSs in the same intron of *Bsn* exhibit low (DHS #1,  
1002 green) vs. high (DHS #2, pink) activity in the cortex. (D) Truncation mutation analysis of a brain  
1003 DHS. A gradual decrease in activity was observed within the ~150 bp critical region (pink),  
1004 corresponding to a phylogenetically conserved peak. Within this critical region, a smaller region  
1005 (vertical blue stripe) was identified that contained an E-box consensus motif ('CANNTG') and a  
1006 motif for a bZIP protein, based on AP-1 ChIP-seq data (Heinz et al. 2010). All browser images  
1007 are from the UCSC Genome Browser (mm9) (Karolchik et al. 2014). DNase-seq data are from  
1008 Mouse ENCODE (Yue et al. 2014). PhastCons depict 30-way vertebrate phylogenetic  
1009 conservation (Siepel et al. 2005). The heat map scale shown in (B) is the same as that used in  
1010 (C) and (D).  
1011

1012 **Figure 7. Validation of individual loci by fluorescence reporter assays.** (A) Critical regions  
1013 (pink areas) identified by CRE-seq truncation mutation analysis at three retinal DHSs (R64, R28,  
1014 and R62) were validated by testing of individual constructs with fluorescence reporter assays.  
1015 Depicted CRE-seq data are based on expression scores averaged across retinal replicates.  
1016 Note that R64 is the same locus as in Figure 6B. For each locus, a 'long' construct containing  
1017 the critical region (CR), a 'short' construct without the critical region, and a 'mutant' construct  
1018 with point mutations (red font) in predicted CRX sites (blue font) were synthesized. Sequences  
1019 are shown for the plus strand of DNA in all cases. For R62, one CRX site fell within the critical  
1020 region, and a second CRX site was immediately adjacent (yellow area). Individual test  
1021 constructs were directionally cloned upstream of the minimal promoter-GFP cassette in a non-  
1022 AAV vector. The test constructs were coelectroporated into explant retinas with *Rho*-CBR3-  
1023 DsRed (Corbo et al. 2010) as a loading control. (B) Representative whole mount images of  
1024 electroporated retinas are shown (exposure times are the same for all images). (C)  
1025 Quantification of the GFP levels normalized to DsRed levels. Error bar represents SEM (n = 10-  
1026 12 retinas per test construct). \*\*\*P-value < 10<sup>-6</sup> (two-tailed student's t test).



1027 **REFERENCES**

- 1028
- 1029 RetNet, <http://www.sph.uth.tmc.edu/RetNet/>.
- 1030 Albert FW, Krugiyak L. 2015. The role of regulatory variation in complex traits and disease. *Nature*
- 1031 *reviews Genetics* **16**(4): 197-212.
- 1032 Allocca M, Mussolino C, Garcia-Hoyos M, Sanges D, Iodice C, Petrillo M, Vandenberghe LH, Wilson JM,
- 1033 Marigo V, Surace EM et al. 2007. Novel adeno-associated virus serotypes efficiently transduce
- 1034 murine photoreceptors. *Journal of virology* **81**(20): 11372-11380.
- 1035 Altrock WD, tom Dieck S, Sokolov M, Meyer AC, Sigler A, Brakebusch C, Fassler R, Richter K, Boeckers TM,
- 1036 Potschka H et al. 2003. Functional inactivation of a fraction of excitatory synapses in mice
- 1037 deficient for the active zone protein bassoon. *Neuron* **37**(5): 787-800.
- 1038 Altschul SF, Gish W, Miller W, Myers EW, Lipman DJ. 1990. Basic local alignment search tool. *Journal of*
- 1039 *molecular biology* **215**(3): 403-410.
- 1040 Andersson R, Gebhard C, Miguel-Escalada I, Hoof I, Bornholdt J, Boyd M, Chen Y, Zhao X, Schmidl C,
- 1041 Suzuki T et al. 2014. An atlas of active enhancers across human cell types and tissues. *Nature*
- 1042 **507**(7493): 455-461.
- 1043 Andre E, Gawlas K, Becker-Andre M. 1998. A novel isoform of the orphan nuclear receptor RORbeta is
- 1044 specifically expressed in pineal gland and retina. *Gene* **216**(2): 277-283.
- 1045 Arnold CD, Gerlach D, Stelzer C, Boryn LM, Rath M, Stark A. 2013. Genome-wide quantitative enhancer
- 1046 activity maps identified by STARR-seq. *Science* **339**(6123): 1074-1077.
- 1047 Arvey A, Agius P, Noble WS, Leslie C. 2012. Sequence and chromatin determinants of cell-type-specific
- 1048 transcription factor binding. *Genome research* **22**(9): 1723-1734.
- 1049 Aschauer DF, Kreuz S, Rumpel S. 2013. Analysis of transduction efficiency, tropism and axonal transport
- 1050 of AAV serotypes 1, 2, 5, 6, 8 and 9 in the mouse brain. *PLoS one* **8**(9): e76310.
- 1051 Ashburner M, Ball CA, Blake JA, Botstein D, Butler H, Cherry JM, Davis AP, Dolinski K, Dwight SS, Eppig JT
- 1052 et al. 2000. Gene ontology: tool for the unification of biology. The Gene Ontology Consortium.
- 1053 *Nature genetics* **25**(1): 25-29.
- 1054 Aurnhammer C, Haase M, Muether N, Hausl M, Rauschhuber C, Huber I, Nitschko H, Busch U, Sing A,
- 1055 Ehrhardt A et al. 2012. Universal real-time PCR for the detection and quantification of adeno-
- 1056 associated virus serotype 2-derived inverted terminal repeat sequences. *Human gene therapy*
- 1057 *methods* **23**(1): 18-28.
- 1058 Bae BI, Jayaraman D, Walsh CA. 2015. Genetic Changes Shaping the Human Brain. *Developmental cell*
- 1059 **32**(4): 423-434.
- 1060 Baker M. 2011. Microarrays, megasynthesis. *Nat Meth* **8**(6): 457-460.
- 1061 Blatti C, Kazemian M, Wolfe S, Brodsky M, Sinha S. 2015. Integrating motif, DNA accessibility and gene
- 1062 expression data to build regulatory maps in an organism. *Nucleic acids research*.
- 1063 Buenrostro JD, Giresi PG, Zaba LC, Chang HY, Greenleaf WJ. 2013. Transposition of native chromatin for
- 1064 fast and sensitive epigenomic profiling of open chromatin, DNA-binding proteins and
- 1065 nucleosome position. *Nature methods* **10**(12): 1213-1218.
- 1066 Butler JE, Kadonaga JT. 2002. The RNA polymerase II core promoter: a key component in the regulation
- 1067 of gene expression. *Genes & development* **16**(20): 2583-2592.
- 1068 Carbon S, Ireland A, Mungall CJ, Shu S, Marshall B, Lewis S, Ami GOH, Web Presence Working G. 2009.
- 1069 AmiGO: online access to ontology and annotation data. *Bioinformatics* **25**(2): 288-289.
- 1070 Cetin A, Komai S, Eliava M, Seeburg PH, Osten P. 2006. Stereotaxic gene delivery in the rodent brain.
- 1071 *Nature protocols* **1**(6): 3166-3173.
- 1072 Chen S, Wang QL, Nie Z, Sun H, Lennon G, Copeland NG, Gilbert DJ, Jenkins NA, Zack DJ. 1997. Crx, a
- 1073 novel Otx-like paired-homeodomain protein, binds to and transactivates photoreceptor cell-
- 1074 specific genes. *Neuron* **19**(5): 1017-1030.

- 1075 Chen S, Zack DJ. 1996. Ret 4, a positive acting rhodopsin regulatory element identified using a bovine  
1076 retina in vitro transcription system. *The Journal of biological chemistry* **271**(45): 28549-28557.
- 1077 Clancy B, Cauller LJ. 1998. Reduction of background autofluorescence in brain sections following  
1078 immersion in sodium borohydride. *Journal of neuroscience methods* **83**(2): 97-102.
- 1079 Clark AJ, Bissinger P, Bullock DW, Damak S, Wallace R, Whitelaw CB, Yull F. 1994. Chromosomal position  
1080 effects and the modulation of transgene expression. *Reproduction, fertility, and development*  
1081 **6**(5): 589-598.
- 1082 Cleary MA, Killian K, Wang Y, Bradshaw J, Cavet G, Ge W, Kulkarni A, Paddison PJ, Chang K, Sheth N et al.  
1083 2004. Production of complex nucleic acid libraries using highly parallel in situ oligonucleotide  
1084 synthesis. *Nature methods* **1**(3): 241-248.
- 1085 Clotman F, Jacquemin P, Plumb-Rudewiez N, Pierreux CE, Van der Smissen P, Dietz HC, Courtoy PJ,  
1086 Rousseau GG, Lemaigre FP. 2005. Control of liver cell fate decision by a gradient of TGF beta  
1087 signaling modulated by Onecut transcription factors. *Genes & development* **19**(16): 1849-1854.
- 1088 Corbo JC, Lawrence KA, Karlstetter M, Myers CA, Abdelaziz M, Dirkes W, Weigelt K, Seifert M, Benes V,  
1089 Fritsche LG et al. 2010. CRX ChIP-seq reveals the cis-regulatory architecture of mouse  
1090 photoreceptors. *Genome research* **20**(11): 1512-1525.
- 1091 Dalkara D, Byrne LC, Lee T, Hoffmann NV, Schaffer DV, Flannery JG. 2012. Enhanced gene delivery to the  
1092 neonatal retina through systemic administration of tyrosine-mutated AAV9. *Gene therapy* **19**(2):  
1093 176-181.
- 1094 Davidson EH. 2001. *Genomic regulatory systems : development and evolution*. Academic Press, San  
1095 Diego.
- 1096 Day TP, Byrne LC, Schaffer DV, Flannery JG. 2014. Advances in AAV vector development for gene therapy  
1097 in the retina. *Advances in experimental medicine and biology* **801**: 687-693.
- 1098 Daya S, Berns KI. 2008. Gene therapy using adeno-associated virus vectors. *Clinical microbiology reviews*  
1099 **21**(4): 583-593.
- 1100 Dermitzakis ET, Clark AG. 2002. Evolution of transcription factor binding sites in Mammalian gene  
1101 regulatory regions: conservation and turnover. *Molecular biology and evolution* **19**(7): 1114-  
1102 1121.
- 1103 Dickel DE, Zhu Y, Nord AS, Wylie JN, Akiyama JA, Afzal V, Plajzer-Frick I, Kirkpatrick A, Gottgens B,  
1104 Bruneau BG et al. 2014. Function-based identification of mammalian enhancers using site-  
1105 specific integration. *Nature methods* **11**(5): 566-571.
- 1106 Dror I, Golan T, Levy C, Rohs R, Mandel-Gutfreund Y. 2015. A widespread role of the motif environment  
1107 in transcription factor binding across diverse protein families. *Genome research* **25**(9): 1268-  
1108 1280.
- 1109 Duttke SH, Lacadie SA, Ibrahim MM, Glass CK, Corcoran DL, Benner C, Heinz S, Kadonaga JT, Ohler U.  
1110 2015. Human promoters are intrinsically directional. *Molecular cell* **57**(4): 674-684.
- 1111 Edmondson DG, Lyons GE, Martin JF, Olson EN. 1994. Mef2 gene expression marks the cardiac and  
1112 skeletal muscle lineages during mouse embryogenesis. *Development* **120**(5): 1251-1263.
- 1113 Fenouil R, Cauchy P, Koch F, Descostes N, Cabeza JZ, Innocenti C, Ferrier P, Spicuglia S, Gut M, Gut I et al.  
1114 2012. CpG islands and GC content dictate nucleosome depletion in a transcription-independent  
1115 manner at mammalian promoters. *Genome research* **22**(12): 2399-2408.
- 1116 Freund CL, Gregory-Evans CY, Furukawa T, Papaioannou M, Looser J, Ploder L, Bellingham J, Ng D,  
1117 Herbrick JA, Duncan A et al. 1997. Cone-rod dystrophy due to mutations in a novel  
1118 photoreceptor-specific homeobox gene (CRX) essential for maintenance of the photoreceptor.  
1119 *Cell* **91**(4): 543-553.
- 1120 Gisselbrecht SS, Barrera LA, Porsch M, Aboukhalil A, Estep PW, 3rd, Vedenko A, Palagi A, Kim Y, Zhu X,  
1121 Busser BW et al. 2013. Highly parallel assays of tissue-specific enhancers in whole Drosophila  
1122 embryos. *Nature methods* **10**(8): 774-780.

