

Systematic analysis of intrinsic enhancer-promoter compatibility in the mouse genome.

¹Miguel Martinez-Ara, ^{1,2}Federico Comoglio, ^{1,3}Joris van Arensbergen,
¹Bas van Steensel*

¹Division of Gene Regulation and OncoCode Institute, Netherlands Cancer Institute, Amsterdam, the Netherlands

²Present address: enGene Statistics GmbH, Basel, Switzerland

³Present address: Annogen B.V., Science park 406, Amsterdam, the Netherlands

*correspondence: b.v.steensel@nki.nl

Abstract

Gene expression is in part controlled by cis-regulatory elements (CREs) such as enhancers and repressive elements. Anecdotal evidence has indicated that a CRE and a promoter need to be biochemically compatible for promoter regulation to occur, but this compatibility has remained poorly characterised in mammalian cells. We used high-throughput combinatorial reporter assays to test thousands of CRE – promoter pairs from three Mb-sized genomic regions in mouse cells. This revealed that CREs vary substantially in their promoter compatibility, ranging from striking specificity for a single promoter to quantitative differences in activation across a broad set of promoters. More than half of the tested CREs exhibit significant promoter selectivity. Housekeeping promoters tend to have similar CRE preferences, but other promoters exhibit a wide diversity of compatibilities. Higher-order TF motif combinations may account for compatibility. CRE–promoter selectivity does not correlate with looping interactions in the native genomic context, suggesting that chromatin folding and compatibility are two orthogonal mechanisms that confer specificity to gene regulation.

Keywords

Enhancer, promoter, cis-regulatory element, repressor, massively parallel reporter assay, combinatorial, gene regulation, transcription.

37 INTRODUCTION

38

39 How genes are regulated by cis-regulatory elements (CREs) such as enhancers and repressor
40 elements is a long-standing topic in molecular biology [1-10]. One conundrum is how CREs
41 'choose' their target promoters. Some enhancers can activate multiple promoters *in cis* over short
42 and long genomic distances [11-13], while others show remarkable specificity, regulating only
43 one of its neighbouring promoters or even skipping one or more promoters to activate more distal
44 ones. In part, 3D folding and compartmentalisation of the chromatin fibre help to establish this
45 specificity, by facilitating certain enhancer-promoter contacts and curbing others [12-14].

46 However, there is also substantial evidence that biochemical (in)compatibility between
47 CREs and promoters contributes to the specificity of their regulatory interactions. This is akin to
48 a lock-and-key mechanism: proteins bound to the CRE and the promoter must be compatible in
49 order to form a productive complex. Examples of such intrinsic selectivity have been documented
50 particularly in *Drosophila*, and in some instances could be attributed to a specific sequence motif
51 in the promoter [15-19]. Data obtained with massively parallel reporter assays (MPRAs) in
52 *Drosophila* cells have suggested a general separation of enhancer-promoter compatibility into
53 housekeeping and tissue-specific classes [20]. Some of this specificity may be determined by the
54 recruitment of co-factors [21]. However, a thorough understanding of the underlying mechanisms
55 is still lacking.

56 While several studies of individual enhancer-promoter combinations indicate that
57 biochemical compatibility also plays a role in mammals (e.g., [22-26]), systematic studies of this
58 mechanism have so far been lacking in mouse or human cells. Thus, it is still unknown how
59 widespread such intrinsic compatibility is in mammalian cells, and what drives this compatibility.

60 In order to address this issue, we systematically tested the compatibility of thousands of
61 combinations of candidate CREs (cCREs) and promoters using MPRAs. We used plasmid-based
62 MPRAs because they are highly scalable [27-29], and because episomal plasmids provide an
63 isolated context that minimises confounding effects of variable chromatin environments and
64 differences in 3D folding. However, so far MPRAs have mostly been used to assess the activity
65 of single elements, either as enhancers or as promoters [27, 30-34], except for one recent study
66 that tested combinations of synthetic elements [29]. To be able to dissect compatibility between
67 enhancers and promoters systematically, we designed cloning strategies that allowed us to test
68 thousands of pairwise cCRE–promoter combinations in different positions and orientations in a
69 reporter plasmid.

70 As models, we chose three genomic loci of 1-3 Mb in mouse embryonic stem cells
71 (mESCs). From these loci, which each encompass ~20 genes, we tested a large fraction of all
72 possible pairwise cCRE–promoter (cCRE-P) combinations. We found that more than half of the
73 active cCREs exhibit significant selectivity for specific subsets of promoters. We dissected some
74 of the underlying sequence determinants. Furthermore, we provide evidence suggesting that 3D
75 folding and intrinsic compatibility are independent mechanisms. Our experimental strategy and
76 datasets provide novel insights into the logic and mechanisms of cCRE-promoter specificity.

77

78

79 RESULTS

80

81 Experimental design

82 To maximise the probability of testing biologically relevant enhancer-promoter pairs, we
83 combined cCREs and promoters coming from the same region in the genome. We selected three
84 loci of 1-3 Mb in size, each roughly centred around a gene (*Nanog*, *Tfcp2l1* or *Klf2*) that is key to
85 the control of pluripotency of mESCs. The regulation of these genes is still incompletely
86 understood. In addition, each locus contains about 20 other genes ([Figure 1A-C](#)).

87 For promoters in the regions of interest we included approximately the -350 to +50 bp
88 segments around all GENCODE-annotated [35] transcription start sites (TSSs). The choice to
89 focus on the range -350 to +50 bp was motivated by our previous study of human promoters,
90 which indicated that most of the relevant information for promoter function is generally contained
91 within this range [30]. This definition of promoters is longer than that of core promoters (which
92 are usually only ~100 bp long) as was used in most previous enhancer reporter assays [21, 27,
93 29, 32-34, 36]. We considered this to be important, because the extra regulatory information
94 contained in those additional sequences may be relevant for interactions of the promoters with
95 CREs.

96 Compared to promoters, the annotation of cCREs is much less accurate. However, most
97 cCREs are centred around DNase I hypersensitive sites (DHS) [5, 37, 38]. We therefore selected
98 fragments of ~400 bp centred around all detected DHS peaks in each locus ([Figure 1A-C](#)). This
99 definition of cCREs within the range of typical enhancer definitions [39]. Some authors consider
100 enhancers combinations of multiple DHSs or longer stretches of DNA sequences. However, other
101 studies have shown that the activity of these long enhancers can be reproduced by shorter
102 versions of ~500 bp [40, 41]. Coordinates of all tested genomic fragments are provided in
103 [Supplementary Dataset 1](#).

104 We designed two MPRA variants to test many cCRE-P combinations ([Figure 1D-E](#)). In
105 the first variant, which we will refer to as Upstream assay, we obtained 82-192 individual cCREs
106 and 18-25 P elements per locus by PCR amplification ([Table 1](#)). We pooled all of these fragments
107 and randomly ligated them to form dimer fragments, which we then cloned *en masse* into a
108 reporter vector, *upstream* of a randomly barcoded transcription unit that lacked a promoter itself.
109 This resulted into highly complex libraries of cCRE-P, cCRE-cCRE, P-P and P-cCRE pairs, with
110 each individual element in two possible orientations. We then sequenced the libraries to identify
111 the paired fragments, their orientations in the reporter vector, and their linked barcodes. Owing
112 to the simple random ligation step, libraries with tens of thousands of cCRE-P combinations can
113 be obtained with this approach ([Table 1](#) and [Supplementary Table 1](#)). Here, we focus on the
114 analysis of cCRE-P pairs, but data from all other configurations are also provided as
115 [Supplementary Dataset 2](#).

116 In a second and complementary approach, we constructed a library in which the cCREs
117 are placed *downstream* of the reporter gene, i.e., separated ~1kb from the promoter ([Figure 1E](#)).
118 This was done in two steps: we first cloned a selection of 10 promoters upstream of the barcoded
119 transcription unit, resulting in a set of reporters with different promoters. Next, we inserted a pool
120 of cCREs into this set, downstream of the barcoded reporter unit and in both possible orientations.

121 We will refer to the assays done with the resulting library as Downstream assay. Due to the two-
122 step cloning protocol, the Downstream assay is less scalable than the Upstream assay, but
123 nevertheless allows for testing of hundreds of cCRE–P combinations (**Table 1**).

124 We used all P and cCRE DNA fragments from each of the three loci in separate Upstream
125 assays, whereas we focused on ten promoters and all cCREs from the *Klf2* locus in the
126 Downstream assay. **Table 1** provides summary statistics of the individual library compositions.
127 Due to the random nature of the combinatorial cloning, we did not recover all possible pairs.
128 Nevertheless, in the three Upstream assays combined we tested a total of 10,678 cCRE–P pairs,
129 or 3,747 pairs if we do not take orientations into account. For the Downstream assay these
130 numbers were 1,364 and 752, respectively. From the *Klf2* locus 847 and 676 pairs, respectively,
131 overlapped between the Upstream and Downstream assay. As references, we also inserted each
132 P and cCRE individually (i.e., unpaired) in the upstream position.

133

134 **Boost indices estimate promoter-specific activity of cCREs**

135 We then transiently transfected each of these libraries into mESCs. Twenty-four hours
136 after transfection we collected mRNA from the cells, and counted the transcribed barcodes by
137 reverse transcription followed by PCR amplification and high-throughput sequencing. In parallel,
138 barcodes were counted in the plasmid libraries. For each barcode we then normalised the counts
139 in cDNA over the counts detected in the plasmid DNA. Further data processing is described in
140 the Methods. We performed 3 biological replicates per library, which correlated with an average
141 Pearson $r=0.87$ (0.83 to 0.90) for the Upstream assay and $r=0.98$ (0.98 to 0.99) for the
142 Downstream assay. (**Figure S1 A-C**)

143 We first analysed the transcriptional activities of all singlet (unpaired) P and cCREs in the
144 upstream position. For promoters, these basal activities varied over a ~100-fold dynamic range
145 (**Figure 2A; Figure S2A**). Of all cCREs, 40.4% showed detectable transcriptional activity in the
146 upstream position without any P (**Figure 2A; Figure S2A**). Such autonomous transcriptional
147 activity is a frequently observed property of enhancers [30, 42, 43], and hence these elements
148 are likely to be enhancers. For a few cCREs this activity was as high as some of the strongest
149 promoters, suggesting that they may in fact be un-annotated promoters or very strong enhancers.

150 We then determined the ability of each cCRE to alter the activity of each linked P. For
151 this, we calculated a *boost index* for each cCRE–P pair, defined as the \log_2 -fold change in activity
152 of the cCRE–P pair compared to the P element alone. Unexpectedly, 20 negative controls that
153 we included in the *Klf2* libraries, consisting of randomly generated DNA sequences of similar size
154 and G/C content as the cCREs, showed a modestly negative boost index (median value -0.45
155 when inserted upstream) (**Figure S1D**). This is possibly because lengthening of the reporter
156 constructs alters the topology, supercoiling, transfection efficiency or a combination of these
157 parameters. We therefore corrected all cCRE–P boost indices for this non-specific negative bias
158 (see Methods). After this correction the negative controls had a marginal residual bias (median
159 \log_2 value -0.19), which we deemed acceptable (**Supplementary Figure S1D**).

160

161 **Identification of activating and repressive cCREs**

162 For each of the three genomic loci, the matrix of corrected boost indices shows a wide
163 diversity of patterns across the cCREs. We observed this both in the Upstream and Downstream
164 assays (**Figure 2B-D, Supplementary Figure S2B-D**). For example, in the *Klf2* locus Upstream
165 assay, cCRE E097 activates most of the tested promoters, while E046 (**Figure 2B**) and E057
166 (arrow in **Figure 2C**) only activate a distinct subset of promoters. Several elements are primarily
167 acting as repressors (e.g, E030 (**Figure 2B**) and E040, (arrow in **Figure 2C**)), and some seem
168 neither activating nor repressive (e.g., E070 (**Figure 2B**) and E085 (arrow in **Figure 2C**)).

169 We broadly classified the cCREs according to their overall effects on the linked promoters
170 (**Figure S3A**). In the Upstream assays, 21% of cCREs showed positive boost indices that were
171 significantly higher than the rest of cCREs across all tested promoters, indicating that they can
172 act as enhancer elements. About 17% of the cCREs showed negative boost indices significantly
173 below the rest of cCREs, and hence are putative repressor elements. For the remaining 62% of
174 cCREs the boost indices across their linked promoters were not significantly higher or lower than
175 the rest; these "ambiguous" elements either have no regulatory effects at all, or they have a mixed
176 repressive/activating/inactive effect that depends on the linked P (see below).

177 We were somewhat surprised to identify similar numbers of putative enhancers and
178 repressors, because most annotated cCREs in mammalian genomes are predicted to be
179 enhancers rather than repressive elements [5, 44]. In some cases this repression may be
180 underestimated in our analysis, as the estimates of negative boost indices for lowly active
181 promoters are less reliable due to the higher noise-to-mean ratios at low expression levels
182 (**Figure S3B**).

183 For activating elements, the boost indices varied in part according to the basal activities
184 of the cCRE and promoters. Strong boosting occurred primarily at promoters with low basal
185 activities, while highly active promoters were more difficult to boost (**FigS3C**). This suggests a
186 saturation effect, or it could indicate that promoters with high basal activity are less dependent
187 on distal enhancers. For cCREs, their basal activity is generally a strong positive predictor of their
188 enhancer potency (**Fig S3D**). However, exceptions to this rule occur, as some cCRE-P pairs
189 show high boost indices even though the basal activity of the cCRE is low (**Fig S3D**, upper left
190 quadrant).

191

192 **cCRE effects are predominantly orientation- and position-independent**

193 Next, we asked whether the ability of cCREs to regulate the linked promoters was generally
194 independent of their orientation and position. This was originally posited for enhancers [1], and
195 in some cases also reported for repressive elements [10]. Indeed, in the Upstream assays we
196 found a general positive correlation of the boost indices between the two orientations of the
197 cCREs (Pearson's $r=0.68$) (**Figure S4A**). These results are similar to those recently obtained
198 with a minimal core promoter [32]. In the Downstream assay the correlation between orientations
199 was somewhat lower (Pearson's $r=0.47$) (**Figure S4B**). This may be due to the lower dynamic
200 range of the Downstream assay data (**Figure S1C**). To simplify, for all other analyses we
201 combined the boost indices of + and - orientations of the cCREs by averaging.

202 We then investigated the degree of position-independence, by comparing the overlapping
203 P-cCRE pairs from the *Klf2* locus Downstream and Upstream assays. This showed an overall

204 Pearson correlation of 0.64 (**Figure S4C**). We conclude that repressive and activating effects of
205 cCREs are substantially but not completely position-independent, at least for the ten tested
206 promoters from the *Klf2* locus.

207 **Extensive selectivity of cCREs for promoters**

208 Visual inspection of the boost index matrices suggested that some cCREs alter the expression
209 of most promoters to similar degrees, while others selectively alter the expression of a subset of
210 promoters. In addition to the examples in **Figure 2B** from the *Klf2* locus, strikingly specific
211 promoter responses to some cCREs are illustrated for the *Tfcp2l1* locus in **Figure 3A**. For
212 example, E060, which forms part of an annotated super-enhancer [45], activates most of the
213 tested promoters, but with boost indices that can vary >50-fold between promoters. Two other
214 remarkable examples from the *Tfcp2l1* locus are E091 and E096, which each activate only a
215 single, distinct promoters out of the 11-12 promoters that were tested in each instance. Much
216 broader specificity is observed for E064, E073, E074 and E090 from the *Nanog* locus, which are
217 part of previously identified super-enhancers [46] (**Figure S2D**).

218 We investigated the degrees of selectivity more systematically. **Figure 4A-B** depicts the
219 distribution of the boost indices for each cCRE. Clearly, some cCREs have a much broader range
220 of boost indices than others. We used an ANOVA approach with Welch F-test to systematically
221 identify cCREs for which the variance of boost indices was larger than could be explained by
222 experimental noise (see methods). Strikingly, out of 233 cCREs with more than 5 tested cCRE-
223 P combinations, a total of 139 (59.9%) (**Figure 4B-C**) showed significant unexplained variance
224 at an estimated false-discovery rate (FDR) cutoff of 5%. Thus, at least throughout the three loci
225 that we tested, cCRE-P selectivity is widespread, ranging from strong specificity for one or a few
226 promoters to low specificity as seen in quantitative differences in the regulation of a broad set of
227 promoters.

