

UCSF

UC San Francisco Previously Published Works

Title

A Compendium of Promoter-Centered Long-Range Chromatin Interactions in the Human Genome

Permalink

<https://escholarship.org/uc/item/9ng2w70x>

Author

Chan, Marilyn

Publication Date

2019-10-01

Data Availability

The data associated with this publication are available at:

<https://www.ncbi.nlm.nih.gov/geo/>

Peer reviewed

1Title:

**2A Compendium of Promoter-Centered Long-Range
3Chromatin Interactions in the Human Genome**

4Authors:

5Inkyung Jung^{1*†}, Anthony Schmitt^{2,3*}, Yarui Diao^{2,4*}, Andrew J. Lee¹, Tristin
6Liu², Dongchan Yang⁵, Catherine Tan², Junghyun Eom¹, Marilyn Chan⁶,
7Sora Chee², Zachary Chiang⁷, Changyoun Kim^{8,9}, Eliezer Masliah^{8,9,10}, Cathy
8L. Barr¹¹, Bin Li¹, Samantha Kuan², Dongsup Kim⁵, Bing Ren^{2,12†}

9

10Affiliations:

11¹Department of Biological Sciences, KAIST, Daejeon 34141, Korea

12²Ludwig Institute for Cancer Research, La Jolla, CA 92093, USA

13³UCSD Biomedical Sciences Graduate Program, La Jolla, CA 92093, USA

14⁴Departments of Cell Biology and Orthopaedic Surgery, Regenerative Next
15Initiative, Duke University School of Medicine. Durham, NC 27710

16⁵Department of Bio and Brain engineering, KAIST, Daejeon 34141, Korea

17⁶University of California San Francisco, San Francisco, CA 94158, USA

18⁷Department of Bioengineering, UCSD, La Jolla, CA 92093, USA

19⁸Molecular Neuropathology Section, Laboratory of Neurogenetics, National
20Institute on Aging, National Institutes of Health, Bethesda, MD 20892, USA

21⁹Department Neurosciences, School of Medicine, University of California,

22San Diego, La Jolla, CA 92093, USA

23¹⁰Department of Pathology, School of Medicine, University of California,
24San Diego, La Jolla, CA 92093, USA

25¹¹Krembil Research Institute, University Health Network, Toronto, and The
26Hospital for Sick Children, Ontario M5T 2S8, Canada

27¹²Department of Cellular and Molecular Medicine, Institute of Genomic
28Medicine, and Moores Cancer Center, University of California at San Diego,
29La Jolla, CA 92093, USA

30

31*These authors contributed equally to this work

32[†]Correspondence to Inkyung Jung (ijung@kaist.ac.kr) and Bing Ren
33(biren@ucsd.edu)

34

35Note:

36All raw and processed data are deposited into GEO database under
37accession number GSE86189. Reviewer's access token is
38mtsxkwgunlipvqb.

39A genome browser session has been set up for visualization of the
40promoter-centered chromatin interactions described in the current study -

41<http://epigenomegateway.wustl.edu/browser/?>

42[genome=hg19&session=IR82F6oIpo&statusId=446157315](http://epigenomegateway.wustl.edu/browser/?genome=hg19&session=IR82F6oIpo&statusId=446157315) (Remy the
43chef)

44

45

46

47

48

49

50**Abstract:**

51**A large number of putative *cis*-regulatory sequences have been**
52**annotated in the human genome, but the genes they control**
53**remain to be defined. To bridge this gap, we generate maps of**
54**long-range chromatin interactions centered on 18,943 well-**
55**annotated promoters for protein-coding genes in 27 human**
56**cell/tissue types. We use this information to infer the target**
57**genes of 70,329 candidate regulatory elements, and suggest**
58**potential regulatory function for 27,325 non-coding sequence**
59**variants associated with 2,117 physiological traits and diseases.**
60**Integrative analysis of these promoter-centered interactome**
61**maps reveals widespread enhancer-like promoters involved in**
62**gene regulation and common molecular pathways underlying**
63**distinct groups of human traits and diseases.**

64

65**Main Text:**

66Genome-Wide Association Studies (GWAS) have uncovered thousands of
67genetic variants associated with human diseases and phenotypic traits¹,
68but molecular characterization of these genetic variants has been
69challenging because they are mostly non-coding and lack clear functional
70annotation. Recent studies have shown that these non-coding variants are
71frequently marked by chromatin signatures of *cis*-regulatory elements
72(cREs) in cells, leading to the hypothesis that a substantial fraction of
73variants may act by affecting transcriptional regulation^{2,3}. To formally test

74this hypothesis, it is critical to define the target genes of cREs in the
75human genome. However, inferring target genes of cREs based on linear
76genomic sequences is not straightforward, since cREs can regulate non-
77adjacent genes over large genomic distances⁴⁻⁷. Such long-range
78regulation can take place because chromatin fibers are folded into a
79higher-order structure in which distant DNA fragments can be juxtaposed
80in space⁸. Consequently, mapping spatial contacts between DNA has the
81potential to uncover target genes of cREs. To this end, Chromosome
82Conformation Capture (3C) techniques such as 4C-seq, ChIA-PET and Hi-C
83have been developed to determine chromatin interactions in a high
84throughput manner⁹⁻¹⁵. More recently, Hi-C combined with targeted
85capture and sequencing (capture Hi-C) has emerged as a cost-effective
86method to map chromatin interactions for specific regions at high-
87resolution¹⁶⁻²⁵.

88

89In order to systematically annotate candidate target genes for the cREs in
90the human genome, we performed capture Hi-C experiments (Fig. 1a;
91Extended Data Fig. 1) to interrogate chromatin interactions centered at
92well-annotated human gene promoters for 19,462 protein-coding genes
93(see Methods). We carried out these experiments with 27 human
94cell/tissue types including embryonic stem cells, four early embryonic
95lineages (mesendoderm, mesenchymal stem cell, neural progenitor cells,
96and trophoblast), two primary cell lines (fibroblast cells and
97lymphoblastoid cells), and 20 primary tissue types (hippocampus,

98dorsolateral prefrontal cortex, esophagus, lung, liver, pancreas, small
99bowel, sigmoid colon, thymus, bladder, adrenal gland, aorta, gastric
100tissue, left heart ventricle, right heart ventricle, right heart atrium, ovary,
101psoas, spleen, and fat) for which reference epigenome maps have
102previously been produced as part of the Epigenome Roadmap project
103(Extended Data Fig. 2a; Supplementary Table 1)²⁶. We designed and
104synthesized 12 capture probes for each promoter, six for each of the
105nearest *HindIII* restriction sites upstream and downstream of the
106transcription start site (TSS). Among 16,720 promoter-containing *HindIII*
107restriction DNA fragments, 14,357 (86%) contain a single promoter, but
108the 2,363 remaining *HindIII* fragments harbor multiple promoters
109(Extended Data Fig. 2b; see Methods). The robustness and the coverage
110of capture probe synthesis were validated by sequencing (Extended Data
111Fig. 2c-f). On average, each capture Hi-C experiment produced 65 million
112unique, on-target paired-end reads, yielding a total of 1.8 billion valid read
113pairs, ~30% of which were between DNA fragments >15kb apart
114(Supplementary Table 2).

115

116To identify the long-range chromatin interactions from the capture Hi-C
117data, two normalization steps were introduced. First, the biases in capture
118efficiency of each promoter (Extended Data Fig. 2g, h) were calibrated
119with the variable “capturability” for each DNA fragment, defined as the
120fraction of total read counts mapped to the region, using a β -spline
121regression model (see Methods). Second, significant chromatin

122interactions were then identified after normalizing against the distance-
123dependent background signals (9% and 5% FDR for P-O and P-P
124interactions, respectively) (see Methods). Focusing on the *HindIII*
125fragments over 15kb away and within 2Mbp of each promoter, we
126determined a total of 892,013 chromatin interactions (431,141 unique
127interacting pairs) in one or more of the 27 human cell/tissue types (Fig.
1281b; Extended Data Fig. 3a; Supplementary Table 3-5). The median
129distance between the interacting DNA pairs was 158kb, which is within a
130similar range of previously reported chromatin loops and eQTL
131associations (Fig. 1c; Supplementary Table 6)^{10,12,13}. The slight discrepancy
132between pHi-C interactions and eQTL-associations may be attributed to
133different experimental approaches, but nevertheless, the two methods
134give complementary information to each other. Between 13% and 45%
135pHi-C interactions detected in a cell or tissue type were unique to that
136cell/tissue type (Extended Data Fig. 3b). As expected, most of the
137detected chromatin interactions were within Topologically Associating
138Domains (TADs) defined in the corresponding tissue/cell type (Extended
139Data Fig. 3c, d)^{27,11}.

140

141To demonstrate that pHi-C could effectively and reproducibly capture
142long-range chromatin interactions as detected by whole genome *in situ*
143Hi-C, we compared the pHi-C data with the *in situ* Hi-C data obtained
144from four distinct biosamples, including two cell lines (IMR90 lung
145fibroblast cell line and GM12878 lymphoblastoid cell line¹³) and two

146 primary tissues - dorsolateral prefrontal cortex and hippocampus (see
147 Methods). Results of pHi-C experiments accurately recapitulated
148 chromatin loops identified from *in situ* Hi-C assays in all samples, with the
149 area under the receiver operating curve (ROC) ranging between 0.84 and
150 0.91 (Extended Data Fig. 4a-e) (see Methods). Additionally, we found high
151 reproducibility of pHi-C chromatin interactions between different donors
152 (average ROC score = 0.85; the average Spearman's rank correlation
153 between replicates = 0.4; Extended Data Fig. 4f-j; Supplementary Table 7;
154 see Methods), and between two independent studies (Extended Data Fig.
155 4k). The observation that interactions identified in both replicates
156 exhibited the strongest interaction signals, while interactions identified in
157 one replicate were moderately strong in one replicate, but moderately
158 weak in the other replicate (Extended Data Fig. 4l-m), suggests that the
159 interactions that are specific to one replicate may be due to under-
160 sampling of the other replicate.

161

162 The chromatin interactome maps allowed us to assign candidate target
163 genes for 70,329 putative cREs, defined based on histone H3K27ac
164 signals in each tissue/cell type profiled previously²⁶, for 17,295 promoters.
165 Each promoter was putatively assigned to 25 cREs on average (Extended
166 Data Fig. 5a), while 45% of cREs were assigned to one candidate target
167 gene (Extended Data Fig. 5b), similar to the previous observation with
168 DNase I hypersensitivity analysis across diverse human cell types²⁸. We
169 took advantage of the existing chromatin datasets collected for the same

170 tissue/cell types²⁶, and examined the relationship of the chromatin states
171 between the cREs and the target promoters (see Methods). As expected,
172 the fragments that extensively interact with multiple promoters were
173 often found at active chromatin regions, such as TF binding clusters or
174 super enhancer regions (Extended Data Fig. 5c-i; Supplementary Table 8-
175 10; see Methods)²⁹. Furthermore, integrative analysis with ChromHMM
176 model revealed that active promoters interact three times more
177 frequently with DNA fragments harboring active enhancers than the
178 bivalent promoters (Fig. 1d). On the other hand, the bivalent promoters
179 interact five times more frequently with genomic regions associated with
180 Polycomb Repressor Complexes than the active promoters (Fig. 1d).
181 Further analysis based on a refined 50-chromatin-state ChromHMM model
182 for 5 cell lines also supports our conclusion (Extended Data Fig. 6).

183

184 Three lines of evidence demonstrate that the above promoter-centered
185 chromatin interactions contain information on regulatory interactions at
186 each promoter in the corresponding cell/tissue types. First, we compared
187 the chromatin interactions at promoters with regulatory relationships
188 inferred from expression quantitative trait loci (eQTL) in 14 matched
189 tissue-types that were recently reported by the GTEx consortium (see
190 Methods) (Fig. 2a; Extended Data Fig. 7a-c)³⁰. For each tissue and cell
191 type, the previously reported eQTLs were highly enriched in the chromatin
192 interactions identified in the corresponding tissue, with enrichment up to
193 five-fold (ovary) (Fig. 2b; Extended Data Fig. 7d and e). A total of 42,627

194eQTL associations were detected by pcHi-C chromatin interactions, while
195only 21,362 were expected by random chance after controlling for linear
196genomic distances (Supplementary Table 11 and 12). Second, there is
197significant correlation between activities of *cis*-regulatory sequences and
198the assigned candidate target gene expression across multiple tissues and
199cell types, consistent with the purported regulatory relationships.
200Specifically, the histone modification status of H3K27ac of these cREs was
201significantly correlated with the promoter H3K27ac levels (KS-test P value
202< 2.2e-16; Extended Data Fig. 7f) and transcription levels of the predicted
203target genes (KS-test P value < 2.2e-16; Extended Data Fig. 7g) across
204these tissues/cell types. For example, the *POU3F3* gene expression
205(second column in Fig. 2c) was highly correlated with H3K27ac signals in
206the distal cRE (first column in Fig. 2c) connected by a tissue-specific
207chromatin interaction (last column in Fig. 2c). Lastly, cell/tissue-specific
208cRE-promoter pairs connected by pcHi-C interactions are significantly
209associated with active cREs and genes that are specific to the same
210cell/tissue types. For example, hippocampus specific cRE-promoter
211chromatin interactions are significantly associated with active cREs (Fig.
2122d) and highly expressed genes, albeit modest, (Fig. 2e) in hippocampus.
213Significant associations of cell/tissue-specific pcHi-C interactions in active
214cREs and highly expressed genes are found in other cell/tissue types as
215well (Fig. 2f-h, KS-test P value < 2.2e-16, see Methods). The above results,
216taken together, strongly suggest that the predicted cRE-promoter pairs

217could uncover regulatory relationships between the cRE and target genes
218in diverse tissues and cell types.

