

# 1 Improved multi-ancestry fine-mapping identifies *cis*-regulatory variants 2 underlying molecular traits and disease risk

3  
4 Zeyun Lu<sup>1</sup>, Xinran Wang<sup>1</sup>, Matthew Carr<sup>1</sup>, Artem Kim<sup>1</sup>, Steven Gazal<sup>1,2,3</sup>, Pejman Mohammadi<sup>4,5,6</sup>, Lang Wu<sup>7</sup>,  
5 Alexander Gusev<sup>8</sup>, James Pirruccello<sup>9</sup>, Linda Kachuri<sup>10,11</sup>, Nicholas Mancuso<sup>1,2,3,12</sup>

- 6 1. Center for Genetic Epidemiology, Keck School of Medicine, University of Southern California, Los Angeles, CA,  
7 USA  
8 2. Department of Population and Public Health Sciences, Keck School of Medicine, University of Southern  
9 California, Los Angeles, CA, USA  
10 3. Department of Quantitative and Computational Biology, University of Southern California, Los Angeles, CA  
11 4. Center for Immunity and Immunotherapies, Seattle Children's Research Institute, Seattle, WA, USA  
12 5. Department of Pediatrics, University of Washington School of Medicine, Seattle, WA, USA  
13 6. Department of Genome Sciences, University of Washington, Seattle, WA, USA  
14 7. Cancer Epidemiology Division, Population Sciences in the Pacific Program, University of Hawai'i Cancer Center,  
15 University of Hawai'i at Mānoa, Honolulu, HI, USA  
16 8. Harvard Medical School and Dana-Farber Cancer Institute, Boston, MA, USA  
17 9. Division of Cardiology, University of California San Francisco, San Francisco, CA, USA  
18 10. Department of Epidemiology and Population Health, Stanford University School of Medicine, Stanford, CA,  
19 USA  
20 11. Stanford Cancer Institute, Stanford University School of Medicine, Stanford, CA, USA  
21 12. Corresponding Author

## 22 23 Contacts:

- 24 1. Zeyun Lu ([zeyunlu@usc.edu](mailto:zeyunlu@usc.edu))  
25 2. Nicholas Mancuso ([Nicholas.Mancuso@med.usc.edu](mailto:Nicholas.Mancuso@med.usc.edu))

## 26 Abstract

27 Multi-ancestry statistical fine-mapping of *cis*-molecular quantitative trait loci (*cis*-molQTL) aims to improve the  
28 precision of distinguishing causal *cis*-molQTLs from tagging variants. However, existing approaches fail to reflect  
29 shared genetic architectures. To solve this limitation, we present the Sum of Shared Single Effects (SuShiE) model,  
30 which leverages LD heterogeneity to improve fine-mapping precision, infer cross-ancestry effect size correlations,  
31 and estimate ancestry-specific expression prediction weights. We apply SuShiE to mRNA expression measured in  
32 PBMCs (n=956) and LCLs (n=814) together with plasma protein levels (n=854) from individuals of diverse  
33 ancestries in the TOPMed MESA and GENOA studies. We find SuShiE fine-maps *cis*-molQTLs for 16% more genes

34 compared with baselines while prioritizing fewer variants with greater functional enrichment. SuShiE infers highly  
35 consistent *cis*-molQTL architectures across ancestries on average; however, we also find evidence of  
36 heterogeneity at genes with predicted loss-of-function intolerance, suggesting that environmental interactions  
37 may partially explain differences in *cis*-molQTL effect sizes across ancestries. Lastly, we leverage estimated *cis*-  
38 molQTL effect-sizes to perform individual-level TWAS and PWAS on six white blood cell-related traits in AOU  
39 Biobank individuals (n=86k), and identify 44 more genes compared with baselines, further highlighting its benefits  
40 in identifying genes relevant for complex disease risk. Overall, SuShiE provides new insights into the *cis*-genetic  
41 architecture of molecular traits.

## 42 Introduction

43 Characterizing the functional consequences of genetic variation remains a crucial task in deciphering the  
44 mechanisms underlying complex disease risk<sup>1,2</sup>. To this end, *cis*-molecular quantitative trait loci (*cis*-molQTL)  
45 mapping seeks to identify genetic variants associated with genomically proximal molecular features measured  
46 across diverse cellular, tissue, and environmental contexts<sup>3-14</sup>. However, due to linkage disequilibrium (LD), it is  
47 challenging to distinguish causal *cis*-molQTLs from tagging variants within a genomic region<sup>3,5</sup>. Statistical fine-  
48 mapping aims to resolve precisely this issue<sup>15-19</sup>, yet pervasive LD signals limit the resolution of these approaches.  
49 Previous efforts have demonstrated that leveraging the heterogeneity of LD and minor allele frequency (MAF)  
50 across diverse ancestries improves the precision of statistical fine-mapping and therefore enhances our biological  
51 understanding of complex diseases<sup>20-25</sup> and molecular traits<sup>26-32</sup>.

52 While existing multi-ancestry fine-mapping frameworks have been proposed for the analysis of complex traits and  
53 diseases<sup>30,33-41</sup>, they have several limitations in the context of large-scale *cis*-molQTL data. First, many approaches  
54 do not model the correlation of causal variant effect sizes across ancestries or assume that they are a-priori  
55 independent across ancestries, which fails to reflect shared or similar genetic architectures<sup>33,35,37,38</sup>. Second,  
56 existing multi-ancestry approaches scale poorly, which precludes their application to thousands of molecular traits

57 commonly measured in *cis*-molQTL studies<sup>33,35,40</sup>. Third, current fine-mapping approaches lack ancestry-specific  
58 effect size estimates<sup>33,35,37</sup>, which neglects their potential use in post-Genome-wide Association Studies (GWASs)  
59 frameworks (e.g., Transcriptome- and Proteome-wide Association Studies (TWASs/PWASs)<sup>42-47</sup>. Last, while recent  
60 approaches address some of these limitations, existing software implementations are capable of analyzing only  
61 two ancestries, which excludes datasets consisting of ever-increasing diverse ancestries<sup>39</sup>.

62 Here, we describe the Sum of Shared Single Effects (SuShiE) approach to fine-map genetic variants shared across  
63 diverse ancestries for thousands of molecular traits. SuShiE integrates genotypic and molecular data from multiple  
64 ancestries to identify *cis*-molQTLs while modeling and learning the covariance structures of shared/non-shared  
65 signals. SuShiE leverages four key insights. First, SuShiE improves fine-mapping precision of the shared *cis*-molQTLs  
66 by leveraging LD across different ancestries. Second, it estimates ancestry-specific effect sizes at shared *cis*-  
67 molQTLs. Third, it infers the prior effect size correlation across ancestries to shed light on genetic similarities and  
68 differences. Lastly, SuShiE is implemented using a scalable variational inference algorithm that runs seamlessly on  
69 CPUs, GPUs, or tensor-processing units (TPUs).

70 Through extensive simulations, we show that SuShiE outputs higher posterior inclusion probabilities (PIPs) at  
71 causal *cis*-molQTLs, outputs smaller credible set sizes, and exhibits better calibration compared with current  
72 approaches<sup>15,38</sup>. Using bulk mRNA expression levels measured in peripheral blood mononuclear cells (PBMCs) and  
73 lymphoblastoid cell lines (LCLs) together with protein abundance measured in plasma, we fine-map 36,911  
74 molecular phenotypes across American European, African, and Hispanic ancestries from TOPMed-MESA<sup>48,49</sup>  
75 ( $n_{\text{mRNA}}=956$  and  $n_{\text{protein}}=814$ ) and GENOA<sup>26</sup> ( $n_{\text{mRNA}}=854$ ). SuShiE fine-maps significantly more *cis*-molQTLs with  
76 smaller credible sets and greater enrichment in relevant functional annotations compared with existing methods.  
77 In addition, SuShiE infers shared genetic architecture of *cis*-molQTL in significantly heritable genes and shows the  
78 heterogeneity across ancestries of signals associated with multiple measures of loss-of-function (LOF) intolerance.  
79 Last, we integrate ancestry-specific *cis*-molQTL effects inferred by SuShiE with six white blood cell-related traits  
80 to perform individual-level TWAS and PWAS in the All of Us Biobank (average  $n=86,345$ )<sup>50</sup> and observe that SuShiE-

81 based prediction models identified 44 additional associated genes compared with the baseline approach. Overall,  
82 our approach sheds light on understanding the genetic *cis*-architecture of molecular data across multiple  
83 ancestries.

## 84 Results

### 85 SuShiE overview

86 Here, we briefly introduce the SuShiE model (for a detailed description, see **Methods** and **Supplementary Note**).  
87 SuShiE assumes *cis*-molQTLs are present in *all* ancestries (defined as shared *cis*-molQTLs) while allowing for effect  
88 sizes at causal *cis*-molQTLs to covary across ancestries a-priori, in contrast to previous multi-ancestry  
89 approaches<sup>15,33,35,37,38</sup>. These assumptions provide enough flexibility to model a variety of *cis*-genetic architectures  
90 across ancestries, including cases when effects are present only in a subset of ancestries. For instance, when  
91 effects are observed only in a subset of ancestries, prior variances can be shrunk towards zero to effectively allow  
92 for *ancestry-specific* causal *cis*-molQTLs.

93 Focusing on the  $i^{\text{th}}$  out of  $k$  ancestries, SuShiE models the normalized levels of a molecular trait  $\mathbf{g}_i$  measured in  
94  $n_i$  individuals as a linear combination of  $p$  genotyped variants  $\mathbf{X}_i$  as

$$95 \quad \mathbf{g}_i = \mathbf{X}_i \left( \sum_{l=1}^L \boldsymbol{\gamma}_l \cdot b_{i,l} \right) + \boldsymbol{\epsilon}_i,$$

96 where  $L$  is the number of shared effects,  $\boldsymbol{\gamma}_l$  is a  $p \times 1$  binary vector selecting the  $l^{\text{th}}$  causal *cis*-molQTL shared  
97 across ancestries,  $b_{i,l}$  is the  $l^{\text{th}}$  effect size in the  $i^{\text{th}}$  ancestry, and environmental noise distributed as  
98  $\boldsymbol{\epsilon}_i \sim N(\mathbf{0}, \sigma_{i,e}^2 \mathbf{I}_{n_i})$  (**Fig. 1A**). Following previous work<sup>15,51,52</sup>, we place a  $\text{Multi}(1, \boldsymbol{\pi})$  prior over  $\boldsymbol{\gamma}_l$  where  $\boldsymbol{\pi}$  is a  
99  $p \times 1$  vector representing prior probability for each SNP to be shared *cis*-molQTLs; however, unlike existing  
100 approaches<sup>33,35,37,38</sup>, we organize ancestry-specific effect sizes under a multivariate normal prior  $[b_{1,l}, \dots, b_{k,l}] \sim$   
101  $MVN(\mathbf{0}, \mathbf{C}_l)$  where  $\mathbf{C}_l$  is the  $l^{\text{th}}$   $k \times k$  prior effect size covariance matrix. To perform scalable inference, we use a

102 variational Bayesian approach and compute, for each of the  $L$  shared effects, the posterior probability of a shared  
103 causal *cis*-molQTL ( $\alpha_i$ ), the ancestry-specific posterior effect sizes, and covariances, in addition to prior effect-size  
104 correlations (**Fig. 1B**) inferred through a procedure analogous to Empirical Bayes. Through learning prior effect-  
105 size correlations, SuShiE can quantify genes' heterogeneity in *cis*-molQTL effects across ancestries.  
106 SuShiE constructs a 90% credible set for each of the  $L$  effects along with a posterior inclusion probability (PIP) for  
107 each SNP to be putative causal *cis*-molQTL (see **Methods**). SuShiE is implemented in an open-source command-  
108 line Python software with JAX (see **Methods** and **Code Availability**) using *Just-In-Time* compilation to achieve high-  
109 speed inference that runs seamlessly on CPUs, GPUs, or TPUs at <https://github.com/mancusolab/sushie>.

## 110 SuShiE outperforms other methods in realistic simulations

111 First, to recapitulate the benefits of multi-ancestry study design<sup>33,35,37–41</sup>, we performed simulations varying the  
112 number of contributing ancestries under a fixed total sample size (see **Methods**). As the number of ancestries  
113 increased, SuShiE produced higher PIPs at causal *cis*-molQTLs, smaller credible set sizes, and better calibration  
114 (**Fig. S1**), reaffirming that increasing genetic diversity refines fine-mapping results compared with expanding the  
115 sample size of a single ancestry. Next, we evaluated the performance of SuShiE in simulations by varying different  
116 parameters and compared against three baselines: SuShiE-Indep (i.e., SuShiE assuming no a-priori correlation of  
117 effect sizes across ancestries), meta-SuSiE (i.e., a meta-analysis on single-ancestry SuSiE), and SuSiE (i.e., SuSiE  
118 performed over data aggregated across ancestries; see **Methods**). For all simulations, SuShiE output higher PIPs  
119 at causal *cis*-molQTLs ( $\sim 0.06$  on average; all  $P < 3.1e-4$ ; **Fig. 2A, S2**), smaller credible set sizes ( $\sim 0.73$  on average; 2  
120 out of 3 comparisons  $P < 0.05$ ; **Fig. 2B, S3**), and better calibration ( $\sim 0.08$  on average; all  $P < 1.51e-7$ ; **Fig. 2C, S4**).  
121 SuShiE similarly outperformed competing methods under simulations with differential power (**Fig. S5**) and genetic  
122 architectures across ancestries (**Fig. S6**). Next, we evaluated the ability of SuShiE to infer prior effect size  
123 correlations from data (see **Methods**). SuShiE accurately estimated primary effect size correlations (**Fig. 2D**) with

124 higher-order effects having diminishing accuracies. This result was likely due to decreasing statistical power, as  
125 evidenced by simulations under increased sample sizes (**Fig. 2D, S7**).

126 Next, we assessed the robustness of SuShiE when there exist genetic variants causal for only a subset of ancestries  
127 in addition to shared causal *cis*-molQTLs (see **Methods**). As the number of ancestry-specific *cis*-molQTLs increased,  
128 the performance of all approaches decreased compared with previous simulations. However, SuShiE continued  
129 producing higher PIPs at shared causal *cis*-molQTLs (**Fig. S8A**), smaller credible set sizes (**Fig. S8B**), and better  
130 calibrated credible sets (**Fig. S8C**), demonstrating SuShiE's robustness when ancestry-specific *cis*-molQTLs are  
131 present. We also evaluated performance in simulations where the number of causal effects (i.e.,  $L$ ) differs from  
132 the number specified at inference and observed that SuShiE similarly outperformed alternative approaches (**Fig.**  
133 **S9**).

134 Last, we evaluated the use of SuShiE-derived ancestry-specific effect sizes in *cis*-molQTL data as a means to predict  
135 the genetic component of gene expression for downstream TWAS<sup>42-44</sup>. Briefly, we performed simulations under a  
136 model in which gene expression mediates disease risk and compared SuShiE predictions with commonly used  
137 approaches for prediction-based TWAS (e.g., LASSO<sup>53</sup>, Elastic Net<sup>54</sup>, and gBLUP<sup>55</sup>) to identify susceptibility genes  
138 (see **Methods**). SuShiE-derived prediction models more accurately recapitulated gene expression levels compared  
139 with existing approaches and exhibited higher statistical power for TWAS with various study sample sizes and  
140 proportion of trait heritability mediated by gene expression (**Fig. 2E-F, S10**).

141 Overall, SuShiE outperforms existing approaches in realistic parameter settings, remains robust under model  
142 misspecifications, and improves statistical power in post-GWAS analyses.

### 143 **SuShiE identifies more functionally relevant *cis*-molQTL signals**

144 Having verified that SuShiE outperforms other methods under realistic simulations, we next sought to perform  
145 fine-mapping on 36, 911 molecular phenotypes from diverse ancestries. Specifically, from the Trans-Omics for  
146 Precision Medicine program Multi-Ethnic Study of Atherosclerosis<sup>48,49</sup> (TOPMed-MESA), we analyzed mRNA

147 expression data of 21,747 genes measured in PBMCs (visit-1; n=956) and protein expression data of 1,274 genes  
148 measured in plasma (visit-1; n=854) for American European, African, and Hispanic ancestries (EUR, AFR, and HIS),  
149 together with mRNA expression data of 13,890 genes measured in LCLs (n=814) for EUR and AFR from the Genetic  
150 Epidemiology Network of Arteriopathy study<sup>26</sup> (GENOA; see **Methods**; **Table S1**).

151 Focusing on 1Mb windows for each gene (i.e., *cis*-region), SuShiE fine-mapped *cis*-molQTLs for 21,088 phenotypes  
152 (e/pGenes), representing an average increase of 3,378 (16%) compared with existing methods (i.e., SuShiE-Indep,  
153 Meta-SuSiE, and SuSiE; all  $P < 2.94 \times 10^{-110}$ ; see **Methods**). For example, SuShiE fine-mapped 21% more e/pGenes  
154 compared to single-ancestry SuSiE followed by meta-analysis (i.e., Meta-SuSiE;  $P = 7.01 \times 10^{-238}$ ), again highlighting  
155 the benefit of multi-ancestry study design. SuShiE-based credible sets maintained higher average PIPs (~0.07 on  
156 average) and higher frequency of *cis*-molQTLs with PIPs  $> 0.9$  (~0.02 on average), as well as smaller credible sets  
157 in most cases (~6.24 on average; **Table S2**). We found the performance advantage slightly diminished in TOPMed-  
158 MESA protein and GENOA mRNA datasets, likely due to lower statistical power. Using the number of credible sets  
159 identified after purity pruning (see **Methods**), SuShiE estimated most (90.4%) molecular phenotypes to exhibit 1-  
160 3 *cis*-molQTL signals (**Fig. 3A**) with PIPs localizing near the transcription start site (TSS; **Fig. 3B**), consistent with  
161 previous studies<sup>3,4,26,56,57</sup>.