228 Intersection of the ANOVA-based classification of selective/unselective cCREs with the
229 above broad classification into enhancers and repressors indicates that almost all (94%) general
230 enhancer elements exhibit significant P selectivity. In contrast, only 34% of the repressors are
231 detectably biased towards a subset of promoters (**Figure 4D**). However, we note that this
232 percentage may be underestimated, because at low expression levels the noise levels are higher
233 (**Figure S3B**). Interestingly, among the "ambiguous" cCREs, 55% are in fact selective. Such
234 elements mostly activate or repress only very few promoters (e.g., E091 and E096 from the
235 *Tfcp2l1* locus; **Figure 3**) and leave all other promoters unaffected. The remainder of the
236 ambiguous cCREs are probably not functional (e.g., E70 from the *Klf2* locus, **Figure 2B**). In
237 summary, these results indicate that more than half of all tested cCREs exhibits significant
238 preference for specific promoters.

239 Promoters of housekeeping and developmental genes in *Drosophila* were reported to
240 have distinct specificities toward cCREs [47]. To investigate whether such a dichotomy could also
241 be observed in our data, we focused on the *Klf2* locus, which has roughly equal numbers of
242 housekeeping and non-housekeeping promoters [48] (the *Tfcp2l1* and *Nanog* loci have only three
243 and zero housekeeping genes, respectively). Indeed, hierarchical clustering of the boost index
244 matrix showed a rough separation of the two classes of promoters (**Figure S5A**). However, this

245 is largely due to the highly similar cCRE specificities among the housekeeping promoters,
246 whereas the cCRE specificities of the non-housekeeping promoters are much more diverse and
247 generally as distinct from each other as from the housekeeping promoters (**Figure S5B**). To test
248 whether a housekeeping versus non-housekeeping dichotomy may largely explain our
249 identification of cCREs with significant selectivity (**Figure 4B-C**), we repeated this analysis after
250 removing all housekeeping promoters. This yielded highly similar results (123 of 221 cCREs are
251 significantly selective at 5% FDR cutoff, **Figure S5C**). We conclude that housekeeping promoters
252 may be similarly regulated, but cCRE selectivity goes beyond a simple distinction between
253 housekeeping and non-housekeeping promoters.

254

255 **Selectivity may be mediated by combinations of multiple TF motifs**

256 Taken together, these results point to a broad spectrum of cCRE specificities for promoters,
257 ranging from largely indiscriminate to highly selective. We searched for sequence motifs that may
258 account for these effects, focusing on binding motifs of transcription factors (TFs) that are
259 expressed in mESCs.

260 We first searched for TF motifs in the cCREs that correlate with boost indices across all
261 promoters. This yielded several dozens of TFs that are candidate activators or repressors (**Figure**
262 **S6A**). Several of these, such as Sox2, Nanog, ETV4 and GABPA are known key regulators in
263 mESC cells [49-51]. These TFs may broadly contribute to enhancer activity.

264 Next, we searched for motifs associated with cCRE-P selectivity. We reasoned that
265 selectivity may be due to certain combinations of TFs bound to cCRE and P. First, we asked
266 whether for any TF the simultaneous presence of its motif at cCRE and P correlated with boost
267 indices (**Figure 6SB**). This only yielded a weak association of FOXO motifs (at a 5% FDR cutoff).
268 Possibly this is due to FOXO1, a known regulator in mESCs [52]. We then asked if selectivity
269 may be mediated by multiple TFs rather than single TFs. For this purpose, we took the TF motifs
270 associated with enhancer activity with effect sizes >0.1 ($n=66$) and searched for combinations of
271 motifs that would be associated with higher boost indices if present at both the cCRE and the P
272 (**Figure 5A-B**). This yielded a few dozen stronger associations (at a 1% FDR cutoff). Some of
273 these associations may be redundant either because of motif similarity or because of motif co-
274 occurrence. For example, the 5 associations between Sox2 and Klf motifs may represent the
275 Klf4-Sox2 pair (**Figure 5B**) which are known to cooperate in mESCs [53]. These results indicate
276 that selectivity may be mediated by combinations of multiple TF motifs. Our dataset does not
277 provide sufficient statistical power for an exhaustive search of such combinations.

278 **Chromatin looping is independent of compatibility.**

279 Finally, we considered that certain pairs of cCREs and promoters frequently contact each other
280 in the nucleus, as is indicated by focal or stripe-like enrichment patterns in high-resolution Hi-C
281 maps [54, 55]. While long-range contacts are irrelevant in our MPRA because the tested
282 elements are directly linked, we asked whether such physical contacts in the native genomic
283 context are related to the selectivity of cCREs for certain promoters according to our MPRA. We
284 considered two models. In one model, the biochemical interactions that underlie cCRE-P
285 selectivity may promote or stabilise cCRE-P looping interactions. Alternatively, looping

286 interactions and cCRE-P selectivity may be independent aspects of cCRE-P interplay that each
287 work by different mechanisms.

288 To discriminate between these two models, we investigated whether the boost indices of
289 cCRE-P pairs correlate with their contact frequencies in Micro-C, a high-resolution variant of Hi-
290 C [55]. Remarkably, we found no correlation between these two quantities (**Figure 6A**). We also
291 found an extremely weak, although statistically significant, correlation between higher boost
292 indices and longer linear distances of cCRE-P pairs along the genome (**Figure 6B**).

293 We conclude that cCRE-P contacts in the nucleus may be independent of their functional
294 compatibility as detected in our reporter assays, raising the interesting possibility that chromatin
295 looping and compatibility are two orthogonal mechanisms of gene regulation.

296

297 **DISCUSSION**

298

299 Only a few other studies have so far attempted to analyse cCRE-P compatibility systematically.
300 An early survey of 27 cCRE-P combinations in human cells did not find evidence for specificity
301 [56], but the assay employed may have been insufficiently quantitative, and the choice of tested
302 elements may have been biased. In contrast, testing of ~200 cCRE-P pairs in zebrafish pointed
303 to extensive specificity [57]. An MPRA study in *Drosophila* cells using seven different promoters
304 and genome-wide cCREs suggested that cCRE-P specificity broadly separates between
305 housekeeping and tissue-specific promoters [47]. To our knowledge, our systematic
306 combinatorial testing of cCRE-P combinations in mESCs is the first large-scale study in
307 mammalian cells. The results reveal a broad spectrum of specificities: some cCREs are
308 promiscuous, others are highly specific for certain promoters, and in many instances the
309 specificity is quantitative rather than qualitative. By statistical analysis we found that more than
310 half of the cCREs exhibit a degree of specificity that cannot be explained by experimental noise.

311 It is likely that cCRE-P compatibility is governed by a complex grammar of TF
312 combinations. Underlying this grammar may be a diversity of molecular mechanisms, including
313 direct and indirect TF-TF interactions [e.g., 53], local concentration of activating factors [33, 58],
314 or functional bridging by cofactors [21, 59]. Due to the complexity of this grammar, its elucidation
315 may require much larger cCRE-P combinatorial datasets than generated here, as well as
316 systematic mutational analysis [60, 61] of individual cCRE-P combinations. Nevertheless, our
317 statistical analysis highlights several candidate combinations of TF motifs that may contribute to
318 the compatibility of some cCRE-P pairs.

319 Our data indicate that some of the cCREs tested may be repressive elements rather than
320 enhancers even though they were selected from DHSs. This is similar to a recent screen of
321 cCREs in human cells, which identified a large set of candidate repressive elements [62] and to
322 another screen in *Drosophila* [63]. It will be interesting to further explore the physiological
323 regulatory role of these elements. Particularly to understand their influence on close genes and
324 how repression works in open regions of the genome.

325 Surprisingly, we found that the boost indices of cCRE-P pairs generally do not correlate
326 with their contact frequencies in the native chromatin context. This suggests that 3D genome

327 organisation and compatibility are regulated by different mechanisms. We envision that
328 compatibility and 3D organisation may be two independent layers necessary for correct selective
329 gene regulation: 3D organisation such as the formation of chromatin loops and compartments
330 may determine whether CREs and promoters are able to interact, while compatibility may
331 determine whether such an interaction is functional, i.e., gives rise to a change in P activity.

332 Our current data were generated with transiently transfected plasmids. Advantages of this
333 approach are that it largely eliminates possible confounding effects of chromatin packaging and
334 3D folding, and that thousands of cCRE–P combinations could be tested. Even higher throughput
335 combinatorial MPRA's will be useful in order to fully dissect the rules behind compatibility either
336 by testing more cCRE-P combinations or mutagenised cCRE-P pairs. However, further studies
337 are needed to verify and analyse the impact of the observed specificities in the native genomic
338 context. Due to genomic confounding factors, such as chromatin context, 3D organisation,
339 regulatory element redundancy/synergy, and poor scalability, such studies will be challenging
340 and may require the development of new technologies.

341

342 MATERIALS AND METHODS

343

344 ***Selection of cCREs and promoters***

345 For the design of the libraries we selected the cCREs and promoters from three TADs centered
346 around each of the *Klf2*, *Nanog* and *Tfcp2l1* genes, using TAD coordinates from [54]. cCREs
347 were selected based on DNase hypersensitivity mapping data from mESCs in both 2i+LIF [38]
348 and serum [5] culturing conditions, which we reprocessed and aligned to the mm10 genome
349 build. DNase hypersensitivity sites (DHSs) were called using Homer v4.10 with default
350 parameters and peak style “factor”. We defined cCREs as 450 bp windows centered on each
351 peak. For promoters we used the Gencode mouse TSS annotation [35]. From each TSS we
352 defined as promoters the -375 +75 bp region. If the promoter regions overlapped with any
353 cCRE then the promoter was redefined as the 450 bp region surrounding the center of the
354 intersection of both elements. PCR primers were designed for each cCRE and promoter using
355 the batch version of Primer3 (BatchPrimer3 v1.0) [64] allowing for primers to be designed on
356 the 50 bps of each end. This yielded PCR products of ~400 bp for each element.

357

358 ***Upstream assay library generation***

359 For each locus, cCREs and promoters were amplified from mouse genomic DNA (extracted
360 from E14TG2a mESCs, ATCC CRL-1821) by PCR using My-Taq Red mix (#BIO-25044; Biorun)
361 in 384 well plates using automated liquid handling (Hamilton Microlab® STAR). PCRs were
362 checked on gel and had a success rate between 60 and 90% depending on the locus. Equal
363 volumes (10ul) of the resulting PCR products were mixed, and the resulting pool was purified
364 by phenol-chloroform extraction followed by gel purification (BIO-52059; Biorun). The purified
365 DNA fragments were then blunted and phosphorylated using End-It DNA End-Repair Kit
366 (#ER0720; Epicentre). Part of the repaired pool was set apart for cloning of singlet libraries.
367 The remainder was self-ligated using Fast-link ligase (LK0750H; Lucigen), after which duplets
368 of ~800bp were excised from agarose gel and purified (BIO-52059; Biorun). Singlet and duplet
369 pools were A-tailed using using Klenow HC 3'→5' exo- (#M0212L; NEB).

370 The SuRE barcoded vector was prepared as described [30]. Then singlet and duplet
371 pools were separately ligated overnight into the SuRE barcoded vector using Takara ligation kit
372 version 2.1 (#6022; Takara). Ligation products were purified using magnetic bead purification
373 (#CPCR-0050; CleanNA). Next, 2 µl of the purified ligation products were electroporated into 20
374 µl of electrocompetent e. coli 10G supreme (#60081-1; Lucigen). Each library was grown
375 overnight in 500 ml of standard Luria Broth (LB) with 50 µg/ml of kanamycin and purified using a
376 maxiprep kit (K210016, Invitrogen).

377

378 ***Downstream assay library generation***

379 The Downstream assay vector was based on a pSMART backbone (Addgene plasmid # 49157;
380 a gift from James Thomson). It was constructed using standard molecular biology techniques
381 and contains a green fluorescent protein (GFP) open reading frame followed by a barcode, and
382 a psiCheck polyadenylation signal (PAS) introduced during barcoding, followed by the cloning
383 site for inserts and a triple polyadenylation site (SV40+bGH+psiCheckPAS).

384 The 10 highest expressing promoters of the *Klf2* Upstream library were selected to be
385 cloned into the Downstream assay vector at the promoter position. These Promoters were
386 amplified by PCR and individually inserted by Gibson assembly (#E2611S; NEB) into the
387 Downstream assay vector. Then each of the 10 constructs were transformed into standard
388 DH5 α competent bacteria (#C2987; NEB) grown overnight in in 500 ml of standard Luria
389 Broth(LB) with 50 μ l/ml of kanamycin and purified.

390 Each of these promoter-containing vectors was then barcoded similarly as the SuRE
391 vector [30]. For this, we digested 10 μ g of each vector with AvrII (#ER1561; Thermo Fischer)
392 and XcmI (#R0533; NEB) and performed a gel-purification. Barcodes were generated by
393 performing 10 PCR reactions of 100 μ l each containing 5 μ l of 10 μ M primer 275JvA, 5 μ l of 10
394 μ M primer 465JvA and 1 μ l of 0.1 μ M template 274JvA (see [Supplementary Table 2](#) for
395 oligonucleotide sequences). A total of 14 PCR cycles were performed using MyTaq Red Mix
396 (#BIO-25043; Bioline), yielding ~30 μ g barcodes. Barcodes were purified by phenol-chloroform
397 extraction and isopropanol precipitation after which they were digested overnight with 80 units
398 of NheI (#R0131S; NEB) and purified using magnetic bead purification (#CPCR-0050;
399 CleanNA). Each vector variant and the barcodes were then ligated in one 100 μ l reaction
400 containing 3 μ g digested vector and 2.7 μ g digested barcodes, 20 units NheI (#R0131S; NEB),
401 20 units AvrII, 10 μ l of 10 \times CutSmart buffer, 10 μ l of 10 mM ATP, 10 units T4 DNA ligase
402 (#10799009001 Roche). A cycle-ligation of six cycles was performed (10 min at 22 $^{\circ}$ C and 10
403 min at 37 $^{\circ}$ C), followed by 20 min heat-inactivation at 80 $^{\circ}$ C. The ligation reaction was purified
404 by magnetic beads and digested with 40 units of XcmI (#R0533S; NEB) for 3 h, and size-
405 selected by gel-purification, yielding ~1 μ g barcoded vector for each variant.

406

407 ***Inverse PCR and sequencing to link inserted elements to barcodes***

408 We identified barcode–insert combinations in the plasmid libraries by inverse-PCR followed by
409 sequencing as described [30]. In brief, the combination of barcode and element(s) was excised
410 from the plasmid by digestion with I-ceul; this fragment was circularised; remaining linear
411 fragments were destroyed; and circular fragments were linearised again with I-sceI. These
412 linear fragments were amplified by PCR with sequencing adaptors. The final product was
413 sequenced on an Illumina MiSeq platform using 150 bp paired-end reads. This process was
414 done separately for each of the libraries. In the singlet libraries the barcodes should be
415 associated to only one insert and in the combinatorial libraries the barcodes should be
416 associated with duplets.

417

418 ***Linking barcodes to element singlets or duplets***

419 For each library the iPCR data was locally aligned using bowtie (version 2.3.4) [65] with very
420 sensitive parameters (--very-sensitive-local) on a custom bowtie genome. This custom genome
421 was generated using bowtie. It consists of virtual chromosomes corresponding to each cCRE or
422 a P from each locus. Bam alignment files were processed using a custom python script that
423 identifies from read 1 the barcode and cCRE or P element, and from read 2 the cCRE or P
424 element. In case of singlet libraries both reads should identify the same element, whereas in

425 combinatorial libraries read 1 is derived from the barcode-proximal element and read 2 from the
426 barcode distal element. In the combinatorial libraries we can not distinguish between a
427 combination of one element with itself in the same orientation or a single element, therefore
428 these were removed from combinatorial libraries. In the Downstream Assay both reads identify
429 the only element cloned in the downstream position. If no element was found, the barcode was
430 assigned as empty vector. The resulting barcode-to-element(s) lists were clustered using
431 Starcode (version 1.1) [66] to remove errors from barcode sequencing. Finally, barcodes
432 present in multiple libraries or matched with multiple element combinations were removed from
433 the data.

