

1 **An Ancestry Based Approach for Detecting Interactions**
2

3 Danny S. Park^{1*}, Itamar Eskin², Eun Yong Kang³, Eric R. Gamazon^{4,5}, Celeste Eng⁶, Christopher R.
4 Gignoux^{1,7}, Joshua M. Galanter⁶, Esteban Burchard^{1,6}, Chun J. Ye⁸, Hugues Aschard⁹, Eleazar Eskin³,
5 Eran Halperin², Noah Zaitlen^{1,6*}
6

7 Affiliations:
8

- 1 1. Department of Bioengineering and Therapeutic Sciences. University of California San
Francisco. San Francisco, CA.
- 2 2. The Blavatnik School of Computer Science. Tel-Aviv University. Tel Aviv, Israel.
- 3 3. Department of Computer Science. University of California Los Angeles. Los Angeles, CA.
- 4 4. Division of Genetic Medicine, Department of Medicine. Vanderbilt University. Nashville, TN.
- 5 5. Academic Medical Center, University of Amsterdam, Amsterdam, The Netherlands
- 6 6. Department of Medicine. University of California San Francisco. San Francisco, CA.
- 7 7. Department of Genetics. Stanford University. Palo Alto, CA.
- 8 8. Institute of Human Genetics. University of California San Francisco. San Francisco, CA.
- 9 9. Department of Epidemiology. Harvard School of Public Health. Boston, MA.

10 * Corresponding Author
11

12 Email: danny.park@ucsf.edu, noah.zaitlen@ucsf.edu
13

14

24 **I. Abstract**

25 Background: Epistasis and gene-environment interactions are known to contribute significantly to
26 variation of complex phenotypes in model organisms. However, their identification in human
27 association studies remains challenging for myriad reasons. In the case of epistatic interactions, the
28 large number of potential interacting sets of genes presents computational, multiple hypothesis
29 correction, and other statistical power issues. In the case of gene-environment interactions, the lack
30 of consistently measured environmental covariates in most disease studies precludes searching for
31 interactions and creates difficulties for replicating studies.

32

33 Results: In this work, we develop a new statistical approach to address these issues that leverages
34 genetic ancestry in admixed populations. We applied our method to gene expression and methylation
35 data from African American and Latino admixed individuals respectively, identifying nine
36 interactions that were significant at $p < 5 \times 10^{-8}$, we show that two of the interactions in methylation
37 data replicate, and the remaining six are significantly enriched for low p-values ($p < 1.8 \times 10^{-6}$).

38

39 Conclusion: We show that genetic ancestry can be a useful proxy for unknown and unmeasured
40 covariates in the search for interaction effects. These results have important implications for our
41 understanding of the genetic architecture of complex traits.

42

43 Keywords: Gene-environment interaction, gene-gene interactions, admixture

44

45

46 II. Background

47 Genetic association studies in humans have focused primarily on the identification of
48 additive SNP effects through marginal tests of association. There is growing evidence that both
49 epistatic and gene-environment ($G \times E$) interactions contribute significantly to phenotypic variation
50 in humans and model organisms[1-5]. In addition to explaining additional components of missing
51 heritability, interactions lend insights into biological pathways that regulate phenotypes and improve
52 our understanding of their genetic architectures. However, identification of interactions in human
53 studies has been complicated by the computational and multiple testing burden in the case
54 of epistatic interactions, and the lack of consistently measured environmental covariates in the case
55 of $G \times E$ interactions[6,7].

56 To overcome these challenges, we leverage the unique nature of genomes from recently
57 admixed populations such as African Americans, Latinos, and Pacific Islanders. Admixed genomes are
58 mosaics of different ancestral segments[8] and for each admixed individual it is possible to
59 accurately estimate θ , the proportion of ancestry derived from each ancestral population (e.g. the
60 fraction of European/African ancestry in African Americans)[9]. Ancestry has been previously
61 leveraged to demonstrate that an array of environmental and biomedical covariates are correlated
62 with θ [10-20] and we therefore consider its use as a surrogate for unmeasured and unknown
63 environmental exposures. θ is also correlated with the genotypes of SNPs that are differentiated
64 between the ancestral populations, suggesting that θ may be effectively used as a proxy for detecting
65 multi-way epistatic interactions. Therefore, we propose a new SNP by θ test of interaction in order to
66 detect evidence of interaction in admixed populations.

67 We first investigate the properties of our method through simulated genotypes and
68 phenotypes of admixed populations. In our simulations we demonstrate that differential linkage-
69 disequilibrium (LD) between ancestral populations can produce false positive SNP by θ interactions
70 when local ancestry is ignored. To accommodate differential LD, we include local ancestry in our
71 statistical model and demonstrate that this properly controls this confounding factor. We also show

72 that our approach, the Ancestry Test of Interaction with Local Ancestry (AITL), is well-powered to
73 detect $G \times E$ interactions when θ is correlated with the environmental covariates of interest and
74 multi-way epistatic interactions. The power for detecting pairwise $G \times G$ interactions at highly
75 differentiated SNPs is lower than direct interaction tests even after accounting for the additional
76 multiple testing burden. However, the results of our simulations show that AITL is well powered to
77 detect multi-way epistasis involving tens or hundreds of SNPs of small effects, not detectable by
78 pairwise tests.

79 We first examined molecular phenotypes by applying our method to gene expression data
80 from African Americans, as well as DNA methylation data from Latinos. Gene expression traits have
81 previously been shown to have large-scale differences as a function of genetic ancestry[13]. Other
82 molecular phenotypes, such as LDL levels, have also been shown to be associated with genetic
83 ancestry [13,16,21-24]. For gene expression in particular, Price *et al.* showed that the effects of
84 ancestry on expression are widespread and not restricted to a handful of genes. Additionally,
85 molecular phenotypes are often used in deep phenotyping and Mendelian randomization studies and
86 are thus directly relevant to elucidating disease biology[25,26].

87 We identified one genome-wide significant interaction ($p < 5 \times 10^{-8}$) associated with gene
88 expression in the African Americans and eight significant interactions ($p < 5 \times 10^{-8}$) associated with
89 methylation in the Latinos. Two of the eight interactions associated with DNA methylation in the
90 Latinos also replicated and the remaining six were enriched for low p-values ($p < 1.8 \times 10^{-6}$). To
91 demonstrate that our approach works in larger data sets we also applied AITL to asthma case-control
92 data from Latinos and observed well-calibrated test statistics. Together, these results provide
93 evidence for the existence of interactions regulating expression and methylation and show that our
94 approach is statistically sound.

95 **III. Results**

96 **Simulated Data**

97 To determine the utility of using θ as a proxy for unmeasured and unknown environmental
98 covariates, we applied the AITL to simulated 2-way admixed individuals. We tested θ_1 , the
99 proportion of ancestry from ancestral population 1, for interaction with simulated SNPs (see
100 Simulation Framework). Power was computed over 1,000 simulations, assuming 10,000 SNPs being
101 tested, and using a Bonferroni correction p-value cutoff of 5×10^{-6} . We calculated the power using
102 assumed interaction effect sizes (either $\beta_{G \times G}$ or $\beta_{G \times E}$) of 0.1, 0.2, 0.3, and 0.4 (see Simulation
103 Framework). Although the few interactions reported for human traits and diseases have smaller
104 effects in terms of the phenotypic variance they explain, we simulated large effects because genetic
105 and environmental effect sizes in omics data, such as the expression and methylation data considered
106 here, are known to be of larger magnitude. For example, some cis-eQTL SNPs explain up to 50% of
107 the variance of gene expression[27]. However for most phenotypes, known interactions will explain a
108 very small proportion of the phenotypic variance, mainly due to the fact that so few interactions have
109 been identified and replicated[28].

110

111 *Power When Using θ as a Proxy for Highly Differentiated SNPs*

112 To determine whether using θ as a proxy for highly differentiated SNPs is more powerful
113 than testing all pairs of potentially interacting SNPs directly, we simulated two interacting SNPs in
114 1000 admixed individuals (see Simulation Framework). We then tested for an interaction using AITL
115 by replacing the genotypes at the highly differentiated SNP with $\vec{\theta}_1$. We observed that even with
116 moderate effect sizes, using θ in place of the actual genotypes does not provide any increase in power
117 even after accounting for multiple corrections (see Figure 1a). This is in agreement with recent work
118 showing the limited utility of local ancestry by local ancestry interaction test to identify underlying
119 SNP by SNP interaction when genotype data are available[29]. For the larger effect sizes we
120 simulated, we do see power increasing as the delta between ancestral frequencies increases. The
121 plots show that AITL has little power unless the effect was very strong. Figure 1b reveals that even
122 with the multiple correction penalty, testing all pairwise SNPs directly is always more powerful. We
123 note that when testing the interacting SNPs directly, we used a cutoff p-value of 1×10^{-9} since in
124 theory we were testing all unique pairs of 10,000 SNPs. Based on these results, we would

125 recommend testing for pairs of interacting SNPs directly if pairwise $G \times G$ interactions are a subject of
126 interest in the study.

127 However, when multi-way interactions are considered, AITL may become more powerful
128 since differentiated SNPs across the genome will be correlated with genetic ancestry. These
129 simulations are important as other studies have suggested that higher order interactions may be
130 important for some traits[1,30,31]. To evaluate the ability of θ to serve as a proxy for multiple
131 (independent) differentiated SNPs, we simulated a scenario where a candidate SNP z had
132 interactions with m SNPs (see Simulation Framework). For each interaction, we assumed a small
133 interaction effect size ($\beta_{G \times G} = 0.025$), which would not be detectable using a pairwise approach, as
134 we demonstrated in the pairwise simulation. Figure 2 shows that AITL is better powered to detect
135 the existence of interactions than a pairwise approach in the presence of multiple interacting SNPs
136 with a candidate SNP.

137

138 *Power When Using θ as a Proxy Environmental Covariate*

139 When assessing the utility of θ as a proxy for an environmental covariate E , we simulated
140 3000 individuals. E was simulated such that it was correlated with the global ancestries in varying
141 degrees (see Simulation Framework). Figure 3 shows the power of the AITL as a function of the
142 Pearson correlation between $\vec{\theta}_1$ and E . The power of testing E directly is exactly the power of the
143 AITL when the correlation is equal to 1. As expected, as the correlation increases, the power
144 increases as well. When the effect size is 0.1, the power to detect a $G \times E$ interaction is low whether
145 one uses θ_1 or E . However, both tests are much better powered for effect sizes greater or equal to 0.2,
146 with the AITL's power being dependent on the level of correlation. Note that using θ as a proxy for E
147 is equivalent to testing $G \times E$ in the presence of measurement error. Under the assumption of non-
148 differential error with regard to the outcome (e.g. the correlation between θ and E is equal among
149 cases and control) such a test is underpowered but has a controlled type I error rate under the
150 null[32].

