

1 **TITLE**

2

3 On the cross-population portability of gene expression prediction models

4

5 **AUTHORS**

6

7 Kevin L. Keys<sup>1,2,\*</sup>, Angel C.Y. Mak<sup>1</sup>, Marquitta J. White<sup>1</sup>, Walter L. Eckalbar<sup>1</sup>, Andy Dahl<sup>1</sup>, Joel  
8 Mefford<sup>1</sup>, Anna V. Mikaylova<sup>3</sup>, María G. Contreras<sup>1,4</sup>, Jennifer R. Elhawary<sup>1</sup>, Celeste Eng<sup>1</sup>, Donglei  
9 Hu<sup>1</sup>, Scott Huntsman<sup>1</sup>, Sam S. Oh<sup>1</sup>, Sandra Salazar<sup>1</sup>, Michael A. Lenoir<sup>5</sup>, Jimmie C. Ye<sup>6,7</sup>, Timothy  
10 A. Thornton<sup>3</sup>, Noah Zaitlen<sup>8</sup>, Esteban G. Burchard<sup>1,7,+</sup>, and Christopher R. Gignoux<sup>9,10,+\*</sup>.

11

12 <sup>1</sup> Department of Medicine, University of California, San Francisco, CA, USA

13 <sup>2</sup> Berkeley Institute for Data Science, University of California, Berkeley, California, USA

14 <sup>3</sup> Department of Biostatistics, University of Washington, Seattle, WA, USA

15 <sup>4</sup> San Francisco State University, San Francisco, CA, USA

16 <sup>5</sup> Bay Area Pediatrics, Oakland, CA, USA

17 <sup>6</sup> Department of Epidemiology and Biostatistics, University of California, San Francisco, CA, USA

18 <sup>7</sup> Department of Bioengineering and Therapeutic Biosciences, University of California, San  
19 Francisco, CA, USA

20 <sup>8</sup> Department of Neurology, University of California, Los Angeles, CA, USA

21 <sup>9</sup> Colorado Center for Personalized Medicine, University of Colorado, Anschutz Medical Campus,  
22 Aurora, CO, USA

23 <sup>10</sup> Department of Biostatistics and Informatics, School of Public Health, University of Colorado,  
24 Anschutz Medical Campus, Aurora, CO, USA

25

26 + Shared senior authorship.

27 \* Correspondence should be addressed to Kevin L. Keys ([klkeys@g.ucla.edu](mailto:klkeys@g.ucla.edu)) and Christopher R.  
28 Gignoux ([chris.gignoux@ucdenver.edu](mailto:chris.gignoux@ucdenver.edu))

29

30 **ABSTRACT**

31

32 The genetic control of gene expression is a core component of human physiology. For the past  
33 several years, transcriptome-wide association studies have leveraged large datasets of linked  
34 genotype and RNA sequencing information to create a powerful gene-based test of association  
35 that has been used in dozens of studies. While numerous discoveries have been made, the  
36 populations in the training data are overwhelmingly of European descent, and little is known  
37 about the portability of these models to other populations. Here, we test for cross-population  
38 portability of gene expression prediction models using a dataset of African American individuals  
39 with RNA-Seq data in whole blood. We find that the default models trained in large datasets such  
40 as GTEx and DGN fare poorly in African Americans, with a notable reduction in prediction  
41 accuracy when compared to European Americans. We replicate these limitations in cross-  
42 population portability using the five populations in the GEUVADIS dataset. Via simulations of both  
43 populations and gene expression, we show that accurate cross-population portability of  
44 transcriptome imputation only arises when eQTL architecture is substantially shared across

45 populations. In contrast, models with non-identical eQTL showed patterns similar to real-world  
46 data. Therefore, generating RNA-Seq data in diverse populations is a critical step towards multi-  
47 ethnic utility of gene expression imputation.

48

#### 49 **KEYWORDS**

50 TWAS, gene expression, admixed populations, GTEx, PrediXcan

51

#### 52 **MANUSCRIPT TEXT**

53

54 In the last decade, large-scale genome-wide genotyping projects have enabled a revolution in our  
55 understanding of complex traits.<sup>1-4</sup> This explosion of genome sequencing data has spurred the  
56 development of new methods that integrate large genotype sets with additional molecular  
57 measurements such as gene expression. A recently popular integrative approach to genetic  
58 association analyses, known as a transcriptome-wide association study (TWAS)<sup>5,6</sup>, leverages  
59 reference datasets such as the Genotype-Tissue Expression (GTEx) repository<sup>7</sup> or the Depression  
60 and Genes Network (DGN)<sup>8</sup> to link associated genetic variants with a molecular trait like gene  
61 expression. The general TWAS framework requires previously estimated *cis*-eQTL for all genes in  
62 a dataset with both genotype and gene expression measurements. The resulting eQTL effect sizes  
63 build a predictive model that can impute gene expression in an independently genotyped  
64 population. A TWAS is similar in spirit to the widely-known genome-wide association study  
65 (GWAS) but suffers less of a multiple testing burden and can potentially detect more associations  
66 as a result.<sup>5,6</sup>

67

68 Unlike a normal GWAS, where phenotypes are regressed onto genotypes, in TWAS the phenotype  
69 is regressed onto the imputed gene expression values, thus constituting a new gene-based  
70 association test. TWAS can also link phenotypes to variation in gene expression and provide  
71 researchers with additional biological and functional insights over those afforded by GWAS alone.  
72 While these models are imperfect predictors, imputing gene expression allows researchers to  
73 test phenotype associations to expression levels in existing GWAS datasets without measuring  
74 gene expression directly. In particular, these methods enable analysis of predicted gene  
75 expression in very large cohorts ( $\sim 10^4 - 10^6$  individuals) rather than typical gene expression  
76 studies that measure expression directly ( $\sim 10^2 - 10^3$  individuals). Several methods have been  
77 recently developed to perform TWAS in existing genotyped datasets. PrediXcan<sup>6</sup> uses eQTL  
78 precomputed from paired genotype-expression data, such as those in GTEx, in conjunction with  
79 a new genotype set to predict gene expression. These gene expression prediction models are  
80 freely available online (PredictDB). Related TWAS approaches, such as FUSION<sup>5</sup>, MetaXcan<sup>9</sup>, or  
81 SMR<sup>10</sup>, leverage eQTL with GWAS summary statistics instead of requiring the availability of raw  
82 individual-level genotype data.

83

84 As evidenced by application to numerous disease domains, the TWAS framework is capable of  
85 uncovering new genic associations.<sup>11-17</sup> However, the power of TWAS is inherently limited by the  
86 data used for eQTL discovery. For example, since gene expression varies by tissue type,  
87 researchers must ensure that the prediction weights are estimated using RNA from a tissue  
88 related to their phenotype, whether that be the direct tissue of interest or one with sufficiently

89 correlated gene expression.<sup>18</sup> Furthermore, the ability of predictive models to impute gene  
90 expression from genotypes is limited by the heritability in the cis region around the gene.<sup>6</sup>  
91 Consequently, genes with little or no measurable genetically regulated effect on their expression  
92 in the discovery data would not be good candidates for TWAS.

93

94 A subtler but more troubling issue arises from the lack of genetic diversity present in the datasets  
95 used for predictive model training: most paired genotype-expression datasets consist almost  
96 entirely of data from European-descent individuals. The European overrepresentation in genetic  
97 studies is well documented<sup>19–21</sup> and has severe negative consequences for equity as well as for  
98 gene discovery<sup>22</sup>, fine mapping<sup>23–25</sup>, and applications in personalized medicine.<sup>26–34</sup> Genetic  
99 architecture and genotype frequencies can vary across populations, which presents a potential  
100 problem for the application of predictive models with genotype predictors across multiple  
101 populations.

102

103 The training data for most models in PredictDB are highly biased toward European ancestry: GTEx  
104 version v6p subjects are over 85% European, while the GTEx v7 and DGN subjects are entirely of  
105 European descent. The lack of suitable genotype-expression datasets in non-European  
106 individuals leads to scenarios in which PredictDB models trained in Europeans are used to impute  
107 into non-European or admixed populations. As shown previously in the context of polygenic risk  
108 scores<sup>35</sup>, multi-SNP prediction models trained in one population can suffer from unpredictable  
109 bias and poor prediction accuracy that impair their cross-population portability. Recent analyses  
110 of genotype-expression data from the Multi-Ethnic Study of Atherosclerosis (MESA)<sup>36–38</sup> explore  
111 cross-population transcriptome imputation and conclude that predictive accuracy is highest  
112 when training and testing populations match in ancestry. These results dovetail with our  
113 experience analyzing diverse populations, but offer little insight into the mechanisms underlying  
114 the cross-population portability of transcriptome prediction models, particularly when eQTL  
115 architecture is known.

116

117 Here, we investigate the cross-population portability of gene expression models using paired  
118 genotype and gene expression data and using simulations derived from real genotypic data and  
119 realistic models of gene expression. We analyze prediction quality from currently available  
120 PrediXcan prediction weights using a pilot subset of paired genotype and whole blood  
121 transcriptome data from the Study of African Americans, Asthma, Genes, and Environment  
122 (SAGE).<sup>39–42</sup> SAGE is a pediatric cohort study of childhood-onset asthma and pulmonary  
123 phenotypes in African American subjects of 8 to 21 years of age. To tease apart cross-population  
124 prediction quality, we turn to GEUVADIS and the 1000 Genomes Project datasets.<sup>4,43</sup> The  
125 GEUVADIS dataset has been used extensively to validate PrediXcan models.<sup>6,38</sup> However, recent  
126 analyses suggest that GTEx and DGN PrediXcan models behave differently on the constituent  
127 populations in GEUVADIS.<sup>44</sup> To our knowledge, nobody has investigated cross-population  
128 portability within GEUVADIS. GEUVADIS provides us an opportunity to investigate predictive  
129 models with an experimentally homogeneous dataset: the GEUVADIS RNA-Seq data were  
130 produced in the same environment under the same protocol, from lymphoblastoid cell lines  
131 (LCLs) derived from similar sampling efforts, providing a high degree of technical harmonization.  
132 We train, test, and validate predictive models wholly within GEUVADIS with a nested cross-

133 validation scheme. Finally, to understand the consequences of eQTL architecture on TWAS, we  
134 use existing 1000 Genomes data to simulate two ancestral populations and an admixed  
135 population and then apply the same “train-test-validate” scheme with various simulated eQTL  
136 models to study cross-population prediction efficacy when a gold standard is known.

137

138 We compared transcriptome imputation accuracy in SAGE whole blood RNA using three  
139 PredictDB prediction weight sets for whole blood RNA: GTEx v6p, GTEx v7, and DGN. We also  
140 evaluated expression prediction with all four MESA monocyte weight sets: MESA\_ALL  
141 (populations combined), MESA\_AFA (African Americans), MESA\_AFHI (combined African  
142 Americans and Hispanic Americans), and MESA\_CAU (Caucasians). For each gene where both  
143 measured RNA-Seq gene expression and predictions are available in SAGE, we compute both the  
144 coefficient of determination ( $R^2$ ) and Spearman correlation to analyze the direction of prediction.  
145 As we are primarily interested in describing the relationship between predicted outcome and real  
146 outcome, we prefer Spearman’s  $\rho$  to describe correlations, while for determining prediction  
147 accuracy, we use the standard regression  $R^2$ , corresponding to the squared Pearson correlation,  
148 to facilitate comparisons to prior work. We then benchmark these against the out-of-sample  $R^2$   
149 and correlations from GTEx v7 and MESA as found in PredictDB. Prediction results in SAGE were  
150 available for 11,545 genes with a predictive model from at least one weight set. Not all sets  
151 derived models at the same genes: the prediction results across all weight sets overlapped at 273  
152 genes, of which 39 genes had predictions with positive correlation to measurements. Since the  
153 estimation of these prediction models requires both high quality expression data and inferred  
154 eQTL, each weight set may have a different number of gene models. Therefore, intersecting  
155 seven different weight sets reduces the overall number of models available for comparison. This  
156 small number of genes in common is largely driven by MESA\_AFA, the repository with the  
157 smallest number of predictive models. MESA\_AFA contains the models that should best reflect  
158 the genetic ancestry in SAGE (Supplementary Table 1). We note that MESA\_AFA also has the  
159 smallest training sample size among our weight sets ( $N = 233$ )<sup>38</sup>, so the small number of predicted  
160 genes from MESA\_AFA probably results from the small training sample size and not from any  
161 feature of the underlying MESA\_AFA training data.

162

163 The concordance between predicted and measured gene expression over the 273 genes in  
164 common to all seven weight sets, with corresponding training metrics from PredictDB as  
165 benchmarks, shows worse performance than expected for  $R^2$  (Figure 1) and correlations (Figure  
166 2). The highest mean  $R^2$  of 0.0336 was observed in DGN. Here, we highlight the intersection of  
167 genes across model sets for investigation, but the overall patterns for all genes are similar; results  
168 for the 11,545 total genes (Supplementary Figure 1) and the 39 genes with positive correlations  
169 (Supplementary Figure 2) showed little appreciable deviation from  $R^2$  shown in Figure 1. Because  
170 SAGE is an independent validation set for the training populations, we would expect to observe  
171 some deterioration in imputation  $R^2$  due to differences in population structure and linkage  
172 disequilibrium. However, Figure 1 shows a marked difference in model performance.

