

# Interpreting SNP heritability in admixed populations

Jinguo Huang<sup>1,2</sup>, Nicole Kleman<sup>3</sup>, Saonli Basu<sup>4</sup>, Mark D. Shriver<sup>2</sup>, and  
Arslan A. Zaidi<sup>†3,5</sup>

<sup>1</sup>Bioinformatics and Genomics, Huck Institutes of the Life Sciences, Pennsylvania State University

<sup>2</sup>Department of Anthropology, Pennsylvania State University

<sup>3</sup>Department of Genetics, Cell Biology, and Development, University of Minnesota

<sup>4</sup>Department of Biostatistics, University of Minnesota

<sup>5</sup>Institute of Health Informatics, University of Minnesota

<sup>†</sup>Correspondence to A.A.Z (aazaidi@umn.edu)

August 2, 2024

## Abstract

SNP heritability ( $h_{snp}^2$ ) is defined as the proportion of phenotypic variance explained by genotyped SNPs and is believed to be a lower bound of heritability ( $h^2$ ), being equal to it if all causal variants are known. Despite the simple intuition behind  $h_{snp}^2$ , its interpretation and equivalence to  $h^2$  is unclear, particularly in the presence of population structure and assortative mating. It is well known that population structure can lead to inflation in  $\hat{h}_{snp}^2$  estimates because of confounding due to linkage disequilibrium (LD) or shared environment. Here we use analytical theory and simulations to demonstrate that  $h_{snp}^2$  estimates can be biased in admixed populations, even in the absence of confounding and even if all causal variants are known. This is because admixture generates LD, which contributes to the genetic variance, and therefore to heritability. Genome-wide restricted maximum likelihood (GREML) does not capture this contribution leading to under- or over-estimates of  $h_{snp}^2$  relative to  $h^2$ , depending on the genetic architecture. In contrast, Haseman-Elston (HE) regression exaggerates the LD contribution leading to biases in the opposite direction. For the same reason, GREML and HE estimates of local ancestry heritability ( $h_{\gamma}^2$ ) are also biased. We describe this bias in  $\hat{h}_{snp}^2$  and  $\hat{h}_{\gamma}^2$  as a function of admixture history and the genetic architecture of the trait and show that it can be recovered under some conditions. We clarify the interpretation of  $\hat{h}_{snp}^2$  in admixed populations and discuss its implication for genome-wide association studies and polygenic prediction.

## 28 Introduction

29 The ability to estimate (narrow-sense) heritability ( $h^2$ ) from unrelated individuals was a major advance in  
30 genetics. Traditionally,  $h^2$  was estimated from family-based studies in which the phenotypic resemblance  
31 between relatives could be modeled as a function of their expected genetic relatedness [1]. However, this  
32 approach was limited to analysis of closely related individuals where pedigree information is available and  
33 the realized genetic relatedness is not too different from expectation [2]. With the advent of genome-wide  
34 association studies (GWAS), we hoped that many of the variants underlying this heritability would be  
35 uncovered. However, when genome-wide significant SNPs explained a much smaller fraction of the phe-  
36 notypic variance, it became important to explain the missing heritability – were family-based estimates  
37 inflated or were GWAS just underpowered, limited by variant discovery?

38 Yang *et al.* (2010) [3] made the key insight that one could estimate the portion of  $h^2$  tagged by  
39 genotyped SNPs, regardless of whether or not they were genome-wide significant, by exploiting the  
40 subtle variation in the realized genetic relatedness among apparently unrelated individuals [3–5]. This  
41 quantity came to be known colloquially as ‘SNP heritability’ ( $h_{snp}^2$ ) and it is believed to be equal to  
42  $h^2$  if all causal variants are included among genotyped SNPs [3]. Indeed, estimates of  $h_{snp}^2$  explain a  
43 much larger fraction of trait heritability than GWAS SNPs [3], approaching family-based estimates of  
44  $h^2$  when whole genome sequence data, which captures rare variants, are used [6]. This has made it clear  
45 that GWAS have yet to uncover more variants with increasing sample size. Now,  $h_{snp}^2$  has become an  
46 important aspect of the design of genetic studies and is often used to define the power of variant discovery  
47 in GWAS and the upper limit of polygenic prediction accuracy.

48 Despite the utility and simple intuition of  $h_{snp}^2$ , there is much confusion about its interpretation and  
49 equivalence to  $h^2$ , particularly in the presence of population structure and assortative mating [7–12].  
50 But much of the discussion of heritability in structured populations has focused on biases in  $\hat{h}_{snp}^2$  – the  
51 estimator – due to confounding effects of shared environment and linkage disequilibrium (LD) with other  
52 variants [7, 9–11, 13]. There is comparatively little discussion, at least in human genetics, on the fact  
53 that LD due to population structure also contributes to genetic variance, and therefore, is a component  
54 of heritability [1] (but see [14–16] for a rigorous discussion). We think this is at least partly due to the  
55 fact that most studies are carried out in cohorts with primarily European ancestry, where the degree of  
56 population structure is minimal and large effects of LD can be ignored. However, that is not the case  
57 for diverse, multi-ethnic cohorts, which have historically been underrepresented in genetic studies, but  
58 thanks to a concerted effort in the field, are now becoming increasingly common [17–23]. The complex  
59 structure in these cohorts also brings unique methodological challenges and it is imperative that we  
60 understand whether existing methods, which have largely been evaluated in more homogeneous groups,  
61 generalize to more diverse cohorts.

62 Our goal in this paper is to study the behavior of  $\hat{h}_{snp}^2$  in admixed populations. What is its inter-  
63 pretation in the ideal situation where causal variants are known? Is it an unbiased estimate of  $h^2$ ? To  
64 answer these questions, we derived a general expression for the genetic variance in admixed populations,  
65 decomposing it in terms of the contribution of population structure, which influences both the genotypic  
66 variance at individual loci and the LD across loci. We used theory and simulations to show that  $\hat{h}_{snp}^2$

67 estimated with genome-wide restricted maximum likelihood (GREML) [3, 5] and Haseman-Elston (HE)  
 68 regression [24] – two widely used approaches – can be biased in admixed and other structured popula-  
 69 tions, even in the absence of confounding and when all causal variants are known. We explain this in  
 70 terms of the discrepancy between the model assumed in  $\hat{h}_{snp}^2$  estimation and the generative model from  
 71 which the genetic architecture of the trait in the population may have been sampled. We describe the  
 72 bias in  $\hat{h}_{snp}^2$  as a function of admixture history and genetic architecture and discuss its implications for  
 73 GWAS and polygenic prediction accuracy.

## 74 Model

### 75 Genetic architecture

76 We begin by describing a generative model for the phenotype. Let  $y = g + e$ , where  $y$  is the phenotypic  
 77 value of an individual,  $g$  is the genotypic value, and  $e$  is random error. We assume additive effects such  
 78 that  $g = \sum_{i=1}^m \beta_i x_i$  where  $\beta_i$  is the effect size of the  $i^{th}$  biallelic locus and  $x_i \in \{0, 1, 2\}$  is the number  
 79 of copies of the trait-increasing allele. Importantly, the effect sizes are fixed quantities and differences in  
 80 genetic values among individuals are due to random variation in genotypes. Note, that this is different  
 81 from the model assumed by GREML where genotypes are fixed and effect sizes are random [14].

82 We denote the mean, variance, and covariance with  $\mathbb{E}(\cdot)$ ,  $\mathbb{V}(\cdot)$ , and  $\mathbb{C}(\cdot, \cdot)$ , respectively, where the  
 83 expectation is measured over random draws from the population rather than random realizations of the  
 84 evolutionary process. We can express the additive genetic variance of a quantitative trait as follows:

$$V_g = \mathbb{V}\left(\sum_{i=1}^m \beta_i x_i\right) = \sum_{i=1}^m \beta_i^2 \mathbb{V}(x_i) + \sum_{j \neq i} \beta_i \beta_j \mathbb{C}(x_i, x_j)$$

85 Here the first term represents the contribution of individual loci (genic variance) and the second term  
 86 is the contribution of linkage disequilibrium (LD contribution). We make the assumption that loci are  
 87 unlinked and therefore, the LD contribution is entirely due to population structure. We describe the  
 88 behavior of  $V_g$  in a population that is a mixture of two previously isolated populations A and B that  
 89 diverged from a common ancestor. To do this, we denote  $\theta$  as the fraction of the genome of an individual  
 90 with ancestry from population A. Thus,  $\theta = 1$  if the individual is from population A, 0 if they are from  
 91 population B, and  $\theta \in (0, 1)$  if they are admixed. Then,  $V_g$  can be expressed in terms of ancestry as  
 92 (Appendix):

$$V_g = 2 \mathbb{E}(\theta) \sum_{i=1}^m \beta_i^2 f_i^A (1 - f_i^A) + 2 \{1 - \mathbb{E}(\theta)\} \sum_{i=1}^m \beta_i^2 f_i^B (1 - f_i^B) \quad (1.1)$$

$$+ 2 \mathbb{E}(\theta) \{1 - \mathbb{E}(\theta)\} \sum_{i=1}^m \beta_i^2 (f_i^A - f_i^B)^2 \quad (1.2)$$

$$+ 2 \mathbb{V}(\theta) \sum_{i=1}^m \beta_i^2 (f_i^A - f_i^B)^2 \quad (1.3)$$

$$+ 4 \mathbb{V}(\theta) \sum_{i \neq j} \beta_i \beta_j (f_i^A - f_i^B) (f_j^A - f_j^B) \quad (1.4)$$

93 where  $f_i^A$  and  $f_i^B$  are the allele frequencies in populations A and B, and  $\mathbb{E}(\theta)$  and  $\mathbb{V}(\theta)$  are the mean  
94 and variance of individual ancestry. The sum of the first three terms represents the genic variance and  
95 the last term represents the LD contribution.

## 96 Demographic history

97 From Eq. 1, it is clear that, conditional on the genetic architecture in the source populations ( $\beta, f^A, f^B$ ),  
98  $V_g$  is a function of the mean,  $\mathbb{E}(\theta)$ , and variance,  $\mathbb{V}(\theta)$ , of individual ancestry in the admixed population.  
99 We consider two demographic models that affect  $\mathbb{E}(\theta)$  and  $\mathbb{V}(\theta)$  in qualitatively different ways. In the  
100 first model, the source populations meet once  $t$  generations ago (we refer to this as  $t = 0$ ) in proportions  
101  $p$  and  $1 - p$ , after which there is no subsequent admixture (Fig. 1A). In the second model, there is  
102 continued gene flow in every generation from one of the source populations such that the mean overall  
103 amount of ancestry from population A is the same as in the first model (Fig. 1A). For brevity, we refer  
104 to these as the hybrid-isolation (HI) and continuous gene flow (CGF) models, respectively, following  
105 Pfaff *et al.* (2001) [25].  $\mathbb{V}(\theta)$  is also affected by ancestry-based assortative mating, where individuals are  
106 more likely to partner with others of similar ancestry. We refer to this simply as assortative mating for  
107 brevity and model this following Zaitlen *et al.* (2017) using a parameter  $P \in (0, 1)$ , which represents the  
108 correlation of the ancestry of individuals across mating pairs in the population [26].

109 Under these conditions, the behavior of  $\mathbb{E}(\theta)$  and  $\mathbb{V}(\theta)$  has been described previously [26, 27] (Fig. 1B  
110 and C). Briefly, in the HI model,  $\mathbb{E}(\theta)$  remains constant at  $p$  in the generations after admixture as there  
111 is no subsequent gene flow.  $\mathbb{V}(\theta)$  is at its maximum at  $t = 0$  when each individual carries chromosomes  
112 either from population A or B, but not both. This genome-wide correlation in ancestry breaks down  
113 in subsequent generations as a function of mating, independent assortment, and recombination, leading  
114 to a decay in  $\mathbb{V}(\theta)$ , the rate depending on the strength of assortative mating (Fig. 1C). In the CGF  
115 model, both  $\mathbb{E}(\theta)$  and  $\mathbb{V}(\theta)$  increase with time as new chromosomes are introduced from the source  
116 populations. But while  $\mathbb{E}(\theta)$  continues to increase monotonically,  $\mathbb{V}(\theta)$  will plateau and decrease due to  
117 the countervailing effects of independent assortment and recombination which redistribute ancestry in  
118 the population, reaching equilibrium at zero if there is no more gene flow and the population is mating  
119 randomly.  $\mathbb{V}(\theta)$  provides an intuitive and quantitative measure of the degree of population structure  
120 (along the axis of ancestry) in admixed populations.

## 121 Results

### 122 Genetic variance in admixed populations

123 To understand the expectation of genetic variance in admixed populations, it is first worth discussing  
124 its behavior in the source populations. In Eq. 1, the first term represents the within-population com-  
125 ponent ( $V_{gw}$ ) and the last three terms altogether represent the component of genetic variance between  
126 populations A and B ( $V_{gb}$ ). Note that  $V_{gb} = \frac{(\bar{g}_A - \bar{g}_B)^2}{2}$  is positive only if there is a difference in the mean  
127 genotypic values (Fig. 2). This variance increases with genetic divergence since the expected values  
128 of both  $(f_i^A - f_i^B)^2$  and  $(f_i^A - f_i^B)(f_j^A - f_j^B)$  are functions of  $F_{ST}$ . While  $\beta_i^2 (f_i^A - f_i^B)^2$  is expected to

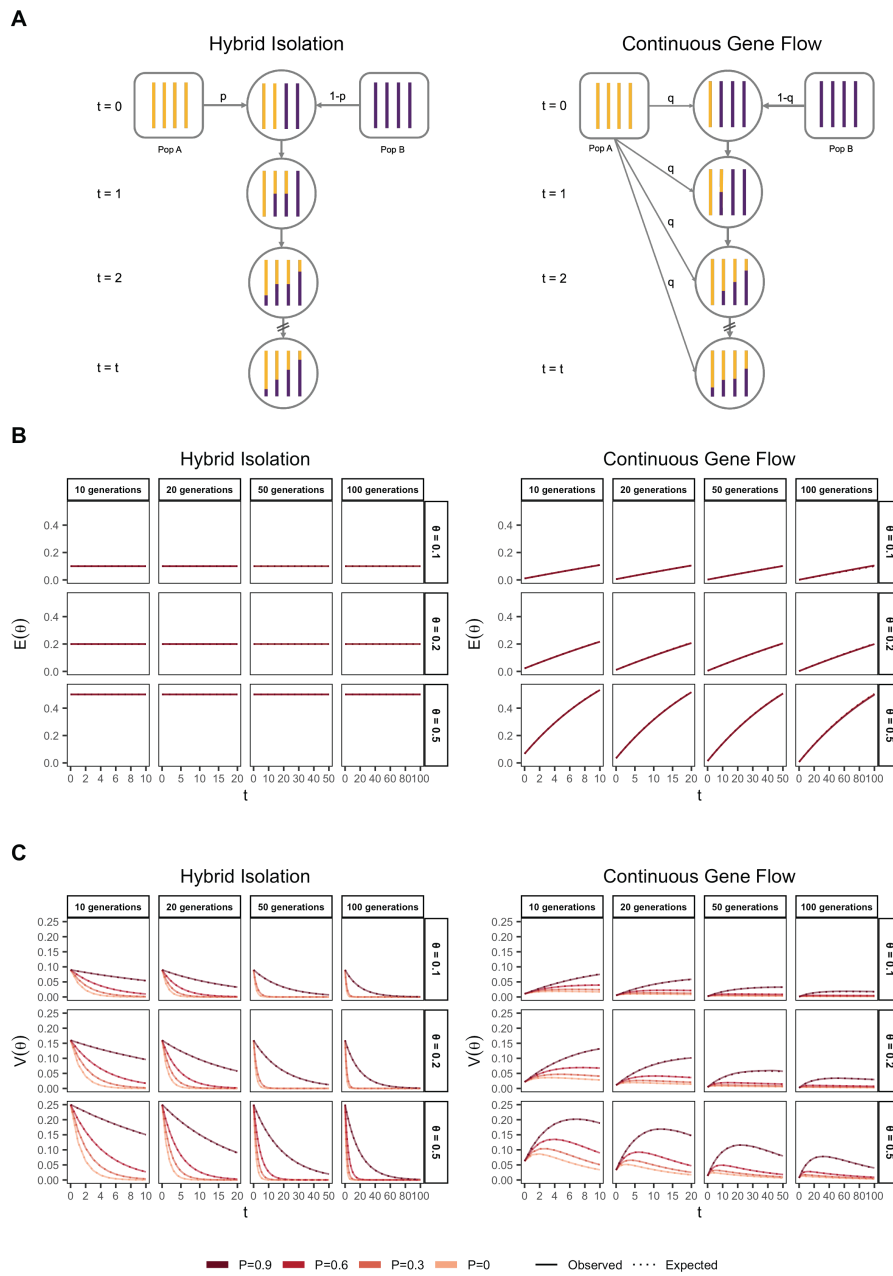


Figure 1: The behavior of mean and variance of individual ancestry as a function of admixture history. (A) Shows the demographic models under which simulations were carried out. Admixture might occur once (Hybrid Isolation, HI, left column) or continuously (Continuous Gene Flow, CGF, right column). (B) The mean individual ancestry,  $\mathbb{E}(\theta)$  remains constant over time in the HI model and increases in the CGF model with continued gene flow. (C) The variance in individual ancestry,  $\mathbb{V}(\theta)$  is maximum at  $t = 0$  in the HI model, decaying subsequently.  $\mathbb{V}(\theta)$  increases with gene flow in the CGF model and will subsequently decrease with time.  $P$  measures the strength of assortative mating, which slows the decay of  $\mathbb{V}(\theta)$ .  $P=0.6$  is missing for simulations run for 50 and 100 generations and  $\theta \in \{0.1, 0.2\}$  due to the difficulty in finding mate pairs (Methods).

129 increase monotonically with increasing divergence,  $\beta_i \beta_j (f_i^A - f_i^B)(f_j^A - f_j^B)$  is expected to be zero under  
 130 neutrality because the direction of frequency change will be uncorrelated across loci. In this case, the LD  
 131 contribution, i.e., (1.4), is expected to be zero and  $V_{gb} = (1.1) + (1.2) + (1.3)$ . However, this is true only  
 132 in expectation over the evolutionary process and the realized LD contribution may be non-zero even for  
 133 neutral traits.

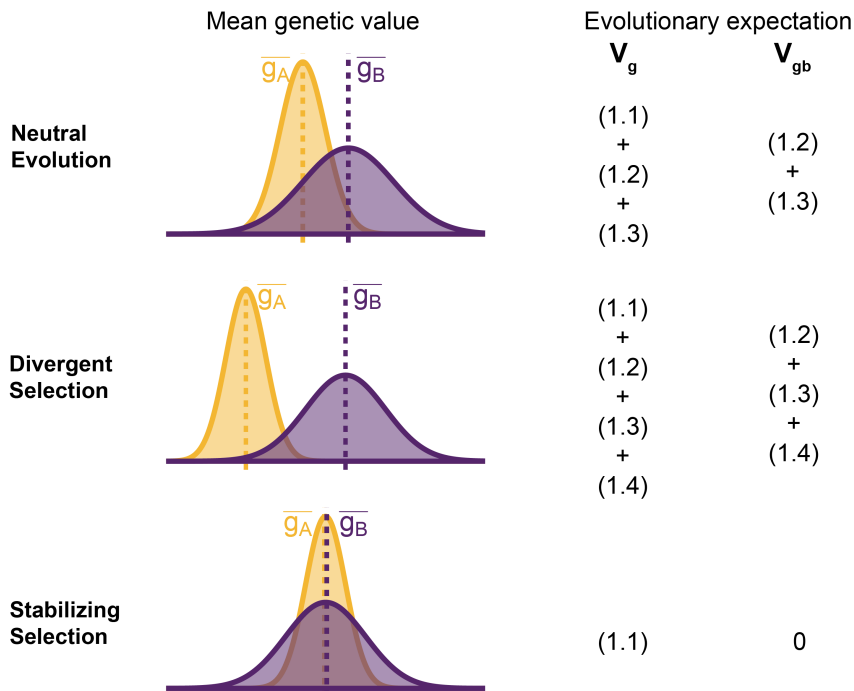


Figure 2: Decomposing genetic variance in a two-population system. The plot illustrates the expected distribution of genetic values in two populations under different selective pressures and the terms on the right list the total ( $V_g$ ) and between-population genetic variance ( $V_{gb}$ ) expected over the evolutionary process. For neutrally evolving traits (top row), we expect there to be an absolute difference in the mean genetic values ( $|\bar{g}_A - \bar{g}_B|$ ) that is proportional to  $F_{ST}$ . For traits under divergent selection (middle),  $|\bar{g}_A - \bar{g}_B|$  is expected to be greater than that expected under genetic drift. For traits under stabilizing selection,  $|\bar{g}_A - \bar{g}_B|$  will be less than that expected under genetic drift, and zero in the extreme case.

