

1 **Colocalization of blood cell traits GWAS associations and variation in PU.1 genomic occupancy**  
2 **prioritizes causal noncoding regulatory variants**

3  
4 Raehoon Jeong<sup>1,2</sup> and Martha L. Bulyk<sup>1,2,3†</sup>

5 <sup>1</sup> Division of Genetics, Department of Medicine, Brigham and Women's Hospital and Harvard Medical  
6 School, Boston, MA 02115, USA.

7 <sup>2</sup> Bioinformatics and Integrative Genomics Graduate Program, Harvard University, Cambridge, MA  
8 02138, USA.

9 <sup>3</sup> Department of Pathology, Brigham and Women's Hospital and Harvard Medical School, Boston, MA  
10 02115, USA.

11

12 † Correspondence: [mlbulyk@genetics.med.harvard.edu](mailto:mlbulyk@genetics.med.harvard.edu)

13

14 **Abstract**

15 Genome-wide association studies (GWAS) have uncovered numerous trait-associated loci across the  
16 human genome, most of which are located in noncoding regions, making interpretations difficult.  
17 Moreover, causal variants are hard to statistically fine-map at many loci because of widespread linkage  
18 disequilibrium. To address this challenge, we present a strategy utilizing transcription factor (TF) binding  
19 quantitative trait loci (bQTLs) for colocalization analysis to identify trait associations likely mediated by  
20 TF occupancy variation and to pinpoint likely causal variants using motif scores. We applied this  
21 approach to PU.1 bQTLs in lymphoblastoid cell lines and blood cell traits GWAS data. Colocalization  
22 analysis revealed 69 blood cell trait GWAS loci putatively driven by PU.1 occupancy variation. We  
23 nominate PU.1 motif-altering variants as the likely shared causal variants at 51 loci. Such integration of  
24 TF bQTL data with other GWAS data may reveal transcriptional regulatory mechanisms and causal  
25 noncoding variants underlying additional complex traits.

26 A recurring challenge in genome-wide association studies (GWAS) is the difficulty of identifying causal  
27 variants, as well as formulating corresponding variant-to-function (V2F) hypotheses<sup>1</sup>. Pinpointing causal  
28 variants is important as it guides subsequent validation experiments<sup>2-4</sup> and development of potential  
29 therapies<sup>5</sup>. More precise identification of causal variants (*e.g.*, fine-mapping) also leads to better genetic  
30 risk predictions across various traits and diseases<sup>6,7</sup>. However, widespread linkage disequilibrium (LD)  
31 typically prevents effective statistical fine-mapping, especially for common variants<sup>1,8</sup>. Moreover, most of  
32 the genome-wide significant loci are noncoding and likely have regulatory functions; in practice, noncoding  
33 variants are much harder to interpret than coding variants because predicting the effects of noncoding  
34 variants on transcription factor (TF) binding *in vivo* is challenging. Since variants predicted to affect TF  
35 binding across the genome have been shown to explain a large proportion of genetic associations to traits  
36 (*i.e.*, heritability enrichment)<sup>9,10</sup>, many studies have examined whether trait-associated variants overlap a  
37 TF binding site motif within the corresponding TF ChIP-seq peak<sup>8,11</sup>, but data demonstrating the variants'  
38 effects on *in vivo* TF binding are necessary to imply causality of the variant on TF binding. Therefore, an  
39 approach to effectively pinpoint regulatory variants and their effects on *in vivo* TF binding at individual  
40 GWAS loci is essential.

41 Previous studies have utilized expression quantitative trait loci (eQTL) or methylation QTL (mQTL)  
42 colocalization to learn about regulatory mechanisms (*e.g.*, causal genes) at GWAS loci<sup>12-16</sup>. Statistical  
43 colocalization specifically tests the hypothesis that genetic signals are shared between a pair of traits (*e.g.*  
44 eQTL and GWAS), whereas positional overlap of associations to two traits alone leads to many false  
45 positives<sup>13,17</sup>. However, a key weakness in eQTL or mQTL colocalization analysis is the inability to  
46 pinpoint a causal regulatory variant effectively because colocalization analyses are not inherently aimed at  
47 identifying the causal variant, and LD typically prevents statistical fine-mapping at single-variant  
48 resolution<sup>14</sup>.

49 Here, we have developed a strategy 1) to analyze colocalization of TF binding QTLs (bQTLs) (*i.e.*,  
50 genomic loci where TF occupancy level, as measured by ChIP-seq, is significantly associated with a genetic  
51 variant) at GWAS loci to highlight TF binding sites that potentially mediate the GWAS associations<sup>18</sup>, and  
52 2) to utilize TF motif models to nominate variants altering a motif of the corresponding TF at those binding  
53 sites as likely shared causal regulatory variants underlying both TF binding variation and the GWAS traits  
54 (Fig. 1a). TF bQTLs are fundamentally different from eQTLs and mQTLs in that TF bQTLs point to likely  
55 causal variants because they are often driven by the corresponding TF motif-altering variants<sup>19,20</sup>. To our  
56 knowledge, this is the first attempt to perform TF bQTL colocalization analysis with GWAS data to fine-  
57 map putative causal variants that affect *in vivo* TF binding.

58 We carried out this strategy with blood cell trait GWAS<sup>21</sup> and bQTL data for the hematopoietic  
59 master regulator PU.1 from lymphoblastoid cell lines (LCLs)<sup>19,22</sup>, which are immortalized B cell lines. PU.1

60 bQTLs in neutrophils have been found previously to colocalize with immune disease susceptibility loci but  
61 were not used to fine-map the causal variants<sup>18</sup>. Blood cell traits (*e.g.*, lymphocyte counts, hemoglobin  
62 concentrations) are indicators of various diseases; for instance, individuals with low lymphocyte counts are  
63 more susceptible to infections, including severe COVID-19<sup>23–25</sup>. Consistent with PU.1’s role in specifying  
64 myeloid and lymphoid lineages during hematopoiesis<sup>26,27</sup> and its expression throughout progenitor cell  
65 types<sup>28</sup> (Supplementary Fig. 1), a recent fine-mapping analysis of blood cell trait GWAS reported that PU.1  
66 was the TF with the highest number of fine-mapped noncoding variants altering its DNA binding site  
67 motif<sup>1</sup>, suggesting that PU.1 motif-altering variants might drive many blood cell trait association signals.

68 In order to identify blood cell trait associations that may be driven by a variant altering PU.1 binding,  
69 we analyzed publicly available PU.1 ChIP-seq data from LCLs across 49 individuals<sup>19,22</sup> and identified  
70 1497 PU.1 bQTLs. Next, PU.1 bQTLs colocalized with at least one blood cell trait association at 69 loci;  
71 for 51 of these loci, we identified PU.1 motif-altering variants as the likely causal variants. Thus, our  
72 approach allowed us to overcome the limitations of statistical fine-mapping in resolving these GWAS  
73 signals to single causal variants. By incorporating chromatin accessibility, histone mark, and transcriptome  
74 data for LCLs, we identified several putative causal genes for traits, including lymphocyte and monocyte  
75 counts. More broadly, our results illustrate the utility of TF bQTL datasets for fine-mapping trait-associated  
76 noncoding loci and in generating mechanistic, V2F models of gene dysregulation for traits of biomedical  
77 importance.

78

## 79 **Results**

80

### 81 **PU.1 motif-altering variants are likely causal for PU.1 bQTL associations**

82

83 First, we reanalyzed available PU.1 ChIP-seq data for LCLs from 49 individuals<sup>19,22</sup>. These individuals  
84 are all of European ancestry, and their genotypes are available through the 1000 Genomes Project<sup>29</sup>  
85 (Supplementary Table 1). After peak calling and normalization of the PU.1 ChIP-seq read counts, we  
86 tested for significant genetic associations with common variants (minor allele frequency (MAF) > 0.05)  
87 within 100 kb of each ChIP-Seq peak. In total, we identified 1497 significant PU.1 bQTLs (FDR < 5%).

88 We next inspected the contribution of PU.1 motif-altering variants to PU.1 bQTLs. First, we  
89 verified that PU.1-occupied regions were enriched for a match to the PU.1 binding site motif, identified  
90 by a position weight matrix (PWM), near the center of the ChIP-Seq peaks (Extended Data Fig. 1a),  
91 suggesting that most of these sites are bound directly by PU.1. Next, we evaluated whether PU.1 motif-  
92 altering variants affect PU.1 binding by training a motif score model gkm-SVM<sup>30,31</sup> to learn gapped *k*-  
93 mers that are overrepresented in PU.1-occupied sequences. The model captured both PU.1 and PU.1:IRF

94 composite motifs (Extended Data Fig. 1b), the latter of which reflects PU.1 binding to DNA as a  
95 heterodimer with either IRF4 or IRF8<sup>32</sup>. Changes in gkm-SVM scores have been shown to predict effects  
96 of variants on TF binding better than PWMs<sup>33</sup>, which imprecisely assume each nucleotide to affect  
97 binding independently. Consistent with our expectations, the predicted change in gkm-SVM scores for  
98 single nucleotide polymorphism (SNP) within PU.1 motifs were significantly correlated with estimated  
99 PU.1 bQTL effect sizes (Pearson  $r = 0.80$ ,  $p = 3.6 \times 10^{-310}$ ) (Fig. 1b, Supplementary Table 2). This strong  
100 positive correlation supports the model that PU.1 motif-altering variants, if present, are likely causal for  
101 those PU.1 bQTLs. Furthermore, significant PU.1 bQTLs with a motif-altering variant (determined by  
102 gkm-SVM) showed that such variants are more concentrated towards the peak centers compared to PU.1  
103 bQTLs without one (Fig. 1c, two-sided Fisher's exact test  $p = 3.1 \times 10^{-18}$ ), consistent with the expectation  
104 that PU.1 motif-altering variants directly affect PU.1 occupancy. Hence, we considered that PU.1 bQTLs  
105 colocalized with blood cell traits association would likely be driven by PU.1 motif-altering variants, if  
106 present (Fig. 1a).

107

#### 108 **PU.1 binding sites and PU.1 bQTLs in LCLs are enriched for blood cell trait association**

109

110 To verify the relevance of these PU.1 bQTLs for investigations of blood cell traits, we evaluated whether  
111 the PU.1 bQTLs are more likely to be significantly associated with each of the 28 blood cell traits  
112 (Supplementary Table 3) than expected by chance. We analyzed blood cell traits GWAS data from UK  
113 Biobank<sup>21</sup>. As a background expectation, we constructed 250 sets of null variants matched with PU.1  
114 bQTL lead variants for allele frequency, number of tagging variants ( $LD r^2 > 0.5$ ), and distance to the  
115 closest transcription start site (TSS). The significant PU.1 bQTLs were more likely to tag lead variants  
116 associated (*i.e.*,  $p < 5 \times 10^{-8}$ ) with myeloid lineage traits (*e.g.* monocyte and neutrophil count) and  
117 lymphoid lineage traits (*e.g.* lymphocyte count) than the sets of null variants (empirical adjusted  $p < 0.05$ )  
118 (Fig. 2a), which is consistent with the known role of PU.1 in myeloid and lymphoid differentiation<sup>26,27</sup>. In  
119 contrast, PU.1 bQTLs were not enriched for other traits like type 2 diabetes or height (Extended Data Fig.  
120 1c).

121

#### 122 **PU.1 bQTL colocalization with blood cell trait associations**

123

124 To identify candidate loci to test for potential colocalization of PU.1 bQTL and blood cell trait  
125 associations, we filtered all significant PU.1 bQTLs for loci with at least one blood cell trait association at  
126  $p < 10^{-6}$ . This resulted in a total of 1621 such PU.1 bQTL-trait pairs, comprising 367 unique loci. We then  
127 applied two distinct colocalization methods – JLIM<sup>13</sup> and Coloc<sup>12</sup> – to test for robust colocalization

128 (Supplementary Table 4). Chun and colleagues showed with simulated data that each method can show  
129 different performance depending on the LD structure of the loci<sup>13</sup>; therefore, we reasoned that requiring  
130 significant colocalization by both methods would enrich true positive cases. We used a significance  
131 threshold of  $p < 0.01172$  (FDR < 5%) for JLIM and posterior probability of colocalization  
132 (PP(Colocalization)) > 0.5 for Coloc.

133 The statistically significant colocalization of PU.1 bQTL-trait pairs identified by JLIM and Coloc  
134 were overall consistent (Pearson  $r = 0.73$ ,  $p = 6.8 \times 10^{-270}$ ; Fig. 2b). We identified a total of 190 (11.7%)  
135 PU.1-trait pairs, spanning 69 unique loci, that were significant by both methods (Fig. 3). Across the blood  
136 cell traits, those related to white blood cells (*e.g.* white blood cell count, lymphocyte count, neutrophil  
137 count) showed a higher proportion of the tested loci showing high-confidence colocalization than red  
138 blood cell or platelet traits (Fig. 3a), similar to the enrichment of tagging variants observed in Fig. 1b. We  
139 also found 1196 (73.8%) cases where a variant that was significant for both PU.1 bQTL and blood cell  
140 traits did not exhibit significant colocalization by either JLIM or Coloc, highlighting the importance of  
141 performing colocalization analysis to distinguish loci with statistical evidence of shared causal variants  
142 from those where the variants associated with each trait are merely in LD with each other<sup>17</sup>. The  
143 remaining 235 (14.5%) pairs showed discordant results between the two methods, which could potentially  
144 stem from lack of statistical power due to weak association signals or many variants showing high LD  
145 with the lead variant (Supplementary Fig. 2, Supplementary Note). This discrepancy justifies the rationale  
146 of applying both methods to identify high-confidence colocalization.