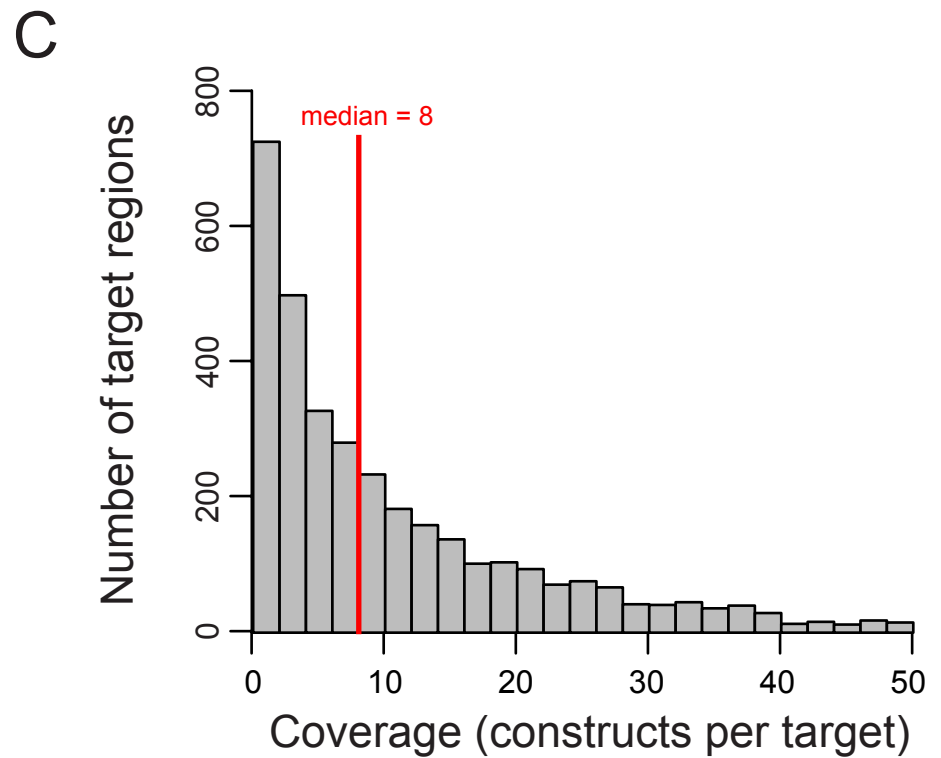
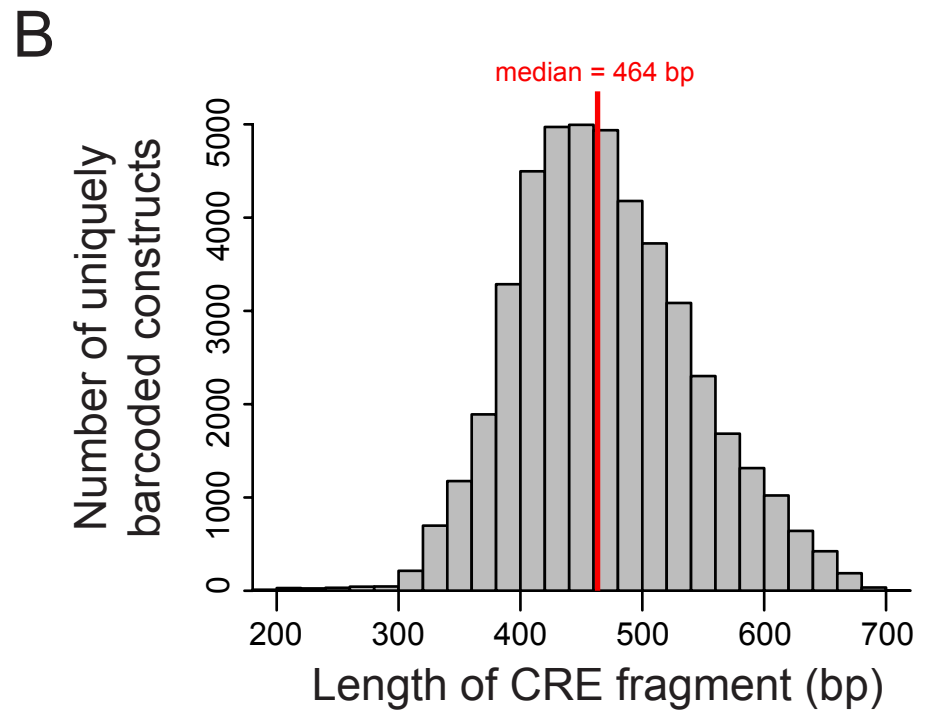
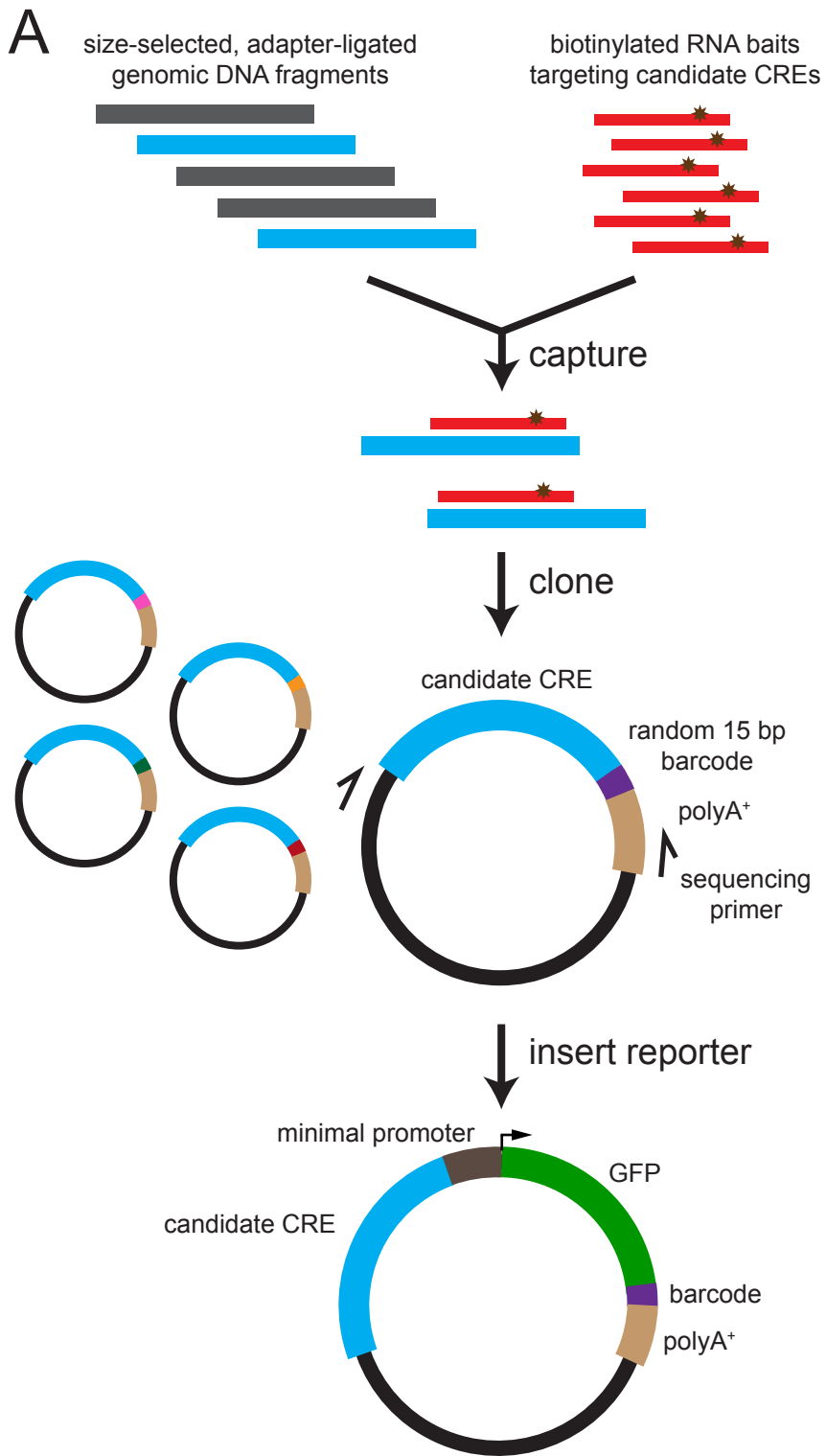
- 1123 Gnirke A, Melnikov A, Maguire J, Rogov P, LeProust EM, Brockman W, Fennell T, Giannoukos G, Fisher S,  
 1124 Russ C et al. 2009. Solution hybrid selection with ultra-long oligonucleotides for massively  
 1125 parallel targeted sequencing. *Nature biotechnology* **27**(2): 182-189.
- 1126 Graur D, Zheng Y, Price N, Azevedo RB, Zufall RA, Elhaik E. 2013. On the immortality of television sets:  
 1127 "function" in the human genome according to the evolution-free gospel of ENCODE. *Genome*  
 1128 *biology and evolution* **5**(3): 578-590.
- 1129 Grieger JC, Choi VW, Samulski RJ. 2006. Production and characterization of adeno-associated viral  
 1130 vectors. *Nature protocols* **1**(3): 1412-1428.
- 1131 Gross DS, Garrard WT. 1988. Nuclease hypersensitive sites in chromatin. *Annual review of biochemistry*  
 1132 **57**: 159-197.
- 1133 Heinz S, Benner C, Spann N, Bertolino E, Lin YC, Laslo P, Cheng JX, Murre C, Singh H, Glass CK. 2010.  
 1134 Simple combinations of lineage-determining transcription factors prime cis-regulatory elements  
 1135 required for macrophage and B cell identities. *Molecular cell* **38**(4): 576-589.
- 1136 Heinz S, Romanoski CE, Benner C, Glass CK. 2015. The selection and function of cell type-specific  
 1137 enhancers. *Nature reviews Molecular cell biology* **16**(3): 144-154.
- 1138 Herweijer H, Wolff JA. 2007. Gene therapy progress and prospects: hydrodynamic gene delivery. *Gene*  
 1139 *therapy* **14**(2): 99-107.
- 1140 Hsiao TH, Diaconu C, Myers CA, Lee J, Cepko CL, Corbo JC. 2007. The cis-regulatory logic of the  
 1141 mammalian photoreceptor transcriptional network. *PLoS one* **2**(7): e643.
- 1142 Hughes AL, Rando OJ. 2014. Mechanisms underlying nucleosome positioning in vivo. *Annual review of*  
 1143 *biophysics* **43**: 41-63.
- 1144 Jeon CJ, Strettoi E, Masland RH. 1998. The major cell populations of the mouse retina. *The Journal of*  
 1145 *neuroscience : the official journal of the Society for Neuroscience* **18**(21): 8936-8946.
- 1146 Karolchik D, Barber GP, Casper J, Clawson H, Cline MS, Diekhans M, Dreszer TR, Fujita PA, Guruvadoo L,  
 1147 Haeussler M et al. 2014. The UCSC Genome Browser database: 2014 update. *Nucleic acids*  
 1148 *research* **42**(Database issue): D764-770.
- 1149 Karra D, Dahm R. 2010. Transfection techniques for neuronal cells. *The Journal of neuroscience : the*  
 1150 *official journal of the Society for Neuroscience* **30**(18): 6171-6177.
- 1151 Kautzmann MA, Kim DS, Felder-Schmittbuhl MP, Swaroop A. 2011. Combinatorial regulation of  
 1152 photoreceptor differentiation factor, neural retina leucine zipper gene NRL, revealed by in vivo  
 1153 promoter analysis. *The Journal of biological chemistry* **286**(32): 28247-28255.
- 1154 Kellis M, Wold B, Snyder MP, Bernstein BE, Kundaje A, Marinov GK, Ward LD, Birney E, Crawford GE,  
 1155 Dekker J et al. 2014. Defining functional DNA elements in the human genome. *Proceedings of*  
 1156 *the National Academy of Sciences of the United States of America* **111**(17): 6131-6138.
- 1157 Kim EJ, Battiste J, Nakagawa Y, Johnson JE. 2008. Ascl1 (Mash1) lineage cells contribute to discrete cell  
 1158 populations in CNS architecture. *Molecular and cellular neurosciences* **38**(4): 595-606.
- 1159 Kumar M, Keller B, Makalou N, Sutton RE. 2001. Systematic determination of the packaging limit of  
 1160 lentiviral vectors. *Human gene therapy* **12**(15): 1893-1905.
- 1161 Kwasniewski JC, Fiore C, Chaudhari HG, Cohen BA. 2014. High-throughput functional testing of ENCODE  
 1162 segmentation predictions. *Genome research* **24**(10): 1595-1602.
- 1163 Kwasniewski JC, Mogno I, Myers CA, Corbo JC, Cohen BA. 2012. Complex effects of nucleotide variants in a  
 1164 mammalian cis-regulatory element. *Proceedings of the National Academy of Sciences of the*  
 1165 *United States of America* **109**(47): 19498-19503.
- 1166 Lancaster MA, Renner M, Martin CA, Wenzel D, Bicknell LS, Hurles ME, Homfray T, Penninger JM,  
 1167 Jackson AP, Knoblich JA. 2013. Cerebral organoids model human brain development and  
 1168 microcephaly. *Nature* **501**(7467): 373-379.
- 1169 Langmead B, Salzberg SL. 2012. Fast gapped-read alignment with Bowtie 2. *Nature methods* **9**(4): 357-  
 1170 359.