173

174 More noteworthy is the substantial proportion of predictions in SAGE with negative correlations  
175 to the real data. All seven weight sets produced negative mean correlations. The least negative  
176 mean correlation (-0.0044) was observed with GTEx\_v6p, while the most negative mean

177 correlation (-0.020) was observed with MESA\_AFA (Supplementary Table 1). The fact that  
178 correlations to SAGE measurements are negative on average suggests that some large  $R^2$  values  
179 in Figure 1 may result from gene models with incorrect direction of prediction. While there are  
180 some fluctuations in prediction accuracy, no prediction weight set produces practically  
181 meaningfully better correlations to data than the others (-0.020 to -0.044). In contrast, the  
182 published PredictDB models for these genes show positive correlations to their training data,  
183 indicating no obvious incapacity for accurate prediction. However, available predictions into  
184 SAGE from otherwise valid prediction models are uniformly limited in power to capture true  
185 genotype-expression relationships.

186  
187 To analyze genes with ostensibly high imputation  $R^2$ , we focus on genes in GTEx v7 with cross-  
188 validated  $R^2 > 0.2$  in the reference population. Figure 3 compares PredictDB testing  $R^2$  against the  
189 empirical  $R^2$  from regressing predictions onto observations in SAGE. In this case, even the better-  
190 imputed gene models derived from PredictDB have limited ability to capture gene expression  
191 accurately in SAGE (mean  $R^2$  0.031, IQR [0.0027, 0.037]).

192  
193 It is important to note that real-world comparisons of RNA-Seq datasets can be subject to  
194 numerous sources of heterogeneity besides differential ancestry. Possible confounders include  
195 technical differences in sequencing protocols, differences in the age of participants<sup>45</sup>, and the  
196 postmortem interval to tissue collection (for GTEx).<sup>46-48</sup> To investigate cross-population  
197 portability in an experimentally homogeneous context, we turn to GEUVADIS.<sup>43</sup> The GEUVADIS  
198 data include two continental population groups from the 1000 Genomes Project: the Europeans  
199 (EUR373), composed of 373 unrelated individuals from four subpopulations (Utahns (CEU), Finns  
200 (FIN), British (GBR), Toscani (TSI)), and the Africans (AFR) composed of 89 unrelated Yoruba (YRI)  
201 individuals. In light of the known bottleneck in Finnish population history, we analyze EUR373  
202 both as one population and as two independent subgroups: the 95 Finnish individuals (FIN) and  
203 the 278 non-Finnish Europeans (EUR278). We used matched RNA-Seq, generated and  
204 harmonized together by the GEUVADIS Consortium and whole-genome genotype data in the  
205 resulting four populations (EUR373, EUR278, FIN, and AFR) to train predictive models for gene  
206 expression in a nested cross-validation scheme<sup>6</sup> and perform cross-population tests of  
207 imputation accuracy.

208  
209 Table 1 shows  $R^2$  from three training sets (EUR373, EUR278 and AFR) into the four testing  
210 populations (EUR373, EUR278, FIN, and AFR) for genes with positive correlation between  
211 prediction and measurement. While the number of genes with applicable models including  
212 genetic data varies in each train-test scenario (see Supplementary Table 3), we note that not all  
213 predictive models are trained on equal sample sizes, so the resulting  $R^2$  only provide a general  
214 idea of how well one population imputes into another. Analyses within a population use out-of-  
215 sample imputation  $R^2$  to avoid overfitting across train-test scenarios. Predicting from a  
216 population into itself yields  $R^2$  ranging from 0.079 – 0.098 (Table 1) consistent with the smaller  
217 sample sizes in GEUVADIS versus GTEx and DGN. In contrast, predicting across populations yields  
218 more variable predictions, with  $R^2$  ranging from 0.029 – 0.087. At the lower range of  $R^2$  (0.029 –  
219 0.039) are predictions from AFR into European testing groups (EUR373, EUR278, and FIN).  
220 Alternatively, when predicting from European training groups into AFR, the  $R^2$  are noticeably



221 higher (0.051 – 0.054). Prediction from EUR278 into FIN ( $R^2 = 0.087$ ) is better than prediction  
222 from EUR278 into AFR ( $R^2 = 0.051$ ), suggesting that imputation  $R^2$  may deteriorate with increased  
223 genetic distance. A comparison of the 564 genes in common across all train-test scenarios (Table  
224 2) yields a more equal basis of comparison between populations, albeit from a subset of genes  
225 with potentially more consistent gene expression levels. In this case involving better-predicted  
226 genes, we see that imputation quality between the European groups improves noticeably, with  
227  $R^2$  ranging between 0.183 to 0.216, while  $R^2$  between Europeans and Africans ranges from 0.095  
228 to 0.147. In general, populations seem to predict better when imputing into themselves, and less  
229 well when imputing into other populations.

230  
231 Combining all European subpopulations obscures population structure and can complicate  
232 analysis of cross-population imputation performance. To that end, we divide the GEUVADIS data  
233 into its five constituent populations and randomly subsample each of them to the smallest  
234 population size ( $n = 89$ ). We then estimate models from each subpopulation and predict into all  
235 five subpopulations. Table 3 shows average  $R^2$  from each population into itself and others. The  
236 populations consistently impute well into themselves, with imputation  $R^2$  ranging from 0.104 –  
237 0.136. However, a notable difference exists between the EUR subpopulations and YRI. The cross-  
238 population  $R^2$  between CEU, TSI, GBR, and FIN ranges from 0.103 to 0.137, while cross-population  
239  $R^2$  from these populations into YRI ranges from 0.062 to 0.084. Imputation between YRI and the  
240 EUR populations taken together is consistently lower than within the EUR populations  
241 (Supplementary Figure 3) and statistically significant ( $p$ -value  $< 1.36 \times 10^{-4}$ , Dunn test; see  
242 Supplementary Table 6). The cross-population differences remain for the 142 genes with positive  
243 correlation in all train-test scenarios (Table 4), where  $R^2$  for imputation into YRI ranges from 0.166  
244 to 0.244, while imputation within EUR populations ranges from 0.239 to 0.331. These results  
245 clearly suggest problems for prediction models that impute gene expression across populations,  
246 in similar regimes to those tested with linear predictive models and datasets of size consistent  
247 with current references. In addition, since AFR is genetically more distant from the EUR  
248 subpopulations than they are to each other, we interpret these results to imply that structure in  
249 populations can potentially exacerbate cross-population imputation quality (Supplementary  
250 Figure 4).

251  
252 The unresolved question is the extent to which these results hold with oracle knowledge of eQTL  
253 architecture, something impossible to investigate in real data when the causal links between  
254 eQTL and gene expression can only be estimated. To investigate genomic architectures giving rise  
255 to gene expression, and in particular to investigate behavior in admixed populations, we simulate  
256 haplotypes from HapMap3<sup>50</sup> CEU and YRI using HAPGEN2<sup>51</sup> and then sample haplotypes in  
257 proportions consistent with observed admixture proportions (80% YRI, 20% CEU)<sup>52</sup> to construct  
258 a simulated African-American (AA) admixed population. We simulate eQTL architectures under  
259 an additive model of size  $k$  causal alleles ( $k = 1, 5, 10, \text{ and } 20$ ) and a phenotype with *cis*-heritability  
260  $h^2 = 0.15$  (recapitulating average  $h^2$  in GTEx) using the genomic background of genic regions on  
261 chromosome 22, thus testing various model sizes and LD patterns. To tease apart the effect of  
262 shared eQTL architecture, we allow the populations to share eQTL with fixed effects in various  
263 proportions (0%, 10%, 20%, ..., 100%). With these simulations providing known architectures for  
264 comparison, we then apply the train-test-validate scheme as before.

265

266 Figure 4 shows the cross-population Spearman correlations between predicted and simulated  
267 phenotypes in our simulated AA, CEU, and YRI, partitioned by proportion of shared eQTL, for  $k =$   
268 10 causal eQTL, with  $k = 5$  and  $k = 20$  showing similar effects (Supplementary Figure 5 and  
269 Supplementary Figure 6). Imputation within a population produced similar correlations in all  
270 cases, ranging from 0.323 to 0.329 (Supplementary Table 4). Secondly, the case of 100% shared  
271 eQTL architecture (where eQTL positions and effects are exactly the same across populations)  
272 models provide predictions with no loss in cross-population portability, with correlations ranging  
273 from 0.299 to 0.336 even when imputing across populations (Supplementary Table 5). This case  
274 suggests that eQTL that are causal in all populations can impute gene expression reliably  
275 regardless of the population in which they were ascertained, provided that the eQTL can be  
276 correctly mapped and genotyped in all populations, that the eQTL effects are identical across  
277 populations, and that a linear model of eQTL is assumed. For cases where eQTL architecture is  
278 not fully shared across populations, we see that imputation from each population into the other  
279 improves as the proportion of shared eQTL increases (Figure 4). The cross-population correlation  
280 between predicted gene expression versus measurement is highest between YRI and AA (0.037  
281 to 0.088), intermediate between CEU and AA (0.0083 to 0.0245), and lowest between CEU and  
282 YRI (0.0017 to 0.0290). When imputing between two populations, the choice of which population  
283 is used to train predictive models produces no obvious difference in imputation quality. More  
284 explicitly, imputation quality between AA to CEU and CEU to AA is not significantly different ( $p$ -  
285 value  $\sim 1$ , Dunn test). The same applies between AA to YRI and YRI to AA ( $p$ -value  $\sim 1$ ) and  
286 between CEU to YRI and YRI to CEU ( $p$ -value  $< 0.12$ ). All other train/test scenarios are significantly  
287 different from each other as expected (Supplementary Table 7). The results for  $k = 5, 10$ , and 20  
288 eQTL are consistent with the higher overall ancestral similarity of AA to YRI versus AA to CEU or  
289 CEU to YRI ( $k = 10$ , Figure 4, similar plots in Supplementary Figure 5 and Supplementary Figure  
290 6). Although less realistic for most genes<sup>5,6,18</sup>, we also analyzed models with a single causal eQTL.  
291 Trends for single-eQTL models are difficult to analyze due to simplicity in architecture  
292 (Supplementary Figure 7) and binary inference as to whether the causal is identified or not.

293

294 Overall, these results highlight two points: firstly, since prediction within populations is better  
295 than prediction between populations, our results reaffirm prior investigations<sup>38</sup> that population  
296 matching matters for optimally imputing gene expression. This is consistent with our results of  
297 impaired transcriptome imputation performance in SAGE with currently available resources.  
298 Secondly, despite decreased prediction accuracy when imputing between different populations,  
299 the populations that are more closely genetically related demonstrate better cross-population  
300 prediction. Imputation results from both GTEx and DGN into SAGE suggest that current predictive  
301 models, even for genes with greater heritability, perform worse than expected despite matching  
302 tissue types. Focusing on imputation  $R^2$ , as previous studies have done, may hide the observation  
303 of a substantial proportion of negative correlations between predictions and gene expression  
304 measurements in cross-population scenarios. Our investigation into cross-population imputation  
305 accuracy with GEUVADIS data replicates this lack of cross-population portability as observed with  
306 current GTEx and DGN predictive models. Since transcriptome prediction models use  
307 multivariate genotype predictors trained on a specific outcome, the impaired cross-population

308 application can be viewed as an analogous observation to that seen previously in polygenic  
309 scores.<sup>35</sup>

310

311 It is important to note that our observations do not reflect shortcomings of either the initial  
312 PrediXcan or TWAS frameworks. Nor do our findings affect the positive discoveries made using  
313 these frameworks over the past several years. These methods fully rely on the data used as input  
314 for training, and the most commonly used datasets for model training are overwhelmingly of  
315 European descent. Here we note that the current models fail to capture the complexity of the  
316 cross-population genomic architecture of gene expression for populations of non-European  
317 descent. Failing to account for this could lead researchers to draw incorrect conclusions from  
318 their genetic data, particularly as these models would lead to false negatives.

319

320 To this end, our simulations strongly suggest that imputing gene expression in a target population  
321 is improved by using predictive models constructed in a genetically similar training population. If  
322 populations share the exact same eQTL architecture, then they are essentially interchangeable  
323 for the purposes of gene expression imputation so long as eQTL are genotyped and accurately  
324 estimated, which remains a technological and statistical challenge. As the proportion of shared  
325 eQTL architecture decreases between two populations, the cross-population imputation quality  
326 decreases as well, and often dramatically. In both SAGE and GEUVADIS, we observe cross-  
327 population patterns consistent with an imperfect overlap of eQTL across populations. Ensuring  
328 representative eQTL architecture for all populations in genotype-expression repositories will  
329 require a solid understanding of true cross-population and population-specific eQTL. However,  
330 expanding the amount of global genetic architecture represented in genotype-expression  
331 repositories, which can be accomplished by sampling more populations, provides the most  
332 desirable course for improving gene expression prediction models. Additionally, this presents an  
333 opportunity for future research in methods that could improve cross-population portability,  
334 particularly when one population is over-represented in reference data. Tools from transfer  
335 learning could facilitate porting TWAS eQTL models from reference populations to target  
336 populations using little or no RNA-Seq data.

