

# Assortative mating and the dynamical decoupling of genetic admixture levels from phenotypes that differ between source populations

Jaehee Kim<sup>1</sup>, Michael D. Edge<sup>2</sup>, Amy Goldberg<sup>3</sup>, and Noah A. Rosenberg<sup>1</sup>

<sup>1</sup>*Department of Biology, Stanford University, Stanford, CA 94305, USA*

<sup>2</sup>*Department of Evolution and Ecology, University of California, Davis, Davis, CA 95616, USA*

<sup>3</sup>*Department of Evolutionary Anthropology, Duke University, Durham, NC 27708, USA*

## Abstract

Source populations for an admixed population can possess distinct patterns of genotype and phenotype at the beginning of the admixture process. Such differences are sometimes taken to serve as markers of ancestry—that is, phenotypes that are initially associated with the ancestral background in one source population are taken to reflect ancestry in that population. Examples exist, however, in which genotypes or phenotypes initially associated with ancestry in one source population have decoupled from overall admixture levels, so that they no longer serve as proxies for genetic ancestry. We develop a mechanistic model for describing the joint dynamics of admixture levels and phenotype distributions in an admixed population. The approach includes a quantitative-genetic model that relates a phenotype to underlying loci that affect its trait value. We consider three forms of mating. First, individuals might assort in a manner that is independent of the overall genetic admixture level. Second, individuals might assort by a quantitative phenotype that is initially correlated with the genetic admixture level. Third, individuals might assort by the genetic admixture level itself. Under the model, we explore the relationship between genetic admixture level and phenotype over time, studying the effect on this relationship of the genetic architecture of the phenotype. We find that the decoupling of genetic ancestry and phenotype can occur surprisingly quickly, especially if the phenotype is driven by a small number of loci. We also find that positive assortative mating attenuates the process of dissociation in relation to a scenario in which mating is random with respect to genetic admixture and with respect to phenotype. The mechanistic framework suggests that in an admixed population, a trait that initially differed between source populations might be a reliable proxy for ancestry for only a short time, especially if the trait is determined by relatively few loci. The results are potentially relevant in admixed human populations, in which phenotypes that have a perceived correlation with ancestry might have social significance as ancestry markers, despite declining correlations with ancestry over time.

## Author Summary

Admixed populations are populations that descend from two or more populations that had been separated for a long time at the beginning of the admixture process. The source populations typically possess distinct patterns of genotype and phenotype. Hence, early in the admixture process, phenotypes of admixed individuals can provide information about the extent to which these individuals possess ancestry in a specific source population. To study correlations between admixture levels and phenotypes that differ between source populations, we construct a genetic and phenotypic model of the dynamical process of admixture. Under the model, we show that correlations between admixture levels and these phenotypes dissipate over time—especially if the genetic architecture of the phenotypes involves only a small number of loci, or if mating in the admixed population is random with respect to both the admixture levels and the phenotypes. The result has the implication that a trait that once reflected ancestry in a specific source population might lose this ancestry correlation. As a consequence, in human

44 populations, after a sufficient length of time, salient phenotypes that can have social meaning as ancestry  
45 markers might no longer bear any relationship to genome-wide genetic ancestry.

## 46 1 Introduction

47 Admixed populations descend from two or more source groups that have long been separated and that likely  
48 possessed distinct patterns of genotype and phenotype at the beginning of the admixture process. Among  
49 individuals in an admixed population in the generations immediately after its founding, admixture levels  
50 from any specific source population are highly heterogeneous [1,2]. For admixed individuals, measurements  
51 of specific genotypes and phenotypes that differ in frequency or distribution between source populations can  
52 often provide reasonable estimates of individual levels of genetic ancestry in the particular source popula-  
53 tions [3,4]—and for some phenotypes, such measurements might even be commonly regarded by researchers,  
54 societies, or admixed individuals themselves as proxies for overall genetic ancestry [5–7].

55 Examples exist, however, in which genotypes or phenotypes initially associated with ancestry in one  
56 source population are decoupled from overall admixture levels, so that they no longer serve as tight prox-  
57 ies for ancestry [5, 6, 8–13]. For example, in human genetics, consider skin pigmentation and eye color,  
58 observable traits for which the phenotypic distribution differs substantially between sub-Saharan African  
59 and European populations. In the Cape Verdean admixed population, descended from European and West  
60 African sources, measurements of skin pigmentation and eye color are correlated with sub-Saharan African  
61 genetic ancestry [11]. At the same time, the correlations between phenotype and ancestry are imperfect;  
62 many individuals with a high proportion of sub-Saharan African genetic ancestry have skin pigmentation and  
63 eye color traits in a range more typical of individuals with higher European genetic ancestry, and vice versa.  
64 Similar patterns of incomplete correlation with overall genetic ancestry hold for genotypes that underlie  
65 these phenotypes [11].

66 How does ancestry level become decoupled from genotype and phenotype in an admixed population?  
67 Parra et al. [8] proposed one scenario for this decoupling, using an example of assortative mating by a phe-  
68 notype correlated with ancestry in Brazil. Parra et al. suggested that in Brazil, assortative mating is largely  
69 dependent on “color,” a phenotypic measure based to a large extent on skin pigmentation. According to  
70 this hypothesis, in a population descended from source groups with substantially different skin pigmentation  
71 distributions (say, sub-Saharan Africans and Europeans), similarity according to a phenotype correlated  
72 with genetic ancestry (say, *color*) increases the probability that a pair is a mating pair. Mating probabilities  
73 for pairs of individuals are more closely related to the phenotype than to overall sub-Saharan African or  
74 European genetic admixture levels *per se*. Whereas in the early generations of such a process, the pheno-  
75 type would strongly reflect genetic ancestry, after a sufficient length of time with assortative mating by the  
76 phenotype, phenotypic variation would be maintained, but with similar genetic ancestry distributions for  
77 individuals with substantially different phenotype (Figure 1). Only at genes associated with the phenotype  
78 and their nearby linked genomic regions would genetic ancestry and the phenotype be associated.

79 Could genetic ancestry in an admixed population become almost entirely decoupled from the phenotypes  
80 that differ between its source populations? This scenario is intriguing, as it would eliminate any connection  
81 between visible phenotypic markers of genetic ancestry and the genetic ancestry itself; the phenotype of an  
82 individual on a trait such as skin pigmentation would reveal little information about the genetic ancestry  
83 of molecular characters in an individual—other than for skin pigmentation genes and their closest genomic

84 neighbors—nor about the total genomic ancestry of the individual.

85 We develop a mechanistic model describing the joint dynamics of admixture levels and phenotype dis-  
86 tributions in an admixed population. The approach includes a quantitative-genetic model that relates a  
87 phenotype to underlying loci that affect its trait value. We consider three forms of mating. First, individuals  
88 might mate randomly, or assort independently of the overall genetic admixture level. Second, individuals  
89 might assort by a phenotype initially correlated with the genetic admixture level, but that is not identical  
90 to it. Third, following studies that have detected evidence of assortative mating by genetic admixture levels  
91 or phenotypes that are tightly connected to them [14,15], individuals might assort by the genetic admixture  
92 level itself. Under the model, we explore the relationship between genetic admixture level and phenotype  
93 over time, studying the effect of the genetic architecture of the phenotype. We find that the decoupling of  
94 genetic ancestry and phenotype is not only possible, but it can occur surprisingly quickly, especially if the  
95 quantitative phenotype is driven by a small number of loci. Moreover, assortative mating is not required for  
96 genotype and phenotype to become decoupled. Indeed, assortative mating attenuates the process compared  
97 with a scenario in which mating is random with respect to admixture and with respect to phenotype.

## 98 2 Model

### 99 2.1 Population Model

100 Our mechanistic admixture model closely follows the model of Verdu, Goldberg, and Rosenberg [1, 16,  
101 17], building on earlier related models [18–20]. We start with individuals in each of two isolated source  
102 populations,  $S_1$  and  $S_2$ . At the founding of an admixed population ( $g = 0$ ), a founding parental pool  $H_0^{\text{par}}$   
103 is formed, containing fraction  $s_{1,0}$  from population  $S_1$  and  $s_{2,0}$  from population  $S_2$ . That is, a random  
104 individual in  $H_0^{\text{par}}$  originates from population  $S_1$  with probability  $s_{1,0}$  and from  $S_2$  with probability  $s_{2,0}$ .  
105 This choice requires  $s_{1,0} + s_{2,0} = 1$  and  $0 \leq s_{1,0}, s_{2,0} \leq 1$ . The individuals in the founding parental pool  
106 mate according to a mating model (Section 2.3) and produce generation  $g = 1$  of admixed offspring ( $H_1$ ).

107 In subsequent generations ( $g \geq 1$ ), in forming an admixed population  $H_{g+1}$  at generation  $g + 1$ , three  
108 populations contribute to its parental pool  $H_g^{\text{par}}$ : the source populations ( $S_1$  and  $S_2$ ) and the admixed  
109 population ( $H_g$ ) of the previous generation, with fractional contributions  $s_{1,g}$ ,  $s_{2,g}$ , and  $h_g$ , respectively.  
110 Here,  $s_{1,g}$ ,  $s_{2,g}$ , and  $h_g$  represent probabilities for a random individual in  $H_g^{\text{par}}$  to originate from populations  
111  $S_1$ ,  $S_2$ , and  $H_g$ , with constraints  $s_{1,g} + s_{2,g} + h_g = 1$  and  $0 \leq s_{1,g}, s_{2,g}, h_g \leq 1$ . Offspring resulting from  
112 mating among individuals in the parental pool  $H_g^{\text{par}}$  define the admixed population  $H_{g+1}$ . A schematic of  
113 the admixture model appears in Figure 2.

114 The total admixture fraction represents the proportion of the genome of an individual originating from a  
115 specific ancestral population,  $S_1$  or  $S_2$ . We denote an individual’s admixture fraction from source population  
116  $S_1$  at generation  $g$  by  $H_{A,g}$ , with the  $A$  indicating consideration of autosomal genetic loci. Given a pair of  
117 individuals with admixture fractions  $H_{A,g}^{(1)}$  and  $H_{A,g}^{(2)}$ , the ancestry of their offspring is deterministically set  
118 to the mean of the admixture fractions of the parents:  $H_{A,g+1} = \frac{1}{2} (H_{A,g}^{(1)} + H_{A,g}^{(2)})$ . The possible values for  
119 the admixture fraction at generation  $g$  are  $0, 1/2^g, 2/2^g, \dots, (2^g - 1)/2^g, 1$ .

## 120 2.2 Quantitative Trait Model

To model a phenotype, we adopt the quantitative trait model of Edge and Rosenberg [21,22]. We assume each individual is diploid and that  $k$  biallelic autosomal loci, each with the same effect size, additively determine the value of a quantitative trait. At each trait locus, we denote the allelic type more prevalent in population  $S_1$  than in population  $S_2$  as allelic type “1”, and the other allelic type as “0”. The choice is arbitrary in case the allele frequency is the same in the two populations. A diploid individual’s genotype at locus  $i$ ,  $i \in \{1, 2, \dots, k\}$ , and allele  $j$ ,  $j \in \{1, 2\}$ , is represented by a random indicator variable  $L_{ij}$ :

$$L_{ij} = \begin{cases} 1 & \text{if the allele has type “1”} \\ 0 & \text{if the allele has type “0”} . \end{cases}$$

Let  $M$  be a random variable representing an individual’s population membership, considering individuals only from the source populations  $S_1$  and  $S_2$ , and define allele frequencies for allelic type “1” at each locus given the population membership:

$$P(L_{ij} = 1 \mid M = S_1) = p_i$$

$$P(L_{ij} = 1 \mid M = S_2) = q_i.$$

121 Here,  $j$  can be either 1 or 2. Because we define allelic type “1” to be more common in population  $S_1$  than  
 122 in population  $S_2$ ,  $0 \leq q_i \leq p_i \leq 1$ .

An individual’s trait value is determined by a sum of contributions across loci. At each locus, we denote an allele that increases the trait value by “+” and the other allele by “-”. The total quantitative trait value  $T$  of an individual given a multilocus genotype is an individual’s total number of “+” alleles. Whether the “1” allelic type or “0” type is the “+” allele at locus  $i$  is determined by a random variable  $X_i$ , following Edge and Rosenberg [21,22]:

$$X_i = \begin{cases} 1 & \text{if allelic type “1” is “+” allele at locus } i \\ 0 & \text{if allelic type “0” is “+” allele at locus } i. \end{cases}$$

123 For a given set of values  $\{X_1, X_2, \dots, X_k\}$  for  $k$  quantitative trait loci, the total trait value for a diploid  
 124 individual is equal to the total number of “+” alleles carried by the individuals, or:

$$T \mid_{\{X_1, X_2, \dots, X_k\}} = \sum_{\{i: X_i=1\}} \sum_{j=1}^2 L_{ij} + \sum_{\{i: X_i=0\}} \sum_{j=1}^2 (1 - L_{ij}). \quad (1)$$

125 This quantity takes values in  $\{0, 1, \dots, 2k\}$ . An example of the quantitative trait model appears in Figure 3.

126 We adopt two scenarios for the  $X_i$ . We first consider an idealized case in which the number of “1” alleles  
 127 is perfectly correlated with the trait value:  $P(X_i = 1) = 1$  and  $P(X_i = 0) = 0$  for all  $i = 1, 2, \dots, k$ , so  
 128 that allelic type “1” is the “+” allele and allelic type “0” is the “-” allele for all loci. Because we define  
 129 “1” to be the more frequent allelic type in source population  $S_1$ , individuals from  $S_1$  are more likely to have  
 130 a larger trait value than are individuals from  $S_2$ . This scenario considers a case in which the phenotype is  
 131 systematically different between populations 1 and 2, and is depicted in the diagram in Figure 1.

132 The second case is detailed in Edge and Rosenberg [21, 22], where the trait is selectively neutral during  
 133 the population divergence and neither allele is preferentially correlated with the trait allele:  $P(X_i = 1) =$   
 134  $P(X_i = 0) = \frac{1}{2}$  for all  $i = 1, 2, \dots, k$ , so that at each locus, allelic type “1” and allelic type “0” have equal  
 135 probability of being the “+” allele.

### 136 2.3 Mating Model

137 We next describe three different mating models: (1) random mating, in which any two individuals have the  
 138 same chance of reproducing, independent of phenotype or ancestry; (2) assortative mating by ancestry, in  
 139 which the probability that two individuals reproduce depends on their ancestries; and (3) assortative mating  
 140 by phenotype, in which the probability that two individuals reproduce depends on their trait values.

141 In each generation  $g$ , the parental pool  $H_g^{\text{par}}$  contains  $2N$  individuals,  $N$  female and  $N$  male. The  
 142 admixture fraction from source population  $S_1$  and trait value of a female individual  $i$  are denoted by  $H_{A,g}^{(i),f}$   
 143 and  $T_g^{(i),f}$ , respectively. Analogous quantities for a male  $j$  are  $H_{A,g}^{(j),m}$  and  $T_g^{(j),m}$ . We construct an  $N \times N$   
 144 mating probability matrix  $M$ , whose entry  $m_{ij}$  represents the probability that a female  $i$  and a male  $j$  mate:

$$\begin{matrix}
 (H_{A,g}^{(1),f}, T_g^{(1),f}) \\
 (H_{A,g}^{(2),f}, T_g^{(2),f}) \\
 \vdots \\
 (H_{A,g}^{(N-1),f}, T_g^{(N-1),f}) \\
 (H_{A,g}^{(N),f}, T_g^{(N),f})
 \end{matrix}
 \begin{pmatrix}
 (H_{A,g}^{(1),m}, T_g^{(1),m}) & (H_{A,g}^{(2),m}, T_g^{(2),m}) & \cdots & (H_{A,g}^{(N-1),m}, T_g^{(N-1),m}) & (H_{A,g}^{(N),m}, T_g^{(N),m}) \\
 m_{11} & m_{12} & \cdots & m_{1(N-1)} & m_{1N} \\
 m_{21} & m_{22} & \cdots & m_{2(N-1)} & m_{2N} \\
 \vdots & \vdots & \ddots & \vdots & \vdots \\
 m_{(N-1)1} & m_{(N-1)2} & \cdots & m_{(N-1)(N-1)} & m_{(N-1)N} \\
 m_{N1} & m_{N2} & \cdots & m_{N(N-1)} & m_{NN}
 \end{pmatrix}.
 \tag{2}$$

145 In the absence of selection, every individual in the population must have the same expected number  
 146 of offspring irrespective of ancestry or phenotype. We assume that the expected number of offspring of an  
 147 individual is proportional to the expected number of matings of the individual. This quantity is the sum of  
 148 mating probabilities across all mates available for an individual. Therefore, the equal-offspring requirement  
 149 translates into an assumption of equal row sums for females and equal column sums for males in the mating  
 150 matrix in Eq. 2. Note that this assumption of equal numbers of offspring independent of ancestry and  
 151 phenotype accords with a standard property of assortative mating models that assortative mating on its own  
 152 does not alter allele frequencies over time [23–28].