- 1171 Lee H, O'Connor BD, Merriman B, Funari VA, Homer N, Chen Z, Cohn DH, Nelson SF. 2009. Improving the  
 1172 efficiency of genomic loci capture using oligonucleotide arrays for high throughput resequencing.  
 1173 *BMC genomics* **10**: 646.
- 1174 Lee J, Myers CA, Williams N, Abdelaziz M, Corbo JC. 2010. Quantitative fine-tuning of photoreceptor cis-  
 1175 regulatory elements through affinity modulation of transcription factor binding sites. *Gene*  
 1176 *therapy* **17**(11): 1390-1399.
- 1177 Levo M, Segal E. 2014. In pursuit of design principles of regulatory sequences. *Nature reviews Genetics*  
 1178 **15**(7): 453-468.
- 1179 Li H, Handsaker B, Wysoker A, Fennell T, Ruan J, Homer N, Marth G, Abecasis G, Durbin R, Genome  
 1180 Project Data Processing S. 2009. The Sequence Alignment/Map format and SAMtools.  
 1181 *Bioinformatics* **25**(16): 2078-2079.
- 1182 Livesey FJ, Cepko CL. 2001. Vertebrate neural cell-fate determination: lessons from the retina. *Nature*  
 1183 *reviews Neuroscience* **2**(2): 109-118.
- 1184 London A, Benhar I, Schwartz M. 2013. The retina as a window to the brain-from eye research to CNS  
 1185 disorders. *Nature reviews Neurology* **9**(1): 44-53.
- 1186 Massari ME, Murre C. 2000. Helix-loop-helix proteins: regulators of transcription in eucaryotic organisms.  
 1187 *Molecular and cellular biology* **20**(2): 429-440.
- 1188 McCarty DM. 2008. Self-complementary AAV vectors; advances and applications. *Molecular therapy :  
 1189 the journal of the American Society of Gene Therapy* **16**(10): 1648-1656.
- 1190 McCarty DM, Young SM, Jr., Samulski RJ. 2004. Integration of adeno-associated virus (AAV) and  
 1191 recombinant AAV vectors. *Annual review of genetics* **38**: 819-845.
- 1192 McLean CY, Bristor D, Hiller M, Clarke SL, Schaar BT, Lowe CB, Wenger AM, Bejerano G. 2010. GREAT  
 1193 improves functional interpretation of cis-regulatory regions. *Nature biotechnology* **28**(5): 495-  
 1194 501.
- 1195 Mears AJ, Kondo M, Swain PK, Takada Y, Bush RA, Saunders TL, Sieving PA, Swaroop A. 2001. Nrl is  
 1196 required for rod photoreceptor development. *Nature genetics* **29**(4): 447-452.
- 1197 Melnikov A, Murugan A, Zhang X, Tesileanu T, Wang L, Rogov P, Feizi S, Gnirke A, Callan CG, Jr., Kinney JB  
 1198 et al. 2012. Systematic dissection and optimization of inducible enhancers in human cells using a  
 1199 massively parallel reporter assay. *Nature biotechnology* **30**(3): 271-277.
- 1200 Mertes F, Elsharawy A, Sauer S, van Helvoort JM, van der Zaag PJ, Franke A, Nilsson M, Lehrach H,  
 1201 Brookes AJ. 2011. Targeted enrichment of genomic DNA regions for next-generation sequencing.  
 1202 *Briefings in functional genomics* **10**(6): 374-386.
- 1203 Mi H, Lazareva-Ulitsky B, Loo R, Kejariwal A, Vandergriff J, Rabkin S, Guo N, Muruganujan A, Doremiex  
 1204 O, Campbell MJ et al. 2005. The PANTHER database of protein families, subfamilies, functions  
 1205 and pathways. *Nucleic acids research* **33**(Database issue): D284-288.
- 1206 Mingozzi F, High KA. 2011. Therapeutic in vivo gene transfer for genetic disease using AAV: progress and  
 1207 challenges. *Nature reviews Genetics* **12**(5): 341-355.
- 1208 Mongrain V, La Spada F, Curie T, Franken P. 2011. Sleep loss reduces the DNA-binding of BMAL1, CLOCK,  
 1209 and NPAS2 to specific clock genes in the mouse cerebral cortex. *PLoS one* **6**(10): e26622.
- 1210 Montana CL, Lawrence KA, Williams NL, Tran NM, Peng GH, Chen S, Corbo JC. 2011a. Transcriptional  
 1211 regulation of neural retina leucine zipper (Nrl), a photoreceptor cell fate determinant. *The*  
 1212 *Journal of biological chemistry* **286**(42): 36921-36931.
- 1213 Montana CL, Myers CA, Corbo JC. 2011b. Quantifying the activity of cis-regulatory elements in the  
 1214 mouse retina by explant electroporation. *Journal of visualized experiments : JoVE*(52).
- 1215 Mortimer I, Tam P, MacLachlan I, Graham RW, Saravolac EG, Joshi PB. 1999. Cationic lipid-mediated  
 1216 transfection of cells in culture requires mitotic activity. *Gene therapy* **6**(3): 403-411.
- 1217 Mullen RJ, Buck CR, Smith AM. 1992. NeuN, a neuronal specific nuclear protein in vertebrates.  
 1218 *Development* **116**(1): 201-211.