337

338 In light of the surging interest in gene expression imputation, we see a pressing need for freely  
339 distributed predictive models of gene expression estimated from coupled transcriptome-genome  
340 data sampled in a variety of populations and tissues. The recently published predictive models  
341 with multi-ethnic MESA data constitute a crucial first step in this direction for researchers  
342 working with admixed populations. However, the clinical and biomedical research communities  
343 must push for more diverse genotype-expression resources to ensure that the fruits of genomic  
344 studies benefit all populations.



345 **Online Resources**

346

347 PredictDB: <http://predictdb.org/>

348 GTEx: <http://gtexportal.org/>

349 DGN: <http://dags.stanford.edu/dgn/>

350 GEUVADIS: <https://www.ebi.ac.uk/Tools/geuvadis-das/>

351 Source code: <https://github.com/asthmacollaboratory/sage-geuvadis-predixcan>

352 Results and simulation data: <https://ucsf.box.com/v/sage-geuvadis-predixcan>

353

## 354 SUPPLEMENTAL MATERIALS AND METHODS

### 355 Genotype and RNA-Seq data

356 RNA-Seq (RNA sequencing) data generation and cleaning protocols for 39 SAGE subjects analyzed  
357 here were initially described in (Mak, White, Eckalbar, et al. 2018).<sup>39</sup> Genotypes were generated  
358 on the Affymetrix Axiom array as described previously.<sup>53</sup> Genotypes were then imputed on the  
359 Michigan Imputation Server<sup>54</sup> with EAGLE v2.3<sup>55</sup> and the 1000 Genomes panel phase 3 v5<sup>56</sup> and  
360 then subjected to the following filters: <5% missing sample, <5% missing genotypes, >1% MAF,  
361 >1e-4 HWE, and >0.3 imputation R<sup>2</sup>. The choice of the 1000 Genomes panel follows GTEx  
362 protocol, though GTEx used the smaller 1000 Genomes phase 1 panel.<sup>4</sup> Gene expression counts  
363 were processed through the GTEx v6p eQTL quality control pipeline and as described  
364 previously.<sup>18</sup> This filtering process kept 20,985 genes with Ensembl identifiers for analysis, of  
365 which 20,268 were autosomal genes. We then quantile normalized the remaining gene  
366 expression values across samples as our gene expression measurements.

367  
368 GEUVADIS genotype VCF files and normalized gene expression data (filename  
369 GD462.GeneQuantRPKM.50FN.sampleName.resk10.txt.gz) were downloaded directly from  
370 the EMBL-EBI GEUVADIS Data Browser. Genotypes were filtered similarly to SAGE subjects. No  
371 manipulation was performed on expression data. This process yielded 23,722 genes for analysis.

372

### 373 Running PrediXcan models

374

375 We ran PrediXcan on SAGE subjects using PredictDB prediction weights from three paired  
376 genotype-expression datasets from PredictDB: GTEx, DGN, and MESA.<sup>6,9,38,57</sup> For GTEx, we used  
377 both GTEx v6p and GTEx v7 weights. For MESA, we used all weight sets from the freeze dated  
378 2018-05-30: African Americans (MESA\_AFA), African Americans and Hispanics (MESA\_AFHI),  
379 Caucasians (MESA\_CAU), and all MESA samples (MESA\_ALL). Overall, the analysis included  
380 10,161 genes, of which only 273 had *both* normalized RNA-Seq measures and predictions from  
381 *all* weight sets. Of these, 126 had positive correlation between prediction and measurement. We  
382 assessed imputation quality by comparing PrediXcan predictions to normalized gene expression  
383 from SAGE using linear regression and correlation tests.

384

### 385 Building prediction models

386

387 We trained prediction models in GEUVADIS on genotypes in a 500Kb window around each of  
388 23,723 genes with measured and normalized gene expression. GEUVADIS subjects were  
389 partitioned into various groups: the Europeans (EUR373), the non-Finnish Europeans (EUR278),  
390 the Yoruba (AFR), and the constituent 1000 Genomes populations (CEU, GBR, TSI, FIN, and YRI).  
391 For each training set, we performed nested cross-validation. The external cross-validation for all  
392 populations used leave-one-out cross-validation (LOOCV). The internal cross-validation used 10-  
393 fold cross-validation for EUR373 and EUR278 and LOOCV for the five constituent GEUVADIS  
394 populations in order to fully utilize the smaller sample size ( $n = 89$ ) compared to EUR278 ( $n = 278$ )

395 and EUR373 ( $n = 373$ ). Internal cross-validation used elastic net regression with mixing parameter  
396  $\alpha = 0.5$  as implemented in the `glmnet` package in R. The nonzero weights for each SNP from each  
397 LOOCV were compiled and averaged for each gene, yielding a single set of prediction weights for  
398 each gene. Predictions were computed by parsing genotype dosages from the target population  
399 corresponding to the nonzero SNP predictors, and then multiplying dosages against the  
400 prediction weights. The resulting predictions were compared to normalized gene expression  
401 measurements downloaded from the GEUVADIS data portal. The comparison of predictive  
402 models cannot easily differentiate predictions of 0 (no gene expression) and NA (missing  
403 expression). We addressed this with two additional filters. Firstly, we removed genes that did not  
404 have any eQTL in their predictive models. Secondly, genes where fewer than half of the  
405 individuals had nonmissing predictions were removed from further analysis. Coefficients of  
406 determination ( $R^2$ ) were computed with the `lm` function in R. Spearman correlations were  
407 computed with the `cor.test` function in R.

408  
409

## 410 **Simulation**

411

412 We downloaded a sample of 20,085 HapMap 3 SNPs<sup>50</sup> from each of CEU and YRI on chromosome  
413 22 as provided by HAPGEN2.<sup>51</sup> The data include 234 phased haplotypes for CEU and 230 phased  
414 haplotypes for YRI. We forward-simulated from these haplotypes to obtain two populations of  $n$   
415 = 1000 individuals each. We then sampled haplotypes in proportions of 80% YRI and 20% CEU to  
416 obtain a mixture of CEU and YRI where the ancestry patterns roughly mimic those of African  
417 Americans. For computational simplicity, and in keeping with the high ancestry LD present in  
418 African Americans<sup>58,59</sup>, for each gene we assumed local ancestry was constant for each haplotype.  
419 For each of the three simulated populations, we applied the same train-test-validate scheme  
420 used for cross-population analysis in GEUVADIS. Genetic data for model simulation were  
421 downloaded from Ensembl 89 and included the largest 100 genes from chromosome 22. We  
422 defined each gene as the start and end positions corresponding to the canonical transcript, plus  
423 1 megabase in each direction. Two genes, PPP6R2 and MOV10L1, spanned no polymorphic  
424 markers in our simulated data, resulting in 98 gene models used for analysis. To simulate  
425 predictive eQTL models, we tested multiple parameter configurations for each gene: we varied  
426 the number of causal eQTL ( $k = 1, 5, 10, \text{ and } 20$ ), the positions (all same or not all same), and the  
427 proportion of shared positions ( $p = 0.0, 0.1, 0.2, \dots, 0.9, 1$ ). Each model included a simulated gene  
428 expression phenotype with *cis*-heritability set to 0.15. For each parameter configuration, we ran  
429 100 different random instantiations of the model simulations.

430

## 431 **Analysis tools**

432

433 Analyses used GNU parallel<sup>60</sup> and the tidyverse bundle of R packages.<sup>61</sup> All plots were  
434 generated with ggplot2.<sup>62</sup>

435

## 436 **CONFLICT OF INTEREST**

437

438 C.R.G. owns stock in 23andMe, Inc. The remaining authors declare no potential conflicts of  
439 interest.

#### 440 **ACKNOWLEDGEMENTS**

441 This work was supported in part by the Sandler Family Foundation, the American Asthma  
442 Foundation, the RWJF Amos Medical Faculty Development Program, the Harry Wm. and Diana V.  
443 Hind Distinguished Professor in Pharmaceutical Sciences II, the National Heart, Lung, and Blood  
444 Institute (NHLBI) R01HL117004, R01HL128439, R01HL135156, X01HL134589, R01HL141992, and  
445 R01HL104608, the National Human Genome Research Institute (NHGRI) U01HG007419, National  
446 Institute of Environmental Health Sciences R01ES015794, R21ES24844, the National Institute on  
447 Minority Health and Health Disparities P60MD006902, R01MD010443, RL5GM118984 and the  
448 Tobacco-Related Disease Research Program under Award Numbers 24RT-0025, 27IR-0030.  
449 Research reported in this article was funded by the National Institutes of Health Common Fund  
450 and Office of Scientific Workforce Diversity under three linked awards RL5GM118984,  
451 TL4GM118986, 1UL1GM118985 administered by the National Institute of General Medical  
452 Sciences.

453  
454 The authors wish to acknowledge the following SAGE co-investigators for subject recruitment,  
455 sample processing and quality control: Luisa N. Borrell, DDS, PhD, Emerita Brigino-Buenaventura,  
456 MD, Adam Davis, MA, MPH, Michael A. LeNoir, MD, Kelley Meade, MD, Fred Lurmann, MS and  
457 Harold J. Farber, MD, MSPH. The authors also wish to thank the staff and participants who  
458 contributed to the SAGE study.

459  
460 K.L.K was additionally supported by a diversity supplement of NHLBI R01HL135156, the UCSF  
461 Bakar Computational Health Sciences Institute, the Gordon and Betty More Foundation grant  
462 GBMF3834, and the Alfred P. Sloan Foundation grant 2013-10-27 to UC Berkeley through the  
463 Moore-Sloan Data Sciences Environment initiative at the Berkeley Institute for Data Science  
464 (BIDS). The logistical space, technical support, administrative assistance, and indefatigable good  
465 humor of the members and staff at BIDS is gratefully acknowledged.

466  
467 M.J.W. was additionally supported by a diversity supplement of NHLBI R01HL117004, an  
468 Institutional Research and Academic Career Development Award K12GM081266, and an NHLBI  
469 Research Career Development (K) Award K01HL140218.

470  
471 M.G.C was additionally supported by NIH MARC U-STAR grant T34GM008574 at San Francisco  
472 State University.

473  
474 C.R.G. was additionally supported by grants R56HG010297 and T32HG00044.  
475

476 **REFERENCES**

477

478 1. Sudlow, C., Gallacher, J., Allen, N., Beral, V., Burton, P., Danesh, J., Downey, P., Elliott, P.,  
479 Green, J., Landray, M., et al. (2015). UK Biobank: An Open Access Resource for Identifying the  
480 Causes of a Wide Range of Complex Diseases of Middle and Old Age. *PLoS Med.* *12*,

481 2. NHLBI Trans-Omics for Precision Medicine.

482 3. NHGRI Genome Sequencing Program (GSP).

483 4. The 1000 Genomes Consortium An integrated map of genetic variation from 1,092 human  
484 genomes | *Nature*.

485 5. Gusev, A., Ko, A., Shi, H., Bhatia, G., Chung, W., Penninx, B.W.J.H., Jansen, R., de Geus,  
486 E.J., Boomsma, D.I., Wright, F.A., et al. (2016). Integrative approaches for large-scale  
487 transcriptome-wide association studies. *Nat. Genet.* *48*, 245–252.

488 6. Gamazon, E.R., Wheeler, H.E., Shah, K.P., Mozaffari, S.V., Aquino-Michaels, K., Carroll,  
489 R.J., Eyler, A.E., Denny, J.C., Nicolae, D.L., Cox, N.J., et al. (2015). A gene-based association  
490 method for mapping traits using reference transcriptome data. *Nat. Genet.* *47*, 1091–1098.

491 7. GTEx Consortium (2013). The Genotype-Tissue Expression (GTEx) project. *Nat. Genet.* *45*,  
492 580–585.

493 8. Battle, A., Mostafavi, S., Zhu, X., Potash, J.B., Weissman, M.M., McCormick, C.,  
494 Haudenschild, C.D., Beckman, K.B., Shi, J., Mei, R., et al. (2014). Characterizing the genetic  
495 basis of transcriptome diversity through RNA-sequencing of 922 individuals. *Genome Res.* *24*,  
496 14–24.

497 9. Barbeira, A.N., Dickinson, S.P., Torres, J.M., Bonazzola, R., Zheng, J., Torstenson, E.S.,  
498 Wheeler, H.E., Shah, K.P., Edwards, T., Garcia, T., et al. (2017). Exploring the phenotypic  
499 consequences of tissue specific gene expression variation inferred from GWAS summary  
500 statistics. *BioRxiv* 045260.

501 10. Zhu, Z., Zhang, F., Hu, H., Bakshi, A., Robinson, M.R., Powell, J.E., Montgomery, G.W.,  
502 Goddard, M.E., Wray, N.R., Visscher, P.M., et al. (2016). Integration of summary data from  
503 GWAS and eQTL studies predicts complex trait gene targets. *Nat. Genet.* *48*, 481–487.

504 11. Mostafavi, S., Gaiteri, C., Sullivan, S.E., White, C.C., Tasaki, S., Xu, J., Taga, M., Klein, H.-  
505 U., Patrick, E., Komashko, V., et al. (2018). A molecular network of the aging human brain  
506 provides insights into the pathology and cognitive decline of Alzheimer's disease. *Nat. Neurosci.*  
507 *21*, 811.

