

1 **Parallel Characterization of *cis*-Regulatory Elements for Multiple Genes Using**
2 **CRISPRpath**

3

4 **Authors:** Xingjie Ren¹, Mengchi Wang², Bingkun Li¹, Kirsty Jamieson¹, Lina Zheng², Ian R.
5 Jones¹, Bin Li³, Maya Asami Takagi¹, Jerry Lee¹, Lenka Maliskova¹, Tsz Wai Tam¹, Miao Yu³,
6 Rong Hu³, Lindsay Lee⁴, Armen Abnousi⁴, Gang Li⁵, Yun Li^{6,7,8}, Ming Hu⁴, Bing Ren^{3,9,10}, Wei
7 Wang^{2,11}, Yin Shen^{1,12,#}

8

9 **Affiliation:**

10 ¹Institute for Human Genetics, University of California, San Francisco, San Francisco, CA,
11 USA.

12 ²Bioinformatics and Systems Biology Graduate Program, University of California, San Diego,
13 La Jolla, CA, USA.

14 ³Ludwig Institute for Cancer Research, La Jolla, California, USA.

15 ⁴Department of Quantitative Health Sciences, Lerner Research Institute, Cleveland Clinic
16 Foundation, Cleveland, OH, USA.

17 ⁵Department of Statistics and Operations Research, University of North Carolina, Chapel Hill,
18 NC, USA.

19 ⁶Department of Biostatistics, University of North Carolina, Chapel Hill, NC, USA.

20 ⁷Department of Genetics, University of North Carolina, Chapel Hill, NC, USA.

21 ⁸Department of Computer Science, University of North Carolina, Chapel Hill, NC, USA.

22 ⁹Department of Cellular and Molecular Medicine, University of California, San Diego, La Jolla,
23 CA, USA.

24 ¹⁰Moore's Cancer Center, University of California, San Diego, La Jolla, California, USA.

25 ¹¹Department of Chemistry and Biochemistry, University of California, San Diego, La Jolla, CA,
26 USA.

27 ¹²Department of Neurology, University of California, San Francisco, San Francisco, CA, USA.

28

29 #Corresponding Author: yin.shen@ucsf.edu

30

31 **Abstract**

32 Current pooled CRISPR screens for *cis*-regulatory elements (CREs) can only accommodate
33 one gene based on its expression level. Here, we describe CRISPRpath, a scalable screening
34 strategy for parallelly characterizing CREs of genes linked to the same biological pathway and
35 converging phenotypes. We demonstrate the ability of CRISPRpath for simultaneously
36 identifying functional enhancers of six genes in the 6-thioguanine-induced DNA mismatch
37 repair pathway using both CRISPR interference (CRISPRi) and CRISPR nuclease (CRISPRn)
38 approaches. 60% of the identified enhancers are known promoters with distinct epigenomic
39 features compared to other active promoters, including increased chromatin accessibility and
40 interactivity. Furthermore, by imposing different levels of selection pressure, CRISPRpath can
41 distinguish enhancers exerting strong impact on gene expression from those exerting weak
42 impact. Our results offer a nuanced view of *cis*-regulation and demonstrate that CRISPRpath
43 can be leveraged for understanding the complex gene regulatory program beyond
44 transcriptional output at scale.

45

46 **Main**

47 *Cis*-regulatory elements (CREs) are key regulators for spatial-temporal control of gene
48 expression. Mutations in CREs can contribute to complex diseases by modulating gene
49 expression over long genomic distances¹⁻³. Thus, functionally characterizing CREs can provide
50 important insight into gene regulation mechanisms and enable us to better interpret non-coding
51 genetic variants associated with diseases. Despite the fact that tremendous numbers of
52 candidate CREs have been mapped by biochemical signature⁴, our knowledge of whether,
53 how, and how much these putative CREs are functional on gene expression remain scarce in
54 the human genome. Pooled CRISPR screens have been developed for testing CREs in their
55 native chromatin context by monitoring the transcriptional levels for the gene of interest⁵⁻¹¹.
56 Although results from these studies have made significant contributions to the annotation of
57 functional DNA elements, challenges remain in pooled CRISPR screens of CREs. First,
58 CRISPR screens for enhancers based on gene expression levels largely depend on
59 generating reporter knock-in cell lines⁷ or using FlowFISH signals⁸. These procedures,
60 involving generation of reporter lines and selection of cells with positive hits by flow cytometry,
61 are time-consuming and difficult to scale up to multiple genes in the same experiment. Second,

62 the approaches of using gene expression as the screening phenotype^{9, 10} fail to connect the
63 functions of DNA elements from transcriptional regulation at the molecular level to interpretable
64 cellular and physiological functions. Third, in cases of CRE screens using phenotypes such as
65 cell proliferation and survival^{11, 12}, they fail to quantify the effect sizes of enhancers on
66 transcriptional output.

67

68 To address these limitations, we developed CRISPRpath, a pooled CRISPR screening
69 approach to simultaneously characterizing CREs for multiple target genes involved in the same
70 biological pathway. CRISPRpath allows us to screen functional DNA elements based on
71 phenotypes associated with well-defined biological pathways. We demonstrate the capacity of
72 CRISPRpath by performing CRISPR interference (CRISPRi) and nuclease (CRISPRn) screens
73 for six genes in human induced pluripotent stem cells (iPSCs), and reveal different strengths of
74 enhancer functions by imposing varying levels of selection pressure on the cells.

75

76 **Results**

77 **Leveraging CRISPRpath for parallel characterization of CREs for multiple genes in** 78 **iPSCs.**

79 To characterize candidate *cis*-regulatory elements (CREs) for multiple genes within the same
80 pooled CRISPR screening, we designed and applied CRISPRpath to six genomic loci
81 containing six genes (*HPRT1*, *MSH2*, *MSH6*, *MLH1*, *PMS2*, *PCNA*) involved in the 6-
82 thioguanine (6TG)-induced mismatch repair (MMR) (**Fig. 1a**). The MMR pathway is highly
83 conserved and essential for the maintenance of genome stability¹³. The MMR pathway
84 recognizes DNA mismatches caused by 6TG treatment and induces cell apoptosis^{14, 15}. On the
85 other hand, cells with a malfunctioning MMR pathway, due to aberrant expression levels of
86 6TG metabolism genes or MMR genes, may survive during 6TG treatment. Employing the
87 properties of the MMR pathway, we used cell survival for selecting cells with the reduced
88 expression of MMR genes due to defects in enhancer activities (**Fig. 1b**). To design the
89 screening library, we first identified open chromatin regions by performing Assay for
90 Transposase Accessible Chromatin using sequencing (ATAC-Seq) in WTC11 iPSCs. We
91 included all open chromatin regions defined by ATAC-seq peaks located 1Mb upstream and
92 1Mb downstream of each of the six genes (spanning a total of 10.6 Mb genomic regions) as

93 candidate CREs for functional characterization (**Supplementary Fig. 1a, b, Supplementary**
94 **Table 1**). We then designed a sgRNA library with 32,383 distal sgRNAs targeting 294 distal
95 ATAC-seq peaks, 2,755 proximal sgRNAs targeting 81 ATAC-seq peaks overlapped with
96 transcription start site (TSS) and coding regions of the six genes, and 625 non-targeting
97 sgRNAs with genomic sequences in the same genomic loci but are not followed by PAM
98 sequences (**Supplementary Fig. 1c, Supplementary Table 2**). In total, we included 35,763
99 sgRNAs in the library with an average of 110 sgRNA per ATAC-seq peak (**Supplementary**
100 **Fig. 1d, Supplementary Fig. 2a, b**). We generated a lentiviral library expressing these
101 sgRNAs and transduced this library into two engineered WTC11 iPSC lines, one expressing
102 doxycycline-inducible dCas9-KRAB (CRISPRi) and the other doxycycline-inducible Cas9
103 (CRISPRn)¹⁶, both at a multiplicity of infection (MOI) of 0.5 (**Fig. 1b**).

104

105 To carry out the screening, we pre-determined the minimal lethal concentration of 6TG at 80
106 ng/mL for CRISPRi and CRISPRn iPSC lines (see **Methods** for more details), and applied
107 three different 6TG concentrations (1X: 80 ng/mL, 2X: 160 ng/mL, 3X: 240 ng/mL) in both
108 CRISPRi and CRISPRn screens. We extracted and sequenced DNA samples from the survival
109 cells seven days after 6TG treatment to determine enriched sgRNAs by comparing the results
110 to that of the control cells taken after sgRNA library infection before the 6TG treatment (**Fig.**
111 **1b**). To avoid confounding signals generated by off-target effects of low-quality sgRNAs¹⁷, we
112 only used sgRNAs with high specificity (defined as specificity score > 0.2¹⁸, and without any
113 off-target sites with sequence similarity of ≤2 mismatches) for data analysis. This led to the use
114 of a total of 12,702 high-quality sgRNAs with an average of 38 sgRNAs per ATAC-seq peak for
115 analysis (**Supplementary Fig. 1e, f, Supplementary Fig. 2c, d**). We performed each screen
116 in two biological replicates with each pair of replicates exhibiting high reproducibility
117 (**Supplementary Fig. 2e**). We compared the abundance of each sgRNA between the 6TG
118 treated population and the control population using a negative binomial model, and computed
119 the fold change and *P* value to quantify the effect size and the significance of enrichment of
120 each sgRNA. We used the 5% percentile of the *P* values from non-targeting control sgRNAs
121 as the empirical significance threshold to achieve a false discovery rate of 5%. sgRNAs with *P*
122 value less than the empirical significance threshold and with fold change > 2 were defined as
123 enriched. (**Supplementary Fig. 3**). As expected, sgRNA targeting TSS and coding region

124 were identified as positive hits from both CRISPRi and CRISPRn screens exhibiting greater
125 fold change in CRISPRn screens compared to the CRISPRi screens (**Fig. 1e, Supplementary**
126 **Fig. 4a**). We also observed enrichment of sgRNA bias towards coding regions over TSS
127 regions for CRISPRn screen (**Supplementary Fig. 4b**). These results are consistent with
128 CRISPRi functioning best near TSS by inhibiting transcription, and CRISPRn can disrupt gene
129 function by generating indels downstream of TSS^{19, 20}.

130

131 Further sgRNA fold-change ranking analysis revealed strong positive correlation between the
132 screens with 2X and 3X 6TG treatment for both CRISPRi and CRISPRn screen (Spearman
133 correlation, CRISPRi = 0.97, CRISPRn = 0.84) (**Fig. 1c**) with the correlations for proximal
134 sgRNAs being higher than for distal sgRNAs (**Supplementary Fig. 4c**). On the contrary,
135 results from the 1X screen correlated poorly with either 2X or 3X screens (**Fig. 1c**), suggesting
136 more substantial selection pressure (2X and 3X) can reduce background noise in CRISPRpath
137 screens. Thus, we used sgRNAs enriched from 2X and 3X screens data for identifying active
138 enhancers in the following section (**Fig. 1d**).

139

140 **CRISPRi is more efficient than CRISPRn in pooled CRISPR screens of CREs.**

141 Performing CRISPRpath with CRISPRi and CRISPRn in the same genetic background with an
142 identical sgRNA library offers a unique opportunity for comparing the efficacies of CRISPRi
143 and CRISPRn in pooled CRISPR screens of CREs. We noticed that CRISPRn screens
144 recovered fewer enriched distal sgRNAs than CRISPRi screens (**Fig. 1f**). This is possibly due
145 to the fact that CRISPRi-mediated heterochromatin formation can more effectively perturb
146 CREs compared to CRISPRn-mediated genetics perturbations. We then called an candidate
147 element as an enhancer if there are at least 3 enriched sgRNAs in that CRE. Based on this
148 criterion, we identified 62 and 33 enhancers from the 2X and 3X CRISPRi screen, respectively,
149 and 19 enhancers from the 2X CRISPRn screen. (**Fig. 1g, Supplementary Table 3**).
150 However, no enhancer was identified from the 3X CRISPRn screen, indicating either the
151 CRISPRn induced mutations did not lead to any strong effect on gene expression to make the
152 cells survive the 3X 6TG treatment or there are insufficient numbers of sgRNAs exhibiting
153 deleterious effects on the tested DNA elements to satisfy our criterion of calling functional
154 enhancers. In total, 66 unique enhancers were identified for the six target genes with

155 CRISPRpath under different 6TG treatments (**Fig. 1g**). Together, we demonstrate
156 CRISPRpath can simultaneously identify enhancers for multiple target genes with CRISPRi
157 outperforming CRISPRn. For the following analysis, we focused on the 63 enhancers identified
158 from the 2X and 3X CRISPRi screens (**Fig. 1g**).

159

160 **Genomic feature of CRISPRpath identified enhancers.**

161 To determine the genomic feature of the enhancers, we plotted all the tested elements by their
162 genomic locations and enrichment scores (average of \log_2 (fold change) of enriched sgRNAs of
163 each element) (**Fig. 2a**). Not surprisingly, our data suggest that each gene can be regulated by
164 multiple enhancers with the identified functional enhancers having no position bias relative to
165 the TSS. The average distance between an enhancer and its paired TSS is about 530 Kb (**Fig.**
166 **2b**) with an average of 10 interval genes between an identified enhancer and its target gene
167 pairs (**Fig. 2c**). Interestingly, we observed a weak negative correlation between the enhancer
168 enrichment score and the distance between an enhancer and its paired TSS (**Fig. 2d**, Pearson
169 correlation, $\rho = -0.36$, $P = 0.01$), suggesting enhancers near to TSS tend to have higher
170 regulatory activity compared to enhancers further away from their target genes. It is worth
171 noting, the relative positions for the enriched sgRNAs exhibited no preference relative to
172 ATAC-seq peaks (**Fig. 2e, Supplementary Fig. 5a**) and no preference for the strand on which
173 the sgRNAs were designed (**Supplementary Fig. 5b**), consistent with our knowledge that
174 CRISPRi mediated heterochromatin spreads over hundreds of base pairs in distance²¹.