434

435 **Cell culture and transfection**

436 All experiments were conducted in E14TG2a mouse embryonic stem cells (mESC) (ATCC
437 CRL-1821) cultured in 2i+LIF culturing media. 2i+LIF was made according to the 4DN nucleome
438 protocol for culturing mESCs (<https://data.4dnucleome.org/protocols/cb03c0c6-4ba6-4bbe-9210-c430ee4fdb2c/>). The reagents used were Neurobasal medium (#21103-049, Gibco),
439 DMEM-F12 medium (#11320-033, Gibco), BSA (#15260-037; Gibco), N27 (#17504-044;
440 Gibco), B2 (#17502-048; Gibco), LIF (#ESG1107; Sigma-Aldrich), CHIR-99021 (#HY-10182;
441 MedChemExpress) and PD0325901 (#HY-10254; MedChemExpress), monothioglycerol
442 (#M6145-25ML; Sigma) and glutamine (#25030-081, Gibco). Monthly tests (#LT07-318; Lonza)
443 confirmed that the cells were not contaminated by mycoplasma. Cells were transiently
444 transfected using Amaxa nucleofector II, program A-30, and Mouse Embryonic Stem Cell
445 Nucleofector™ Kit (#VPH-1001, Lonza). *Klf2* and *Nanog* loci Upstream assay libraries were
446 mixed and transfected together, *Tfcp2l1* Upstream Assay libraries were transfected in separate
447 experiments. All the Downstream assay sub-libraries were transfected as a mix. Three
448 independent biological replicates were done for each library mix. For each biological replicate
449 16 million cells were transfected (4 million cells with 4 µg plasmid per cuvette)

451

452 **RNA extraction and cDNA sequencing**

453 RNA was extracted and processed for sequencing as described [30] with a few modifications.
454 Cells were harvested 24 h after transfection, resuspended in Trisure (#BIO-38032; Bioline) and
455 frozen at -80 °C until further processing. From the Trisure suspension, the aqueous phase
456 containing the RNA was extracted and loaded into RNA extraction columns (#K0732, Thermo
457 Scientific). Total RNA was divided into 10 µl reactions containing 5 µg of RNA and was treated
458 for 30 mins with 10 units of DNase I (#04716728001; Roche). Then DNase I was inactivated by
459 addition of 1 µl of 25 mM EDTA and incubation at 70°C for 10 min.

460 For the Upstream Assay the cDNA was produced and amplified by PCR as described
461 [30]. Per biological replicate 8 to 10 reactions were carried out in parallel in order to cover
462 enough barcode complexity of the library. For the Downstream Assay the RNA was extracted
463 and processed the same way until cDNA production Here, cDNA was produced using a specific
464 primer (304JvA sequence in [Supplementary Table 2](#) for oligonucleotide sequences). Primer
465 304JvA introduces an adaptor sequence 5' to the primer sequence which is targeted in the first

466 PCR (see below) to ensure strand specific amplification of barcodes. Then cDNA was amplified
467 in 2 steps (nested PCRs) in order to make the reaction strand-specific. The first PCR reaction
468 was run for 10 cycles (1 min 96 °C, 10 times (15 s 96 °C, 15 s 60 °C, 15 s 72 °C)) using (index
469 variants of) primers 285JvA (containing the S2, index and p7 adaptor) and 305JvA (targeting
470 the adapter introduced by 304JvA). Each 20 µl RT reaction was amplified in a 100-µl PCR
471 reaction with MyTaq Red mix. The second PCR reaction was performed using 10ul of the
472 product of the previous reaction in 100 µl reactions (1 min 96 °C, 8×(15 s 96 °C, 15 s 60 °C, 15
473 s 72 °C)) using the same index variant primer and primer 437JvA (containing the S1, and p5
474 adaptor). For both Upstream and Downstream assays, the resulting PCR products were
475 sequenced on an Illumina 2500 HiSeq platform with 65bp single end reads.

476

477 ***Plasmid DNA (pDNA) barcode sequencing***

478 For normalisation purposes, barcodes in the plasmid pools were counted as follows. For both
479 assays the process was the same. For each library 1 µg of plasmid was digested with I-sceI in
480 order to linearise the plasmid. Then, barcodes were amplified by PCR from 50 ng of material
481 using the same primers and reaction conditions as in the amplification of cDNA in the Upstream
482 assay, but only 9 cycles of amplification were used (1 min 96 °C, 9 times (15 s 96 °C, 15 s 60
483 °C, 15 s 72 °C)). For each library, two technical replicates were carried out by using different
484 index primers for each replicate. Samples were sequenced on an Illumina 2500 HiSeq platform
485 with 65bp single end reads.

486

487 ***Pre-Processing of cDNA and pDNA reads***

488 For each replicate of each library pool transfection barcodes were extracted from the single end
489 reads by using a custom python script that identifies the constant region after the barcode.
490 Near-identical barcodes were pooled using Starcode (version 1.1) [66] to remove errors from
491 barcode sequencing, and barcode counts were summarised. The process was the same for
492 cDNA and pDNA counts and for Upstream and Downstream data.

493

494 ***Post processing of cDNA and pDNA counts***

495 For each transfection, barcodes identified in the cDNA were matched to the barcodes in the
496 iPCR data, and all barcodes were counted in cDNA and pDNA replicates. Barcode counts were
497 normalised to the total number of barcode reads from each sample. Activity per barcode was
498 then calculated as a cDNA:pDNA ratio of normalised counts. Next, activities from multiple
499 barcodes belonging to the same element singlet or combination were averaged, requiring a
500 minimum of 5 barcodes per singlet or combination and at least 8 pDNA counts per barcode.
501 The mean activity of each singlet or combination across replicates was calculated as the
502 geometric mean of the three replicates.

503

504 ***Calculation of boost indices***

505 We initially calculated raw boost indices simply as a log₂ ratio of the activity of each cCRE–P
506 pair over the activity of the corresponding P alone. However, 20 negative controls that we
507 included in the *Klf2* libraries, consisting of randomly generated DNA sequences of similar size

508 and G/C content as the cCREs (**Supplementary dataset 1**), generally showed a negative
509 boost index by this measure (median value -0.45 when inserted upstream) (**Figure S1D**). We
510 therefore calculated corrected boost indices as the \log_2 ratio of cCRE-P activity over the median
511 cCRE-P activity per promoter (**Figure S1D**). Importantly, in the *Klf2* library data this largely
512 removed the negative bias that we observed with the negative controls; we thus assume that
513 this correction is adequate and therefore also applied it to the boost indices obtained with the
514 other libraries. For the analyses in **Figures 2-6** and **Supplementary figures 2-6** except **4A-B**
515 the boost indices of cCREs were averaged over both orientations of the cCREs.

516

517 ***Analysis of selectivity***

518 We performed a Welch's ANOVA (or Welch F-test) to assess the selectivity of each cCRE with
519 more than 5 cCRE-P combinations. For this purpose, each replicate of each orientation of the
520 cCRE-P was used as a datapoint and each cCRE-P combination was used as a group. P-
521 values were corrected for multiple hypothesis testing using the Benjamini-Hochberg method
522 and an FDR cutoff of 5% was chosen. The Welch F-test was chosen over the classic ANOVA
523 due to heteroscedasticity of the data.

524

525 ***TF motif Survey***

526 We used a custom TF motif database provided by the lab of Gioacchino Natoli containing 2,448
527 TF motifs which was built on top of a previously published version [67] (Dataset composition
528 and sources available at [##GitHub-url](#)). TF motifs were filtered for expression of TFs in mESCs
529 cultured in 2i+LIF according to published RNA-seq (higher expression than 1RPM) [38]. We
530 scored presence or absence of a TF motif in each cCRE using FIMO (MEME suite, version
531 5.0.2). We then searched for motifs associated with (1) general enhancer activity, (2) self-
532 compatibility and (3) duplets of self-compatible motifs. In (1), for each TF motif we compared
533 the general cCRE-P population to combinations where the TF motif was present at the cCRE.
534 In (2), for each TF motif we compared the cCRE-P combinations where the TF motif was
535 present at the cCRE to the combinations where it was present at both the cCRE and the
536 promoter. In (3), we took all the significant TF motifs at a 1% FDR and an effect size higher
537 than 0.1 ($n=66$). Then we tested all pairwise non-repeated TF motif duplets. Per TF motif duplet
538 we compared the cCRE- promoters where both TF motif were present at the cCRE to the
539 combinations where both were present at both the cCRE and the promoter. In all comparisons
540 a Wilcoxon test was applied to the boost indices of each group and the effect size was
541 calculated a difference of median boost indices. In each analysis p-values were corrected for
542 multiple hypothesis testing using the Benjamini-Hochberg method. We required a minimum of
543 50 combinations per group.

544

545 ***Micro-C data correlation***

546 Micro-C data was obtained from [55]. Contact scores between cCRE-P pairs were averaged
547 across bins overlapping a ± 500 bp window from the location of each element using 400 bp
548 bins.

549

550 ***Data analysis and data availability.***

551 All data analysis was performed in R [68]. Code of data processing pipelines and analysis
552 scripts are available at [##Github-url](#). Raw and processed data are available at GEO (accession
553 nr GSE186265). Processed datasets and pipeline output files are available at OSF ([##OSF-ur](#)).

554

555

556

557

558 **Author contributions**

559 M.M.A, F.C. and B.v.S. designed the study. M.M.A and F.C. developed computational methods
560 and performed analyses. M.M.A. and J.v.A. developed experimental methods. M.M.A.
561 performed experiments. B.v.S. and M.M.A. wrote the manuscript, with input from F.C. and
562 J.v.A. B.v.S. supervised the study.

563

564 **Acknowledgements**

565 We thank Tao Chen for initial input on the study design, the NKI Genomics, Robotics and
566 Research High Performing Computing facilities for technical support, Barak Cohen and his lab
567 for insightful discussions, and Gioacchino Natoli and his lab for providing the TF motif
568 database. Supported by ERC Advanced Grant 694466 (B.v.S.) and Swiss National Science
569 Foundation postdoctoral fellowship P2EZP3_165206 (F.C). Oncode is partly funded by the
570 Dutch Cancer Society KWF.

571

572 **Competing Interests**

573 J.v.A. is founder of Gen-X B.V. and Annogen B.V. F.C. is a co-founder of enGene Statistics
574 GmbH.

575 FIGURE LEGENDS

576

577 **Figure 1.** Regulatory element selection and library construction. **A-C)** Representations of
578 *Nanog*, *Tfcp2l1*, and *Klf2* loci, respectively. In **C)** the zoom-in displays a DNase I sensitivity
579 track [38] where peaks overlap with cCREs. **D)** Cloning strategy for the Upstream assay.
580 cCREs and promoters were amplified by PCR from genomic DNA and pooled. Fragments in
581 this pool were then randomly ligated to generate duplets. Singlets and duplets were cloned into
582 the same barcoded vector to generate two libraries per locus, a singlet library and a
583 combinatorial library. **E)** Cloning strategy for the Downstream assay. The singlet pool from the
584 *Klf2* locus was cloned into ten vectors, each of them carrying a different promoter. The resulting
585 ten sub-libraries were combined into one Downstream assay library.

586

587 **Figure 2.** Singlet and combinatorial activities of cCREs and promoters from the *Klf2* locus. **A)**
588 Transcription activities of singlet cCREs and promoters. Each dot represents the mean activity
589 of one singlet. Horizontal lines represent the average background activity of empty vectors
590 (black line) plus or minus two standard deviations (grey lines). Elements with activities more
591 than two standard deviations above the average background signal are defined as active. **B)**
592 Examples of Upstream assay cCRE-P combinations for cCREs E097, E046, E030 and E070 of
593 the *Klf2* locus. Barplots represent the mean boost index of each combination, vertical lines
594 represent the standard deviations. Crosses mark missing data. **C-D)** Boost index matrices of
595 cCRE-P combinations from the *Klf2* locus according to Upstream (**C)** and Downstream (**D)**
596 assays. White tiles indicate missing data. Barplots on the right and top of each panel show
597 basal activities of each tested P or cCRE, respectively, with the black line indicating the
598 background activity of the empty vector. All data are averages over 3 independent biological
599 replicates.

600

601 **Figure 3.** Examples of selective cCREs from the *Tfcp2l1* locus. Boost indices obtained in the
602 Upstream assay are shown for cCRE-P combinations of cCREs E060, E091 and E096 of the
603 *Tfcp2l1* locus. Barplots indicate the mean boost index of each combination, vertical lines
604 indicate standard deviations. All data are averages over 3 independent biological replicates.

605

606 **Figure 4.** Promoter selectivity of cCREs. **A)** Plot showing the broad diversity of boost indices of
607 many cCREs. Data are from Upstream assays of *Klf2*, *Nanog* and *Tfcp2l1* loci combined.
608 Vertical axis indicates boost indices of all tested cCRE-P pairs, which are horizontally ordered
609 by the mean boost index of each cCRE. **B)** Boost index distributions for each cCRE from the
610 *Klf2* locus (Upstream assay). Each dot represents one cCRE-P combination; black bar
611 represents the mean. Turquoise colouring marks cCREs that have a larger variance of their
612 boost indices than may be expected based on experimental noise, according to the Welch F-
613 test after multiple hypothesis correction (5% FDR cutoff). **C)** Summary of Welch F-test
614 selectivity analysis results for all cCREs from the three loci with more than 5 cCRE-P
615 combinations. Each dot represents one cCRE; the size of the dots indicates the number of
616 cCRE-P pairs. Significantly selective cCREs (5% FDR cutoff) are highlighted in turquoise. **D)**

617 Proportion of significantly selective (turquoise) cCRE in the three categories as shown in
618 **Figure S3A**. All data are averages over 3 independent biological replicates.

619
620 **Figure 5**. Association of TF motif Duos with higher boost indices. **A)** Results of TF survey for
621 self-compatible TF motif Duos. TF motif duos associated with higher or lower boost indices at a
622 1% FDR cutoff are highlighted. **B)** Association of Sox2+Klf4 motifs at both cCRE and P with
623 higher boost indices. cCRE-P combinations are split into 3 groups according to presence or
624 absence of Sox2+Klf4 motifs both at the cCRE and the promoter, or only the cCRE. Numbers at
625 the top of horizontal brackets are the p-values obtained from comparing the different groups
626 boost index distributions using a Wilcoxon rank-sum test. Boxplots represent median and
627 interquartile ranges. Barplots at the top represent the number of combinations in each group.

628
629 **Figure 6**. Absent or very weak correlation between boost indices and **(A)** contact frequencies
630 according to micro-C [55] or **(B)** linear genomic distance, for all cCRE-P pairs from the three
631 loci combined. All boost index data are averages over 3 independent biological replicates.

632
633
634

635 SUPPLEMENTARY FIGURE LEGENDS

636
637 **Figure S1**. Reproducibility of data and boost index calculation. **(A-C)** Correlograms of the three
638 biological replicates of each library pool. Lower left panels show pairwise scatterplots of the
639 activities of all cCRE-P pairs per replicate. Middle panels show the density of data distribution in
640 each replicate and upper right panels show the Pearson correlation coefficients. **A)** *Klf2* and
641 *Nanog* Upstream libraries. **B)** *Tfcp2l1* Upstream library. **C)** *Klf2* Downstream libraries. **D)**
642 Upstream assay boost index distributions for cCRE-P and negative controls – promoter (NC-P)
643 combinations. Left panel: raw boost indices; right panel: boost indices after correction for
644 negative bias (see Methods).