153 The mating probability  $m_{ij}$  in Eq. 2 between female  $i$  and male  $j$  can be expressed as:

$$m_{ij} = \alpha_{ij} \psi(H_{A,g}^{(i),f}, T_g^{(i),f}, H_{A,g}^{(j),m}, T_g^{(j),m}), \tag{3}$$

154 where  $\psi$  is a function that quantifies the dependence of the mating probability  $m_{ij}$  on the ancestry and  
 155 trait values of the individuals in a pair, and  $\alpha_{ij}$  is a normalization constant specific to the mating pair  $(i, j)$ .  
 156 The constant  $\alpha_{ij}$  is included in order to permit the matrix entries to satisfy the constant row and column  
 157 sum constraints. Without loss of generality, we choose the constant row and column sums to be 1 so that  
 158 the mating probability matrix  $M$  is doubly stochastic. Procedures to evaluate the  $\alpha_{ij}$  appear later in the  
 159 section. The properties of  $\psi$  are determined by a mating model.

160 In random mating, the mating probability is independent of individual ancestry and trait values, so that

Jaehee Kim

161  $\psi(H_{A,g}^{(i),f}, T_g^{(i),f}, H_{A,g}^{(j),m}, T_g^{(j),m})$  is constant across all  $i$  and  $j$ , and  $m_{ij}$  has the same value for all mating pairs.  
 162 Therefore, for all pairs  $(i, j)$ , each taken from  $\{1, 2, \dots, N\}$ ,  $m_{ij} = \alpha_{ij} = \alpha$  for some constant  $\alpha \in [0, 1]$ .

163 In assortative mating by ancestry, the mating probability depends only on the ancestries of poten-  
 164 tial mates and not on the phenotypes:  $\psi(H_{A,g}^{(i),f}, T_g^{(i),f}, H_{A,g}^{(j),m}, T_g^{(j),m}) = \psi(H_{A,g}^{(i),f}, H_{A,g}^{(j),m})$  and  $m_{ij} =$   
 165  $\alpha_{ij}\psi(H_{A,g}^{(i),f}, H_{A,g}^{(j),m})$ . For positive assortment, the mating function  $\psi$  has higher values if two individuals  
 166 have similar ancestries and lower values as the ancestries become more different. For example, in complete  
 167 assortment,  $\psi$  is 1 if the two input parameters have the same value and 0 if the values differ. For negative  
 168 assortment, the behavior of  $\psi$  is reversed compared with the positive assortment case:  $\psi$  increases as the  
 169 difference between the ancestries of the two individuals in a mating pair increases.

170 In assortative mating by phenotype, the mating probability depends only on the trait values of po-  
 171 tential mates and not on the ancestries:  $\psi(H_{A,g}^{(i),f}, T_g^{(i),f}, H_{A,g}^{(j),m}, T_g^{(j),m}) = \psi(T_g^{(i),f}, T_g^{(j),m})$  and  $m_{ij} =$   
 172  $\alpha_{ij}\psi(T_g^{(i),f}, T_g^{(j),m})$ . The qualitative requirements for the function  $\psi$  are the same as with assortative mating  
 173 by ancestry, but with the trait values of the mating pair as arguments instead of the ancestries.

174 We adopt the following form for the mating function:

$$\psi(X_g^{(i),f}, X_g^{(j),m}) = e^{-\frac{c|X_g^{(i),f} - X_g^{(j),m}|}{\sigma_{X_g}}} \quad (4)$$

175 The finite constant  $c$  quantifies the strength of the assortative mating. For a given pair of values  $(X_g^{(i),f}, X_g^{(j),m})$ ,  
 176 where  $X_g = H_{A,g}$  or  $X_g = T_g$ , a larger  $c$  value results in a lower mating probability, which gives stronger  
 177 positive assortative mating compared with a smaller  $c$  value. For a positive value of  $c$ , the function takes a  
 178 value of 1 if two potential mates have the same ancestry level (or phenotype), and  $\psi$  decreases exponentially  
 179 as the difference between the two individuals increases. A negative value of  $c$  indicates negative assortative  
 180 mating, where two individuals with different ancestry (or phenotype) have a higher probability of mating  
 181 than do two individuals with similar ancestry. We focus on positive assortative mating.

182 At each generation  $g$ , the admixture fraction  $H_{A,g}$  takes values in  $\{0, 1/2^g, 2/2^g, \dots, (2^g - 1)/2^g, 1\}$   
 183 (Section 2.1), and the phenotype  $T_g$  takes values in  $\{0, 1, \dots, 2k\}$  (Section 2.2). To compare statistics from  
 184 different mating schemes, we consider variables that are standardized by dividing  $X_g$  ( $H_{A,g}$  or  $T_g$ ) by its  
 185 standard deviation  $\sigma_{X_g}$  based on its distribution in  $H_g^{\text{par}}$  at each generation  $g$ . For the unstandardized  
 186 variables, because  $T_g$  takes a higher value than  $H_{A,g}$ , the effect of assortative mating by phenotype at the  
 187 same assortative mating strength  $c$  is artificially inflated compared to the effect of assortative mating by  
 188 admixture fraction.

189 Having specified the mating function  $\psi$ , we now formally state the normalization condition for the mating  
 190 matrix  $M$ : the sum across potential mates of the mating probabilities of a random individual in the parental  
 191 pool must be 1. Recalling that each entry  $m_{ij}$  in  $M$  represents the probability that individuals  $i$  and  $j$  mate,  
 192 the condition requires the row and column sums in the mating matrix to be 1 for each row and column. We  
 193 start with an unnormalized mating matrix  $\tilde{M} = [\tilde{m}_{ij}]$  whose entries are:

$$\tilde{m}_{ij} = \begin{cases} 1 & \text{random mating} \\ \psi(H_{A,g}^{(i),f}, H_{A,g}^{(j),m}) & \text{assortative mating by ancestry} \\ \psi(T_g^{(i),f}, T_g^{(j),m}) & \text{assortative mating by phenotype.} \end{cases} \quad (5)$$

194 We must obtain  $N^2$  normalizing constants  $\alpha_{ij}$  such that the mating matrix  $M = [m_{ij}] = [\alpha_{ij}\tilde{m}_{ij}]$  satisfies  
 195 the double stochasticity requirement. This requirement gives  $2N$  constraints, one for each row and one for  
 196 each column. Because each entry in the mating matrix represents the probability for two individuals to  
 197 mate, we also require  $m_{ij}$  to be in  $[0, 1]$  for all  $(i, j) \in \{1, 2, \dots, N\}^2$ .

198 Infinitely many matrices satisfy the constraints, as the set of  $2N$  equations with  $N^2$  variables is underde-  
 199 termined. We choose the matrix  $M$  by identifying the matrix that satisfies the set of constraints and that is  
 200 closest to our model matrix  $\tilde{M}$  according to the principle of minimum discrimination information (pp. 36-43  
 201 in [29]). Here, the ‘‘closeness’’ of a pair of matrices is measured by the Kullback-Leibler divergence  $D_{KL}$   
 202 (pp. 1-11 in [29]), which is nonnegative and is equal to zero if and only if the two matrices are identical.

203 The problem of identifying  $M$  can be formally written as a convex optimization problem. The objective  
 204 function that we seek to minimize is

$$\min_{\{m_{ij}\}} D_{KL}(M||\tilde{M}) = \min_{\{m_{ij}\}} \sum_{i=1}^N \sum_{j=1}^N m_{ij} \log \frac{m_{ij}}{\tilde{m}_{ij}}, \quad (6)$$

and we have constraints

$$\begin{aligned} \sum_{j=1}^N m_{ij} &= 1 \quad \text{for each } i \text{ from } 1 \text{ to } N, \\ \sum_{i=1}^N m_{ij} &= 1 \quad \text{for each } j \text{ from } 1 \text{ to } N, \\ 0 \leq m_{ij} &\leq 1 \text{ for all } (i, j) \in \{1, 2, \dots, N\}^2. \end{aligned} \quad (7)$$

We use the interior-point method [30,31], which iteratively traverses within the feasible region to obtain the optimal solution numerically, as implemented in `mosek` function of R package `Rmosek` [32]. For fixed  $\tilde{M}$ , the Hessian of the KL divergence has

$$\frac{\partial^2 D_{KL}(M||\tilde{M})}{\partial m_{ij} \partial m_{kl}} = \frac{1}{m_{ij}} \delta_{ik} \delta_{jl},$$

205 where  $\delta$  is the Kronecker delta. Because  $\nabla^2 D_{KL} > 0$  for all  $m_{ij} \in (0, 1)$ , the KL divergence function is  
 206 strictly convex (Section 3.1.4 in [33]) in each of the  $N^2$  variables in  $M$  for fixed  $\tilde{M}$ , and thus, the optimal  
 207 solution found by numerical minimization is the unique global minimum (Section 4.2.1 in [33]).

## 208 2.4 Expectation and Variance of the Admixture Fraction

209 To interpret our simulations of admixture dynamics, we will need a series of results concerning the mean and  
 210 variance of the admixture fraction in the admixed population. In particular, we derive a relationship between  
 211 the variance of the admixture fraction and the correlation in admixture levels for members of mating pairs.

212 Let  $H_{A,g}$  be a random variable representing the admixture fraction of an individual chosen at random  
 213 in the admixed population  $H_g$  at generation  $g \geq 1$ . We denote by  $(H_{A,g}^{f,p}, H_{A,g}^{m,p})$  the admixture fractions of  
 214 the members of a mating pair chosen at random from a parental pool  $H_g^{\text{par}}$  in generation  $g \geq 0$ . Here, the  
 215 superscript  $p$  denotes that the individual is from the parental pool. The parental pool  $H_g^{\text{par}}$ , from which  
 216 the admixed population  $H_{g+1}$  at generation  $g + 1$  is formed, consists of populations  $S_1$ ,  $S_2$ , and  $H_g$ , with



fractional contributions  $s_{1,g}$ ,  $s_{2,g}$ , and  $h_g$  respectively (Section 2.1 and Figure 2). Because we assume each population ( $S_1$ ,  $S_2$ , and  $H_g$ ) has equally many males and females, the numbers of males and females each remain constant at  $N$ , every individual has the same expected number of offspring, and no sex bias by population of origin exists in parental pairings (Section 2.3),  $H_{A,g}^{f,p}$  and  $H_{A,g}^{m,p}$  are identically distributed.

Let the random variable  $Y$  indicate the population membership of a random individual in  $H_g^{\text{par}}$ . Then

$$H_{A,g}^{f,p}, H_{A,g}^{m,p} = \begin{cases} H_{A,g} & \text{with } P(Y = H_g) = h_g \\ 1 & \text{with } P(Y = S_1) = s_{1,g} \\ 0 & \text{with } P(Y = S_2) = s_{2,g}. \end{cases} \quad (8)$$

For the expectation of admixture in the parental pool, we have

$$\begin{aligned} E[H_{A,g}^{f,p}] &= E[H_{A,g}^{m,p}] = E_Y \left[ E[H_{A,g}^{f,p} | Y] \right] = \sum_{y \in \{S_1, S_2, H_g\}} P(Y = y) E[H_{A,g}^{f,p} | Y = y] \\ &= s_{1,g} + h_g E[H_{A,g}] = s_{1,g} + h_g \mu_g. \end{aligned} \quad (9)$$

As a consequence of Eq. 9, we also have

$$E[(H_{A,g}^{f,p})^2] = E[(H_{A,g}^{m,p})^2] = s_{1,g} + h_g E[H_{A,g}^2] = s_{1,g} + h_g \mu_g^2 + h_g \text{Var}[H_{A,g}] \quad (10)$$

$$\text{Var}[H_{A,g}^{f,p}] = \text{Var}[H_{A,g}^{m,p}] = s_{1,g} + h_g \mu_g^2 + h_g \text{Var}[H_{A,g}] - (s_{1,g} + h_g \mu_g)^2. \quad (11)$$

Here,  $\mu_g = E[H_{A,g}]$  indicates the expectation of the admixture fraction of a random individual in the admixed population  $H_g$  at generation  $g \geq 1$ .

The ancestry of an offspring individual is deterministically set to the mean of the admixture fractions of the parents. This choice gives:

$$E[H_{A,g+1}] = E \left[ \frac{1}{2} (H_{A,g}^{f,p} + H_{A,g}^{m,p}) \right] = s_{1,g} + h_g E[H_{A,g}] = s_{1,g} + h_g \mu_g. \quad (12)$$

We obtain the recursion for the variance of the admixture fraction over a single generation as follows:

$$\begin{aligned} \text{Var}[H_{A,g+1}] &= E[H_{A,g+1}^2] - (E[H_{A,g+1}])^2 \\ &= \frac{1}{4} E[(H_{A,g}^{f,p} + H_{A,g}^{m,p})(H_{A,g}^{f,p} + H_{A,g}^{m,p})] - (E[H_{A,g+1}])^2 \\ &= \frac{1}{2} \left( E[(H_{A,g}^{f,p})^2] + E[H_{A,g}^{f,p} H_{A,g}^{m,p}] \right) - (E[H_{A,g+1}])^2 \\ &= \frac{1}{2} \left( E[(H_{A,g}^{f,p})^2] + \text{Cor}[H_{A,g}^{f,p}, H_{A,g}^{m,p}] \text{Var}[H_{A,g}^{f,p}] + (E[H_{A,g}^{f,p}])^2 \right) - (E[H_{A,g+1}])^2 \\ &= \frac{1}{2} (1 + r_{H_{A,g}}) \left[ h_g \text{Var}[H_{A,g}] + \mu_g^2 h_g (1 - h_g) - 2\mu_g h_g s_{1,g} + s_{1,g} (1 - s_{1,g}) \right], \end{aligned} \quad (13)$$

where  $r_{H_{A,g}} = \text{Cor}[H_{A,g}^{f,p}, H_{A,g}^{m,p}]$  denotes the correlation of the admixture fractions in a mating pair. The last step is obtained from Eqs. 9–12. As we will see, the time-varying  $r_{H_{A,g}}$  value in general depends on the parameters of the population model, the quantitative trait model, and the mating model.

For a special case of a single admixture event in which source populations  $S_1$  and  $S_2$  do not contribute to the admixed population after its founding ( $s_{1,g} = s_{2,g} = 0$  and  $h_g = 1$  for all  $g \geq 1$ ), the expectation of



231 the admixture fraction stays constant in time (Eq. 12), and the variance reduces to a simple formula:

$$\text{Var}[H_{A,g+1}] = \frac{1}{2}(1 + r_{H_{A,g}})\text{Var}[H_{A,g}]. \quad (14)$$

232 Under random mating in an infinite population with no ongoing contributions from the source populations,  
233 with  $r_{H_{A,g}} = 0$  for all  $g \geq 0$ , Eq. 14 reduces to the formula  $\text{Var}[H_{A,g}] = s_{1,0}(1 - s_{1,0})/2^g$  of Verdu and  
234 Rosenberg [1]. Eq. 14 has also been derived by Zaitlen et al. [27] but under assumptions of constant mating  
235 correlation ( $r_{H_{A,g}} = r$ ) across all generations, no migration, and infinite population size.

## 236 3 Simulation

### 237 3.1 Simulation Procedure

238 Having specified the populations of interest, the properties of trait values in the populations, and the mating  
239 probabilities for pairs of individuals, we now describe how we simulate populations under the model. At the  
240 first time step ( $g = 0$ ),  $s_{1,0}N$  and  $s_{2,0}N$  males are randomly generated from the source populations  $S_1$  and  
241  $S_2$ , respectively, with  $s_{1,0} + s_{2,0} = 1$ . The corresponding numbers of females  $s_{1,0}N$  and  $s_{2,0}N$  are randomly  
242 drawn from source populations,  $S_1$  and  $S_2$ , respectively, constituting the founding parental pool  $H_0^{\text{par}}$  of  
243  $2N$  individuals, with  $N$  males and  $N$  females. All individuals in source population  $S_1$  have an admixture  
244 fraction value of 1, and all individuals in source population  $S_2$  have an admixture fraction value of 0, by  
245 definition. For each individual in populations  $S_1$  and  $S_2$ , genotypes at each of  $k$  quantitative trait loci are  
246 then randomly generated on the basis of pre-specified allele frequencies  $p_i$  and  $q_i$ .

247 We consider two different distributions for the  $p_i$  and  $q_i$ . First, we assume that the two source populations  
248 display fixed differences at all trait loci, so  $p_i = 1$  and  $q_i = 0$  for all  $k$  loci. In this case, every individual in  
249 population  $S_1$  has the “1” allele at all trait loci, and every individual in population  $S_2$  has the “0” allele at  
250 all trait loci. In subsequent generations, allele “1” can be traced back to population  $S_1$ , and allele “0” to  $S_2$   
251 (Figure 1). This choice for the  $p_i$  and  $q_i$  models a case in which trait-influencing alleles are initially entirely  
252 predictive of ancestry and vice versa.