147 Most (56/69) loci showing high-confidence colocalization had some biologically plausible  
148 putative causal variants (*i.e.*, directly affecting a PU.1 binding sequence) (Fig. 2c, Extended Data Fig. 2a,  
149 Supplementary Table 5). 43 (62.3%) loci had a SNP altering a PU.1 motif, while 7 (10.1%) had a short  
150 insertion or deletion (indel) variant. In addition, there was one locus where two adjacent SNPs were in  
151 perfect LD ( $r^2 = 1$ ) and altered a single PU.1 motif sequence (Extended Data Fig. 2a and Supplementary  
152 Table 6). These SNPs and short indels showed a balance of gained and lost PU.1 binding (two-sided  
153 binomial test  $p = 0.67$ ), and changes in gkm-SVM motif scores were highly correlated with the estimated  
154 PU.1 bQTL effect sizes (Pearson  $r = 0.89$ ,  $p = 5.2 \times 10^{-18}$ ) (Extended Data Fig. 2b). The PU.1 motif-  
155 altering SNPs at colocalized loci are distributed within the PU.1 or PU.1:IRF motif, with the highest  
156 frequencies at the core “GGAAG” positions (Fig. 2d and Supplementary Table 7). We retrieved fine-  
157 mapping results for 25 colocalized loci with a PU.1 motif-altering variant (*i.e.* SNP or indel) from a recent  
158 blood cell trait GWAS study<sup>8</sup> (Supplementary Note). 19 of these 25 (76%) loci had more than 10 variants  
159 in the 95% credible set (*i.e.*, minimal set of variants that have 95% posterior probability of containing the  
160 causal variant), none of which was fine-mapped to a single variant (Fig. 2e and Supplementary Table 8).  
161 Despite difficulty in fine-mapping due to LD structure, we were able to pinpoint putative causal variants

162 in these loci using a specific TF's (*i.e.*, PU.1) motif information. There were also 5 loci with large  
163 deletions that completely removed the PU.1 binding site, which we were able to uncover because the  
164 1000 Genomes Project (1KGP)<sup>29</sup> genotypes included structural variants (Extended Data Fig. 2c); whether  
165 the deletions are true causal variants will need to be tested experimentally in future studies.

166 Pinpointing likely causal regulatory variants allowed us to derive specific hypotheses about gene  
167 regulatory mechanisms that are perturbed by the variants, as described below. We show one example  
168 where a PU.1 motif-altering SNP (rs12517864) represents a secondary expression QTL (eQTL) (*i.e.*, a  
169 weaker signal independent from the strongest, primary eQTL) to *ZNF608* in LCLs, and only this  
170 secondary signal colocalizes with lymphocyte count association (Fig. 4); an eQTL-centric analysis in  
171 LCLs would have missed this locus without accounting for multiple independent signals, highlighting the  
172 power of the use of TF bQTL data in colocalization analysis with GWAS data. Two other examples show  
173 reporter assay results corroborating the regulatory effects of PU.1 motif-altering variants identified in  
174 colocalized loci (Fig. 5 and 6).

175

#### 176 **bQTL colocalization reveals a putative causal variant that is not the primary eQTL**

177

178 Causal genes at a trait-associated locus frequently have been identified using eQTL data for nearby  
179 genes<sup>14,34</sup>. However, eQTLs can often have multiple independent signals<sup>14</sup>, and these signals detected in  
180 any one cell type may not all be associated with a GWAS trait, such as if the regulatory effects manifest  
181 themselves only in certain cellular contexts. This complicates colocalization analyses that often assume a  
182 single shared causal variant at a locus<sup>12,13</sup>. In contrast, TF bQTLs capture regulatory effects of individual  
183 regulatory elements. Therefore, TF bQTL colocalization analysis can isolate the effects of variants on  
184 specific regulatory elements, lowering the probability of multiple causal variants compared to that of  
185 eQTLs.

186 For example, the *ZNF608* locus shows significant colocalization of PU.1 bQTL and lymphocyte  
187 count association (Fig. 4a, Extended Data Fig. 3a). Although the molecular function of *ZNF608* remains  
188 unclear, a study of follicular lymphoma (FL), a type of cancer in which B lymphocytes divide  
189 uncontrollably, found this gene to be among the 39 genes significantly enriched for missense or predicted-  
190 loss-of-function (pLOF) somatic mutations in FL patients<sup>35</sup>, suggesting it plays a role in B lymphocyte  
191 development. The associated PU.1 binding site is located about 257 kilobases (kb) upstream of the  
192 *ZNF608* promoter, and the SNP rs12517864 that increases the PU.1 binding motif score (0.68→2.69) is  
193 located near the center of the PU.1 occupancy site (Fig. 4b,g).

194 Multiple lines of evidence support the regulatory effect of rs12517864. We reanalyzed ATAC-seq  
195 and histone mark ChIP-seq data for LCLs<sup>36,37</sup> and found that rs12517864 is significantly associated with

196 each of these molecular phenotypes that overlap the PU.1 binding site, suggesting that the variant, if  
197 causal, affects gene regulation (Fig. 4f, h). Furthermore, the variant falls within a fragment that physically  
198 interacts only with the *ZNF608* promoter in primary B cells according to promoter-capture Hi-C (PCHi-  
199 C) data<sup>38</sup>, supporting the model that rs12517864 directly regulates *ZNF608* (Fig. 4g).

200 Surprisingly, initial inspection of *ZNF608* eQTL signals in LCLs<sup>39</sup> seemed contradictory because  
201 the lead variant for this eQTL (rs2028854) is located elsewhere, 200 kb upstream of the *ZNF608*  
202 promoter, and is not strongly associated with lymphocyte count<sup>21</sup> ( $p = 0.04$ ) (Fig. 4c, g). We therefore  
203 examined the possibility of multiple independent *ZNF608* eQTL signals in LCLs by performing  
204 conditional analysis on the lead variant, as well as fine-mapping using SuSiE<sup>40</sup>, which can detect multiple  
205 signals. Once conditioned on the lead eQTL SNP rs2028854, association of rs12517864 to *ZNF608*  
206 expression became much stronger ( $p = 6.98 \times 10^{-7}$ ) (Fig. 4c). Moreover, the fine-mapping analysis  
207 identified two independent credible sets for *ZNF608* eQTL signal, one of which contained rs12517864 as  
208 the variant with the highest posterior inclusion probability (PIP = 0.07), demonstrating that this variant is  
209 likely to be causally associated with *ZNF608* expression level (Fig. 4d).

210 Since only one of the two independent *ZNF608* eQTL signals in LCLs is associated with  
211 lymphocyte count, we hypothesized that even though both SNPs are significant eQTLs in LCLs, only  
212 rs12517864 (*i.e.*, the secondary eQTL signal), and not rs2028854 (*i.e.*, the primary eQTL signal),  
213 modulates *ZNF608* expression in the causal cell type. Analysis of RNA-seq data for various blood cells<sup>28</sup>  
214 revealed that *ZNF608* is highly expressed in common lymphoid progenitors and B cells (Fig. 4i).  
215 Inspection of eQTL data for B cells in the eQTL Catalogue<sup>41,42</sup> showed that only rs12517864, and not  
216 rs2028854 ( $p = 0.25$ ), is significantly associated with *ZNF608* expression ( $p = 4.39 \times 10^{-5}$ ) (Fig. 4e).  
217 Although we cannot unambiguously conclude that B cells are the causal cell type, rs12517864 is likely  
218 the only variant increased lymphocyte count through increased *ZNF608* expression (Fig. 4h and Extended  
219 Data Fig. 3a).

220

## 221 **Blood cell trait-associated PU.1 motif-altering variants show regulatory effects in reporter assays**

222

223 To verify that the nominated PU.1 motif-altering variants are indeed regulatory variants, we inspected  
224 massively parallel reporter assay (MPRA) studies data<sup>43,44</sup>, which measured the regulatory effects of two  
225 such variants: (1) rs5827412, a PU.1 motif-altering short deletion associated with monocyte percentage  
226 affects expression levels of *LRRC25*, a gene previously shown to be necessary for granulocyte  
227 differentiation<sup>45</sup>, in monocytes (Fig. 5); and (2) rs3808619, a PU.1 motif-altering SNP at the promoter of  
228 *ZC2HC1A*, a functionally uncharacterized gene, as the regulatory causal variant for association with  
229 lower lymphocyte count (Fig. 6).

230 *LRRC25*, also called monocyte and plasmacytoid-activated protein (MAPA), is a gene shown to  
231 impair differentiation of granulocytes, which share lineages with monocytes, if knocked down or knocked  
232 out<sup>45</sup>. At this locus, we found that the PU.1 bQTL signal showed significant colocalization with monocyte  
233 count and percentage, neutrophil count and percentage, and white blood cell count association signals<sup>8,21</sup>  
234 (Fig. 5b and Extended Data Fig. 4a). The corresponding PU.1 binding site contains a short deletion  
235 rs5827412 that lowers the PU.1 motif score and is associated with reduced PU.1 binding, as well as  
236 chromatin accessibility, active histone mark levels, and *LRRC25* expression<sup>36,37,39</sup> (Fig. 5a and Extended  
237 Data Fig. 4b). This deletion significantly reduced regulatory activity in a reporter assay<sup>44</sup> (two-sided *t*-test  
238  $p = 6.9 \times 10^{-5}$ ) (Fig. 5f); data from another study suggested concordant direction of effect despite not being  
239 statistically significant<sup>43</sup> (negative binomial regression  $p = 0.26$ ) (Fig. 5c). Next, we analyzed available  
240 ATAC-seq data from *SP11*, the gene encoding PU.1, knockout pro-B cell lines (RS4;11) to verify whether  
241 PU.1 is likely to be the trans factor for the regulatory variant<sup>46</sup>, and determined that *SP11* knockout  
242 resulted in significantly reduced chromatin accessibility at sites of PU.1 occupancy genome-wide<sup>47</sup> (chi  
243 square test  $p < 1 \times 10^{-300}$ ) (Supplementary Fig. 3). Indeed, the activity of the regulatory element that  
244 contains rs5827412 is likely dependent on PU.1 binding as *SP11* knockout cell lines showed reduced  
245 chromatin accessibility at this region (DESeq2 adjusted  $p = 8.73 \times 10^{-5}$ ) (Fig. 5d). RNA-seq data for 13  
246 blood cell types<sup>28</sup> indicates that *LRRC25* is specifically expressed in monocytes at a much higher level  
247 than in other blood cell types and is sharply upregulated as progenitor cells differentiate to monocytes  
248 (Fig. 5e and Extended Data Fig. 4c). Consistent with the variant's strongest effect on monocyte  
249 percentage ( $p = 1.3 \times 10^{-96}$ ) and monocyte-specific expression of *LRRC25*, we found that rs5827412 is also  
250 significantly associated with reduced *LRRC25* expression in monocytes<sup>16</sup> ( $p = 3.78 \times 10^{-22}$ ) (Fig. 5f) and is  
251 in a regulatory element that is accessible throughout monocyte differentiation (Fig. 5g). Altogether, our  
252 results provide strong support for rs5827412 reducing *LRRC25* gene expression levels in monocytes and  
253 decreasing monocyte percentage while increasing neutrophil percentage.

254 The *ZC2HC1A* locus, which is primarily associated with lymphocyte count and percentage<sup>21</sup> (Fig.  
255 6b and Extended Data Fig. 5a,b), represents a challenging locus for fine-mapping. Here, 44 variants  
256 comprise the 95% credible set (*i.e.*, a minimal set of putative causal variants), based on a UK Biobank  
257 fine-mapping study<sup>48</sup> (Fig. 6c). Among the candidate causal variants at the *ZC2HC1A* locus, rs3808619 is  
258 the only PU.1 motif-altering variant found within the associated PU.1 binding site at the *ZC2HC1A*  
259 promoter; rs3808619 increases the strength of a PU.1 motif, resulting in a higher affinity DNA binding  
260 site (Fig. 6a). Of multiple tagging variants in this locus that were tested for reporter activity (59 variants  
261 in Abell et al.<sup>43</sup> and 30 variants in Tewhey et al.<sup>44</sup>), only rs3808619 showed a significantly increased  
262 reporter activity that is concordant in direction with that of the variant's associations to elevated  
263 chromatin accessibility, active histone mark levels, and *ZC2HC1A* expression in LCLs<sup>36,37,39</sup> (Fig. 6d,e,f).



264 Finally, as for rs5827412, we detected significantly reduced chromatin accessibility levels at the  
265 *ZC2HC1A* promoter in *SP11* knockout cell lines<sup>46</sup> (DESeq2 adjusted  $p = 1.76 \times 10^{-13}$ ), supporting the likely  
266 role of PU.1 at this promoter (Fig. 6g). rs3808619 is also associated with multiple sclerosis<sup>49</sup> ( $p = 1.1 \times 10^{-9}$ )  
267 (Extended Data Fig. 5c,d), suggesting it plays a multifactorial role in IMDs. Our results suggest that a  
268 direct consequence of rs3808619, which is associated with lower lymphocyte count, is likely *ZC2HC1A*  
269 upregulation (Supplementary Note).

270

## 271 Discussion

272

273 Our results with PU.1 binding and blood cell trait GWAS data demonstrate the utility of TF bQTL data in  
274 identifying which of many variants in LD are the likely causal regulatory variants underlying GWAS trait  
275 associations, as the presence of motif-altering variants suggests that they directly affect binding of the  
276 corresponding TF (Fig. 2c). Incorporating PU.1 bQTLs in our colocalization analysis conferred two key  
277 advantages: 1) identification of trait-associated regulatory elements and 2) identification of putatively  
278 causal PU.1 motif-altering variants. Together, they highlight a likely transcriptional regulatory  
279 mechanism underlying the trait association. In contrast, eQTL colocalization cannot assist fine-mapping  
280 in this way because there is no prior expectation that a specific noncoding region regulates the associated  
281 gene and that a regulatory variant would alter a certain TF binding site.

282 For instance, at the *ZNF608* locus, pinpointing the putative causal variant and associated  
283 regulatory element would have been difficult without PU.1 bQTLs, especially because there is another  
284 stronger eQTL signal, which did not colocalize with the lymphocyte count association for *ZNF608* in  
285 LCLs (Fig. 4). Such a situation may partially explain the observation that many significant eQTL signals  
286 failed to colocalize with the GWAS associations using existing colocalization methods<sup>13</sup>; however, this  
287 locus was the only such example in our study. Nevertheless, this example motivates applying TF bQTL  
288 colocalization to isolate independent eQTL signals, generating eQTL data in trait-relevant cell types<sup>50</sup>,  
289 and applying colocalization methods that allow multiple causal variants to eQTLs<sup>51</sup>, if accurate LD  
290 matrices or individual genotypes are available for both traits, which is often not the case for GWAS data.