645
646 **Figure S2**. Element activities and boost indices obtained with *Nanog* and *Tfcp2l1* Upstream
647 libraries. **A)** Transcriptional activities of cCREs and promoters. Each dot represents the mean
648 activity of one singlet. Horizontal lines represent the average background activity of empty
649 vectors (black line) plus or minus two standard deviations (grey lines). Elements with activities
650 more than two standard deviations above the average background signal are defined as active.
651 **B-C)** Boost index matrices for cCRE–P pairs from *Nanog* and *Tfcp2l1* loci (both Upstream
652 assays). White tiles indicate missing data. Barplots on the right and top of each panel show
653 basal activities of each tested P or cCRE, respectively, with the black line indicating the
654 background activity of the empty vector. **D)** Examples of cCRE-P combinations for cCREs
655 E064, E073, E074 and E090 of the *Nanog* locus. Barplots represent the mean boost index of
656 each combination, vertical lines represent the standard deviation of each boost index. All data
657 are averages over 3 independent biological replicates.

658

659 **Figure S3.** cCRE functional classification and activity influence on Boost indices. **A)** Volcano
660 plot of cCREs associated with activation or repression across promoters. A Wilcoxon test is
661 performed per cCRE comparing the boost indices of all the cCRE-P combinations of that cCRE
662 against the rest of cCRE-P combinations. A minimum of 6 combinations is required per cCRE.
663 P-values are corrected for multiple hypothesis testing using the Benjamini-Hochberg method
664 (FDR). **B)** Relationship between noise-to-mean ratio (Standard Deviation/mean Activity) and
665 mean activity of cCRE-Ps. Horizontal lines represent noise-to-mean ratios of 1 and of 4 in log₂
666 scale. **C)** Relationship between boost indices and basal (singlet) P activity. Each column of dots
667 shows the data of cCRE–P pairs for one P. Data are from Upstream assays of all three loci
668 combined. **D)** Relationship between boost indices and basal (singlet) cCRE activity. All data are
669 averages over 3 independent biological replicates.

670
671
672 **Figure S4.** Orientation and position independence of cCREs. **(A-B)** Correlation between boost
673 indices of both cCRE orientations of the same cCRE-P combination, in the **(A)** Upstream assay
674 and **(B)** Downstream assay. Data are from the *Klf2* locus libraries. Note that "+" and "-"
675 orientations are arbitrary labels, because cCREs do not have an intrinsic orientation. **(C)**
676 Correlation between boost indices of cCRE-P combinations shared between the Upstream and
677 Downstream assays of the *Klf2* locus. In all panels R is the Pearson correlation coefficient. All
678 data are averages over 3 independent biological replicates. In C Boost indices are averaged
679 over cCRE orientations.

680
681
682 **Figure S5.** Housekeeping promoters show a distinct pattern of cCRE compatibility. **A)**
683 Hierarchical clustering of the Upstream assay boosting matrix of the *Klf2* locus. In order to
684 facilitate hierarchical clustering the matrix has been restricted to almost complete cases
685 (cCREs >15 combinations) **B)** Density plot of pairwise Pearson correlation coefficients of the
686 boost indices of *Klf2* locus promoters classified as either housekeeping or non-housekeeping
687 [48]. Blue: correlations between all pairs of housekeeping promoters; red: all correlations
688 between pairs of non-housekeeping promoters; grey: all correlations between one
689 housekeeping and one non-housekeeping promoter. Vertical lines represent the median of
690 each group. Unlike in (A), all promoters in the Upstream assay were included in this analysis.
691 **C)** Results of selectivity analysis as performed in **Figure 4C**, but excluding housekeeping
692 promoters. All data are averages over 3 independent biological replicates.

693
694
695 **Figure S6.** Identification of single TF motifs that correlate with boost indices. **(A)** TF motifs in
696 cCREs associated (at 1% FDR cutoff) with activation (turquoise) or repression (red). **(B)** Motifs
697 of putative self-compatible TFs, i.e. motifs that predict increased or reduced boosting indices
698 when present both at the cCRE and P, compared to being present only at the cCRE. TF motifs
699 associated with higher or lower boost indices at a 1% FDR cutoff are highlighted. We note that

700 TF motifs with multiple hits from the same family, such as for ELK, FOXO and ELF factors, may
 701 in fact be due to the activity of one TF motif of that family [69].

702

703

704 **TABLES**

705

706 **Table 1. Numbers of tested Promoters (Ps), cCREs and cCRE–P pairs in each combinatorial**
 707 **MPRA library.**

Library	Ps present	cCREs present	cCRE–P pairs tested	cCRE–P pairs (orientation-independent)
Klf2 Upstream	23	82	3758	1400
Nanog Upstream	18	88	1321	595
Tfcp2l1 Upstream	25	198	5599	2490
Klf2 Downstream	10	84	1364	752

708

709

710 **SUPPLEMENTARY TABLES**

711

712 **Supplementary table 1. Other combinations of cCRE and P elements in each MPRA library.**

Library	cCRE-cCRE	cCRE-cCRE (orientation-independent)	P-P	P-P (orientation-independent)	P-cCRE	P-cCRE (orientation-independent)
Klf2 Upstream	10626	4284	1335	441	4067	1439
Nanog Upstream	10536	4769	155	82	1511	713
Tfcp2l1 Upstream	44515	21149	626	274	5239	2386
Klf2 Downstream	0	0	420	225	0	0

713

714 **Supplementary table 2. Oligonucleotide and plasmid sequences**
 715 **(supplementary file)**

716

717 **SUPPLEMENTARY DATASETS**

718

719 **Data Set 1 Coordinates and sequences of cCREs and Promoters**

720 **Data Set 2 Activities of all cCRE-cCRE, cCRE-P, P-cCRE and P-P combinations Upstream**
 721 **assay**

722 **Data Set 3 Boost indices of cCRE-P combinations Upstream assay**

723 **Data Set 4 Boost indices of cCRE-P combinations Downstream assay**

724

725

726 **REFERENCES**

- 727
- 728 1. Banerji, J., S. Rusconi, and W. Schaffner, Expression of a beta-globin gene is enhanced
729 by remote SV40 DNA sequences. *Cell*, 1981. 27(2 Pt 1): p. 299-308.
 - 730 2. Tuan, D., et al., The "beta-like-globin" gene domain in human erythroid cells. *Proc Natl*
731 *Acad Sci U S A*, 1985. 82(19): p. 6384-8.
 - 732 3. Fiering, S., et al., Targeted deletion of 5'HS2 of the murine beta-globin LCR reveals that
733 it is not essential for proper regulation of the beta-globin locus. *Genes Dev*, 1995. 9(18):
734 p. 2203-13.
 - 735 4. Lettice, L.A., et al., A long-range Shh enhancer regulates expression in the developing
736 limb and fin and is associated with preaxial polydactyly. *Hum Mol Genet*, 2003. 12(14):
737 p. 1725-35.
 - 738 5. Consortium, E.P., An integrated encyclopedia of DNA elements in the human genome.
739 *Nature*, 2012. 489(7414): p. 57-74.
 - 740 6. van Arensbergen, J., B. van Steensel, and H.J. Bussemaker, In search of the
741 determinants of enhancer-promoter interaction specificity. *Trends Cell Biol*, 2014.
742 24(11): p. 695-702.
 - 743 7. Zabidi, M.A. and A. Stark, Regulatory Enhancer-Core-Promoter Communication via
744 Transcription Factors and Cofactors. *Trends Genet*, 2016. 32(12): p. 801-814.
 - 745 8. Farley, E.K., K.M. Olson, and M.S. Levine, Regulatory Principles Governing Tissue
746 Specificity of Developmental Enhancers. *Cold Spring Harb Symp Quant Biol*, 2015. 80:
747 p. 27-32.
 - 748 9. Robson, M.I., A.R. Ringel, and S. Mundlos, Regulatory Landscaping: How Enhancer-
749 Promoter Communication Is Sculpted in 3D. *Mol Cell*, 2019. 74(6): p. 1110-1122.
 - 750 10. Segert, J.A., S.S. Gisselbrecht, and M.L. Bulyk, Transcriptional Silencers: Driving Gene
751 Expression with the Brakes On. *Trends Genet*, 2021. 37(6): p. 514-527.
 - 752 11. Shlyueva, D., G. Stampfel, and A. Stark, Transcriptional enhancers: from properties to
753 genome-wide predictions. *Nat Rev Genet*, 2014. 15(4): p. 272-86.
 - 754 12. Schoenfelder, S. and P. Fraser, Long-range enhancer-promoter contacts in gene
755 expression control. *Nat Rev Genet*, 2019. 20(8): p. 437-455.
 - 756 13. Furlong, E.E.M. and M. Levine, Developmental enhancers and chromosome topology.
757 *Science*, 2018. 361(6409): p. 1341-1345.
 - 758 14. Lupianez, D.G., et al., Disruptions of topological chromatin domains cause pathogenic
759 rewiring of gene-enhancer interactions. *Cell*, 2015. 161(5): p. 1012-1025.
 - 760 15. Li, X. and M. Noll, Compatibility between enhancers and promoters determines the
761 transcriptional specificity of gooseberry and gooseberry neuro in the *Drosophila* embryo.
762 *EMBO J*, 1994. 13(2): p. 400-6.
 - 763 16. Merli, C., et al., Promoter specificity mediates the independent regulation of neighboring
764 genes. *Genes Dev*, 1996. 10(10): p. 1260-70.
 - 765 17. Butler, J.E. and J.T. Kadonaga, Enhancer-promoter specificity mediated by DPE or
766 TATA core promoter motifs. *Genes Dev*, 2001. 15(19): p. 2515-9.
 - 767 18. Juven-Gershon, T., J.Y. Hsu, and J.T. Kadonaga, Caudal, a key developmental
768 regulator, is a DPE-specific transcriptional factor. *Genes Dev*, 2008. 22(20): p. 2823-30.
 - 769 19. Kwon, D., et al., Enhancer-promoter communication at the *Drosophila* engrailed locus.
770 *Development*, 2009. 136(18): p. 3067-75.
 - 771 20. Arnold, C.D., et al., Genome-wide assessment of sequence-intrinsic enhancer
772 responsiveness at single-base-pair resolution. *Nat Biotechnol*, 2017. 35(2): p. 136-144.

- 773 21. Haberle, V., et al., Transcriptional cofactors display specificity for distinct types of core
774 promoters. *Nature*, 2019. 570(7759): p. 122-126.
- 775 22. Bertolino, E. and H. Singh, POU/TBP cooperativity: a mechanism for enhancer action
776 from a distance. *Mol Cell*, 2002. 10(2): p. 397-407.
- 777 23. Vakoc, C.R., et al., Proximity among distant regulatory elements at the beta-globin locus
778 requires GATA-1 and FOG-1. *Mol Cell*, 2005. 17(3): p. 453-62.
- 779 24. Jing, H., et al., Exchange of GATA factors mediates transitions in looped chromatin
780 organization at a developmentally regulated gene locus. *Mol Cell*, 2008. 29(2): p. 232-
781 42.
- 782 25. Deng, W., et al., Controlling long-range genomic interactions at a native locus by
783 targeted tethering of a looping factor. *Cell*, 2012. 149(6): p. 1233-44.
- 784 26. Chang, T.H., et al., An enhancer directs differential expression of the linked *Mrf4* and
785 *Myf5* myogenic regulatory genes in the mouse. *Dev Biol*, 2004. 269(2): p. 595-608.
- 786 27. Inoue, F. and N. Ahituv, Decoding enhancers using massively parallel reporter assays.
787 *Genomics*, 2015. 106(3): p. 159-164.
- 788 28. van Arensbergen, J., et al., High-throughput identification of human SNPs affecting
789 regulatory element activity. *Nat Genet*, 2019. 51(7): p. 1160-1169.
- 790 29. Sahu, B., et al., Sequence determinants of human gene regulatory elements. *bioRxiv*,
791 2021: p. 2021.03.18.435942.
- 792 30. van Arensbergen, J., et al., Genome-wide mapping of autonomous promoter activity in
793 human cells. *Nat Biotechnol*, 2017. 35(2): p. 145-153.
- 794 31. Arnold, C.D., et al., Genome-wide quantitative enhancer activity maps identified by
795 STARR-seq. *Science*, 2013. 339(6123): p. 1074-7.
- 796 32. Klein, J.C., et al., A systematic evaluation of the design and context dependencies of
797 massively parallel reporter assays. *Nat Methods*, 2020. 17(11): p. 1083-1091.
- 798 33. Davis, J.E., et al., Dissection of c-AMP Response Element Architecture by Using
799 Genomic and Episomal Massively Parallel Reporter Assays. *Cell Syst*, 2020. 11(1): p.
800 75-85 e7.
- 801 34. King, D.M., et al., Synthetic and genomic regulatory elements reveal aspects of cis-
802 regulatory grammar in mouse embryonic stem cells. *Elife*, 2020. 9.
- 803 35. Frankish, A., et al., GENCODE reference annotation for the human and mouse
804 genomes. *Nucleic Acids Res*, 2019. 47(D1): p. D766-D773.
- 805 36. Ohler, U., et al., Computational analysis of core promoters in the *Drosophila* genome.
806 *Genome Biol*, 2002. 3(12): p. RESEARCH0087.
- 807 37. Groudine, M., et al., Human fetal to adult hemoglobin switching: changes in chromatin
808 structure of the beta-globin gene locus. *Proc Natl Acad Sci U S A*, 1983. 80(24): p.
809 7551-5.
- 810 38. Joshi, O., et al., Dynamic Reorganization of Extremely Long-Range Promoter-Promoter
811 Interactions between Two States of Pluripotency. *Cell Stem Cell*, 2015. 17(6): p. 748-
812 757.
- 813 39. Long, H.K., S.L. Prescott, and J. Wysocka, Ever-Changing Landscapes: Transcriptional
814 Enhancers in Development and Evolution. *Cell*, 2016. 167(5): p. 1170-1187.
- 815 40. Barakat, T.S., et al., Functional Dissection of the Enhancer Repertoire in Human
816 Embryonic Stem Cells. *Cell Stem Cell*, 2018. 23(2): p. 276-288 e8.
- 817 41. Agrawal, P., et al., Genome editing demonstrates that the -5 kb *Nanog* enhancer
818 regulates *Nanog* expression by modulating RNAPII initiation and/or recruitment. *J Biol*
819 *Chem*, 2021. 296: p. 100189.
- 820 42. Djebali, S., et al., Landscape of transcription in human cells. *Nature*, 2012. 489(7414):
821 p. 101-8.

