

1 **Promoter-anchored chromatin interactions predicted from genetic analysis of**
2 **epigenomic data**

3
4 Yang Wu^{1,7}, Ting Qi^{1,7}, Huanwei Wang¹, Futao Zhang¹, Zhili Zheng^{1,2}, Jennifer E. Phillips-Cremins³,
5 Ian J. Deary^{4,5}, Allan F. McRae¹, Naomi R. Wray^{1,6}, Jian Zeng¹, Jian Yang^{1,2,6,*}

6
7 ¹ Institute for Molecular Bioscience, The University of Queensland, Brisbane, Queensland 4072,
8 Australia

9 ² Institute for Advanced Research, Wenzhou Medical University, Wenzhou, Zhejiang 325027,
10 China

11 ³ Department of Bioengineering, University of Pennsylvania, Philadelphia, PA 19104, USA

12 ⁴ Centre for Cognitive Ageing and Cognitive Epidemiology, University of Edinburgh, Edinburgh
13 EH8 9JZ, UK

14 ⁵ Department of Psychology, University of Edinburgh, Edinburgh EH8 9JZ, UK

15 ⁶ Queensland Brain Institute, The University of Queensland, Brisbane, Queensland 4072,
16 Australia

17 ⁷ These authors contributed equally to this work.

18
19 * Correspondence: Jian Yang (jian.yang@uq.edu.au)

20

21

22 **Abstract**

23 Promoter-anchored chromatin interactions (PAIs) play a pivotal role in transcriptional regulation.
24 Current high-throughput technologies for detecting PAIs, such as promoter capture Hi-C, are often
25 limited in sample size due to the complexity of the experiments. Here, we present an analytical
26 approach that uses summary-level data from DNA methylation (DNAm) quantitative trait locus
27 (mQTL) studies to predict PAIs. Using mQTL data from human peripheral blood ($n=1,980$), we
28 predicted 34,797 PAIs which showed strong overlap with the chromatin contacts identified by
29 experimental assays. The promoter-interacting DNAm sites were enriched in enhancers or near
30 expression QTLs. Genes whose promoters were involved in PAIs were more actively expressed,
31 and gene pairs with promoter-promoter interactions were enriched for co-expression.
32 Integration of the predicted PAIs with GWAS data highlighted interactions among 601 DNAm sites
33 associated with 15 complex traits. This study demonstrates the use of mQTL data to predict PAIs
34 and provide insights into the role of PAIs in complex trait variation.

35

36 **Introduction**

37 Genome-wide association studies (GWASs) in the past decade have identified tens of thousands
38 of genetic variants associated with human complex traits (including common diseases) at a
39 stringent genome-wide significance level^{1,2}. However, most of the trait-associated variants are
40 located in non-coding regions^{3,4}, and the causal variants as well as their functional roles in trait
41 etiology are largely unknown. One hypothesis is that the genetic variants affect the trait through
42 genetic regulation of gene expression⁴. Promoter-anchored chromatin interaction (PAI)^{5,6} is a key
43 regulatory mechanism whereby non-coding genetic variants alter the activity of cis-regulatory
44 elements and subsequently regulate the expression levels of the target genes. Therefore, a
45 genome-wide map of PAIs is essential to understand transcriptional regulation and the genetic
46 regulatory mechanisms underpinning complex trait variation.

47
48 High-throughput experiments, such as Hi-C⁷ and ChIA-PET (chromatin interaction analysis by
49 paired-end tag sequencing)⁸, have been developed to detect chromatin interactions by a massively
50 parallelized assay of ligated DNA fragments. Hi-C is a technique based on chromosome
51 conformation capture (3C)⁹ to quantify genome-wide interactions between genomic loci that are
52 close in three-dimensional (3D) space. ChIA-PET is a method that combines ChIP-based methods¹⁰
53 and 3C. The ChIA-PET method uses particular immune-precipitated proteins (e.g., transcription
54 factors) that can bind to specific genomic regions to identify protein-specific chromatin
55 interactions⁸. However, these high-throughput assays are currently not scalable to population-
56 based cohorts with large sample sizes because of the complexity of generating a DNA library for
57 each individual (tissue or cell line) and the extremely high sequencing depth needed to achieve
58 high detection resolution¹¹. On the other hand, recent technological advances have facilitated the
59 use of epigenomic marks to infer the chromatin state of a specific genomic locus and further to
60 predict the transcriptional activity of a particular gene^{12,13}. There have been increasing interests
61 in the use of epigenomic data (e.g., DNA methylation (DNAm) and/or histone modification) to
62 infer chromatin interactions¹⁴⁻¹⁷. These analyses, however, rely on individual-level chromatin
63 accessibility data often only available in small samples^{14,16}, and it is not straightforward to use the
64 predicted chromatin interactions to interpret the variant-trait associations identified by GWAS.

65
66 In this study, we proposed an analytical approach to predict chromatin interaction by detecting
67 the association between DNAm levels of two CpG sites due to the same set of genetic variants (i.e.,
68 pleiotropic association between DNAm sites). This can be achieved because if the methylation
69 levels (unmethylated, partly methylated or fully methylated) of a pair of relatively distal CpG sites
70 covary across individuals and such covariation is not (or at least not completely) caused by
71 environmental or experimental factors (evidenced by the sharing of a common set of causal

72 genetic variants in cis) (**Fig. 1a**), it is very likely that the two genomic regions interact (having
73 contacts or functional links because of their close physical proximity in 3D space). Our analytical
74 approach was based on two recently developed methods, i.e., the summary-data-based Mendelian
75 randomization (SMR) test and the test for heterogeneity in deendent instruments (HEIDI)¹⁸,
76 which can be used in combination to detect pleiotropic associations between a molecular
77 phenotype (e.g. gene expression or DNA methylation) and a complex trait¹⁸ or between two
78 molecular phenotypes¹⁹. The SMR and HEIDI approaches only require summary-level data from
79 DNA methylation quantitative trait locus (mQTL) studies, providing the flexibility of using mQTL
80 data from studies with large sample sizes to ensure efficient power. Since the proposed method is
81 based on cohort-based genetic data, it also allows us to integrate the predicted chromatin
82 interactions with GWAS results to understand the genetic regulatory mechanisms for complex
83 traits. In this study, we analyzed mQTL summary data from a meta-analysis of studies on 1,980
84 individuals with DNAm levels measured by Illumina 450K methylation arrays and SNP data from
85 SNP-array-based genotyping followed by imputation to the 1000 Genome Project (1KGP)
86 reference panels^{19,20}. For the ease of computation and to control the number of tests, we limited
87 the analysis to predict the interactions between promoters and genomic regions (of
88 approximately 4 Mb) centered around the focal promoters.

89

90 **Results**

91 **Predicting promoter-anchored chromatin interactions using mQTL data**

92 As described above, our underlying hypothesis was that if the variation between people in DNAm
93 levels of two relatively distal CpG sites are associated due to the same set of causal genetic variants
94 (**Fig. 1a**), then it is very likely that these two chromatin regions have contacts or functional links
95 because of their close physical proximity in 3D space. Hence, we set out to predict the promoter-
96 anchored chromatin interactions (PAIs) from mQTL data. Our approach was to apply the SMR and
97 HEIDI approaches (both implemented in the SMR software tool)¹⁸ to test for pleiotropic
98 associations between a DNAm site in the promoter region of a gene and all the other DNAm sites
99 within 2 Mb of the gene (excluding the DNAm sites in the same promoter region as the target
100 DNAm site) using mQTL summary data from peripheral blood samples (**Fig. 1, Fig. S1** and
101 **Methods**). Therefore, our analysis was not limited to detect chromatin looping interactions but a
102 scan for genomic regions that are functionally associated with promoter regions likely because of
103 chromatin contacts or close physical proximity in 3D space. In the SMR analysis, the promoter
104 DNAm site was used as the “exposure” and each of the other DNAm sites in the region was used
105 as the “outcome” (**Fig. 1**). The mQTL summary data were generated from a meta-analysis of the
106 mQTL data sets available in McRae et al. ($n = 1,980$)^{19,20}. The mQTL effects were in standard
107 deviation (SD) units of DNAm levels. For exposure probes, we included in the SMR analysis only

108 the DNAm sites with at least one cis-mQTL (SNPs within 2 Mb of the CpG site associated with
109 variation in DNAm level) at $P_{\text{mQTL}} < 5 \times 10^{-8}$ (note that a basic assumption of Mendelian
110 randomization is that the SNP instrument needs to be strongly associated with the exposure^{21,22}).
111 There were 90,749 DNAm probes with at least one cis-mQTL at $P_{\text{mQTL}} < 5 \times 10^{-8}$, 28,732 of which
112 were located in promoters annotated based on data from blood samples of the Roadmap
113 Epigenomics Mapping Consortium (REMC)¹³. We used the 1KGP-imputed Health and Retirement
114 Study (HRS)²³ data as a reference sample for linkage disequilibrium (LD) estimation to perform
115 the HEIDI test, which rejects SMR association between DNAm sites that are not driven by the same
116 set of causal variants (called linkage model in Zhu et al.¹⁸). In total, we identified 34,797 PAIs
117 between pairwise DNAm sites that passed the SMR test ($P_{\text{SMR}} < 1.69 \times 10^{-9}$ based on Bonferroni
118 correction for multiple tests) and were not rejected by the HEIDI test ($P_{\text{HEIDI}} > 0.01$; see Wu et al.¹⁹
119 for the justification of the use of this HEIDI threshold p-value). The significant PAIs comprised of
120 21,787 unique DNAm sites, among which 10,249 were the “exposure” probes in promoter regions
121 of 4,617 annotated genes. Most of the DNAm sites in promoters showed pleiotropic associations
122 with multiple DNAm sites (mean = 4) (**Fig. S2a**). The distances between 95% of the pairwise
123 interacting DNAm sites were less than 500 Kb (mean = 79 Kb and median = 23 Kb). Approximately
124 0.7% of the predicted PAIs were between DNAm sites greater than 1 Mb apart (**Fig. S2b**).

