

1 Powerful mapping of *cis*-genetic effects on gene expression across diverse populations 2 reveals novel disease-critical genes

3 Kai Akamatsu^{1,2,3}, Stephen Golzari^{2,3,4}, and Tiffany Amariuta^{2,3}

- 4
- 5 1. School of Biological Sciences, UC San Diego, La Jolla, CA, USA.
- 6 2. Department of Medicine, Division of Biomedical Informatics, UC San Diego, La Jolla, CA, USA.
- 7 3. Halicioğlu Data Science Institute, UC San Diego, La Jolla, CA, USA.
- 8 4. Shu Chien-Gen Lay Department of Bioengineering, UC San Diego, La Jolla, CA, USA.
- 9

10 Abstract

11

12 While disease-associated variants identified by genome-wide association studies (GWAS) most
13 likely regulate gene expression levels, linking variants to target genes is critical to determining
14 the functional mechanisms of these variants. Genetic effects on gene expression have been
15 extensively characterized by expression quantitative trait loci (eQTL) studies, yet data from non-
16 European populations is limited. This restricts our understanding of disease to genes whose
17 regulatory variants are common in European populations. While previous work has leveraged
18 data from multiple populations to improve GWAS power and polygenic risk score (PRS)
19 accuracy, multi-ancestry data has not yet been used to better estimate *cis*-genetic effects on gene
20 expression. Here, we present a new method, Multi-Ancestry Gene Expression Prediction
21 Regularized Optimization (MAGEPRO), which constructs robust genetic models of gene
22 expression in understudied populations or cell types by fitting a regularized linear combination
23 of eQTL summary data across diverse cohorts. In simulations, our tool generates more accurate
24 models of gene expression than widely-used LASSO and the state-of-the-art multi-ancestry PRS
25 method, PRS-CSx, adapted to gene expression prediction. We attribute this improvement to
26 MAGEPRO's ability to more accurately estimate causal eQTL effect sizes ($p < 3.98 \times 10^{-4}$,
27 two-sided paired t-test). With real data, we applied MAGEPRO to 8 eQTL cohorts representing 3
28 ancestries (average $n = 355$) and consistently outperformed each of 6 competing methods in
29 gene expression prediction tasks. Integration with GWAS summary statistics across 66 complex
30 traits (representing 22 phenotypes and 3 ancestries) resulted in 2,331 new gene-trait associations,
31 many of which replicate across multiple ancestries, including *PHTF1* linked to white blood cell
32 count, a gene which is overexpressed in leukemia patients. MAGEPRO also identified
33 biologically plausible novel findings, such as *PIGB*, an essential component of GPI biosynthesis,
34 associated with heart failure, which has been previously evidenced by clinical outcome data.
35 Overall, MAGEPRO is a powerful tool to enhance inference of gene regulatory effects in
36 underpowered datasets and has improved our understanding of population-specific and shared
37 genetic effects on complex traits.

38 Introduction

39

40 Many genetic variants drive complex traits by regulating gene expression¹⁻⁸. Confident
41 characterization of genetic effects on gene expression is required for the functional interpretation
42 of disease-associated variants from genome-wide association studies (GWAS)⁹⁻¹¹. For example,
43 transcriptome-wide association studies (TWAS) integrate GWAS and gene expression data to
44 enable the identification of gene-disease associations, which can reveal genes underpinning
45 disease susceptibility, nominate candidate biomarkers for clinical use, or propel therapeutic
46 development¹²⁻¹⁴. Despite the potential to unravel the functional mechanisms of diseases, our
47 current understanding of disease-critical genes has been limited by variant-to-gene linking
48 strategies that rely heavily on sample size.

49

50 Although there is widespread availability of expression quantitative trait loci (eQTL) summary
51 statistics, such as across different human tissues from the Genotype-Tissue Expression (GTEx)¹⁵
52 project or from single cell RNA-sequencing data generated by eQTLGen¹⁶, datasets from non-
53 European populations are severely limited. Differences in allele frequency, linkage
54 disequilibrium (LD), and potentially causal variants reduce the applicability of genetic models
55 (of gene expression and complex traits alike) trained in European populations to non-European
56 populations¹⁷⁻²¹ and therefore limit the relevance of disease-gene associations detected by
57 European TWAS to other global populations. Therefore, there is an urgent need to more
58 accurately infer which genetic variants regulate gene expression and by how much, specifically
59 in understudied populations. Orthogonal to cross-ancestry fine-mapping of TWAS associations²²,
60 there also exists an opportunity to prune dense genomic loci with multiple gene-disease
61 associations to effects that are shared across ancestries, as causal genes are expected to be shared
62 across ancestries, more so in fact than causal variants.

63

64 Efforts to include diverse groups of individuals in genetic studies have yielded a modest number
65 of publicly available eQTL summary statistics from non-European populations²³⁻²⁷. Although the
66 statistical power of the eQTL studies performed in non-European populations remains
67 considerably weaker than that of European studies (6.5- and 2.6-fold difference in sample size
68 between European and African American individuals in GTEx¹⁵ and the Multi-Ethnic Study of
69 Atherosclerosis (MESA)²⁴, respectively), these data provide a unique opportunity to capture
70 varying genetic effects on gene expression across diverse ancestries. However, current gene
71 expression prediction models (such as LASSO, elastic net, and the best linear unbiased predictor
72 (BLUP) used in TWAS) can only model the limited individual-level genotype and gene
73 expression data from a single population to compute noisy estimates of variant-gene effect sizes.
74 Previous studies have proven the feasibility of leveraging data from multiple populations to
75 enhance GWAS association power²⁸, polygenic risk score (PRS) accuracy²⁹⁻³¹ and GWAS fine-
76 mapping^{32,33}. Thus, we hypothesized that multi-ancestry data would enhance the construction of
77 *cis*-genetic models of gene expression by improving the estimation of variant-level effects and

78 overall expression prediction accuracy. Current multi-ancestry TWAS approaches do not tackle
79 the issue of large uncertainty of inferred *cis*-genetic effects on gene expression in small non-
80 European cohorts. For example, TESLA improves association power by colocalizing a single
81 eQTL dataset with a cross-population meta-analysis of GWAS summary statistics, producing
82 results with mixed or uncertain relevance to each ancestry³⁴. Another approach called METRO
83 models the uncertainty of gene expression models across multiple cohorts to maximize
84 colocalization with GWAS³⁵, resulting in findings that are highly driven by European data when
85 other gene models are derived from smaller non-European datasets. To date, multi-ancestry data
86 has not been used to reduce uncertainty and improve accuracy of population-specific genetic
87 models of gene expression.

88

89 Here, we introduce a new method, Multi-Ancestry Gene Expression Prediction Regularized
90 Optimization (MAGEPRO), that improves gene expression prediction accuracy in underpowered
91 ancestries or undersampled tissues by optimally combining eQTL summary statistics from
92 ancestrally and functionally diverse datasets. We evaluate the robustness of our method in
93 various simulated genetic architectures and compare the predictive performance of MAGEPRO
94 to alternative methods of gene expression prediction, including an adaptation of a multi-ancestry
95 complex trait PRS method called PRS-CSx³⁰, using 8 different eQTL cohorts representing 3
96 ancestries. We additionally applied MAGEPRO gene models to perform TWAS with 15 blood-
97 cell traits and 7 immune-mediated diseases, each represented by GWAS cohorts of individuals of
98 African, European, and Hispanic ancestries, to identify novel disease-gene associations and
99 interrogate the population-specificity of these putative disease genes.

100 **Results**

101

102 ***Overview of MAGEPRO***

103

104 MAGEPRO maximizes our ability to infer gene regulatory effects in small sample size eQTL
105 datasets and constructs robust *cis*-genetic models of gene expression that are specific to an
106 ancestry. Given individual-level genotype and gene expression data of the target cohort and
107 external eQTL data from diverse ancestries and tissues, MAGEPRO first estimates effect sizes
108 for single nucleotide polymorphism (SNP)-gene pairs in *cis* that are specific to the target
109 population via a LASSO (L1 norm)-regularized linear regression (**Figure 1**, green box). This
110 step constitutes the conventional TWAS gene expression prediction model. Next, MAGEPRO
111 applies the Sum of Single Effects (SuSiE)^{36,37} regression model to each set of external eQTL
112 summary statistics to identify putative causal variants and estimate posterior effect size estimates
113 for all *cis*-variants (**Figure 1**, blue box). Assuming most causal variants are shared, this step is
114 critical to maximizing the cross-population transferability of information from external datasets
115 to the target cohort. Causal variants are more likely to possess predictive power in the target
116 population compared to variants that merely tag the causal variant; specifically, the causal variant
117 may not be sufficiently tagged in the target population if there are differences in linkage
118 disequilibrium and allele frequency between training and target populations. Finally, our
119 approach finds an optimal ridge (L2 norm)-regularized linear combination of posterior effect size
120 estimates from SuSiE and the target population SNP-gene weights to produce the final gene
121 expression prediction model (**Figure 1**, white box). By utilizing existing fine-mapping
122 frameworks and regularizing the combination of SNP-gene weights across datasets, MAGEPRO
123 is designed to include only information that is potentially relevant to the target population, as
124 opposed to other strategies such as METRO (see above) or a meta-analysis approaches where
125 inferred effect sizes are driven by the largest (European) datasets in the analysis.

126

127 Throughout this study, we compare MAGEPRO to several methods for gene expression
128 prediction. These include single-ancestry methods commonly used in TWAS, such as LASSO
129 regression^{12,14,38,39}, and multi-ancestry approaches, such as a cross-population meta-analysis of
130 eQTL summary statistics. We also utilized methods that are conventionally applied to gene
131 expression or GWAS data, like SuSiE^{36,37} and pruning and threshold (P+T)^{17,40,41}. Notably, we
132 benchmarked our tool against a variation of MAGEPRO that we refer to as Multipop, which does
133 not use SuSiE, but rather fits a ridge (L2 norm)-regularized linear combination of raw eQTL
134 summary statistics. Lastly, we benchmarked MAGEPRO against PRS-CSx³⁰, a state-of-the-art
135 multi-ancestry PRS method for genome-wide complex trait/disease data. PRS-CSx is a Bayesian
136 framework that models LD heterogeneity across datasets and infers a shared shrinkage parameter
137 to enforce sparsity, which assumes that causal effects are shared, a common assumption of most
138 multi-ancestry fine-mapping models^{32,42}. While PRS-CSx is a popular choice for PRS using
139 ancestrally diverse GWAS data⁴³⁻⁴⁹, this method has not yet been applied to integrate cross-

140 population eQTL summary statistics to create more predictive models of gene expression. In our
141 study, we compare gene expression prediction accuracy ($\frac{R_{CV}^2}{\hat{h}_{ge}^2}$) between methods, which is defined
142 as the fraction of gene expression variance explained by the model in cross-validation (R_{CV}^2),
143 normalized by the upper limit of the prediction: the *cis*-heritability estimated by GCTA⁵⁰ (\hat{h}_{ge}^2).
144 Each competing method is described in further detail in Methods.

145

146 **Simulations**

147

148 We performed extensive simulations to compare the performance of MAGEPRO to the most
149 popular approaches, LASSO for single-ancestry and PRS-CSx for multi-ancestry, under various
150 genetic architectures, using code adapted from the Mancuso Lab TWAS simulator (Code
151 Availability)^{51,52}. We used real genotypes from the 1000 Genomes Project¹⁸ as LD reference
152 panels to simulate genotypes and *cis*-regulated gene expression data across African, European,
153 and American ancestries (Methods). We compared the 5-fold cross-validation accuracy of each
154 model in predicting *cis*-regulated gene expression in African individuals (target), using simulated
155 European and American summary statistics (external) for both PRS-CSx and MAGEPRO. In our
156 primary analysis, we simulated genes with four causal *cis*-eQTLs shared across populations with
157 correlated true effect sizes ($r = 0.8$); we varied target population sample sizes, the heritability of
158 gene expression, and the number of causal *cis*-eQTLs. In secondary analyses, we varied whether
159 or not eQTL effects were correlated across ancestries, changed whether or not there were
160 ancestry-specific causal *cis*-eQTLs in high LD with the causal variant of the target ancestry, and
161 lastly, evaluated if MAGEPRO can still improve the accuracy of gene models when SuSiE fails
162 to identify a likely causal variant. More details on our simulation framework are described in
163 Methods and the Supplementary Note.

164

165 Within our primary analyses, we first compared the prediction accuracy of the three methods,
166 calculated as $\frac{R_{CV}^2}{\hat{h}_{ge}^2}$ (see above) across target population sample sizes ranging from 80 to 500
167 individuals and gene expression heritability ranging from 5% to 40%. Across 1,000
168 independently simulated genes, MAGEPRO outperformed both LASSO and PRS-CSx in each of
169 20 different sample size and *cis*-heritability settings with an average improvement of 5.7% and
170 4.5% in accuracy, respectively (**Figure 2A, Supplementary Tables 1-2**). Generally, larger
171 sample sizes of the target population resulted in more accurate predictions for a given
172 heritability; and, accuracy notably increased and began to approach 100% for each method
173 within the most heritable genes (40%), thanks to the larger and more easily identifiable eQTL
174 effects. The utility of MAGEPRO is most clearly demonstrated at smaller sample sizes and
175 higher gene expression heritability (**Supplementary Figures 1-2**), enhancing accuracy by > 9%
176 compared to LASSO ($p < 1.4 \times 10^{-56}$) and by > 7% compared to PRS-CSx ($p < 2.3 \times 10^{-49}$)
177 when the sample size of the target cohort is 80 individuals and the heritability of the gene is \geq
178 20%. For lowly heritable genes, MAGEPRO demonstrates an increasing margin of advantage

179 over the other two methods as sample sizes grow (**Supplementary Figures 1-2**), suggesting that
180 MAGEPRO may be especially useful for modeling the genetic architecture of disease-critical
181 genes whose regulatory effects are flattened by natural selection and thus have lower *cis*-
182 heritability^{53,54}.

183
184 We further hypothesized that MAGEPRO would achieve superior prediction accuracy by
185 estimating more accurate eQTL effect sizes. Indeed, when we compare the squared difference
186 between simulated (true) and estimated causal eQTL effect sizes, MAGEPRO produces smaller
187 errors compared to both competing methods across the five different sample sizes at 10% gene
188 expression heritability (all $p < 3.98 \times 10^{-4}$, **Figure 2B, Supplementary Tables 3-4**). Although
189 the accuracy of causal eQTL effect sizes is not a requirement for prediction methods (e.g.,
190 prediction can be achieved with strong tagging variants), we believe this characteristic of
191 MAGEPRO may lead to more accurate results from downstream gene-based association analysis
192 like TWAS.

193
194 We also evaluated each method across genetic architectures with varying numbers of causal *cis*-
195 eQTLs while maintaining a constant 10% *cis*-heritability and target sample size of 240, which is
196 synonymous with decreasing the per-SNP heritability ($\frac{h_{ge}^2}{m \text{ causal eQTLs}}$). Overall, as the per-SNP
197 heritability decreases, the prediction accuracy of all methods decreases due to the difficulty of
198 capturing larger quantities of smaller effects (**Figure 2C, Supplementary Tables 5-6**),
199 exemplifying the challenge of modeling the genetic regulation of disease-critical genes, which
200 are more likely to have lower *cis*-heritability (see above). Despite this challenge, MAGEPRO
201 outperformed both LASSO and PRS-CSx in each per-SNP heritability setting (all $p <$
202 1.5×10^{-5}), while PRS-CSx notably surpassed the accuracy of LASSO for the two lower per-
203 SNP heritability settings. This indicates that at current eQTL study sample sizes, leveraging
204 multi-ancestry data is a useful tool for accurately modeling the genetic regulation of potentially
205 disease-relevant genes and may help more confidently identify which diseases they influence via
206 gene-based association tests.

