# Genotype inference from aggregated chromatin accessibility data reveals genetic regulatory mechanisms

Brandon M. Wenz*[1], Yuan He*[2], Nae-Chyun Chen[3], Joseph K. Pickrell[4], Jeremiah H. Li[4], Max F. Dudek[5], Taibo Li[2], Rebecca Keener[2], Benjamin F. Voight[6,7,8], Christopher D. Brown[6], Alexis Battle[2,3,9,10,11]

1. Genetics and Epigenetics Program, Cell and Molecular Biology Graduate Group, Biomedical Graduate Studies, University of Pennsylvania - Perelman School of Medicine, Philadelphia PA 19104
2. Department of Biomedical Engineering, Johns Hopkins University; Baltimore, MD, 21218
3. Department of Computer Science, Johns Hopkins University; Baltimore, MD, 21218
4. Gencove, Inc., New York, NY, 11101
5. Graduate Group in Genomics and Computational Biology, University of Pennsylvania, Philadelphia, PA 19104
6. Department of Genetics, University of Pennsylvania - Perelman School of Medicine, Philadelphia, PA, 19104
7. Department of Systems Pharmacology and Translational Therapeutics, University of Pennsylvania - Perelman School of Medicine, Philadelphia PA, 19104
8. Institute for Translational Medicine and Therapeutics, University of Pennsylvania – Perelman School of Medicine, Philadelphia, PA, 19104
9. Department of Genetic Medicine, Johns Hopkins University; Baltimore, MD, 21218
10. Malone Center for Engineering in Healthcare, Johns Hopkins University, Baltimore, MD, 21218
11. Data Science and AI Institute, Johns Hopkins University, Baltimore, MD, 21218

*: These authors contributed jointly to the work.

Correspondence to:

Alexis Battle, PhD

ajbattle@jhu.edu

28

# Abstract

30

**Background**

Understanding the genetic causes for variability in chromatin accessibility can shed light on the molecular mechanisms through which genetic variants may affect complex traits. Thousands of ATAC-seq samples have been collected that hold information about chromatin accessibility across diverse cell types and contexts, but most of these are not paired with genetic information and come from diverse distinct projects and laboratories.

**Results**

We report here joint genotyping, chromatin accessibility peak calling, and discovery of quantitative trait loci which influence chromatin accessibility (caQTLs), demonstrating the capability of performing caQTL analysis on a large scale in a diverse sample set without pre-existing genotype information. Using 10,293 profiling samples representing 1,454 unique donor individuals across 653 studies from public databases, we catalog 23,381 caQTLs in total. After joint discovery analysis, we cluster samples based on accessible chromatin profiles to identify context-specific caQTLs. We find that caQTLs are strongly enriched for annotations of gene regulatory elements across diverse cell types and tissues and are often strongly linked with genetic variation associated with changes in expression (eQTLs), indicating that caQTLs can mediate genetic effects on gene expression. We demonstrate sharing of causal variants for chromatin accessibility and diverse complex human traits, enabling a more complete picture of the genetic mechanisms underlying complex human phenotypes.

51

52    **Conclusions**

53    Our work provides a proof of principle for caQTL calling from previously ungenotyped samples,

54    and represents one of the largest, most diverse caQTL resources currently available, informing

55    mechanisms of genetic regulation of gene expression and contribution to disease.


56    # Introduction

57    Genome wide association studies (GWAS) have identified thousands of loci and

58    common human genetic variants that are associated with a wide range of complex human traits,

59    diseases, and risk factors[1]. GWAS variants are often found in noncoding regions, where they

60    are likely to be involved in gene regulation[2,3]. However, a full picture of the causal regulatory

61    elements that underlie these associations remains incomplete for most loci[4]. Characterizing

62    the genetic effects of variants on gene expression as revealed by expression quantitative trait

63    locus (eQTL) mapping has provided insights into the molecular basis of phenotypes[3,5–7].

64    Although some eQTL variants directly affect open-reading frames, the vast majority are in non-

65    coding regions, as has been described for GWAS variants. Connecting causal variants to the

66    regulatory elements and the genes of action that they perturb remains a central goal of the post-

67    GWAS era.

68    Accessibility of chromatin regions to transcriptional machinery is a key factor in gene

69    regulation[8,9], and genetic variants can affect complex traits through changes in gene

70    expression levels that are mediated by chromatin accessibility[10,11]. Improved understanding

71    of the mechanisms involved in chromatin accessibility, revealed by genetic variants that

72    modulate chromatin accessibility (i.e., caQTLs), has the potential to illuminate the molecular

73    mechanisms and genetic regulatory architecture of complex traits. caQTLs have been

74    measured in a variety of tissue and cell types, at both bulk[12–16] and single-cell

75    resolutions[17]. caQTLs have been used in a variety of studies to characterize gene expression

76    regulation[18], and to propose mechanisms for risk loci identified through GWAS[19]. caQTLs

77    may co-occur with eQTLs together, thus describing a more complete picture of the genetic

78    mechanism underlying GWAS-associated signals. However, relevant caQTLs may be

79    discovered even in the absence of any established eQTL, as eQTL studies may not include the

80    relevant cell type or environmental context to reveal the change to gene expression. Analysis of

81    the contribution of caQTLs to complex human traits can help us better understand the molecular

82    impact of these variants and the mechanism(s) driving GWAS signals. To date, caQTL studies

83    have mostly been performed in analyses restricted to single tissue/cell types, a majority of which

84    have assayed a limited number of samples.

85        The Assay for Transposase-Accessible Chromatin using sequencing (ATAC-seq)

86    technology has been widely used to capture chromatin accessibility in various cell types and

87    experimental conditions[20–22]. There is a rapidly accumulating trove of ATAC-seq data

88    generated from various experiments, labs, and conditions. This wealth of information has the

89    potential to boost power for caQTL analysis. Unfortunately, many of these samples do not have

90    matched genotype information, a necessary component for QTL analyses. ATAC-seq reads,

91    however, naturally carry the sequence information at nucleotide resolution, providing the

92    possibility of inferring sample genotypes from these data directly.

93        Here, we have selected and evaluated pipelines to uniformly process ATAC-seq

94    samples, including peak calling and genetic variant calling directly from ATAC-seq reads. We

95    called genotypes using a pipeline incorporating Gencove's low-pass sequencing methods

96    applied to ATAC-seq reads in accessible chromatin, which utilizes imputation to infer genotype

97    for variants that are located outside of regions covered by observed reads in accessible

98    regions[23,24]. We benchmarked this pipeline, using gold standard genotype information

99    available for a subset of samples, and compared it with other methods. Because large-scale

100   public data often contains multiple samples from the same donor or even the same cell line, we

101     also developed a method to automatically infer donor assignment based on genotype from the

102     called variants. Peak calling from thousands of diverse samples presents challenges of

103     identifying true, distinct regions of chromatin accessibility rather than low-signal false positives,

104     or large regions merged from what should be distinct peaks[25,26]. Based on comparisons

105     across various peak-calling approaches, we finalized a pipeline based on an Genrich, an ATAC-

106     seq specific method[27] for collectively calling peaks across large, diverse data sets and

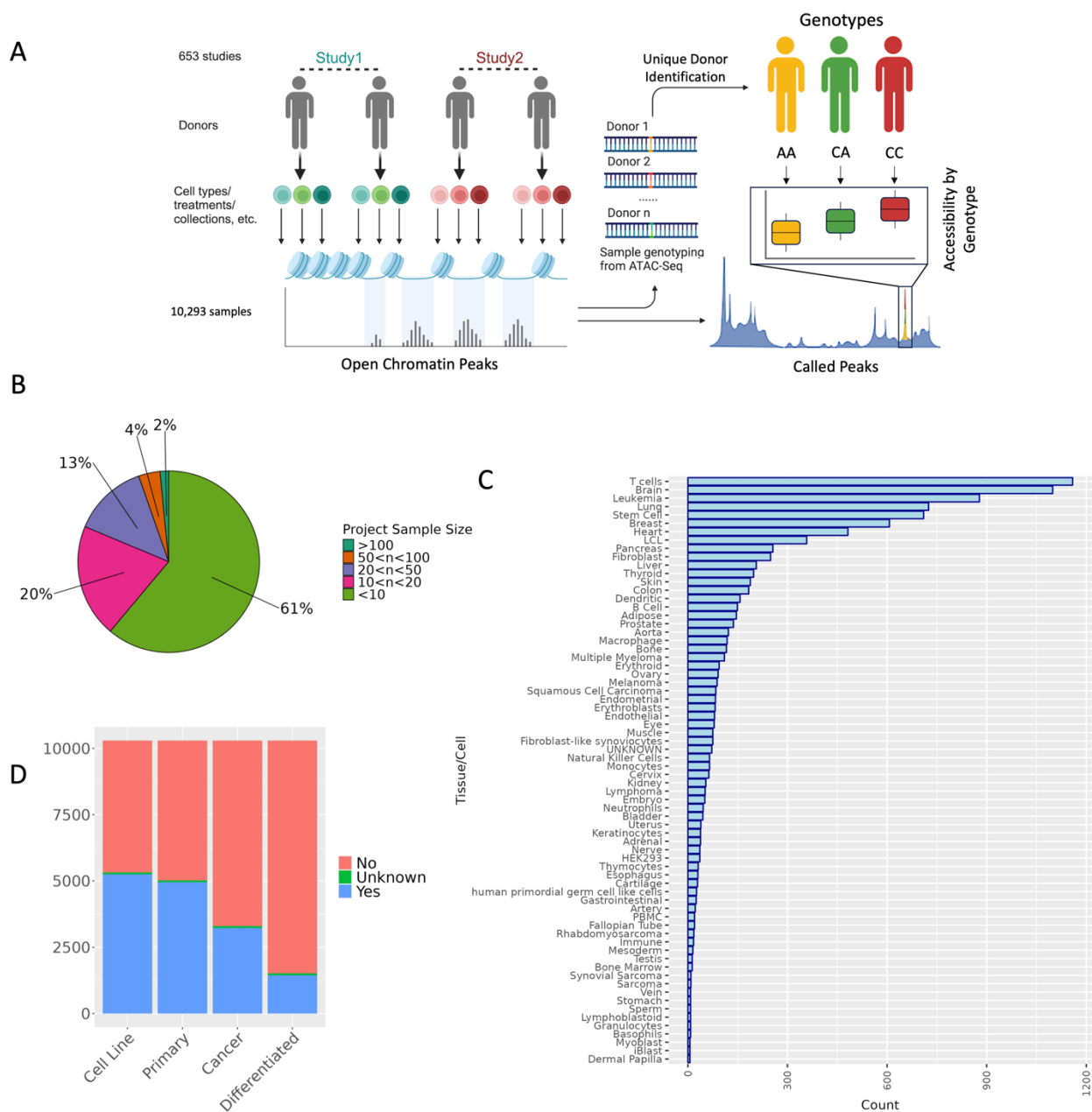107     quantifying accessibility in each peak.

108        Using our ATAC-seq derived genotypes and accessibility estimates across peaks and

109     samples, we then called caQTLs from this collection of publicly available ATAC-seq data. We

110     identified thousands of caQTLs that share a causal signal with GWAS signals, many of which

111     are not explained by known eQTLs. Additionally, we identified many GWAS signals that appear

112     to share a causal signal with both eQTLs and caQTLs, enabling a more comprehensive analysis

113     predicting target gene, gene regulatory element and even potential transcription factors that are

114     driving GWAS signals for a variety of complex human traits. Furthermore, to capture context-

115     specific caQTLs, we inferred clusters of samples with similar accessibility profiles, mostly

116     reflecting cell or tissue type, and identified cluster-specific caQTLs. With the captured global and

117     cluster-specific caQTLs, we investigated potential mechanisms involving transcription factors

118     and their role in target gene regulation.

119  # Results

120  **Accurate genotyping and imputation based on ATAC-seq reads from public**

121  **repositories**

122        We established a workflow to collect a diverse set of publicly available ATAC-seq

123     datasets and ascertain donor genotype from ATAC-seq reads, with the overall objective of

124     mapping genetic variants that are associated with differences in chromatin accessibility for

125    diverse tissues and contexts on a large scale (Figure 1A). We collected 10,293 human samples

126    from 653 projects from the Gene Expression Omnibus (GEO) data repository, where most

127    projects were comprised of 10 or fewer samples (Figure 1B, Supplementary Table 1). The

128    aggregated data includes samples from a wide variety of tissues or cell types (Figure 1C),

129    labeled based on a manual curation of project abstracts, sample labels, and project methods,

130    with the most common cell/tissue types including T cells and brain.  Additionally, based on our

131    metadata review, both cancer and normal primary tissue are well represented, along with cell

132    lines and experimentally differentiated cell types (Figure 1D). The diversity of samples highlights

133    the value of a workflow that can aggregate data and genotype samples from ATAC-seq reads,

134    providing an overall large sample size, but also tissue-specific sample sizes larger than any

135    existing genotyped chromatin accessibility study for several individual tissues including lung,

136    breast, heart, and pancreas[12,28–30].