134 For traits under selection, the LD contribution is expected to be greater or less than zero, depending  
 135 on the type of selection. Under divergent selection, trait-increasing alleles will be systematically more  
 136 frequent in one population over the other, inducing positive LD across loci [28, 29], increasing the  
 137 LD contribution, i.e., term (1.4). Stabilizing selection, on the other hand, induces negative LD [30,  
 138 31]. In the extreme case, the mean genetic values of the two populations are exactly equal and  $V_{gb} =$   
 139  $(1.2) + (1.3) + (1.4) = 0$ . For this to be true, (1.4) has to be negative and equal to  $(1.2) + (1.3)$ , which  
 140 are both positive, and the total genetic variance is reduced to the within-population variance, i.e., term  
 141 (1.1) (Fig. 2). This is relevant because, as we show in the following sections, the behavior of the genetic  
 142 variance in admixed populations depends on the magnitude of  $V_{gb}$  between the source populations.

143 We illustrate this by tracking the genetic variance in admixed populations for two traits, both with  
 144 the same mean  $F_{ST}$  at causal loci but with different LD contributions (term 1.4): one where the LD

145 contribution is positive (Trait 1) and the other where it is negative (Trait 2). Thus, traits 1 and 2  
146 can be thought of as examples of phenotypes under divergent and stabilizing selection, respectively, and  
147 we refer to them as such from hereon. To simulate the genetic variance of such traits, we drew the  
148 allele frequencies ( $f^A$  and  $f^B$ ) in populations A and B for 1,000 causal loci with  $F_{ST} \sim 0.2$  using the  
149 Balding-Nichols model [32]. We drew their effects ( $\beta$ ) from  $\mathcal{N}(0, \frac{1}{2mf(1-f)})$  where  $\bar{f}$  is the mean allele  
150 frequency between the two populations,  $m$  is the number of loci. To simulate positive and negative  
151 LD, we permuted the effect signs across variants 100 times and selected the combinations that gave the  
152 most positive and negative LD contribution to represent the genetic architecture of traits that might  
153 be under directional (Trait 1) and stabilizing (Trait 2) selection, respectively (Methods). We simulated  
154 the genotypes of 10,000 individuals under the HI and CGF models for  $t \in \{10, 20, 50, 100\}$  generations  
155 post-admixture and calculated genetic values for both traits using  $g = \sum_{i=1}^m \beta_i x_i$ , where  $m = 1,000$   
156 (Method). The observed genetic variance at any time can then be calculated simply as the variance in  
157 genetic values, i.e.  $V_g = \mathbb{V}(g)$ .

158 In the HI model,  $\mathbb{E}(\theta)$  does not change (Fig. 1B) so terms (1.1) and (1.2) are constant through time.  
159 Terms (1.3) and (1.4) decay towards zero as the variance in ancestry goes to zero and  $V_g$  ultimately  
160 converges to (1.1) + (1.2) (Fig. 3). This equilibrium value is equal to the  $\mathbb{E}(V_g|\theta)$  (Appendix) and the  
161 rate of convergence depends on the strength of assortative mating, which slows the rate at which  $\mathbb{V}(\theta)$   
162 decays.  $V_g$  approaches equilibrium from a higher value for traits under divergent selection and lower value  
163 for traits under stabilizing selection because of positive and negative LD contributions, respectively, at  
164  $t = 0$  (Fig. 3). In the CGF model,  $V_g$  increases initially for both traits with increasing gene flow (Fig.  
165 3). This might seem counter-intuitive at first because gene flow increases admixture LD, which leads  
166 to more negative values of the LD contribution for traits under stabilizing selection (Fig. S1). But this  
167 is outweighed by positive contributions from the genic variance – terms (1.1) + (1.2) + (1.3) – all of  
168 which initially increase with gene flow (Fig. S1). After a certain point, the increase in  $V_g$  slows down as  
169 any increase in  $\mathbb{V}(\theta)$  due to gene flow is counterbalanced by recombination and independent assortment.  
170 Ultimately,  $V_g$  will decrease if there is no more gene flow, reaching the same equilibrium value as in the  
171 HI model, i.e.,  $\mathbb{E}(V_g|\theta) = (1.1) + (1.2)$ . Because the loci are unlinked, we refer to the sum (1.3) + (1.4)  
172 as the contribution of population structure.

## 173 GREML estimation

174 In their original paper, Yang *et al.* (2010) defined  $h_{snp}^2$  as the variance explained by genotyped SNPs and  
175 not as heritability [3]. This is because  $h^2$  is the genetic variance explained by causal variants, which are  
176 unknown. Genotyped SNPs may not overlap with or tag all causal variants and thus,  $h_{snp}^2$  is understood  
177 to be a lower bound of  $h^2$ , both being equal if causal variants are known [3]. Our goal is to demonstrate  
178 that this may not be true in structured populations and quantify the bias in  $\hat{h}_{snp}^2$ , even in the ideal  
179 situation when causal variants are known.

180 We used GREML, implemented in GCTA [3, 5], to estimate the genetic variance for our simulated  
181 traits. GCTA assumes the following model:  $\mathbf{y} = \mathbf{Z}\mathbf{u} + \boldsymbol{\epsilon}$  where  $\mathbf{Z}$  is an  $n \times m$  standardized genotype  
182 matrix such that the genotype of the  $k^{th}$  individual at the  $i^{th}$  locus is  $z_{ik} = \frac{x_{ik} - 2f_i}{\sqrt{2f_i(1-f_i)}}$ ,  $f_i$  being the

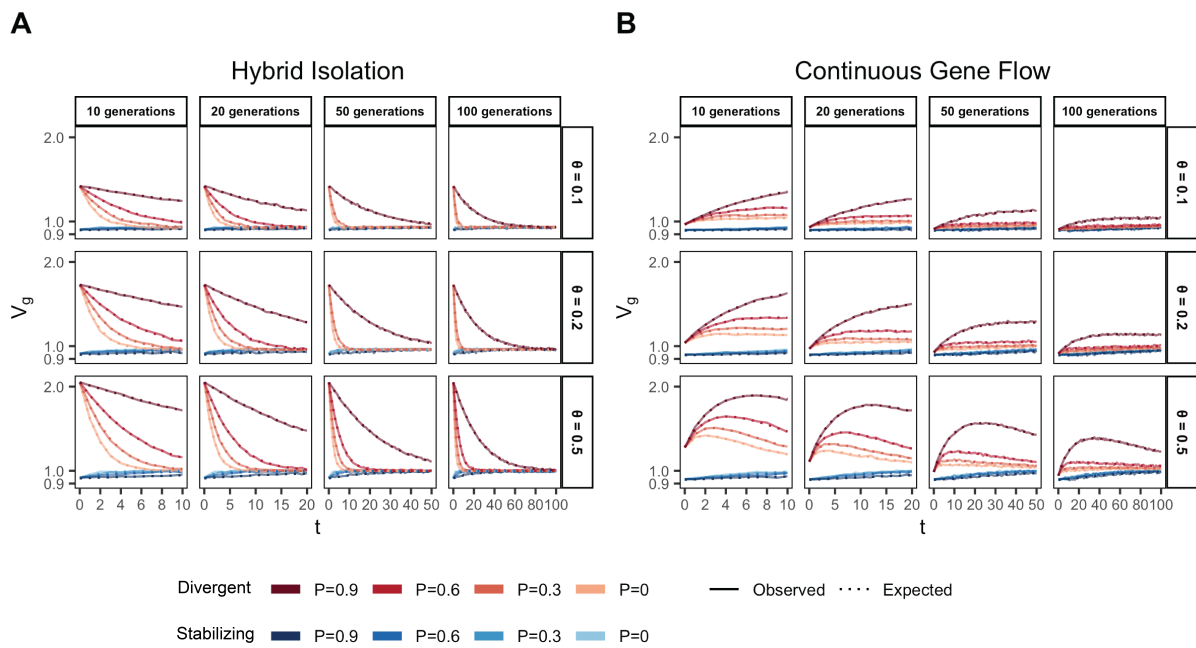


Figure 3: Genetic variance in admixed populations under the (A) HI and (B) CGF models. Dotted lines represent the expected genetic variance based on Eq. (1) and solid lines represent results of simulations averaged over ten replicates. Red and blue lines represent traits under divergent and stabilizing selection, respectively.  $P = 0.6$  is missing for simulations run for 50 and 100 generations and  $\theta \in \{0.1, 0.2\}$  due to the difficulty in finding mate pairs (Methods)



183 allele frequency. The SNP effects corresponding to the scaled genotypes are assumed to be random and  
 184 independent such that  $\mathbf{u} \sim \mathcal{N}(0, \mathbf{I} \frac{\sigma_u^2}{m})$  and  $\epsilon \sim \mathcal{N}(0, \mathbf{I} \sigma_\epsilon^2)$  is random environmental error. Then, the  
 185 phenotypic variance can be decomposed as:

$$\begin{aligned} \mathbb{V}(\mathbf{y}) &= \mathbb{V}(\mathbf{Z}\mathbf{u}) + \mathbb{V}(e) \\ &= \frac{\mathbf{Z}\mathbf{Z}'}{m} \sigma_u^2 + \sigma_\epsilon^2 \end{aligned}$$

186 where  $\frac{\mathbf{Z}\mathbf{Z}'}{m}$  is the genetic relationship matrix (GRM), the variance components  $\sigma_u^2$  and  $\sigma_\epsilon^2$  are estimated  
 187 using restricted maximum likelihood, and  $\hat{h}_{snp}^2$  is calculated as  $\frac{\hat{\sigma}_u^2}{\hat{\sigma}_u^2 + \hat{\sigma}_\epsilon^2}$ . We are interested in asking  
 188 whether  $\hat{\sigma}_u^2$  is an unbiased estimate of  $V_g$ . To answer this, we constructed the GRM with causal variants  
 189 and estimated  $\hat{\sigma}_u^2$  using GCTA [3, 4].

190 GCTA under- and over-estimates the genetic variance in admixed populations for traits under diver-  
 191 gent (Trait 1) and stabilizing selection (Trait 2), respectively, when there is population structure, i.e.,  
 192 when  $\mathbb{V}(\theta) > 0$  (Fig. 4A). One reason for this bias is that the GREML model assumes that the effects  
 193 are independent, and therefore the LD contribution is zero. This, as discussed in the previous section, is  
 194 not true for traits under divergent or stabilizing selection between the source populations, and only true  
 195 for neutral traits in expectation. Because of this,  $\hat{\sigma}_u^2$  does not capture the LD contribution, i.e. term  
 196 (1.4) (Fig. 4A). But  $\hat{\sigma}_u^2$  can be biased even if the LD contribution is zero if the genotypes are scaled with  
 197  $\sqrt{2f_i(1-f_i)}$  – the standard practice – where  $f_i$  is the frequency of the allele in the population. This  
 198 scaling assumes that  $\mathbb{V}(x_i) = 2f_i(1-f_i)$ , which is true only if the population were mating randomly.  
 199 In an admixed population  $\mathbb{V}(x_i) = 2f_i(1-f_i) + 2\mathbb{V}(\theta)(f_i^A - f_i^B)^2$ , where  $f_i$ ,  $f_i^A$ , and  $f_i^B$  correspond  
 200 to frequency in the admixed population, and source populations, A and B, respectively (Appendix).  
 201 Alternatively, if the genotypes are scaled,  $\mathbb{V}(z_i) = 1 + 2\mathbb{V}(\theta)F_{st}^{(i)}$  where  $F_{st}^{(i)}$  is the  $F_{st}$  at the  $i^{th}$  locus.  
 202 We show that this assumption biases  $\hat{\sigma}_u^2$  downwards by a factor of  $2\mathbb{V}(\theta)(f_i^A - f_i^B)^2$  (or  $2\mathbb{V}(\theta)F_{st}^{(i)}$  if  
 203 genotypes are scaled) – term (1.3) (Fig. 4B, Appendix). Thus, with the standard scaling,  $\hat{\sigma}_u^2$  gives a  
 204 biased estimate in the presence of population structure, even of the genic variance.

205 The overall bias in  $\hat{\sigma}_u^2$  is determined by the relative magnitude and direction of terms (1.3) and  
 206 (1.4), both of which are functions of  $\mathbb{V}(\theta)$ , and therefore, of the degree of structure in the population.  
 207 The contribution of term (1.3) will be modest, even in highly structured populations (Fig. S1) and  
 208 therefore, the overall bias is largely driven by the LD contribution. If there is no more gene flow,  $\mathbb{V}(\theta)$   
 209 will ultimately go to zero and  $V_g$  will converge towards  $\hat{\sigma}_u^2$ . Thus,  $\hat{\sigma}_u^2$  is more accurately interpreted  
 210 as the genetic variance expected if the LD contribution were zero and if the population were mating  
 211 randomly. In other words,  $\mathbb{E}(\hat{\sigma}_u^2) = (1.1) + (1.2) \neq V_g$  (Fig. 4B).

212 In principle, we can recover the missing components of  $V_g$  by scaling the genotypes appropriately.  
 213 For example, we can recover term (1.3) by scaling the genotype at each variant  $i$  by its sample variance,  
 214 i.e.,  $z_{ik} = \frac{x_{ik} - 2f_i}{\sqrt{\mathbb{V}(x_i)}}$  (Fig. 4C) (Appendix). We can also recover term (1.4) by scaling the genotypes with  
 215 the covariance between SNPs, i.e., the LD matrix, as previously proposed [33, 34] (Methods). In matrix  
 216 form, the ‘LD-scaled’ genotypes can be written as  $\mathbf{Z} = (\mathbf{X} - 2\mathbf{P})\mathbf{U}^{-1}$  where  $\mathbf{P}$  is an  $n \times m$  matrix such  
 217 that all elements of the  $i^{th}$  column contain the frequency of the  $i^{th}$  SNP and  $\mathbf{U}$  is the (upper triangular)  
 218 square root matrix of the LD matrix, i.e.,  $\mathbf{\Sigma} = \mathbf{U}'\mathbf{U}$  [33]. GREML recovers the LD contribution under

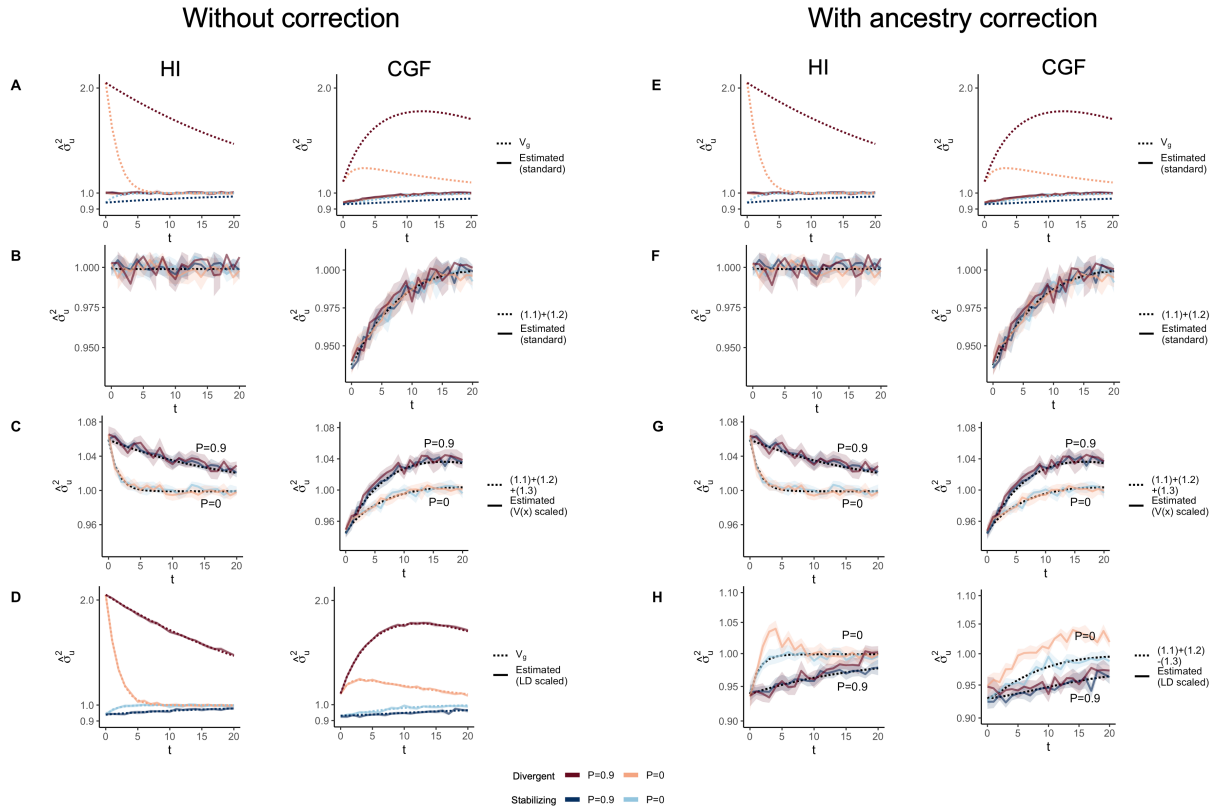


Figure 4: The behavior of GREML estimates of the genetic variance ( $\hat{\sigma}_u^2$ ) in admixed populations under the HI (left column) and CGF (right column) models either without (A-D) or with (E-H) individual ancestry as a fixed effect. The solid lines represent estimates from simulated data averaged across ten replicates with red and blue colors representing estimates for traits under divergent and stabilizing selection, respectively.  $P$  indicates the strength of assortative mating. The shaded area represents the 95% confidence bands generated by bootstrapping (sampling with replacement 100 times) the point estimate reported by GCTA. The dotted lines either represent the expected variance in the population based on Eq. 1 (A & B) or the expected estimate for three different ways of scaling genotypes (B-D & F-H). (A-B & E-F) show the behavior of  $\hat{\sigma}_u^2$  for the default scaling, (C, G) shows  $\hat{\sigma}_u^2$  when the genotype at a locus is scaled by its sample variance ( $\mathbb{V}(x)$  scaled), and (D, H) when it is scaled by the sample covariance (LD scaled).

219 this scaling, resulting in unbiased estimates of  $V_g$  for both traits (Fig. 4D, Appendix).

220 In practice, however, the LD contribution may not be fully recoverable for two reasons. One, the  
 221 LD-scaled GRM requires computing the inverse of  $\Sigma$  or  $U$  which may not exist, especially if the number  
 222 of markers is greater than the sample size – the case for most human genetic studies. Second, it is  
 223 common to include individual ancestry or principal components of the GRM as fixed effects in the model  
 224 to account for inflation in heritability estimates due to shared environment. This should also have the  
 225 effect of removing the components of genetic variance along the ancestry axes, the residual variance being  
 226 equal to  $\mathbb{E}\{\mathbb{V}(g|\theta)\} = (1.1) + (1.2) - (1.3)$  (Appendix). Indeed, this is what we observe in Fig. 4H. Thus,  
 227 if ancestry is included as a fixed effect, we expect  $V_g$  to be underestimated in the presence of population  
 228 structure, regardless of genetic architecture.