253 Second, departing from the idealized model, we simulate sets of  $k$  allele frequency pairs  $(p_i, q_i)$ ,  $i \in$   
254  $\{1, \dots, k\}$  following Edge and Rosenberg [22]. Allele frequencies  $\pi_i$  for derived alleles in the “ancestral”  
255 population of  $S_1$  and  $S_2$  are drawn based on the neutral site frequency spectrum:  $P[\pi_i = j/(2N_a)] \propto 1/j$ ,  
256 where  $N_a$  indicates the size of the ancestral population (Eq. B6.6.1 in [34]). We use  $2N_a = 20,000$ . We  
257 assume each locus  $i$  in  $S_1$  and  $S_2$  undergoes independent genetic drift following a split. We add random  
258 numbers  $\epsilon_{i,1}$  and  $\epsilon_{i,2}$  drawn from a  $\text{Normal}(0, \gamma\pi_i(1 - \pi_i))$  distribution to  $\pi_i$  to simulate derived allele  
259 frequencies at locus  $i$  in populations  $S_1$  and  $S_2$ , respectively. The parameter  $\gamma$  represents the amount of  
260 variance introduced by drift into the allele frequencies of the divergent populations. Following Edge and  
261 Rosenberg [22], we choose  $\gamma = 0.3$  so that the overall degree of genetic differentiation between  $S_1$  and  $S_2$   
262 at a group of simulated loci approximates worldwide human  $F_{ST}$  estimates. If  $\epsilon_{i,1} \geq \epsilon_{i,2}$ , then we assign  
263  $p_i = \pi_i + \epsilon_{i,1}$  and  $q_i = \pi_i + \epsilon_{i,2}$ . If  $\epsilon_{i,1} < \epsilon_{i,2}$ , then we assign  $p_i = 1 - (\pi_i + \epsilon_{i,1})$  and  $q_i = 1 - (\pi_i + \epsilon_{i,2})$ . Note  
264 that if this procedure produces  $p_i > 1$  or  $q_i < 0$ , then we assign  $p_i = 1$  and  $q_i = 0$  so that  $0 \leq q_i \leq p_i \leq 1$ .

265 Using the mating function in Eq. 4, we compute an unnormalized mating matrix for every pair containing

266 a male and a female from the parental pool:

$$\widetilde{M} = \begin{pmatrix} 1 & e^{-c\Delta_{1,2}} & \dots & e^{-c\Delta_{1,N-1}} & e^{-c\Delta_{1,N}} \\ e^{-c\Delta_{2,1}} & 1 & \dots & e^{-c\Delta_{2,N-1}} & e^{-c\Delta_{2,N}} \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ e^{-c\Delta_{N-1,1}} & e^{-c\Delta_{N-1,2}} & \dots & 1 & e^{-c\Delta_{N-1,N}} \\ e^{-c\Delta_{N,1}} & e^{-c\Delta_{N,2}} & \dots & e^{-c\Delta_{N,N-1}} & 1 \end{pmatrix}, \quad (15)$$

267 where  $\Delta_{i,j} = |H_{A,0}^{(i),f} - H_{A,0}^{(j),m}| / \sigma_{H_{A,0}}$  for assortative mating by ancestry and  $\Delta_{i,j} = |T_0^{(i),f} - T_0^{(j),m}| / \sigma_{T_0}$   
 268 for assortative mating by trait. For random mating,  $c = 0$  and all entries equal 1. The matrix  $\widetilde{M}$  is  
 269 normalized using the procedure of Section 2.3, producing the mating probability matrix  $M$ .

270 Considering all  $N^2$  potential mating pairs, we randomly draw  $N$  mating pairs with replacement from  
 271 the parental pool, weighting mate choices by the mating probabilities in  $M$ . Once mates are chosen, each  
 272 mating pair produces two children, one male and one female, in order to keep the population size of the  
 273 offspring generation constant at  $N$  males and  $N$  females. An admixture fraction for an offspring individual  
 274 is then assigned as the mean of its parental admixture fractions. Assuming no linkage disequilibrium and no  
 275 mutation, the genotype of the offspring at the quantitative trait loci is then determined by independently  
 276 selecting at each locus one random allele from one parent and one from the other. The resulting  $2N$  offspring  
 277 form the admixed population  $H_1$  at generation 1.

278 In subsequent generations  $g \geq 1$ , we randomly select  $s_{1,g}N$ ,  $s_{2,g}N$ , and  $h_gN$  males and  $s_{1,g}N$ ,  $s_{2,g}N$ ,  
 279 and  $h_gN$  females from  $S_1$ ,  $S_2$ , and  $H_g$ , respectively, forming a  $g$ th generation parental pool  $H_g^{\text{par}}$  of  $2N$   
 280 individuals, consisting of  $N$  males and  $N$  females. The procedure to generate the offspring population  $H_{g+1}$   
 281 from  $H_g^{\text{par}}$  is the same as the procedure for generating  $H_1$  from  $H_0^{\text{par}}$ .

282 Throughout the simulation, we keep the population size parameter  $N$  constant at 1,000 for computational  
 283 efficiency in the matrix normalization step. Here, the admixed population size ( $N$ ) is not necessarily identical  
 284 to the source population sizes ( $N_a$ ). For each set of parameters,  $(k, p_1, p_2, \dots, p_k, q_1, q_2, \dots, q_k, X_1, X_2, \dots, X_k,$   
 285  $c, s_{1,0}, s_{1,g}, s_{2,g})$ , we propagated the population to  $G = 40$  generations. We generated 100 independent tra-  
 286 jectories for each parameter set. For each trajectory, we computed statistics of interest, averaging them over  
 287 all 100 trajectories in the simulation given the fixed set of parameters.

### 288 3.2 Base Case

289 We start with an idealized base case. First, we specify the parameters involving the population model  
 290 (Section 2.1). We assume an equal influx from each source population at founding  $g = 0$ :  $s_{1,0} = 0.5$ ,  
 291  $s_{2,0} = 1 - s_{1,0} = 0.5$ . We also assume no additional contributions from the source populations in the  
 292 subsequent generations,  $s_{1,g} = s_{2,g} = 0$ , and  $h_g = 1 - s_{1,g} - s_{2,g} = 1$  for all  $g \geq 1$ .

293 Next, we choose parameter values for the quantitative trait model (Section 2.2). We consider  $k = 10$   
 294 trait loci. Across the  $k$  loci, all “1” alleles come from source population  $S_1$  and all “0” alleles come from  $S_2$ :  
 295  $p_i = 1$  and  $q_i = 0$  for all  $i = 1, 2, \dots, k$ . For each locus  $i$  contributing to the quantitative trait, we define “1”  
 296 to be the “+” allele and “0” to be the “-” allele:  $X_i = 1$  for all  $i = 1, 2, \dots, k$ .

297 Finally, for the mating model (Section 2.3), we set the assortative mating strength  $c$  in Eq. 4 to 0.5.

### 3.3 Statistics Measured

In each simulated admixed population, in each generation  $g$ , we computed the following statistics: correlation between admixture fraction and trait ( $\text{Cor}[H_A, T]$ ), variance of the admixture fraction in the population ( $\text{Var}[H_A]$ ), and variance of the trait value in the population ( $\text{Var}[T]$ ). In the following section, we discuss how these statistics of interest change as we modify the simulation parameters.

## 4 Results

### 4.1 Base Case

#### 4.1.1 Correlation between genetic ancestry and phenotype ( $\text{Cor}[H_A, T]$ )

In the base case, each individual from  $S_1$  has admixture fraction  $H_A = 1$  and trait value  $T = k$ , and each individual from  $S_2$  has admixture fraction  $H_A = 0$  and trait value  $T = 0$ . Therefore, in the founding parental pool,  $H_0^{\text{par}}$ , admixture fraction and trait value are perfectly correlated:  $\text{Cor}[H_A, T] = 1$ . In subsequent generations, however, the correlation between the admixture fraction and the trait values starts to decouple, as illustrated in Figure 1. With all parameters involving the population model and the quantitative trait model fixed, the rate of decay in  $\text{Cor}[H_A, T]$  depends on the mating model.

A comparison of  $\text{Cor}[H_A, T]$  under the three mating models using base-case parameters appears in Figure 4E. Irrespective of the mating model, the founding parental pool has a perfect correlation between ancestry and phenotype. Even if the population starts with perfect correlation between admixture fractions and trait values, however, then random mating rapidly decouples them (red curve). It takes 6 generations of random mating for the correlation to decrease below 0.5 ( $\text{Cor}[H_A, T] = 0.490$ ). After  $g = 20$  generations, the correlation becomes 0.137, and it is near zero at  $g = 40$  ( $-0.003$ ).

Compared to random mating, positive assortative mating slows the decoupling of admixture fractions and trait values. Assortative mating by phenotype (green curve in Figure 4E) maintains the correlation longer than assortative mating by admixture fraction (blue curve in Figure 4E). It takes 11 generations under assortative mating by phenotype for the correlation to drop below  $\frac{1}{2}$  ( $\text{Cor}[H_A, T] = 0.490$ ), and 10 generations under assortative mating by admixture ( $\text{Cor}[H_A, T] = 0.443$ ). Across the 40 generations we simulated,  $\text{Cor}[H_A, T]$  is consistently higher under assortative mating by phenotype than under assortative mating by admixture fraction. The correlation decreases to 0.227 at  $g = 20$  and 0.043 at  $g = 40$  under assortative mating by phenotype. The corresponding values under assortative mating by admixture are 0.065 at  $g = 20$  and 0.009 at  $g = 40$ , both considerably lower than under assortative mating by phenotype.

The speed of the decoupling between admixture fraction and trait value increases with two factors: Mendelian noise in generating admixed individuals from the admixed population and decreasing contributions from the source populations. Compared to random mating, both assortative mating models have higher probabilities for matings within source populations, and thus, the proportion of individuals produced in the admixed population at  $g = 1$  that are genetically admixed is smaller (blue and green lines in the marginal plots for  $H_A$  in Figure S1A). Over time, as displayed in Figure S1, random mating pulls individuals away from the source populations, pushing the  $H_A$  and  $T$  distributions towards the mean values rapidly. On the other hand, both assortative mating models maintain individuals with  $H_A$  and  $T$  values near the source population values for longer, and thus, they retain higher  $\text{Cor}[H_A, T]$  than random mating.

The difference in  $\text{Cor}[H_A, T]$  between the two assortative mating models arises from the difference between  $\text{Var}[H_A]$  and  $\text{Var}[T]$ .  $\text{Cov}[H_A, T]$  is similar under the two models. Given the similar covariance,

$$\frac{\text{Cor}[H_A, T]_{\text{phen}}}{\text{Cor}[H_A, T]_{\text{gen}}} \approx \sqrt{\frac{\text{Var}[H_A]_{\text{gen}}}{\text{Var}[H_A]_{\text{phen}}} \times \frac{\text{Var}[T]_{\text{gen}}}{\text{Var}[T]_{\text{phen}}}},$$

where the subscripts “gen” and “phen” indicate the property on which mating pairs assort. As we show in the next section, both assortative mating models increase  $\text{Var}[H_A]$  and  $\text{Var}[T]$  compared to random mating, and the increase in variance is the largest for the property on which mating assorts:  $\text{Var}[H_A]_{\text{gen}} > \text{Var}[H_A]_{\text{phen}}$  and  $\text{Var}[T]_{\text{phen}} > \text{Var}[T]_{\text{gen}}$ . However, we will see that the increase in  $\text{Var}[H_A]$  due to assortative mating by admixture fraction exceeds the increase in  $\text{Var}[T]$  due to assortative mating by trait:

$$\frac{\text{Var}[H_A]_{\text{gen}}}{\text{Var}[H_A]_{\text{phen}}} > \frac{\text{Var}[T]_{\text{phen}}}{\text{Var}[T]_{\text{gen}}}.$$

336 This result leads to higher  $\text{Cor}[H_A, T]$  under assortative mating by trait compared to that under assortative  
337 mating by admixture fraction.

#### 338 4.1.2 Variance of Ancestry and Phenotype ( $\text{Var}[H_A]$ and $\text{Var}[T]$ )

339 Each individual in  $S_1$  has admixture fraction 1, and each individual in  $S_2$  has admixture fraction 0. In  
340 the founding parental pool,  $\text{Var}[H_A] = 0.250$  for all three mating models. As discussed in Section 2.4, the  
341 variance of the admixture fraction can be understood in relation to the correlation coefficient  $\text{Cor}[H_{A,g}^f, H_{A,g}^m]$   
342 of the admixture fractions of members of mating pairs. Figure S2 shows this correlation coefficient for the  
343 simulations of Figure 4, and Figure S3 shows the analogous correlation  $\text{Cor}[T_g^f, T_g^m]$  of trait values.

344 Figure 5E then shows the variance of the admixture fraction over time under the three mating models,  
345 for the same simulations from Figure 4E with the base case parameters. The  $\text{Var}[H_A]$  curves in Figure 5E  
346 under the three mating models follow Eq. 13, using the time-varying  $r_{H_A,g}$  in Figure S2.

347 Among the three mating models,  $\text{Var}[H_A]$  decreases fastest for random mating. After one generation,  
348  $\text{Var}[H_A]$  falls in half (0.125), and it continues to decrease monotonically by half. After 40 generations, the  
349 value decreases to  $2.118 \times 10^{-13}$ . The distribution of the admixture fraction concentrates around  $H_A = \frac{1}{2}$   
350 at each generation. Because the offspring admixture fraction is the mean of those of its parents, without  
351 additional influx from the source populations after the founding event, random mating rapidly drives the  
352 admixture fraction away from extreme values (0 or 1) toward the mean value of the parental pool ( $\frac{1}{2}$ ).

353 Under assortative mating by admixture, pairs with similar admixture fraction have higher mating proba-  
354 bilities than under random mating. The fraction of offspring that are admixed is smaller than under random  
355 mating, and the admixture fraction distribution remains close to the extreme values (0 or 1) for longer  
356 (Figure S1). Hence,  $\text{Var}[H_A]$  is larger under assortative mating by admixture fraction (Figure 5E). Without  
357 influx from the source populations,  $\text{Var}[H_A]$  eventually decreases to zero, but the decrease is slower than  
358 for random mating.  $\text{Var}[H_A] = 0.184$  after one generation of assortative mating by admixture fraction, and  
359  $\text{Var}[H_A] = 1.118 \times 10^{-8}$  after 40 generations. This result can also be seen in Eq. 13. From generation  $g$   
360 to  $g + 1$ ,  $\text{Var}[H_A]_g$  decreases by a factor of  $(1 + r_{H_A,g})/2$ . With positive assortative mating by admixture  
361 ( $r_{H_A,g} > 0$ ),  $\text{Var}[H_A]$  in the next generation is increased compared to the case of random mating ( $r_{H_A,g} = 0$ ).

362 Under assortative mating by phenotype,  $\text{Var}[H_A] = 0.183$  after one generation of assortative mating by

363 phenotype, and  $\text{Var}[H_A] = 4.910 \times 10^{-12}$  after 40 generations. For the first few generations ( $g < 5$ ), because  
364  $\text{Cor}[H_A, T]$  is high, the correlation between the admixture fraction of mating pairs, and thus  $\text{Var}[H_A]$ , is  
365 similar under the two assortative mating models, as shown in the comparison of the green and blue curves  
366 in Figures S2 and 5. However, because the admixture fraction and phenotype decouple by introduction of  
367 the Mendelian noise, mating assortatively by phenotype results in lower  $r_{H_A, g}$  than mating assortatively  
368 by admixture fraction. In accord with Eq. 13, assortative mating by phenotype produces faster decay in  
369  $\text{Var}[H_A]$  with its lower  $r_{H_A, g}$  at each generation than assortative mating by admixture fraction.

370 For the variance of the phenotype, using Eq. 1, all individuals in  $S_1$  and  $S_2$  have trait values of 20  
371 and 0, respectively. Therefore, in the founding parental pool,  $H_0^{\text{par}}$ , noting that  $S_1$  and  $S_2$  each have 1,000  
372 individuals, this variance has the same constant value of  $\frac{2,000}{1,999} \cdot 10^2 = 100.050$  irrespective of the mating model.  
373 Figure 6E displays the variance of the phenotype, which decreases most rapidly under random mating. After  
374 one generation of random mating,  $\text{Var}[T]$  decreases by half (50.025), and it approaches its steady-state value  
375 of  $\approx 4.957$  after 13 generations. Opposite to what was seen for  $\text{Var}[H_A]$ , however, assortative mating by  
376 trait retains  $\text{Var}[T]$  higher for longer than assortative mating by admixture fraction. Similar to the case with  
377  $\text{Var}[H_A]$  under assortative mating by admixture fraction, assortative mating by phenotype keeps the trait  
378 values close to extreme values for longer than the other two mating models.

379 Having examined the behavior of  $\text{Cor}[H_A, T]$ ,  $\text{Var}[H_A]$ , and  $\text{Var}[T]$  in the base case, we now explore the  
380 effect of the assortative mating strength  $c$  and the parameters involving the quantitative trait—the number  
381 of loci  $k$ , allele frequencies  $p_i$  and  $q_i$ , and trait contribution at each locus  $X_i$ —on these quantities.

