

Statistical aspects of evolution under natural selection, with implications for the advantage of sexual reproduction[☆]

D. J. M. Crouch*

Department of Oncology, University of Oxford, Oxford, UK, OX3 7DQ

Abstract

The prevalence of sexual reproduction remains mysterious, as it poses clear evolutionary drawbacks compared to reproducing asexually. Several possible explanations exist, with one of the most likely being that finite population size causes linkage disequilibria to randomly generate and impede the progress of natural selection, and that these are eroded by recombination via sexual reproduction. Previous investigations have either analysed this phenomenon in detail for small numbers of loci, or performed population simulations for many loci. Here we present a quantitative genetic model for fitness, based on the Price Equation, in order to examine the theoretical consequences of randomly generated linkage disequilibria when there are many loci. In addition, most previous work has been concerned with the long-term consequences of deleterious linkage disequilibria for population fitness. The expected change in mean fitness between consecutive generations, a measure of short-term evolutionary success, is shown under random environmental influences to be related to the autocovariance in mean fitness between the generations, capturing the effects of stochastic forces such as genetic drift. Interaction between genetic drift and natural selection, due to randomly generated linkage disequilibria, is demonstrated to be one possible source of mean fitness autocovariance. This suggests

[☆]Preprint submitted to the Journal of Theoretical Biology

*Corresponding author

Email address: daniel.crouch@oncology.ox.ac.uk (D. J. M. Crouch)

a possible role for sexual reproduction in reducing the negative effects of genetic drift, thereby improving the short-term efficacy of natural selection.

Keywords: Sexual reproduction, Linkage disequilibrium, Hill-Robertson effect

1. Introduction

Sexual reproduction is by far the most prevalent mating system among the animals and plants, despite appearing to confer substantial disadvantages. For instance, when males contribute genetic material but no economic resources e.g. food or protection, individuals that clone themselves asexually ought to have 2^{n-1} as many descendants, after n generations, as those reproducing sexually. In the absence of any intrinsic benefit to reproducing sexually, a population of asexuals producing on average two offspring each will double in size per generation, while a sexually reproducing variety would stay at a constant size, as half of each female's resources are spent on male offspring that only reproduce by utilising female resources. This observation is referred to as the two-fold cost of sex (Maynard Smith, 1978). The effect is diminished when males invest economic resources in their offspring, but persists to some extent so long as they invest less than females, as is typically the case.

The preponderance of sexual reproduction is thus perplexing on the face of it, and various theories aspire to explain why it persists (Barton, 2010; Otto and Gerstein, 2006; Otto, 2009). Popular ideas are centred on the capacity of sexual reproduction to clear away the otherwise inexorable accumulation of deleterious mutations in finite-sized populations of asexual organisms (Muller, 1932), or similarly to combine favourable mutations within individuals more efficiently than under asexuality; referred to as the Fisher-Muller model (Muller, 1932; Fisher, 1930). When some degree of recombination is present, the Hill-Robertson effect (Hill and Robertson, 1966) operates via a related mechanism, whereby, under selection, loci that are linked are on average more susceptible to the effects genetic drift than unlinked loci, making recombination favourable. This is because

linkage disequilibrium (LD) that is negative between beneficial alleles, though just as likely to occur through genetic drift as positive LD, persists for longer by reducing the pace of natural selection. It has been argued that the Fisher-Muller and Hill-Robertson arguments are fundamentally equivalent, as both are consequences of finite population size (Felsenstein, 1974). In an effectively infinitely sized population there is still a possible advantage to recombination, as selection can generate negative LD between beneficial loci so long as epistasis is common, weak and producing fitnesses lower than those expected based on the marginal effects of alleles (Barton, 1995; Kondrashov, 1988) (known as negative epistasis). Another theory is based on interactions between species, and emphasises the role of host-parasite dynamics. When parasites evolve to exploit particular genotypes, the capacity of the host population to rapidly bring together rare or unique combinations of previously disparate alleles, via sexual recombination, may benefit its overall health (Hamilton, 1980). Similar arguments apply to random fluctuations in the abiotic environment (Charlesworth, 1976).

Although these are all productive theories, none are yet universally accepted (Otto, 2009). Empirical investigations mostly suggest that the preponderance of negative epistasis required to drive selection for sex is unlikely (de Visser and Elena, 2007), though this has been recently challenged (Sohail et al., 2017). The parasite avoidance theory is persuasive, but requires relatively strict constraints on the nature of host-parasite interactions (Otto and Gerstein, 2006; Otto, 2009; Iles et al., 2003), for instance that the parasites must have very strong selective effects on their hosts. The Fisher-Muller and Hill-Robertson mechanisms, predicated on certain consequences of finite population sizes, enjoy considerable theoretical support (Iles et al., 2003; Barton and Otto, 2005; Keightley and Otto, 2006), but do rely on genetic drift being a significant force.