291 A prior study that performed colocalization analysis of PU.1 bQTLs in neutrophils and immune  
292 diseases GWAS found that the majority (>50%) of colocalized variants altered the binding site motifs of  
293 other TFs<sup>18</sup>; in contrast, we found that the majority (87%) of the colocalized blood cell trait GWAS loci  
294 had a variant that altered a PU.1 motif (Fig. 2c), even though only a minority (34%) of all PU.1 bQTLs,  
295 colocalized or not, did overall (Fig. 1c). The increased proportion of PU.1 motif-altering variants present  
296 in this study may be due to PU.1's central role in blood cell traits<sup>26</sup> and highlights the increased likelihood  
297 that PU.1 binding is mediating the genetic effects on blood cell traits.

298 We observed that only a minority of the tested GWAS loci (69 / 367) showed significant  
299 colocalization. This is not surprising because we selected candidate loci solely based on the marginal  
300 association to PU.1 binding and blood cell traits<sup>13</sup>, without filtering for high LD between the two lead  
301 variants<sup>13</sup> to ‘cast a wide net’ for discovery. This observation is a testament to the importance of  
302 performing colocalization analysis to distinguish loci with a single causal variant for the two phenotypes  
303 (here, PU.1 binding and a particular blood cell trait) from those with distinct variants responsible for the  
304 different phenotypes. Furthermore, even though PU.1 bQTLs were enriched for blood cell traits  
305 association (Fig. 2a), they explain only a subset of all associated loci, likely indicating that other TFs are  
306 mediating genetic effects at other associated loci.

307 We offer guidelines for broad application of colocalization analysis with TF bQTLs. First, high-  
308 quality ChIP-grade antibodies<sup>52</sup> or, alternatively, cell lines in which the TF has been epitope-tagged, are  
309 essential. Second, TFs for bQTL analysis, as well as the cell type for the ChIP experiments, must be  
310 selected to be relevant to the trait or disease of interest. The feasibility of our analysis relied on the  
311 relevance of PU.1, a known hematopoietic master regulator, and LCLs, a model of mature B cells, to  
312 specific blood cell traits, such as lymphocyte count and monocyte count. Future studies will need to  
313 validate the regulatory functions of the variants in the relevant primary cell types. Third, sufficiently large  
314 sample sizes for both GWAS and TF bQTL are necessary for discovery, as colocalization can return false  
315 negative results due to limited statistical power<sup>53</sup>; although the sample size of 49 for the PU.1 bQTL data  
316 led to 69 robustly colocalized loci, we anticipate that a larger sample size could increase the power to  
317 detect weaker colocalization signals.

318 Future studies could use TF bQTL data in colocalization analysis to elucidate the ever-increasing  
319 number of trait-associated loci<sup>1</sup>. Where TFs important for a trait are known, TF bQTLs identified in the  
320 relevant cell type(s) could mediate a subset of trait associations, shedding light on putative causal  
321 variants, as well as the pathogenic mechanisms. Such colocalization analysis with TF bQTL data uniquely  
322 provides a path to pinpointing causal regulatory elements and variants, and thus a smaller set of  
323 mechanistic hypotheses to test experimentally to verify the underlying causes of the disease.

324 **References**

- 325 1. Claussnitzer, M. *et al.* A brief history of human disease genetics. *Nature* **577**, 179–189 (2020).
- 326 2. Claussnitzer, M., Dankel, S. N., Kim, K.-H., Hauner, H. & Kellis, M. FTO obesity variant  
327 circuitry and adipocyte browning in humans. *New England Journal of Medicine* vol. 6 895–907  
328 (2015).
- 329 3. Nasser, J. *et al.* Genome-wide enhancer maps link risk variants to disease genes. *Nature* **593**, 238–  
330 243 (2021).
- 331 4. International Common Disease Alliance. *International Common Disease Alliance White Paper*  
332 *v1.0*. <https://www.icda.bio/> (2020).
- 333 5. Visscher, P. M. *et al.* 10 Years of GWAS Discovery: Biology, Function, and Translation. *Am. J.*  
334 *Hum. Genet.* **101**, 5–22 (2017).
- 335 6. Amariuta, T. *et al.* Improving the trans-ancestry portability of polygenic risk scores by prioritizing  
336 variants in predicted cell-type-specific regulatory elements. *Nat. Genet.* **52**, 1346–1354 (2020).
- 337 7. Weissbrod, O. *et al.* Leveraging fine-mapping and multipopulation training data to improve cross-  
338 population polygenic risk scores. *Nat. Genet.* **54**, 450–458 (2022).
- 339 8. Vuckovic, D. *et al.* The Polygenic and Monogenic Basis of Blood Traits and Diseases. *Cell* **182**,  
340 1214–1231 (2020).
- 341 9. Amariuta, T. *et al.* IMPACT: Genomic Annotation of Cell-State-Specific Regulatory Elements  
342 Inferred from the Epigenome of Bound Transcription Factors. *Am. J. Hum. Genet.* **104**, 879–895  
343 (2019).
- 344 10. van de Geijn, B. *et al.* Annotations capturing cell type-specific TF binding explain a large fraction  
345 of disease heritability. *Hum. Mol. Genet.* **29**, 1057–1067 (2020).
- 346 11. Ulirsch, J. C. *et al.* Interrogation of human hematopoiesis at single-cell and single-variant  
347 resolution. *Nat. Genet.* **51**, 683–693 (2019).
- 348 12. Giambartolomei, C. *et al.* Bayesian test for colocalisation between pairs of genetic association  
349 studies using summary statistics. *PLoS Genet.* **10**, e1004383 (2014).
- 350 13. Chun, S. *et al.* Limited statistical evidence for shared genetic effects of eQTLs and autoimmune-  
351 disease-associated loci in three major immune-cell types. *Nat. Genet.* **49**, 600–605 (2017).
- 352 14. GTEx Consortium. The GTEx Consortium atlas of genetic regulatory effects across human tissues.  
353 *Science* **369**, 1318–1330 (2020).
- 354 15. Barbeira, A. N. *et al.* Exploiting the GTEx resources to decipher the mechanisms at GWAS loci.  
355 *Genome Biol.* **22**, 49 (2021).
- 356 16. Chen, L. *et al.* Genetic Drivers of Epigenetic and Transcriptional Variation in Human Immune

- 357 Cells. *Cell* **167**, 1398–1414.e24 (2016).
- 358 17. Liu, B., Gloudemans, M. J., Rao, A. S., Ingelsson, E. & Montgomery, S. B. Abundant associations  
359 with gene expression complicate GWAS follow-up. *Nat. Genet.* **51**, 768–769 (2019).
- 360 18. Watt, S. *et al.* Genetic perturbation of PU.1 binding and chromatin looping at neutrophil enhancers  
361 associates with autoimmune disease. *Nat. Commun.* **12**, 1–12 (2021).
- 362 19. Kilpinen, H. *et al.* Coordinated effects of sequence variation on DNA binding, chromatin  
363 structure, and transcription. *Science* **342**, 744–747 (2013).
- 364 20. Deplancke, B., Alpern, D. & Gardeux, V. The Genetics of Transcription Factor DNA Binding  
365 Variation. *Cell* **166**, 538–554 (2016).
- 366 21. Canela-Xandri, O., Rawlik, K. & Tenesa, A. An atlas of genetic associations in UK Biobank. *Nat.*  
367 *Genet.* **50**, 1593–1599 (2018).
- 368 22. Waszak, S. M. *et al.* Population Variation and Genetic Control of Modular Chromatin Architecture  
369 in Humans. *Cell* **162**, 1039–1050 (2015).
- 370 23. Guan, W.-J. *et al.* Clinical Characteristics of Coronavirus Disease 2019 in China. *N. Engl. J. Med.*  
371 **382**, 1708–1720 (2020).
- 372 24. Terpos, E. *et al.* Hematological findings and complications of COVID-19. *Am. J. Hematol.* **95**,  
373 834–847 (2020).
- 374 25. Wang, S., Sheng, Y., Tu, J. & Zhang, L. Association between peripheral lymphocyte count and the  
375 mortality risk of COVID-19 inpatients. *BMC Pulm. Med.* **21**, 55 (2021).
- 376 26. Fisher, R. C. & Scott, E. W. Role of PU.1 in hematopoiesis. *Stem Cells* **16**, 25–37 (1998).
- 377 27. Rothenberg, E. V., Hosokawa, H. & Ungerback, J. Mechanisms of Action of Hematopoietic  
378 Transcription Factor PU.1 in Initiation of T-Cell Development. *Front. Immunol.* **10**, 228 (2019).
- 379 28. Corces, M. R. *et al.* Lineage-specific and single-cell chromatin accessibility charts human  
380 hematopoiesis and leukemia evolution. *Nat. Genet.* **48**, 1193–1203 (2016).
- 381 29. Auton, A. *et al.* A global reference for human genetic variation. *Nature* **526**, 68–74 (2015).
- 382 30. Ghandi, M., Lee, D., Mohammad-Noori, M. & Beer, M. A. Enhanced Regulatory Sequence  
383 Prediction Using Gapped k-mer Features. *PLoS Comput. Biol.* **10**, (2014).
- 384 31. Lee, D. *et al.* A method to predict the impact of regulatory variants from DNA sequence. *Nat.*  
385 *Genet.* **47**, 955–961 (2015).
- 386 32. Escalante, C. R. *et al.* Crystal structure of PU.1/IRF-4/DNA ternary complex. *Mol. Cell* **10**, 1097–  
387 1105 (2002).
- 388 33. Yan, J. *et al.* Systematic analysis of binding of transcription factors to noncoding variants. *Nature*  
389 **591**, 147–151 (2021).
- 390 34. Hormozdiari, F. *et al.* Colocalization of GWAS and eQTL Signals Detects Target Genes. *Am. J.*

- 391 *Hum. Genet.* **99**, 1245–1260 (2016).
- 392 35. Krysiak, K. *et al.* Recurrent somatic mutations affecting B-cell receptor signaling pathway genes  
393 in follicular lymphoma. *Blood* **129**, 473–483 (2017).
- 394 36. Kumasaka, N., Knights, A. J. & Gaffney, D. J. High-resolution genetic mapping of putative causal  
395 interactions between regions of open chromatin. *Nat. Genet.* **51**, 128–137 (2019).
- 396 37. Delaneau, O. *et al.* Chromatin three-dimensional interactions mediate genetic effects on gene  
397 expression. *Science* **364**, (2019).
- 398 38. Javierre, B. M. *et al.* Lineage-Specific Genome Architecture Links Enhancers and Non-coding  
399 Disease Variants to Target Gene Promoters. *Cell* **167**, 1369–1384 (2016).
- 400 39. Lappalainen, T. *et al.* Transcriptome and genome sequencing uncovers functional variation in  
401 humans. *Nature* **501**, 506–511 (2013).
- 402 40. Wang, G., Sarkar, A., Carbonetto, P. & Stephens, M. A simple new approach to variable selection  
403 in regression, with application to genetic fine mapping. *J. R. Stat. Soc. Ser. B Stat. Methodol.* **82**,  
404 1273–1300 (2020).
- 405 41. Kerimov, N. *et al.* A compendium of uniformly processed human gene expression and splicing  
406 quantitative trait loci. *Nat. Genet.* **53**, 1290–1299 (2021).
- 407 42. Schmiedel, B. J. *et al.* Impact of Genetic Polymorphisms on Human Immune Cell Gene  
408 Expression. *Cell* **175**, 1701–1715 (2018).
- 409 43. Abell, N. S. *et al.* Multiple causal variants underlie genetic associations in humans. *Science* **375**,  
410 1247–1254 (2022).
- 411 44. Tewhey, R. *et al.* Direct identification of hundreds of expression-modulating variants using a  
412 multiplexed reporter assay. *Cell* **165**, 1519–1529 (2016).
- 413 45. Liu, W. *et al.* LRRC25 plays a key role in all-trans retinoic acid-induced granulocytic  
414 differentiation as a novel potential leukocyte differentiation antigen. *Protein Cell* **9**, 785–798  
415 (2018).
- 416 46. Coz, C. Le *et al.* Constrained chromatin accessibility in PU.1-mutated agammaglobulinemia  
417 patients. *J. Exp. Med.* **218**, (2021).
- 418 47. Wu, J. N. *et al.* Functionally distinct patterns of nucleosome remodeling at enhancers in  
419 glucocorticoid-treated acute lymphoblastic leukemia. *Epigenetics Chromatin* **8**, 53 (2015).
- 420 48. Kanai, M. *et al.* Insights from complex trait fine-mapping across diverse populations. Preprint at  
421 <https://www.medrxiv.org/content/10.1101/2021.09.03.21262975v1> (2021).
- 422 49. International Multiple Sclerosis Genetics Consortium (IMSGC) *et al.* Analysis of immune-related  
423 loci identifies 48 new susceptibility variants for multiple sclerosis. *Nat. Genet.* **45**, 1353–60  
424 (2013).