219

220Widespread promoter-promoter (P-P) interactions have been reported in
221cultured mammalian cells and a few primary tissues^{14,21,31}. The promoter-
222centered interaction maps obtained from 27 diverse tissues and cell types
223allowed us to test whether this is a general phenomenon. Indeed,
224consistent with previous reports, a significant fraction of the chromatin
225interactions was found between two promoters (9%, $n = 79,989$, Fisher's
226Exact test p value $< 2.2e-16$, Extended Data Fig. 8a). The physical
227proximity of these promoters is accompanied by a strikingly high
228correlation in chromatin modification state between the pair of promoters
229across diverse cell/tissue types (Fig. 3a, b). Previously, several promoter
230loci have been shown to function as enhancers to regulate distal genes³²⁻
231³⁴. In support of the functional significance of enhancer-like promoters
232identified in the current study, 6,127 eQTLs match P-P interaction pairs,
233while only 2,722 eQTLs were expected by random chance (Fig. 3c;
234Extended Data Fig. 8b-d; Supplementary Table 13 and 14; see Methods).
235For instance, strong chromatin interactions were found between the
236*DACT3* and *AP2S1* gene promoter regions, and one significant eQTL
237(rs78730097) for *DACT3* gene was located in the *AP2S1* promoter in the
238dorsolateral prefrontal cortex (Fig. 3d). Notably, this eQTL does not show
239any meaningful genetic association with the adjacent downstream gene
240(*AP2S1*) or nearby genes, but is exclusively associated with *DACT3*

241(Extended Data Fig. 8e), suggesting regulatory potential of the *AP2S1*
242promoter region in distal *DACT3* gene regulation. To validate the function
243of enhancer-like promoters, we deleted 2 core promoter regions, where
244the downstream gene is not expressed but the promoter region shows
245active chromatin marks, using CRISPR-mediated system (Extended Data
246Fig. 8f, g; see Methods). Deletion of the *ARIH2OS* core promoter resulted
247in marked down-regulation of the distal target gene (FDR adjusted p-value
248= 0.02), *NCKIPSD*, identified by long-range chromatin interactions (Fig. 3e)
249with no significant or moderate effect on nearby genes (Extended Data
250Fig. 8h). Importantly, sgRNA-induced mutations in selected eQTLs
251proximal to transcriptional start sites demonstrated significant down-
252regulation effect on distal target genes but no significant effect on nearby
253gene expression in H1-hESC (Fig. 3f; Extended Data Fig. 8i; see Methods).
254Our results strongly suggest genome-wide presence of enhancer-like
255promoters in the human genome and provide additional insight into their
256potential function in distal gene regulation.

257

258The above promoter-centered chromatin interaction maps allowed us to
259infer the target genes of sequences harboring disease-associated variants
260and understand the molecular basis of human disease. We focused on
26142,633 putative disease/trait-associated genetic variants from a recent
262public repository of GWAS catalog¹. Consistent with previous reports^{2,35}, a
263significant portion of SNPs (30%, Fisher's Exact test p value < 2.2e-16)
264were found in putative cREs, emphasizing the importance of target gene

265identification of cREs in functional interpretation of disease associated
266genetic variants. Since the causal SNPs are unknown in most cases, we
267also included SNPs that lie outside the previously defined cREs for further
268analysis. In total, we were able to assign target genes for 27,325 SNPs in
269the list. On average, each SNP was assigned to between 1 and 3
270candidate target genes in each cell/tissue type, with the caveat that the
271precise number of target genes could potentially be affected by the
272modest resolution of our promoter capture strategy and the heterogeneity
273of tissue samples (Extended Data Fig. 9a; Supplementary Table 15; see
274Methods). The above maps therefore provided many more predictions of
275disease-associated genes than using the nearest neighbor gene
276predictions alone (one example is provided for the Parkinson disease in
277Extended Data Fig. 9b, c), with only about 8% of the putative target genes
278inferred from our promoter-centered chromatin interaction maps were
279found to be the closest gene to the sequence variant (Extended Data Fig.
2809d). To evaluate the validity of target predictions based on the promoter-
281centered chromatin interaction maps, we focused on 7 GWAS variants
282that overlap with previously annotated cREs and eQTLs in human
283lymphoblastoid cell line GM12878 cells. We introduced deletions to these
284elements in GM12878 using CRISPR-Cas9 genome editing tools and
285examined the expression of predicted target genes using RT-qPCR in the
286mutant cells and controls. For 5 of the 7 tested cREs, genetic perturbation
287led to down regulation of the predicted distal target genes (Fig. 4a and

288Extended Data Fig. 9e-f). This result supports the target gene predictions
289based on the pHi-C interactions.

290

291Many diseases and traits could be linked to common molecular pathways,
292and the identification of these shared molecular pathways can be
293beneficial in understanding disease pathogenesis and developing
294treatment. To uncover the common molecular pathways underlying
295different diseases and physiological traits, we first determined the
296diseases/traits that share a significant number of common target genes
297predicted from their respective GWAS-associated SNPs. We grouped 687
298traits and diseases into 40 clusters (Fig. 4b; Extended Data Fig. 10a-c;
299Supplementary Table 16; see Methods). Many physiological traits with
300known connections are found to be clustered together. For examples, C5
301clusters oxygen transport related traits together, C6 groups together traits
302related to renal functions, and C20 includes vascular function associated
303traits (Fig. 4b). The above grouping is made possible thanks to the
304promoter-centered chromatin interactome maps, because the similarities
305among related traits observed in Fig. 4b were much less evident when we
306used either GWAS SNPs or nearest genes of the GWAS SNPs to compute
307the similarities as control experiments (Fig. 4c, d, Extended Data Fig.
30810d). Our result suggests the power of target gene identification of GWAS
309variants to uncover trait-trait associations.

310

311To further understand the common molecular pathways affected in
312various human diseases, we carried out gene ontology (GO) analysis for
313the predicted target genes of the GWAS SNPs within each cluster
314(Supplementary Table 17; see Methods). The enriched GO biological
315processes suggest potential shared molecular pathways for disease and
316trait types in each cluster (Fig. 4e, Extended Data Fig. 10e,
317Supplementary Table 18), including unexpected connections between
318specific traits. For example, C39 exposes a link between the susceptibility
319to infectious and autoimmune diseases and the risk of chemotherapeutic
320toxicity by carboplatin and cisplatin. In support of such link, a putative
321target gene associated with the response to carboplatin and cisplatin is
322*ABCF1*, which is involved in inflammatory response³⁶. While speculative,
323the shared molecular pathways uncovered by our analyses may provide
324new leads for investigation of the molecular basis of complex traits and
325disease phenotypes.

326

327In summary, we have generated promoter-centered chromatin
328interactome maps across diverse human cell/tissue types. Our analysis
329covers a broad range of human tissue types and provides prediction of
330target genes for over 70,000 putative *cis*-regulatory elements and 27,000
331GWAS SNP variants. This resource enables a new approach to
332understanding the molecular pathways dysregulated in distinct diseases
333and traits²¹. In future studies, delineation of disease-specific chromatin
334interactions with clinical samples by comparing our reference chromatin

335interaction maps could greatly improve the functional interpretation of
336many disease and trait associated genetic variants.

337

338It should be noted that the current study only surveys a limited number of
339human tissues and cell types, and assigned target genes for a small
340fraction of the putative *cis*-regulatory elements annotated in the human
341genome. Furthermore, the heterogeneous nature of the tissue samples
342used in this study prevents us from accessing the cell types in which the
343identified chromatin interactions occur, except for a few cell lines.

344Nevertheless, this resource lays the ground for further understanding of
345human disease pathogenesis and development of new treatment
346strategies.

347

348

349**Methods**

350**Human tissue samples**

351Esophagus, lung, liver, pancreas, small bowel, sigmoid colon, thymus,
352bladder, adrenal gland, aorta, gastric, left heart ventricle, right heart
353ventricle, right heart atrium, ovary, psoas, spleen, and fat tissues were
354obtained from deceased donors at the time of organ procurement at
355Barnes-Jewish Hospital (St. Louis, USA) as described in our previous
356study²⁶. The same tissue types from different donors were combined
357together during downstream data analysis. Human dorsolateral prefrontal
358cortex (DLPFC rep1) and hippocampus (HC rep1) tissues were obtained
359from the National Institute of Child Health and Human Development
360(NICHHD) Brain Bank for Developmental Disorders. These two samples were
361from a healthy 31-year-old male donor. Ethics approval was obtained from
362the University Health Network and The Hospital for Sick Children for the
363use of these tissues. Another set of human dorsolateral prefrontal cortex
364(DLPFC rep2) and hippocampus (HC rep2) tissues were obtained from the
365Shiley-Marcos Alzheimer's Disease Research Center (ADRC). These two
366samples were from a healthy 80-year-old female donor. Institutional
367Review Board (IRB) approval was obtained from KAIST for the use of these
368tissues.

369

370**Hi-C library on human tissue samples and early embryonic cell**

371**types**

372 Human tissue samples were flash frozen and pulverized prior to
373 formaldehyde cross-linking. Fibroblasts (IMR90) and lymphoblastoid cell
374 lines (GM12878 and GM19240) were cultured and 5 million cells were
375 formaldehyde cross-linked for each Hi-C library. Hi-C was then conducted
376 on the samples as previously described, using *HindIII* for Hi-C library
377 preparation³⁷. Previously constructed Hi-C libraries¹¹ were used for human
378 ES cells (H1) and early embryonic cell types including mesendoderm,
379 mesenchymal stem cell, neural progenitor cells, and trophoblast-like cells.
380

381 **Generation of capture RNA probes**

382 In order to perform Promoter Capture Hi-C, we computationally designed
383 RNA probes that capture promoter regions of previously annotated human
384 protein coding genes. Capture regions were selected for 19,462 well-
385 annotated protein coding gene promoters across 22 autosomes and X
386 chromosome according to GENCODE v19 annotation with confidence level
387 1 and 2. The annotation confidence level 1 and 2 comprise of genes that
388 are accurately annotated with sufficient validation and manual annotation
389 by combining the manual gene annotation from the Human and
390 Vertebrate Analysis and Annotation (HAVANA) group, automatic gene
391 annotation from Ensembl, and validating by CAGE. Due to the variability
392 of capture efficiency, 19,328 promoter regions (99%) were captured in
393 this study. Among them, 18,943 promoter regions were involved in pHi-C
394 interactions in one or more cell/tissue types analyzed in this study. For
395 each transcription start site, the two nearest left hand- and right hand-

396side *HindIII* restriction sites were selected. Six capture oligos were
397designed to be of 120 nucleotide (nt) length and to have 30nt tiling
398overhang. Oligos were designed +/- 300bp upstream and downstream of
399each restriction site. As two restriction sites were chosen for each
400transcription start site, a total of 12 capture oligos were designed to
401target each promoter region. Capture sequences that overlap with directly
402adjacent *HindIII* restriction sites were removed. GC contents of 94%
403capture sequences ranged from 25% to 65%. Some promoters shared the
404same *HindIII* fragment with at least one other, while 14,357 *HindIII*
405fragments (86%) were uniquely assigned to one promoter. The effect of
406the DNA fragments harboring multiple promoters on the quality of our
407analytical findings is modest because only 15% of pHi-C interactions
408emanated from the promoter sharing DNA fragments, and eliminating
409these fragments results in no significant changes in our conclusion for
410both eQTL enrichment test and gene set enrichment analysis. Further,
411strong correlation of GWAS trait associations remains even after excluding
412unresolvable promoters. In total, our capture oligo design generated
413280,445 unique probe sequences including randomly selected capture
414regions (i.e. gene deserts). Single-stranded DNA oligos were then
415synthesized by CustomArray Inc. Single-stranded DNA oligos contained
416universal forward and reverse primer sequences (total length 31nt),
417whereby the forward priming sequence contained a truncated SP6
418recognition sequence that was completed by the overhanging forward
419primer during PCR amplification of the oligos. After PCR, double-stranded

420DNA was converted into biotinylated RNA probes through *in vitro*
421transcription with the SP6 Megascript kit and in the presence of a
422biotinylated UTP, as previously described¹¹.

423

424**Promoter Capture Hi-C library construction**

425Promoter Capture Hi-C library was constructed by performing target
426enrichment protocol (enriching target promoter-centered proximity
427ligation fragments from Hi-C library using capture RNA probes). Briefly, we
428incubated 500ng Hi-C library for 24h at 65°C in a humidified hybridization
429chamber with 2.5ug human Cot-1 DNA (Life Technologies), 2.5ug salmon
430sperm DNA (Life Technologies), and p5/p7 blocking oligos with
431hybridization buffer mix (10X SSPE, 10mM EDTA, 10X Denhardts solution,
432and 0.26% SDS) and 500ng RNA probes. RNA:DNA hybrids were enriched
433using 50ul T1 streptavidin beads (Invitrogen) through 30min incubation at
434RT. Next, bead-bound hybrids were washed through a 15min incubation in
435wash buffer1 (1X SSC and 0.1% SDS) with frequent vortexing, and then
436washed three times with 500ul of pre-warmed (65 °C) wash buffer2 (0.1X
437SSC and 0.1% SDS), then finally resuspended in nuclease-free water. The
438resulting capture Hi-C libraries were amplified while bound to T1 beads,
439and purified using AMPure XP beads, followed by sequencing.