## 229 HE estimation

230 Haseman-Elston (HE) regression also assumes a random-effects model but uses a method-of-moments  
 231 approach, as opposed to GREML, which maximizes the likelihood to estimate  $V_g$ . Previous work has  
 232 shown that as long as all causal variants are included in the GRM calculation, the HE estimator will  
 233 not be biased, even if they are in LD with each other [35]. We show that in the presence of positive  
 234 and negative LD between causal loci, as exemplified by traits under divergent and stabilizing selection,  
 235 respectively, the HE estimates of  $V_g$  are biased upwards and downwards, respectively (Fig. 5A-B). To  
 236 understand this discrepancy and the source of bias in our simulations, recall that HE estimates  $V_g$  from  
 237 the regression of the (pairwise) phenotypic covariance between individuals on their genotypic covariance  
 238 [24]. More specifically, if we denote  $Y_{kl} = y_k y_l$  as the product of the (centered) phenotypes of  $k^{th}$  and  
 239  $l^{th}$  individuals, and  $\psi_{kl}$  as the  $k^{th}$  and  $l^{th}$  entry of the GRM, then the HE estimator can be written as:

$$\begin{aligned}
 \hat{V}_g &= \frac{Cov(Y_{kl}, \psi_{kl})}{Var(\psi_{kl})} \\
 &= \frac{\mathbb{E}(y_k y_l \sum_{w=1}^M z_{wk} z_{wl})}{\mathbb{E}(\sum_{i=1}^M z_{ik} z_{il} \sum_{w=1}^M z_{wk} z_{wl})} \\
 &= \frac{\mathbb{E}\{(g_k + e_k)(g_l + e_l) \sum_{w=1}^M z_{wk} z_{wl}\}}{\mathbb{E}(\sum_{i=1}^M z_{ik} z_{il} \sum_{w=1}^M z_{wk} z_{wl})} \\
 &= \frac{\mathbb{E}(g_k g_l \sum_{w=1}^M z_{wk} z_{wl})}{\mathbb{E}(\sum_{i=1}^M z_{ik} z_{il} \sum_{j=1}^M z_{wk} z_{wl})} \\
 &= \frac{\mathbb{E}(\sum_{i=1}^M \sum_{j=1}^M u_i u_j z_{ik} z_{jl} \sum_{w=1}^M z_{wk} z_{wl})}{\mathbb{E}(\sum_{i=1}^M z_{ik} z_{il} \sum_{w=1}^M z_{wk} z_{wl})} \\
 &= \frac{\mathbb{E}(\sum_{i=1}^M \sum_{j=1}^M u_i u_j \sum_{w=1}^M z_{ik} z_{jl} z_{wk} z_{wl})}{\mathbb{E}(\sum_{i=1}^M \sum_{w=1}^M z_{ik} z_{il} z_{wk} z_{wl})} \\
 &= \frac{\mathbb{E}(\sum_{i=1}^M u_i^2 \sum_{w=1}^M z_{ik} z_{jl} z_{wk} z_{wl})}{\mathbb{E}(\sum_{i=1}^M \sum_{w=1}^M z_{ik} z_{il} z_{wk} z_{wl})} + \frac{\mathbb{E}(\sum_{i=1}^M \sum_{j \neq i} u_i u_j \sum_{w=1}^M z_{ik} z_{jl} z_{wk} z_{wl})}{\mathbb{E}(\sum_{i=1}^M \sum_{w=1}^M z_{ik} z_{il} z_{wk} z_{wl})} \quad (2)
 \end{aligned}$$

240 Where the first and second terms represent the genic and LD components, respectively, of the esti-  
 241 mate. Population structure induces correlations between the alleles at a given locus as well as across  
 242 loci (i.e., LD). But the LD may not be directional, i.e., trait-increasing alleles may be as likely to be  
 243 co-inherited with each other as they are to trait-decreasing alleles, and vice versa – implicit under the

244 standard random-effects model. Thus, in the absence of directional LD, the second term is zero and the  
245 first term is unaffected as long as all causal variants are included in the GRM, because the increase in  
246 the numerator due to population structure is proportional to the denominator [35]. Directional LD does  
247 not affect the first term but exaggerates the contribution from the second term, i.e., the LD component  
248 (see Appendix section A3.2). Consequently, HE regression over- and under-estimates  $V_g$  for traits with  
249 positive and negative LD, respectively. Note that this bias is in the opposite direction of the bias observed  
250 with GREML, which fails to capture the LD contribution. Scaling the genotype at a locus by its LD with  
251 other loci, as discussed in the previous section, corrects for the bias in HE regression regardless of genetic  
252 architecture, yielding estimates consistent with GREML (Fig. 5C). Thus, GREML and HE regression  
253 are guaranteed to yield the same estimates only if the underlying model specifying the distribution of  
254 effects is consistent with the true architecture of the trait.

255 The practice of including individual ancestry as a covariate in HE regression to account for shared  
256 environment [11] reduces the bias from exaggerated LD contributions (Fig. 5D-F). But, as with GREML,  
257 this also removes any genetic variance that may exist along the ancestry axis, yielding underestimates  
258 of  $V_g$ , regardless of genetic architecture.

## 259 Local ancestry heritability

260 A related quantity of interest in admixed populations is local ancestry heritability ( $h_\gamma^2$ ), which is defined  
261 as the proportion of phenotypic variance that can be explained by local ancestry. Zaitlen *et al.* (2014)  
262 [36] showed that this quantity is related to, and can be used to estimate,  $h^2$  in admixed populations.  
263 The advantage of this approach is that local ancestry segments shared between individuals are identical  
264 by descent and are therefore, more likely to tag causal variants compared to array markers, allowing  
265 one to potentially capture the contributions of rare variants [36]. Here, we show that in the presence of  
266 population structure, (i) the relationship between  $h_\gamma^2$  and  $h^2$  is not straightforward and (ii)  $\hat{h}_\gamma^2$  may be a  
267 biased estimate of local ancestry heritability under the random effects model for the same reasons that  
268  $\hat{h}_{snp}^2$  is biased.

269 We define local ancestry  $\gamma_i \in \{0, 1, 2\}$  as the number of alleles at locus  $i$  that trace their ancestry  
270 to population A. Thus, ancestry at the  $i^{th}$  locus in individual  $k$  is a binomial random variable with  
271  $\mathbb{E}(\gamma_{ik}) = 2\theta_k$ ,  $\theta_k$  being the ancestry of the  $k^{th}$  individual. Similar to genetic value, the ‘ancestry value’  
272 of an individual can be defined as  $\sum_{i=1}^m \phi_i \gamma_i$ , where  $\phi_i = \beta_i(f_i^A - f_i^B)$  is the effect size of local ancestry  
273 (Appendix). Then, the genetic variance due to local ancestry can be expressed as (Appendix):

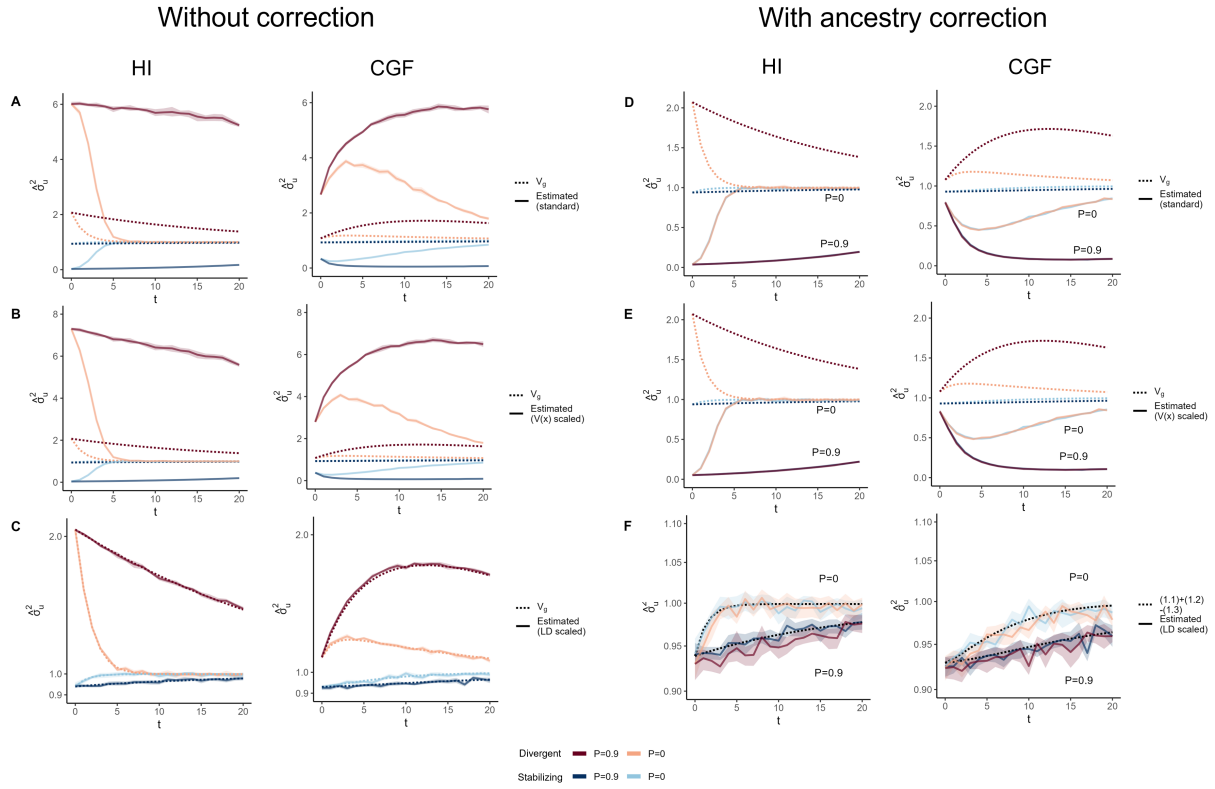


Figure 5: Genetic variance ( $\hat{V}_g$ ) estimated with HE regression in admixed populations under the HI (left column) and CGF (right column) models either without (A-C) or with (D-F) adjustment for individual ancestry. The solid lines represent estimates from simulated data averaged across ten replicates with red and blue colors representing estimates for traits under divergent and stabilizing selection, respectively.  $P$  indicates the strength of assortative mating. (A & D) show the behavior of  $\hat{V}_g$  for the default scaling, (B, E) shows  $\hat{V}_g$  when the genotype at a locus is scaled by its sample variance ( $V(x)$  scaled), and (C, F) when it is scaled by the sample covariance (LD scaled). The dotted lines in A-E represent the expected  $V_g$  in the population based on Eq. 1 and in F, represent the expected  $V_g$  after removing any genetic variance along the ancestry axis. The shaded areas represent the 95% bootstrapped confidence bands of the estimate.

$$\begin{aligned}
 V_\gamma &= \mathbb{V} \left( \sum_{i=1}^m \phi_i \gamma_i \right) = \sum_{i=1}^m \phi_i^2 \mathbb{V}(\gamma_i) + \sum_{i=1}^m \sum_{j \neq i}^m \phi_i \phi_j \mathbb{C}(\gamma_i, \gamma_j) \\
 &= 2 \mathbb{E}(\theta) \{1 - \mathbb{E}(\theta)\} \sum_{i=1}^m \phi_i^2 + 2 \mathbb{V}(\theta) \sum_{i=1}^m \phi_i^2 + 4 \mathbb{V}(\theta) \sum_{i=1}^m \sum_{j \neq i}^m \phi_i \phi_j \\
 &= 2 \mathbb{E}(\theta) \{1 - \mathbb{E}(\theta)\} \sum_{i=1}^m \beta_i^2 (f_i^A - f_i^B)^2 \\
 &\quad + 2 \mathbb{V}(\theta) \sum_{i=1}^m \beta_i^2 (f_i^A - f_i^B)^2 \\
 &\quad + 4 \mathbb{V}(\theta) \sum_{i=1}^m \sum_{j \neq i}^m \beta_i \beta_j (f_i^A - f_i^B)(f_j^A - f_j^B)
 \end{aligned}$$

274 and heritability explained by local ancestry is simply the ratio of  $V_\gamma$  and the phenotypic variance. Note  
 275 that  $V_\gamma = (1.2) + (1.3) + (1.4)$  and therefore its behavior is similar to  $V_g$  in that the terms (1.3) and (1.4)  
 276 decay towards zero as  $\mathbb{V}(\theta) \rightarrow 0$ , and  $V_\gamma$  converges to (1.2) (Fig. S2). Additionally, the dependence of  
 277  $V_\gamma$  on both  $\mathbb{E}(\theta)$  and  $\mathbb{V}(\theta)$  precludes a straightforward derivation between local ancestry heritability and  
 278  $h^2$ .

279 GREML estimation of  $\hat{h}_\gamma^2$  is similar to that of  $\hat{h}_{snp}^2$ , the key difference being that the former involves  
 280 constructing the GRM using local ancestry instead of genotypes [36]. The following model is assumed:  
 281  $\mathbf{y} = \mathbf{W}\mathbf{v} + \boldsymbol{\xi}$  where  $\mathbf{W}$  is an  $n \times m$  standardized local ancestry matrix,  $\mathbf{v} \sim \mathcal{N}(0, \mathbf{I} \frac{\sigma_v^2}{m})$  are local  
 282 ancestry effects, and  $\boldsymbol{\xi} \sim \mathcal{N}(0, \mathbf{I} \sigma_\xi^2)$ . Note that  $\sigma_\xi^2$  captures both environmental noise as well as any  
 283 genetic variance independent of local ancestry. The phenotypic variance is decomposed as  $\mathbb{V}(\mathbf{y}) =$   
 284  $\mathbb{V}(\mathbf{W}\mathbf{v}) + \mathbb{V}(\boldsymbol{\xi}) = \frac{\mathbf{W}\mathbf{W}'}{m} \sigma_v^2 + \sigma_\xi^2$  where  $\frac{\mathbf{W}\mathbf{W}'}{m}$  is the local ancestry GRM and  $\sigma_v^2$  is the parameter of  
 285 interest, which is believed to be equal to  $V_\gamma$  – the genetic variance due to local ancestry.

286 We show that, in the presence of population structure, i.e., when  $\mathbb{V}(\theta) > 0$ , GREML  $\hat{\sigma}_v^2$  is biased  
 287 downwards relative to  $V_\gamma$  for traits under divergent selection and upwards for traits under stabilizing  
 288 selection because it does not capture the contribution of LD (Fig. 6A). But there is another source of  
 289 bias in  $\hat{\sigma}_v^2$ , which tends to be inflated in the presence of population structure if individual ancestry is  
 290 not included as a covariate, even with respect to the expectation of  $V_\gamma$  under equilibrium (seen more  
 291 clearly in Fig. 6B-C). We suspect this inflation is because of strong correlations between local ancestry  
 292 – local ancestry disequilibrium – across loci that inflates  $\hat{\sigma}_v^2$  in a way that is not adequately corrected  
 293 even when all causal variants are included in the model [4, 10]. Scaling local ancestry by its covariance  
 294 removes this bias and recovers the contribution of LD (Fig. 6D) presumably because this accounts for  
 295 the correlation in genotypes across loci. Including individual ancestry as a fixed effect also corrects for  
 296 the inflation in  $\hat{\sigma}_v^2$  (Fig. 6E-H). But as with  $\hat{\sigma}_u^2$ , this practice will underestimate the genetic variance  
 297 due to local ancestry in the presence of population structure because it removes the variance along the  
 298 ancestry axis (Fig. 6E-H).

299 Based on the above, GREML  $\hat{h}_\gamma^2$  and corresponding estimates of  $h^2$  are more accurately interpreted as  
 300 the heritability due to local ancestry and heritability, respectively, expected in the absence of population  
 301 structure. We believe  $\hat{h}_\gamma^2$  is still useful in that, because it should capture the effects of rare variants, it  
 302 can be used to estimate the upper bound of  $\hat{h}_{snp}^2$ .

303 In a previous paper, we suggested that local ancestry heritability could potentially be used to estimate  
 304 the genetic variance between populations [37]. Our results suggest this is not possible for two reasons.  
 305 First, the GREML estimator of local ancestry heritability, as we show in this section is biased and  
 306 does not capture the LD contribution. But even if we were able to recover the LD component, our  
 307 decomposition shows that local ancestry is equal to the genetic variance between populations ( $V_{gb}$ ) only  
 308 when  $\mathbb{E}(\theta) = 0.5$  and  $\mathbb{V}(\theta) = \mathbb{E}(\theta)\{1 - \mathbb{E}(\theta)\} = 0.25$ , which is only possible at  $t = 0$  in the HI model. After  
 309 admixture,  $\mathbb{V}(\theta)$  decays and the equivalence between  $V_\gamma$  and  $V_{gb}$  is lost, making it impossible to estimate  
 310 the latter from admixed populations, especially for traits under divergent or stabilizing selection, even  
 311 if the environment is randomly distributed with respect to ancestry. We note that this conclusion was  
 312 recently reached independently by Schraiber and Edge (2023) [38].

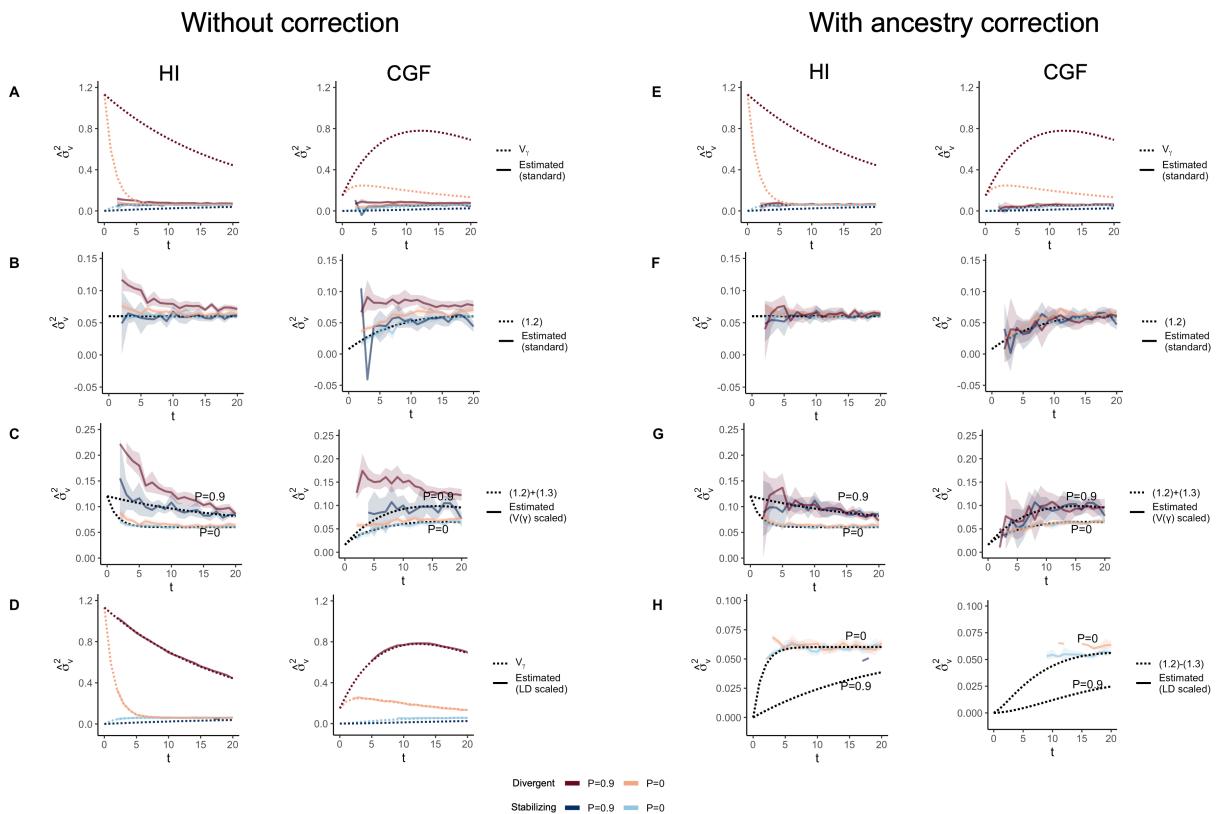


Figure 6: The behavior of GREML estimates of the variance due to local ancestry ( $\hat{\sigma}_v^2$ ) in admixed populations under the HI (left column) and CGF (right column) models either without (A-D) or with (E-H) individual ancestry included as a fixed effect. The solid lines represent estimates from simulated data averaged across ten replicates with red and blue colors representing estimates for traits under divergent and stabilizing selection, respectively. P indicates the strength of assortative mating. The dotted lines either represent the expected variance in the population (A & B) or the expected estimate for three different ways of scaling local ancestry (B-D & F-H). (A-B & E-F) show the behavior of  $\hat{\sigma}_v^2$  for the default scaling, (C, G) shows  $\hat{\sigma}_v^2$  when local ancestry is scaled by the sample variance, and (D, H) when it is scaled by the sample covariance. Shaded regions represent the 95% confidence bands. Some runs in (D & H) failed to converge as seen by the missing segments of the solid lines because the expected variance in such cases was too small.

### 313 How much does LD contribute to $V_g$ in practice?

314 In the previous sections, we showed theoretically that  $\hat{h}_{snp}^2$  may be biased in admixed populations even  
315 if the causal variants are known and in the absence of confounding by shared environment. GREML  
316 fails to capture the LD contribution whereas HE regression overestimates it. The extent to which  $\hat{h}_{snp}^2$  is  
317 biased because of this reason in practice is ultimately an empirical question, which is difficult to answer  
318 because the true genetic architecture – the LD contribution in particular – is unknown. In this section,  
319 we develop some intuition for this contribution among individuals with mixed African and European  
320 ancestry using a combination of simulations and empirical data analysis.

321 First, we simulated a neutral trait using genotype data from the African Americans (ASW) from the  
322 1,000 Genomes Project (1KGP) [39]. To do this, we sampled  $m \in \{10, 100, 1,000\}$  causal loci from a set of  
323 common (MAF > 0.01), LD pruned variants and assigned them effects such that  $\beta_i \sim \mathcal{N}\left(0, \frac{1}{\sqrt{m \mathbb{V}(x_i)}}\right)$ ,  
324 i.e., the expected *genic* variance is  $\mathbb{E}\{\sum_{i=1}^m \beta_i^2 \text{Var}(x_i)\} = 1$  (Methods). We computed the genic and  
325 LD contributions and repeated this process 1,000 times where each replicate can be thought of as an  
326 independent realization of the genetic architecture of a neutrally evolving trait. We show that the LD  
327 contribution may be zero in expectation but can be substantial for a given trait (up to 50% of the genic  
328 variance, Fig. S4), even in the absence of selection.