162 To characterize the regulatory function of identified *cis*-molQTL signals, we performed enrichment analysis using  
163 PIPs with 89 genomic functional annotations (see **Methods**). We observed that PIPs inferred by SuShiE were  
164 enriched in 83/89 annotations across all three datasets, with the highest enrichment occurring in promoter  
165 regions (**Table S3**). For example, PIPs were enriched in 4/5 candidate *cis*-regulatory elements (cCREs) from  
166 ENCODE Registry v3<sup>58</sup> (**Fig. 3C**) and in all 10 cell-type/tissue-specific cCREs using single-nucleus(sn) or single-cell(sc)  
167 ATAC-Seq<sup>59,60</sup> (**Fig. S11**). Importantly, PIPs inferred by SuShiE were more enriched across functional annotations  
168 compared with those computed from existing fine-mapping methods (all  $P < 8.13 \times 10^{-3}$ ; **Table S4**), highlighting  
169 SuShiE's ability to better prioritize functionally relevant *cis*-molQTLs. Next, to explore how potential regulatory  
170 function may differ among *cis*-molQTLs contributing to the same gene, we repeated the above analyses using per-

171 effect posterior probabilities ( $\alpha_l$ ), rather than overall inclusion probabilities (i.e., PIPs). First, the initial three  
172 shared effects were similarly localized near the TSS (**Fig. S12**) and were more enriched in promoter regions  
173 compared to the PIP-based analyses (**Fig. S13; Table S5**), echoing the previous finding that most genes are  
174 regulated by 1-3 *cis*-molQTLs<sup>3,4,26,57</sup>. Second, we found *cis*-molQTLs with weaker effects were further away from  
175 the TSS on average (**Fig. S14**), likely due to statistical power. For example, we observed the expected distance to  
176 TSS for the initial three shared effects was 84.7kb compared with 144.5 kb for the remaining shared effects (i.e.,  
177 from L=6 to L=10; P=8.39e-113).

178 Last, we sought to validate our fine-mapping results by applying SuShiE on molecular phenotypes from three  
179 independent datasets: mRNA expression measured in PBMCs of EUR, AFR, and HIS ancestries from TOPMed-  
180 MESA<sup>48,49</sup> (visit-5, ten-year after visit-1; n=875), mRNA expression measured in LCLs (n=462) of EUR and Yoruba  
181 (YRI) ancestries from GEUVADIS study<sup>61</sup>, and protein expression measured in plasma of EUR ancestry (N=3,301;  
182 single-ancestry SuSiE) from INTERVAL study<sup>5</sup> (see **Methods; Table S1**). First, we confirmed SuShiE identifies 4,361  
183 (21%; all P<2.89e-112) more e/pGenes on average compared with existing methods while obtaining higher  
184 average PIPs (~0.07 on average), smaller credible set sizes (~6.54 on average), and more *cis*-molQTLs with PIPs >  
185 0.9 (~0.04 on average) for TOPMed-MESA visit-5 and GEUVADIS (**Table S6**). Second, focusing on 20,502 e/pGenes  
186 identified by SuShiE that were also measured in validation datasets, 34% (41%, 32%, and 13% for TOPMed-MESA  
187 visit-5, INTERVAL, and GEUVADIS, respectively) *cis*-molQTLs replicated in the validation datasets with an average  
188 cosine similarity of 0.70 (0.72, 0.63, and 0.45 for the three mentioned studies; P<2e-200 for all), which increased  
189 to 73% and 0.75 respectively after conditioning on significantly heritable genes and the primary signal (see  
190 **Methods**). The diminished replication performance of GEUVADIS likely resulted from a combination of  
191 significantly reduced sample sizes, admixture differences between African YRI and American Africans in GENOA,  
192 and genotyping differences (see **Methods**). Furthermore, SuShiE exhibited similar replication ratios and cosine  
193 similarities compared to existing methods, suggesting the higher number of e/pGenes identified by SuShiE were  
194 not likely due to false positives (**Table S7; see Methods**).



195 Overall, by jointly modeling multi-ancestry data, SuShiE identifies additional *cis*-regulatory mechanisms for  
196 molecular traits.

### 197 SuShiE identifies putative eQTL for *URGCP*

198 Here, we showcase a putative eQTL for *URGCP*, a gene on chromosome 7 that has been implicated in tumor growth  
199 and progression<sup>62–66</sup>. SuShiE fine-mapped a single SNP in TOPMed-MESA mRNA (*rs2528382*; GRCh38: 7:43926148;  
200 PIP=0.94; **Fig. 4A**), while alternative methods did not produce credible sets for this gene. Importantly, SuShiE  
201 replicated *rs2528382* in TOPMed-MESA visit-5 mRNA data. We found *rs2528382* was reported as significant in  
202 whole blood eQTL data from the eQTLGen Consortium<sup>4</sup>, the Study of African Americans, Asthma, Genes, and  
203 Environments (SAGE), and the Genes-Environments and Admixture in Latino Asthmatics (GALA II) study<sup>31</sup>, further  
204 supporting its role in regulating *URGCP* expression levels. Investigating the functional consequences of *rs2528382*  
205 using genomic annotations, we found *rs2528382* represents a non-coding exon variant within the 5' UTR<sup>67</sup>, and  
206 localizes within a proximal enhancer region (pELS), as evidenced by strong signals of H3K27ac in PBMCs<sup>58</sup> falling  
207 within 2kb of the TSS (**Fig. 4B**). Lastly, through snATAC-seq<sup>59</sup> and scATAC-seq<sup>60</sup>, we found *rs2528382* localizes  
208 within an open chromatin accessibility region measured in different cell types, such as PBMCs, naive T cells, naive  
209 B cells, cytotoxic NK cells, and monocytes. Altogether, these results suggest that *rs2528382* regulates *URGCP*  
210 expression levels in PBMCs through disruption of regulatory activity.

### 211 SuShiE reveals heterogeneity of *cis*-molQTL effect sizes at the loss-of-function intolerant genes

212 After validating *cis*-molQTLs identified by SuShiE, we next sought to characterize genetic architectures of  
213 molecular traits across ancestries. First, we computed *cis*-SNP heritability for all e/pGenes of each ancestry and  
214 observed 87% significant heritable genes (in at least one ancestry) across studies (**Fig. S15**), which resulted in  
215 highly correlated estimates across ancestries (**Fig. S16**). Next, using SuShiE-derived estimates of *cis*-molQTL  
216 correlation across ancestries (see **Methods**), we found highly consistent effect-size correlations on average (0.81,

217 0.86, and 0.87 for EUR-AFR, EUR-HIS, and AFR-HIS, respectively), which further increased when focusing on genes  
218 whose heritabilities are significant in all ancestries (0.94, 0.98 and 0.99, respectively; 9,885 genes; 46.9%; **Figs.**  
219 **S17-S18**). Altogether, these results further affirm previous results<sup>20,21,23,68-74</sup> demonstrating primarily shared  
220 genetic architectures for molecular traits across ancestries.

221 Despite this evidence, we observed a long tail of heterogeneous effect sizes (i.e., SuShiE-estimated effect size  
222 correlation <1), suggesting the presence of ancestry-specific *cis*-molQTL effects (**Fig. S19**), which is consistent with  
223 previous multi-ancestry *cis*-molQTL studies<sup>27,31,72</sup>. To characterize this apparent heterogeneity across ancestries,  
224 we correlated the estimated correlation signals with multiple measures of constraint (pLI<sup>75</sup>, LOEUF<sup>76</sup>, EDS<sup>77</sup>, RVIS<sup>78</sup>,  
225 and  $s_{\text{het}}$ <sup>79</sup>) and found highly significant associations (**Table 1**; see **Methods**). Overall, genes with lower effect-size  
226 correlations across ancestries exhibited higher intolerance to loss-of-function mutations on average. For example  
227 using TOPMed-MESA mRNA dataset, we observed an average *cis*-molQTL effect size correlation of 0.81 (when L=1;  
228 SE=0.02) between EUR and AFR individuals at genes that exhibited pLI >0.9, which increased to 0.86 (when L=1;  
229 SE=0.01) when focusing on genes with pLI <0.1. Genes with high constraint exhibited lower estimates of *cis*-SNP  
230 heritability on average (**Table S8**), which may result in apparent heterogeneity arising from low statistical power.  
231 Given this, we re-analyzed putative relationships using estimated covariances, only primary signals (L=1), and  
232 bootstrapped standard errors and found broadly consistent results (**Table 1**). In addition, we observed our results  
233 were robust to adjusting for Wright's fixation index ( $F_{\text{st}}$ ; **Table 1**; see **Methods**), suggesting  
234 heterogeneity/constraint associations are not driven solely by allele frequency differences across ancestries.

235 To investigate the relationship between *cis*-molQTLs identified by SuShiE and gene constraint, we first observed  
236 inverse associations between the number of fine-mapped *cis*-molQTLs per gene and constraint (**Fig. S20**),  
237 consistent with several previous studies showing the depletion of *cis*-molQTLs for high constraint genes<sup>56,77,80</sup>.  
238 However, we also observed positive associations between expected *cis*-molQTLs' distance to TSS and constraint,  
239 affirming previous results that high constraint genes tended to have more complex regulatory regions<sup>56,77</sup> (**Fig.**  
240 **S21**; see **Methods**). In addition, we correlated gene enrichment scores from ENCODE<sup>58</sup> cCREs with constraint

241 scores. We found that putative causal *cis*-molQTLs for high constraint genes tended to be enriched for distal  
242 enhancers (dELS) and depleted for promoter (PLS) and proximal enhancers (pELS) compared with weakly  
243 constrained genes, consistent with several previous studies<sup>56,77</sup> (**Fig. S22**). We found these associations remained  
244 significant after accounting for  $F_{st}$ , suggesting average allele frequency differences across ancestries cannot solely  
245 explain the observed heterogeneity.

246 Overall, SuShiE recapitulates the findings of primarily shared genetic architectures of molecular traits and show  
247 that effect size heterogeneity is consistent with gene LOF intolerance.

#### 248 Posterior *cis*-molQTL effect sizes improve T/PWAS power in white blood cell traits

249 Lastly, to showcase the downstream benefits of SuShiE, we performed TWAS and PWAS<sup>42-44</sup> on six white blood-  
250 cell-related traits in AOU biobank<sup>50</sup> (average  $n=86,336$ ; **Table S9**). First, we assessed the predictive performance  
251 of SuShiE-based weights compared to alternative expression-prediction methods. Specifically, SuShiE obtained  
252 better cross-validation estimates ( $cv-r^2$ ) compared to SuShiE-Indep, Meta-SuSiE, SuSiE, Elastic Net and gBLUP (2  
253 out of 5 comparisons  $P<0.05$ ) and comparable estimates relative to LASSO ( $P=0.64$ ; **Fig. S23A**). When focusing on  
254 genes with estimated *cis*-molQTL effect size correlation  $<0.9$  across ancestries, we find SuShiE consistently  
255 outperformed other approaches (4 out of 6 comparisons  $P<0.05$ ; **Fig. S23B**), suggesting the benefits in modeling  
256 and learning the prior effect size covariances. We observed significantly decreased prediction performance when  
257 evaluating cross-ancestry prediction (e.g., predicting mRNA expression of AFR using EUR weights; see **Methods**;  
258  $P=1.71e-53$ ; **Fig. S24**), consistent with previous works<sup>22,27,36,81</sup> and further motivating ancestry-matched analyses.  
259 Given this, we predicted the expression levels of 20,515 genes (mRNA) and 573 proteins using ancestry-matched  
260 SuShiE *cis*-molQTL prediction weights from the above analyses and AOU genotypes. Overall, we identified 221  
261 T/PWAS significant associations in white blood count (WBC), eosinophil count (EOS), and monocyte count (MON;  
262 **Table S10**; **Fig. S25**). Of these associations, ~90% were identified in WBC due to substantially increased statistical  
263 power (i.e., 21,476 more participants on average). We found no significant associations in lymphocyte count (LYM),

264 neutrophil count (NEU), and basophil count (BAS), likely due to low detected cell counts, similar to previous  
265 studies<sup>36,82</sup> that identified fewer associations compared to models based on WBC.

266 Consistent with our simulation results (**Fig. 2F**), SuShiE demonstrated higher T/PWAS chi-square statistics and  
267 identified 44 more T/PWAS associations compared to results driven by SuSiE prediction weights (**Fig. 5A**). In  
268 addition, we observed that the SuShiE T/PWAS signals associated with multiple measures of LOF intolerance  
269 (**Table S11**), in contrast to previous work demonstrating that high LOF intolerance genes are typically depleted in  
270 TWAS models due to weak eQTL signals<sup>56,77,80</sup> (**Fig. 5B**; see **Methods**). We found less support for a relationship  
271 between SuSiE-based TWAS signals and LOF intolerance ( $P=9.21e-10$ ; **Table S11**), further demonstrating SuShiE's  
272 advantage. To validate our results, we compared our TWAS statistics with multiple independent white blood cell-  
273 related TWASs<sup>31,36,82-84</sup>. Overall, we found SuShiE-based TWAS replicated at rates similar to SuSiE, suggesting that  
274 its improved power is unlikely due to false positives and further highlighting its benefit in identifying disease-  
275 related genes.

276 Overall, our work has shown that by jointly modeling the molecular data across different ancestries while allowing  
277 effect sizes to differ, SuShiE outputs more accurate *cis*-molQTL prediction weights, thus boosting the downstream  
278 statistical power for integrative analyses with GWASs.

## 279 Discussion

280 Here, we present the Sum of Shared Single Effect approach (SuShiE), a novel approach for multi-ancestry SNP fine-  
281 mapping of molecular traits using a scalable variational approach. SuShiE assumes the joint *cis*-molQTL effects  
282 arise as a linear combination of per-ancestry effect sizes across shared causal variants. Through extensive  
283 simulations, SuShiE first improved the fine-mapping precision in disentangling the causal *cis*-molQTLs from tagging  
284 SNPs by leveraging LD heterogeneity across diverse ancestries. Second, SuShiE accurately learned prior effect size  
285 correlations across ancestries employing a procedure analogous to Empirical Bayes. Third, SuShiE estimated  
286 ancestry-specific *cis*-molQTL prediction weights, boosting findings in the post-GWAS framework (e.g., TWAS and

287 PWAS), compared to the baselines that did not model effect size covariance across ancestries or ignored ancestry  
288 altogether. We applied SuShiE to 36,911 molecular phenotypes of diverse ancestries from three datasets: mRNA  
289 and protein expression from TOPMed-MESA and mRNA expression from GENOA. SuShiE fine-mapped 16% more  
290 genes on average compared to the existing methods, exhibiting smaller credible set sizes and higher enrichment  
291 in relevant functional annotations. SuShiE inferred highly correlated *cis*-molQTL effect sizes across ancestries on  
292 average in significantly heritable genes, reflecting primarily shared *cis*-molQTL architectures. In addition, we  
293 observed *cis*-molQTL effect size heterogeneity across ancestries associated with multiple constraint  
294 measurements, consistent with environmental interactions may partially drive differences in effect sizes across  
295 ancestries. Last, we performed TWAS and PWAS on six white blood cell-related traits from AOU biobank using  
296 SuShiE-derived ancestry-specific *cis*-molQTL prediction weights and identified 44 more significant genes compared  
297 to the existing method. We also observed that SuShiE T/PWAS signals are associated with multiple measures of  
298 LOF intolerance, further showing the benefit of multi-ancestry approaches in identifying genes relevant to  
299 complex disease risk.

300 Next, we describe caveats in our real data analysis. First, SuShiE approximates ancestry as a discrete category,  
301 allowing us to model *cis*-molQTL effect sizes using a multivariate normal distribution (see **Methods**). While this  
302 simplifies modeling and inference tasks, we emphasize that this is a heuristic approach that neglects the complex  
303 and shared demographic histories underlying all humans. Indeed, recent work has demonstrated the importance  
304 of viewing genetic ancestries as a continuous spectrum rather than discrete categories<sup>85</sup>. Relatedly, previous  
305 studies<sup>33,35,37-41</sup> and our simulation results (**Fig. S1**) have shown that increasing the number of ancestries within a  
306 multi-ancestry framework improves fine-mapping precision. However, SuShiE and similar frameworks perform  
307 inference on variants present after filtering on MAF thresholds (e.g., 1%) within each ancestry. As a result, this  
308 requirement can exclude *cis*-molQTLs from analysis due to small sample sizes within an ancestry, suggesting a  
309 trade-off in practice between increasing overall sample size versus excluding informative genetic variants. For  
310 instance, we obtained mRNA expression data measured in EUR (n=402), AFR (n=175), HIS (n=277), and East Asian

311 ancestry (EAS; n=96) individuals from TOPMed MESA study visit-1. From two-ancestry fine-mapping (EUR and AFR)  
312 to three-ancestry (+HIS), we filtered an additional 29 SNPs per gene on average. However, this number increased  
313 to 501 SNPs by including the additional 96 participants of EAS ancestry. As a result, we opted to not include EAS  
314 participants in our analysis in order to maximize the genetic variants analyzed. Modeling genetic ancestry  
315 continuously can potentially avoid this type of *cis*-molQTL loss, thus improving the fine-mapping precision with a  
316 larger sample size.

317 Second, we note that our data consist of African- (AFR) and Hispanic-American (HIS) individuals, which contain  
318 recent admixture events. To account for complex diversity within ancestries, we included genotyping PCs as a  
319 covariate in our models. Several works have suggested that admixture can be sufficiently corrected for using global  
320 ancestry information (i.e., genotyping PCs) in association testing<sup>73,86-91</sup>, especially when causal effect sizes are  
321 largely consistent across ancestries<sup>86,87,89</sup> (**Fig. S16-S18**). On the other hand, accounting for local ancestry may  
322 increase the associating testing power when causal effects are highly different across ancestries<sup>86,87,92</sup> or aid fine-  
323 mapping in post-GWAS analysis<sup>87,89,93</sup>, which can be one of the future directions for SuShiE.