## 382 4.2 Assortative Mating Strength ( $c$ )

### 383 4.2.1 $\text{Cor}[H_A, T]$

384 Each row of Figure 4 illustrates the influence of the assortative mating strength  $c$  on  $\text{Cor}[H_A, T]$  with a fixed  
385 number of trait loci  $k$ , and each column depicts the effect of the number of loci  $k$  on  $\text{Cor}[H_A, T]$  with fixed  
386 assortative mating strength  $c$ . All parameters other than  $c$  and  $k$  are held constant at the base case values.

387 With different assortative mating strengths and numbers of trait loci,  $c = 0.1, 0.5, 1.0$  and  $k = 1, 10, 100$ ,  
388 the qualitative behavior of  $\text{Cor}[H_A, T]$  over time remains the same as in the base case. As before, we observe  
389 decay in  $\text{Cor}[H_A, T]$  under all three mating models, with random mating decoupling ancestry and trait values  
390 the most rapidly.  $\text{Cor}[H_A, T]$  remains higher for longer under assortative mating by phenotype than under  
391 assortative mating by admixture fraction. The rate of decay and the degree to which the patterns differ  
392 across the three mating models depend on the assortative mating strength and the number of loci.

393 If assortative mating is weak ( $c = 0.1$  in Figure 4A, D, G), then the effect of assortative mating is small,  
394 and thus,  $\text{Cor}[H_A, T]$  under assortative mating by admixture and by phenotype closely follows that under  
395 random mating. This pattern is seen irrespective of the number of loci. Note that in the limit of  $c = 0$ , the  
396 assortative mating and random mating models are identical because the mating function in Eq. 4 becomes  
397 a constant, the same value for all three mating models.

398 Comparing panels within rows of Figure 4, the results from random mating are identical, as the assortative  
399 mating strength does not affect the random mating model. Under both assortative mating models, however,  
400  $\text{Cor}[H_A, T]$  increases with the assortative mating strength. In an extreme case of complete assortment  
401 ( $c \rightarrow \infty$ ), the correlation would stay constant at 1 across all generations: we start with complete correlation

402 between admixture and phenotype, and  $c \rightarrow \infty$  implies that only identical individuals can mate, so that the  
403 correlation persists unchanged.

404 The difference among the three models increases with the assortative mating strength given a fixed  
405 number of trait loci. The difference is the greatest if  $k = 1$  and  $c = 1.0$  (Figure 4C). Even after 40  
406 generations, assortative mating by trait retains a high correlation at 0.788, whereas the corresponding values  
407 under random mating and assortative mating by admixture are 0.006 and 0.010, respectively.

#### 408 **4.2.2 Var[ $H_A$ ] and Var[ $T$ ]**

409 The plots of Var[ $H_A$ ] in Figure 5 and Var[ $T$ ] in Figure 6 consider the same simulations that appear for  
410 Cor[ $H_A, T$ ] in Figure 4. As is seen in classical work [26, 35, 36], compared to random mating, assortative  
411 mating increases the variance of the property on which assortment takes place. Thus, the variance of the  
412 admixture fraction is increased to a greater extent under mating by admixture fraction than under mating  
413 by phenotype. Similarly, the variance of the phenotype is increased to a greater extent under mating by  
414 phenotype than under mating by admixture fraction. Both types of assortative mating increase both Var[ $H_A$ ]  
415 and Var[ $T$ ] compared with random mating.

416 The variance-increasing effect of the assortative mating is visible when comparing panels within each  
417 row. For low assortative mating strength ( $c = 0.1$ ), panels A, D, and G in Figures 5 and 6 show that  
418 minimal differences in Var[ $H_A$ ] and Var[ $T$ ] exist between mating models. As  $c$  increases, for a given number  
419 of loci, Figures 5 and 6 display increased differences between random and assortative mating, with maximal  
420 separation at the largest assortative mating strength simulated,  $c = 1$  (panels C, F, I). The random mating  
421 model is unaffected by the assortative mating strength  $c$ , as was seen with Cor[ $H_A, T$ ] in Section 4.2.1.

### 422 **4.3 Number of Trait Loci ( $k$ )**

#### 423 **4.3.1 Cor[ $H_A, T$ ]**

424 A comparison of panels within columns of Figure 4 shows that under random mating, with more loci as-  
425 sociated with the phenotype, the ancestry-phenotype correlation is higher and stays high for longer. In  
426 other words, it takes longer for  $H_A$  and  $T$  to become decoupled. In particular, under random mating, the  
427 correlation between trait and ancestry falls below 0.5 at  $g = 3$  if  $k = 1$ ,  $g = 6$  if  $k = 10$ , and  $g = 10$  if  
428  $k = 100$ , independent of the assortative mating strength.

429 As the number of loci increases, results from the models with assortative mating by phenotype and by  
430 admixture become similar. If  $(c, k) = (1, 100)$  (Figure 4I), then it takes 24 generations for Cor[ $H_A, T$ ] values  
431 under the two models to differ by more than 0.1. Corresponding times for  $(c, k) = (1, 1)$  (Figure 4C) and  
432  $(c, k) = (1, 10)$  (Figure 4F) are  $g = 6$  and  $g = 15$ , respectively. Recall that the admixture fraction represents  
433 the probability that a random allele at a random autosomal genetic locus originates from source population  
434  $S_1$ , assuming infinitely many loci. In the  $k \rightarrow \infty$  limit, with the whole genome contributing to the trait, the  
435 assortative mating models by admixture and by phenotype would behave in exactly the same way.

#### 436 **4.3.2 Var[ $H_A$ ] and Var[ $T$ ]**

437 Comparing panels within columns in Figure 5, for a given assortative mating strength, Var[ $H_A$ ] under  
438 assortative mating by admixture follows the same curve irrespective of the number of loci. Because the



439 mating probability is independent of trait values if mating assortatively by admixture,  $k$  has no effect.

440 As in the base case (Section 4.1.2), both assortative mating models have higher  $\text{Var}[H_A]$  and  $\text{Var}[T]$  than  
441 random mating. Of the two assortative mating models, assortative mating by admixture fraction has greater  
442  $\text{Var}[H_A]$  than assortative mating by trait at each generation. For  $\text{Var}[T]$ , assortative mating by trait has  
443 greater values than assortative mating by admixture fraction. As was seen with  $\text{Cor}[H_A, T]$  (Section 4.3.1),  
444 for  $\text{Var}[H_A]$  and  $\text{Var}[T]$ , the difference between random mating and both assortative mating models increases  
445 with  $k$ , and the difference between the two assortative mating models diminishes as  $k$  increases.

#### 446 4.4 Allele Frequencies ( $p_i$ and $q_i$ )

447 Departing further from the base case, we next evaluate the effect of the allele frequencies,  $p_i$  and  $q_i$ , on  
448 the quantities of interest. Instead of treating the two source populations as fixed for different alleles, the  
449 frequencies  $p_i$  and  $q_i$  are now sampled according to the simulation procedure described in Section 3.1.  
450 Because our results show a monotonic trend across the number of loci we examined (Section 4.3), we focus  
451 this analysis on a single value of  $k = 10$ , the number of loci corresponding to the base case.

##### 452 4.4.1 $\text{Cor}[H_A, T]$

453 Figure 7A displays  $\text{Cor}[H_A, T]$  under the model with simulated rather than fixed allele frequencies.  $\text{Cor}[H_A, T]$   
454 starts from a lower correlation value at time  $g = 0$ , 0.456, compared to the base case (Figure 4E) value of  
455 1. If all loci have  $X_i = 1$ , as shown in Eq. 1, then an individual's trait value is determined by the number  
456 of "1" alleles across the trait loci. Because the allele "1" is randomly drawn at each locus  $i = 1, 2, \dots, k$   
457 with probabilities  $P(L_{ij} = 1 \mid M = S_1) = p_i$  and  $P(L_{ij} = 1 \mid M = S_2) = q_i$  with  $j = 1, 2$  (Section 2.2)  
458 and the mean absolute difference between simulated  $p_i$  and  $q_i$  across  $k$  loci is small, some individuals in the  
459 source population  $S_1$  have lower trait values than some individuals in  $S_2$ , and vice versa. However, due to  
460 the constraint  $p_i \geq q_i$  across all trait loci, individuals from  $S_1$  have higher probability of having a larger trait  
461 value than those from  $S_2$ . This property accounts for the nonzero correlation between ancestry and trait  
462 present in the source populations outside the base case setting.

463 The qualitative differences between the three mating models remain similar to the base case, as shown  
464 in Figure 7A. All three mating models, however, show an increased rate of decoupling between the admix-  
465 ture fraction and the trait, in that the correlation decreases more rapidly. For random mating, it takes  
466 only 4 generations for  $\text{Cor}[H_A, T]$  to drop to below half of its starting value, reaching 0.170. The corre-  
467 sponding values under assortative mating by admixture fraction and assortative mating by trait are  $g = 5$   
468 ( $\text{Cor}[H_A, T] = 0.240$ ) and  $g = 10$  ( $\text{Cor}[H_A, T] = 0.220$ ), respectively. Compared to the base case, the cor-  
469 relation between ancestry and trait in the source population is weaker if the allele frequencies are drawn  
470 from the simulation, and thus, the "1" allele does not necessarily trace back to the source population  $S_1$ .  
471 Under this setting, the effect of the Mendelian noise in decoupling of admixture fraction and phenotype in  
472 producing admixed individuals becomes more significant than the base case.

##### 473 4.4.2 $\text{Var}[H_A]$ and $\text{Var}[T]$

474 The admixture fraction values at the source populations are not affected by the allele frequencies:  $H_A = 1$   
475 and  $H_A = 0$  for all individuals in  $S_1$  and  $S_2$ , respectively. If  $s_{1,0} = s_{2,0} = 0.5$ , then  $\text{Var}[H_A]$  starts at 0.25 in  
476 the founding parental pool, irrespective of the allele frequencies. Comparing Figure 7B and 5E, the  $\text{Var}[H_A]$



477 curves under random mating (red) and assortative mating by admixture (blue) are not affected by the change  
478 in allele frequencies  $p_i$  and  $q_i$ , holding other parameters fixed. Under random mating and assortative mating  
479 by admixture, mate choice is independent of the parameters that affect the quantitative trait, and thus, the  
480 change in  $p_i$  and  $q_i$  does not alter the admixture fraction distribution at each generation.

481 By contrast, under assortative mating by phenotype,  $\text{Var}[H_A]$  (green) is affected by the change in the  
482 nature of the allele frequencies.  $\text{Var}[H_A]$  under assortative mating by phenotype closely follows that under  
483 random mating. The simulated allele frequencies have relatively small differences ( $\bar{\delta} \approx 0.0509$ ) between  
484 source populations  $S_1$  and  $S_2$ . With  $X_i = 1$  for all loci, the between-group difference in trait values is small  
485 as well, whereas all individuals in  $S_1$  and  $S_2$  still have  $H_A = 1$  and  $H_A = 0$ , respectively. Therefore, with  
486 the simulated allele frequencies, the effect on the admixture fraction of assortative mating by phenotype  
487 is similar to that in the random mating case. This scenario contrasts with the base case, where allele “1”  
488 can be associated with the source population  $S_1$  with certainty, and  $\text{Var}[H_A]$  under assortative mating by  
489 phenotype behaves similarly to the case of assortative mating by admixture fraction.

490 With the simulated allele frequencies,  $\text{Var}[T] = 0.877$  in the founding parental pool. At  $g = 1$ ,  $\text{Var}[T]$   
491 values under random mating and under assortative mating by admixture are 0.784 and 0.825, respectively.  
492 Assortative mating by admixture maintains higher  $\text{Var}[T]$  than random mating until  $g = 8$  and then follows  
493 the  $\text{Var}[T]$  curve for random mating. By contrast,  $\text{Var}[T]$  under assortative mating by trait gradually  
494 increases until  $g = 13$ , at which it achieves its maximum of 0.977, and then decreases to 0.935 at  $g = 40$ .

## 495 4.5 Trait Contributions of Individual Loci ( $X_i$ )

496 Returning to the case with fixed allele frequencies of 1 and 0 in the source populations, we next examine the  
497 case in which the trait has the property that both alleles have equal probability of being the “+” allele, as  
498 described in Section 2.2:  $P(X_i = 1) = P(X_i = 0) = \frac{1}{2}$  for all  $i = 1, 2, \dots, k$ . Figure 8 displays the results  
499 using the number of trait loci from the base case,  $k = 10$ . The qualitative behavior of the result does not  
500 depend on the number of loci with the other parameters fixed.

### 501 4.5.1 $\text{Cor}[H_A, T]$

502 If we let the number of loci with  $X_i = 1$  be  $z$ , then the number of loci with  $X_i = 0$  is  $k - z$ . Because  
503  $p_i = 1$  and  $q_i = 0$  across all loci in the base case, the trait value is  $2z$  for every individual in  $S_1$  and  
504  $2(k - z)$  for every individual in  $S_2$  (Eq. 1). For a randomly generated set of  $X_i$ ,  $\{i = 1, 2, \dots, k\}$ , under  
505  $P(X_i = 1) = P(X_i = 0) = \frac{1}{2}$ , if  $z \neq k - z$ , then  $\text{Cor}[H_A, T] = 1$  in the founding parental pool  $H_0^{\text{par}}$ , as  
506 shown in Figure 8A. However, compared with the base case (Figure 3.2), the correlation decays much more  
507 rapidly. With the  $P(X_i = 1) \neq \frac{1}{2}$  setting, the ancestry and trait are not as tightly coupled in the source  
508 populations. However, as in Section 4.1-4.4, assortative mating by phenotype preserves the correlation for  
509 the longest, and random mating decouples the correlation the fastest of the three mating models.

510 If the numbers of loci with  $X_i = 1$  and  $X_i = 0$  are equal ( $z = k - z$ ), then all individuals in the source  
511 populations have trait value  $k$  irrespective of their origin, and thus, no correlation exists between trait and  
512 ancestry in the source population. Hence,  $\text{Cor}[H_A, T]$  is 0 in the founding parental pool  $H_0^{\text{par}}$ , and the  
513 correlation remains at 0 throughout the time simulated, irrespective of the mating type (Figure 8D).

#### 514 4.5.2 $\text{Var}[H_A]$ and $\text{Var}[T]$

515 The panels of Figure 8 display results under  $P(X_i = 1) = P(X_i = 0) = \frac{1}{2}$ , fixing other parameters as in  
516 base case. By the same reasoning as in Section 4.4.2, the change in parameters involving the quantitative  
517 trait does not affect  $\text{Var}[H_A]$  under random mating and under assortative mating by admixture fraction. A  
518 comparison of Figure 8B and 8E with Figure 5E shows that  $\text{Var}[H_A]$  values under assortative mating by  
519 admixture fraction are not affected by the change in the  $X_i$ .

520 For  $z \neq k - z$ , results with  $z = 6$  and  $k - z = 4$  appear in Figure 8B. Some correlation between the  
521 admixture fraction and allele “1” exists in the source populations, and thus,  $\text{Var}[H_A]$  for assortative mating  
522 by trait (green) somewhat follows that for assortative mating by admixture (blue). However, compared to  
523 the base case (Figure 5E), where the “1” allele can be traced back to  $S_1$  with certainty, a more noticeable  
524 deviation from the blue curve is observed. As  $z$  increases from 6 to 10, the pattern is similar; the quantitative  
525 behavior of  $\text{Var}[H_A]$  would approach the base case, equivalent to  $z = 10$  (green curve in Figure 5E).

526 In the founding parental pool,  $\text{Var}[T] = 4.002$  in our example with  $z \neq k - z$ . After one generation of  
527 mating,  $\text{Var}[T]$  drops to 1.997, 2.934, and 2.923, under random mating, assortative mating by admixture  
528 fraction, and assortative mating by trait, respectively. From  $g = 2$ ,  $\text{Var}[T]$  gradually increases and achieves  
529 steady state values for the three models near 4.942 at  $g = 5$ , 4.934 at  $g = 8$ , and 6.929 at  $g = 14$ , respectively.  
530 In accord with Section 4.1-4.4, assortative mating by trait has the highest  $\text{Var}[T]$  values across generations.

531 If  $z = k - z$  (Figure 8E), then allele “1” has equal probability of traced back to either source population.  
532 In this scenario, the  $\text{Var}[H_A]$  curve from assortative mating follows the  $\text{Var}[H_A]$  curve from random mating.  
533 Because all individuals have the same trait value in the founding parental pool irrespective of their origin,  
534  $\text{Var}[T] = 0$  at  $g = 0$ . For all three mating models,  $\text{Var}[T]$  gradually increases from  $g = 1$  to achieve steady  
535 state values that are the same as those from the  $z \neq k - z$  case.

## 536 5 Discussion

537 In this paper, we have devised a mechanistic admixture model in which an admixed population is formed  
538 from contributions of a pair of mutually isolated source populations under assortative mating, either by  
539 admixture level or by a quantitative phenotype. The approach includes a quantitative-genetic model that  
540 relates a quantitative phenotype to underlying loci affecting its trait value and, ultimately, to mate choices.  
541 The admixture level and the quantitative phenotype are studied using a discrete-time recursion that describes  
542 the evolution of the admixed population. Under this model, we have examined the correlation between genetic  
543 ancestry and phenotype in the admixed population as a function of time under three mating models: random  
544 mating, assortative mating by admixture fraction, and assortative mating by phenotype.