125

126 **Overlap of the predicted PAIs with Hi-C data**

127 We first examined whether the predicted PAIs are consistent with chromatin contacts identified
128 by experimental assays, such as Hi-C²⁴ and promoter captured Hi-C (PChi-C)⁵. While the majority
129 of experimental assays are measured in primary cell lines, topological associated domains (TADs)
130 annotated from Hi-C are relatively conserved across cell types²⁵. We therefore tested the overlap
131 of our predicted PAIs with the TADs identified from recent Hi-C and PChi-C studies^{5,24,26}. We found
132 that 22,024 (22,024/34,797 = 63.3%) of the predicted PAIs were between DNAm sites located in
133 the TADs identified by Rao et al. using Hi-C in the GM12878 cell lines²⁴, 27,200 (27,200/34,797 =
134 78.2%) in those by Dixon et al. using Hi-C in embryonic stem cells²⁶, and 27,716 (27,716/34,797
135 = 79.7%) in those by Javierre et al. using PChi-C in primary hematopoietic cells⁵. These numbers
136 of overlaps with Hi-C and PChi-C data were significantly higher than those for the same number
137 of DNAm pairs randomly sampled (repeated 1,000 times) from distance-matched DNAm pairs
138 tested in the SMR analysis ($P < 0.001$ for all the three Hi-C or PChi-C data sets; note that the p-
139 value was truncated at 0.001 due to the finite number of resampling) (**Fig. 2a, Fig. 2b, Fig. 2c** and
140 **Methods**). One example was the *MAD1L1* locus (a ~450 Kb region) on chromosome 7 (**Fig. 2d**
141 and **Fig. 2e**) where there were a large number of predicted PAIs highly consistent with TADs
142 identified by Hi-C from the Rao et al. study²⁴. There were also scenarios where the predicted PAIs
143 were not aligned well with the TAD data. For example, 58.5% of the predicted PAIs at the *RPS6KA2*

144 gene locus did not overlap with the TADs identified by Hi-C from the Rao et al. study²⁴ (**Fig. S3a**).
145 These predicted interactions, however, are very likely to be functional as indicated by our
146 subsequent analysis with GWAS and omics data (see below). Additionally, the predicted PAIs were
147 slightly enriched for chromatin loops identified by Hi-C²⁴ (1.49-fold, $P < 0.001$), although the
148 number of overlaps was small ($m = 130$).

149

150 **Enrichment of the predicted PAIs in functional annotations**

151 To investigate the functional role of the DNAm sites that showed significant interactions with the
152 DNAm sites in promoter regions (called promoter-interacting DNAm sites or PIDSs hereafter; i.e.,
153 the “outcome” probes of the significant PAIs), we conducted an enrichment analysis of the PIDSs
154 ($m = 14,361$) in 14 main functional annotation categories derived from the REMC blood samples
155 (**Methods**). Note that the PIDS of a PAI can also be located in the promoter region of another gene
156 (i.e., promoter-promoter interaction; **Fig. 1b**). The fold-enrichment was computed as the
157 proportion of PIDSs in a functional category divided by that for the same number of variance-
158 matched “control” probes randomly sampled from all the “outcome” probes used in the SMR
159 analysis. The standard deviation of the estimate of fold-enrichment was computed by repeatedly
160 sampling the control set 1,000 times. We found a significant enrichment of PIDSs in enhancers
161 (fold-enrichment=2.17 and $P_{\text{enrichment}} < 0.001$), repressed Polycomb regions (fold-
162 enrichment=1.50 and $P_{\text{enrichment}} < 0.001$), primary DNase (fold-enrichment=1.37 and $P_{\text{enrichment}} <$
163 0.001) and bivalent promoters (fold-enrichment=1.21 and $P_{\text{enrichment}} < 0.001$) and a significant
164 underrepresentation in transcription starting sites (fold-enrichment=0.25 and $P_{\text{enrichment}} < 0.001$),
165 quiescent regions (fold-enrichment=0.74 and $P_{\text{enrichment}} < 0.001$), promoters around transcription
166 starting sites (fold-enrichment=0.83 and $P_{\text{enrichment}} < 0.001$), and transcribed regions (fold-
167 enrichment=0.86 and $P_{\text{enrichment}} < 0.001$) in comparison with the “control” probes (**Fig. 3a** and **Fig.**
168 **3b**). Although the PIDSs are underrepresented in promoters, it is of note that a large proportion
169 (~21%) of the predicted PAIs were promoter-promoter interactions (PmPmI), consistent with
170 the result from a previous study^{5,27} that PmPmI were widespread and may play an important role
171 in transcriptional regulation.

172

173 **Relevance of the predicted PAIs with gene expression**

174 We then tested whether pairwise genes with significant PmPmI were enriched for co-expression.
175 We used gene expression data (measured by Transcript Per Kilobase Million mapped reads (TPM))
176 from the blood samples of the Genotype-Tissue Expression (GTEx) project²⁸ and computed the
177 Pearson correlation of expression levels across individuals between pairwise genes (r_p). We
178 randomly sampled the same number of distance-matched gene pairs ($m = 2,236$) from all the pairs,
179 whose promoters were tested for interaction in the PAI analysis, and repeated the sampling 1,000

180 times to generate a distribution of mean correlations of expression levels between a set of “control”
181 gene pairs. The mean correlation for the significant PmPmI gene pairs (\bar{r}_p) was 0.375, significantly
182 ($P < 0.001$) higher than that computed from the control gene pairs (mean of $\bar{r}_p = 0.317$) (**Fig. 3c**),
183 suggesting that pairwise genes with PmPmI are more likely to be co-expressed.

184
185 We also tested whether genes whose promoters were involved in significant PAI (called Pm-PAI
186 genes hereafter, **Fig. 1**) were expressed more actively than the same number of “control” genes
187 whose promoter DNAm sites were included in the SMR analysis. Similar to the analysis above, we
188 used the gene expression data (measured by TPM) from the blood samples of the GTEx project
189 and tested the significance of enrichment of Pm-PAI genes in different expression level groups.
190 We found that in comparison to the “control” genes, Pm-PAI genes were significantly
191 overrepresented ($P < 0.001$) among the group of genes with the highest expression levels and
192 significantly underrepresented ($P < 0.001$) among genes that were not actively expressed (median
193 TPM < 0.1) (**Fig. 3d** and **Methods**), implicating the regulatory role of the PIDSs in transcription
194 and their asymmetric effects on gene expression.

195 196 **Enrichment of eQTLs in the PIDS regions**

197 We have shown that the PIDSs are located in regions enriched with regulatory elements (e.g.,
198 enhancers) (**Fig. 3b**) and that the Pm-PAI genes tend to have higher expression levels (**Fig. 3d**).
199 We next investigated if genomic regions near PIDS are enriched for genetic variants associated
200 with expression levels of Pm-PAI genes using data from an expression quantitative trait locus
201 (eQTL) study in blood²⁹. There were 11,204 independent cis-eQTLs at $P_{\text{eQTL}} < 5 \times 10^{-8}$ for 9,967
202 genes, among which 2,053 were Pm-PAI genes (**Methods**). We mapped cis-eQTLs to a 10 Kb
203 region centered around each PIDS (5 Kb on either side) and counted the number of cis-eQTLs
204 associated with expression levels of the corresponding Pm-PAI gene for each PIDS. There were
205 591 independent eQTLs located in the PIDS regions of the Pm-PAI genes, significantly higher than
206 ($P < 0.001$) that from a “control” sample (mean = 454), where the number of independent eQTL
207 was computed from the same number of 10 Kb regions around distance-matched pairs of DNAm
208 probes randomly sampled from the SMR “exposure” and “outcome” probes (**Fig. 4a**). These results
209 again imply the regulatory role of the PIDSs in transcription through eQTLs and provide evidence
210 supporting the functional role of the predicted PAIs.

211
212 There were examples where a cis-eQTL was located in a PIDS region predicted to interact with
213 the promoters of multiple genes. For instance, our result showed that a cis-eQTL was located in
214 an enhancer region that was predicted to interact with the promoters of three genes (i.e., *ABCB9*,
215 *ARL6IP4*, and *MPHOSPH9*) (**Fig. S4**), and the predicted interactions were consistent with the TADs

216 identified by Hi-C from Rao et al.²⁴ (**Fig. S3b**). Furthermore, the predicted interactions between
217 promoter regions of *ARL6IP4* and *MPHOSPH9* are consistent with the chromatin contact loops
218 identified by Hi-C in the the GM12878 cells²⁴ (**Fig. S4**). The eQTL association signals were highly
219 consistent for the three genes, and the pattern was also consistent with the SNP association
220 signals for schizophrenia (SCZ) and years of education (EY) as shown in our previous work¹⁹,
221 suggesting a plausible mechanism whereby the SNP effects on SCZ and EY are mediated by the
222 expression levels of at least one of the three co-regulated genes through the interactions of the
223 enhancer and three promoters (**Fig. S4**).

224

225 We have shown previously that the functional association between a DNAm site and a gene nearby
226 can be inferred by the pleiotropic association analysis using SMR and HEIDI considering the
227 DNAm level of a CpG site as the exposure and gene expression level as the outcome¹⁹. We further
228 tested if the PIDSs are enriched among the DNAm sites showing pleiotropic associations with the
229 expression levels of the neighboring Pm-PAI genes. We found that approximately 10% of the
230 PIDSs were the gene-associated DNAm sites identified in our previous study¹⁹, significantly higher
231 ($P < 0.001$) than that computed from the distance-matched control probe pairs (1.1%) described
232 above (**Fig. 4b**).