151

152 *Differential LD*

153 To demonstrate that differential LD has the potential to cause inflated test-statistics, we ran
154 10,000 simulations of 1000 admixed individuals. For each individual we simulated 2 SNPs, a causal
155 SNP and a tag SNP. The LD between the tag SNP and causal SNP was different based on the ancestral
156 background the SNPs were on (see Simulation Framework). Over 10,000 simulations, we computed
157 the mean χ^2_1 test-statistic for the AIT and the AITL. We note that the phenotypes for these
158 simulations were generated under a model that assumed no interaction. We observed a mean $\chi^2_1 =$
159 0.996 with a standard deviation of 1.53 for AITL. AIT, which does not condition on local ancestry, had
160 a mean $\chi^2_1 = 3.59$ with a standard deviation of 3.60. We also looked at genomic control λ_{GC} , the ratio
161 of the observed median χ^2 over the expected median χ^2 under the null[33]. λ_{GC} compares the
162 median observed χ^2 test-statistic versus the true median under the null. In our simulations, we
163 observed $\lambda_{GC} = 5.81$ for AIT and $\lambda_{GC} = 0.980$ for AITL (see Supplementary Figure S1). Last, we
164 computed the proportion of test-statistics that passed a p-value threshold of .05 and .01 in our
165 simulations. The AIT had 3687 statistics passing a p-value of .05 and 1687 at a threshold of .01,
166 whereas AITL had 464 and 96 at the same p-value thresholds. The results for AITL are as expected
167 under a true null. The results from our simulations show that not accounting for local ancestry can
168 result in inflated test-statistics and can potentially lead to false positive findings.
169

170 **Real Data**

171 *Coriell Gene Expression Results*

172 We first applied our method to the Coriell gene expression dataset[34]. The Coriell cohort is
173 composed of 94 African-American individuals and the gene expression values of ~8800 genes in
174 lymphoblastoid cell lines (LCLs). Since African Americans derive their genomes from African and
175 European ancestral backgrounds, we tested for interaction between a given SNP and the proportion
176 of European ancestry, θ_{EUR} . Each SNP by θ_{EUR} term was tested once for association with the
177 expression of the gene closest to the SNP. We observed well-calibrated statistics with a λ_{GC} equal to
178 1.04 (see Supplementary Figure S2). In the LCLs, we found that interaction of rs7585465 with θ_{EUR}
179 was associated with ERBB4 expression (AITL $p = 2.95 \times 10^{-8}$, marginal $p = 0.404$) at a genome-wide
180 significant threshold ($p \leq 5 \times 10^{-8}$). rs7585465 has a 'C' allele frequency of 0.218 in the Coriell data

181 and appears to be differentiated between CEU and YRI with allele frequencies of 0.619 and 0.097 in
182 the respective populations.

183 Given that the gene expression values come from LCLs (all cultured according to the same
184 standards), the SNPs may be interacting with epigenetic alterations due to environmental exposures
185 that have persisted since transformation into LCLs. This scenario is unlikely, and we believe that
186 signals are driven by multi-way epistatic interactions. In our simulations, we showed that using θ as
187 a proxy for a single highly differentiated SNP is underpowered compared to testing all pairs of
188 potentially interacting SNPs directly. However, there are many SNPs that are highly differentiated
189 across the genome with which θ will be correlated. It is therefore possible that θ is capturing the
190 interaction between the aggregate of many differentiated trans-SNPs (i.e. global genetic background)
191 and the candidate SNP. This is consistent with a recently reported finding, conducted in human iPS
192 cell lines, that genetic background accounts for much of the transcriptional variation[2,35].

193 Although we believe the ERBB4 result to be representative of multi-way epistasis, we
194 performed a standard pairwise interaction test (see Methods) to check for interaction between
195 rs7585465 and other SNPs genome-wide. Interestingly, we found that the standard interaction test
196 (see Methods) showed substantial departure from the null with a λ_{GC} equal to 1.8 (see
197 Supplementary Figure S3). Since the interaction of rs7585465 by θ was significant, the pairwise
198 interaction test-statistics of rs7585465 by any SNP j can be inflated if j is correlated with θ . We found
199 that including the original significant SNP by θ term in the null (see Methods) brought the λ_{GC} down
200 to 1.05, and controlled for such scenarios in this dataset (See Supplementary Figure S3). As we had
201 previously anticipated, identifying the exact interactions driving the SNP by θ interaction proved to
202 be difficult. We found one borderline significant SNP (rs4839709, $p = 3.08 \times 10^{-7}$) but no
203 interactions that passed genome-wide significance. These results are consistent with what we have
204 observed in simulations, in which even though a standard pairwise interaction test is underpowered
205 to detect interactions, AITL is able to identify the main locus involved in a multi-way interaction.

206
207

208 *GALA II Case-Control*

209 To determine if our method is biased in large structured GWAS data, we applied AITL to
210 case-control data from a study of asthmatic Latino individuals called the Genes-environments and
211 Admixture in Latino Americans (GALA II)[36]. The dataset includes 1158 Mexicans and 1605 Puerto
212 Ricans, which were analyzed separately. Case status was assigned to individuals if they were
213 between the ages of 8 and 40 years with a physician-diagnosed mild to moderate-to-severe asthma.
214 Additionally, they had to have experienced 2 or more asthma related symptoms in the previous 2
215 years at the time of recruitment[37]. In the Mexicans and Puerto Ricans there were 548 and 797
216 cases, respectively. In our analysis, we also included BMI, age, and sex as additional covariates. We
217 observed well-calibrated statistics with a λ_{GC} equal to 1.00 and 0.98 in the Mexicans and Puerto
218 Ricans, respectively (see Supplementary Figure S5). In contrast to the molecular phenotype data,
219 searches for interactions in these phenotypes did not yield any findings passing genome-wide
220 significance. This is consistent with previous disease studies that have failed to find many replicable
221 interactions in disease studies[28]. In the data here, the lack of any findings may be due to the
222 relatively small sample size or because the effects of the interactions are extremely small (if they
223 exist for covariates correlated with θ_{EUR}).

224

225 *GALA II Methylation Results*

226 We searched for interactions in methylation data derived from a study of GALA II asthmatic
227 Latino individuals[36]. The methylation data is composed of 141 Mexicans and 184 Puerto Ricans. As
228 the phenotype, we used DNA methylation measurements on ~300,000 markers from peripheral
229 blood. As we had done with gene expression, we tested for interaction between a given SNP and θ_{EUR}
230 using AITL. All SNPs within a 1 MB window centered around the methylation probe were tested. We
231 used the European component of ancestry because it is the component shared most between
232 Mexicans and Puerto Ricans (see Table 1). We observed well-calibrated test-statistics with λ_{GC} equal
233 to 1.06 in the Mexicans and 0.96 in the Puerto Ricans (see Supplementary Figure S6). We tested
234 128,794,325 methylation-SNP pairs, which result in a Bonferroni corrected p-value cutoff of
235 3.88×10^{-10} . However, this cutoff is extremely conservative given the tests are not independent. We
236 therefore report all results that are significant at 5×10^{-8} in either set as an initial filter. We found 5

237 interactions in the Mexicans and 3 in the Puerto Ricans that are significant at this threshold (see
238 Table 2).

239 Unlike the Coriell individuals, who are 2-way admixed, the GALA II Latinos are 3-way
240 admixed and derive their ancestries from European, African, and Native American ancestral groups.
241 Consequently, to confirm that incomplete modeling or better tagging on one of the non-European
242 ancestries was not driving the results, we retested all significant interactions including a second
243 component of ancestry for AITL. In the case of the Mexicans, we included African and European
244 ancestry, and in the case of the Puerto Ricans, we included European and Native American ancestry.
245 Even after adjusting for the second ancestry the interactions between SNP and θ_{EUR} remained highly
246 significant (see Supplementary Table 1).

247 As we did for the gene expression data, we attempted to identify pairwise interactions
248 involved in the methylation data results. For each genome-wide significant result, we performed a
249 standard pairwise interaction test of all SNPs with the original SNP found to be significant with AITL.
250 We were unable to identify any significant interactions after applying genomic control to the results.
251 For all tests, we included the significant SNP by θ term (see Methods) in the null. For this dataset,
252 unlike the gene expression data, we observed substantial remaining departure from the null (see
253 Supplementary Table S2) even after including the original significant SNP by θ term, suggesting there
254 may be other factors that need to be accounted for when testing for interactions in admixed
255 populations. The results from our pairwise scan are what we would anticipate, given that in
256 simulations only AITL (not the standard pairwise interaction test) was able to identify the main locus
257 involved in the multi-way interaction.

258 We then performed a replication study of the significant Puerto Rican associations in the
259 Mexican cohort and vice versa. To account for the fact that we are replicating eight total results
260 across both populations, we used a Bonferroni corrected p-value threshold equal to $.05/8 =$
261 6.25×10^{-3} . The interaction of rs4312379 and rs4312379 with ancestry in the Puerto Ricans
262 replicated in the Mexicans. Furthermore, there was a highly significant enrichment of low p-values in
263 the replication study among the discovery results (permutation $p < 1 \times 10^{-4}$). Furthermore, 5 out of
264 the 6 non-replicating results have a p-value less than 0.05 (binomial test $p < 1.8 \times 10^{-6}$). The results

265 of the permutation and binomial test suggests that the interactions that did not replicate are likely to
266 do so with bigger sample sizes. It is important to note that replicated interactions and the enrichment
267 for low p-values do not necessarily indicate that the same genetic or environmental covariates are
268 interacting with the genetic locus in both populations. The covariates correlated with θ_{EUR} in one
269 population are not necessarily those correlated with θ_{EUR} in the other population. There may be
270 correlations which exist in both populations but θ_{EUR} serves as a proxy for all such correlated
271 covariates and therefore should not be necessarily viewed as a proxy for any specific one. Overall,
272 our results from the GALA II (methylation) cohort suggest there are both genetic and environmental
273 variables contributing to epistasis that have yet to be discovered in admixed individuals.

274

275 **IV. Discussion and Conclusions**

276 For many disease architectures, interactions are believed to be a major component of
277 missing heritability[38]. Finding new interactions has proven to be difficult for logistical, statistical,
278 biological, and computational reasons. In this study, we have demonstrated that in admixed
279 populations, testing for $G \times \theta$ interactions can be leveraged to overcome some of the difficulties
280 typically encountered when searching for interactions. The computational cost is minimal and has
281 the same order as running a standard GWAS.