545 Initially, ancestry and phenotype are coupled, as the source populations differ in phenotype. Random  
546 mating then decouples the correlation between ancestry and trait faster than is seen in both assortative  
547 mating models (Figure 4), and assortative mating by phenotype maintains the correlation to a greater extent  
548 than does assortative mating by admixture (Figure 4). Compared with random mating, in a similar manner  
549 to classic assortative mating models [26, 35, 36], the assortative mating increases the population variance of  
550 the property on which the assortment is based (Figures 5 and 6). In fact, our Eq. 13 multiplies the variance  
551 of admixture in a model without assortment [1] by a factor that increases with positive assortative mating.

552 Increasing the strength of assortative mating magnifies the difference observed among the models in the

553 speed at which the correlation declines (Figure 4). Generally similar qualitative patterns are observed if  
554 the source populations are regarded as having fixed differences (Figure 4) or merely frequency differences in  
555 alleles affecting the quantitative trait (Figures 7 and 8).

556 A key observation is that, as the number of loci underlying the quantitative trait increases, the difference  
557 in trajectories between the two assortative mating models decreases. Because assortative mating by admix-  
558 ture fraction affects all loci, whereas assortative mating by trait affects only trait loci and their genomic  
559 neighbors, as the trait is determined by an increasing number of loci, the behavior of assortative mating by  
560 trait increasingly follows that of assortative mating by ancestry (Figure 4). As the number of loci affecting  
561 the phenotype increases, assortative mating by trait increasingly reflects assortative mating by admixture,  
562 and the correlation between ancestry and trait value persists for longer (Figure 4).

563 Our study was motivated partly by a hypothesis of Parra et al. [8] claiming that assortative mating by  
564 the *color* phenotype in Brazil could eventually decouple *color* from genetic ancestry, so that largely separate  
565 subpopulations with distinct *color* could eventually possess similar levels of African genetic ancestry. We  
566 have seen not only that assortative mating by a quantitative trait that differs between source populations can  
567 decouple the phenotype from the genetic ancestry, but that random mating can decouple the phenotype from  
568 genetic ancestry as well. Moreover, the decoupling is *slower* with assortative mating than it is with random  
569 mating. In an admixed population with assortative mating that is heavily influenced by a salient phenotype  
570 (such as *color* in the scenario of Parra et al. [8]), mating by other genetically influenced phenotypes is  
571 random, or perhaps less strongly assortative. Thus, in an admixed population, we might expect that among  
572 all the traits to which genotypes contribute, traits that have little influence on mating behavior will decouple  
573 from ancestry most rapidly. Those traits on which assortment does occur, such as *color* in Brazil, will be the  
574 slowest to decouple from ancestry—but under our model, they eventually will do so. Thus, in an admixed  
575 population, phenotypes that once reflected ancestry in the source populations might no longer be predictive  
576 of genetic ancestry after a sufficient length of time has passed.

577 The focus of our simulations has been on understanding demographic phenomena, but the model is  
578 relevant to efforts to investigate determinants of disease traits in admixed populations. For example, in  
579 admixture-mapping studies and in studies of health disparities involving admixed populations, correlations  
580 of phenotypes and admixture levels are often computed [37–40]. The mechanistic model can potentially  
581 provide insights into the way in which these correlations change over time in scenarios in which specific trait  
582 architectures are of interest.

583 We note that we have examined a model with a single admixture event at the founding of the admixed  
584 population. Under this idealized model, even if the founding admixed population starts with a perfect  
585 correlation between admixture fraction and trait value, the correlation decreases over time and eventually  
586 approaches zero in the absence of further influx from source populations. In principle, the framework can  
587 account for continuous influx from the source populations. If ancestry–trait correlation exists in source  
588 populations, then such influx would be expected to slow the decoupling between admixture and phenotype  
589 in the admixed populations under all three mating models, while qualitatively maintaining their relative  
590 order in the rate of decoupling.

591 Although the theoretical framework we have developed can incorporate various genetic architectures and  
592 population admixture processes (including disassortative mating), our model has a number of limitations.  
593 First, it does not include sex bias during the admixture event, a phenomenon that often occurs in admixture

594 processes [17, 41, 42]. A recent genetic model in a scenario with continuing contributions from the sources  
595 does allow for sex bias with assortative mating, but with no phenotype and with the assortative mating  
596 occurring by population membership—in the admixed population or in one of the sources—rather than by  
597 the admixture level itself [28]. Next, we have chosen to model admixture in an individual as the mean of  
598 parental admixture levels; this approach does not account for stochasticity during genetic transmission. Our  
599 quantitative trait model does not incorporate dominance, spatial positioning of trait loci along a genome,  
600 epistasis, environmental effects on the phenotype, or genotype-by-environment interaction. The latter pair  
601 of limitations might be particularly important in using the results in human data, as numerous studies  
602 have shown that assortative mating often operates on sociocultural traits [43–46]. It will be important to  
603 extend our theoretical framework to include additional features of the quantitative-genetic model, such as by  
604 incorporating dominance, linkage, variable effect sizes across the loci contributing to the quantitative trait,  
605 and varying heritability of the phenotype.

606 **Acknowledgments.** We acknowledge support from National Institutes of Health grants R01 HG005855  
607 and F32 GM130050 and National Science Foundation grant BCS-1515127.

## 608 References

- 609 [1] Verdu P, Rosenberg NA. A general mechanistic model for admixture histories of hybrid populations. *Genetics*.  
610 2011;189(4):1413–1426. doi:10.1534/genetics.111.132787.
- 611 [2] Gravel S. Population genetics models of local ancestry. *Genetics*. 2012;191(2):607–619. doi:10.1534/genetics.112.139808.
- 612 [3] Shriver MD, Smith MW, Jin L, Marcini A, Akey JM, Deka R, et al. Ethnic-affiliation estimation by use of population-  
613 specific DNA markers. *The American Journal of Human Genetics*. 1997;60(4):957–964.
- 614 [4] Parra EJ, Marcini A, Akey J, Martinson J, Batzer MA, Cooper R, et al. Estimating African American admixture  
615 proportions by use of population-specific alleles. *The American Journal of Human Genetics*. 1998;63(6):1839–1851.  
616 doi:10.1086/302148.
- 617 [5] Parra EJ, Kittles RA, Shriver MD. Implications of correlations between skin color and genetic ancestry for biomedical  
618 research. *Nature Genetics*. 2004;36:S54–S60. doi:10.1038/ng1440.
- 619 [6] Ruiz-Linares A, Adhikari K, Acuña-Alonzo V, Quinto-Sanchez M, Jaramillo C, Arias W, et al. Admixture in Latin  
620 America: geographic structure, phenotypic diversity and self-perception of ancestry based on 7,342 individuals. *PLOS*  
621 *Genetics*. 2014;10(9):1–13. doi:10.1371/journal.pgen.1004572.
- 622 [7] Algee-Hewitt BFB. Population inference from contemporary American craniometrics. *American Journal of Physical*  
623 *Anthropology*. 2016;160(4):604–624. doi:10.1002/ajpa.22959.
- 624 [8] Parra FC, Amado RC, Lambertucci JR, Rocha J, Antunes CM, Pena SDJ. Color and genomic ancestry in Brazil-  
625 ians. *Proceedings of the National Academy of Sciences of the United States of America*. 2003;100(1):177–182.  
626 doi:10.1073/pnas.0126614100.
- 627 [9] Pimenta JR, Zuccherato LW, Debes AA, Maselli L, Soares RP, Moura-Neto RS, et al. Color and genomic ancestry in  
628 Brazilians: a study with forensic microsatellites. *Human Heredity*. 2006;62(4):190–195. doi:10.1159/000096872.
- 629 [10] Leite TKM, Fonseca RMC, de França NM, Parra EJ, Pereira RW. Genomic ancestry, self-reported “color” and  
630 quantitative measures of skin pigmentation in Brazilian admixed siblings. *PLOS ONE*. 2011;6(11):e27162–e27162.  
631 doi:10.1371/journal.pone.0027162.
- 632 [11] Beleza S, Johnson NA, Candille SI, Absher DM, Coram MA, Lopes J, et al. Genetic architecture of skin and eye color in  
633 an African-European admixed population. *PLOS Genetics*. 2013;9(3):1–15. doi:10.1371/journal.pgen.1003372.
- 634 [12] Durso DF, Bydlowski SP, Hutz MH, Suarez-Kurtz G, Magalhães TR, Junho Pena SD. Association of genetic variants with  
635 self-assessed color categories in Brazilians. *PLOS ONE*. 2014;9(1):1–8. doi:10.1371/journal.pone.0083926.

- 636 [13] Magalhães da Silva T, Sandhya Rani MR, de Oliveira Costa GN, Figueiredo MA, Melo PS, Nascimento JaF, et al. The  
637 correlation between ancestry and color in two cities of Northeast Brazil with contrasting ethnic compositions. *European*  
638 *Journal Of Human Genetics*. 2014;23:984–989. doi:10.1038/ejhg.2014.215.
- 639 [14] Risch N, Choudhry S, Via M, Basu A, Sebros R, Eng C, et al. Ancestry-related assortative mating in Latino populations.  
640 *Genome Biology*. 2009;10(11):1–16. doi:10.1186/gb-2009-10-11-r132.
- 641 [15] Zou JY, Park DS, Burchard EG, Torgerson DG, Pino-Yanes M, Song YS, et al. Genetic and socioeconomic study of mate  
642 choice in Latinos reveals novel assortment patterns. *Proceedings of the National Academy of Sciences of the United States*  
643 *of America*. 2015;112(44):13621–13626. doi:10.1073/pnas.1501741112.
- 644 [16] Goldberg A, Verdu P, Rosenberg NA. Autosomal admixture levels are informative about sex bias in admixed populations.  
645 *Genetics*. 2014;198(3):1209–1229. doi:10.1534/genetics.114.166793.
- 646 [17] Goldberg A, Rosenberg NA. Beyond 2/3 and 1/3: the complex signatures of sex-biased admixture on the X Chromosome.  
647 *Genetics*. 2015;201(1):263–279. doi:10.1534/genetics.115.178509.
- 648 [18] Long JC. The genetic structure of admixed populations. *Genetics*. 1991;127(2):417–428.
- 649 [19] Ewens WJ, Spielman RS. The transmission/disequilibrium test: history, subdivision, and admixture. *The American*  
650 *Journal of Human Genetics*. 1995;57(2):455–464.
- 651 [20] Guo W, Fung WK, Shi N, Guo J. On the formula for admixture linkage disequilibrium. *Human Heredity*. 2005;60(3):177–  
652 180. doi:10.1159/000090119.
- 653 [21] Edge MD, Rosenberg NA. Implications of the apportionment of human genetic diversity for the apportionment of human  
654 phenotypic diversity. *Studies in History and Philosophy of Science Part C: Studies in History and Philosophy of Biological*  
655 *and Biomedical Sciences*. 2015;52:32–45. doi:10.1016/j.shpsc.2014.12.005.
- 656 [22] Edge MD, Rosenberg NA. A general model of the relationship between the apportionment of human ge-  
657 netic diversity and the apportionment of human phenotypic diversity. *Human Biology*. 2015;87(4):313–337.  
658 doi:10.13110/humanbiology.87.4.0313.
- 659 [23] Jennings HS. The numerical results of diverse systems of breeding. *Genetics*. 1916;1(1):53–89.
- 660 [24] Fisher RA. XV. The correlation between relatives on the supposition of Mendelian inheritance. *Transactions of the Royal*  
661 *Society of Edinburgh*. 1919;52(2):399–433. doi:10.1017/S0080456800012163.
- 662 [25] Wright S. Systems of mating. III. Assortative mating based on somatic resemblance. *Genetics*. 1921;6(2):144–161.
- 663 [26] Crow JF, Kimura M. *An introduction to population genetics theory*. Blackburn Press; 2009.
- 664 [27] Zaitlen N, Huntsman S, Hu D, Spear M, Eng C, Oh SS, et al. The effects of migration and assortative mating on admixture  
665 linkage disequilibrium. *Genetics*. 2017;205(1):375–383. doi:10.1534/genetics.116.192138.
- 666 [28] Goldberg A, Rastogi A, Rosenberg NA. Assortative mating by population of origin in a mechanistic model of admixture.  
667 *Theoretical Population Biology*. 2019;.
- 668 [29] Kullback S. *Information theory and statistics*. Mineola, NY: Dover Publications; 1997.
- 669 [30] Nesterov Y, Nemirovskii A. *Interior-point polynomial algorithms in convex programming*. Society for Industrial and  
670 *Applied Mathematics*; 1994.
- 671 [31] Forsgren A, Gill PE, Wright MH. Interior methods for nonlinear optimization. *SIAM Review*. 2002;44(4):525–597.  
672 doi:10.1137/S0036144502414942.
- 673 [32] MOSEK ApS. *MOSEK Rmosek Package 8.0.0.78*. 2017;.
- 674 [33] Boyd S, Vandenberghe L. *Convex optimization*. Cambridge, UK: Cambridge University Press; 2004.
- 675 [34] Charlesworth B, Charlesworth D. *Elements of evolutionary genetics*. Greenwood Village, CO: Roberts and Co. Publishers;  
676 2010.
- 677 [35] Crow JF, Felsenstein J. The effect of assortative mating on the genetic composition of a population. *Eugenics Quarterly*.  
678 1968;15(2):85–97. doi:10.1080/19485565.1968.9987760.
- 679 [36] Felsenstein J. Continuous-genotype models and assortative mating. *Theoretical Population Biology*. 1981;19(3):341–357.  
680 doi:10.1016/0040-5809(81)90025-3.
- 681 [37] Tang H, Jorgenson E, Gadde M, Kardina SLR, Rao DC, Zhu X, et al. Racial admixture and its impact on BMI and blood  
682 pressure in African and Mexican Americans. *Human Genetics*. 2006;119(6):624–633. doi:10.1007/s00439-006-0175-4.

- 683 [38] Peralta CA, Ziv E, Katz R, Reiner A, Burchard EG, Fried L, et al. African ancestry, socioeconomic status, and kidney  
684 function in elderly African Americans: a genetic admixture analysis. *Journal of the American Society of Nephrology*.  
685 2006;17(12):3491–3496. doi:10.1681/ASN.2006050493.
- 686 [39] Gravlee CC, Non AL, Mulligan CJ. Genetic ancestry, social classification, and racial inequalities in blood pressure in  
687 Southeastern Puerto Rico. *PLOS ONE*. 2009;4(9):e6821–e6821. doi:10.1371/journal.pone.0006821.
- 688 [40] Non AL, Gravlee CC, Mulligan CJ. Education, genetic ancestry, and blood pressure in African Americans and Whites.  
689 *American Journal of Public Health*. 2012;102(8):1559–1565. doi:10.2105/AJPH.2011.300448.
- 690 [41] Wilkins JF. Unraveling male and female histories from human genetic data. *Current Opinion in Genetics & Development*.  
691 2006;16(6):611–617. doi:10.1016/j.gde.2006.10.004.
- 692 [42] Adhikari K, Chacón-Duque JC, Mendoza-Revilla J, Fuentes-Guajardo M, Ruiz-Linares A. The genetic diversity of the  
693 Americas. *Annual Review of Genomics and Human Genetics*. 2017;18(1):277–296. doi:10.1146/annurev-genom-083115-  
694 022331.
- 695 [43] Kalmijn M. Inter-marriage and homogamy: causes, patterns, trends. *Annual Review of Sociology*. 1998;24(1):395–421.  
696 doi:10.1146/annurev.soc.24.1.395.
- 697 [44] Watson D, Klohnen EC, Casillas A, Nus Simms E, Haig J, Berry DS. Match makers and deal breakers: analyses of assort-  
698 ative mating in newlywed couples. *Journal of Personality*. 2004;72(5):1029–1068. doi:10.1111/j.0022-3506.2004.00289.x.
- 699 [45] Rosenfeld MJ. Racial, educational and religious endogamy in the United States: a comparative historical perspective.  
700 *Social Forces*. 2008;87(1):1–31. doi:10.1353/sof.0.0077.
- 701 [46] Schwartz CR. Trends and variation in assortative mating: causes and consequences. *Annual Review of Sociology*.  
702 2013;39(1):451–470. doi:10.1146/annurev-soc-071312-145544.