137

**Figure 1. Study overview and characteristics of specimens utilized in this study. (A)**

Overview of study design to jointly call genotype and caQTLs across studies. Human ATAC-seq

datasets were obtained from GEO. After variant-calling (Methods), we identified the unique

donors in the dataset (Methods) for use in caQTL mapping. **(B)** The distribution of the number of

samples collected across all n=653 studies. **(C)** Frequency of the Cell/Tissue types present in
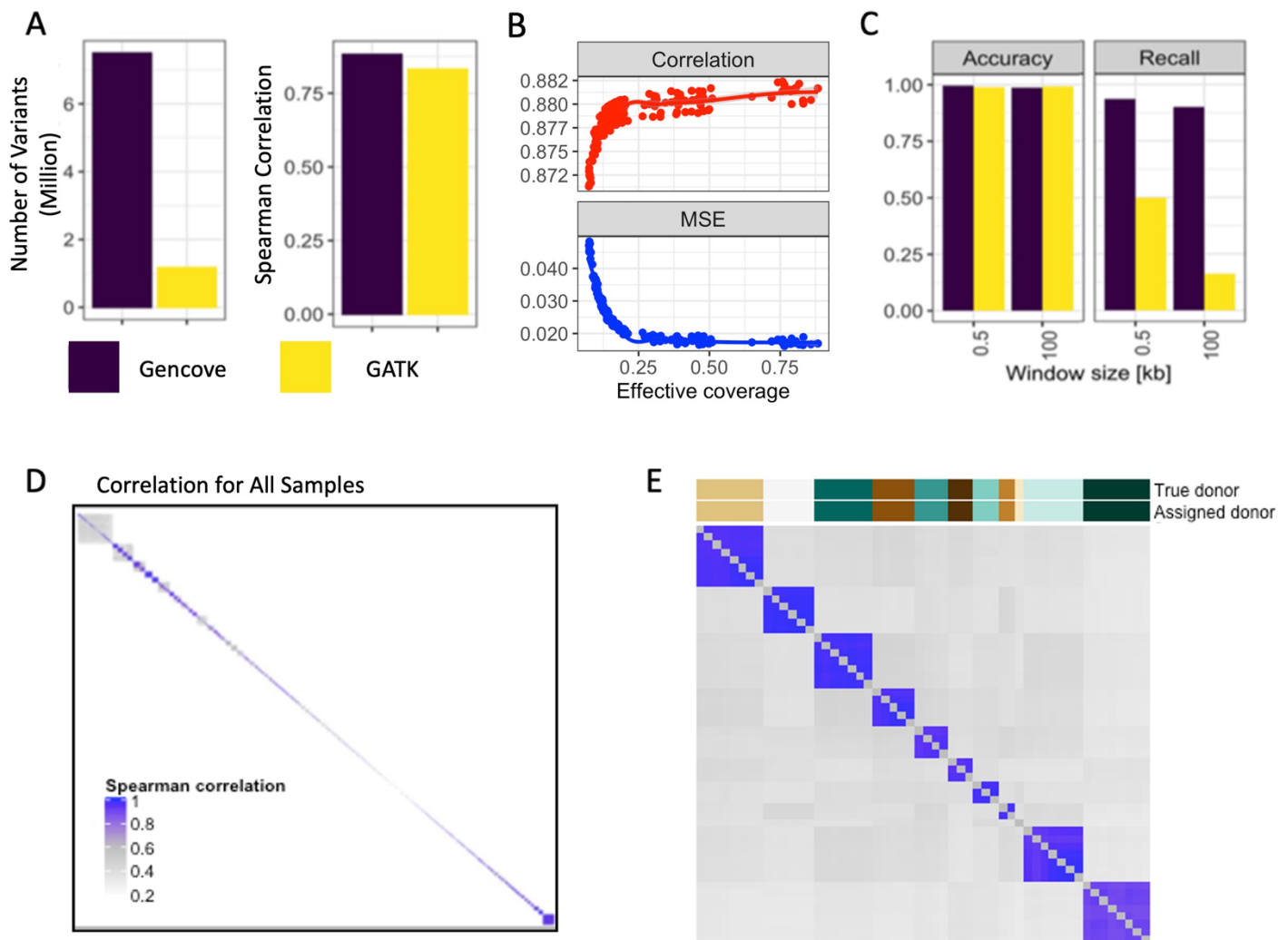
143    samples collected across studies based on manual metadata curation **(D)** Frequencies of

144    cancer, non-cancer, primary tissues, and cell-line samples included in our study based on our

145    metadata review. For each category, samples were assigned a "Yes" if they belonged to that

146    category (e.g. cell line samples for 'Cell Line' category), a "No" if they did not belong (e.g.

147    primary tissue samples for 'Cell Line' category), or an "Unknown" if it was not clear from the

148    metadata.

149

150         QTL mapping requires paired genotype and molecular phenotype information for each

151    sample. In standard QTL studies, genotyping arrays or whole genome sequencing (WGS) are

152    used to ascertain sample genotype information[31]. Unfortunately, for most of the ATAC-seq

153    data in public repositories that has already been collected, genotype data is not readily

154    available. However, ATAC-Seq directly captures genomic DNA fragments from accessible

155    chromatin regions; thus, we surmised that it might instead be possible extract genotype

156    information for these samples directly from the ATAC-seq reads. To obtain genotyping from

157    ATAC-sequencing and evaluate the performance of variant calling using ATAC-seq reads, we

158    applied two approaches: a pipeline incorporating genotyping from Gencove, which optimizes

159    genotyping and imputation for low-pass sequencing data[23,24,32,33], and a standard GATK

160    variant calling pipeline[32,33](Methods). To benchmark the performance of our workflow, we

161    used a published dataset of 71 HapMap lymphoblastoid cell lines (LCL) samples with paired

162    ATAC-seq and WGS data [34]. We observed that, compared to the standard GATK variant

163    calling pipeline, the Gencove pipeline with imputation greatly increased the number of variants

164    called and resulted in a median correlation of over 0.88 between true and called donor genotype

165    (Figure 2A). To quantify the effects of read coverage on the performance of variant calling, we

166    randomly subselected ATAC-seq reads at varying total read counts for use with the Gencove

167    pipeline. We observed a marginal increase in accuracy with deeper coverage, however, variant-

168    calling accuracy remained high at effective coverage as low as 0.04 (Figure 2B). In our full

169    dataset, the distribution of effective coverage in the full sample set was within the range

170    previously tested with the gold standard HapMap LCL samples, verifying the accuracy of

171    genotype calling in this larger data set. These analyses demonstrate the capabilities of accurate

172    inference of genome-wide genotypes directly from ATAC-seq data.

173        As a proof of concept, we next performed caQTL mapping using genotypes called from

174    ATAC-seq reads, comparing the results to the caQTLs identified using the full set of gold

175    standard genotypes in these 71 HapMap LCL samples. We observed that caQTL calling using

176    ATAC-seq reads and the Gencove pipeline performed better than the GATK pipeline with 99%

177    accuracy and over 90% recall compared to caQTL calling using WGS data. The increased recall

178    is due to the Gencove pipeline's imputation step and sacrifices very little in accuracy (Figure

179    2C). The performance of the Gencove pipeline had substantially greater benefit when testing

180    variants in larger caQTL mapping window sizes where recall remained above 90% for the

181    Gencove pipeline but dropped to 16% for the GATK pipeline at 100 kb (Figure 2C). Overall, we

182    conclude that genotype calling from ATAC-seq reads leads to highly accurate caQTL calling

183    with relatively high recall with a low rate of false positives. Given the diverse samples collected

184    and varying study designs, an individual donor will likely have multiple ATAC-seq samples

185    represented. As such, we next developed a pipeline to infer unique donors based on the

186    correlation between inferred sample genotypes across different samples and projects (Figure

187    2D-E, Methods). Applying this pipeline to all samples, we identified 1,454 unique donors across

188    our entire dataset (Supplementary Table 2). The majority of donors (~82%) are found within a

189    single project only. As expected, the occurrence of multiple samples per donor was especially

190    common amongst cell lines, which is reflected in the reduced proportion of cell line samples in

191    the final unique donor sample set (Supplementary Figure 1).

**Figure 2. High quality genotyping with unique donor information is inferable directly from reads obtained by ATAC-Seq. (A)** Variants called for the HapMap samples using two pipelines - Gencove, and GATK HaplotypeCaller. **(B)** Accuracy of variant genotype called by Gencove pipeline using a random subset of sample reads. Spearman correlation and mean squared error (MSE) are computed between the called genotype and genotype from WGS. **(C)** caQTLs called using ATAC-seq derived genotypes across the HapMap samples. **(D)** Spearman correlation of called genotypes between all samples. **(E)** Spearman correlation of called genotypes between samples in study PRJNA388006. On the top the "True donor" indicates the donor assignment obtained from metadata information for this study, and "Assigned donor" indicates the donor

10

201    assignment derived from called genotypes **(Methods)**.

202

203    **Peak calling across all samples identifies a plethora of open chromatin regions**

204    **with regulatory potential**

205    The next step in our pipeline was to identify open chromatin regions. We called

206    chromatin accessibility peaks based on evidence across all samples using Genrich, a peak

207    caller optimized for ATAC-seq reads[27]. Genrich assigns p-values to genomic positions within

208    each sample, then combines p-values across samples using Fisher's method to call peaks. We

209    compared this Genrich pipeline to strategies which called peaks in individual samples followed

210    by merging. The Genrich strategy alone produced peaks that are likely derived from

211    nucleosome-free and mono-nucleosome fragments, as seen by enrichment around 100 bp and

212    200 bp in the observed peak length distribution (Figure 3A).

213    Across 10,293 samples, we identified 1,659,379 autosomal peaks with a median peak

214    length of 250 base pairs, covering approximately 27% of the genome (Figure 3A). Chromatin

215    accessibility is influenced by a variety of regulatory processes[35–37], and we would expect to

216    see chromatin accessibility peaks in regions associated with gene regulation. To verify the

217    quality of our ATAC-seq peaks, we annotated our peaks, along with length-matched, randomly

218    selected controls, with various genomic features that included transcript annotations and

219    enhancer annotations as defined by the FANTOM5 enhancer atlas[38,39] (Methods). We found

220    that relative to controls, our ATAC-seq peaks were enriched for genomic regions annotated as

221    enhancers and all transcript annotations but depleted for gene intergenic regions

222    (Supplementary Figure 2, Supplementary Table 3). Similarly, we would expect our ATAC-seq

223    peaks to be enriched for histone modifications associated with gene regulatory regions[40–42].

224    The ENCODE Roadmap Epigenomics Mapping Consortium[43] provides chromatin

225    immunoprecipitation with sequencing (ChIP-seq) data representing eight different histone marks

226    from 556 cell line, tissue, and primary cell samples derived from a variety of biological origins.

227    Using these data, the highest enrichment of our ATAC-seq peaks and chromatin histone marks

228    was for H3K4me1, a histone mark that has been linked to enhancers (Supplementary Table

229    4)[40]. In contrast, our ATAC-seq peaks were depleted for overlap with the histone mark

230    H3K9me3, which is associated with gene repression and heterochromatin[44]. Together, these

231    data suggest that our ATAC-seq peaks are enriched for cis-regulatory regions, as expected for

232    genomic sequences implicated in regulatory activity and indicating high quality peak calls.

233

234    **Inferred genotypes support high-powered caQTL mapping across samples**

235          Next, we sought to identify genetic variants that are associated with differences in

236    measured chromatin accessibility in ATAC-seq peaks, i.e., caQTLs. We tested a 10 kilobase

237    (kb) window in *cis* flanking each chromatin accessibility peak, as we anticipate that genetically

238    altered active transcription factor binding sites are likely to be found within or very nearby

239    regions of chromatin accessibility[45,46]. Utilizing our peak calling and genotyping pipelines, we

240    identified 23,381 chromatin accessibility peaks with a significant caQTL at FDR 5% across

241    1,454 unique donor samples (Figure 3B, Methods, Supplementary Tables 5-6). To mitigate

242    potential confounding from population stratification, we estimated variation in similarity across

243    donors generated by our genotyping via principal components analysis (PCA), including 3 PCs

244    as covariates in discovery analysis. In addition, we also included 200 PCs generated from the

245    donor chromatin accessibility peak read count matrix to mitigate potential latent confounders in

246    QTL mapping [47] (Methods).

247          We examined the quality of our caQTL variants by determining whether they were

248    enriched for expected functional characteristics. First, we confirmed that the distribution of

249    positions for lead caQTL variants was centered within the open chromatin peak tested, as

250    expected (Figure 3C). In addition, we observed that peaks with a mapped caQTL were the most

251    strongly enriched for gene 5' UTRs and enhancer regions while depleted in gene intergenic

252    regions (Supplementary Figure 3, Supplementary Table 7). Interestingly, caQTL peaks were

253    further enriched in enhancer regions compared to all chromatin accessibility peaks, suggesting

254    that caQTLs we mapped may be found at genomic elements involved in distal gene regulation.