Here, a quantitative genetic model with an arbitrary number of loci is analysed in order to examine the interference of selection by randomly generated LD (a consequence of genetic drift), when acting upon standing genetic variation.

This provides an understanding of how reduction of LD provides short-term gains in fitness to sexual populations. Previously, simulations of many loci have shown that recombination confers long-term advantages even for large populations, in which genetic drift is expected to be less powerful, provided that the number of loci under selection is sufficiently large (Iles et al., 2003). Detailed analysis of theoretical models has demonstrated that recombination is most effective at increasing mean fitness in small populations (Bodmer, 1970; Felsenstein, 1974), but has not considered cases where there are many loci under selection. It also remains unclear what precise mechanisms provide advantages to sex in the short term, i.e. the increase in mean fitness from one generation to the next, under this theory, as the majority of work focuses on long-term consequences of randomly generated LD.

The most widely accepted explanations for sex are centred on its capacity to break down deleterious LD (i.e. negative LD between beneficial alleles), increasing the additive genetic variance in fitness among individuals in the population. Fisher's Fundamental Theorem of Natural Selection (Fisher, 1930) states that this variance is proportional to the rate of increase in mean fitness due to natural selection. As part of this work, we derive an expanded version of the Fundamental Theorem. When there is genetic drift, the expected rate of increase in mean fitness is affected by randomly-induced correlations, more formally the autocovariance, between mean fitnesses of consecutive generations. The expanded theorem is then shown to encapsulate certain finite-population based advantages of sexual reproduction, as randomly generated LD induces autocovariance of mean fitness between generations by interfering systematically with natural selection. In contrast with the classic formulation of the Hill-Robertson effect, which describes the long-term detriment to fitness caused by the persistence of deleterious LD, this is shown to slow the rate of change of mean fitness in the short term, immediately after LD is generated.

2. The model

We assume that the fitness w_i for individual i is the linear combination of their genotype scores; $x_{ij} \in \{0, 1, 2\}$ for 0, 1 or 2 major alleles respectively (where j is a locus identifier and all loci are biallelic), and the *average effects* on fitness $\hat{\beta}_j$, plus an intercept $\hat{\beta}_0$, and a residual ϵ_i that captures variation in reproductive success that cannot be attributed to the genotypes:

$$w_i = \hat{\beta}_0 + \sum_{j=1}^M \hat{\beta}_j x_{ij} + \epsilon_i. \quad (1)$$

The hats on the average effects indicate that they are fitted least squares regression coefficients, treating all individuals in the population as the 'data'. There are M polymorphic loci in total. The breeding value for fitness, $\hat{g}_i = \hat{\beta}_0 + \sum_{j=1}^M \hat{\beta}_j x_{ij}$, represents its heritable component, and is the main subject of this analysis. The mean breeding value is written \bar{g} , and is equal to the mean fitness, $\bar{w} = \sum_i^N w_i/N$, where N is the number of individuals. This is because, due to properties of least squares regression, $\sum_{i=1}^N \epsilon_i = 0$. In a deterministic system, i.e. where all the variables in Equation 1 are known, the change in mean breeding value between parental and offspring generations can be written using the Price Equation (Price, 1970; Robertson, 1966; Frank, 1998) as

$$\Delta \bar{g} = \frac{\text{var}_i[\hat{g}_i] + E_i[w_i \Delta \hat{g}_i]}{\bar{w}}, \quad (2)$$

where \hat{g}_i represents the breeding value for a given individual and the subscript i on the variance and expectation operators ($\text{var}_i[\]$ and $E_i[\]$) indicates that these are taken with respect to the individuals within the population (i.e. summing the expression in parentheses over i and dividing by N). In the terminology of quantitative genetics, $\text{var}_i[\hat{g}_i]$ is the additive genetic variance in fitness. The mean difference in breeding value between a given parent and their offspring is represented by $\Delta \hat{g}_i$, and can be influenced both by environmental effects and by genetic differences between the parent and offspring that are due to sexual recombination. *De novo* mutations, though also a cause of

parent-offspring differences, are sufficiently rare to be unimportant, and are thus excluded from our analysis for simplicity. Equation 2 partitions the causes of breeding value evolution into two components: the first term represents the action of natural selection, and the second term the transition between parents and offspring (Frank, 2012). Extending to incorporate the stochastic effects of the environment, we can write:

$$E_e[\bar{w}\Delta\bar{g}] = E_e[\text{var}_i[\hat{g}_i]] + E_e[E_i[w_i\Delta\hat{g}_i]], \quad (3)$$

where the newly added expectations are taken with respect to the random environmental effects, represented by e . Here, random environmental factors are considered to be any that affect the distribution of genetic (x_{ij}) or non-genetic (ϵ_i) variables, or those that are affected by both e.g. the 'fitted' average effects $\hat{\beta}_j$. Genetic variables are considered random variables with respect to the environment, as random environmental forces can affect the distribution of genotypes in subsequent generations, e.g. through genetic drift. The subscript e therefore refers to quite general random phenomena and is intended mainly to differentiate from expectation, variance and covariance operations over i , which describe statistical relationships between the individuals within a given population. Equation 3 is a formulation similar to that of Grafen (Grafen, 2000), who refer to the random environmental effects as "states of nature". In the following section, the model will be rearranged so as to partition the causes of breeding value evolution into separate statistical processes.

3. Results

Individual variance in fitness

First, we separately derive expressions for the two terms on the right hand side of Equation 3 that permit their further analysis in terms of population genetic parameters. For a given environmental state initially, the first term, which is the between-individual variance in breeding values, can be expressed in terms of LD coefficients and average effects using the variance sum law:

$$\begin{aligned}
\text{var}_i[\hat{g}_i] &= \frac{1}{N} \sum_{i=1}^N \hat{g}_i^2 - \bar{g}^2 \\
&= \frac{1}{N} \sum_{i=1}^N \sum_{j=0}^M \sum_{k=0}^M \hat{\beta}_j \hat{\beta}_k x_{ij} x_{ik} - 4 \sum_{j=0}^M \sum_{k=0}^M \hat{\beta}_j \hat{\beta}_k p_j p_k \\
&= \sum_{j=0}^M \sum_{k=0}^M \hat{\beta}_j \hat{\beta}_k (E_i[x_{ij} x_{ik}] - 4p_j p_k) \\
&= 2 \sum_{j=0}^M \sum_{k=0}^M \hat{\beta}_j \hat{\beta}_k D_{jk}.
\end{aligned} \tag{4}$$

Recall that each x_{ij} is the number of major alleles present in a diploid individual, ranging from 0 to 2, so its expected value over i is $2p_j$ where p_j is the allele frequency. Summation is performed from $j = 0$ to account for the effect of the intercept term, $\hat{\beta}_0$, so $x_{i0} = 2$ and $p_0 = 1$. The coefficient of LD between loci j and k , D_{jk} , is defined as standard, and D_{jj} is used to denote the genotype variance of allele j ; $p_j(1 - p_j)$. The last line assumes that there is no inbreeding in the population, i.e. that the maternal and paternal chromosomes are uncorrelated. Therefore, where x_{ijA} and x_{ijB} are the genotypes on the maternal and paternal chromosomes of individual i , and $x_{ij} = x_{ijA} + x_{ijB}$, independence of the chromosomes implies that

$$\begin{aligned}
E_i[x_{ij} x_{ik}] - 4p_j p_k &= E_i[(x_{ijA} + x_{ijB})(x_{ikA} + x_{ikB})] - 4p_j p_k \\
&= E_i[x_{ijA} x_{ikA}] + E_i[x_{ijB} x_{ikB}] \\
&\quad + E_i[x_{ijA}] E_i[x_{ikB}] + E_i[x_{ikA}] E_i[x_{ijB}] - 4p_j p_k \\
&= E_i[x_{ijA} x_{ikA}] + E_i[x_{ijB} x_{ikB}] - 2p_j p_k \\
&= 2D_{jk}.
\end{aligned} \tag{5}$$

This also holds in diploid asexuals, as under the additive fitness model each organism can be counted as two combined independent haploid genomes. Then, taking the expectation over random environmental effects, and assuming that there are no gene-environment correlations (Appendix A),

$$\begin{aligned}
E_e[\text{var}_i[\hat{g}_i]] &= 2 \sum_{j=0}^M \sum_{k=0}^M E_e[\hat{\beta}_j \hat{\beta}_k D_{jk}] \\
&= 2 \sum_{j=0}^M \sum_{k=0}^M \text{cov}_e[\hat{\beta}_j, \hat{\beta}_k D_{jk}] \\
&\quad + 2 \sum_{j=0}^M \sum_{k=0}^M \beta_j \beta_k E_e[D_{jk}],
\end{aligned} \tag{6}$$

where $\beta_j \equiv E_e[\hat{\beta}_j]$, which is the average effect of the major allele at locus j that one would find in a population of infinite size, and expect in a randomly chosen population. This is assumed to be equal for all generations, i.e. $\beta_j = \beta'_j$. Appendix A shows that $E_e[\hat{\beta}_k D_{jk}] = \beta_k E_e[D_{jk}]$ which is necessary to produce the final term.