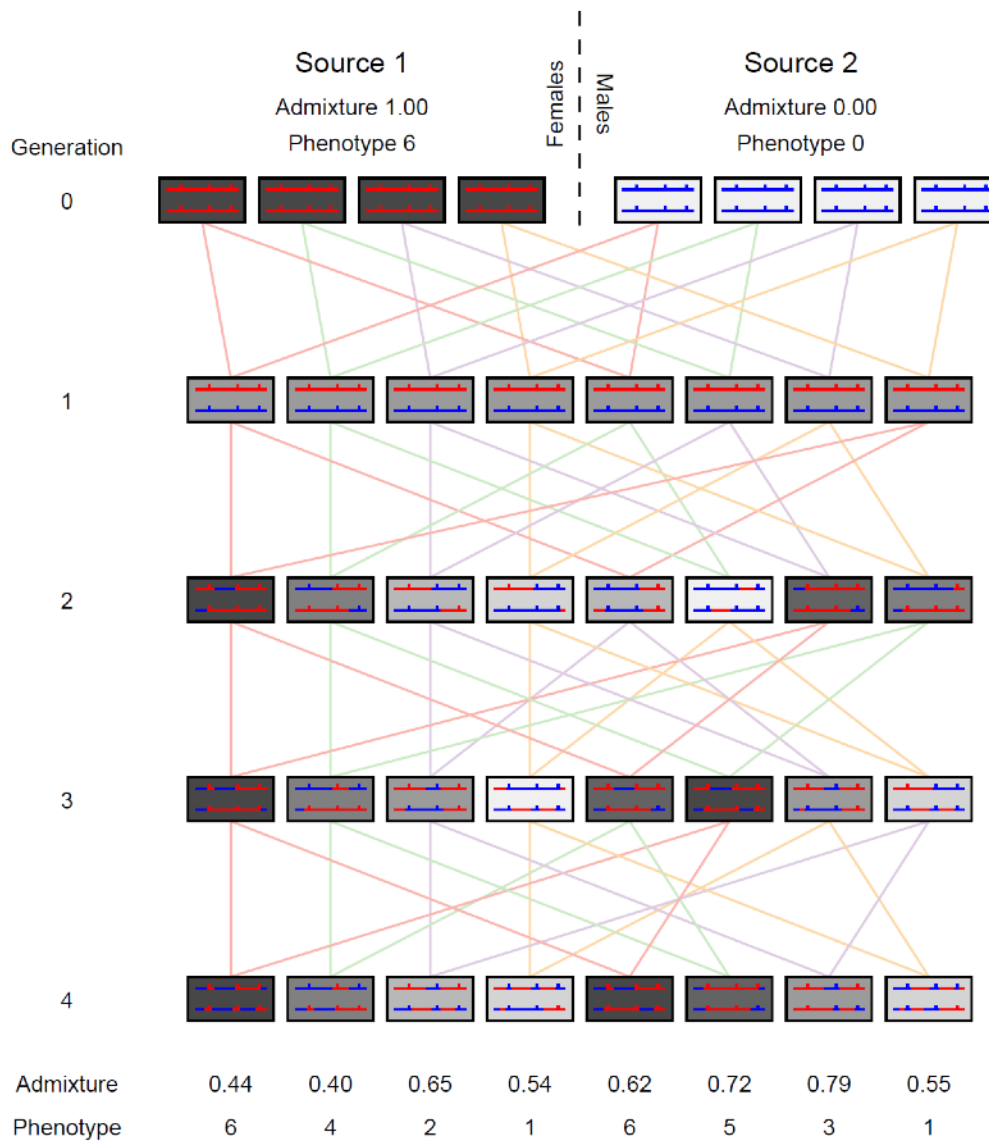


Figure 1: A schematic of an admixture process with positive assortative mating by a phenotype initially correlated with admixture levels. In generation 0, an admixture process begins with females from one population (source 1, left) and males from another (source 2, right). For a quantitative phenotype, source population 1 begins with a high trait value of 6 and source population 2 has a low trait value of 0. Three loci contribute additively to the genetic architecture of the phenotype; each allele derived from source population 1 contributes a value of 1 to the phenotype. The phenotype is represented by the shading of a box. Individuals are depicted as pairs of chromosomes with the ancestral sources of those chromosomes; short vertical lines along the chromosome indicate the three loci that contribute to the phenotype. After generation 1, positive assortative mating by phenotype proceeds in the admixed population. Lines connecting generations are displayed in four colors, representing four mating pairs. Initially, in generation 2, a strong correlation exists between admixture and phenotype ( $r = 0.96$ ). By generation 4, however, owing to recombination events that stochastically dissociate the trait loci from the overall genetic admixture, the genetic admixture has been decoupled from the phenotype, so that some of the individuals with the highest trait values have among the lowest admixture coefficients for source population 1, and the correlation between phenotype and overall genetic admixture has dissipated ( $r = -0.09$ ).



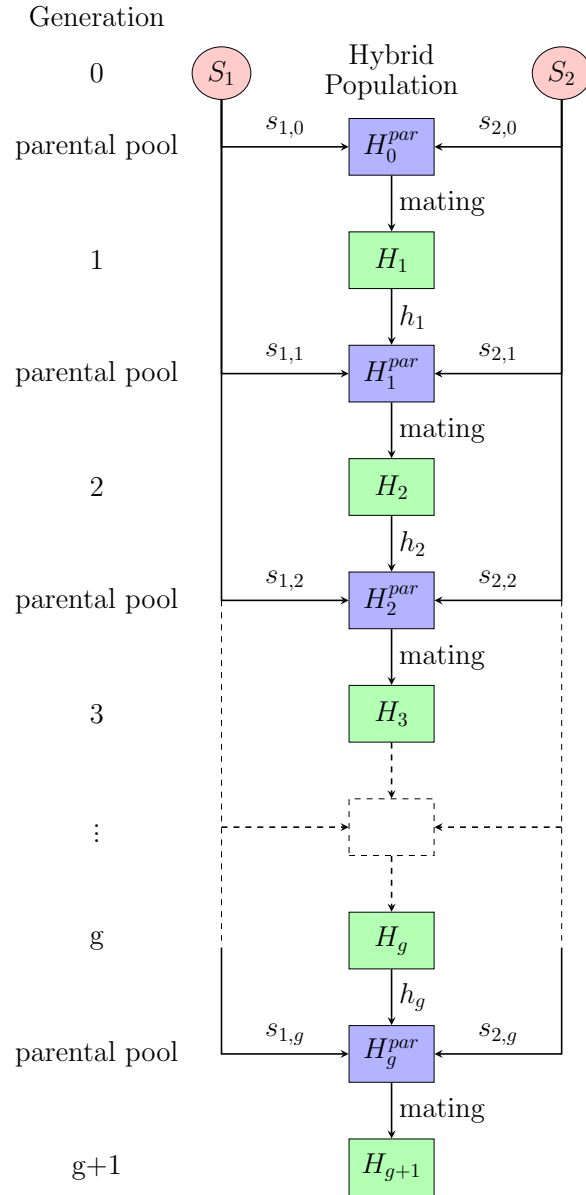


Figure 2: A schematic diagram of the admixture process. At the founding of the population ( $g = 0$ ), two isolated source populations produce the first generation of a admixed population ( $H_1$ ). In the subsequent generations ( $g \geq 1$ ), populations from  $S_1$ ,  $S_2$ , and  $H_g$  constitute a parental pool  $H_g^{\text{par}}$  at generation  $g$  from which the admixed population  $H_{g+1}$  at generation  $g + 1$  is produced. Fractional contributions from three populations in forming the parental pool are  $s_{1,g}$ ,  $s_{2,g}$ , and  $h_g$ , respectively. Individuals in the parental pool mate based on mating models described in Section 2.3.

	<b>Locus</b>							
	1	2	3	4	5	6	7	8
$L_{i1}$	1	0	0	1	1	1	0	1
$L_{i2}$	0	0	0	1	1	1	1	1
$X_i$	1	1	0	0	0	1	1	0
<b>Contribution to <math>T</math></b>	1	0	2	0	0	2	1	0
<b><math>T = 6</math></b>								

Figure 3: An example of our quantitative trait model. Here, a diploid individual with  $k = 8$  trait loci is shown. At each locus  $i$ , an allele  $L_{ij}$  contributes to the overall trait value if and only if  $L_{ij} = X_i$  where  $X_i$  is a variable indicating which of two alleles, “0” or “1” increases the trait value. The total trait value of an individual equals the number of alleles satisfying  $L_{ij} = X_i$  across the  $k$  trait loci. In this example, the individual has  $T = 6$ .

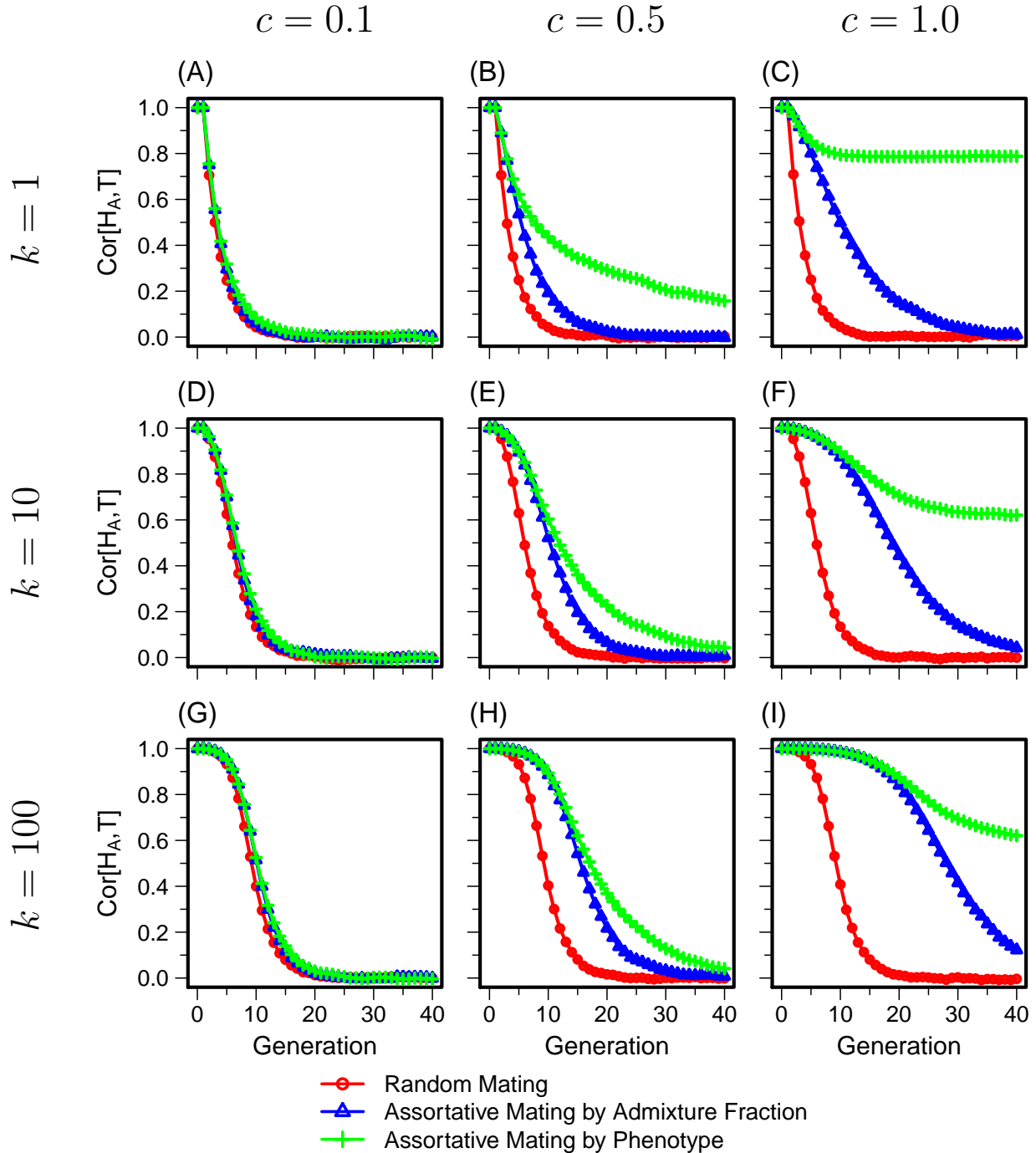


Figure 4: Correlation between admixture fraction and quantitative trait value ( $\text{Cor}[H_A, T]$ ) as a function of time. All parameter values in panel (E) follow the base case described in Section 3.2; the number of quantitative trait loci  $k$  and the assortative mating strength  $c$  vary across panels. In each panel, for a given  $(k, c)$  pair, for each mating scheme, the mean of 100 simulated trajectories is plotted. The red, blue, and green curves represent results from random mating, assortative mating by admixture fraction, and assortative mating by phenotype, respectively. (A)  $k = 1$ ,  $c = 0.1$ . (B)  $k = 1$ ,  $c = 0.5$ . (C)  $k = 1$ ,  $c = 1.0$ . (D)  $k = 10$ ,  $c = 0.1$ . (E)  $k = 10$ ,  $c = 0.5$ . (F)  $k = 10$ ,  $c = 1.0$ . (G)  $k = 100$ ,  $c = 0.1$ . (H)  $k = 100$ ,  $c = 0.5$ . (I)  $k = 100$ ,  $c = 1.0$ .

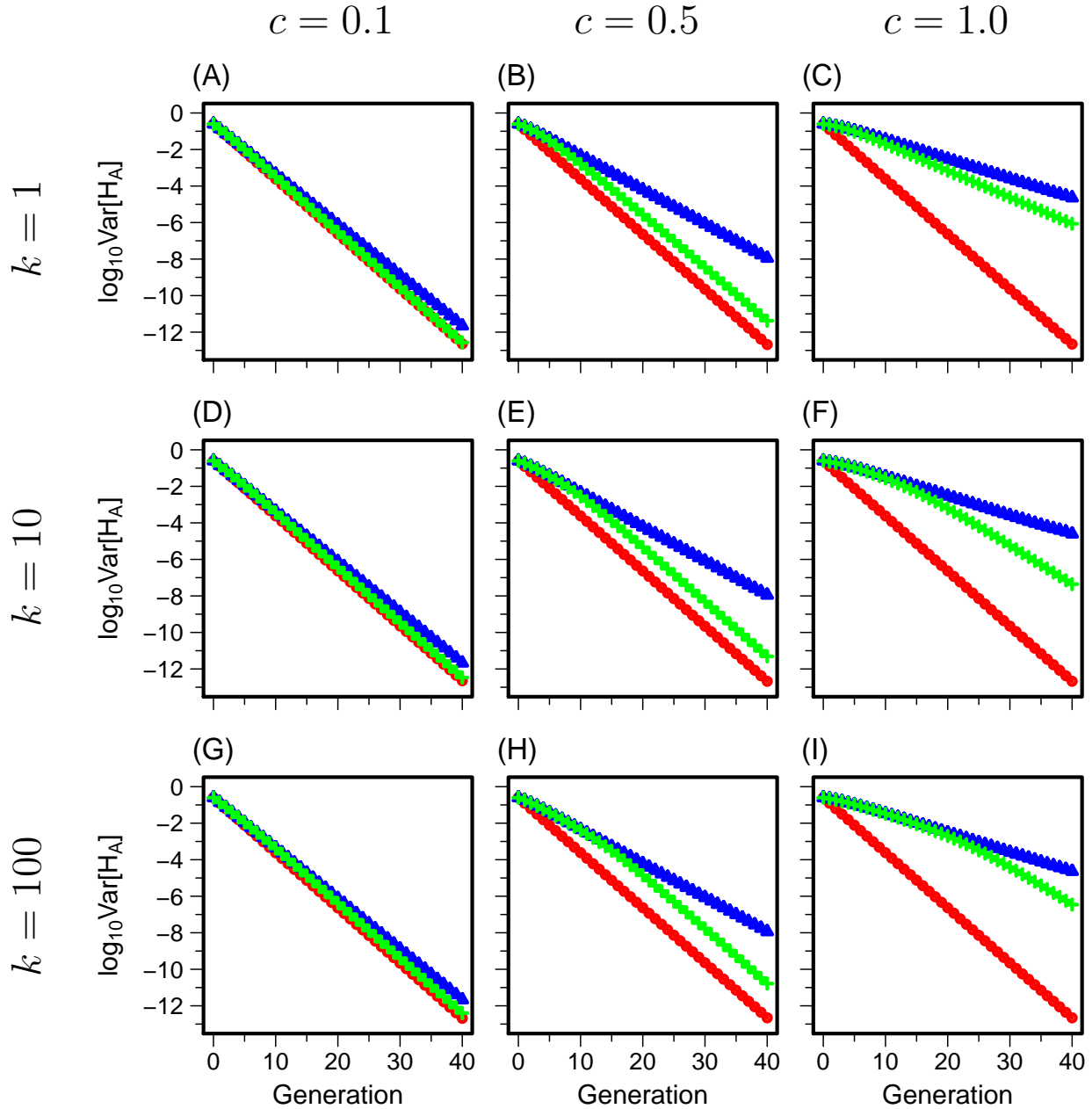


Figure 5: Variance of admixture fraction ( $\text{Var}[H_A]$ ) as a function of time. The simulations shown are the same ones from Figure 4. (A)  $k = 1$ ,  $c = 0.1$ . (B)  $k = 1$ ,  $c = 0.5$ . (C)  $k = 1$ ,  $c = 1.0$ . (D)  $k = 10$ ,  $c = 0.1$ . (E)  $k = 10$ ,  $c = 0.5$ . (F)  $k = 10$ ,  $c = 1.0$ . (G)  $k = 100$ ,  $c = 0.1$ . (H)  $k = 100$ ,  $c = 0.5$ . (I)  $k = 100$ ,  $c = 1.0$ . Colors and symbols follow Figure 4. The y-axis is plotted on a logarithmic scale.