255    This could potentially arise due to selective pressure reducing functional variation in promoters

256    and other proximal elements.

257         Additionally, we examined whether our caQTL peaks were enriched for transcription

258    factor binding sites in the ENCODE transcription factor ChIP-seq data from 129 cell types and

259    340 transcription factors[48]. As expected, caQTL peaks, compared to length-matched random

260    controls, were enriched for binding sites for all transcription factors except for SRSF9, which is

261    depleted in caQTL peaks (Supplementary Table 8). Enrichment of these functional

262    characteristics support the conclusion that our caQTLs are high quality, reflect enrichment in

263    expected regulatory elements, and can help identify genetic mechanisms relevant to regulation

264    of gene expression. We sought further evidence that caQTL variants were enriched for

265    functional roles in gene expression regulation by intersecting them with eQTLs. Across all 49

266    Genotype-Tissue Expression (GTEx) v8 tissues, we observed caQTL/eQTL enrichments

267    ranging from 2.1 to 4.8-fold per tissue and a total of 2,859 (~13% of unique caQTLs) unique

268    overlapping lead caQTL/lead eQTL variants found across all tissues, for an enrichment of

269    approximately 1.8-fold (Supplementary Table 9).

270         Finally, to further demonstrate that our catalog represents reproducible peaks and

271    caQTLs, we compared our findings here to a recent caQTL study that identified variants

272    associated with chromatin accessibility in African LCL samples[49] not included in our discovery

273    effort. Lead caQTLs and peaks identified in our study resulted in a replication rate ($\pi_1$

274    value[50,51]) of 0.62 with this orthogonal study (Figure 3D). Together, these analyses further

275    demonstrate that on average, our catalog of caQTLs are high quality and provide insight into

276    how genetic variation may affect gene regulation and complex traits.
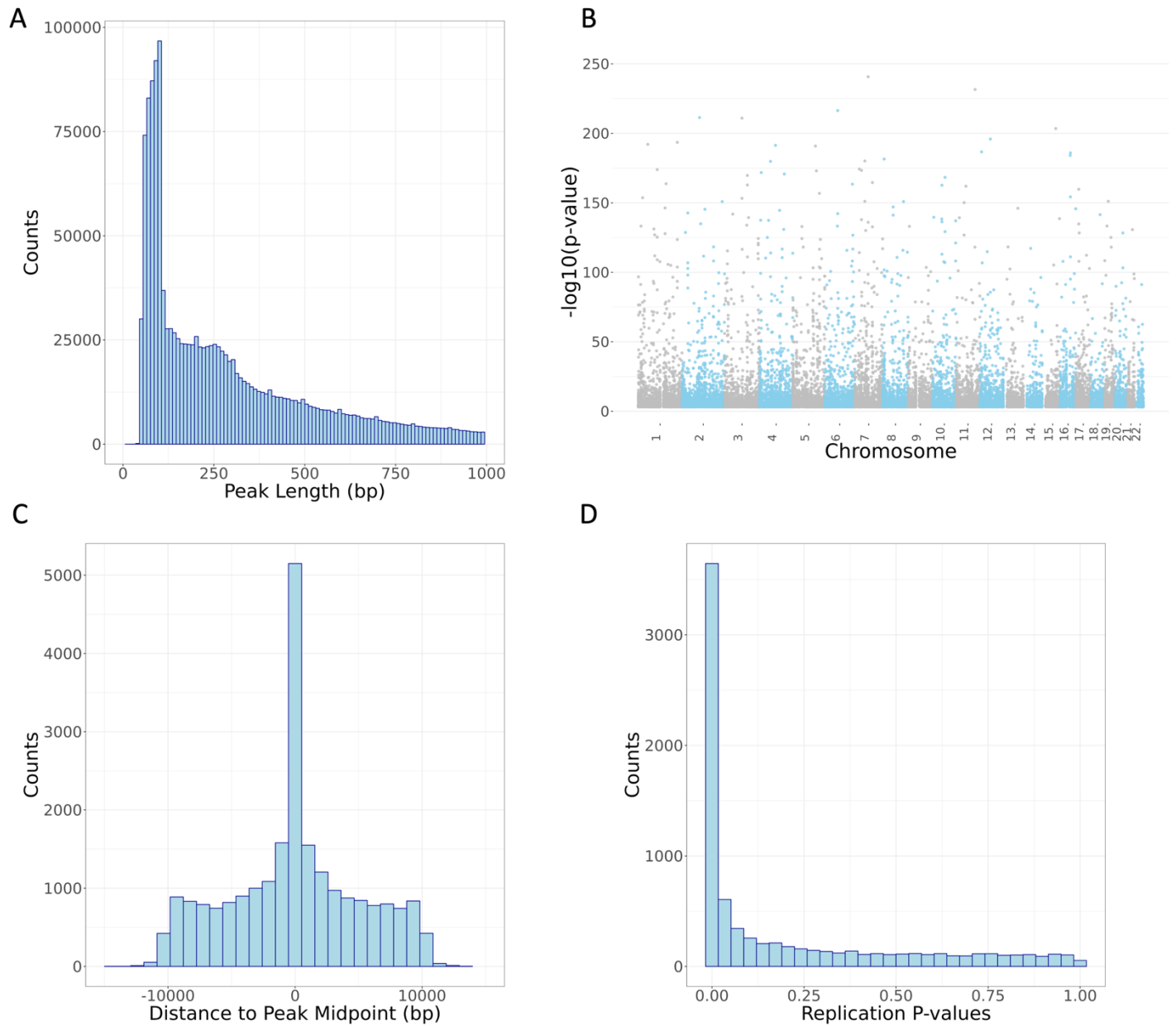
277

278

279

280

281

282



**Figure 3. Characteristics of chromatin accessibility peaks and caQTL variants identified**

14

284     **in this study. (A)** Distribution of peak length across 1,659,379 called peaks (peaks under 1000

285     bp shown). **(B)** Manhattan plot of lead variant for 23,381 caQTL peaks. **(C)** Distance from lead

286     caQTL variant to midpoint of caQTL peak showing elevation of caQTL variant within the

287     identified chromatin accessibility peak. **(D)** Lead variants for 23,381 caQTL peaks were matched

288     in external caQTL mapping dataset of African LCLs[49]; p-values from the replication study are

289     plotted here.

290

291

292     **Colocalization suggests shared causality between chromatin accessibility,**

293     **complex traits, and expression QTLs**

294            To gain further insight into the molecular mechanisms underlying GWAS signals, we

295     sought to link GWAS association signals, expression QTLs (eQTLs), and caQTLs together via

296     statistical colocalization (**Methods**). Colocalization analysis discerns if an association signal is

297     likely shared between two traits, suggestive of a common underlying genetic mechanism. First,

298     we examined which caQTL signals are shared with GWAS signals across a variety of complex

299     human traits. We obtained GWAS summary statistics from a subset of the UK Biobank (UKBB)

300     study, selecting 78 traits of interest with high confidence of significant heritability (**Methods**)[52].

301     We then performed colocalization analysis (**Methods**) for any caQTL peak that was located

302     within 1 Mb of a genome-wide significant lead GWAS signal (**Methods**). We observed that 67

303     traits had a caQTL/GWAS colocalization event (PP3+PP4 > 0.8 and PP4/(PP3+PP4) > 0.9.) for

304     a total of 12,882 colocalization events across all traits, involving 4,351 (~19%) unique caQTL

305     peaks and 4,706 (~34%) unique tested GWAS signals (Supplementary Table 10).

306            Regulatory variants do not always affect the nearest gene and assigning a GWAS signal

307     to a causal gene is not a trivial procedure[53,54]. Furthermore, comparison of the overlap

308     between lead variants of GWAS signals and the lead variant of eQTLs can suggest the incorrect

15

309    causal gene[55]. Given the prominence of long-range gene expression regulation, colocalization

310    of cis regulatory elements with eGenes can suggest a shared causal variant[56,57]. We

311    performed colocalization analyses between caQTLs and 49 GTEx v8 eQTL tissues. Across all

312    tissues, between 358 (Kidney) and 5,427 (Thyroid) eGenes colocalized with our caQTLs.

313    Colocalized caQTLs/eQTLs were shared across a median of three tissues and 17,471 unique

314    eGenes colocalized with caQTLs in any GTEx tissue (Supplementary Figure 4, Supplementary

315    Table 11). We found that only 13% of eQTL/caQTL colocalizations involve the gene nearest to

316    the lead caQTL and that there was a median of 6 genes closer to the lead caQTL than the

317    colocalizing gene (Supplementary Figure 4). Additionally, the putative regulated gene

318    transcription start site (TSS) was a median of 80,798 base pairs away from the colocalizing

319    caQTL (Supplementary Figure 4). These results suggest that caQTLs may often be found

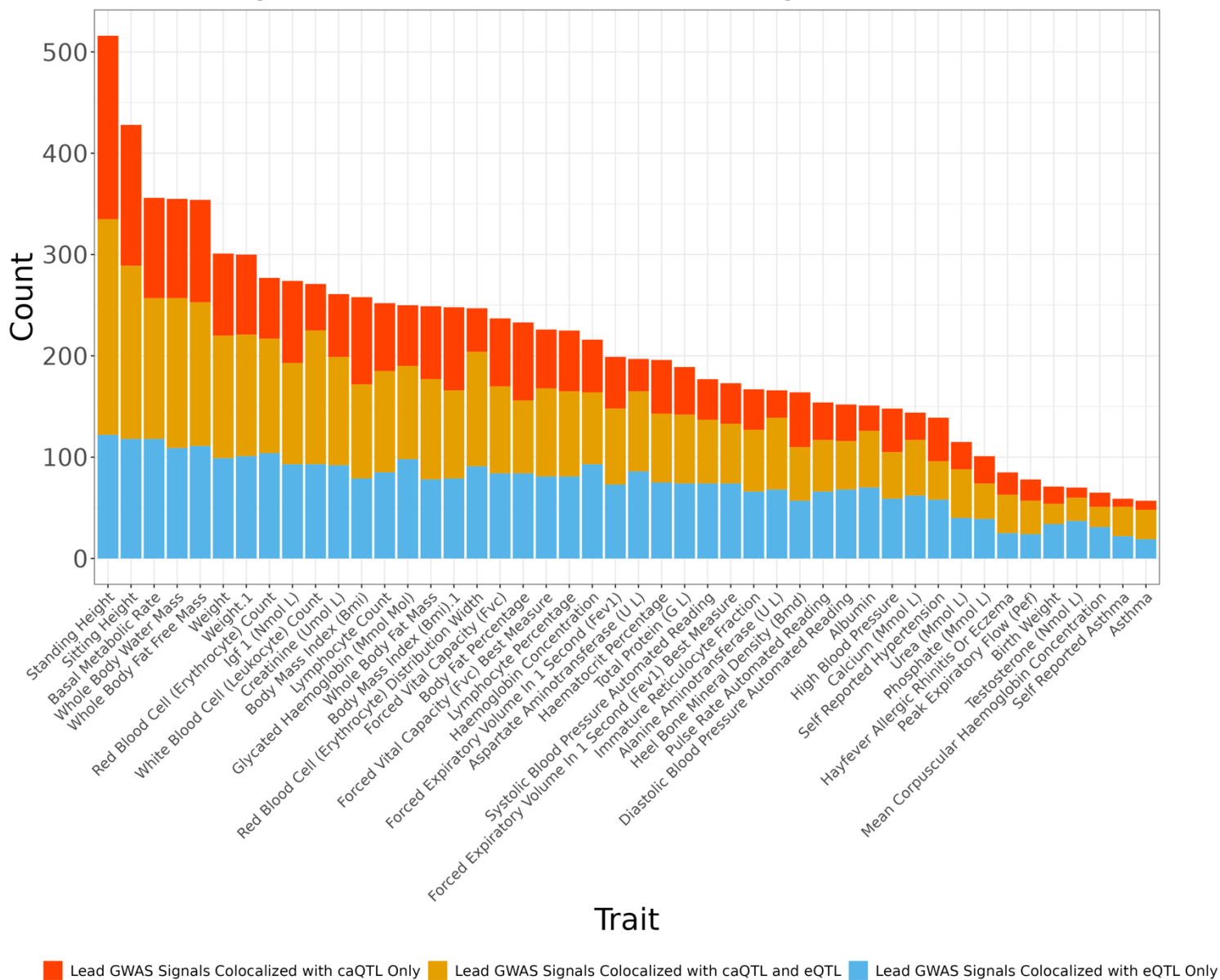320    tagging and potentially modifying the behavior of distal gene regulatory elements.