329 Second, we estimated the LD contribution of genome-wide significant SNPs for 26 quantitative traits  
330 from the GWAS catalog [40]. To do this, we decomposed the variance explained in ASW into the four  
331 components in Equation 1 using allele frequencies ( $f^A$  and  $f^B$ ) from the YRI and CEU and the mean  
332 ( $\mathbb{E}(\theta) \approx 0.77$ ) and variance ( $\mathbb{V}(\theta) \approx 0.02$ ) of individual ancestry from ASW (Methods). We show that  
333 for skin pigmentation – a trait under strong divergent selection – the LD contribution, i.e. term (1.4),  
334 is positive and accounts for  $\approx 40 - 50\%$  of the total variance explained. This is because of large allele  
335 frequency differences between Africans and Europeans that are correlated across skin pigmentation loci,  
336 consistent with strong polygenic selection favoring alleles for darker pigmentation in regions with high UV  
337 exposure and vice versa [37, 41–44]. But for most other traits, LD contributes relatively little, explaining  
338 a modest, but non-negligible proportion of the genetic variance in height, LDL and HDL cholesterol, mean  
339 corpuscular hemoglobin (MCH), neutrophil count (NEU), and white blood cell count (WBC) (Fig. 7).  
340 Because we selected independent associations for this exercise (Methods), the LD contribution is driven  
341 entirely due to population structure in ASW. The contribution of population structure to the genic  
342 variance, i.e., term (1.3) is also small even for traits like skin pigmentation and neutrophil count with  
343 large effect alleles that are highly diverged in frequency between Africans and Europeans [42, 43, 45–  
344 47]. Overall, this suggests that population structure contributes relatively little, at least to the variance  
345 explained by GWAS SNPs.

## 346 Discussion

347 Despite the growing size of GWAS and discovery of thousands of variants for hundreds of traits [40], the  
348 heritability explained by GWAS SNPs remains a fraction of twin-based heritability estimates. Yang *et*  
349 *al.* (2010) introduced the concept of SNP heritability ( $h_{snp}^2$ ) that does not depend on the discovery of



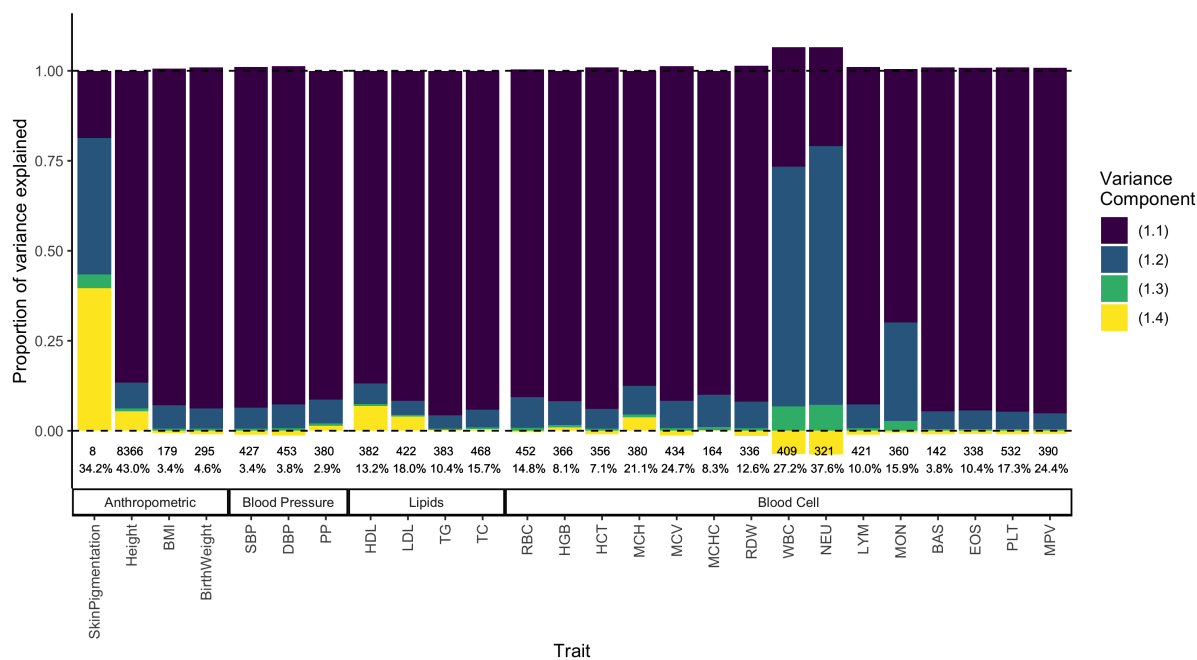


Figure 7: Decomposing the genetic variance explained by GWAS SNPs in the 1000 Genomes ASW (African Americans from Southwest). We calculated the four variance components listed in Eq. 1, their values shown on the y-axis as a fraction of the total variance explained (shown as percentage at the bottom). The LD contribution, which can be positive or negative, is shown in yellow. The number of variants used to calculate variance components for each trait is also shown at the bottom.

350 causal variants but assumes that they are numerous and are more or less uniformly distributed across the  
 351 genome (the infinitesimal model), their contributions to the genetic variance ‘tagged’ by genotyped SNPs  
 352 [3].  $h_{snp}^2$  is now routinely estimated in most genomic studies and at least for some traits (e.g. height and  
 353 BMI), these estimates now approach twin-based heritability [6]. But despite the widespread use of  $\hat{h}_{snp}^2$ ,  
 354 its interpretation remains unclear, particularly in the presence of admixture and population structure.  
 355 It is generally accepted that  $\hat{h}_{snp}^2$  can be biased in structured populations because of confounding effects  
 356 of unobserved environmental factors and LD between causal variants [4, 7, 9–11, 48]. But  $\hat{h}_{snp}^2$  may be  
 357 biased even in the absence of confounding because of misspecification of the underlying random-effects  
 358 model, i.e., if the model does not represent the genetic architecture from which the trait is sampled  
 359 [14–16, 49, 50].

360 Under the standard GREML model, SNP effects are assumed to be uncorrelated and the total genetic  
 361 variance can be represented as the sum of the variance explained by individual loci, i.e. the genic variance  
 362 [14–16]. In admixed populations, there is substantial LD, which can contribute to the genetic variance,  
 363 and can persist for a number of generations, despite recombination, due to continued gene flow and/or  
 364 ancestry-based assortative mating. GREML does not capture this LD contribution [12, 15], and therefore,  
 365 may lead to biased estimates of  $h_{snp}^2$ . The LD contribution can be negative for traits under stabilizing  
 366 selection, and positive for traits under divergent selection between the source populations, leading to  
 367 over- or under-estimates, respectively. Thus, GREML estimates of  $h_{snp}^2$ , assuming genotypes are scaled  
 368 properly (see below), is better interpreted as the proportion of phenotypic variance explained by the

369 *genic* variance. Estimates of local ancestry heritability ( $\hat{h}_\gamma^2$ ) [36, 51] should be interpreted similarly.

370 We show that with GREML,  $\hat{h}_{snp}^2$  can be biased even when the LD contribution is zero if the geno-  
371 types are scaled by  $\sqrt{2f(1-f)}$  – the standard approach, which implicitly assumes a randomly mating  
372 population. In the presence of population structure, the variance in genotypes can be higher and  $\hat{h}_{snp}^2$   
373 does not capture this additional variance, which we show can be recovered by scaling genotypes by the  
374 SNP variance ( $\sqrt{Var(x)}$ ). In principle, the LD contribution can also be recovered by scaling genotypes  
375 by the SNP covariance, i.e., the LD matrix, as previously suggested [33, 34]. But this approach is limited  
376 to situations where the sample size is much larger than the number of markers.

377 We also investigated the behavior of another widely used approach to estimate  $h_{snp}^2$  – Haseman-  
378 Elston regression. We show that  $\hat{h}_{snp}^2$  estimated with HE regression is also biased, but for different  
379 reasons and in the opposite direction of the bias observed with GREML. HE regression exaggerates  
380 the LD contribution, leading to over- and under-estimates of  $h_{snp}^2$  for traits where the causal loci are  
381 in positive and negative LD, respectively. Approaches that correct for population structure [35] should  
382 remove this source of bias but would also remove any genetic variance in the trait along the ancestry axis,  
383 including the LD contribution. This results in underestimates of  $h_{snp}^2$ , regardless of trait architecture.

384 One limitation of this paper is that we have focused on random-effects estimators of  $h_{snp}^2$  because of  
385 their widespread use. Estimators of  $h_{snp}^2$  can be broadly grouped into random- and fixed effect estimators  
386 based on how they treat SNP effects [35]. Fixed effect estimators make fewer distributional assumptions  
387 but they are not as widely used because they require conditional estimates of all variants – a high-  
388 dimensional problem where the number of markers is often far larger than the sample size [52]. This is  
389 one reason why random effect estimators, such as GREML, are popular – because they reduce the number  
390 of parameters that need to be estimated by assuming that the effects are drawn from some distribution  
391 where the variance is the only parameter of interest. Fixed effects estimators, in principle, should be able  
392 to capture the LD contribution but this is not obvious in practice since the simulations used to evaluate  
393 the accuracy of such estimators still assume uncorrelated effects [35, 52, 53]. Further research is needed  
394 to clarify the interpretation of the different estimators of  $h_{snp}^2$  in structured populations under a range  
395 of genetic architectures.

396 Does the LD contribution to the genetic variance have practical implications? The answer to this  
397 depends on the context in which SNP heritability is used.  $\hat{h}_{snp}^2$  can be useful in quantifying the power  
398 to detect variants in GWAS where the quantity of interest is the genic variance. But  $\hat{h}_{snp}^2$  can lead to  
399 misleading conclusions if used to measure the extent to which genetic variation contributes to phenotypic  
400 variation, in predicting the response to selection, or in defining the upper limit of polygenic prediction  
401 accuracy [2] – applications where the LD contribution is important.

402 Ultimately, the discrepancy between  $\hat{h}_{snp}^2$  and  $h^2$  in practice is an empirical question, the answer to  
403 which depends on the degree of population structure (which we can measure) and the genetic architecture  
404 of the trait (which we do not know *a priori*). We show that for most traits, the contribution of population  
405 structure to the variance explained by GWAS SNPs is modest among African Americans. Thus, if we  
406 assume that the genetic architecture of GWAS SNPs represents that of all causal variants, then despite  
407 incorrect assumptions, the discrepancy between  $\hat{h}_{snp}^2$  and  $h^2$  should be fairly modest. But this assumption  
408 is unrealistic given that GWAS SNPs are common variants that in most cases cumulatively explain a

409 fraction of trait heritability. What is the LD contribution of the rest of the genome, particularly rare  
410 variants? This is not obvious and will become clearer in the near future through large sequence-based  
411 studies [54]. While these are underway, theoretical studies are needed to understand how different  
412 selection regimes influence the directional LD between causal variants – clearly an important aspect of  
413 the genetic architecture of complex traits.

## 414 Methods

### 415 Simulating genetic architecture

416 We first drew the allele frequency ( $f^0$ ) of 1,000 biallelic causal loci in the ancestor of populations A and  
417 B from a uniform distribution,  $U(0.001, 0.999)$ . Then, we simulated their frequency in populations A and  
418 B ( $f^A$  and  $f^B$ ) under the Balding-Nichols model [32], such that  $f^A, f^B \sim \text{Beta}(\frac{f^0(1-F)}{F}, \frac{(1-f^0)(1-F)}{F})$   
419 where  $F = 0.2$  is the inbreeding coefficient. We implemented this using code adapted from Lin *et al.*  
420 (2021) [55]. To avoid drawing extremely rare alleles, we continued to draw  $f^A$  and  $f^B$  until we had 1,000  
421 loci with  $f^A, f^B \in (0.01, 0.99)$ .

422 We generated the effect size ( $\beta$ ) of each locus by sampling from  $\mathcal{N}(0, \frac{1}{2mf(1-f)})$ , where  $m$  is the  
423 number of loci and  $\bar{f}$  is the mean allele frequency across populations A and B. Thus, rare variants have  
424 larger effects than common variants and the total genetic variance sums to 1. Given these effects, we  
425 simulated two different traits, one with a large difference in means between populations A and B (Trait  
426 1) and the other with roughly no difference (Trait 2). This was achieved by permuting the signs of the  
427 effects 100 times to get a distribution of  $V_{gb}$  – the genetic variance between populations. This has the  
428 effect of varying the LD contribution without changing the  $F_{ST}$  at causal loci. We selected the maximum  
429 and minimum of  $V_{gb}$  to represent Traits 1 and 2.

### 430 Simulating admixture

431 We simulated the genotypes, local ancestry, and phenotype for 10,000 admixed individuals per generation  
432 under the hybrid isolation (HI) and continuous gene flow (CGF) models by adapting the code from Zaitlen  
433 *et al.* (2017) [26]. We denote the ancestry of a randomly selected individual  $k$  with  $\theta$ , the fraction of their  
434 genome from population A. At  $t = 0$  under the HI model, we set  $\theta$  to 1 for individuals from population A  
435 and 0 if they were from population B such that  $\mathbb{E}(\theta) = p \in \{0.1, 0.2, 0.5\}$  with no further gene flow from  
436 either source population. In the CGF model, population B receives a constant amount  $q$  from population  
437 A in every generation starting at  $t = 0$ . The mean overall proportion of ancestry in the population is  
438 kept the same as the HI model by setting  $q = 1 - (1 - p)^{\frac{1}{t}}$  where  $t$  is the number of generations of gene  
439 flow from A. In every generation, we simulated ancestry-based assortative mating by selecting mates  
440 such that the correlation between their ancestries is  $P \in \{0, 0.3, 0.6, 0.9\}$  in every generation. We do this  
441 by repeatedly permuting individuals with respect to each other until  $P$  falls within  $\pm 0.01$  of the desired  
442 value. It becomes difficult to meet this criterion when  $\mathbb{V}(\theta)$  is small (Fig.1C). To overcome this, we  
443 relaxed the threshold up to 0.04 for some conditions, i.e., when  $\theta \in \{0.1, 0.2\}$  and  $t \geq 50$ . We generated  
444 expected variance in individual ancestry using the expression in Zaitlen *et al.* (2017) [26]. At time  $t$

445 since admixture,  $\mathbb{V}(\theta_t) = \mathbb{V}(\theta_{t-1})\frac{(1+P)}{2}$  under the HI model where  $P$  measures the strength of assortative  
446 mating, i.e, the correlation between the ancestry between individuals in a mating pair. Under the CGF  
447 model,  $\mathbb{V}(\theta_t) = q(1-q)\mathbb{E}(\theta_{t-1})^2 + q(1-q)\{1-2\mathbb{E}(\theta_{t-1})\} + (1-q)\mathbb{V}(\theta_{t-1})\frac{(1+P)}{2}$  (Appendix).

448 We sampled the local ancestry at each  $i^{th}$  locus as  $\gamma_i = \gamma_{if} + \gamma_{im}$  where  $\gamma_{im} \sim Bin(1, \theta_m)$ ,  $\gamma_{if} \sim$   
449  $Bin(1, \theta_f)$  and  $\theta_m$  and  $\theta_f$  represent the ancestry of the maternal and paternal chromosome, respectively.  
450 The global ancestry of the individual is then calculated as  $\theta_k = \sum_{i=1}^m \frac{\gamma_{im} + \gamma_{if}}{2m}$ , where  $m$  is the number of  
451 loci. We sample the genotypes  $x_{im}$  and  $x_{if}$  from a binomial distribution conditioning on local ancestry.  
452 For example, the genotype on the maternal chromosome is  $x_{im} \sim Bin(1, f_i^A)$  if  $\gamma_{im} = 1$  and  $x_{im} \sim$   
453  $Bin(1, f_i^B)$  if  $\gamma_{im} = 0$  where  $f_i^A$  and  $f_i^B$  represent the allele frequency in populations A and B, respectively.  
454 Then, the genotype can be obtained as the sum of the maternal and paternal genotypes:  $x_i = x_{im} + x_{ip}$ .  
455 We calculate the genetic value of each individual as  $g = \sum_{i=1}^m \beta_i x_i$  and the genetic variance as  $\mathbb{V}(g)$ .

## 456 Heritability estimation with GREML

457 We used the `--reml` and `--reml-no-constrain` flags in GCTA [5] to estimate  $\sigma_u^2$  and  $\sigma_v^2$ , the genetic variance  
458 due to genotypes and local ancestry, respectively. We could not run GCTA without noise in the genetic  
459 values so we simulated individual phenotypes with a heritability of  $h^2 = 0.8$  by adding random noise  
460  $e \sim \mathcal{N}(0, V_g \frac{1-h^2}{h^2})$ . We computed three different GRMs, which correspond to different transformations  
461 of the genotypes: (i) standard, (i) Variance or  $V(x)$  scaled, and (ii) LD-scaled.

462 For the standard GRM, the genotypes at the  $i^{th}$  SNP are standardized such that  $z_i = \frac{x_i - 2f_i}{\sqrt{2f_i(1-f_i)}}$ .  
463 For the variance scaled GRM, we computed  $z_i = \frac{x_i - 2f_i}{\sqrt{\mathbb{V}(x_i)}}$  where  $\mathbb{V}(x_i)$  is the sample variance of the  
464 genotypes at the  $i^{th}$  SNP. The LD-scaled GRM conceptually corresponds to standardizing the genotypes  
465 by the SNP covariance, rather than its variance. Let  $\mathbf{X}$  represent the  $n \times m$  *unstandardized* matrix of  
466 genotypes and  $\mathbf{P}$  represent an  $n \times m$  matrix where the  $i^{th}$  column contains the allele frequency of that  
467 SNP. Let  $\mathbf{U}$  be the upper triangular ‘square root’ matrix of the  $m \times m$  SNP covariance matrix  $\mathbf{\Sigma}$  such  
468 that  $\mathbf{\Sigma} = \mathbf{U}'\mathbf{U}$ . Then, the standardized genotypes are computed as  $\mathbf{Z} = (\mathbf{X} - 2\mathbf{P})\mathbf{U}^{-1}$  and the GRM  
469 becomes  $(\mathbf{X} - 2\mathbf{P})\mathbf{\Sigma}^{-1}(\mathbf{X} - 2\mathbf{P})'$  [33]. Similarly, the three GRMs for local ancestry were computed by  
470 scaling local ancestry with (i)  $\sqrt{2\bar{\gamma}_i(1-\bar{\gamma}_i)}$  where we denote  $\bar{\gamma}_i$  as the mean local ancestry at the  $i^{th}$   
471 SNP, or with the (ii) variance, or (iii) covariance of local ancestry, respectively. We estimated  $\sigma_u^2$  and  
472  $\sigma_v^2$  with and without individual ancestry as a fixed effect to correct for any confounding due to genetic  
473 stratification. This was done by using the `--qcovar` flag.

## 474 Heritability estimation with HE regression

475 Haseman-Elston regression with and without ancestry correction was implemented using custom scripts  
476 in R [56]. To estimate  $V_g$  without ancestry correction, we first computed the cross-product of the centered  
477 phenotypes ( $\mathbf{y}$ ), resulting in an  $n \times n$  matrix  $\mathbf{yy}'$ . We stacked the upper-triangular matrix of  $\mathbf{yy}'$  into a  
478 vector and regressed it on the corresponding elements of the GRM ( $\mathbf{\psi}$ ), taking the slope as an estimate  
479 of  $V_g$ :

$$\hat{V}_g = \frac{\sum_{k=1} \sum_{l < k} y_k y_l \psi_{kl}}{\sum_k \sum_{l < k} \psi_{kl}^2}$$

480 To correct for individual ancestry, we followed the approach of Min et al. (2022) [35]. To do this, we  
481 first regressed out the effect of individual ancestry ( $\theta$ ) on phenotype. The regression coefficient can be  
482 expressed as  $\theta(\theta'\theta)^{-1}\theta'$  and the residuals as  $\mathbf{y}^* = (\mathbf{I} - \theta(\theta'\theta)^{-1}\theta')\mathbf{y}$ . Then, we fit the following model:

$$\mathbb{E}(\mathbf{y} * \mathbf{y}^*) = V_g \boldsymbol{\psi} + V_e \mathbf{I} + \delta \boldsymbol{\theta} \boldsymbol{\theta}'$$

483 where  $\boldsymbol{\theta} \boldsymbol{\theta}'$  represents the cross-product of individual ancestry,  $\delta$  represents its corresponding regression  
484 coefficient, and  $V_g$  represents the parameter of interest, i.e., the genetic variance and  $V_e$ , the residual  
485 variance.

486 To demonstrate that the bias in HE estimates arises because of a bias in the estimate of LD contri-  
487 bution, not the genic variance, we carried out a simple simulation where half of the individuals in the  
488 population derive their ancestry from population A and the rest from population B. This is equivalent to  
489 the meta-population under the HI model at  $t = 0$  where  $\mathbb{E}(\theta) = 0.5$ . We simulated genotypes for 1,000  
490 individuals for  $m = 100$  loci where the allele frequencies in populations A and B were set to  $f_A = 0.1$   
491 and  $f_B = 0.8$ , respectively. We standardized the genotypes at each locus  $i$  using the square-root of the  
492 sample variance and assigned effect sizes such that the total genetic variance explained by all loci is equal  
493 to 1, i.e., the effect of the scaled genotype at the  $i^{th}$  locus is  $u_i = \frac{1}{\sqrt{m}}$ . This is equivalent to the effect  
494 size of the unscaled genotypes being  $\beta_i = \frac{1}{\sqrt{m \mathbb{V}(x_i)}}$  where  $\mathbb{V}(x_i)$  is the sample variance at the  $i^{th}$  locus.  
495 We introduced randomness in the direction of the effect by assigning a negative or positive sign to each  
496 locus uniformly at random 100 times to generate 100 traits with the same genic variance but varying LD  
497 contributions. Then, for each trait we computed the two terms in Eq. 2, which should converge to the  
498 genic variance and LD contributions, which represent the genic and LD components to the HE regression  
499 estimate. Fig. S5 shows that in the presence of directional LD, the overall bias is in the HE regression  
500 estimate is due to an exaggerated estimate of the LD contribution.