207
208 In secondary analyses, we tested the performance of MAGEPRO when the effect sizes of shared
209 causal *cis*-eQTLs are drawn independently across ancestries and are thus uncorrelated. Although
210 MAGEPRO achieves larger improvements relative to LASSO and PRS-CSx when effect sizes
211 are correlated across ancestries, our tool robustly improves prediction accuracy even when effect
212 sizes are independent (**Supplementary Figure 3**) and trends across sample sizes and
213 heritabilities are largely shared with simulations with correlated eQTL effects. Recent work
214 shows that effect size correlations across ancestries are lower for loss-of-function intolerant
215 genes³⁹ and variants with ancestry-specific disease effects may reside closer to genes interacting
216 with the environment, such as immune responses⁵⁵. This suggests that MAGEPRO will continue
217 to improve gene model accuracy, even when causal eQTL effect sizes are independent, which
218 could potentially lead to the discovery of novel gene-disease associations. In a related

219 framework, we simulated gene expression prediction models based on a single causal eQTL in
220 the target African population. In this analysis, the single causal eQTL is not shared across any
221 ancestries, but the two causal variants from the European and American populations are in high
222 LD with the causal variant of the target population (**Supplementary Figure 4**). Overall, we
223 observed highly similar trends with that of **Figure 2**; in fact, the accuracies across sample sizes
224 and heritabilities were greater than in **Figure 2** due to the fact that per-SNP heritability was
225 proportionally higher thanks to simulating a single causal variant.

226

227 Lastly, we explored whether the improvement in accuracy provided by MAGEPRO depends on
228 the ability of SuSiE to identify causal *cis*-eQTLs in external datasets. The enhancement of
229 prediction accuracy relative to LASSO is nominally larger when SuSiE identifies at least 1
230 causal *cis*-eQTL ($PIP \geq 0.95$) across the external datasets and this difference is only statistically
231 significant at the largest target population sample size of 500 ($p = 0.03$) (**Supplementary**
232 **Figure 5**). This implies that although isolating the causal regulatory variants contributes to
233 improved prediction, MAGEPRO does not rely on fine-mapped SNPs with high PIPs, but rather
234 on posterior effect size estimates.

235

236 ***Benchmarking MAGEPRO against alternative gene expression prediction methods***

237

238 In real data analysis, we employed MAGEPRO to create *cis*-genetic models of gene expression
239 for 8 eQTL cohorts across 3 different ancestries (average $n = 355$) using up to 5 external
240 summary statistic datasets as features in the MAGEPRO model (**Table 1**)^{15,16,24,25,27}. For each
241 gene, we performed variable selection, e.g., eQTL fine-mapping, applying SuSiE to each
242 summary statistic dataset (Methods). We explored the possibility of leveraging IMPACT, a tool
243 we have previously developed to estimate the probability that a variant participates in cell-type-
244 specific gene regulation⁵⁶, as Bayesian SNP-selection priors in SuSiE have been shown to
245 improve fine-mapping power⁵⁷. Although this increased the number of genes with at least 1
246 putatively causal eQTL (posterior inclusion probability (PIP) ≥ 0.95), increased average PIPs in
247 credible sets, and decreased average credible set size, it did not substantially affect the accuracy
248 of MAGEPRO gene models (**Supplementary Figures 6-7**). Even random priors seemed to
249 improve fine-mapping metrics, likely by randomly pruning high PIP variants in high LD; but,
250 ultimately the predictive capacity of posterior effect size estimates do not strictly depend on
251 reduced credible set size and high PIP SNPs, thus the gene model accuracy is not necessarily
252 affected (**Supplementary Figure 6**). These results are consistent with our simulations that
253 indicated MAGEPRO need not find a putatively causal eQTL to enhance prediction accuracy
254 relative to LASSO. Therefore, we elected to not use IMPACT priors in the default
255 implementation of MAGEPRO.

256

257 Next, we applied GCTA to each target eQTL cohort to estimate the *cis*-heritability (h_{ge}^2) of each
258 gene. For genes with larger *cis*-heritability estimates, SuSiE detected a larger number of

259 putatively causal eQTLs on average ($PIP \geq 0.95$) (**Supplementary Figure 8**). We also observed
260 that the estimated *cis*-heritabilities of gene expression were highly correlated across ancestries,
261 consistent with previous work²² (Pearson correlation (r) ranging from 0.32 to 0.83 in
262 comparisons between European, Hispanic/Latino, and African American populations)
263 (**Supplementary Figure 9**). However, we observed similar heterogeneity of heritability
264 estimates even across cohorts within the same ancestry ($r = 0.34$, 95% CI [0.311, 0.367])
265 between European individuals in GEUVADIS and GENOA cohorts), suggesting that cross-
266 cohort variation may limit out-of-cohort prediction accuracy.

267

268 We next compared the performance of various methods in predicting expression levels of
269 significantly *cis*-heritable genes in each target cohort ($GCTA \hat{h}_{ge}^2 > 0; p < 0.01$). These
270 methods, introduced above and in more detail in Methods, comprise a cross-population meta-
271 analysis, pruning and thresholding (P+T) of target marginal *cis*-eQTL, LASSO of the target
272 population, SuSiE applied to the target population, a ridge (L2 norm) regression of full external
273 *cis*-eQTL summary statistics (which we refer to as “Multipop”), PRS-CSx, and MAGEPRO. We
274 note that not all external summary statistics contain associations for all genes, and thus
275 MAGEPRO utilizes only relevant external datasets available to each gene.

276

277 First, we applied each method to predict lymphoblastoid cell line (LCL) gene expression in the
278 Genetic Epidemiology Network of Arteriopathy (GENOA) African American (AA) cohort ($n =$
279 346). MAGEPRO outperformed all competing methods (all paired one sided t-test $p <$
280 3×10^{-10}) and improved prediction accuracy by 10.4% relative to LASSO averaged across
281 4,141 *cis*-heritable genes (**Figure 3A, Supplementary Table 7**). MAGEPRO’s accuracy
282 exceeded that of Multipop ($p = 8 \times 10^{-32}$), suggesting that the posterior effect sizes estimated
283 by SuSiE are prioritizing variants that are critical in predicting gene expression. Notably, our
284 model increased prediction accuracy relative to LASSO by over 20% for 1,177 genes and
285 introduced 204 new genes with an R_{cv}^2 significantly greater than 0 ($p < 0.05$). We then down-
286 sampled the GENOA AA cohort ($n = 100$) to challenge MAGEPRO in a small sample size
287 setting (one that is similar to the number of African American individuals in GTEx). We found
288 that MAGEPRO maintains improved accuracy compared to all methods when target population
289 genotype and gene expression data is extremely limited ($p < 0.01$ across all comparisons,
290 **Figure 3B, Supplementary Table 8**). At this sample size, we achieved a 4.4% improvement in
291 accuracy relative to PRS-CSx ($p = 4 \times 10^{-10}$), suggesting that the layers of regularization in
292 our framework minimize overfitting even with small training cohorts.

293

294 We observed similar trends across all 8 target eQTL cohorts (10 including down-sampled
295 cohorts). In predicting monocyte gene expression in the Multi-Ethnic Study of Atherosclerosis
296 (MESA) Hispanic/Latino (HIS) cohort, MAGEPRO again outperformed all competing methods
297 (all $p < 6 \times 10^{-21}$), improving prediction accuracy relative to LASSO by over 20% for 942
298 genes and creating 191 new gene models with significantly positive R_{cv}^2 (**Supplementary Figure**

299 **10)** MAGEPRO improved prediction accuracy relative to LASSO by 14.7% in the GTEx AA
300 cohort ($n = 80$, Whole Blood) and by 13.5% in the down-sampled GEUVADIS European (EUR)
301 cohort ($n = 100$, LCL), suggesting that our method provides the largest relative improvement
302 when the target cohort sample size is limited (**Figure 3C, Supplementary Table 9**).

303
304 Next, we aimed to characterize the genes for which MAGEPRO is most useful for capturing the
305 *cis*-genetic component of expression. We observed that the change in accuracy between
306 MAGEPRO and LASSO ($\text{MAGEPRO } \frac{R_{cv}^2}{\hat{h}_{ge}^2} - \text{LASSO } \frac{R_{cv}^2}{\hat{h}_{ge}^2}$) is negatively correlated with *cis*-
307 heritability estimates ($r = -0.14$, $p = 3.7 \times 10^{-20}$ and $r = -0.17$, $p = 4.7 \times 10^{-23}$ for
308 GENOA AA and MESA HIS respectively; **Figure 3D, Supplementary Table 10,**
309 **Supplementary Figure 11**). This indicates that MAGEPRO offers the greatest modeling
310 improvements to low heritability genes, which are more likely to be disease-critical, as natural
311 selection restricts the magnitude of *cis*-genetic effects (and thus heritability) on disease-critical
312 genes. For example, we found that loss-of-function intolerant genes ($pLI > 0.9$)⁵⁸ indeed have the
313 lowest gene expression heritability estimates (**Supplementary Figure 12**, $p < 7.0 \times 10^{-8}$).
314 Additionally, we found that MAGEPRO offers the greatest advantage over PRS-CSx when the
315 per-SNP heritability of the gene ($\frac{\hat{h}_{ge}^2}{\# \text{ SNPs with } PIP \geq 0.95}$), which is proportional to the power to detect
316 *cis*-genetic effects³¹, is low (**Supplementary Figure 13**).

317
318 We also evaluated the generalizability of each model to individuals from a different study cohort
319 in the same target ancestry. To this end, we compared out-of-cohort prediction accuracy. We
320 trained gene expression prediction models in GENOA AA and GEUVADIS EUR cohorts, each at
321 two different sample sizes, and then applied these models to predict LCL gene expression in
322 GEUVADIS Yoruba (YRI) and GENOA EUR cohorts, respectively. MAGEPRO and SuSiE
323 consistently outperformed the other methods (LASSO, Multipop, PRS-CSx) in out-of-cohort
324 prediction, suggesting that frameworks which prioritize putative causal eQTL may result in more
325 generalizable predictive models (**Supplementary Figure 14**). We note that we did not assess
326 cross-population meta-analysis or P+T in this analysis, as they performed much more poorly in
327 within-cohort cross-validation tasks. However, the performance of MAGEPRO relative to SuSiE
328 (applied directly to the training population) was highly variable. For example, the SuSiE model
329 trained in the down-sampled GENOA AA cohort ($n = 100$) achieved a higher out-of-cohort R^2
330 than MAGEPRO ($p = 0.006$, **Supplementary Figure 14**), possibly due to the different extent of
331 admixture between African American (training) and Yoruba individuals (testing)
332 (**Supplementary Figure 15**) or due to the inherent cross-cohort variation in the genetic
333 architecture of gene expression that we previously observed (**Supplementary Figure 9**). In
334 contrast, the MAGEPRO model trained in the down-sampled GEUVADIS EUR cohort ($n =$
335 100) exceeded SuSiE in out-of-cohort prediction by 6% ($p = 8.6 \times 10^{-9}$, **Supplementary**
336 **Figure 14**). MAGEPRO generally excels in out-of-cohort prediction when the genetic ancestry
337 of the training and testing cohorts are closely related (**Supplementary Figures 14-15**),

338 highlighting the population-specific nature of MAGEPRO models. In other words, SuSiE applied
339 to the target training population is effective at assaying causal variants that are likely to be shared
340 across populations, but more population-specific effects may be identified by MAGEPRO, which
341 is tailored to the training population.

342
343 We found that MAGEPRO is consistently most useful when the target population genotype and
344 gene expression data is limited. We hypothesized that this may include situations where the
345 target tissue is less accessible and/or data is scarce. Therefore, we explored if genetic models of
346 gene expression in tissues that are seemingly unrelated to blood can be improved by integrating
347 widely available blood-derived eQTL summary statistics. To this end, we applied MAGEPRO to
348 create Lung gene models in GTEx using blood-related external *cis*-eQTL summary statistics
349 (**Table 1**). MAGEPRO produced impressively accurate gene models (59% on average) while
350 outperforming all competing methods (all $p < 1 \times 10^{-46}$), likely owing to the correlation of *cis*-
351 genetic regulation of gene expression across tissues¹⁵ (**Supplementary Figure 16**), not unlike the
352 cross-population sharing of causal effects. Moreover, this suggests that MAGEPRO successfully
353 identifies regulatory effects from blood tissue that are transferable to lung tissue, notably
354 resulting in an 8.4% average improvement over the lung-specific LASSO model ($p =$
355 9×10^{-307}).

356
357 We implemented MAGEPRO as a publicly available pipeline on GitHub (Code Availability),
358 leveraging multiple threads on both high-performance computing (HPC) clusters⁵⁹ and personal
359 devices to enhance computational efficiency (**Supplementary Figure 17**).

360
361 ***Transcriptome-wide association studies are sensitive to cis-genetic models of gene expression***

362
363 We hypothesized that one of the most compelling applications of MAGEPRO would be to make
364 the inference of disease-critical genes more powerful for underrepresented populations. To this
365 end, we applied LASSO, SuSiE, PRS-CSx, and MAGEPRO models trained in 7 blood-related
366 eQTL cohorts (MESA AA Monocyte, GENOA AA LCL, GTEx AA Whole Blood, MESA EUR
367 Monocyte, GEUVADIS EUR LCL, GTEx EUR Whole Blood, MESA HIS Monocyte) to perform
368 TWAS for 15 blood cell traits and 7 immune-mediated diseases using ancestry matched GWAS
369 summary statistics from Chen and colleagues⁶⁰ (AFR $N = 13,391$, EUR $N = 516,979$, HIS $N =$
370 $6,849$) and the Global Biobank Meta-analysis Initiative (GBMI)⁶¹ (AFR $N = 26,052$, EUR $N =$
371 $1,024,298$, Native American ancestry (AMR) $N = 15,490$), respectively (**Supplementary Table**
372 **11**). We note that we did not have access to AMR eQTL data and, therefore, we used HIS gene
373 expression prediction models as proxies to perform TWAS in the AMR population. To avoid
374 complicated notation, we refer to subsequent TWAS analysis involving HIS eQTL data and
375 AMR GWAS as HIS. Generally, we observed two main phenomena. In one case, MAGEPRO
376 models led to more accurate *cis*-genetic models of gene expression (relative to LASSO), and this
377 subsequently eliminated the statistically significant TWAS association observed for LASSO. In

378 the other case, MAGEPRO generated predictive gene expression models (significantly positive
379 R^2) even though LASSO failed to do so; this resulted in many new gene-trait/disease
380 associations, exemplifying the utility of MAGEPRO to enhance disease inference in
381 underpowered cohorts and underrepresented populations. Ultimately, both of these scenarios
382 allowed us to explore the sensitivity of TWAS to slight variations in *cis*-genetic gene models.
383 We explore examples of both cases below in more depth.

384

385 First, we observed that the average change in gene expression prediction R^2 (MAGEPRO R^2 –
386 LASSO R^2) does not correlate with the average change in TWAS chi-square statistic (χ^2)
387 (MAGEPRO TWAS χ^2 - LASSO TWAS χ^2) across significantly *cis*-heritable genes
388 (**Supplementary Figure 18**). This result is not surprising as few genes play critical roles for any
389 one disease, and MAGEPRO is able to improve the mapping of *cis*-genetic effects for both
390 disease-critical and non-critical genes. However, this observation led us to understand that
391 sometimes an improved gene expression prediction model may actually produce a weaker TWAS
392 association, implying that less accurate gene models were only spuriously correlated with
393 disease. In other words, MAGEPRO provides an additional utility of enhancing the confidence in
394 TWAS association results by increasing the gene expression prediction accuracy. While TWAS is
395 most well-powered to identify genes with large *cis*-genetic effects that colocalize with disease,
396 our observation here does not invalidate the compelling nature of our previous finding that
397 MAGEPRO produces the largest improvements in model accuracy for low heritability genes,
398 which due to natural selection may be more disease-critical. Therefore, by learning more
399 accurate *cis*-genetic models of gene expression, MAGEPRO may be additionally poised to help
400 derive disease-critical effects on gene expression in frameworks beyond TWAS.