321

322    **Multiple molecular QTL datasets provide insight into regulatory mechanisms underlying**

323    **GWAS associations**

324            eQTLs have been shown to provide a regulatory mechanistic hypothesis for GWAS

325    associated signals, yet only an estimated ~25-43% of GWAS signals colocalize with known

326    eQTLs[6,58], implying that more than half of GWAS loci may lack an obvious functional,

327    mechanistic hypothesis[6,59–61]. caQTL mapping could help close that gap if, for example, the

328    effects of the eQTL are only apparent in certain cellular contexts, during specific developmental

329    stages, or in the presence of external stimuli[62–64], whereas chromatin accessibility may be

330    primed and reveal effects in a wider range of context. Across all traits and GTEx tissues, we find

331    that lead GWAS signals colocalize with a median of 5 eQTLs and 2 caQTLs (Supplementary

332    Figure 5). For each GWAS trait, we then considered whether independent GWAS lead signals

333    colocalize only with eQTLs, colocalize with both caQTLs and eQTLs, or colocalize only with

334    caQTLs. Across all GWAS, a median of 35 unique signals colocalized with a caQTL only, a

16

335    median of 66 unique signals colocalized with an eQTL only, and a median of 53 unique signals

336    colocalized with both a caQTL and an eQTL (Figure 4, Supplementary Table 13). These

337    differences may reflect context-specific behavior of gene regulation that is not well captured by

338    steady-state, adult gene expression data, but may still be reflected in chromatin accessibility.

339    These results demonstrate that incorporating both caQTLs and eQTLs nominates putative

340    causal mechanisms for approximately 29% more GWAS signals than using eQTLs alone.

341    Furthermore, 57% of GWAS signals we tested were linked with either a caQTL, eQTL, or both

342    (Supplementary Figure 6). Instances where GWAS signals colocalized with both caQTLs and

343    eQTLs may also allow for a better delineation of the mechanism at these loci by nominating a

344    candidate caQTL-associated gene regulatory element to a target eGene[65].

**Figure 4. caQTLs map to regions tagged by GWAS and eQTL variation.** For each GWAS trait, independent lead GWAS variant signals were checked for colocalization with caQTL and eQTL signals across all GTEx tissues. Plotted is the number of unique lead GWAS signals per colocalization group, as multiple caQTL peaks, eGenes, etc. can colocalize with the same lead GWAS signal. Traits with greater than 50 colocalizing lead variants shown.

352     To gain insight into molecular mechanisms that may be unique to caQTLs as compared

353     to eQTLs, we calculated the enrichment of colocalizing caQTLs and lead eQTLs for diverse

354     genomic annotations. caQTLs and eQTLs involved in colocalizations with GWAS signals were

355     both significantly enriched for all tested genomic annotation categories except for intergenic

356     regions, where they were significantly depleted, compared to matched random controls

357     (Supplementary Figure 7, **Methods**). However, caQTLs from GWAS/caQTL and

358     caQTL/GWAS/eQTL colocalization events were further enriched for enhancer regions and less

359     depleted in intergenic regions than eQTLs from GWAS/eQTL colocalizations alone

360     (Supplementary Figures 8-9). In contrast, lead variants of eQTLs that colocalized with a GWAS

361     signal only were further enriched for gene promoters and other gene proximal categories, less

362     enriched in enhancer regions, and showed greater depletion for intergenic regions, consistent

363     with previous reports (Supplementary Figure 10)[6,66]. These differences in enrichment may be

364     due to systematic differences in GWAS signals that are explained by eQTLs compared to those

365     explained by potentially distal regulatory mechanisms captured by caQTLs[67].

366     While our caQTLs were called from heterogeneous cell/tissue samples, they are

367     predominantly from brain and whole blood (Figure 1). To reflect this, we also performed an

368     analysis of caQTL/GWAS colocalizations compared to eQTL/GWAS colocalizations from brain

369     cortex and whole blood only. Across 69 GWAS, each trait has at least 1 GWAS signal that

370     colocalizes only with a caQTL, and one trait, standing height, had 360 lead GWAS variants that

371     colocalize exclusively with caQTLs compared to brain eQTLs. In contrast, we identify a

372     maximum of 66 lead GWAS variants that colocalize only with eQTLs for a given trait. Across all

373     GWAS, a median of 76 unique signals colocalized with a caQTL only, a median of 15 unique

374     signals colocalized with an eQTL only in Whole Blood, and a median of 11 unique signals

375     colocalized with both a caQTL and a Whole Blood eQTL (Supplementary Figure 11,

376     Supplementary Table 14). Furthermore, across all GWAS, a median of 83 unique signals

377     colocalized with a caQTL only, a median of 10 unique signals colocalized with an eQTL only in

378    Brain Cortex, and a median of 7 unique signals colocalized with both a caQTL and a Brain

379    Cortex eQTL (Supplementary Figure 12, Supplementary Table 15). Compared to the analysis

380    considering eQTLs across all tissues, we find that caQTL/GWAS only colocalizations occur with

381    a larger proportion of GWAS signals in single tissue eQTL analysis colocalizations. This

382    discrepancy provides further evidence that using caQTLs can provide molecular insight into

383    GWAS association signals beyond eQTLs when restricting to a single eQTL tissue.

384

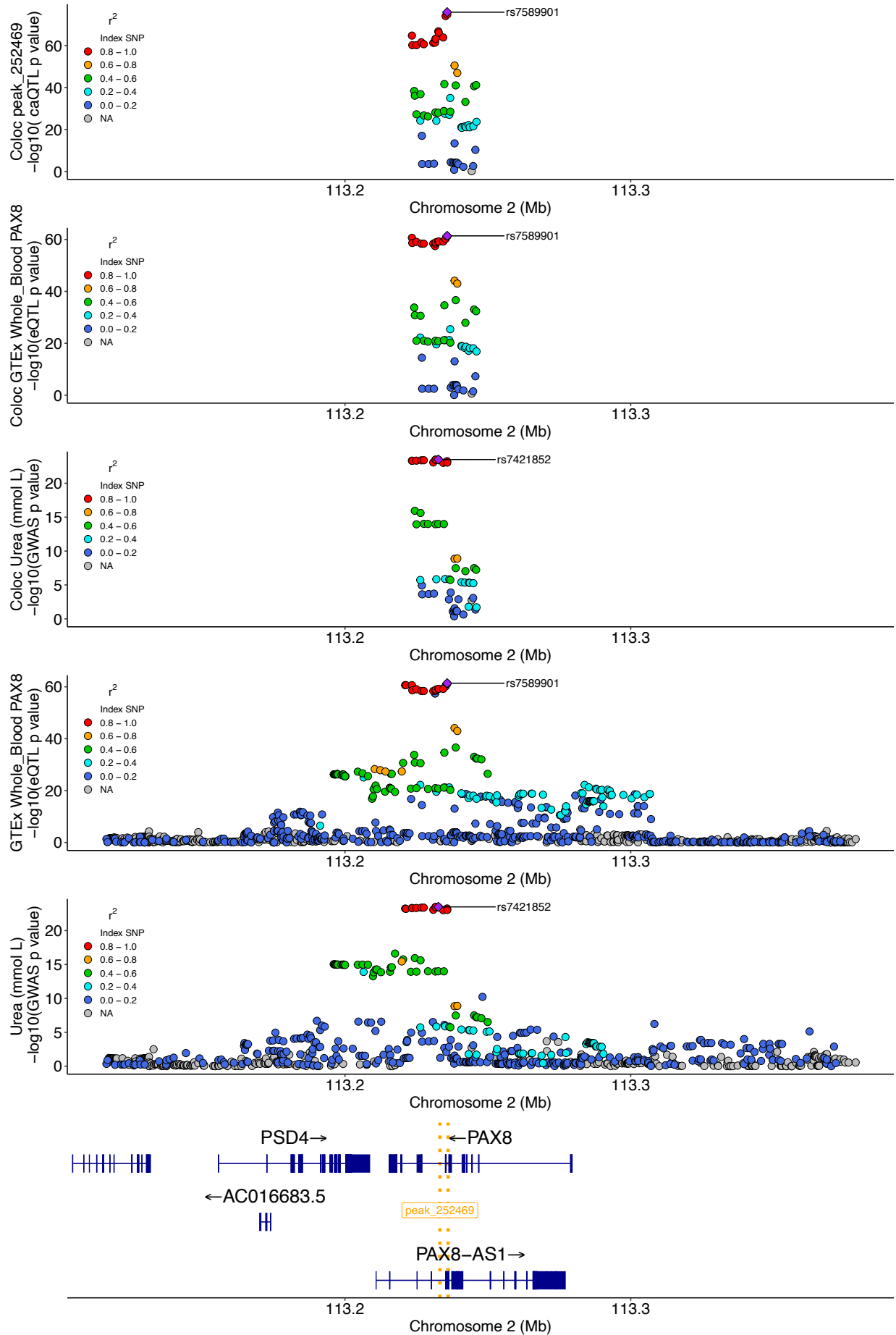385    **Integration of caQTLs informs mechanistic interpretation at many GWAS loci**

386    Colocalization analysis with QTL datasets across multiple modalities, such as

387    expression and chromatin accessibility, has previously been shown to nominate putative target

388    genes underlying more GWAS signals than a single modality alone[65,68]signals that

389    colocalized separately with both caQTLs and eQTLs and quantified how many of the GWAS-

390    colocalizing caQTLs and eQTLs also colocalized with each other. We identified 43,005 unique

391    colocalization events involving a GWAS trait, caQTL peak, eGene, and eGene tissue

392    (Supplementary Table 16). These were comprised of 2,177 unique eGenes and 1,695 unique

393    caQTL peaks.

394    In cases where caQTLs colocalize with both GWAS signals and eQTLs, they provide a

395    more complete picture of the mechanisms likely driving the association signal. First, we provide

396    an instructive example of a well-characterized GWAS locus strongly associated with plasma

397    low-density lipoprotein cholesterol (LDL-C) at the 1p13 locus. eQTL colocalization analyses at

398    this locus, followed by functional characterization in vitro and in vivo, suggest that the causal

399    gene at this locus is *SORT1*, with expression differences observed in the liver [46]. We find a

400    caQTL at this locus that colocalizes with both the *SORT1* eQTL in liver, and the GWAS trait self-

401    reported high cholesterol (Supplementary Figure 13). This caQTL peak contains a well-studied

402    noncoding variant that creates a C/EBP (CCAAT/enhancer binding protein) TF binding site,

403    altering hepatic expression of *SORT1* and plasma LDL-C levels[46]. This highlights the ability of

20

404    our analyses to identify verified mechanisms underlying GWAS signals.

405         In a second example, we identified a compelling locus where a caQTL peak, a whole

406    blood eQTL for *PAX8*, and a GWAS signal for blood urea levels colocalized (Figure 5). The

407    shared lead caQTL and eQTL variant, rs7589901, is an intronic variant within the *PAX8* gene.

408    The reference allele of rs7589901-A is associated with increased chromatin accessibility in the

409    associated peak (Supplementary Figure 14). Based on motif analysis, ZNF135 is predicted to

410    bind to a motif overlapping rs7589901, with the alternate C allele strongly favored for binding

411    (PWM value=0.8, Supplementary Figure 15). In GTEx, the rs7589901 eQTL direction of effect is

412    concordant with the caQTL direction of effect, suggesting that increased accessibility at this

413    locus is associated with increased *PAX8* gene expression in whole blood. The lead GWAS

414    variant at this locus, rs7421852, is associated with increased blood urea levels, is ~3,000 bp

415    from rs7589901, and is in strong LD ($r^2$=0.85) with rs7589901 in our caQTL sample genotypes.

416    These results suggest a potential mechanism where ZNF135 is acting as a transcriptional

417    repressor at this locus, a functional role that has been implicated in a different context[69]. The

418    culmination of evidence suggests a mechanism where decreased ZNF135 binding leads to

419    increased chromatin accessibility, increased expression of the *PAX8* gene, and lower blood

420    urea levels. Such examples demonstrate the power of integrating multiple molecular QTL

421    datasets to nominate mechanistic hypotheses that may be further validated experimentally.

422

423 **Figure 5**. **Change in chromatin accessibility and expression implicate *PAX8* in serum**

424 **urea levels.** The top three plots are the colocalization windows (10kb + caQTL peak) for the

425 caQTL, eQTL, and GWAS, respectively. The following two plots are showing a larger window to

426 illustrate the eQTL and GWAS signals, respectively, at this locus at a different scale. The

427 bottom gene track highlights the position of genes at this locus, as well as the location of the

428 caQTL peak (gold dotted lines).

429

430

431 **Sample heterogeneity enables identification of context-specific clusters**

432 Because profiles of chromatin accessibility often segregate context or cell-type specific

433 information, we next grouped our samples by their profiles of chromatin[70]. We performed

434 dimensionality reduction[71] and applied a semi-supervised clustering method[71] to identify

435 groups of similar samples, identifying 11 clusters (Figure 6A). We used sample metadata to

436 assign a label to each cluster, denoting biological origin. Overall, clustering appears to be

437 mainly driven by the tissue or cell type from which the sample is derived (Supplementary

438 Figures 16-17). For example, blood cell types appear to be grouped together or near each other

439 in separate, but related clusters. In addition, we found other examples of clusters where nearly

440 half of the samples are derived from a single tissue, such as pancreas. Annotating samples with

441 other aspects of metadata, such as primary sample vs. cell line, or cancer vs. non-cancer

442 samples, did not appear to explain clustering results (Supplementary Figure 18).