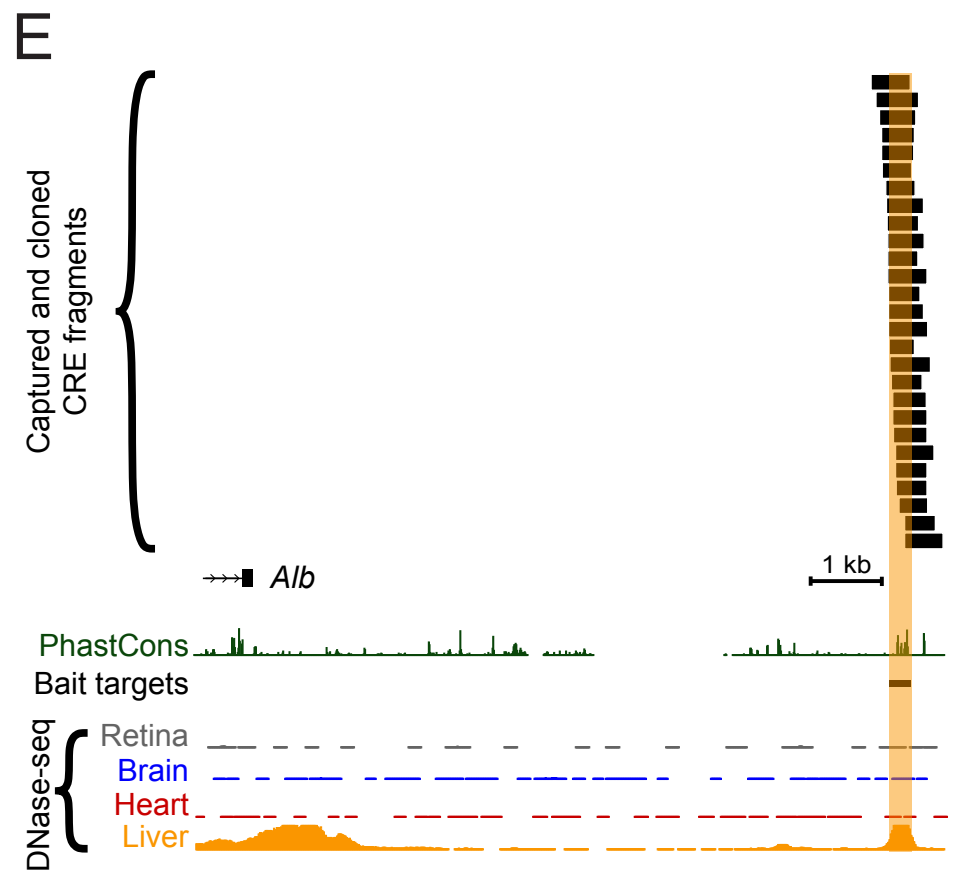
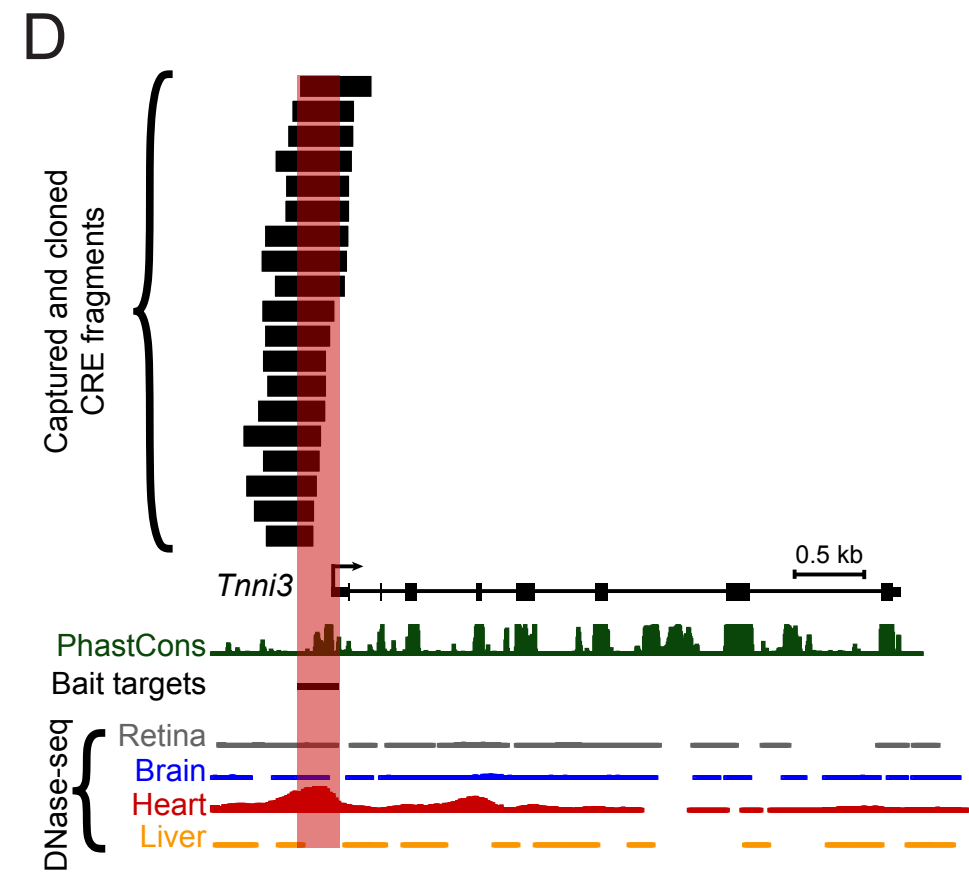
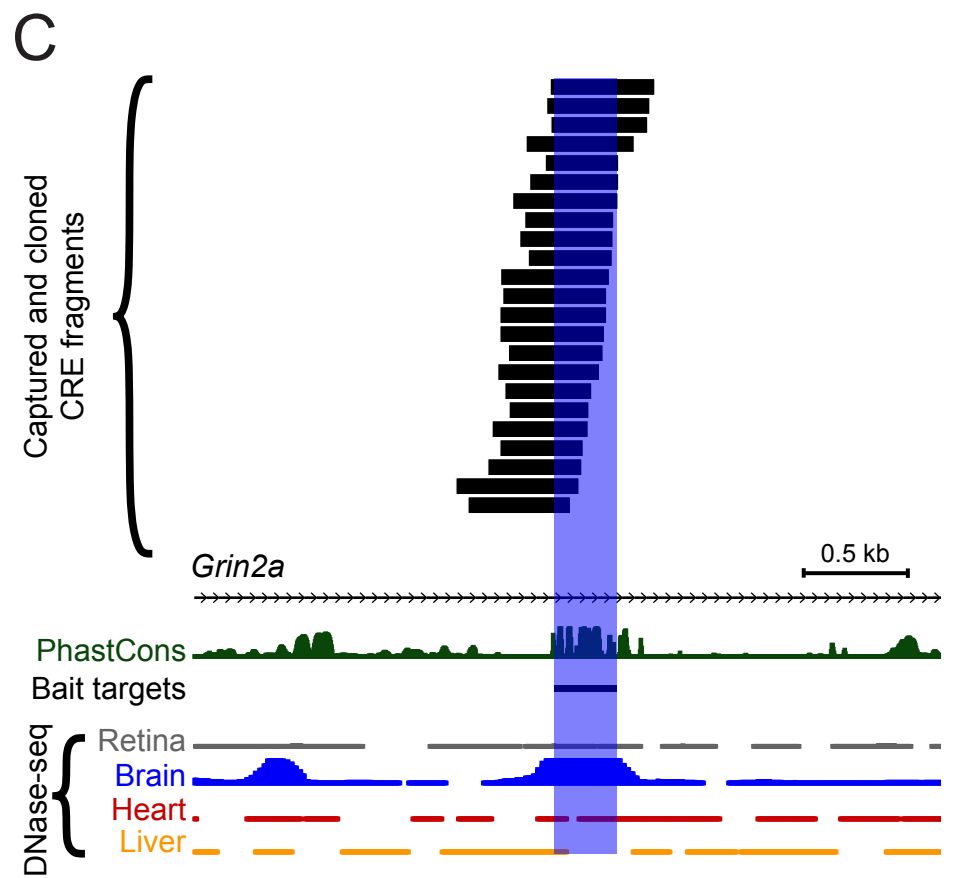
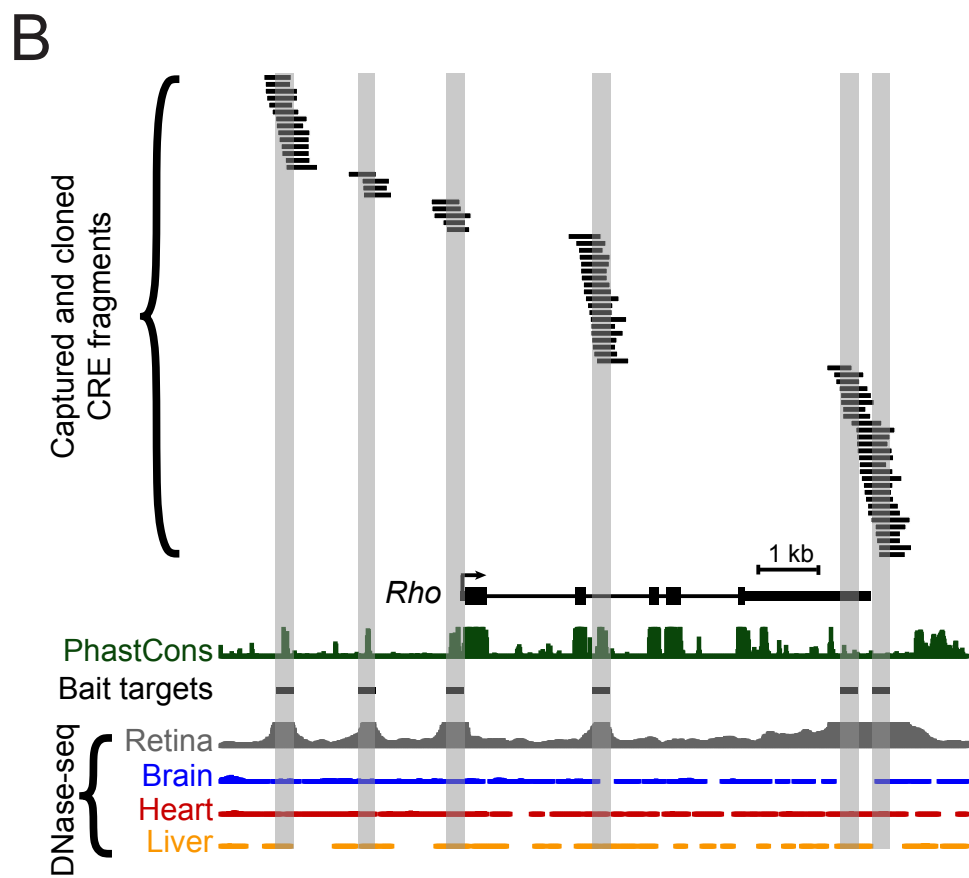
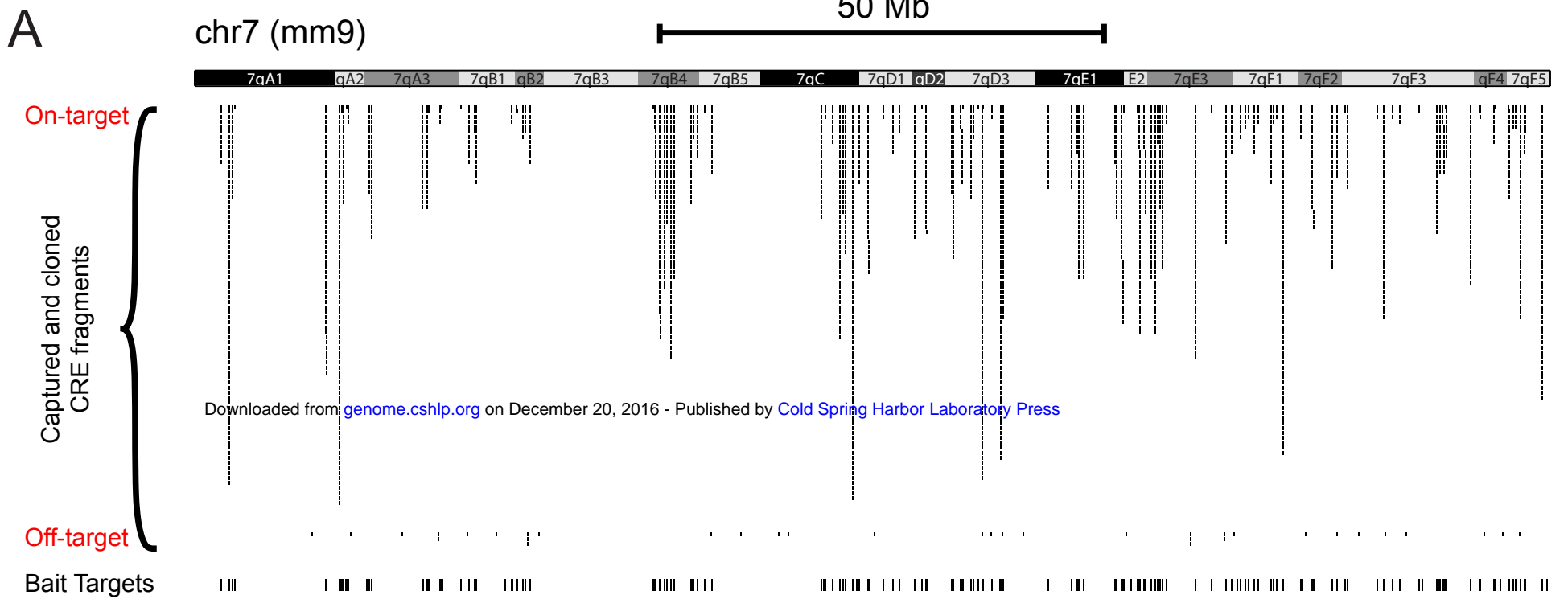
- 1219 Murtha M, Tokcaer-Keskin Z, Tang Z, Strino F, Chen X, Wang Y, Xi X, Basilico C, Brown S, Bonneau R et al.  
 1220 2014. FIREWACH: high-throughput functional detection of transcriptional regulatory modules in  
 1221 mammalian cells. *Nature methods* **11**(5): 559-565.
- 1222 Natarajan A, Yardimci GG, Sheffield NC, Crawford GE, Ohler U. 2012. Predicting cell-type-specific gene  
 1223 expression from regions of open chromatin. *Genome research* **22**(9): 1711-1722.
- 1224 Nord AS, Blow MJ, Attanasio C, Akiyama JA, Holt A, Hosseini R, Phouanavong S, Plajzer-Frick I, Shoukry  
 1225 M, Afzal V et al. 2013. Rapid and pervasive changes in genome-wide enhancer usage during  
 1226 mammalian development. *Cell* **155**(7): 1521-1531.
- 1227 Nord AS, Pattabiraman K, Visel A, Rubenstein JL. 2015. Genomic Perspectives of Transcriptional  
 1228 Regulation in Forebrain Development. *Neuron* **85**(1): 27-47.
- 1229 Okaty BW, Sugino K, Nelson SB. 2011. Cell type-specific transcriptomics in the brain. *The Journal of*  
 1230 *neuroscience : the official journal of the Society for Neuroscience* **31**(19): 6939-6943.
- 1231 Ostuni R, Piccolo V, Barozzi I, Polletti S, Termanini A, Bonifacio S, Curina A, Prosperini E, Ghisletti S,  
 1232 Natoli G. 2013. Latent enhancers activated by stimulation in differentiated cells. *Cell* **152**(1-2):  
 1233 157-171.
- 1234 Patwardhan RP, Hiatt JB, Witten DM, Kim MJ, Smith RP, May D, Lee C, Andrie JM, Lee SI, Cooper GM et  
 1235 al. 2012. Massively parallel functional dissection of mammalian enhancers in vivo. *Nature*  
 1236 *biotechnology* **30**(3): 265-270.
- 1237 Patwardhan RP, Lee C, Litvin O, Young DL, Pe'er D, Shendure J. 2009. High-resolution analysis of DNA  
 1238 regulatory elements by synthetic saturation mutagenesis. *Nature biotechnology* **27**(12): 1173-  
 1239 1175.
- 1240 Penaud-Budloo M, Le Guiner C, Nowrouzi A, Toromanoff A, Cherel Y, Chenuaud P, Schmidt M, von Kalle  
 1241 C, Rolling F, Moullier P et al. 2008. Adeno-associated virus vector genomes persist as episomal  
 1242 chromatin in primate muscle. *Journal of virology* **82**(16): 7875-7885.
- 1243 Quinlan AR, Hall IM. 2010. BEDTools: a flexible suite of utilities for comparing genomic features.  
 1244 *Bioinformatics* **26**(6): 841-842.
- 1245 Raivich G, Behrens A. 2006. Role of the AP-1 transcription factor c-Jun in developing, adult and injured  
 1246 brain. *Progress in neurobiology* **78**(6): 347-363.
- 1247 Rands CM, Meader S, Ponting CP, Lunter G. 2014. 8.2% of the Human genome is constrained: variation  
 1248 in rates of turnover across functional element classes in the human lineage. *PLoS genetics* **10**(7):  
 1249 e1004525.
- 1250 Reynolds N, O'Shaughnessy A, Hendrich B. 2013. Transcriptional repressors: multifaceted regulators of  
 1251 gene expression. *Development* **140**(3): 505-512.
- 1252 Romanoski CE, Glass CK, Stunnenberg HG, Wilson L, Almouzni G. 2015. Epigenomics: Roadmap for  
 1253 regulation. *Nature* **518**(7539): 314-316.
- 1254 Sandelin A, Carninci P, Lenhard B, Ponjavic J, Hayashizaki Y, Hume DA. 2007. Mammalian RNA  
 1255 polymerase II core promoters: insights from genome-wide studies. *Nature reviews Genetics* **8**(6):  
 1256 424-436.
- 1257 Selever J, Kong JQ, Arenkiel BR. 2011. A rapid approach to high-resolution fluorescence imaging in semi-  
 1258 thick brain slices. *Journal of visualized experiments : JoVE*(53).
- 1259 Shen Y, Yue F, McCleary DF, Ye Z, Edsall L, Kuan S, Wagner U, Dixon J, Lee L, Lobanenkov VV et al. 2012.  
 1260 A map of the cis-regulatory sequences in the mouse genome. *Nature* **488**(7409): 116-120.
- 1261 Shlyueva D, Stampfel G, Stark A. 2014. Transcriptional enhancers: from properties to genome-wide  
 1262 predictions. *Nature reviews Genetics* **15**(4): 272-286.
- 1263 Shu W, Chen H, Bo X, Wang S. 2011. Genome-wide analysis of the relationships between DNaseI HS,  
 1264 histone modifications and gene expression reveals distinct modes of chromatin domains. *Nucleic*  
 1265 *acids research* **39**(17): 7428-7443.