175

176 Previous studies have revealed that promoters can function as enhancers^{7, 22}. Indeed, 60% (38
177 out of 63) of the functional enhancers identified in CRISPRi screens overlapped with annotated
178 promoters, providing an excellent opportunity to further explore the genomic features of these
179 enhancer-like promoters. To validate whether these promoters function as bona fide
180 enhancers, we targeted three enhancer-like promoters with CRISPRi. We confirmed significant
181 downregulation of their target genes including *MSH6*, *MSH2*, and *PCNA* (**Fig. 3a-c**). In
182 contrast, shRNAs against the transcripts from these promoters (*SOC5*, *FOXN2*, and
183 *TMEM230*) only led to a significant downregulation of its own transcripts and did not affect their
184 target gene expression (**Fig. 3a-c**). These results confirm that these promoter sequences
185 identified by CRISPRpath can function as enhancers. Although it has been shown that

186 enhancer-like protomors are enriched with active chromatin marks and physically close to
187 target genes⁷, it is not clear whether enhancer-like promoters have unique genome features
188 that can differentiate them from other regular active promoters. To this end, we compared
189 chromatin accessibility, occupancy of histone 3 lysine 4 trimethylation (H3K4me3), histone 3
190 lysine 27 acetylation (H3K27ac) and CTCF, transcription, and chromatin interactivity levels
191 between enhancer-like promoters and all other active promoters that did not show enhancer
192 activity in our CRISPRi screens. We show that enhancer-like promoters exhibit higher
193 chromatin accessibility, higher level of transcription, stronger H3K4me3 and H3K27ac signals
194 than those at other active promoters (**Fig. 3d**). On the other hand, we did not observe a
195 significant difference for CTCF binding signals between enhancer-like promoters and control
196 promoters (**Fig. 3d**). Furthermore, by evaluating chromatin interaction data using H3K4me3
197 Proximity Ligation-Assisted ChIP-seq (PLAC-seq), we show enhancer-like promoters have
198 significantly more and stronger interactions compared to control promoters (**Fig. 3e**).

199

200 **CRISPRpath is capable of distinguishing enhancers with distinct effect sizes.**

201 Gene expression often is a result of combinatorial regulatory effects from multiple *cis*-
202 regulatory elements^{11, 23}. Understanding how individual enhancers contribute to gene
203 expression in a quantitative manner is an important first step in dissecting how enhancers
204 orchestrate precise transcriptional control. We seek a new strategy to differentiate enhancers
205 based on their effect sizes on gene expression using CRISPRpath. We hypothesized that cells
206 with drastic down-regulation of MMR genes have a fitness advantage under higher 6TG
207 concentration than cells with modest down-regulation of MMR genes. Consistent with this
208 hypothesis, proximal sgRNAs exhibit larger fold changes than distal sgRNAs (**Fig. 1e**) because
209 perturbing proximal regions has more profound effects on gene down-regulation than
210 perturbing distal regulatory regions. Based on these observations, we hypothesize that
211 enhancers identified under different selection pressure represent distinct regulatory strengths
212 on transcriptional activation. We noticed that enhancers identified under strong selection
213 pressure (3X) have higher enrichment scores compared to enhancers uniquely identified under
214 weak selection pressure (2X), with the TSS regions manifesting the highest enrichment scores
215 (**Fig. 4a**). Thus, enhancers identified in the 3X screen are strong enhancers (n=33), while
216 enhancers uniquely identified in the 2X screen are weak enhancers (n=30) (**Fig. 1g**). To

217 confirm the quantitative effect of enhancers on target gene expression, we tested 11 strong
218 and 10 weak enhancers using CRISPRi followed by RT-qPCR measurement of the
219 corresponding target gene expression (**Fig. 4b, Supplementary Fig. 6a**). We show
220 perturbations of strong enhancers led to significantly more down-regulation of target gene
221 expression (mean down-regulation of target gene by 21%) than perturbations of weak
222 enhancers (mean down-regulation of target gene by 6%), with the perturbations of TSS
223 regions achieving the strongest down-regulation of target genes, by an average of 68%
224 reduction in gene expression (**Fig. 4b, Supplementary Fig. 6a, b**). These quantitative effects
225 on target gene expression are consistent with the enrichment scores from our CRISPRpath
226 screens (**Supplementary Fig. 6c**) and demonstrate the capacity of distinguishing enhancers
227 with different effect sizes by imposing different levels of selection pressures.

228

229 We further explored chromatin features of strong and weak enhancers by analyzing chromatin
230 accessibility, H3K4me3, H3K27ac, and CTCF binding signals in these regions. At individual
231 chromatin mark level, while CRISPRpath-identified enhancers were more accessible, and
232 enriched with active chromatin marks, such as H3K4me3 and H3K27ac, and CTCF binding
233 compared to negative elements or random elements (**Fig. 4c**), we did not observe significant
234 differences between strong and weak enhancers in the chromatin features we individually
235 examined. However, strong enhancers tend to have more active chromatin signatures than
236 weak enhancers (**Fig. 4d**), suggesting combined signatures of active chromatin can be a better
237 indicator of enhancer strength. Strong enhancers tend to have higher distance normalized
238 PLAC-seq contact frequencies with their target promoters than weak enhancers, though not
239 statistically significant, possibly due to the small sample size in this study (**Fig. 4e**). We
240 obtained similar results by expanding this analysis for characterized enhancers in K562 cells
241 and mouse embryonic stem cells^{6, 8} (**Supplementary Fig. 7**), which reinforces the idea that
242 enhancers with larger effects on gene expression tend to have higher chromatin interactions
243 with their cognate promoters. To explore the possible mechanisms that drive enhancer
244 activities in a quantitative manner, we evaluated potential transcription factors (TFs) binding
245 motifs in strong and weak enhancer sequences. Both strong and weak enhancers are enriched
246 with CTCF binding motif (**Fig. 4f**). Indeed, most of strong and weak enhancers are bound by
247 CTCF (**Fig. 4d**), consistent with the notion that CTCF-mediated chromatin loops are essential

248 for gene activation²⁴. Furthermore, strong enhancers and weak enhancers have differential
249 enrichment with TFs binding motifs. For example, the binding motifs for SP/KLF family²⁵ and
250 E2F family^{26, 27} appear more frequently in strong enhancers compared to weak enhancers,
251 suggesting these strong enhancers could be major docking sites for master regulators in
252 iPSCs (**Fig. 4f**).

253

254 **Discussion**

255 CRISPR-mediated high throughput screening using bulk cells allows for the functional
256 characterization of regulatory elements in their native genomic context. However, current
257 approaches are limited to validating a small number of regulatory elements for a single gene^{5, 7,}
258 ^{9, 12, 28, 29}. To overcome this bottleneck, we developed CRISPRpath, a strategy for functional
259 characterization of enhancers for multiple genes simultaneously by leveraging the genes
260 involved in the same biological pathway so that the effects can be measured via a defined
261 phenotype. For example, alpha-toxin resistance phenotype can be used to identify CREs for
262 17 genes in glycosylphosphatidylinositol (GPI)-anchor synthesis pathway³⁰. CRISPRpath can
263 also be leveraged to identify CREs for protein folding regulators that contribute to the
264 endoplasmic reticulum stress-response pathway³¹ using UPRE reporter in mammalian cells.
265 Since CRISPR screen technology is widely used, the CRISPRpath strategy is readily
266 applicable to simultaneously identifying enhancers for genes converging in the defined
267 biological processes and pathways across different cell types. Compared to the existing pooled
268 CRISPR screens of CREs^{5, 7, 8, 10-12, 28, 29, 32}, CRISPRpath is scalable with additional benefits of
269 connecting DNA elements to cellular function, beyond the most standard molecular phenotype
270 of gene expression.

271

272 Promoters can function as enhancers more widespread than expected, with more than half of
273 the enhancers identified for MMR genes in our study being previously annotated promoters.
274 This is consistent with previous reports that enhancer-like promoters are more prevalent for
275 ubiquitously expressed genes³³. Enhancer-like promoters are more accessible compared to
276 other promoters, possibly because these regions are required to be more open to
277 accommodate additional transcriptional machinery such as TF for activating target gene
278 expression besides their own transcription³⁴. Enhancer-like promoters also exhibit significantly

279 higher levels of chromatin interactions with distal regions compared to other active promoters.
280 This observation can be explained by the fact that enhancer-like promoters will not only form
281 chromatin loops with their distal target genes, but also with CREs for controlling the expression
282 of their own genes.

283

284 Genomic studies of chromatin marks have revealed hundreds of thousands candidate CREs in
285 the human genome but with very little quantitative information regarding how CREs contribute
286 to gene regulation^{35, 36}. Using CRISPRpath, we can systematically classify enhancers based
287 on their effect sizes on transcription. Identifying and characterizing the effect size for each
288 individual enhancer is the critical first step to future studies of their combinatory effects on
289 target gene expression. Interestingly, strong and weak enhancers can not be distinguished by
290 individual epigenetic marks we examined. One possible explanation for this observation is that
291 chromatin features only mark enhancer's identity but do not quantify enhancer activity. On the
292 other hand, the strong and weak enhancers we identified may regulate other genes differently
293 from regulating the MMR gene. Interestingly, strong enhancers tend to harbor more than one
294 active chromatin signature, which indicates that enhancer activities are regulated by multiple
295 epigenetic factors, for example, TF mediated transcriptional regulation. Differential TFs binding
296 motifs observed within strong and weak enhancers suggest that enhancer strength is
297 modulated by TF binding. Future studies that further integrating TF binding datasets with
298 functional data of enhancers will shed light on the molecular mechanisms that drives
299 enhancers' effect sizes on gene regulation.

300

301 **Methods**

302

303 **Cell culture**

304 Doxycycline inducible CRISPRi and CRISPRn WTC11 iPSC lines were purchased from
305 Gladstone Stem Cell Core. Both CRISPRi and CRISPRn WTC11 iPS cells were cultured on
306 Matrigel-coated (Corning, 354277) plates with Essential 8™ Medium (Life Technologies,
307 A1517001). iPSCs were passaged using Accutase (STEMCELL Technologies, 07922) and 10
308 μM ROCK inhibitor Y-27632 (STEMCELL Technologies, 72302). HEK293T cells were cultured
309 in Dulbecco's modified Eagle's medium (Gibco, 11995065) with 10% fetal bovine serum

310 (CPSSerum, FBS-500). HEK293T cells were passaged with Trypsin-EDTA (Gibco, 25200072).
311 All the cells were grown with 5% CO₂ at 37°C and verified mycoplasma free using the
312 MycoAlert Mycoplasma Detection Kit (Lonza, LT07-218).

313

314 **sgRNA library design**

315 CRISPRpath sgRNA library was designed to screen cis-regulatory elements for *HPRT1*,
316 *MSH2*, *MSH6*, *MLH1*, *PMS2* and *PCNA*. ATAC-seq peaks within the region of 1 Mb upstream
317 and 1 Mb downstream of each target gene including TSS and coding regions were selected as
318 targeting regions for the sgRNA library design (**Supplementary Table 1**). We generated a
319 genome-wide sgRNA database containing all the available unique sgRNAs, each followed by a
320 'NGG' PAM sequence. All the designed unique sgRNAs in the target regions were added in
321 the sgRNA library, excluding sgRNAs containing AATAAA, AAAAA, TTTTT or TTTTTT
322 sequences. Unique 20-bp sequences in the target regions that were not followed by the 'NGG'
323 or 'NAG' PAM sequences were taken as non-targeting control sgRNAs, excluding non-
324 targeting sgRNAs containing TTT, TTNTT, AATAAA, AAAAA, TTTTT or TTTTTT sequences.
325 Then, a guanine nucleotide was added to all the sgRNAs if the sequence did not start with G to
326 increase efficiency of transcription from U6 promoter. Final sgRNA oligos adhered to the
327 following template: 5'-ATATCTTGTGGAAAGGACGAAACACC-[20- or 21-bp sgRNA
328 sequence]-GTTTTAGAGCTAGAAATAGCAAGTTAAAATAAGGC-3'. In total, 35,763 sgRNAs
329 were included in the library (**Supplementary Fig. 1 and Supplementary Table 2**). We
330 retrieved specificity score and off-target site for each sgRNA from GuideScan¹⁸
331 (www.guidescan.com) and assigned the specificity score of sgRNAs not existed in the
332 GuideScan database to 0. The high-quality sgRNAs were filtered with specificity score >0.2
333 and without perfectly matched or 1-2 mismatches off-target sites.

334

335 **Oligo synthesis and library cloning**

336 sgRNA library oligos were synthesized by TWIST BIOSCIENCE and amplified with the forward
337 primer 5'- TCGATTTCTTGGCTTTATATATCTTGTGGAAAGGACGAAACAC-3' and the
338 reverse primer
339 5'-AACGGACTAGCCTTATTTTAACTTGCTATTTCTAGCTCTAAAAC-3'. We replaced the
340 Cas9 sequence in lentiCRISPR v2 plasmid (Addgene, 52961) with blasticidin S deaminase

341 sequence to construct the lentiCRISPR-v2-Blast-Puro plasmid (Addgene, 167186). The PCR
342 products were purified via gel excision and column purification (Promega, A9282), and then
343 inserted into the BsmBI-digested lentiCRISPR-v2-Blast-Puro vector by Gibson assembly (New
344 England Biolabs, E2621L). The assembled products were transformed into NEB 5-alpha
345 electrocompetent *E. coli* cells (New England Biolabs, C2989K) by electroporation. About 40
346 million independent bacterial colonies were cultured, and sgRNA library plasmids were
347 extracted with the Qiagen EndoFree Plasmid Mega Kit (Qiagen, 12381). The recovery rate and
348 distribution of the sgRNA library were checked with next generation sequencing
349 **(Supplementary Fig. 2a-d)**.

350

351 **Lentivirus production and titration**

352 To make the lentiviral library, 5 µg of sgRNA plasmid library was co-transfected with 3 µg of
353 psPAX (Addgene, 12260) and 1 µg of pMD2.G (Addgene, 12259) lentivirus packaging
354 plasmids into 8 million HEK293T cells in a 10-cm dish with PolyJet (SigmaGen Laboratories,
355 SL100688). For each individual sgRNA, 3.75 µg of sgRNA plasmid was co-transfected with
356 2.25 µg of psPAX (Addgene, 12260) and 0.75 µg of pMD2.G (Addgene, 12259) plasmids into
357 4 million HEK293T cells in a T25 flask with PolyJet (SigmaGen Laboratories, SL100688).
358 Media was replaced 12 h after transfection, and harvested every 24 h for a total of three
359 harvests. Harvested media containing the desired virus were filtered through Millex-HV 0.45-
360 µm PVDF filters (Millipore, SLHV033RS) and further concentrated with 100,000 NMWL Amicon
361 Ultra-15 centrifugal filter units (Amicon, UFC910008).

362

363 The titer of lentivirus was determined by transducing 500,000 cells with varying amount (0, 0.5,
364 1.0, 2.0, 4.0 and 8.0 µL) of concentrated virus and polybrene (Millipore, TR-1003-G, 8 µg/mL).
365 Viral transduction was performed by centrifuging the lentivirus and cell combination at 1000
366 RCF for 90 min at 37°C. 3 to 4 h later, virus containing media was replaced with fresh media.
367 24 h after the transduction, transduced cells were dissociated with Accutase, and seeded as
368 duplicates. One replicate was treated with blasticidin (Gibco, A1113903, 4 µg/mL), and the
369 other replicate was not treated with blasticidin. Four days later, the blasticidin resistant cells
370 and control cells were counted to calculate the ratio of infected cells and the viral titer.