401

402 There were several genes for which the conventional single population TWAS model produced a
403 significant TWAS association that was ablated when the gene model was improved with
404 MAGEPRO. For example, the association between *ZNF213-AS1* and red blood cell count in the
405 African American population diminished as MAGEPRO improved the accuracy of gene
406 expression prediction (**Figure 4A, Supplementary Table 12**). Investigating how the *cis*-genetic
407 model of gene expression colocalizes with GWAS summary statistics reveals that the
408 MAGEPRO model captured a new eQTL signal (“MAGEPRO-specific” in teal), improving gene
409 expression prediction accuracy (from 24% with SuSiE or 33% with LASSO to 45% with
410 MAGEPRO) but providing conflicting evidence against the negative association with the GWAS
411 phenotype (**Figure 4A**). *ZNF213-AS1* is a noncoding antisense RNA gene which controls breast
412 cancer progression by modulating estrogen receptor signaling^{62,63}, but links to blood-related
413 phenotypes have not been reported in the literature. Additionally, this association was not found
414 in the European TWAS ($z = -2.8$, not significant [n.s.]), although the gene model achieved near
415 perfect accuracy. To summarize, while TWAS does not account for the uncertainty of gene
416 expression models, our findings suggest that considering association statistics across different

417 models for the same gene can reveal unstable gene-disease associations and potentially false
418 positives.

419
420 Second, we observed that modest changes to *cis*-genetic models of gene expression can also give
421 rise to biologically plausible new disease-gene associations. For instance, *RGS14* was not
422 analyzed in the European TWAS using LASSO because the model produced an R^2 that was not
423 significantly greater than 0 (**Figure 4B, Supplementary Table 13**). The MAGEPRO model
424 introduced a new eQTL signal (teal dotted line), which helped the model achieve a significantly
425 positive R^2 ($p < 0.05$) and provided additional evidence to the negative association with asthma
426 (**Figure 4B**). The estimated heritability (\hat{h}_{ge}^2) of *RGS14* was only 0.03 (se = 0.015), reflecting the
427 inherent difficulty in modeling genetic effects on genes with low heritability and the utility of
428 MAGEPRO for detecting putative disease-critical genes that could not previously be reliably
429 analyzed. *RGS14* belongs to a family of proteins that regulate G protein signaling, which plays a
430 significant role in asthma^{64,65}. Current asthma therapies include G protein signaling agonists and
431 antagonists, which relax airway smooth muscles and reduce airway inflammation, respectively⁶⁶.
432 Our finding suggests that regulatory variants modulating G protein signaling may carry genetic
433 risk for asthma.

434 435 ***MAGEPRO recapitulates gene-disease associations across diverse ancestries and reveals*** 436 ***ancestry-specific findings***

437
438 Now that we understand the dominant mechanisms by which MAGEPRO can inform gene-
439 disease association studies (e.g., by ablating the significant association produced by less accurate
440 models, or by producing significant associations for genes that previously lacked predictive
441 models), we sought to apply our models across diverse ancestries to characterize population-
442 specific or population-shared gene-level effects on complex traits and diseases. We organized our
443 analysis into two disjoint sets of genes: those with fairly accurate predictive models ($R^2 > 0$,
444 $p < 0.05$) across all methods (LASSO, SuSiE, PRS-CSx, MAGEPRO) and those that lacked a
445 predictive LASSO model.

446
447 We first analyzed all genes with a gene expression prediction R^2 significantly greater than 0 in
448 all methods. Aggregating results across 7 blood-related eQTL cohorts and 66 GWAS summary
449 statistics (accounting for 22 unique diseases/traits and 3 ancestries), MAGEPRO identified 2,521
450 gene-trait associations ($p < \frac{0.05}{\# \text{ genes tested in dataset}}$) that were not found by LASSO
451 (**Supplementary Table 14**). Considering all four methods, we found that MAGEPRO identified
452 1,350 significant gene-trait associations that are not identified by any other model
453 (**Supplementary Table 15**), showcasing the benefit of MAGEPRO in augmenting current gene
454 expression prediction models in the TWAS framework. However, MAGEPRO gene models do
455 not necessarily generate more significant gene-trait associations than other methods
456 (**Supplementary Figure 19**). This is because improving genetic models of gene expression

457 yields TWAS results that are more reliable, but not necessarily stronger in association as we
458 discussed previously (**Figure 4A, Supplementary Figure 18**). When we applied Monocyte gene
459 models trained in MESA African American individuals to TWAS, MAGEPRO found 8
460 significant associations not identified by LASSO (6 of them as a result of larger gene model R^2)
461 (**Figure 5A, Supplementary Table 16**) and 20 significant associations not found by PRS-CSx
462 (10 of them as a result of larger gene model R^2) (**Figure 5B, Supplementary Table 17**). In
463 contrast, when we applied our LCL gene models trained in GENOA African American
464 individuals, PRS-CSx identified 16 associations not found by MAGEPRO (**Supplementary**
465 **Figure 19**). However, MAGEPRO produced a more accurate genetic model of gene expression
466 for 9 of these 16 genes, suggesting that a majority of the gene-trait associations undetected by
467 MAGEPRO may be false positives, or at the least, unreliable associations. We found similar
468 patterns when comparing TWAS associations across MAGEPRO, LASSO, and PRS-CSx in
469 Hispanic/Latino individuals (**Supplementary Figure 20**), although the limited GWAS sample
470 size for this population greatly reduced our power to assess patterns of gene-trait associations
471 across methods. Reflecting on our results, our suggested best practice is to use the most accurate
472 *cis*-genetic model of gene expression for each gene, as similarly implemented in FUSION.
473 Although it does not always lead to more statistically significant gene-trait associations
474 (**Supplementary Figure 19**), TWAS results will be more credible when the gene expression
475 prediction models are more accurate.

476
477 Second, we explored how improving genetic models of gene expression in underpowered
478 ancestries can help us challenge or recapitulate results from European TWAS studies. To this
479 end, we investigated TWAS results for white blood cell (WBC) count using Monocyte gene
480 models developed for European, African American, and Hispanic populations; we focus on 4
481 associations that were consistent across at least two populations: *PHTF1*, *LAMTOR2*, *PTPN22*,
482 and *LMNA* (**Figure 5C, Supplementary Table 18**). *PHTF1* was not evaluated in African-
483 ancestry TWAS with LASSO because the gene model R^2 was not significantly greater than 0.
484 However, MAGEPRO improved this gene expression prediction model and identified a positive
485 association with WBC count, recapitulating findings from the European population (**Figure 5C,**
486 **Supplementary Figure 21**). *PHTF1* has been associated with other immune-mediated diseases,
487 such as type 1 diabetes in early genetic studies⁶⁷. Additionally, differential expression analysis
488 has shown that this gene is overexpressed in patients with acute lymphoblastic leukemia⁶⁸, a
489 condition characterized by the overproduction of immature white blood cells. This indicates that
490 *PHTF1* is a plausible candidate for regulating white blood cell count and extreme dysregulation
491 of this gene may be linked to forms of leukemia. Furthermore, leveraging MAGEPRO to
492 improve the genetic model of gene expression for *LAMTOR2* by 54% resulted in a new
493 association for individuals of African ancestry, which is consistent with findings from European
494 TWAS (**Figure 5C, Supplementary Figure 21**). Previous work shows that experimental
495 knockout of *LAMTOR2* results in an expansion of conventional dendritic cells in mice⁶⁹ and the
496 deficiency of this gene causes immunodeficiency syndromes in humans^{70,71}. The replication

497 across ancestries and the layers of evidence in the literature suggest that *LAMTOR2* is another
498 candidate regulator of white blood cell count in humans. *PTPN22*, a well-known regulator of
499 immune signaling^{72–76}, and *LMNA*, a major component of the mammalian lamina with important
500 functions in immune cells⁷⁷, was also identified by TWAS for both African and European
501 ancestries using either LASSO or MAGEPRO models. Our findings demonstrate that applying
502 MAGEPRO to improve genetic models of gene expression in understudied populations can help
503 identify potentially causal disease/trait-associated genes that replicate across different ancestries.
504

505 Third, we evaluated MAGEPRO's capacity to identify ancestry-specific gene-trait associations.
506 To achieve this, we analyzed genes with a gene expression prediction R^2 significantly greater
507 than 0 in both LASSO and MAGEPRO and used the better-performing model for TWAS. We
508 identified 137 associations in African or Hispanic populations which were not found in European
509 TWAS (**Supplementary Table 19**). Among these, 13 genes were exclusively identified by
510 MAGEPRO, 5 by LASSO, and 119 by both methods. Notably, MAGEPRO improved the
511 predictive performance of the *UBAP2L* Monocyte gene model in the African American
512 population, modestly raising the R^2 from 0.10 (LASSO) to 0.11. As a result, MAGEPRO
513 detected an association between *UBAP2L* and neutrophil count (NEU) ($z = -6.02$), which was not
514 found by any European model across monocyte, LCL and whole blood tissues. Previous
515 experimental studies have demonstrated that *UBAP2L* plays a crucial role in the regulation of
516 long-term hematopoietic stem cells⁷⁸, supporting its potential as a candidate regulator of
517 neutrophil counts.
518

519 Lastly, we sought to use MAGEPRO to identify disease-critical roles specifically for genes that
520 lacked a predictive LASSO model (R^2 not significantly positive), and thus could not be
521 previously analyzed by TWAS. In this category, MAGEPRO offered 3,195 new gene models
522 across 7 eQTL cohorts. The *cis*-genetic effects of these genes were inherently difficult to model
523 due to the low heritability of gene expression (average $\hat{h}_{ge}^2 = 0.095$, lowest quantile in **Figure**
524 **3D**). Nevertheless, MAGEPRO enhanced the average R^2 of these models from 0.0047 with
525 LASSO to 0.031 (a 560% increase). Applying these newly modeled genes to TWAS across all
526 66 traits yielded 981 associations at Bonferroni significance, where a different threshold was
527 determined for each of 7 eQTL cohorts (**Figure 6, Supplementary Table 20**). Several of these
528 associations recapitulate existing results from colocalization analysis using European GWAS.
529 For example, European MAGEPRO models identified an association of *IRF8*⁷⁹ to monocyte
530 count (MON) and *RCCD1*^{80,81} to red blood cell distribution width (RDW), which are consistent
531 with European colocalization analyses^{82,83} (**Figure 6**). Additionally, some of these associations
532 replicate previously reported European TWAS results in an understudied ancestry. For instance,
533 the relationship between *FAM234* and mean corpuscular hemoglobin concentration (MCHC) has
534 been established in European TWAS^{84–86}, but to our knowledge has not been reported using
535 genetic associations from individuals of African ancestry until now (**Figure 6**).
536

537 The new MAGEPRO gene models also resulted in biologically plausible novel findings. For
538 example, African American MAGEPRO models for whole blood identified an association
539 between *SH2D1B* and *SLAMF8* to both neutrophil count (NEU) and white blood cell (WBC)
540 count (**Figure 6**). Multiple lines of evidence support that the proteins encoded by these two
541 genes interact to control immune response^{87,88,89}, and some studies have promoted SLAM
542 receptors as potential therapeutic targets for immune-mediated diseases⁹⁰. Improved European
543 genetic models of gene expression for whole blood also revealed an association between *PIGB*
544 and heart failure, as well as *NOC3L* and asthma (**Figure 6**). Genetic variation in *PIGB* causes
545 defects in glycosylphosphatidylinositol (GPI) biosynthesis⁹¹, which has been linked to
546 cardiomyopathy from clinical outcome data⁹². The mammalian homolog of *NOC3L*, called
547 *FAD24*, regulates the development of adipocytes⁹³, which release adiponectin, a hormone that
548 controls inflammation and is linked to asthma⁹⁴.

549

550 Overall, our study has demonstrated several compelling applications and utilities of MAGEPRO.
551 First, applying MAGEPRO gene expression prediction models to TWAS flags unstable
552 disease/trait-associated genes by sometimes ablating significant associations generated by less
553 accurate gene models. Second, MAGEPRO can help replicate European TWAS results in
554 understudied ancestries, confirming population-shared gene-level effects on disease which has
555 the potential to inform which European findings may be most clinically relevant to other
556 populations. Third, utilizing MAGEPRO to perform TWAS in non-European populations can
557 reveal population-specific gene-level disease effects. Fourth, MAGEPRO identifies biologically
558 plausible novel connections between disease and putative gene-level risk factors, which
559 previously could not be identified due to the lack of an available predictive *cis*-genetic gene
560 model.

561

562 **Discussion**

563

564 We developed a new method, MAGEPRO, that enhances population-specific gene expression
565 prediction models by leveraging eQTL summary statistics from diverse ancestries and cell types.
566 Briefly, MAGEPRO utilizes SuSiE to prioritize putative causal variants in external eQTL
567 datasets, which are likely more informative than tagging variants when applied to the target
568 population. We applied MAGEPRO to 8 eQTL cohorts representing 3 different ancestries,
569 improving prediction accuracy by an average of 11% relative to LASSO and consistently
570 outperforming all competing methods, including the state-of-the-art tool for genome-wide
571 complex trait PRS using multi-ancestry data, PRS-CSx. The advantages offered by MAGEPRO
572 were exemplified in small training cohorts (maximized improvement over conventional LASSO
573 models), in low *cis*-heritable genes – which are more likely to be disease-critical, and in out-of-
574 cohort prediction tasks for genetically similar populations.

575

576 When we applied MAGEPRO models to the TWAS framework, we identified 2,331 novel
577 disease/trait-associated genes, including 1,350 as a result of improving (or adjusting) existing
578 gene-trait associations and 981 that could not be identified by LASSO due to the lack of a
579 predictive *cis*-genetic gene model. MAGEPRO identified several genes associated with white
580 blood cell count that replicate across multiple ancestries, such as *PHTF1*, which is differentially
581 expressed in leukemia patients. MAGEPRO also identified biologically plausible new
582 associations, such as *PIGB* linked to heart failure, which has been evidenced by clinical outcome
583 data.

584

585 We note several limitations to our work. First, MAGEPRO relies on the availability of target
586 population genotype and gene expression data, which may be scarce for some ancestries (such as
587 South Asians, South Americans, and others) and less accessible tissues. Second, MAGEPRO
588 applies SuSiE to each external dataset independently, which may not be as powerful as modeling
589 cross-ancestry or cross-tissue effect size correlations while fine-mapping. Third, MAGEPRO
590 models are population-specific by design, which may complicate downstream analysis and limit
591 generalizability when there are slight mismatches between the population structure of the
592 training eQTL cohort and the target population (i.e., if the GWAS cohort has higher degrees of
593 admixture). Fourth, while MAGEPRO definitively improves the accuracy of *cis*-genetic models
594 of gene expression, limited availability of large ancestrally diverse GWAS continues to restrict
595 the power of gene-disease association studies like TWAS. Despite these limitations, MAGEPRO
596 is a powerful and robust method for creating population-specific *cis*-genetic models of gene
597 expression and has provided clarifying and new insights related to the underlying risk factors of
598 blood cell complex traits and immune-mediated diseases.