508 12. Ferreira, M.A.R., Jansen, R., Willemsen, G., Penninx, B., Bain, L.M., Vicente, C.T., Revez,  
509 J.A., Matheson, M.C., Hui, J., Tung, J.Y., et al. (2017). Gene-based analysis of regulatory  
510 variants identifies four putative novel asthma risk genes related to nucleotide synthesis and  
511 signaling. *J. Allergy Clin. Immunol.* *139*, 1148–1157.



- 512 13. Lamontagne, M., Bérubé, J.-C., Obeidat, M., Cho, M.H., Hobbs, B.D., Sakornsakolpat, P., de  
513 Jong, K., Boezen, H.M., International COPD Genetics Consortium, Nickle, D., et al. (2018).  
514 Leveraging lung tissue transcriptome to uncover candidate causal genes in COPD genetic  
515 associations. *Hum. Mol. Genet.* 27, 1819–1829.
- 516 14. Thériault, S., Gaudreault, N., Lamontagne, M., Rosa, M., Boulanger, M.-C., Messika-  
517 Zeitoun, D., Clavel, M.-A., Capoulade, R., Dagenais, F., Pibarot, P., et al. (2018). A  
518 transcriptome-wide association study identifies PALMD as a susceptibility gene for calcific  
519 aortic valve stenosis. *Nat. Commun.* 9, 988.
- 520 15. Porcu, E., Rüeger, S., Consortium, eQTLGen, Santoni, F.A., Reymond, A., and Kutalik, Z.  
521 (2018). Mendelian Randomization integrating GWAS and eQTL data reveals genetic  
522 determinants of complex and clinical traits. *BioRxiv* 377267.
- 523 16. Gusev, A., Lawrenson, K., Segato, F., Fonseca, M., Kar, S., Lee, J., Pejovic, T., Consortium,  
524 O.C.A., Karlan, B., Freedman, M., et al. (2018). Multi-Tissue Transcriptome-Wide Association  
525 Studies Identify 21 Novel Candidate Susceptibility Genes for High Grade Serous Epithelial  
526 Ovarian Cancer. *BioRxiv* 330613.
- 527 17. Huckins, L.M., Dobbyn, A., Ruderfer, D., Hoffman, G., Wang, W., Pardini, A.F.,  
528 Rajagopal, V.M., Als, T.D., Hoang, H.T., Girdhar, K., et al. (2017). Gene expression imputation  
529 across multiple brain regions reveals schizophrenia risk throughout development. *BioRxiv*  
530 222596.
- 531 18. GTEx Consortium (2017). Genetic effects on gene expression across human tissues. *Nature*  
532 550, 204–213.
- 533 19. Bustamante, C.D., Burchard, E.G., and De la Vega, F.M. (2011). Genomics for the world.  
534 *Nature* 475, 163–165.
- 535 20. Popejoy, A.B., and Fullerton, S.M. (2016). Genomics is failing on diversity. *Nature* 538,  
536 161–164.
- 537 21. Bentley, A.R., Callier, S., and Rotimi, C.N. (2017). Diversity and inclusion in genomic  
538 research: why the uneven progress? *J. Community Genet.* 8, 255–266.
- 539 22. Hindorff, L.A., Bonham, V.L., Brody, L.C., Ginoza, M.E.C., Hutter, C.M., Manolio, T.A.,  
540 and Green, E.D. (2018). Prioritizing diversity in human genomics research. *Nat. Rev. Genet.* 19,  
541 175–185.
- 542 23. Asimit, J.L., Hatzikotoulas, K., McCarthy, M., Morris, A.P., and Zeggini, E. (2016). Trans-  
543 ethnic study design approaches for fine-mapping. *Eur. J. Hum. Genet.* 24, 1330–1336.
- 544 24. Wang, X., Cheng, C.-Y., Liao, J., Sim, X., Liu, J., Chia, K.-S., Tai, E.-S., Little, P., Khor, C.-  
545 C., Aung, T., et al. (2016). Evaluation of transethnic fine mapping with population-specific and  
546 cosmopolitan imputation reference panels in diverse Asian populations. *Eur. J. Hum. Genet.* 24,  
547 592–599.

- 548 25. Li, Y.R., and Keating, B.J. (2014). Trans-ethnic genome-wide association studies:  
549 advantages and challenges of mapping in diverse populations. *Genome Med.* 6, 91.
- 550 26. Kumar, R., Seibold, M.A., Aldrich, M.C., Williams, L.K., Reiner, A.P., Colangelo, L.,  
551 Galanter, J., Gignoux, C., Hu, D., Sen, S., et al. (2010). Genetic ancestry in lung-function  
552 predictions. *N. Engl. J. Med.* 363, 321–330.
- 553 27. Yang, J.J., Cheng, C., Devidas, M., Cao, X., Fan, Y., Campana, D., Yang, W., Neale, G.,  
554 Cox, N.J., Scheet, P., et al. (2011). Ancestry and pharmacogenomics of relapse in acute  
555 lymphoblastic leukemia. *Nat. Genet.* 43, 237–241.
- 556 28. Acuña-Alonzo, V., Flores-Dorantes, T., Kruit, J.K., Villarreal-Molina, T., Arellano-Campos,  
557 O., Hünemeier, T., Moreno-Estrada, A., Ortiz-López, M.G., Villamil-Ramírez, H., León-Mimila,  
558 P., et al. (2010). A functional ABCA1 gene variant is associated with low HDL-cholesterol  
559 levels and shows evidence of positive selection in Native Americans. *Hum. Mol. Genet.* 19,  
560 2877–2885.
- 561 29. Adeyemo, A., and Rotimi, C. (2010). Genetic variants associated with complex human  
562 diseases show wide variation across multiple populations. *Public Health Genomics* 13, 72–79.
- 563 30. Manrai, A.K., Funke, B.H., Rehm, H.L., Olesen, M.S., Maron, B.A., Szolovits, P.,  
564 Margulies, D.M., Loscalzo, J., and Kohane, I.S. (2016). Genetic Misdiagnoses and the Potential  
565 for Health Disparities. *N. Engl. J. Med.* 375, 655–665.
- 566 31. Petrovski, S., and Goldstein, D.B. (2016). Unequal representation of genetic variation across  
567 ancestry groups creates healthcare inequality in the application of precision medicine. *Genome*  
568 *Biol.* 17, 157.
- 569 32. Oh, S.S., White, M.J., Gignoux, C.R., and Burchard, E.G. (2016). Making Precision  
570 Medicine Socially Precise. Take a Deep Breath. *Am. J. Respir. Crit. Care Med.* 193, 348–350.
- 571 33. Oh, S.S., Galanter, J., Thakur, N., Pino-Yanes, M., Barcelo, N.E., White, M.J., de Bruin,  
572 D.M., Greenblatt, R.M., Bibbins-Domingo, K., Wu, A.H.B., et al. (2015). Diversity in Clinical  
573 and Biomedical Research: A Promise Yet to Be Fulfilled. *PLoS Med.* 12,.
- 574 34. Belbin, G.M., Nieves-Colón, M.A., Kenny, E.E., Moreno-Estrada, A., and Gignoux, C.R.  
575 (2018). Genetic diversity in populations across Latin America: implications for population and  
576 medical genetic studies. *Curr. Opin. Genet. Dev.* 53, 98–104.
- 577 35. Martin, A.R., Gignoux, C.R., Walters, R.K., Wojcik, G.L., Neale, B.M., Gravel, S., Daly,  
578 M.J., Bustamante, C.D., and Kenny, E.E. (2017). Human Demographic History Impacts Genetic  
579 Risk Prediction across Diverse Populations. *Am. J. Hum. Genet.* 100, 635–649.
- 580 36. Bild, D.E., Bluemke, D.A., Burke, G.L., Detrano, R., Diez Roux, A.V., Folsom, A.R.,  
581 Greenland, P., Jacob, D.R., Kronmal, R., Liu, K., et al. (2002). Multi-Ethnic Study of  
582 Atherosclerosis: objectives and design. *Am. J. Epidemiol.* 156, 871–881.

- 583 37. Liu, Y., Ding, J., Reynolds, L.M., Lohman, K., Register, T.C., De La Fuente, A., Howard,  
584 T.D., Hawkins, G.A., Cui, W., Morris, J., et al. (2013). Methylomics of gene expression in  
585 human monocytes. *Hum. Mol. Genet.* 22, 5065–5074.
- 586 38. Mogil, L.S., Andaleon, A., Badalamenti, A., Dickinson, S.P., Guo, X., Rotter, J.I., Johnson,  
587 W.C., Im, H.K., Liu, Y., and Wheeler, H.E. (2018). Genetic architecture of gene expression traits  
588 across diverse populations. *PLOS Genet.* 14, e1007586.
- 589 39. Mak, A.C.Y., White, M.J., Eckalbar, W.L., Szpiech, Z.A., Oh, S.S., Pino-Yanes, M., Hu, D.,  
590 Goddard, P., Huntsman, S., Galanter, J., et al. (2018). Whole-Genome Sequencing of  
591 Pharmacogenetic Drug Response in Racially Diverse Children with Asthma. *Am. J. Respir. Crit.*  
592 *Care Med.* 197, 1552–1564.
- 593 40. Thakur, N., Oh, S.S., Nguyen, E.A., Martin, M., Roth, L.A., Galanter, J., Gignoux, C.R.,  
594 Eng, C., Davis, A., Meade, K., et al. (2013). Socioeconomic status and childhood asthma in  
595 urban minority youths. The GALA II and SAGE II studies. *Am. J. Respir. Crit. Care Med.* 188,  
596 1202–1209.
- 597 41. Borrell, L.N., Nguyen, E.A., Roth, L.A., Oh, S.S., Tcheurekdjian, H., Sen, S., Davis, A.,  
598 Farber, H.J., Avila, P.C., Brigino-Buenaventura, E., et al. (2013). Childhood Obesity and Asthma  
599 Control in the GALA II and SAGE II Studies. *Am. J. Respir. Crit. Care Med.* 187, 697–702.
- 600 42. Nishimura, K.K., Galanter, J.M., Roth, L.A., Oh, S.S., Thakur, N., Nguyen, E.A., Thyne, S.,  
601 Farber, H.J., Serebrisky, D., Kumar, R., et al. (2013). Early-life air pollution and asthma risk in  
602 minority children. The GALA II and SAGE II studies. *Am. J. Respir. Crit. Care Med.* 188, 309–  
603 318.
- 604 43. Lappalainen, T., Sammeth, M., Friedländer, M.R., 't Hoen, P.A., Monlong, J., Rivas, M.A.,  
605 González-Porta, M., Kurbatova, N., Griebel, T., Ferreira, P.G., et al. (2013). Transcriptome and  
606 genome sequencing uncovers functional variation in humans. *Nature* 501, 506–511.
- 607 44. Mikhaylova, A.V., and Thornton, T.A. (2019). Accuracy of gene expression prediction from  
608 genotype data with PrediXcan varies across diverse populations. *BioRxiv* 524728.
- 609 45. Viñuela, A., Brown, A.A., Buil, A., Tsai, P.-C., Davies, M.N., Bell, J.T., Dermitzakis, E.T.,  
610 Spector, T.D., and Small, K.S. (2018). Age-dependent changes in mean and variance of gene  
611 expression across tissues in a twin cohort. *Hum. Mol. Genet.* 27, 732–741.
- 612 46. McCall, M.N., Illei, P.B., and Halushka, M.K. (2016). Complex Sources of Variation in  
613 Tissue Expression Data: Analysis of the GTEx Lung Transcriptome. *Am. J. Hum. Genet.* 99,  
614 624–635.
- 615 47. Zhu, Y., Wang, L., Yin, Y., and Yang, E. (2017). Systematic analysis of gene expression  
616 patterns associated with postmortem interval in human tissues. *Sci. Rep.* 7, 5435.
- 617 48. Ferreira, P.G., Muñoz-Aguirre, M., Reverter, F., Godinho, C.P.S., Sousa, A., Amadoz, A.,  
618 Sodaiei, R., Hidalgo, M.R., Pervouchine, D., Carbonell-Caballero, J., et al. (2018). The effects of  
619 death and post-mortem cold ischemia on human tissue transcriptomes. *Nat. Commun.* 9, 490.