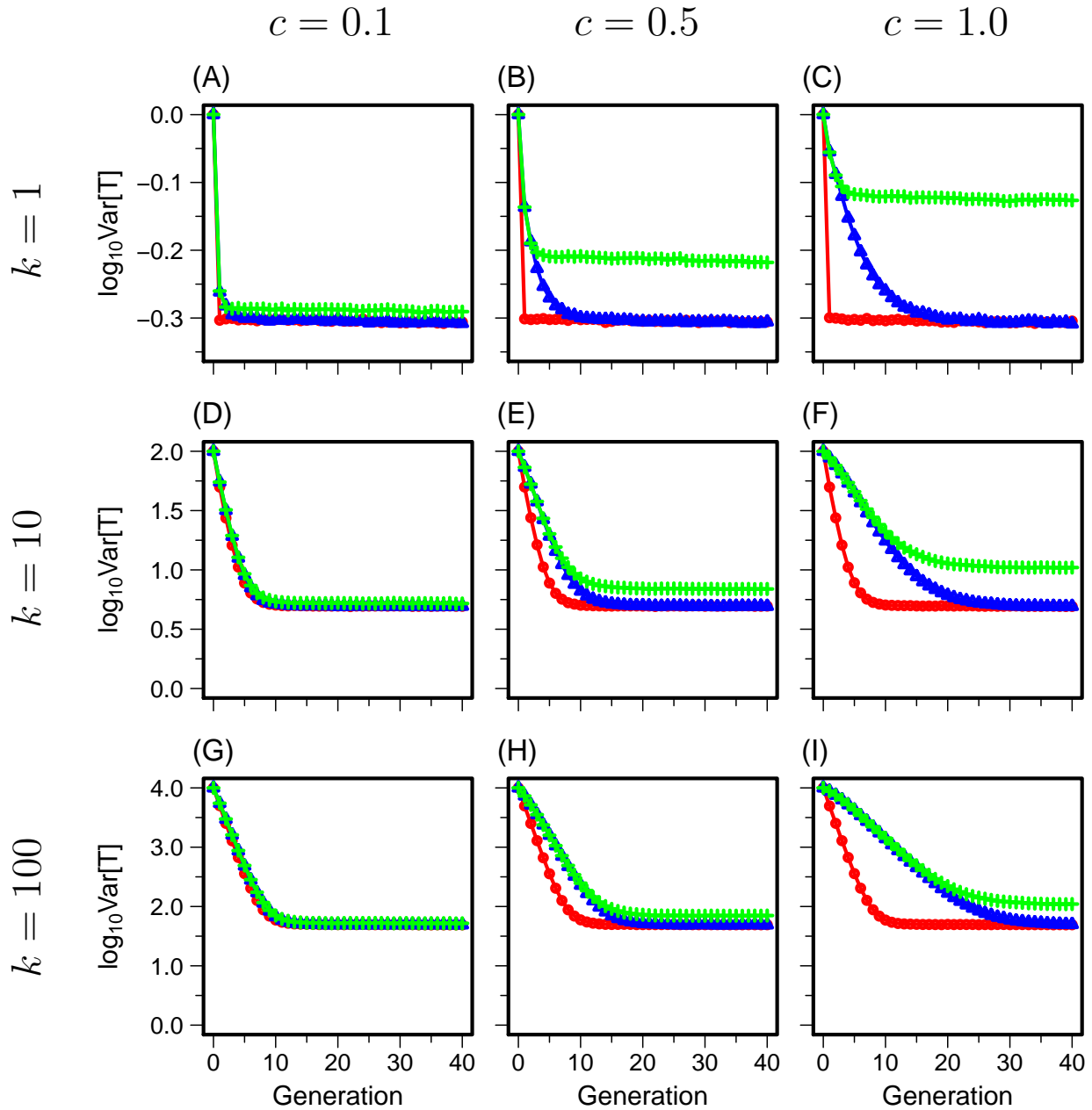


Figure 6: Variance of the phenotype ( $\text{Var}[T]$ ) as a function of time. The simulations shown are the same ones from Fig. 4. (A)  $k = 1$ ,  $c = 0.1$ . (B)  $k = 1$ ,  $c = 0.5$ . (C)  $k = 1$ ,  $c = 1.0$ . (D)  $k = 10$ ,  $c = 0.1$ . (E)  $k = 10$ ,  $c = 0.5$ . (F)  $k = 10$ ,  $c = 1.0$ . (G)  $k = 100$ ,  $c = 0.1$ . (H)  $k = 100$ ,  $c = 0.5$ . (I)  $k = 100$ ,  $c = 1.0$ . Colors and symbols follow Figure 4. The y-axis is plotted on a logarithmic scale.

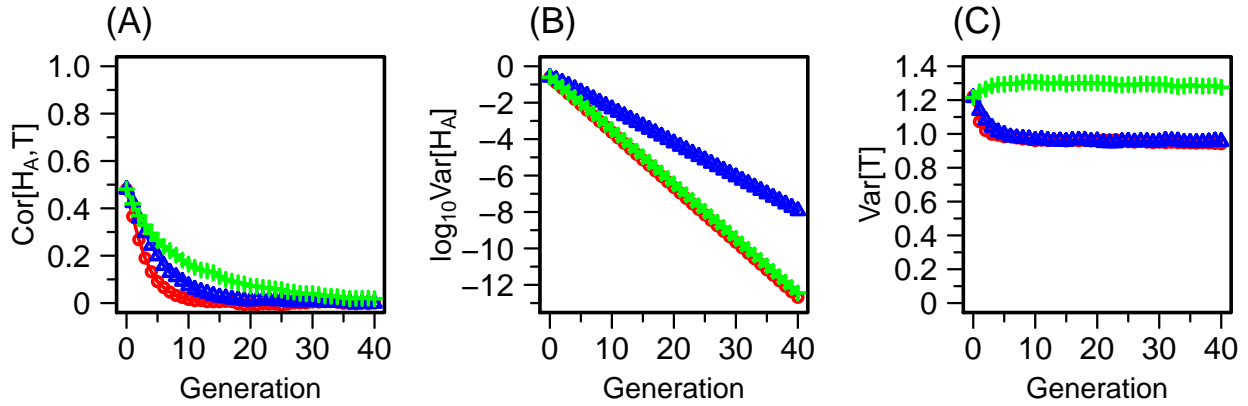


Figure 7:  $\text{Cor}[H_A, T]$ ,  $\text{Var}[H_A]$ , and  $\text{Var}[T]$  in a model in which the allele frequencies  $p_i$  and  $q_i$  in the source populations  $S_1$  and  $S_2$  are drawn from a simulation rather than being treated as fixed at 1 and 0, respectively. All other parameters are kept at the values of the base case (Section 3.2). (A) Correlation between admixture fraction and trait ( $\text{Cor}[H_A, T]$ ). (B) Variance of the admixture fraction ( $\text{Var}[H_A]$ ). (C) Variance of the phenotype ( $\text{Var}[T]$ ). Colors and symbols follow Figure 4. The figure relies on a single replicate of simulated allele frequencies  $p_i$  and  $q_i$  following a genetic drift model in which  $S_1$  and  $S_2$  descend from a common ancestral population, as described in Section 3.1. The simulated allele frequencies across  $k = 10$  loci have mean values of  $\bar{p} \approx 0.502$  and  $\bar{q} \approx 0.449$  and variance  $s_p^2 \approx 0.214$  and  $s_q^2 \approx 0.235$ . If we let  $\delta_i = p_i - q_i$ , with  $\delta_i > 0$ , then the mean of the allele frequency difference across the 10 loci is  $\bar{\delta} \approx 5.310 \times 10^{-2}$ , with  $\bar{\delta}^2 \approx 7.865 \times 10^{-3}$ . Across  $k = 10$  loci,  $F_{ST} \approx 0.075$ , as computed using Eq. 14 of [22]. The y-axis of  $\text{Var}[H_A]$  is plotted on a logarithmic scale. Results using other replicates of simulated allele frequencies with  $k = 10$  are shown in Figure S4.

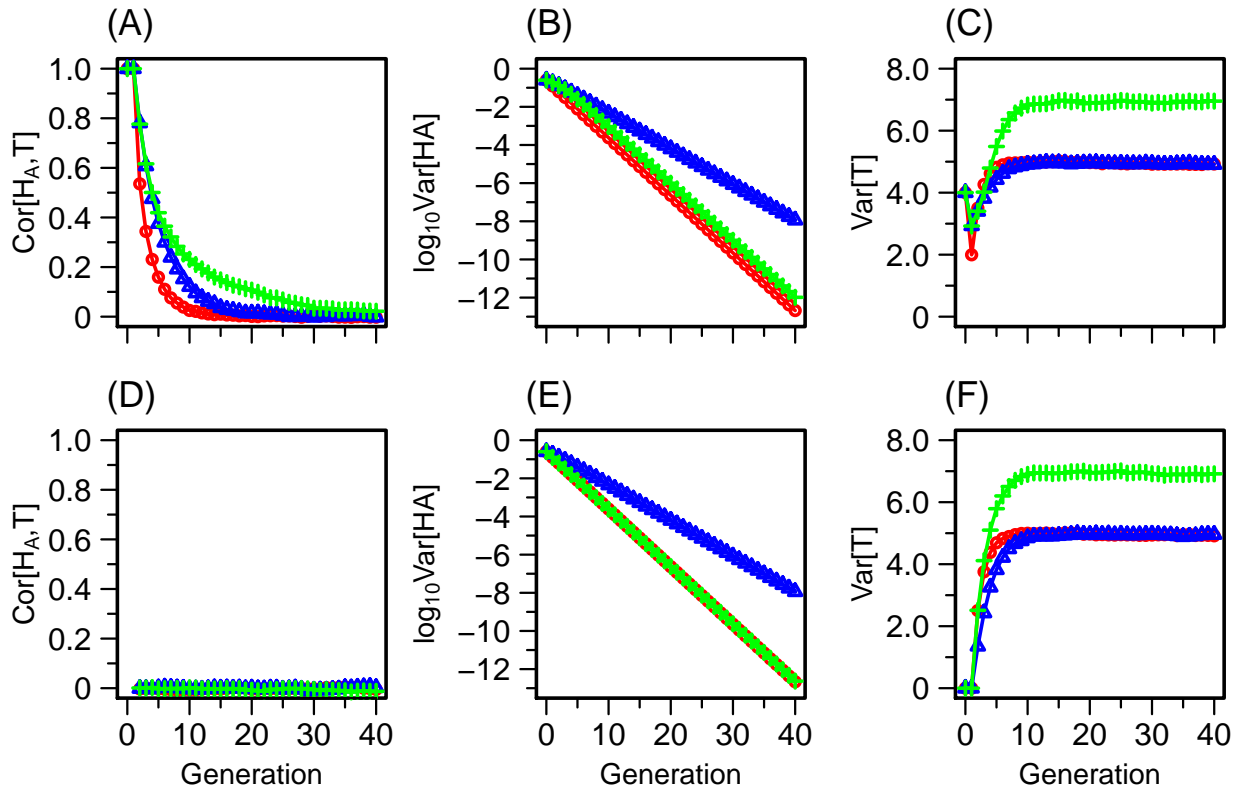


Figure 8:  $\text{Cor}[H_A, T]$ ,  $\text{Var}[H_A]$ , and  $\text{Var}[T]$  under a trait model in which trait loci do not systematically have greater values in one source population:  $P(X_i = 1) = P(X_i = 0) = 0.5$ . All other parameters are kept at the values of the base case (Section 3.2). Of  $k = 10$  trait loci, we denote the number of randomly selected loci to have  $X_i = 1$  by  $z$ . (A)  $\text{Cor}[H_A, T]$ ,  $z = 6$ . (B)  $\text{Var}[H_A]$ ,  $z = 6$ . (C)  $\text{Var}[T]$ ,  $z = 6$ . (D)  $\text{Cor}[H_A, T]$ ,  $z = 5$ . (E)  $\text{Var}[H_A]$ ,  $z = 5$ . (F)  $\text{Var}[T]$ ,  $z = 5$ . Colors and symbols follow Figure 4. Panels A-C and D-F each relies on a single replicate of a set of  $X_i$  obtained by sampling the  $X_i$  from a  $\text{Binomial}(10, \frac{1}{2})$  distribution and retaining those with the specified value of  $z = 6$  (top panels) and  $z = 5$  (bottom panels). The y-axis of  $\text{Var}[H_A]$  is plotted on a logarithmic scale. Results using other replicates of simulated allele frequencies with  $k = 10$  are shown in Figures S5 (for  $z = 6$ ) and S6 (for  $z = 5$ ).



703 **Supplementary Material**

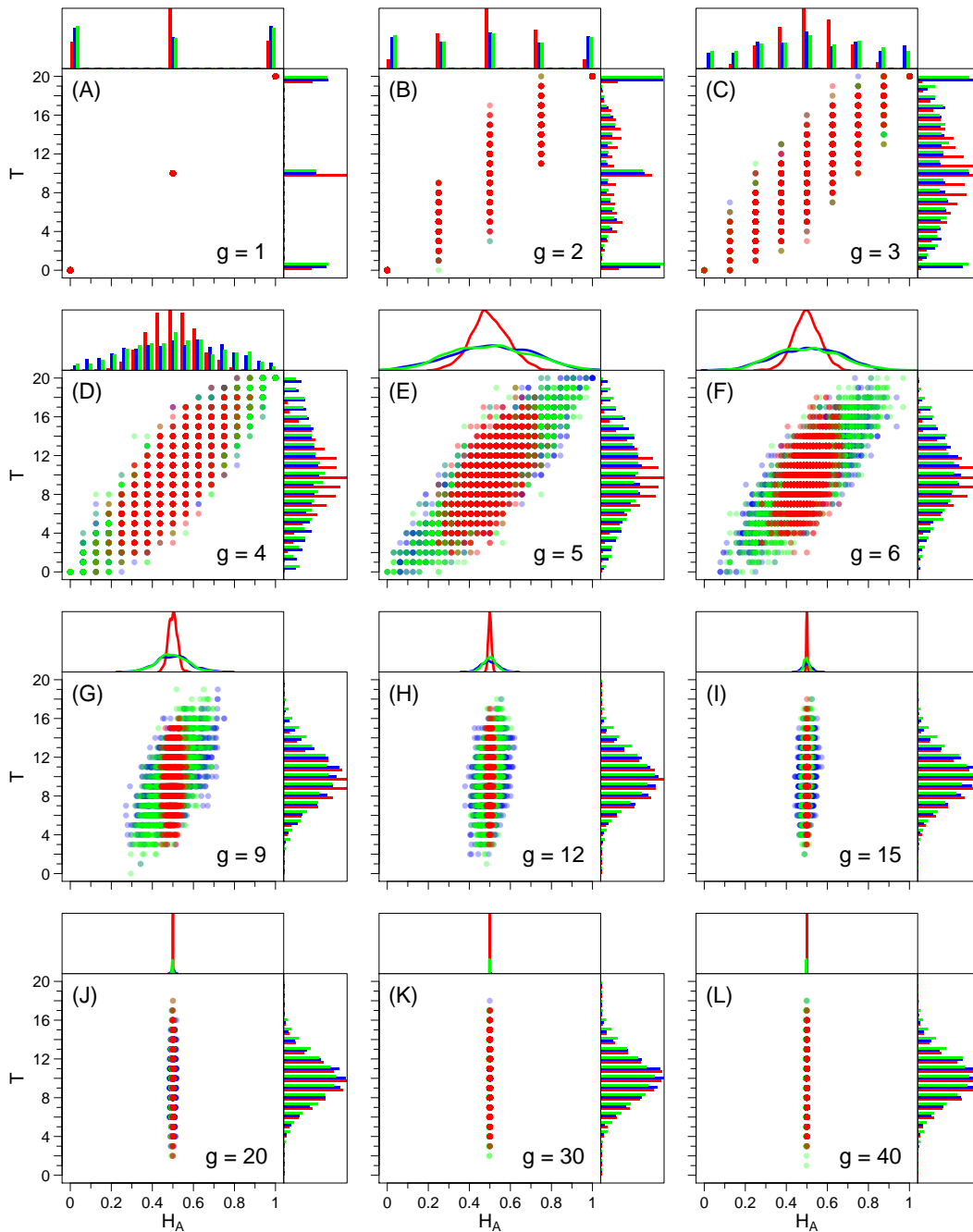


Figure S1: Joint distribution of  $H_A$  and  $T$  as a function of time. The simulation shown is the same one from Figures 4E, 5E, and 6E, using the parameters from base case (Section 3.2). As described in Section 2.1 and 2.2, the possible values for the admixture fraction at generation  $g$  are  $0, 1/2^g, 2/2^g, \dots, (2^g - 1)/2^g, 1$ , whereas the possible values for the trait are  $0, 1, \dots, 2k$  across all generations. In each panel, the top, right, and center plots display a marginal distribution of  $H_A$ , a marginal distribution of  $T$ , and a joint distribution of  $H_A$  and  $T$ , respectively. Colors follow Figure 4.