599 **References**

600

- 601 1. Smemo, S. *et al.* Obesity-associated variants within FTO form long-range functional
602 connections with IRX3. *Nature* **507**, 371–375 (2014).
- 603 2. Li, Y. I. *et al.* RNA splicing is a primary link between genetic variation and disease.
604 *Science* **352**, 600–604 (2016).
- 605 3. Gusev, A. *et al.* Partitioning heritability of regulatory and cell-type-specific variants
606 across 11 common diseases. *Am. J. Hum. Genet.* **95**, 535–552 (2014).
- 607 4. Finucane, H. K. *et al.* Partitioning heritability by functional annotation using genome-
608 wide association summary statistics. *Nat. Genet.* **47**, 1228–1235 (2015).
- 609 5. Roadmap Epigenomics Consortium *et al.* Integrative analysis of 111 reference human
610 epigenomes. *Nature* **518**, 317–330 (2015).
- 611 6. Meuleman, W. *et al.* Index and biological spectrum of human DNase I hypersensitive
612 sites. *Nature* **584**, 244–251 (2020).
- 613 7. Maurano, M. T. *et al.* Systematic localization of common disease-associated variation in
614 regulatory DNA. *Science* **337**, 1190–1195 (2012).
- 615 8. Albert, F. W. & Kruglyak, L. The role of regulatory variation in complex traits and
616 disease. *Nat. Rev. Genet.* **16**, 197–212 (2015).
- 617 9. Hormozdiari, F. *et al.* Colocalization of GWAS and eQTL Signals Detects Target Genes.
618 *Am. J. Hum. Genet.* **99**, 1245–1260 (2016).
- 619 10. Hormozdiari, F. *et al.* Leveraging molecular quantitative trait loci to understand the
620 genetic architecture of diseases and complex traits. *Nat. Genet.* **50**, 1041–1047 (2018).
- 621 11. Giambartolomei, C. *et al.* Bayesian test for colocalisation between pairs of genetic
622 association studies using summary statistics. *PLoS Genet.* **10**, e1004383 (2014).
- 623 12. Gusev, A. *et al.* Integrative approaches for large-scale transcriptome-wide association
624 studies. *Nat. Genet.* **48**, 245–252 (2016).
- 625 13. Wainberg, M. *et al.* Opportunities and challenges for transcriptome-wide association
626 studies. *Nat. Genet.* **51**, 592–599 (2019).
- 627 14. Gamazon, E. R. *et al.* A gene-based association method for mapping traits using
628 reference transcriptome data. *Nat. Genet.* **47**, 1091–1098 (2015).
- 629 15. GTEx Consortium. The GTEx Consortium atlas of genetic regulatory effects across
630 human tissues. *Science* **369**, 1318–1330 (2020).
- 631 16. Võsa, U. *et al.* Large-scale cis- and trans-eQTL analyses identify thousands of genetic
632 loci and polygenic scores that regulate blood gene expression. *Nat. Genet.* **53**, 1300–1310
633 (2021).
- 634 17. Martin, A. R. *et al.* Clinical use of current polygenic risk scores may exacerbate health
635 disparities. *Nat. Genet.* **51**, 584–591 (2019).
- 636 18. 1000 Genomes Project Consortium *et al.* A global reference for human genetic variation.
637 *Nature* **526**, 68–74 (2015).

- 638 19. Cavazos, T. B. & Witte, J. S. Inclusion of variants discovered from diverse populations
639 improves polygenic risk score transferability. *HGG Adv.* **2**, (2021).
- 640 20. Keys, K. L. *et al.* On the cross-population generalizability of gene expression prediction
641 models. *PLoS Genet.* **16**, e1008927 (2020).
- 642 21. Weissbrod, O. *et al.* Leveraging fine-mapping and multipopulation training data to
643 improve cross-population polygenic risk scores. *Nat. Genet.* **54**, 450–458 (2022).
- 644 22. Lu, Z. *et al.* Multi-ancestry fine-mapping improves precision to identify causal genes in
645 transcriptome-wide association studies. *Am. J. Hum. Genet.* **109**, 1388–1404 (2022).
- 646 23. Kachuri, L. *et al.* Gene expression in African Americans, Puerto Ricans and Mexican
647 Americans reveals ancestry-specific patterns of genetic architecture. *Nat. Genet.* **55**, 952–
648 963 (2023).
- 649 24. Mogil, L. S. *et al.* Genetic architecture of gene expression traits across diverse
650 populations. *PLoS Genet.* **14**, e1007586 (2018).
- 651 25. Shang, L. *et al.* Genetic Architecture of Gene Expression in European and African
652 Americans: An eQTL Mapping Study in GENOA. *Am. J. Hum. Genet.* **106**, 496–512
653 (2020).
- 654 26. Ota, M. *et al.* Dynamic landscape of immune cell-specific gene regulation in immune-
655 mediated diseases. *Cell* **184**, 3006–3021.e17 (2021).
- 656 27. Lappalainen, T. *et al.* Transcriptome and genome sequencing uncovers functional
657 variation in humans. *Nature* **501**, 506–511 (2013).
- 658 28. Okada, Y. *et al.* Genetics of rheumatoid arthritis contributes to biology and drug
659 discovery. *Nature* **506**, 376–381 (2014).
- 660 29. Márquez-Luna, C., Loh, P.-R., South Asian Type 2 Diabetes (SAT2D) Consortium,
661 SIGMA Type 2 Diabetes Consortium & Price, A. L. Multiethnic polygenic risk scores
662 improve risk prediction in diverse populations. *Genet. Epidemiol.* **41**, 811–823 (2017).
- 663 30. Ruan, Y. *et al.* Improving polygenic prediction in ancestrally diverse populations. *Nat.*
664 *Genet.* **54**, 573–580 (2022).
- 665 31. Hoggart, C. J. *et al.* BridgePRS leverages shared genetic effects across ancestries to
666 increase polygenic risk score portability. *Nat. Genet.* **56**, 180–186 (2024).
- 667 32. Yuan, K. *et al.* Fine-mapping across diverse ancestries drives the discovery of putative
668 causal variants underlying human complex traits and diseases. *Nat. Genet.* (2024)
669 doi:10.1038/s41588-024-01870-z.
- 670 33. Zaitlen, N., Paşaniuc, B., Gur, T., Ziv, E. & Halperin, E. Leveraging genetic variability
671 across populations for the identification of causal variants. *Am. J. Hum. Genet.* **86**, 23–33
672 (2010).
- 673 34. Chen, F. *et al.* Multi-ancestry transcriptome-wide association analyses yield insights into
674 tobacco use biology and drug repurposing. *Nat. Genet.* **55**, 291–300 (2023).
- 675 35. Li, Z. *et al.* METRO: Multi-ancestry transcriptome-wide association studies for powerful
676 gene-trait association detection. *Am. J. Hum. Genet.* **109**, 783–801 (2022).

- 677 36. Wang, G., Sarkar, A., Carbonetto, P. & Stephens, M. A simple new approach to variable
678 selection in regression, with application to genetic fine mapping. *J. R. Stat. Soc. Series B*
679 *Stat. Methodol.* **82**, 1273–1300 (2020).
- 680 37. Zou, Y., Carbonetto, P., Wang, G. & Stephens, M. Fine-mapping from summary data
681 with the “Sum of Single Effects” model. *PLoS Genet.* **18**, e1010299 (2022).
- 682 38. Siewert-Rocks, K. M., Kim, S. S., Yao, D. W., Shi, H. & Price, A. L. Leveraging gene
683 co-regulation to identify gene sets enriched for disease heritability. *Am. J. Hum. Genet.*
684 **109**, 393–404 (2022).
- 685 39. Lu, Z. *et al.* Improved multi-ancestry fine-mapping identifies cis-regulatory variants
686 underlying molecular traits and disease risk. *medRxiv* (2024)
687 doi:10.1101/2024.04.15.24305836.
- 688 40. Amariuta, T. *et al.* Improving the trans-ancestry portability of polygenic risk scores by
689 prioritizing variants in predicted cell-type-specific regulatory elements. *Nat. Genet.* **52**,
690 1346–1354 (2020).
- 691 41. International Schizophrenia Consortium *et al.* Common polygenic variation contributes to
692 risk of schizophrenia and bipolar disorder. *Nature* **460**, 748–752 (2009).
- 693 42. Rossen, J. *et al.* MultiSuSiE improves multi-ancestry fine-mapping in All of Us whole-
694 genome sequencing data. *medRxiv* (2024) doi:10.1101/2024.05.13.24307291.
- 695 43. Tsuo, K. *et al.* Multi-ancestry meta-analysis of asthma identifies novel associations and
696 highlights the value of increased power and diversity. *Cell Genomics* **2**, 100212 (2022).
- 697 44. Jee, Y. H. *et al.* Multi-ancestry polygenic risk scores for venous thromboembolism. *Hum.*
698 *Mol. Genet.* **33**, 1584–1591 (2024).
- 699 45. Tsuo, K. *et al.* All of Us diversity and scale improve polygenic prediction contextually
700 with greatest improvements for under-represented populations. *BioRxiv* (2024)
701 doi:10.1101/2024.08.06.606846.
- 702 46. Pham, D. *et al.* Assessing polygenic risk score models for applications in populations
703 with under-represented genomics data: an example of Vietnam. *Brief. Bioinformatics* **23**,
704 (2022).
- 705 47. Ge, T. *et al.* Development and validation of a trans-ancestry polygenic risk score for type
706 2 diabetes in diverse populations. *Genome Med.* **14**, 70 (2022).
- 707 48. Zhao, Z., Fritsche, L. G., Smith, J. A., Mukherjee, B. & Lee, S. The construction of
708 cross-population polygenic risk scores using transfer learning. *Am. J. Hum. Genet.* **109**,
709 1998–2008 (2022).
- 710 49. Sun, Q. *et al.* Improving polygenic risk prediction in admixed populations by explicitly
711 modeling ancestral-differential effects via GAUDI. *Nat. Commun.* **15**, 1016 (2024).
- 712 50. Yang, J., Lee, S. H., Goddard, M. E. & Visscher, P. M. GCTA: a tool for genome-wide
713 complex trait analysis. *Am. J. Hum. Genet.* **88**, 76–82 (2011).
- 714 51. Wang, X., Lu, Z., Bhattacharya, A., Pasaniuc, B. & Mancuso, N. *twas_sim*, a Python-
715 based tool for simulation and power analysis of transcriptome-wide association analysis.
716 *Bioinformatics* **39**, (2023).

- 717 52. Mancuso, N. *et al.* Probabilistic fine-mapping of transcriptome-wide association studies.
718 *Nat. Genet.* **51**, 675–682 (2019).
- 719 53. O'Connor, L. J. *et al.* Extreme polygenicity of complex traits is explained by negative
720 selection. *Am. J. Hum. Genet.* **105**, 456–476 (2019).
- 721 54. Mostafavi, H., Spence, J. P., Naqvi, S. & Pritchard, J. K. Systematic differences in
722 discovery of genetic effects on gene expression and complex traits. *Nat. Genet.* **55**, 1866–
723 1875 (2023).
- 724 55. Wang, J., Zhang, Z., Lu, Z., Mancuso, N. & Gazal, S. Genes with differential expression
725 across ancestries are enriched in ancestry-specific disease effects likely due to gene-by-
726 environment interactions. *Am. J. Hum. Genet.* (2024) doi:10.1016/j.ajhg.2024.07.021.
- 727 56. Amariuta, T. *et al.* IMPACT: Genomic Annotation of Cell-State-Specific Regulatory
728 Elements Inferred from the Epigenome of Bound Transcription Factors. *Am. J. Hum.*
729 *Genet.* **104**, 879–895 (2019).
- 730 57. Zhang, X., Jiang, W. & Zhao, H. Integration of expression QTLs with fine mapping via
731 SuSiE. *PLoS Genet.* **20**, e1010929 (2024).
- 732 58. Lek, M. *et al.* Analysis of protein-coding genetic variation in 60,706 humans. *Nature*
733 **536**, 285–291 (2016).
- 734 59. Boerner, T. J., Deems, S., Furlani, T. R., Knuth, S. L. & Towns, J. ACCESS: advancing
735 innovation: nsf's advanced cyberinfrastructure coordination ecosystem: services &
736 support. in *Practice and Experience in Advanced Research Computing* (eds. Sinkovits,
737 R., Romanella, A., Knuth, S., Hackworth, K. & Pummill, J.) 173–176 (ACM, 2023).
738 doi:10.1145/3569951.3597559.
- 739 60. Chen, M.-H. *et al.* Trans-ethnic and Ancestry-Specific Blood-Cell Genetics in 746,667
740 Individuals from 5 Global Populations. *Cell* **182**, 1198-1213.e14 (2020).
- 741 61. Zhou, W. *et al.* Global Biobank Meta-analysis Initiative: Powering genetic discovery
742 across human disease. *Cell Genomics* **2**, 100192 (2022).
- 743 62. Yang, H. *et al.* ZNF213 facilitates ER alpha signaling in breast cancer cells. *Front.*
744 *Oncol.* **11**, 638751 (2021).
- 745 63. Liu, Y. *et al.* ZNF213 negatively controls triple negative breast cancer progression via
746 Hippo/YAP signaling. *Cancer Sci.* **112**, 2714–2727 (2021).
- 747 64. Billington, C. K. & Penn, R. B. Signaling and regulation of G protein-coupled receptors
748 in airway smooth muscle. *Respir. Res.* **4**, 2 (2003).
- 749 65. Fuentes, N., McCullough, M., Panettieri, R. A. & Druey, K. M. RGS proteins, GRKs, and
750 beta-arrestins modulate G protein-mediated signaling pathways in asthma. *Pharmacol.*
751 *Ther.* **223**, 107818 (2021).
- 752 66. Johnson, E. N. & Druey, K. M. Heterotrimeric G protein signaling: role in asthma and
753 allergic inflammation. *J. Allergy Clin. Immunol.* **109**, 592–602 (2002).
- 754 67. Douroudis, K., Kisand, K., Nemvalts, V., Rajasalu, T. & Uibo, R. Allelic variants in the
755 PHTF1-PTPN22, C12orf30 and CD226 regions as candidate susceptibility factors for the
756 type 1 diabetes in the Estonian population. *BMC Med. Genet.* **11**, 11 (2010).

- 757 68. Huang, X. *et al.* Analysis of the expression of PHTF1 and related genes in acute
758 lymphoblastic leukemia. *Cancer Cell Int.* **15**, 93 (2015).
- 759 69. Scheffler, J. M. *et al.* LAMTOR2 regulates dendritic cell homeostasis through FLT3-
760 dependent mTOR signalling. *Nat. Commun.* **5**, 5138 (2014).
- 761 70. Sparber, F. *et al.* The late endosomal adaptor molecule p14 (LAMTOR2) represents a
762 novel regulator of Langerhans cell homeostasis. *Blood* **123**, 217–227 (2014).
- 763 71. Bohn, G. *et al.* A novel human primary immunodeficiency syndrome caused by
764 deficiency of the endosomal adaptor protein p14. *Nat. Med.* **13**, 38–45 (2007).
- 765 72. Bottini, N. & Peterson, E. J. Tyrosine phosphatase PTPN22: multifunctional regulator of
766 immune signaling, development, and disease. *Annu. Rev. Immunol.* **32**, 83–119 (2014).
- 767 73. Stanford, S. M. & Bottini, N. PTPN22: the archetypal non-HLA autoimmunity gene. *Nat.*
768 *Rev. Rheumatol.* **10**, 602–611 (2014).
- 769 74. Plenge, R. M. *et al.* Replication of putative candidate-gene associations with rheumatoid
770 arthritis in >4,000 samples from North America and Sweden: association of susceptibility
771 with PTPN22, CTLA4, and PADI4. *Am. J. Hum. Genet.* **77**, 1044–1060 (2005).
- 772 75. Wellcome Trust Case Control Consortium. Genome-wide association study of 14,000
773 cases of seven common diseases and 3,000 shared controls. *Nature* **447**, 661–678 (2007).
- 774 76. Gregersen, P. K. *et al.* REL, encoding a member of the NF-kappaB family of
775 transcription factors, is a newly defined risk locus for rheumatoid arthritis. *Nat. Genet.*
776 **41**, 820–823 (2009).
- 777 77. Saez, A. *et al.* Lamin A/C and the immune system: one intermediate filament, many
778 faces. *Int. J. Mol. Sci.* **21**, (2020).
- 779 78. Bordeleau, M.-E. *et al.* UBAP2L is a novel BMI1-interacting protein essential for
780 hematopoietic stem cell activity. *Blood* **124**, 2362–2369 (2014).
- 781 79. Yáñez, A., Ng, M. Y., Hassanzadeh-Kiabi, N. & Goodridge, H. S. IRF8 acts in lineage-
782 committed rather than oligopotent progenitors to control neutrophil vs monocyte
783 production. *Blood* **125**, 1452–1459 (2015).
- 784 80. Hoffman, J. D. *et al.* Cis-eQTL-based trans-ethnic meta-analysis reveals novel genes
785 associated with breast cancer risk. *PLoS Genet.* **13**, e1006690 (2017).
- 786 81. Takeuchi, H. *et al.* Elevated red cell distribution width to platelet count ratio predicts
787 poor prognosis in patients with breast cancer. *Sci. Rep.* **9**, 3033 (2019).
- 788 82. Ghousaini, M. *et al.* Open Targets Genetics: systematic identification of trait-associated
789 genes using large-scale genetics and functional genomics. *Nucleic Acids Res.* **49**, D1311–
790 D1320 (2021).
- 791 83. Mountjoy, E. *et al.* An open approach to systematically prioritize causal variants and
792 genes at all published human GWAS trait-associated loci. *Nat. Genet.* **53**, 1527–1533
793 (2021).
- 794 84. Tapia, A. L. *et al.* A large-scale transcriptome-wide association study (TWAS) of 10
795 blood cell phenotypes reveals complexities of TWAS fine-mapping. *Genet. Epidemiol.*
796 **46**, 3–16 (2022).