233

234 **Replication of the predicted PAIs across tissues**

235 To investigate the robustness of the predicted PAIs across tissues, we performed the PAI analysis
236 using brain mQTL data from the Religious Orders Study and Memory and Aging Project
237 (ROSMAP)³⁰ ($n = 468$). Of the 11,082 PAIs with $P_{\text{SMR}} < 1.69 \times 10^{-9}$ and $P_{\text{HEIDI}} > 0.01$ in blood and
238 available in brain, 2,940 (26.5%) showed significant PAIs in brain after Bonferroni correction for
239 multiple testing ($P_{\text{SMR}} < 4.51 \times 10^{-6}$ and $P_{\text{HEIDI}} > 0.01$). If we use a less stringent threshold for
240 replication, e.g., the nominal P value of 0.05, 66.31% of PAIs predicted in blood were replicated in
241 brain. Here, the replication rate is computed based on a p-value threshold, which is dependent of
242 the sample size of the replication data. Alternatively, we can estimate the correlation of PAI effects
243 (i.e., the effect of the exposure DNAm site on the outcome site of a predicted PAI) between brain
244 and blood using the r_b method³¹. This method does not rely on a p-value threshold and accounts
245 for estimation errors in the estimated effects, which is therefore not dependent of the replication
246 sample size. The estimate of r_b was 0.527 (SE = 0.0051) for 11,082 PAIs between brain and blood,
247 suggesting a relatively strong overlap in PAI between brain and blood.

248

249 It is of note that among the 2,940 blood PAIs replicated at $P_{\text{SMR}} < 4.51 \times 10^{-6}$ and $P_{\text{HEIDI}} > 0.01$ in
250 brain, there were 268 PAIs for which the PAI effects in blood were in opposite directions to those
251 in brain (**Supplementary Table 1**). For example, the estimated PAI effect between the *SORT1* and

252 *SYPL2* loci was 0.49 in blood and -0.86 in brain. This tissue-specific effect is supported by the
253 differences in gene expression correlation (correlation of expression levels between *SORT1* and
254 *SYPL2* was -0.07 in whole blood and -0.37 in brain frontal cortex; $P_{\text{difference}} = 0.0018$) and the
255 chromatin state of the promoter of *SYPL2* (bivalent promoter in blood and active promoter in
256 brain; **Fig. S5**) between brain and blood. Taken together, while there are tissue-specific PAIs, a
257 substantial proportion of the predicted PAIs in blood are consistent with those in brain.

258

259 **Putative target genes of the disease-associated PIDs**

260 We have shown above the potential functional roles of the predicted PAIs in transcriptional
261 regulation. We then turned to ask how the predicted PAIs can be used to infer the genetic and
262 epigenetic regulatory mechanisms at the GWAS loci for complex traits and diseases. We have
263 previously reported 1,203 pleiotropic associations between 1,045 DNAm sites and 15 complex
264 traits and diseases by an integrative analysis of mQTL, eQTL and GWAS data using the SMR and
265 HEIDI approaches¹⁹. Of the 1,045 trait-associated DNAm sites, 601 (57.5%) sites were involved in
266 the predicted PAIs related to 299 Pm-PAI genes (**Supplementary Table 2**). We first tested the
267 functional enrichment of the Pm-PAI genes of the trait-associated PIDs by a gene set enrichment
268 analysis using FUMA³². For the 15 complex traits analysed in Wu et al.¹⁹, our FUMA analyses
269 identified enrichment in multiple GO and KEGG pathways relevant to the corresponding
270 phenotypes such as the inflammatory response pathway for Crohn's disease (CD) and steroid
271 metabolic process for body mass index (BMI) (**Supplementary Table 3**), demonstrating the
272 regulatory role of the trait-associated PIDs in biological processes and tissues relevant to the
273 trait or disease.

274

275 There were a number of examples where the predicted PAIs provided important insights to the
276 functional genes underlying the GWAS loci and the underlying mechanisms by which the DNA
277 variants affect the trait through genetic regulation of gene expression. One notable example was
278 a PID (cg00271210) in an enhancer region predicted to interact in 3D space with the promoter
279 regions of two genes (i.e., *RNASET2* and *RPS6KA2*), the expression levels of both of which were
280 associated with ulcerative colitis (UC) and CD as reported in our previous study¹⁹ (**Fig. 5**). The
281 SNP-association signals were consistent across CD GWAS, eQTL, and mQTL studies, suggesting
282 that the genetic effect on CD is likely to be mediated through epigenetic regulation of gene
283 expression. Our predicted PAIs further implicated a plausible mechanism whereby the expression
284 levels of *RNASET2* and *RPS6KA2* are co-regulated through the interactions of their promoters with
285 a shared enhancer (**Fig. 5**), although only 41.5% of the predicted PAIs in this region overlapped
286 with the TADs identified by Hi-C from the Rao et al. study²⁴ (**Fig. S3a**) as mentioned above.
287 According to the functional annotation data derived from the REMC samples, it appears that this

288 shared enhancer is highly tissue-specific and present only in B cell and digestive system that are
289 closely relevant to CD (**Fig. 5**). The over-expression of *RNASET2* in spleen (**Fig. S6**) is an additional
290 piece of evidence supporting the functional relevance of this gene to CD. Another interesting
291 example is the *ATG16L1* locus (**Fig. S7**). We have shown previously that five DNAm sites are in
292 pleiotropic associations with CD and the expression level of *ATG16L1*¹⁹. Of these five DNAm sites,
293 three were in an enhancer region and predicted to interact in 3D space with two DNAm sites in
294 the promoter region of *ATG16L1* (**Fig. S7**), suggesting a plausible mechanism that the genetic
295 effect on CD at this locus is mediated by genetic and epigenetic regulation of the expression level
296 of *ATG16L1* through promoter–enhancer interactions.

297

298 **Discussion**

299 We have presented an analytical approach on the basis of the recently developed SMR and HEIDI
300 methods to predict promoter-anchored chromatin interactions using mQTL summary data. The
301 proposed approach uses DNAm level of a CpG site in the promoter region of a gene as the bait to
302 detect its pleiotropic associations with DNAm levels of the other CpG sites (**Fig. 1**) within 2 Mb
303 distance of the promoter in either direction. In contrast to experimental assays, such as Hi-C and
304 PChi-C, our approach is cost-effective (because of the reuse of data available from experiments
305 not originally designed for this purpose) and scalable to data with large sample sizes. Our method
306 utilises a genetic model to perform a Mendelian randomization analysis so that the detected
307 associations are not confounded by non-genetic factors, which is also distinct from the methods
308 that predict chromatin interactions from the correlations of chromatin accessibility measures^{14,16}.

309

310 Using mQTL summary-level data from human peripheral blood ($n = 1,980$), we predicted 34,797
311 PAIs for the promoter regions of 4,617 genes. We showed that the predicted PAIs were enriched
312 in TADs detected by published Hi-C and PChi-C assays and that the PIDS regions were enriched
313 with eQTLs of target genes. We also showed that the PIDSs were enriched in enhancers and that
314 the Pm-PAI genes tended to be more actively expressed than matched control genes. These results
315 demonstrate the functional relevance of the predicted PAIs to transcriptional regulation and the
316 feasibility of using data from genetic studies of chromatin status to infer three-dimensional
317 chromatin interactions. The proposed approach is applicable to data from genetic studies of other
318 chromatin features such as histone modification (i.e., hQTL)³³ or chromatin accessibility (caQTL)³⁴.
319 The flexibility of the method also allowed us to analyse data from different tissues or cell types.
320 Using summary data from a brain mQTL study ($n = 468$), we replicated 26.5% of blood PAIs in
321 brain at a very stringent threshold ($P_{\text{SMR}} < 0.05 / m$ with m being the number of tests in the
322 replication set and $P_{\text{HEIDI}} > 0.01$) and 66.31% at a less stringent threshold ($P_{\text{SMR}} < 0.05$). Together
323 with an estimate of r_b of 0.527 for the correlation of PAI effects between brain and blood, we

324 demonstrated a substantial overlap of the predicted PAIs between blood and brain, in line with
325 the finding from a recent study that cis-mQTLs are largely shared between brain and blood³¹.

326

327 The use of a genetic model to detect PAIs also facilitated the integration of the predicted PAIs with
328 GWAS data. In a previous study, Wu et al.¹⁹ mapped DNAm sites to genes and then to a trait by
329 checking the consistency of pleiotropic association signals across all the three layers. This strategy
330 is robust but conservative because such a three-layer model will be rejected if the pleiotropic
331 association is not significant in any of the layers due to the lack of power and/or the stringent of
332 HEIDI threshold (note that if there are errors in the summary data and/or heterogeneity in LD
333 among the mQTL, eQTL, GWAS and LD reference samples, the pleiotropic association will be
334 rejected by the HEIDI test). In this study, we have shown examples of how to integrate the
335 predicted PAIs with GWAS, eQTL and functional annotation data to better understand the genetic
336 and epigenetic regulatory mechanisms underlying the GWAS loci for complex traits (**Figs. 5, S4,**
337 **and S7**). The pleiotropic associations between DNAm sites involved in PAIs and a complex trait
338 are also helpful to link genes to the trait at GWAS loci even in the absence of eQTL data. This can
339 be achieved by testing for pleiotropic associations of both DNAm sites of a PAI with the trait. Of
340 the 1,045 DNAm sites that showed pleiotropic associations with 15 complex traits as reported in
341 Wu et al.¹⁹, 601 sites were involved in the PAIs for 299 Pm-PAI genes identified in this study. In
342 this case, these Pm-PAI genes are very likely to be the functionally relevant genes at the GWAS
343 loci. In comparison with 66 gene targets identified in Wu et al.¹⁹ (34/66 overlapped with 299 Pm-
344 PAI genes), integration of PAIs with GWAS facilitates the discovery of more putative gene targets
345 for complex traits.