- 620 49. Martin, A.R., Karczewski, K.J., Kerminen, S., Kurki, M.I., Sarin, A.-P., Artomov, M.,  
621 Eriksson, J.G., Esko, T., Genovese, G., Havulinna, A.S., et al. (2018). Haplotype Sharing  
622 Provides Insights into Fine-Scale Population History and Disease in Finland. *Am. J. Hum. Genet.*  
623 *102*, 760–775.
- 624 50. The International HapMap 3 Consortium (2010). Integrating common and rare genetic  
625 variation in diverse human populations. *Nature* *467*, 52–58.
- 626 51. Su, Z., Marchini, J., and Donnelly, P. (2011). HAPGEN2: simulation of multiple disease  
627 SNPs. *Bioinformatics* *27*, 2304–2305.
- 628 52. Baharian, S., Barakatt, M., Gignoux, C.R., Shringarpure, S., Errington, J., Blot, W.J.,  
629 Bustamante, C.D., Kenny, E.E., Williams, S.M., Aldrich, M.C., et al. (2016). The Great  
630 Migration and African-American Genomic Diversity. *PLOS Genet.* *12*, e1006059.
- 631 53. Hoffmann, T.J., Zhan, Y., Kvale, M.N., Hesselson, S.E., Gollub, J., Iribarren, C., Lu, Y.,  
632 Mei, G., Purdy, M.M., Quesenberry, C., et al. (2011). Design and coverage of high throughput  
633 genotyping arrays optimized for individuals of East Asian, African American, and Latino  
634 race/ethnicity using imputation and a novel hybrid SNP selection algorithm. *Genomics* *98*, 422–  
635 430.
- 636 54. Das, S., Forer, L., Schönherr, S., Sidore, C., Locke, A.E., Kwong, A., Vrieze, S.I., Chew,  
637 E.Y., Levy, S., McGue, M., et al. (2016). Next-generation genotype imputation service and  
638 methods. *Nat. Genet.* *48*, 1284–1287.
- 639 55. Loh, P.-R., Danecek, P., Palamara, P.F., Fuchsberger, C., Reshef, Y.A., Finucane, H.K.,  
640 Schoenherr, S., Forer, L., McCarthy, S., Abecasis, G.R., et al. (2016). Reference-based phasing  
641 using the Haplotype Reference Consortium panel. *Nat. Genet.* *48*, 1443–1448.
- 642 56. Sudmant, P.H., Rausch, T., Gardner, E.J., Handsaker, R.E., Abyzov, A., Huddleston, J.,  
643 Zhang, Y., Ye, K., Jun, G., Hsi-Yang Fritz, M., et al. (2015). An integrated map of structural  
644 variation in 2,504 human genomes. *Nature* *526*, 75–81.
- 645 57. Wheeler, H.E., Shah, K.P., Brenner, J., Garcia, T., Aquino-Michaels, K., Consortium, Gte.,  
646 Cox, N.J., Nicolae, D.L., and Im, H.K. (2016). Survey of the Heritability and Sparse Architecture  
647 of Gene Expression Traits across Human Tissues. *PLOS Genet.* *12*, e1006423.
- 648 58. Gravel, S. (2012). Population genetics models of local ancestry. *Genetics* *191*, 607–619.
- 649 59. Price, A.L., Tandon, A., Patterson, N., Barnes, K.C., Rafaels, N., Ruczinski, I., Beaty, T.H.,  
650 Mathias, R., Reich, D., and Myers, S. (2009). Sensitive detection of chromosomal segments of  
651 distinct ancestry in admixed populations. *PLoS Genet.* *5*, e1000519.
- 652 60. Tange, O. (2018). GNU Parallel 2018 (Ole Tange).
- 653 61. Wickham, Hadley, and Gromm, Garrett (2017). R for Data Science (O’Reilly Media,  
654 Inc.).

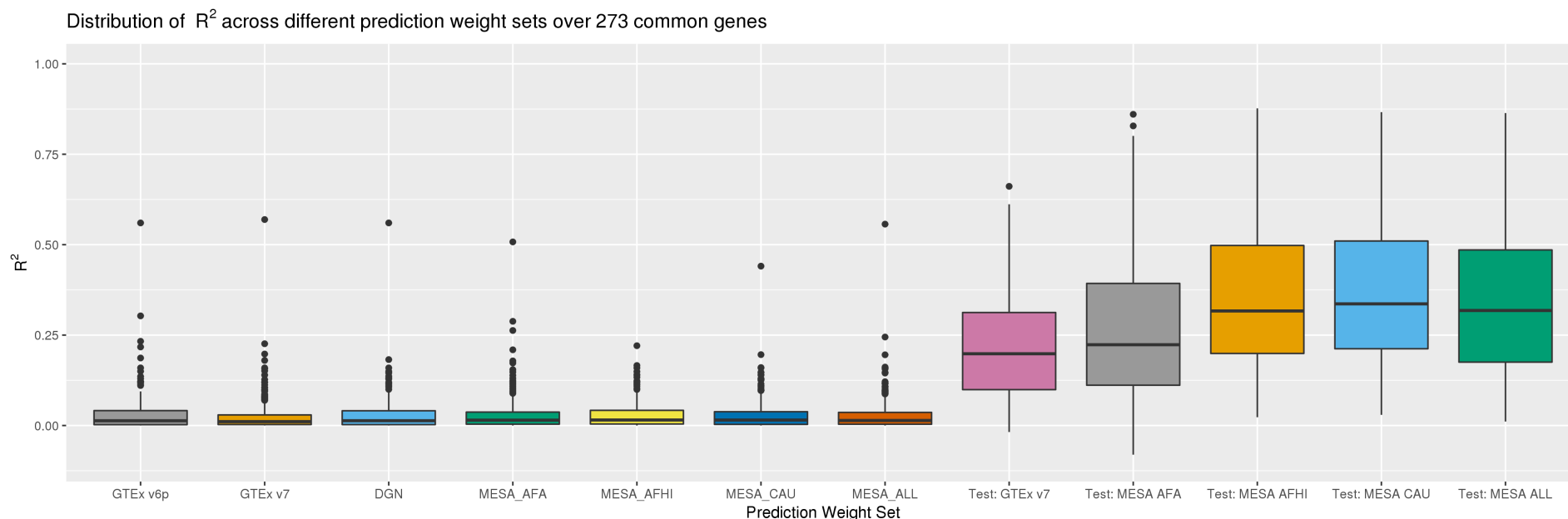
655 62. Wickham, Hadley (2016). *ggplot2: Elegant Graphics for Data Analysis* (Springer-Verlag  
656 New York).

657 63. The 1000 Genomes Project Consortium (2010). A map of human genome variation from  
658 population-scale sequencing. *Nature* 467, 1061–1073.

659



660



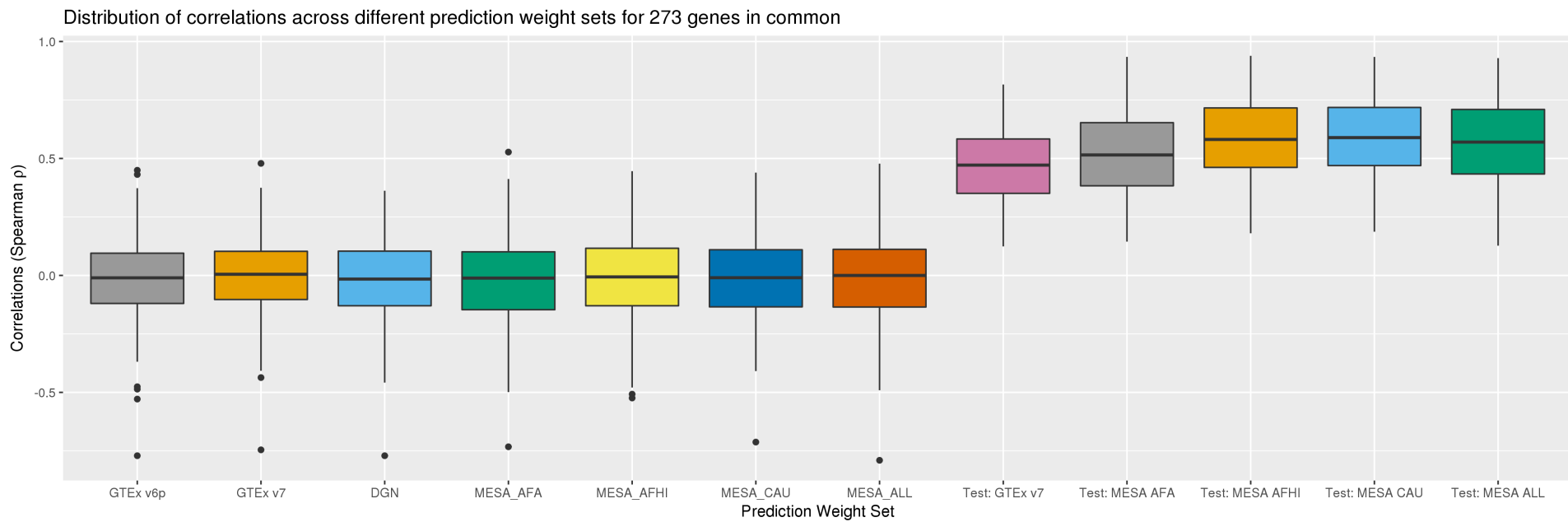
661

662 *Figure 1:  $R^2$  of measured gene expression versus predictions from PrediXcan. The prediction weights used here are, from left to right:*

663 *GTEx v6p, GTEx v7, DGN, MESA African Americans, MESA African Americans and Hispanics, MESA Caucasians, and all MESA subjects.*

664 *Test  $R^2$  from model training in GTEx 7 and MESA appear on the right and provide a performance baseline.*

665

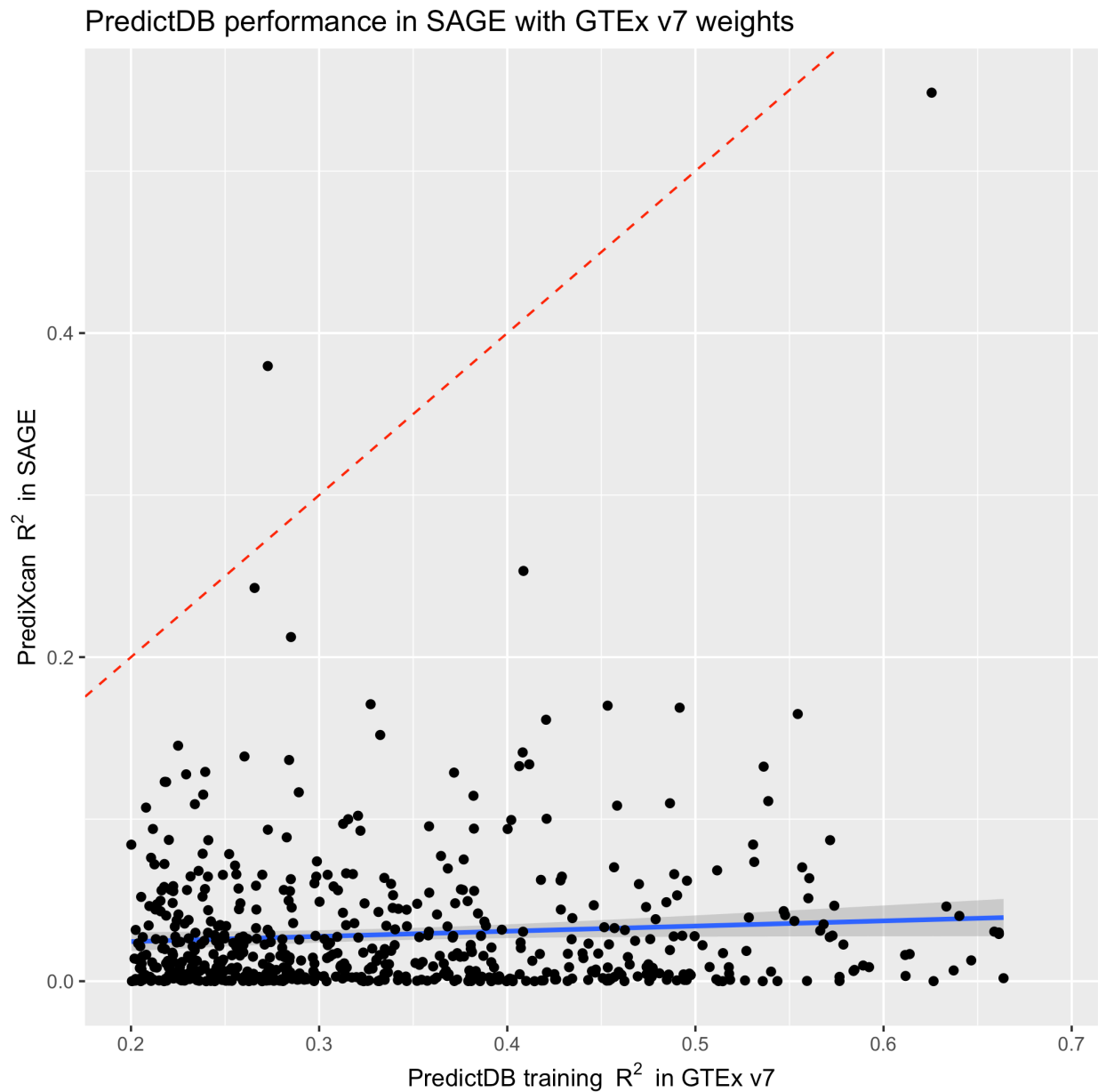


666

667 *Figure 2: Spearman correlations of measured gene expression versus predicted expression from PrediXcan. The order of the weight*  
 668 *sets matches Figure 1.*

669

670



671

672 *Figure 3: A comparison of  $R^2$  from SAGE and GTEx v7 training diagnostics. The SAGE  $R^2$  are*  
673 *computed from regressing PrediXcan predictions onto gene expression measurements. The GTEx*  
674 *v7  $R^2$  are taken from PredictDB. The red dotted line marks where  $R^2$  between the two groups*  
675 *match, while the blue line denotes the best linear fit.*