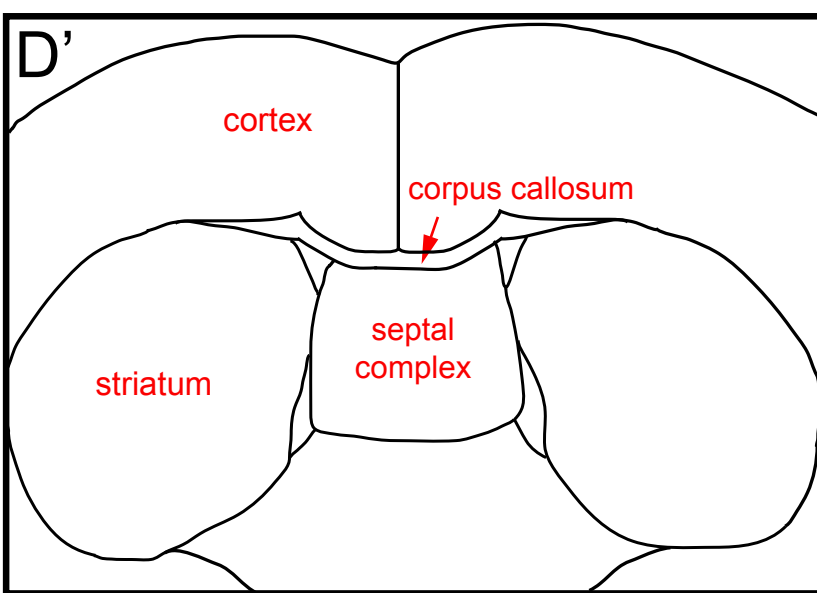
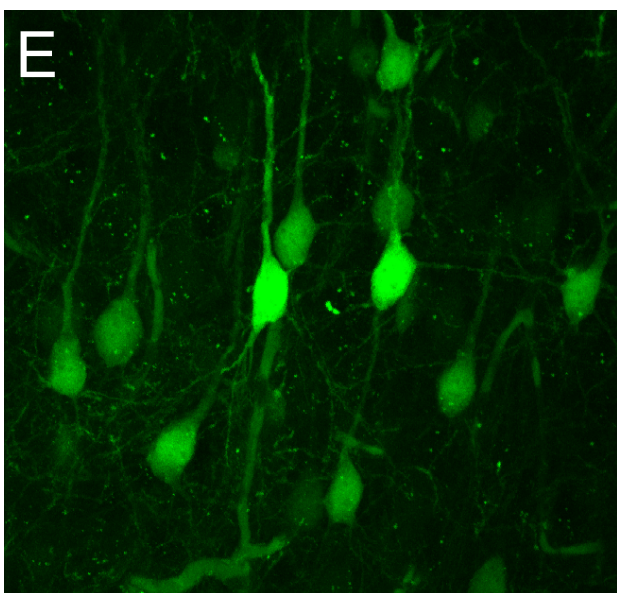
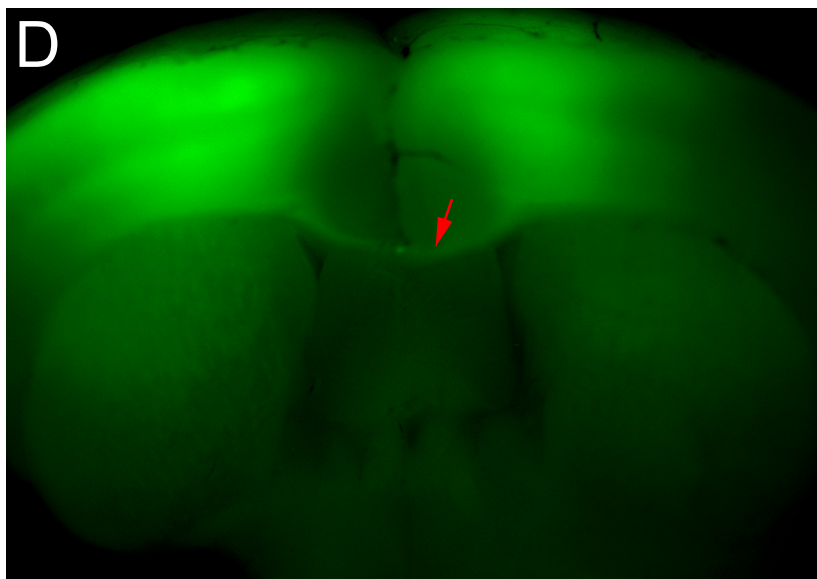
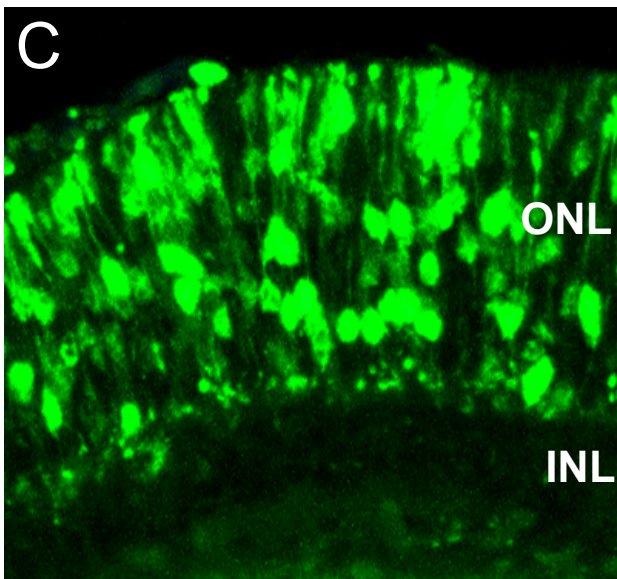
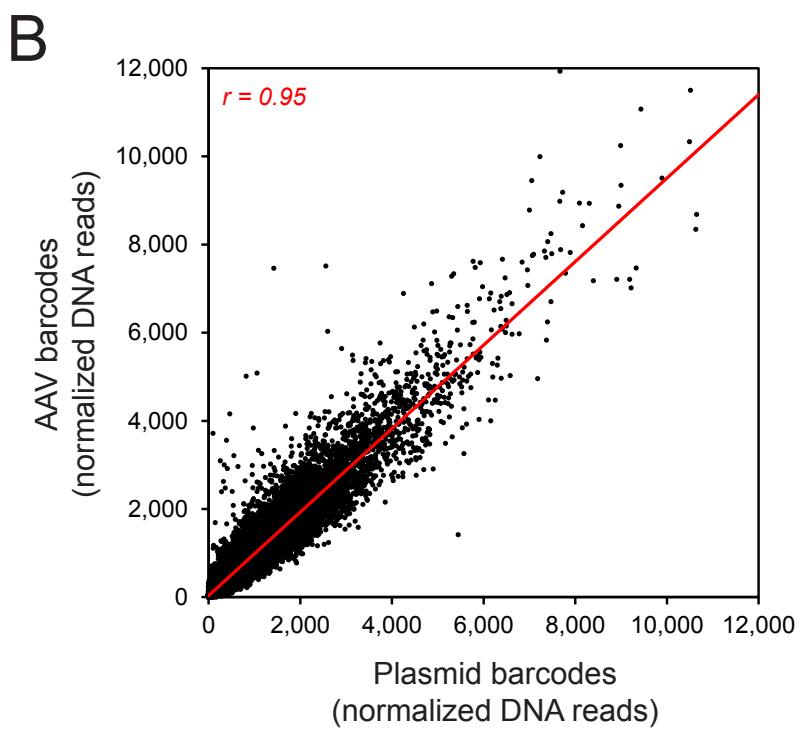
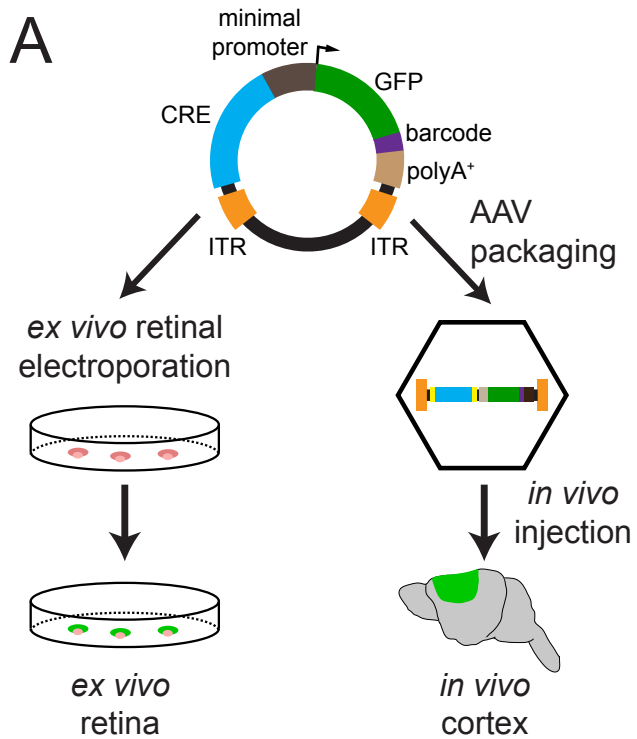
- 1266 Siepel A, Bejerano G, Pedersen JS, Hinrichs AS, Hou M, Rosenbloom K, Clawson H, Spieth J, Hillier LW,  
1267 Richards S et al. 2005. Evolutionarily conserved elements in vertebrate, insect, worm, and yeast  
1268 genomes. *Genome research* **15**(8): 1034-1050.
- 1269 Sing T, Sander O, Beerenwinkel N, Lengauer T. 2005. ROCr: visualizing classifier performance in R.  
1270 *Bioinformatics* **21**(20): 3940-3941.
- 1271 Swaroop A, Kim D, Forrest D. 2010. Transcriptional regulation of photoreceptor development and  
1272 homeostasis in the mammalian retina. *Nature reviews Neuroscience* **11**(8): 563-576.
- 1273 The ENCODE Project Consortium. 2012. An integrated encyclopedia of DNA elements in the human  
1274 genome. *Nature* **489**(7414): 57-74.
- 1275 Thurman RE, Rynes E, Humbert R, Vierstra J, Maurano MT, Haugen E, Sheffield NC, Stergachis AB, Wang  
1276 H, Vernot B et al. 2012. The accessible chromatin landscape of the human genome. *Nature*  
1277 **489**(7414): 75-82.
- 1278 Tillo D, Hughes TR. 2009. G+C content dominates intrinsic nucleosome occupancy. *BMC bioinformatics*  
1279 **10**: 442.
- 1280 Tillo D, Kaplan N, Moore IK, Fondufe-Mittendorf Y, Gossett AJ, Field Y, Lieb JD, Widom J, Segal E, Hughes  
1281 TR. 2010. High nucleosome occupancy is encoded at human regulatory sequences. *PloS one* **5**(2):  
1282 e9129.
- 1283 van Arensbergen J, van Steensel B, Bussemaker HJ. 2014. In search of the determinants of enhancer-  
1284 promoter interaction specificity. *Trends in cell biology* **24**(11): 695-702.
- 1285 Vanhille L, Griffon A, Maqbool MA, Zacarias-Cabeza J, Dao LT, Fernandez N, Ballester B, Andrau JC,  
1286 Spicuglia S. 2015. High-throughput and quantitative assessment of enhancer activity in  
1287 mammals by CapStarr-seq. *Nature communications* **6**: 6905.
- 1288 Vierstra J, Rynes E, Sandstrom R, Zhang M, Canfield T, Hansen RS, Stehling-Sun S, Sabo PJ, Byron R,  
1289 Humbert R et al. 2014. Mouse regulatory DNA landscapes reveal global principles of cis-  
1290 regulatory evolution. *Science* **346**(6212): 1007-1012.
- 1291 Visel A, Blow MJ, Li Z, Zhang T, Akiyama JA, Holt A, Plajzer-Frick I, Shoukry M, Wright C, Chen F et al.  
1292 2009. ChIP-seq accurately predicts tissue-specific activity of enhancers. *Nature* **457**(7231): 854-  
1293 858.
- 1294 Visel A, Taher L, Girgis H, May D, Golonzhka O, Hoch RV, McKinsey GL, Pattabiraman K, Silberberg SN,  
1295 Blow MJ et al. 2013. A high-resolution enhancer atlas of the developing telencephalon. *Cell*  
1296 **152**(4): 895-908.
- 1297 Wang J, Zhuang J, Iyer S, Lin X, Whitfield TW, Greven MC, Pierce BG, Dong X, Kundaje A, Cheng Y et al.  
1298 2012. Sequence features and chromatin structure around the genomic regions bound by 119  
1299 human transcription factors. *Genome research* **22**(9): 1798-1812.
- 1300 Ward LD, Kellis M. 2012. Interpreting noncoding genetic variation in complex traits and human disease.  
1301 *Nature biotechnology* **30**(11): 1095-1106.
- 1302 White MA, Myers CA, Corbo JC, Cohen BA. 2013. Massively parallel in vivo enhancer assay reveals that  
1303 highly local features determine the cis-regulatory function of ChIP-seq peaks. *Proceedings of the*  
1304 *National Academy of Sciences of the United States of America* **110**(29): 11952-11957.
- 1305 Wilken MSB, J.A.; La Torre, A.; Siebenthall, K.; Thurman R.; Sabo, P.; Sandstrom, R.S.; Vierstra, J.;  
1306 Canfield, T.K.; Hansen, R.S.; Bender, M.A.; Stamatoyannopoulos, J.; Reh, T.A. 2015. DNase I  
1307 hypersensitivity analysis of the mouse brain and retina identifies region-specific regulatory  
1308 elements. *Epigenetics & Chromatin* **8**(8).
- 1309 Wray GA. 2007. The evolutionary significance of cis-regulatory mutations. *Nature reviews Genetics* **8**(3):  
1310 206-216.
- 1311 Wright AF, Chakarova CF, Abd El-Aziz MM, Bhattacharya SS. 2010. Photoreceptor degeneration: genetic  
1312 and mechanistic dissection of a complex trait. *Nature reviews Genetics* **11**(4): 273-284.