371

372 **Determining 6TG concentration via killing curve titration**

373 Both CRISPRi and CRISPRn WTC11 iPSCs were used to determine the minimal lethal
374 concentration of 6TG. Cells were seeded in 24 well plates. When the cells reached around
375 50% confluence (Day 0), they were treated with 6TG concentrations of 0 (control), 20, 40, 60,
376 80, 100, 120, 140, and 160 ng/mL. Two wells were allocated for each condition. The cells were
377 examined daily and cultured for 7 days. The media was replaced daily with the specified 6TG
378 concentration. After 3 days, wells with 6TG concentration greater than or equal to 100 ng/mL
379 had no surviving cells. On Day 4 of treatment, the wells with 80 ng/mL 6TG treatment had no
380 surviving cells. On the last day of treatment, the wells with 40 and 60 ng/mL treatments had
381 very few surviving cells, while the 20 ng/mL treatment had many surviving cells. Based on
382 these results, we set 80 ng/mL as the minimal lethal concentration for 6TG.

383

384 **CRISPRpath screening and sequencing library preparation**

385 CRISPRpath screens were carried out with 72 million doxycycline inducible CRISPRi or
386 CRISPRn iPSCs in biological replicate. The cells for lentiviral transduction were seeded into 6
387 well plates with 1 million cells per well, and the lentiviral library (MOI = 0.5) was transduced
388 into the iPSCs with 8 µg/mL Polybrene (Millipore, TR-1003-G) and spun at 1000 RCF at 37°C
389 for 90 min. The transduced cells were treated with doxycycline (Sigma, D9891, 2 µM) and
390 blasticidin (Gibco, A1113903, 4 µg/mL) for 4 days. After this doxycycline and blasticidin
391 treatment, 10 million cells were reserved as a control population, and 100 million cells were
392 used for CRISPRpath screen with doxycycline and 6TG (Sigma, A4660) treatment for 7 days.
393 Finally, survival cells were collected from the 6TG treated population.

394

395 The genomic DNA was extracted from each sample via cell lysis and digestion (100 mM Tris-
396 HCl pH8.5, 5 mM EDTA, 200 mM NaCl, 0.2% SDS, 100 µg/mL proteinase k),
397 phenol:chloroform (Thermo Scientific, 17908) extraction and isopropanol (Fisher Scientific,
398 BP2618500) precipitation. To amplify the sgRNA sequences from each sample, thirty-two 50 µl
399 PCR reactions were performed using 500 ng genomic DNA for each reaction and NEBNext®
400 High-Fidelity 2X PCR Master Mix (New England Biolabs, M0541S). The purified libraries were
401 sequenced on the NovaSeq 6000 with 150-bp paired-end sequencing. The detail protocol is

402 available on ENCODE portal (<https://www.encodeproject.org/documents/2e6451a9-3b98-4d95-922e-a3d8d2100ddf/>).

404

405 **CRISPRpath data analysis**

406 The sequence files were down-sampled to the same amount of total reads, and then mapped
407 to the sgRNA library with the requirement of exact match of designed sgRNA sequences in the
408 following pattern 5'-CCG-[N19 or N20]-GTT-3'. Only the highly specific sgRNAs (specificity
409 score >0.2, without perfectly matched or 1-2 mismatches off-target sites) were used for
410 downstream data analysis. The sgRNA enrichment for each screen was calculated by
411 comparing 6TG treated samples with the associated control samples with edgeR and TMM
412 normalization. We first used edgeR³⁷ to calculate P value based on negative binomial model
413 for both targeting sgRNAs and non-targeting control sgRNAs. To achieve empirical false
414 discovery rate less than 5%, we then selected a P value cutoff corresponding to the 5%
415 percentile of P values from non-targeting control sgRNAs. Finally, we defined enriched
416 sgRNAs with P value less than the selected P value cutoff, and fold change >2. The ATAC-seq
417 peaks were identified as functional enhancers for the six MMR genes by having at least 3
418 significant enriched sgRNAs. Analysis scripts are available at
419 <https://github.com/MichaelMW/crispy>.

420

421 **Analysis of genomic feature and chromatin signature of identified enhancers**

422 Genomic distances between enhancer and TSS pairs were calculated based on the distance
423 from the center of enhancers to the transcription start sites of the target genes. The number of
424 interval genes is the number of all the RefSeq annotated genes between each enhancer and
425 paired target gene. The signal of chromatin signatures, including ATAC-seq, H3K27ac,
426 H3K4me3, CTCF binding and RNA-seq, were calculated by deeptools (v3.4.3)³⁸. The
427 enhancer-like promoters are the enhancers overlap with the region 500 bp upstream and
428 downstream of a RefSeq annotated TSS.

429

430 **Validation of identified enhancers using CRISPRi**

431 We cloned lentiCRISPR-v2-HygR-EGFP (Addgene, 167188) and lentiCRISPR-v2-HygR-
432 mCherry (Addgene, 167189) vectors by replacing the Cas9 and puromycin N-acetyltransferase

433 sequences in lentiCRISPR v2 plasmid (Addgene, 52961) with hygromycin B
434 phosphotransferase and EGFP or mCherry sequences. To validate the identified enhancers,
435 individual sgRNAs targeting identified enhancers were cloned into the lentiCRISPR-v2-HygR-
436 GFP or lentiCRISPR-v2-HygR-mCherry vector. The doxycycline inducible CRISPRi WTC11
437 iPSCs were infected with the lentivirus expressing sgRNAs for three replicates per sgRNA.
438 The sgRNA infected cells were grown with hygromycin (Gibco, 10687010, 150 µg/mL) and
439 doxycycline (Sigma, D9891, 2 µM) containing media. Seven days later, the cells were collected
440 and total RNA was extracted from the cells using the Qiagen RNeasy® Plus Kit (Qiagen,
441 74134). One µg of RNA was then used to synthesize cDNA using the Bio-RAD iScript cDNA
442 Synthesis Kit (Bio-RAD, 1708840). qPCR reactions for targeted genes were performed with
443 the Luminaris HiGreen qPCR Master Mix (Thermo Scientific, K0993) on the Roche LightCycler
444 96 System. The qPCR primers are listed in **Supplementary Table 4** and the sgRNA
445 sequences are listed in **Supplementary Table 6**. For each tested element in **Fig. 3a-c** and
446 **Supplementary Fig. 6b**, we performed CRISPRi experiments with two independent sgRNAs
447 and used the results from the sgRNA with stronger transcriptional repression in **Fig. 4b**.

448

449 **shRNA mediated RNA interference**

450 shRNAs were designed by using DSIR tool (<http://biodev.extra.cea.fr/DSIR/DSIR.html>)
451 targeting *SOCS5*, *FOXN2* and *TMEM230*. The sequences of shRNAs are listed in
452 **Supplementary Table 5**. The shRNAs were cloned into lentiCRISPR-v2-HygR-mCherry
453 vector under the control of human U6 promoter and packaged into lentivirus for cell
454 transduction. The WTC11 iPSCs transduced with shRNA lentivirus were treated with
455 hygromycin (Gibco, 10687010, 150 µg/mL) for 7 days and then collected for RNA extraction
456 and qPCR.

457

458 **ATAC-seq**

459 ATAC-seq was carried out using the Nextera DNA Library Prep Kit (Illumina, FC-121-1030) as
460 previously described³⁹. The detailed protocol is available on the ENCODE portal
461 (<https://www.encodeproject.org/documents/0317894c-5a42-4f03-b865-c2a2d08708ef/>). Briefly,
462 each library started with 100,000 fresh iPSCs, and the cells were incubated with ice cold nuclei
463 extraction buffer (10 mM TrisHCl pH 7.5, 10 mM NaCl, 3 mM MgCl₂, 0.1% Igepal CA630, and

464 1x protease inhibitor) for 5 min on ice, then centrifuged at 500 RCF for 5 min. 50,000 resulting
465 nuclei were treated with tagmentation buffer (25 μ L Buffer TD with 50,000 nuclei, 22.5 μ L
466 water, 2.5 μ L TDE1) for 30 min at 37°C. The transposed DNA was purified using Qiagen
467 MinElute PCR purification kit (Qiagen, 28006), amplified using Nextera primers, then size-
468 selected for fragments between 150 and 1000 bp using SPRIselect beads (Beckman Coulter,
469 B23319). Libraries were sent for single-end sequencing on the HiSeq 4000 (50 bp single-end
470 reads). Reads were mapped to hg38/GRCh38 and processed using the ENCODE pipeline
471 (https://github.com/kundajelab/atac_dnase_pipelines, V1.8.0), which ran on the default
472 settings. The ATAC-seq peaks were filtered with FDR cutoff of 0.1%, and adjacent peaks were
473 merged if they are less than 1 kb apart.

474

475 **RNA-seq**

476 RNA was extracted from fresh cells using the RNeasy Plus Mini Kit (Qiagen, 74134).
477 Approximately 1000 ng of extracted RNA was used to prepare libraries for sequencing using
478 the TruSeq Stranded mRNA Library Prep Kit (Illumina, 20020594). Libraries were sent for
479 paired-end sequencing on the NovaSeq 6000 (100 bp paired-end reads). Reads were aligned
480 to hg38/GRCh38 using STAR 2.7.0f⁴⁰ with the standard ENCODE settings, and transcript
481 quantification was performed in a strand-specific manner using RSEM 1.3.1⁴¹ with the
482 annotation from GENCODE v32. Only the first read was used, and all reads were trimmed to
483 51bp using TrimGalore 0.4.5 running the following options: -q 20 --length 20 -- stringency 3 --
484 trim-n. The edgeR package in R (3.20.9)³⁷ was used to calculate TMM-normalized FPKM
485 values for each gene based on the expected counts and gene lengths for each library. The
486 mean gene expression across all replicates was used for analysis.

487

488 **ChIP-seq**

489 ChIP-seq libraries were constructed from 2 million WTC11 iPSCs. Cells were crosslinked in
490 1% formaldehyde at room temperature for 20 min and then quenched with 2.5 M glycine at
491 room temperature for 5 min. Fixed cells were lysed and chromatin was sonicated by Covaris
492 with the following parameters: Duty Factor: 2%, Peak Incident Power: 105W, Cycles per Burst:
493 200, for 30 min. Input chromatin was removed and stored at -20°C for later processing.
494 Magnetic beads (Invitrogen, Dynabeads Protein A, 10001D) were preincubated with H3K27ac

495 antibody (Active Motif, 39133, Lot 22618011) for 2 hours at 4°C before being added to sheared
496 chromatin. Samples were incubated overnight at 4°C. Beads were washed 3 times and
497 chromatin was then eluted. Samples were incubated at 65°C overnight to reverse the
498 crosslinking. DNA was treated with RNase A for 1 hr at 37°C and Proteinase K (New England
499 Biolabs, 8107) for 1 hr at 55°C. DNA was purified by phenol-chloroform extraction and ethanol
500 precipitation. Libraries were prepared using Tru-seq adapters and size-selected using
501 SPRIselect beads prior to amplification and paired-end sequencing. Libraries were sent for
502 paired-end sequencing on the NovaSeq 6000 (150 bp paired-end reads). Sequencing reads
503 were trimmed to 50 bp and mapped to hg38 using bowtie2 with the following options: --local --
504 very-sensitive-local --no-unal --no-mixed --no-discordant --phred33 -l 10 -X 700. Picard Tools
505 was used to remove blacklisted regions and duplicate reads and MACS2 was used to call
506 peaks on merged replicates at an FDR cutoff of 1%.

507

508 **CUT&Tag**

509 CUT&Tag libraries were constructed from 150,000 WTC11 iPSCs according to previously
510 described methods⁴². Cells were lysed in nuclei extraction buffer (20 mM HEPES-KOH pH 7.9,
511 10 mM MgCl₂, 0.1% Triton X-100, 20% glycerol and 1x protease inhibitor) on ice for 10 min.
512 The samples were spun and resuspended in 100 µl nuclei extraction buffer. Meanwhile, 10 µl
513 of BioMag Plus Concanavalin A (Bangs Laboratories, BP531) were equilibrated in binding
514 buffer (1x PBS, 1 mM CaCl₂, 1 mM MgCl₂ and 1 mM MnCl₂). The equilibrated beads were
515 added to the samples and incubated with rotation for 15 min at 4°C. Nuclei-bound beads were
516 washed with Buffer 1 (20 mM HEPES-KOH pH 7.9, 150 mM NaCl, 2 mM EDTA, 0.5 mM
517 spermidine, 0.1% BSA and 1x protease inhibitor) and Buffer 2 (20 mM HEPES-KOH pH 7.9,
518 150 mM NaCl, 0.5 mM spermidine, 0.1% BSA and 1x protease inhibitor). After washing nuclei-
519 bound beads were resuspended in 50 µl Buffer 2 with 0.5 µl antibody (H3K4me3 from
520 Millipore, 04-745, Lot 3543820 and CTCF from Millipore, 07-729, Lot 3059608) and incubated
521 with rotation overnight at 4°C. Samples were washed twice with Buffer 2 and resuspended in
522 50 µl Buffer 2 with antibody (antibodies-online Inc., Guinea Pig anti-Rabbit IgG, ABIN101961,
523 Lot 42323) and incubated for 1 hr at room temperature with rotation. Samples were washed
524 again with Buffer 2 and resuspended in 100 µl Buffer 3 (20 mM HEPES-KOH pH 7.9, 300 mM
525 NaCl, 0.5 mM Spermidine, 0.1% BSA and 1x proteinase inhibitor) containing 0.04 µM pA-Tn5.

526 Samples were incubated for 1 hr at room temperature, washed three times with Buffer 3 and
527 resuspended in tagmentation buffer (20 mM HEPES-KOH pH 7.9, 300 mM NaCl, 0.5 mM
528 Spermidine, 10 mM MgCl₂, 0.1% BSA and 1x proteinase inhibitor). Samples were incubated
529 for 1 hr at 37°C. Samples were treated with Proteinase K (New England Biolabs, 8107) for 1 hr
530 at 50°C. DNA was purified by phenol-chloroform extraction and ethanol precipitation. Libraries
531 were prepared using Tru-seq adapters and size-selected using SPRIselect beads prior to
532 amplification and paired-end sequencing. Libraries were sent for paired-end sequencing on the
533 Mini-seq (37 bp paired-end reads, H3K4me3 libraries) or NovaSeq 6000 (150 bp paired-end
534 reads, CTCF libraries). Sequencing reads (CTCF libraries were trimmed to 50 bp) were
535 mapped to hg38 using bowtie2 with the following options: --local --very-sensitive-local --no-
536 unal --no-mixed --no-discordant --phred33 -l 10 -X 700. Picard Tools was used to remove
537 blacklisted regions and duplicate reads and SEACR⁴³ was used to call peaks on merged
538 replicates.