282 One drawback of our method is that it does not identify which covariate is interacting with a
283 genetic locus. Nevertheless, the approach can show whether an interaction effect exists in a given
284 dataset and if it does exist, our method ensures that an underlying genetic or environmental
285 covariate(s) is correlated with ancestry. Additionally, in the case where there is no marginal effect,
286 our approach identifies new loci and shows that the genetic locus influences the phenotype and
287 exerts its effects through interactions, which has important implications for the genetic architecture
288 of the phenotype. The relative contribution of additive and non-additive genetic effects to variability
289 in molecular phenotypes and disease risk is an important area of investigation, and our approach
290 provides a direct test for detecting non-additive contributions[39].

291 Environmental covariates are often not consistently measured across cohorts whereas
292 genetic ancestry is nearly perfectly replicable. Testing for the presence of interaction using a nearly

293 perfectly reproducible covariate may enhance our understanding of the genetic basis of disease and
294 other traits. Our method also provides the additional benefit of not being confounded by interactions
295 between unaccounted-for covariates[40].

296 Association testing for interaction effects involving continuous environmental exposures in
297 the context of mixed-models remains an open problem. For binary environmental exposures, it has
298 been shown that mixed-models control for population structure nominally better than including
299 genetic ancestry (or principal components) as a covariate[41]. Because it is unclear how mixed-
300 models perform with continuous environmental exposures, especially those correlated with
301 ancestry, in our analyses we took the standard approach of filtering related individuals and including
302 ancestry as a covariate.

303 It has been shown that 2-step analyses may be more powerful for detecting interactions
304 when exposures are binary [42-44]. However, these studies have primarily been done in a single
305 homogeneous population, and the correct null distribution for the interaction effect must assume
306 that the 2nd stage procedure is independent of the marginal effect test-statistic. In real data, using a 2-
307 step approach in conjunction with AITL to test for interactions may be problematic because the
308 interaction effect size will not necessarily be independent of the marginal effect size, as the allele
309 frequency at any SNP will be a function of ancestry in an admixed population. Additionally, only 1 of
310 the interaction results that we report here had a marginal effect ($p < 0.05$) and thus would have been
311 missed by a 2-step approach. Thus, our approach can serve to complement or extend the frequently
312 used 2-step procedure for detecting interaction effects.

313 Results from our multi-way epistasis simulation analyses and empirical data in cell lines
314 suggest that genetic ancestry is a good proxy for genetic background, since all highly differentiated
315 SNPs across the genome will be correlated with genetic ancestry. Our simulations also demonstrated
316 that genetic ancestry can be a good proxy for an environmental covariate depending on the
317 correlation between the two. However, it may be the case that there are multiple environmental
318 factors interacting with a genetic locus, all of which are correlated with θ in differing degrees and
319 effect sizes. Such a situation would mirror what we saw in our multi-way $G \times G$ simulations where a
320 single interaction may not be detectable by using a traditional $G \times E$ test, but because θ aggregates the

321 effects of all interacting covariates, AITL would be able to detect it. There are also other contexts in
322 which modeling SNP by θ may be useful, such as using variance components. For example, SNP by θ
323 interaction terms can be used in a mixed-model framework to test for interaction effects because
324 genetic ancestry is correlated with many genetic markers and environmental covariates[45].

325 For some traits, there may be systematic differences between ancestral populations in the
326 genetic effects on the trait. In admixed individuals with these ancestral populations, the effect of
327 genetic variation on phenotype will be reflected in the correlation between phenotype and θ , thereby
328 affecting epistatic and $G \times E$ interactions. It will be interesting to see how much of the phenotype-
329 ancestry correlations are due to epistatic and $G \times E$ interactions.

330 In our analysis of real data, we discovered gene by θ interactions associated with genes that
331 have known interactions. In the GALA II Mexicans, the interaction of rs925736 with ancestry was
332 associated with the methylation of HDAC4, a known histone deacetylase (HDAC). In concert with
333 DNA methylases, HDACs function to regulate gene expression by altering chromatin state[46]. In
334 Europeans, HDACs have been shown to be associated with lung function through direct genetic
335 effects and through environmental interactions[47,48]. For the GALA II Puerto Ricans, rs17091085
336 showed an interaction associated with the methylation state of SERPINA6. Of note, interaction
337 between birth weight and SERPINA6 has been previously associated with Hypothalamic-Pituitary-
338 Adrenal axis function[49]. Further investigations of our interaction findings are thus warranted.

339 In the GALA II (methylation) dataset, two of the eight significant associations replicated and,
340 in general, the results had an enrichment of low p-values in the replication dataset. However, we note
341 that if the interactions detected by AITL are multi-way epistasis it is more likely that the results will
342 replicate. This is because most SNPs differentiated in the Mexicans will still be differentiated in the
343 Puerto Ricans, and thus still be correlated with θ . If the interactions detected by AITL are $G \times E$
344 interactions, then the interactions are less likely to replicate because the same environmental
345 covariate(s) will need to be correlated with ancestry in both groups.

346 Another caveat is that the Mexicans and Puerto Ricans, though independent, are part of the
347 same study and occasionally technical artifacts, such as issues with genotyping or measuring
348 methylation, can affect downstream analyses of both populations. For our analyses, we have taken

349 careful quality-control steps to ensure that this is not the case and there is no apparent inflation of
350 test-statistics as demonstrated by our values for genomic control. Future research of interactions
351 using AITL should keep such caveats in mind.

352 We investigated in detail the potential of single SNP-SNP interactions driving the results that
353 were found both in the gene expression and methylation datasets. As demonstrated by the wide
354 range of λ_{GC} values, we observed that non-linear effects can cause substantial departure from the null
355 when testing for pairwise SNP-SNP interactions. This is especially true when testing for interaction
356 between SNPs s and j , where s has a significant interaction with θ and j is correlated covariates that
357 are also correlated with θ . As we saw in the gene expression data, including the significant SNP by θ
358 term can properly control for such situations, but its use in standard pairwise interaction tests
359 warrants further investigation.

360 Our analysis revealed the existence of interactions but does not provide a direct way to
361 determine the covariate that is interacting with a SNP. Further methodological work is required to
362 uncover the exact environmental exposures or genetic loci with which SNPs are interacting. The
363 existence of gene by θ interactions in GALA II underscores why modeling interactions should be
364 considered for future association studies and for heritability estimation in admixed populations.

365

366 V. Materials and Methods

367 Our approach is best illustrated with an example. First consider testing a SNP s for
368 interaction with an environmental covariate E . θ can serve as a proxy for E if the two are correlated,
369 even if E is unknown or unmeasured (see Figure 4a). Now consider testing s for interaction with a
370 SNP $j \neq s$ that is highly differentiated in terms of ancestral allele frequencies. For example, a SNP that
371 has a high allele frequency in one ancestral population and a low allele frequency in the other
372 ancestral population. θ can be used as a proxy for j because θ and the genotypes of SNP j will be
373 correlated. Consider the case where j has a frequency of 0.9 in population 1 and frequency of 0.1 in
374 population 2. Individuals with large values of θ_1 are more likely to have derived j from population 1
375 and on average have greater genotype values at j . Similarly, individuals with small values of θ_1 are

376 more likely to have derived j from population 2 and on average have smaller genotype values. Thus, θ
377 will be correlated with the genotypes of the individuals for highly differentiated SNPs and can serve
378 as a proxy for detecting interactions (see Figure 4b).

379 Consider an admixed individual i who derives his or her genome from k ancestral
380 populations. We denote individual i 's global ancestry proportion as $\theta_i =$
381 $\langle \theta_{i1}, \theta_{i2}, \dots, \theta_{ik} \rangle$, where $\sum_k \theta_{ik} = 1$. The local ancestry of individual i at a SNP s is denoted as $\gamma_{ais} \in$
382 $\{0, 1, 2\}$ and is equal to the number of alleles from ancestry $a \in \{1 \dots k\}$ inherited at SNP s . Current
383 methods allow us to estimate ancestry directly from genotype data both globally and at specific
384 SNPs[9,50,51]. We denote the genotype of an individual i at SNP s as $g_{is} \in \{0, 1, 2\}$ and the
385 corresponding phenotype as y_i .

386 In this work, we model continuous phenotypes in an additive linear regression framework.
387 Assuming n (unrelated) individuals, define \vec{y} to be the vector of all individuals' phenotypes. The
388 model for the phenotype is then

389
$$\vec{y} = \mathbf{X}\vec{\beta} + \vec{\varepsilon}$$

390 where $\vec{\varepsilon} \sim \mathcal{N}(0, \sigma)$ is a $n \times 1$ vector of error terms, \mathbf{X} is a $n \times v$ matrix of v covariates, and $\vec{\beta}$ is a $v \times 1$
391 vector of the covariate effect sizes. We note that in our notation $\vec{v}^2 = \vec{v}^T \vec{v}$ for a vector \vec{v} . Assuming
392 independence, the likelihood under this model is:

393
$$L = \left(\frac{1}{\sigma\sqrt{2\pi}} \right)^n \exp \left(-\frac{1}{2\sigma^2} (\vec{y} - \mathbf{X}\beta)^2 \right)$$

394 We can compute the log-likelihood ratio statistic (D) using a maximum likelihood approach:

395
$$D = -2 (\log L_1 - \log L_0) = -2 \left(n \log(\hat{\sigma}_{L_1}) + \frac{(\vec{y} - \mathbf{X}\hat{\beta}_{L_1})^2}{2\hat{\sigma}_{L_1}^2} \right) + 2 \left(n \log(\hat{\sigma}_{L_0}) + \frac{(\vec{y} - \mathbf{X}\hat{\beta}_{L_0})^2}{2\hat{\sigma}_{L_0}^2} \right)$$

396 We note that for a case-control phenotype we would use the following likelihood and log-likelihood
397 ratio statistic:

398
$$L = \prod_{i=1}^n \left[\frac{1}{1 + e^{-X_i \beta}} \right]^{y_i} \left[1 - \frac{1}{1 + e^{-X_i \beta}} \right]^{1-y_i}$$

399
$$D = -2 (\log L_1 - \log L_0)$$

400
$$= -2 \left(\sum_{i=1}^n -\log(1 + e^{-X_i \hat{\beta}_{L1}}) + \sum_{i=1}^n y_i (X_i \hat{\beta}_{L1}) \right)$$

401
$$+ 2 \left(\sum_{i=1}^n -\log(1 + e^{-X_i \hat{\beta}_{L0}}) + \sum_{i=1}^n y_i (X_i \hat{\beta}_{L0}) \right)$$

402

403 where X_i is the i -th row of the matrix \mathbf{X} , which correspond to the covariates of individual i .