324 Third, we observed significant associations between gene LOF intolerance and several SuShiE-estimated metrics,  
325 including effect size heterogeneity across ancestries, the number of *cis*-molQTLs, *cis*-molQTL distance to TSS, and  
326 functional enrichments. The relationship remained significant after adjusting for  $F_{st}$ , suggesting allele frequency  
327 differences across ancestries are not sufficient to fully explain estimated heterogeneity. As a result, we  
328 hypothesized that *cis*-molQTL effect size heterogeneity could be in part due to gene-by-environment (GxE)  
329 interactions<sup>69,77,94-96</sup>. Highly constrained genes exhibit more complex regulatory landscapes with fewer *cis*-  
330 molQTLs (or apparent *cis*-molQTLs due to smaller effect sizes)<sup>56,77</sup>. As a result, these genes may be less resilient to  
331 environmental perturbations<sup>77</sup>, which may induce effect-size heterogeneity across different ancestries. On the  
332 other hand, it is possible that our  $F_{st}$  estimates are underpowered to detect subtle allele frequency differences  
333 across ancestries. Therefore, these associations may provide indirect evidence for natural selection partially  
334 driving *cis*-molQTL effect size heterogeneity across ancestries. To explicitly investigate the role of selection in

335 molecular differences across ancestries, we likely require a more principled modeling procedure based in  
336 population genetics together with higher-resolution molecular data measured in diverse ancestries<sup>56,80,97–100</sup>. For  
337 instance, recent work has shown the promise of using single-cell data to demonstrate how selection impacts genes  
338 expressed differentially across ancestries<sup>101</sup>.

339 Fourth, SuShiE assumes causal *cis*-molQTLs are shared across ancestries. Our simulations show that SuShiE  
340 remains robust when ancestry-specific *cis*-molQTLs are present (**Fig. S8**). However, in situations where there exist  
341 shared *cis*-molQTLs but ancestries have different sample sizes, SuShiE may prioritize shared *cis*-molQTLs along  
342 with SNPs tagged in LD of the ancestry with larger sample sizes, evidenced through simulations (**Fig. S5B**). However,  
343 through our case study in *URGCP* (**Fig. 4**), we observed relatively higher signals in AFR but not in EUR and HIS,  
344 despite AFR having the smallest sample size, suggesting this limitation may be minimal overall.

345 Lastly, in our T/PWAS analysis, we selected six white blood-cell related traits to best match PBMC and LCL contexts.  
346 However, alternative cell-types not included in our analyses may better capture relevant contexts. For example,  
347 PBMCs and LCLs do not contain neutrophils, basophils, and eosinophils, and LCLs additionally do not include  
348 monocytes, which may result in a loss in statistical power. As single-cell RNA-seq datasets become more  
349 available<sup>102</sup>, one possible direction would be to perform TWAS in fine-grained cellular contexts and backgrounds.  
350 In addition, after predicting expression levels using ancestry-matched weights for each individual, we performed  
351 individual-level T/PWAS by concatenating the predicted expression levels across ancestries rather than perform  
352 ancestry-specific TWAS followed by meta-analysis<sup>103</sup>. The premise of the meta-analysis approach is that  
353 researchers obtain ancestry-specific GWAS and then integrate with corresponding eQTL weights. Because the  
354 causal genes for complex traits are likely shared across ancestries<sup>20,21,23,36,68–74</sup>, a regression framework with  
355 individual-level data concatenated across ancestries (the largest sample size) can maximize power.

356 We briefly discuss potential directions for future work. First, recent studies have shown that incorporating  
357 functional annotation in the prior distribution can improve the fine-mapping precision<sup>33,79,104</sup>. SuShiE currently  
358 employs a uniform distribution for prior causal probability. Including functionally-informed priors is likely to

359 improve its performance further. Second, SuShiE fine-maps individual-level molecular and genotypic data in a  
360 prespecified locus flanking the TSS and TES regions of a gene. In theory, users can apply SuShiE on individual-level  
361 complex trait data, however, this likely will require additional analyses (e.g., pre-specifying GWAS significant loci)  
362 and care in controlling for genome-wide backgrounds and population structure. In addition, the limited  
363 accessibility to the individual-level complex trait data allows the extension of SuSiE-like models to be compatible  
364 with summary statistics<sup>38,39,41,51</sup>, which typically requires external LD reference panels. As more *cis*-molQTL  
365 summary statistics are available to the community<sup>4,102</sup>, we foresee a potential demand to implement this  
366 compatibility in SuShiE. Last, SuShiE currently cannot model molecular data in their original read-count format,  
367 which is usually transformed to a continuous scale (i.e., inverse normal transformation). Extending SuShiE to a  
368 GLM-like model naturally would encompass this scenario and present an exciting direction for SuShiE.  
369 Overall, SuShiE, together with its application on large-scale molecular data of diverse ancestries, identifies more  
370 *cis*-regulatory mechanisms and reveals its genetic architecture. We anticipate considerable demand for our  
371 approach in the genetics field characterized by forthcoming multi-ancestry and multi-omics research.

## 372 Online Methods

### 373 Sum of Shared Single Effects Model

374 Here, we describe the statistical model underlying SuShiE (see **Supplementary Note** for a detailed description).  
375 SuShiE assumes *cis*-molQTLs are present in *all* ancestries, defined as shared *cis*-molQTLs while allowing for effect  
376 sizes at causal *cis*-molQTLs to covary across ancestries a-priori. For the  $i^{\text{th}}$  of total  $k$  ancestries, SuShiE models the  
377 centered and standardized levels of a molecular trait  $\mathbf{g}_i$  measured in  $n_i$  individuals as a linear combination of  $p$   
378 genotyped variants  $\mathbf{X}_i$  as

$$379 \quad \mathbf{g}_i = \mathbf{X}_i \boldsymbol{\beta}_i + \boldsymbol{\epsilon}_i$$



380 where  $\boldsymbol{\beta}_i$  is a  $p \times 1$  vector of ancestry-specific *cis*-molQTL effects, and  $\boldsymbol{\epsilon}_i \sim N(0, \sigma_{i,e}^2 \mathbf{I}_{n_i})$  is environmental noise.

381 In addition, we model  $\boldsymbol{\beta}_i = \sum_{l=1}^L \boldsymbol{\beta}_{i,l}$  as the sum of  $L$  effects  $\boldsymbol{\beta}_{i,l} = \boldsymbol{\gamma}_l \cdot b_{i,l}$  where  $\boldsymbol{\gamma}_l$  is a  $p \times 1$  binary vector

382 indicating which variant is the shared *cis*-molQTL for the  $l^{\text{th}}$  effect while allowing ancestry-specific effect sizes  $b_{i,l}$ .

383 Furthermore, we model  $\boldsymbol{\gamma}_l \sim \text{Multi}(1, \boldsymbol{\pi})$  where  $\boldsymbol{\pi}$  is a  $p \times 1$  vector representing prior probability for each SNP

384 to be a *cis*-molQTL, and model  $b_l = [b_{1,l} \cdots b_{i,l} \cdots b_{k,l}] \sim N(\mathbf{0}, \mathbf{C}_l)$  where

$$385 \quad \mathbf{C}_{i,i',l} = \begin{cases} \sigma_{i,b,l}^2 & \text{if } i = i' \\ \rho_{i,i',l} \sigma_{i,b,l} \sigma_{i',b,l} & \text{otherwise} \end{cases}$$

386  $\mathbf{C}_l$  is the  $l^{\text{th}}$  prior  $k \times k$  effect size covariance matrix with  $\sigma_{i,b,l}^2$  as variance, and  $\rho_l$  as correlation.

### 387 Variational inference of model parameters

388 To infer the *cis*-molQTL effects, we seek to estimate the posterior distribution of  $\Pr(\boldsymbol{\beta}_1, \dots, \boldsymbol{\beta}_k | \mathbf{Data}) =$

389  $\Pr(\boldsymbol{\beta}_1, \dots, \boldsymbol{\beta}_k | \mathbf{g}_1, \dots, \mathbf{g}_k, \mathbf{X}_1, \dots, \mathbf{X}_k, \mathbf{C}, \boldsymbol{\pi}, \sigma_{1,e}^2, \dots, \sigma_{k,e}^2)$  where  $\mathbf{C} = \{\mathbf{C}_1, \dots, \mathbf{C}_L\}$ . We regard  $\boldsymbol{\beta}_*$  as latent

390 variables,  $\mathbf{g}_*$ , and  $\mathbf{X}_*$ , as observed data, and  $\mathbf{C}$ ,  $\boldsymbol{\pi}$ , and  $\sigma_{*,e}^2$  are the hyperparameters. However, inferring the exact

391 distributions of latent variables is computationally intractable due to non-conjugacy with the prior distribution.

392 Therefore, we seek a surrogate distribution  $Q(\boldsymbol{\beta}_1, \dots, \boldsymbol{\beta}_k)$ , which minimizes the Kullback–Leibler (KL) divergence

393 with  $\Pr(\boldsymbol{\beta}_1, \dots, \boldsymbol{\beta}_k | \mathbf{Data})$ . Specifically, we have:

$$394 \quad D_{KL}(Q(\boldsymbol{\beta}_1, \dots, \boldsymbol{\beta}_k) || \Pr(\boldsymbol{\beta}_1, \dots, \boldsymbol{\beta}_k | \mathbf{Data})) \\ 395 \quad = \log \Pr(\mathbf{g}_1, \dots, \mathbf{g}_k | \mathbf{X}_1, \dots, \mathbf{X}_k, \mathbf{C}, \boldsymbol{\pi}, \sigma_{1,e}^2, \dots, \sigma_{k,e}^2) \\ 396 \quad - E[\log \Pr(\boldsymbol{\beta}_1, \dots, \boldsymbol{\beta}_k, \mathbf{g}_1, \dots, \mathbf{g}_k | \mathbf{C}, \boldsymbol{\pi}, \sigma_{1,e}^2, \dots, \sigma_{k,e}^2) - \log Q(\boldsymbol{\beta}_1, \boldsymbol{\beta}_2)]$$

397 where the first term is the log evidence, and the expectation term is the evidence lower bound (ELBO). Since the

398 log evidence is constant with respect to model variables, minimizing the KL divergence is equivalent to maximizing

399 the ELBO. Furthermore, to limit the universe of possible forms that the surrogate distribution  $Q(\boldsymbol{\beta}_1, \dots, \boldsymbol{\beta}_k)$  may

400 take, we impose an additional mean-field assumption<sup>105</sup>. Namely, SuShiE assumes that each of the  $L$  shared

401 effects  $\beta_l$  are mutually independent under  $Q$ :

402 
$$Q(\boldsymbol{\beta}_1, \dots, \boldsymbol{\beta}_k) = \prod_{l=1}^L Q(\boldsymbol{\beta}_{1,l}, \dots, \boldsymbol{\beta}_{k,l}) = \prod_{l=1}^L Q(\mathbf{b}_l | \boldsymbol{\gamma}_l) Q(\boldsymbol{\gamma}_l).$$

403 Therefore, to approximate the posterior distributions  $Q(\cdot)$  for latent variables  $\mathbf{b}_{l,j}$  (a  $k \times 1$  vector) and  $\boldsymbol{\gamma}_{l,j}$  (a  
404 scalar) at SNP  $j \in [1, p]$  of  $l^{\text{th}}$  shared effect, we need to compute the expectation of complete data log-likelihood  
405  $L(\boldsymbol{\beta}_1, \dots, \boldsymbol{\beta}_k, \mathbf{g}_1, \dots, \mathbf{g}_k | \mathbf{C}, \boldsymbol{\pi}, \sigma_{1,e}^2, \dots, \sigma_{k,e}^2)$  (i.e., the joint distribution) while holding other variables constant.  
406 Through the principles of coordinate-ascent variational inference (CAVI)<sup>105</sup>, we can identify each  $Q(\cdot)$  surrogate  
407 as,

408 
$$Q(\mathbf{b}_{l,j} | \boldsymbol{\gamma}_{l,j} = 1) = N(\mathbf{b}_{l,j} | \boldsymbol{\mu}_{l,j}, \boldsymbol{\Sigma}_{l,j})$$
  
409 
$$Q(\boldsymbol{\gamma}_{l,j} = 1) \propto \text{softmax}(\log \boldsymbol{\pi}_j - \log N(\boldsymbol{\mu}_{l,j} | \mathbf{0}, \boldsymbol{\Sigma}_{l,j}))$$
  
410 
$$Q(\boldsymbol{\gamma}_l) = \text{Multi}(\boldsymbol{\gamma}_l | \mathbf{1}, \boldsymbol{\alpha}_l)$$

411 where  $\boldsymbol{\mu}_{l,j} \in \mathbb{R}^{k \times 1}$  and  $\boldsymbol{\Sigma}_{l,j} \in \mathbb{R}^{k \times k}$  are the corresponding posterior mean and covariance, and  $\boldsymbol{\alpha}_l \in \mathbb{R}^{p \times 1}$  is  
412 each SNP's posterior probability to explain the  $l^{\text{th}}$  effect. We provide the complete mathematical derivations,  
413 inference algorithms, and detailed definitions in the **Supplementary Note**.

#### 414 [Computing posterior inclusion probability and \$\eta\$ -credible sets](#)

415 We define the posterior inclusion probability (PIP) for SNP  $j$  with  $\boldsymbol{\alpha}_1, \dots, \boldsymbol{\alpha}_L$  as  $\text{PIP}_j := 1 - \prod_{l=1}^L (1 - \alpha_{l,j})$ . To  
416 compute an  $\eta$ -credible set for each  $L$ , where  $\eta$  represents the desired probability that the set contains *cis*-  
417 molQTLs, we decreasingly sort  $\boldsymbol{\alpha}_l$  and take a greedy approach to include SNPs until their cumulative sum exceeds  
418  $\eta$ .

419 In the case that the inferred number of effects  $L$  surpasses the actual number of *cis*-molQTLs, the unnecessary  
420 credible sets will contain most SNPs with low posterior probability close to  $\alpha_{l,j} = 1/p$ , where  $p$  is the number of  
421 SNPs. To refine the final inference results, we remove the credible sets whose lowest absolute pairwise correlation,  
422 which is defined as “purity”<sup>15</sup> and weighted by sample size across all ancestries, among SNPs is less than 0.5. In  
423 practice, following the previous work<sup>15</sup>, we empirically specify  $L$  as 10.

## 424 Inferring cross-ancestry effect size correlations

425 SuShiE features the capability to estimate the correlation of *cis*-molQTL effect sizes across multiple ancestries. For  
426 some gene  $t$ , SuShiE by default outputs  $L$  estimates of the effect size correlation  $\hat{\rho}_{t,1}, \dots, \hat{\rho}_{t,L}$  for each credible set.  
427 If we apply SuShiE to  $T$  genes in total, we empirically recommend computing effect size correlation across  
428 ancestries with  $\hat{\rho} = \frac{1}{T} \sum_{t=1}^T \hat{\rho}_{t,1}$ .

## 429 Simulating genotypes and quantitative molecular traits

430 To evaluate SuShiE's performance in simulations, we first simulated genotypes and quantitative molecular traits  
431 to mimic the real-world scenarios using our previous simulation frameworks<sup>36,106,107</sup>. To simulate genotype data,  
432 we used LD estimates from individuals of European (EUR;  $n=489$ ), African (AFR;  $n=639$ ), and East Asian (EAS;  $n=481$ )  
433 ancestries from the 1000 Genomes Project (1000G) phase three data<sup>108</sup>. We limited LD to biallelic HapMap SNPs<sup>109</sup>,  
434 discarded those with missingness ( $>1\%$ ), MAF ( $<1\%$ ), and violated Hardy-Weinberg equilibrium (HWE mid-adjusted  
435  $P < 1e-6$ ). We obtained chromosome, transcription start site (TSS), and transcription end site (TES) information for  
436 19,279 protein-coding autosomal genes using GENCODE release 26 (GRCh37)<sup>110</sup>. We extended each gene 500,000  
437 base pairs (bp) upstream of TSS and 500,000 bp downstream of TES, and randomly selected 500 genes that have  
438 at least 500 common SNPs across EUR, AFR, and EAS genotypes.

439 We first focused on simulations using EUR and AFR ( $k = 2$ ). At each gene, we simulated centered and standardized  
440 genotype matrix  $\mathbf{X}_i \in \mathbb{R}^{n_i \times p}$  for  $i^{\text{th}}$  ancestry using a multivariate normal distribution  $N(0, \mathbf{V}_i)$  where  $n_i \in$   
441  $\{200, 400, 600, 800\}$  is the *cis*-molQTL study sample size,  $p$  is the number of common SNPs across ancestries in  
442 the locus, and  $\mathbf{V}_i \in \mathbb{R}^{p \times p}$  is the ancestry-specific LD matrix estimated from 1000G genotypes<sup>108</sup>. Next, we  
443 uniformly chose  $m \in \{1, 2, 3\}$  out of  $p$  common SNPs as *cis*-molQTLs and simulated their ancestry-specific effect  
444 sizes  $\tilde{\boldsymbol{\beta}}_1, \tilde{\boldsymbol{\beta}}_2 \in \mathbb{R}^{m \times 1}$  under a bivariate normal distribution as

445 
$$(\tilde{\boldsymbol{\beta}}_1, \tilde{\boldsymbol{\beta}}_2) \sim N \left( \mathbf{0}, \begin{bmatrix} h_{g,1}^2 & \rho \cdot \sqrt{h_{g,1}^2 \cdot h_{g,2}^2} \\ \rho \cdot \sqrt{h_{g,1}^2 \cdot h_{g,2}^2} & h_{g,2}^2 \end{bmatrix} / m \right) \otimes I_m,$$

446 where  $h_{g,i}^2 \in \{0.01, 0.05, 0.1, 0.2\}$  is the proportion of variance in gene expression explained by *cis*-molQTLs (i.e.,  
 447 *cis*-SNP heritability of the molecular trait) and  $\rho \in \{0.01, 0.4, 0.8, 0.99\}$  is the effect size correlation. Then, we  
 448 constructed effect-size vectors  $\boldsymbol{\beta}_1$  and  $\boldsymbol{\beta}_2$ , where  $\tilde{\boldsymbol{\beta}}_1$  and  $\tilde{\boldsymbol{\beta}}_2$  are the  $m$  non-zero entries at the same index  
 449 representing shared *cis*-molQTLs and the rest  $p - m$  entries are zero representing the null SNPs. Next, we  
 450 computed the quantitative molecular traits  $\boldsymbol{g}_i$  using  $\boldsymbol{X}_i \boldsymbol{\beta}_i + \boldsymbol{\epsilon}_i$  where  $\boldsymbol{\epsilon}_i \sim N \left( \mathbf{0}, s_{g,i}^2 \left( \frac{1}{h_{g,i}^2} - 1 \right) \boldsymbol{I}_{n_i} \right)$  is the random  
 451 environmental noise and  $s_{g,i}^2 = \boldsymbol{\beta}_i^T \boldsymbol{V}_i \boldsymbol{\beta}_i$  is the genetic variance after accounting for LD. To reflect cases where  
 452 heterogeneity exists in the genetic architecture of molecular traits across ancestries<sup>31,72</sup>, we allowed *cis*-SNP  
 453 heritability to be ancestry-specific with  $h_{g,1}^2 = 0.05$  and  $h_{g,2}^2 \in \{0.01, 0.05, 0.1, 0.2\}$ ; we also evaluated the  
 454 performance under different statistical power where  $n_1 = 400$  and  $n_2 \in \{200, 400, 600, 800\}$ . To determine  
 455 whether incorporating additional ancestry improves SuShiE's performance, we simulated the genotypic and  
 456 phenotypic data for EAS with the same total sample sizes and genetic architecture.