- 425 50. Umans, B. D., Battle, A. & Gilad, Y. Where Are the Disease-Associated eQTLs? *Trends Genet.*  
426 **37**, 109–124 (2021).
- 427 51. Wallace, C. A more accurate method for colocalisation analysis allowing for multiple causal  
428 variants. *PLoS Genet.* **17**, 1–11 (2021).
- 429 52. Baker, M. Reproducibility crisis: Blame it on the antibodies. *Nature* **521**, 274–6 (2015).
- 430 53. Hukku, A. *et al.* Probabilistic Colocalization of Genetic Variants from Complex and Molecular  
431 Traits: Promise and Limitations. *Am. J. Hum. Genet.* **108**, 25–35 (2020).
- 432 54. Langmead, B. & Salzberg, S. L. Fast gapped-read alignment with Bowtie 2. *Nat. Methods* **9**, 357–  
433 359 (2012).
- 434 55. Van De Geijn, B., Mcvicker, G., Gilad, Y. & Pritchard, J. K. WASP: Allele-specific software for  
435 robust molecular quantitative trait locus discovery. *Nat. Methods* **12**, 1061–1063 (2015).
- 436 56. Wu, T. D. & Nacu, S. Fast and SNP-tolerant detection of complex variants and splicing in short  
437 reads. *Bioinformatics* **26**, 873–881 (2010).
- 438 57. Zhang, Y. *et al.* Model-based analysis of ChIP-Seq (MACS). *Genome Biol.* **9**, R137 (2008).
- 439 58. Liao, Y., Smyth, G. K. & Shi, W. FeatureCounts: An efficient general purpose program for  
440 assigning sequence reads to genomic features. *Bioinformatics* **30**, 923–930 (2014).
- 441 59. Robinson, M. D. & Oshlack, A. A scaling normalization method for differential expression  
442 analysis of RNA-seq data. *Genome Biol.* (2010).
- 443 60. Stegle, O., Parts, L., Piipari, M., Winn, J. & Durbin, R. Using probabilistic estimation of  
444 expression residuals (PEER) to obtain increased power and interpretability of gene expression  
445 analyses. *Nat. Protoc.* **7**, 500–507 (2012).
- 446 61. Das, S. *et al.* Next-generation genotype imputation service and methods. *Nat. Genet.* **48**, 1284–  
447 1287 (2016).
- 448 62. Delaneau, O. *et al.* A complete tool set for molecular QTL discovery and analysis. *Nat. Commun.*  
449 **8**, 15452 (2017).
- 450 63. McCarthy, S. *et al.* A reference panel of 64,976 haplotypes for genotype imputation. *Nat. Genet.*  
451 **48**, 1279–1283 (2016).
- 452 64. Pers, T. H., Timshel, P. & Hirschhorn, J. N. SNPsnip: a Web-based tool for identification and  
453 annotation of matched SNPs. *Bioinformatics* **31**, 418–20 (2015).
- 454 65. Mahajan, A. *et al.* Fine-mapping type 2 diabetes loci to single-variant resolution using high-  
455 density imputation and islet-specific epigenome maps. *Nat. Genet.* **50**, 1505–1513 (2018).
- 456 66. Wood, A. R. *et al.* Defining the role of common variation in the genomic and biological  
457 architecture of adult human height. *Nat. Genet.* **46**, 1173–1186 (2014).
- 458 67. Ambrosini, G., Groux, R. & Bucher, P. PWMScan: a fast tool for scanning entire genomes with a

- 459 position-specific weight matrix. *Bioinformatics* **34**, 2483–2484 (2018).
- 460 68. Quinlan, A. R. & Hall, I. M. BEDTools: a flexible suite of utilities for comparing genomic  
461 features. *Bioinformatics* **26**, 841–2 (2010).
- 462 69. Love, M. I., Huber, W. & Anders, S. Moderated estimation of fold change and dispersion for  
463 RNA-seq data with DESeq2. *Genome Biol.* **15**, 1–21 (2014).
- 464

## 465 **Methods**

### 466 **PU.1 ChIP-seq data processing**

467 We downloaded PU.1 ChIP-seq fastq files from EMBL-EBI ArrayExpress under accession E-MTAB-  
468 3657<sup>22</sup> ( $n=45$ ) and E-MTAB-1884<sup>19</sup> ( $n=4$ ). The list of samples is provided in Supplementary Table 1. We  
469 mapped the reads to the hg19 reference genome supplemented with the Epstein-Barr virus (EBV) using  
470 Bowtie 2<sup>54</sup>. In order to eliminate reference allele bias in read mapping, we applied WASP<sup>55</sup> to filter reads  
471 that mapped to a different position when variants were added, and used GSNAP<sup>56</sup>, which is a SNP-  
472 tolerant read alignment method, to remap filtered out reads.

473 PU.1 ChIP-seq peaks were called using MACS2<sup>57</sup>. For equal representation, we subsampled 5  
474 million reads from each sample and performed peak calling on the aggregate alignment file. To account  
475 for the size of the merged read set, we downloaded 8 available control ChIP-seq samples in GM12878  
476 from ENCODE (File ID: ENCFF032WUR, ENCFF426WJH, ENCFF508HCX,  
477 ENCFF537DAJ, ENCFF812HUT, ENCFF837IOW, ENCFF849LYY, ENCFF892TNJ). To define 200 bp  
478 sequences occupied by PU.1, we took the summits and extended them 100 bp in each direction. In total,  
479 there were 78720 peaks.

480

### 481 **PU.1 binding quantitative trait loci**

482 First, we quantified the PU.1 binding levels at identified occupancy sites. We counted the number of  
483 reads overlapping each 200 bp peak using featureCounts<sup>58</sup>. The read counts were normalized for library  
484 size using trimmed mean of M-values<sup>59</sup> and further normalized to follow a standard normal distribution  
485 across the samples, using quantile normalization. Finally, in order to eliminate the effect of variables,  
486 such as batch, gender, and ancestry, we used PEER<sup>60</sup> to residualize the phenotype values, correcting for  
487 batch (*i.e.*, which publication), sex, and 3 genotype principal components, as well as 10 PEER factors.

488 Second, we obtained the genotypes of the LCL samples from the 1000 Genomes Project data<sup>29</sup>. 4  
489 out of 49 samples only had microarray genotype data from Illumina Omni2.5 chips, and these genotypes  
490 were phased and imputed using the European samples of the 1000 Genomes project phase 3 data<sup>29</sup> on the  
491 Michigan Imputation Server<sup>61</sup>. Genotypes of all samples were converted to biallelic form and aggregated.  
492 Afterwards, variants with minor allele frequency less than 5% were removed from the PU.1 binding  
493 quantitative trait loci analysis.

494 Finally, we tested for genetic associations to PU.1 binding levels using the phenotype matrix and  
495 the genotype data. We utilized QTLtools<sup>62</sup> to approximate linear regression efficiently while also  
496 correcting for multiple hypotheses tested with permutations and false discovery rate estimation. For each



497 PU.1 occupancy site, variants within 100 kb were included in the QTL analysis. In the end, there were  
498 1497 significant PU.1 bQTLs.

499

### 500 **UK Biobank blood cell trait GWAS summary statistics**

501 We downloaded 28 blood cell trait GWAS summary statistics from UK Biobank<sup>21</sup> for the colocalization  
502 analysis. The authors performed a linear mixed model-based regression analysis on 452,264 White British  
503 individuals using rank-normalized phenotypes. The 28 blood cell traits are listed in Supplementary Table  
504 3. One limitation of these summary statistics is that the authors used the Haplotype Reference Consortium  
505 imputation panel, which only included SNPs by design, for imputation<sup>63</sup> (Supplementary Note). Thus,  
506 short deletions like rs5827412 were missing in these summary statistics. For Figure 5, we verified that the  
507 variant is associated with decreased monocyte percentage and increased neutrophil percentage in  
508 summary statistics from another analysis of the UK Biobank data<sup>8</sup>, and utilized these data for  
509 visualization.

510

### 511 **Fold Enrichment of GWAS signal in PU.1 bQTLs**

512 We first generated 250 sets of null variants matched with the significant PU.1 bQTL lead variants for  
513 allele frequency, number of tagging SNPs ( $LD\ r^2 > 0.5$ ), and distance to the closest transcription start site  
514 (TSS), using SNPsnap<sup>64</sup>. 250 sets of null variants were successfully generated for 1292 of the PU.1 bQTL  
515 lead variants, so we restricted the downstream analysis within them. Using the distribution of number of  
516 variants tagging ( $r^2 > 0.8$ ) trait-associated lead variants as the background, we computed the fold  
517 enrichment of the number of PU.1 bQTLs tagging those variants. The empirical  $p$  values are derived for  
518 each blood cell trait by counting how many sets had SNPs tagging ( $r^2 > 0.8$ ) trait-associated variants more  
519 than or equal to the number of PU.1 bQTLs tagging them and dividing by 251. The  $p$  values were  
520 adjusted using *qvalue* package in R. For non-blood traits, lead SNPs from GWAS of type 2 diabetes<sup>65</sup> and  
521 height<sup>66</sup> were used.

522

### 523 **Position weight matrix and gkm-SVM PU.1 motif models**

524 To initially scan for the position of PU.1 motif sequences within occupancy sites, we used PWMScan<sup>67</sup>.  
525 With a PU.1 (SPI1) motif position weight matrix (PWM) selected within the tool (CISBP: M6119\_1) we  
526 scanned for the motif ( $p < 10^{-5}$ ) within PU.1 occupancy sites, which resulted in a total of 30812 instances.  
527 To determine the relative location of PU.1 motifs within the PU.1 occupancy sites, we subtracted the start  
528 or end position of the motif from the center position of the 200 bp PU.1 peak, depending on the strand  
529 (Extended Data Fig. 1a).

530           Afterwards, we trained a PU.1 motif model using gkm-SVM<sup>31</sup>, as a more sophisticated  
531 counterpart to PWM. We used the 200 bp sequences detected to be PU.1 occupancy sites for positive  
532 sequences in the training set. We left out PU.1 occupancy sites with a variant overlapping PU.1 motifs  
533 identified using PWMs (*i.e.*, one of the alleles with log-likelihood score > 8) from the training set so that  
534 the model effectively captures the motif sequences and excludes potentially causal PU.1 bQTLs. We  
535 generated negative sequences using the ‘genNullSeqs’ function in the gkmSVM R package. Then, we  
536 trained the model using default parameters with LS-GKM<sup>31</sup>, which is a faster implementation from the  
537 developers. Throughout the study, we defined PU.1 motif-altering variants as those where one of the  
538 alleles shows a gkm-SVM score greater than 0 for a 30 bp sequence centered at the variant, and the  
539 variant induces a non-zero change.

540

#### 541 **Colocalization analysis using JLIM and Coloc**

542 We selected 1621 PU.1-trait pairs at loci where the significant PU.1 bQTLs also show at least one blood  
543 cell trait association at  $p < 10^{-6}$  to perform colocalization. For JLIM<sup>13</sup>, we used the default parameters.  $p$   
544 values were derived by permuting the PU.1 binding level matrix. For Coloc<sup>12</sup>, we used the prior  
545 parameters  $p_1=10^{-4}$ ,  $p_2=10^{-4}$ , and  $p_{12}=10^{-6}$ , which is more conservative than the default, and ran Coloc on  
546 the summary statistics. For both analyses, we considered variants within a 200 kb window around the  
547 GWAS lead variant. We used a significance threshold of  $p < 0.01172$  (FDR < 5%) for JLIM and posterior  
548 probability of colocalization (PP(Colocalization)) > 0.5. The FDR cutoff for JLIM was determined by the  
549 equation:

$$550 \quad FDR(p_{cutoff}) = \frac{p_{cutoff} N}{\#\{P_{JLIM} \leq p_{cutoff}\}}$$

551 where  $p_{cutoff}$  is the  $p$  value cutoff,  $N$  is the number of PU.1-trait loci tested, and  $P_{JLIM}$  is the JLIM  $p$  value.

552

#### 553 **Chromatin accessibility, histone mark, and expression QTLs in LCLs**

554 ATAC-seq<sup>36</sup> ( $n=100$ ), histone mark ChIP-seq ( $n=158$ <sup>13</sup> and  $n=2$ <sup>34</sup>, respectively), and RNA-seq<sup>39</sup> ( $n=373$ )  
555 data were downloaded from European Nucleotide Archive (ERP110508), EMBL-EBI ArrayExpress (E-  
556 MTAB-3657 and E-GEUV-1), respectively. ATAC-seq data were only available as bam files, so we used  
557 bamtofastq command from bedtools<sup>68</sup> to extract reads. We processed ATAC-seq and histone mark ChIP-  
558 seq read data similarly to PU.1 ChIP-seq data (*i.e.*, alignment, duplicate removal, peak calling,  
559 quantification, and then PEER<sup>60</sup> normalization). The processed gene expression matrix derived from  
560 RNA-seq was downloaded directly.

561           We obtained the genotypes of the LCL samples from the 1000 Genomes Project data. We  
562 imputed 9 out of 100, 9 out of 160, and 15 out of 373 samples, respectively, from available microarray

563 data to the 1000 Genomes Project phase 3 data<sup>29</sup> on the Michigan Imputation Server<sup>61</sup>. Common variants  
564 (MAF > 5%) from the merged genotypes and the prepared phenotype matrices were used to test genetic  
565 associations to the corresponding molecular phenotypes with QTLtools<sup>62</sup>.

566

### 567 **Chromatin accessibility and gene expression levels across blood cell types**

568 ATAC-seq and RNA-seq data from multiple blood cell types throughout hematopoiesis were downloaded  
569 from GEO series GSE74912 and GSE74246, respectively<sup>28</sup>. We aligned ATAC-seq read data to the hg19  
570 reference genome, and merged data from each cell type for visualization. The genome tracks in Fig. 5  
571 were generated with fold enrichment over average genome coverage to account for library size  
572 differences. We downloaded the count matrix for RNA-seq and converted them to counts per million for  
573 comparison across cell types.

574

### 575 **MPRA data analysis**

576 We downloaded MPRA analysis tables from the two studies<sup>43,44</sup>. We extracted statistics for rs5827412  
577 and rs3808619, which were the only two putative causal PU.1 motif-altering variants at colocalized loci  
578 with MPRA data. For rs3808619, we also extracted the statistics for the other 29 and 58 variants tagging  
579 rs3808619 from Tewhey et al. and Abell et al., respectively. From Tewhey et al. data, we referred to the  
580 combined LCL analysis statistics, and from Abell et al. data, we referred to the allele effect statistics to  
581 measure the regulatory effects of variants.