346

347 There are some limitations of this study. First, we restricted the PAI analysis to the cis-region (+/-
348 2 Mb of the DNAm site in the promoter region of a gene) so that PAIs between DNAm sites more
349 than 2 Mb are beyond the scope of this study. Compared to a genome-wide scan, this strategy
350 substantially decreased the computational burden and gained power for the PAIs in cis-regions
351 because of the much less stringent significance threshold owing to the order of magnitude of the
352 smaller number of tests. Second, chromatin interactions are likely to be tissue- and temporal-
353 specific whereas our PAI analyses were limited to mQTL data from blood and brain owing to data
354 availability and thus were unable to detect PAIs in specific tissues or at different developmental
355 stages. Third, the sample size of our blood mQTL summary data is large ($n = \sim 2,000$). However,
356 even with such a relatively large sample size, the PAI analysis could be underpowered if the
357 proportion of variance in exposure or outcome explained by the top associated cis-mQTL is small.
358 Fourth, the mQTL data sets used in this study were all from assays based on the Illumina 450K
359 methylation array. Although the 450K array has a genome-wide coverage, the probes only cover

360 a limited proportion of the regulatory elements. Our predicted PAIs are likely to be sparse as
361 illustrated in **Fig. 2d**. In addition, the SMR analysis requires an ascertainment on probes with at
362 least one mQTL at a stringent significance level ($P_{\text{mQTL}} < 5 \times 10^{-8}$), resulting in a further loss of
363 power. Fifth, the functional annotation data derived from the REMC samples could potentially
364 include noise due to the small sample sizes, leading to uncertainty in defining the bait promoter
365 regions. Sixth, if the DNAm levels of two CpG sites are affected by two sets of causal variants in
366 very high LD, these two DNAm sites will appear to be associated in the SMR analysis and the power
367 of the HEIDI test to reject such an SMR association will be limited because of the high LD^{18,19}.
368 However, this phenomenon is likely to be rare given that most of the promoter-anchored DNAm
369 sites were predicted to interact with multiple DNAm sites which are very unlikely to be all caused
370 by distinct sets of causal variants in high LD. Despite these limitations, our study provides a novel
371 computational paradigm to predict PAIs from genetic effects on epigenetic markers with high
372 resolution. Integrating of the predicted PAIs with GWAS, gene expression, and functional
373 annotation data provides novel insights into the regulatory mechanisms underlying GWAS loci for
374 complex traits. The computational framework is general and applicable to other types of
375 chromatin and histone modification data, to further decipher the functional organisation of the
376 genome.
377
378

379 **Methods**

380 **Predicting PAIs from mQTL data by the SMR and HEIDI analyses**

381 We used summary-level mQTL data to test whether the variation between people in DNAm levels
382 of two CpG sites are associated because of a set of shared causal variants. Mendelian
383 Randomization (MR) is an approach developed to test for the causal effect of an exposure and an
384 outcome using a genetic variant as the instrumental variable^{21,22}. Summary-data-based
385 Mendelian Randomization (SMR) is a variant of MR, originally designed to test for association
386 between the expression level of a gene and a complex trait using summary-level data from GWAS
387 and eQTL studies¹⁸ and subsequently applied to test for associations between DNAm and gene
388 expression and between DNAm and complex traits¹⁹. Here, we applied the SMR analysis to detect
389 associations between DNAm sites. We specified the DNAm level of a probe within the promoter
390 region of a gene as the “exposure” and tested its associations with the DNAm levels of other probes
391 (“outcomes”) within 2 Mb of the exposure probe (**Fig. 1** and **Fig. S1**). Probe pairs in the same
392 promoter region were not included in the analysis. For a pair of probes in two different promoter
393 regions, the one with higher variance explained by its top associated cis-mQTL was used as the
394 exposure and the other one was used as the outcome. The associations passed the SMR test could
395 possibly be due to linkage (i.e., distinct sets of causal variants in LD, one set affecting the exposure
396 and the other set affecting the outcome), which is less of biological interest in comparison with
397 pleiotropy (i.e., the same set of causal variants affecting both the exposure and the outcome). We
398 then applied the HEIDI (heterogeneity in dependent instruments) test to distinguish pleiotropy
399 from linkage. The HEIDI test uses multiple cis-mQTL SNPs in LD with the top cis-mQTL for the
400 exposure to detect whether the SNP associations with the exposure and those with the outcome
401 are due to the same set of causal variants, with pairwise LD between SNPs estimated from the
402 Health and Retirement Study (HRS)²³ with SNP data imputed to the 1000 Genomes Project
403 (1KGP)³⁵. We rejected the SMR associations with $P_{\text{HEIDI}} < 0.01$. All these analyses have been
404 implemented in the SMR software tool (**URLs**). Because the mQTL data for the exposure and the
405 outcome were obtained from the same sample, we investigated whether the SMR and HEIDI test-
406 statistics were biased by the sample overlap. To this end, we computed the phenotypic correlation
407 between each pair of exposure and outcome probes as well as the variance explained by the top
408 associated cis-mQTL of each exposure probe, and performed the simulation based on these
409 observed distributions (**Supplementary Note 1**). The simulation results showed that P values
410 from both SMR and HEIDI test were evenly distributed under the null model without inflation or
411 deflation (**Fig. S8**).

412

413 **Data used for the PAI analysis**

414 The peripheral blood mQTL summary data were from the Brisbane Systems Genetics Study
415 (BSGS)³⁶ ($n=614$) and Lothian Birth Cohorts (LBC) of 1921 and 1936³⁷ ($n=1,366$). We performed
416 a meta-analysis of the two cohorts and identified 90,749 DNAm probes with at least a cis-mQTL
417 at $P_{\text{mQTL}} < 5 \times 10^{-8}$ (excluding the probes in the major histocompatibility complex (MHC) region
418 because of the complexity of this region), of which 28,732 DNAm probes were in the promoter
419 regions defined by the annotation data derived from 23 REMC blood samples (T-cell, B-cell, and
420 Hematopoietic stem cells). The prefrontal cortex mQTL summary data were from the Religious
421 Orders Study and Memory and Aging Project (ROSMAP)³⁰ ($n=468$), comprising 419,253 probes
422 and approximate 6.5 million genetic variants. In the ROSMAP data, there were 67,995 DNAm
423 probes with at least a cis-mQTL at $P_{\text{mQTL}} < 5 \times 10^{-8}$ (not including the probes in the MHC region), of
424 which 22,285 DNAm probes were in the promoter regions defined by the annotation data derived
425 from 10 REMC brain samples. For all the DNAm probes, enhanced annotation data from Price *et*
426 *al.*³⁸ (**URLs**) were used to annotate the closest gene of each DNAm probe.

427
428 We included in the analysis 15 complex traits (including disease) as analysed in Wu *et al.*¹⁹. They
429 are height³⁹, body mass index (BMI)⁴⁰, waist-hip-ratio adjusted by BMI (WHRadjBMI)⁴¹, high-
430 density lipoprotein (HDL)⁴², low-density lipoprotein (LDL)⁴², thyroglobulin (TG)⁴², educational
431 years (EY)⁴³, rheumatoid arthritis (RA)⁴⁴, schizophrenia (SCZ)⁴⁵, coronary artery disease (CAD)⁴⁶,
432 type 2 diabetes (T2D)⁴⁷, Crohn's disease (CD)⁴⁸, ulcerative colitis (UC)⁴⁸, Alzheimer's disease
433 (AD)⁴⁹ and inflammatory bowel disease (IBD)⁴⁸. The GWAS summary data were from the large
434 GWAS meta-analyses (predominantly in samples of European ancestry) with sample sizes of up
435 to 339,224. The number of SNPs varied from 2.5 to 9.4 million across traits.

436 437 **Annotations of the chromatin state**

438 The epigenomic annotation data used in this study were from the Roadmap Epigenomics Mapping
439 Consortium (REMC), publicly available at <http://compbio.mit.edu/roadmap/>. We used these data
440 to annotate the functional relevance of the DNAm sites and their cell type or tissue specificity. The
441 chromatin state annotations from the Roadmap Epigenomics Project¹³ were predicted by
442 ChromHMM¹² based on the imputed data of 12 histone-modification marks. It contains 25
443 functional categories for 127 epigenomes in a wide range of primary tissue and cell types (**URLs**).
444 The 25 chromatin states were further combined into 14 main functional annotations (as shown
445 in **Fig. 3B**), as described in Wu *et al.*¹⁹ study.

446 447 **Overlap of the predicted PAIs with Hi-C and PCHi-C data**

448 To test the overlap between our predicted PAIs and chromatin contacts detected by Hi-C or PCHi-
449 C, we used chromatin contact loops and topological associated domains (TADs) data from the Rao

450 et al. study called in the GM12812 cells²⁴ and the Dixon et al. study in embryonic stem cells²⁶, and
451 PChI-C interaction data generated from human primary hematopoietic cells⁵. To demonstrate the
452 significance of enrichment, we generated a null distribution by random sampling of 1,000 sets of
453 “control” probe pairs (with the same number as that of the significant pairs) from all the distance-
454 matched probe pairs tested in the SMR analysis. We mapped both the predicted PAIs and the
455 control probe pairs to the TAD regions or chromatin contact loops detected by experimental
456 assays and quantified the number of overlapping pairs. We estimated the fold enrichment by the
457 ratio of the overlapping number for the predicted PAIs to the mean of the null distribution and
458 computed the empirical *P*-value by comparing the overlapping number for the predicted PAIs
459 with the null distribution.

460

461 **Enrichment of the PIDSs in functional annotations**

462 To conduct an enrichment test of the promoter interacting DNAm sites (PIDSs) in different
463 functional annotation categories, we first extracted chromatin state data of 23 blood samples from
464 the REMC samples. We then mapped the PIDSs to 14 main functional categories based on the
465 physical positions, and counted the number of PIDSs in each functional category. Again, we
466 generated a null distribution by randomly sampling the same number of control probes (with
467 variance in DNAm level matched with the PIDSs) from all the probes tested in the PAI analysis and
468 repeated the random sampling 1,000 times. The fold enrichment was calculated by the ratio of the
469 observed value to the mean of the null distribution, and an empirical *P*-value was computed by
470 comparing the observed value with the null distribution.