- 1313 Wu Z, Asokan A, Samulski RJ. 2006. Adeno-associated virus serotypes: vector toolkit for human gene  
1314 therapy. *Molecular therapy : the journal of the American Society of Gene Therapy* **14**(3): 316-327.
- 1315 Wu Z, Yang H, Colosi P. 2010. Effect of genome size on AAV vector packaging. *Molecular therapy : the*  
1316 *journal of the American Society of Gene Therapy* **18**(1): 80-86.
- 1317 Wurmbach E, Gonzalez-Maeso J, Yuen T, Ebersole BJ, Mastaitis JW, Mobbs CV, Sealfon SC. 2002.  
1318 Validated genomic approach to study differentially expressed genes in complex tissues.  
1319 *Neurochemical research* **27**(10): 1027-1033.
- 1320 Yan Z, Zak R, Zhang Y, Engelhardt JF. 2005. Inverted terminal repeat sequences are important for  
1321 intermolecular recombination and circularization of adeno-associated virus genomes. *Journal of*  
1322 *virology* **79**(1): 364-379.
- 1323 Yue F Cheng Y Breschi A Vierstra J Wu W Ryba T Sandstrom R Ma Z Davis C Pope BD et al. 2014. A  
1324 comparative encyclopedia of DNA elements in the mouse genome. *Nature* **515**(7527): 355-364.
- 1325 Zabidi MA, Arnold CD, Schernhuber K, Pagani M, Rath M, Frank O, Stark A. 2015. Enhancer-core-  
1326 promoter specificity separates developmental and housekeeping gene regulation. *Nature*  
1327 **518**(7540): 556-559.
- 1328 Zhang G, Gurtu V, Kain SR. 1996. An enhanced green fluorescent protein allows sensitive detection of  
1329 gene transfer in mammalian cells. *Biochemical and biophysical research communications* **227**(3):  
1330 707-711.
- 1331 Zhang Y, Liu T, Meyer CA, Eeckhoutte J, Johnson DS, Bernstein BE, Nusbaum C, Myers RM, Brown M, Li W  
1332 et al. 2008. Model-based analysis of ChIP-Seq (MACS). *Genome biology* **9**(9): R137.
- 1333 Zhong L, Li B, Mah CS, Govindasamy L, Agbandje-McKenna M, Cooper M, Herzog RW, Zolotukhin I,  
1334 Warrington KH, Jr., Weigel-Van Aken KA et al. 2008. Next generation of adeno-associated virus 2  
1335 vectors: point mutations in tyrosines lead to high-efficiency transduction at lower doses.  
1336 *Proceedings of the National Academy of Sciences of the United States of America* **105**(22): 7827-  
1337 7832.
- 1338 Zincarelli C, Soltys S, Rengo G, Rabinowitz JE. 2008. Analysis of AAV serotypes 1-9 mediated gene  
1339 expression and tropism in mice after systemic injection. *Molecular therapy : the journal of the*  
1340 *American Society of Gene Therapy* **16**(6): 1073-1080.
- 1341 Zolotukhin S, Potter M, Zolotukhin I, Sakai Y, Loiler S, Fraites TJ, Jr., Chiodo VA, Phillipsberg T, Muzyczka  
1342 N, Hauswirth WW et al. 2002. Production and purification of serotype 1, 2, and 5 recombinant  
1343 adeno-associated viral vectors. *Methods* **28**(2): 158-167.
- 1344



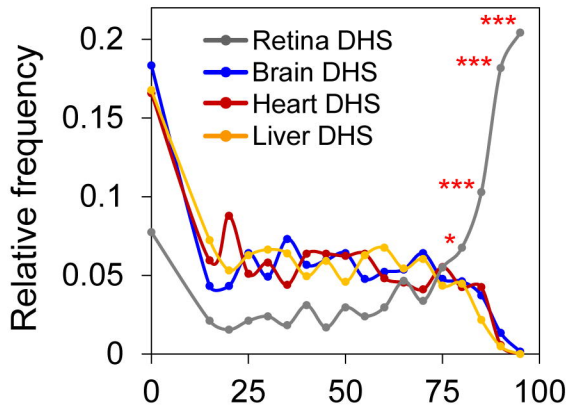




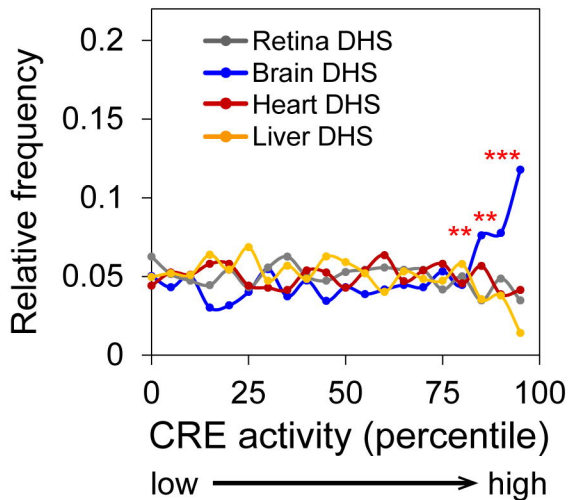


**A**

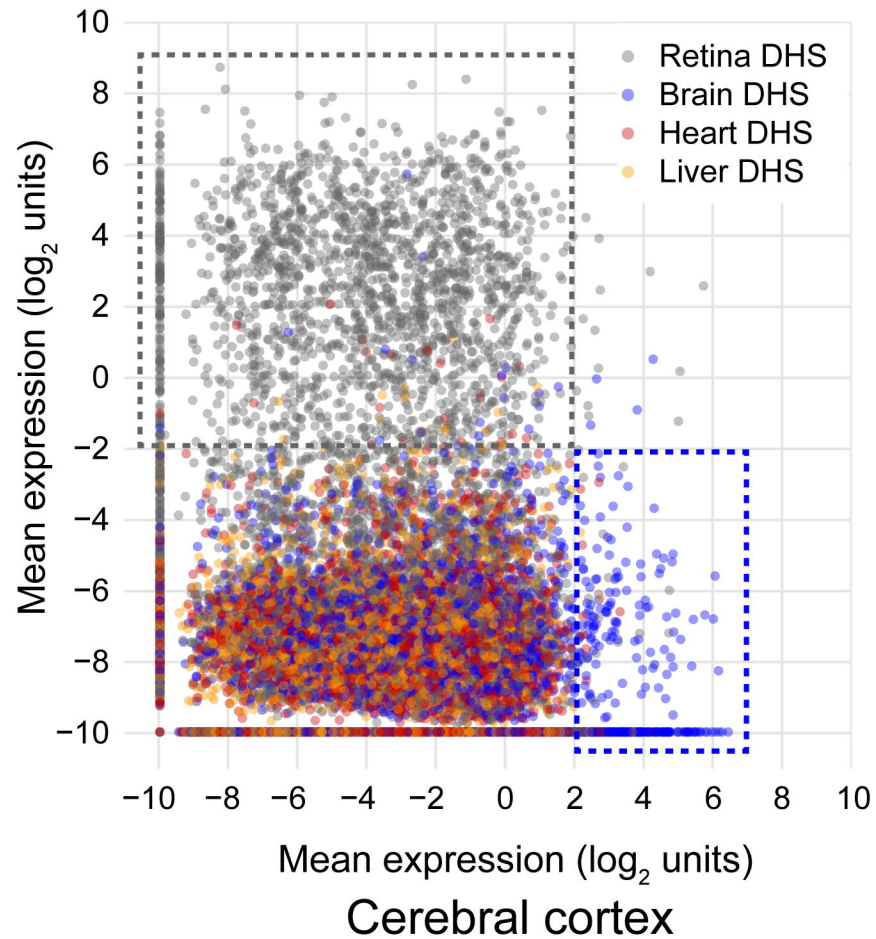
Retina



Cerebral cortex

**B**

Retina



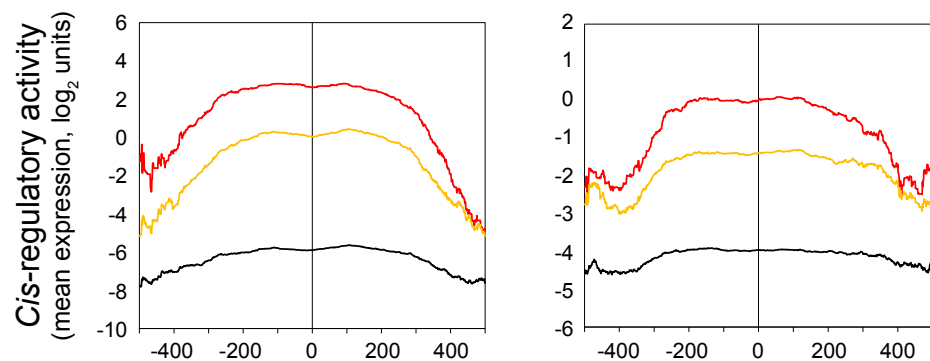
## Retina

Top 100 retina DHSs in retina  
 Top 200 retina DHSs in retina  
 All retina DHSs in retina

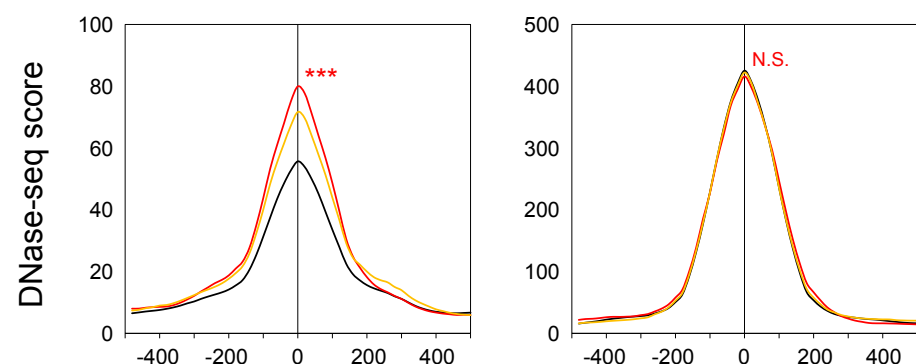
## Cerebral cortex

Top 100 brain DHSs in cortex  
 Top 200 brain DHSs in cortex  
 All brain DHSs in cortex