404 For linear regression, the maximum likelihood estimator (MLE) of the effect sizes is $\hat{\beta} =$

405 $(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \vec{y}$, and the MLE of the error variance is $\hat{\sigma}^2 = \frac{1}{n} (\vec{y} - \mathbf{X} \hat{\beta})^2$. Here, L_1 is the likelihood under

406 the alternative and L_0 is the likelihood under the null. $(\hat{\beta}_{L1}, \hat{\sigma}_{L1}^2)$ and $(\hat{\beta}_{L0}, \hat{\sigma}_{L0}^2)$ are the effect sizes and

407 error variance estimates that maximize the respective likelihoods. D is distributed as χ^2 with k

408 degrees of freedom (df), where k is the number of parameters constrained under the null.

409

410 **1-df Ancestry Interaction Test (AIT)**

411 The first test we present is the standard direct test of interaction. We test for a SNP's

412 interaction with θ instead of an environmental covariate or another genotype. Let $\vec{g}_s = \langle g_{1s} \dots g_{ns} \rangle$ be

413 the vector of the individuals' genotypes at SNP s , $\vec{\theta}_a = \langle \theta_{1a} \dots \theta_{na} \rangle$ be the vector of their global

414 ancestries for ancestry a , and $\vec{g}_s \times \vec{\theta}_a$ be the vector of interaction terms which result from the

415 component-wise multiplication of the genotype and global ancestry vectors. We test the alternative

416 hypothesis $(\hat{\beta}_{G \times \theta} \neq 0)$ against the null hypothesis $(\hat{\beta}_{G \times \theta} = 0)$.

417

418
$$H_1: \vec{y} = \vec{g}_s + \vec{g}_s \times \vec{\theta}_a + \vec{\theta}_a$$

419
$$H_0: \vec{y} = \vec{g}_s + \vec{\theta}_a$$

420

421 In this test of interaction, we test a single ancestry versus the other ancestries that may be present in
422 the population of interest. One parameter is constrained under the null which results in a statistic
423 with $k=1$ df. Let $\hat{\beta}_{L_{\{0,1\}}(s)}$, $\hat{\beta}_{L_{\{0,1\}}(G \times \theta)}$, and $\hat{\beta}_{L_{\{0,1\}}(\theta)}$ denote the effect sizes of genotype, interaction, and
424 global ancestry under a given hypothesis respectively. The statistic is given below.

425

$$426 D = -2 \left(n \log(\hat{\sigma}_{L_1}) + \frac{[\vec{y} - \mathbf{X} \langle \hat{\beta}_{L_1(s)}, \hat{\beta}_{L_1(G \times \theta)}, \hat{\beta}_{L_1(\theta)} \rangle]^2}{2\hat{\sigma}_{L_1}^2} \right) + 2 \left(n \log(\hat{\sigma}_{L_0}) + \frac{[\vec{y} - \mathbf{X} \langle \hat{\beta}_{L_0(s)}, 0, \hat{\beta}_{L_0(\theta)} \rangle]^2}{2\hat{\sigma}_{L_0}^2} \right)$$

427 where \mathbf{X} is an $n \times 3$ matrix composed of \vec{g}_s , $\vec{\theta}_a$, and $\vec{g}_s \times \vec{\theta}_a$ as columns.

428

429 **1-df Ancestry Interaction Test with Local Ancestry (AITL)**

430 Given that the individuals we analyze in this work are assumed to be admixed, there is
431 potential for confounding due to differential LD. An interaction that is not driven by biology could
432 occur due to the possibility that a causal variant may be better tagged by a SNP being tested on one
433 ancestral background versus another (See Figure 4c). We account for the different LD patterns on
434 varying ancestral backgrounds by including local ancestry as an additional covariate in AITL. By
435 including local ancestry, we assume that the SNP being tested is on the same local ancestry block as
436 the causal SNP that it may be tagging. Such an assumption is reasonable because admixture in
437 populations such as Latinos and African Americans are relatively recent events and their genomes
438 have not undergone many recombination events. As a result, local ancestry blocks on average stretch
439 for several hundred kilobases[52,53].

440 Let $\vec{\gamma}_{as} = \langle \gamma_{a1s} \dots \gamma_{a1s} \rangle$ be the vector of local ancestry calls for all individuals for ancestry a
441 and let $\vec{g}_s \times \vec{\gamma}_{as}$ be the interaction terms from piecewise multiplication of the two vectors. We use the
442 following alternative and null hypotheses:

443

$$444 H_1: \vec{y} = \vec{g}_s + \vec{g}_s \times \vec{\theta}_a + \vec{\theta}_a + \vec{\gamma}_{as} + \vec{g}_s \times \vec{\gamma}_{as}$$

$$445 H_0: \vec{y} = \vec{g}_s + \vec{\theta}_a + \vec{\gamma}_{as} + \vec{g}_s \times \vec{\gamma}_{as}$$

446

447 Here we are testing for an interaction effect, i.e. $\hat{\beta}_{G \times \theta} \neq 0$, and constrain one parameter under the
448 null resulting in a statistic with $k=1$ df. Let $\hat{\beta}_{L_{\{0,1\}}(G \times \gamma)}$ and $\hat{\beta}_{L_{\{0,1\}}(\gamma)}$ denote the effect sizes of the
449 interaction between genotype and local ancestry and just local ancestry, respectively. The log
450 likelihood ratio statistic is given by

451

452
$$D = -2 \left(n \log(\hat{\sigma}_{L_1}) + \frac{[\vec{y} - \mathbf{X} \langle \hat{\beta}_{L_1(s)}, \hat{\beta}_{L_1(G \times \theta)}, \hat{\beta}_{L_1(\theta)}, \hat{\beta}_{L_1(\gamma)}, \hat{\beta}_{L_1(G \times \gamma)} \rangle]^2}{2\hat{\sigma}_{L_1}^2} \right)$$

453
$$+ 2 \left(n \log(\hat{\sigma}_{L_0}) + \frac{[\vec{y} - \mathbf{X} \langle \hat{\beta}_{L_0(s)}, 0, \hat{\beta}_{L_0(\theta)}, \hat{\beta}_{L_0(\gamma)}, \hat{\beta}_{L_0(G \times \gamma)} \rangle]^2}{2\hat{\sigma}_{L_0}^2} \right)$$

454 where \mathbf{X} is an $n \times 5$ matrix composed of \vec{g}_s , $\vec{\theta}_a$, $\vec{g}_s \times \vec{\theta}_a$, $\vec{\gamma}_{as}$, and $\vec{g}_s \times \vec{\gamma}_{as}$ as columns. All of these test-
455 statistics are straightforwardly modified to jointly incorporate several ancestries in the case of multi-
456 way admixed populations.

457

458 **Standard Pairwise Test of Interaction and Controlling Confounding in Admixed Populations**

459 Here we present the standard approach for testing for interaction between two SNPs. We
460 use the following alternative and null hypotheses.

461
$$H_1: \vec{y} = \vec{g}_1 + \vec{g}_2 + \vec{g}_1 \times \vec{g}_2 + \vec{\theta}_a$$

462
$$H_0: \vec{y} = \vec{g}_1 + \vec{g}_2 + \vec{\theta}_a$$

463 If AITL is significant for a given SNP s , then any SNP j tested for interaction with s may be biased if j is
464 correlated with covariates that are also correlated with θ . Furthermore, if the effects of the covariates
465 correlated with θ are non-linear then controlling for the main effects of the SNPs and ancestry will
466 account for the non-linear effects. We thus, propose the use of the following alternative and null
467 hypotheses.

468
$$H_1: \vec{y} = \vec{g}_s + \vec{g}_j + \vec{g}_s \times \vec{g}_j + \vec{\theta}_a + \vec{g}_s \times \vec{\theta}_a$$

469
$$H_0: \vec{y} = \vec{g}_s + \vec{g}_j + \vec{\theta}_a + \vec{g}_s \times \vec{\theta}_a$$

470 We note that the utility of this test will require further investigation (see Discussion).

471

472

473 **Simulation Framework**

474 For all our simulations, we simulated 2-way admixed individuals. Global ancestry for
475 ancestral population 1 (θ_1) was drawn from a normal distribution with $\mu = 0.7$ and $\sigma = 0.2$.
476 Individuals with $\theta_1 > 1$ or $\theta_1 < 0$ were assigned a value of 1 or 0, respectively. We simulated
477 phenotypes of individuals to investigate our method in four different scenarios: $G \times E$ interactions,
478 pairwise epistatic interactions, multi-way epistatic interactions, and false positive interactions due to
479 local differential tagging.

480

481 To simulate phenotypes under the situation of a $G \times E$ interaction, we simulated a single SNP.
482 For each individual i , we assigned the local ancestry or the number of alleles derived from population
483 1 (γ_{ai}) for each haplotype by performing two binomial trials with the probability of success equal to
484 θ_{i1} . We then drew ancestry specific allele frequencies following the Balding-Nichols model by
485 assuming a $F_{ST} = 0.16$ and drawing two ancestral frequencies, p_1 and p_2 , from the following beta
486 distribution[54].

487

488
$$p_1, p_2 \sim Beta \left(\frac{p(1 - F_{ST})}{F_{ST}}, \frac{(1 - p)(1 - F_{ST})}{F_{ST}} \right)$$

489

490 where p is the underlying MAF in the entire population and is set to 0.2. Genotypes were drawn using
491 a binomial trial for each local ancestry haplotype with the probability of success equal to p_1 or p_2 for
492 values of $\gamma_{ai} = 0$ or 1, respectively. Environmental covariates correlated with θ_1 , E_i , were generated
493 for each individual i by drawing from a normal distribution $\mathcal{N}(\mu = \theta_{i1}, \sigma_E)$, where σ_E is the standard
494 deviation of the environmental covariates. σ_E was varied from 0 to 5 in increments of 0.005 to create
495 E_i 's that were correlated with individuals' global ancestries in varying degrees. We generated
496 phenotypes for individuals assuming only an interaction effect by drawing from a normal
497 distribution, $\mathcal{N}(\mu = \beta_{G \times E} \times g_{i1} \times E_i, \sigma = 1)$ for a given interaction effect size ($\beta_{G \times E}$).