471

472 **Quantifying the expression levels of Pm-PAI genes**

473 To quantify the expression levels of genes whose promoters were involved in the predicted PAIs
474 (Pm-PAI genes), we used gene expression data (measured by Transcript Per Kilobase Million
475 mapped reads (TPM)) from blood samples of the Genotype-Tissue Expression (GTEx) project²⁸.
476 We classified all the genes into two groups based on their expression levels in GTEx blood, i.e.,
477 active and inactive (TPM < 0.1). For the active genes, we further divided them into four quartiles
478 based on their expression levels in GTEx blood, and counted the number of Pm-PAI genes in each
479 of the five groups. To generate the null distribution, we randomly sampled the same number of
480 “control” genes whose promoter DNAm sites were included in the SMR analysis, and repeated the
481 random sampling 1,000 times. We computed the number of Pm-PAI genes and “control” genes in
482 each group and assessed the significance by comparing the number of Pm-PAI genes with the null
483 distribution in each group.

484

485 **Enrichment of eQTLs and gene-associated DNAm in the PIDS regions**

486 The eQTL enrichment analysis was conducted using all the independent cis-eQTLs ($m=11,204$)
487 from CAGE²⁹ study. The independent cis-eQTLs were from SNP-probe associations ($P < 5 \times 10^{-8}$)
488 after clumping analysis in PLINK⁵⁰ followed by a conditional and joint (COJO) analysis in GCTA⁵¹.
489 We only retained the cis-eQTLs whose target genes had at least a PIDS and mapped the cis-eQTL
490 to a 10 Kb region centred around each corresponding PIDS of a Pm-PAI gene. To assess the
491 significance of the enrichment, we generated a null distribution by mapping the cis-eQTLs to the
492 same number of “control” gene-DNA_m pairs (strictly speaking, it is the bait DNA_m probe in the
493 promoter of a gene together with another non-promoter DNA_m probe) randomly sampled (with
494 1,000 repeats) from those included in the PAI analysis with the distance between a control pair
495 matched with that between a Pm-PAI gene and the corresponding PIDS. In addition, we have
496 identified a set of DNA_m sites that showed pleiotropic associations with gene expressions in a
497 previous study¹⁹. We used the same approach as described above to test the significance of
498 enrichment of the gene-associated DNA_m sites in the PIDSs.

499

500 **Supplemental information**

501 Supplemental data include 8 supplemental figures and 3 supplemental tables.

502

503 **URLs**

504 M2Mdb, <http://cnsgenomics.com/shiny/M2Mdb/>

505 SMR, <http://cnsgenomics.com/software/smr>

506 GTEx, <http://www.gtexportal.org/home/>

507 Annotation file for the Illumina HumanMethylation450 BeadChip,

508 <https://www.ncbi.nlm.nih.gov/geo/query/acc.cgi?acc=GPL16304>

509

510 **Acknowledgements**

511 We thank Peter Visscher for helpful discussion. This research was supported by the Australian
512 Research Council (DP160101343, DP160101056 and FT180100186), the Australian National
513 Health and Medical Research Council (1107258, 1083656, 1078901 and 1113400), and the Sylvia
514 & Charles Viertel Charitable Foundation. The Lothian Birth Cohorts (LBC) are supported by Age
515 UK (Disconnected Mind programme). Methylation typing was supported by Centre for Cognitive
516 Ageing and Cognitive Epidemiology (Pilot Fund award), Age UK, The Wellcome Trust Institutional
517 Strategic Support Fund, The University of Edinburgh, and The University of Queensland. The LBC
518 resource is prepared in the Centre for Cognitive Ageing and Cognitive Epidemiology, which is
519 supported by the Medical Research Council and Biotechnology and Biological Sciences Research
520 Council (MR/K026992/1), and which supports I.J.D.. This study makes use of data from dbGaP

521 (accessions: phs000428.v1.p1 and phs000424.v1.p1) and EGA (accession: EGAS00001000108).
522 A full list of acknowledgements to these data sets can be found in the **Supplementary Note 2**.

523

524 **Author Contributions**

525 J.Y. conceived the study. Y.W. and J.Y. designed the experiment. Y.W. and T.Q. performed
526 simulations and statistical analyses under the assistance or guidance from J.Y., J.Z., H.W., F.Z., and
527 Z.Z.. I.J.D., N.R.W. and A.F.M. contributed the blood DNA methylation data. J.E.P.C. provided critical
528 advice that significantly improved the interpretation of the results. N.R.W. and J.Y. contributed
529 funding and resources. Y.W., T.Q., J.Z. and J.Y. wrote the manuscript with the participation of all
530 authors.

531

532 **Declaration of Interests**

533 We declare that all authors have no competing interests.

534

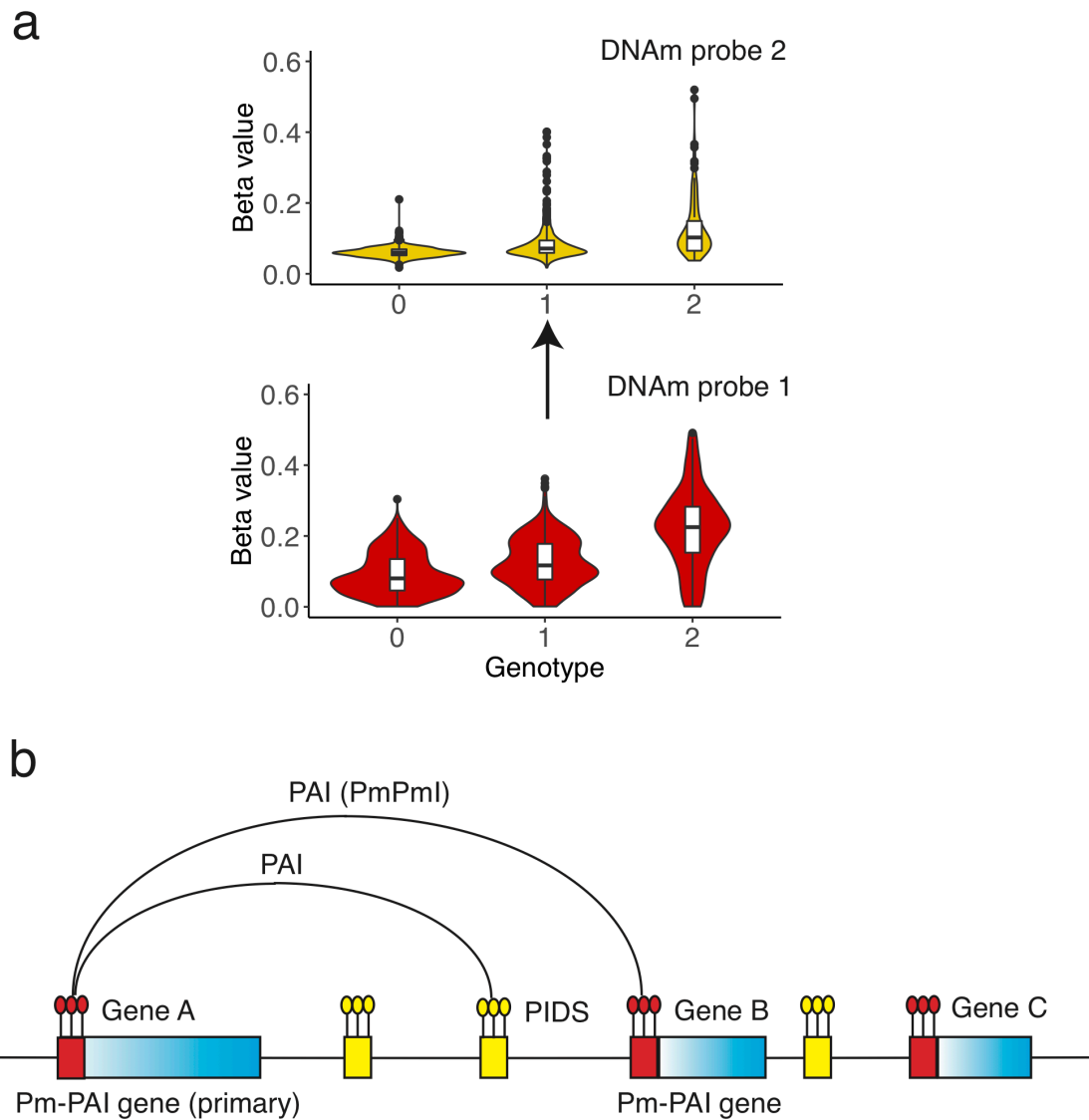
535 **References**

- 536 1. MacArthur, J. *et al.* The new NHGRI-EBI Catalog of published genome-wide association
537 studies (GWAS Catalog). *Nucleic Acids Research* **45**, D896-D901 (2017).
- 538 2. Visscher, P.M. *et al.* 10 Years of GWAS Discovery: Biology, Function, and Translation. *Am J*
539 *Hum Genet* **101**, 5-22 (2017).
- 540 3. Farh, K.K.-H. *et al.* Genetic and epigenetic fine mapping of causal autoimmune disease
541 variants. *Nature* **518**, 337-343 (2015).
- 542 4. Claussnitzer, M. *et al.* FTO Obesity Variant Circuitry and Adipocyte Browning in Humans.
543 *New England Journal of Medicine* **373**, 895-907 (2015).
- 544 5. Javierre, B.M. *et al.* Lineage-Specific Genome Architecture Links Enhancers and Non-
545 coding Disease Variants to Target Gene Promoters. *Cell* **167**, 1369-1384 e19 (2016).
- 546 6. Khurana, E. *et al.* Role of non-coding sequence variants in cancer. *Nat Rev Genet* **17**, 93-
547 108 (2016).
- 548 7. Lieberman-Aiden, E. *et al.* Comprehensive Mapping of Long-Range Interactions Reveals
549 Folding Principles of the Human Genome. *Science* **326**, 289-293 (2009).
- 550 8. Fullwood, M.J. *et al.* An oestrogen-receptor-alpha-bound human chromatin interactome.
551 *Nature* **462**, 58-64 (2009).
- 552 9. Wit, E.d. & Laat, W.d. A decade of 3C technologies: insights into nuclear organization. *Genes*
553 *& Development* **26**, 11-24 (2012).
- 554 10. Kuo, M.-H. & Allis, C.D. In Vivo Cross-Linking and Immunoprecipitation for Studying
555 Dynamic Protein:DNA Associations in a Chromatin Environment. *Methods* **19**, 425-433
556 (1999).