- 797 85. Lu, M. *et al.* TWAS Atlas: a curated knowledgebase of transcriptome-wide association
798 studies. *Nucleic Acids Res.* **51**, D1179–D1187 (2023).
- 799 86. Mancuso, N. *et al.* Integrating Gene Expression with Summary Association Statistics to
800 Identify Genes Associated with 30 Complex Traits. *Am. J. Hum. Genet.* **100**, 473–487
801 (2017).
- 802 87. Morra, M. *et al.* Structural basis for the interaction of the free SH2 domain EAT-2 with
803 SLAM receptors in hematopoietic cells. *EMBO J.* **20**, 5840–5852 (2001).
- 804 88. Nagy, N. *et al.* SH2D1A and slam protein expression in human lymphocytes and derived
805 cell lines. *International Journal of Cancer* (2000).
- 806 89. Calpe, S. *et al.* The SLAM and SAP gene families control innate and adaptive immune
807 responses. *Adv. Immunol.* **97**, 177–250 (2008).
- 808 90. Dragovich, M. A. & Mor, A. The SLAM family receptors: Potential therapeutic targets
809 for inflammatory and autoimmune diseases. *Autoimmun. Rev.* **17**, 674–682 (2018).
- 810 91. Murakami, Y. *et al.* Mutations in PIGB Cause an Inherited GPI Biosynthesis Defect with
811 an Axonal Neuropathy and Metabolic Abnormality in Severe Cases. *Am. J. Hum. Genet.*
812 **105**, 384–394 (2019).
- 813 92. Bayat, A. *et al.* GPI-anchoring disorders and the heart: Is cardiomyopathy an overlooked
814 feature? *Clin. Genet.* **104**, 598–603 (2023).
- 815 93. Johmura, Y., Osada, S., Nishizuka, M. & Imagawa, M. FAD24, a regulator of
816 adipogenesis, is required for the regulation of DNA replication in cell proliferation. *Biol.*
817 *Pharm. Bull.* **31**, 1092–1095 (2008).
- 818 94. Otelea, M. R., Arghir, O. C., Zugravu, C. & Rascu, A. Adiponectin and asthma: knowns,
819 unknowns and controversies. *Int. J. Mol. Sci.* **22**, (2021).

820 **Methods**

821

822 ***Baseline genetic model of gene expression: LASSO***

823

824 We used the FUSION tool to build the standard gene expression prediction model, which uses
825 individual-level genotype and gene expression data from a single target population (**Figure 1**,
826 green box). In this baseline model, a single gene’s expression is modeled with standardized
827 genotypes of *cis*-variants (within 1 Mb of the gene’s transcription start site (TSS)) in a
828 multivariate linear regression:

$$829 \quad y_i = \sum_j X_{ij}\beta_j + \epsilon_i$$

830 where for each individual i , y_i is the gene expression of one gene, j indexes *cis*-variants, X_{ij} is
831 the standardized genotype of individual i at SNP j , β_j is the true unobserved eQTL effect size,
832 and ϵ_i is the residual of gene expression not explained by modeled *cis*-genetic effects. We used
833 LASSO (L1 norm) regularized linear regression from PLINK⁹⁵ to estimate $\hat{\beta}_j$ for each *cis*-variant
834 such that we minimize the penalized sum of squares:

$$835 \quad \min_{\hat{\beta}_j} \left(\sum_i (y_i - \sum_j X_{ij}\hat{\beta}_j)^2 + \lambda \sum_j |\hat{\beta}_j| \right)$$

836 where λ is the sparsity parameter which is tuned via cross-validation. L1 regularization avoids
837 overfitting by shrinking coefficients of less informative features (e.g., SNPs) to 0 and assigns
838 nonzero coefficients to potentially predictive SNPs. When LASSO regression fails to find any
839 meaningful predictors and pushes all coefficients to zero (potentially due to the limited sample
840 size of the target population), we employ the “top 1” model as is done in the FUSION
841 framework. The “top 1” model uses a single predictor SNP, specifically the SNP with the largest
842 squared effect size from marginal *cis*-eQTL analysis. This approach systematically enables us to
843 build a standard gene model for every gene in the analysis, to which we can compare
844 MAGEPRO models informed by multiple ancestries.

845

846 ***MAGEPRO (Multi-Ancestry Gene Expression Prediction Regularized Optimization)***

847

848 MAGEPRO takes a three-step approach. First, it learns noisy estimates of SNP-gene effect sizes
849 in the target population with a LASSO-regularized linear regression, identical to the baseline
850 model described above (**Figure 1**, green box). Second, we apply the Sum of Single Effects
851 (SuSiE) linear regression to each set of external eQTL summary statistics and we retain the
852 posterior effect size estimates (**Figure 1**, blue box). SuSiE serves as a variable selection step,
853 prioritizing potentially causal eQTLs which are more likely to be informative to the target
854 population (see “Sum of Single Effects to prioritize variants from external summary statistics”
855 section for more details regarding SuSiE). Finally, MAGEPRO models the gene expression of

856 the target population as a function of the baseline LASSO-regularized model and the SuSiE
857 posterior eQTL effect size estimates for each external dataset (**Figure 1**, white box):

858
$$y_i \sim \sum_{D \in t,d} (\alpha_D \sum_j X_{ij} \hat{\beta}_{jD})$$

859 where for each individual i , y_i is the gene expression of one gene, D indexes target (t) and
860 external datasets (d), j indexes *cis*-variants, X_{ij} is the standardized genotype of individual i at
861 SNP j , $\hat{\beta}_d$ is a vector of posterior eQTL effect size estimates from external dataset d , and $\hat{\beta}_t$ is a
862 vector of estimated effect sizes from applying the baseline model described above to the target
863 dataset. We used ridge (L2 norm) regression to fit $\hat{\alpha}_t$ and $\hat{\alpha}_d$; the dataset-specific mixing weights
864 represent the relative contribution of each dataset to the prediction of gene expression, such that
865 we minimize the loss function:

866
$$\min_{\hat{\alpha}_D} \sum_i (y_i - (\sum_{D \in t,d} (\hat{\alpha}_D \sum_j X_{ij} \hat{\beta}_{jD})))^2 + \lambda \sum_{D \in t,d} \hat{\alpha}_D^2$$

867 where λ is the sparsity parameter, which is tuned by ten-fold cross-validation⁹⁶. We applied ridge
868 regression to constrain the coefficients when two or more vectors are collinear, which may be
869 common given that causal eQTL architecture is at least partially shared across populations.

870

871 **Simulations**

872

873 We conducted simulations with various sample sizes and gene expression *cis*-heritability values
874 to assess the robustness of MAGEPRO. We applied MAGEPRO, PRS-CSx, and LASSO to four
875 predetermined levels of heritability (0.05, 0.1, 0.2, 0.4), which we confirmed using GCTA
876 (**Supplementary Figure 22**). These heritability values were chosen based on the average
877 estimated heritability values in quartiles of significantly heritable genes in LCL gene expression
878 data from the GENOA African American (AA) population (0.088, 0.139, 0.202, 0.382). For each
879 heritability value, we simulated 1,000 random genes and investigated the performance of each
880 model across five target population (African) sample sizes (80, 160, 240, 400, 500). Simulated
881 genotypes and gene expression levels for 500 EUR individuals (based on LD from the 1000
882 Genomes European ancestry group) and 500 AMR individuals (based on LD from the 1000
883 Genomes American ancestry group) were used to compute summary statistics, which we used as
884 external datasets to apply MAGEPRO and PRS-CSx. Many of the functions that we used for our
885 simulations are adopted from the Mancuso Lab TWAS simulator.

886

887 We assessed the performance of MAGEPRO in various simulated genetic architectures of gene
888 expression: (1) the causal *cis*-eQTLs are the same across populations (same genomic position but
889 not necessarily correlated in effect size), (2) the causal *cis*-eQTLs are different variants across
890 populations but in high LD ($r^2 > 0.8$), (3) true effect sizes of all shared causal *cis*-eQTLs are
891 drawn independently across populations, and (4) true effect sizes of all shared causal *cis*-eQTLs

892 are correlated across populations with effect size correlation set to 0.8, following recent work
893 which estimated cis-molQTL (molecular quantitative trait loci) effect size correlations across
894 ancestries³⁹. The performances of LASSO, PRS-CSx, and MAGEPRO in simulations are
895 evaluated with the prediction accuracy defined as $\frac{R_{CV}^2}{\hat{h}_{ge}^2}$. Please see the Supplementary Note
896 section called “Simulation framework” for more details.

897

898 ***Competing methods of gene expression prediction***

899

900 We compare the performance of MAGEPRO against six different methods, capturing
901 conventional methods applied to genome-wide complex trait data and gene expression data:
902 meta-analysis, P+T, LASSO, SuSiE, Multipop, and PRS-CSx (see “*Baseline genetic model of*
903 *gene expression: LASSO*” for more information on the LASSO model). We note that we do not
904 compare the performance of elastic net or BLUP as recent work has shown that neither
905 significantly outperform LASSO³⁹.

906

907 The meta-analysis model refers to a sample-size weighted meta-analysis of all datasets, including
908 the LASSO gene model which was developed using the training split of the target cohort. This
909 strategy is commonly applied to GWAS data to maximize association power and identify shared
910 effects.

911

912 P+T (pruning and thresholding) is an LD-informed pruning and p-value thresholding method⁹⁷,
913 also referred to as clumping and thresholding. Briefly, we iterate through SNPs in order of
914 increasing p-value below a chosen threshold; p-values are computed from a marginal *cis*-eQTL
915 analysis with the target cohort data. All variants in LD with the current SNP are removed until
916 the iteration finishes. We performed a small grid-search across several LD r^2 thresholds (0.2, 0.5,
917 0.8) and p-value thresholds (0.001, 0.01, 0.1, 0.5) to identify the pair of parameters that result in
918 the best prediction result in 5-fold cross-validation. We performed P+T using PLINK and we
919 used the target population genotypes as the in-sample LD reference panel.

920

921 SuSiE is the Sum of Single Effects regression model applied to the individual-level target
922 population genotype and gene expression data. We used default parameters to run SuSiE
923 (including a maximum number of allowed credible sets: $L = 10$, up to 100 iterative Bayesian
924 stepwise selection (IBSS) iterations, and setting the estimated residual variance flag to TRUE if
925 in-sample LD files were available and FALSE otherwise) and retained the resulting posterior
926 effect size estimates to predict gene expression.

927

928 Multipop refers to a variation of MAGEPRO without the variable selection step using SuSiE. In
929 this model, the raw external marginal *cis*-eQTL summary statistics are combined with the target
930 population LASSO model using ridge regression. Benchmarking against this method allows us to

931 evaluate if using SuSiE to prioritize potentially causal variants helps us create more accurate
932 predictive models.

933

934 PRS-CSx is a Bayesian framework that improves cross-population polygenic prediction by
935 learning an optimal linear combination of GWAS summary statistics from multiple ancestry
936 groups to produce the final PRS. PRS-CSx employs a shared continuous shrinkage prior to SNP
937 effects across populations (which assumes shared effects across populations) and leverages LD
938 diversity across samples to enhance accuracy in effect size estimates. Although this method was
939 originally designed to improve PRS for genome-wide complex traits and polygenic diseases in
940 ancestrally diverse populations, we applied their command line tool to gene expression
941 prediction to benchmark MAGEPRO. We utilized the shared shrinkage prior from PRS-CSx on
942 the same datasets employed in MAGEPRO. Then, we learned an optimal linear combination of
943 the post-shrinkage external datasets. To ensure that PRS-CSx utilizes the same features as
944 MAGEPRO, we also added the LASSO gene model for the target population as one of the
945 features in the linear combination. The authors of PRS-CSx recommend that the global shrinkage
946 parameter, Φ , is adjusted based on the polygenicity of the phenotype. Since we expected the *cis*-
947 genetic component of gene expression to be much less polygenic (involve fewer causal variants)
948 than a genome-wide trait, we considered values of $[10^{-5}, 10^{-6}, 10^{-7}, 10^{-8}, 10^{-9}]$. We applied
949 PRS-CSx with these shrinkage parameters for 200 random genes with $\hat{h}_{ge}^2 > 0$ and $\hat{h}_{ge}^2 p <$
950 0.05 . We observed that gene model accuracy was robust across all values of Φ , and thus we
951 selected the intermediate value (10^{-7}) for the remaining analyses, which assumes that the
952 polygenicity of *cis*-genetic gene expression regulation was well-represented by these 200
953 randomly selected genes (**Supplementary Figure 23**).

954

955 We note that BridgePRS³¹ is a recently published multi-ancestry PRS method that we considered
956 for our study. However, their study demonstrated that BridgePRS only nominally outperforms
957 PRS-CSx under highly polygenic genetic architectures, such as genome-wide complex traits.
958 Therefore, we benchmarked MAGEPRO against PRS-CSx because we believed it was the best
959 candidate among multi-ancestry PRS frameworks that are applicable to gene expression
960 prediction.

961

962 ***Preparing external summary statistics for MAGEPRO***

963

964 We downloaded eQTL summary statistics from 5 publicly available datasets from 3 different
965 ancestries including European, Latino/Hispanic and African American cohorts. For each dataset,
966 we extracted full *cis*-eQTL summary statistics and filtered for 1,034,897 HapMap 3 SNPs
967 included in GTEx. If the effect allele and alternate allele of the eQTLs were flipped in
968 comparison to the target cohort SNPs, we multiplied the effect size of the eQTL from the
969 external dataset by -1. We split each dataset into gene-specific files to facilitate downstream
970 analysis with MAGEPRO. Dataset-specific preprocessing details are described in the

971 Supplementary Note. To avoid overfitting, we utilized different combinations of external
972 summary statistics depending on the target population to build the predictive model (**Table 1**).
973

974 *Sum of Single Effects model to prioritize variants from external summary statistics*

975
976 We utilized the Sum of Single Effects regression model (SuSiE), specifically “SuSiE-RSS”
977 (Regression with Summary Statistics), for variable selection from eQTL summary statistics data.
978 SuSiE is a variable selection method that quantifies the uncertainty in which variables are
979 selected by expressing the regression coefficients as a sum of single effects where only one of
980 the variables has a nonzero coefficient. The model is fit with the IBSS procedure and produces
981 posterior inclusion probabilities (PIPs) and posterior effect sizes for each SNP. The original
982 SuSiE method requires individual-level phenotype and genotype data. In our MAGEPRO
983 pipeline, external datasets only contain summary data, hence, we use SuSiE-RSS, which employs
984 the “IBSS-ss” algorithm that relies only on sufficient statistics that can be approximated from the
985 summary statistics. Within our pipeline, we conduct fine-mapping separately for each gene in
986 each eQTL dataset. When available, we utilize in-sample correlation matrices (e.g., for MESA or
987 GENOA datasets). In cases where in-sample matrices are not available, we employ out-of-cohort
988 ancestry-matched alternatives (e.g., we used LD from the 1000 Genomes European population to
989 fine-map the European eQTLGen dataset).