539

540 **H3K4me3 PLAC-seq**

541 H3K4me3 PLAC-seq data in WTC11 cells were generated as previously described⁴⁴ in
542 biological replicates (clone 6 and clone 28) ([https://data.4dnucleome.org/experiment-set-
543 replicates/4DNESDRL4ZKM/](https://data.4dnucleome.org/experiment-set-replicates/4DNESDRL4ZKM/) and [https://data.4dnucleome.org/experiment-set-
544 replicates/4DNESIZ5TTHO/](https://data.4dnucleome.org/experiment-set-replicates/4DNESIZ5TTHO/)). We combined the two biological replicates, and applied the
545 MAPS pipeline⁴⁵ to identify significant long-range chromatin interactions at 5 kb bin resolution
546 for the genomic distance 10 kb ~ 1 Mb. The reference genome is GRCh38. In addition, for
547 each 5 kb bin pair anchored at H3K4me3 peaks, the MAPS pipeline outputs the normalized
548 contact frequency, which adjusts for the biases from effective fragment length, GC content,
549 sequence mappability, H3K4me3 enrichment level and 1D genomic distance effect.

550

551 **Comparison between strong enhancers and weak enhancers using H3K4me3 PLAC-seq 552 data**

553 For 4 genes *HPRT1*, *MLH1*, *PMS2* and *PCNA*, there are 23 enhancer-promoter pairs between
554 strong enhancers and their target genes, and 21 enhancer-promoter pairs between weak
555 enhancers and their target genes. We mapped each enhancer and promoter of target gene
556 into 5 kb bins, and obtained the distance normalized H3K4me3 PLAC-seq contact frequency

557 for 5 kb bin pairs containing the enhancer-promoter pairs. Since *MSH2* and *MSH6* are located
558 within 407 kb linear genomic distance with each other and we can't assign enhancers to either
559 gene reliably, enhancers identified near *MSH2* and *MSH6* were excluded from this analysis.

560

561 **Comparison between enhancer-like promoters and control promoters using H3K4me3** 562 **PLAC-seq data**

563 For this analysis, control promoters are active promoter regions with annotated ATAC-seq
564 peaks and tested negative as enhancers for the for MMR genes. We mapped each promoter
565 into a 5 kb bin that was used in the PLAC-seq analysis. We only choose the bins with one
566 annotated active promoter, which gave us 31 enhancer-like promoters and 43 control
567 promoters in this analysis. We counted the number of significant H3K4me3 PLAC-seq
568 interactions anchored at the 5 kb bins with these promoter sequences. In addition, as
569 described in our previous study²³, for promoters with at least one significant interaction, we
570 calculated the summation of $-\log_{10}$ FDR of significant interactions, which is a measure of the
571 overall interaction strength.

572

573 **Chromatin contact frequency comparison between strong enhancers and weak** 574 **enhancers in K562 cells and mESCs**

575 For the chromatin contact frequency comparison of enhancers in K562 cells and mouse
576 embryonic stem cells (mESCs), we downloaded the identified enhancers from each
577 publication^{6, 8}, and defined strong enhancer with cutoff of $50\% \leq$ transcriptional contribution \leq
578 100% , weak enhancer with cutoff of $0\% <$ transcriptional contribution $\leq 20\%$. H3K27ac HiChIP
579 data in K562 cells⁴⁶ and H3K4me3 PLAC-seq data in mESCs⁴⁵ were used for comparison. The
580 comparisons were performed in 10 kb resolution.

581

582 **Motif scan and Transcription factor identification**

583 The fasta files were first generated in the hg38 genome for the identified strong enhancers and
584 weak enhancers separately. For each strong enhancer and weak enhancer, the FIMO software
585 (version 5.1.0)⁴⁷ with human motif database HOCOMOCO (v11 FULL)⁴⁸ was used to scan the
586 motifs. All the FIMO motif scans were in default settings. We then filtered the transcription
587 factors (TFs) in each strong and weak enhancer loci by FDR cutoff of 0.05 and p-value cutoff

588 of 0.0001 and gene expressions cutoff of FPKM > 1. By taking the TFs with TF motif appeared
589 in more than 80% enhancers, 47 TFs were considered as commonly appearing in the strong
590 enhancers, and 35 TFs were in the weak enhancers.

591

592 **Acknowledgements**

593 This work was supported by the National Institutes of Health grants UM1HG009402 (to Y.S.,
594 B.R., and W.W.) and U54DK107977 (to B.R. and M.H.).

595

596 **Data availability**

597 The CRISPRpath screen datasets used in this study are available at the ENCODE portal
598 (www.encodeproject.org, accession number: ENCSR617AZY (sgRNA plasmid library),
599 ENCSR427OPP (CRISPRi control), ENCSR900AXT (CRISPRi 1X), ENCSR254SJU (CRISPRi
600 2X), ENCSR793DSE (CRISPRi 3X), ENCSR250ZWC (CRISPRn control), ENCSR117YGQ
601 (CRISPRn 1X), ENCSR071ZGB (CRISPRn 2X), ENCSR482PHH (CRISPRn 3X)). WTC11
602 iPSCs H3K4me3 PLAC-seq datasets are available at the 4DN data portal
603 (data.4dnucleome.org, accession number: 4DNESIZ5TTTHO and 4DNESDRL4ZKM). ATAC-
604 seq, ChIP-seq, CUT&Tag, and RNA-seq datasets in WTC11 iPSCs are available at the Gene
605 Expression Omnibus under the accession number GSE166839. Data can be visualized on the
606 WashU Epigenome Browser using the session at the following link
607 ([https://epigenomegateway.wustl.edu/browser/?genome=hg38&sessionFile=https://shen-
608 xren.s3-us-west-1.amazonaws.com/CRISPRpath/eg-session-QRXJ0218-4d710b60-6ea7-
609 11eb-8d8d-03c7189570c0.json](https://epigenomegateway.wustl.edu/browser/?genome=hg38&sessionFile=https://shen-xren.s3-us-west-1.amazonaws.com/CRISPRpath/eg-session-QRXJ0218-4d710b60-6ea7-11eb-8d8d-03c7189570c0.json)). Tracks include ATAC-seq, H3K27ac, H3K4me3, and CTCF
610 signals, and the identified enhancers from CRISPRi 2X and 3X screens. RefGene 38 genes
611 are also displayed. The plasmids generated in this study are available from Addgene
612 (#167186, #167188, #167189).

613

614 **Code availability**

615 The computer code used for analyzing CRISPRpath datasets is available
616 at <https://github.com/MichaelMW/crispy>.

617

618 **Competing interests**

619 B.R. is co-founder and shareholder of Arima Genomics and Epigenome Technologies. The
620 other authors declare that they have no competing interests.

621

622 **Author contributions**

623 X.R. and Y.S. designed the study. X.R. and B.L. designed the sgRNA library under the
624 supervision of Y.S. and B.R. X.R., K.J., I.R.J., M.A.T., J.L., L.M., and T.W.T. performed the
625 experiments. M.W. designed the CRISPY under the supervision of W.W. X.R., B.L., L.Z., G.L.,
626 and Y.L. performed data analysis. M.Y. and R.H. constructed the H3K4me3 PLAC-seq
627 libraries. L.L., A.A. and M.H. analyzed PLAC-seq and HiChIP data. X.R. and Y.S. prepared the
628 manuscript with input from all other authors.

629 **Figure 1. CRISPRpath for identifying enhancers of multiple genes.**

630 (a) Six genes (*HPRT1*, *MSH2*, *MSH6*, *MLH1*, *PMS2* and *PCNA*) in the 6TG-induced mismatch
631 repair process were used for CRISPRpath screen in this study. (b) Schematic of the
632 CRISPRpath screening strategy with 6TG treatment in iPSCs. Cell survival was used as
633 readout for the screen. (c) Spearman correlation analysis of sgRNA ranking based on fold
634 change for CRISPRpath screens with different 6TG concentrations (1X, 2X and 3X). (d) Venn
635 diagram shows the overlapping enriched sgRNAs identified from the screens with 2X and 3X
636 6TG treatments. (e) Box plots show the fold change of the enriched distal and proximal
637 sgRNAs from 2X, 3X CRISPRi and CRISPRn screens. Star indicates no enriched distal
638 sgRNA was identified from 3X CRISPRn screen. Boxplots indicate the median, interquartile
639 range (IQR), $Q1 - 1.5 \times IQR$ and $Q3 + 1.5 \times IQR$. (f) Bar plot shows the number of enriched
640 distal and proximal sgRNAs from 2X, 3X CRISPRi and CRISPRn screens. Star indicates no
641 enriched distal sgRNA was identified from 3X CRISPRn screen. (g) Venn diagram shows the
642 identified enhancers from each CRISPRpath screen.

643

644 **Figure 2. Genomic features of identified enhancers from CRISPRpath using CRISPRi.**

645 (a) Genomic locations of identified enhancers relative to TSS. Circles indicate enhancers
646 identified from the CRISPRi 3X screen (red), enhancers uniquely identified from the CRISPRi
647 2X screen (blue), tested CREs that are not identified as enhancers (grey). Purple lines label
648 the location of each target gene. (b) Histogram shows the distance distribution between
649 identified enhancers and their paired TSS. (c) Histogram shows the number of interval genes
650 between enhancers and their target gene TSS. Mean is indicated with an orange dashed line
651 and only the enhancers for *MLH1*, *PMS2*, *PCNA*, *HPRT1* are included in b and c. (d) A weak
652 negative correlation is observed between enrichment score and genomic distance between
653 enhancers and their target genes (Pearson correlation, $r = -0.36$, $P = 0.01$). Black circles
654 indicate promoters. The red and blue circles are enhancers showing in a. (e) Density plot
655 shows no significant difference (two-tailed two-sample Kolmogorov–Smirnov test) for the
656 distribution of all distal sgRNAs (gray), enriched distal sgRNAs from 2X (blue) and 3X screen
657 (red) CRISPRi screens.

658

659 **Figure 3. Enhancer-like promoters act as functional enhancers.**

660 (a, b, c) Three examples of promoters function as enhancers. CRISPRi silencing of the
661 promoter region of *SOCS5*, *FOXN2* and *TMEM230* results in significant downregulation of
662 *MSH6*, *MSH2* and *PCNA*, respectively. shRNA knockdown of *SOCS5*, *FOXN2* and *TMEM230*
663 can only downregulate *SOCS5*, *FOXN2* and *TMEM230* expression. Three independent
664 replicates per condition and two independent sgRNAs or shRNAs per replicate were used for
665 each experiment. *P* values are from two-tailed two-sample t-test. (d) Average signal
666 enrichment of ATAC-seq, gene transcription, H3K4me3, H3K27ac and CTCF binding for
667 enhancer-like promoters (n = 38) and control promoters (n = 47). *P* values are from Wilcoxon
668 test. Boxplots indicate the median, IQR, Q1 - 1.5 × IQR and Q3 + 1.5 × IQR. (e) Number of
669 H3K4me3 mediated chromatin interactions and cumulative interaction score for enhancer-like
670 promoters (n = 31) and control promoters (n = 43). Boxplots indicate median, IQR, Q1 - 1.5 ×
671 IQR and Q3 + 1.5 × IQR. *P* values are calculated from Wilcoxon test.

672

673 **Figure 4. CRISPRpath can distinguish weak and strong enhancers by imposing**
674 **different selection pressures.**

675 (a) Box plots show the enrichment score of the tested elements. TSS regions (black circles)
676 show highest enrichment scores. Enhancers uniquely identified from the lower selection
677 pressure (CRISPRi 2X, blue circles) exhibit lower enrichment scores compared to the
678 enhancers identified from the higher selection pressure (CRISPRi 3X, red circles). *P* values
679 are from Wilcoxon test. (b) Box plots show the CRISPRi perturbation at enhancers induced
680 various degrees of transcriptional repression of target genes measured with RT-qPCR. Each
681 dot represents the average value from three biological replicates. CRISPRi targeting TSS
682 regions (dark gray) achieved the highest transcriptional repression. CRISPRi targeting strong
683 enhancers (pink) leads a more substantial transcription silencing on target gene compared to
684 CRISPRi targeting weak enhancers (cyan). *P* values are from Wilcoxon test. (c) Enrichment
685 analysis of ATAC-seq, H3K27ac, H3K4me3, CTCF binding signals for strong (n = 33) and
686 weak (n = 30) enhancers. Boxplots indicate the median, IQR, Q1 - 1.5 × IQR and Q3 + 1.5 ×
687 IQR. *P* values for the difference between strong and weak enhancers are from Wilcoxon test;
688 see Supplementary Table 7 for *P* values of all pairwise comparisons. (d) Intersection of
689 genomic features for weak enhancers (blue bar) and strong enhancers (red bar). (e) Distance
690 normalized H3K4me3 PLAC-seq contact frequency for strong (n = 23) and weak (n = 21)

691 enhancers. Only the enhancers for *MLH1*, *PMS2*, *PCNA*, *HPRT1* are included (see Methods).
692 Boxplots indicate the median, IQR, $Q1 - 1.5 \times IQR$ and $Q3 + 1.5 \times IQR$. *P* value is from
693 Wilcoxon test. (f) Heatmap shows the frequency of transcription factor motifs found in strong
694 and weak enhancers.

695

696 **Supplementary Figure 1. Features of sgRNA library for CRISPRpath screen.**

697 (a) Bar graph shows the number of distal ATAC-seq peaks used as candidate CREs for six
698 target genes. (b) Histogram shows size distribution of distal ATAC-seq peaks. The average
699 size is 488 bp (blue dash line). (c, e) The composition of the sgRNA library. In total, 35,763
700 sgRNAs were included in the library (c), and 12,702 sgRNAs are high quality sgRNAs (e). (d,
701 f) Distribution of the number of sgRNAs per distal ATAC-seq peak. Average numbers of
702 sgRNA per ATAC-seq peak are indicated with blue dash lines.

703

704 **Supplementary Figure 2. Quality of the sgRNA library and CRISPRpath screen libraries.**

705 (a) Distribution of sgRNA oligo read counts in the sgRNA library. (b) Cumulative frequency of
706 sgRNAs in the sgRNA library. (c) Distribution of high quality sgRNAs read counts in the sgRNA
707 library. (d) Cumulative frequency of high quality sgRNAs in the sgRNA library. The constructed
708 sgRNA plasmid library recovered all the designed sgRNAs with the copy number difference
709 less than five fold for at least 97% designed sgRNAs. (e) PCA analysis shows the high
710 reproducibility of the CRISPRpath screen libraries between biological replicates.

711

712 **Supplementary Figure 3. *P* value cutoff used for identifying enriched sgRNAs from each**
713 **screens.**

714 (a-f) Distribution of *P* value for tested distal, proximal and non-targeting control sgRNA groups.
715 Orange dash lines indicate 5% percentile of the *P* values from non-targeting control sgRNAs to
716 achieve a false discovery rate of 5%.