## 501 Estimating variance explained by GWAS SNPs

502 To decompose the variance explained by GWAS SNPs in African Americans, we needed four quantities:  
503 (i) effect sizes of GWAS SNPs, (ii) their allele frequencies in Africans and Europeans, and (iii) the mean  
504 and variance of global ancestry in African Americans (Equation 1).

505 We retrieved the summary statistics of 26 traits from GWAS catalog [40]. Full list of traits and the  
506 source papers [44, 57–64] are listed in Table S1. To maximize the number of variants discovered, we chose  
507 summary statistics from studies that were conducted in both European and multi-ancestry samples and  
508 that reported the following information: effect allele, effect size, p-value, and genomic position. For birth  
509 weight, we downloaded the data from the Early Growth Genetics (EGG) consortium website [61] since the  
510 version reported on the GWAS catalog is incomplete. For skin pigmentation, we chose summary statistics

511 from the UKB [65] released by the Neale Lab (<http://www.nealelab.is/uk-biobank>) and processed by Ju  
512 and Mathieson [44] to represent effect sizes estimated among individuals of European ancestry. We  
513 also selected summary statistics from Lona-Durazo *et al.* (2019) where effect sizes were meta-analyzed  
514 across four admixed cohorts [57]. Lona-Durazo *et al.* provide summary statistics separately with and  
515 without conditioning on rs1426654 and rs35397 – two large effect variants in *SLC24A5* and *SLC45A2*.  
516 We used the ‘conditioned’ effect sizes and added in the effects of rs1426654 and rs35397 to estimate  
517 genetic variance.

518 We selected independent hits for each trait by pruning and thresholding with PLINK v1.90b6.21 [66]  
519 in two steps as in Ju *et al.* (2020) [44]. We used the genotype data of GBR from the 1000 genome  
520 project [39] as the LD reference panel. We kept only SNPs (indels were removed) that passed the  
521 genome-wide significant threshold (`--clump-p1 5e-8`) with a pairwise LD cutoff of 0.05 (`--clump-r2 0.05`)  
522 and a physical distance threshold of 250Kb (`--clump-kb 250`) for clumping. Second, we applied a second  
523 round of clumping (`--clump-kb 100`) to remove SNPs within 100kb.

524 When GWAS was carried out separately in different ancestry cohorts in the same study, we used  
525 inverse-variance weighting to meta-analyze effect sizes for variants that were genome-wide significant  
526 (p-value  $< 5 \times 10^{-8}$ ) in at least one cohort. This allowed us to maximize the discovery of variants such  
527 as the Duffy null allele that are absent among individuals of European ancestry but polymorphic in other  
528 populations [47].

529 We used allele frequencies from the 1000 Genomes CEU and YRI to represent the allele frequencies  
530 of GWAS SNPs in Europeans and Africans, respectively, making sure that the alleles reported in the  
531 summary statistics matched the alleles reported in the 1000 Genomes. We estimated the global ances-  
532 try of ASW individuals ( $N = 74$ ) with CEU and YRI individuals from 1000 genome (phase 3) using  
533 ADMIXTURE 1.3.0 [67] with  $k=2$  and used it to calculate the mean (proportion of African ancestry =  
534 0.767) and variance (0.018) of global ancestry in ASW. With the effect sizes, allele frequencies, and the  
535 mean and variance in ancestry, we calculated the four components of genetic variance using Equation 1  
536 and expressed them as a fraction of the total genetic variance.

537 Initially, the multi-ancestry summary statistics for a few traits (NEU, WBC, MON, MCH, BAS)  
538 yielded values  $> 1$  for the proportion of variance explained. This is likely because, despite LD pruning,  
539 some of the variants in the model are not independent and tag large effect variants under divergent  
540 selection such as the Duffy null allele, leading to an inflated contribution of LD. We checked this by  
541 calculating the pairwise contribution, i.e.,  $\beta_i \beta_j (f_i^A - f_i^B)(f_j^A - f_j^B)$ , of all SNPs in the model and show  
542 long-range positive LD between variants on chromosome 1 for NEU, WBC, and MON, especially with the  
543 Duffy null allele (Fig. S6A-C). A similar pattern was observed on chromosome 16 for MCH, confirming  
544 our suspicion. This also suggests that for certain traits, pruning and thresholding approaches are not  
545 guaranteed to yield independent hits. To get around this problem, we retained only one association with  
546 the lowest p-value, each from chromosome 1 (rs2814778 for NEU, WBC, and MON) and chromosome 16  
547 (rs13331259 for MCH) (Fig. S6D). For BAS, we observed that the variance explained was driven by a  
548 rare variant (rs188411703, MAF = 0.0024) of large effect ( $\beta = -2.27$ ). We believe this effect estimate  
549 to be inflated and therefore, we removed it from our calculation.

550 As a sanity check, we independently estimated the genetic variance as the variance in polygenic

551 scores, calculated using `--score` flag in PLINK, [66] in ASW individuals. We compared the first estimate  
552 of the genetic variance to the second (Fig. S7) to confirm two things: (i) the allele frequencies, and  
553 mean and variance in ancestry are estimated correctly, and (ii) the variants are more or less independent  
554 in that they do not absorb the effects of other variants in the model. We show that the two estimates  
555 of the genetic variance are strongly correlated ( $r \sim 0.85$ , Fig. S7). The 95% confidence intervals were  
556 calculated by sampling individuals with replacement 10,000 times.

### 557 Code availability

558 We carried out all analyses in R version 4.2.3 [56], PLINK v1.90b6.21 and PLINK 2.0 [66, 68], and GCTA  
559 version 1.94.1 [5]. All code is freely available on [https://github.com/jinguohuang/admix\\_heritability.git](https://github.com/jinguohuang/admix_heritability.git).

### 560 Acknowledgements

561 We thank Iain Mathieson and Doc Edge for helpful comments on the manuscript. This study was funded  
562 by National Institute of General Medical Sciences award R00GM137076 to A.A.Z. The content of this  
563 paper is solely the responsibility of the authors and does not necessarily represent the official views of  
564 the National Institutes of Health.

## 565 References

- 566 1. Lynch, M. & Walsh, B. *Genetics and analysis of quantitative traits* 1–980 (Sinauer Associates, Inc,  
567 1998).
- 568 2. Visscher, P. M., Hill, W. G. & Wray, N. R. Heritability in the genomics era—concepts and miscon-  
569 ceptions. *Nature reviews. Genetics* **9**, 255–66 (4 2008).
- 570 3. Yang, J., Benyamin, B., *et al.* Common SNPs explain a large proportion of the heritability for  
571 human height. *Nature Genetics* **42**, 565–569 (7 2010).
- 572 4. Yang, J., Manolio, T. A., *et al.* Genome partitioning of genetic variation for complex traits using  
573 common SNPs. *Nature Genetics* *2011 43:6* **43**, 519–525 (6 2011).
- 574 5. Yang, J., Lee, S. H., Goddard, M. E. & Visscher, P. M. GCTA: a tool for genome-wide complex  
575 trait analysis. *American journal of human genetics* **88**, 76–82 (1 2011).
- 576 6. Wainschein, P., Jain, D., *et al.* Assessing the contribution of rare variants to complex trait heri-  
577 tability from whole-genome sequence data. *Nature Genetics* *2022 54:3* **54**, 263–273 (3 2022).
- 578 7. Browning, S. R. & Browning, B. L. Population structure can inflate SNP-based heritability esti-  
579 mates. *American Journal of Human Genetics* **89**, 191–193 (1 2011).
- 580 8. Goddard, M. E., Lee, S. H., Yang, J., Wray, N. R. & Visscher, P. M. Response to browning and  
581 browning. *American Journal of Human Genetics* **89**, 193–195 (1 2011).
- 582 9. Kumar, S. K., Feldman, M. W., Rehkopf, D. H. & Tuljapurkar, S. Limitations of GCTA as a solution  
583 to the missing heritability problem. *Proceedings of the National Academy of Sciences of the United*  
584 *States of America* **113**, E61–E70 (1 2016).
- 585 10. Yang, J., Leec, S. H., Wraya, N. R., Goddardd, M. E. & Visscher, P. M. GCTA-GREML accounts  
586 for linkage disequilibrium when estimating genetic variance from genome-wide SNPs. *Proceedings*  
587 *of the National Academy of Sciences* **113**, E4579–E4580 (32 2016).
- 588 11. Lin, Z., Seal, S. & Basu, S. Estimating SNP heritability in presence of population substructure in  
589 biobank-scale datasets. *Genetics* **220** (4 2022).
- 590 12. Border, R., O’Rourke, S., *et al.* Assortative mating biases marker-based heritability estimators.  
591 *Nature Communications* *2022 13:1* **13**, 1–10 (1 2022).
- 592 13. Visscher, P. M., Yang, J. & Goddard, M. E. A Commentary on ‘Common SNPs Explain a Large  
593 Proportion of the Heritability for Human Height’ by Yang *et al.* (2010). *Twin Research and Human*  
594 *Genetics* **13**, 517–524 (6 2010).
- 595 14. De los Campos, G., Sorensen, D. & Gianola, D. Genomic Heritability: What Is It? *PLOS Genetics*  
596 **11**, e1005048 (5 2015).
- 597 15. Rawlik, K., Canela-Xandri, O., Oolliams, J. W. & Tenesa, A. SNP heritability: What are we esti-  
598 mating? *bioRxiv*, 2020.09.15.276121 (2020).



- 599 16. Lara, L. A. C., Pocrnic, I., de P. Oliveira, T., Gaynor, R. C. & Gorjanc, G. Temporal and genomic  
600 analysis of additive genetic variance in breeding programmes. *Heredity* 2021 128:1 **128**, 21–32 (1  
601 2021).
- 602 17. Wojcik, G. L., Graff, M., *et al.* Genetic analyses of diverse populations improves discovery for  
603 complex traits. *Nature* **570**, 514–518 (7762 2019).
- 604 18. Ben-Eghan, C., Sun, R., *et al.* Don't ignore genetic data from minority populations. *Nature* 2021  
605 585:7824 **585**, 184–186 (7824 2020).
- 606 19. Verma, A., Damrauer, S. M., *et al.* The Penn Medicine BioBank: Towards a Genomics-Enabled  
607 Learning Healthcare System to Accelerate Precision Medicine in a Diverse Population. *Journal of*  
608 *Personalized Medicine* **12**, 1974 (12 2022).
- 609 20. Fatumo, S., Chikowore, T., *et al.* A roadmap to increase diversity in genomic studies. *Nature*  
610 *Medicine* 2022 28:2 **28**, 243–250 (2 2022).
- 611 21. The All of Us Research Program Investigators. The “All of Us” Research Program. *New England*  
612 *Journal of Medicine* **381**, 668–676 (7 2019).
- 613 22. Sohail, M., Chong, A. Y., *et al.* Nationwide genomic biobank in Mexico unravels demographic  
614 history and complex trait architecture from 6,057 individuals. *bioRxiv*, 2022.07.11.499652 (2022).
- 615 23. Johnson, R., Ding, Y., *et al.* The UCLA ATLAS Community Health Initiative: Promoting precision  
616 health research in a diverse biobank. *Cell Genomics* **3**, 100243 (1 2023).
- 617 24. Haseman, J. K. & Elston, R. C. The investigation of linkage between a quantitative trait and a  
618 marker locus. *Behavior genetics* **2**, 3–19 (1 1972).
- 619 25. Pfaff, C. L., Parra, E. J., *et al.* Population structure in admixed populations: effect of admixture  
620 dynamics on the pattern of linkage disequilibrium. *American journal of human genetics* **68**, 198–207  
621 (1 2001).
- 622 26. Zaitlen, N., Huntsman, S., *et al.* The Effects of Migration and Assortative Mating on Admixture  
623 Linkage Disequilibrium. *Genetics* **205**, 375–383 (1 2017).
- 624 27. Verdu, P. & Rosenberg, N. A. A General Mechanistic Model for Admixture Histories of Hybrid  
625 Populations. *Genetics* **189**, 1413–1426 (4 2011).
- 626 28. Corre, V. L. & Kremer, A. Genetic Variability at Neutral Markers, Quantitative Trait Loci and  
627 Trait in a Subdivided Population Under Selection. *Genetics* **164**, 1205–1219 (3 2003).
- 628 29. Berg, J. J. & Coop, G. A Population Genetic Signal of Polygenic Adaptation. *PLoS Genetics* **10**,  
629 e1004412 (8 2014).
- 630 30. Bulmer, M. G. The Effect of Selection on Genetic Variability. *The American Naturalist* **105**, 201–  
631 211 (943 1971).
- 632 31. Yair, S. & Coop, G. Population differentiation of polygenic score predictions under stabilizing  
633 selection. *Philosophical Transactions of the Royal Society B* **377** (1852 2022).

- 634 32. Balding, D. J. & Nichols, R. A. A method for quantifying differentiation between populations at  
635 multi-allelic loci and its implications for investigating identity and paternity. *Genetica* **96**, 3–12  
636 (1-2 1995).
- 637 33. Mathew, B., Léon, J. & Sillanpää, M. J. A novel linkage-disequilibrium corrected genomic rela-  
638 tionship matrix for SNP-heritability estimation and genomic prediction. *Heredity* *2017 120:4* **120**,  
639 356–368 (4 2017).
- 640 34. Ma, R. & Dicker, L. H. The Mahalanobis kernel for heritability estimation in genome-wide associ-  
641 ation studies: fixed-effects and random-effects methods (2019).
- 642 35. Min, A., Thompson, E. & Basu, S. Comparing heritability estimators under alternative structures  
643 of linkage disequilibrium. *G3 Genes/Genomes/Genetics* **12** (8 2022).
- 644 36. Zaitlen, N., Pasaniuc, B., *et al.* Leveraging population admixture to characterize the heritability of  
645 complex traits. *Nature Genetics* **46**, 1356–1362 (12 2014).
- 646 37. Zaidi, A. A., Mattern, B. C., *et al.* Investigating the case of human nose shape and climate adap-  
647 tation. *PLoS Genetics* **13**, 2017 (3 2017).
- 648 38. Schraiber, J. G. & Edge, M. D. Heritability within groups is uninformative about differences among  
649 groups: cases from behavioral, evolutionary, and statistical genetics. *bioRxiv*, 2023.11.06.565864  
650 (2023).
- 651 39. Auton, A., Abecasis, G. R., *et al.* A global reference for human genetic variation. *Nature* *2015*  
652 *526:7571* **526**, 68–74 (7571 2015).
- 653 40. Sollis, E., Mosaku, A., *et al.* The NHGRI-EBI GWAS Catalog: knowledgebase and deposition re-  
654 source. *Nucleic Acids Research* **51**, D977–D985 (D1 2023).
- 655 41. Jablonski, N. G. the Evolution of Human Skin and Skin Color. *Annual Review of Anthropology* **33**,  
656 585–623 (1 2004).
- 657 42. Lamason, R. L., Mohideen, M. A. P., *et al.* SLC24A5, a putative cation exchanger, affects pigmen-  
658 tation in zebrafish and humans. *Science* **310**, 1782–1786 (5755 2005).
- 659 43. Beleza, S., Johnson, N. A., *et al.* Genetic architecture of skin and eye color in an African-European  
660 admixed population. *PLoS genetics* **9** (ed Spritz, R. A.) e1003372 (3 2013).
- 661 44. Ju, D. & Mathieson, I. The evolution of skin pigmentation-associated variation in West Eurasia.  
662 *Proceedings of the National Academy of Sciences of the United States of America* **118**, e2009227118  
663 (1 2020).
- 664 45. Nalls, M. A., Wilson, J. G., *et al.* Admixture Mapping of White Cell Count: Genetic Locus Re-  
665 sponsible for Lower White Blood Cell Count in the Health ABC and Jackson Heart Studies. *The*  
666 *American Journal of Human Genetics* **82**, 81–87 (1 2008).
- 667 46. Reich, D., Nalls, M. A., *et al.* Reduced Neutrophil Count in People of African Descent Is Due To  
668 a Regulatory Variant in the Duffy Antigen Receptor for Chemokines Gene. *PLoS Genetics* **5** (ed  
669 Visscher, P. M.) e1000360 (1 2009).

- 670 47. McManus, K. F., Taravella, A. M., *et al.* Population genetic analysis of the DARC locus (Duffy)  
671 reveals adaptation from standing variation associated with malaria resistance in humans. *PLOS*  
672 *Genetics* **13**, e1006560 (3 2017).
- 673 48. Kumar, S. K., Feldman, M. W., Rehkop, D. H. & Tuljapurkar, S. Reply to Yang *et al.*: GCTA  
674 produces unreliable heritability estimates. *Proceedings of the National Academy of Sciences* **113**,  
675 E4581–E4581 (32 2016).
- 676 49. Speed, D., Hemani, G., Johnson, M. R. & Balding, D. J. Improved heritability estimation from  
677 genome-wide SNPs. *American journal of human genetics* **91**, 1011–21 (6 2012).
- 678 50. Speed, D., Cai, N., Johnson, M. R., Nejentsev, S. & Balding, D. J. Reevaluation of SNP heritability  
679 in complex human traits. *Nature Genetics* **49**, 986–992 (7 2017).
- 680 51. Chan, T. F., Rui, X., *et al.* Estimating heritability explained by local ancestry and evaluating strat-  
681 ification bias in admixture mapping from summary statistics. *bioRxiv*, 2023.04.10.536252 (2023).
- 682 52. Schwartzman, A., Schork, A. J., Zablocki, R. & Thompson, W. K. A simple, consistent estimator  
683 of SNP heritability from genome-wide association studies. <https://doi.org/10.1214/19-AOAS1291>  
684 **13**, 2509–2538 (4 2019).
- 685 53. Hou, K., Ding, Y., *et al.* Causal effects on complex traits are similar for common variants across  
686 segments of different continental ancestries within admixed individuals. *Nature Genetics* **2023** *55*:4  
687 **55**, 549–558 (4 2023).
- 688 54. Backman, J. D., Li, A. H., *et al.* Exome sequencing and analysis of 454,787 UK Biobank participants.  
689 *Nature* **2021** *599*:7886 **599**, 628–634 (7886 2021).
- 690 55. Lin, M., Park, D. S., Zaitlen, N. A., Henn, B. M. & Gignoux, C. R. Admixed Populations Improve  
691 Power for Variant Discovery and Portability in Genome-Wide Association Studies. *Frontiers in*  
692 *Genetics* **12**, 673167 (2021).
- 693 56. R Core Team. *R: A Language and Environment for Statistical Computing* R Foundation for Sta-  
694 tistical Computing (Vienna, Austria, 2023).
- 695 57. Lona-Durazo, F., Hernandez-Pacheco, N., *et al.* Meta-analysis of GWA studies provides new insights  
696 on the genetic architecture of skin pigmentation in recently admixed populations. *BMC Genetics*  
697 **20**, 1–16 (1 2019).
- 698 58. Yengo, L., Vedantam, S., *et al.* A saturated map of common genetic variants associated with human  
699 height. *Nature* **610**, 704–712 (7933 2022).
- 700 59. Hoffmann, T. J., Choquet, H., *et al.* A large multiethnic genome-wide association study of adult  
701 body mass index identifies novel loci. *Genetics* **210**, 499–515 (2 2018).
- 702 60. Pulit, S. L., Stoneman, C., *et al.* Meta-Analysis of genome-wide association studies for body fat  
703 distribution in 694 649 individuals of European ancestry. *Human Molecular Genetics* **28**, 166–174  
704 (1 2019).
- 705 61. Warrington, N. M., Beaumont, R. N., *et al.* Maternal and fetal genetic effects on birth weight and  
706 their relevance to cardio-metabolic risk factors. *Nature Genetics* **2019** *51*:5 **51**, 804–814 (5 2019).