440

441**Promoter Capture Hi-C library sequencing, read alignment, and** 442**off-target read filtering**

443 Promoter Capture Hi-C library sequencing procedures were carried out as
444 previously described according to Illumina HiSeq2500 or HiSeq4000
445 protocols with minor modifications (Illumina, San Diego, CA). Read pairs
446 from Promoter Capture Hi-C library were independently mapped to human
447 genome hg19 using BWA-mem and manually paired with in house script.
448 Unmapped, non-uniquely mapped, and PCR duplicate reads were
449 removed. Trans-chromosomal read pairs and putative self-ligated
450 products (<15kb read pairs) were also removed. Off-target reads were
451 removed when both read pairs did not match the capture probe
452 sequences. The resulting on-target rates in Promoter Capture Hi-C library
453 ranged from 17% to 44% after removing PCR duplicate reads.

454

455 **Promoter Capture Hi-C normalization**

456 Interaction frequencies obtained from Promoter Capture Hi-C were
457 normalized in terms of DNA fragment resolution restricted by *HindIII*. We
458 defined DNA fragments that spanned each *HindIII* restriction site. The
459 start and the end of DNA fragments were defined by taking the midpoint
460 of adjacent upstream and downstream restriction sites, respectively. We
461 merged adjacent DNA fragments if the total length of the DNA fragments
462 was less than 3kb. As a result, 510,045 DNA fragments were defined with
463 a median length of 4.8kb. After that, we calculated raw interaction
464 frequencies at DNA fragment resolution and performed normalization to
465 remove experimental biases caused by intrinsic DNA sequence biases (GC
466 contents, mappability, and effective fragment lengths), RNA probe

467synthesis efficiency bias, and RNA probe hybridization efficiency bias.
468Highly variable RNA probe synthesis efficiency would greatly complicate
469the control of experimental bias. However, if the efficiency bias was
470reproducible, the bias can be computationally removed. To prove such
471bias reproducibility, we performed RNA-seq with two sets of RNA probes
472that were synthesized independently. The RNA-seq results can
473quantitatively measure the amount of synthesized RNA probes, which is
474an indicator of the probe synthesis efficiency. We observed highly
475reproducible RNA-seq results (Pearson Correlation Coefficient = 0.98),
476indicating reproducible probe synthesis efficiency. To address the high
477complexity of different types of experimental biases, we defined a new
478term named “Capturability”, which refers to the probability of the region
479being captured. We assumed that “Capturability” represents all combined
480experimental biases and can be estimated by the total number of capture
481reads spanning a given DNA fragment divided by the total number of
482captured reads in *cis*. We found that “Capturability” in each DNA fragment
483is highly reproducible across samples with 0.95 Pearson correlation
484coefficient between samples on average. Therefore, we defined universal
485“Capturability” as the summation of all “Capturability” defined in each
486sample and normalized raw interaction frequencies by considering
487“Capturability” of two DNA fragments. During normalization, we processed
488promoter-promoter interactions and promoter-other interactions
489independently because promoter regions tend to show very high
490“Capturability” as our capture probes were designed to target promoter

491regions. Also, we only considered promoter-centered long-range
492interactions over 15kb and within 2Mb from TSS of each gene. We
493denoted Y_{ij} to represent the raw interaction frequency between DNA
494fragment i and j and C_i to represent “Capturability” defined in DNA
495fragment i . We assumed Y_{ij} to follow a negative binomial distribution with
496mean μ and variance $\mu + \alpha\mu^2$. Here, $\alpha > 0$ is a parameter to measure the
497magnitude of over-dispersion. We then fitted a negative binomial
498regression model as follows: $\log u_{ij} = \beta_0 + \beta_1 BS(C_i) + \beta_2 BS(C_j)$, where u_{ij} is
499an raw interaction frequency between DNA fragment i and j with coverage
500 C_i and C_j and defined the residual $R_{ij} = Y_{ij} / \exp(\dots)$ as a normalized
501interaction frequency between DNA fragments i and j . BS represents a
502basis vector obtained from B -spline regression, which applied to a vector
503of values of input variable, C , during negative binomial regression model
504fitting for robustness and memory efficient calculation.

505

506 **Identification of P-P and P-O pcHi-C long-range chromatin**

507 **interactions**

508 To identify significant pcHi-C chromatin interactions, we removed distance
509 dependent background signals from normalized interaction frequencies.

510 Here, we assumed that normalized interaction frequency R_{ij} follows a

511 negative binomial distribution with mean μ and variance $\mu + \alpha\mu^2$. Similar to

512 the interaction frequency normalization step above, we calculated the

513 expected interaction frequency at a given distance by fitting it to a

514 negative binomial regression model with basis vectors obtained from B -

4523

46

515 spline regression of distance between two DNA fragments. We denoted E_d
516 to represent the expected interaction frequency at a given distance d
517 calculated from a negative binomial regression model. Distance
518 dependent background signals were removed by taking signal to
519 background ratio as follows: $(R_{ij} + \text{avg}(R)) / (E_d + \text{avg}(R))$, where d
520 indicates distance between DNA fragment i and j . We confirm that the
521 average of normalized interaction frequencies against distance dependent
522 background signals are close to one in all distance, indicating the
523 successful elimination of distance dependent background signals using
524 our method. Next, using 'fitDistr' function in propagate R package we
525 found that 3-parameter Weibull distribution well follows the values of
526 normalized interaction frequencies. Thus, we modeled background
527 distribution of distance normalized interaction frequencies with 3-
528 parameter Weibull distribution. Based on this, significant long-range
529 chromatin interactions are defined when observed interaction frequencies
530 show lower than 0.01 p-value thresholds by fitting distance background
531 removed interaction frequencies with 3-parameter Weibull distribution. To
532 eliminate false pHi-C interactions caused by experimental noise, we
533 applied the criteria of minimum raw interaction frequencies (having more
534 than 5 raw interaction frequencies), which is chosen by investigating
535 reproducibility between biological replicates using lymphoblastoid and
536 mesenchymal stem cell. Note that as the interaction frequencies in pHi-C
537 are mostly zeros or close to zero, the distribution of p-values does not
538 follow the uniform distribution, violating the basic assumption of FDR

539 calculation, which assumes that the null distribution follows uniform (0,1)
540 distribution. Thus, we simulated normalized interaction frequencies that
541 follow 3-parameter Weibull distribution in a sample specific manner, and
542 computed the estimated FDR through multiple permutations. The
543 estimated FDR through multiple permutation (n=1,000) for P-O and P-P
544 pcHi-C interactions is 9% and 5% on average, respectively

545

546 ***in situ* Hi-C experiments and validation of pcHi-C long-range** 547 **chromatin interactions**

548 The visual inspection of normalized interaction frequencies between
549 IMR90 Promoter Capture Hi-C and high resolution IMR90 Hi-C showed high
550 consistency based on manual inspection despite pcHi-C having only 10%
551 sequencing depth compared to high resolution Hi-C (Extended Data Fig.
552 4a). Next, we compared the identified pcHi-C interactions with “loops”
553 defined from IMR90, GM12878, dorsolateral prefrontal cortex, and
554 hippocampus tissues using *in situ* Hi-C experiments (Extended Data Fig.
555 4b-e). Although there is a huge discrepancy between the number of *in situ*
556 Hi-C loops and pcHi-C interactions, we may consider ‘loops’ are a subset
557 of high confident long-range chromatin interactions that involve ‘loop’
558 domains but cannot cover all promoter-mediated long-range chromatin
559 interactions. Loops of IMR90 and GM12878 *in situ* Hi-C result were
560 obtained from previous publication¹³. Loops of dorsolateral prefrontal
561 cortex and hippocampus were identified using HiCCUPS, distributed with
562 Juicer v1.7.6¹³. The loops were called from Knight-Ruiz normalized 5kb,

56310kb, and 25kb resolution data, as these parameters were suggested for a
564medium resolution Hi-C map by the authors of HiCCUPS. As a result, 7,722
565and 8,040 loops were identified from dorsolateral prefrontal cortex and
566hippocampus, respectively. We compared the identified pHi-C long-range
567chromatin interactions to loops of *in situ* Hi-C data and measured the
568reproducibility in terms of ROC curve (receiver operating characteristic
569curve), a plot of the true positive rate against the false positive rate at
570various threshold settings. Here, we set loops as true interactions. We
571ranked all tested pHi-C DNA fragment pairs in terms of p-values and then
572calculated the fraction of true positive and false positive to draw ROC
573curve. We only considered “loops” emanating from promoter-containing
574DNA fragments defined in our Promoter Capture Hi-C result. Each point on
575the ROC curve indicates the true and false positive rate for each 1,000
576false positive interactions. The area under the ROC curve is defined as an
577ROC score and an ROC score of 1 indicates that the rank of DNA fragment
578pairs matched by loops are always higher than all other tested DNA
579fragment pairs according to pHi-C interaction p-values.

580

581 **Reproducibility of pHi-C chromatin interactions between** 582 **biological replicates**

583The reproducibility of pHi-C chromatin interactions between biological
584replicates were measured in terms of ROC curve (Extended Data Fig. 4f).
585Here, we set pHi-C interactions identified in one replicate as true
586interactions. For the other replicate, we ranked all tested DNA fragment

587pairs in terms of p-values and then calculated the fraction of true positive
588and false positive to draw ROC curve. The area under the ROC curve is
589defined as an ROC score and an ROC score of 1 indicates that the rank of
590all pHi-C interactions identified in one replicate is always higher than all
591other tested DNA fragment pairs in another replicate. Due to different
592sequencing depths in each replicate, we first defined true interaction sets
593with one replicate that identified fewer number of pHi-C interactions than
594the other replicate, then tested how these true interactions were well
595detected in the other replicate. Both P-P and P-O interactions were
596combined together for calculating ROC scores. Each dot in ROC curve
597indicates the true positive rate at the corresponding false positive rate
598with increment of 1% of false positive rate. We tested biological replicates
599in the following 12 tissue/cell types: aorta (AO2/AO3, ROC score=0.79),
600lung (LG1/LG2, ROC score=0.80), small bowel (SB1/SB2, ROC
601score=0.83), spleen (SX1/SX3, ROC score=0.80), dorsolateral prefrontal
602cortex (FC_rep1/FC_rep2, ROC score=0.92), left ventricle (LV1/LV3, ROC
603score=0.85), mesenchymal stem cell (MSC_rep1/MSC_rep2, ROC
604Score=0.99), hippocampus (HC_rep1/HC_rep2, ROC score=0.81), gastric
605(GA2/GA3, ROC score=0.91), lymphoblastoid cell lines
606(GM12878/GM19240, ROC score=0.98), right ventricle (RV1/RV3, ROC
607Score=0.83), and pancreas (PA2/PA3, ROC score=0.73). Indeed, we
608calculated Spearman's rank correlation of p values between replicates and
609found that the average Spearman's rank correlation was around 0.40.

610

611 **Enrichment of pHi-C interactions regarding TAD, boundary, and** 612 **unorganized regions**

613 The TAD annotations for 22 samples by DomainCaller¹⁴ with 2MB windows
614 size were downloaded from the 3DIV database³⁸. The regions between
615 TADs were classified as “unorganized” when the gap is longer than 400kb,
616 otherwise, the remaining regions were classified as “boundary”. Then, the
617 types of pHi-C interactions were classified based on the location of DNA
618 fragment’s centroid.

- 619 1. “Within TAD”, if both fragments’ centroids are located in the
620 identical TAD.
- 621 2. “Within unorganized region”, if both fragments’ centroids are
622 located in the identical unorganized region.
- 623 3. “Between different TADs”, if one fragment’s centroid is located in a
624 TAD while another fragment’s centroid is located in a different TAD.
- 625 4. “Between TAD and boundary”, if one fragment’s centroid is located
626 in a TAD while another fragment’s centroid is covered by boundary
627 region.
- 628 5. “Between TAD and unorganized region”, if one fragment’s centroid
629 is located in a TAD while another fragment’s centroid is located in
630 an unorganized region.

631

632 **Annotation of ChromHMM 18-chromatin state to DNA fragments**

633 The pre-calculated chromatin state annotations were downloaded from
634 the 18-state ChromHMM model established by Roadmap Epigenomics
635 Project. As the genomic proportion of promoter and enhancer regions are
636 relatively low, we assigned the chromatin states to DNA fragments based
637 on the following priority order (TssA-EnhA1-EnhA2-TssFlnk-TssFlnkU-

638TssflnkD-EnhG1/G2-EnhWk-TssBiv-Enhbiv). For instance, the chromatin
639state of a fragment was assigned as TssFlnkU, if the fragment contained
640two annotations TssFlnkU and EnhWk. EnhG1 and EnhG2 annotations
641were merged because of their low occurrence percentage. We considered
642two promoter types (TssA and TssBiv) according to ChromHMM
643annotations and investigated the preference of their interacting partners.
644For each promoter type, the occurrence of each chromatin status at
645interacting DNA fragments was divided by the total number of interacting
646DNA fragments. This fraction value of each chromatin status was
647normalized against the genomic fraction of each chromatin status. KS-test
648was performed to measure the statistical significance of each chromatin
649status at interacting DNA fragments between TssA and TssBiv promoters.
650

651**Analysis with a 50-chromatin-state ChromHMM model**

652To supplement our analysis with the ChromHMM 18-chromatin state
653model, we conducted in-depth investigations with 5 samples, including H1
654embryonic stem cell, mesendoderm, mesenchymal stem cell, trophoblast,
655and IMR90, using a 50-state ChromHMM model produced by the Roadmap
656Epigenomics Project³⁵. The ChromHMM model utilized combination of 29
657chromatin marks to generate a 50-state ChromHMM model. To be
658consistent with the 18-state ChromHMM model, we used the same
659definition for TssA and TssBiv promoter containing fragments, but
660chromatin state of their interacting partners was further refined using the

66150-state ChromHMM model. The statistical test was performed as
662described in the analysis with the 18-chromatin-state ChromHMM model.
663

664**Identification of extensively interacting DNA fragments**

665In order to identify DNA fragments that showed extensive long-range
666interactions with multiple promoters, we systematically defined these
667promiscuously interacting DNA fragments from P-P pHi-C interaction
668maps and P-O pHi-C interaction maps, respectively. For each cell or
669tissue-type, we selected frequently interacting DNA fragments with
670multiple promoters in terms of 0.01 Poisson p value cutoff.