717

718 **Supplementary Figure 4. Enriched proximal sgRNAs and sgRNA ranking analysis.**

719 (a) Number and fold change of the enriched proximal sgRNAs for the six target genes from
720 CRISPRi and CRISPRn screens. The color indicates fold changes, and the size of circle
721 indicates the number of enriched sgRNAs. (b) Enrichment analysis shows the enriched

722 proximal sgRNAs bias towards to the TSS region for CRISPRi screens and the protein coding
723 region (CDS) for CRISPRn screens. Color represents the number of enriched sgRNAs. (c)
724 Spearman correlation analysis of the distal and proximal sgRNAs ranking shows proximal
725 sgRNAs exhibiting higher correlation between each screen compared to distal sgRNAs.

726

727 **Supplementary Figure 5. Enriched sgRNAs identified from CRISPRi screens exhibit no**
728 **position and strand preference.**

729 (a) Enriched sgRNAs from CRISPRi 2X (red dots, n=448) and CRISPRi 3X (red dots, n=260)
730 screens showed similar distributions across candidate CREs. (b) Odds ratio analysis of the
731 fold change of enriched sgRNAs shows enriched sgRNAs have no strand preference.
732 Enhancers with enriched sgRNAs only targeting one strand were exclude for the analysis.
733 Odds ratio was calculated for each element with the equation of $\text{ave}(\log_2(\text{fold change of}$
734 $\text{sgRNA targeting plus strand})) / \text{ave}(\log_2(\text{fold change of sgRNA targeting minus strand}))$. Violin
735 plots show the distributions of odds ratio values within each screen, and boxplots indicate the
736 median, IQR, $Q1 - 1.5 \times \text{IQR}$ and $Q3 + 1.5 \times \text{IQR}$.

737

738 **Supplementary Figure 6. Validation of CIRSPRpath identified enhancers.**

739 (a) Validation of the strong (black) and weak (grey) enhancers with CRISPRi followed by RT-
740 qPCR. Three independent replicates per condition. The significance was calculated with two-
741 tailed two-sample *t*-test. Data are mean and s.d. * $P < 0.05$, ** $P < 0.01$, *** $P < 0.001$. (b)
742 CRISPRi-mediated transcriptional repression of six target genes by targeting TSS of each
743 gene. Three independent replicates per condition. The significance was calculated with two-
744 tailed two-sample *t*-test. Data are mean and s.d. * $P < 0.05$, ** $P < 0.01$, *** $P < 0.001$. (c)
745 Pearson correlation analysis reveals element enrichment score from CRISPRpath screens
746 correlates with element effect size on transcription from CRISPRi (Pearson correlation, $\text{PCC} =$
747 -0.68 , $P = 3.2 \times 10^{-5}$).

748

749 **Supplementary Figure 7. Chromatin contract frequency analysis for the enhancers in**
750 **K562 cells and mESCs.**

751 (a) Distance normalized H3K27ac HiChIP contact frequency for strong (n = 34) and weak (n =
752 82) enhancers identified with crisprQTL mapping in K562 cells. (b) Distance normalized

753 H3K27ac HiChIP contact frequency for strong (n = 2) and weak (n = 20) enhancers identified
754 with CRISPRi-FlowFISH screen in K562 cells. (c) Distance normalized H3K4me3 PLAC-seq
755 contact frequency for strong (n = 10) and weak (n = 3) enhancers identified in mouse
756 embryonic stem cells. Boxplots indicate the median, IQR, $Q1 - 1.5 \times IQR$ and $Q3 + 1.5 \times IQR$.
757 *P* values are from Wilcoxon test.

758

759 **Supplementary Table 1. ATAC-seq peak regions used to design the sgRNA library.**

760

761 **Supplementary Table 2. List of sgRNA sequences used for CRISPRpath screen.**

762

763 **Supplementary Table 3. Enhancers identified from CRISPRn 2X, CRISPRi 2X and**
764 **CRISPRi 3X screens.**

765

766 **Supplementary Table 4. List of primers used for RT-qPCR.**

767

768 **Supplementary Table 5. List of shRNA sequences used for RNA interference**
769 **experiments.**

770

771 **Supplementary Table 6. List of sgRNA sequences used for CRISPRi-mediated enhancer**
772 **validation experiments.**

773

774 **Supplementary Table 7. Pairwise comparisons of data in Figure 4c.**

775

776 **References**

- 777 1. Soldner, F. et al. Parkinson-associated risk variant in distal enhancer of alpha-synuclein
778 modulates target gene expression. *Nature* **533**, 95-99 (2016).
- 779 2. Visel, A., Rubin, E.M. & Pennacchio, L.A. Genomic views of distant-acting enhancers. *Nature*
780 **461**, 199-205 (2009).
- 781 3. Long, H.K. et al. Loss of Extreme Long-Range Enhancers in Human Neural Crest Drives a
782 Craniofacial Disorder. *Cell Stem Cell* **27**, 765-783 e714 (2020).
- 783 4. Consortium, E.P. et al. Expanded encyclopaedias of DNA elements in the human and mouse
784 genomes. *Nature* **583**, 699-710 (2020).
- 785 5. Sanjana, N.E. et al. High-resolution interrogation of functional elements in the noncoding
786 genome. *Science* **353**, 1545-1549 (2016).
- 787 6. Gasperini, M. et al. A Genome-wide Framework for Mapping Gene Regulation via Cellular
788 Genetic Screens. *Cell* **176**, 377-390 e319 (2019).
- 789 7. Diao, Y. et al. A tiling-deletion-based genetic screen for cis-regulatory element identification
790 in mammalian cells. *Nat Methods* **14**, 629-635 (2017).
- 791 8. Fulco, C.P. et al. Activity-by-contact model of enhancer-promoter regulation from
792 thousands of CRISPR perturbations. *Nat Genet* **51**, 1664-1669 (2019).
- 793 9. Simeonov, D.R. et al. Discovery of stimulation-responsive immune enhancers with CRISPR
794 activation. *Nature* **549**, 111-115 (2017).
- 795 10. Klann, T.S. et al. CRISPR-Cas9 epigenome editing enables high-throughput screening for
796 functional regulatory elements in the human genome. *Nat Biotechnol* **35**, 561-568 (2017).
- 797 11. Fulco, C.P. et al. Systematic mapping of functional enhancer-promoter connections with
798 CRISPR interference. *Science* **354**, 769-773 (2016).
- 799 12. Gasperini, M. et al. CRISPR/Cas9-Mediated Scanning for Regulatory Elements Required for
800 HPRT1 Expression via Thousands of Large, Programmed Genomic Deletions. *Am J Hum*
801 *Genet* **101**, 192-205 (2017).
- 802 13. Pecina-Slaus, N., Kafka, A., Salamon, I. & Bukovac, A. Mismatch Repair Pathway, Genome
803 Stability and Cancer. *Front Mol Biosci* **7**, 122 (2020).
- 804 14. Yan, T., Berry, S.E., Desai, A.B. & Kinsella, T.J. DNA mismatch repair (MMR) mediates 6-
805 thioguanine genotoxicity by introducing single-strand breaks to signal a G2-M arrest in
806 MMR-proficient RKO cells. *Clin Cancer Res* **9**, 2327-2334 (2003).
- 807 15. Li, G.M. Mechanisms and functions of DNA mismatch repair. *Cell Res* **18**, 85-98 (2008).
- 808 16. Mandegar, M.A. et al. CRISPR Interference Efficiently Induces Specific and Reversible Gene
809 Silencing in Human iPSCs. *Cell Stem Cell* **18**, 541-553 (2016).
- 810 17. Tycko, J. et al. Mitigation of off-target toxicity in CRISPR-Cas9 screens for essential non-
811 coding elements. *Nat Commun* **10**, 4063 (2019).
- 812 18. Perez, A.R. et al. GuideScan software for improved single and paired CRISPR guide RNA
813 design. *Nat Biotechnol* **35**, 347-349 (2017).
- 814 19. Radzishewska, A., Shlyueva, D., Muller, I. & Helin, K. Optimizing sgRNA position markedly
815 improves the efficiency of CRISPR/dCas9-mediated transcriptional repression. *Nucleic*
816 *Acids Res* **44**, e141 (2016).
- 817 20. Rosenbluh, J. et al. Complementary information derived from CRISPR Cas9 mediated gene
818 deletion and suppression. *Nat Commun* **8**, 15403 (2017).
- 819 21. Li, K. et al. Interrogation of enhancer function by enhancer-targeting CRISPR epigenetic
820 editing. *Nat Commun* **11**, 485 (2020).

- 821 22. Engreitz, J.M. et al. Local regulation of gene expression by lncRNA promoters, transcription
822 and splicing. *Nature* **539**, 452-455 (2016).
- 823 23. Song, M. et al. Cell-type-specific 3D epigenomes in the developing human cortex. *Nature*
824 **587**, 644-649 (2020).
- 825 24. Kubo, N. et al. Promoter-proximal CTCF binding promotes distal enhancer-dependent gene
826 activation. *Nat Struct Mol Biol* (2021).
- 827 25. Fernandez-Zapico, M.E. et al. A functional family-wide screening of SP/KLF proteins
828 identifies a subset of suppressors of KRAS-mediated cell growth. *Biochem J* **435**, 529-537
829 (2011).
- 830 26. Ren, B. et al. E2F integrates cell cycle progression with DNA repair, replication, and G(2)/M
831 checkpoints. *Genes Dev* **16**, 245-256 (2002).
- 832 27. Polager, S., Kalma, Y., Berkovich, E. & Ginsberg, D. E2Fs up-regulate expression of genes
833 involved in DNA replication, DNA repair and mitosis. *Oncogene* **21**, 437-446 (2002).
- 834 28. Canver, M.C. et al. BCL11A enhancer dissection by Cas9-mediated in situ saturating
835 mutagenesis. *Nature* **527**, 192-197 (2015).
- 836 29. Rajagopal, N. et al. High-throughput mapping of regulatory DNA. *Nat Biotechnol* **34**, 167-
837 174 (2016).
- 838 30. Koike-Yusa, H., Li, Y., Tan, E.P., Velasco-Herrera Mdel, C. & Yusa, K. Genome-wide recessive
839 genetic screening in mammalian cells with a lentiviral CRISPR-guide RNA library. *Nat*
840 *Biotechnol* **32**, 267-273 (2014).
- 841 31. Adamson, B. et al. A Multiplexed Single-Cell CRISPR Screening Platform Enables Systematic
842 Dissection of the Unfolded Protein Response. *Cell* **167**, 1867-1882 e1821 (2016).
- 843 32. Diao, Y. et al. A new class of temporarily phenotypic enhancers identified by CRISPR/Cas9-
844 mediated genetic screening. *Genome Res* **26**, 397-405 (2016).
- 845 33. Zabidi, M.A. et al. Enhancer-core-promoter specificity separates developmental and
846 housekeeping gene regulation. *Nature* **518**, 556-559 (2015).
- 847 34. Dao, L.T.M. & Spicuglia, S. Transcriptional regulation by promoters with enhancer function.
848 *Transcription* **9**, 307-314 (2018).
- 849 35. Mikhaylichenko, O. et al. The degree of enhancer or promoter activity is reflected by the
850 levels and directionality of eRNA transcription. *Genes Dev* **32**, 42-57 (2018).
- 851 36. Thomas, H.F. et al. Temporal dissection of an enhancer cluster reveals distinct temporal
852 and functional contributions of individual elements. *Mol Cell* (2021).
- 853 37. Robinson, M.D., McCarthy, D.J. & Smyth, G.K. edgeR: a Bioconductor package for differential
854 expression analysis of digital gene expression data. *Bioinformatics* **26**, 139-140 (2010).
- 855 38. Ramirez, F. et al. deepTools2: a next generation web server for deep-sequencing data
856 analysis. *Nucleic Acids Res* **44**, W160-165 (2016).
- 857 39. Buenrostro, J.D., Giresi, P.G., Zaba, L.C., Chang, H.Y. & Greenleaf, W.J. Transposition of native
858 chromatin for fast and sensitive epigenomic profiling of open chromatin, DNA-binding
859 proteins and nucleosome position. *Nat Methods* **10**, 1213-1218 (2013).
- 860 40. Dobin, A. et al. STAR: ultrafast universal RNA-seq aligner. *Bioinformatics* **29**, 15-21 (2013).
- 861 41. Li, B. & Dewey, C.N. RSEM: accurate transcript quantification from RNA-Seq data with or
862 without a reference genome. *BMC Bioinformatics* **12**, 323 (2011).
- 863 42. Kaya-Okur, H.S. et al. CUT&Tag for efficient epigenomic profiling of small samples and
864 single cells. *Nat Commun* **10**, 1930 (2019).
- 865 43. Meers, M.P., Tenenbaum, D. & Henikoff, S. Peak calling by Sparse Enrichment Analysis for
866 CUT&RUN chromatin profiling. *Epigenetics Chromatin* **12**, 42 (2019).

- 867 44. Fang, R. et al. Mapping of long-range chromatin interactions by proximity ligation-assisted
868 ChIP-seq. *Cell Res* **26**, 1345-1348 (2016).
- 869 45. Juric, I. et al. MAPS: Model-based analysis of long-range chromatin interactions from PLAC-
870 seq and HiChIP experiments. *PLoS Comput Biol* **15**, e1006982 (2019).
- 871 46. Mumbach, M.R. et al. Enhancer connectome in primary human cells identifies target genes
872 of disease-associated DNA elements. *Nat Genet* **49**, 1602-1612 (2017).
- 873 47. Grant, C.E., Bailey, T.L. & Noble, W.S. FIMO: scanning for occurrences of a given motif.
874 *Bioinformatics* **27**, 1017-1018 (2011).
- 875 48. Kulakovskiy, I.V. et al. HOCOMOCO: towards a complete collection of transcription factor
876 binding models for human and mouse via large-scale ChIP-Seq analysis. *Nucleic Acids Res*
877 **46**, D252-D259 (2018).
878

Figure 1

bioRxiv preprint doi: <https://doi.org/10.1101/2021.02.19.431931>; this version posted May 20, 2021. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted bioRxiv a license to display the preprint in perpetuity. It is made available under a [CC-BY-NC-ND 4.0 International license](#).