457 In addition, we simulated two cases under model misspecification. We first evaluated SuShiE's performance when  
 458 ancestry-specific *cis*-molQTLs exist, we simulated  $m_{i,AS} \in \{1, 2, 3\}$  *cis*-molQTLs for both ancestries in addition to  
 459 shared *cis*-molQTLs  $m = 2$  while fixing  $h_{g,i}^2 = 0.05$  for ancestry  $i$ . Second, to reflect cases where the number of  
 460 shared *cis*-molQTL ( $m$ ) is different from inferred  $L$  by fixing  $m = 2$  and varying the inferred  $L \in \{2, 5, 10\}$ .

## 461 Default parameters and performance metrics

462 We performed SNP fine-mapping using SuShiE on simulated genotypes and molecular data across EUR and AFR  
 463 individuals. In terms of variational inference parameters, we specified  $L \in \{1, 2, 3\}$  to match the actual number of  
 464 simulated effects and initialized *cis*-molQTL effects  $\hat{\boldsymbol{b}}_{l,j}$  as  $\begin{bmatrix} 0 \\ 0 \end{bmatrix}$ , their covariance matrix  $\hat{C}_l$  as  $\begin{bmatrix} 0.001 & 0.1 \\ 0.1 & 0.001 \end{bmatrix}$ , the

465 prior estimates of environmental noises  $\hat{\sigma}_{i,e}^2$  as 0.001, the prior probability for SNPs to be *cis*-molQTLs as  $1/p$   
466 where  $p$  is the number of common SNPs.

467 To evaluate the gain in parametrizing the effect size correlation across ancestries, we compared our method  
468 SuShiE to “SuShiE-Indep” which assumes the *cis*-molQTL effect sizes are independent across ancestries; that is,  
469 we fixed the effect size correlation prior  $\rho = 0$ , and did not learn it through the Empirical-Bayes-like procedure.

470 To demonstrate that SuShiE’s improvement does not result from the accumulation of samples across ancestries,  
471 we compared SuShiE’s performance to two “baseline” methods: first, we performed single-ancestry SuSiE and  
472 then meta-analyzed the resulting PIPs by  $PIP_{\text{meta}} = 1 - (1 - PIP_{\text{EUR}}) \cdot (1 - PIP_{\text{AFR}})$ ; we refer to this method as  
473 “meta-SuSiE”. Second, we row-stacked the genotype matrices and molecular trait vectors across ancestries and  
474 then performed single-ancestry SuSiE as “SuSiE.” Overall, we performed four methods (SuShiE, SuShiE-Indep,  
475 meta-SuSiE, and SuSiE) on 500 genes’ simulated genotypes and molecular traits to output corresponding PIPs,  
476 credible sets, and ancestry-specific effect size estimates. We varied four parameters: per-ancestry *cis*-molQTL  
477 study sample size ( $n_i$ ), the number of *cis*-molQTLs ( $m$ ), the *cis*-SNP heritability of molecular traits ( $h_g^2$ ) for each  
478 ancestry, and the effect size correlation ( $\rho$ ). To reflect a practical study design, the default parameters were fixed  
479 at  $n_i = 400$ ,  $m = 2$ ,  $h_g^2 = 0.05$ , and  $\rho = 0.8$  unless stated otherwise. Furthermore, we evaluated the fine-  
480 mapping performance with three metrics across 500 simulated genes: PIPs at causal *cis*-molQTLs, credible set size,  
481 and frequency that causal *cis*-molQTLs are contained in 90% credible sets (calibration). We computed the metrics  
482 of meta-SuSiE based on the union of the credible sets across two single-ancestry SuSiE. As different methods may  
483 or may not prune credible sets at the same simulated gene, to show a fair comparison, we computed the credible  
484 set size metric only using the credible set that none of the four methods pruned out. To compare metrics across  
485 methods, we ran linear regression adjusted for relevant simulation parameters and reported one-sided Wald test  
486 P values.

## 487 Simulating GWAS and TWAS

488 Transcriptome-wide Association Studies (TWASs) leverage GWAS summary statistics, eQTL prediction weights,  
489 and LD reference to identify genes whose predicted expression levels are associated with complex traits<sup>42-44</sup>. A  
490 more accurate eQTL prediction weight will increase the power of the TWAS framework. Therefore, we compared  
491 the prediction weights inferred by SuShiE to other methods: SuShiE-Indep, Meta-SuSiE, SuSiE, least absolute  
492 shrinkage and selection operator (LASSO)<sup>53</sup>, elastic net regularization (Elastic Net)<sup>54</sup>, and genomic best linear  
493 unbiased prediction (gBLUP)<sup>55</sup>. We simulated the expression and genotype data for the training and testing set  
494 separately, with the same method mentioned in the previous sections. For the training set, we varied the per-  
495 ancestry sample size  $n_t \in \{200, 400, 600, 800\}$  and set the out-of-sample testing set sample size  $n_v = 200$ . Then,  
496 we predicted the expressions using ancestry-matched fitted weights on testing genotype data, and computed the  
497 coefficients of determination ( $r^2$ ) between the predicted and simulated expression. For Meta-SuSiE, we trained  
498 the prediction weights for each ancestry using per-ancestry sample size. For SuSiE, LASSO, Elastic Net, and gBLUP,  
499 we trained the prediction weights after concatenating data across ancestries to guarantee that the total sample  
500 sizes were the same as SuShiE as fair comparisons.

501 To showcase that SuShiE's prediction weights introduce more power in TWAS, we simulated GWAS summary  
502 statistics and computed TWAS statistics using different prediction weights. First, because individuals in GWASs  
503 are usually different from ones in the eQTL studies, we re-simulated the genotype matrix  $\mathbf{X}_{\text{GWAS},i} \in \mathbb{R}^{n_{\text{GWAS},i} \times p}$   
504 where  $n_{\text{GWAS},i}$  is the GWAS sample size for ancestry  $i$  using the same generating approach above. Then, we used  
505 the eQTL effect size vectors  $\boldsymbol{\beta}_i$  generated in the previous section to simulate a complex trait  $\mathbf{y}_i \in \mathbb{R}^{n_{\text{GWAS},i} \times 1}$  as  
506 a linear combination of expression levels  $\mathbf{g}_i \in \mathbb{R}^{n_{\text{GWAS},i} \times 1}$  as

$$507 \quad \mathbf{y}_i = \mathbf{g}_i \delta + \mathbf{e}_i = \mathbf{X}_{\text{GWAS},i} \boldsymbol{\beta}_i \delta + \mathbf{e}_i,$$

508 where  $\delta \sim N(0,1)$  is the gene expression effect on the complex trait,  $\mathbf{e}_i \sim N(0, s_i^2 \left( \frac{1}{h_{GE}^2} - 1 \right) I_{n_{\text{GWAS},i}})$  is the random  
509 noises for the complex traits,  $s_i^2 = \boldsymbol{\beta}_i^T \mathbf{V}_i \boldsymbol{\beta}_i \delta^2$ ,  $\mathbf{V}_i$  is the LD matrix generated from 1000G<sup>108</sup>, and  $h_{GE}^2 \in$

510  $\{6 \times 10^{-5}, 1.5 \times 10^{-4}, 3 \times 10^{-4}, 6 \times 10^{-4}\}$  is the proportion of variation of the complex trait explained by the  
511 expression of a single gene<sup>111</sup>. Then, we regressed the complex trait  $\mathbf{y}_i$  on each SNP in  $\mathbf{X}_{\text{GWAS},i}$  marginally to  
512 compute the GWAS summary statistics  $\mathbf{z}_{\text{GWAS},i} \in \mathbb{R}^p \times 1$ . Last, we computed TWAS summary statistics with  
513  $\mathbf{z}_{\text{TWAS},*i} = \frac{\mathbf{w}_{*,i}^T \mathbf{z}_{\text{GWAS},i}}{\mathbf{w}_{*,i}^T \mathbf{V}_i \mathbf{w}_{*,i}}$  along with its P value where  $\mathbf{w}_{*,i}$  is the prediction weights fitted by different methods. We  
514 define the TWAS power as the frequency of the Bonferroni-corrected P value is less than 0.05.

## 515 Overview of real-data analyses

516 We applied SuShiE and other methods (e.g., SuShiE-Indep, Meta-SuSiE, and SuSiE) to three datasets: mRNA (visit-  
517 1) measured in peripheral blood mononuclear cells (PBMCs) and protein expression measured in plasma of three  
518 EUR, AFR, and HIS ancestries from Trans-Omics for Precision Medicine program Multi-Ethnic Study of  
519 Atherosclerosis (TOPMed MESA)<sup>48,49</sup> and mRNA expression measured in lymphoblastoid cell lines (LCLs) of EUR  
520 and AFR ancestries from the Genetic Epidemiology Network of Arteriopathy (GENOA) study<sup>26</sup>. We excluded the  
521 mRNA expression levels data measured in T cells and monocytes from TOPMed MESA study due to relatively  
522 smaller sample sizes. We explain the detailed quality control (QC) procedure in the sections below. We conducted  
523 pairwise comparisons of methods on four basic summary statistics, focusing on the genes for which both methods  
524 output credible sets; the summary statistics included the number of genes identified with *cis*-molQTLs (*e/p*Genes),  
525 the average PIPs of the SNPs in the credible sets, the average single-effect-specific credible set sizes, and the  
526 frequency of having genes whose credible sets contained SNPs with PIPs greater than 0.9. We defined the number  
527 of *cis*-molQTLs as the number of credible sets output after pruning for purity (see previous section for the  
528 definition). Next, we performed enrichment analyses using 89 functional annotations and a case study focusing  
529 on a gene that was only identified by SuShiE, and missed by other methods. Last, using SuShiE-derived ancestry-  
530 specific *cis*-molQTL effect sizes, we performed individual-level TWAS and PWAS with All Of Us (AOU) biobank<sup>50</sup>  
531 individuals and compared to the results derived from SuSiE.

532 To validate SuShiE's results on the three main datasets mentioned above, we applied SuShiE and other methods  
533 to three separate datasets: mRNA expression (visit-5) measured in PBMC of EUR, AFR, and HIS ancestries from  
534 TOPMed MESA, protein expression measured in plasma of EUR ancestry from INTERVAL study<sup>5</sup>, and the mRNA  
535 expression measured in LCL of EUR and Yoruba in Ibadan (YRI) ancestries from the GEUVADIS study<sup>61</sup>. We  
536 computed two statistics to evaluate validation performance: first, focusing on the *cis*-molQTLs of e/pGenes  
537 identified by SuShiE, the percentage for which SuShiE identified *cis*-molQTLs in the validation datasets. Second,  
538 focusing on the credible sets for which we identified the same *cis*-molQTLs in both main and validation studies,  
539 we computed the cosine similarity of posterior probabilities ( $\alpha_l$ ) to see whether they prioritized the same SNPs.  
540 For SNPs that are not in the overlap between main and validation studies, we manually assigned them a value of  
541 0 for cosine similarity calculation. For each credible set, we randomly shuffled the  $\alpha_l$  in validation studies 500  
542 times to construct the null distribution of the cosine similarity and compute its z score. We computed the average  
543 cosine similarity and z scores across all credible sets as an aggregation estimate and its corresponding significance.  
544 For all the fine-mapping analysis, we used the SNPs that are shared across ancestries on the genomic window of  
545 each gene that is 500,000 bp upstream and downstream of each gene's TSS and TES (one million bp in total),  
546 respectively, based on the GENCODE v34<sup>110,112</sup>. In addition, we only included genes that are located on the  
547 autosomes, do not overlap with the major histocompatibility complex (MHC) region, have more than 100 SNPs on  
548 the genomic window present in all ancestries, and whose ENSEMBL gene IDs match the records in GENCODE  
549 v34<sup>110,112</sup>. We adjusted for covariates by regressing them from both mRNA/protein levels and each SNP. In addition,  
550 we computed the *cis*-SNP heritability using the limix python package (see **Code Availability**) for each analyzed  
551 molecule within each ancestry. We used PLINK2.0, vcftools, and bcftools for genotype manipulation<sup>113-116</sup>.

## 552 Genotype data in the TOPMed MESA study

553 We obtained the whole-genome sequencing (WGS) data (freeze 9; GRCh 38) of 5,379 individuals from the  
554 TOPMed MESA<sup>48,49</sup>. Specifically, we removed the SNPs with the following criteria: both duplicate genotype



555 discordance and mendelian genotype discordance are greater than 2%, genotype missing rate at depth 10 is  
556 greater than 2%, Milk-SVM score for variant quality is less than -0.5, variants that overlap with centromeric regions,  
557 HWE p-value is less than  $1e-6$ , and MAF is less than 1%, resulting in a total of 125,089,612 SNPs. In addition, we  
558 computed the genotype principal components (PCs) with SNPs that are pruned for LD using PLINK2.0 (--indep-  
559 pairwise 200 1 0.3)<sup>113,114,117</sup>. Last, we retained individuals who are self-identified as EUR, AFR, or HIS ancestries and  
560 have measurements in mRNA (both visits 1 and 5) and protein datasets, resulting in a total of 1,292 individuals.

### 561 mRNA expression data in the TOPMed MESA study

562 We obtained RNA-seq data in gene-level read counts and reads per kilobase of transcript per million mapped  
563 reads (RPKM) of 57,615 genes for 2,137 samples (both visits-1 and visit-5) measured in PBMC using RNA-SeqQC  
564 v2.0.0 from the TOPMed MESA study. The data was pre-processed based on the TOPMed RNA-seq harmonization  
565 pipeline (see **Code Availability**). We first calculated the gene expression PCs on all samples' read counts using the  
566 PCA function of the scikit-learn package<sup>118</sup>, and normalized it across all samples within each PC. Then, focusing on  
567 the samples measured in visit-1, we followed the GTEx<sup>3</sup> eQTL analysis preparation script to select gene whose  
568 transcript per million (TPM) is  $>0.1$  and raw read counts  $>6$  reads in at least 20% of samples (see **Code Availability**).  
569 For individuals with replicate samples, we only kept one sample with the greatest sum of reads across all genes;  
570 we also removed individuals with whom we did not have self-identified ancestry information, resulting in 402 EUR,  
571 175 AFR, and 277 HIS individuals. Then, within each ancestry, we normalized expression levels between samples  
572 using `edger_cpm` function in the `pyqtl` package, (see **Code Availability**) with `normalized_lib_sizes=True`, which is  
573 a Python implementation of `edgeR`<sup>119</sup> ; we next performed inverse-rank normalization using the  
574 `inverse_normal_transform` function. Last, focusing on 21,747 genes filtered based on inclusion criteria and using  
575 SNPs whose MAF  $>1\%$  and HWE mid-adjusted  $P > 1e-6$  within each ancestry, we ran SuShiE and other methods  
576 using SNPs on the genomic windows of each gene, adjusting for 15 gene expression PCs, 10 genotype PCs, age,  
577 sex, and the assay lab. We did not include individuals who self-identified as East Asian in TOPMed MESA study due

578 to the small sample size (n=96). We removed SNPs based on MAF <1%, and including EAS participants would  
579 exclude 501 more SNPs on average per gene from downstream analyses.

## 580 Protein expression data in the TOPMed MESA study

581 We obtained the protein expression levels of 1,317 target proteins for 1,966 samples (both visits-1 and -5) from  
582 the TOPMed MESA study using SOMAscan, an aptamer-based technology. First, we computed the protein  
583 expression PCs on all samples using the PCA function of the scikit-learn package<sup>118</sup>, and normalized it across all  
584 samples for each PC. Then, focusing on the samples measured in visit-1, we removed individuals with whom we  
585 did not have self-identified ancestry information, resulting in 398 EUR, 297 AFR, and 261 HIS individuals. Within  
586 each ancestry, we inverse-rank normalized the protein expression data using the `inverse_normal_transform`  
587 function in the `pyqtl` package (see **Code Availability**). As some proteins may be targeted by multiple aptamers,  
588 which correspond to different isoforms of proteins<sup>120</sup>, we regarded each isoform as a unique protein. As a result,  
589 we obtained 1,274 proteins based on gene inclusion criteria and performed fine-mapping using SuShiE and other  
590 methods on the genomic windows adjusted for 15 protein expression PCs, 10 genotype PCs, sex, and age, using  
591 SNPs whose MAF > 1% and HWE mid-adjusted  $P > 1e-6$  within each ancestry.