- 707 62. Surendran, P., Feofanova, E. V., *et al.* Discovery of rare variants associated with blood pressure  
708 regulation through meta-analysis of 1.3 million individuals. *Nature Genetics* **52**, 1314–1332 (12  
709 2020).
- 710 63. Graham, S. E., Clarke, S. L., *et al.* The power of genetic diversity in genome-wide association studies  
711 of lipids. *Nature* **600**, 675–679 (7890 2021).
- 712 64. Chen, M. H., Raffield, L. M., *et al.* Trans-ethnic and Ancestry-Specific Blood-Cell Genetics in  
713 746,667 Individuals from 5 Global Populations. *Cell* **182**, 1198–1213.e14 (5 2020).
- 714 65. Bycroft, C., Freeman, C., *et al.* The UK Biobank resource with deep phenotyping and genomic  
715 data. *Nature* **562**, 203–209 (7726 2018).
- 716 66. Chang, C. C., Chow, C. C., *et al.* Second-generation PLINK: rising to the challenge of larger and  
717 richer datasets. *GigaScience* **4**, s13742–015–0047–8 (1 2015).
- 718 67. Alexander, D. H., Novembre, J. & Lange, K. Fast model-based estimation of ancestry in unrelated  
719 individuals. *Genome research* **19**, 1655–64 (9 2009).
- 720 68. Purcell, S., Neale, B., *et al.* PLINK: a tool set for whole-genome association and population-based  
721 linkage analyses. *American journal of human genetics* **81**, 559–75 (3 2007).
- 722 69. Chen, G. B. Estimating heritability of complex traits from genome-wide association studies using  
723 IBS-based Haseman-Elston regression. *Frontiers in Genetics* **5**, 72296 (APR 2014).

## 724 Appendix

### 725 A1 Variance in ancestry

726 We denote variance and covariance with  $\mathbb{V}(\cdot)$  and  $\mathbb{C}(\cdot)$  and used the expressions in [26] to generate the  
 727 expected value for the variance in ancestry, i.e.,  $\mathbb{V}(\theta)$ . This is straightforward for the HI model, where at  
 728 time  $t$   $\mathbb{V}(\theta_t) = \mathbb{V}(\theta_{t-1}) \frac{(1+P_{t-1})}{2}$ .  $P_t = \text{Cor}(\theta_m, \theta_f)$  measures the strength of assortative mating, i.e, the  
 729 correlation between the ancestry across mating pairs  $(\theta_m, \theta_f)$  at time  $t$ . For simplicity, we assumed this to  
 730 be constant in every generation, i.e.  $P_t = P_{t-1} = P$  following [26]. Since our notation slightly differs from  
 731 [26], we re-derived the expression for  $V(\theta_t)$  for the CGF model where population B receives a constant  
 732 amount  $q$  of gene flow from population A in every generation. Note, that  $\mathbb{E}(\theta_t) = q + (1 - q)\mathbb{E}(\theta_{t-1})$ .  
 733 Then,

$$\begin{aligned}
 \mathbb{V}(\theta_t) &= \mathbb{E}(\theta_t^2) - \mathbb{E}(\theta_t)^2 \\
 &= q + (1 - q) \mathbb{E} \left[ \left( \frac{\theta_{t-1}^m + \theta_{t-1}^f}{2} \right) \left( \frac{\theta_{t-1}^m + \theta_{t-1}^f}{2} \right) \right] - \{q + (1 - q) \mathbb{E}(\theta_{t-1})\}^2 \\
 &= q + \frac{(1 - q)}{4} \{2 \mathbb{E}(\theta_{t-1}^2) + 2 \mathbb{E}(\theta_{t-1}^m \theta_{t-1}^f)\} - \{q^2 + 2q(1 - q) \mathbb{E}(\theta_{t-1}) + (1 - q)^2 \mathbb{E}(\theta_{t-1})^2\} \\
 &= q + \frac{1 - q}{2} \mathbb{E}(\theta_{t-1}^2) + \frac{1 - q}{2} \mathbb{E}(\theta_{t-1}^m \theta_{t-1}^f) - q^2 - 2q(1 - q) \mathbb{E}(\theta_{t-1}) - (1 - q)^2 \mathbb{E}(\theta_{t-1})^2 \\
 &= q(1 - q) + \frac{1 - q}{2} \{\mathbb{V}(\theta_{t-1}) + \mathbb{E}(\theta_{t-1})^2\} + \frac{1 - q}{2} \{\mathbb{C}(\theta_{t-1}^m, \theta_{t-1}^f) + \mathbb{E}(\theta_{t-1})^2\} - 2q(1 - q) \mathbb{E}(\theta_{t-1}) - \mathbb{E}(\theta_{t-1})^2 \\
 &= q(1 - q) + \frac{1 - q}{2} \mathbb{V}(\theta_{t-1}) + \frac{1 - q}{2} \mathbb{E}(\theta_{t-1})^2 + \frac{1 - q}{2} P_{t-1} \mathbb{V}(\theta_{t-1}) + \frac{1 - q}{2} \mathbb{E}(\theta_{t-1})^2 - 2q(1 - q) \mathbb{E}(\theta_{t-1}) - \mathbb{E}(\theta_{t-1})^2 \\
 &= q(1 - q) + \frac{1 - q}{2} \mathbb{V}(\theta_{t-1}) \{1 + P_{t-1}\} + (1 - q) \mathbb{E}(\theta_{t-1})^2 - 2q(1 - q) \mathbb{E}(\theta_{t-1}) - (1 - q)^2 \mathbb{E}(\theta_{t-1})^2 \\
 &= q(1 - q) \mathbb{E}(\theta_{t-1})^2 + q(1 - q) \{1 - 2 \mathbb{E}(\theta_{t-1})\} + \frac{1 - q}{2} \mathbb{V}(\theta_{t-1}) \{1 + P_{t-1}\}
 \end{aligned}$$

### 734 A2 Genetic variance

735 Let  $y = g + e$ , where  $y$  is the phenotypic value of an individual,  $g$  is the genotypic value, and  $e$  is random  
 736 error. We assume additive effects such that  $g = \sum_{i=1}^m \beta_i x_i$  where  $\beta_i$  is the effect size of the  $i^{\text{th}}$  biallelic  
 737 locus and  $x_i \in \{0, 1, 2\}$  is the number of copies of the trait-increasing allele. Then, the genetic variance  
 738  $V_g$  is:

$$V_g = \mathbb{V} \left( \sum_{i=1}^m \beta_i x_i \right) = \sum_{i=1}^m \beta_i^2 \mathbb{V}(x_i) + \sum_{j \neq i} \beta_i \beta_j \mathbb{C}(x_i, x_j)$$

739 In the following sections, we decompose  $\mathbb{V}(x_i)$  and  $\mathbb{C}(x_i, x_j)$  further as functions of ancestry.

#### 740 A2.1 $\mathbb{V}(x_i)$

741 We first derive  $\mathbb{V}(x_i)$  as a function of ancestry ( $\theta$ ) using the law of total variance:

$$\mathbb{V}(x_i) = \mathbb{E}_{\theta}\{\mathbb{V}(x_i|\theta)\} + \mathbb{V}\{\mathbb{E}_{\theta}(x_i|\theta)\}$$

742 where  $\mathbb{E}_{\theta}$  represents the expectation taken over  $\theta$ .

743 **A2.1.1**  $\mathbb{E}_{\theta}\{\mathbb{V}(x_i|\theta)\}$

744 We derive  $\mathbb{V}(x_i|\theta)$  by further conditioning on the local ancestry at each locus.

$$\mathbb{V}(x_i|\theta) = \mathbb{E}_{\gamma}\{\mathbb{V}(x_i|\gamma, \theta)\} + \mathbb{V}\{\mathbb{E}_{\gamma}(x_i|\gamma, \theta)\}$$

745 where  $\mathbb{E}_{\gamma}$  represents expectation taken over local ancestry. Since we are interested in the variance at

746 a single locus, we will ignore the subscript  $i$  and denote the frequency of the trait-increasing allele in

747 populations A and B with  $f^A$  and  $f^B$ , respectively.

$$\begin{aligned} \mathbb{E}_{\gamma}\{\mathbb{V}(x_i|\gamma, \theta)\} &= \mathbb{V}(x_i|\gamma = 0, \theta) \mathbb{P}(\gamma = 0|\theta) + \mathbb{V}(x_i|\gamma = 1, \theta) \mathbb{P}(\gamma = 1|\theta) + \mathbb{V}(x_i|\gamma = 2, \theta) \mathbb{P}(\gamma = 2|\theta) \\ &= 2f^B(1-f^B)(1-\theta)^2 + \{f^A(1-f^A) + f^B(1-f^B)\}2\theta(1-\theta) + 2f^A(1-f^A)\theta^2 \\ &= (2f^B - 2f^{A^2})(1-2\theta + \theta^2) + (f^A - f^{A^2} + f^B - f^{B^2})(2\theta - 2\theta^2) + (2f^A - 2f^{A^2})\theta^2 \\ &= 2f^B - 2\theta f^B - 2f^{B^2} + 2\theta f^{B^2} + 2\theta f^A - 2\theta f^{A^2} \\ &= 2f^B(1-\theta) - 2f^{B^2}(1-\theta) + 2\theta f^A(1-f^A) \\ &= 2f^B(1-f^B)(1-\theta) + 2\theta f^A(1-f^A) \end{aligned}$$

748 To derive  $\mathbb{V}\{\mathbb{E}_{\gamma}(x|\gamma, \theta)\}$ , note that

$$\begin{aligned} \mathbb{E}_{\gamma}(x|\gamma, \theta) &= \mathbb{E}_{\gamma}\{\mathbb{E}(x|\theta)\} \\ &= \mathbb{E}(x|\gamma = 0, \theta) \mathbb{P}(\gamma = 0|\theta) + \mathbb{E}(x|\gamma = 1, \theta) \mathbb{P}(\gamma = 1|\theta) + \mathbb{E}(x|\gamma = 2, \theta) \mathbb{P}(\gamma = 2|\theta) \\ &= 2\theta f^A + 2(1-\theta) f^B \end{aligned}$$

749 And,

$$\begin{aligned} \mathbb{V}\{\mathbb{E}_{\gamma}(x|\gamma, \theta)\} &= [\mathbb{E}(x|\gamma = 0, \theta) - \mathbb{E}(x|\theta)]^2 \mathbb{P}(\gamma = 0|\theta) \\ &\quad + [\mathbb{E}(x|\gamma = 1, \theta) - \mathbb{E}(x|\theta)]^2 \mathbb{P}(\gamma = 1|\theta) \\ &\quad + [\mathbb{E}(x|\gamma = 2, \theta) - \mathbb{E}(x|\theta)]^2 \mathbb{P}(\gamma = 2|\theta) \\ &= \theta^2 [2f^A - \{2\theta f^A + 2(1-\theta) f^B\}]^2 \\ &\quad + 2\theta(1-\theta) [f^A + f^B - \{2\theta f^A + 2(1-\theta) f^B\}]^2 \\ &\quad + (1-\theta)^2 [2f^B - \{2\theta f^A + 2(1-\theta) f^B\}]^2 \\ &= 2\theta(1-\theta)(f^A - f^B)^2 \end{aligned}$$

**750** Putting this together,

$$\begin{aligned}
 \mathbb{E}\{\mathbb{V}(x_i|\theta)\} &= \mathbb{E}\{2f^B(1-f^B)(1-\theta) + 2\theta f^A(1-f^A) + 2\theta(1-\theta)(f^A-f^B)^2\} \\
 &= 2f^B(1-f^B)\{1-\mathbb{E}(\theta)\} + 2\mathbb{E}(\theta)f^A(1-f^A) + 2\mathbb{E}(\theta-\theta^2)(f^A-f^B)^2 \\
 &= 2f^B(1-f^B)\{1-\mathbb{E}(\theta)\} + 2\mathbb{E}(\theta)f^A(1-f^A) + 2\{\mathbb{E}(\theta)-\mathbb{E}(\theta^2)\}(f^A-f^B)^2 \\
 &= 2f^B(1-f^B)\{1-\mathbb{E}(\theta)\} + 2\mathbb{E}(\theta)f^A(1-f^A) + 2\{\mathbb{E}(\theta)-\mathbb{V}(\theta)-\mathbb{E}(\theta)^2\}(f^A-f^B)^2 \\
 &= 2f^B(1-f^B)\{1-\mathbb{E}(\theta)\} + 2\mathbb{E}(\theta)f^A(1-f^A) + 2\mathbb{E}(\theta)(1-\mathbb{E}(\theta))(f^A-f^B)^2 - 2\mathbb{V}(\theta)(f^A-f^B)^2
 \end{aligned}$$

**751** **A2.1.2**  $\mathbb{V}\{\mathbb{E}_\theta(x_i|\theta)\}$

**752** Recall from the previous section that  $\mathbb{E}_\theta(x_i|\theta) = 2\theta f^A + 2(1-\theta)f^B$ . Then,

$$\begin{aligned}
 \mathbb{V}\{\mathbb{E}_\theta(x_i|\theta)\} &= \mathbb{V}\{2\theta f^A + 2(1-\theta)f^B\} \\
 &= 4\mathbb{V}(\theta)f^{A^2} + 4\mathbb{V}(1-\theta)f^{B^2} + 2\mathbb{C}(2\theta f^A, 2(1-\theta)f^B) \\
 &= 4\mathbb{V}(\theta)f^{A^2} + 4\mathbb{V}(1-\theta)f^{B^2} - 8f^A f^B \mathbb{V}(\theta) \\
 &= 4\mathbb{V}(\theta)(f^A - f^B)^2
 \end{aligned}$$

**753** We are now ready to express  $\mathbb{V}(x_i)$ :

$$\begin{aligned}
 \mathbb{V}(x_i) &= 2f^B(1-f^B)\{1-\mathbb{E}(\theta)\} + 2\mathbb{E}(\theta)f^A(1-f^A) + 2\mathbb{E}(\theta)(1-\mathbb{E}(\theta))(f^A-f^B)^2 \\
 &\quad - 2\mathbb{V}(\theta)(f^A-f^B)^2 + 4\mathbb{V}(\theta)(f^A-f^B)^2 \\
 &= 2\mathbb{E}(\theta)f_i^A(1-f_i^A) + 2\{1-\mathbb{E}(\theta)\}f_i^B(1-f_i^B) \\
 &\quad + 2\mathbb{E}(\theta)\{1-\mathbb{E}(\theta)\}(f_i^A-f_i^B)^2 - 2\mathbb{V}(\theta)(f_i^A-f_i^B)^2
 \end{aligned}$$

**754** Note, that we can also express  $V(x_i)$  as:

$$\mathbb{V}(x_i) = 2f_i(1-f_i) + 2\mathbb{V}(\theta)(f_i^A - f_i^B)^2$$

**755** where the second term is the contribution of population structure to the genetic variance at locus  $i$ .

**756** **A2.2**  $\mathbb{C}(x_i, x_j)$

**757** We can derive  $\mathbb{C}(x_i, x_j)$  using the law of total covariance:

$$\begin{aligned}
\mathbb{C}(x_i, x_j) &= \mathbb{E}_\theta\{\mathbb{C}(x_i, x_j|\theta)\} + \mathbb{C}\{\mathbb{E}_\theta(x_i|\theta), \mathbb{E}_\theta(x_j|\theta)\} \\
&= 0 + \mathbb{C}\{2f_i^A\theta + 2f_i^B(1-\theta), 2f_j^A\theta + 2f_j^B(1-\theta)\} \\
&= \mathbb{C}(2f_i^A\theta, 2f_j^A\theta) + \mathbb{C}(2f_i^A\theta, 2f_j^B(1-\theta)) + \\
&\quad \mathbb{C}(2f_i^B(1-\theta), 2f_j^A\theta) + \mathbb{C}(2f_i^B(1-\theta), 2f_j^B(1-\theta)) \\
&= 4\mathbb{V}(\theta)(f_i^A - f_i^B)(f_j^A - f_j^B)
\end{aligned}$$

758  $\mathbb{E}_\theta\{\mathbb{C}(x_i, x_j|\theta)\} = 0$  because we assume that the loci are unlinked and therefore,  $x_i$  and  $x_j$  are condi-  
759 tionally independent. Putting this all together, we get the genetic variance in admixed populations as  
760 presented in the main text:

$$\begin{aligned}
V_g &= \sum_{i=1}^m \beta_i^2 \mathbb{V}(x_i) + \sum_{j \neq i} \beta_i \beta_j \mathbb{C}(x_i, x_j) \\
&= \sum_{i=1}^m \beta_i^2 2\mathbb{E}(\theta) f_i^A (1 - f_i^A) + \sum_{i=1}^m \beta_i^2 2\{1 - \mathbb{E}(\theta)\} f_i^B (1 - f_i^B) \\
&\quad + \sum_{i=1}^m \beta_i^2 2\mathbb{E}(\theta)\{1 - \mathbb{E}(\theta)\} (f_i^A - f_i^B)^2 + \\
&\quad + \sum_{i=1}^m \beta_i^2 2\mathbb{V}(\theta) (f_i^A - f_i^B)^2] \\
&\quad + \sum_{j \neq i} \beta_i \beta_j 4\mathbb{V}(\theta) (f_i^A - f_i^B) (f_j^A - f_j^B)
\end{aligned}$$

761 The only difference being that in the main text we use  $\mathbb{E}$  instead of  $\mathbb{E}_\theta$  for simplicity. With two ‘unad-  
762 mixed’ source populations with equal number of individuals,  $\mathbb{E}(\theta) = 0.5$  and  $\mathbb{V}(\theta) = \mathbb{E}(\theta)\{1 - \mathbb{E}(\theta)\} = 0.25$   
763 and  $V_g$  reduces to:

$$\begin{aligned}
V_g &= \mathbb{V}\left(\sum_{i=1}^m \beta_i x_i\right) = \sum_{i=1}^m \beta_i^2 \mathbb{V}(x_i) + \sum_{j \neq i} \beta_i \beta_j \mathbb{C}(x_i, x_j) \\
&= \sum_{i=1}^m \beta_i^2 [f_i^A (1 - f_i^A) + f_i^B (1 - f_i^B)] \\
&\quad + \sum_{i=1}^m \beta_i^2 (f_i^A - f_i^B)^2 \\
&\quad + \sum_{i \neq j} \beta_i \beta_j (f_i^A - f_i^B) (f_j^A - f_j^B)
\end{aligned}$$

### 764 A3 The effect of genotype scale on $\hat{V}_g$

765 In the main text, we showed that both GREML and Haseman-Elston regression estimates of  $V_g$  depend  
766 on how the genotypes are scaled. We provide an explanation of this behavior using the Haseman-Elston  
767 (HE) regression estimator, which is asymptotically equivalent to the GREML estimator if effects are  
768 uncorrelated [69] but which, unlike GREML, has a closed-form solution.



## 769 A3.1 No directional LD

### 770 A3.1.1 Scaling by $2f_i(1 - f_i)$

771 First, let's assume a genetic architecture where all loci contribute equally to the genetic variance and there  
 772 is no LD contribution. With the standard scaling, the genotype at a given locus  $i$  is  $z_i = \frac{x_i - 2f_i}{2f_i(1 - f_i)}$  where  
 773  $f_i$  is the frequency of the allele in the population. Under the random-effects model, this is equivalent  
 774 to saying that the unscaled effects are:  $\beta_i \sim \mathcal{N}(0, \frac{\sigma_u^2}{2mf_i(1 - f_i)})$ ,  $\sigma_u^2$  being the parameter of interest. In a  
 775 panmictic population,

$$\begin{aligned} V_g &= \mathbb{V} \left( \sum_{i=1}^m \beta_i x_i \right) = \sum_{i=1}^m \beta_i^2 \mathbb{V}(x_i) \\ &= \sum_{i=1}^m \frac{\sigma_u^2}{2mf_i(1 - f_i)} 2f_i(1 - f_i) \\ &= \sigma_u^2 \end{aligned}$$

776 In an admixed population,

$$\begin{aligned} V_g &= \sum_{i=1}^m \beta_i^2 \{2f_i(1 - f_i) + 2\mathbb{V}(\theta)(f_i^A - f_i^B)^2\} \\ &= \sum_{i=1}^m \frac{\sigma_u^2}{2mf_i(1 - f_i)} \{2f_i(1 - f_i) + 2\mathbb{V}(\theta)(f_i^A - f_i^B)^2\} \\ &= \frac{\sigma_u^2}{m} \sum_{i=1}^m \left\{ 1 + \mathbb{V}(\theta) \frac{(f_i^A - f_i^B)^2}{f_i(1 - f_i)} \right\} \\ &= \sigma_u^2 + \underbrace{\mathbb{V}(\theta) \frac{\sigma_u^2}{m} \sum_{i=1}^m \frac{(f_i^A - f_i^B)^2}{f_i(1 - f_i)}}_{\text{contribution of population structure to the genic variance}} \end{aligned}$$

777 The HE estimator of  $V_g$  is based on the regression of products of (centered) phenotypes  $y_k y_l$  for all pairs  
 778 of individuals  $k \neq l$  on the corresponding entries of the GRM ( $\psi$ ) where  $\psi_{kl} = \frac{\sum_{i=1}^m z_{ik} z_{il}}{m}$  and  $z_{ik}$  is the  
 779 centered and scaled genotype of individual  $k$  for locus  $i$ :

$$\begin{aligned} \hat{V}_g &= \frac{\mathbb{C}(y_k y_l, \psi_{kl})}{\mathbb{V}(\psi_{kl})} \\ &= \frac{\mathbb{E}_{kl}(y_k y_l \psi_{kl}) - \mathbb{E}_{kl}(y_k y_l) \mathbb{E}_{kl}(\psi_{kl})}{\mathbb{E}_{kl}(\psi_{kl}^2) - \mathbb{E}(\psi_{kl})^2} \\ &= \frac{\mathbb{E}_{kl}(y_k y_l \psi_{kl})}{\mathbb{E}_{kl}(\psi_{kl}^2)} \end{aligned}$$

780 Where  $\mathbb{E}_{kl}$  represents the expectation over all  $k \times l$  pairwise comparisons between individuals. It is  
 781 simpler to express the HE estimator in terms of the scaled effects  $u_i \sim \mathcal{N}(0, \frac{\sigma_u^2}{m})$ .

$$\begin{aligned}
 \hat{V}_g &= \frac{\mathbb{E}_{kl}(y_k y_l \psi_{kl})}{\mathbb{E}_{kl}(\psi_{kl}^2)} = \frac{\mathbb{E}_{kl}(\sum_{i=1}^m u_i z_{ik} \sum_{i=1}^m u_i z_{il} \psi_{kl})}{\mathbb{E}_{kl}(\psi_{kl}^2)} \\
 &= \frac{\mathbb{E}_{kl}(\sum_{i=1}^m u_i^2 z_{ik} z_{il} \psi_{kl})}{\mathbb{E}_{kl}(\psi_{kl}^2)} + \frac{\mathbb{E}_{kl}(\sum_{i=1}^m \sum_{j \neq i} u_i u_j z_{ik} z_{jl} \psi_{kl})}{\mathbb{E}_{kl}(\psi_{kl}^2)} \\
 &= \frac{\mathbb{E}(u_i^2) \mathbb{E}_{kl}(\sum_{i=1}^m z_{ik} z_{il} \psi_{kl})}{\mathbb{E}_{kl}(\psi_{kl}^2)} + \frac{\mathbb{E}(u_i u_j) \mathbb{E}_{kl}(\sum_{i=1}^m \sum_{j \neq i} z_{ik} z_{jl} \psi_{kl})}{\mathbb{E}_{kl}(\psi_{kl}^2)} \\
 &= \frac{\mathbb{E}(u_i^2) \mathbb{E}_{kl}(m \psi_{kl}^2)}{\mathbb{E}_{kl}(\psi_{kl}^2)} + 0 = \sigma_u^2
 \end{aligned}$$

782 Where  $\mathbb{E}(u_i)$  and  $\mathbb{E}(u_i u_j)$  represent expectations over random realizations of effect sizes. Thus, the last  
 783 line follows from our assumption that the effect sizes are independent in expectation, i.e.,  $\mathbb{E}(u_i u_j) = 0$ .  
 784 Note, that the estimate is still biased since it does not capture the contribution of population structure.