- 822 43. Andersson, R., et al., An atlas of active enhancers across human cell types and tissues.
823 Nature, 2014. 507(7493): p. 455-461.
- 824 44. Consortium, E.P., et al., Expanded encyclopaedias of DNA elements in the human and
825 mouse genomes. Nature, 2020. 583(7818): p. 699-710.
- 826 45. Khan, A. and X. Zhang, dbSUPER: a database of super-enhancers in mouse and
827 human genome. Nucleic Acids Res, 2016. 44(D1): p. D164-71.
- 828 46. Blinka, S., et al., Super-Enhancers at the Nanog Locus Differentially Regulate
829 Neighboring Pluripotency-Associated Genes. Cell Rep, 2016. 17(1): p. 19-28.
- 830 47. Zabidi, M.A., et al., Enhancer-core-promoter specificity separates developmental and
831 housekeeping gene regulation. Nature, 2015. 518(7540): p. 556-9.
- 832 48. Hounkpe, B.W., et al., HRT Atlas v1.0 database: redefining human and mouse
833 housekeeping genes and candidate reference transcripts by mining massive RNA-seq
834 datasets. Nucleic Acids Res, 2021. 49(D1): p. D947-D955.
- 835 49. Kim, J., et al., An extended transcriptional network for pluripotency of embryonic stem
836 cells. Cell, 2008. 132(6): p. 1049-61.
- 837 50. Akagi, T., et al., ETS-related transcription factors ETV4 and ETV5 are involved in
838 proliferation and induction of differentiation-associated genes in embryonic stem (ES)
839 cells. J Biol Chem, 2015. 290(37): p. 22460-73.
- 840 51. Kinoshita, K., et al., GABPalph regulates Oct-3/4 expression in mouse embryonic stem
841 cells. Biochem Biophys Res Commun, 2007. 353(3): p. 686-91.
- 842 52. Zhang, X., et al., FOXO1 is an essential regulator of pluripotency in human embryonic
843 stem cells. Nat Cell Biol, 2011. 13(9): p. 1092-9.
- 844 53. Wei, Z., et al., Klf4 interacts directly with Oct4 and Sox2 to promote reprogramming.
845 Stem Cells, 2009. 27(12): p. 2969-78.
- 846 54. Bonev, B., et al., Multiscale 3D Genome Rewiring during Mouse Neural Development.
847 Cell, 2017. 171(3): p. 557-572 e24.
- 848 55. Hsieh, T.S., et al., Resolving the 3D Landscape of Transcription-Linked Mammalian
849 Chromatin Folding. Mol Cell, 2020. 78(3): p. 539-553 e8.
- 850 56. Kermekchiev, M., et al., Every enhancer works with every promoter for all the
851 combinations tested: could new regulatory pathways evolve by enhancer shuffling?
852 Gene Expr, 1991. 1(1): p. 71-81.
- 853 57. Gehrig, J., et al., Automated high-throughput mapping of promoter-enhancer
854 interactions in zebrafish embryos. Nat Methods, 2009. 6(12): p. 911-6.
- 855 58. Tak, Y.E., et al., Augmenting and directing long-range CRISPR-mediated activation in
856 human cells. Nat Methods, 2021. 18(9): p. 1075-1081.
- 857 59. El Khattabi, L., et al., A Pliable Mediator Acts as a Functional Rather Than an
858 Architectural Bridge between Promoters and Enhancers. Cell, 2019. 178(5): p. 1145-
859 1158 e20.
- 860 60. Fuqua, T., et al., Dense and pleiotropic regulatory information in a developmental
861 enhancer. Nature, 2020. 587(7833): p. 235-239.
- 862 61. Kircher, M., et al., Saturation mutagenesis of twenty disease-associated regulatory
863 elements at single base-pair resolution. Nat Commun, 2019. 10(1): p. 3583.
- 864 62. Pang, B. and M.P. Snyder, Systematic identification of silencers in human cells. Nat
865 Genet, 2020. 52(3): p. 254-263.
- 866 63. Gisselbrecht, S.S., et al., Transcriptional Silencers in Drosophila Serve a Dual Role as
867 Transcriptional Enhancers in Alternate Cellular Contexts. Mol Cell, 2020. 77(2): p. 324-
868 337 e8.
- 869 64. You, F.M., et al., BatchPrimer3: a high throughput web application for PCR and
870 sequencing primer design. BMC Bioinformatics, 2008. 9: p. 253.

- 871 65. Langmead, B. and S.L. Salzberg, Fast gapped-read alignment with Bowtie 2. *Nat*
872 *Methods*, 2012. 9(4): p. 357-9.
- 873 66. Zorita, E., P. Cusco, and G.J. Filion, Starcode: sequence clustering based on all-pairs
874 search. *Bioinformatics*, 2015. 31(12): p. 1913-9.
- 875 67. Diaferia, G.R., et al., Dissection of transcriptional and cis-regulatory control of
876 differentiation in human pancreatic cancer. *EMBO J*, 2016. 35(6): p. 595-617.
- 877 68. Team, R.C., R: A Language and Environment for Statistical Computing. R Foundation
878 for Statistical Computing, 2021. 2021.
- 879 69. Weirauch, M.T., et al., Determination and inference of eukaryotic transcription factor
880 sequence specificity. *Cell*, 2014. 158(6): p. 1431-1443.
- 881

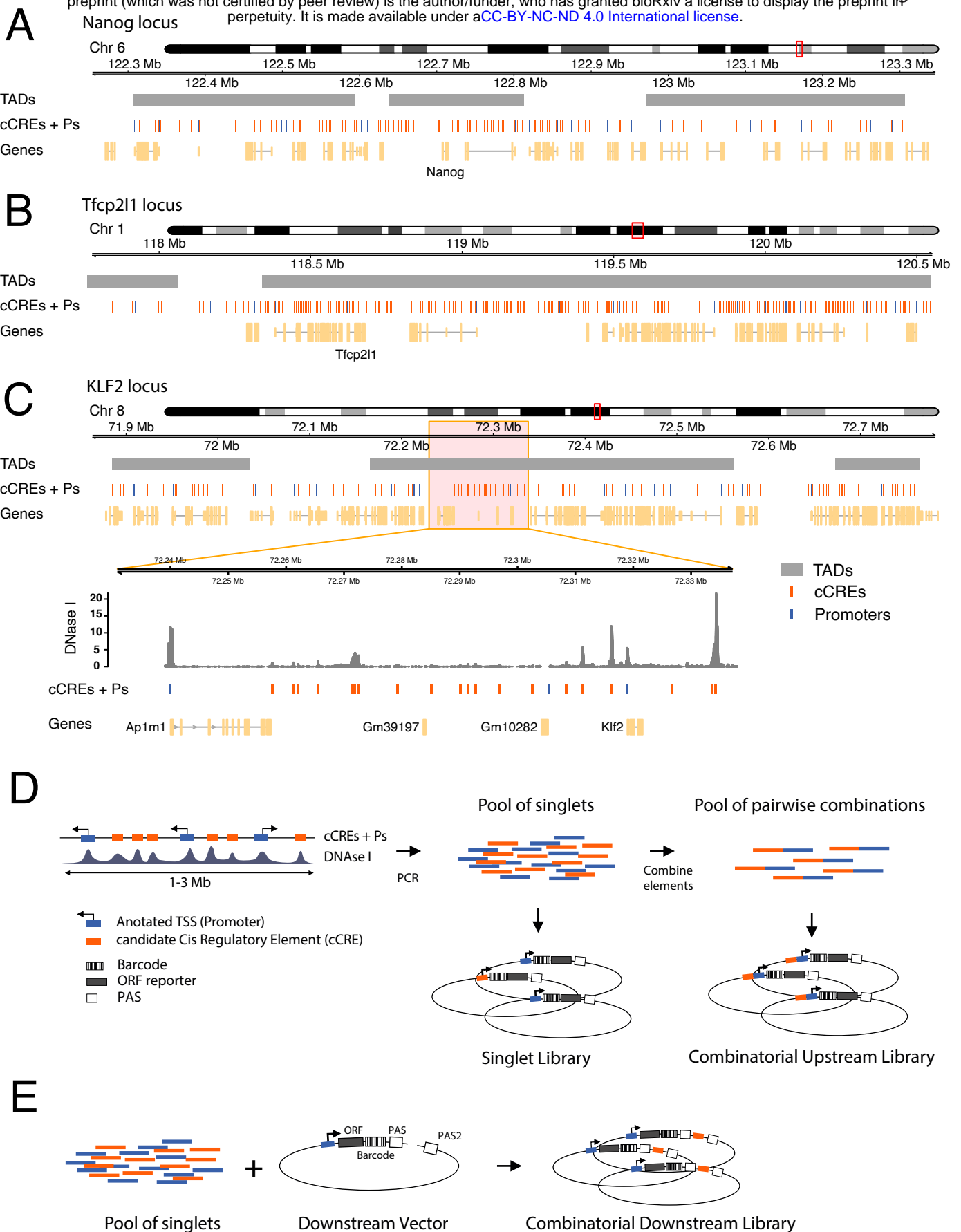
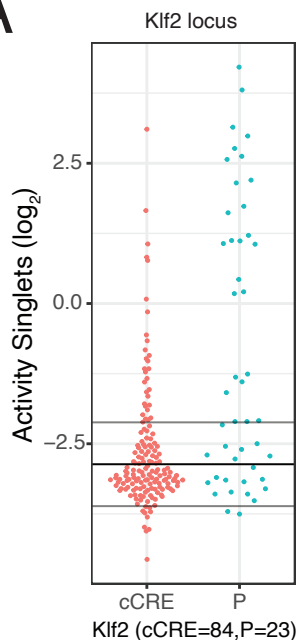
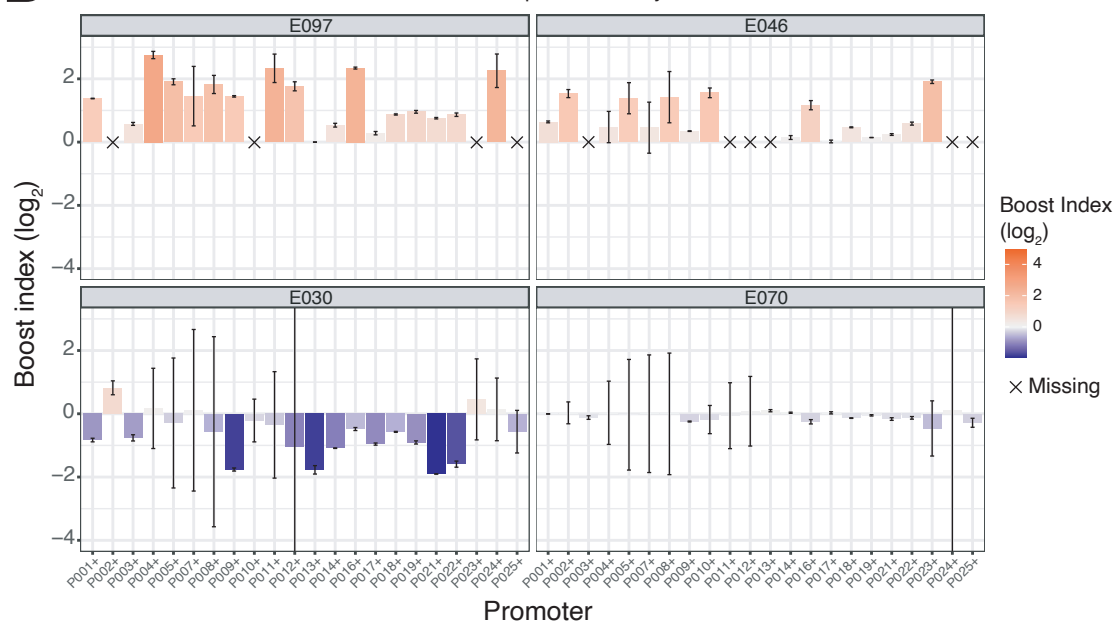


Figure 1. Regulatory element selection and library construction. **A-C)** Representations of Nanog, Tfcp2l1, and Klf2 loci, respectively. In **C)** the zoom-in displays a DNase I sensitivity track [38] where peaks overlap with cCREs. **D)** Cloning strategy for the Upstream assay. cCREs and promoters were amplified by PCR from genomic DNA and pooled. Fragments in this pool were then randomly ligated to generate duplets. Singlets and duplets were cloned into the same barcoded vector to generate two libraries per locus, a singlet library and a combinatorial library. **E)** Cloning strategy for the Downstream assay. The singlet pool from the Klf2 locus was cloned into ten vectors, each of them carrying a different promoter. The resulting ten sub-libraries were combined into one Downstream assay library.

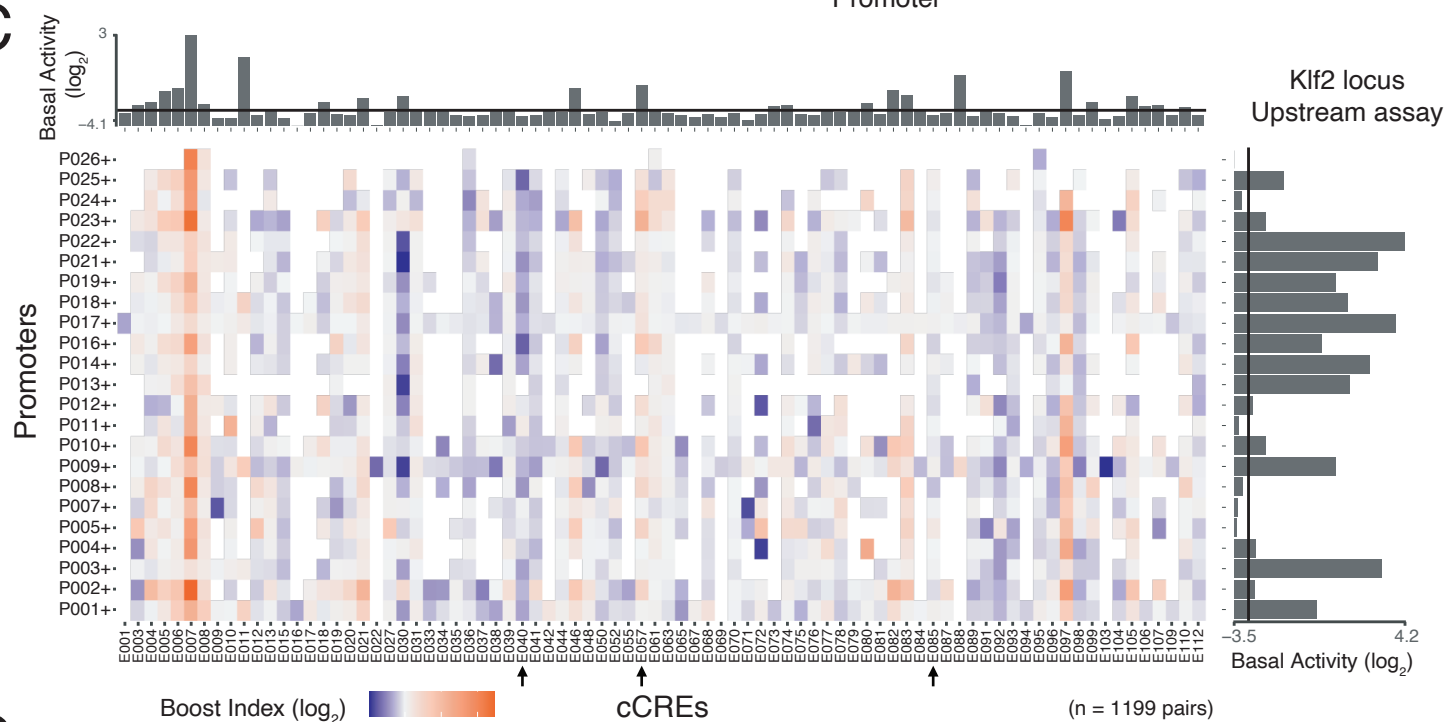
A



B



C



D

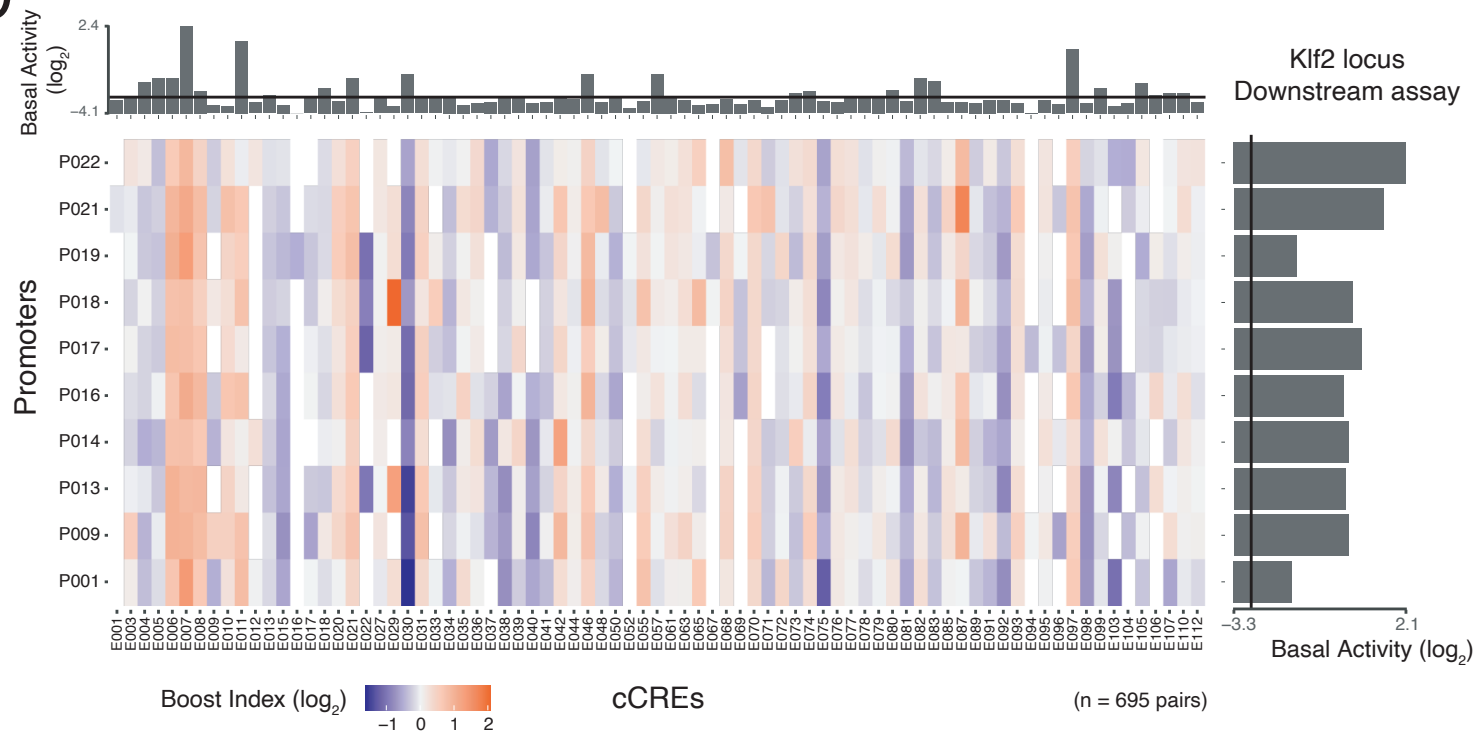


Figure 2. Singlet and combinatorial activities of cCREs and promoters from the *Klf2* locus. **A)** Transcription activities of singlet cCREs and promoters. Each dot represents the mean activity of one singlet. Horizontal lines represent the average background activity of empty vectors (black line) plus or minus two standard deviations (grey lines). Elements with activities more than two standard deviations above the average background signal are defined as active. **B)** Examples of Upstream assay cCRE-P combinations for cCREs E097, E046, E030 and E070 of the *Klf2* locus. Barplots represent the mean boost index of each combination, vertical lines represent the standard deviations. Crosses mark missing data. **C-D)** Boost index matrices of cCRE–P combinations from the *Klf2* locus according to Upstream **(C)** and Downstream **(D)** assays. White tiles indicate missing data. Barplots on the right and top of each panel show basal activities of each tested P or cCRE, respectively, with the black line indicating the background activity of the empty vector. All data are averages over 3 independent biological replicates.