671

672**Identification of TF clusters in H1-hESC and GM12878**

673Transcription factor ChIP-seq datasets on human lymphoblastoid cell lines
674(GM12878) and human embryonic stem cell (H1-hESC) were collected
675from ENCODE. These ChIP-seq reads were aligned against human genome
676hg19 using BWA-mem with default parameters. Non-uniquely mapped,
677low quality (MAPQ<10), and PCR duplicate reads were removed. Peak
678calling of individual ChIP-seq experiments was performed with MACS2
679callpeak with default parameters³⁹. We defined TF clusters by calling
680peaks from combined bed files of TF peaked regions using MACS2
681bdgpeakcall. The regions occupied by multiple TF peaks were recognized
682as TF clusters. To remove parameter dependent bias, we retrieved TF
683clusters 40 times with various parameter sets as following; minimum # of
684TFs within cluster (5 or 10), minimum length of cluster from 100bp to

6851600bp, and maximum gap length within cluster from 100bp to 51.2kb.

686Final TF clusters were defined when the region was detected as TF

687clusters more than 50 times from 100 different parameter sets.

688

689**Enrichment analysis of TF clusters and super-enhancers**

690In order to calculate the enrichment of TF clusters or super-enhancers at

691extensively interacting DNA fragments (EIF), we counted the number of

692matched TF clusters and super-enhancers. The list of super-enhancers

693was obtained from the 3DIV database³⁸. Permutation test was performed

694to calculate the expected values. Using Bedtools shuffleBed, we

695generated random genomic locations that resemble actual TF clusters

696with the same size but different genomic coordinates. Bedtools

697intersectBed identified any overlap between EIF and TF clusters or random

698genomic locations. Standard deviations of error bars in the random

699genomic locations were calculated from 10,000 random data sets. In order

700to test the enrichment of TF clusters compared to typical TF peaks, we

701generated random genomic locations that resemble actual TF clusters

702with the same size but different genomic coordinates matched to typical

703TF peaks. Standard deviations of error bars in the typical TF peaks were

704calculated from 10,000 random data sets. Similarly, enrichment analysis

705of super-enhancers was conducted by generating random genomic

706locations of the same size as super-enhancers but at different genomic

707coordinates. We also conducted the enrichment test with typical

708enhancers. We revealed that P-O EIFs highly co-exist with super-enhancer

709regions, rather than typical enhancers and genomic background for most
710of the samples, except two samples, lymphoblastoid cell lines and gastric
711tissue. Note that half of lymphoblastoid P-O EIFs are co-occupied with
712typical enhancers that are classified as super-enhancers in other
713cell/tissue types.

714

715**Comparison between eQTL associations and P-O interactions**

716In order to test the enrichment for P-O pHi-C chromatin interactions in
717significant eQTL associations, we compared P-O pHi-C interactions to
718significant eQTL associations in the matching tissue types. The eQTL
719associations were downloaded directly from GTEx Portal (downloaded on
720Nov. 10th, 2017) for all matching tissue types (n=14, adrenal gland, aorta,
721dorsolateral prefrontal cortex, brain hippocampus, sigmoid colon,
722esophagus, left heart ventricle, liver, lung, ovary, pancreas, small
723intestine terminal ileum for small bowel, spleen, and stomach for gastric).
724First, the significant eQTLs defined by GTEx (q-value ≤ 0.05) were filtered
725so that only the eQTL variants within the fragments that involve P-O pHi-
726C interactions remain for comparison. Then, we removed pHi-C
727interactions beyond 1Mb in distance to match the range of eQTL
728association, and discarded eQTL associations with distance below 15kb to
729match the valid interaction cutoff. The filtered, significant eQTL
730associations were compared with pHi-C and randomized interactions in
731the same condition. Here, we only considered P-O pHi-C interactions with
732DNA fragments that do not harbor multiple promoters. For the random

733 expectation, we generated a simulated pHi-C interaction pool by creating
734 all possible combinations of DNA fragments with no TSS and the protein
735 coding genes that exist within the distance range. The pHi-C interactions
736 that exist in any of the tissue/cell type were removed from the control
737 interaction pool for the enrichment analysis. To avoid variation caused by
738 the difference in distance between pHi-C interactions and eQTL
739 associations, we created distance matched control, in which the number
740 of pHi-C interactions was stored at the interval of 40kb, and the same
741 number of interactions was drawn randomly from the control interaction
742 pool. The number of randomized interactions drawn from each
743 chromosome was matched to the pHi-C interactions. The standard
744 deviation was obtained by permuting the random expectation with 1,000
745 iterations and was used to calculate the statistical confidence

746

747 To illustrate the filtering process of the eQTL data, for example, the
748 549,763 significant eQTLs in adrenal gland were reduced to 237,181 after
749 collecting eQTLs located in the DNA fragments without TSS and discarding
750 eQTL association with the distance below 15kb and with a pseudogene
751 target. This filtered set of significant eQTL associations was used for
752 enrichment test for both pHi-C and randomized interactions. The number
753 of total tested significant eQTL association, 19,996 in case of adrenal
754 gland, in Supplementary Table 11, indicates the number of significant
755 eQTLs located in the DNA fragments that are associated with the pHi-C
756 interactions in the corresponding cell/tissue type.

757

**758 Variations in H3K27ac signals at promoters and cREs connected
759 by P-O interactions**

760 We conducted correlation analysis of H3K27ac signals across all available
761 cell/tissue types for each promoter-cRE pair connected by P-O interactions
762 in at least one cell/tissue type analyzed. First, we defined putative distal
763 *cis*-regulatory elements (cREs) marked by H3K27ac peaks across all 27
764 cell/tissue types. We merged these elements if the peaks are within 3kb of
765 each other, then we defined cRE-containing DNA fragment when the DNA
766 fragment harbors at least one *cis*-regulatory element. When a DNA
767 fragment contained both TSS and cRE, we defined the fragment as a
768 promoter-containing DNA fragment instead of a cRE-containing DNA
769 fragment because our experiment is designed to target promoter regions.
770 We used input normalized H3K27ac RPKM values by taking log₂
771 transformation as H3K27ac signals at promoters and cREs. Pearson
772 correlation coefficient values were calculated for each promoter-cRE pair
773 connected by pHi-C interactions after excluding cREs spanning adjacent
774 DNA fragments and visualized as a box plot. Random expectation values
775 were calculated after randomization of H3k27ac signals at promoters and
776 cREs. Distance matching random expectation values were calculated after
777 random selection of cRE-promoter pairs by controlling distance
778 information as same as identified cRE-promoter pairs.

779

780 Analysis of H3K27ac signals at cREs and expression of target
781 genes connected by cell/tissue specific cRE-promoter pairs

782 In order to investigate cell/tissue-specific cRE-promoter pairs, for each
783 cell/tissue-type unique cRE-promoter pairs were collected and then
784 distance normalized pHi-C interaction frequencies of corresponding P-O
785 pHi-C interactions were obtained across all cell/tissue types. We only
786 considered a unique P-O interaction pair when multiple cREs are located in
787 a same DNA fragment and target a same promoter. The cell/tissue-
788 specific cRE-promoter pairs exhibit strong enrichment of pHi-C
789 interaction frequencies in the corresponding cell/tissue type but depleted
790 in other cell/tissue types, validating the cell/tissue-specificity of these cRE-
791 promoter pairs. Statistical significance of pHi-C interaction frequencies
792 was tested by conducting KS-test between mean of pHi-C interaction
793 frequencies in the matched cell/tissue types (values in diagonal in Fig. 2f)
794 and those in other cell/tissue types (values in off diagonal in Fig. 2f).

795

796 In order to investigate cell/tissue-specific activity of cREs connected by
797 cell/tissue-specific cRE-promoter pairs, we identified group of cREs that
798 are connected by unique cRE-promoter pairs for each cell/tissue type.
799 After that, H3K27ac signals were calculated for each cRE across all
800 cell/tissue types and these values were normalized by taking z-score
801 transformation to obtain relative H3K27ac enrichment signals. The mean
802 values of normalized H3K27ac signals were calculated for each group of
803 cREs in each cell/tissue type. KS-test was performed between the mean

804 values in the corresponding cell/tissue types (values in diagonal in Fig. 2g)
805 and those in other cell/tissue types (values in off diagonal in Fig. 2g).

806

807 In order to investigate expression levels of target genes connected by
808 cell/tissue-specific cRE-promoter pairs, we first defined a group of genes
809 that are connected by unique cRE-promoter pairs more than twice for
810 each cell/tissue-type. After that, gene expression levels (FPKM) were
811 calculated for each gene across all cell/tissue types. Relative gene
812 expression levels were obtained by taking z-score transformation for each
813 gene across all cell/tissue types. The mean values of z-score transformed
814 FPKM values were calculated for each group of genes in each cell/tissue
815 type. KS-test was performed between the mean values in the
816 corresponding cell/tissue types (values in diagonal in Fig. 2h) and those in
817 other cell/tissue types (values in off diagonal in Fig. 2h).

818

819 **Comparison between eQTL associations and P-P chromatin** 820 **interactions**

821 In order to assess the enrichment for promoter-promoter pHi-C
822 interactions in the significant eQTL associations, we computed the number
823 of P-P pHi-C interactions matched to the significant eQTL associations
824 (downloaded on Nov. 10th, 2017). For the tested tissue types (n=13,
825 adrenal gland, aorta, dorsolateral prefrontal cortex BA9, hippocampus,
826 sigmoid colon, left ventricle, liver, lung, ovary, pancreas, small intestine
827 terminal ileum for small bowel, spleen, and stomach for gastric), we

828 considered only the eQTLs that are located within 2.5kb from a TSS of a
829 protein coding gene. For accurate comparison, we removed P-P chromatin
830 interactions beyond 1Mb in distance to match the range of eQTL
831 association, and discarded eQTL associations with distance below 15kb to
832 match the valid interaction cutoff. Finally, the significant eQTLs were
833 filtered by collecting only the eQTLs within the fragments that involve P-P
834 pcHi-C interactions in the corresponding cell/tissue and by removing
835 eQTLs that target pseudogenes. Then, the number of filtered significant
836 eQTLs that match P-P pcHi-C interactions was counted. The DNA
837 fragments that harbor multiple promoters were removed from the
838 analysis. For the random expectation, we created a control pool of all P-P
839 pairs within the range of 15kb to 1Mb, selected the same number of
840 random P-P pairs as used in significant eQTL comparison, and counted the
841 matched number of random P-P pairs with P-P pcHi-C interactions. The P-P
842 pcHi-C interactions that exist in any of the tissue/cell type were removed
843 from the control interaction pool for the enrichment analysis. In addition,
844 to avoid variation caused by the difference in distance between pcHi-C
845 interactions and eQTL associations, we created distance matched control,
846 in which the number of pcHi-C interactions was stored at the interval of
847 40kb, and the same number of interactions was drawn randomly from the
848 randomized interaction pool. In addition, the number of randomized
849 interactions drawn from each chromosome was matched to the pcHi-C
850 interactions. The standard deviation was obtained by permuting the

851 random expectation with 1,000 iterations and was used to calculate the
852 statistical significance.

853

854 **Visualization of eQTL-supported P-P and P-O chromatin**

855 **interactions**

856 The pcHi-C interactions that matched significant eQTLs were visualized by
857 LocusZoom⁴⁰. We collected and merged significant and all tested eQTLs
858 for each tissue type and extracted the relevant p-values and SNP IDs for a
859 queried gene. Then, LocusZoom was run with default parameters to show
860 the pcHi-C interaction and its eQTL associations surrounding the region.

861

862 **Experimental validation of enhancer-like function of promoters**

863 H1-hESC was cultured in mTeSR1 medium on Matrigel coated plates³³. To
864 knockout the core promoter regions of *ZNF891* (chr12:133706994) and
865 *ARIH2OS* (chr3:48956862) in H1-hESC, we utilized CRISPR/Cas9 RNP
866 method as previously described by Diao, et al.³³. Briefly, we used *in vitro*
867 synthesized CRISPR crRNA and CRISPR tracrRNA (IDT) with the sequences
868 specified below.