443

444 **Clustering allows for identification of caQTLs in specific clusters**

445 To determine whether clustering samples of similar biological origin enables the

446 discovery of additional caQTL signals, we next performed caQTL mapping within each cluster.

447 Each cluster is composed of a different number of samples, with varying contributions from cell

448    types and projects, which is reflected in the number of caQTLs identified in each cluster. Cluster

449    sample size ranged from 80-220 samples (Supplementary Table 17) and resulted in 174-15,277

450    (FDR<5%) caQTLs identified in a single cluster. As in the global analysis, cluster caQTLs

451    showed similar patterns of genomic region annotation enrichments (Supplementary Figure 19)

452    and lead caQTLs were centered within the open chromatin peak tested (Supplementary Figure

453    20). Across all clusters, cluster caQTLs rediscovered 34-94% of caQTL peaks observed in the

454    global analysis (Figure 6B) with median global caQTL replication rate of 0.99 ($\pi_1$ value) across

455    all clusters (Supplementary Figure 21). Analysis comparing cluster caQTL peak discoveries to

456    other clusters resulted in a range of caQTL peak rediscovery (Supplementary Figure 22) but

457    high replication rate across clusters ($\pi_1$ value 0.91-0.99) (Figure 6C, Supplementary Table 18).

458    This suggests that clusters are capturing common global signals, but some clusters are better

459    powered at identifying caQTLs that might be cell/tissue-specific. For example, cluster 9, which

460    identified the largest number of cluster caQTLs, is comprised of more than 50% LCL samples,

461    many of which are from a single study (Supplemental Table). Approximately 2/3 of the caQTL

462    peaks identified in cluster 9 are not identified as caQTL peaks in the global analysis performed

463    across all tissues/cell types, suggesting that cluster 9 may be better powered to discover

464    caQTLs more prevalent in LCLs and related blood cell samples. As a measure of reproducibility

465    across experiments, we found that Cluster 9 caQTL lead variants were enriched for evidence of

466    caQTL peak causality in the original study[34] that the majority of cluster 9 samples originate

467    from (Supplementary Figure 23). These results suggest that as with eQTLs, future work

468    increasing the sample size to examine cell/tissue-specific caQTLs is likely to capture novel

469    caQTLs that will be useful for elucidating molecular mechanisms underlying GWAS signals.
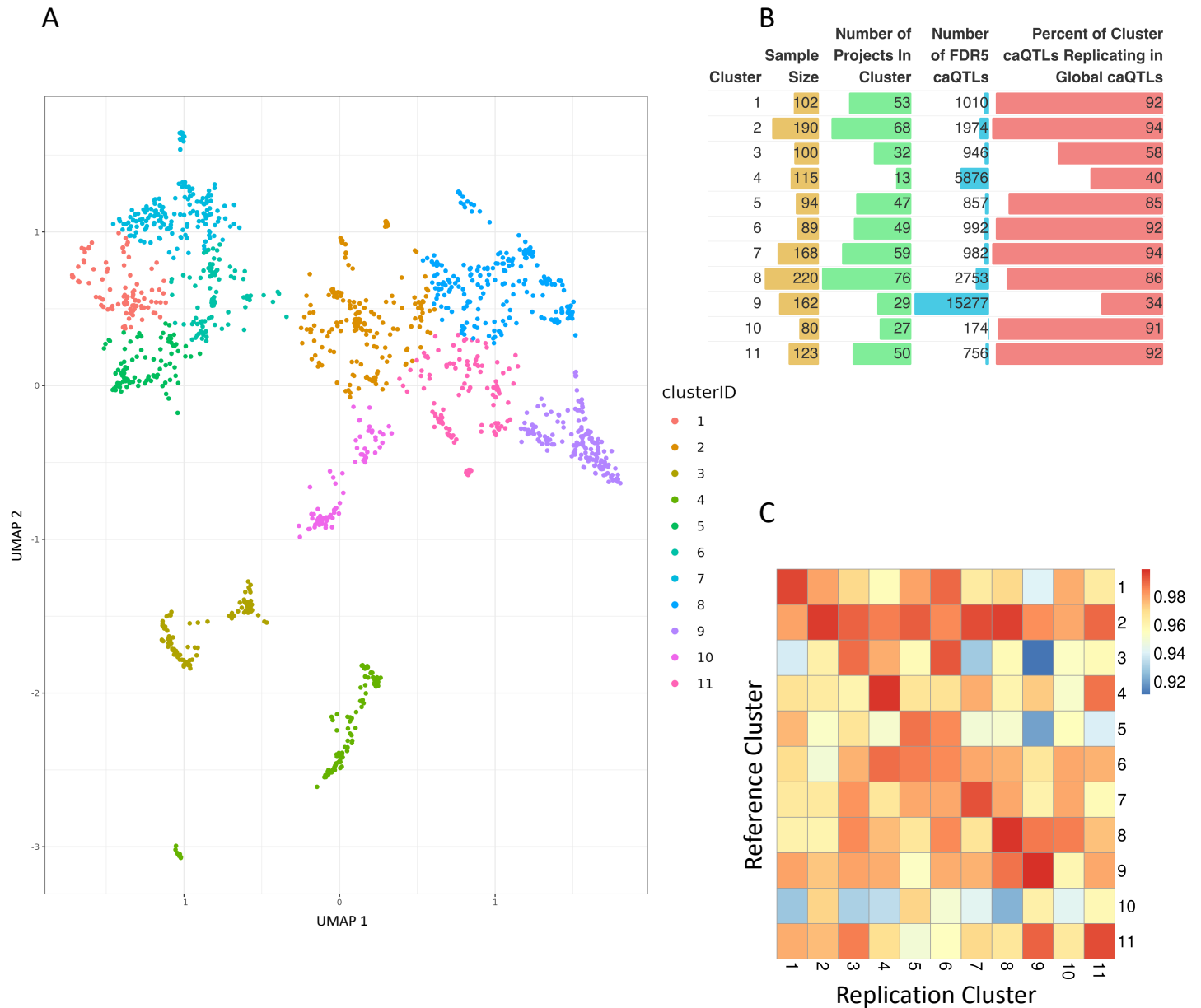
470            Mapping caQTLs in clusters highlights the increase in caQTL discovery power of

471    aggregating all samples across experiments, particularly for caQTLs that might be found across

472    cell types. In our global analysis we identified 23,381 caQTL peaks, with a maximum of 5,169 of

24

473    those also identified in a single cluster caQTL mapping experiment. This suggests that by

474    considering all samples, we achieve greater than a 4.5X increase in caQTL discovery power for

475    global caQTLs. Across all clusters, we identify 8,610 (37% of global) caQTL peaks that were

476    also found in the global analysis and 14,795 caQTL peaks that were not found in the global

477    analysis.

478    **Figure 6. Clustering and discovery of cluster caQTLs across ATAC-Seq samples.**



| Cluster | Sample Size | Number of Projects In Cluster | Number of FDR5 caQTLs | Percent of Cluster caQTLs Replicating in Global caQTLs |
|---|---|---|---|---|
| 1 | 102 | 53 | 1010 | 92 |
| 2 | 190 | 68 | 1974 | 94 |
| 3 | 100 | 32 | 946 | 58 |
| 4 | 115 | 13 | 5876 | 40 |
| 5 | 94 | 47 | 857 | 85 |
| 6 | 89 | 49 | 992 | 92 |
| 7 | 168 | 59 | 982 | 94 |
| 8 | 220 | 76 | 2753 | 86 |
| 9 | 162 | 29 | 15277 | 34 |
| 10 | 80 | 27 | 174 | 91 |
| 11 | 123 | 50 | 756 | 92 |

479    **(A)** UMAP followed by k-means clustering to identify groups of related samples based

480    on chromatin accessibility profiles across all peaks. **(B)** Cluster characteristics, caQTLs

26

481      identified, and replication with respect to global caQTL mapping. **(C)** Replication rate ($\pi_1$

482      value) of caQTLs identified in each cluster compared to those found in all other clusters.

483

484

485      **Cluster-specific caQTLs can explain additional gene regulation and GWAS signal**

486      **causality**

487      We next performed colocalization analysis between GTEx eQTLs and the caQTLs

488      identified within each cluster to determine if cluster-specific caQTLs appear to be involved in

489      gene regulation as well. As in the cluster caQTL analysis, we find that the number of

490      colocalizations found per cluster was commensurate with the number of caQTLs identified in

491      each cluster. We find a maximum of 13,688 unique eGenes colocalizing in a single cluster, and

492      a total of 16,833 unique eGenes colocalize when considering all clusters (Supplementary

493      Tables 19-20). Compared to the global analysis, which identified a total of 17,471 unique

494      colocalizing eGenes, 14,017 of which also colocalized in the cluster analyses, suggesting that

495      the majority of colocalizing eGenes are identified across both analyses. As in the cluster caQTL

496      analyses, we find that colocalizing eGenes are often shared across clusters (Supplementary

497      Figure 24). Considering all cluster colocalization events, 7,789 total eGenes were found to

498      uniquely colocalize in a single cluster, with 5,532 (71%) of these in cluster 9. Overall, we find a

499      variable number of cluster-specific caQTL/eQTL colocalizations per cluster, many of which are

500      shared across clusters.

501      Our previous analyses assessed the benefit of utilizing global caQTLs in GWAS

502      colocalizations compared to eQTLs. In this analysis, we considered eQTLs that were discovered

503      in experiments performed in single tissues, experiments that are much more likely to identify

504      variants with tissue-specific effects compared to our multi-tissue, global caQTL mapping

505      strategy. Cluster-specific caQTLs might more closely mimic these single-tissue eQTL datasets,

506    as these caQTLs were mapped in clusters of samples that likely shared a similar biological

507    origin. To better compare the contribution of eQTLs and caQTLs to GWAS signals, we

508    considered caQTLs identified in both global and cluster-specific analyses to assess

509    colocalization improvement. Across all GWAS traits and eQTL tissues tested, we find that

510    combining global and cluster-specific caQTLs results in an increase of the contribution of

511    caQTLs to GWAS colocalizations. Specifically, we find a median of 41 GWAS signals

512    colocalizing with caQTLs only and a median of 67.5 GWAS signals colocalizing with both

513    caQTLs and eQTLs (Supplementary Figure 25, Supplementary Table 21). Both measurements

514    are increases compared to the global analysis only. In contrast, the median number of GWAS

515    signals that colocalize with eQTLs only decreased to 39 (Supplementary Figure 25,

516    Supplementary Table 21). Leveraging both global and cluster caQTLs, together with eQTLs, we

517    explained a median of 62% of GWAS signals tested (Supplementary Figure 26). Overall, we find

518    that both global and cluster-specific caQTLs can contribute to the causal mechanisms

519    underlying GWAS signals not captured by eQTLs.

520    **Discussion**

521        We developed a pipeline to discover caQTLs on a large scale by aggregating and

522    genotyping large-scale ATAC-seq data across many studies. We collected 10,293 human

523    ATAC-seq samples, representing 1,454 unique donors, from public databases that come from a

524    diversity of cell types and conditions, demonstrating that genotype data can be accurately called

525    from ATAC-seq data, and identified unique sample donors, both within and across projects.

526    Combining accessibility and genotype information, we performed caQTL analysis and were able

527    to capture global and cluster-specific caQTLs. caQTL studies are often limited by sample size

528    constraints. We show that amassing public-domain project data allows for identification of a

529    greater number of caQTLs than smaller individual studies alone. We demonstrated that caQTLs

530    are enriched for various regulatory elements and likely underlie gene expression differences

531   and complex human traits. We provide our large catalog of global and cluster caQTLs as a

532   resource.

533   Our study does have limitations and opportunities for further development. Naturally, as

534   more ATAC-seq data are generated, a similar study could be repeated on a larger scale.

535   Additionally, the clustering performed in our study was coarse, and may have grouped multiple

536   cell types or contexts together. With a larger sample size from new studies or more extensive

537   exploration of clustering methods or cell type prediction approaches, these grouping could be

538   further refined and made more homogeneous, which would be expected to boost statistical

539   power for discovery. Although we analyzed a large and diverse set of samples and experiments,

540   many GWAS signals were not tagged by one of our caQTLs (and/or by eQTLs). One

541   explanation for this is that we are missing many cluster/context-specific caQTLs that may

542   underlie the remaining GWAS signals. One limitation of this study is that while the sample

543   contexts were diverse, we still do not have sufficient sample size across some disease-relevant

544   contexts to fully examine context-specific caQTLs. Further work, perhaps using single cell

545   ATAC-seq data, is necessary to gain insight into tissue/cell context specific caQTLs. Other

546   types of molecular QTLs may underlie some unexplained GWAS signals[60]. Incorporating

547   additional data modalities, such as those reflecting chromosome conformation changes, may

548   identify additional QTLs underlying GWAS loci. A recent study has shown that genetic variants

549   in enhancer regions affect gene expression changes via enhancer-promoter touching and

550   looping processes[72]. Integrating HiC or HiChIP datasets with ATAC-seq data can provide

551   insight into this process. These datasets may also help identify target genes or resolve

552   situations where multiple eGenes are implicated as causal genes at a locus[73]. Furthermore,

553   other mechanisms, such as DNA methylation (meQTLs)[74,75] or post-transcriptional processes

554   such as splicing (sQTLs)[66] or protein concentrations (pQTLs)[76] could underlie GWAS

555   signals that have yet to be explained.