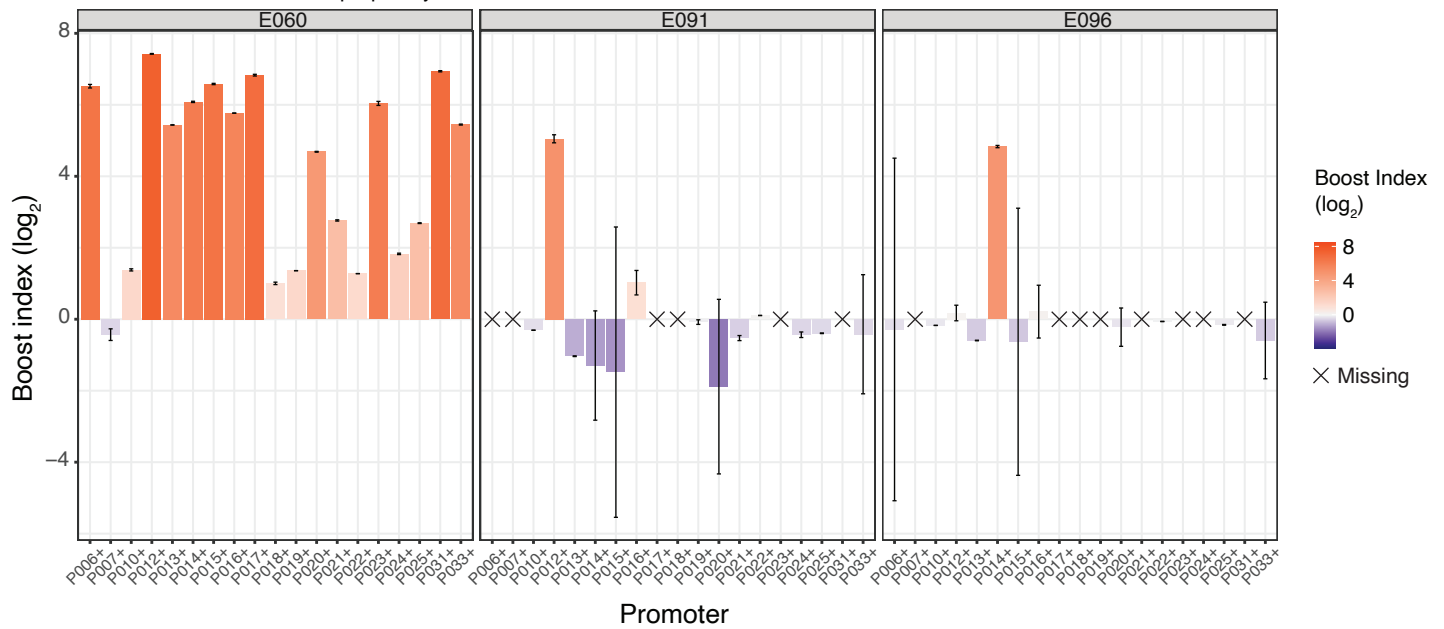
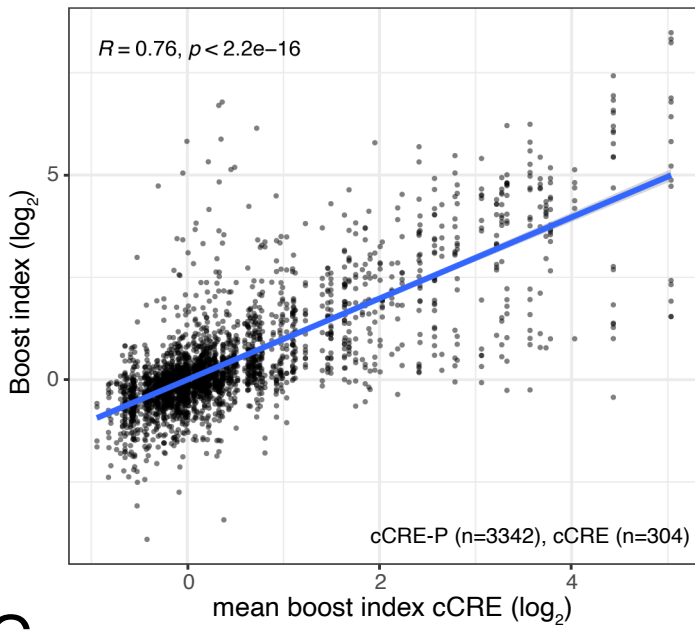
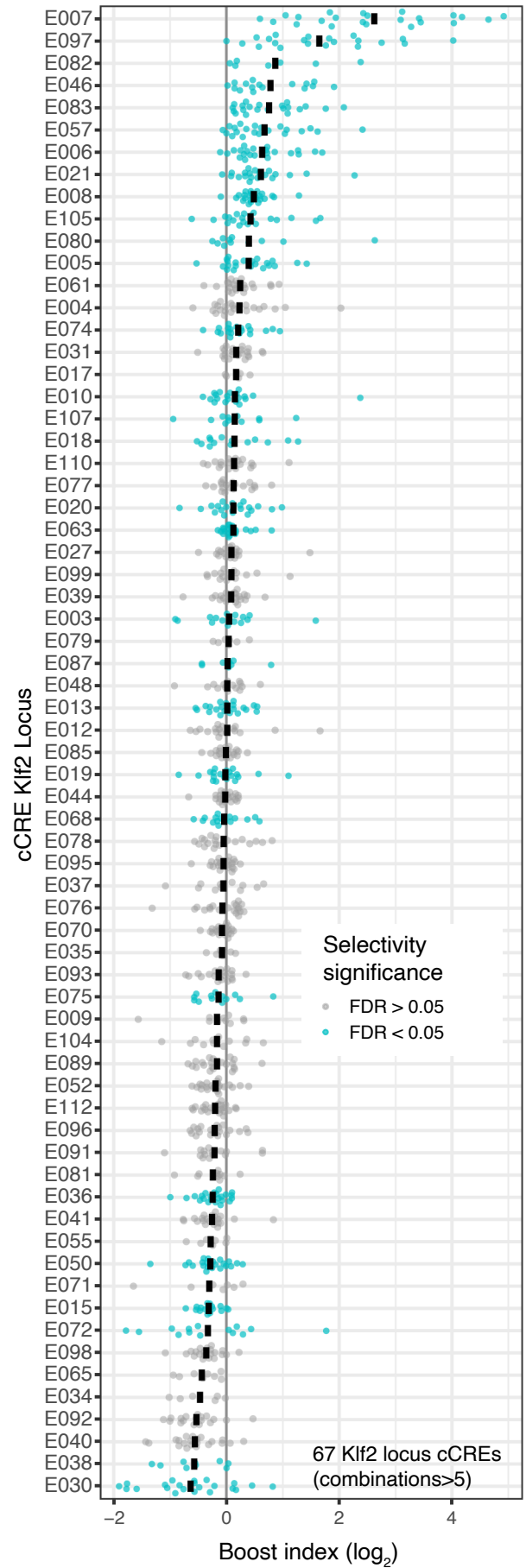


Figure 3. Examples of selective cCREs from the *Tfcp2l1* locus. Boost indices obtained in the Upstream assay are shown for cCRE-P combinations of cCREs E060, E091 and E096 of the *Tfcp2l1* locus. Barplots indicate the mean boost index of each combination, vertical lines indicate standard deviations. All data are averages over 3 independent biological replicates.

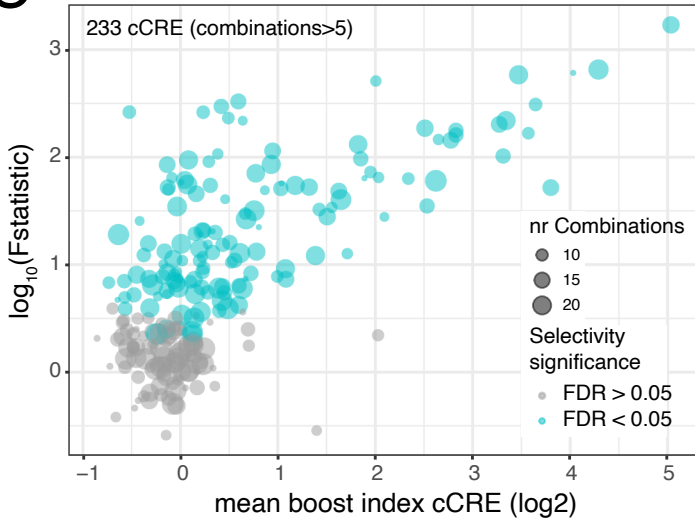
A



B



C



D

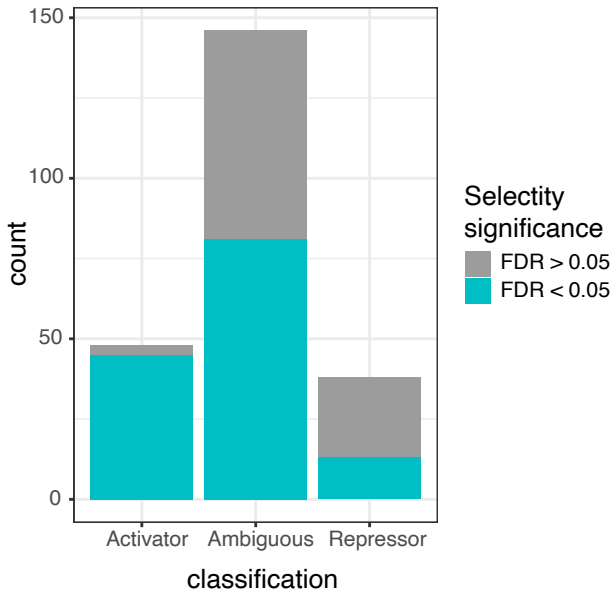
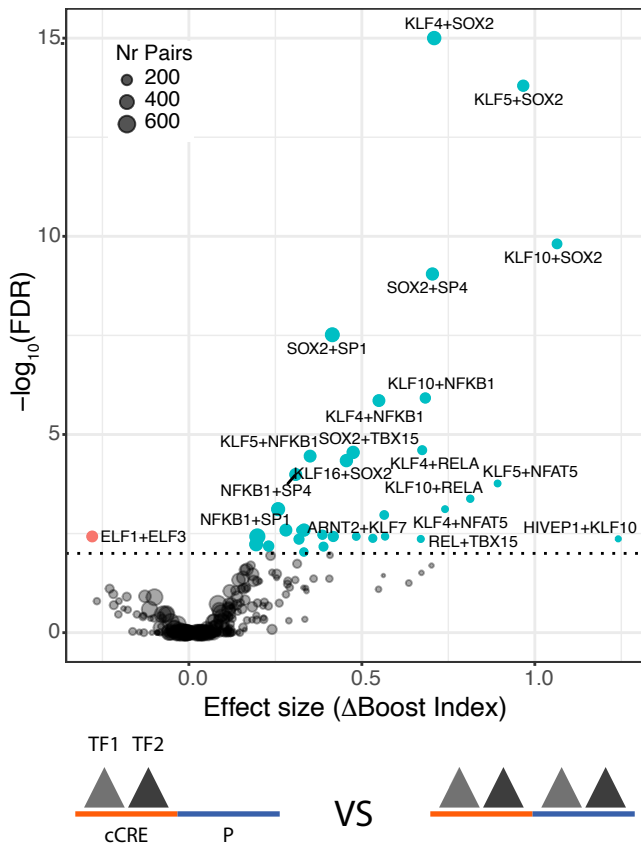


Figure 4. Promoter selectivity of cCREs. **A)** Plot showing the broad diversity of boost indices of many cCREs. Data are from Upstream assays of *Klf2*, *Nanog* and *Tfcp2l1* loci combined. Vertical axis indicates boost indices of all tested cCRE–P pairs, which are horizontally ordered by the mean boost index of each cCRE. **B)** Boost index distributions for each cCRE from the *Klf2* locus (Upstream assay). Each dot represents one cCRE–P combination; black bar represents the mean. Turquoise colouring marks cCREs that have a larger variance of their boost indices than may be expected based on experimental noise, according to the Welch F-test after multiple hypothesis correction (5% FDR cutoff). **C)** Summary of Welch F-test selectivity analysis results for all cCREs from the three loci with more than 5 cCRE–P combinations. Each dot represents one cCRE; the size of the dots indicates the number of cCRE–P pairs. Significantly selective cCREs (5% FDR cutoff) are highlighted in turquoise. **D)** Proportion of significantly selective (turquoise) cCRE in the three categories as shown in [Figure S3A](#). All data are averages over 3 independent biological replicates.

A



B

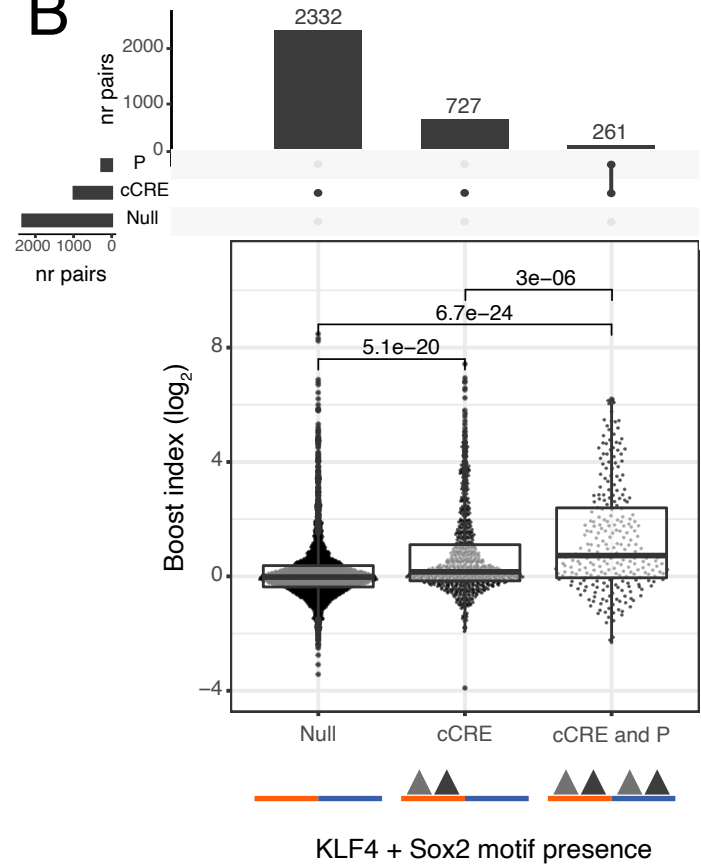


Figure 5. Association of TF motif Duos with higher boost indices. **A)** Results of TF survey for self-compatible TF motif Duos. TF motif duos associated with higher or lower boost indices at a 1% FDR cutoff are highlighted. **B)** Association of Sox2+Klf4 motifs at both cCRE and P with higher boost indices. cCRE-P combinations are split into 3 groups according to presence or absence of Sox2+Klf4 motifs both at the cCRE and the promoter, or only the cCRE. Numbers at the top of horizontal brackets are the p-values obtained from comparing the different groups boost index distributions using a Wilcoxon rank-sum test. Boxplots represent median and interquartile ranges. Barplots at the top represent the number of combinations in each group.