## 592 Genotype and mRNA expression data in the GENOA study

593 From the GENOA study<sup>26</sup>, we obtained paired genotype and LCL mRNA expression data of 373 EUR and 441 AFR  
594 individuals, together with corresponding covariates, processed by previous works<sup>26,36</sup>. Briefly, we restricted  
595 TOPMed-imputed<sup>121</sup> genotype data on biallelic SNPs with imputation score  $r^2 > 0.6$ , MAF >1%, and HWE mid-  
596 adjusted  $P > 1e-6$  within each ancestry. Focusing on 14,797 genes based on gene inclusion criteria, we performed  
597 fine-mapping on the genomic window, adjusted for 30 gene expression PCs, five genotype PCs, age, sex, and  
598 genotyping platform.

## 599 Genotype and molecular data in three validation datasets

600 To validate SuShiE's results of PBMC mRNA expression (visit-1) in TOPMed MESA<sup>48,49</sup>, we used the mRNA  
601 expression data measured in PBMC of the same study but collected from visit-5, a 10-year-later follow-up visit.  
602 We performed the identical pipeline mentioned in the previous section, resulting in 21,695 genes (21,240  
603 overlapped with visit-1) from 422 EUR, 168 AFR, and 285 HIS individuals.

604 To validate the plasma protein expression results in TOPMed MESA, we obtained the inverse-rank normalized  
605 protein expression levels of 3,301 EUR individuals measured in plasma from the INTERVAL study<sup>5</sup>. The genotype  
606 data was pre-processed, imputed, and annotated with dbSNP v153 by previous studies<sup>5,122,123</sup>. We obtained 3,187  
607 ENSEMBLE-UniProt-SOMAmer ID triplets (1,313 overlapped with the TOPMed MESA) based on gene selection  
608 criteria and performed single-ancestry SuSiE fine-mapping on the genomic window, adjusted for sex, age, duration  
609 between blood draw and process, 3 genotype PCs, and subcohort, and 5 expression PCs, using SNPs whose MAF >1%  
610 and HWE mid-adjusted  $P > 1e-6$ .

611 To validate the mRNA expression data measured in LCLs from the GENOA study, we obtained paired genotype and  
612 gene expression data measured in LCLs in gene-level RPKM of 23,722 genes for 373 EUR and 89 YRI individuals  
613 from the GEUVADIS study<sup>61</sup>. First, we computed the expression PCs on all the individuals using the PCA function  
614 of the scikit-learn package<sup>118</sup>. Then, we kept high-expressed genes whose TPM >0.1 in at least 20% of all the  
615 individuals<sup>3</sup> and filtered based on gene selection criteria, resulting in a total of 19,882 genes (10,439 overlapped  
616 with GENOA). Last, using SNPs whose MAF >1% and HWE mid-adjusted  $P > 1e-6$  within each ancestry, we  
617 performed SuShiE fine-mapping on the genomic window, adjusted for sex, five expression PCs, and five genotype  
618 PCs, which is calculated on the LD-pruned pipeline defined in the previous section.

## 619 Functional enrichment analyses and case study

620 We ran functional enrichment analysis only on the genes identified with *cis*-molQTLs (i.e., SuShiE outputs credible  
621 sets; e/pGenes). To visualize the relationship between the PIPs inferred by SuShiE and their distance to the TSS,

622 we grouped fine-mapped SNPs into 2,000 bins that are 500 bp long to cover the one-million-bp window around  
623 the TSS for each gene and computed the average PIPs within each bin. To visualize the relationship between single  
624 effects' posterior probabilities and their distance to the TSS, we performed the same procedure focusing on the  
625 shared effects that had credible set output (i.e., passed the purity threshold; see previous method section).

626 We performed enrichment analysis using 89 functional annotations. First, we downloaded 5 candidate cis-  
627 regulatory elements (cCREs) from ENCODE Registry v3<sup>58</sup>. Then, we obtained 9 cell-type specific cCREs measured  
628 in PBMC using snATAC-Seq<sup>59</sup> and one cCRE measured in frozen PBMC using scATAC-seq<sup>60</sup>. Last, we obtained the  
629 74 categorical functional annotations from LDSC baseline annotations v2.2<sup>124,125</sup>, and remapped to GRCh38 using  
630 LiftOver (see **Code Availability**). To compute the functional enrichment scores, we employed an approach that is  
631 similar to TORUS<sup>126</sup>. Briefly, for each functional annotation and each gene, we performed the logistic regression  
632  $g(\mathbf{P}) = \mathbf{a}\omega$  where  $g(\cdot)$  is the logit link function,  $\mathbf{P}$  is the vector for the PIPs of all the SNPs,  $\mathbf{a}$  is the binary vector  
633 indicating whether the SNPs fall into the annotation, and  $\omega$  is the desired log-enrichment scores. After removing  
634 the genes on which logistic regression does not converge, we meta-analyzed the log-enrichment scores across  
635 genes by  $\omega_{\text{meta}} = \frac{\sum \phi_i \omega_i}{\sum \phi_i}$  and  $z_{\omega_{\text{meta}}} = \frac{\sum \phi_i \omega_i}{\sqrt{\sum \phi_i}}$  where  $\phi_i$  is the inverse of the squared standard error for gene  $i$ .

636 When comparing enrichment results across methods, we focused on e/pGenes fine-mapped by both methods.  
637 We computed the comparison z score as  $\frac{\omega_{\text{meta},j} - \omega_{\text{meta},j'}}{\sqrt{se^2_{\omega_{\text{meta},j}} + se^2_{\omega_{\text{meta},j}'}}}$  for method  $j$  and  $j'$ . For the enrichment analyses  
638 focusing on individual shared effect using  $\alpha_L$ , rather than PIPs, we limited analyses to those single effects that had  
639 corresponding credible sets (i.e., were not pruned).

640 To perform a case study, we selected *URGCP*, which was fine-mapped by SuShiE, but missed by other methods.  
641 To annotate the genomic region around *URGCP*, we downloaded the ChIP-Seq H3K27ac data of ENCODE<sup>58</sup> from  
642 WashU Epigenome Browser<sup>127</sup> (see **Code Availability**) and proximal enhancer (pELS) cCREs from ENCODE Registry  
643 v3, PBMC annotation using scATAC-seq in Satpathy et al.<sup>60</sup>, naive T cells, naive B cells, cytotoxic natural killer (cNK)  
644 cells, and monocytes annotations using snATAC-seq in Chiou et al.<sup>59</sup>

## 645 Prior *cis*-molQTL correlation analyses

646 To shed light on the relationship between heterogeneity of effect-sizes across ancestries and genes' constraint,  
647 using all the credible sets output by SuShiE, we tested for association between SuShiE-inferred effect size  
648 correlations across ancestries ( $\rho_l$ ) and five measures of constraint ( $s$ ) using all the fine-mapped e/pGenes:  
649 probability of being Loss-of-Function Intolerant (pLI)<sup>75</sup>, loss-of-function observed/expected upper bound fraction  
650 (LOEUF)<sup>76</sup>, enhancer-domain score (EDS)<sup>77</sup>, the Residual Variation Intolerance Score (RVIS)<sup>78</sup>, and  $s_{\text{het}}$ <sup>79</sup>. We  
651 downloaded pLI and LOEUF from gnomAD browser v4.0 (see **Code Availability**), we downloaded EDS, RVIS, and  
652  $s_{\text{het}}$  from their original papers. Our base model is according to:

$$653 \quad E(s) = \mathbf{v}_0 + \rho_l v_1 + \mathbf{L}v_2 + \mathbf{d}v_3 + \mathbf{r}v_4$$

654 where  $\mathbf{v}_0$  is the intercept term,  $\mathbf{L}$  is the ordered and categorical single effect index representing the order of  
655 variance explained,  $\mathbf{d}$  is the corresponding ancestry pair indicator (e.g., the correlation of EUR-AFR, EUR-HIS, or  
656 HIS-AFR),  $\mathbf{r}$  is the study indicator (e.g., TOPMed MESA mRNA, TOPMed MESA proteins, or GENOA mRNA),  $v_i$ s are  
657 the corresponding coefficients. We test the significance of  $v_1$  in a linear regression framework. A negative value  
658 of  $v_1$  for pLI, EDS, and  $s_{\text{het}}$  is taken to indicate stronger associations between *cis*-molQTL effect size heterogeneity  
659 across ancestries and gene constraint, while a lower value of LOEUF and RVIS is suggestive of stronger associations.  
660 In addition, to show robustness, we re-tested these associations using estimated covariance by replacing  $\rho_l$  by  
661  $\sigma_b^2$ . We also only focused on correlations estimated only from the primary effect (i.e.,  $L=1$ ); in this case, we  
662 removed  $\mathbf{L}$  from the base model. We also re-computed the standard error using bootstrap. Specifically, for each  
663 study, each ancestry pair, and each  $L$ , we sampled the genes with replacement and computed the  $v_1$ . We  
664 repeated 100 times to construct the null distributions for  $v_1$  and used its standard deviation as a new standard  
665 error. In addition, to adjust for allele frequency differences across ancestries, we added Wright's fixation index  
666 ( $F_{\text{st}}$ ) as an additional term. To compute  $F_{\text{st}}$ , we only used the fine-mapped SNPs to compute the  $F_{\text{st}}$  using  
667 PLINK2<sup>113,114</sup> with the "Hudson" method<sup>128,129</sup> for each gene. To investigate the relationship between expected *cis*-

668 molQTLs's distance to TSS and genes' constraint, we computed the expected distance to TSS for each gene  
669 according to  $\frac{\sum PIP_i * D_i}{\sum PIP_i}$  where  $D_i$  is the distance (absolute value) to the TSS for SNP  $i$ .

## 670 TWAS and PWAS analyses in All Of Us biobank

671 We performed individual-level Transcriptome- and Proteome-wide Association Studies (TWASs and PWASs)<sup>42-44,47</sup>  
672 on 6 white blood cell-related traits: basophil count (BAS), eosinophil count (EOS), lymphocyte count (LYM),  
673 monocyte count (MON), neutrophil count (NEU), and white blood cell count (WBC; **Table S9**) measured in AOU  
674 biobank<sup>50</sup>. We excluded individuals who had acute abdomen, acute appendicitis, acute cholangitis, acute  
675 cholecystitis, acute pancreatitis, anemia due to and following chemotherapy, bone marrow transplant present,  
676 chemotherapy-induced nausea and vomiting, cirrhosis of liver, clostridium difficile colitis, complication of  
677 chemotherapy, congenital anemia, congenital hemolytic anemia, convalescence after chemotherapy, dermatosis  
678 resulting from cytotoxic therapy, diverticulitis of intestine, end-stage renal disease, fatigue due to chemotherapy,  
679 hereditary hemolytic anemia, human immunodeficiency virus infection, leukemia, mucositis following  
680 chemotherapy, myelodysplastic syndrome (clinical), neutropenia due to and following chemotherapy,  
681 pancytopenia due to antineoplastic chemotherapy, peripheral neuropathy due to and following antineoplastic  
682 therapy, post-splenectomy disorder and post-splenectomy thrombocytosis. For WBC, we only included  
683 measurements <200e9/L. For all the traits, we also excluded measurements that were 3 standard deviations away  
684 from the mean, resulting in a total of 86,345 individuals on average. We identified individual ancestry information  
685 based on AOU precomputed information (i.e., "eur", "afr", and "amr" labels), resulting in 53,268 EUR, 16,748 AFR,  
686 and 16,329 HIS individuals on average.

687 From our previous analysis, we obtained the eQTL prediction weights of EUR, AFR, and HIS in the TOPMed MESA  
688 mRNA dataset, the pQTL prediction weights of EUR, AFR, and HIS in the TOPMed MESA protein dataset, and the  
689 eQTL prediction weights of EUR and AFR in the GENOA mRNA dataset. We evaluated the prediction accuracy for  
690 SuShiE SuShiE-Indep, Meta-SuSiE, SuSiE, LASSO<sup>53</sup>, Elastic Net<sup>54</sup>, and gBLUP<sup>55</sup> with five-fold cross-validation. For

691 Meta-SuSiE, we trained the prediction weights for each ancestry. For SuSiE, LASSO, Elastic Net, and gBLUP, we  
692 trained the prediction weights after concatenating genotype and phenotype data across ancestries to ensure the  
693 equal sample sizes as SuShiE (i.e., the same prediction weights for all ancestries). We computed cross validation  
694  $r^2$  ( $cv-r^2$ ) between the measured expression levels and predicted expression levels concatenated across each fold  
695 and each ancestry. We also used the SuShiE-based ancestry-specific prediction weights to evaluate the prediction  
696 performance using cross-ancestry weights. Specifically, we predicted the expression levels of EUR individuals using  
697 AFR weights (of AFR individuals using HIS weights and of HIS individuals using EUR weights).

698 To perform T/PWAS, we first predicted expression levels (either mRNA or proteins) for EUR, AFR, and HIS  
699 individuals in AOU using each ancestry-matched e/pQTL prediction weights with the score function in PLINK2<sup>113,114</sup>.  
700 Then, we standardized the expression vector (centered by mean and scaled by standard deviation) within each  
701 ancestry and then concatenated them into a single vector across ancestries. Then, we regressed out sex, age,  
702 squared age, and ten genotype PCs from the trait measurements. Last, we regressed the inverse-rank normalized  
703 residuals on the predicted expression levels to compute the TWAS or PWAS statistics. We re-performed the  
704 procedure using SuSiE-derived e/pQTL prediction weights as comparisons. We applied the Bonferroni correction  
705 to adjust the reported P-values with  $n=23,000$ . To validate our TWAS results, we compared them to five  
706 independent TWAS studies: Lu and Gopalan et al.<sup>36</sup>, Kachuri et al.<sup>31</sup>, Tapia et al.<sup>82</sup>, Rowland et al.<sup>84</sup>, and Wen et  
707 al.<sup>83</sup> We released our *cis*-molQTL prediction weights to the public, which can be found at the Data Availability  
708 section. To test the association between T/PWAS chi-square statistics and genes' constraint scores: pLI<sup>75</sup>, LOEUF<sup>76</sup>,  
709 EDS<sup>77</sup>, RVIS<sup>78</sup>, and  $S_{het}$ <sup>79</sup>, we used linear regression adjusted for phenotype and study and reported one-sided P  
710 values. To compare significance of these associations between SuShiE and SuSiE, we computed the z score as  
711  $\frac{\sum_{i=1}^I X_{i,SuShiE}^2 - \sum_{i=1}^I X_{i,SuSiE}^2}{\sqrt{2*(2I)}}$  where  $X_{i,*}^2$  is the chi-square statistics for constraint score  $i$ . We classified genes into three  
712 groups: Low, Middle, and High based on different scores, respectively. For pLI, we labeled genes with pLI >0.9 as  
713 High, <0.1 as Low, and otherwise middle. For other scores, we labeled genes whose value is greater than 90%  
714 quantile as High, smaller than 10% quantile as Low, and otherwise middle.

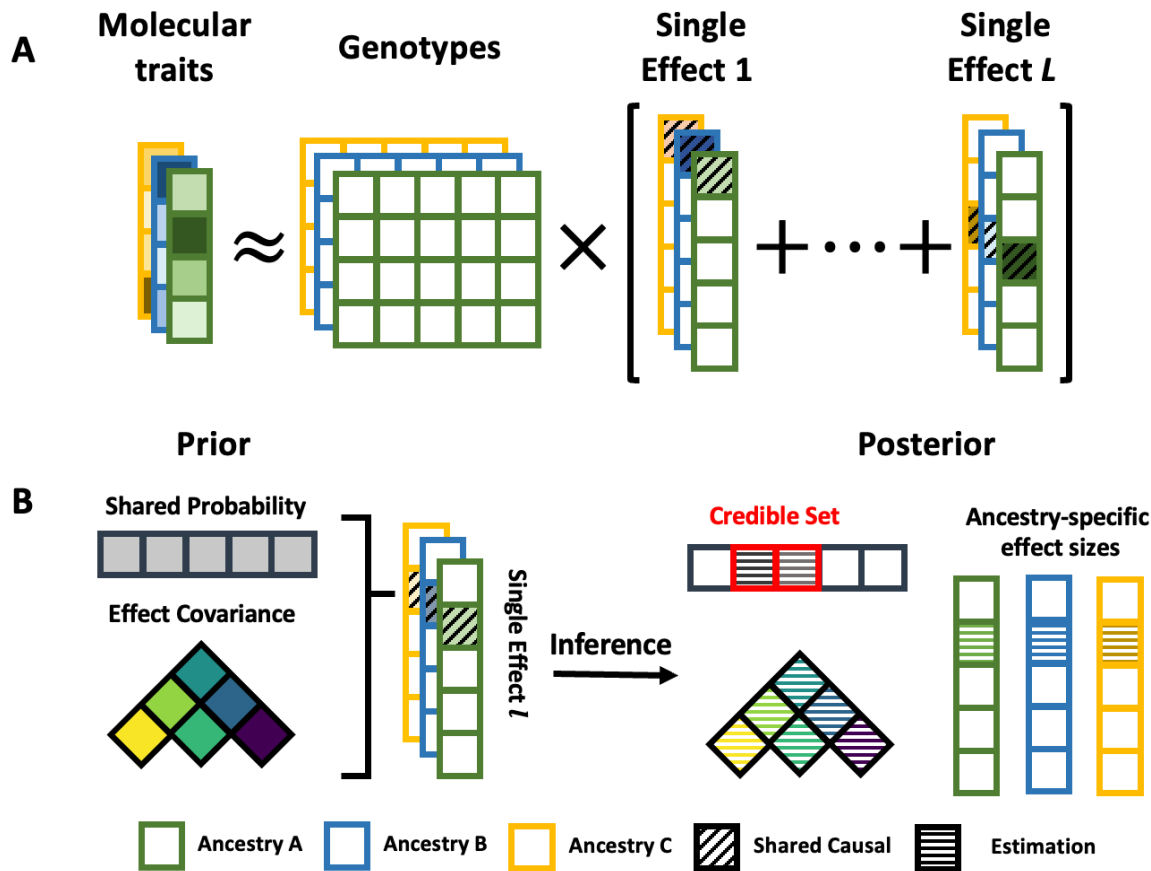
## 715 High-speed inference of SuShiE using JAX

716 We implemented SuShiE in an open-sourced command-line Python software *sushie*, which can read individual-  
717 level genotype data in three formats: PLINK1.9<sup>113,114</sup>, bgen1.3<sup>130</sup>, and vcf<sup>116</sup>, together with phenotypic and  
718 covariates data in tab-separated-values format. We leveraged *Just In Time* compilation in JAX (see **Code**  
719 **Availability**) to facilitate high-speed inference on CPUs, GPUs, or TPUs. This technique allows users to process, in  
720 a scalable fashion, thousands of molecular phenotypes with the backgrounds of diverse ancestries specified by  
721 the user. Not only can *sushie* perform our method, but it can also perform single-ancestry SuSiE<sup>15</sup>, effect size  
722 correlation estimation, *cis*-SNP heritability estimation, cross-validation for the *cis*-molQTL prediction weights, and  
723 contain the script to convert the *cis*-molQTL prediction results to FUSION format<sup>42</sup>, thus can be used in TWAS  
724 framework. We also implemented basic QC on the input data. Users can also customize the *sushie* inference  
725 function according to their preferences. We have compiled comprehensive documentation about the software at  
726 <https://mancusolab.github.io/sushie/>.