582

### 583 **Differential accessibility analysis in *SP11* knockout RS4;11 lines**

584 ATAC-seq data from wild type and *SP11* knockout RS4;11 cell lines were downloaded from EMBL-EBI  
585 ArrayExpress under accession E-MTAB-8676<sup>46</sup>. We aligned the reads using Bowtie2<sup>54</sup> and removed  
586 duplicate alignments using scripts from WASP<sup>55</sup>. Then, we pooled the three replicates per genotype to  
587 call accessible regions using MACS2<sup>57</sup> with  $q < 0.05$  cutoff, and the two sets of accessible regions were  
588 merged using bedtools<sup>68</sup>. After counting the number of reads from each region using featureCount<sup>58</sup>, we  
589 tested for differential accessibility using DESeq2<sup>69</sup>. PU.1 ChIP-seq and input DNA data from  
590 unstimulated RS4;11 cell lines were downloaded from GEO series GSE71616<sup>47</sup>. After alignment using  
591 Bowtie2<sup>54</sup> and duplicate removal<sup>55</sup>, we called peaks using MACS2<sup>57</sup>. Accessible regions were stratified by  
592 whether they intersect identified PU.1 occupancy sites. The significance of observing reduced  
593 accessibility in *SP11* knockout lines was tested using a chi square test.

### 594 **Data availability**

595 Processed data for generating the figures presented in the manuscript are available at  
596 <https://github.com/BulykLab/PU1-colocalization-manuscript>. PU.1 and Histone mark ChIP-seq data are  
597 available from EMBL-EBI ArrayExpress under accession [E-MTAB-3657](https://www.ebi.ac.uk/ena/arrayexpress/experiments/E-MTAB-3657) and [E-MTAB-1884](https://www.ebi.ac.uk/ena/arrayexpress/experiments/E-MTAB-1884). ATAC-seq  
598 data in LCLs are available under European Nucleotide Archive accession [ERP110508](https://www.ebi.ac.uk/ena/arrayexpress/experiments/ERP110508). Processed RNA-  
599 seq data in LCLs are available under EMBL-EBI ArrayExpress under accession [E-GEUV-1](https://www.ebi.ac.uk/ena/arrayexpress/experiments/E-GEUV-1). The 1000  
600 Genomes Project Phase 3 genotype data are available at  
601 <ftp://ftp.1000genomes.ebi.ac.uk/vol1/ftp/release/20130502>. UK Biobank blood cell traits GWAS data  
602 from Canela-Xandri et al.<sup>21</sup> are available at <http://geneatlas.roslin.ed.ac.uk/>, and those from Vuckovic et  
603 al.<sup>8</sup> are available  
604 at [ftp://ftp.sanger.ac.uk/pub/project/humgen/summary\\_statistics/UKBB\\_blood\\_cell\\_traits](ftp://ftp.sanger.ac.uk/pub/project/humgen/summary_statistics/UKBB_blood_cell_traits). Monocyte  
605 eQTL data from BLUEPRINT<sup>16</sup> are available at <http://blueprint-dev.bioinfo.cnio.es/WP10/qtls>. Naïve B  
606 cell eQTL data from the eQTL Catalogue<sup>41</sup> are available at  
607 [ftp://ftp.ebi.ac.uk/pub/databases/spot/eQTL/sumstats/Schmiedel\\_2018/ge/Schmiedel\\_2018\\_ge\\_monocyte.](ftp://ftp.ebi.ac.uk/pub/databases/spot/eQTL/sumstats/Schmiedel_2018/ge/Schmiedel_2018_ge_monocyte_all.tsv.gz)  
608 [all.tsv.gz](https://www.ebi.ac.uk/ena/arrayexpress/experiments/E-MTAB-8676). ATAC-seq data from control and *SPI1* knockout RS4;11 cell lines are available under EMBL-  
609 EBI ArrayExpress accession [E-MTAB-8676](https://www.ebi.ac.uk/ena/arrayexpress/experiments/E-MTAB-8676), and PU.1 ChIP-seq data from RS4;11 cell line are available  
610 under GEO series accession [GSE71616](https://www.ncbi.nlm.nih.gov/geo/query/acc.cgi?acc=GSE71616). Fine-mapping results for blood cell trait GWAS are available at  
611 [https://github.com/bloodcellgwas/manuscript\\_code/tree/master/data/finemap\\_bedfiles/ukbb\\_v2](https://github.com/bloodcellgwas/manuscript_code/tree/master/data/finemap_bedfiles/ukbb_v2) and  
612 <https://www.finucanelab.org/data>.

### 613 Code availability

614 Codes for generating the figures are available at [https://github.com/BulykLab/PU1-colocalization-](https://github.com/BulykLab/PU1-colocalization-manuscript)  
615 [manuscript](https://github.com/BulykLab/PU1-colocalization-manuscript). We trained a PU.1 motif gkm-SVM model using LS-GKM ([https://github.com/Dongwon-](https://github.com/Dongwon-Lee/lsgkm)  
616 [Lee/lsgkm](https://github.com/Dongwon-Lee/lsgkm)). We performed genotype imputation using the Michigan Imputation Server  
617 (<https://imputationserver.sph.umich.edu/>). We processed genotype data using BCFtools  
618 (<https://samtools.github.io/bcftools/bcftools>) and PLINK (<https://www.cog-genomics.org/plink2/>). We  
619 estimated hidden factors for QTL analyses using PEER (<https://github.com/PMBio/peer>). We generated  
620 sets of null variants for PU.1 bQTL enrichment analysis using SNPsnap  
621 (<https://data.broadinstitute.org/mpg/snpsnap/>). We performed colocalization analysis using JLIM  
622 (<https://github.com/cotsapaslab/jlim>) and Coloc (<https://chr1swallace.github.io/coloc/>). The Fuji plot (Fig.  
623 3b) was made using code from <https://github.com/mkanai/fujiplot>. We adapted codes from  
624 LocusCompareR (<https://github.com/boxiangliu/locuscomparer>) to create association plots. We  
625 performed fine-mapping analysis for *ZNF608* eQTL in LCLs using SuSiE  
626 (<https://stephenslab.github.io/susieR/>).

627 **Acknowledgments**

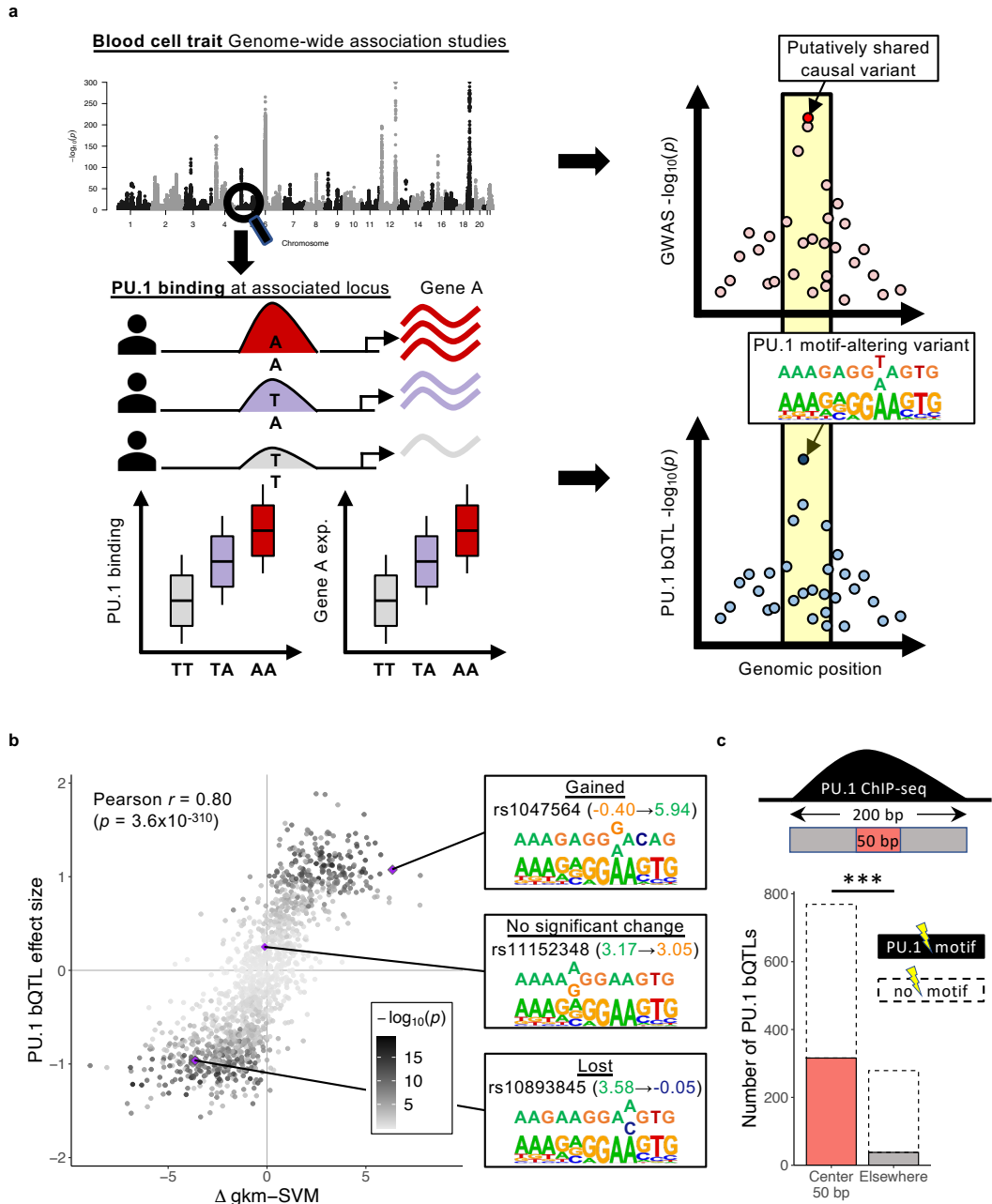
628 We thank members of the Bulyk lab, Vijay Sankaran, Shamil Sunyaev, Alexander Gusev, and members  
629 of the Raychaudhuri lab, including, but not limited to, Soumya Raychaudhuri, Kazuyoshi Ishigaki, Saori  
630 Sakaue, Tiffany Amariuta, Yang Luo, and Samira Asgari for helpful discussion and Shubham Khetan and  
631 Shamil Sunyaev for critical reading of the manuscript. This work was funded by a grant from the Brigham  
632 Research Institute Fund to Sustain Research Excellence and NIH grant R01 HG010501.

633 **Author Contributions**

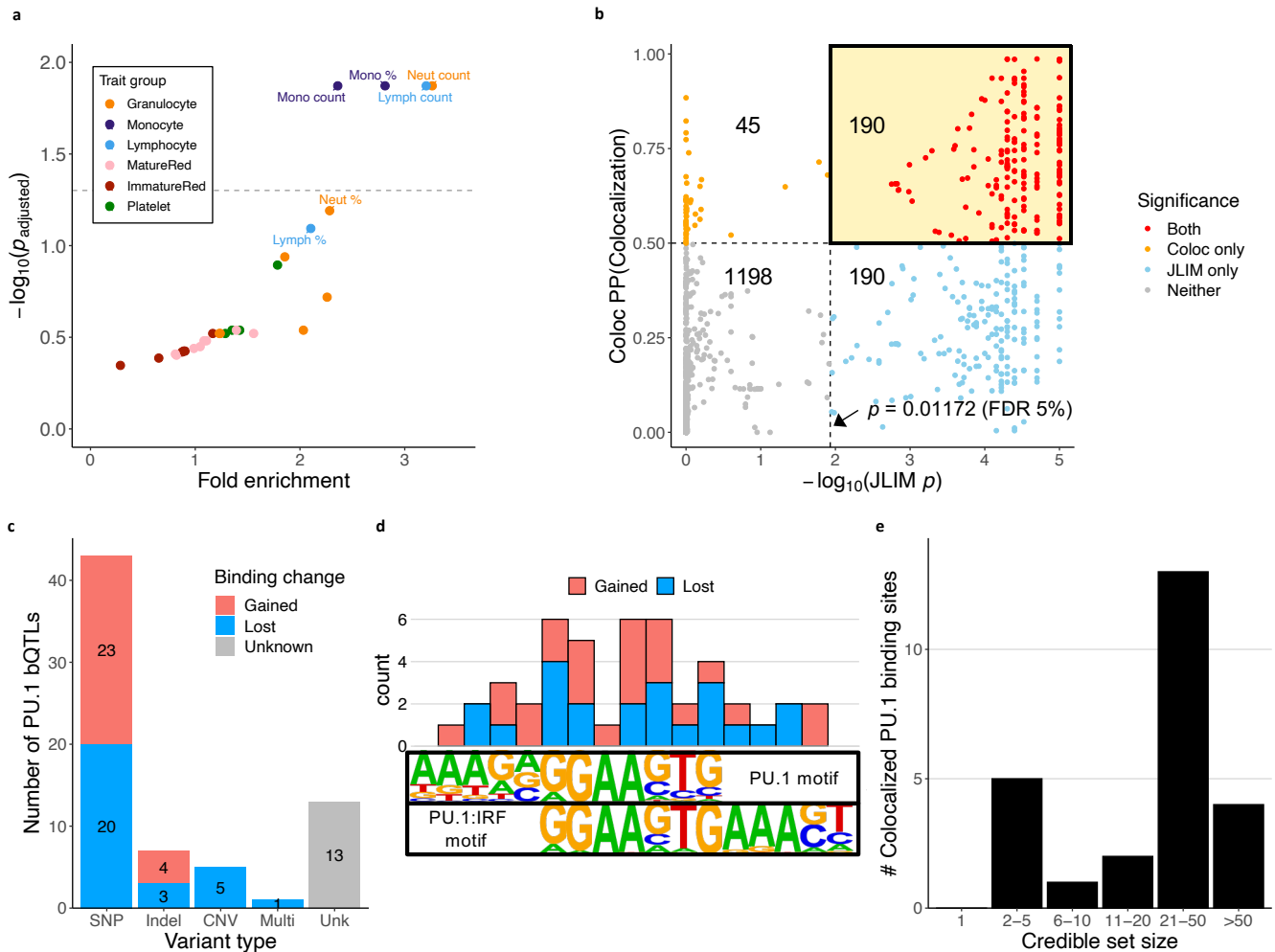
634 R.J. and M.L.B. conceived and designed the research project. R.J. performed all analyses and prepared  
635 the figures. M.L.B. supervised the research. R.J. and M.L.B. wrote the manuscript.