498

499 To simulate phenotypes based on pairwise epistatic interactions, we simulated two SNPs. At
500 both SNPs, we assigned the local ancestry values as described for the $G \times E$ case. We assigned
501 genotypes for individuals at the first SNP assuming an allele frequency of 0.5 for both populations
502 and drawing from two binomial trials. We assigned genotypes at the second SNP over a wide range of
503 ancestry specific allele frequencies to simulate different levels of SNP differentiation. Ancestry
504 specific allele frequencies were initially $p_1 = p_2 = 0.5$ and iteratively increasing p_1 by 0.005 while
505 simultaneously decreasing p_2 by 0.005 until $p_1 = 0.05$ and $p_2 = 0.95$. Genotypes at the second SNP
506 were drawn using the same approach described for $G \times E$. Using the simulated genotypes, phenotypes
507 were drawn from a normal distribution, $\mathcal{N}(\mu = \beta_{G \times G} \times g_{i1} \times g_{i2}, \sigma = 1)$, where g_{is} is the genotype for
508 individual i at the simulated SNP s .

509 To simulate phenotypes based on multi-way epistatic interactions, we simulated a SNP z and
510 m (independent) SNPs with pairwise interactions with z . Genotypes for individuals at SNP z were
511 assigned assuming an allele frequency of 0.5 for both populations and drawing from two binomial
512 trials. Genotypes at the m interacting SNPs were assigned in the same manner as the 2nd SNP in the
513 pairwise interaction simulations. Using the simulated genotypes, phenotypes were drawn from a
514 normal distribution, $\mathcal{N}(\mu = \sum_{x=1}^m \beta_{G \times G} \times g_{iz} \times g_x, \sigma = 1)$ where g_{is} is the genotype for individual i at
515 the simulated SNP s .

516 To simulate the scenario of differential LD on different ancestral backgrounds leading to
517 false positives, we simulated phenotypes based on a single causal SNP that was tagged by another
518 SNP. At both SNPs, local ancestries were assigned as described previously and genotypes were drawn
519 using ancestry specific allele frequencies. Ancestral allele frequencies were assigned such that the
520 average r^2 between the causal and tag SNP was 0.272 on the background of ancestral population 1
521 and 0.024 on the background of ancestral population 2. Thus, the tag SNP was only a tag on the
522 population 1 background and not on the population 2 background. Phenotypes were drawn from a
523 normal distribution, $\mathcal{N}(\mu = \beta_{Causal} \times g_{ic}, \sigma = 1)$, assuming no interaction and $\beta_{Causal} = 0.7$, where
524 g_{ic} is the genotype of individual i at the causal variant.

525

526 We implemented our approach in an R package (GxTheta), which is available for download
527 at <http://www.scandb.org/newinterface/GxTheta.html>

528

529

530 **Ancestry Inference**

531 Global ancestry inference was done using ADMIXTURE [9] and local ancestry inference was
532 done using LAMP-LD [55]. CEU and YRI from 1000 Genomes Phase 3 [56] were used as the European
533 and African reference panels. For the Native American reference panels, 95 Native Americans
534 genotyped on the Axiom LAT1 array were used[57].

535 **Filtering for Related Individuals**

536 All analyses in real data were filtered for related individuals due to the possibility of cryptic
537 relatedness causing false positives. To filter for related individuals, we estimated kinship coefficients
538 between all pairs of individuals using REAP [58]. We defined two individuals as related if they had a
539 kinship coefficient greater than 0.025. For a pair of related individuals, we removed the one with a
540 greater number of other individuals to whom he or she was related. In the case of a tie, we removed
541 one of the pair at random.

542 **Data Normalization**

543 *Gene Expression Normalization*

544 Gene expression data (see Results) were first standardized for each gene such that mean
545 expression was 0 and variance was 1. We then computed a covariance matrix of individual's
546 expression values and performed PCA on the covariance matrix. Residuals were computed for all
547 expression values by adjusting for the top 10 principal components and the mean for each gene was
548 added back to the residuals. Due to the high dynamic range of gene expression compared to
549 methylation we conservatively chose to additionally perform quantile normalization. We then sorted
550 the gene expression residuals and used the quantiles of their rank order to draw new expression
551 values from a normal distribution, $\mathcal{N}(\mu = 0, \sigma = 1)$, by using the inverse cumulative density
552 function^{24,25}.

553

554 *Methylation Data Normalization*

555 Raw methylation values (see Results) were first normalized using Illumina's control probe
556 scaling procedures. All probes with median methylation less than 1% or greater than 99% were
557 removed and the remaining probes were logit-transformed as previously described[59]. To control
558 for extreme outliers, we truncated the distribution of methylation values. For a given probe, we first
559 computed the mean and standard deviation of the methylation values. We then set any methylation
560 values deviating more than 2.58 standard deviations from the mean to the methylation value
561 corresponding to the 99.5th quantile.

562

563 **Availability of Supporting Data**

564 The Coriell data is available from dbGAP under accession number phs000211.v1.p1. The GALA
565 and SAGE data is available by emailing the study organizers at <https://pharm.ucsf.edu/gala/contact>.
566

567 **Competing Interests**

568 The authors declare that they have no competing interests.
569

570 **Authors' Contributions**

571 DSP, IE, EK, EE, EH and NZ designed research. DSP, IE, EK, ERG, and NZ performed research. DSP, IE,
572 EK, EE, CE, CRG, JMG, EG, HA, CJY, EE, EH, and NZ contributed new reagents/analytic tools. DSP, ERG,
573 and NZ wrote the manuscript. All authors read and approved the final manuscript.
574

575 **Description of Additional Data Files**

576 The following data are available with the online version of this paper. The Supplemental contains QQ-
577 plots for the simulations and real analyses performed as well as a table containing p-values for the 2-
578 component ancestry analysis of the GALA methylation data.
579

580 **Acknowledgements**

581 We would like to thank Lancelote Leong for his helpful manuscript comments.
582

583

584 **References**

- 585 1. Hemani G, Shakhsbazov K, Westra H-J, Esko T, Henders AK, Mcrae AF, et al. Detection and
586 replication of epistasis influencing transcription in humans. *Nature*. Nature Publishing Group;
587 2014 Apr 10;508(7495):249–53.
- 588 2. Rouhani F, Kumasaka N, de Brito MC, Bradley A, Vallier L, Gaffney D. Genetic Background
589 Drives Transcriptional Variation in Human Induced Pluripotent Stem Cells. Gibson G, editor.
590 *PLoS Genet*. 2014;10(6):e1004432.
- 591 3. Kang EY, Han B, Furlotte N, Joo JWJ, Shih D, Davis RC, et al. Meta-Analysis Identifies Gene-by-
592 Environment Interactions as Demonstrated in a Study of 4,965 Mice. Gibson G, editor. *PLoS*
593 *Genet*. Public Library of Science; 2014 Jan 9;10(1):e1004022.
- 594 4. Jemal A, Bray F, Center MM, Ferlay J, Ward E, Forman D. Global cancer statistics. *CA Cancer J*
595 *Clin*. 2011 Mar;61(2):69–90.
- 596 5. Lee M, Raj T, Castillo IW. ImmVar Project: Genetic architecture of leukocyte gene expression
597 in healthy humans. *JOURNAL OF ...*; 2012.
- 598 6. Manolio TA, Collins FS, Cox NJ, Goldstein DB, Hindorff LA, Hunter DJ, et al. Finding the missing
599 heritability of complex diseases. *Nature*. Nature Publishing Group; 2009 Oct
600 8;461(7265):747–53.
- 601 7. Eichler EE, Flint J, Gibson G, Kong A, Leal SM. Missing heritability and strategies for finding the
602 underlying causes of complex disease. *Nature Reviews* 2010.
- 603 8. Seldin MF, Pasaniuc B, Price AL. New approaches to disease mapping in admixed populations.
604 *Nature Reviews Genetics*. Nature Publishing Group; 2011 Aug 1;12(8):523–8.
- 605 9. Alexander DH, Novembre J, Lange K. Fast model-based estimation of ancestry in unrelated
606 individuals. *Genome Res*. Cold Spring Harbor Lab; 2009 Sep 1;19(9):1655–64.
- 607 10. Burchard EG, Ziv E, Coyle N, Gomez SL. The importance of race and ethnic background in
608 biomedical research and clinical practice. *New England Journal* 2003.
- 609 11. Kumar R, Seibold MA, Aldrich MC, Williams LK, Reiner AP, Colangelo L, et al. Genetic Ancestry
610 in Lung-Function Predictions. *N Engl J Med*. 2010 Jul 22;363(4):321–30.
- 611 12. Kumar R, Nguyen EA, Roth LA, Oh SS, Gignoux CR, Huntsman S, et al. Factors associated with
612 degree of atopy in Latino children in a nationwide pediatric sample: The Genes-environments
613 and Admixture in Latino Asthmatics (GALA II) study. *J Allergy Clin Immunol*. Elsevier; 2013
614 May 16;132(4):896–905.e1.
- 615 13. Price AL, Patterson N, Hancks DC, Myers S, Reich D, Cheung VG, et al. Effects of cis and trans
616 Genetic Ancestry on Gene Expression in African Americans. Gibson G, editor. *PLoS Genet*.
617 Public Library of Science; 2008 Dec 5;4(12):e1000294.
- 618 14. Shaffer JR, Kammerer CM, Reich D, McDonald G, Patterson N, Goodpaster B, et al. Genetic
619 markers for ancestry are correlated with body composition traits in older African Americans.
620 *Osteoporos Int*. Springer-Verlag; 2007;18(6):733–41.
- 621 15. Florez JC, Price AL, Campbell D, Riba L, Parra MV, Yu F, et al. Strong Association of
622 Socioeconomic Status and Genetic Ancestry in Latinos: Implications for Admixture Studies of