### 785 A3.1.2 Scaling by $\mathbb{V}(x_i)$

786 Next, we consider the case where the genotypes are standardized instead by the sample variance, i.e.,  
 787  $z_{kl} = \frac{x_{ik} - 2f_i}{\sqrt{\mathbb{V}(x_i)}}$  such that  $\mathbb{V}(z_i) = 1$ . We can derive  $\mathbb{E}(u_i^2)$  corresponding to this scaling by noting that the  
 788 genetic variance is invariant under linear transformations of the genotype [14]:

$$\begin{aligned}
 \sum_{i=1}^m \beta_i^2 \mathbb{V}(x_i) &= \sum_{i=1}^m u_i^2 \mathbb{V}(z_i) \\
 m \mathbb{E}(u_i^2) &= \sigma_u^2 + \mathbb{V}(\theta) \frac{\sigma_u^2}{m} \sum_{i=1}^m \frac{(f_i^A - f_i^B)^2}{f_i(1 - f_i)} \\
 \mathbb{E}(u_i^2) &= \frac{\sigma_u^2}{m} + \mathbb{V}(\theta) \frac{\sigma_u^2}{m^2} \sum_{i=1}^m \frac{(f_i^A - f_i^B)^2}{f_i(1 - f_i)}
 \end{aligned}$$

789 Then, the HE estimator becomes:

$$\begin{aligned}
 \hat{V}_g &= m \mathbb{E}(u_i^2) \\
 &= m \left( \frac{\sigma_u^2}{m} + \mathbb{V}(\theta) \frac{\sigma_u^2}{m^2} \sum_{i=1}^m \frac{(f_i^A - f_i^B)^2}{f_i(1 - f_i)} \right) \\
 &= \sigma_u^2 + \mathbb{V}(\theta) \frac{\sigma_u^2}{m} \sum_{i=1}^m \frac{(f_i^A - f_i^B)^2}{f_i(1 - f_i)}
 \end{aligned}$$

790 Which provides an unbiased estimate of the genic variance. It's important to note that even though we  
 791 assumed effect sizes under a random-effect model, the above result holds under a fixed-effect model as  
 792 long as there is no directional LD. We discuss the implications of directional LD in the following section.

### 793 A3.2 Directional LD

794 Under the standard random-effect model, the effect sizes are assumed to be independent *in expectation*.  
 795 We discussed in the main text how certain processes (e.g. selection and assortative mating) can induce  
 796 directional LD across causal loci. But directional LD might arise even for neutral traits and under the  
 797 random-effects model for any given realization of effects. This can lead to biases both HE and GREML

798 estimates of  $V_g$ , though the direction and reason for bias is different for the two methods. GREML does  
 799 not have a closed-form solution so the exact estimand is difficult to derive. Instead, we develop some  
 800 intuition based on HE regression.

### 801 A3.2.1 Scaling by $\mathbb{V}(x_i)$

802 To do this, let  $\mathbf{u}' = [u_1, u_2, \dots, u_m]$  represent the vector of a given realization of (fixed) effects correspond-  
 803 ing to the standardized genotypes such that each locus contributes equally to  $\sigma_u^2$ , the genic variance, i.e.,  
 804  $u_i^2 = \frac{\sigma_u^2}{m}$ . Let there be positive LD across loci such that all cross-product terms are  $u_i u_j = \frac{\sigma_u^2}{m}$ . Then,  
 805 the genetic variance explained by all loci is:

$$\begin{aligned} V_g &= \sum_{i=1}^m u_i^2 \mathbb{V}(z_i) + \sum_{j \neq i} u_i u_j \mathbb{C}(z_i, z_j) \\ &= \sum_{i=1}^m u_i^2 + \sum_{j \neq i} u_i u_j \mathbb{C}(z_i, z_j) \\ &= \sigma_u^2 + \frac{\sigma_u^2}{m} \sum_{j \neq i} \mathbb{C}(z_i, z_j) \end{aligned} \quad (3)$$

806 where  $\mathbb{C}(z_i, z_j)$  is the LD between the  $i^{\text{th}}$  and  $j^{\text{th}}$  loci that ranges from 0 (no LD) to 1 (perfect LD).  
 807 Thus, the LD contribution to  $V_g$  ranges from 0 to  $(m-1)\sigma_m^2$ . In comparison, the HE estimator is:

$$\begin{aligned} \hat{V}_g &= \frac{\mathbb{E}_{kl} \left( \sum_{i=1}^m u_i^2 z_{ik} z_{il} \psi_{kl} \right)}{\mathbb{E}_{kl}(\psi_{kl}^2)} + \frac{\mathbb{E}_{kl} \left( \sum_{i=1}^m \sum_{j \neq i} u_i u_j z_{ik} z_{jl} \psi_{kl} \right)}{\mathbb{E}_{kl}(\psi_{kl}^2)} \\ &= \frac{\mathbb{E}_{kl} \left( \sum_{i=1}^m \frac{\sigma_u^2}{m} z_{ik} z_{il} \sum_{w=1}^m z_{wk} z_{wl} / m \right)}{\mathbb{E}_{kl} \left( \sum_{i=1}^m z_{ik} z_{il} / m \sum_{w=1}^m z_{wk} z_{wl} / m \right)} + \frac{\mathbb{E}_{kl} \left( \sum_{i=1}^m \sum_{j \neq i} \frac{\sigma_u^2}{m} z_{ik} z_{jl} \sum_{w=1}^m z_{wk} z_{wl} / m \right)}{\mathbb{E}_{kl} \left( \sum_{i=1}^m z_{ik} z_{il} / m \sum_{w=1}^m z_{wk} z_{wl} / m \right)} \\ &= \sigma_u^2 + \sigma_u^2 \frac{\mathbb{E}_{kl} \left( \sum_{i=1}^m \sum_{j \neq i} z_{ik} z_{jl} z_{wk} z_{wl} \right)}{\mathbb{E}_{kl} \left( \sum_{i=1}^m z_{ik} z_{il} \sum_{w=1}^m z_{wk} z_{wl} \right)} \end{aligned} \quad (4)$$

808 This shows that the bias due to directional LD in the HE estimate of  $V_g$  does not come from the genic,  
 809 but the LD component. When there is no LD, e.g. if the population has reached equilibrium after  
 810 generations of random mating, this component goes to zero and both the estimate and  $V_g$  converge to  
 811 the same value – the genic variance. The LD component is maximum when the  $i^{\text{th}}$  and  $j^{\text{th}}$  loci are  
 812 in perfect LD. In this case,  $i$  and  $j$  are exchangeable and the second term of the estimator reduces  
 813 to  $(m-1)\sigma_m^2$ . Thus, HE regression should give an unbiased estimate of  $V_g$ , even in the presence of  
 814 directional LD, but only when LD is perfect. For any other value  $0 < C(z_i, z_j) < 1$ , the estimate is  
 815 biased (Fig. A1). An interpretable, analytical derivation of the second term in Eq. 4 is complicated but  
 816 we illustrate the bias with simulations below.

817 For unlinked markers,  $\mathbb{C}(z_i, z_j)$  is a function of  $4\mathbb{V}(\theta)(f_i^A - f_i^B)(f_j^A - f_j^B)$  (see A2.1). Perfect LD arises  
 818 when (i) both  $f_i^A - f_i^B = 1$  and  $f_j^A - f_j^B = 1$  and (ii)  $\mathbb{V}(\theta)$  is maximum, which occurs at the time of  
 819 admixture when source populations mix equally, i.e.  $\mathbb{E}(\theta) = 0.5$ . To generate a range of LD, we simulated  
 820 an admixed population ( $N = 1,000$ ) with equal number of individuals from populations A and B. Thus,

821  $4\mathbb{V}(\theta) = 4\mathbb{E}(\theta)\{1 - \mathbb{E}(\theta)\} = 1$ . We simulated genotypes for each individual at 50 ‘causal’ loci where  
 822 the difference between the frequencies in the source populations,  $f_i^A - f_i^B \in [0, 1]$  with the condition  
 823 that  $\frac{f_i^A + f_i^B}{2} = 0.5$ . We assigned each locus the same effect size (on the variance-standardized scale) of  
 824  $+1/\sqrt{m}$  summing up to a genic variance of 1. The positive sign ensures positive LD across loci, i.e.,  
 825 all off-diagonal elements of  $\mathbf{u}\mathbf{u}'$  are set to  $1/m$ . For each simulation, we computed the expected and  
 826 estimated LD component using the second terms in Eqs. 3 and 4, respectively, and averaged the results  
 827 over 100 replications.

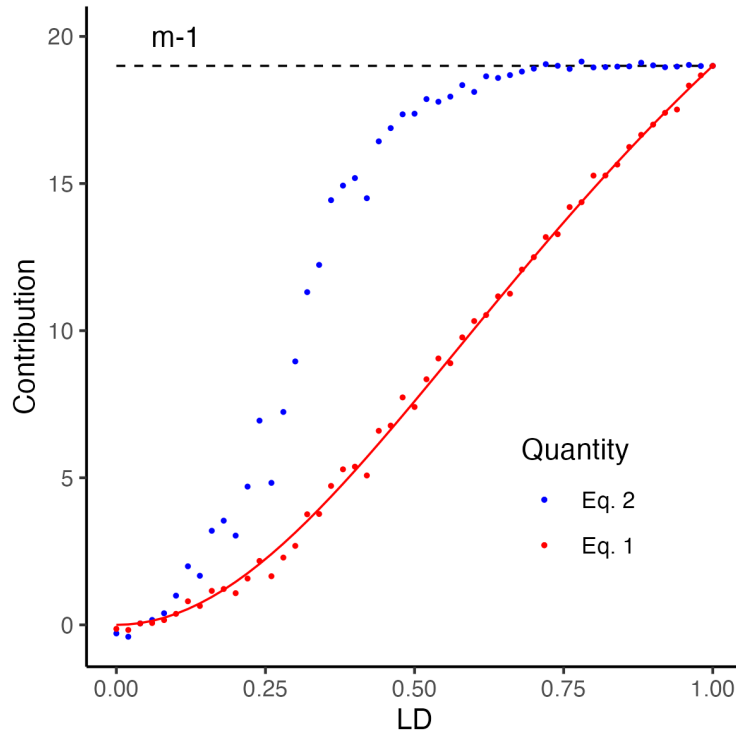


Figure A1: The behavior of the LD contribution (y-axis) to the genetic variance (red) and the Haseman-Elston regression estimate (blue) as a function of LD, i.e.  $\mathbb{C}(z_i, z_j)$  (x-axis). Each point represents the contribution calculated from a random draw of genotypes, given  $\mathbb{C}(z_i, z_j) \propto (f_i^A - f_i^B)(f_j^A - f_j^B)$ . The red line represents the expected LD contribution and the black dashed line represents the contribution expected in the case of perfect LD.

### 828 A3.2.2 Scaling by LD

829 In the main text, we showed that standardizing the genotypes at a locus by its covariance with other loci  
 830 accounts for the bias for GREML and HE estimators. More specifically, the ‘LD-scaled’ genotypes can  
 831 be written as  $\mathbf{Z} = (\mathbf{X} - 2\mathbf{P})\mathbf{U}^{-1}$  where  $\mathbf{P}$  is an  $n \times m$  matrix such that all elements of the  $i^{th}$  column  
 832 contain the frequency of the  $i^{th}$  SNP and  $\mathbf{U}$  is the (upper triangular) square root of the LD matrix, i.e.,  
 833  $\mathbf{\Sigma} = \mathbf{U}'\mathbf{U}$ . Under this scheme, the standardized genotypes are uncorrelated and therefore, the second  
 834 term in Eqs. 3 and 4 are zero. This reduces the estimator to the first term, representing the sum of  
 835 squares of effect sizes, i.e.  $\mathbf{u}'\mathbf{u} = \sum_{i=1}^m u_i^2$ . The effect sizes corresponding to the LD scaled genotypes  
 836 are  $\mathbf{u} = \mathbf{U}\boldsymbol{\beta}$  and the sum of squares is:

$$\mathbf{u}'\mathbf{u} = (\mathbf{U}\boldsymbol{\beta})'(\mathbf{U}\boldsymbol{\beta}) = \boldsymbol{\beta}'\mathbf{U}'\mathbf{U}\boldsymbol{\beta} = \boldsymbol{\beta}'\boldsymbol{\Sigma}\boldsymbol{\beta} = \sum_{i=1}^m \sum_{j=1}^m \beta_i \beta_j \mathbb{C}(x_i, x_j)$$

837 Which captures both the genic and LD contributions and therefore, provides an unbiased estimate of  $V_g$ .

#### 838 A4 Genetic variance after correction for individual ancestry

839 In the main text, we stated that including individual ancestry as a fixed effect in GREML can lead to an  
 840 underestimate of  $V_g$  in the presence of population structure. Mixed effect models deal with fixed effects  
 841 (ancestry in our case) by projecting them out of the phenotypes, and estimating the residual variance.  
 842 This is conceptually equivalent to measuring the residual variance of the regression between phenotype  
 843 and ancestry. As a result, any variance in the phenotype that is explained by ancestry is removed. To  
 844 understand this quantitatively, it is helpful to decompose  $V_g$  into components of variance explained by  
 845 and variance orthogonal to ancestry:

$$\mathbb{V}(g) = \underbrace{\mathbb{V}\{\mathbb{E}(g|\theta)\}}_{\text{variance along ancestry axis}} + \underbrace{\mathbb{E}\{\mathbb{V}(g|\theta)\}}_{\text{variance orthogonal to ancestry axis}}$$

846 We can express the residual variance as:

$$\begin{aligned} \mathbb{E}_{\theta}\{\mathbb{V}(g|\theta)\} &= \mathbb{E}_{\theta}\{\mathbb{V}(\sum_{i=1}^M \beta_i^2 x_i|\theta)\} \\ &= \mathbb{E}_{\theta}\{\sum_{i=1}^M \beta_i^2 \mathbb{V}(x_i|\theta)\} + \mathbb{E}_{\theta}\{\sum_{i \neq j} \beta_i \beta_j \mathbb{C}(x_i, x_j|\theta)\} \\ &= \sum_{i=1}^M \beta_i^2 \mathbb{E}_{\theta}\{\mathbb{V}(x_i|\theta)\} + 0 \\ &= 2\mathbb{E}(\theta) \sum_{i=1}^M \beta^2 f_i^A (1 - f_i^A) + 2\{1 - \mathbb{E}(\theta)\} \sum_{i=1}^M \beta^2 f_i^B (1 - f_i^B) \\ &\quad + 2\mathbb{E}(\theta) \sum_{i=1}^M \beta^2 \{1 - \mathbb{E}(\theta)\} (f_i^A - f_i^B)^2 - 2\mathbb{V}(\theta) \sum_{i=1}^M \beta^2 (f_i^A - f_i^B)^2 \end{aligned}$$

847 Note, that this represents the following components of  $V_g$  from the main text: (1.1) + (1.2) - (1.3). Note,  
 848 that (1.3), which is subtracted out, is always positive and depends on  $\mathbb{V}(\theta)$ . Thus, the residual genetic  
 849 variance will be underestimated, regardless of trait architecture, in the presence of population structure,  
 850 i.e. when  $\mathbb{V}(\theta) > 0$ .

#### 851 A5 Effect size of local ancestry

852 We define local ancestry  $\gamma_i \in \{0, 1, 2\}$  as the number of alleles at locus  $i$  that trace their ancestry to  
 853 population A. Thus, the local ancestry at locus  $i$  in individual  $k$  is a Binomial random variable with

**854**  $\mathbb{E}(\gamma_{i,k}) = 2\theta_k$ . We define the ancestry value of an individual as the weighted sum of their local ancestry:

**855**  $\sum_{i=1}^m \phi_i \gamma_i$  where  $\phi_i = \beta_i(f_i^B - f_i^A)$ .

**856** To show this, note that  $\phi = \mathbb{E}(y|\gamma = 1) - \mathbb{E}(y|\gamma = 0)$  where  $\mathbb{E}(y|\gamma = 1) = \int_{-\infty}^{\infty} yh(y|\gamma = 1) dy$  and  $h$  is

**857** a density function. Our goal is to express  $\phi$  in terms of  $\beta$ , which is equal to  $\mathbb{E}(y|x = 1) - \mathbb{E}(y|x = 0)$ .