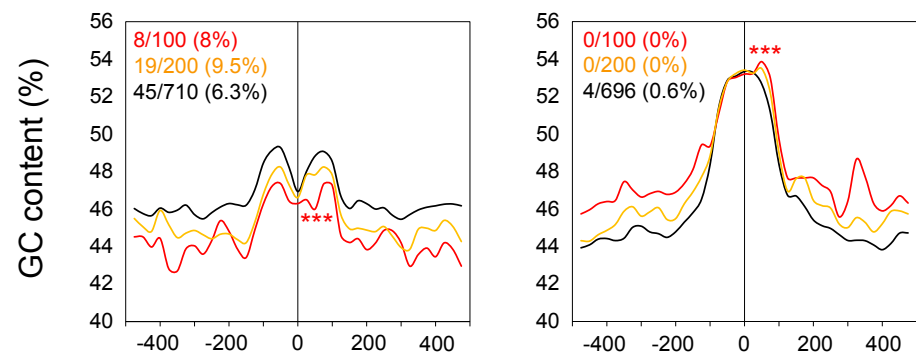
### A



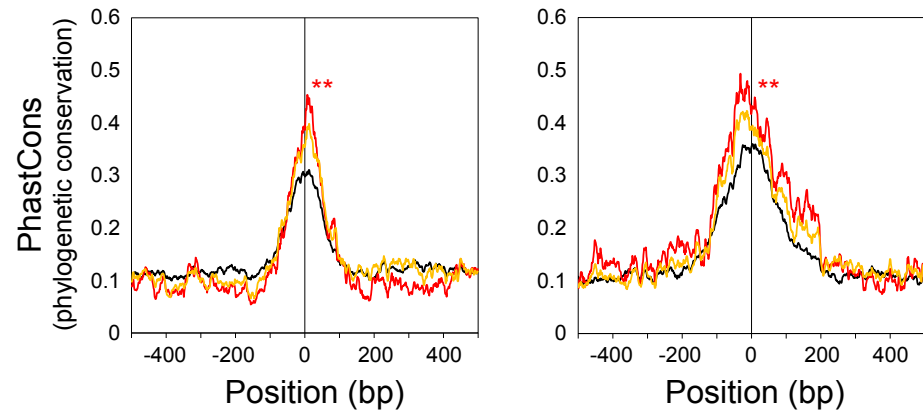
### B



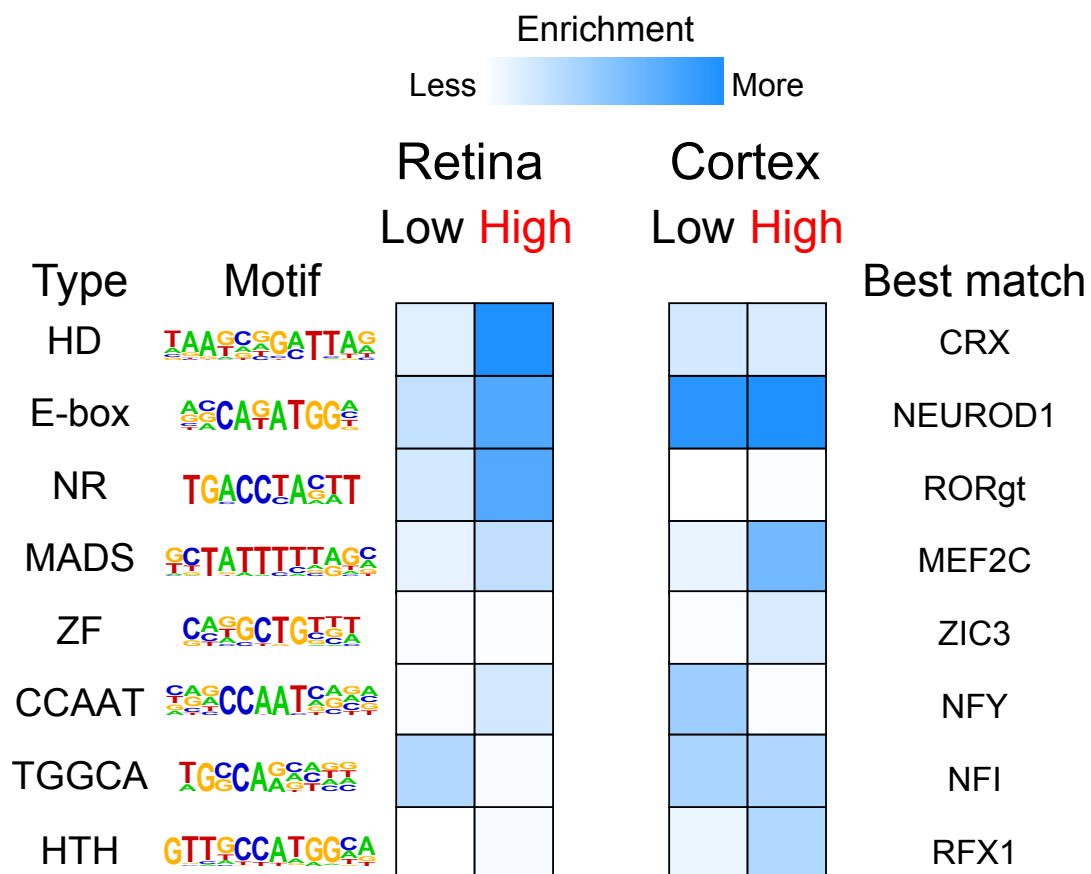
### C



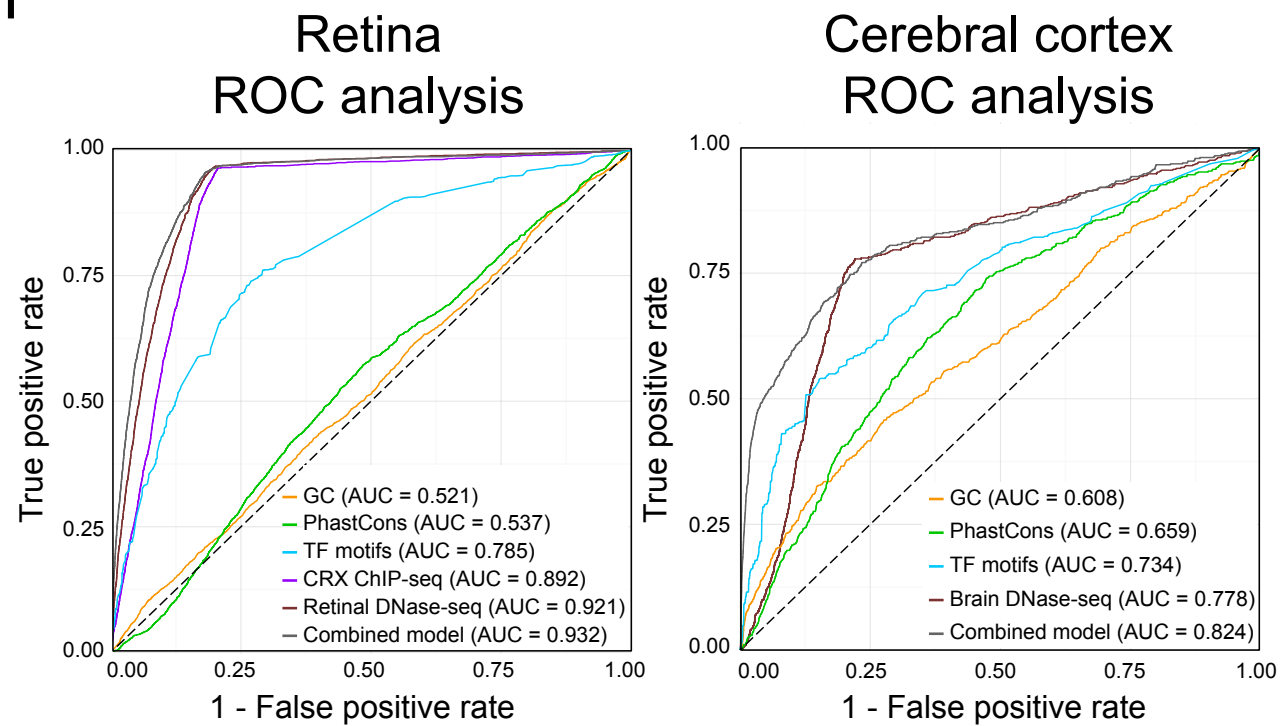
### D

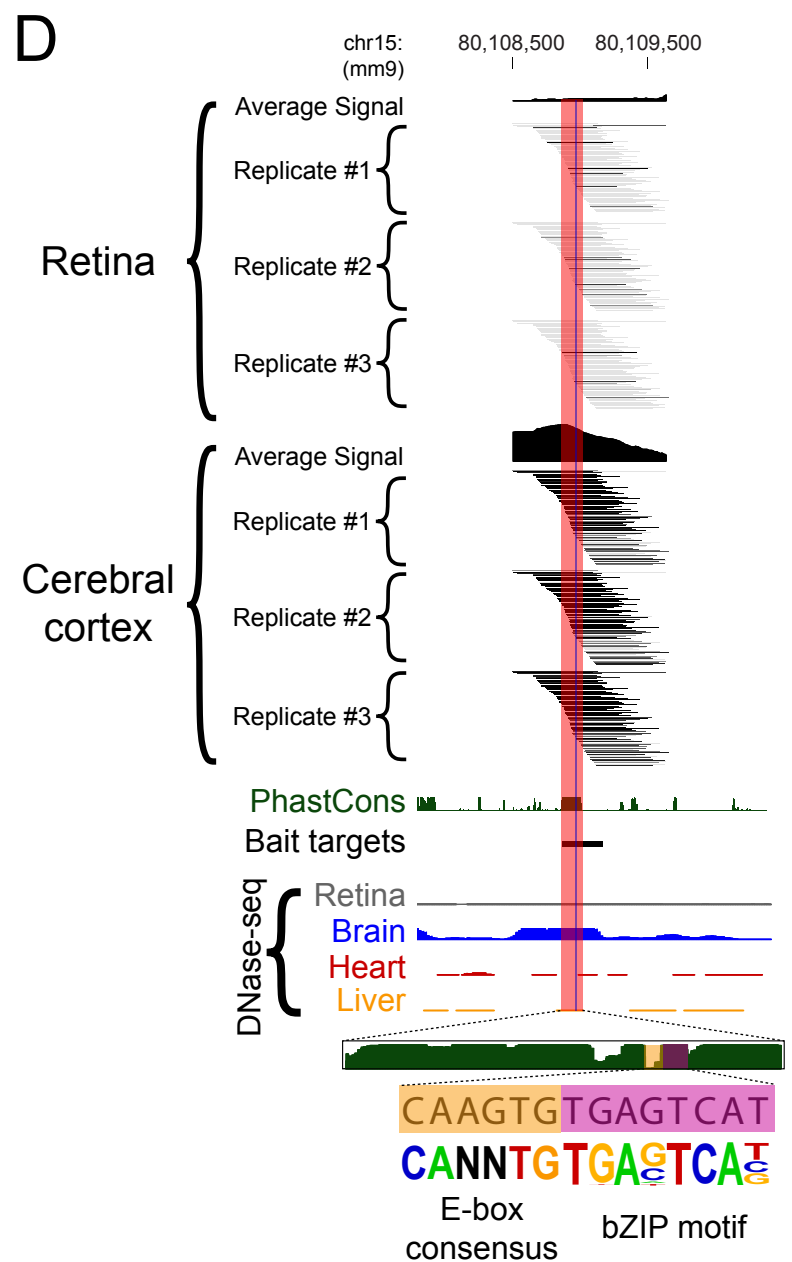
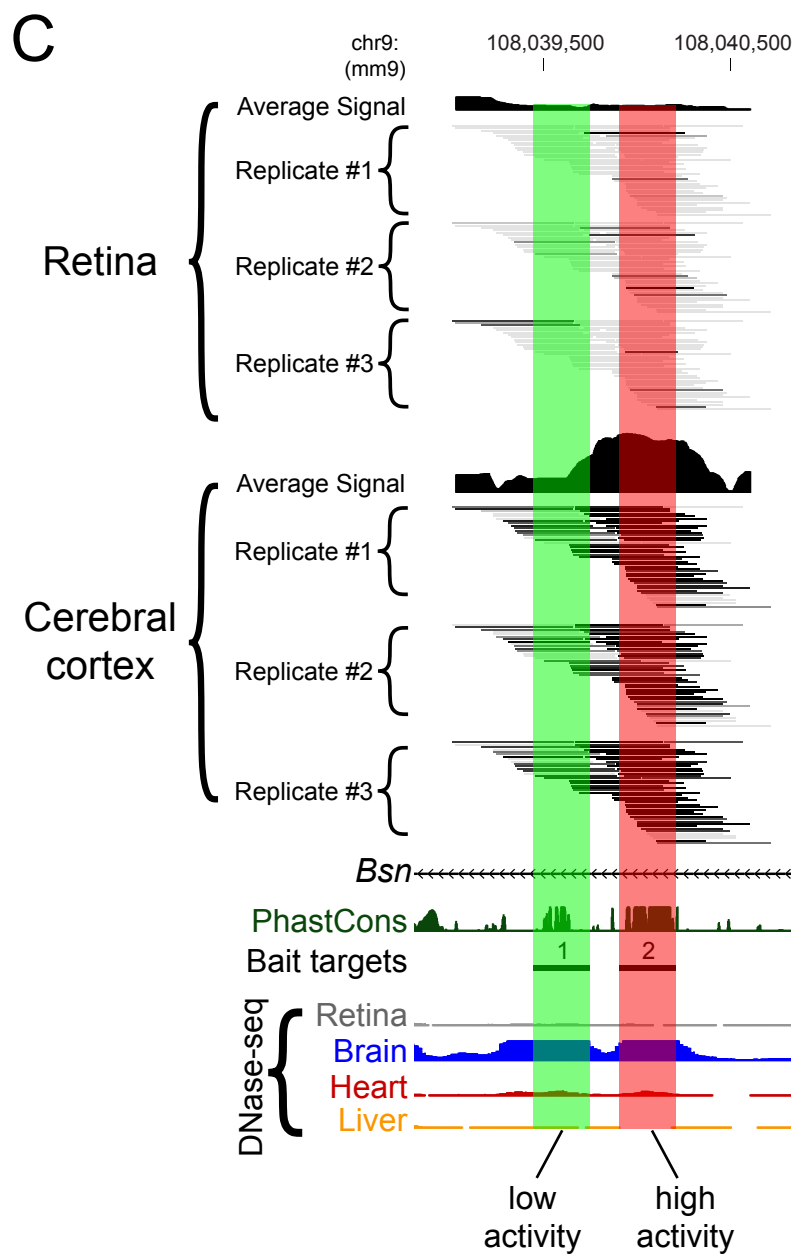
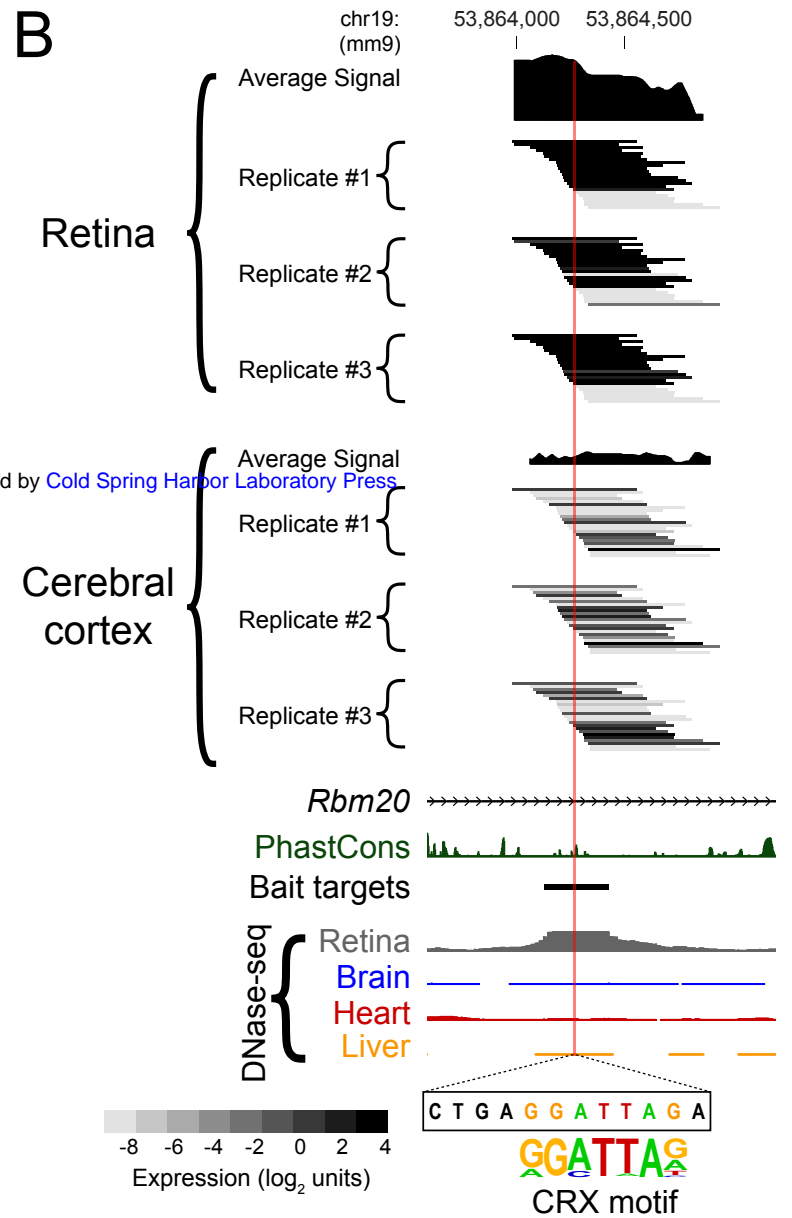
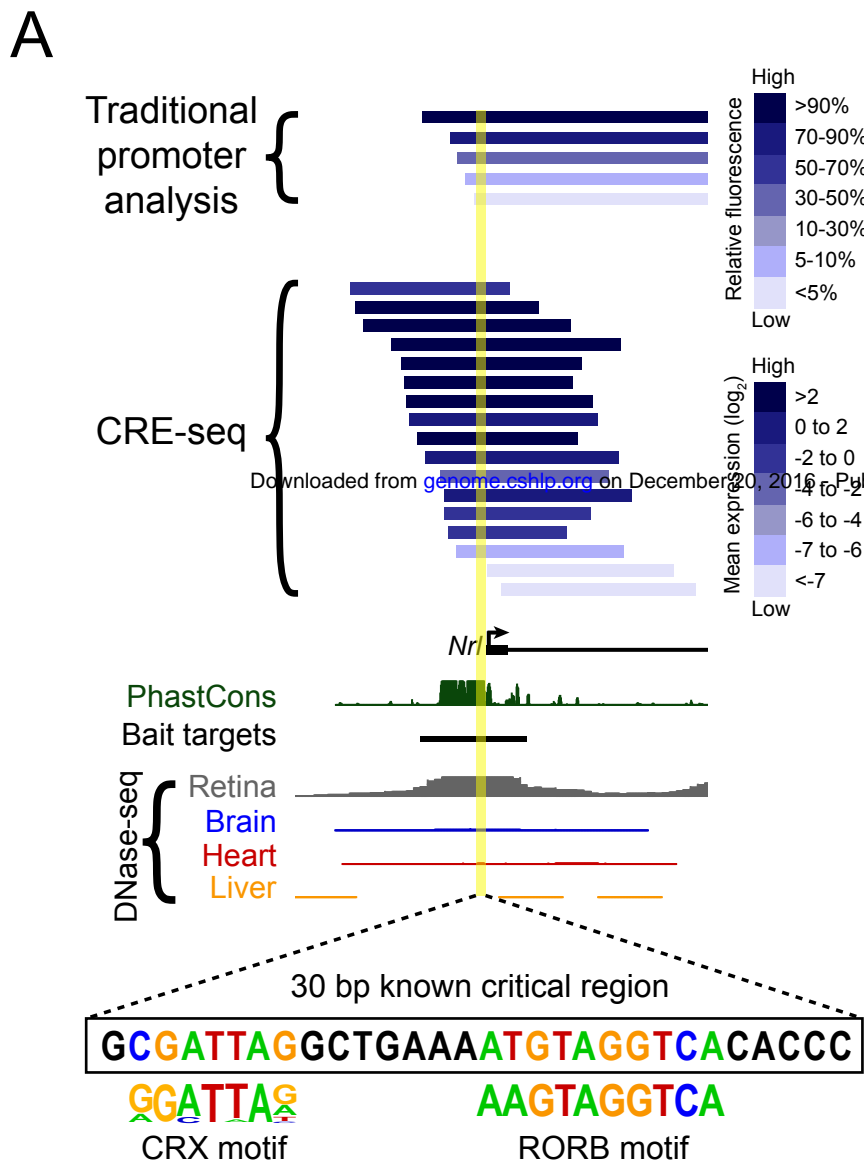