676

$R^2$		Train Pop		
		EUR373	EUR278	AFR
Test Pop	EUR373	0.098	n/a	0.029
	EUR278	n/a	0.096	0.030
	FIN	n/a	0.087	0.039
	AFR	0.054	0.051	0.079

677 *Table 1: Imputation  $R^2$  between populations in GEUVADIS for genes with positive correlation*  
 678 *between predictions and measurements. Scenarios where the training sample is contained in*  
 679 *the testing sample cannot be accurately tested and are marked with “n/a”. EUR373 includes all*  
 680 *Europeans, EUR278 includes only non-Finnish Europeans, FIN includes only the Finnish, and AFR*  
 681 *includes only the Yoruba.*

682

$R^2$		Train Pop		
		EUR373	EUR278	AFR
Test Pop	EUR373	0.201	n/a	0.096
	EUR278	n/a	0.183	0.095
	FIN	n/a	0.216	0.111
	AFR	0.147	0.141	0.130

683 *Table 2: Imputation  $R^2$  between populations in GEUVADIS for 564 gene models that show*  
 684 *positive correlation between prediction and measurement in all 9 train-test scenarios that were*  
 685 *analyzed. Scenarios that were not tested are marked with “n/a”. As before, EUR373 includes all*  
 686 *Europeans, EUR278 includes only non-Finnish Europeans, FIN includes only the Finnish, and AFR*  
 687 *includes only the Yoruba.*

R2 Mean (Std Err)		Training population				
		CEU	TSI	GBR	FIN	YRI
Testing Pop	CEU	0.115 (0.139)	0.106 (0.139)	0.107 (0.134)	0.103 (0.133)	0.069 (0.116)
	TSI	0.124 (0.158)	0.121 (0.151)	0.124 (0.149)	0.118 (0.145)	0.083 (0.13)
	GBR	0.132 (0.16)	0.137 (0.155)	0.136 (0.156)	0.133 (0.155)	0.087 (0.132)
	FIN	0.128 (0.158)	0.130 (0.155)	0.130 (0.153)	0.130 (0.152)	0.084 (0.134)
	YRI	0.065 (0.108)	0.069 (0.112)	0.063 (0.1)	0.062 (0.102)	0.104 (0.138)

688 *Table 3: Cross-population prediction performance across all five constituent GEUVADIS*  
689 *populations over genes with positive correlation between predictions and measurements. All*  
690 *populations were subsampled to N = 89 individuals. The number of genes represented varies by*  
691 *training sample (CEU: N = 1029, FIN: N = 1320, GBR: 1436, TSI: 1250, YRI: 914).*

692

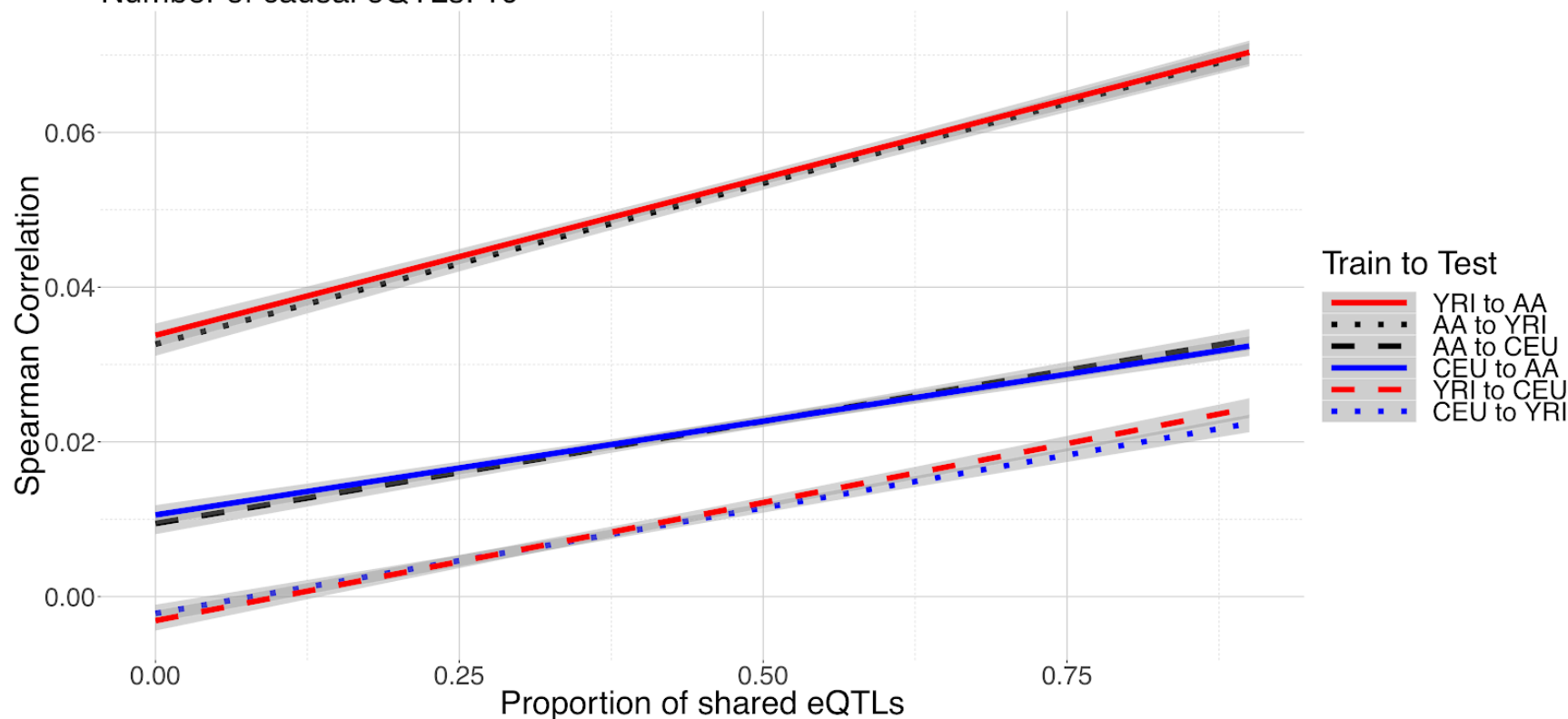


R2 Mean (Std Err)		Training population				
		CEU	TSI	GBR	FIN	YRI
Testing Pop	CEU	0.239 (0.18)	0.269 (0.177)	0.291 (0.166)	0.297 (0.168)	0.201 (0.164)
	TSI	0.307 (0.188)	0.294 (0.21)	0.331 (0.182)	0.322 (0.185)	0.227 (0.185)
	GBR	0.320 (0.175)	0.326 (0.181)	0.318 (0.191)	0.350 (0.178)	0.235 (0.183)
	FIN	0.318 (0.191)	0.320 (0.198)	0.343 (0.182)	0.323 (0.201)	0.244 (0.192)
	YRI	0.166 (0.164)	0.205 (0.163)	0.195 (0.157)	0.189 (0.156)	0.213 (0.177)

693 *Table 4: Cross-population prediction performance across all five subsampled GEUVADIS*  
694 *populations over the 142 genes with positive correlation between prediction and measurement*  
695 *in all 25 train-test scenarios.*

# Crosspopulation correlations of predicted versus simulated gene expression

Number of causal eQTLs: 10



696  
697 *Figure 4: Correlations between predictions and simulated gene expression measurements from simulated populations across various*  
698 *proportions of shared eQTL architecture with 10 causal cis-eQTL. Here YRI is simulated from the 1000 Genomes Yoruba, CEU is*  
699 *simulated from the Utahns, and AA is constructed from YRI and CEU. Each trend line represents a linear interpolation of correlation*  
700 *versus shared eQTL proportion. Gray areas denote 95% confidence regions of LOESS-smoothed mean correlations conditional on the*  
701 *proportion of shared eQTL.*

702

<b>Weight set</b>	<b>Gene models</b>	<b>Genes predicted in SAGE</b>	<b>Genes both predicted and measured</b>	<b>Genes with positively correlated predictions and measurements</b>	<b>Mean Correlation (273 common genes)</b>
GTEEx v6p	6588	5773	5348	2730	-0.0044
GTEEx v7	6297	2742	2570	1319	-0.0113
DGN	13171	4033	3678	1819	-0.0124
MESA_AFA	3551	995	982	497	-0.0204
MESA_AFHI	5556	1889	1862	969	-0.0049
MESA_CAU	4674	1654	1633	837	-0.0082
MESA_ALL	6217	2443	2408	1201	-0.0107

*Supplementary Table 1: Summary statistics for analyzing gene expression prediction in SAGE for all seven weight sets in PredictDB. SAGE has measurements for 20,985 genes, of which 20,268 are autosomal. The intersection of genes with both predictions and measurements in SAGE across all seven weight sets is 273, of which 39 produce predictions positively correlated to data in all comparisons.*

703  
704

<b>Pop</b>	<b>Measured genes</b>	<b>Predictive Models</b>	<b>With &gt;50% samples predicted</b>	<b>Analyzed prediction v. measurement</b>	<b>Positive correlation</b>
EUR373	23723	20418	11917	11914	5586
EUR278	23723	20182	11043	11043	4817
YRI89	23723	20699	11180	11179	4867

705

*Supplementary Table 2: Summary statistics for each filtering step in the analysis of gene expression models from GEUVADIS for the 3 training populations EUR373, EUR278, and AFR. The analysis of prediction vs. measurement contains 5038 genes in common between all three populations. Of these genes, 1476 genes demonstrate positive correlation between predictions and measurements.*

Training Pop	Testing Pop	R <sup>2</sup>	Correlation	Transcripts
AFR	AFR	0.079	0.2329	2562
AFR	EUR278	0.030	0.1122	2996
AFR	EUR373	0.029	0.1072	3043
AFR	FIN	0.039	0.1377	2908
EUR278	AFR	0.051	0.1632	3079
EUR278	EUR278	0.096	0.2291	2857
EUR278	FIN	0.087	0.2171	3994
EUR373	AFR	0.054	0.1683	3105
EUR373	EUR373	0.098	0.2325	3132

706 *Supplementary Table 3: Summary statistics from training and testing results with continental*  
707 *GEUVADIS populations for gene models with positive correlations. The R<sup>2</sup> correspond to Table 1.*  
708 *The column “Correlation” lists the Spearman correlations for each scenario, while “Transcripts”*  
709 *gives the number of gene models used to compute the R<sup>2</sup> and correlation summaries.*

710

<b>Training Pop</b>	<b>Testing Pop</b>	<b>Shared eQTL Proportion</b>	<b>Correlation (Mean)</b>	<b>Correlation (StdErr)</b>
AA	AA	0	0.323	0.0350
AA	AA	0.1	0.323	0.0357
AA	AA	0.2	0.323	0.0356
AA	AA	0.3	0.324	0.0355
AA	AA	0.4	0.323	0.0361
AA	AA	0.5	0.323	0.0355
AA	AA	0.6	0.323	0.0358
AA	AA	0.7	0.321	0.0364
AA	AA	0.8	0.323	0.0361
AA	AA	0.9	0.322	0.0358
CEU	CEU	0	0.329	0.0345
CEU	CEU	0.1	0.329	0.0345
CEU	CEU	0.2	0.329	0.0345
CEU	CEU	0.3	0.329	0.0345
CEU	CEU	0.4	0.329	0.0345
CEU	CEU	0.5	0.329	0.0345
CEU	CEU	0.6	0.329	0.0345
CEU	CEU	0.7	0.329	0.0346
CEU	CEU	0.8	0.329	0.0345
CEU	CEU	0.9	0.329	0.0345
YRI	YRI	0	0.325	0.0354
YRI	YRI	0.1	0.325	0.0355
YRI	YRI	0.2	0.324	0.0351
YRI	YRI	0.3	0.324	0.0354
YRI	YRI	0.4	0.325	0.0354
YRI	YRI	0.5	0.325	0.0352
YRI	YRI	0.6	0.324	0.0351
YRI	YRI	0.7	0.322	0.0354
YRI	YRI	0.8	0.324	0.0352
YRI	YRI	0.9	0.324	0.0350

*Supplementary Table 4: Spearman correlations between prediction versus simulated measurement from simulated populations to themselves across various shared eQTL proportions for  $k = 10$  causal eQTL.*

Correlation Mean (Std Err)		Train-test direction		
		AA	CEU	YRI
Training Pop	AA	0.321 (0.0071)	0.308 (0.0058)	0.336 (0.0052)
	CEU	0.334 (0.006)	0.329 (0.0069)	0.326 (0.0063)
	YRI	0.336 (0.0051)	0.299 (0.007)	0.325 (0.0063)

712 *Supplementary Table 5: Prediction performance under fully shared eQTL architecture for  $k = 10$*   
713 *eQTL yields reliable cross-population gene expression imputation. Results for other sizes of eQTL*  
714 *models are similar.*

<b>R<sup>2</sup></b>	<b>AFR to AFR</b>	<b>AFR to EUR</b>	<b>EUR to AFR</b>
<b>AFR to EUR</b>	1.222 x 10 <sup>-12</sup>		
<b>EUR to AFR</b>	1.705 x 10 <sup>-24</sup>	6.636 x 10 <sup>-06</sup>	
<b>EUR to EUR</b>	1.357 x 10 <sup>-04</sup>	1.487 x 10 <sup>-112</sup>	1.753 x 10 <sup>-228</sup>