990
991 We note that the incorporation of the recently developed multi-ancestry statistical fine-mapping
992 method, Sum of Shared Single Effects (SuShiE), may enhance the MAGEPRO framework by
993 leveraging LD heterogeneity and modeling cross-ancestry effect size correlations to improve
994 variable selection and effect size estimates in external eQTL datasets³⁹. However, a version of
995 SuShiE that is compatible with summary statistics was not released at the time of this study.
996 Additionally, fine-mapping methods that are most compatible with MAGEPRO may also benefit
997 from modeling cross-cell-type correlations to enable the sharing of information across eQTL
998 datasets from different ancestries and cell types.

999

1000 *Processing individual-level genotype and gene expression data*

1001

1002 We used the same variant and relatedness filtering for all genotyping data, regardless of cohort.
1003 All genotype data processing was done using PLINK v1.9 and bcftools⁹⁸. For the GENOA and
1004 MESA cohort, we imputed genotype data on the TOPMed server. Each ancestry/dataset assayed
1005 on different genotype platforms were imputed separately. The imputation was run using
1006 Minimac4 (1.8.0-beta4), using the TOPMed r3 reference panel and Eagle v2.4 phasing. We kept
1007 biallelic SNPs with high imputation quality ($r^2 > 0.9$) for each imputed dataset and removed
1008 SNPs with MAF < 1%, Hardy Weinberg Equilibrium (HWE) $p < 1 \times 10^{-6}$, and genotyping rate
1009 < 1. We used plink (--rel-cutoff) to remove one individual of a pair that exhibited a relatedness
1010 greater than 0.05. When fitting the gene expression prediction models, we subset to HapMap 3

1011 SNPs present in the dataset. Compared to keeping all SNPs in the genotype data, utilizing only
1012 HapMap 3 SNPs produces heritability estimates with smaller standard errors (**Supplementary**
1013 **Figure 24**).

1014

1015 The gene expression data for each cohort was inverse-normal transformed across individuals
1016 before fitting the gene expression prediction models¹⁵. We defined the *cis*-window of each gene
1017 as [start – 500 kilobases (Kb), end + 500 Kb]. The start and end positions were defined by
1018 gencode v26 gene annotations.

1019

1020 ***Fitting gene expression prediction models***

1021

1022 To calculate gene expression weights from real data, we used genotypes and gene expression
1023 data from whole blood and lung tissues of the GTEx cohort (EUR and AA populations), LCL
1024 gene expression data from GEUVADIS (EUR) and GENOA (AA), and monocyte gene
1025 expression data from MESA (EUR, AA, HIS) (**Table 1**). After extracting samples with both
1026 genotype and gene expression data, we performed imputation, variant-based filtering, and
1027 individual-level filtering steps described above. We regressed out the appropriate covariates from
1028 the gene expression data before fitting the gene expression prediction models. These covariates
1029 generally included 5 genotype PCs, genotype platform / site of data collection, sex, age, and gene
1030 expression PCs (depending on the sample size of the cohort). Please see the Supplementary Note
1031 for dataset-specific information.

1032

1033 The performance of gene expression prediction models in this paper are evaluated with R_{cv}^2 from
1034 a 5-fold cross validation. In each iteration of the cross-validation, we use the training split (4
1035 folds) to learn a noisy estimate of *cis*-variant weights in a model identical to the standard gene
1036 expression prediction models described above. We include these weights from the training fold
1037 in a regularized linear combination with the other external datasets (consisting of SuSiE posterior
1038 effect sizes), and use the training split again to estimate the mixing weights (\hat{a}_D). Finally, we
1039 extract the estimated coefficients and predict gene expression on the remaining testing split (5th
1040 fold).

1041

1042 MAGEPRO computes both the target population SNP-gene weights ($\hat{\beta}_{target}$) and the dataset
1043 mixture weights (\hat{a}_D) using the same training split. Therefore, we tested two potential training
1044 approaches: (1) the MAGEPRO training approach described above and (2) a training approach
1045 adopted from Márquez-Luna and colleagues²⁹. In this second approach, we iteratively split the
1046 training samples (4 folds in 5-fold cross validation) into a 90% set used to estimate $\hat{\beta}_{target}$ and
1047 computed the predicted gene expression for the 10% set (for each of the 10 folds). We then
1048 performed ridge regression across all training samples to estimate \hat{a}_D and finally re-estimated
1049 $\hat{\beta}_{target}$ with the entire training split. We evaluated the two training approaches in predicting LCL
1050 gene expression at two different sample sizes ($n = 100$ and $n = 346$). We found that our

1051 MAGEPRO training approach outperformed the nested cross-validation approach in cross-
1052 validation prediction ($p = 2 \times 10^{-90}$ and $p < 1 \times 10^{-200}$ at $n = 100$ and $n = 346$,
1053 respectively, **Supplementary Figure 25**). While this could result from overfitting by
1054 MAGEPRO, we further compared the two approaches via an out-of-cohort prediction task in the
1055 GEUVADIS Yoruba (YRI) cohort. The gene models trained using the MAGEPRO approach
1056 exhibited higher accuracy ($p = 0.01$ and $p = 4.9 \times 10^{-15}$ at $n = 100$ and $n = 346$,
1057 respectively, **Supplementary Figure 25**). Therefore, we concluded that our training approach
1058 that utilizes the same training split to estimate both $\hat{\beta}_{target}$ and \hat{a}_D is valid.

1059

1060 *Validation of MAGEPRO models out-of-cohort*

1061

1062 We validate the improved MAGEPRO models by training our models in one cohort and applying
1063 them to a different cohort of a similar ancestry and cell type. To facilitate the application of gene
1064 expression prediction models across datasets, we subset to SNPs in common between the two
1065 datasets within each ancestry. Without this additional SNP-based filtering step, we risk creating
1066 predictive models that assign a non-zero effect size to SNPs that are not present in the out-of-
1067 cohort validation set.

1068

1069 To validate the LCL gene models in the European population, we built predictive models in the
1070 GEUVADIS population and validated them in the GENOA population. We worked with 718,414
1071 HapMap 3 SNPs that are present among GEUVADIS European individuals and GENOA
1072 European American individuals.

1073

1074 For individuals of African American descent, we built predictive models in the GENOA
1075 population and validated them in the GEUVADIS YRI (Yoruba) population. We worked with
1076 718,838 HapMap 3 SNPs that are present among GENOA African American individuals and
1077 GEUVADIS YRI individuals.

1078

1079 *TWAS using GWAS summary statistics*

1080

1081 We collected GWAS summary statistics for 15 blood cell traits from a previous study⁶⁰ (AFR N
1082 = 13,391, EUR N = 516,979, HIS N = 6,849) and 7 immune-mediated diseases from the Global
1083 Biobank Meta-analysis Initiative (GBMI) (AFR N = 26,052, EUR N = 1,024,298, AMR N =
1084 15,490). We updated the variant identifiers to dbSNP v151 and used the `munge_sumstats.py`
1085 script from LD score regression⁹⁹ to perform quality control and filtering. We evaluated the
1086 TWAS results for the union of significantly heritable genes across populations (LCL: 6,872
1087 genes, Monocyte: 5,920 genes, Lung: 8,807 genes) that have gene models that explain some
1088 proportion of variance in gene expression ($R^2 > 0$, $p < 0.05$). TWAS p-values were subjected
1089 to a Bonferroni significance threshold to account for multiple hypothesis testing.

1090

1091 **Statistics and Reproducibility**

1092
1093 First, as described above, we filtered each external eQTL dataset and target cohort genotypes to
1094 HapMap 3 SNPs. Second, we evaluated the performance of MAGEPRO on significantly
1095 heritable genes ($\hat{h}_{ge}^2 > 0$, $p < 0.01$) with eQTL data from at least 1 external dataset. Third, as
1096 described above, we performed random down-sampling of certain cohorts to test MAGEPRO at
1097 smaller sample sizes. Fourth, as described above, we evaluated TWAS results from gene models
1098 that explain some proportion of variance in gene expression ($R^2 > 0$, $p < 0.05$) to prevent
1099 spurious associations from estimated eQTL effect sizes that poorly capture gene expression
1100 regulation. Fifth, as described above, 1,000 random genes were simulated for each genetic
1101 architecture to robustly evaluate MAGEPRO performance. Randomization and blinding were not
1102 pertinent to our study.

1103

1104 **Data Availability**

1105

1106 Blood trait GWAS summary statistics are available at [http://www.mhi-](http://www.mhi-humangenetics.org/en/resources/)
1107 [humangenetics.org/en/resources/](http://www.mhi-humangenetics.org/en/resources/). Immune-related disease GWAS summary statistics are
1108 available at <https://www.globalbiobankmeta.org/resources>. GTEx gene expression and genotype
1109 data were acquired from dbGaP accession phs000424.v9.p2. MESA genotype data was acquired
1110 from dbGaP accession phs000209.v13.p3 (file names:
1111 phg000071.v2.NHLBI_SHARE_MESA.genotype-calls-matrixfmt.c1 and
1112 phg000071.v2.NHLBI_SHARE_MESA.genotype-calls-matrixfmt.c2), GENOA genotype data
1113 was acquired from dbGaP accession phs001238.v2.p1, and GEUVADIS genotype data is
1114 publicly available at <https://www.ebi.ac.uk/biostudies/arrayexpress/studies/E-GEUV-1>. MESA
1115 gene expression data was acquired from NCBI GEO accession GSE56045, GENOA gene
1116 expression data was acquired from NCBI GEO accessions GSE138914 (African American
1117 individuals) and GSE49531 (European individuals), and GEUVADIS gene expression data is
1118 publicly available at <https://uchicago.app.box.com/s/ewnrqs31ivobz2sn6462cq2eb423dvpr>. 1000
1119 Genomes LD reference files were acquired from https://www.bridgeprs.net/guide_input/.

1120

1121 As described in https://github.com/kaiakamatsu/MAGEPRO/tree/main/PROCESS_DATASET,
1122 all eQTL summary statistics were publicly available: eQTLGen
1123 ([https://molgenis26.gcc.rug.nl/downloads/eqtlgen/cis-eqtl/SMR_formatted/cis-eQTL-](https://molgenis26.gcc.rug.nl/downloads/eqtlgen/cis-eqtl/SMR_formatted/cis-eQTL-SMR_20191212.tar.gz)
1124 [SMR_20191212.tar.gz](https://molgenis26.gcc.rug.nl/downloads/eqtlgen/cis-eqtl/SMR_formatted/cis-eQTL-SMR_20191212.tar.gz)), GTEx ([https://console.cloud.google.com/storage/browser/gtex-](https://console.cloud.google.com/storage/browser/gtex-resources;tab=objects?prefix=&forceOnObjectsSortingFiltering=false)
1125 [resources;tab=objects?prefix=&forceOnObjectsSortingFiltering=false](https://console.cloud.google.com/storage/browser/gtex-resources;tab=objects?prefix=&forceOnObjectsSortingFiltering=false)), GENOA
1126 (http://www.xzlab.org/data/AA_summary_statistics.txt.gz), and MESA
1127 (<https://www.dropbox.com/sh/f6un5evevyvvy19/AAA3sfa1DgqY67tx4q36P341a?dl=0>).

1128

1129 **Code Availability**

1130

1131 MAGEPRO software including documentation and tutorial is publicly available at
1132 <https://github.com/kaiakamatsu/MAGEPRO> [DOI [10.5281/zenodo.13765893](https://doi.org/10.5281/zenodo.13765893)]. The Mancuso
1133 Lab TWAS Simulator is available at https://github.com/mancusolab/twas_sim. The FUSION
1134 software is available at <http://gusevlab.org/projects/fusion>. PRS-CSx is available at
1135 <https://github.com/getian107/PRScsx>. SuSiE is available as an R package and it is described at
1136 <https://stephenslab.github.io/susieR/index.html>. The munge_sumstats.py script is available in the
1137 LDSC github at <https://github.com/bulik/ldsc/tree/master>. To improve the runtime of
1138 MAGEPRO, we utilized GNU Parallel available at <https://zenodo.org/records/10901541>.

1139

1140 **Methods-only References**

1141

- 1142 95. Purcell, S. *et al.* PLINK: a tool set for whole-genome association and population-based
1143 linkage analyses. *Am. J. Hum. Genet.* **81**, 559–575 (2007).
- 1144 96. Friedman, J., Hastie, T. & Tibshirani, R. Regularization Paths for Generalized Linear
1145 Models via Coordinate Descent. *J. Stat. Softw.* **33**, 1–22 (2010).
- 1146 97. Privé, F., Vilhjálmsson, B. J., Aschard, H. & Blum, M. G. B. Making the most of
1147 clumping and thresholding for polygenic scores. *Am. J. Hum. Genet.* **105**, 1213–1221
1148 (2019).
- 1149 98. Danecek, P. *et al.* Twelve years of SAMtools and BCFtools. *Gigascience* **10**, (2021).
- 1150 99. Bulik-Sullivan, B. *et al.* LD Score regression distinguishes confounding from
1151 polygenicity in genome-wide association studies. *Nat. Genet.* **47**, 291–295 (2015).

1152

1153 **Acknowledgements**

1154

1155 This work was supported by funding from the National Science Foundation (NSF) (Award
1156 #2336469 awarded to T.A.) and the National Institutes of Health (NIH) (NHGRI R01HG013671
1157 awarded to T.A.). The funders played no role in study design, data collection and analysis,
1158 decision to publish or preparation of the manuscript. Support for GENOA was provided by the
1159 National Heart, Lung and Blood Institute (HL054457, HL054464, HL054481, HL119443, and
1160 HL087660) of the National Institutes of Health. We would like to thank the Mayo Clinic
1161 Genotyping Core, the DNA Sequencing and Gene Analysis Center at the University of
1162 Washington, and the Broad Institute for their genotyping and sequencing services. We would
1163 also like to thank the GENOA participants. This manuscript was not prepared in collaboration
1164 with investigators from the Genetic Epidemiology Network of Arteriopathy and does not
1165 necessarily reflect the opinions or views of the Genetic Epidemiology Network of Arteriopathy
1166 or NHLBI. MESA and the MESA SHARe project are conducted and supported by the National
1167 Heart, Lung, and Blood Institute (NHLBI) in collaboration with MESA investigators. Support
1168 for MESA is provided by contracts N01-HC95159, N01-HC-95160, N01-HC-95161, N01-HC-
1169 95162, N01-HC-95163, N01-HC-95164, N01-HC-95165, N01-HC95166, N01-HC-95167, N01-
1170 HC-95168, N01-HC-95169 and CTSA UL1-RR-024156. The Genotype-Tissue Expression

1171 (GTEx) Project was supported by the Common Fund of the Office of the Director of the National
1172 Institutes of Health, and by NCI, NHGRI, NHLBI, NIDA, NIMH, and NINDS. This work used
1173 the Expanse HPC server at the San Diego Supercomputer Center (SDSC) through allocation
1174 BIO230210 from the Advanced Cyberinfrastructure Coordination Ecosystem: Services &
1175 Support (ACCESS) program, which is supported by National Science Foundation grants
1176 #2138259, #2138286, #2138307, #2137603, and #2138296.