Change in average effects between generations

Turning to the second term on the right hand side of Equation 3, which describes on average how offspring breeding values differ from those of their parents; in an asexual population

$$\begin{aligned}
E_i[w_i \Delta \hat{g}_i] &= \frac{1}{N} \sum_{i=1}^N \sum_{j=0}^M w_i (\hat{\beta}'_j - \hat{\beta}_j) x_{ij} \\
&= \sum_{j=0}^M \sum_{k=0}^M \hat{\beta}_k (\hat{\beta}'_j - \hat{\beta}_j) E_i[x_{ij} x_{ik}],
\end{aligned} \tag{7}$$

as $\Delta \hat{g}_i = \sum_{j=0}^M (\hat{\beta}'_j x'_{ij} - \hat{\beta}_j x_{ij})$, and $x'_{ij} = x_{ij}$ as we have chosen to ignore *de novo* mutations. Prime symbols here and throughout refer to variables in the subsequent (i.e. offspring) generation, and x'_{ij} is the mean genotype among the offspring produced by individual i at locus j . Appendix B demonstrates that Equation 7 generalises to sexual organisms under the additional assumptions of fair meiosis and absence of selection acting on gametes. To go from the first to the second line, w_i can be replaced by \hat{g}_i , as residuals and predictor variables do

not correlate (i.e. $E_i[\epsilon_i x_{ij}] = E_i[\epsilon_i]E_i[x_{ij}] = 0$, so $E_i[(\hat{g}_i + \epsilon_i)x_{ij}] = E_i[\hat{g}_i x_{ij}]$) due to properties of linear regression. Then, averaging over the environmental effects as in Equation 6,

$$\begin{aligned}
E_e[E_i[w_i \Delta \hat{g}_i]] &= \sum_{j=0}^M \sum_{k=0}^M E_e[\hat{\beta}_k (\hat{\beta}'_j - \hat{\beta}_j) E_i[x_{ij} x_{ik}]] \\
&= \sum_{j=0}^M \sum_{k=0}^M \text{cov}_e[\hat{\beta}'_j, \hat{\beta}_k E_i[x_{ij} x_{ik}]] \\
&\quad - \sum_{j=0}^M \sum_{k=0}^M \text{cov}_e[\hat{\beta}_j, \hat{\beta}_k E_i[x_{ij} x_{ik}]],
\end{aligned} \tag{8}$$

as we assume that the expected effect of an allele is the same regardless of the generation, implying $E_e[\hat{\beta}'_j]E_e[\hat{\beta}_j E_i[x_{ij} x_{ik}]] = E_e[\hat{\beta}_j]E_e[\hat{\beta}_j E_i[x_{ij} x_{ik}]]$, so

$$\begin{aligned}
&E_e[\hat{\beta}'_j \hat{\beta}_k E_i[x_{ij} x_{ik}]] - E_e[\hat{\beta}_j \hat{\beta}_k E_i[x_{ij} x_{ik}]] \\
&= \text{cov}_e[\hat{\beta}'_j, \hat{\beta}_k E_i[x_{ij} x_{ik}]] - \text{cov}_e[\hat{\beta}_j, \hat{\beta}_k E_i[x_{ij} x_{ik}]].
\end{aligned} \tag{9}$$

Assuming no gene-environment correlations, the covariance involving average effects in each generation, $\text{cov}_e[\hat{\beta}'_j, \hat{\beta}_k E_i[x_{ij} x_{ik}]]$, is zero (Appendix C). Then, we find that

$$E_e[E_i[w_i \Delta \hat{g}_i]] = - \sum_{j=0}^M \sum_{k=0}^M \text{cov}_e[\hat{\beta}_j, \hat{\beta}_k E_i[x_{ij} x_{ik}]], \tag{10}$$

meaning that if the effects of finite population size on the estimation of the average effects is non-negligible, the breeding value of offspring will on average decrease relative to their parents. This is because, although errors in offspring breeding values are expected to be distributed symmetrically around zero, parents with high breeding value errors will have disproportionately many offspring.