- 557 11. Belton, J.M. *et al.* Hi-C: a comprehensive technique to capture the conformation of genomes.
558 *Methods* **58**, 268-76 (2012).
- 559 12. Ernst, J. & Kellis, M. ChromHMM: automating chromatin-state discovery and
560 characterization. *Nat Methods* **9**, 215-6 (2012).
- 561 13. Roadmap Epigenomics Consortium *et al.* Integrative analysis of 111 reference human
562 epigenomes. *Nature* **518**, 317-30 (2015).
- 563 14. Zhu, Y. *et al.* Constructing 3D interaction maps from 1D epigenomes. *Nat Commun* **7**, 10812
564 (2016).
- 565 15. Huang, J., Marco, E., Pinello, L. & Yuan, G.-C. Predicting chromatin organization using
566 histone marks. *Genome Biology* **16**, 162 (2015).
- 567 16. Fortin, J.-P. & Hansen, K.D. Reconstructing A/B compartments as revealed by Hi-C using
568 long-range correlations in epigenetic data. *Genome Biology* **16**, 180 (2015).
- 569 17. Kumasaka, N., Knights, A.J. & Gaffney, D.J. High-resolution genetic mapping of putative
570 causal interactions between regions of open chromatin. *Nature Genetics* **51**, 128-137
571 (2019).
- 572 18. Zhu, Z. *et al.* Integration of summary data from GWAS and eQTL studies predicts complex
573 trait gene targets. *Nat Genet* **48**, 481-7 (2016).
- 574 19. Wu, Y. *et al.* Integrative analysis of omics summary data reveals putative mechanisms
575 underlying complex traits. *Nature Communications* **9**, 918 (2018).
- 576 20. McRae, A.F. *et al.* Identification of 55,000 Replicated DNA Methylation QTL. *Scientific*
577 *Reports* **8**, 17605 (2018).
- 578 21. Davey Smith, G. & Hemani, G. Mendelian randomization: genetic anchors for causal
579 inference in epidemiological studies. *Human Molecular Genetics* **23**, R89-R98 (2014).
- 580 22. Davey Smith, G. & Ebrahim, S. 'Mendelian randomization': can genetic epidemiology
581 contribute to understanding environmental determinants of disease? *International*
582 *Journal of Epidemiology* **32**, 1-22 (2003).
- 583 23. Sonnega, A. *et al.* Cohort Profile: the Health and Retirement Study (HRS). *Int J Epidemiol*
584 **43**, 576-85 (2014).
- 585 24. Rao, Suhas S.P. *et al.* A 3D Map of the Human Genome at Kilobase Resolution Reveals
586 Principles of Chromatin Looping. *Cell* **159**, 1665-1680 (2014).
- 587 25. Pombo, A. & Dillon, N. Three-dimensional genome architecture: players and mechanisms.
588 *Nat Rev Mol Cell Biol* **16**, 245-57 (2015).
- 589 26. Dixon, J.R. *et al.* Topological domains in mammalian genomes identified by analysis of
590 chromatin interactions. *Nature* **485**, 376-80 (2012).
- 591 27. Li, G. *et al.* Extensive Promoter-centered Chromatin Interactions Provide a Topological
592 Basis for Transcription Regulation. *Cell* **148**, 84-98 (2012).

- 593 28. The GTEx Consortium. The Genotype-Tissue Expression (GTEx) pilot analysis: Multitissue
594 gene regulation in humans. *Science* **348**, 648-660 (2015).
- 595 29. Lloyd-Jones, L.R. *et al.* The Genetic Architecture of Gene Expression in Peripheral Blood.
596 *Am J Hum Genet* **100**, 371 (2017).
- 597 30. Ng, B. *et al.* An xQTL map integrates the genetic architecture of the human brain's
598 transcriptome and epigenome. *Nature Neuroscience* **20**, 1418 (2017).
- 599 31. Qi, T. *et al.* Identifying gene targets for brain-related traits using transcriptomic and
600 methylomic data from blood. *bioRxiv* (2018).
- 601 32. Watanabe, K., Taskesen, E., van Bochoven, A. & Posthuma, D. Functional mapping and
602 annotation of genetic associations with FUMA. *Nature Communications* **8**, 1826 (2017).
- 603 33. Chen, L. *et al.* Genetic Drivers of Epigenetic and Transcriptional Variation in Human
604 Immune Cells. *Cell* **167**, 1398-1414 e24 (2016).
- 605 34. Gate, R.E. *et al.* Genetic determinants of co-accessible chromatin regions in activated T
606 cells across humans. *Nat Genet* **50**, 1140-1150 (2018).
- 607 35. Abecasis, G.R., Altshuler, D., Auton, A., Brooks, L.D. & Durbin, R.M. A map of human genome
608 variation from population-scale sequencing. *Nature* **467**(2010).
- 609 36. Powell, J.E. *et al.* The Brisbane Systems Genetics Study: genetical genomics meets complex
610 trait genetics. *PLoS One* **7**, e35430 (2012).
- 611 37. Chen, B.H. *et al.* DNA methylation-based measures of biological age: meta-analysis
612 predicting time to death. *Aging (Albany NY)* **8**, 1844-1865 (2016).
- 613 38. Price, M.E. *et al.* Additional annotation enhances potential for biologically-relevant
614 analysis of the Illumina Infinium HumanMethylation450 BeadChip array. *Epigenetics*
615 *Chromatin* **6**, 4 (2013).
- 616 39. Wood, A.R. *et al.* Defining the role of common variation in the genomic and biological
617 architecture of adult human height. *Nat Genet* **46**, 1173-86 (2014).
- 618 40. Locke, A.E. *et al.* Genetic studies of body mass index yield new insights for obesity biology.
619 *Nature* **518**, 197-206 (2015).
- 620 41. Shungin, D. *et al.* New genetic loci link adipose and insulin biology to body fat distribution.
621 *Nature* **518**, 187-96 (2015).
- 622 42. Global Lipids Genetics Consortium *et al.* Discovery and refinement of loci associated with
623 lipid levels. *Nat Genet* **45**, 1274-83 (2013).
- 624 43. Okbay, A. *et al.* Genome-wide association study identifies 74 loci associated with
625 educational attainment. *Nature* **533**, 539-42 (2016).
- 626 44. Okada, Y. *et al.* Genetics of rheumatoid arthritis contributes to biology and drug discovery.
627 *Nature* **506**, 376-381 (2014).

- 628 45. Schizophrenia Working Group of the Psychiatric Genomics Consortium. Biological insights
629 from 108 schizophrenia-associated genetic loci. *Nature* **511**, 421-7 (2014).
- 630 46. Nikpay, M. *et al.* A comprehensive 1,000 Genomes-based genome-wide association meta-
631 analysis of coronary artery disease. *Nat Genet* **47**, 1121-30 (2015).
- 632 47. Morris, A.P. *et al.* Large-scale association analysis provides insights into the genetic
633 architecture and pathophysiology of type 2 diabetes. *Nature genetics* **44**, 981 (2012).
- 634 48. Liu, J.Z. *et al.* Association analyses identify 38 susceptibility loci for inflammatory bowel
635 disease and highlight shared genetic risk across populations. *Nat Genet* **47**, 979-986
636 (2015).
- 637 49. Lambert, J.C. *et al.* Meta-analysis of 74,046 individuals identifies 11 new susceptibility loci
638 for Alzheimer's disease. *Nat Genet* **45**, 1452-8 (2013).
- 639 50. Purcell, S. *et al.* PLINK: a tool set for whole-genome association and population-based
640 linkage analyses. *Am J Hum Genet* **81**, 559-75 (2007).
- 641 51. Yang, J. *et al.* Conditional and joint multiple-SNP analysis of GWAS summary statistics
642 identifies additional variants influencing complex traits. *Nature Genetics* **44**, 369 (2012).
- 643