715 *Supplementary Table 6: A Dunn test shows statistically significant differences when imputing*  
716 *between AFR and EUR populations versus imputing between EUR populations.*

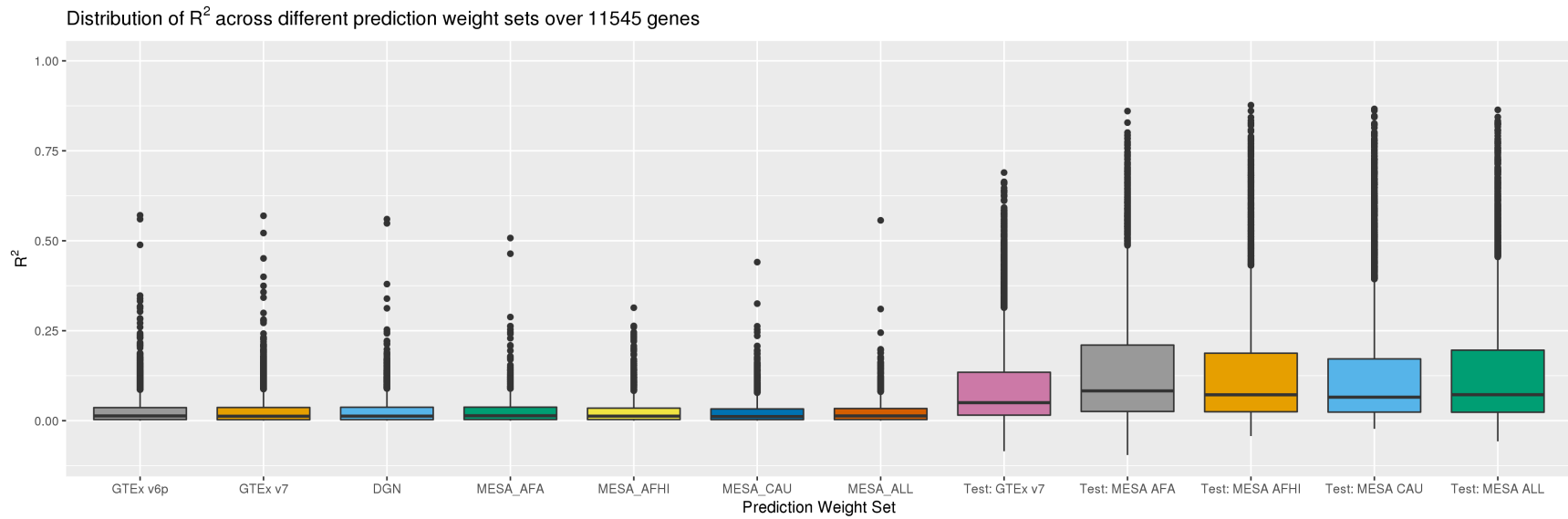
717



Z-score ( <i>p</i> -value)		Train-test direction				
		AA to CEU	AA to YRI	CEU to AA	CEU to YRI	YRI to AA
Train-test direction	AA to YRI	-28.029 ( <i>p</i> < 5.6 x 10 <sup>-172</sup> )	n/a n/a	n/a n/a	n/a n/a	n/a n/a
	CEU to AA	-0.244 ( <i>p</i> ~ 1)	27.784 ( <i>p</i> < 5.0 x 10 <sup>-169</sup> )	n/a n/a	n/a n/a	n/a n/a
	CEU to YRI	13.373 ( <i>p</i> < 6.5 x 10 <sup>-40</sup> )	41.403 ( <i>p</i> ~ 0)	13.618 ( <i>p</i> < 2.3 x 10 <sup>-41</sup> )	n/a n/a	n/a n/a
	YRI to AA	-28.725 ( <i>p</i> < 1.4 x 10 <sup>-180</sup> )	-0.695 ( <i>p</i> ~ 1)	-28.480 ( <i>p</i> < 1.5 x 10 <sup>-177</sup> )	-42.099 ( <i>p</i> ~ 0)	n/a n/a
	YRI to CEU	12.508 ( <i>p</i> ~ 0)	40.538 ( <i>p</i> ~ 0)	12.753 ( <i>p</i> ~ 0)	-0.865 ( <i>p</i> ~ 0)	41.234 ( <i>p</i> ~ 0)

718 *Supplementary Table 7: Differences in cross-population imputation performance are statistically significant, with a few notable*  
719 *exceptions. Imputation between AA and YRI, between AA and CEU, and between CEU and YRI is essentially the same, indicating that*  
720 *the direction of imputation does not matter.*

721



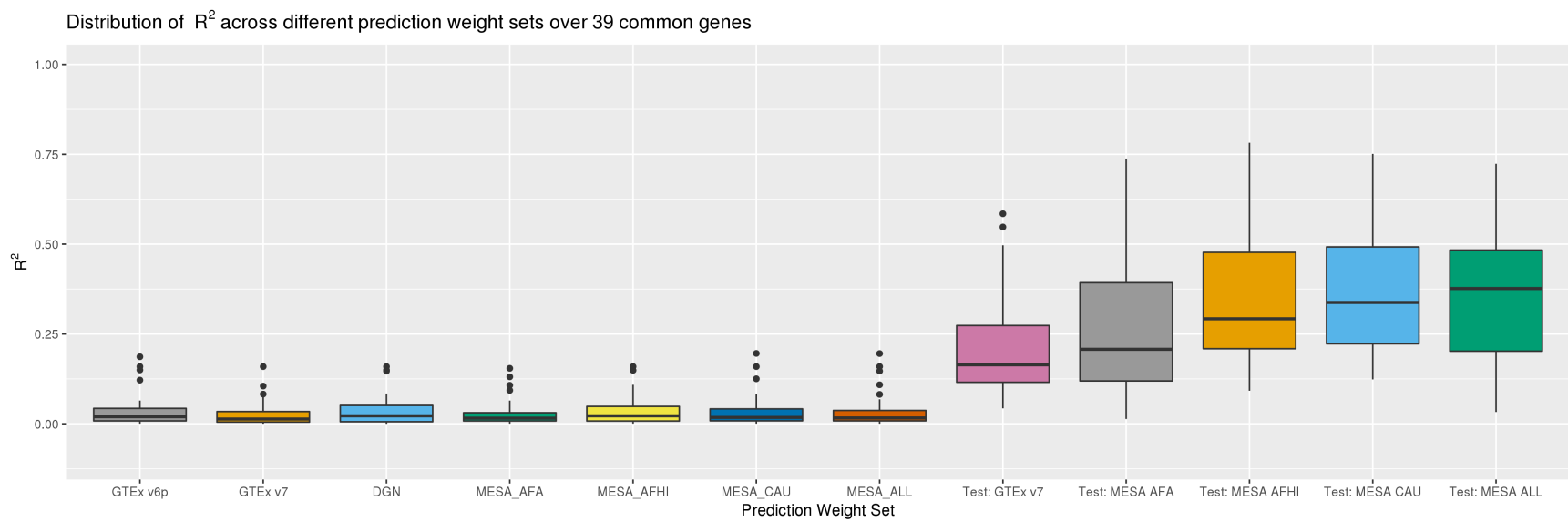
722

723

724

Supplementary Figure 1: A comparison of  $R^2$  between prediction and measurement in SAGE, with PredictDB test metrics as benchmarks, for 11,545 genes total. The number of genes per weight set varies; see Supplementary Table 1.

725

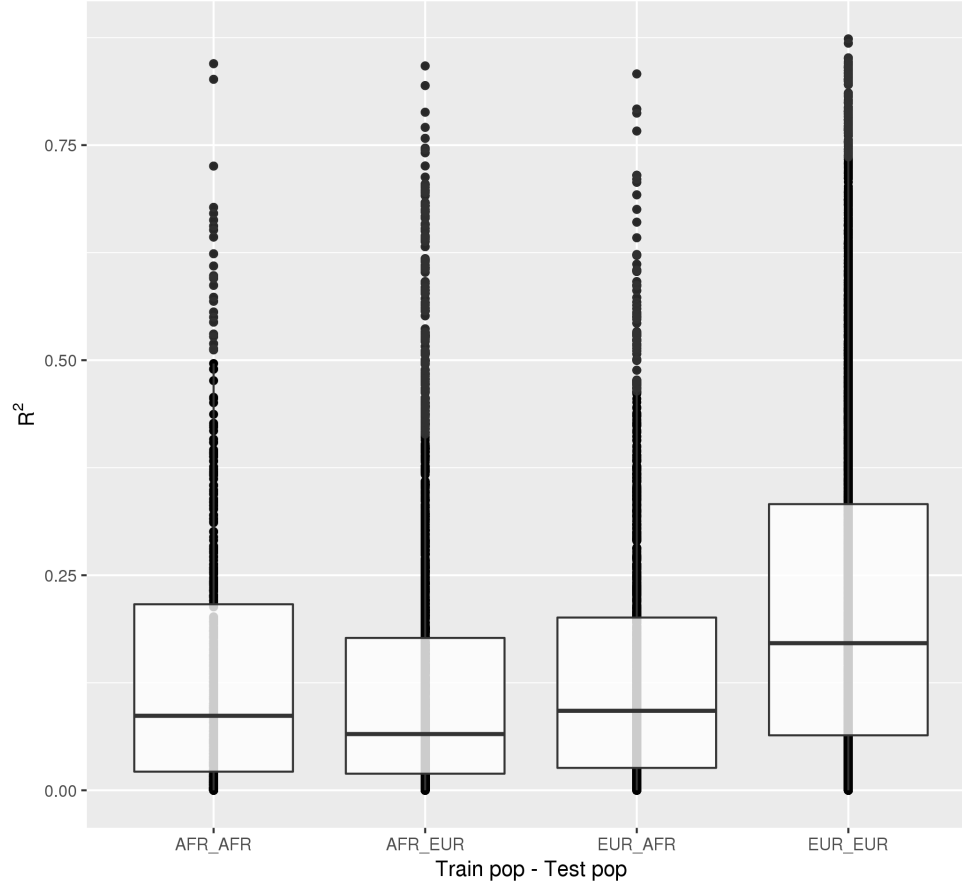


726  
727  
728

*Supplementary Figure 2:  $R^2$  between prediction and measurement in SAGE only using the 39 genes with positive correlation between prediction and measurement in all weight sets and benchmarks.*

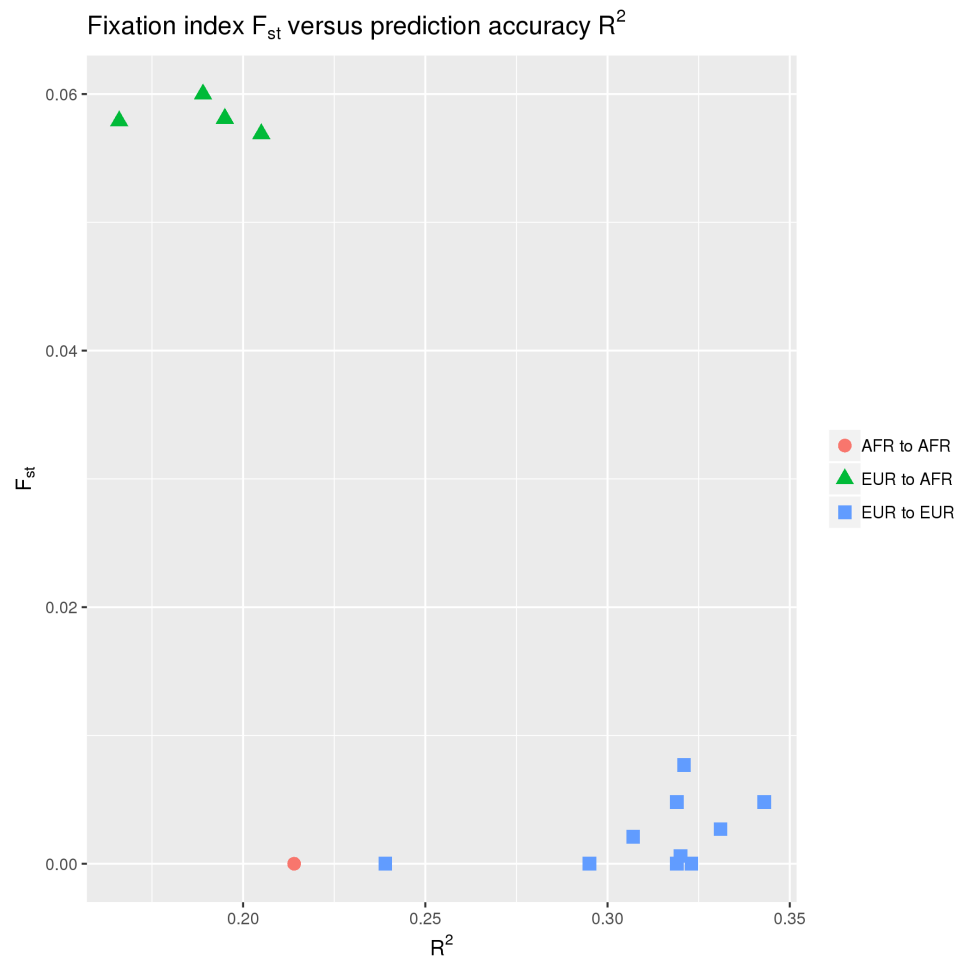
Comparison of  $R^2$  between continental GEUVADIS populations