869 *ZNF891* sgRNAs 5p-1: GCGTCCGTGACGCACAGACC

870 *ZNF891* sgRNAs 5p-2: GACCAGGCCCTCTGCGGGG

871 *ZNF891* sgRNAs 3p-1: AGGCTGGGGCGCGTGCCTAA

872 *ZNF891* sgRNAs 3p-2: GTGCGTAACGGTGTGTGTTG

873 *ZNF891* genotyping primer 5p: GTCCTCAGTGCCTGCCTC

874 *ZNF891* genotyping primer 3p: CAGCAACAGCAAACAGAGAAC

875 *ARIH2OS* sgRNAs 5p-1: GCTCCCAAAGATGACTCGAG

876 *ARIH2OS* sgRNAs 5p-2: GACTCGAGTGGTGAGCCCCG

877 *ARIH2OS* sgRNAs 3p-1: GGAGAAGTCATCCAAGAACG

7538

76

878ARIH2OS sgRNAs 3p-2: CGCTATGACAGAAAGTTCTA
879ARIH2OS genotyping primer 5p: CATCTAGGCCCTCTCTCCCT
880ARIH2OS genotyping primer 3p: TCAGCAATTTTCGTTTCAAATC
881

882Each of the core promoter was knocked-out by two sets of sgRNA pairs to
883avoid the potential off-target effect caused by CRISPR/Cas9 genome
884editing. The Cas9 recombinant protein was purchased from NEB (Cat
885M0386M) and the Cas9/crRNA/tracrRNA was assembled *in vitro* by
886following the previously described protocol³³. The RNP complex was
887electro-transfected into POU5F1-eGFP hESC reporter line with Neon
888Transfection System 10 μ l kit (ThermoFisher Scientific, Cat#: MPK1096)
889with the default electrotransfection protocol. Seven days after post-
890transfection, individual colonies were picked and expanded, followed by
891genotyping and in-depth analysis. After genotype validation, we
892performed RNA-seq using Ovation[®] RNA-Seq System V2 (NuGEN) as
893previously described⁴¹.

894

895**RNA-seq data analysis between WT and mutants upon promoter** 896**deletion**

897Raw RNA-seq fastq files were aligned to the reference genome (hg19)
898using BWA-mem. Duplicate reads were discarded with Picard to avoid any
899artifact caused by the amplification step originated from Ovation[®] RNA-
900Seq System V2 (NuGEN). Then, FPKM values were calculated using
901Cufflinks with GENCODE v19 annotation. Reproducibility between
902biological replicates were measured (PCC of FPKM for WT = 0.98, *ZNF891*

903promoter deletion clone #1 =0.99, *ZNF891* promoter deletion clone #2 =
9040.99, *ZNF891* promoter deletion clone #3=0.99 , and *ARIH2OS* promoter
905deletion clone #1 =0.98). FPKM values of *ZNF84* and *NCKIPSD* were
906investigated as distal target genes of *ZNF891* and *ARIH2OS*, respectively,
907between mutant and WT to test the effect of deletion of core promoters
908on distal target genes.

909

910**Experimental validation of promoter-proximal eQTL distal target** 911**genes**

912In order to validate the distal target genes of promoter-proximal eQTLs
913identified by pHi-C results, we designed sgRNA sequences targeting +/-
9145bp of the eQTLs in H1-hESC and cloned the sgRNAs into lentiCRISPRv2
915backbone, followed by lentiviral preparation, infection and Puromycin
916selection as previously described³³. Two weeks after the infection, single
917clones were selected and genotyped to confirm the mutations on the
918targeted eQTL sites (genotyping PCR results are listed in the oligo file).
919Total RNA was purified from each single clone and subjected to RT-qPCR
920analysis as previously described (Genotyping PCR primers are listed in the
921oligo file)³³. To conduct statistical analysis, two separate sgRNAs were
922generated, which target the same eQTL. Then, three clones were isolated
923and cultured for a single sgRNA in order to induce the knockout, and each
924of these clones was considered as a biological replicate. Each clone was
925consisted of technical triplicates for the stable measurement of the
926expression during RT-qPCR experiment.

927

928chr9:139305041_1 sgRNA in H1-hESC: GCCTTGGGCCGTCGGCGAGGGGG
929chr9:139305041_2 sgRNA in H1-hESC: TGGGCCGTCGGCGAGGGGGAGGG
930chr17:18128865_1 sgRNA in H1-hESC: GCGGGGCCGGGCCTGCACGGGGG
931chr17:18128865_2 sgRNA in H1-hESC: CGCGCGGGGCCGGGCCTGCACGG
932chr14:104029246_1 sgRNA in H1-hESC: CGAAGCCCGAGGAAGCGCGGGCGG
933chr14:104029246_2 sgRNA in H1-hESC: CGGCAGGGTCGCGAAGCCCGAGG
934chr3:184032262_1 sgRNA in H1-hESC: GGCAAATCCCATGTGCTCGGCGG
935chr3:184032262_2 sgRNA in H1-hESC: GGGGGCAAATCCCATGTGCTCGG

936

937chr9:139305041_F genotyping primer: CGCTGGTAGCCCGACATC
938chr9:139305041_R genotyping primer: CCCCCTTCAGTCGTCAC
939chr17:18128865_F genotyping primer: CCCAGTTCACCATTGTCTGG
940chr17:18128865_R genotyping primer: AACCGAACTTCATCATCTTGC
941chr14:104029246_F genotyping primer: GAGGCAGCCTGGAGTGAC
942chr14:104029246_R genotyping primer: GAGAAAGGTCTTCTTCCCGG
943chr3:184032262_F genotyping primer: AATGAACTAAAGAATCGCGGAA
944chr3:184032262_R genotyping primer: CACAGACGTAGTCCACAACCAT

945

946**Experimental validation of distal target genes for disease-**

947**associated genetic variants**

948In order to validate the distal target genes of disease-associated genetic
949variants (GWAS-SNPs) identified by pHi-C results, we designed sgRNA
950sequences targeting +/- 5bp of the GWAS-SNPs in lymphoblastoid cells,
951and cloned the sgRNAs into lentiCRISPRv2 backbone as described above.
952Two weeks after the infection, single clones were selected and genotyped
953to confirm the mutations on the targeted GWAS-SNP sites (genotyping
954PCR results are listed in the oligo file). Total RNA was purified from each
955single clone and subjected to RT-qPCR analysis as previously described
956(Genotyping PCR primers are listed in the oligo file)³³. To conduct

8141

82

957 statistical analysis, two separate sgRNAs were generated, which target
958 the same GWAS SNP. Then, two clones were isolated and cultured for a
959 single sgRNA in order to induce the knockout, and each of these clones
960 was considered as a biological replicate. Each clone was consisted of
961 technical triplicates for the stable measurement of the expression during
962 RT-qPCR experiment.

963

964 chr5.96297527 sgRNA in GM12878: TGCCATTCAGTCTATAGATCTGG
965 chr17.38032460 sgRNA in GM12878: TGGGCTTTGGCTGGGCGCAGTGG
966 chr17.38023745 sgRNA in GM12878: GGGCTCCATCCCTACAGAAAAGG
967 chr3.52707026 sgRNA in GM12878: GAGTTTTGCTCTTATTGTCCAGG
968 chr3.52703615 sgRNA in GM12878: AGTTATTACAAATAACATCATGG
969 chr3.52728804 sgRNA in GM12878: TCCTGGAAGATAGCATGCGTGGG
970 chr3.52706724 sgRNA in GM12878: GGTCTCGAACTCCTGCACTCAGG

971

972 chr5: 96297527_F genotyping primer: ACCAGTTTACACGAATCATCCC
973 chr17:38032460_F genotyping primer: TAGAGACAGAGTTTCGCCCTGT
974 chr17:38023745_F genotyping primer: TGGGCTCTCTCTACTAACCAGC
975 chr3:52707026_F genotyping primer: TGACAGCAAGAGAGGAAAGATG
976 chr3:52703615_F genotyping primer: TCAAATGAAGTTCCAGGAGACA
977 chr3:52728804_F genotyping primer: ACTTGTAAGGCAGATGGAGAC
978 chr3:52706724_F genotyping primer: GTTCAAGTGATTCTCCTGCCTC
979 chr5: 96297527_R genotyping primer: ACTTCATCATGGGCAGTAAACC
980 chr17:38032460_R genotyping primer: AGGACCATTCTGTTTTCTTCA
981 chr17:38023745_R genotyping primer: GTGACCTTGCTTTAAAATGGG
982 chr3:52707026_R genotyping primer: AGGTGGGAGAATTGCTTGAAC
983 chr3:52703615_R genotyping primer: AACCTGTCAGCTAAGGTTCCAA
984 chr3:52728804_R genotyping primer: GCAAATTCAACCTAATCCGAAG
985 chr3:52706724_R genotyping primer: ATGCCTGTAATCCCAACACTTT

986

987 **Extended GWAS-SNPs with high LD structure**

988GWAS-SNPs were obtained from GWAS catalogue database (version1.0.1,
989downloaded on February 2018) and selected with p-value cutoff of 10^{-5}
990with minor manual curations. As GWAS-SNPs obtained from GWAS catalog
991database contain tag SNP information only, we extended the GWAS-SNP
992information using linkage disequilibrium (LD) structure. LD scores were
993calculated using PLINK for five different populations obtained from 1000
994genome phase 3 data. For each tag SNP, we included all associated SNPs
995that had tight LD scores (>0.8) across all five populations (AFR, AMR, EAS,
996EUR, and SAS). With the p-value cutoff of 10^{-5} , we collected 42,674
997significant GWAS-SNPs across 2,310 GWAS mapped traits and expanded
998this list to 180,893 by including LD information. Then, putative target
999genes of GWAS-SNPs were identified by aggregating all unique pHi-C
1000interactions. We noted that the cutoff value of high LD association is
1001arbitrarily determined by considering a stringent cutoff value presented in
1002a set of previous studies to minimize additional noise in the data analysis.

1003

1004**Enrichment test of disease genes in putative GWAS-SNP target** 1005**genes**

1006The list of putative disease associated genes was downloaded from
1007GeneCard database, obtaining 9,989 disease-associated genes. Then, we
1008defined putative target genes of GWAS-SNPs associated with Parkinson
1009disease by using pHi-C interactions or the nearest gene information,
1010respectively. Then, we counted the number of matched disease-
1011associated genes in each set of putative GWAS-SNP target genes.

1012

1013 **Clustering of GWAS mapped traits based on putative target gene** 1014 **similarities**

1015 The “mapped traits” were obtained from GWAS catalog database
1016 (version 1.0.1, downloaded on February 2018), and paired with their
1017 corresponding GWAS SNPs. Then, putative target genes for each GWAS
1018 SNP were obtained by the unique and aggregated pHi-C interactions.
1019 After defining putative target genes and their target frequency for each
1020 trait, we constructed a 1442 by 1442 correlation matrix where each entry
1021 indicates a similarity score between the mapped traits in terms of the
1022 Pearson correlation coefficient (PCC), for which only the traits with a total
1023 gene count greater than 5 were considered. The correlation matrix was
1024 subjected to K-means clustering (n=30) using Euclidean distance, and the
1025 cluster containing ungrouped terms was excluded in further analysis to
1026 eliminate miscellaneous terms. To avoid having a predetermined number
1027 of clusters, the remaining 687 traits were rearranged in a correlation
1028 matrix in terms of their hierarchical relationship (Pearson uncentered and
1029 complete linkage). The final hierarchically clustered correlation matrix
1030 showed a clear organization of 40 clusters with a threshold of dendrogram
1031 height, 0.9. Fig. 4c was drawn by using the nearest gene of GWAS SNPs.
1032 After defining the list of nearest genes for each mapped trait, we again
1033 measured the similarity between the mapped traits by calculating the
1034 Pearson correlation coefficient. We presented similarity values between
1035 the mapped traits as in the same order of mapped traits in Fig. 4b.

1036 Similarly, Fig. 4d was drawn by using the GWAS SNPs alone. We measured
1037 the similarity of the mapped traits by calculating the Pearson correlation
1038 coefficient between GWAS SNPs of each trait, and presented the values as
1039 in the same order of mapped traits in Fig. 4b.

1040

1041 **Analysis of functional enrichment using DAVID**

1042 To identify the enriched biological pathways in the GWAS mapped traits
1043 for the clusters, we extracted putative target genes associated with each
1044 cluster. Then, we performed Gene Ontology (GO) analysis using DAVID
1045 (6.8 version) to obtain the list of enriched biological pathways for each
1046 cluster with a cutoff p-value of 10^{-3} by using the GO_BP annotation
1047 selection. After that we constructed 40 (number of clusters) by 126
1048 (number of GO_BP annotations) matrix where each entry indicates
1049 $-\log_{10}(\text{p-value})$ of corresponding GO_BP annotation. Next, we performed
1050 hierarchical clustering in respect to the enriched biological pathways by
1051 Pearson correlation matrix and average linkage parameter. In
1052 Supplementary Table 17, we presented GO_BP annotation information.

1053

1054 To see the effect of multiple TSS co-existing in a DNA fragment during
1055 gene set enrichment analysis, we calculated the number of genes that are
1056 located in the defined DNA fragments for all genes and the genes in
1057 cluster 38. To see the effect of fragment-sharing TSS of genes on the
1058 enriched biological pathways, we submitted the genes in cluster 38 for
1059 enriched pathway analysis using three different queries; 1) total genes in

1060the cluster, 2) random selection of genes in case of fragment-sharing, and
10613) after removal of the fragment-sharing genes, as illustrated in
1062Supplementary Table 18. We did not observe any significant effect on
1063gene set enrichment analysis caused by promoters shared by the same
1064*Hind*III fragment with at least one other promoter.