1177

1178 **Author Contributions**

1179

1180 K.A. and T.A. conceived and designed the study. K.A. conducted simulation analyses. K.A. and
1181 and S.G. conducted real data analysis. T.A. managed GTEx, GENOA, and MESA data through
1182 dbGaP. K.A., S.G., and T.A. wrote the initial draft of the manuscript and contributed to the final
1183 manuscript.

1184

1185 **Competing Interests**

1186

1187 The authors declare no competing interests.

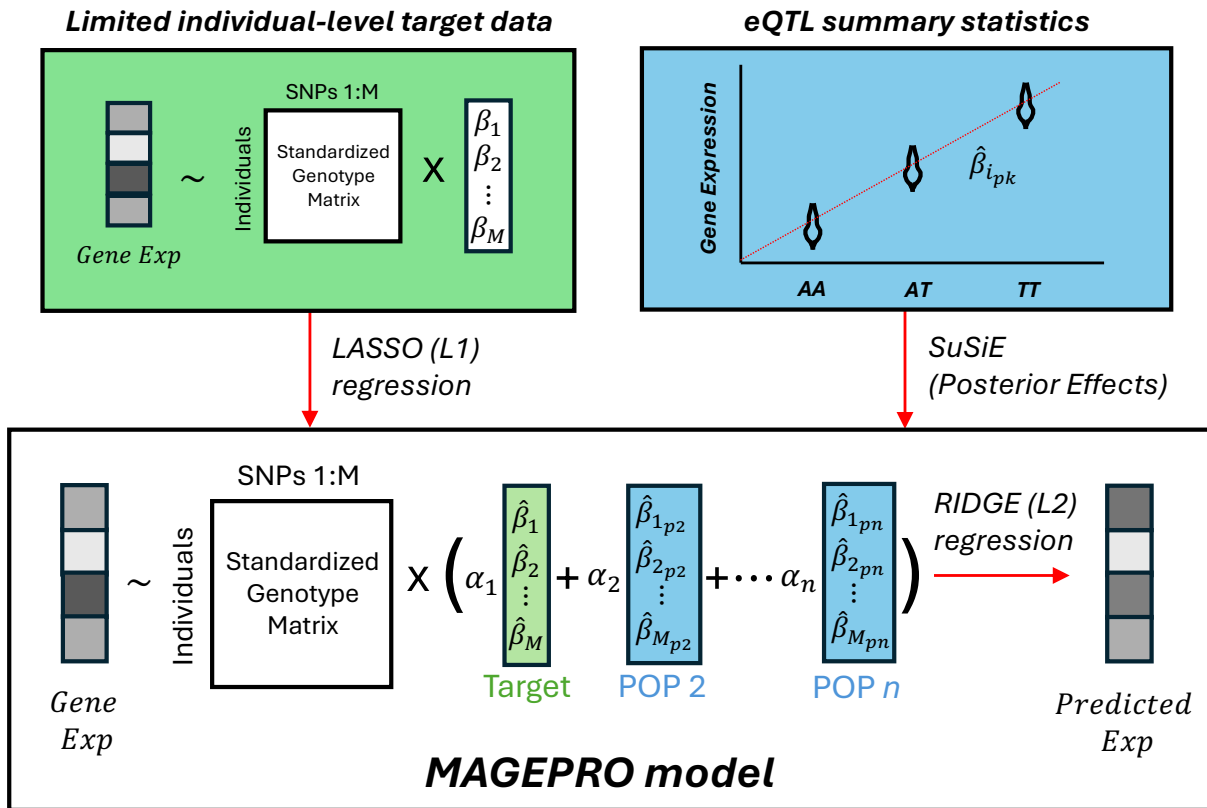
1188 **Tables**

Cohort	Population	Tissue or Cell type	Sample Size	eQTLGen EUR WB	GTE _x EUR WB	MESA HIS MONO	MESA AA MONO	GENOA AA LCL
MESA	AA	Monocyte (MONO)	224	✓	✓	✓		✓
MESA	HIS	Monocyte (MONO)	242	✓	✓		✓	✓
MESA	EUR	Monocyte (MONO)	574	✓	✓	✓	✓	✓
GTE _x	AA	Whole Blood (WB)	80	✓		✓	✓	✓
GTE _x	EUR	Whole Blood (WB)	568	✓		✓	✓	✓
GTE _x	EUR	Lung	440	✓		✓	✓	✓
GENOA	AA	LCL	346	✓	✓	✓	✓	
GEUVADIS	EUR	LCL	364	✓	✓	✓	✓	

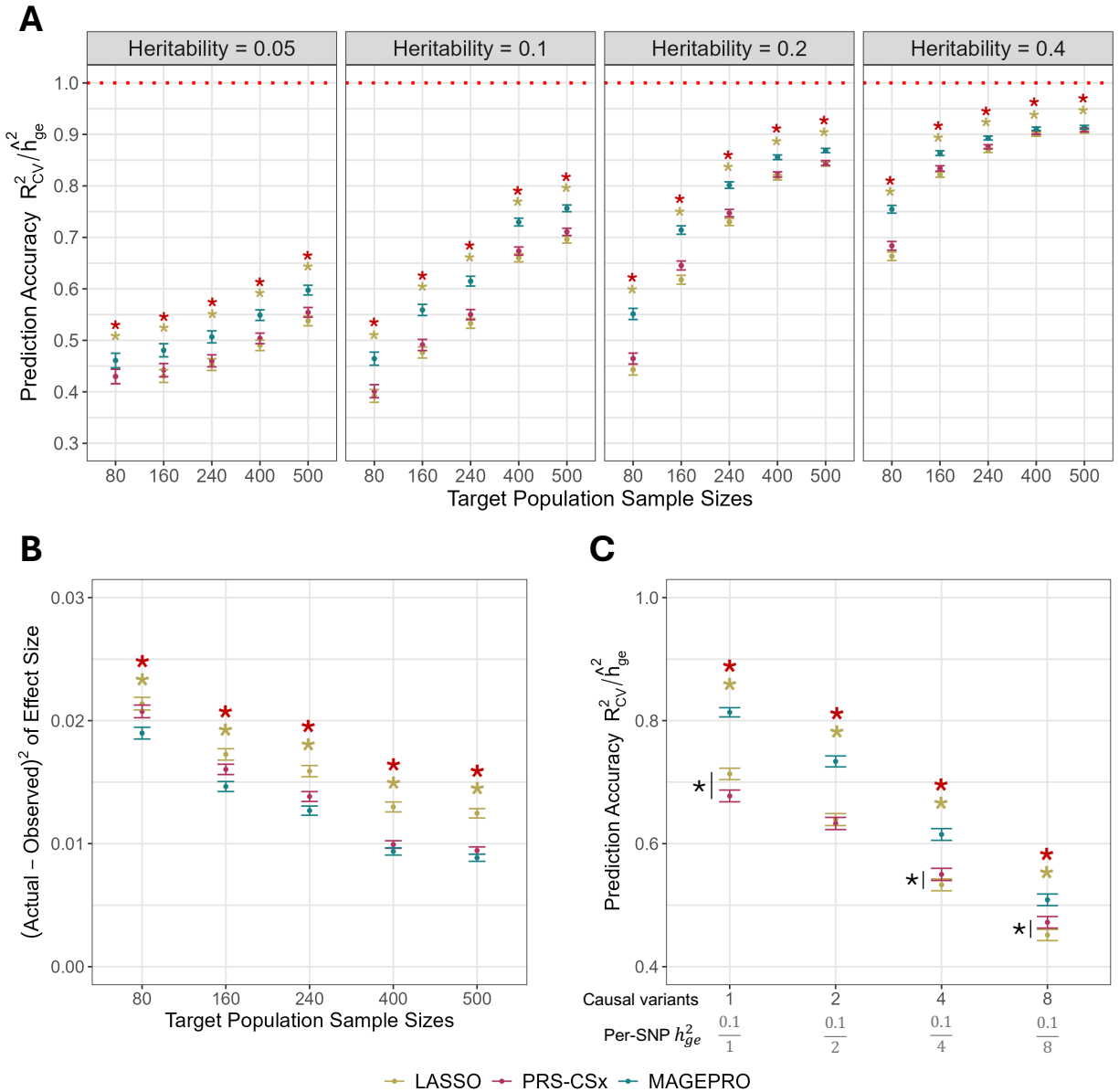
1189

1190 **Table 1. External eQTL summary statistics used for each target cohort.** Rows correspond to
 1191 each target cohort for which individual-level gene expression and genotype data were used to
 1192 create genetic models of gene expression. The last five columns correspond to external eQTL
 1193 summary statistics used as inputs to MAGEPRO. We avoided using external summary statistics
 1194 that contain the same individuals as the target cohort to prevent over-fitting and inflation of
 1195 cross-validation results. Sample sizes indicate the number of individuals in a target cohort after
 1196 relatedness-based filtering (Methods). AA, African American; HIS, Hispanic/Latino; EUR,
 1197 European; LCL: lymphoblastoid cell line.

1198 **Figures**

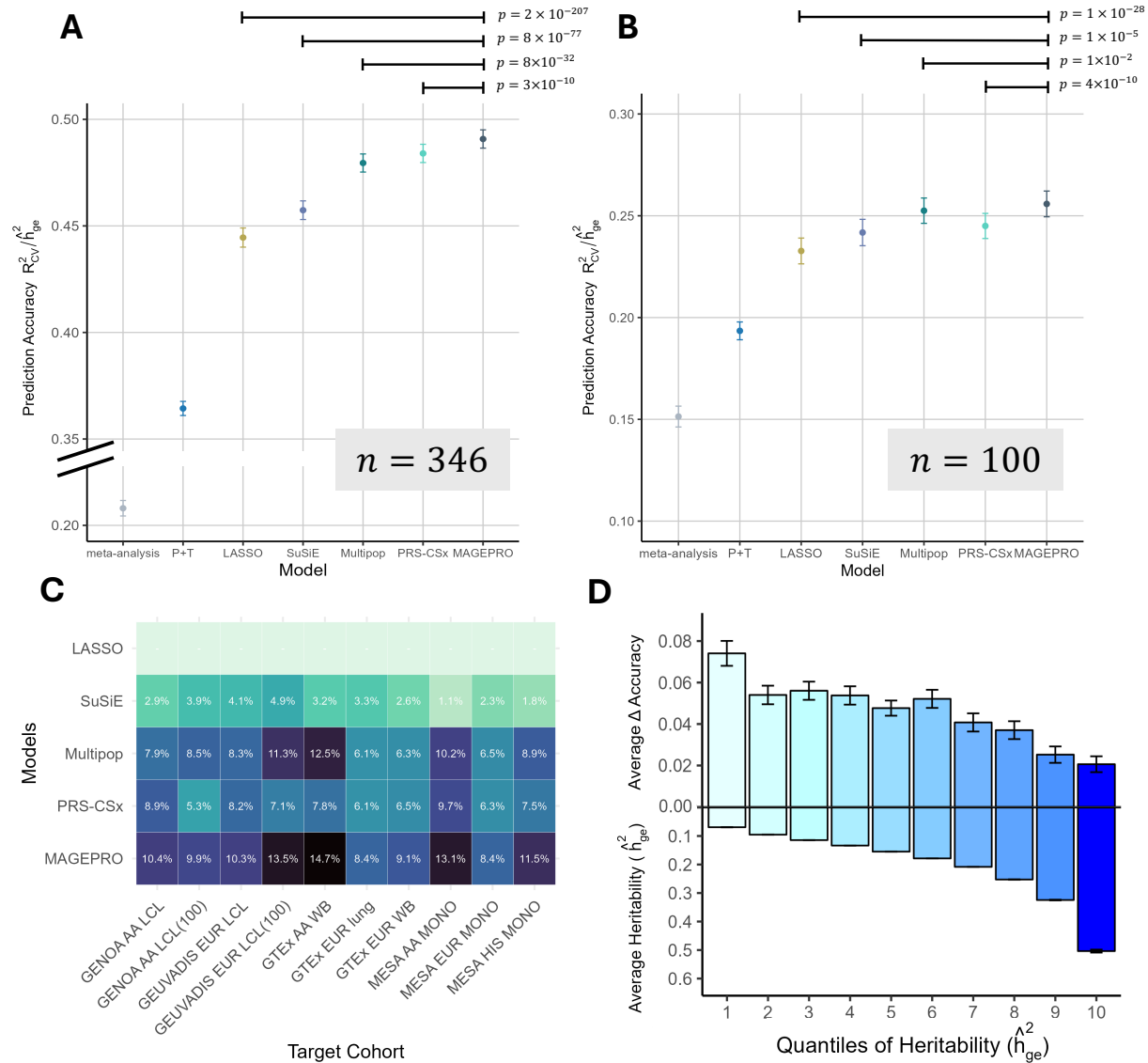


1199
 1200 **Figure 1. Overview of the MAGEPRO model.** Schema of the MAGEPRO model for one gene.
 1201 MAGEPRO takes limited individual-level target data (green) and external eQTL summary
 1202 statistics (blue) as input. Red arrows indicate the three main operations of MAGEPRO. First,
 1203 individual-level gene expression and standardized genotypes are used to estimate noisy effect
 1204 sizes for the target population ($\hat{\beta}_i$ for SNP i) using an L1-regularized linear regression. Next, we
 1205 estimate the posterior effect size estimates for each set of external eQTL summary statistics
 1206 using SuSiE, designated by $\hat{\beta}_{i,pk}$ for SNP i and population k . Finally, we estimate optimal mixing
 1207 weights of effect sizes across all populations, including the target, using L2-regularized linear
 1208 regression (α_k for population k). The *cis*-heritability of the gene expression (\hat{h}_{ge}^2) is estimated
 1209 using the limited individual-level target data and is used to normalize the prediction accuracy
 1210 $\left(\frac{R_{cv}^2}{\hat{h}_{ge}^2}\right)$ to allow comparisons across genes with different heritabilities.