N = 521 common genes



729

730 *Supplementary Figure 3: Imputation  $R^2$  between AFR (YRI) and EUR (CEU, TSI, GBR, and FIN). Imputing into and from AFR produces*  
731 *consistently lower  $R^2$  than imputing within EUR, suggesting a potential decrease in prediction accuracy when imputing across*  
732 *continental population groups.*

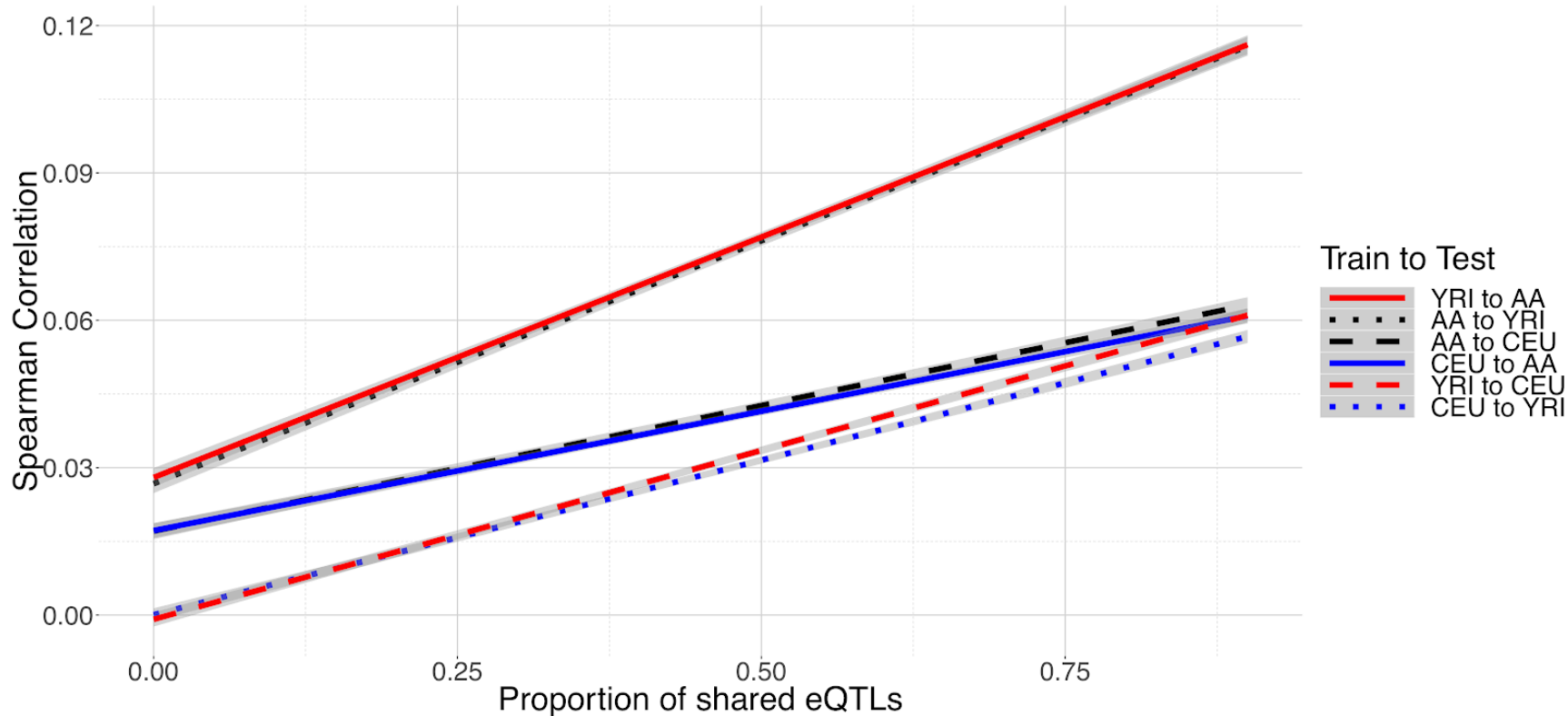


734

735 *Supplementary Figure 4: Genetic distance versus imputation accuracy over 142 genes with positive correlation across all train-test*  
 736 *scenarios. Here the GEUVADIS populations are arranged into three groups. AFR to AFR includes imputation from YRI into itself; EUR*  
 737 *to AFR includes imputation into YRI from CEU, GBR, TSI, and FIN; and EUR to EUR includes imputation within and between all*  
 738 *European populations in GEUVADIS. Clustering by genetic distance separates imputation between European populations from*  
 739 *imputation between European populations and AFR.  $F_{ST}$  are taken from the 1000 Genomes Project (Table S11).<sup>63</sup>*

740  
741  
742

### Crosspopulation correlations of predicted versus simulated gene expression Number of causal eQTLs: 5



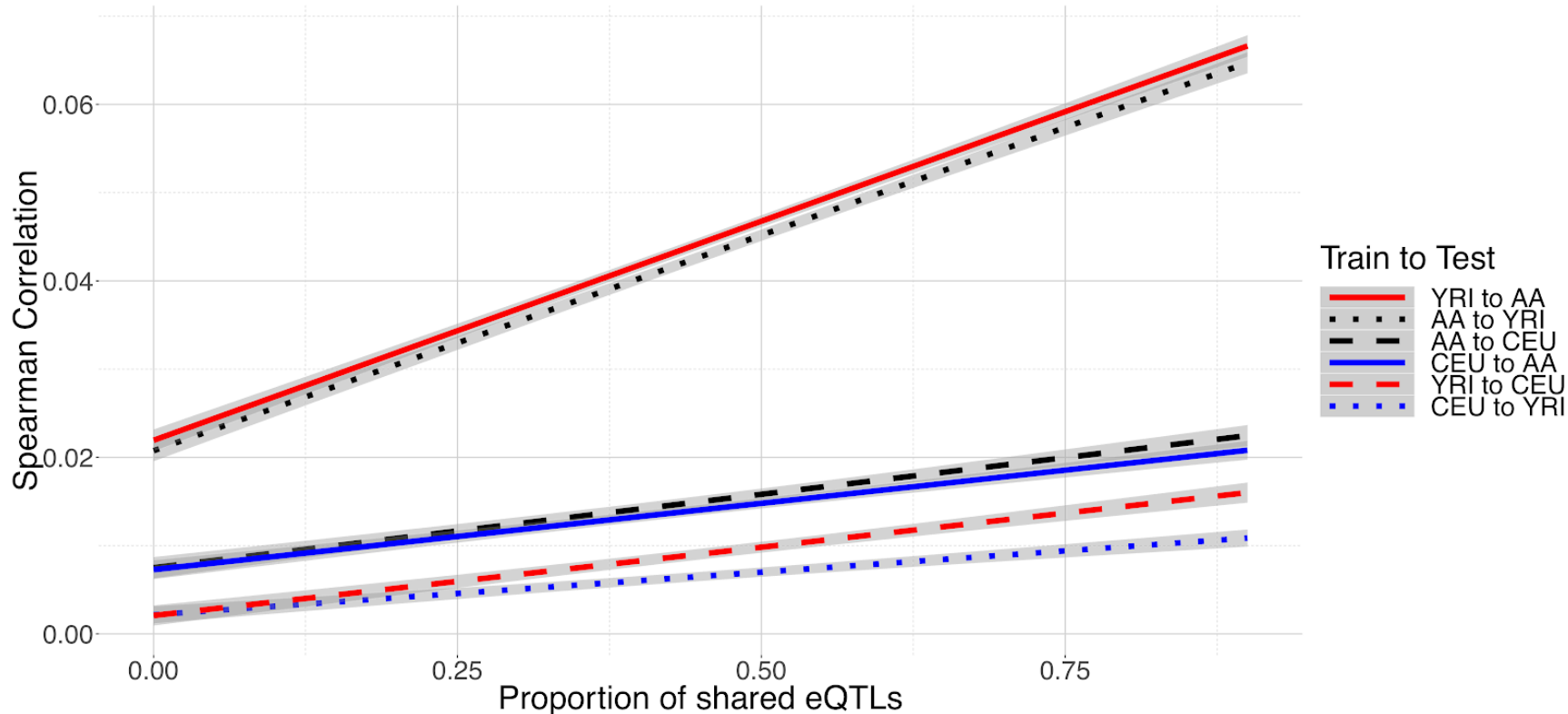
743  
744 *Supplementary Figure 5: Correlations between predictions and simulated gene expression measurements from simulated populations*  
745 *across various proportions of shared eQTL architecture with 5 causal cis-eQTL.*

746

747

748

### Crosspopulation correlations of predicted versus simulated gene expression Number of causal eQTLs: 20



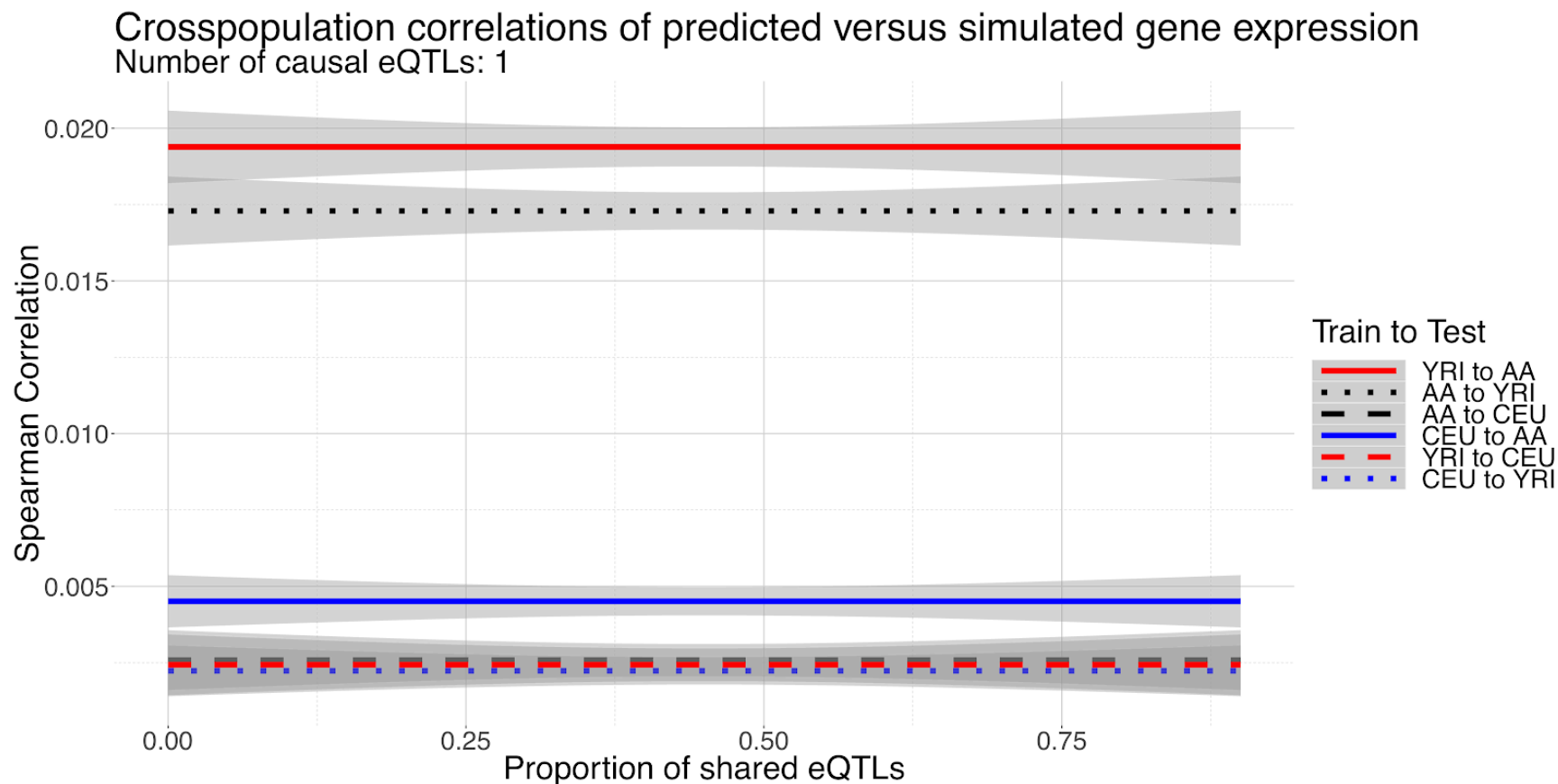
749

750 *Supplementary Figure 6: Correlations between predictions and simulated gene expression measurements from simulated populations*  
751 *across various proportions of shared eQTL architecture with 20 causal cis-eQTL.*

752



753  
754



755  
756  
757

Supplementary Figure 7: Correlations between predictions and simulated gene expression measurements from simulated populations across various proportions of shared eQTL architecture with a single causal cis-eQTL.