556    Although we observed colocalization analysis between our caQTLs and GWAS signals

557    on par with previous studies [77], experimental validation is necessary to determine whether

558    putative causal variants underlying these QTLs directly mediate disease risk[78,79]. Previous

559    studies have shown that this type of analysis has led to the correct identification of molecular

560    mechanisms underlying disease. For example, regulatory mapping has successfully identified

561    gene targets that can be experimentally modulated to produce a phenotypic effect both in vitro

562    and in vivo[80]. Furthermore, caQTL analyses have been used to predict mechanisms

563    underlying GWAS signals with follow-up functional experiment results supporting these

564    predictions[15]. Ultimately, regulatory elements and gene targets that we identify as implicated

565    at GWAS loci will need additional support from low-throughput experimental techniques to

566    confirm our findings, such as using base editing to dissect variant function[81]. Toward the goal

567    of understanding molecular mechanisms underlying GWAS signals, molecular QTLs generate

568    hypotheses and our work has demonstrated that including caQTLs in these experiments

569    increases the number of GWAS signals for which a putative molecular mechanisms may be

570    identified.

571    **Conclusions**

572    In summary, we have deployed a pipeline to call a set of consensus peaks from thousands of

573    publicly available ATAC-seq samples and genotype these samples directly from the

574    experimental sequencing reads. We leveraged these data to identify caQTLs that likely share

575    causal variants with eQTLs and GWAS signals. We show that caQTLs can improve our

576    understanding of the mechanisms underlying GWAS signals and we provide this dataset as a

577    resource for use in further fine-mapping experiments.

578

579

580    **METHODS**

581    **Sample Collection**

582    ATAC-seq samples were identified through the Gene Expression Omnibus (GEO) database and

583    downloaded. Collected sample metadata is found in Supplementary Table 1.

584

585    **Benchmarking on HapMap samples**

586    We downloaded ATAC-seq for 71 HapMap samples from ENA project PRJEB28318[34]. We

587    aligned the sequencing reads to GRCh38 using bowtie2 and retained only autosomal

588    chromosomes. Duplicated reads tagged by Picard were removed and Base Quality Score

589    Recalibration (BQSR) was performed using GATK tools. Variant calling was done using GATK

590    HaplotypeCaller. Loci with less than 2 reads were filtered out and variants were mapped to

591    GRCh37 using Picard LiftoverVcf. Minimac4 was utilized to run imputation with reference panel

592    derived 1000G Phase 3

593    (https://csg.sph.umich.edu/abecasis/mach/download/1000G.Phase3.v5.html). We kept only the

594    genotype for common variants derived from 1000G with MAF > 0.05. The gold standard variants

595    were obtained from https://www.internationalgenome.org/data-portal/data-collection/grch38. For

596    the ATAC-seq data, we converted cram files to bam files, and removed the reads that map to

597    mitochondrial genome. We obtained the genotype from the 1000 Genome Project on the

598    GRCh38 genome assembly[82].

599

600    **Benchmarking for caQTLs in HapMap samples**

601    We first obtained caQTLs using ATAC-seq reads with BH corrected P-value < 0.05, then ran

602    QTL analysis using gold standard genotype and obtained caQTLs with BH corrected P-value <

603    0.05. The precision is computed as the percentage of replicated caQTLs at FDR < 0.05 using

604    the gold standard genotype. Similarly, we first obtained caQTLs using gold standard genotypes

605    with BH corrected P-value < 0.05, then ran QTL analysis using ATAC-seq reads and obtained

606    caQTLs with BH corrected P-value < 0.05. The recall is computed as the percentage of

607    replicated caQTLs at FDR < 0.05 using the ATAC-seq reads.

608

609    **Variant calling**

610    For the ATAC-seq data, we performed two pipelines of variant calling, one using GATK

611    HaplotypeCaller, and the other with Gencove's low-pass sequencing pipeline. Using the GATK

612    HaplotypeCaller, we performed alignment using Bowtie2, and removed duplicated reads and

613    applied base quality score recalibration, followed by GATK HaplotypeCaller[33,83,84]. Variants

614    with at least 3 reads were extracted. We then compared the called genotype dosage to the gold

615    standard genotype by computing the Spearman correlation and mean squared error (MSE).

616

617    **Peak Calling**

618    Genrich[27] (v0.6.1) was used to call peaks. A slightly modified version of Genrich was applied

619    to allow peak calling across a large number of samples (https://github.com/maxdudek/Genrich).

620    Genrich assigns p-values to genomic positions within each sample followed by combining p-

621    values across samples using Fisher's method to call peaks. Bam files were filtered using:

622    'samtools view -S -b -q 10'. Bam files were name sorted using: 'samtools sort -n

623    /path/to/q10_filtered_bams/sample.bam | samtools view -h -o

624    /path/to/nameSortedBams/sample.bam'. Peak calling parameters were: 'Genrich -t

625    /path/to/nameSortedBams/sample1.bam, path/to/nameSortedBams/sample2.bam,

626    path/to/nameSortedBams/sampleN.bam, -j , -o /path/to/outputFile -v -E

627    /path/to/blacklistRegions.bed -r -q 0.05'.

628

629    **Genomic Annotation Enrichment**

630    Genomic annotation enrichment analyses were performed using the R package annotatr

631    (v.1.28.0) (https://bioconductor.org/packages/release/bioc/html/annotatr.html). 100 iterations of

632    random, matched background data using bedtools shuffle with flags "-chrom -excl

633    /path/to/blacklistRegions.bed -g /path/to/chrSizes.txt". P values were calculated by quantifying

634    the number of random data iterations that were more extreme than the true data values for each

635    category.

636

**Encode Roadmap Enrichment**

638    Histone ChIP-seq data derived from adult human samples were downloaded from

639    https://www.encodeproject.org/search/?type=Experiment&status=released&award.project=Road

640    map . ATAC-seq peaks that overlapped histone mark data were identified using bedtools

641    intersect -wo -a /path/to/encodeData.bed -b /path/to/peakCoords.txt. 100 iterations of random,

642    matched background data using bedtools shuffle with flags "-chrom -excl

643    /path/to/blacklistRegions.bed -g /path/to/chrSizes.txt". P values were calculated by quantifying

644    the number of random data iterations that were more extreme than the true data values for each

645    histone mark.

646

**caQTL Mapping**

648    Sample peak counts were generated for all samples. To remove potential outlier peak regions,

649    peaks with mean count <1 and max count > 100,000 were removed. Peaks were also removed

650    if >5000 samples had a read count of zero in that peak. Given that a single individual might

651    contribute multiple samples to the 10,293 sample pool, we identified each sample that can be

652    attributed to each individual and averaged sample peak CPM values to calculate a single CPM

653    value per peak for each individual donor. This workflow results in 1454 individual donor samples

654    for caQTL mapping. Code available in file

655    "Post_peakCalling_CountMatrixGeneration_Pipeline.txt". tensorQTL (v.1.0.9) [85] was used to

656    identify caQTLs using a linear model with 3 genotype PCs and 200 principal components as

657    covariates. PCs generated from each cluster's chromatin accessibility peak read count data

658     sample matrix was used to map caQTLs on chromosome 1 over a large range of included PCs.

659     The optimized PC covariate number was chosen based on the elbow of the PCs included vs.

660     caQTL discovery plot on chromosome 1 (Supplementary Table 23). We tested all genotyped

661     biallelic genetic variants with MAF > 0.05 within 10 kilobases of all open chromatin peak

662     boundaries detected by Genrich from the ATAC-Seq data[35]. Empirical p-values were

663     estimated by tensorQTL to get peak-level p-values and q-values [86]. caQTL mapping code

664     available in file "caQTL_mapping_code_pipeline.txt".

665

666     **Lead caQTL/eQTL Enrichment**

667     Significant lead eQTL variants were downloaded for 49 tissues from GTEx v8 publicly available

668     data. Unique global sample analysis lead caQTLs (n= 21,647) were intersected with lead eQTL

669     variants to assess overlap within each GTEx tissue. The unique intersection of overlaps across

670     all tissues was considered to determine the total number of caQTL lead variants that were found

671     to be a lead eQTL variant in at least one tissue. Background variants were selected to perform

672     enrichment analyses. Background variants were chosen by randomly sampling non-lead caQTL

673     genetic variants that were matched, +/- 10%, to the allele frequency and distance to nearest

674     gene transcription start site of true lead caQTL variants. Enrichment of caQTLs/eQTLs in each

675     tissue was calculated as the ratio of the overlap of true lead caQTL/eQTL compared to the

676     overlap of background variants/eQTL across 100 iterations.

677

678     **Replication Analysis**

679     An external dataset was identified that was not included in our peak calling or caQTL mapping

680     workflow[49]. Global FDR5 caQTL peaks with any overlap with the external study and variants

681     tested in both analyses against these shared peaks were identified. External study p values

682     were used for $\pi_1$ replication rate calculation and plotted.

683

684 **GWAS Trait/Signal Selection**

685 GWAS summary stats for traits were downloaded February 2021 from the UKBB Neale Lab

686 repository and selected for relevant traits based on the following filters: h2 > 0.05, z > 7,

687 confidence == high. Independent significant GWAS signals from 78 traits were chosen to

688 prevent counting a single GWAS signal multiple times. This was done by selecting GWAS

689 signals with a minimum p-value of 5e-08, considering a window of 50 kb on either side of these

690 variants, clumping all variants with R2 > 0.01, and selecting the variant with the most significant

691 p-value as the lead GWAS signal for this locus.

692

693 **Colocalization Analyses**

694 Colocalization was performed using coloc[59] (v.5.2.3). All reported colocalizations utilized a

695 previously published approach to define significance[87]. This approach consists of considering

696 whether the colocalization is sufficiently powered, PP3+PP4 > 0.8. For those events that

697 surpass this threshold, we assessed whether the colocalization is significant, PP4/(PP3+PP4) >

698 0.9. GTEx v8 data were downloaded from https://www.gtexportal.org/home/downloads/adult-

699 gtex/bulk_tissue_expression.

700

701 **Colocalization Genome Annotations**

702 Genomic annotation enrichment analyses were performed using the R package annotatr

703 (v.1.28.0)(https://bioconductor.org/packages/release/bioc/html/annotatr.html). For each type of

704 colocalization, caQTL peaks involved in the colocalization were labeled with genomic

705 annotations they overlap. To perform an enrichment analysis, true data results were compared

706 with the median of 1000 iterations of random genomic regions matched to the true data using

707 bedtools shuffle with flags "-chrom -excl /path/to/blacklistRegions.bed -g /path/to/chrSizes.txt".

708 Summaries were produced by identifying significant enrichments (annotation category

709 enriched/depleted p value <= 0.05) across all traits or trait/tissue pairs and calculating the mean

710    and median enrichment/depletion values.

711

**Clustering Analyses**

713    To reduce the dimensions of the data, Uniform Manifold Approximation and Projection (UMAP)

714    was performed on the normalized sample CPM count matrix across all peaks. Kmeans

715    clustering was performed on UMAP coordinates 1 and 2. 11 outlier samples were removed from

716    analysis. The number of clusters was optimized using several clustering metrics

717    (Supplementary Table 22) and samples were assigned to a cluster based on the results of the

718    clustering algorithm.

719

**Cluster-specific caQTL mapping**

721    caQTL mapping was performed as in the global analysis. In this analysis, peaks identified in the

722    global analysis were included if at least 50% of cluster samples had non-zero CPMs in that

723    feature, resulting in the removal of 5-5920 (0.0003-0.35% of total peaks). All steps of the caQTL

724    mapping pipeline were performed within each cluster. caQTL mapping was performed including

725    3 genotype PCs and an optimized number of principal components based on each cluster. For

726    each cluster, a range of PCs generated from each cluster's chromatin accessibility peak read

727    count data sample matrix was used to map caQTLs on chromosome 1. The optimized PC

728    covariate number was chosen based on the elbow of the PCs included vs. caQTL discovery

729    plot. We tested all genotyped biallelic genetic variants with MAF > 0.05 within 10 kilobases of all

730    open chromatin peak boundaries detected by Genrich from the ATAC-Seq data[35]. Empirical p-

731    values were estimated by tensorQTL to get peak-level p-values and q-values [86]. All

732    colocalizations were performed as described for the global analyses.

733

**Cluster caQTL replication analyses**

735    Cluster caQTL replication of global caQTLs was assessed by extracting global caQTL peak test

36

736   statistics from each cluster and calculating $\pi_1$ replication rate. The reported replication rate for

737   each cluster was calculated by calculating the median $\pi_1$ replication rate after calculating $\pi_1$

738   replication rate with a range of values for the lambda parameter (from=0.1,to=0.9,by=0.05).