1065

1066

1067

1068

1069

1070References

10711. Welter, D. *et al.* The NHGRI GWAS Catalog, a curated resource of SNP-trait
1072 associations. *Nucleic Acids Res* **42**, D1001-6 (2014).
10732. Maurano, M.T. *et al.* Systematic localization of common disease-associated
1074 variation in regulatory DNA. *Science* **337**, 1190-5 (2012).
10753. Hindorf, L.A. *et al.* Potential etiologic and functional implications of genome-wide
1076 association loci for human diseases and traits. *Proc Natl Acad Sci U S A* **106**,
1077 9362-7 (2009).
10784. Lettice, L.A. *et al.* A long-range Shh enhancer regulates expression in the
1079 developing limb and fin and is associated with preaxial polydactyly. *Hum Mol*
1080 *Genet* **12**, 1725-35 (2003).
10815. Uslu, V.V. *et al.* Long-range enhancers regulating Myc expression are required for
1082 normal facial morphogenesis. *Nat Genet* **46**, 753-8 (2014).
10836. Claussnitzer, M. *et al.* FTO Obesity Variant Circuitry and Adipocyte Browning in
1084 Humans. *N Engl J Med* **373**, 895-907 (2015).
10857. Smemo, S. *et al.* Obesity-associated variants within FTO form long-range
1086 functional connections with IRX3. *Nature* **507**, 371-5 (2014).
10878. Yu, M. & Ren, B. The Three-Dimensional Organization of Mammalian Genomes.
1088 *Annu Rev Cell Dev Biol* **33**, 265-289 (2017).
10899. de Wit, E. *et al.* The pluripotent genome in three dimensions is shaped around
1090 pluripotency factors. *Nature* **501**, 227-31 (2013).
109110. Sanyal, A., Lajoie, B.R., Jain, G. & Dekker, J. The long-range interaction landscape
1092 of gene promoters. *Nature* **489**, 109-13 (2012).
109311. Dixon, J.R. *et al.* Chromatin architecture reorganization during stem cell
1094 differentiation. *Nature* **518**, 331-6 (2015).
109512. Jin, F. *et al.* A high-resolution map of the three-dimensional chromatin interactome
1096 in human cells. *Nature* **503**, 290-4 (2013).
109713. Rao, S.S. *et al.* A 3D map of the human genome at kilobase resolution reveals
1098 principles of chromatin looping. *Cell* **159**, 1665-80 (2014).
109914. Dixon, J.R. *et al.* Topological domains in mammalian genomes identified by
1100 analysis of chromatin interactions. *Nature* **485**, 376-80 (2012).
110115. Tang, Z. *et al.* CTCF-Mediated Human 3D Genome Architecture Reveals Chromatin
1102 Topology for Transcription. *Cell* **163**, 1611-27 (2015).
110316. Sahlen, P. *et al.* Genome-wide mapping of promoter-anchored interactions with
1104 close to single-enhancer resolution. *Genome Biol* **16**, 156 (2015).
110517. Jager, R. *et al.* Capture Hi-C identifies the chromatin interactome of colorectal
1106 cancer risk loci. *Nat Commun* **6**, 6178 (2015).
110718. Mifsud, B. *et al.* Mapping long-range promoter contacts in human cells with high-
1108 resolution capture Hi-C. *Nat Genet* **47**, 598-606 (2015).
110919. Dryden, N.H. *et al.* Unbiased analysis of potential targets of breast cancer
1110 susceptibility loci by Capture Hi-C. *Genome Res* **24**, 1854-68 (2014).
111120. Martin, P. *et al.* Capture Hi-C reveals novel candidate genes and complex long-
1112 range interactions with related autoimmune risk loci. *Nat Commun* **6**, 10069
1113 (2015).
111421. Javierre, B.M. *et al.* Lineage-Specific Genome Architecture Links Enhancers and
1115 Non-coding Disease Variants to Target Gene Promoters. *Cell* **167**, 1369-1384 e19
1116 (2016).
111722. Freire-Pritchett, P. *et al.* Global reorganisation of cis-regulatory units upon lineage
1118 commitment of human embryonic stem cells. *Elife* **6**(2017).
111923. Siersbaek, R. *et al.* Dynamic Rewiring of Promoter-Anchored Chromatin Loops
1120 during Adipocyte Differentiation. *Mol Cell* **66**, 420-435 e5 (2017).
112124. Rubin, A.J. *et al.* Lineage-specific dynamic and pre-established enhancer-promoter
1122 contacts cooperate in terminal differentiation. *Nat Genet* **49**, 1522-1528 (2017).
112325. Orlando, G. *et al.* Promoter capture Hi-C-based identification of recurrent
1124 noncoding mutations in colorectal cancer. *Nat Genet* (2018).

112526. Leung, D. *et al.* Integrative analysis of haplotype-resolved epigenomes across
1126 human tissues. *Nature* **518**, 350-354 (2015).
112727. Schmitt, A.D. *et al.* A Compendium of Chromatin Contact Maps Reveals Spatially
1128 Active Regions in the Human Genome. *Cell Rep* **17**, 2042-2059 (2016).
112928. Thurman, R.E. *et al.* The accessible chromatin landscape of the human genome.
1130 *Nature* **489**, 75-82 (2012).
113129. Whyte, W.A. *et al.* Master transcription factors and mediator establish super-
1132 enhancers at key cell identity genes. *Cell* **153**, 307-19 (2013).
113330. Consortium, G.T. Human genomics. The Genotype-Tissue Expression (GTEx) pilot
1134 analysis: multitissue gene regulation in humans. *Science* **348**, 648-60 (2015).
113531. Zhang, Y. *et al.* Chromatin connectivity maps reveal dynamic promoter-enhancer
1136 long-range associations. *Nature* **504**, 306-10 (2013).
113732. Rajagopal, N. *et al.* High-throughput mapping of regulatory DNA. *Nat Biotechnol*
1138 **34**, 167-74 (2016).
113933. Diao, Y. *et al.* A tiling-deletion-based genetic screen for cis-regulatory element
1140 identification in mammalian cells. *Nat Methods* **14**, 629-635 (2017).
114134. Engreitz, J.M. *et al.* Local regulation of gene expression by lncRNA promoters,
1142 transcription and splicing. *Nature* **539**, 452-455 (2016).
114335. Roadmap Epigenomics, C. *et al.* Integrative analysis of 111 reference human
1144 epigenomes. *Nature* **518**, 317-30 (2015).
114536. Richard, M., Drouin, R. & Beaulieu, A.D. ABC50, a novel human ATP-binding
1146 cassette protein found in tumor necrosis factor- α -stimulated synoviocytes.
1147 *Genomics* **53**, 137-45 (1998).
114837. Lieberman-Aiden, E. *et al.* Comprehensive mapping of long-range interactions
1149 reveals folding principles of the human genome. *Science* **326**, 289-93 (2009).
115038. Yang, D. *et al.* 3DIV: A 3D-genome Interaction Viewer and database. *Nucleic Acids*
1151 *Res* **46**, D52-D57 (2018).
115239. Zhang, Y. *et al.* Model-based analysis of ChIP-Seq (MACS). *Genome Biol* **9**, R137
1153 (2008).
115440. Pruim, R.J. *et al.* LocusZoom: regional visualization of genome-wide association
1155 scan results. *Bioinformatics* **26**, 2336-7 (2010).
115641. Dahl, J.A. *et al.* Broad histone H3K4me3 domains in mouse oocytes modulate
1157 maternal-to-zygotic transition. *Nature* **537**, 548-552 (2016).

1158

1159

1160 **Supplementary Information** is linked to the online version of the paper

1161 at www.nature.com/ng.

1162

1163 **Acknowledgement**

1164 We thank members of the Ren laboratory for support and critical

1165 suggestions throughout the course of this work. We are thankful to Dr.

1166 Naoki Nariai (UCSD) for sharing LD information. This work was funded by

1167 in part by the Ludwig Institute for Cancer Research (to B.R.), NIH

1168(U54HG006997 and 1R01ES024984, to B.R.), the Ministry of Science, ICT,
1169and Future Planning through the National Research Foundation in Republic
1170of Korea (2017R1C1B2008838 to I. J.), Korean Ministry of Health and
1171Welfare (HI17C0328 to I. J.), and SUHF Fellowship (to I.J.).

1172

1173**Author Contributions**

1174Ij, AS, YD and BR conceived the study. Ij, AS, and YD performed
1175experiments with assistance from TL, CT, and SC. Ij, AJL, and DY
1176performed data analysis with assistance from JE, MC, ZC, and CLB. DK
1177supervised data analysis by DY. CK, EM, and CLB contributed to provide
1178human brain tissue samples. BL and SK contributed to sequencing and
1179initial data processing. Ij prepared the manuscript with assistance from
1180AS, YD, AJL, JE, and BR. All authors read and commented on the
1181manuscript.

1182

1183**Author Information**

1184Reprints and permissions information is available at
1185www.nature.com/reprints. The authors declare no competing financial
1186interests: details are available in the online version of the paper. Readers
1187are welcome to comment on the online version of the paper.
1188Correspondence and requests for materials should be addressed to B.R.
1189(biren@ucsd.edu) or I.J. (ijung@kaist.ac.kr). All raw and processed data
1190have been deposited in the GEO database under accession number
1191GSE86189.

1192

1193

1194 **Figure Legends**

1195 **Figure 1. Genome-wide mapping of promoter-centered chromatin** 1196 **interactions in diverse human tissues and cell types.**

1197 **a**, A schematic of the pcHi-C procedure. **b**, Barplots of normalized
1198 promoter-centered chromatin interaction frequencies (y-axis) emanating
1199 from the *ADAMTS1* promoter (translucent gray). The identified chromatin
1200 interactions are shown below the axis (purple loops). Highlighted in
1201 translucent yellow are cell/tissue type specific interactions. **c**, Barplots of
1202 the number of chromatin interactions that span a given genomic distance
1203 are shown. Orange line indicates the accumulated fraction of chromatin
1204 interactions from all 27 tissues/cell types. **d**, Boxplots showing the fold
1205 enrichment of the interaction frequencies between the active (colored
1206 dots) or bivalent promoters (gray dots) and each chromatin state. The 17
1207 chromatin states shown were obtained by processing 18-state ChromHMM
1208 model after merging genic enhancer 1 and 2 annotations. KS-tests were
1209 performed between interactions originating from active promoter regions
1210 (colored dots) and those from bivalent promoters (gray dots) (** p value <
1211 0.01 and *** p value < 0.001). The chromatin states that interact more
1212 frequently with active promoters than bivalent promoters were
1213 highlighted in translucent yellow. The chromatin states that interact more
1214 frequently with bivalent promoters than active promoters were
1215 highlighted in translucent blue. Whiskers correspond to the highest and
1216 lowest points within 1.5× the interquartile range.

1217

1218 **Figure 2. Inference of target genes of *cis*-regulatory sequences**
1219 **from pHi-C data.**

1220 **a**, Illustrative LocusZoom plot of eQTLs for *VLDLR* (top) and pHi-C
1221 interactions originating from the *VLDLR* promoter region in aorta tissue
1222 (bottom). Dots along the LocusZoom plot represent the P-values of SNPs'
1223 association with *VLDLR* gene expression levels in the aorta (data obtained
1224 from GTEx). Dots are also color-coded based on their Linkage
1225 Disequilibrium (LD) scores with a tagging SNP. The blue bars indicate the
1226 recombination rate. **b**, Barplots showing fold enrichment between the
1227 number of eQTL-associations matched to P-O pHi-C interactions and that
1228 of distance matched random P-O pHi-C interactions for 12 corresponding
1229 tissue types. P-O interactions in all 12 tissues were significantly enriched
1230 for eQTL associations (empirical p value < 0.01). The dotted line indicates
1231 the expected fold-enrichment (i.e. 1). Error bars indicate standard
1232 deviation obtained by 1,000 random trials. **c**, An illustrative example of
1233 tissue specifically expressed gene, showing positive correlation between
1234 the chromatin state (H3K27ac) at a distal cRE and expression levels (RNA-
1235 seq) of the promoter connected by long-range chromatin interactions. The
1236 significant chromatin interaction between the *POU3F3* promoter and a
1237 distal cRE marked by H3K27ac ~350kb upstream in hippocampus (HC)
1238 tissue is shown at the top. Shown below are H3K27ac signals and
1239 locations of genes. The bar plots at the lower half show the H3K27ac
1240 signals at the distal cRE (left), the transcript levels of the *POU3F3*
1241 (middle), and the normalized pHi-C interaction frequencies between the

1242 *POU3F3* promoter and the distal cRE (right). **d**, Boxplots illustrating the
1243 H3K27ac signals after quantile normalization at the cREs exhibiting
1244 hippocampus specific pHi-C interactions with putative target promoters.
1245 These cREs are marked by higher levels of H3K27ac in hippocampus than
1246 in other cell/tissues types (KS-test p value < 0.005). Whiskers correspond
1247 to the highest and lowest points within 1.5× the interquartile range. **e**,
1248 Boxplots showing transcript levels of the putative target genes predicted
1249 by hippocampus specific pHi-C interactions. Genes are significantly
1250 expressed in hippocampus compared to other cell/tissues types (KS-test p
1251 value < 0.005) except dorsolateral prefrontal cortex (KS-test p value 0.27)
1252 and mesenchymal stem cell (KS-test p value 0.02). Whiskers correspond
1253 to the highest and lowest points within 1.5× the interquartile range. **f-h**,
1254 Heatmaps demonstrate the enrichment of pHi-C interactions for
1255 cell/tissue-specific cRE-promoter pairs (column) in the corresponding
1256 cell/tissue type (row) (f), z-score transformed H3K27ac signals (column) at
1257 the promoter associated cREs (row) (g), and z-score transformed FPKM
1258 values (column) of RNA-seq at the cREs' putative target genes (row) (h).
1259 Color indicates mean values of distance normalized pHi-C interaction
1260 frequencies for H1 (n=5,096), MES (n=3,380), MSC (n=5,188), NPC
1261 (n=1,295), TB (n=5,830), HC (7,100), FC (n=15,733), IMR90 (n=5,313),
1262 LG (n=1,101), LI (n=2,656), PA (n=2,751), SB (n=1,072), TH (n=2,233),
1263 GA (n=1,511), LV (n=1,501), PO (n=865), RV (n=1,049), SX (n=9,228),
1264 AD (n=1,998), AO (n=4,407), and LCL (n=10,283) (f), z-score transformed
1265 H3K27ac signals for H1 (n=5,813), MES (n=3,951), MSC (n=5,790), NPC

1266(n=1,631), TB (n=6,616), HC (7,712), FC (n=15,389), IMR90 (n=6,146),
1267LG (n=1,345), LI (n=3,224), PA (n=3,211), SB (n=1,310), TH (n=2,717),
1268GA (n=1,903), LV (n=1,741), PO (n=1,087), RV (n=1,296), SX (n=10,077),
1269AD (n=2,342), AO (n=5,179), and LCL (n=10,945) (g), and z-score
1270transformed FPKM values for H1 (n=1,589), MES (n=1,024), MSC
1271(n=1,587), NPC (n=450), TB (n=1,920), HC (2,339), FC (n=4,830), IMR90
1272(n=1,743), LG (n=310), LI (n=870), PA (n=845), SB (n=293), TH (n=747),
1273GA (n=460), LV (n=368), PO (n=281), RV (n=295), SX (n=3,054), AD
1274(n=550), AO (n=1,381), and LCL (n=3167) (h). KS-test was performed
1275between pcHi-C interaction frequencies, z-score transformed H3K27ac
1276signals, and z-score transformed FPKM values in the matched cell/tissue
1277types (values in diagonal in each heatmap) and those in other cell/tissue
1278types (values in off diagonal in each heatmap), demonstrating significant
1279association of cRE-promoter pairs with cell/tissue-specific cRE H3K27ac
1280signals and gene expression (KS-test p value < 2.2e-16).