636 **Ethics Declarations**

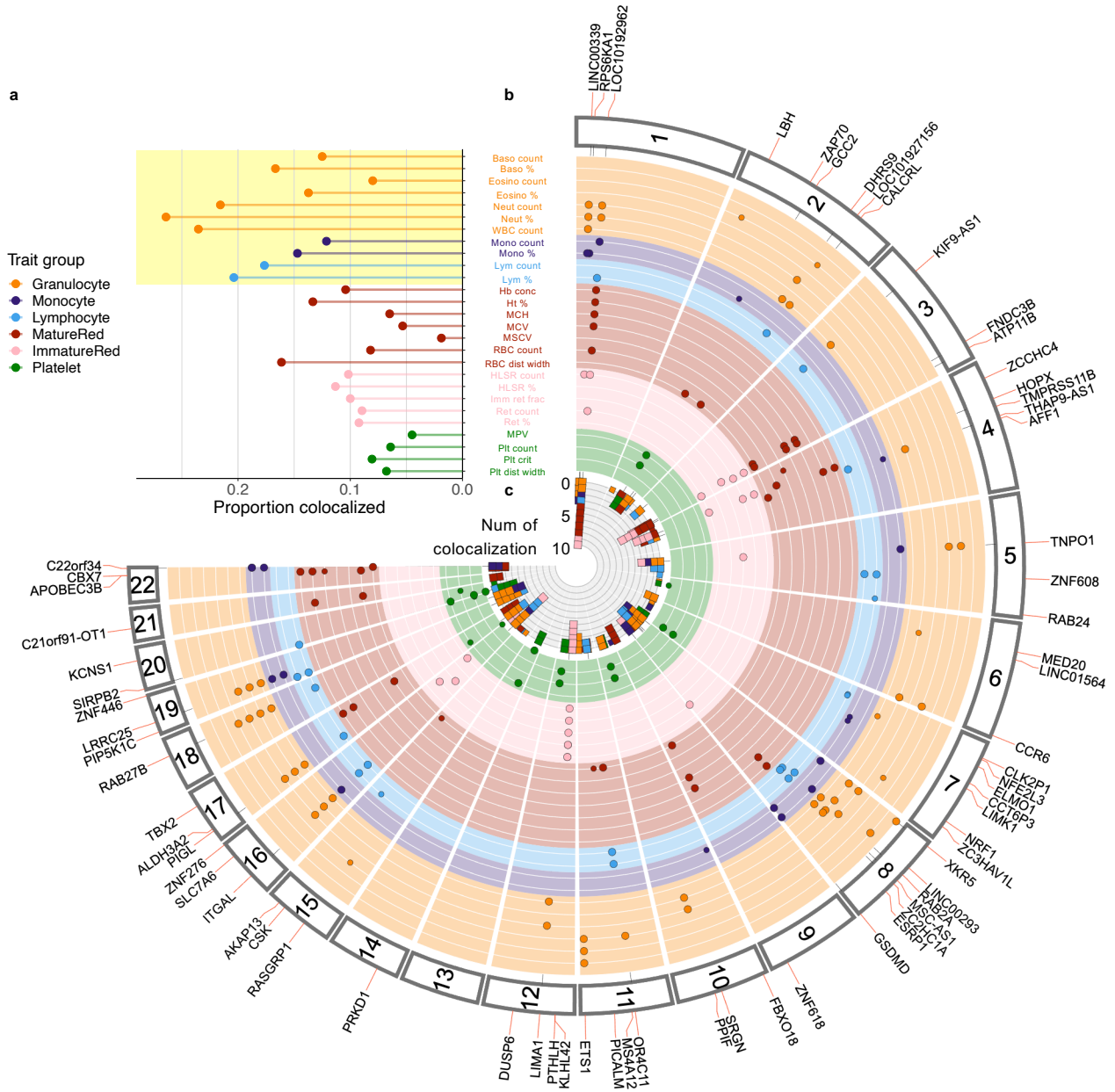
637 The authors declare no competing interests.



**Fig. 1 | Relevance of PU.1 bQTLs in LCLs to blood cell trait associations.** (a) (Left) Blood cell trait-associated loci may have overlapping PU.1 bQTLs and, potentially, expression QTL (eQTL) associations. (Right) Significant colocalization suggests that the causal variants are shared. If there is a PU.1 motif-altering variant at a colocalized PU.1 bQTL, the variant is likely to be the shared causal variant. (b) Comparison of changes in motif score ( $\Delta$  gkm-SVM) and estimated bQTL effect sizes at PU.1 motif-altering variants within 200bp PU.1 ChIP-seq peaks. The color represents the  $-\log_{10}(p)$  of PU.1 bQTL association (linear regression). (c) Number of significant PU.1 bQTLs with PU.1 motif-altering variants at each region within the 200bp PU.1 ChIP-seq peaks. \*\*\*:  $p < 2.2 \times 10^{-16}$  (Fisher's exact test).

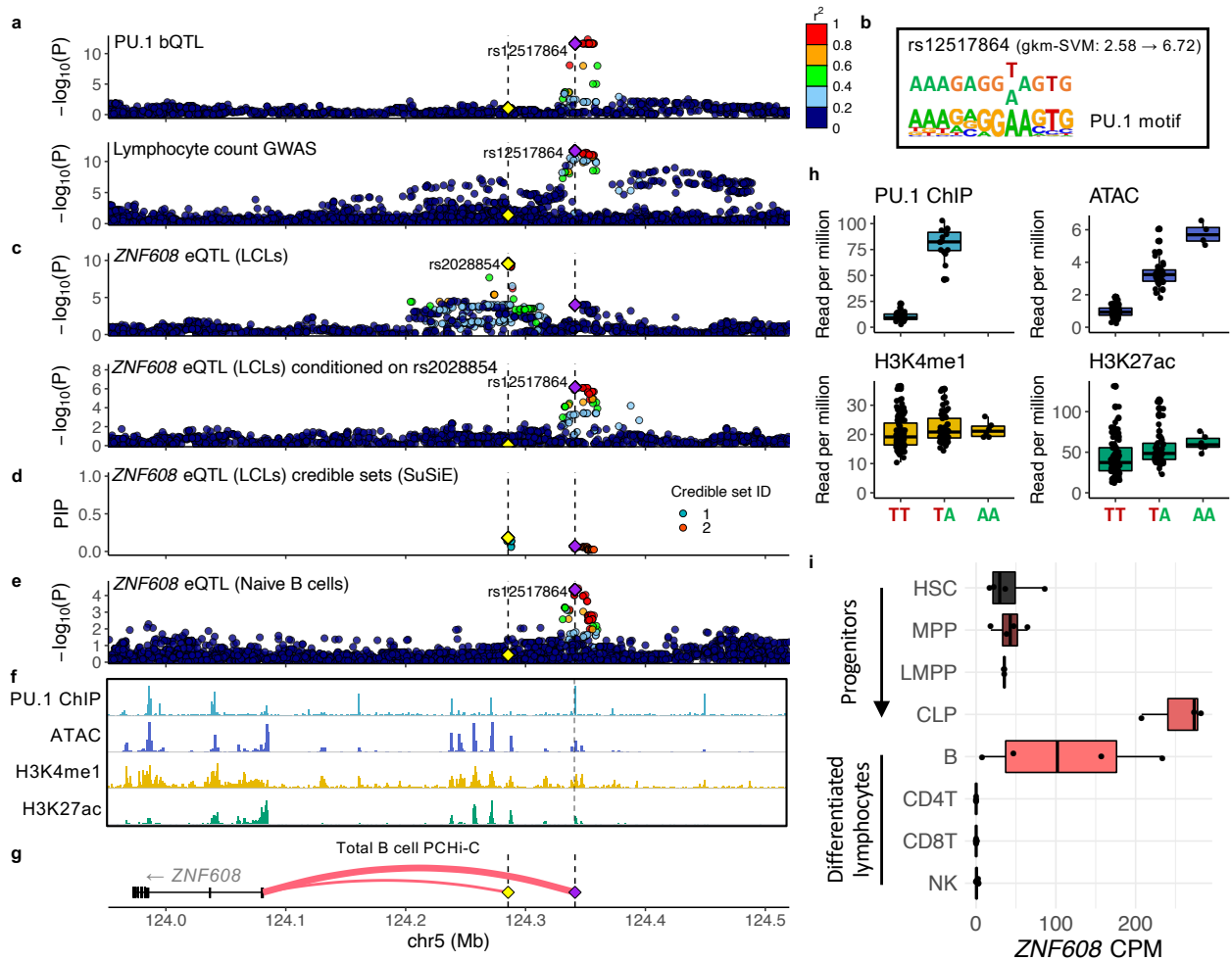


**Fig. 2 | Colocalization of blood cell traits GWAS and PU.1 bQTLs.** (a) Enrichment of PU.1 bQTLs for associations to specific blood cell traits. Traits with empirical adjusted  $p < 0.05$  (above the dashed line) are labeled. Abbreviations of blood cell traits are described in Supplementary Table 3. (b) Colocalization results from JLIM and Coloc. Each point is a PU.1 bQTL - Trait pair. The number shown in each quadrant is the number of points within the significance category. Dashed lines indicate the respective significance thresholds (JLIM:  $p < 0.01172$  (FDR 5%), Coloc: PP(colocalized)  $> 0.5$ ). (c) The types of putative causal variants at colocalized PU.1 bQTLs that alter PU.1 motifs or the copy number of the PU.1 occupancy site. SNPs, indels, and multi-variants alter PU.1 motifs. CNV: copy number variation altering copy number of PU.1 binding sites; Multi: multiple variants in perfect LD ( $r^2 = 1$ ) within a PU.1 motif sequence; Unk (Unknown): No variant altering PU.1 motif sequence or its copy number. (d) Number of PU.1 motif-altering SNPs at each nucleotide position at colocalized PU.1 binding sites. Motif logos are from Homer database. (e) Blood cell trait GWAS credible set size at loci with colocalized PU.1 bQTLs and a PU.1 motif-altering variant. Only 25 loci with fine-mapping result in Vuckovic et al. 2020 are represented.

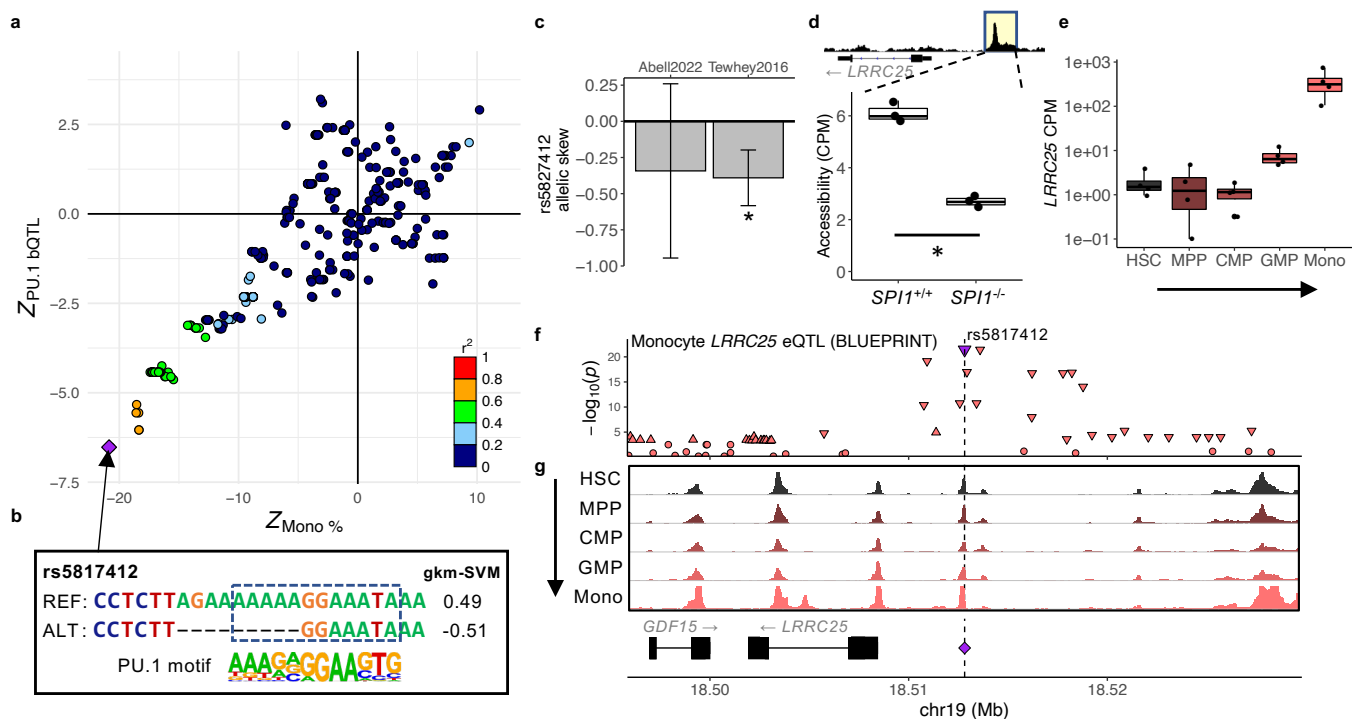


**Fig. 3 | Distribution of colocalized loci across the genome. (a)** Proportion of tested loci with significant colocalization. The colors represent the trait groups. The blood cell traits highlighted in yellow correspond to white blood cell traits. Abbreviations of blood cell traits are described in Supplementary Table 3. **(b)** Fuji plot depicting the genomic distribution of blood cell trait-associated loci that show high-confidence colocalization with PU.1 bQTLs. The colors are as in panel a. **(c)** The stacked bar plot at the center shows the number of traits each PU.1 bQTL colocalizes with.