739   Cluster caQTL replication rate across all other clusters was calculated in a similar fashion. For

740   each cluster, cluster caQTL peak test statistics were extract from all other clusters and $\pi_1$

741   replication rate was calculated. The reported replication rate for each cluster was calculated by

742   calculating the median $\pi_1$ replication rate after calculating $\pi_1$ replication rate with a range of

743   values for the lambda parameter (from=0.1,to=0.9,by=0.05).

744

745   **Declarations**

746

747   **Ethics approval and consent to participate**

748   'Not applicable'

749   **Consent for publication**

750   'Not applicable'

751   **Availability of data and materials**

752   All data generated or analyzed during this study are included in this published article [and its

753   supplementary information files]. Publicly available samples used are listed in Supplementary

754   Table 1. The code used to generate the results and figures and generated data/results are

755   deposited in a Zenodo repository (https://doi.org/10.5281/zenodo.12706263) and will be made

756   public upon publication.

757   **Competing interests**

758   A.B. is a co-founder and equity holder of CellCipher, Inc, a stockholder in Alphabet, Inc, and has

759   consulted for Third Rock Ventures. N.C. is an employee and shareholder of Exai Bio, Inc. J.K.P.

760   and J.H.L. are employees of Gencove, Inc.

37

761 **Funding**

764 **Authors' contributions**

765 A.B., C.D.B., Y.H., B.M.W. designed the study. J.K.P. and J.H.L. generated genotype data. Y.H.

766 and B.M.W. performed computational analyses and prepared figures and tables. N.C. and T.L.

767 assisted with analyses. M.F.D. assisted with software. A.B., R.K., B.F.V., Y.H. and B.M.W.

768 wrote and revised the manuscript. All authors read and approved the final manuscript.

769 **Acknowledgements**

771

772

773

774 **References**

775 [1]  Buniello A, Macarthur JAL, Cerezo M, Harris LW, Hayhurst J, Malangone C, et al. The
776       NHGRI-EBI GWAS Catalog of published genome-wide association studies, targeted
777       arrays and summary statistics 2019. Nucleic Acids Res 2019;47:D1005–12.
778       https://doi.org/10.1093/NAR/GKY1120.
779 [2]  Maurano MT, Humbert R, Rynes E, Thurman RE, Haugen E, Wang H, et al. Systematic
780       localization of common disease-associated variation in regulatory DNA. Science
781       2012;337:1190–5. https://doi.org/10.1126/science.1222794.
782 [3]  Nicolae DL, Gamazon E, Zhang W, Duan S, Eileen Dolan M, Cox NJ. Trait-Associated
783       SNPs Are More Likely to Be eQTLs: Annotation to Enhance Discovery from GWAS. PLoS
784       Genet 2010;6:e1000888. https://doi.org/10.1371/JOURNAL.PGEN.1000888.
785 [4]  Gallagher MD, Chen-Plotkin AS. The Post-GWAS Era: From Association to Function. Am
786       J Hum Genet 2018;102:717–30. https://doi.org/10.1016/J.AJHG.2018.04.002.
787 [5]  Aguet F, Brown AA, Castel SE, Davis JR, He Y, Jo B, et al. Genetic effects on gene
788       expression across human tissues. Nature 2017 550:7675 2017;550:204–13.
789       https://doi.org/10.1038/nature24277.
790 [6]  Aguet F, Barbeira AN, Bonazzola R, Brown A, Castel SE, Jo B, et al. The GTEx
791       Consortium atlas of genetic regulatory effects across human tissues. Science (1979)
792       2020;369. https://doi.org/10.1126/SCIENCE.AAZ1776.

[7]   Battle A, Mostafavi S, Zhu X, Potash JB, Weissman MM, McCormick C, et al. Characterizing the genetic basis of transcriptome diversity through RNA-sequencing of 922 individuals. Genome Res 2014;24. https://doi.org/10.1101/gr.155192.113.

[8]   Boeger H, Griesenbeck J, Kornberg RD. Nucleosome Retention and the Stochastic Nature of Promoter Chromatin Remodeling for Transcription. Cell 2008;133. https://doi.org/10.1016/j.cell.2008.02.051.

[9]   Workman JL, Kingston RE. Alteration of nucleosome structure as a mechanism of transcriptional regulation. Annu Rev Biochem 1998;67. https://doi.org/10.1146/annurev.biochem.67.1.545.

[10]  Kumasaka N, Knights AJ, Gaffney DJ. Fine-mapping cellular QTLs with RASQUAL and ATAC-seq. Nat Genet 2016;48:206. https://doi.org/10.1038/NG.3467.

[11]  Degner JF, Pai AA, Pique-Regi R, Veyrieras JB, Gaffney DJ, Pickrell JK, et al. DNase-I sensitivity QTLs are a major determinant of human expression variation. Nature 2012;482. https://doi.org/10.1038/nature10808.

[12]  Khetan S, Kursawe R, Youn A, Lawlor N, Jillette A, Marquez EJ, et al. Type 2 diabetes-associated genetic variants regulate chromatin accessibility in human islets. Diabetes, vol. 67, 2018. https://doi.org/10.2337/db18-0393.

[13]  Krause MD, Huang RT, Wu D, Shentu TP, Harrison DL, Whalen MB, et al. Genetic variant at coronary artery disease and ischemic stroke locus 1p32.2 regulates endothelial responses to hemodynamics. Proc Natl Acad Sci U S A 2018;115. https://doi.org/10.1073/pnas.1810568115.

[14]  Tehranchi A, Hie B, Dacre M, Kaplow I, Pettie K, Combs P, et al. Fine-mapping cis-regulatory variants in diverse human populations. Elife 2019;8. https://doi.org/10.7554/eLife.39595.

[15]  Currin KW, Erdos MR, Narisu N, Rai V, Vadlamudi S, Perrin HJ, et al. Genetic effects on liver chromatin accessibility identify disease regulatory variants. Am J Hum Genet 2021;108. https://doi.org/10.1016/j.ajhg.2021.05.001.

[16]  Zeng B, Bendl J, Deng C, Lee D, Misir R, Reach SM, et al. Genetic regulation of cell-type specific chromatin accessibility shapes the etiology of brain diseases. BioRxiv 2023.

[17]  Wang J, Cheng X, Liang Q, Owen LA, Lu J, Zheng Y, et al. Single-cell multiomics of the human retina reveals hierarchical transcription factor collaboration in mediating cell type-specific effects of genetic variants on gene regulation. Genome Biol 2023;24. https://doi.org/10.1186/s13059-023-03111-8.

[18]  Keele GR, Quach BC, Israel JW, Chappell GA, Lewis L, Safi A, et al. Integrative QTL analysis of gene expression and chromatin accessibility identifies multi-tissue patterns of genetic regulation. PLoS Genet 2020;16. https://doi.org/10.1371/journal.pgen.1008537.

[19]  Pandey GK, Vadlamudi S, Currin KW, Moxley AH, Nicholas JC, McAfee JC, et al. Liver regulatory mechanisms of noncoding variants at lipid and metabolic trait loci. Human Genetics and Genomics Advances 2024;5. https://doi.org/10.1016/j.xhgg.2024.100275.

[20]  Buenrostro JD, Wu B, Chang HY, Greenleaf WJ. ATAC-seq: A Method for Assaying Chromatin Accessibility Genome-Wide. Current Protocols in Molecular Biology / Edited by Frederick M Ausubel . [et Al] 2015;109:21.29.1. https://doi.org/10.1002/0471142727.MB2129S109.

836  [21]  Buenrostro JD, Giresi PG, Zaba LC, Chang HY, Greenleaf WJ. Transposition of native
837         chromatin for fast and sensitive epigenomic profiling of open chromatin, DNA-binding
838         proteins and nucleosome position. Nat Methods 2013;10.
839         https://doi.org/10.1038/nmeth.2688.
840  [22]  Buenrostro JD, Wu B, Litzenburger UM, Ruff D, Gonzales ML, Snyder MP, et al. Single-
841         cell chromatin accessibility reveals principles of regulatory variation. Nature 2015;523.
842         https://doi.org/10.1038/nature14590.
843  [23]  Wasik K, Berisa T, Pickrell JK, Li JH, Fraser DJ, King K, et al. Comparing low-pass
844         sequencing and genotyping for trait mapping in pharmacogenetics. BMC Genomics
845         2021;22. https://doi.org/10.1186/s12864-021-07508-2.
846  [24]  Li JH, Mazur CA, Berisa T, Pickrell JK. Low-pass sequencing increases the power of
847         GWAS and decreases measurement error of polygenic risk scores compared to
848         genotyping arrays. Genome Res 2021;31. https://doi.org/10.1101/GR.266486.120.
849  [25]  Bailey T, Krajewski P, Ladunga I, Lefebvre C, Li Q, Liu T, et al. Practical Guidelines for
850         the Comprehensive Analysis of ChIP-seq Data. PLoS Comput Biol 2013;9.
851         https://doi.org/10.1371/journal.pcbi.1003326.
852  [26]  Liu H, Li R, Hu K, Ou J, Pak M, Green MR, et al. Best practices for the ATAC-seq assay
853         and its data analysis. Rigor and Reproducibility in Genetics and Genomics: Peer-
854         reviewed, Published, Cited, 2023. https://doi.org/10.1016/B978-0-12-817218-6.00016-4.
855  [27]  Gaspar JM. Genrich. Https://GithubCom/Jsh58/Genrich n.d.
856  [28]  Chiou J, Zeng C, Cheng Z, Han JY, Schlichting M, Miller M, et al. Single-cell chromatin
857         accessibility identifies pancreatic islet cell type– and state-specific regulatory programs of
858         diabetes risk. Nat Genet 2021;53. https://doi.org/10.1038/s41588-021-00823-0.
859  [29]  Wang D, Wu X, Jiang G, Yang J, Yu Z, Yang Y, et al. Systematic analysis of the effects
860         of genetic variants on chromatin accessibility to decipher functional variants in non-coding
861         regions. Front Oncol 2022;12. https://doi.org/10.3389/fonc.2022.1035855.
862  [30]  Turner AW, Hu SS, Mosquera JV, Ma WF, Hodonsky CJ, Wong D, et al. Single-nucleus
863         chromatin accessibility profiling highlights regulatory mechanisms of coronary artery
864         disease risk. Nat Genet 2022;54. https://doi.org/10.1038/s41588-022-01069-0.
865  [31]  Ha NT, Freytag S, Bickeboeller H. Coverage and efficiency in current SNP chips.
866         European Journal of Human Genetics 2014;22. https://doi.org/10.1038/ejhg.2013.304.
867  [32]  De Summa S, Malerba G, Pinto R, Mori A, Mijatovic V, Tommasi S. GATK hard filtering:
868         Tunable parameters to improve variant calling for next generation sequencing targeted
869         gene panel data. BMC Bioinformatics 2017;18. https://doi.org/10.1186/s12859-017-1537-
870         8.
871  [33]  Brouard JS, Schenkel F, Marete A, Bissonnette N. The GATK joint genotyping workflow
872         is appropriate for calling variants in RNA-seq experiments. J Anim Sci Biotechnol
873         2019;10. https://doi.org/10.1186/s40104-019-0359-0.
874  [34]  Kumasaka N, Knights AJ, Gaffney DJ. High-resolution genetic mapping of putative causal
875         interactions between regions of open chromatin. Nat Genet 2019;51.
876         https://doi.org/10.1038/s41588-018-0278-6.
877  [35]  Boyle AP, Song L, Lee BK, London D, Keefe D, Birney E, et al. High-resolution genome-
878         wide in vivo footprinting of diverse transcription factors in human cells. Genome Res
879         2011;21. https://doi.org/10.1101/gr.112656.110.