1281

1282**Figure 3. Enhancer-like promoters involved in regulation of distal**
1283**target genes.**

1284**a**, Browser snapshots of the *TMED4* locus showing H3K27ac signals and
1285promoter-centered chromatin interactions. Shown at the RefSeq genes
1286(top), H3K27ac histone modification signals as measured by ChIP-seq
1287(middle) and promoter-centered chromatin interactions detected from
1288pcHi-C experiments (bottom). Highlighted in translucent blue are
1289promoter-promoter pairs showing highly correlated H3K27ac signal and

1290 significant pcHi-C interactions. Highlighted in gray is an adjacent promoter
1291 of the *TMED4*. Shown below are Pearson correlation coefficient (PCC)
1292 values based on H3K27ac signals and links based on pcHi-C interactions,
1293 with MSC as the acronym for mesenchymal stem cell. **b**, Density plots
1294 showing distributions of PCC values of H3K27ac (blue, median of
1295 PCC=0.45, n=48,893), H3K4me1 (orange, median of PCC=0.67,
1296 n=48,893), and H3K4me3 (green, median of PCC=0.64, n=48,893) signals
1297 for P-P pcHi-C interactions. As a control, a density plot of PCC distributions
1298 of H3K27ac signals for randomly selected promoter-promoter pairs is
1299 shown (gray, median of PCC=0.02, n=48,142). X-axis indicates PCC of
1300 histone modification signals between promoter-promoter pairs across 27
1301 cell/tissue types. **c**, A pie chart showing the fraction of unique P-P
1302 interactions matched by eQTL associations, of which 5.7% are P-P
1303 interactions (n=1,976) in 12 matched tissue types (n=34,880). **d**, An
1304 illustrative LocusZoom plot of eQTLs for *DACT3* gene expression in
1305 dorsolateral prefrontal cortex. Both the *DACT3* gene promoter region and
1306 the *AP2S1* gene promoter that contains significant eQTLs are highlighted
1307 in translucent orange, dots along the LocusZoom plot represent SNPs, and
1308 their significance of association with the *DACT3* gene expression is plotted
1309 along the left y-axis. Dots are also color-coded based on their LD score
1310 with a tag SNP (rs78730097). The blue line indicates the estimated
1311 recombination rate, as plotted along the right y-axis. Gene expression
1312 levels detected by RNA-seq and RefSeq genes are plotted below the
1313 LocusZoom plot. **e**, Illustrative genome browser snapshot of RNA-seq

1314 results between control and mutant clones with deletion of the core
1315 promoter regions of the *ARIH2OS*. In both control and mutant cells, the
1316 *ARIH2OS* gene was not expressed, but the expression of the *NCKIPSD*
1317 gene, which displays chromatin interactions with the *ARIH2OS* gene
1318 promoter, was significantly down-regulated in the mutant clones (FDR
1319 adjusted p value from cuffdiff = 0.02). **f**, Genome browser snapshot
1320 showing the promoter containing an eQTL targeted by sgRNAs and its
1321 distal target gene, *ABCF3*, together with H3K27ac and chromatin
1322 accessibility (DNaseI). The relative mRNA expression levels of the *ABCF3*
1323 quantified by RT-qPCR are shown below, which were significantly down-
1324 regulated in both mutants (***) one-tailed KS-test p value < 0.001). Error
1325 bars indicate standard deviation of three mutant clones with technical
1326 triplicates.

1327

1328 Figure 4. Analysis of human diseases and physiological traits
1329 based on the putative target genes of GWAS-SNPs.

1330 **a**, Genome browser snapshot showing multiple cREs harboring GWAS-
1331 SNPs and their common target gene, *NT5DC2*, together with signals of
1332 H3K27ac (ChIP-seq) and chromatin accessibility (DNaseI). The DNA
1333 fragments containing all these cREs interact with the *NT5DC2* gene
1334 promoter region as evidenced by pHi-C analysis (arcs). The relative
1335 mRNA expression levels of the *NT5DC2* upon induced mutation of GWAS-
1336 SNPs with sgRNA were quantified by RT-qPCR as shown below. Error bars
1337 indicate standard deviation of two mutant clones with technical triplicates

1338(KS-test, ** p value < 0.01, *** p value < 0.001). **b**, Hierarchical clustering
1339of human diseases and traits based on similarities of the putative target
1340genes of trait-associated SNPs and SNPs in LD. The color intensity of each
1341dot indicates Pearson correlation coefficient (PCC) of the putative target
1342genes between two diseases or traits. Color bars on the left and top
1343demarcate the clusters. **c, d**, Shown are similarities, as measured by
1344Pearson correlation coefficient (PCC), between traits in the same order as
1345Fig. 4b, based on either the nearest genes of the GWAS SNPs (c) or the
1346GWAS SNPs alone (d). The color intensity of each dot indicates PCC of
1347target gene similarities between two traits. **e**, Hierarchical clustering of
1348GO biological processes (each column, n=126) for the trait clusters
1349defined in Fig. 4b (each row, n=40). Each entry indicates $-\log_{10}(\text{p-value})$
1350value of GO biological processes in the corresponding cluster. Several
1351representative biological processes are highlighted.

1352

1353 **Extended Data Figure Legends**

1354 **Extended Data Figure 1. Capture Hi-C design, probe synthesis,** 1355 **and target enrichment workflow.**

1356 **a**, Schematic of probe design for Promoter Capture Hi-C experiments. For
1357 each promoter (black rectangle), two flanking *HindIII* cut sites were
1358 identified. A 15bp buffer was then added to each side of the *HindIII* cut
1359 site, followed by allocation of three 120-mer capture probes to the same
1360 sites, with a 30bp shift between the adjacent probes. In total, 12 capture
1361 probes were designed for each promoter and all probes were targeted to
1362 the Watson Strand. **b**, Schematic workflow of custom RNA probe
1363 synthesis. Single stranded DNA (ssDNA) probe synthesis by CustomArray,
1364 Inc., is shown from top to bottom; PCR amplification with SP6 recognition
1365 sequence completion and purification, BsrDI digestion and purification, *in*
1366 *vitro* transcription in the presence of biotinylated UTP and purification, and
1367 pooling of probe batches using equal mass ratios. **c**, Schematic workflow
1368 of target enrichment of Hi-C libraries (Promoter Capture Hi-C). From top to
1369 bottom, preparation of library mix, hybridization buffer, and probe mix,
1370 followed by combining the mixes and overnight incubation to bind probes
1371 to Hi-C template. Then, preparation of streptavidin beads and wash
1372 buffers, followed by binding of RNA:DNA duplexes to streptavidin beads
1373 and rigorous washing to remove off-target binding. And lastly, PCR
1374 amplification of the resulting Promoter Capture Hi-C library.

1375

1376 **Extended Data Figure 2. Overview of samples and capture probe**
1377 **quality control.**

1378 **a**, Schematic overview of cell and tissue types analyzed by Promoter
1379 Capture-Hi-C and note of other datasets available for these samples.
1380 Embryonic or embryonic-derived cell types are on the left and tissues are
1381 tabled on the right according to their developmental origin. **b**, Bar plots
1382 showing the fraction of number of TSS in a DNA fragment. **c**, Scatter plot
1383 showing the reproducibility of probe density from RNA-seq data between
1384 two probe synthesis experiments. Each dot on the scatter plot represents
1385 a single promoter and the value is the aggregated probe density from all
1386 probes assigned to that given promoter. **d**, Venn diagram showing the
1387 number of targeted regions that contain detectable probe density based
1388 on RNA-seq of the capture probes from each replicate of probe synthesis.
1389 **e**, Snapshot of Promoter Capture-Hi-C probe density from RNA-seq
1390 analysis of the capture probes. Two replicates of probe synthesis and
1391 subsequent RNA-seq are shown, followed by GENCODE gene annotations.
1392 **f**, Zoomed-in snapshot of Promoter Capture Hi-C probe density from RNA-
1393 seq analysis of the capture probes. Below the replicate RNA-seq datasets
1394 are the *HindIII* cut sites and GENCODE gene annotations, illustrating that
1395 the vast majority of probe density is only found around *HindIII* restriction
1396 sites flanking promoters. **g, h**, Histograms of the probe densities
1397 measured by RNA-seq (x-axis) in each promoter from replicate 1 (g) and
1398 replicate 2 (h) of probe synthesis.

1399

1400 **Extended Data Figure 3. General characterization of promoter-**
1401 **centered long-range interactions.**

1402 **a**, Identified pHi-C chromatin interactions across multiple cell/tissue
1403 types are plotted in Genome Browser, with the darkness of blue
1404 corresponding to the strength of interactions. RefSeq genes are presented
1405 below the snapshot. **b**, Fraction of pHi-C interactions uniquely detected
1406 in one cell/tissue type (green) or also detected in other cell/tissue types
1407 (orange). The average fraction of cell/tissue-specific interactions is not
1408 over-estimated due to the number of tested samples (at 22 samples the
1409 fraction of cell/tissue-specific interactions reach plateau) and tissue-
1410 heterogeneity (similar trend was observed when we only considered pHi-
1411 C interactions obtained from cell lines). **c**, Snapshot of a locus showing
1412 promoter-centered long-range interactions from pHi-C data in H1-hESC
1413 (bottom, purple loops) in the context of TAD annotations (blue rectangles)
1414 identified from Hi-C data (top, red) in H1-hESC. RefSeq genes are shown
1415 at the bottom. **d**, Fraction of P-O pHi-C chromatin interactions in the
1416 context of TAD annotations with the respective cell/tissue types.

1417

1418 **Extended Data Figure 4. Validation of Promoter Capture Hi-C**
1419 **Interactions.**

1420 **a**, Browser snapshot of the *CCL* gene cluster, highlighting the similarity of
1421 promoter-centered interactions from Promoter Capture Hi-C and high
1422 resolution Hi-C data in IMR90. The top two tracks show histone
1423 modification signals for H3K4me3 and H3K27ac, followed by RefSeq

1424genes. Below are pcHi-C chromatin in IMR90 (blue loops) and promoter-
1425centered chromatin interactions from high-resolution Hi-C data in IMR90
1426(reddish brown loops). **b-e**, ROC plots illustrating the prediction
1427performance of Promoter Capture Hi-C result for *in situ* Hi-C loops
1428anchored at promoters in lymphoblastoid (b), IMR90 (c), hippocampus (d),
1429and dorsolateral prefrontal cortex (e). Promoter centered interactions for
1430*in situ* Hi-C loops were considered as true interactions, and ROC plots
1431were drawn for the corresponding pcHi-C result. ROC scores are shown in
1432the ROC plot. **f**, ROC plots showing the reproducibility of pcHi-C chromatin
1433interactions between biological replicates. pcHi-C interactions from one
1434replicate were used as true interactions, and ROC plots were drawn for the
1435other replicate. **g-k**. Venn diagrams presenting the number of commonly
1436identified pcHi-C interactions between biological replicates for
1437lymphoblastoid (g), dorsolateral prefrontal cortex (h), mesenchymal stem
1438cell (i), lymphoblastoid processed by CHICAGO (j), and GM12878 with
1439previously published pcHi-C data¹⁸ (k). Hypergeometric p-values are
1440shown together. **l-m**, Illustration of interaction intensity in the replicates
1441of lymphoblastoid (l) and mesenchymal stem cells (m), depending on the
1442replicate consistency. Whiskers correspond to the highest and lowest
1443points within 1.5× the interquartile range.