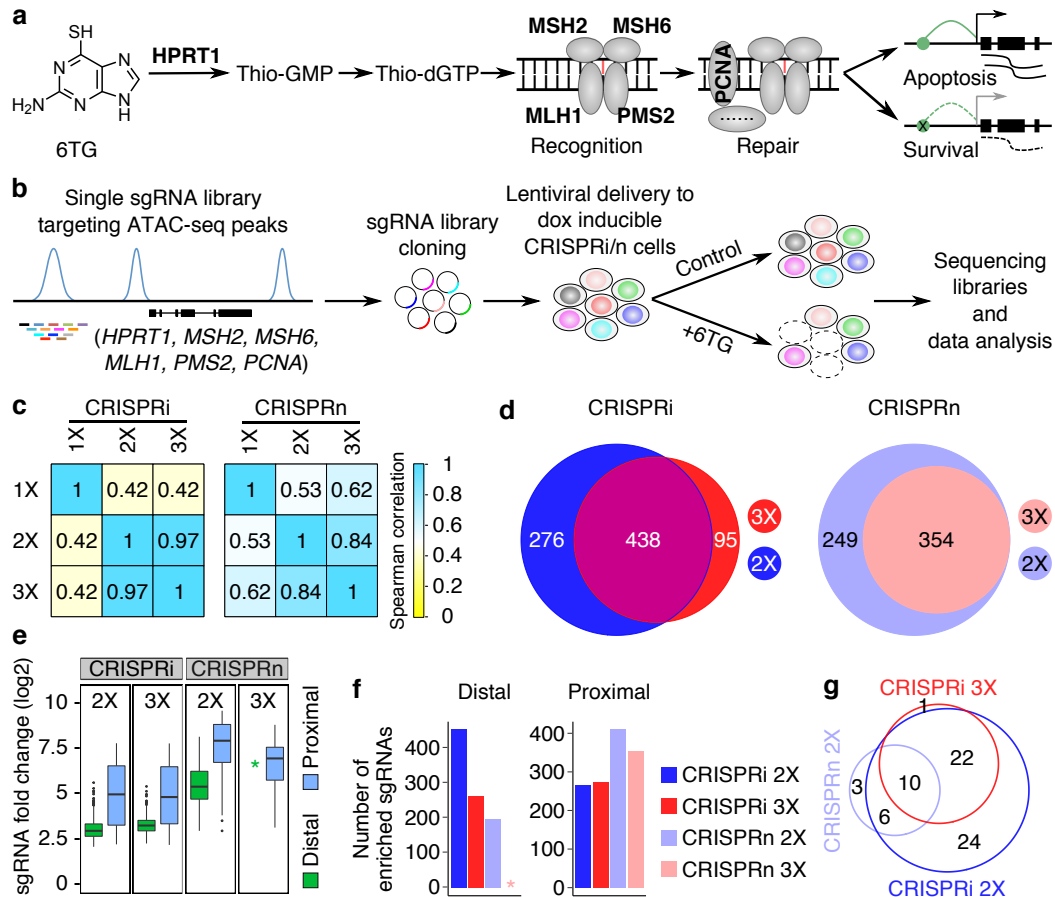


Figure 1. CRISPRpath for identifying enhancers of multiple genes.

(a) Six genes (*HPRT1*, *MSH2*, *MSH6*, *MLH1*, *PMS2* and *PCNA*) in the 6TG-induced mismatch repair process were used for CRISPRpath screen in this study. (b) Schematic of the CRISPRpath screening strategy with 6TG treatment in iPSCs. Cell survival was used as readout for the screen. (c) Spearman correlation analysis of sgRNA ranking based on fold change for CRISPRpath screens with different 6TG concentrations (1X, 2X and 3X). (d) Venn diagram shows the overlapping enriched sgRNAs identified from the screens with 2X and 3X 6TG treatments. (e) Box plots show the fold change of the enriched distal and proximal sgRNAs from 2X, 3X CRISPRi and CRISPRn screens. Star indicates no enriched distal sgRNA was identified from 3X CRISPRn screen. Boxplots indicate the median, interquartile range (IQR), $Q1 - 1.5 \times IQR$ and $Q3 + 1.5 \times IQR$. (f) Bar plot shows the number of enriched distal and proximal sgRNAs from 2X, 3X CRISPRi and CRISPRn screens. Star indicates no enriched distal sgRNA was identified from 3X CRISPRn screen. (g) Venn diagram shows the identified enhancers from each CRISPRpath screen.

Figure 2

bioRxiv preprint doi: <https://doi.org/10.1101/2021.02.19.431931>; this version posted May 20, 2021. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted bioRxiv a license to display the preprint in perpetuity. It is made available under a [CC-BY-NC-ND 4.0 International license](#).

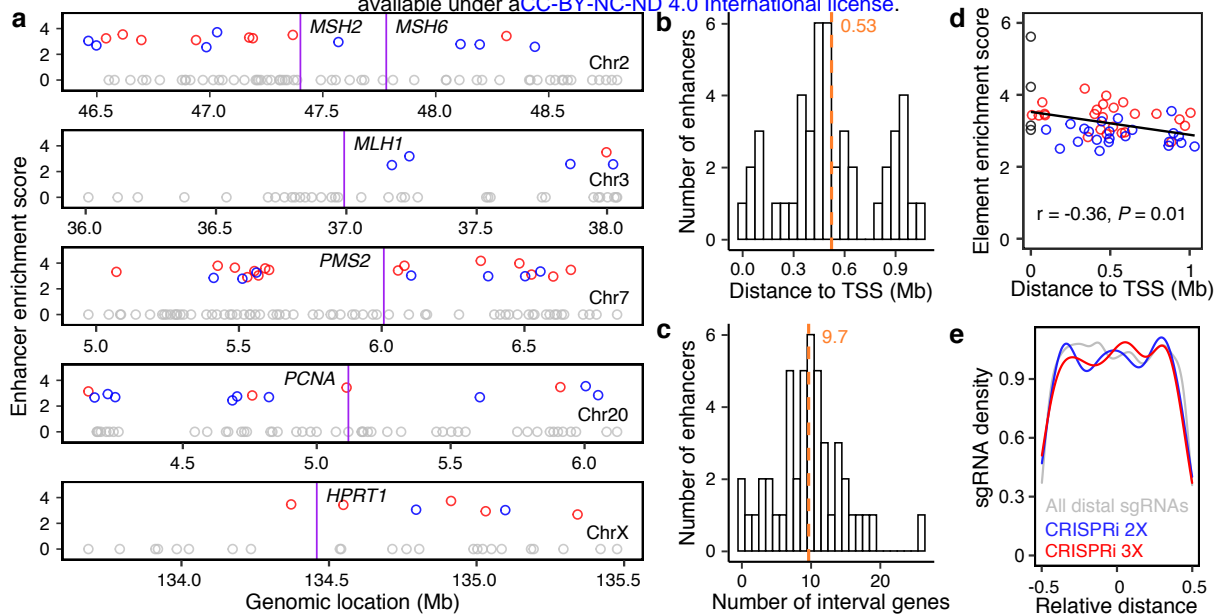


Figure 2. Genomic features of identified enhancers from CRISPRpath using CRISPRi.

(a) Genomic locations of identified enhancers relative to TSS. Circles indicate enhancers identified from the CRISPRi 3X screen (red), enhancers uniquely identified from the CRISPRi 2X screen (blue), tested CREs that are not identified as enhancers (grey). Purple lines label the location of each target gene. (b) Histogram shows the distance distribution between identified enhancers and their paired TSS. Mean is indicated with an orange dashed line and only the enhancers for *MLH1*, *PMS2*, *PCNA*, *HPRT1* are included in b and c. (c) Histogram shows the number of interval genes between enhancers and their target gene TSS. Mean is indicated with an orange dashed line and only the enhancers for *MLH1*, *PMS2*, *PCNA*, *HPRT1* are included in b and c. (d) A weak negative correlation is observed between enrichment score and genomic distance between enhancers and their target genes (Pearson correlation, $r = -0.36$, $P = 0.01$). Black circles indicate promoters. The red and blue circles are enhancers showing in a. (e) Density plot shows no significant difference (two-tailed two-sample Kolmogorov–Smirnov test) for the distribution of all distal sgRNAs (gray), enriched distal sgRNAs from 2X (blue) and 3X screen (red) CRISPRi screens.

Figure 3

bioRxiv preprint doi: <https://doi.org/10.1101/2021.02.19.431931>; this version posted May 20, 2021. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted bioRxiv a license to display the preprint in perpetuity. It is made available under aCC-BY-NC-ND 4.0 International license.

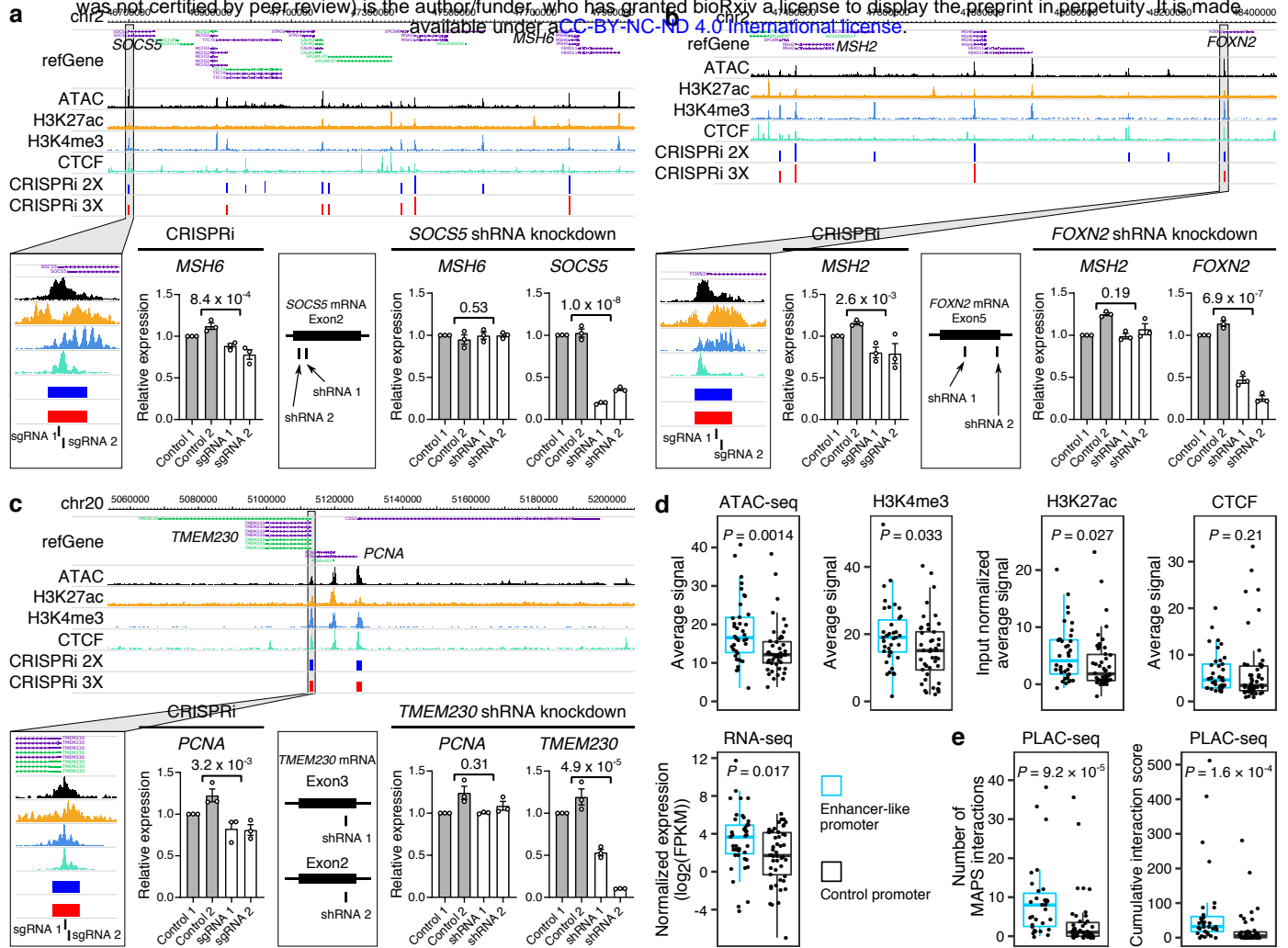


Figure 3. Enhancer-like promoters act as functional enhancers.

(a, b, c) Three examples of promoters function as enhancers. CRISPRi silencing of the promoter region of *SOCS5*, *FOXN2* and *TMEM230* results in significant downregulation of *MSH6*, *MSH2* and *PCNA*, respectively. shRNA knockdown of *SOCS5*, *FOXN2* and *TMEM230* can only downregulate *SOCS5*, *FOXN2* and *TMEM230* expression. Three independent replicates per condition and two independent sgRNAs or shRNAs per replicate were used for each experiment. P values are from two-tailed two-sample t-test. (d) Average signal enrichment of ATAC-seq, gene transcription, H3K4me3, H3K27ac and CTCF binding for enhancer-like promoters ($n = 38$) and control promoters ($n = 47$). P values are from Wilcoxon test. Boxplots indicate the median, IQR, $Q1 - 1.5 \times \text{IQR}$ and $Q3 + 1.5 \times \text{IQR}$. (e) Number of H3K4me3 mediated chromatin interactions and cumulative interaction score for enhancer-like promoters ($n = 31$) and control promoters ($n = 43$). Boxplots indicate median, IQR, $Q1 - 1.5 \times \text{IQR}$ and $Q3 + 1.5 \times \text{IQR}$. P values are calculated from Wilcoxon test.

Figure 4

bioRxiv preprint doi: <https://doi.org/10.1101/2021.02.19.431931>; this version posted May 20, 2021. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted bioRxiv a license to display the preprint in perpetuity. It is made available under aCC-BY-NC-ND 4.0 International license.