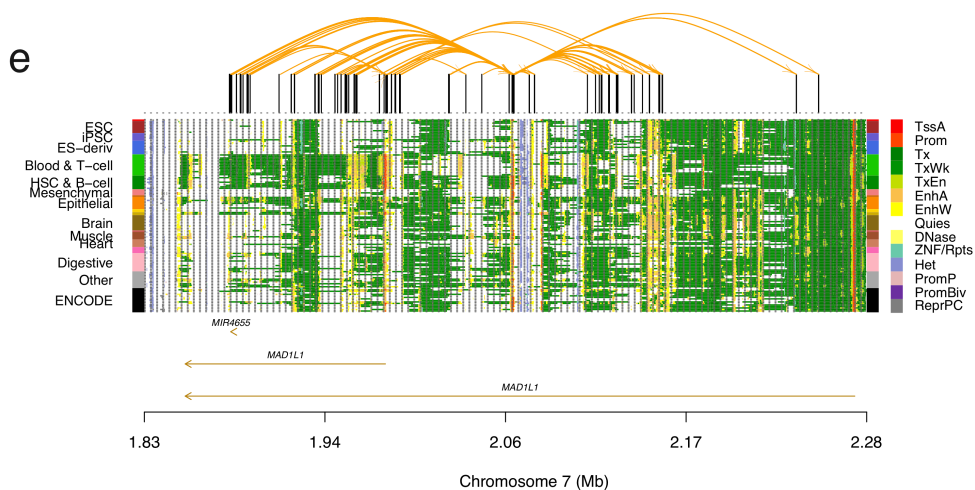
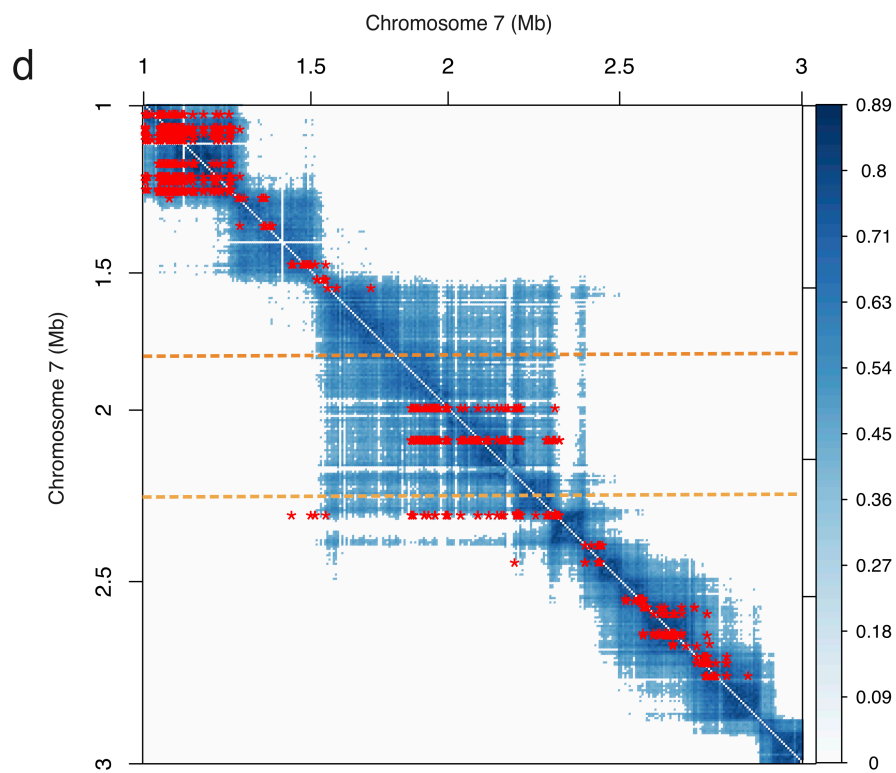
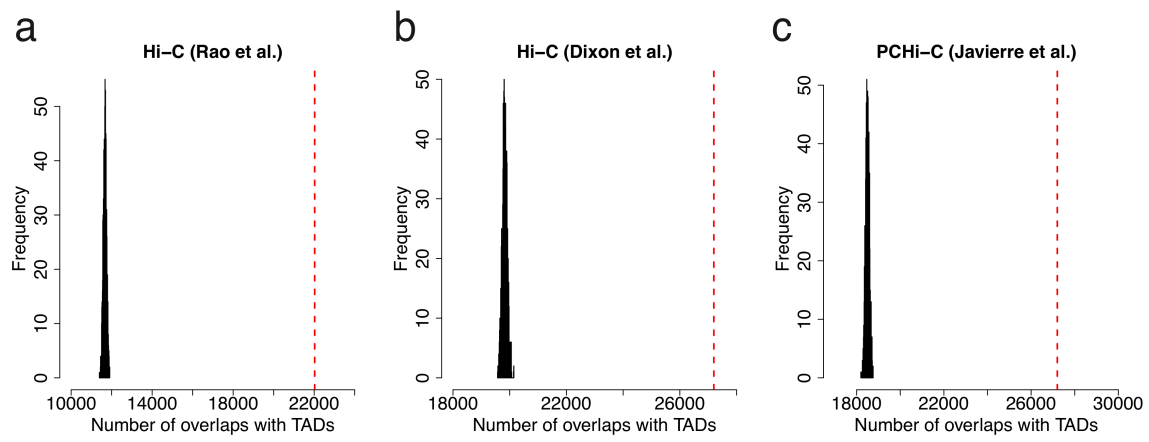


644

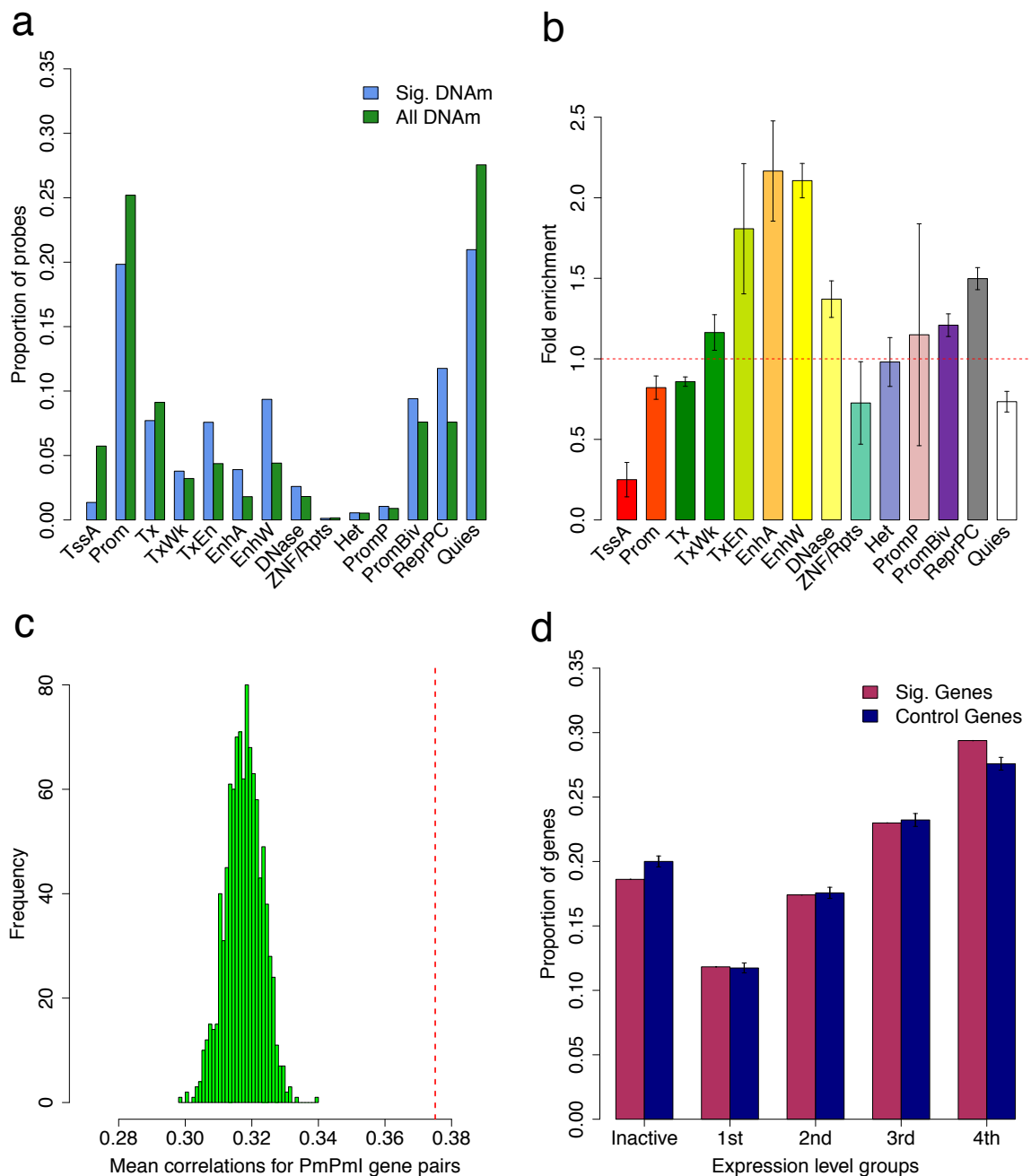
645 **Figure 1** Schematic of the promoter-anchored interaction (PAI) analysis. Panel a): a schematic of
646 the pleiotropic model that variation between people in DNAm levels of two CpG sites are
647 associated because of a shared causal variant. The DNAm level is measured by beta value, ranging
648 from 0 to 1 (with 0 being unmethylated and 1 being fully methylated). It is the ratio of the
649 methylated probe intensity to the overall intensity (sum of methylated and unmethylated probe
650 intensities). Panel b): a schematic of the PAI analysis. The blue rectangles represent genes with
651 their promoter regions color coated in red. The small yellow bars represent other functional
652 regions (e.g., enhancers). In this toy example, the promoter region of Gene A is used as the bait for
653 the PAI analysis. Genes whose promoters are involved in significant PAIs are defined as Pm-PAI
654 genes. PIDS: promoter interacting DNAm site. PmPml: promoter-promoter interaction.