1444

1445**Extended Data Figure 5. Integrative analysis of long-range**
1446**chromatin interactions with epigenome.**

1447**a, b**, Shown are histograms of number of interacting cREs per promoter
1448(a) and number of interacting promoters per cRE (b). Y-axis indicates
1449frequency of the corresponding value in x-axis. **c**, Depiction of identified
1450long-range promoter-centered interactions across a 0.84Mb locus in
1451lymphoblastoid (top). Shown below are histone modification signals
1452obtained from ChIP-seq analyses³⁵, as well as accessible chromatin
1453regions measured from DNaseI hypersensitivity assay. **d**, Depiction of
1454extensively interacting DNA fragments (EIF) from P-P and P-O interactions,
1455and transcription factor (TF) binding clusters identified in GM12878 cells
1456for the same region shown in Extended Data Fig. 5c. Below are 67 TF
1457binding profiles obtained from TF ChIP-seq results performed in GM12878
1458cells. Highlighted in translucent blue are overlapping EIF and TF binding
1459clusters. EIF was defined in each cell/tissue type by selecting frequently
1460interacting DNA fragments with multiple promoters in terms of 0.01
1461Poisson p value cutoff. **e, f**, Bar plots showing the number of P-O EIF
1462overlapping with TF clusters compared to random expectation in
1463lymphoblastoid (e) and H1-hESC (f). Error bars indicate standard deviation
1464of expectation values calculated by using typical TF peaked regions (blue)
1465and generating random genomic regions (green). Empirical p-value shows
1466statistical significance (***) p value < 0.001). **g, h**, Bar plots showing the
1467number of P-P EIF overlapping with TF clusters compared to random
1468expectation in lymphoblastoid (g) and H1-hESC (h). Error bars indicate
1469standard deviation of expectation values calculated by using typical TF
1470peaked regions (blue) and generating random genomic regions (green).

1471 Empirical p-value shows statistical significance (***) p value < 0.001). **i**, An
1472 array of bar plots showing the number of P-O EIF overlapping with super-
1473 enhancers (first bar plot, orange), compared to typical enhancers (middle
1474 bar plot, blue) and random genomic regions (last bar plot, purple). Error
1475 bars indicate standard deviation of expectation values obtained by 10,000
1476 permutations. Empirical p-value showed statistical significance for all
1477 tested cell/tissue types compared to random genomic regions (p value <
1478 0.0001).

1479

1480 **Extended Data Figure 6. Enrichment of long-range chromatin**
1481 **interactions at various chromatin states generated by a 50-state**
1482 **ChromHMM model.**

1483 **a**, Boxplots showing the fold change of interaction frequencies between
1484 active/bivalent promoters and each chromatin state over expected values.
1485 The 50 chromatin states (E01-E50) were obtained from the 50-state
1486 ChromHMM model. KS-tests were performed between active promoters
1487 and bivalent promoters (two adjacent boxplots) (** p value < 0.01 and ***
1488 p value < 0.001). The chromatin states that interact more frequently with
1489 active promoters than bivalent promoters were highlighted in pink
1490 asterisk. The chromatin states that interact more frequently with bivalent
1491 promoters than active promoters were highlighted in blue asterisk.
1492 Whiskers correspond to the highest and lowest points within 1.5× the
1493 interquartile range. **b**, A heatmap showing an emission parameter matrix
1494 of each chromatin state in which each row corresponds to a different

1495chromatin state and each column corresponds to an emission probability
1496of a chromatin mark shown at the top. The pre-calculated emission
1497parameter heatmap was downloaded from the 50-state ChromHMM model
1498established by Roadmap Epigenomics Project.

1499

1500**Extended Data Figure 7. Validation of P-O interactions with eQTL**
1501**associations.**

1502**a-c**, Illustrative LocusZoom plots of eQTLs for the *HS3ST1* (a), the
1503*METTL25* (b), and the *DAAM1* (c) gene expression in left ventricle,
1504dorsolateral prefrontal cortex, and aorta, respectively. RefSeq genes are
1505plotted below the LocusZoom plot. Identified pHi-C interactions are
1506shown as loops (purple) in the bottom. **d**, Array of bar plots showing the
1507number of matched eQTL associations between P-O pHi-C chromatin
1508interactions after exclusion of DNA fragment shared promoters and
1509random expectation across 14 matched tissue types from GTEx database.
1510All P-O pHi-C interactions are significantly enriched by eQTL associations
1511compared to random P-O pHi-C interactions with or without distance
1512match (* empirical p-value <0.05, ** empirical p-value <0.01, ***
1513empirical p-value <0.001). Error bars indicate standard deviation of
1514random expectation values. **e**, Density plots showing the number of
1515unique eQTLs per P-O pHi-C interaction fragment and randomized
1516interactions. No significant difference between pHi-C interactions and
1517randomized interactions (KS-test p value > 0.05) except pancreas (p value
1518= 0.02), gastric (p value = 0.009), and lung (p value = 0.03). **f**, Shown are

1519 boxplots of the distribution of PCC between H3K27ac signals in cRE-
1520 promoter pairs connected by pHi-C interactions after exclusion of
1521 multiple fragment spanning cREs (Orange, n=154,055), compared to the
1522 distribution of random expectation with matched distance (dark gray,
1523 n=154,055) and without matched distance (gray, n=154,055). We only
1524 considered P-O pairs where other DNA fragments are marked by H3K27ac
1525 peaks in at least one cell/tissue type analyzed. We also excluded two
1526 fragments spanning cREs. KS-test was performed between P-O pairs and
1527 random control, demonstrating that P-O pairs showed significant positive
1528 correlation (***) Welch's t-test p value < 2.2e-16). Whiskers correspond to
1529 the highest and lowest points within 1.5× the interquartile range. **g**,
1530 Similar to Extended Data Fig. 7e, but the distribution of PCC between
1531 H3K27ac signals at a cRE and target gene expressions of the cRE
1532 connected by pHi-C interactions. KS-test was performed between P-O
1533 pairs (orange, n=154,055), distance matched random control (dark gray,
1534 n=154,055), and random control (gray, n=154,055), revealing that P-O
1535 pairs showed significant positive correlation (***) Welch's t-test p value <
1536 2.2e-16). Whiskers correspond to the highest and lowest points within
1537 1.5× the interquartile range.

1538

1539 Extended Data Figure 8. Functional analysis of promoter-
1540 promoter interactions.

1541 **a**, Pie chart showing the fraction of promoter-promoter interactions (P-P)
1542 among all pHi-C interactions. The fraction of P-P pHi-C interactions

1543 modestly decrease to 6.5% after excluding fragment that harbor multiple
1544 promoters. **b**, An array of bar plots showing the number of eQTL
1545 associations matched to P-P pcHi-C interactions (left, purple), compared to
1546 random expectation with matched distance (middle, blue) and without
1547 matched distance (right, light blue). Each bar plot represents analysis of a
1548 different tissue. Error bars indicate standard deviation of random
1549 expectation values. Empirical p values are shown at the top (* < 0.05, **
1550 < 0.01, *** < 0.001). **c**, **d**, Illustrative LocusZoom plots of *FHOD1* eQTLs
1551 (c) and *POFUT2* eQTLs (d) in left ventricle and aorta, respectively.
1552 Promoters that contain significant eQTLs and target promoters are
1553 highlighted in translucent orange. Dots along the LocusZoom plot
1554 represent SNPs, and their significance of association with *FHOD1* and
1555 *POFUT2* gene expression is plotted along the left y-axis, respectively. The
1556 blue line traveling across the scatterplot indicates the estimated
1557 recombination rate, as plotted along the right y-axis. RefSeq genes and
1558 RNA-seq are plotted below the LocusZoom plot. pcHi-C interactions are
1559 shown as purple in the bottom. **e**, Bar plot showing the eQTL associations
1560 between the SNP rs78730097 and surrounding genes, showing the most
1561 significant association with the distal gene *DACT3*. Y-axis indicates -
1562 \log_{10} (eQTL association p value). **f**, **g**, Bar plots showing FPKM values of
1563 distal target gene expressions upon deletion of core promoter regions of
1564 the *ARIH2OS* (f) and the *ZNF891* (g). Two biological replicates of one
1565 mutant clone for the *NCKIPSD* and two biological replicates of three
1566 mutant clones for the *ZNF84* were measured using RNA-seq, respectively.

1567 FDR-adjusted p value obtained from cuffdiff is shown together. N.S
1568 indicates statistically non-significant. **h**, Bar plots showing FPKM values of
1569 two nearby genes of the *ARIH2OS* and one nearest gene of the *NCKIPSD*
1570 (y-axis) upon deletion of core promoter regions of the *ARIH2OS*. The
1571 *ARIH2*, a DNA fragment sharing gene with the *ARIH2OS*, is excluded. FDR-
1572 adjusted p value obtained from cuffdiff is shown together. Corresponding
1573 gene name is shown on the top of bar plots. **i**, The relative mRNA
1574 expression levels of distal target genes (orange) and nearby genes (gray)
1575 of promoter-proximal eQTLs quantified by RT-qPCR are shown. Error bars
1576 indicate standard deviation from total six mutant clones for two separate
1577 sgRNAs with technical triplicates. One-sided KS-test p values are shown
1578 together on the top of each bar plot (***) p value < 0.001).

1579

1580 **Extended Data Figure 9. Identification of target genes of disease-**
1581 **associated genetic variants.**

1582 **a**, Illustration of the strategy to identify target genes of each GWAS trait.
1583 An example result is shown for Alzheimer disease. Both known and novel
1584 target genes were identified according to literature search. **b**, Venn
1585 diagram showing number of target genes by pHi-C interactions and by
1586 nearby gene information for the GWAS-SNPs associated with Parkinson
1587 disease. **c**, Number of matched disease-associated genes in each group of
1588 target genes identified in Parkinson disease. **d**, Fraction of distal genes
1589 (blue) and nearby genes (gray) among the identified target genes of
1590 GWAS-SNPs based on pHi-C interactions (left). Expected fraction is shown

1591by calculating the fraction of nearby genes when we consider a nearest
1592gene over 15kb as a GWAS-SNP target gene (right). **e**, Barplot showing
1593the relative mRNA expression levels of *GNL3* upon induced mutation of
1594GWAS-SNPs with sgRNA as quantified by RT-qPCR as a control. Error bars
1595indicate standard deviation of two mutant clones with technical triplicates.
1596**f**, Barplot showing RT-qPCR results of relative target gene expression (y-
1597axis) between mutant and control. Error bars indicate standard deviation
1598of two mutant clones with technical triplicates. The mutants showing
1599significant down regulation of target genes are shown in orange (KS-test,
1600** p value < 0.01, *** p value < 0.001). sgRNA target GWAS-SNP genomic
1601coordinate, rsID, associated disease, distal target gene information, high
1602LD SNP on coding region, and related publication PMID information are
1603shown together.

1604

1605**Extended Data Figure 10. Analysis of disease-disease**
1606**associations.**

1607**a**, Illustration of the strategy to calculate the similarity between GWAS
1608mapped traits using target gene similarity information. **b**, **c**, Shown are
1609similarities, as measured by Pearson correlation coefficient (PCC),
1610between traits in the same order as Fig. 4b based on similarities of the
1611putative GWAS-SNP target genes without shared promoters (b) and
1612without genes located in HLA and HIST locus (c). The color intensity of
1613each dot indicates Pearson correlation coefficient (PCC) of the putative
1614target genes between two diseases or traits. **d**, Shown are similarities, as

1615 measured by Pearson correlation coefficient (PCC), between traits based
1616 on the 5 nearest genes of the GWAS SNPs. The color intensity of each dot
1617 indicates PCC of target gene similarities between GWAS mapped traits. **e,**
1618 Bar plots showing the fraction of number of TSS in a DNA fragment
1619 between all TSS and TSS corresponding genes in cluster 38 of Fig. 4b.

1620

1621

1622 **Supplementary Tables**

1623

1624 **Supplementary Table 1. List of cell/tissue types analyzed in this**
1625 **study**

1626 **Supplementary Table 2. Number of processed pcHi-C reads**

1627 **Supplementary Table 3. List of P-O interactions**

1628 **Supplementary Table 4. List of P-P interactions**

1629 **Supplementary Table 5. Number of significant pcHi-C promoter-**
1630 **centered interactions**

1631 **Supplementary Table 6. The list of mean and median distance of**
1632 **pcHi-C and eQTL associations in each cell/tissue type**

1633 **Supplementary Table 7. The numbers and fractions of overlapped**
1634 **interactions between replicates**

1635 **Supplementary Table 8. Total number of extensively interacting**
1636 **DNA fragments (Poisson P value < 0.01)**

1637 **Supplementary Table 9. List of TF ChIP-seq data used to define**
1638 **GM12878 TF clusters**

1639 **Supplementary Table 10. List of TF ChIP-seq data used to define**
1640 **H1-hESC TF clusters**

1641 **Supplementary Table 11. Summary of matched eQTL-associations**
1642 **with P-O pcHi-C interactions**

1643 **Supplementary Table 12. List of P-O pcHi-C interactions and**
1644 **matched eQTL relationships**

1645 **Supplementary Table 13. Summary of matched eQTL-associations**
1646 **with P-P pcHi-C interactions**

1647 **Supplementary Table 14. List of P-P pcHi-C interactions and**
1648 **matched eQTL relationships**

1649 **Supplementary Table 15. Summary of average number of target**
1650 **genes of GWAS-SNPs**

1651 **Supplementary Table 16. List of putative target genes of GWAS-**
1652 **SNPs**

1653 **Supplementary Table 17. List of GWAS mapped traits and**
1654 **enriched GO biological processes in Fig. 4b**

1655 **Supplementary Table 18. Enriched pathway analysis of Cluster 38**
1656 **in Fig. 4b**