727



728 Figures



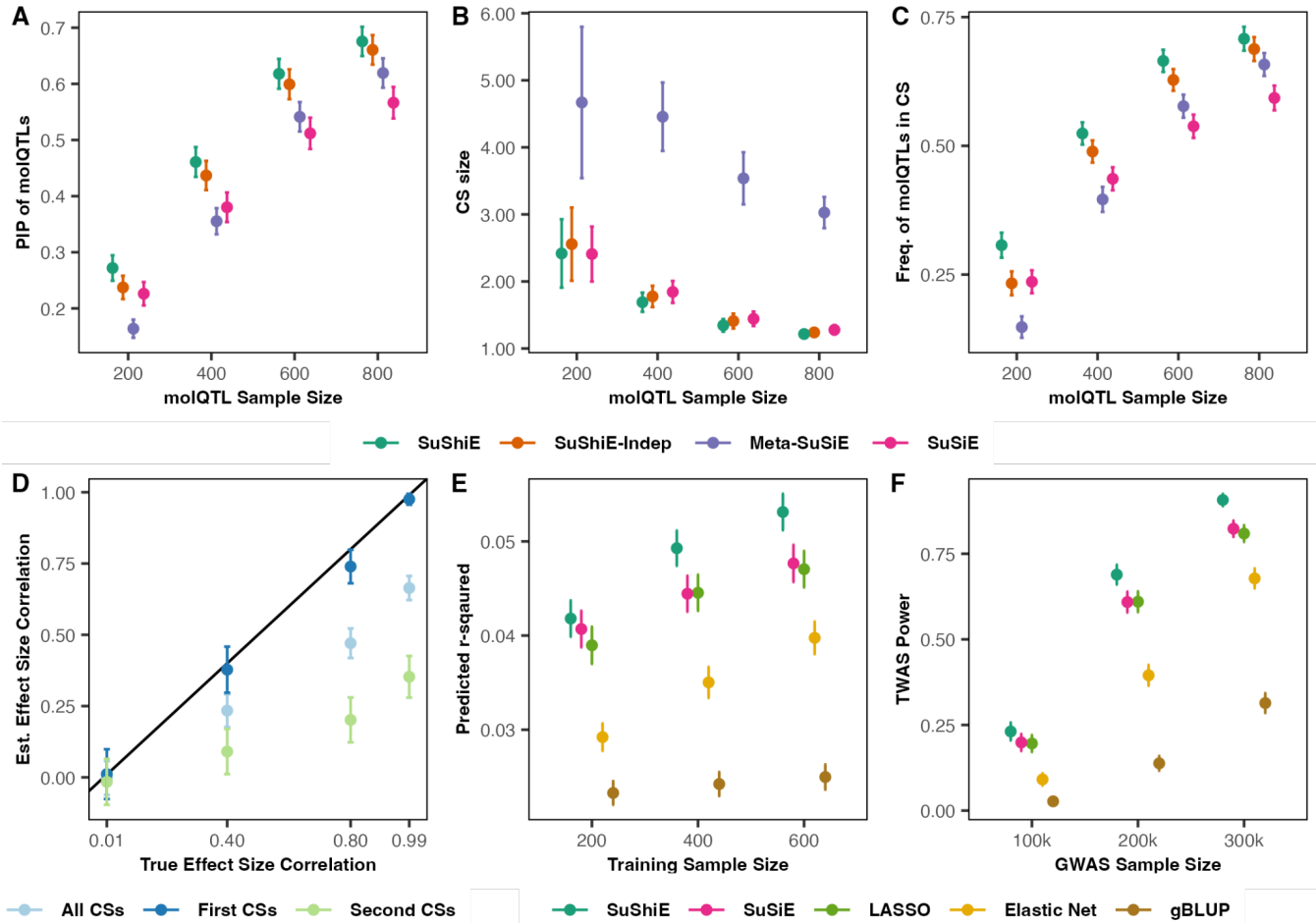
729

730 **Fig. 1: SuShiE infers ancestry-specific effect sizes, PIPs, and credible sets by leveraging shared**  
 731 **genetic architectures and LD heterogeneity.**

732 **A)** SuShiE takes individual-level phenotypic and genotypic data as input and assumes the shared *cis*-molQTL effects  
 733 as a linear combination of single effects.

734 **B)** For each single shared effect, SuShiE models the *cis*-molQTL effect size follows a multivariate normal prior  
 735 distribution with a covariance matrix, and the probability for each SNP to be molQTL follows a uniform prior  
 736 distribution; through the inference, SuShiE outputs a credible set that includes putative causal *cis*-molQTLs, learns  
 737 the effect-size covariance prior, and estimates the ancestry-specific effect sizes.

738



739

740

741

## Fig. 2: SuShiE outperforms other methods, estimates accurate effect-size correlation, and boosts higher power of TWAS in realistic simulations

742

743

744

745

746

747

748

749

750

751

752

753

754

755

756

757

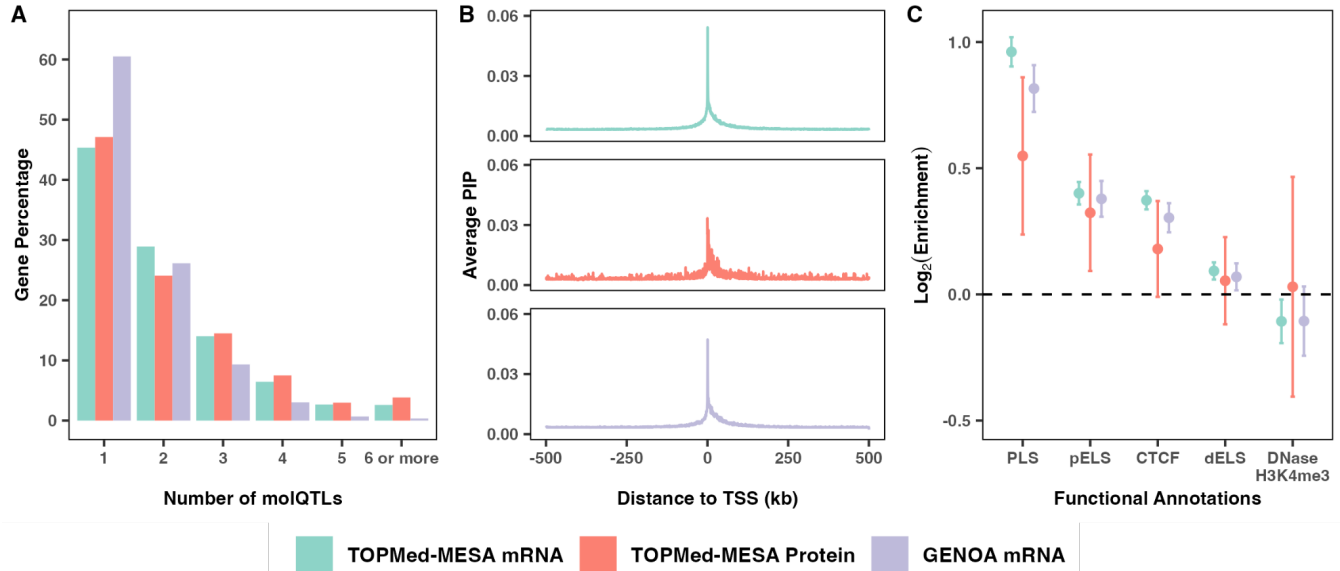
**A-C)** SuShiE outputs higher posterior inclusion probabilities (PIPs; A), smaller credible set sizes (B), and higher frequency of *cis*-molQTLs in the credible sets (calibration; C) compared to SuShiE-Indep ( $2.60 \times 10^{-4}$ ,  $1.5 \times 10^{-1}$ , and  $1.30 \times 10^{-11}$ ), Meta-SuSiE ( $P=9.67 \times 10^{-43}$ ,  $9.35 \times 10^{-231}$ , and  $1.17 \times 10^{-76}$ ), and SuSiE ( $P=6.98 \times 10^{-63}$ ,  $6.65 \times 10^{-2}$ , and  $1.58 \times 10^{-104}$ ).

**D)** SuShiE accurately estimates the true effect-size correlation across ancestries using the primary effect (First credible sets; CSs) while exhibiting an underestimation using the secondary effects (Second CSs) or combined (All CSs) because the variance explained by the secondary effect decreases, thus requiring higher statistical power. The error bar is a 95% confidence interval.

**E)** SuShiE outputs higher ancestry-specific prediction accuracy compared against SuSiE, LASSO, Elastic Net, and gBLUP (all  $P < 9.57 \times 10^{-8}$ ) with the fixed sample size. The plots are aggregation across two ancestries.

**F)** SuShiE induces higher TWAS power compared to SuSiE, LASSO, Elastic Net, and gBLUP (all  $P < 4.34 \times 10^{-14}$ ) with the fixed sample size. The plots are aggregation across two ancestries.

By default, the simulation assumes that there are 2 causal *cis*-molQTLs, the per-ancestry training sample size is 400, and the testing sample size is 200, *cis*-SNP heritability is 0.05, the effect size correlation is 0.8 across ancestries, and the proportion of *cis*-SNP heritability of complex trait explained by gene expression is  $1.5 \times 10^{-14}$ . The error bar is a 95% confidence interval.



758

759

### Fig. 3: SuShiE reveals cis-regulatory mechanisms for mRNA and protein expression

760

**A)** SuShiE identified *cis*-molQTLs for 14,590, 573, and 5,925 genes whose 88%, 86%, and 96% contain 1-3 *cis*-molQTLs for the TOPMed-MESA mRNA, TOPMed-MESA protein, and GENOA mRNA dataset, respectively.

761

**B)** Posterior inclusion probabilities (PIPs) of *cis*-molQTLs inferred by SuShiE are mainly enriched around the TSS region of genes. We grouped SNPs into 500-bp-long bins and computed their PIP average. There are 2,000 bins to cover a one-million-bp-long genomic window around the genes' TSS.

762

**C)** Across all three studies, *cis*-molQTLs identified by SuShiE are enriched in four out of five candidate *cis*-regulatory elements (cCREs) from ENCODE<sup>58</sup>, with the promoter (PLS) as the most enriched category. Specifically, the mRNA expression from TOPMed-MESA and GENOA showed enrichment in the promoter, proximal enhancer (pELS), CTCF, and distal enhancer (dELS) but depletion in DNase-H3K4me3. Protein expression from TOPMed-MESA showed enrichment in PLS and pELS but non-significant enrichment in CTCF and dELS because of the low number of genes identified with pQTLs (n=573). The error bar is a 95% confidence interval.

763

764

765

766

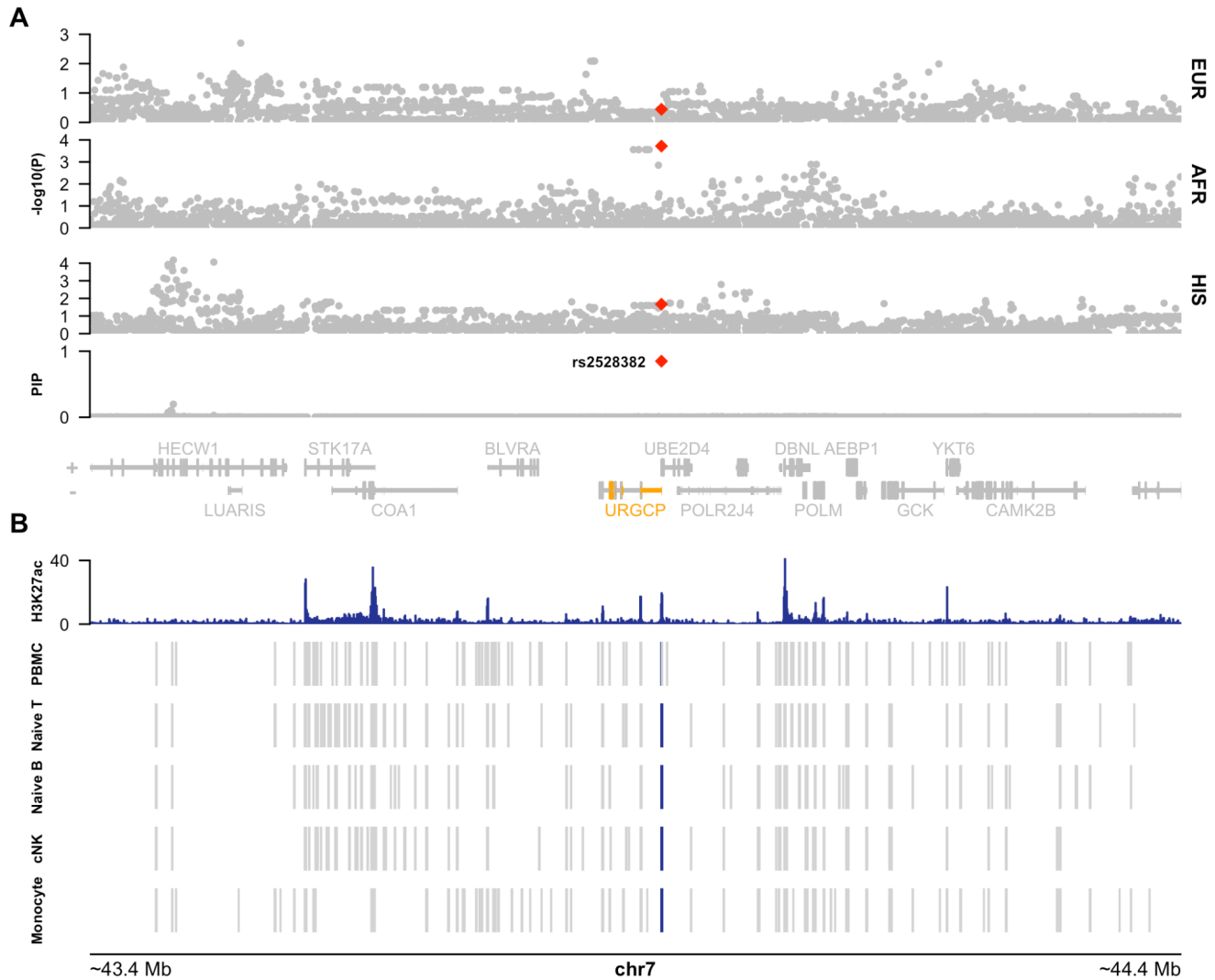
767

768

769

770

771



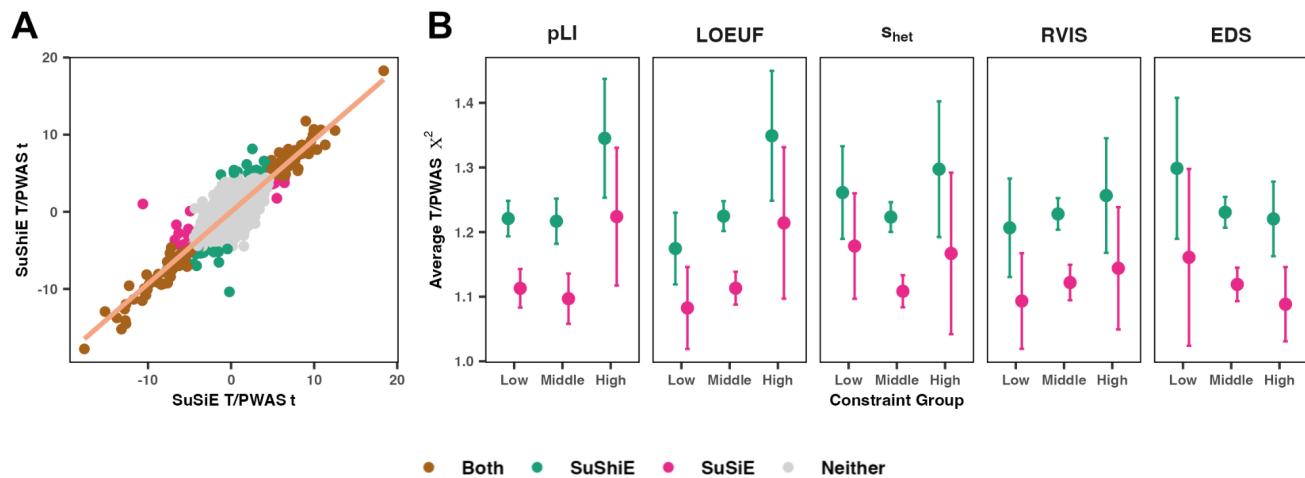
772

773 **Fig. 4: SuShiE identifies eQTL *rs2528382* for *URGCP* with functional support**

774 **A)** Manhattan plot of *cis*-eQTL scans of *URGCP* (denoted in orange) for each ancestry (above) with SuShiE fine-  
 775 mapping results (below). SuShiE was the only method to output credible sets for *URGCP* and prioritized a single  
 776 SNP (*rs2528382*; denoted in red).