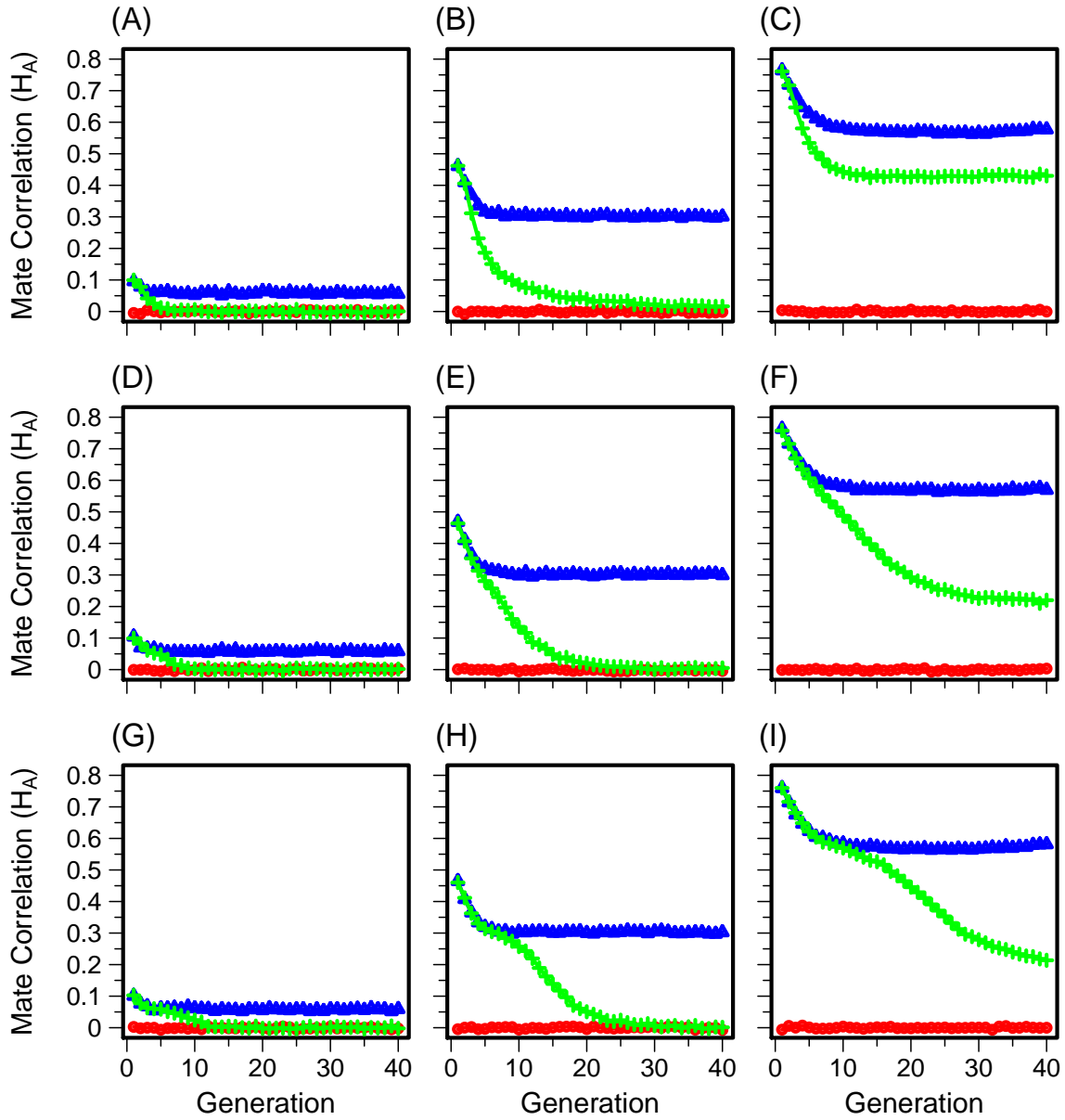


Figure S2: The correlation coefficient  $\text{Cor}[H_{A,g}^f, H_{A,g}^m]$  between the admixture fractions of the members of mating pairs as a function of time. The simulations shown are the same ones from Fig. 4. (A)  $k = 1, c = 0.1$ . (B)  $k = 1, c = 0.5$ . (C)  $k = 1, c = 1.0$ . (D)  $k = 10, c = 0.1$ . (E)  $k = 10, c = 0.5$ . (F)  $k = 10, c = 1.0$ . (G)  $k = 100, c = 0.1$ . (H)  $k = 100, c = 0.5$ . (I)  $k = 100, c = 1.0$ . Colors and symbols follow Figure 4.

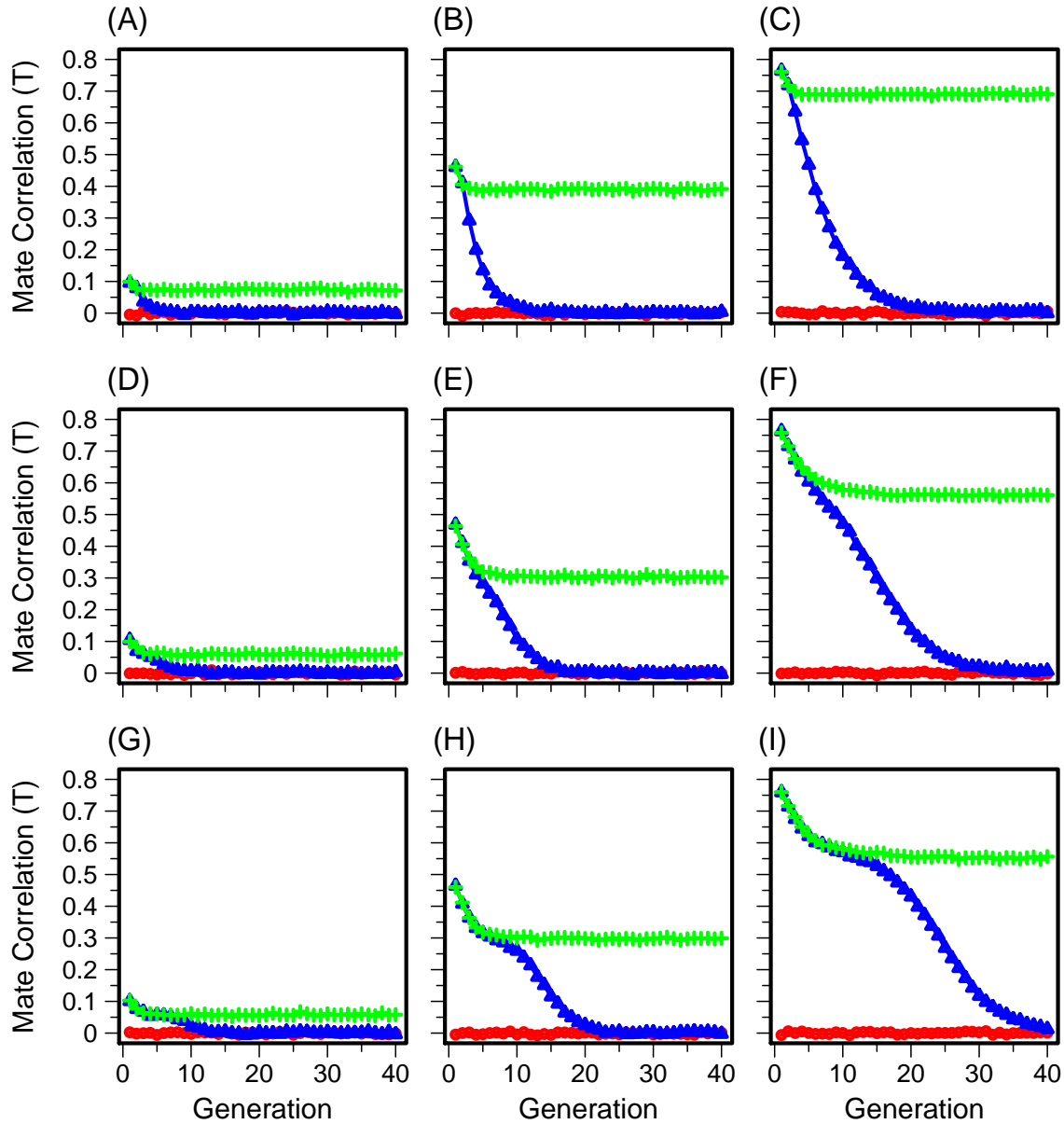


Figure S3: The correlation coefficient  $\text{Cor}[T_g^f, T_g^m]$  between the phenotypes of the members of mating pairs as a function of time. The simulations shown are the same ones from Fig. 4. (A)  $k = 1, c = 0.1$ . (B)  $k = 1, c = 0.5$ . (C)  $k = 1, c = 1.0$ . (D)  $k = 10, c = 0.1$ . (E)  $k = 10, c = 0.5$ . (F)  $k = 10, c = 1.0$ . (G)  $k = 100, c = 0.1$ . (H)  $k = 100, c = 0.5$ . (I)  $k = 100, c = 1.0$ . Colors and symbols follow Figure 4.

Jaehee Kim

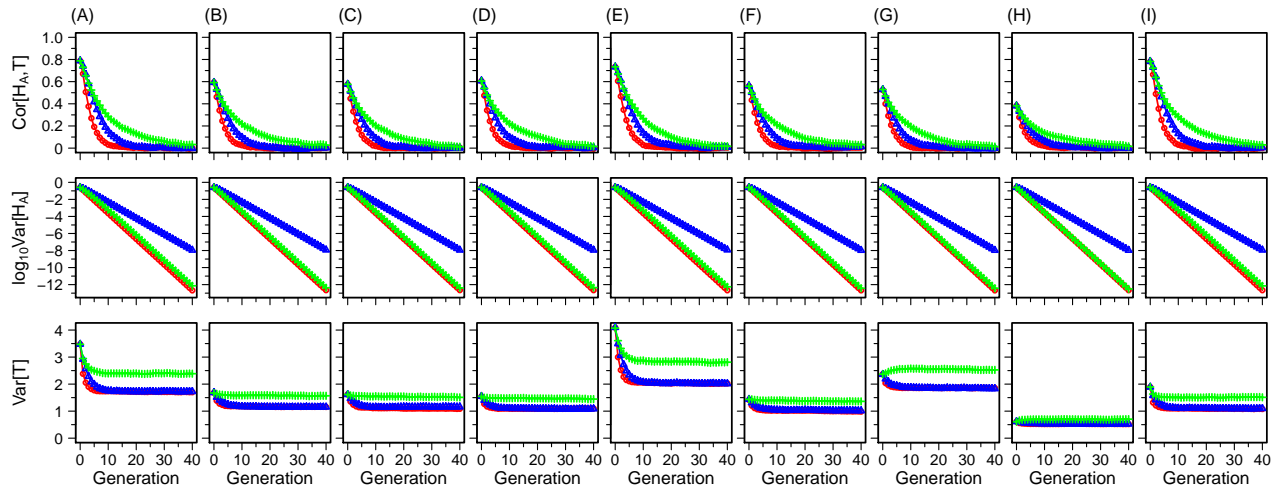


Figure S4:  $\text{Cor}[H_A, T]$  (top row),  $\text{Var}[H_A]$  (middle row), and  $\text{Var}[T]$  (bottom row) using different replicate sets of simulated allele frequencies  $p_i$  and  $q_i$  with  $k = 10$  loci, as described in Section 3.1 and Figure 7. Different columns represent results from different replicates of simulated  $p_i$  and  $q_i$ . Colors and symbols follow Figure 4. The y-axis of  $\text{Var}[H_A]$  is plotted on a logarithmic scale with base 10.

Jaehee Kim

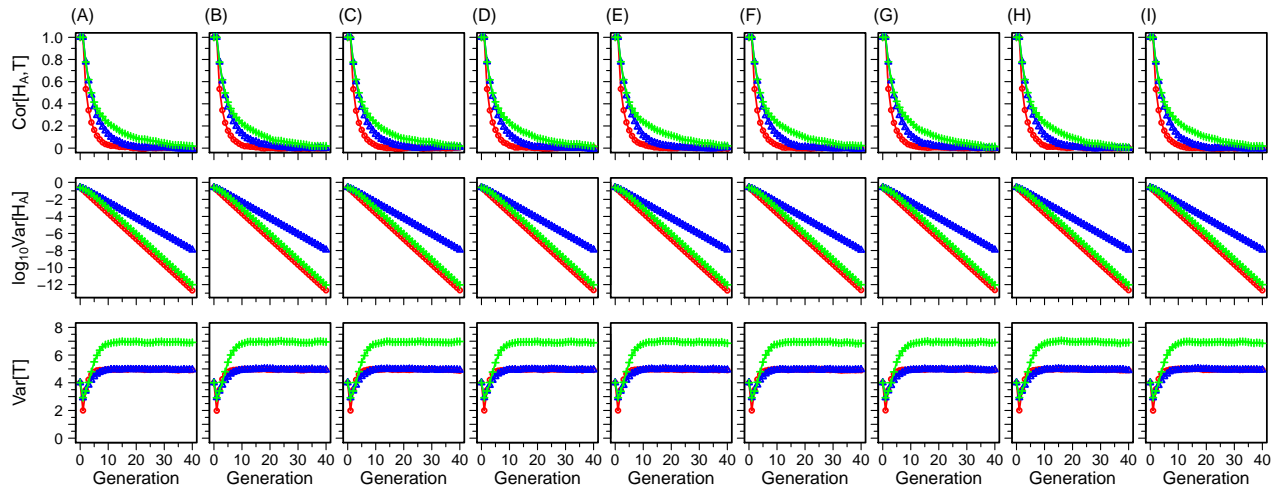


Figure S5:  $\text{Cor}[H_A, T]$  (top row),  $\text{Var}[H_A]$  (middle row), and  $\text{Var}[T]$  (bottom row) using different replicates of a set of  $X_i$  for  $k = 10$  loci sampled from a  $\text{Binomial}(10, \frac{1}{2})$  distribution with a constraint  $z = 6$  (Section 2.2 and Figure 8A-C). Colors and symbols follow Figure 4. The y-axis of  $\text{Var}[H_A]$  is plotted on a logarithmic scale.

Jaehee Kim

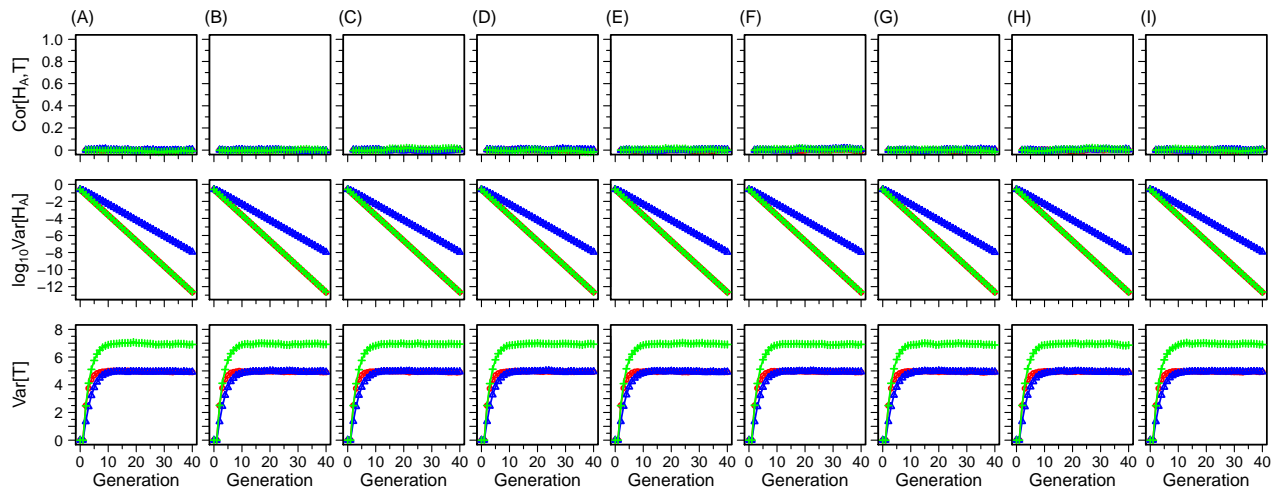


Figure S6:  $\text{Cor}[H_A, T]$  (top row),  $\text{Var}[H_A]$  (middle row), and  $\text{Var}[T]$  (bottom row) using different replicates of a set of  $X_i$  for  $k = 10$  loci sampled from a  $\text{Binomial}(10, \frac{1}{2})$  distribution with a constraint  $z = 5$  (Section 2.2 and Figure 8D-F). Colors and symbols follow Figure 4. The y-axis of  $\text{Var}[H_A]$  is plotted on a logarithmic scale.