- 623 Type 2 Diabetes. In: Racial Identities, Genetic Ancestry, and Health in South America. Palgrave
624 Macmillan US; 2011. pp. 137–53.
- 625 16. Reiner AP, Carlson CS, Ziv E, Iribarren C, Jaquish CE, Nickerson DA. Genetic ancestry,
626 population sub-structure, and cardiovascular disease-related traits among African-American
627 participants in the CARDIA Study. *Hum Genet*. Springer-Verlag; 2007;121(5):565–75.
- 628 17. Sanchez E, Webb RD, Rasmussen A, Kelly JA, Riba L, Kaufman KM, et al. Genetically
629 determined Amerindian ancestry correlates with increased frequency of risk alleles for
630 systemic lupus erythematosus. *Arthritis & Rheumatism*. Wiley Subscription Services, Inc., A
631 Wiley Company; 2010 Dec 1;62(12):3722–9.
- 632 18. Ziv E, John EM, Choudhry S, Kho J, Lorizio W, Pérez-Stable EJ, et al. Genetic Ancestry and Risk
633 Factors for Breast Cancer among Latinas in the San Francisco Bay Area. *Cancer Epidemiol
634 Biomarkers Prev*. American Association for Cancer Research; 2006 Oct 1;15(10):1878–85.
- 635 19. Cheng C-Y, Reich D, Haiman CA, Tandon A, Patterson N, Elizabeth S, et al. African Ancestry and
636 Its Correlation to Type 2 Diabetes in African Americans: A Genetic Admixture Analysis in
637 Three U.S. Population Cohorts. Atkin SL, editor. PLoS ONE. Public Library of Science; 2012
638 Mar 16;7(3):e32840.
- 639 20. Choudhry S, Burchard EG, Borrell LN, Tang H, Gomez I, Naqvi M, et al. Ancestry–Environment
640 Interactions and Asthma Risk among Puerto Ricans. *Am J Respir Crit Care Med* [Internet].
641 American Thoracic Society; 2006 Nov 15;174(10):1088–93. Available from:
642 <http://www.atsjournals.org/doi/abs/10.1164/rccm.200605-5960C>
- 643 21. Spielman RS, Bastone LA, Burdick JT, Morley M, Ewens WJ, Cheung VG. Common genetic
644 variants account for differences in gene expression among ethnic groups. *Nat Genet*. Nature
645 Publishing Group; 2007 Feb 1;39(2):226–31.
- 646 22. Peralta CA, Risch N, Lin F, Shlipak MG, Reiner A, Ziv E, et al. The Association of African
647 Ancestry and Elevated Creatinine in the Coronary Artery Risk Development in Young Adults
648 (CARDIA) Study. *Am J Nephrol*. Karger Publishers; 2009 Dec 21;31(3):202–8.
- 649 23. Galanter JM, Gignoux CR, Oh SS, Torgerson D, Pino-Yanes M, Thakur N, et al. Methylation
650 Analysis Reveals Fundamental Differences Between Ethnicity and Genetic Ancestry. *bioRxiv*.
651 Cold Spring Harbor Labs Journals; 2016 Jan 15;:036822.
- 652 24. Population-specificity of human DNA methylation. 2012.
- 653 25. Delude CM. Deep phenotyping: The details of disease. *Nature*. Nature Publishing Group; 2015
654 Nov 5;527(7576):S14–5.
- 655 26. Vimalesarwan KS, Berry DJ, Lu C, Tikkanen E, Pilz S, Hiraki LT, et al. Causal Relationship
656 between Obesity and Vitamin D Status: Bi-Directional Mendelian Randomization Analysis of
657 Multiple Cohorts. Minelli C, editor. PLOS Med. Public Library of Science; 2013 Feb
658 5;10(2):e1001383.
- 659 27. Grundberg E, Small KS, Hedman ÅK, Nica AC, Buil A, Keildson S, et al. Mapping cis- and trans-
660 regulatory effects across multiple tissues in twins. *Nat Genet*. Nature Publishing Group; 2012
661 Oct 1;44(10):1084–9.
- 662 28. Aschard H, Lutz S, Maus B, Duell EJ, Fingerlin TE, Chatterjee N, et al. Challenges and
663 opportunities in genome-wide environmental interaction (GWEI) studies. *Hum Genet*.

- 664 Springer-Verlag; 2012;131(10):1591–613.
- 665 29. Aschard H, Gusev A, Brown R, Pasaniuc B. Leveraging local ancestry to detect gene-gene
666 interactions in genome-wide data. *BMC Genetics*. BioMed Central Ltd; 2015 Oct 24;16(1):124.
- 667 30. De R, Hu T, Moore JH, Gilbert-Diamond D. Characterizing gene-gene interactions in a
668 statistical epistasis network of twelve candidate genes for obesity. *BioData Mining*. BioMed
669 Central; 2015;8(1):1–16.
- 670 31. Ritchie MD, Hahn LW, Roodi N, Bailey LR, Dupont WD, Parl FF, et al. Multifactor-
671 Dimensionality Reduction Reveals High-Order Interactions among Estrogen-Metabolism
672 Genes in Sporadic Breast Cancer. *The American Journal of Human Genetics*. 2001
673 Jul;69(1):138–47.
- 674 32. Wong MY, Day NE, Luan JA, Chan KP, Wareham NJ. The detection of gene-environment
675 interaction for continuous traits: should we deal with measurement error by bigger studies or
676 better measurement? *Int J Epidemiol*. 2003 Feb;32(1):51–7.
- 677 33. Devlin B, Roeder K. Genomic Control for Association Studies. *Biometrics* [Internet]. Blackwell
678 Publishing Ltd; 2004 May 25;55(4):997–1004. Available from:
679 <http://doi.wiley.com/10.1111/j.0006-341X.1999.00997.x>
- 680 34. Simon-Sanchez J, Scholz S, Fung H-C, Matarin M, Hernandez D, Gibbs JR, et al. Genome-wide
681 SNP assay reveals structural genomic variation, extended homozygosity and cell-line induced
682 alterations in normal individuals. *Hum Mol Genet*. Oxford University Press; 2007 Jan
683 1;16(1):1–14.
- 684 35. Martin AR, Costa HA, Lappalainen T, Henn BM, Kidd JM, Yee M-C, et al. Transcriptome
685 Sequencing from Diverse Human Populations Reveals Differentiated Regulatory Architecture.
686 Gibson G, editor. PLoS Genet. Public Library of Science; 2014 Aug 14;10(8):e1004549.
- 687 36. Borrell LN, Nguyen EA, Roth LA, Oh SS, Tcheurekdjian H, Sen S, et al. Childhood Obesity and
688 Asthma Control in the GALA II and SAGE II Studies. dx.doi.org. American Thoracic Society;
689 2013. 6 p.
- 690 37. Torgerson DG, Gignoux CR, Galanter JM, Drake KA, Roth LA, Eng C, et al. Case-control
691 admixture mapping in Latino populations enriches for known asthma-associated genes. *J*
692 *Allergy Clin Immunol*. 2012 Jul;130(1):76–82.e12.
- 693 38. Eichler EE, Flint J, Gibson G, Kong A, Leal SM, Moore JH, et al. Missing heritability and
694 strategies for finding the underlying causes of complex disease. *Nature Reviews Genetics*.
695 Nature Publishing Group; 2010 Jun 1;11(6):446–50.
- 696 39. Powell JE, Henders AK, Mcrae AF, Kim J, Hemani G, Martin NG, et al. Congruence of Additive
697 and Non-Additive Effects on Gene Expression Estimated from Pedigree and SNP Data. Spector
698 TD, editor. PLoS Genet. Public Library of Science; 2013 May 16;9(5):e1003502.
- 699 40. Keller MC. Gene × Environment Interaction Studies Have Not Properly Controlled for
700 Potential Confounders: The Problem and the (Simple) Solution. *Biological Psychiatry*. 2014
701 Jan;75(1):18–24.
- 702 41. Sul JH, Bilow M, Yang W-Y, Kostem E, Furlotte N, He D, et al. Accounting for Population
703 Structure in Gene-by-Environment Interactions in Genome-Wide Association Studies Using
704 Mixed Models. Schork NJ, editor. PLoS Genet. Public Library of Science; 2016 Mar

- 705 4;12(3):e1005849.
- 706 42. Kooperberg C, LeBlanc M. Increasing the power of identifying gene × gene interactions in
707 genome-wide association studies. *Genet Epidemiol*. Wiley Subscription Services, Inc., A Wiley
708 Company; 2008 Apr 1;32(3):255–63.
- 709 43. Murcray CE, Lewinger JP, Gauderman WJ. Gene-Environment Interaction in Genome-Wide
710 Association Studies. *Am J Epidemiol*. Oxford University Press; 2009 Jan 15;169(2):219–26.
- 711 44. Hsu L, Jiao S, Dai JY, Hutter C, Peters U, Kooperberg C. Powerful Cocktail Methods for
712 Detecting Genome-Wide Gene-Environment Interaction. *Genet Epidemiol*. 2012 Apr
713 1;36(3):183–94.
- 714 45. Yang J, Benyamin B, McEvoy BP, Gordon S, Henders AK, Nyholt DR, et al. Common SNPs
715 explain a large proportion of the heritability for human height. *Nat Genet*. 2010 Jun
716 20;42(7):565–9.
- 717 46. Smith ZD, Meissner A. DNA methylation: roles in mammalian development. *Nature Reviews
718 Genetics*. Nature Publishing Group; 2013 Mar 1;14(3):204–20.
- 719 47. Artigas MS, Loth DW, Wain LV, Gharib SA, Obeidat M, Tang W, et al. Genome-wide association
720 and large-scale follow up identifies 16 new loci influencing lung function. *Nat Genet*. Nature
721 Publishing Group; 2011 Nov 1;43(11):1082–90.
- 722 48. Liao SY, Lin X, Christiani DC. Gene-environment interaction effects on lung function-a
723 genome-wide association study within the Framingham heart study. *Environ Health*. 2013.
- 724 49. Anderson LN, Briollais L, Atkinson HC, Marsh JA, Xu J, Connor KL, et al. Investigation of
725 Genetic Variants, Birthweight and Hypothalamic-Pituitary-Adrenal Axis Function Suggests a
726 Genetic Variant in the SERPINA6 Gene Is Associated with Corticosteroid Binding Globulin in
727 the Western Australia Pregnancy Cohort (Raine) Study. Hsu Y-H, editor. PLoS ONE. Public
728 Library of Science; 2014 Apr 1;9(4):e92957.
- 729 50. Baran Y, Pasaniuc B, Sankararaman S, Torgerson DG, Gignoux C, Eng C, et al. Fast and accurate
730 inference of local ancestry in Latino populations. *Bioinformatics*. Oxford University Press;
731 2012 May 15;28(10):1359–67.
- 732 51. Sankararaman S, Sridhar S, Kimmel G. Estimating local ancestry in admixed populations. *The
733 American Journal of* 2008.
- 734 52. Price AL, Patterson N, Yu F, Cox DR, Waliszewska A, McDonald GJ, et al. A Genomewide
735 Admixture Map for Latino Populations. *The American Journal of Human Genetics*. 2007
736 Jun;80(6):1024–36.
- 737 53. Smith MW, Patterson N, Lautenberger JA, Truelove AL, McDonald GJ, Waliszewska A, et al. A
738 High-Density Admixture Map for Disease Gene Discovery in African Americans. *The American
739 Journal of Human Genetics*. 2004 May;74(5):1001–13.
- 740 54. Balding DJ, Nichols RA. A method for quantifying differentiation between populations at
741 multi-allelic loci and its implications for investigating identity and paternity. *Human
742 Identification: The Use of DNA Markers*. 1995.
- 743 55. Baran Y, Pasaniuc B, Sankararaman S, Torgerson DG, Gignoux C, Eng C, et al. Fast and accurate
744 inference of local ancestry in Latino populations. *Bioinformatics*. 2012 May 15;28(10):1359–

- 745 67.
- 746 56. Consortium T1GP. An integrated map of genetic variation from 1,092 human genomes.
747 Nature. Nature Publishing Group; 2012 Nov 1;491(7422):56–65.
- 748 57. Drake KA, Torgerson DG, Gignoux CR, Galanter JM, Roth LA, Huntsman S, et al. A genome-wide
749 association study of bronchodilator response in Latinos implicates rare variants. Journal of
750 Allergy and Clinical Immunology. 2014 Feb;133(2):370–378.e15.
- 751 58. Thornton T, Tang H, Hoffmann TJ, Ochs-Balcom HM, Caan BJ, Risch N. Estimating kinship in
752 admixed populations. Am J Hum Genet. 2012 Jul 13;91(1):122–38.
- 753 59. Du P, Zhang X, Huang C-C, Jafari N, Kibbe WA, Hou L, et al. Comparison of Beta-value and M-
754 value methods for quantifying methylation levels by microarray analysis. BMC
755 Bioinformatics. BioMed Central Ltd; 2010 Nov 30;11(1):587.
- 756
- 757

758 **Figure Legends**

759 Figure 1. Power Plots for Pairwise Interaction Simulations.

760 Power of testing $G \times \theta$ (a) versus testing pairwise SNPs directly (b) as a function of the difference in
761 the ancestral allele frequencies at a differentiated SNP.