777 **B)** Functional annotations at *URGCP* locus show colocalization of active enhancer activity and chromatin  
 778 accessibility with *rs2528382*. H3K27ac CHIP-seq peaks measured in PBMCs (intensity denoted in blue) and 0/1  
 779 accessibility annotations determined from scATAC-seq measured in PBMCs and snATAC-seq measured in naive T  
 780 cells, naive B cells, cytotoxic NK (cNK) cells, and monocytes. Blue rectangles denote a putative cCRE called from  
 781 sc/snATAC-seq data that colocalize with *rs2528382* (gray no colocalization).

782



783

784 **Fig. 5: SuShiE identifies more T/PWAS genes compared with SuSiE**

785 **A)** Scatter plot of T/PWAS t-statistics between SuShiE (y-axis) and SuSiE (x-axis) across all phenotypes and  
786 contributing *cis*-molQTL studies.

787 **B)** Average T/PWAS chi-square statistics within low, middle, and high constraint scores (see **Methods**). Error bars  
788 represent 95% confidence intervals.

789

## 790 Tables

791

	pLI	LOEUF	S <sub>het</sub>	RVIS	EDS
<b>Base Model</b>	-0.022 (4.13e-33)	0.021 (5.92e-20)	-0.007 (4.15e-40)	0.043 (2.04e-14)	-0.002 (1.25e-02)
<b>Bootstrap SE</b>	-0.022 (5.84e-32)	0.021 (4.92e-20)	-0.007 (3.13e-37)	0.043 (1.68e-17)	-0.002 (1.56e-02)
<b>Primary Effect</b>	-0.034 (3.51e-23)	0.027 (7.45e-11)	-0.011 (7.69e-29)	0.055 (4.09e-09)	-0.004 (1.27e-03)
<b>Effect Covariance</b>	-0.339 (7.59e-177)	0.334 (1.77e-109)	-0.089 (1.33e-154)	0.537 (9.93e-49)	-0.053 (3.10e-25)
<b>Adjusted F<sub>st</sub></b>	-0.022 (2.00e-32)	0.021 (9.90e-20)	-0.007 (2.22e-39)	0.042 (5.63e-14)	-0.002 (1.08e-02)

792

### 793 **Table 1: Across-ancestry *cis*-molQTL effect size correlations are negatively associated with** 794 **gene constraint scores**

795 The estimates and corresponding P-value in the regression framework testing associations between inferred  
796 effect size correlations across ancestries and constraint scores (see **Methods** for the base model). “Bootstrap SE”  
797 is to re-estimate standard error using bootstrap. “Primary Effect” is to only use estimates from L=1. “Effect  
798 Covariance” is to replace estimated correlation with estimated effect size covariance across ancestries. “Adjusted  
799 F<sub>st</sub>” is to additionally adjusted for F<sub>st</sub> from the base model. A higher value of pLI, S<sub>het</sub>, and, EDS is taken to indicate  
800 stronger constraint, while a lower value of LOEUF and RVIS is suggestive of more constraint. The reported P-value  
801 is one-sided.

802

## 803 Data availability

804 SuShiE-derived prediction models for TWAS/PWAS, fine-mapping, and other analyzed results across *cis*-molQTL  
805 datasets can be found at <https://zenodo.org/records/10963034>.

## 806 Code availability

807 SuShiE: <https://github.com/mancusolab/sushie>

808 The analysis codes for simulation and real-data analysis of this manuscript:

809 <https://github.com/mancusolab/sushie-project-codes>

810 TOPMed RNA-seq Harmonization pipeline: [https://github.com/broadinstitute/gtex-](https://github.com/broadinstitute/gtex-pipeline/blob/master/TOPMed_RNAseq_pipeline.md)  
811 [pipeline/blob/master/TOPMed\\_RNAseq\\_pipeline.md](https://github.com/broadinstitute/gtex-pipeline/blob/master/TOPMed_RNAseq_pipeline.md)

812 [gnomAD v4.0: https://gnomad.broadinstitute.org/news/2023-11-gnomad-v4-0/](https://gnomad.broadinstitute.org/news/2023-11-gnomad-v4-0/)

813 GTEx eQTL analysis pipeline: <https://www.gtexportal.org/home/methods>

814 pyqtl software: <https://github.com/broadinstitute/pyqtl>

815 PLINK: <https://www.cog-genomics.org/plink/2.0>

816 BCFTOOLS: <https://samtools.github.io/bcftools/bcftools.html>

817 JAX: <https://github.com/google/jax>

818 scikit-learn: <https://scikit-learn.org/stable/>

819 FUSION: <http://gusevlab.org/projects/fusion/>  
820 limix: <https://github.com/limix/limix>  
821 LiftOver: <https://genome.ucsc.edu/cgi-bin/hgLiftOver>  
822 WashU Epigenome Browser: <https://epigenomegateway.wustl.edu/>

## 823 Acknowledgements

824 The authors would like to thank members of the Mancuso and Gazal labs for fruitful discussions regarding this  
825 manuscript. The authors would also like to specially thank Dr. Michael D. Edge for his thoughtful comments and  
826 suggestions. This work was funded in part by National Institutes of Health (NIH) under awards R01HG012133,  
827 R01CA258808, R01GM140287, R35GM142783, R01GM140287, U54HG013243, R35GM147789, and  
828 K08HL159346.

829 MESA phenotypes (dbGaP: phs000209.v13.p3): MESA and the MESA SHARe project are conducted and supported  
830 by the National Heart, Lung, and Blood Institute (NHLBI) in collaboration with MESA investigators. Support for  
831 MESA is provided by contracts HHSN268201500003I, N01-HC-95159, N01-HC-95160, N01-HC-95161, N01-HC-  
832 95162, N01-HC-95163, N01-HC95164, N01-HC-95165, N01-HC-95166, N01-HC-95167, N01-HC-95168, N01-HC-  
833 95169, UL1-TR-001079, UL1-TR000040, UL1-TR-001420, UL1-TR-001881, and DK063491. Funding for SHARe  
834 genotyping was provided by NHLBI Contract N02-HL-64278. TOPMed MESA WGS genotype, mRNA, and protein  
835 expression data (dbGaP: phs001416.v3.p1): Molecular data for the Trans-Omics in Precision Medicine (TOPMed)  
836 program was supported by the National Heart, Lung and Blood Institute (NHLBI). WGS genotype data for NHLBI  
837 TOPMed: MESA (phs001416.v3.p1) was performed at Broad Genomics (HHSN268201600034I). mRNA expression  
838 data for NHLBI TOPMed: MESA (phs001416.v3.p1) was performed at NWGC (HHSN268201600032I). SOMAscan  
839 proteomics for NHLBI TOPMed: Multi-Ethnic Study of Atherosclerosis (MESA) (phs001416.v1.p1) was performed  
840 at the Broad Institute and Beth Israel Proteomics Platform (HHSN268201600034I). Core support including  
841 centralized genomic read mapping and genotype calling, along with variant quality metrics and filtering were  
842 provided by the TOPMed Informatics Research Center (3R01HL-117626-02S1; contract HHSN268201800002I).  
843 Core support including phenotype harmonization, data management, sample-identity QC, and general program  
844 coordination were provided by the TOPMed Data Coordinating Center (R01HL-120393; U01HL-120393; contract  
845 HHSN268201800001I). We gratefully acknowledge the studies and participants who provided biological samples  
846 and data for TOPMed.

847 GENOA genotype (dbGaP: phs001238.v2.p1) and gene expression (GEO: GSE138914) data were supported by  
848 grants from NIH NHLBI (HL054457, HL054464, HL054481, HL119443, and HL087660). The authors would like to  
849 acknowledge Drs. Sharon Kardia and Jennifer Smith in preparing GENOA eQTL data.

850 The All of Us Research Program is supported by the National Institutes of Health, Office of the Director: Regional  
851 Medical Centers: 1 OT2 OD026549; 1 OT2 OD026554; 1 OT2 OD026557; 1 OT2 OD026556; 1 OT2 OD026550; 1  
852 OT2 OD 026552; 1 OT2 OD026553; 1 OT2 OD026548; 1 OT2 OD026551; 1 OT2 OD026555; IAA #: AOD 16037;  
853 Federally Qualified Health Centers: HHSN 263201600085U; Data and Research Center: 5 U2C OD023196; Biobank:  
854 1 U24 OD023121; The Participant Center: U24 OD023176; Participant Technology Systems Center: 1 U24  
855 OD023163; Communications and Engagement: 3 OT2 OD023205; 3 OT2 OD023206; and Community Partners: 1  
856 OT2 OD025277; 3 OT2 OD025315; 1 OT2 OD025337; 1 OT2 OD025276. In addition, the All of Us Research Program  
857 would not be possible without the partnership of its participants.

## 858 Author contributions

859 Z.L. and N.M. developed the model and study design. Z.L. performed simulations and fine-mapping analyses. Z.L.,  
860 X.W., J.P., and L.K. performed TWAS and AoU analyses. Z.L., M.C., and N.M. developed the model and inference  
861 scheme. Z.L. and A.K. prepared functional genomic annotations and enrichment analyses. Z.L. and N.M. wrote the  
862 initial manuscript. All authors edited the final manuscript.

## 863 Competing interests

864 L.W. provided consulting service to Pupil Bio Inc. and reviewed manuscripts for Gastroenterology Report, not  
865 related to this study, and received honorarium. No potential conflicts of interest were disclosed by the other  
866 authors.



## 867 References

- 868 1. Claussnitzer, M. *et al.* A brief history of human disease genetics. *Nature* **577**, 179–189 (2020).
- 869 2. Cheung, V. G. *et al.* Mapping determinants of human gene expression by regional and genome-wide  
870 association. *Nature* **437**, 1365–1369 (2005).
- 871 3. GTEx Consortium. The GTEx Consortium atlas of genetic regulatory effects across human tissues. *Science*  
872 **369**, 1318–1330 (2020).
- 873 4. Vösa, U. *et al.* Large-scale cis- and trans-eQTL analyses identify thousands of genetic loci and polygenic  
874 scores that regulate blood gene expression. *Nat. Genet.* **53**, 1300–1310 (2021).
- 875 5. Sun, B. B. *et al.* Genomic atlas of the human plasma proteome. *Nature* **558**, 73–79 (2018).
- 876 6. Gate, R. E. *et al.* Genetic determinants of co-accessible chromatin regions in activated T cells across  
877 humans. *Nat. Genet.* **50**, 1140–1150 (2018).
- 878 7. Battle, A. *et al.* Genomic variation. Impact of regulatory variation from RNA to protein. *Science* **347**, 664–  
879 667 (2015).
- 880 8. Li, Y. I. *et al.* RNA splicing is a primary link between genetic variation and disease. *Science* **352**, 600–604  
881 (2016).
- 882 9. McVicker, G. *et al.* Identification of genetic variants that affect histone modifications in human cells.  
883 *Science* **342**, 747–749 (2013).
- 884 10. Oliva, M. *et al.* DNA methylation QTL mapping across diverse human tissues provides molecular links  
885 between genetic variation and complex traits. *Nat. Genet.* **55**, 112–122 (2023).
- 886 11. Wu, L. *et al.* Variation and genetic control of protein abundance in humans. *Nature* **499**, 79–82 (2013).
- 887 12. Aguet, F. *et al.* Molecular quantitative trait loci. *Nat. Rev. Methods Primers* **3**, 1–22 (2023).
- 888 13. Albert, F. W. & Kruglyak, L. The role of regulatory variation in complex traits and disease. *Nat. Rev. Genet.*  
889 **16**, 197–212 (2015).
- 890 14. Suhre, K., McCarthy, M. I. & Schwenk, J. M. Genetics meets proteomics: perspectives for large population-  
891 based studies. *Nat. Rev. Genet.* **22**, 19–37 (2021).
- 892 15. Wang, G., Sarkar, A., Carbonetto, P. & Stephens, M. A simple new approach to variable selection in  
893 regression, with application to genetic fine mapping. *J. R. Stat. Soc. Series B Stat. Methodol.* **82**, 1273–1300  
894 (2020).
- 895 16. Hormozdiari, F., Kichaev, G., Yang, W.-Y., Pasaniuc, B. & Eskin, E. Identification of causal genes for complex  
896 traits. *Bioinformatics* **31**, i206-13 (2015).
- 897 17. Benner, C. *et al.* FINEMAP: efficient variable selection using summary data from genome-wide association  
898 studies. *Bioinformatics* **32**, 1493–1501 (2016).
- 899 18. Schaid, D. J., Chen, W. & Larson, N. B. From genome-wide associations to candidate causal variants by  
900 statistical fine-mapping. *Nat. Rev. Genet.* **19**, 491–504 (2018).
- 901 19. Kichaev, G. *et al.* Integrating functional data to prioritize causal variants in statistical fine-mapping studies.  
902 *PLoS Genet.* **10**, e1004722 (2014).
- 903 20. Wojcik, G. L. *et al.* Genetic analyses of diverse populations improves discovery for complex traits. *Nature*  
904 **570**, 514–518 (2019).
- 905 21. Chen, M.-H. *et al.* Trans-ethnic and Ancestry-Specific Blood-Cell Genetics in 746,667 Individuals from 5  
906 Global Populations. *Cell* **182**, 1198-1213.e14 (2020).
- 907 22. Martin, A. R. *et al.* Clinical use of current polygenic risk scores may exacerbate health disparities. *Nat.*  
908 *Genet.* **51**, 584–591 (2019).
- 909 23. Sakaue, S. *et al.* A cross-population atlas of genetic associations for 220 human phenotypes. *Nat. Genet.*  
910 **53**, 1415–1424 (2021).
- 911 24. Conti, D. V. *et al.* Trans-ancestry genome-wide association meta-analysis of prostate cancer identifies new  
912 susceptibility loci and informs genetic risk prediction. *Nat. Genet.* **53**, 65–75 (2021).
- 913 25. Wang, A. *et al.* Characterizing prostate cancer risk through multi-ancestry genome-wide discovery of 187

- 914 novel risk variants. *Nat. Genet.* (2023) doi:10.1038/s41588-023-01534-4.
- 915 26. Shang, L. *et al.* Genetic Architecture of Gene Expression in European and African Americans: An eQTL  
916 Mapping Study in GENOA. *Am. J. Hum. Genet.* **106**, 496–512 (2020).
- 917 27. Mogil, L. S. *et al.* Genetic architecture of gene expression traits across diverse populations. *PLoS Genet.* **14**,  
918 e1007586 (2018).
- 919 28. Schubert, R. *et al.* Protein prediction for trait mapping in diverse populations. *PLoS One* **17**, e0264341  
920 (2022).
- 921 29. Tehranchi, A. *et al.* Fine-mapping cis-regulatory variants in diverse human populations. *Elife* **8**, (2019).
- 922 30. Wen, X., Luca, F. & Pique-Regi, R. Cross-population joint analysis of eQTLs: fine mapping and functional  
923 annotation. *PLoS Genet.* **11**, e1005176 (2015).
- 924 31. Kachuri, L. *et al.* Gene expression in African Americans, Puerto Ricans and Mexican Americans reveals  
925 ancestry-specific patterns of genetic architecture. *Nat. Genet.* **55**, 952–963 (2023).
- 926 32. Kasela, S. *et al.* Interaction molecular QTL mapping discovers cellular and environmental modifiers of  
927 genetic regulatory effects. *Am. J. Hum. Genet.* **111**, 133–149 (2024).
- 928 33. Kichaev, G. & Pasaniuc, B. Leveraging Functional-Annotation Data in Trans-ethnic Fine-Mapping Studies.  
929 *Am. J. Hum. Genet.* **97**, 260–271 (2015).
- 930 34. Asimit, J. L., Hatzikotoulas, K., McCarthy, M., Morris, A. P. & Zeggini, E. Trans-ethnic study design  
931 approaches for fine-mapping. *Eur. J. Hum. Genet.* **24**, 1330–1336 (2016).
- 932 35. LaPierre, N. *et al.* Identifying causal variants by fine mapping across multiple studies. *PLoS Genet.* **17**,  
933 e1009733 (2021).
- 934 36. Lu, Z. *et al.* Multi-ancestry fine-mapping improves precision to identify causal genes in transcriptome-wide  
935 association studies. *Am. J. Hum. Genet.* **109**, 1388–1404 (2022).
- 936 37. Shen, J. *et al.* Fine-mapping and credible set construction using a multi-population Joint Analysis of  
937 Marginal summary statistics from Genome-wide Association Studies. *bioRxiv* (2022)  
938 doi:10.1101/2022.12.22.521659.
- 939 38. Yuan, K. *et al.* Fine-mapping across diverse ancestries drives the discovery of putative causal variants  
940 underlying human complex traits and diseases. *medRxiv* (2023) doi:10.1101/2023.01.07.23284293.
- 941 39. Cai, M. *et al.* XMAP: Cross-population fine-mapping by leveraging genetic diversity and accounting for  
942 confounding bias. *Nat. Commun.* **14**, 6870 (2023).
- 943 40. Zhou, F. *et al.* Leveraging information between multiple population groups and traits improves fine-  
944 mapping resolution. *Nat. Commun.* **14**, 7279 (2023).
- 945 41. Gao, B. & Zhou, X. MESuSiE enables scalable and powerful multi-ancestry fine-mapping of causal variants in  
946 genome-wide association studies. *Nat. Genet.* (2024) doi:10.1038/s41588-023-01604-7.
- 947 42. Gusev, A. *et al.* Integrative approaches for large-scale transcriptome-wide association studies. *Nat. Genet.*  
948 **48**, 245–252 (2016).
- 949 43. Gamazon, E. R. *et al.* A gene-based association method for mapping traits using reference transcriptome  
950 data. *Nat. Genet.* **47**, 1091–1098 (2015).
- 951 44. Zhu, Z. *et al.* Integration of summary data from GWAS and eQTL studies predicts complex trait gene  
952 targets. *Nat. Genet.* **48**, 481–487 (2016).
- 953 45. Gusev, A. *et al.* Transcriptome-wide association study of schizophrenia and chromatin activity yields  
954 mechanistic disease insights. *Nat. Genet.* **50**, 538–548 (2018).
- 955 46. Mancuso, N. *et al.* Large-scale transcriptome-wide association study identifies new prostate cancer risk  
956 regions. *Nat. Commun.* **9**, 4079 (2018).
- 957 47. Zhang, J. *et al.* Plasma proteome analyses in individuals of European and African ancestry identify cis-pQTLs  
958 and models for proteome-wide association studies. *Nat. Genet.* **54**, 593–602 (2022).
- 959 48. Bild, D. E. *et al.* Ethnic differences in coronary calcification: the Multi-Ethnic Study of Atherosclerosis  
960 (MESA). *Circulation* **111**, 1313–1320 (2005).
- 961 49. Taliun, D. *et al.* Sequencing of 53,831 diverse genomes from the NHLBI TOPMed Program. *Nature* **590**,