655

656



658 **Figure 2** Overlap of the predicted PAIs with Hi-C and PCHi-C data. Panels a), b) and c): overlaps
659 of the predicted PAIs with TADs identified by a) Rao et al.²⁴ and b) Dixon et al.²⁶ using Hi-C and by
660 c) Javierre et al.⁵ using PCHi-C. The red dash lines represent the observed number and histograms
661 represent the distribution of “control” sets. Panel d): a heatmap of the predicted PAIs (red
662 asterisks) and chromatin interactions with correlation score > 0.4 (blue dots) identified by Rao et
663 al.²⁴ using Hi-C in a 2 Mb region on chromosome 7. The heatmap is asymmetric for the PAIs with
664 the x- and y-axes representing the physical positions of “outcome” and “exposure” probes
665 respectively. Panel e): the predicted PAIs at the *MAD1L1* locus, a 450-Kb sub-region of that shown
666 between two orange dashed lines in panel d). The orange curved lines on the top represent the
667 significant PAIs between 14 DNAm sites in the promoter regions of *MAD1L1* (multiple transcripts)
668 and other nearby DNAm sites. The panel on the bottom represents 14 chromatin state annotations
669 (indicated by different colours) inferred from data of 127 REMC samples (one row per sample).
670



671

672 **Figure 3** Enrichment of PIDs and Pm-PAI genes. Panels a) and b): enrichment of PIDs in 14 main

673 functional annotation categories inferred from the 127 REMC samples. Each error bar in panel b)

674 represents the standard deviation of the estimate under the null obtained from 1,000 random

675 samples. The 14 functional categories are: TssA, active transcription start site; Prom,

676 upstream/downstream TSS promoter; Tx, actively transcribed state; TxWk, weak transcription;

677 TxEn, transcribed and regulatory Prom/Enh; EnhA, active enhancer; EnhW, weak enhancer;

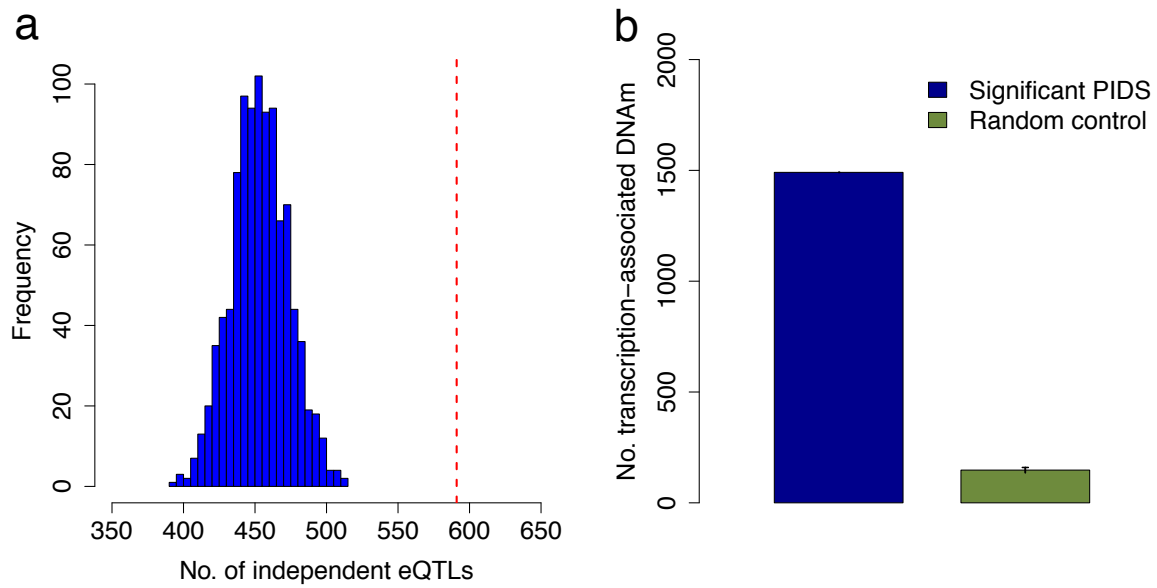
678 DNase, primary DNase; ZNF/Rpts, state associated with zinc finger protein genes; Het,

679 constitutive heterochromatin; PromP, Poised promoter; PromBiv, bivalent regulatory states;

680 ReprPC, repressed Polycomb states; and Quies, a quiescent state. Panel c): mean Pearson

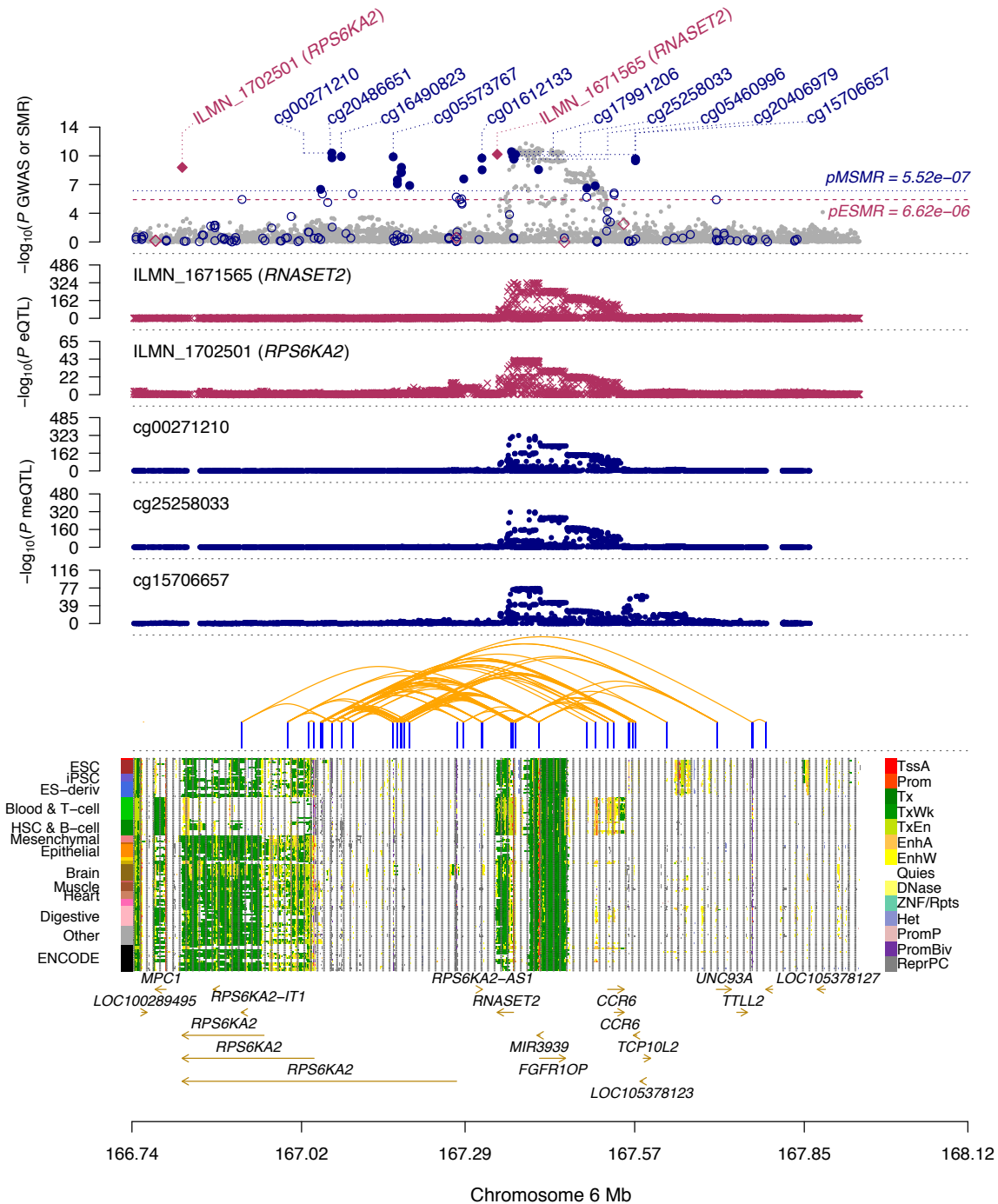
681 correlation of expression levels for gene pairs whose promoters were involved in PmPml. The red

682 dash line represents the observed value and the histogram represents the distribution for the
683 control gene pairs. Panel d): proportion of Pm-PAI genes in five gene activity groups. The five gene
684 activity groups are inactive (TPM <0.1) together with four quartiles defined based on the
685 expression levels of all genes in the GTEx blood samples. Each error bar represents the standard
686 deviation estimated from the control genes in 1,000 random samples.



687

688 **Figure 4** Enrichment of eQTLs or transcription-associated DNAm sites in PIDS regions of the Pm-
689 PAI genes. Panel a): the number of independent cis-eQTLs ($P_{eQTL} < 5 \times 10^{-8}$) located in PIDS regions
690 of the Pm-PAI genes. The red dash line represents the observed number and the blue histogram
691 represents the distribution of 1000 control sets. Panel b): the number of transcription-associated
692 DNAm sites located in PIDS regions of the Pm-PAI genes. The blue bar represents the observed
693 number and the green bar represents the mean of 1000 control sets. The error bar represents the
694 standard deviation estimated from the control sets.



695
 696 **Figure 5** Prioritizing genes and functional regions at the *RPS6KA2* locus for Crohn's disease (CD).
 697 The top plot shows $-\log_{10}(P$ values) of SNPs from the GWAS meta-analysis (grey dots) for CD⁴⁴.
 698 Red diamonds and blue circles represent $-\log_{10}(P$ values) from SMR tests for associations of gene
 699 expression and DNAm probes with CD, respectively. Solid diamonds and circles are the probes not
 700 rejected by the HEIDI test ($P_{\text{HEIDI}} > 0.01$). The second and third plots show $-\log_{10}(P$ values) of SNP
 701 associations for the expression levels of probe ILMN_1671565 (tagging *RNASET2*) and
 702 ILMN_1702501 (tagging *RPS6KA2*), respectively, from the CAGE data. The fourth, fifth and sixth
 703 plots shows $-\log_{10}(P$ values) of SNP associations for the DNAm levels of probes cg00271210,

704 cg25258033, and cg15706657, respectively, from the mQTL meta-analysis. The panel on the
705 bottom shows 14 chromatin state annotations (indicated by colours) inferred from 127 REMC
706 samples (one sample per row) with the predicted PAIs annotated by orange curved lines on the
707 top (see **Fig. S3a** for the overlap of the predicted PAIs with Hi-C data).