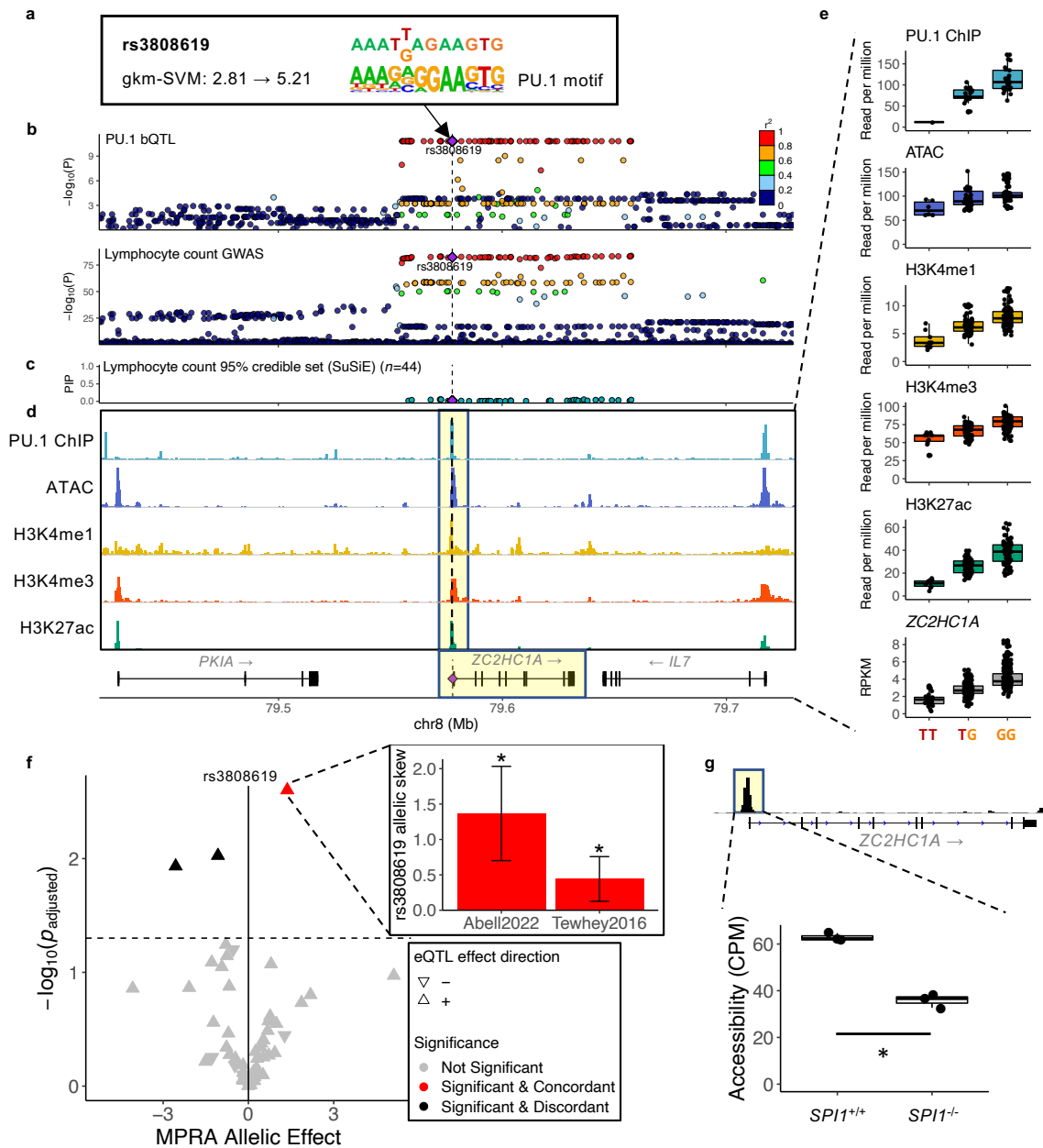




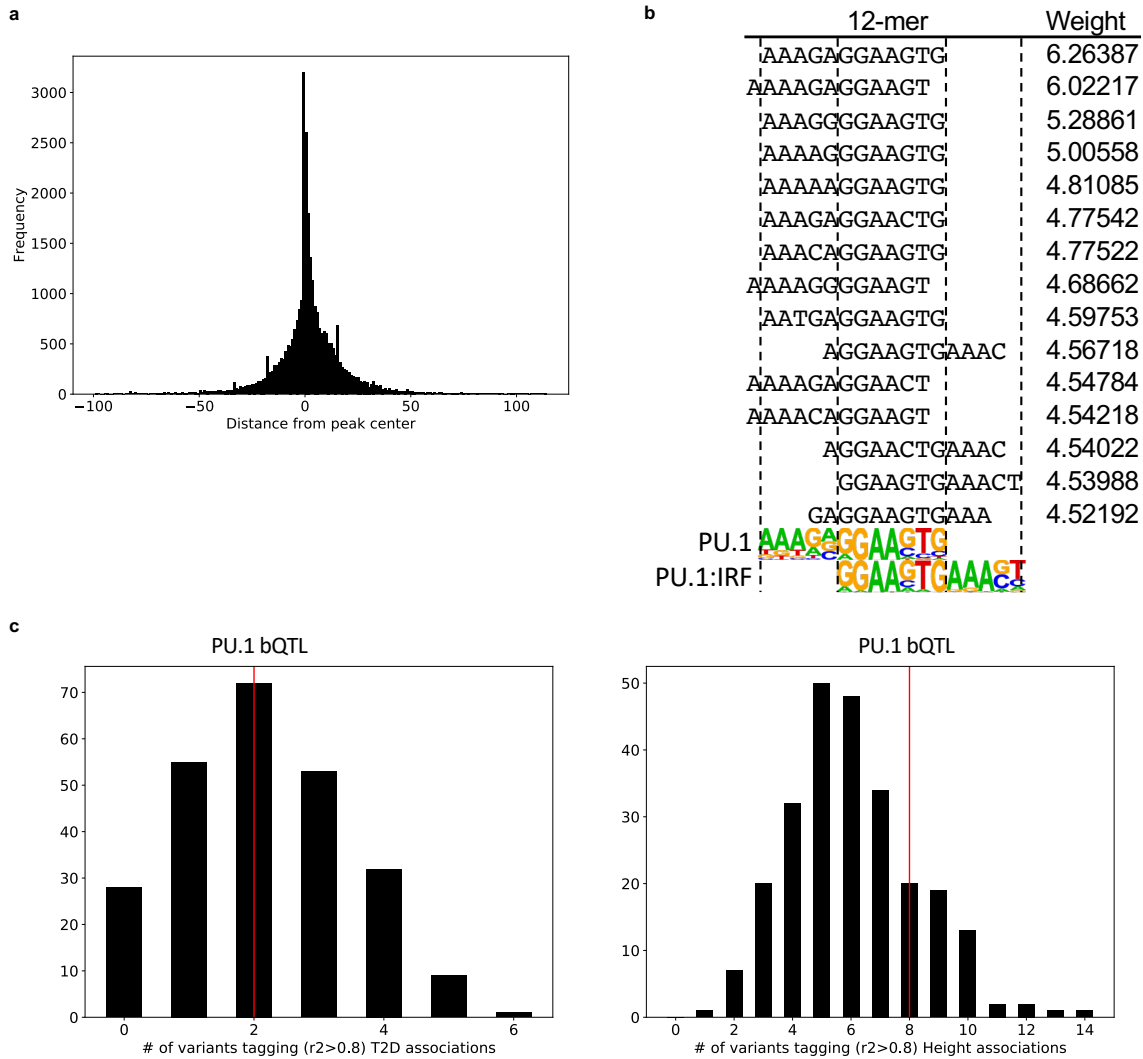
**Fig. 4 | PU.1 motif alteration pinpoints a lymphocyte count-associated variant that is a secondary *ZNF608* eQTL variant.** (a, c-e, g) PU.1 motif-altering variant rs12517864 is shown as a purple diamond, and the *ZNF608* eQTL lead variant rs2028854 is shown as a yellow diamond. Vertical dashed lines mark the position of these two variants. Unless noted otherwise, points are colored by LD  $r^2$  with respect to rs12517864. (a) PU.1 bQTL and lymphocyte count association signals. (b) The effect of rs2028854 on the sequence with respect to the PU.1 binding motif. (c) (Top) Primary *ZNF608* eQTL signals in LCLs. LD  $r^2$  is calculated with respect to rs2028854, the lead variant. (Bottom) *ZNF608* eQTL signals in LCLs conditioned on the rs2028854 dosage. (d) Fine-mapping result of *ZNF608* eQTL signals in LCLs, using SuSiE. Points are colored by the credible set they belong to. PIP: Posterior inclusion probability. (e) *ZNF608* eQTL association signals in naïve B cells (DICE). (f) Genome tracks of PU.1 ChIP-seq, ATAC-seq, H3K4me1 and H3K27ac ChIP-seq assayed in GM12878. (g) Gene track showing *ZNF608* and the two variants. The weights of the red curves indicate the CHiCAGO scores calculated in Javierre et al. 2016. (h-i) On top of the box plots, all the data points are shown. (h) The effect of rs12517864 dosage on various molecular phenotypes shown in panel f. For PU.1 ChIP-seq data, there weren't any individuals with homozygous alternate allele (AA). (i) *ZNF608* expression levels (count per million) through lymphocyte differentiation and across various lymphocyte types. HSC: hematopoietic stem cell, MPP: multipotent progenitor, LMPP: lymphoid-primed multipotent progenitor, CLP: common lymphoid progenitor, B: B cell, CD4T: CD4<sup>+</sup> T cell, CD8T: CD8<sup>+</sup> T cell, NK: natural killer cell.



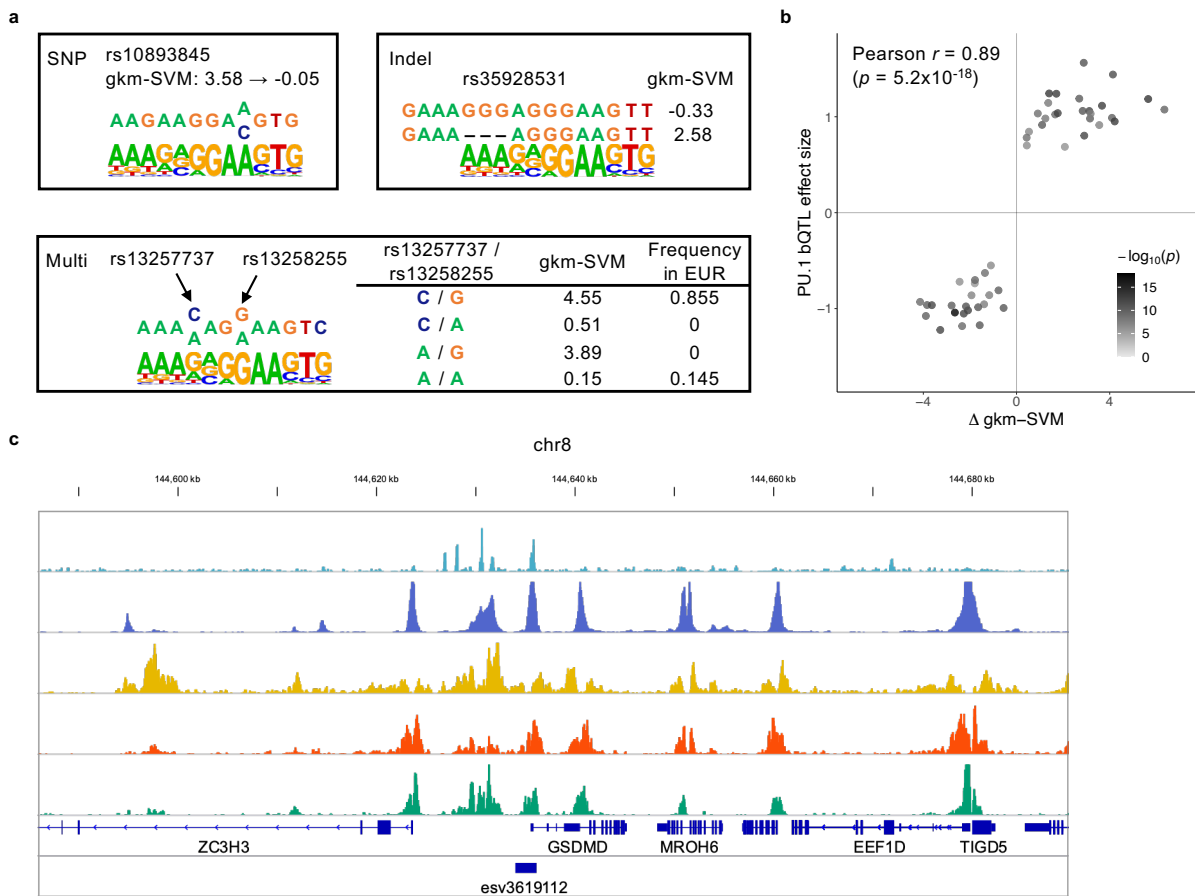
**Fig. 5 | PU.1 motif-altering deletion rs5827412 at *LRRC25* locus associated with lower monocyte counts.** (a) PU.1 bQTL and monocyte percentage association signals colocalize. (b) The effect of rs5827412 on the PU.1 motif. (c) Reduced reporter activity by rs5827412 in log<sub>2</sub> fold change. Error bars indicate 95% confidence intervals. \*: adjusted  $p < 0.05$ . (d-e) Boxplots are formatted as in Fig 4. (d) A boxplot showing PU.1-dependent reduction in chromatin accessibility levels (count per million) at the regulatory element surrounding rs5827412 in control pro-B cell lines ( $SPI1^{+/+}$ ) and counterparts with  $SPI1$  knocked out ( $SPI1^{-/-}$ ). Regions highlighted in yellow marks the accessible region corresponding to the boxplot.  $n = 3$  for each condition. \*: DESeq2 adjusted  $p < 0.05$ . (e) A boxplot showing *LRRC25* expression levels (count per million) through monocyte differentiation. HSC: hematopoietic stem cell, MPP: multipotent progenitor, CMP: common myeloid progenitor, GMP: granulocyte-macrophage progenitor, Mono: monocyte. (f-g) Purple triangle and diamond, as well as the dashed line, mark rs5827412. (f) Monocyte *LRRC25* eQTL association. Downward and upward triangles indicate the direction of effect (down- and up-regulation, respectively) for variants with  $p < 1 \times 10^{-3}$ . (g) ATAC-seq tracks as fold enrichment over average (range 0-40) for various blood cell types through monocyte differentiation.