- 962 290–299 (2021).
- 963 50. All of Us Research Program Genomics Investigators. Genomic data in the All of Us Research Program.
- 964 *Nature* (2024) doi:10.1038/s41586-023-06957-x.
- 965 51. Zou, Y., Carbonetto, P., Wang, G. & Stephens, M. Fine-mapping from summary data with the “Sum of Single
- 966 Effects” model. *PLoS Genet.* **18**, e1010299 (2022).
- 967 52. Zou, Y., Carbonetto, P., Xie, D., Wang, G. & Stephens, M. Fast and flexible joint fine-mapping of multiple
- 968 traits via the Sum of Single Effects model. *bioRxiv* (2023) doi:10.1101/2023.04.14.536893.
- 969 53. Tibshirani, R. Regression Shrinkage and Selection via the Lasso. *J. R. Stat. Soc. Series B Stat. Methodol.* **58**,
- 970 267–288 (1996).
- 971 54. Zou, H. & Hastie, T. Regularization and Variable Selection Via the Elastic Net. *J. R. Stat. Soc. Series B Stat.*
- 972 *Methodol.* **67**, 301–320 (2005).
- 973 55. Clark, S. A. & van der Werf, J. Genomic best linear unbiased prediction (gBLUP) for the estimation of
- 974 genomic breeding values. *Methods Mol. Biol.* **1019**, 321–330 (2013).
- 975 56. Mostafavi, H., Spence, J. P., Naqvi, S. & Pritchard, J. K. Systematic differences in discovery of genetic effects
- 976 on gene expression and complex traits. *Nat. Genet.* **55**, 1866–1875 (2023).
- 977 57. Sun, B. B. *et al.* Plasma proteomic associations with genetics and health in the UK Biobank. *Nature* **622**,
- 978 329–338 (2023).
- 979 58. ENCODE Project Consortium *et al.* Expanded encyclopaedias of DNA elements in the human and mouse
- 980 genomes. *Nature* **583**, 699–710 (2020).
- 981 59. Chiou, J. *et al.* Interpreting type 1 diabetes risk with genetics and single-cell epigenomics. *Nature* **594**, 398–
- 982 402 (2021).
- 983 60. Satpathy, A. T. *et al.* Massively parallel single-cell chromatin landscapes of human immune cell
- 984 development and intratumoral T cell exhaustion. *Nat. Biotechnol.* **37**, 925–936 (2019).
- 985 61. Lappalainen, T. *et al.* Transcriptome and genome sequencing uncovers functional variation in humans.
- 986 *Nature* **501**, 506–511 (2013).
- 987 62. Tufan, N. L. S. *et al.* Hepatitis Bx antigen stimulates expression of a novel cellular gene, URG4, that
- 988 promotes hepatocellular growth and survival. *Neoplasia* **4**, 355–368 (2002).
- 989 63. Song, J. *et al.* Enhanced cell survival of gastric cancer cells by a novel gene URG4. *Neoplasia* **8**, 995–1002
- 990 (2006).
- 991 64. Li, W. & Zhou, N. URG4 upregulation is associated with tumor growth and poor survival in epithelial ovarian
- 992 cancer. *Arch. Gynecol. Obstet.* **286**, 209–215 (2012).
- 993 65. Xie, C. *et al.* Upregulator of cell proliferation predicts poor prognosis in hepatocellular carcinoma and
- 994 contributes to hepatocarcinogenesis by downregulating FOXO3a. *PLoS One* **7**, e40607 (2012).
- 995 66. Cai, J. *et al.* URGCP promotes non-small cell lung cancer invasiveness by activating the NF- $\kappa$ B-MMP-9
- 996 pathway. *Oncotarget* **6**, 36489–36504 (2015).
- 997 67. Cingolani, P. *et al.* A program for annotating and predicting the effects of single nucleotide polymorphisms,
- 998 SnpEff: SNPs in the genome of *Drosophila melanogaster* strain w1118; iso-2; iso-3. *Fly* **6**, 80–92 (2012).
- 999 68. Shi, H., Mancuso, N., Spendlove, S. & Pasaniuc, B. Local Genetic Correlation Gives Insights into the Shared
- 000 Genetic Architecture of Complex Traits. *Am. J. Hum. Genet.* **101**, 737–751 (2017).
- 001 69. Shi, H. *et al.* Population-specific causal disease effect sizes in functionally important regions impacted by
- 002 selection. *Nat. Commun.* **12**, 1098 (2021).
- 003 70. Saitou, M., Dahl, A., Wang, Q. & Liu, X. Allele frequency differences of causal variants have a major impact
- 004 on low cross-ancestry portability of PRS. *bioRxiv* (2022) doi:10.1101/2022.10.21.22281371.
- 005 71. Taylor, D. J. *et al.* Sources of gene expression variation in a globally diverse human cohort. *bioRxiv* (2023)
- 006 doi:10.1101/2023.11.04.565639.
- 007 72. Brown, B. C., Asian Genetic Epidemiology Network Type 2 Diabetes Consortium, Ye, C. J., Price, A. L. &
- 008 Zaitlen, N. Transethnic Genetic-Correlation Estimates from Summary Statistics. *Am. J. Hum. Genet.* **99**, 76–
- 009 88 (2016).

- 010 73. Hou, K. *et al.* Causal effects on complex traits are similar for common variants across segments of different  
011 continental ancestries within admixed individuals. *Nat. Genet.* **55**, 549–558 (2023).
- 012 74. Shi, H. *et al.* Localizing Components of Shared Transethnic Genetic Architecture of Complex Traits from  
013 GWAS Summary Data. *Am. J. Hum. Genet.* **106**, 805–817 (2020).
- 014 75. Lek, M. *et al.* Analysis of protein-coding genetic variation in 60,706 humans. *Nature* **536**, 285–291 (2016).
- 015 76. Karczewski, K. J. *et al.* The mutational constraint spectrum quantified from variation in 141,456 humans.  
016 *Nature* **581**, 434–443 (2020).
- 017 77. Wang, X. & Goldstein, D. B. Enhancer Domains Predict Gene Pathogenicity and Inform Gene Discovery in  
018 Complex Disease. *Am. J. Hum. Genet.* **106**, 215–233 (2020).
- 019 78. Petrovski, S., Wang, Q., Heinzen, E. L., Allen, A. S. & Goldstein, D. B. Genic intolerance to functional  
020 variation and the interpretation of personal genomes. *PLoS Genet.* **9**, e1003709 (2013).
- 021 79. Zeng, T., Spence, J. P., Mostafavi, H. & Pritchard, J. K. Bayesian estimation of gene constraint from an  
022 evolutionary model with gene features. *bioRxiv* (2023) doi:10.1101/2023.05.19.541520.
- 023 80. Berg, J. J. *et al.* Reduced signal for polygenic adaptation of height in UK Biobank. *Elife* **8**, (2019).
- 024 81. Keys, K. L. *et al.* On the cross-population generalizability of gene expression prediction models. *PLoS Genet.*  
025 **16**, e1008927 (2020).
- 026 82. Tapia, A. L. *et al.* A large-scale transcriptome-wide association study (TWAS) of 10 blood cell phenotypes  
027 reveals complexities of TWAS fine-mapping. *Genet. Epidemiol.* **46**, 3–16 (2022).
- 028 83. Wen, J. *et al.* Transcriptome-Wide Association Study of Blood Cell Traits in African Ancestry and  
029 Hispanic/Latino Populations. *Genes* **12**, (2021).
- 030 84. Rowland, B. *et al.* Transcriptome-wide association study in UK Biobank Europeans identifies associations  
031 with blood cell traits. *Hum. Mol. Genet.* **31**, 2333–2347 (2022).
- 032 85. Ding, Y. *et al.* Polygenic scoring accuracy varies across the genetic ancestry continuum. *Nature* **618**, 774–  
033 781 (2023).
- 034 86. Mester, R. *et al.* Impact of cross-ancestry genetic architecture on GWASs in admixed populations. *Am. J.*  
035 *Hum. Genet.* **110**, 927–939 (2023).
- 036 87. Hou, K., Bhattacharya, A., Mester, R., Burch, K. S. & Pasaniuc, B. On powerful GWAS in admixed  
037 populations. *Nature genetics* vol. 53 1631–1633 (2021).
- 038 88. Zhong, Y., Perera, M. A. & Gamazon, E. R. On Using Local Ancestry to Characterize the Genetic Architecture  
039 of Human Traits: Genetic Regulation of Gene Expression in Multiethnic or Admixed Populations. *Am. J.*  
040 *Hum. Genet.* **104**, 1097–1115 (2019).
- 041 89. Zhang, J. & Stram, D. O. The role of local ancestry adjustment in association studies using admixed  
042 populations. *Genet. Epidemiol.* **38**, 502–515 (2014).
- 043 90. Pasaniuc, B. *et al.* Enhanced statistical tests for GWAS in admixed populations: assessment using African  
044 Americans from CARE and a Breast Cancer Consortium. *PLoS Genet.* **7**, e1001371 (2011).
- 045 91. Seldin, M. F., Pasaniuc, B. & Price, A. L. New approaches to disease mapping in admixed populations. *Nat.*  
046 *Rev. Genet.* **12**, 523–528 (2011).
- 047 92. Atkinson, E. G. *et al.* Tractor uses local ancestry to enable the inclusion of admixed individuals in GWAS and  
048 to boost power. *Nat. Genet.* **53**, 195–204 (2021).
- 049 93. Qin, H. *et al.* Interrogating local population structure for fine mapping in genome-wide association studies.  
050 *Bioinformatics* **26**, 2961–2968 (2010).
- 051 94. Aracena, K. A. *et al.* Epigenetic variation impacts individual differences in the transcriptional response to  
052 influenza infection. *Nat. Genet.* **56**, 408–419 (2024).
- 053 95. Randolph, H. E. *et al.* Genetic ancestry effects on the response to viral infection are pervasive but cell type  
054 specific. *Science* **374**, 1127–1133 (2021).
- 055 96. Robinson, M. R. *et al.* Genotype-covariate interaction effects and the heritability of adult body mass index.  
056 *Nat. Genet.* **49**, 1174–1181 (2017).
- 057 97. Durvasula, A. & Lohmueller, K. E. Negative selection on complex traits limits phenotype prediction accuracy

- 058 between populations. *Am. J. Hum. Genet.* **108**, 620–631 (2021).
- 059 98. Yair, S. & Coop, G. Population differentiation of polygenic score predictions under stabilizing selection.  
060 *Philos. Trans. R. Soc. Lond. B Biol. Sci.* **377**, 20200416 (2022).
- 061 99. Agarwal, I., Fuller, Z. L., Myers, S. R. & Przeworski, M. Relating pathogenic loss-of-function mutations in  
062 humans to their evolutionary fitness costs. *Elife* **12**, (2023).
- 063 100. Speidel, L., Forest, M., Shi, S. & Myers, S. R. A method for genome-wide genealogy estimation for  
064 thousands of samples. *Nat. Genet.* **51**, 1321–1329 (2019).
- 065 101. Wang, J. & Gazal, S. Ancestry-specific regulatory and disease architectures are likely due to cell-type-  
066 specific gene-by-environment interactions. *medRxiv* (2023) doi:10.1101/2023.10.20.23297214.
- 067 102. Yazar, S. *et al.* Single-cell eQTL mapping identifies cell type-specific genetic control of autoimmune disease.  
068 *Science* **376**, eabf3041 (2022).
- 069 103. Bhattacharya, A. *et al.* Best practices for multi-ancestry, meta-analytic transcriptome-wide association  
070 studies: Lessons from the Global Biobank Meta-analysis Initiative. *Cell Genom* **2**, (2022).
- 071 104. Selewa, A. *et al.* Single-cell genomics improves the discovery of risk variants and genes of atrial fibrillation.  
072 *Nat. Commun.* **14**, 4999 (2023).
- 073 105. Blei, D. M., Kucukelbir, A. & McAuliffe, J. D. Variational inference: A review for statisticians. *J. Am. Stat.*  
074 *Assoc.* **112**, 859–877 (2017).
- 075 106. Mancuso, N. *et al.* Probabilistic fine-mapping of transcriptome-wide association studies. *Nat. Genet.* **51**,  
076 675–682 (2019).
- 077 107. Wang, X., Lu, Z., Bhattacharya, A., Pasaniuc, B. & Mancuso, N. *twas\_sim*, a Python-based tool for simulation  
078 and power analysis of transcriptome-wide association analysis. *Bioinformatics* (2023)  
079 doi:10.1093/bioinformatics/btad288.
- 080 108. 1000 Genomes Project Consortium *et al.* A global reference for human genetic variation. *Nature* **526**, 68–  
081 74 (2015).
- 082 109. International HapMap Consortium. The International HapMap Project. *Nature* **426**, 789–796 (2003).
- 083 110. Frankish, A. *et al.* GENCODE reference annotation for the human and mouse genomes. *Nucleic Acids Res.*  
084 **47**, D766–D773 (2019).
- 085 111. Yao, D. W., O’Connor, L. J., Price, A. L. & Gusev, A. Quantifying genetic effects on disease mediated by  
086 assayed gene expression levels. *Nat. Genet.* **52**, 626–633 (2020).
- 087 112. Frankish, A. *et al.* GENCODE 2021. *Nucleic Acids Res.* **49**, D916–D923 (2021).
- 088 113. Purcell, S. *et al.* PLINK: A Tool Set for Whole-Genome Association and Population-Based Linkage Analyses.  
089 *The American Journal of Human Genetics* vol. 81 559–575 Preprint at <https://doi.org/10.1086/519795>  
090 (2007).
- 091 114. Chang, C. C. *et al.* Second-generation PLINK: rising to the challenge of larger and richer datasets.  
092 *GigaScience* vol. 4 Preprint at <https://doi.org/10.1186/s13742-015-0047-8> (2015).
- 093 115. Li, H. *et al.* The Sequence Alignment/Map format and SAMtools. *Bioinformatics* vol. 25 2078–2079 Preprint  
094 at <https://doi.org/10.1093/bioinformatics/btp352> (2009).
- 095 116. Danecek, P. *et al.* The variant call format and VCFtools. *Bioinformatics* vol. 27 2156–2158 Preprint at  
096 <https://doi.org/10.1093/bioinformatics/btr330> (2011).
- 097 117. Price, A. L. *et al.* Long-range LD can confound genome scans in admixed populations. *American journal of*  
098 *human genetics* vol. 83 132–5; author reply 135–9 (2008).
- 099 118. Buitinck, L. *et al.* API design for machine learning software: experiences from the scikit-learn project. *arXiv*  
100 *[cs.LG]* (2013).
- 101 119. Robinson, M. D. & Oshlack, A. A scaling normalization method for differential expression analysis of RNA-  
102 seq data. *Genome Biol.* **11**, R25 (2010).
- 103 120. Gold, L. *et al.* Aptamer-based multiplexed proteomic technology for biomarker discovery. *PLoS One* **5**,  
104 e15004 (2010).
- 105 121. Loh, P.-R. *et al.* Reference-based phasing using the Haplotype Reference Consortium panel. Preprint at

- 106 <https://doi.org/10.1101/052308>.
- 107 122. Di Angelantonio, E. *et al.* Efficiency and safety of varying the frequency of whole blood donation  
108 (INTERVAL): a randomised trial of 45 000 donors. *Lancet* **390**, 2360–2371 (2017).
- 109 123. McCarthy, S. *et al.* A reference panel of 64,976 haplotypes for genotype imputation. *Nat. Genet.* **48**, 1279–  
110 1283 (2016).
- 111 124. Finucane, H. K. *et al.* Partitioning heritability by functional annotation using genome-wide association  
112 summary statistics. *Nature Genetics* vol. 47 1228–1235 Preprint at <https://doi.org/10.1038/ng.3404> (2015).
- 113 125. Hujoel, M. L. A., Gazal, S., Hormozdiari, F., van de Geijn, B. & Price, A. L. Disease Heritability Enrichment of  
114 Regulatory Elements Is Concentrated in Elements with Ancient Sequence Age and Conserved Function  
115 across Species. *Am. J. Hum. Genet.* **104**, 611–624 (2019).
- 116 126. Wen, X. MOLECULAR QTL DISCOVERY INCORPORATING GENOMIC ANNOTATIONS USING BAYESIAN FALSE  
117 DISCOVERY RATE CONTROL. *Ann. Appl. Stat.* **10**, 1619–1638 (2016).
- 118 127. Li, D. *et al.* WashU Epigenome Browser update 2022. *Nucleic Acids Res.* **50**, W774–W781 (2022).
- 119 128. Hudson, R. R., Slatkin, M. & Maddison, W. P. Estimation of levels of gene flow from DNA sequence data.  
120 *Genetics* **132**, 583–589 (1992).
- 121 129. Bhatia, G., Patterson, N., Sankararaman, S. & Price, A. L. Estimating and interpreting FST: the impact of rare  
122 variants. *Genome Res.* **23**, 1514–1521 (2013).
- 123 130. Band, G. & Marchini, J. BGEN: a binary file format for imputed genotype and haplotype data. *bioRxiv*  
124 308296 (2018) doi:10.1101/308296.
- 125