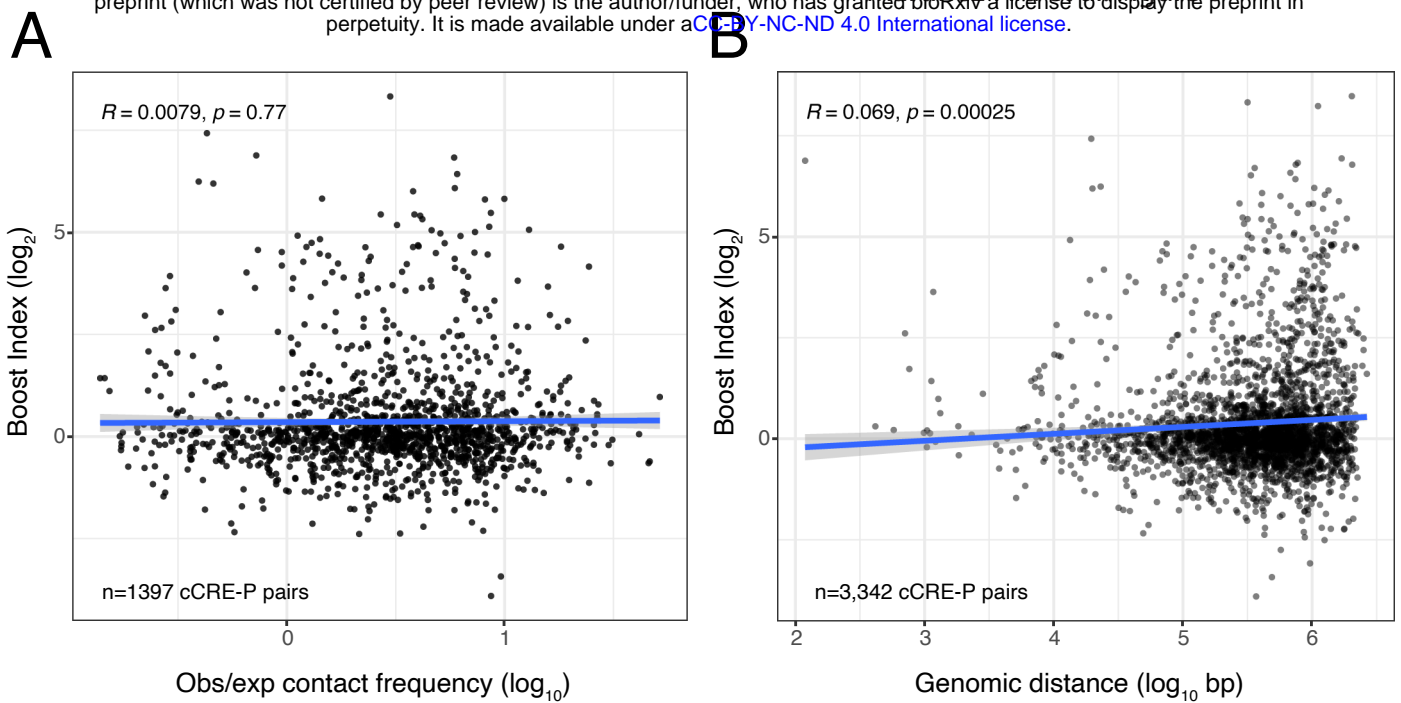


Figure 6. Absent or very weak correlation between boost indices and **(A)** contact frequencies according to micro-C [55] or **(B)** linear genomic distance, for all cCRE-P pairs from the three loci combined. All boost index data are averages over 3 independent biological replicates.

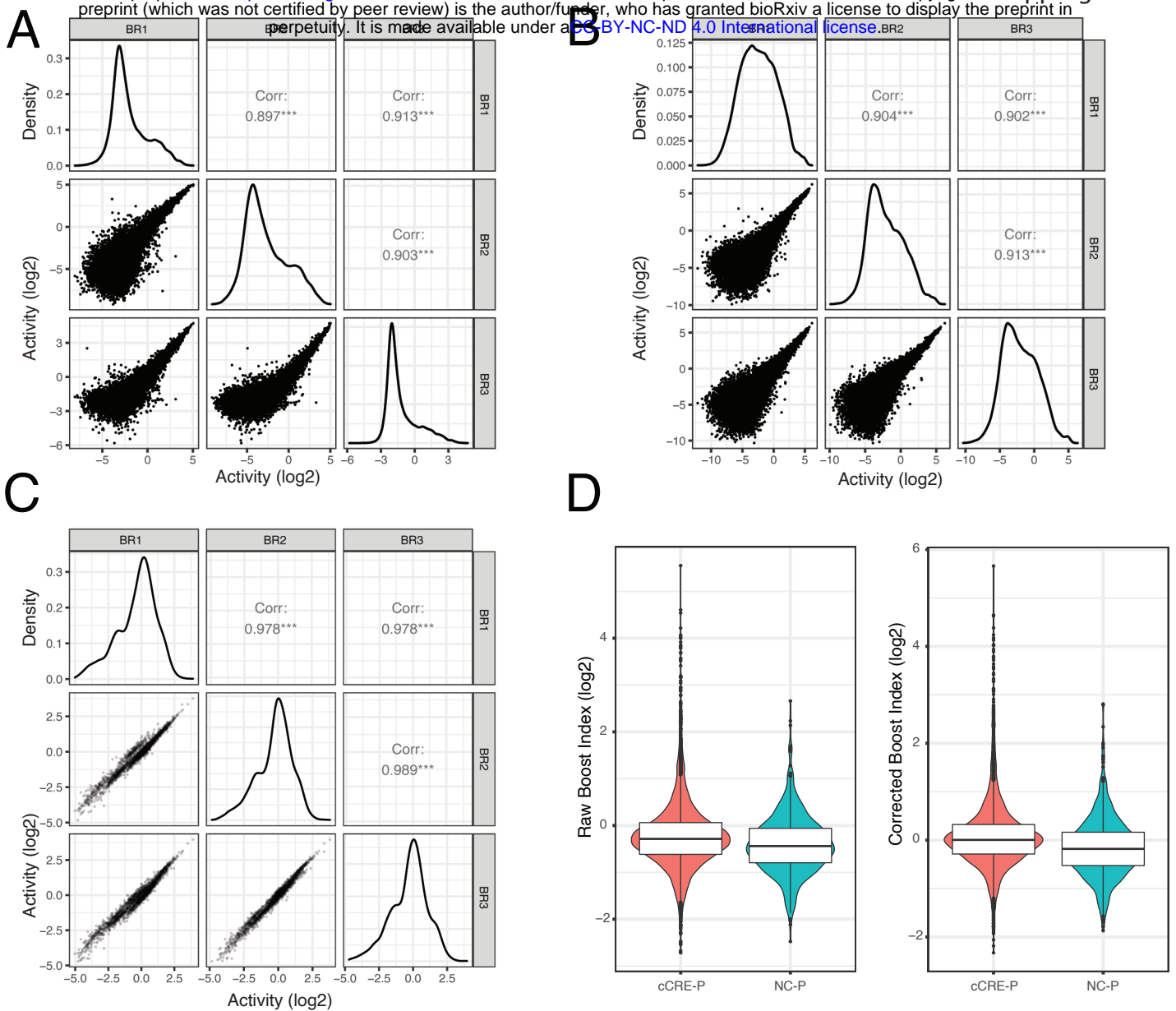


Figure S1. Reproducibility of data and boost index calculation. **(A-C)** Correlograms of the three biological replicates of each library pool. Lower left panels show pairwise scatterplots of the activities of all cCRE-P pairs per replicate. Middle panels show the density of data distribution in each replicate and upper right panels show the Pearson correlation coefficients. **A)** Klf2 and Nanog Upstream libraries. **B)** Tfc211 Upstream library. **C)** Klf2 Downstream libraries. **D)** Upstream assay boost index distributions for cCRE-P and negative controls – promoter (NC-P) combinations. Left panel: raw boost indices; right panel: boost indices after correction for negative bias (see Methods).

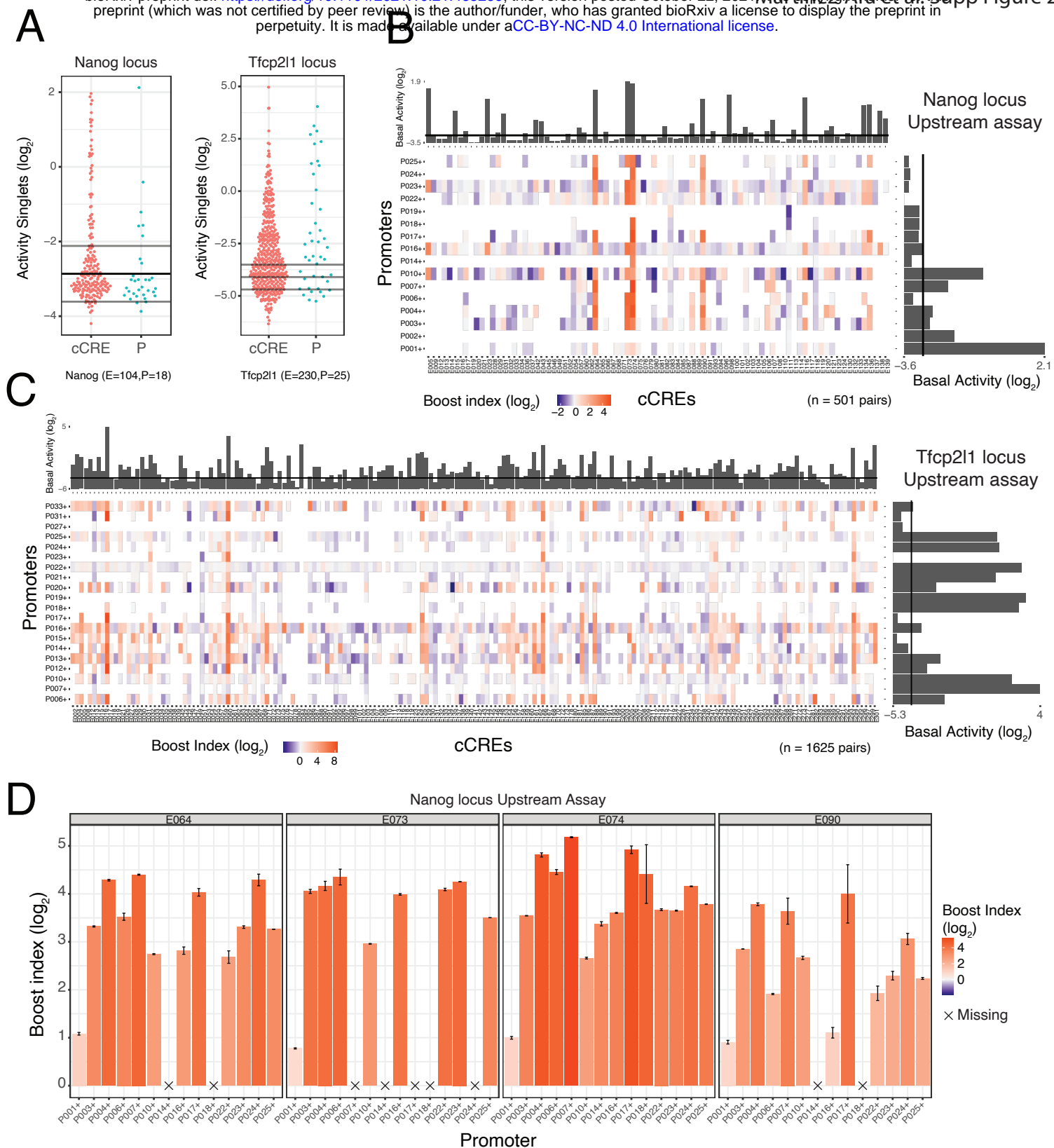


Figure S2. Element activities and boost indices obtained with Nanog and Tfc2p11 Upstream libraries. **A)** Transcriptional activities of cCREs and promoters. Each dot represents the mean activity of one singlet. Horizontal lines represent the average background activity of empty vectors (black line) plus or minus two standard deviations (grey lines). Elements with activities more than two standard deviations above the average background signal are defined as active. **B-C)** Boost index matrices for cCRE–P pairs from Nanog and Tfc2p11 loci (both Upstream assays). White tiles indicate missing data. Barplots on the right and top of each panel show basal activities of each tested P or cCRE, respectively, with the black line indicating the background activity of the empty vector. **D)** Examples of cCRE–P combinations for cCREs E064, E073, E074 and E090 of the Nanog locus. Barplots represent the mean boost index of each combination, vertical lines represent the standard deviation of each boost index. All data are averages over 3 independent biological replicates.