**858** Furthermore,  $\mathbb{E}(y|x = 1) = \int_{-\infty}^{\infty} yh(y|x = 1) dy$ . We can express  $h(y|\gamma)$  in terms of  $h(y|x)$  as follows:

$$\begin{aligned} h(y|\gamma = 1) &= h(y|x = 0) \mathbb{P}(x = 0|\gamma = 1) + h(y|x = 1) \mathbb{P}(x = 1|\gamma = 1) + h(y|x = 2) \mathbb{P}(x = 2|\gamma = 1) \\ &= h(y|x = 0)2(1 - f^A)(1 - f^B) + h(y|x = 1)\{f^A(1 - f^B) + f^B(1 - f^A)\} + h(y|x = 2)2f^A f^B \end{aligned}$$

$$\begin{aligned} \mathbb{E}(y|\gamma = 1) &= \int_{-\infty}^{\infty} yh(y|\gamma = 1) dy \\ &= (1 - f^A)(1 - f^B) \int_{-\infty}^{\infty} yh(y|x = 0) dy \\ &\quad + \{f^A(1 - f^B) + f^B(1 - f^A)\} \int_{-\infty}^{\infty} yh(y|x = 1) dy \\ &\quad + f^A f^B \int_{-\infty}^{\infty} yh(y|x = 2) dy \\ &= (1 - f^A)(1 - f^B) \mathbb{E}(y|x = 0) + \{f^A(1 - f^B) + f^B(1 - f^A)\} \mathbb{E}(y|x = 1) + f^A f^B \mathbb{E}(y|x = 2) \\ &= 0 + \{f^A(1 - f^B) + f^B(1 - f^A)\} \beta + f^A f^B 2\beta \\ &= \beta f^A + \beta f^B \end{aligned}$$

**859** Similarly,  $\mathbb{E}(y|\gamma = 0) = 2\beta f^B$  and  $\phi = \mathbb{E}(y|\gamma = 1) - \mathbb{E}(y|\gamma = 0) = \beta(f^B - f^A)$

## **860 A6 Genetic variance due to local ancestry**

$$\begin{aligned} V_\gamma &= \mathbb{V}\left(\sum_{i=1}^m \phi_i \gamma_i\right) \\ &= \sum_{i=1}^m \phi_i^2 \mathbb{V}(\gamma_i) + \sum_{i=1}^m \sum_{j \neq i}^m \phi_i \phi_j \mathbb{C}(\gamma_i, \gamma_j) \end{aligned} \tag{5}$$

**861** We use the law of total variance and covariance to derive  $\mathbb{V}(\gamma_i)$  and  $\mathbb{C}(\gamma_i, \gamma_j)$ :

$$\begin{aligned} \mathbb{V}(\gamma_i) &= \mathbb{E}\{\mathbb{V}(\gamma_i|\theta)\} + \mathbb{V}\{\mathbb{E}(\gamma_i|\theta)\} \\ &= \mathbb{E}\{2\theta(1 - \theta)\} + \mathbb{V}(2\theta) \\ &= 2\mathbb{E}(\theta) - 2\mathbb{E}(\theta^2) + 4\mathbb{V}(\theta) \\ &= 2\mathbb{E}(\theta) - 2\mathbb{V}(\theta) - 2\mathbb{E}(\theta)^2 + 4\mathbb{V}(\theta) \\ &= 2\mathbb{E}(\theta)\{1 - \mathbb{E}(\theta)\} + 2\mathbb{V}(\theta) \end{aligned}$$

$$\begin{aligned} \mathbb{C}(\gamma_i, \gamma_j) &= \mathbb{E}\{\mathbb{C}(\gamma_i, \gamma_j|\theta)\} + \mathbb{C}\{\mathbb{E}(\gamma_i, \gamma_j|\theta)\} \\ &= 0 + \mathbb{C}(2\theta, 2\theta) = 4\mathbb{V}(\theta) \end{aligned}$$

$$\begin{aligned} V_\gamma &= 2\mathbb{E}(\theta)\{1 - \mathbb{E}(\theta)\} \sum_{i=1}^m \phi_i^2 + 2\mathbb{V}(\theta) \sum_{i=1}^m \phi_i^2 + 4\mathbb{V}(\theta) \sum_{i=1}^m \sum_{j \neq i} \phi_i \phi_j \\ &= 2\mathbb{E}(\theta)\{1 - \mathbb{E}(\theta)\} \sum_{i=1}^m \beta_i^2 (f_i^B - f_i^A)^2 \\ &\quad + 2\mathbb{V}(\theta) \sum_{i=1}^m \beta_i^2 (f_i^B - f_i^A)^2 \\ &\quad + 4\mathbb{V}(\theta) \sum_{i=1}^m \sum_{j \neq i} \beta_i \beta_j (f_i^B - f_i^A)(f_j^B - f_j^A) \end{aligned}$$

## 862 Supplement

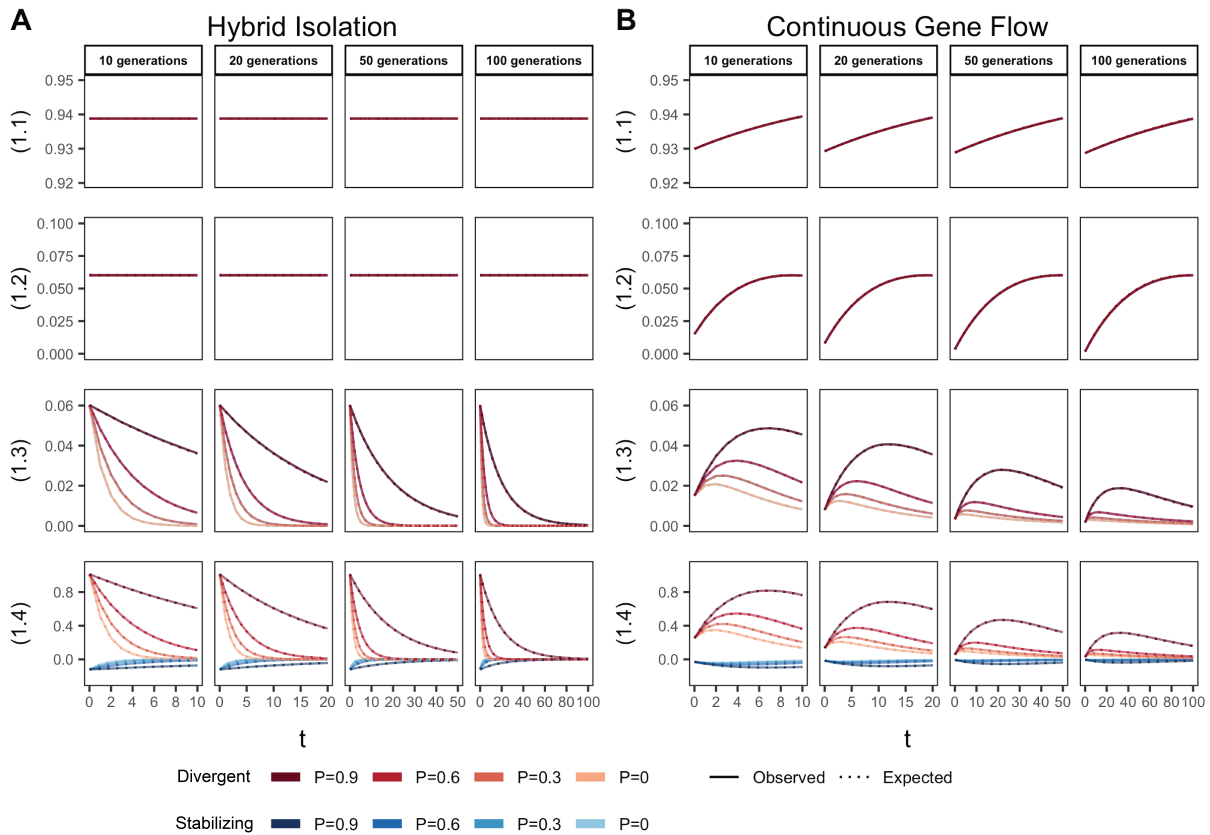


Figure S1: The behavior of the four components of genetic variance in admixed populations under the (A) HI and (B) CGF models. We assume that the mean ancestry proportion in the population is 0.5. The solid lines represent values observed in simulations averaged across ten replicates and the dotted lines represent the expected values based on Eq. 1 of the main text. The red and blue lines represent values for traits 1 and 2, respectively.  $P$  indicates the strength of assortative mating.  $P=0.6$  is missing for simulations run for 50 and 100 generations and  $\theta \in \{0.1, 0.2\}$  due to the difficulty in finding mate pairs (Methods).



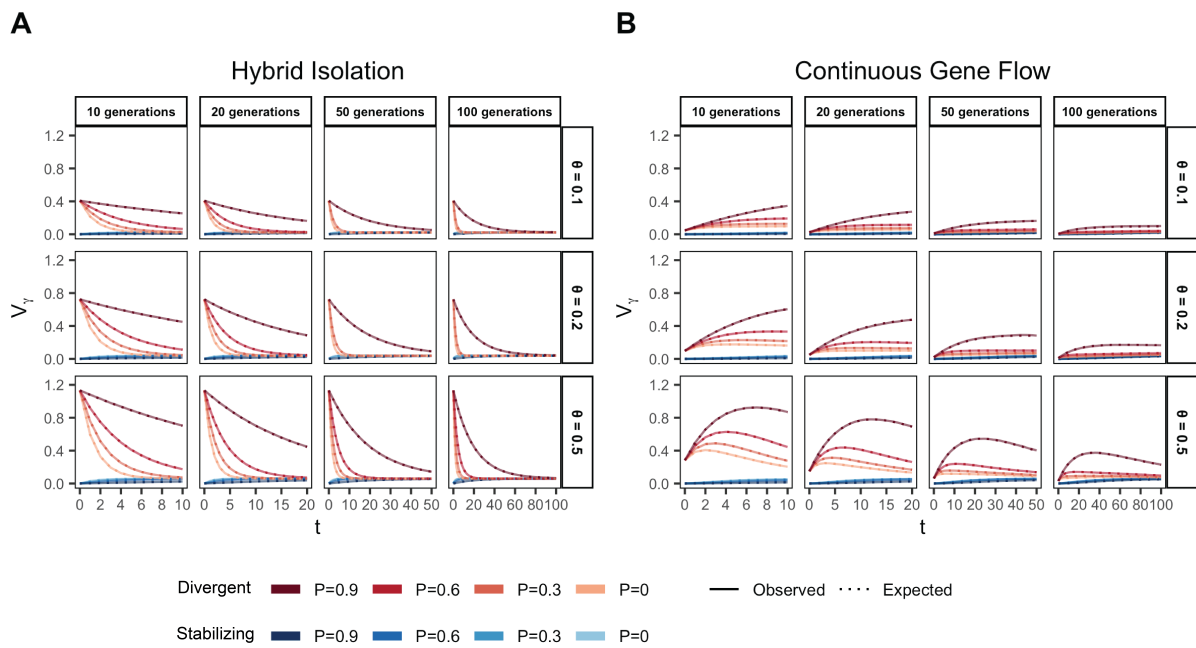


Figure S2: The behavior of the genetic variance due to local ancestry in admixed populations under the (A) HI and (B) CGF models. The solid lines represent values observed in simulations averaged across ten replicates and the dotted lines represent the expected values based on Eq. 1 of the main text. The red and blue lines represent values for traits 1 and 2, respectively.  $P$  indicates the strength of assortative mating.  $P=0.6$  is missing for simulations run for 50 and 100 generations and  $\theta \in \{0.1, 0.2\}$  due to the difficulty in finding mate pairs (Methods).

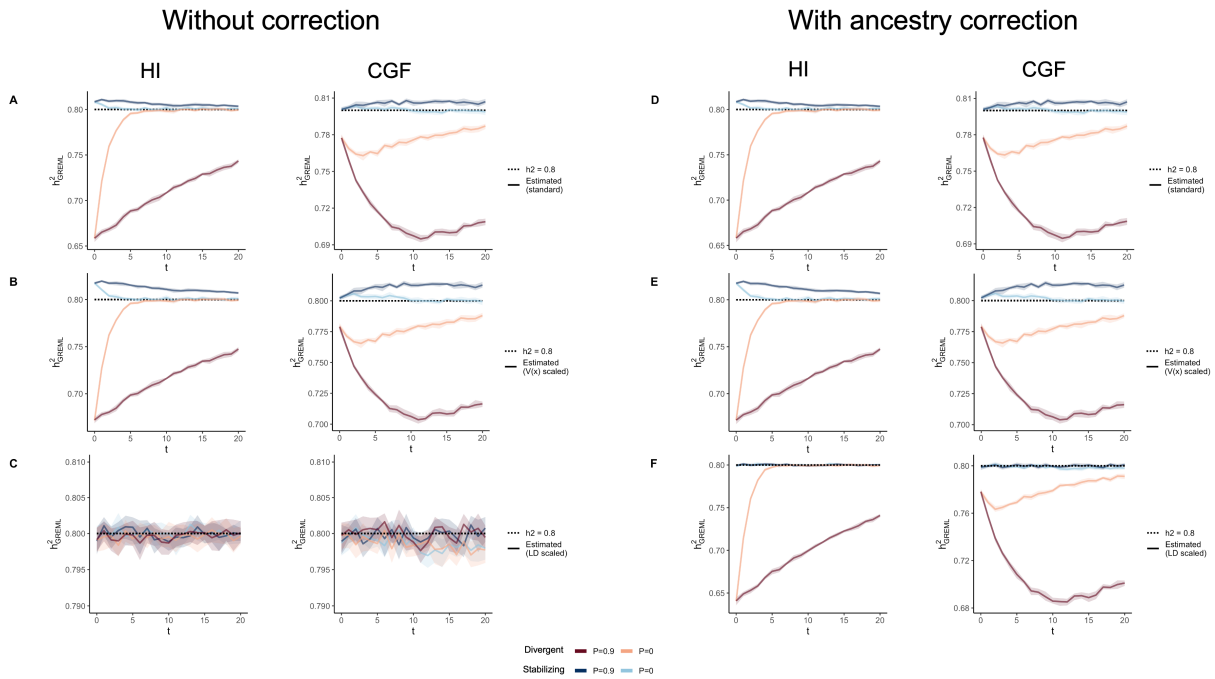


Figure S3: The behavior of GREML estimates of SNP heritability ( $\hat{h}_{snp}^2$ ) in admixed populations under the HI (left column) and CGF (right column) models either without (A-C) or with (D-F) individual ancestry as a fixed effect. The solid lines represent  $\hat{h}_{snp}^2$  averaged across ten replicates, with red and blue colors representing estimates for traits under divergent and stabilizing selection, respectively. (A, D) show the behavior of  $\hat{h}_{snp}^2$  for the default scaling, (B, E) shows  $\hat{h}_{snp}^2$  when the genotype at a locus is scaled by its sample variance ( $\mathbb{V}(x)$  scaled), and (C, F) when it is scaled by the sample covariance (LD scaled). The shaded area represents the 95% confidence bands generated by bootstrapping (sampling with replacement 100 times) the point estimate reported by GCTA. The black dotted lines represent the expected heritability value given the simulation settings ( $h^2 = 0.8$ ).  $P$  indicates the strength of assortative mating

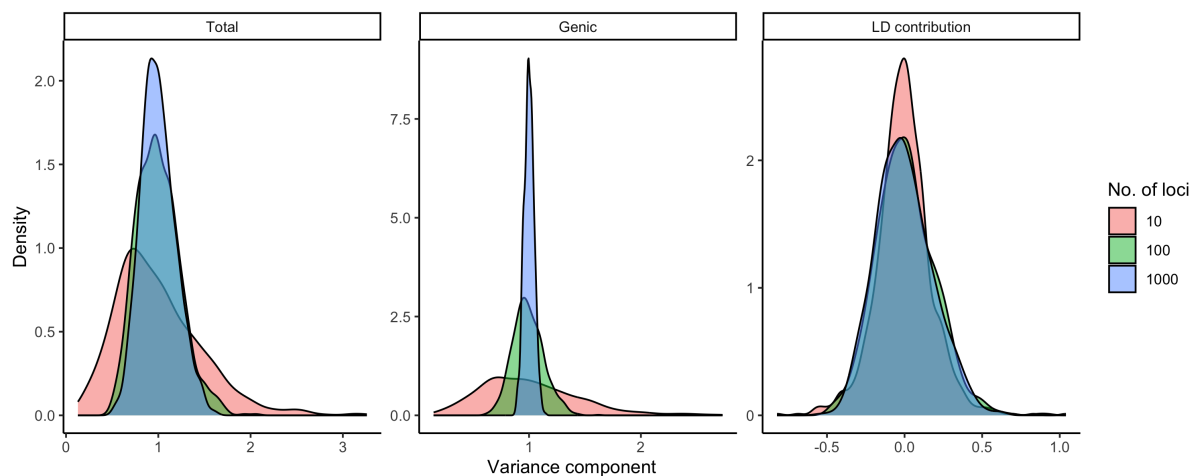


Figure S4: Distribution of the total genetic variance (left), genic variance (middle), and LD component (right) for a neutral trait simulated by drawing effects for 10, 100, or 1,000 causal variants in ASW. The total genetic variance is the sum of the genic and LD components.

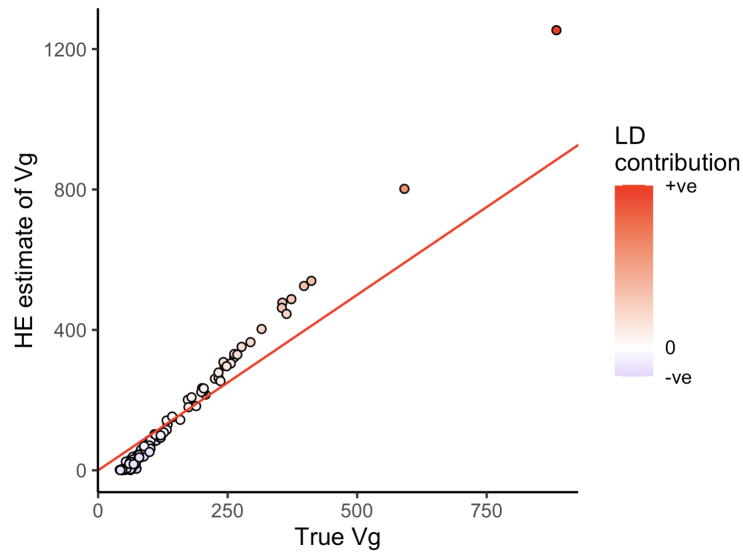


Figure S5: The effect of directional LD on Haseman-Elston estimate of genetic variance ( $V_g$ ). Each individual point is an independent simulation where the effects were drawn from a normal distribution and applied to genotypes from an admixed population (Methods). The solid red line shows the  $y = x$  line and the color of each point represents the contribution of LD to  $V_g$ .

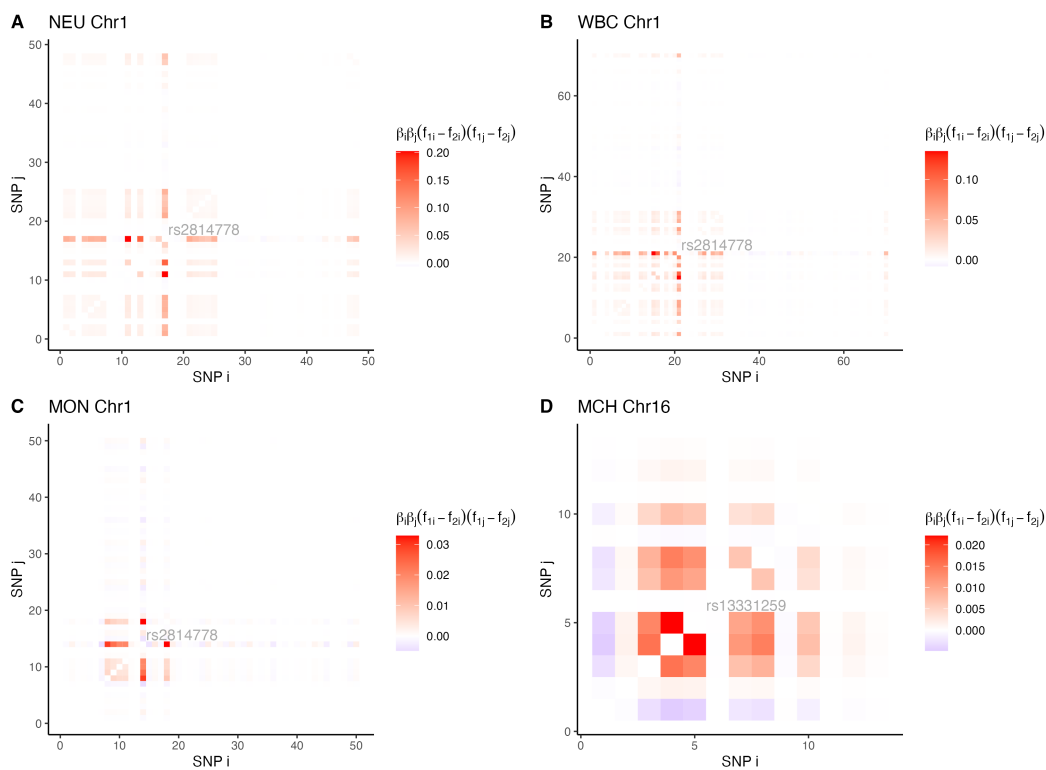


Figure S6: The LD contribution to the variance explained by variant pairs for (A) neutrophil counts (NEU), (B) white blood count (WBC), (C) monocyte count (MON), and (D) mean corpuscular hemoglobin (MCH). Only chromosomes where we suspected there was a disproportionate contribution to the variance explained are shown.

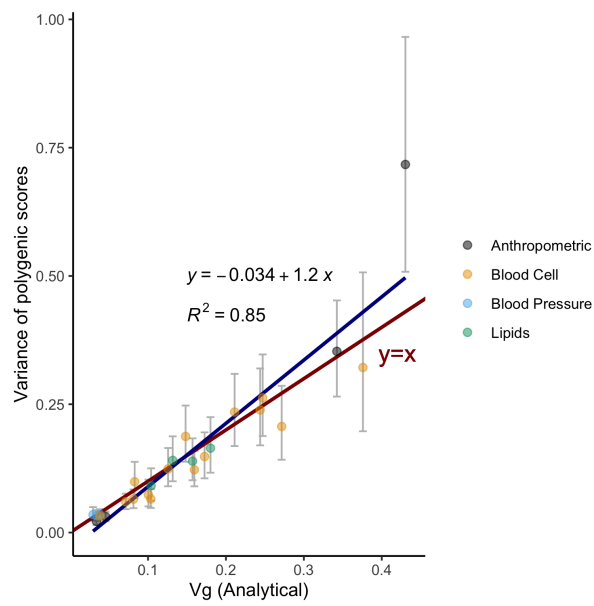


Figure S7: Expected variance explained estimated using Equation 1 (x-axis) vs the observed variance in polygenic scores (y-axis) in ASW are strongly correlated. Confidence intervals were generated by non-parametric bootstrap (Methods).