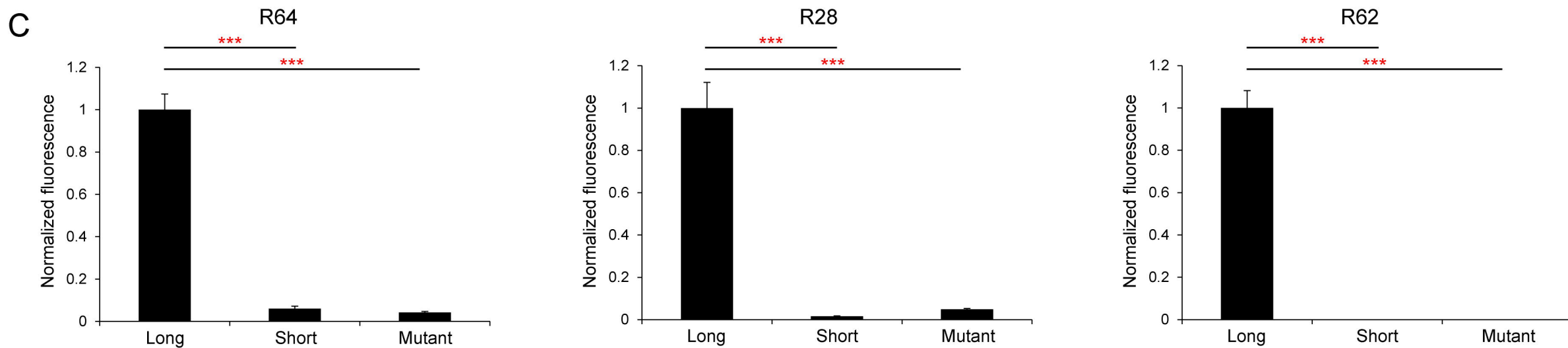
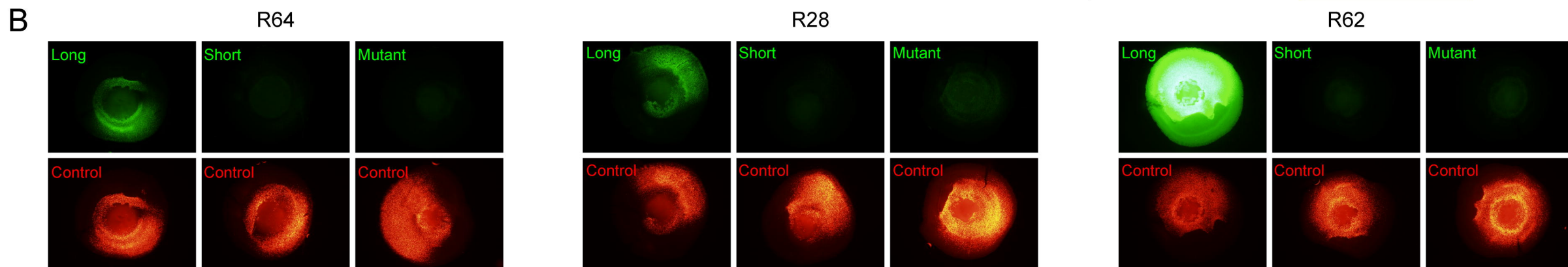
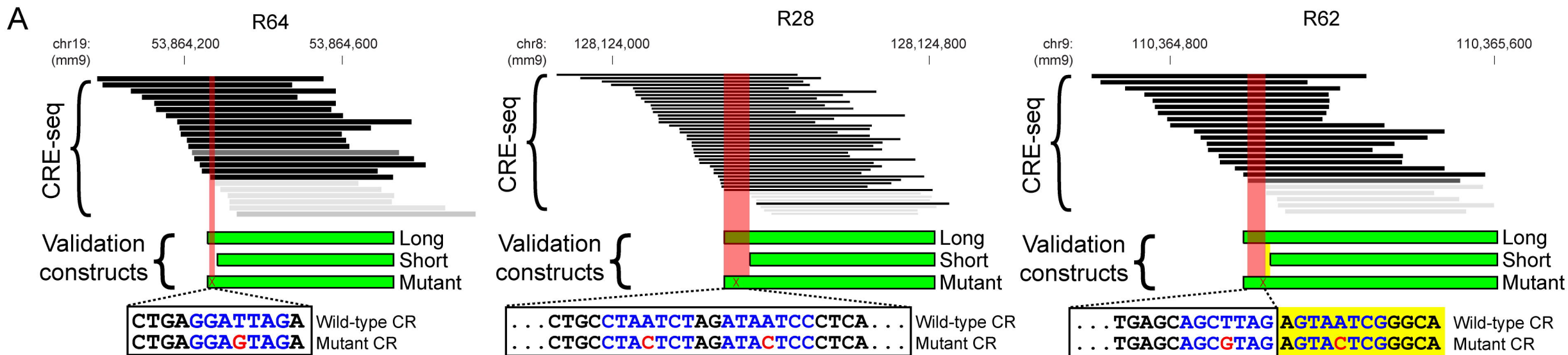
### E



### F









## Massively parallel *cis*-regulatory analysis in the mammalian central nervous system

Susan Q Shen, Connie A Myers, Andrew EO Hughes, et al.

*Genome Res.* published online November 17, 2015  
Access the most recent version at doi:[10.1101/gr.193789.115](https://doi.org/10.1101/gr.193789.115)

---

**Supplemental Material** <http://genome.cshlp.org/content/suppl/2015/11/17/gr.193789.115.DC1.html>

**P<P** Published online November 17, 2015 in advance of the print journal.

**Accepted Manuscript** Peer-reviewed and accepted for publication but not copyedited or typeset; accepted manuscript is likely to differ from the final, published version.

**Creative Commons License** This article is distributed exclusively by Cold Spring Harbor Laboratory Press for the first six months after the full-issue publication date (see <http://genome.cshlp.org/site/misc/terms.xhtml>). After six months, it is available under a Creative Commons License (Attribution-NonCommercial 4.0 International), as described at <http://creativecommons.org/licenses/by-nc/4.0/>.

**Email Alerting Service** Receive free email alerts when new articles cite this article - sign up in the box at the top right corner of the article or [click here](#).

---

---

Advance online articles have been peer reviewed and accepted for publication but have not yet appeared in the paper journal (edited, typeset versions may be posted when available prior to final publication). Advance online articles are citable and establish publication priority; they are indexed by PubMed from initial publication. Citations to Advance online articles must include the digital object identifier (DOIs) and date of initial publication.

---

To subscribe to *Genome Research* go to:  
<http://genome.cshlp.org/subscriptions>

---