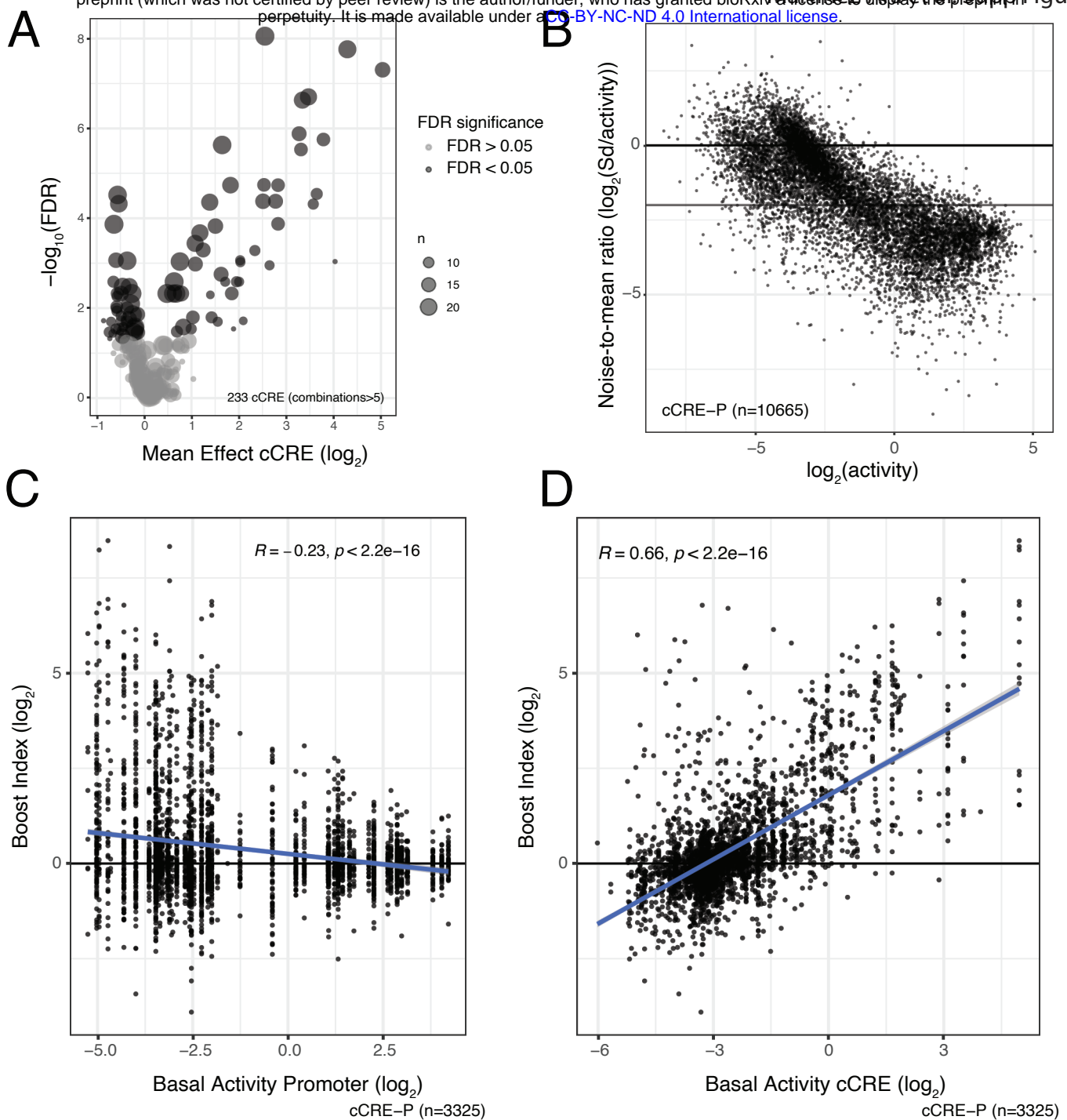


Figure S3. cCRE functional classification and activity influence on Boost indices. **A)** Volcano plot of cCREs associated with activation or repression across promoters. A Wilcoxon test is performed per cCRE comparing the boost indices of all the cCRE-P combinations of that cCRE against the rest of cCRE-P combinations. A minimum of 6 combinations is required per cCRE. P-values are corrected for multiple hypothesis testing using the Benjamini-Hochberg method (FDR). **B)** Relationship between noise-to-mean ratio (Standard Deviation/mean Activity) and mean activity of cCRE-Ps. Horizontal lines represent noise-to-mean ratios of 1 and of 4 in \log_2 scale. **C)** Relationship between boost indices and basal (singlet) P activity. Each column of dots shows the data of cCRE-P pairs for one P. Data are from Upstream assays of all three loci combined. **D)** Relationship between boost indices and basal (singlet) cCRE activity. All data are averages over 3 independent biological replicates.

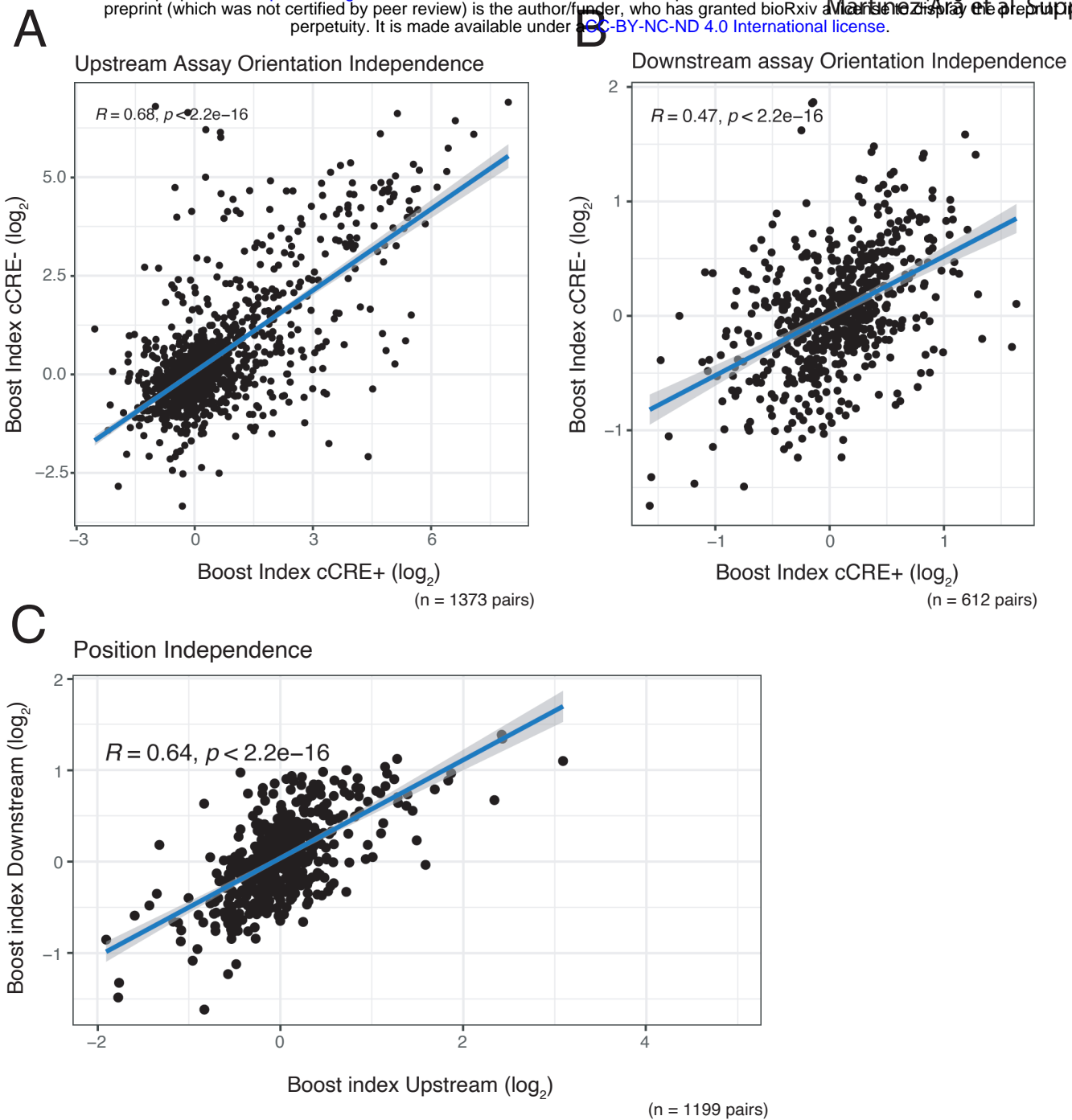


Figure S4. Orientation and position independence of cCREs. **(A-B)** Correlation between boost indices of both cCRE orientations of the same cCRE-P combination, in the **(A)** Upstream assay and **(B)** Downstream assay. Data are from the Klf2 locus libraries. Note that "+" and "-" orientations are arbitrary labels, because cCREs do not have an intrinsic orientation. **(C)** Correlation between boost indices of cCRE-P combinations shared between the Upstream and Downstream assays of the Klf2 locus. In all panels R is the Pearson correlation coefficient. All data are averages over 3 independent biological replicates. In C Boost indices are averaged over cCRE orientations.

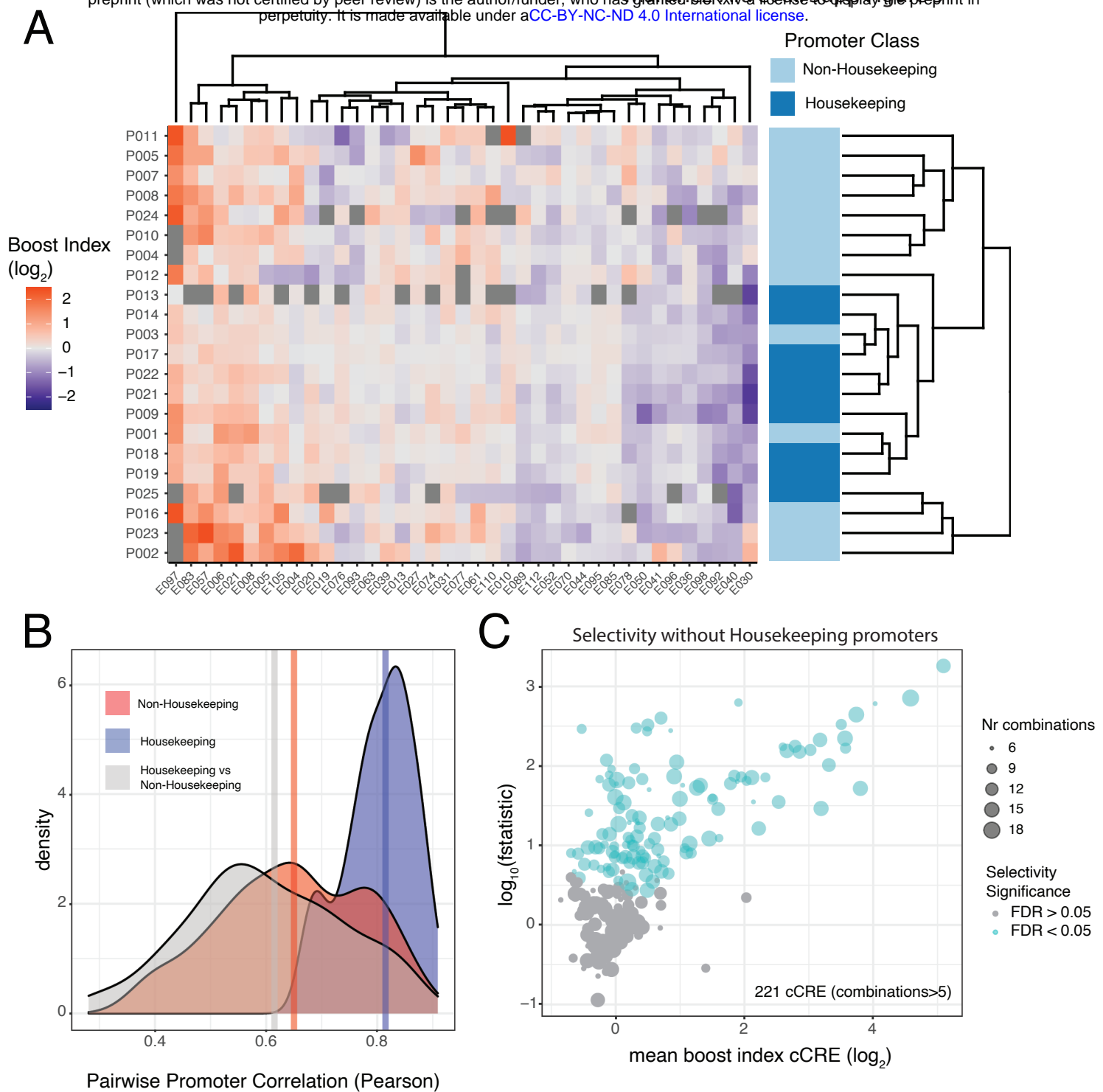
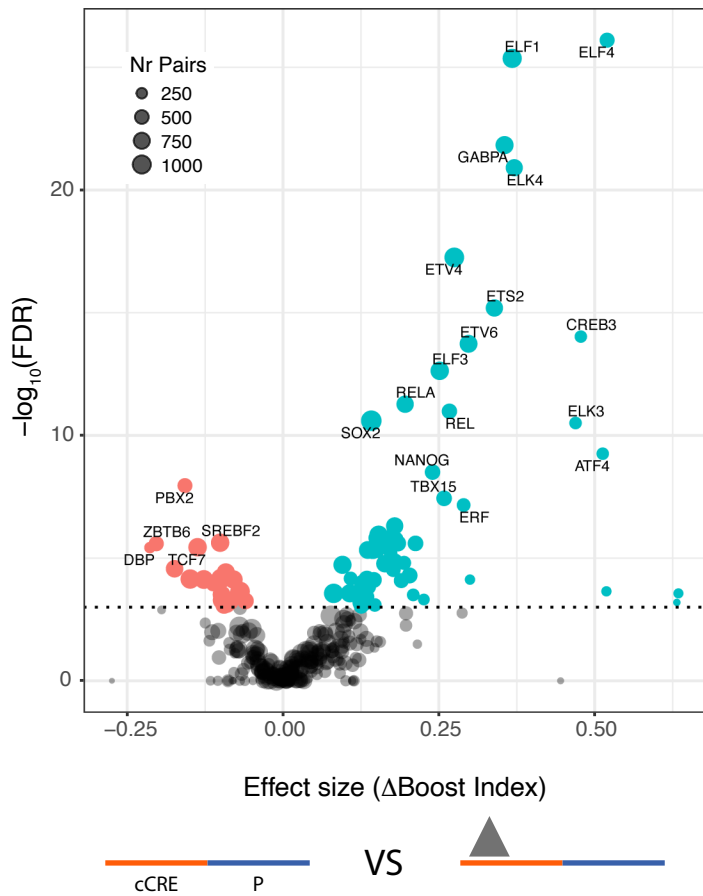


Figure S5. Housekeeping promoters show a distinct pattern of cCRE compatibility. **A)** Hierarchical clustering of the Upstream assay boosting matrix of the Klf2 locus. In order to facilitate hierarchical clustering the matrix has been restricted to almost complete cases (cCREs >15 combinations) **B)** Density plot of pairwise Pearson correlation coefficients of the boost indices of Klf2 locus promoters classified as either housekeeping or non-housekeeping [48]. Blue: correlations between all pairs of housekeeping promoters; red: all correlations between pairs of non-housekeeping promoters; grey: all correlations between one housekeeping and one non-housekeeping promoter. Vertical lines represent the median of each group. Unlike in **(A)**, all promoters in the Upstream assay were included in this analysis. **C)** Results of selectivity analysis as performed in Figure 4C, but excluding housekeeping promoters. All data are averages over 3 independent biological replicates.

A



B

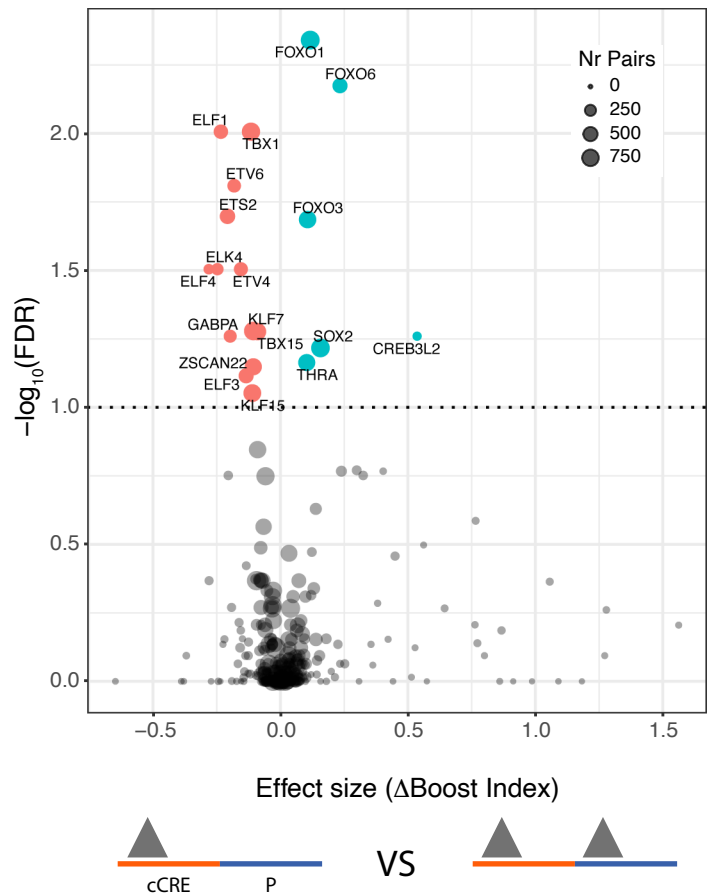


Figure S6. Identification of single TF motifs that correlate with boost indices. **(A)** TF motifs in cCREs associated (at 1% FDR cutoff) with activation (turquoise) or repression (red). **(B)** Motifs of putative self-compatible TFs, i.e. motifs that predict increased or reduced boosting indices when present both at the cCRE and P, compared to being present only at the cCRE. TF motifs associated with higher or lower boost indices at a 1% FDR cutoff are highlighted. We note that TF motifs with multiple hits from the same family, such as for ELK, FOXO and ELF factors, may in fact be due to the activity of one TF motif of that family [69].