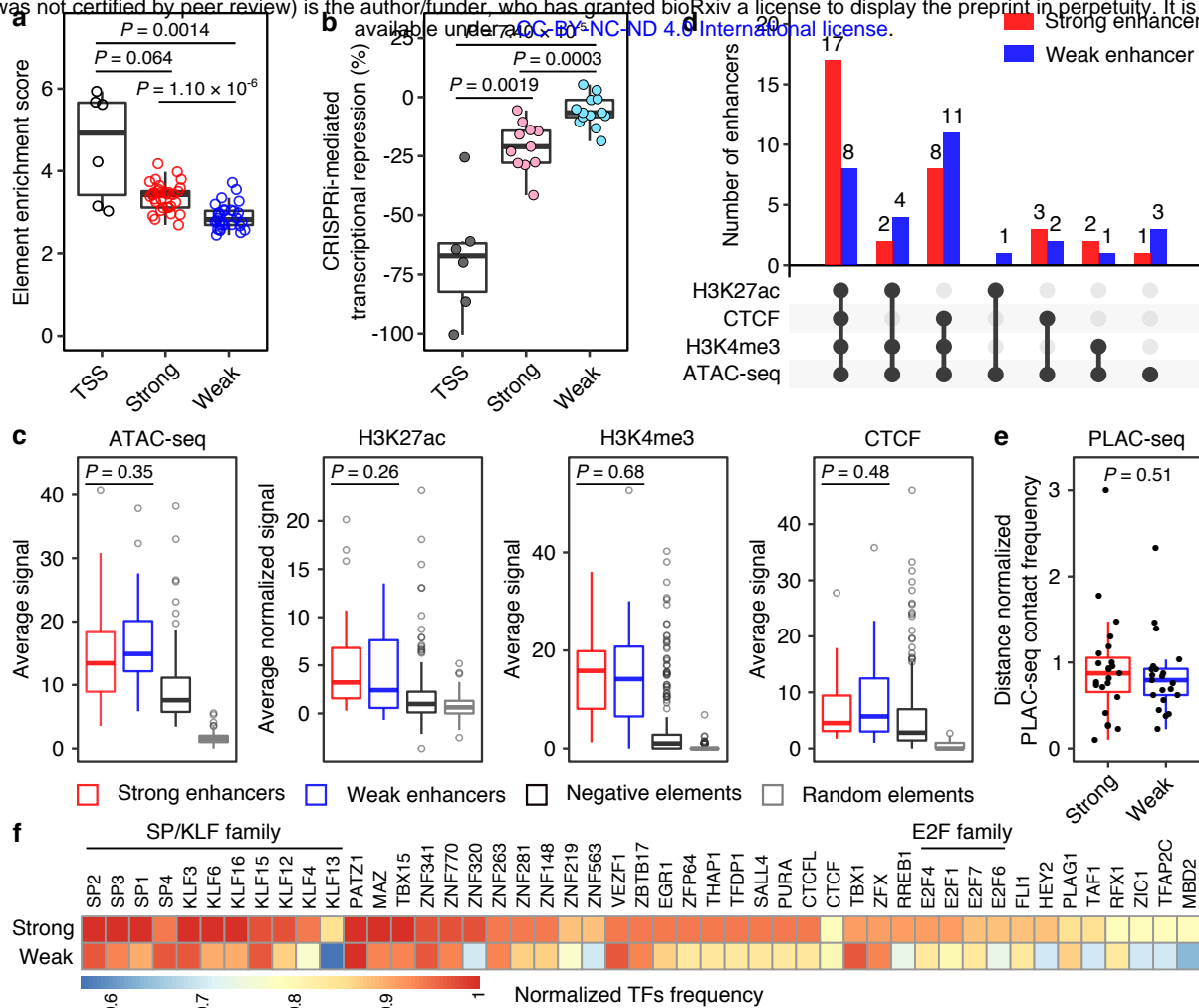
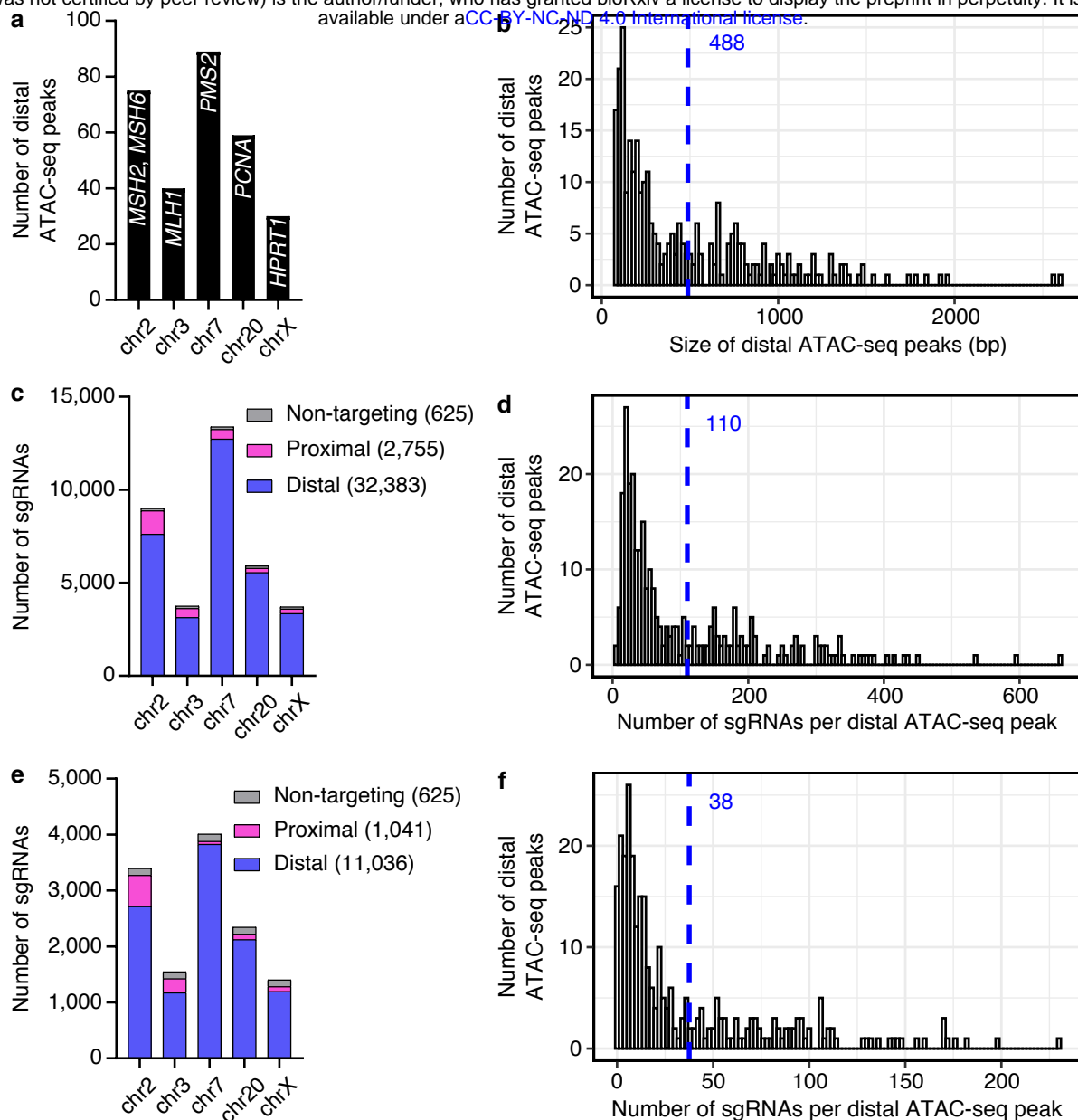


Figure 4. CRISPRpath can distinguish weak and strong enhancers by imposing different selection pressures.

(a) Box plots show the enrichment score of the tested elements. TSS regions (black circles) show highest enrichment scores. Enhancers uniquely identified from the lower selection pressure (CRISPRi 2X, blue circles) exhibit lower enrichment scores compared to the enhancers identified from the higher selection pressure (CRISPRi 3X, red circles). P values are from Wilcoxon test. (b) Box plots show the CRISPRi perturbation at enhancers induced various degrees of transcriptional repression of target genes measured with RT-qPCR. Each dot represents the average value from three biological replicates. CRISPRi targeting TSS regions (dark gray) achieved the highest transcriptional repression. CRISPRi targeting strong enhancers (pink) leads a more substantial transcription silencing on target gene compared to CRISPRi targeting weak enhancers (cyan). P values are from Wilcoxon test. (c) Enrichment analysis of ATAC-seq, H3K27ac, H3K4me3, CTCF binding signals for strong ($n = 33$) and weak ($n = 30$) enhancers. Boxplots indicate the median, IQR, $Q1 - 1.5 \times IQR$ and $Q3 + 1.5 \times IQR$. P values for the difference between strong and weak enhancers are from Wilcoxon test; see Supplementary Table 7 for P values of all pairwise comparisons. (d) Intersection of genomic features for weak enhancers (blue bar) and strong enhancers (red bar). (e) Distance normalized H3K4me3 PLAC-seq contact frequency for strong ($n = 23$) and weak ($n = 21$) enhancers. Only the enhancers for *MLH1*, *PMS2*, *PCNA*, *HPRT1* are included (see Methods). Boxplots indicate the median, IQR, $Q1 - 1.5 \times IQR$ and $Q3 + 1.5 \times IQR$. P value is from Wilcoxon test. (f) Heatmap shows the frequency of transcription factor motifs found in strong and weak enhancers.

Supplementary Figure 1

bioRxiv preprint doi: <https://doi.org/10.1101/2021.02.19.431931>; this version posted May 20, 2021. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted bioRxiv a license to display the preprint in perpetuity. It is made available under aCC-BY-NC-ND 4.0 International license.

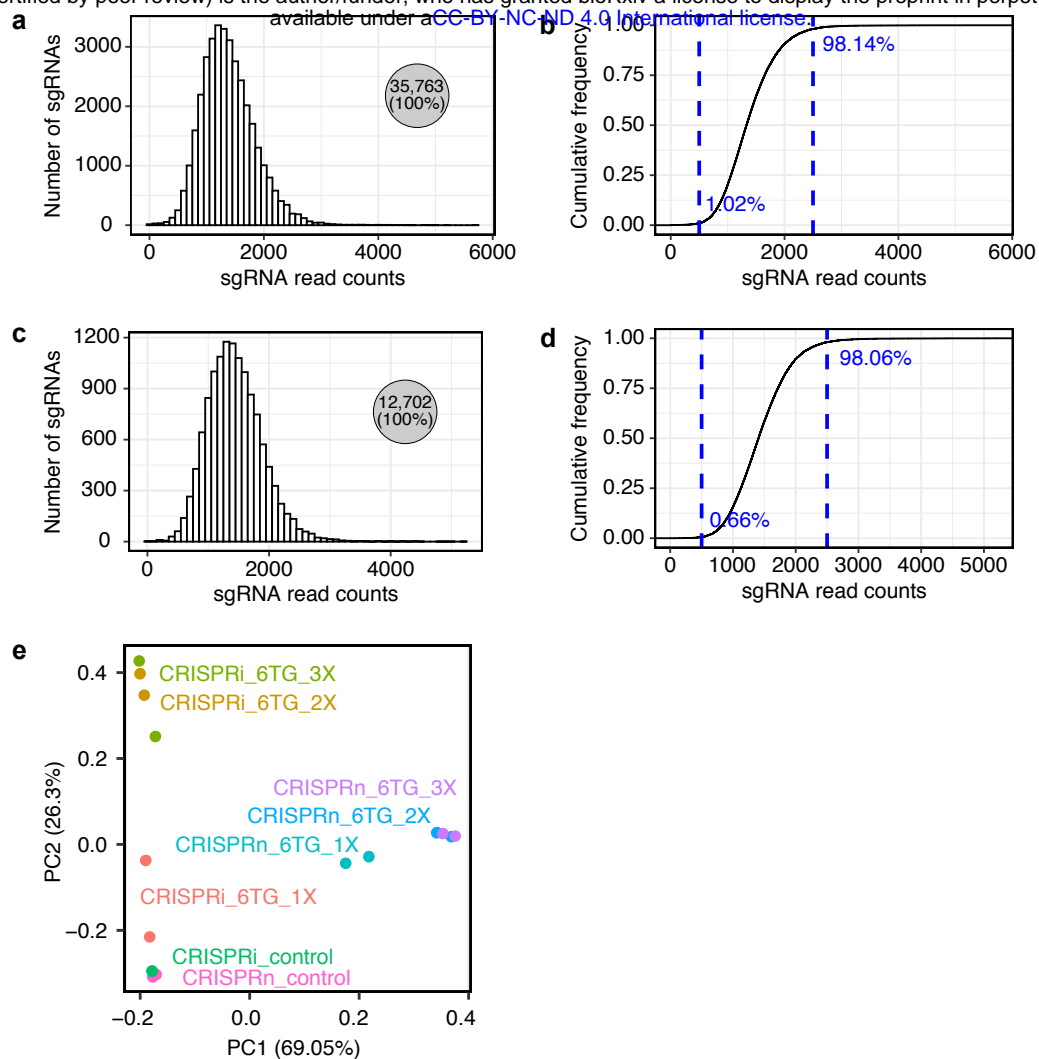


Supplementary Figure 1. Features of sgRNA library for CRISPRpath screen.

(a) Bar graph shows the number of distal ATAC-seq peaks used as candidate CREs for six target genes. (b) Histogram shows size distribution of distal ATAC-seq peaks. The average size is 488 bp (blue dash line). (c, e) The composition of the sgRNA library. In total, 35,763 sgRNAs were included in the library (c), and 12,702 sgRNAs are high quality sgRNAs (e). (d, f) Distribution of the number of sgRNAs per distal ATAC-seq peak. Average numbers of sgRNA per ATAC-seq peak are indicated with blue dash lines.

Supplementary Figure 2

bioRxiv preprint doi: <https://doi.org/10.1101/2021.02.19.431931>; this version posted May 20, 2021. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted bioRxiv a license to display the preprint in perpetuity. It is made available under aCC-BY-NC-ND 4.0 International license.

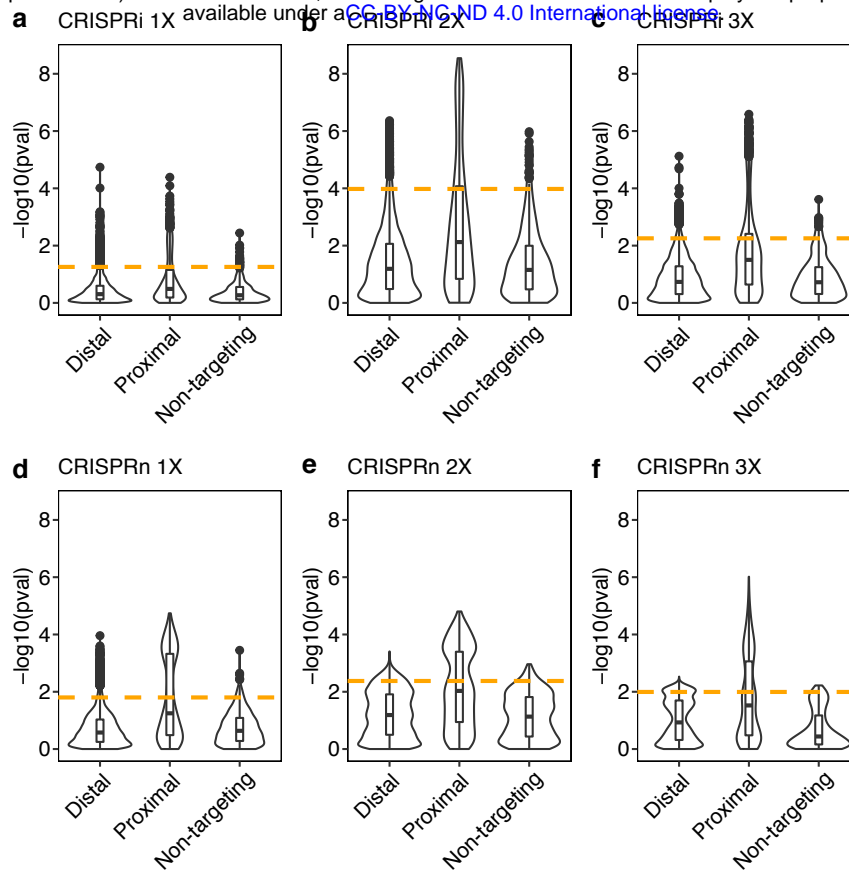


Supplementary Figure 2. Quality of the sgRNA library and CRISPRpath screen libraries.