1211
 1212 **Figure 2. MAGEPRO outperforms alternative gene expression prediction models in**
 1213 **various simulated architectures.** (A) Predictive accuracy of LASSO, PRS-CSx, and
 1214 MAGEPRO across different gene expression heritability and sample size settings. Across all
 1215 settings, genes were simulated with four causal variants. Accuracy is calculated as the ratio of
 1216 the cross-validation R_{cv}^2 and the GCTA-estimated *cis*-heritability of gene expression (\hat{h}_{ge}^2). (B)
 1217 Squared difference between the simulated (actual) and estimated effect sizes of the four causal
 1218 variants per gene. *Cis*-heritability was set to 10%. (C) Predictive accuracy ($\frac{R_{cv}^2}{\hat{h}_{ge}^2}$) of methods
 1219 while varying the number of causal variants and maintaining the total *cis*-heritability (h_{ge}^2) at
 1220 10%. Sample size was set to 240. In all panels, data are presented as mean values across 1,000
 1221 independently simulated genes with confidence intervals representing ± 1 standard error. Yellow

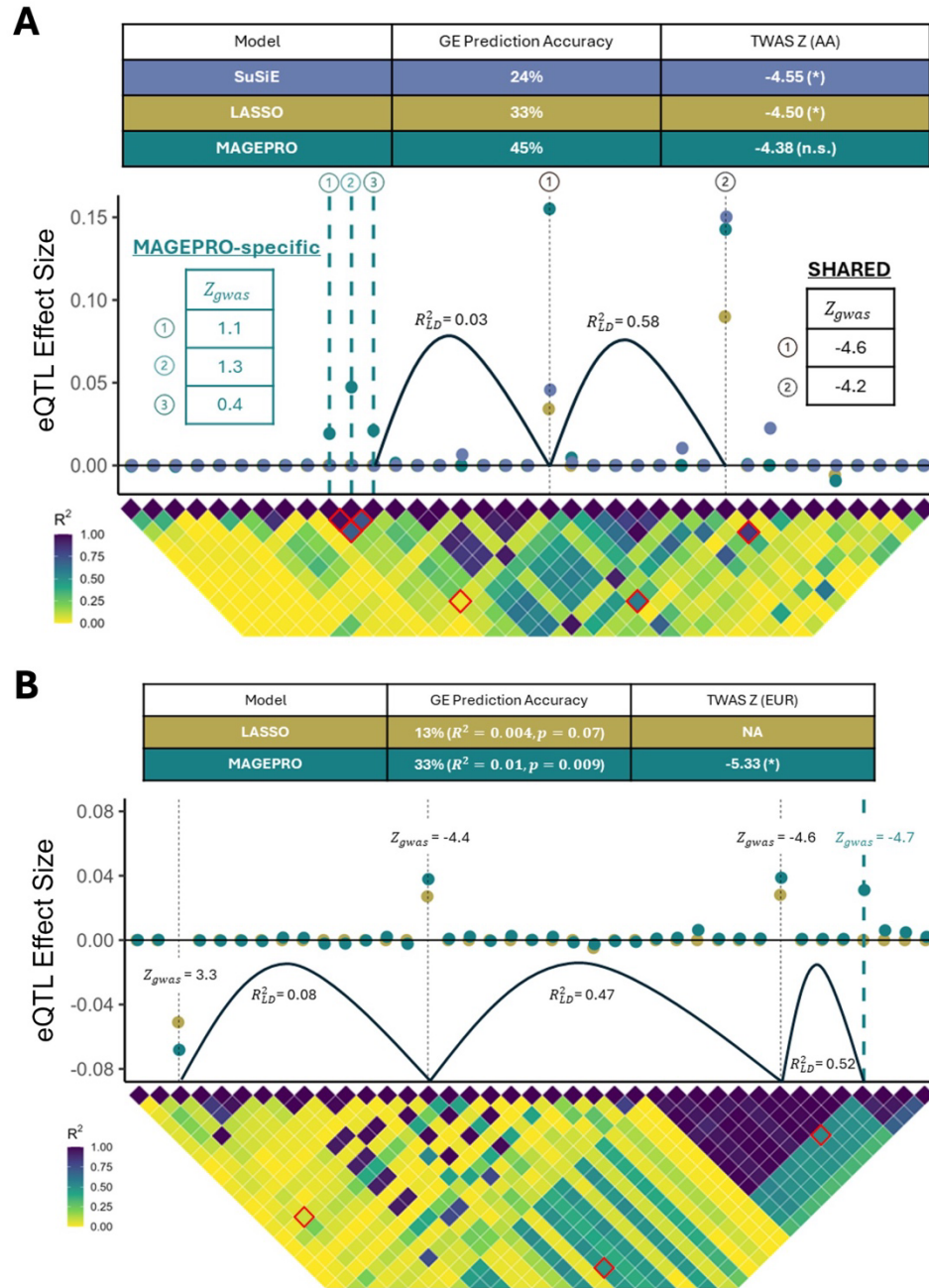
1222 (resp. red) asterisks indicate that the difference between MAGEPRO and LASSO (resp. PRS-
1223 CSx) results is significant. Black asterisks highlight pairwise comparisons. All hypothesis tests
1224 are two-sided paired t-tests. Numerical results are reported in **Supplementary Tables 1-6**.



1225

1226 **Figure 3. MAGEPRO outperforms alternative methods in real data.** (A,B) Comparison of
 1227 the accuracy of different models in predicting LCL gene expression in the GENOA AA
 1228 population at two different sample sizes (A: full cohort, B: random down-sampling to 100
 1229 individuals). P-values are derived from a one-sided paired t-test, testing the alternative
 1230 hypothesis that MAGEPRO produces larger accuracies. Comparisons between MAGEPRO and
 1231 meta-analysis or P+T not annotated due to low precision to estimate such small p-values. (C)
 1232 Performance of the top five gene expression prediction methods across the eight different target
 1233 cohorts from **Table 1** plus two randomly down-sampled cohorts, indicated by (100). Values in
 1234 the heatmap are the percent change in predictive accuracy relative to LASSO regression. All
 1235 percent differences are significant according to one-sided paired t-tests ($p < 0.05$). (D) The
 1236 average change in accuracy between MAGEPRO and LASSO ($\text{MAGEPRO } \frac{R_{CV}^2}{\hat{h}_{ge}^2} - \text{LASSO } \frac{R_{CV}^2}{\hat{h}_{ge}^2}$)
 1237 across 10 quantiles of genes grouped by GCTA-estimated *cis*-heritability. Accuracy and

1238 heritability values were estimated for LCL gene expression in the GENOA AA population. In
1239 panels A, B, and D, data are presented as mean values with confidence intervals representing ± 1
1240 standard error. LCL, lymphoblastoid cell line; AA, African American. Numerical results are
1241 reported in **Supplementary Tables 7-10**.



1242

1243 **Figure 4. Gene-disease associations are sensitive to gene expression prediction models. (A)**

1244 An example of a TWAS association that becomes non-significant after MAGEPRO has

1245 improved the accuracy of the gene model. MAGEPRO introduces three variants to the gene

1246 model for *ZNF213-ASI*, trained on GENOA LCL data from African American individuals, that

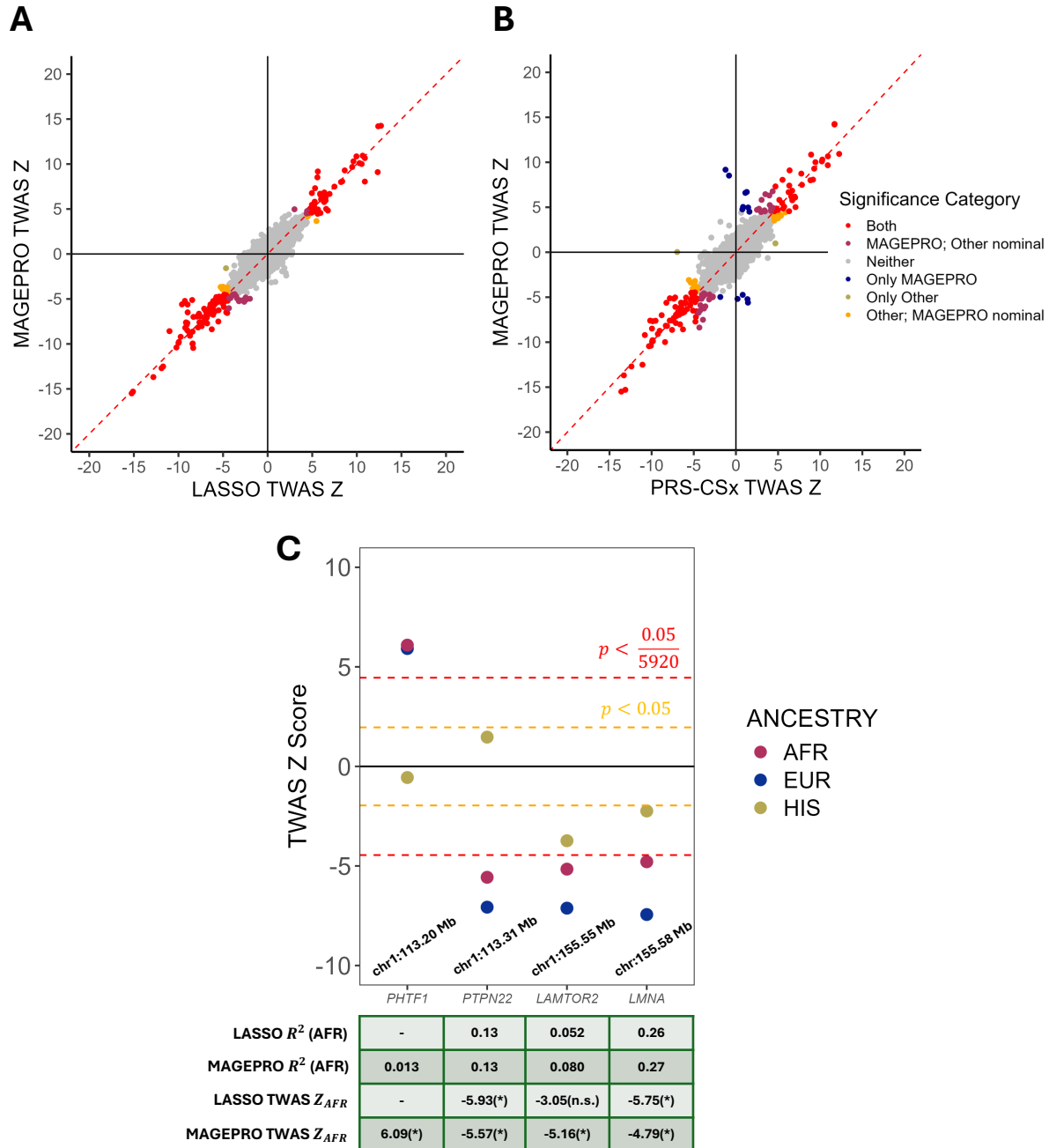
1247 do not colocalize well with a GWAS for red blood cell count (teal). (B) An example of a new

1248 TWAS association introduced by MAGEPRO. *RGS14* is newly associated with asthma based on

1249 a European monocyte model from the MESA cohort. In both panels, the asterisk (*) indicates

1250 significance using a Bonferroni threshold across cohort-specific *cis*-heritable genes ($\frac{0.05}{6872}$ for A,

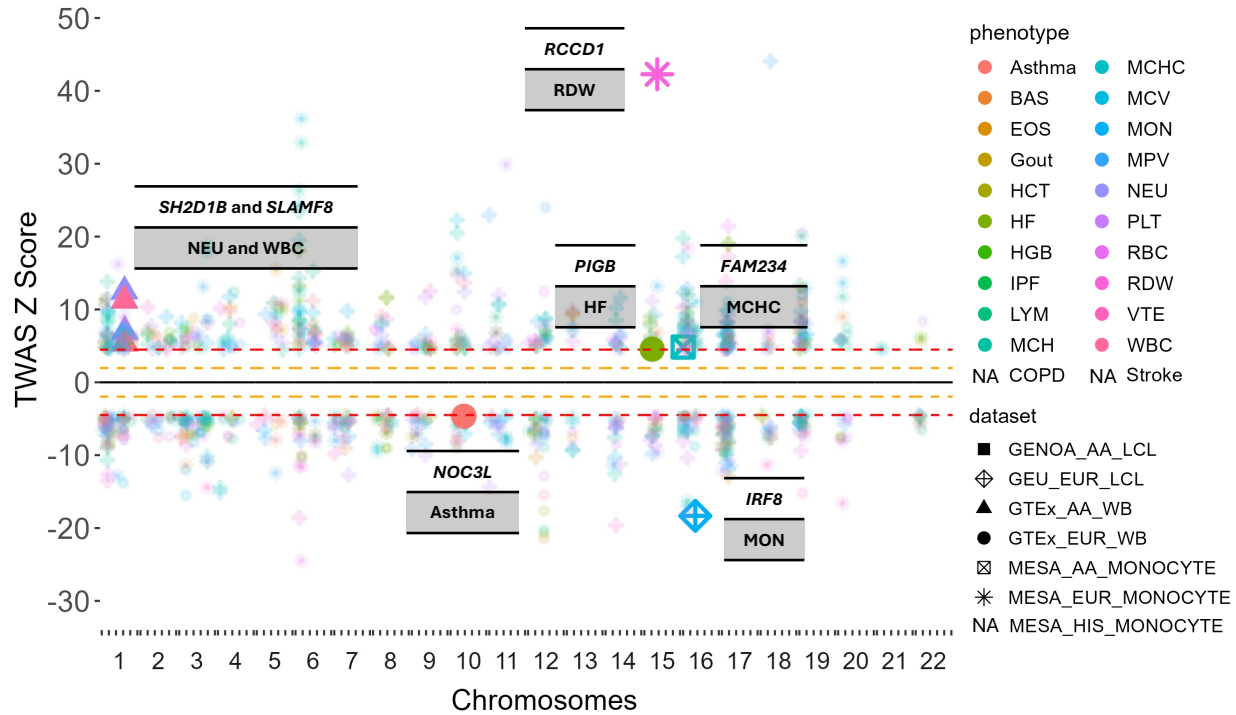
1251 $\frac{0.05}{5920}$ for B). The dot plot shows the effect sizes inferred by the *cis*-genetic model of gene
1252 expression created by each method. Black dotted vertical lines designate eQTL effects identified
1253 by all/both models and teal dotted vertical lines designate effects captured specifically by
1254 MAGEPRO. In the heatmap below, R_{LD}^2 values that are relevant to potential interactions between
1255 variants are boxed in red. Distances between SNPs are not to scale; the x-axis indicates the
1256 indices of the *cis*-SNPs ordered by increasing genomic coordinate. LCL, lymphoblastoid cell
1257 line; GE, gene expression. Numerical results are reported in **Supplementary Tables 12-13**.



1258

1259 **Figure 5. Applying MAGEPRO to improve Monocyte gene expression prediction across**
 1260 **three ancestries identifies novel genes associated with blood cell traits.** (A,B) Comparison of
 1261 TWAS z-scores between MAGEPRO and other gene expression prediction methods for the
 1262 African (AFR) ancestry. Colors correspond to groups of significance described in the legend.
 1263 “Other” refers to the model in comparison on the x-axis. Results are aggregated across 15 blood
 1264 cell traits. (C) Miami plot of TWAS associations with white blood cell counts across three
 1265 different ancestries. Green table display the gene expression prediction R^2 and TWAS z-score in
 1266 the AFR population (statistics for other population are presented in **Supplementary Figure 22**).

1267 Positions indicate the start of the *cis* window, nominal significance threshold is $p < 0.05$ and
1268 Bonferroni significance threshold is $p < \frac{0.05}{5920}$ for all panels. Numerical results are reported in
1269 **Supplementary Tables 16-18.**



1270

1271 **Figure 6. New gene models created by MAGEPRO recapitulate previous findings and**
 1272 **identify several biologically plausible new findings.** Miami plot shows genome-wide signed
 1273 TWAS associations from analysis of 22 unique complex traits/diseases across three ancestries
 1274 (represented by seven independent cohorts). Only gene-trait associations resulting from new
 1275 gene models created by MAGEPRO (MAGEPRO $R^2 > 0$, $p < 0.05$ while LASSO R^2 not
 1276 significantly greater than 0) and passing Bonferroni significance are plotted. Phenotypes and
 1277 datasets are labeled “NA” if there are no such associations. Examples highlighted in the text are
 1278 labeled with the gene symbol, associated phenotype, and enlarged point. Yellow dotted line
 1279 indicates nominal $p < 0.05$ and red dotted line indicates Bonferroni threshold

1280 ($p < \frac{0.05}{\# \text{ genes tested in dataset}}$). BAS, basophil count; EOS, eosinophil count; HCT, hematocrit;
 1281 HF, heart failure; HGB, hemoglobin concentration; IPF, idiopathic pulmonary fibrosis; LYM,
 1282 lymphocyte count; MCH, mean corpuscular hemoglobin; COPD, Chronic obstructive pulmonary
 1283 disease; MCHC, mean corpuscular hemoglobin concentration; MCV, mean corpuscular volume;
 1284 MON, monocyte count; MPV, mean platelet volume; NEU, neutrophil count; PLT, platelet
 1285 count; RBC, red blood cell count; RDW, red blood cell distribution width; VTE, venous
 1286 thromboembolism; WBC, total white blood cell count. Numerical results are reported in
 1287 **Supplementary Table 20.**