762

763 Figure 2. Power Plots for Multi-way Pairwise Interaction Simulations

764 Power of testing $G \times \theta$ as a function of the difference in the ancestral allele frequencies for multiple
765 interacting SNPs.

766

767 Figure 3. Power Plots for $G \times E$ Interaction Simulations.

768 Power of testing $G \times \theta$ as a function of the correlation between an environmental covariate and
769 genetic ancestry.

770

771 Figure 4. Examples of How Genetic Ancestry Can Be A Proxy for Interacting Covariates.

772 (a) Model of how genetic ancestry θ can be correlated with various environmental exposures, some
773 of which affect a phenotype. (b) Example of how the correlation between the probability of an AA
774 genotype (bars 2-4) and values of θ (bar 1) increase with higher levels of SNP allele frequency
775 differentiation. In this plot p_1 and p_2 denote the allele frequency of allele A in ancestral populations 1
776 and 2 respectively. (c) Example of how effect sizes at a tag-SNP may differ due to differential LD on
777 distinct ancestral backgrounds (here, EUR and AFR).

778

779

780 **Tables**

781

782

Table 1. Distribution of Ancestry in Coriell and GALA II.

Dataset	θ_{EUR}	θ_{AFR}	θ_{NAM}
Coriell	$\mu=0.212, \sigma=0.021$	$\mu=0.788, \sigma=0.021$	NA
GALA II MX	$\mu=0.396, \sigma=0.149$	$\mu=0.043, \sigma=0.025$	$\mu=0.561, \sigma=0.159$
GALA II PR	$\mu=0.641, \sigma=0.094$	$\mu=0.246, \sigma=0.101$	$\mu=0.113, \sigma=0.024$

783 Mean and variance of the global ancestry distributions for each dataset.

784

785

Table 2. GALA II DNA Methylation Analysis Results.

GALA II Population	Probe Gene	Probe ID	rsid	Distance of SNP to Probe	Marginal p-value	AITL p-value	AITL Replication p-value
MX	CNFN	cg14327995	rs16975986	280795	2.49E-09	5.69E-09	9.27E-03
MX	C11orf95	cg16678159	rs7106153	249768	2.58E-01	2.52E-08	9.39E-02
MX	NA	cg05697734	rs1560919	13711	1.14E-01	2.21E-08	8.18E-03
MX	TNK2	cg01792640	rs67217828	278866	4.49E-01	6.38E-09	1.43E-02
MX	HDAC4	cg06533788	rs925736	9548	4.51E-01	3.09E-09	2.80E-02
PR	NA	cg07436864*	rs8117083	31813	7.46E-02	1.34E-09	5.34E-03
PR	NA	cg16803083*	rs4312379	63847	3.69E-01	2.29E-08	2.31E-04
PR	SERPINA6	cg10025865	rs17091085	247796	6.83E-01	2.97E-08	8.05E-03

786 P-values for AITL applied to the methylation data in the GALA II Latinos. MX and PR denote Mexicans
787 and Puerto Ricans respectively in the GALA II population columns. The probe gene column shows the
788 gene that the methylation probe lies in. The marginal column is the p-value for standard linear
789 regression of methylation on genotype while controlling for population structure. * indicates results
790 that replicated between the Mexicans and Puerto Ricans.