(a) Distribution of sgRNA oligo read counts in the sgRNA library. (b) Cumulative frequency of sgRNAs in the sgRNA library. (c) Distribution of high quality sgRNAs read counts in the sgRNA library. (d) Cumulative frequency of high quality sgRNAs in the sgRNA library. The constructed sgRNA plasmid library recovered all the designed sgRNAs with the copy number difference less than five fold for at least 97% designed sgRNAs. (e) PCA analysis shows the high reproducibility of the CRISPRpath screen libraries between biological replicates.

Supplementary Figure 3

bioRxiv preprint doi: <https://doi.org/10.1101/2021.02.19.431931>; this version posted May 20, 2021. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted bioRxiv a license to display the preprint in perpetuity. It is made available under aCC-BY-NC-ND 4.0 International license.

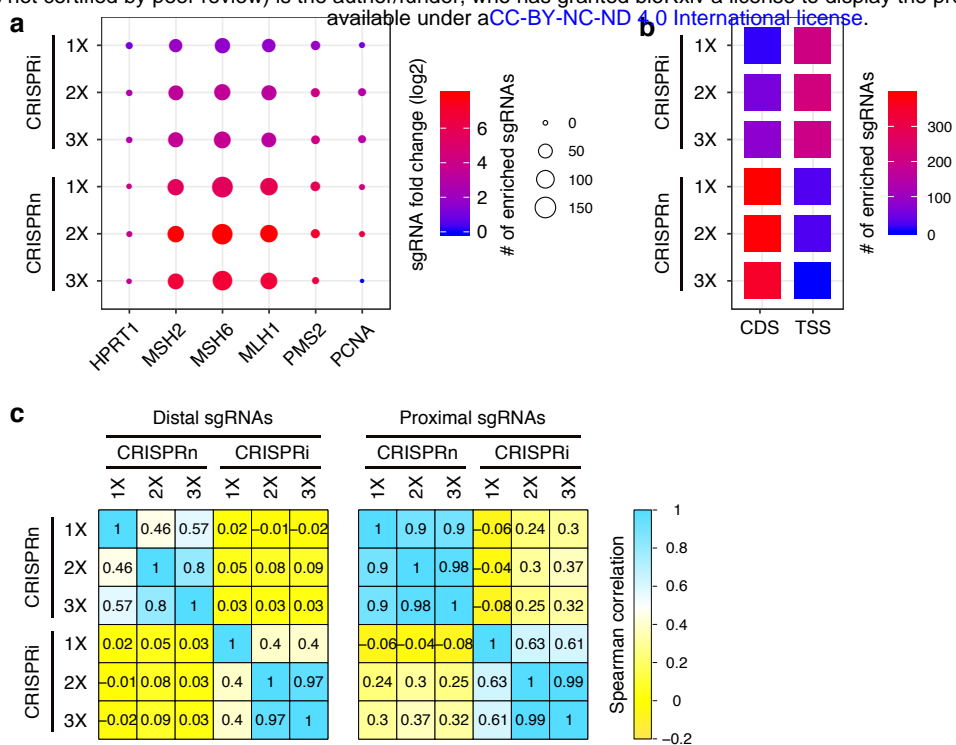


Supplementary Figure 3. P value cutoff used for identifying enriched sgRNAs from each screens.

(a-f) Distribution of P value for tested distal, proximal and non-targeting control sgRNA groups. Orange dash lines indicate 5% percentile of the P values from non-targeting control sgRNAs to achieve a false discovery rate of 5%.

Supplementary Figure 4

bioRxiv preprint doi: <https://doi.org/10.1101/2021.02.19.431931>; this version posted May 20, 2021. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted bioRxiv a license to display the preprint in perpetuity. It is made available under a [CC-BY-NC-ND 4.0 International license](#).

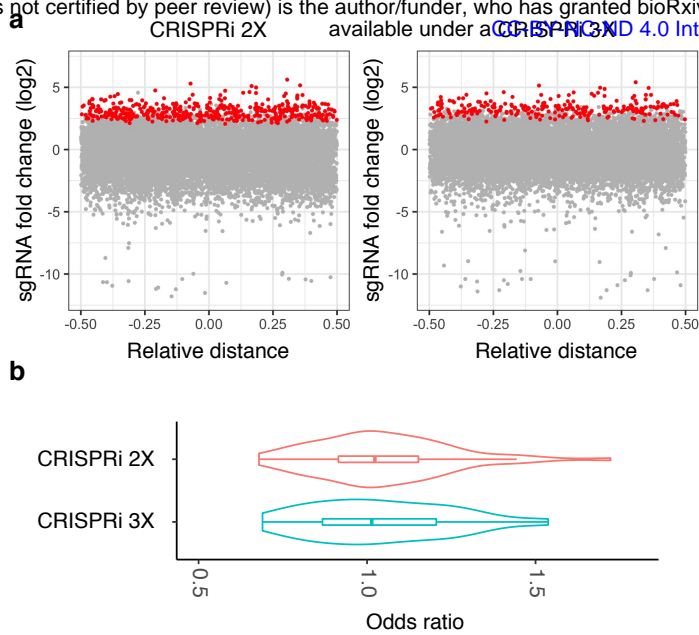


Supplementary Figure 4. Enriched proximal sgRNAs and sgRNA ranking analysis.

(a) Number and fold change of the enriched proximal sgRNAs for the six target genes from CRISPRi and CRISPRn screens. The color indicates fold changes, and the size of circle indicates the number of enriched sgRNAs. (b) Enrichment analysis shows the enriched proximal sgRNAs bias towards the TSS region for CRISPRi screens and the protein coding region (CDS) for CRISPRn screens. Color represents the number of enriched sgRNAs. (c) Spearman correlation analysis of the distal and proximal sgRNAs ranking shows proximal sgRNAs exhibiting higher correlation between each screen compared to distal sgRNAs.

Supplementary Figure 5

bioRxiv preprint doi: <https://doi.org/10.1101/2021.02.19.431931>; this version posted May 20, 2021. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted bioRxiv a license to display the preprint in perpetuity. It is made available under a [CC-BY-NC-ND 4.0 International license](#).

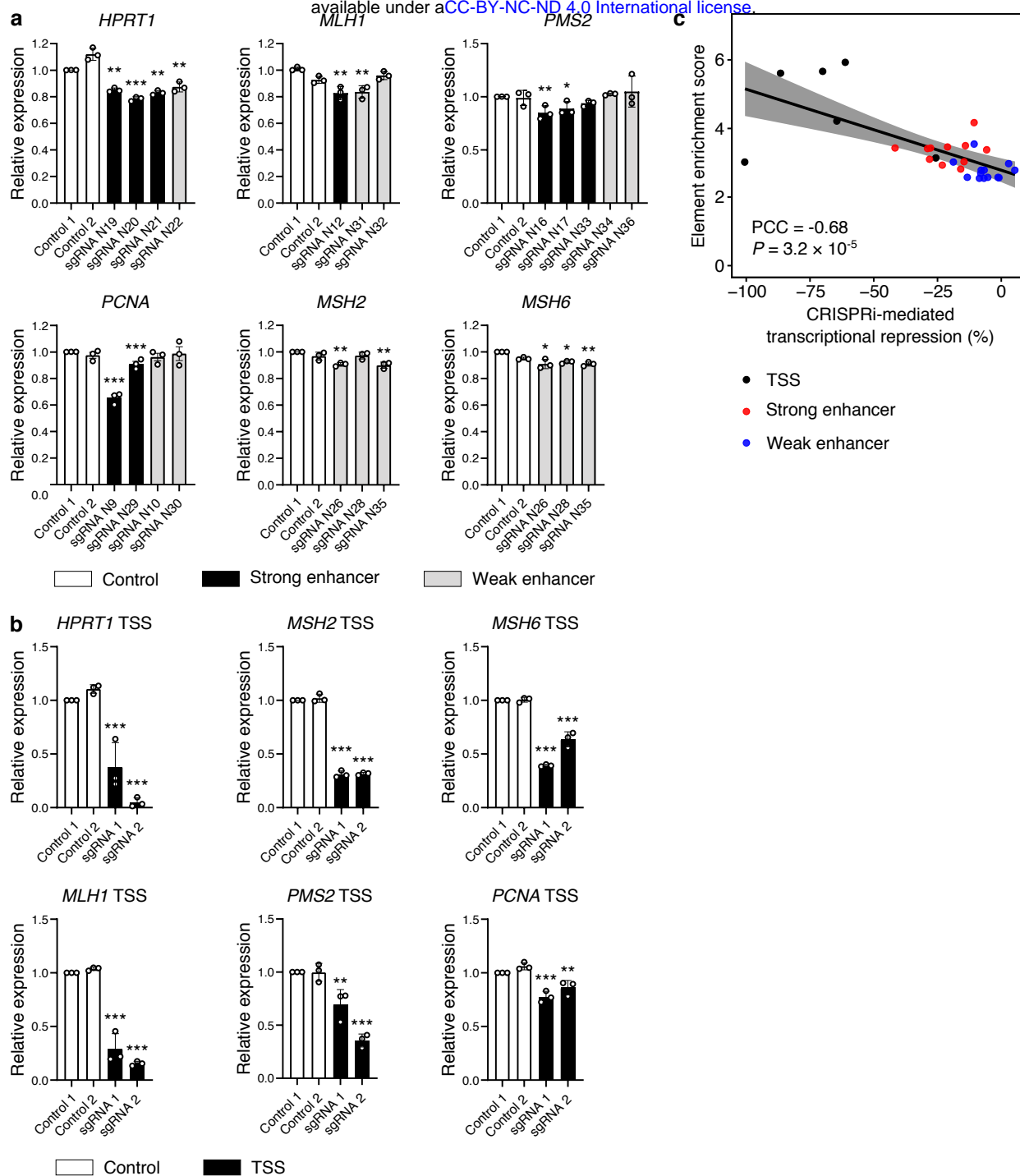


Supplementary Figure 5. Enriched sgRNAs identified from CRISPRi screens exhibit no position and strand preference.

(a) Enriched sgRNAs from CRISPRi 2X (red dots, n=448) and CRISPRi 3X (red dots, n=260) screens showed similar distributions across candidate CREs. (b) Odds ratio analysis of the fold change of enriched sgRNAs shows enriched sgRNAs have no strand preference. Enhancers with enriched sgRNAs only targeting one strand were excluded for the analysis. Odds ratio was calculated for each element with the equation of $\text{ave}(\log_2(\text{fold change of sgRNA targeting plus strand})) / \text{ave}(\log_2(\text{fold change of sgRNA targeting minus strand}))$. Violin plots show the distributions of odds ratio values within each screen, and boxplots indicate the median, IQR, $Q1 - 1.5 \times \text{IQR}$ and $Q3 + 1.5 \times \text{IQR}$.

Supplementary Figure 6

bioRxiv preprint doi: <https://doi.org/10.1101/2021.02.19.431931>; this version posted May 20, 2021. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted bioRxiv a license to display the preprint in perpetuity. It is made available under a [CC-BY-NC-ND 4.0 International license](#).

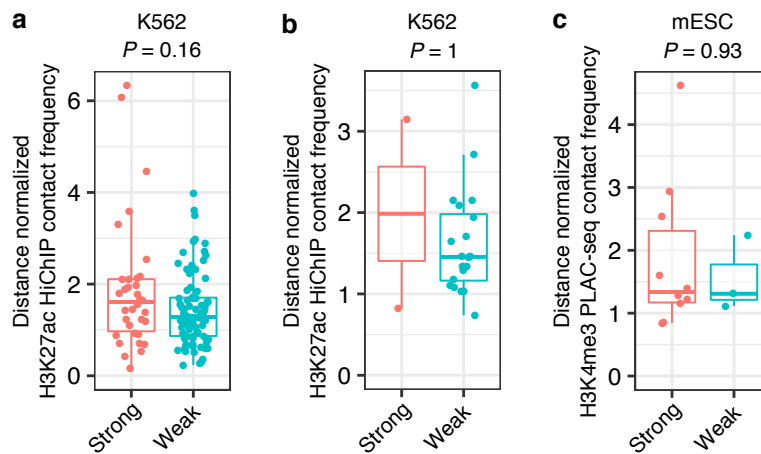


Supplementary Figure 6. Validation of CRISPRpath identified enhancers.

(a) Validation of the strong (black) and weak (grey) enhancers with CRISPRi followed by RT-qPCR. Three independent replicates per condition. The significance was calculated with two-tailed two-sample t-test. Data are mean and s.d. * $P < 0.05$, ** $P < 0.01$, *** $P < 0.001$. (b) CRISPRi-mediated transcriptional repression of six target genes by targeting TSS of each gene. Three independent replicates per condition. The significance was calculated with two-tailed two-sample t-test. Data are mean and s.d. * $P < 0.05$, ** $P < 0.01$, *** $P < 0.001$. (c) Pearson correlation analysis reveals element enrichment score from CRISPRpath screens correlates with element effect size on transcription from CRISPRi (Pearson correlation, PCC = -0.68, $P = 3.2 \times 10^{-5}$).

Supplementary Figure 7

bioRxiv preprint doi: <https://doi.org/10.1101/2021.02.19.431931>; this version posted May 20, 2021. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted bioRxiv a license to display the preprint in perpetuity. It is made available under a [CC-BY-NC-ND 4.0 International license](#).



Supplementary Figure 7. Chromatin contact frequency analysis for the enhancers in K562 cells and mESCs.

(a) Distance normalized H3K27ac HiChIP contact frequency for strong ($n = 34$) and weak ($n = 82$) enhancers identified with crisprQTL mapping in K562 cells. (b) Distance normalized H3K27ac HiChIP contact frequency for strong ($n = 2$) and weak ($n = 20$) enhancers identified with CRISPRi-FlowFISH screen in K562 cells. (c) Distance normalized H3K4me3 PLAC-seq contact frequency for strong ($n = 10$) and weak ($n = 3$) enhancers identified in mouse embryonic stem cells. Boxplots indicate the median, IQR, $Q1 - 1.5 \times IQR$ and $Q3 + 1.5 \times IQR$. P values are from Wilcoxon test.