880   [36]   Lee CK, Shibata Y, Rao B, Strahl BD, Lieb JD. Evidence for nucleosome depletion at
881           active regulatory regions genome-wide. Nat Genet 2004;36.
882           https://doi.org/10.1038/ng1400.
883   [37]   Thurman RE, Rynes E, Humbert R, Vierstra J, Maurano MT, Haugen E, et al. The
884           accessible chromatin landscape of the human genome. Nature 2012 489:7414
885           2012;489:75–82. https://doi.org/10.1038/nature11232.
886   [38]   Lizio M, Abugessaisa I, Noguchi S, Kondo A, Hasegawa A, Hon CC, et al. Update of the
887           FANTOM web resource: Expansion to provide additional transcriptome atlases. Nucleic
888           Acids Res 2019;47. https://doi.org/10.1093/nar/gky1099.
889   [39]   Lizio M, Harshbarger J, Shimoji H, Severin J, Kasukawa T, Sahin S, et al. Gateways to
890           the FANTOM5 promoter level mammalian expression atlas. Genome Biol 2015;16.
891           https://doi.org/10.1186/s13059-014-0560-6.
892   [40]   Heintzman ND, Stuart RK, Hon G, Fu Y, Ching CW, Hawkins RD, et al. Distinct and
893           predictive chromatin signatures of transcriptional promoters and enhancers in the human
894           genome. Nat Genet 2007;39. https://doi.org/10.1038/ng1966.
895   [41]   Creyghton MP, Cheng AW, Welstead GG, Kooistra T, Carey BW, Steine EJ, et al.
896           Histone H3K27ac separates active from poised enhancers and predicts developmental
897           state. Proc Natl Acad Sci U S A 2010;107. https://doi.org/10.1073/pnas.1016071107.
898   [42]   Roadmap Epigenomics Consortium, Kundaje A, Meuleman W, Ernst J, Bilenky M, Yen A,
899           et al. Integrative analysis of 111 reference human epigenomes. Nature 2015;518.
900           https://doi.org/10.1038/nature14248.
901   [43]   Bernstein BE, Stamatoyannopoulos JA, Costello JF, Ren B, Milosavljevic A, Meissner A,
902           et al. The NIH Roadmap Epigenomics Mapping Consortium. Nat Biotechnol
903           2010;28:1045–8. https://doi.org/10.1038/nbt1010-1045.
904   [44]   Ninova M, Tóth KF, Aravin AA. The control of gene expression and cell identity by H3K9
905           trimethylation. Development (Cambridge) 2019;146. https://doi.org/10.1242/dev.181180.
906   [45]   Johnston AD, Simões-Pires CA, Thompson T V., Suzuki M, Greally JM. Functional
907           genetic variants can mediate their regulatory effects through alteration of transcription
908           factor binding. Nat Commun 2019;10. https://doi.org/10.1038/S41467-019-11412-5.
909   [46]   Musunuru K, Strong A, Frank-Kamenetsky M, Lee NE, Ahfeldt T, Sachs K V., et al. From
910           noncoding variant to phenotype via SORT1 at the 1p13 cholesterol locus. Nature 2010
911           466:7307 2010;466:714–9. https://doi.org/10.1038/nature09266.
912   [47]   Aguet F, Alasoo K, Li YI, Battle A, Im HK, Montgomery SB, et al. Molecular quantitative
913           trait loci. Nature Reviews Methods Primers 2023;3. https://doi.org/10.1038/s43586-022-
914           00188-6.
915   [48]   ENCODE Project Consortium. An integrated encyclopedia of DNA elements in the human
916           genome. Nature 2012;489:57–74. https://doi.org/10.1038/nature11247.
917   [49]   DeGorter MK, Goddard PC, Karakoc E, Kundu S, Yan SM, Nachun D, et al.
918           Transcriptomics and chromatin accessibility in multiple African population samples.
919           BioRxiv 2023. https://doi.org/10.1101/2023.11.04.564839.
920   [50]   Storey JD, Tibshirani R. Statistical significance for genomewide studies. Proc Natl Acad
921           Sci U S A 2003;100. https://doi.org/10.1073/pnas.1530509100.
922   [51]   Storey JD. A direct approach to false discovery rates. J R Stat Soc Series B Stat
923           Methodol 2002;64. https://doi.org/10.1111/1467-9868.00346.

924 [52] Karczewski KJ, Gupta R, Kanai M, Lu W, Tsuo K, Wang Y, et al. Pan-UK Biobank GWAS
925 improves discovery, analysis of genetic architecture, and resolution into ancestry-
926 enriched effects. MedRxiv 2024.
927 [53] Zhu Z, Zhang F, Hu H, Bakshi A, Robinson MR, Powell JE, et al. Integration of summary
928 data from GWAS and eQTL studies predicts complex trait gene targets. Nat Genet
929 2016;48. https://doi.org/10.1038/ng.3538.
930 [54] Gamazon ER, Segrè A V., Van De Bunt M, Wen X, Xi HS, Hormozdiari F, et al. Using an
931 atlas of gene regulation across 44 human tissues to inform complex disease- and trait-
932 associated variation. Nat Genet 2018;50. https://doi.org/10.1038/s41588-018-0154-4.
933 [55] Liu B, Gloudemans MJ, Rao AS, Ingelsson E, Montgomery SB. Abundant associations
934 with gene expression complicate GWAS follow-up. Nat Genet 2019;51.
935 https://doi.org/10.1038/s41588-019-0404-0.
936 [56] Kleinjan DA, Van Heyningen V. Long-range control of gene expression: Emerging
937 mechanisms and disruption in disease. Am J Hum Genet 2005;76.
938 https://doi.org/10.1086/426833.
939 [57] Smemo S, Tena JJ, Kim KH, Gamazon ER, Sakabe NJ, Gómez-Marín C, et al. Obesity-
940 associated variants within FTO form long-range functional connections with IRX3. Nature
941 2014;507:371. https://doi.org/10.1038/NATURE13138.
942 [58] Chun S, Casparino A, Patsopoulos NA, Croteau-Chonka DC, Raby BA, De Jager PL, et
943 al. Limited statistical evidence for shared genetic effects of eQTLs and autoimmune-
944 disease-associated loci in three major immune-cell types. Nat Genet 2017;49.
945 https://doi.org/10.1038/ng.3795.
946 [59] Giambartolomei C, Vukcevic D, Schadt EE, Franke L, Hingorani AD, Wallace C, et al.
947 Bayesian Test for Colocalisation between Pairs of Genetic Association Studies Using
948 Summary Statistics. PLoS Genet 2014;10:e1004383.
949 https://doi.org/10.1371/JOURNAL.PGEN.1004383.
950 [60] Umans BD, Battle A, Gilad Y. Where Are the Disease-Associated eQTLs? Trends in
951 Genetics 2021;37. https://doi.org/10.1016/j.tig.2020.08.009.
952 [61] Barbeira AN, Bonazzola R, Gamazon ER, Liang Y, Park YS, Kim-Hellmuth S, et al.
953 Exploiting the GTEx resources to decipher the mechanisms at GWAS loci. Genome Biol
954 2021;22:1–24. https://doi.org/10.1186/S13059-020-02252-4/FIGURES/6.
955 [62] Strober BJ, Elorbany R, Rhodes K, Krishnan N, Tayeb K, Battle A, et al. Dynamic genetic
956 regulation of gene expression during cellular differentiation. Science (1979) 2019;364.
957 https://doi.org/10.1126/science.aaw0040.
958 [63] Zhernakova D V., Deelen P, Vermaat M, Van Iterson M, Van Galen M, Arindrarto W, et
959 al. Identification of context-dependent expression quantitative trait loci in whole blood. Nat
960 Genet 2017;49. https://doi.org/10.1038/ng.3737.
961 [64] Fairfax BP, Humburg P, Makino S, Naranbhai V, Wong D, Lau E, et al. Innate immune
962 activity conditions the effect of regulatory variants upon monocyte gene expression.
963 Science (1979) 2014;343. https://doi.org/10.1126/science.1246949.
964 [65] Çalışkan M, Manduchi E, Rao HS, Segert JA, Beltrame MH, Trizzino M, et al. Genetic
965 and Epigenetic Fine Mapping of Complex Trait Associated Loci in the Human Liver. Am J
966 Hum Genet 2019;105:89. https://doi.org/10.1016/J.AJHG.2019.05.010.

[66]   Li YI, Van De Geijn B, Raj A, Knowles DA, Petti AA, Golan D, et al. RNA splicing is a primary link between genetic variation and disease. Science (1979) 2016;352. https://doi.org/10.1126/science.aad9417.

[67]   Mostafavi H, Spence JP, Naqvi S, Pritchard JK. Systematic differences in discovery of genetic effects on gene expression and complex traits. Nat Genet 2023;55:1866–75. https://doi.org/10.1038/s41588-023-01529-1.

[68]   Mountjoy E, Schmidt EM, Carmona M, Schwartzentruber J, Peat G, Miranda A, et al. An open approach to systematically prioritize causal variants and genes at all published human GWAS trait-associated loci. Nat Genet 2021;53. https://doi.org/10.1038/s41588-021-00945-5.

[69]   Raghuram V, Weber S, Raber J, Chen DH, Bird TD, Maylie J, et al. Assessment of mutations in KCNN2 and ZNF135 to patient neurological symptoms. Neuroreport 2017;28. https://doi.org/10.1097/WNR.0000000000000754.

[70]   Corces MR, Buenrostro JD, Wu B, Greenside PG, Chan SM, Koenig JL, et al. Lineage-specific and single-cell chromatin accessibility charts human hematopoiesis and leukemia evolution. Nat Genet 2016;48. https://doi.org/10.1038/ng.3646.

[71]   McInnes L, Healy J, Saul N, Großberger L. UMAP: Uniform Manifold Approximation and Projection. J Open Source Softw 2018;3. https://doi.org/10.21105/joss.00861.

[72]   Nasser J, Bergman DT, Fulco CP, Guckelberger P, Doughty BR, Patwardhan TA, et al. Genome-wide enhancer maps link risk variants to disease genes. Nature 2021;593. https://doi.org/10.1038/s41586-021-03446-x.

[73]   Shi C, Ray-Jones H, Ding J, Duffus K, Fu Y, Gaddi VP, et al. Chromatin Looping Links Target Genes with Genetic Risk Loci for Dermatological Traits. Journal of Investigative Dermatology 2021;141. https://doi.org/10.1016/j.jid.2021.01.015.

[74]   McClay JL, Shabalin AA, Dozmorov MG, Adkins DE, Kumar G, Nerella S, et al. High density methylation QTL analysis in human blood via next-generation sequencing of the methylated genomic DNA fraction. Genome Biol 2015;16. https://doi.org/10.1186/s13059-015-0842-7.

[75]   Lemire M, Zaidi SHE, Ban M, Ge B, Aïssi D, Germain M, et al. Long-range epigenetic regulation is conferred by genetic variation located at thousands of independent loci. Nat Commun 2015;6. https://doi.org/10.1038/ncomms7326.

[76]   Melzer D, Perry JRB, Hernandez D, Corsi AM, Stevens K, Rafferty I, et al. A genome-wide association study identifies protein quantitative trait loci (pQTLs). PLoS Genet 2008;4. https://doi.org/10.1371/journal.pgen.1000072.

[77]   Jeong R, Bulyk ML. Chromatin accessibility variation provides insights into missing regulation underlying immune-mediated diseases. BioRxiv 2024:2024.04.12.589213. https://doi.org/10.1101/2024.04.12.589213.

[78]   Zeng B, Lloyd-Jones LR, Montgomery GW, Metspalu A, Esko T, Franke L, et al. Comprehensive multiple eQTL detection and its application to GWAS interpretation. Genetics 2019;212. https://doi.org/10.1534/genetics.119.302091.

[79]   Yao DW, O'Connor LJ, Price AL, Gusev A. Quantifying genetic effects on disease mediated by assayed gene expression levels. Nat Genet 2020;52. https://doi.org/10.1038/s41588-020-0625-2.

1010 [80] Claussnitzer M, Dankel SN, Kim K-H, Quon G, Meuleman W, Haugen C, et al. FTO
1011      Obesity Variant Circuitry and Adipocyte Browning in Humans . New England Journal of
1012      Medicine 2015;373. https://doi.org/10.1056/nejmoa1502214.

1013 [81] Martin-Rufino JD, Castano N, Pang M, Grody EI, Joubran S, Caulier A, et al. Massively
1014      parallel base editing to map variant effects in human hematopoiesis. Cell 2023;186.
1015      https://doi.org/10.1016/j.cell.2023.03.035.

1016 [82] Lowy-Gallego E, Fairley S, Zheng-Bradley X, Ruffier M, Clarke L, Flicek P. Variant calling
1017      on the grch38 assembly with the data from phase three of the 1000 genomes project
1018      [version 2; peer review: 1 approved, 1 not approved]. Wellcome Open Res 2019;4.
1019      https://doi.org/10.12688/wellcomeopenres.15126.1.

1020 [83] Hwang S, Kim E, Lee I, Marcotte EM. Systematic comparison of variant calling pipelines
1021      using gold standard personal exome variants. Sci Rep 2015;5.
1022      https://doi.org/10.1038/srep17875.

1023 [84] Langmead B, Salzberg SL. Fast gapped-read alignment with Bowtie 2. Nat Methods
1024      2012;9. https://doi.org/10.1038/nmeth.1923.

1025 [85] Taylor-Weiner A, Aguet F, Haradhvala NJ, Gosai S, Anand S, Kim J, et al. Scaling
1026      computational genomics to millions of individuals with GPUs. Genome Biol 2019;20.
1027      https://doi.org/10.1186/s13059-019-1836-7.

1028 [86] Ongen H, Buil A, Brown AA, Dermitzakis ET, Delaneau O. Fast and efficient QTL mapper
1029      for thousands of molecular phenotypes. Bioinformatics 2016;32.
1030      https://doi.org/10.1093/bioinformatics/btv722.

1031 [87] Alasoo K, Rodrigues J, Mukhopadhyay S, Knights AJ, Mann AL, Kundu K, et al. Shared
1032      genetic effects on chromatin and gene expression indicate a role for enhancer priming in
1033      immune response. Nat Genet 2018;50. https://doi.org/10.1038/s41588-018-0046-7.

1034

1035