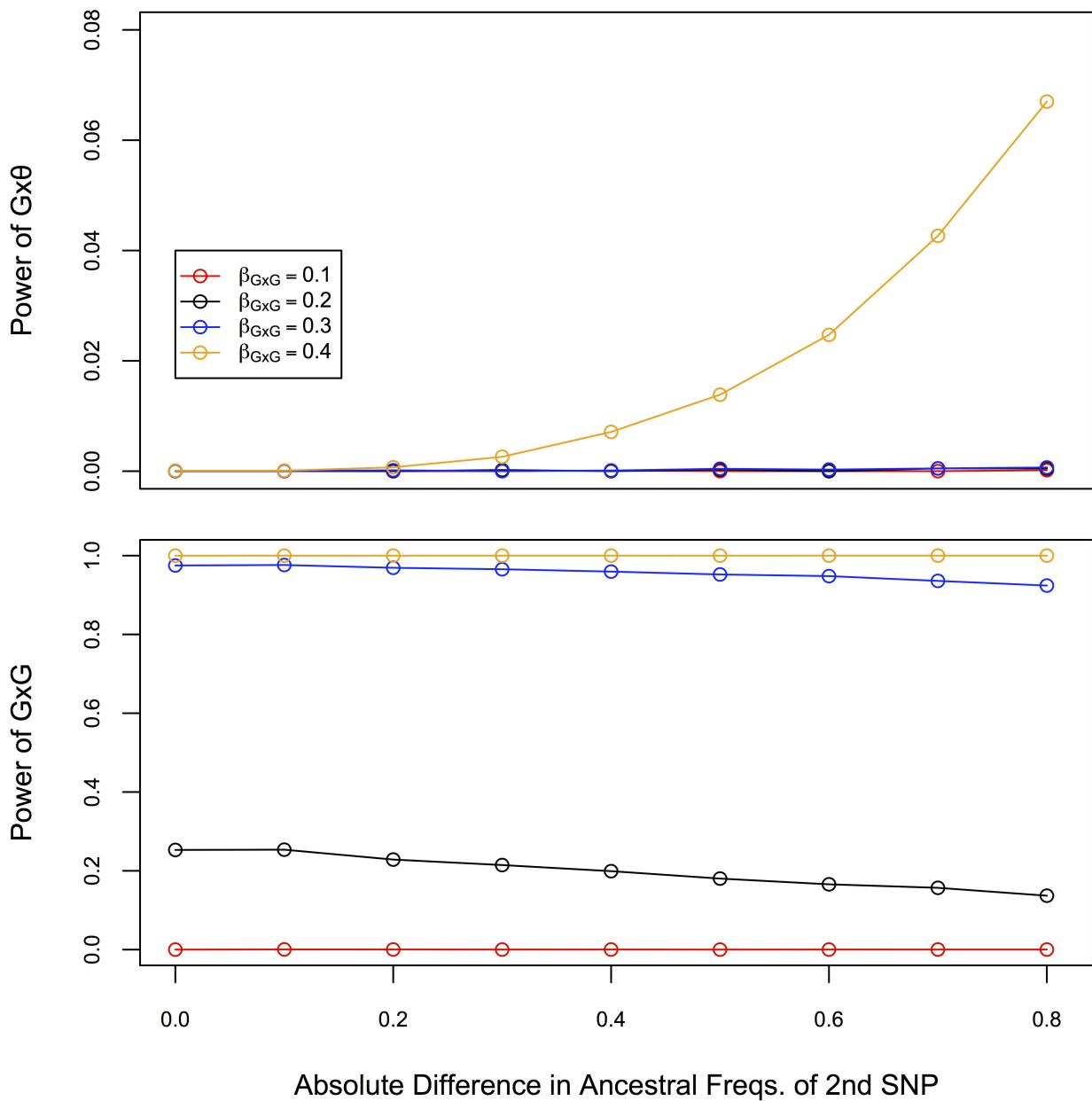
791

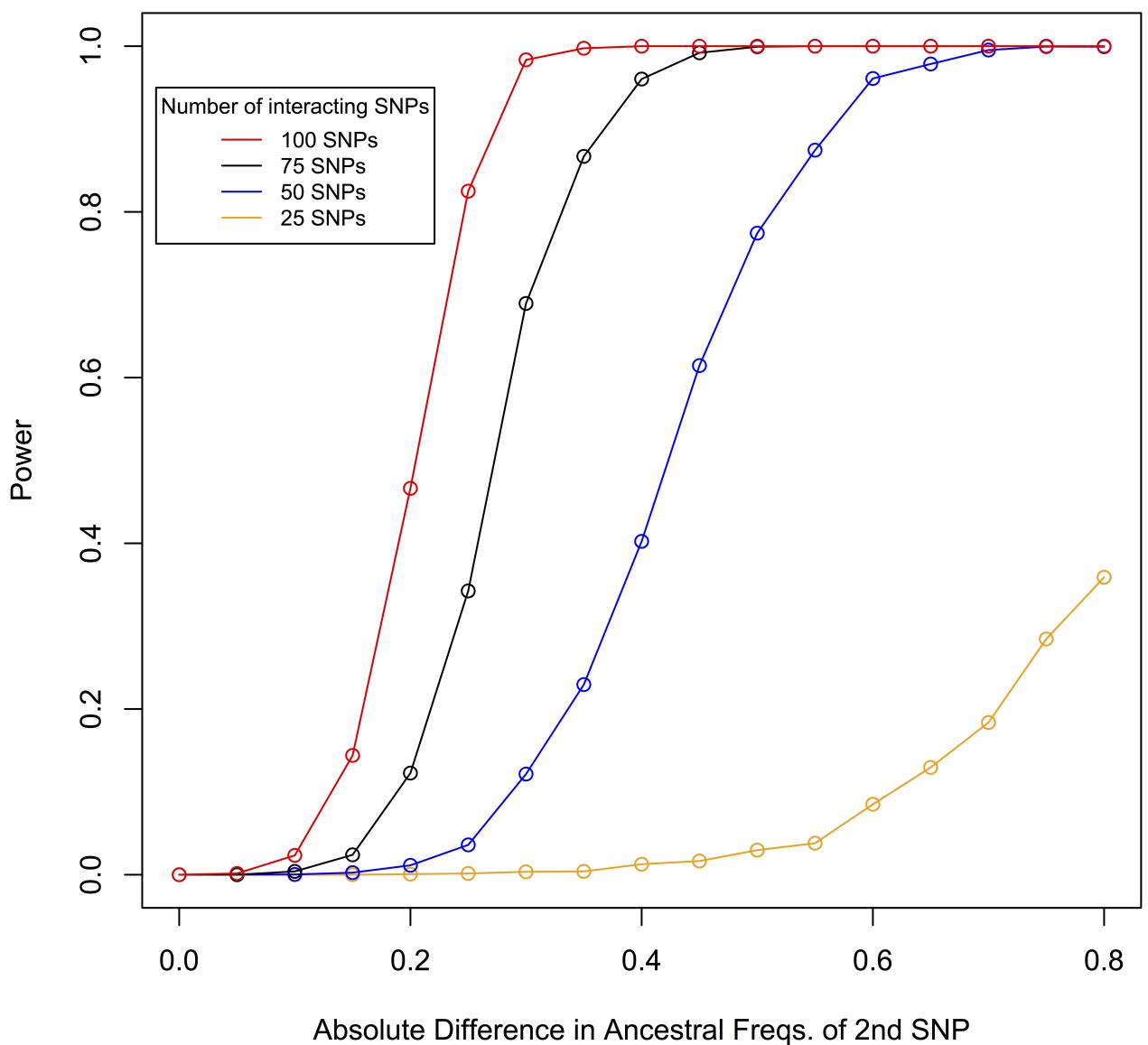
792

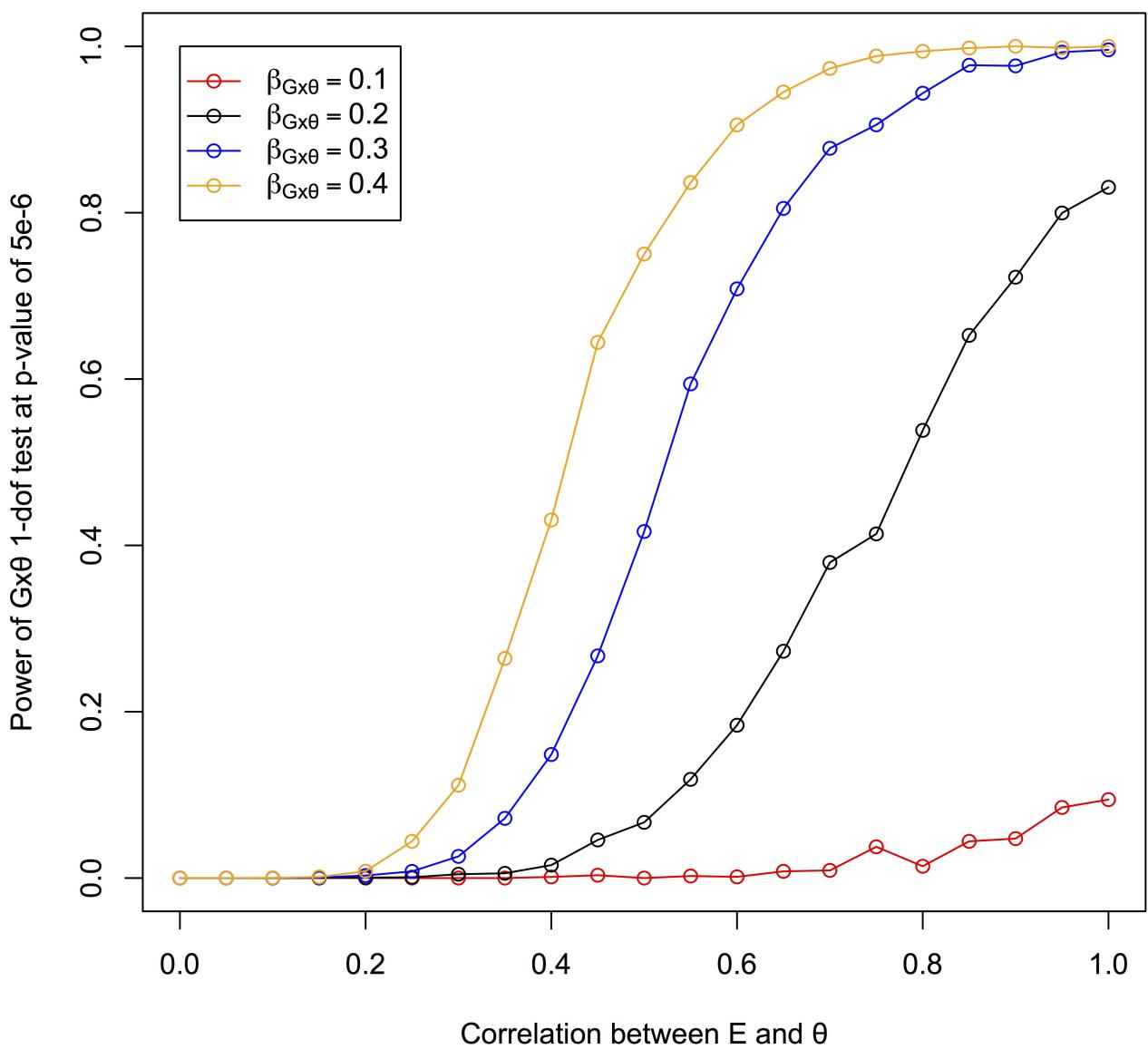
793

794

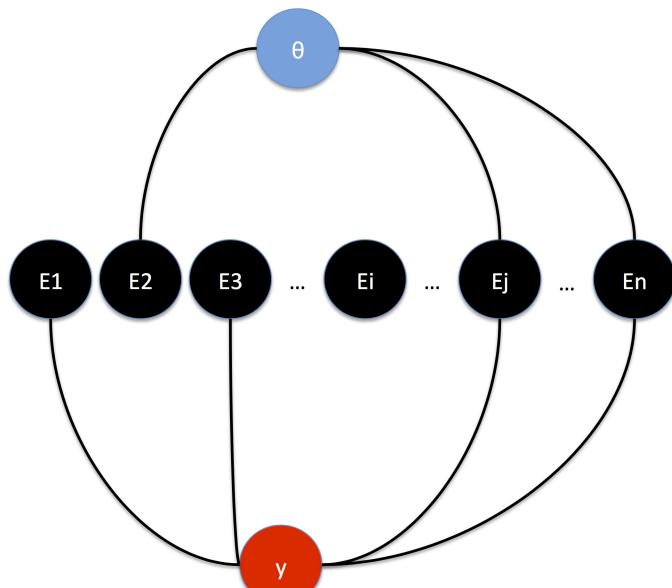
795



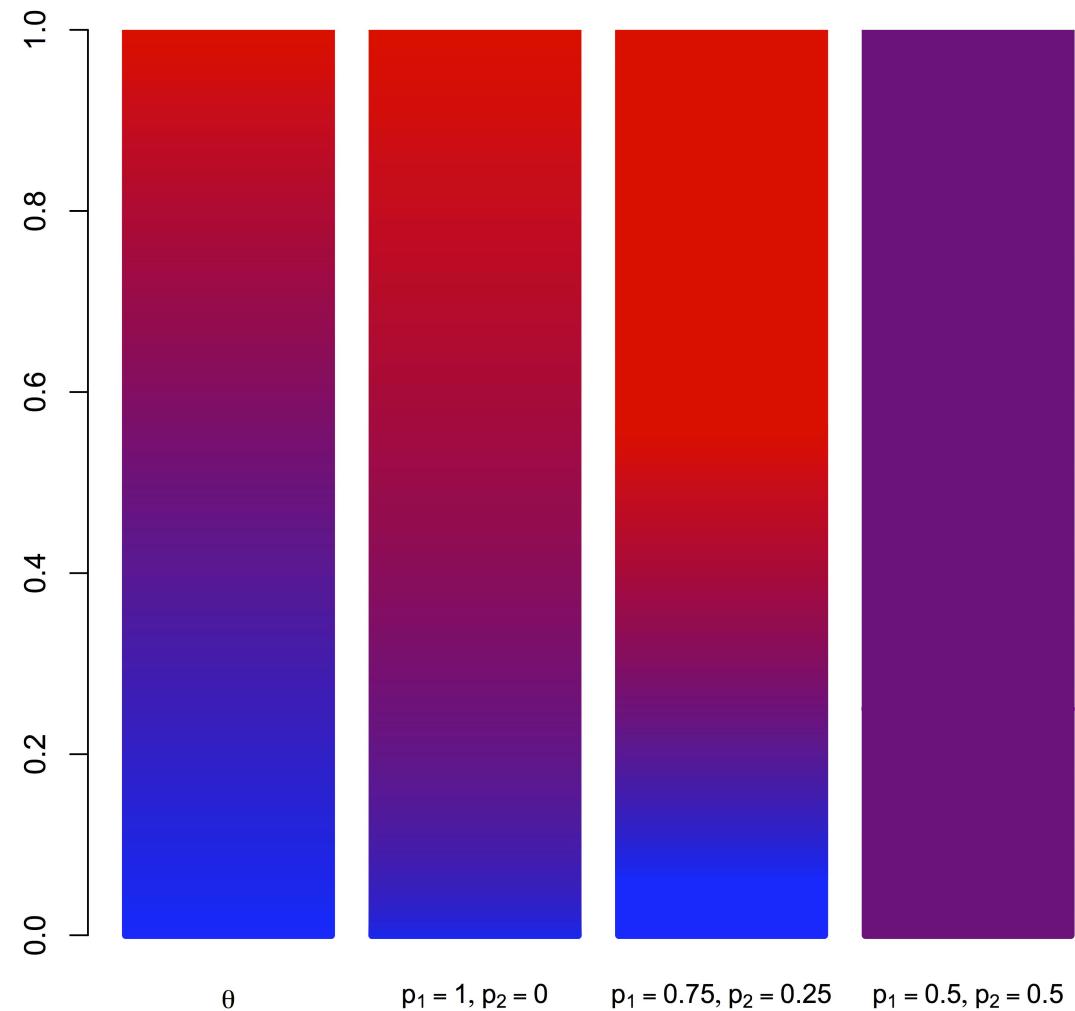




(a)



(b)



(c)