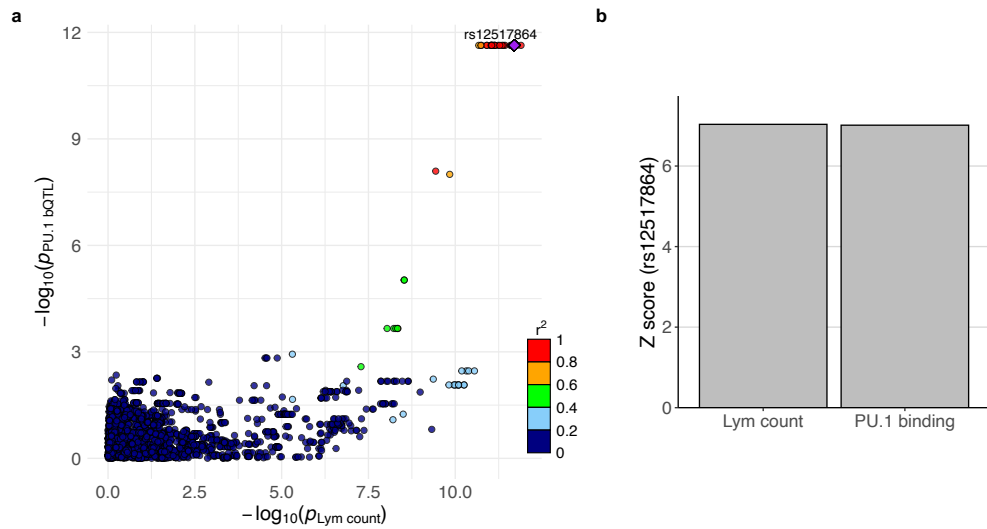
**Fig. 6 | *ZC2HC1A* locus: PU.1 motif-alteration highlights a regulatory variant among those in high LD. (a-d)** PU.1 motif-altering variant rs3808619 is shown as a purple diamond. Vertical dashed line also mark the position of this variant. **(a)** The effect of rs3808619 on the PU.1 composite motif. **(b)** PU.1 bQTL and lymphocyte count association signal at the *ZC2HC1A* locus. **(c)** Posterior inclusion probability (PIP) of variants in the 95% credible set of lymphocyte count association at the *ZC2HC1A* locus. **(d)** Genome tracks of PU.1 ChIP-seq, ATAC-seq, H3K4me1, H3K4me3, H3K27ac ChIP-seq assayed in GM12878. The highlighted regions correspond to molecular phenotypes with QTL associations in **e**. **(e)** The effect of rs3808619 dosage on various molecular phenotypes shown in panel d. Box plots are formatted as in Fig. 4. **(f)** Regulatory effects of rs3808619 and 58 tagging variants in a reporter assay. MPRA allelic effect corresponds to log<sub>2</sub> fold change of regulatory activity of the oligo sequence with the alternate allele over that with the reference allele. The inset shows the allelic skew estimates with 95% confidence intervals from Abell et al. and Tewhey et al. \*: adjusted  $p < 0.05$ . **(g)** PU.1-dependent reduction in chromatin accessibility levels (count per million) at the regulatory element surrounding rs3808619 in control pro-B cell lines (*SPI1*<sup>+/+</sup>) and counterparts with *SPI1* knocked out (*SPI1*<sup>-/-</sup>).  $n = 3$  for each condition. \*: DESeq2 adjusted  $p < 0.05$ . The panel is formatted as in Fig. 5d.



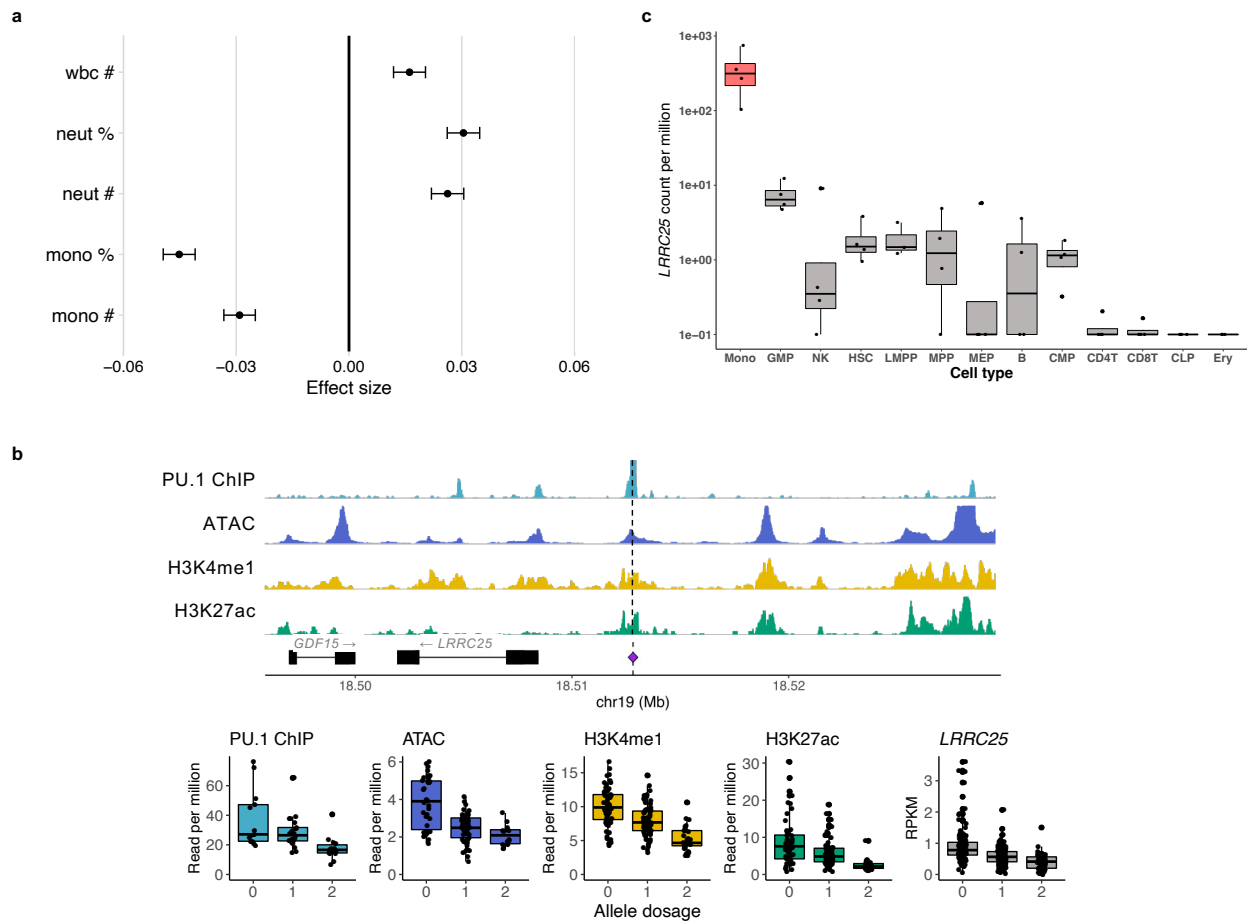
**Extended Data Fig. 1 | Properties of PU.1 binding sites and bQTLs.** (a) Position of PU.1 motifs at PU.1 binding sites. The bp distance is measured from the center of a 200 bp PU.1 ChIP-seq peak. (b) 12-mers with the highest (top 15) gkm-SVM weights aligned to PU.1 motif and PU.1:IRF composite motif. (c) Lack of enrichment in PU.1 bQTL lead variants tagging ( $LD r^2 > 0.8$ ) type 2 diabetes (T2D) and height GWAS associations. The histogram shows the number of variants tagging GWAS associations for each of 250 sets of null variants. The red lines indicate the number of PU.1 bQTL lead variants tagging GWAS associations.



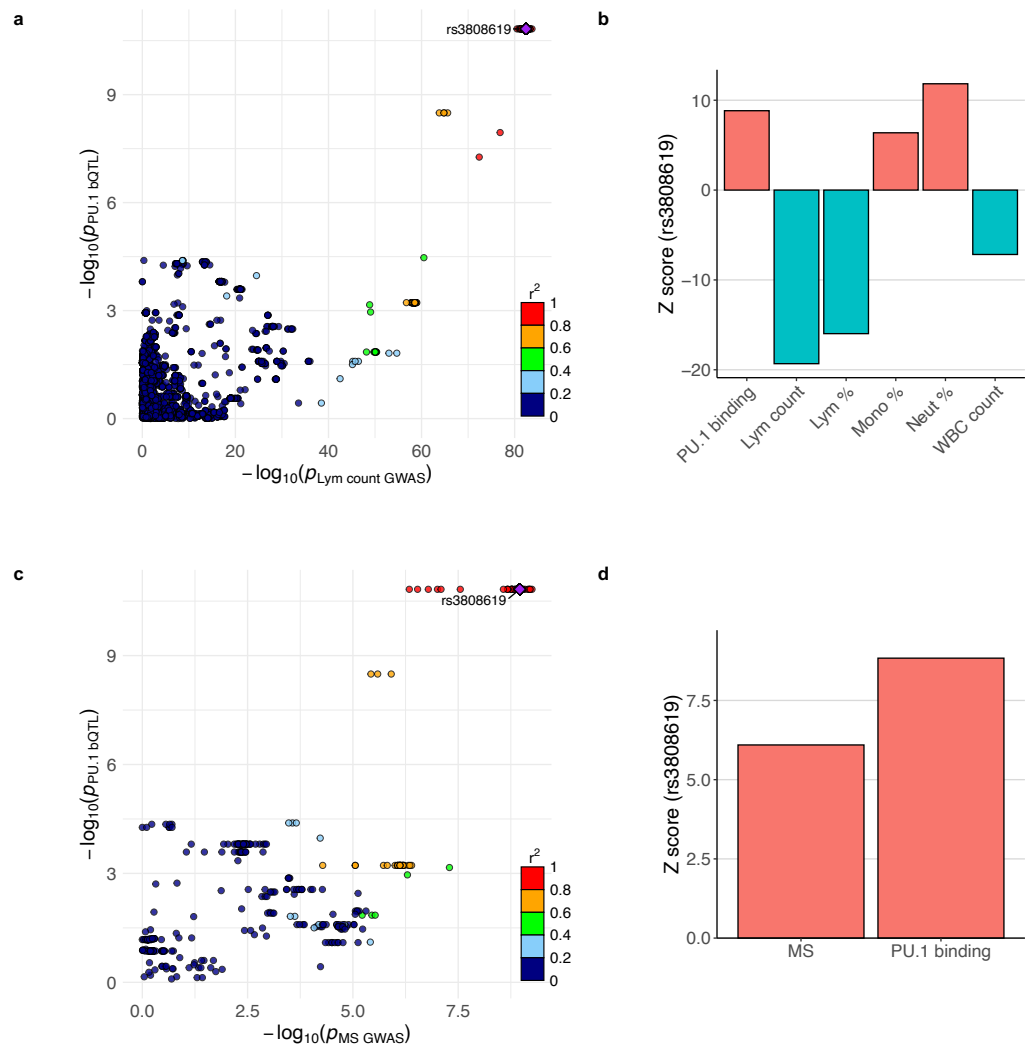
**Extended Data Fig. 2 | Examples of variants affecting PU.1 binding.** (a) Examples of PU.1 motif-altering variants. Categorization of the variants correspond to Fig. 2b. EUR: European ancestry population in the 1000 Genomes Project. (b) Comparison of changes in motif score ( $\Delta$  gkm-SVM) and estimated bQTL effect sizes of PU.1 motif-altering variants (SNPs and indels) at 49 colocalized loci. (c) An example of a copy number variation (esv3619112) affecting a PU.1 binding site.



**Extended Data Fig. 3 | Colocalization of PU.1 bQTL and lymphocyte count association signals at *ZNF608* locus.** (a) Merged association plot for PU.1 bQTL and lymphocyte count association signals. Points are colored by LD  $r^2$  with respect to rs12517864, which is labeled with a purple diamond. (b) Z scores of rs12517864 for lymphocyte count and PU.1 bQTL association.



**Extended Data Fig. 4 | Effects of PU.1 motif-altering deletion rs5827412.** (a) GWAS effect size estimates for rs5827412 on 5 blood cell traits. The error bars indicate 95% confidence interval. Abbreviations of blood cell traits are described in Supplementary Table 2. (b-c) Boxplots are formatted as in Fig 4. (b) Regulatory QTL effects of rs5827412. (top) Genome tracks show PU.1 ChIP-seq, ATAC-seq, and H3K4me1 and H3K27ac ChIP-seq data from LCLs, respectively. (bottom) 4 phenotype values in read per million for each genome track and reads per kilobase million for *LRRC25* expression levels. Allele dosage corresponds to the deletion allele. (c) *LRRC25* expression level across 13 blood cell types. Monocyte is colored red. Cell types abbreviated as in Supplementary Fig. 1.



**Extended Data Fig. 5 | Colocalization of PU.1 bQTL and multiple sclerosis association signals at *ZC2HC1A* locus.** (a,c) Points are colored by LD  $r^2$  in the 1000 Genomes Project European population, with respect to rs3808619, which is labeled with a purple diamond. (a) Merged association plot for PU.1 bQTL and lymphocyte count association signals. (b) Z scores of rs3808619 for PU.1 bQTL and 5 blood cell traits association. (c) Merged association plot for PU.1 bQTL and multiple sclerosis (MS) association signals. (d) Z scores of rs3808619 for MS and PU.1 bQTL association.