Substituting Equations 6 and 10 into Equation 3 to give the expectation of the mean breeding value change multiplied by mean fitness;

$$\begin{aligned}
E_e[\bar{w}\Delta\bar{g}] &= 2 \sum_{j=0}^M \sum_{k=0}^M \beta_j \beta_k E_e[D_{jk}] \\
&\quad + \sum_{j=0}^M \sum_{k=0}^M cov_e[\hat{\beta}_j, \hat{\beta}_k(2D_{jk} - E_i[x_{ij}x_{ik}])] \\
&= 2 \sum_{j=0}^M \sum_{k=0}^M \beta_j \beta_k E_e[D_{jk}] \\
&\quad - 4 \sum_{j=0}^M \sum_{k=0}^M cov_e[\hat{\beta}_j, \hat{\beta}_k p_j p_k]. \\
&= E_e[var_i[g_i]] - 4 \sum_{j=0}^M \sum_{k=0}^M cov_e[\hat{\beta}_j, \hat{\beta}_k p_j p_k],
\end{aligned} \tag{11}$$

where g_i is the 'true' breeding value for a given individual which, in contrast with \hat{g}_i , is a function of the expected, not fitted, average effects.

Partitioning the stochastic and deterministic causes of breeding value evolution

The terms on the right hand side of Equation 11 can now be split into separate components representing processes that are either affected or unaffected by population size (i.e. deterministic or stochastic). The first term, $E_e[var_i[g_i]]$, is

$$2 \sum_{j=0}^M \sum_{k=0}^M \beta_j \beta_k (\theta_{jk} - 2cov_e[p_j, p_k]), \tag{12}$$

where θ_{jk} is the LD coefficient between loci j and k obtained if the population was infinitely sized. In statistical terminology this is the LD parameter that D_{jk} acts to estimate, and approaches as the population size increases towards infinity. Specifically,

$$\begin{aligned}
\theta_{jk} &= \lim_{N \rightarrow \infty} E_e[D_{jk}] = E_e[X_j X_k]/2 - E_e[X_j]E_e[X_k]/2 \\
&= E_e[E_i[x_{ij}x_{ik}]]/2 - 2E_e[p_j]E_e[p_k] \\
&= E_e[E_i[x_{ij}x_{ik}]]/2 - 2E_e[p_j p_k] + 2cov_e[p_j, p_k] \\
&= E_e[D_{jk}] + 2cov_e[p_j, p_k],
\end{aligned} \tag{13}$$

where X_j and X_k are random genotype variables for an arbitrary individual. The equality in the first line is due to D_{jk} being a *consistent* estimator of half the covariance between the genotypes X_j and X_k , represented on the right hand side. Under infinite population size the allele frequencies p_j and p_k will not covary due to random effects; therefore what would be $2E_e[p_j p_k]$ in the expression for the LD coefficient in a finite size population is equal to $2E_e[p_j]E_e[p_k]$. This produces Expression 12, which separates the causes of evolutionary change into two major components. The idea is made clearer by rewriting the first term on the right hand side of Equation 11 as

$$E_e[var_i[g_i]] = var_i[g_i] - 4 \sum_{j=0}^M \sum_{k=0}^M \beta_j \beta_k cov_e[p_j, p_k], \tag{14}$$

where $var_i[g_i]$ is the variance in breeding values that would be expected if the number of individuals was infinite - i.e. under deterministic forces alone:

$$var_i[g_i] = 2 \sum_{j=0}^M \sum_{k=0}^M \beta_j \beta_k \theta_{jk}. \tag{15}$$

Substituting Equation 14 for the first term in Equation 11;

$$E_e[\bar{w}\Delta\bar{g}] = var_i[g_i] - 4 \sum_{j=0}^M \sum_{k=0}^M \beta_j \beta_k cov_e[p_j, p_k] - 4 \sum_{j=0}^M \sum_{k=0}^M cov_e[\hat{\beta}_j, \hat{\beta}_k p_j p_k]. \tag{16}$$

We are primarily interested in the expected change in mean breeding value for fitness, $E_e[\Delta\bar{g}]$, rather than $E_e[\bar{w}\Delta\bar{g}]$. This can be obtained by adding $var_e[\bar{g}] - cov_e[\bar{g}, \bar{g}']$ to both sides of Equation 16, as

$$\begin{aligned}
E_e[\bar{w}\Delta\bar{g}] &= E_e[\bar{g}\bar{g}'] - E_e[\bar{g}^2] \\
&= \left(E_e[\bar{g}]E_e[\bar{g}'] + cov_e[\bar{g}, \bar{g}'] \right) - \left(E_e[\bar{g}]^2 + var_e[\bar{g}] \right) \quad (17) \\
&= E_e[\bar{g}]E_e[\Delta\bar{g}] + cov_e[\bar{g}, \bar{g}'] - var_e[\bar{g}],
\end{aligned}$$

recalling that $\bar{g} = \bar{w}$, so $E_e[\bar{g}]E_e[\Delta\bar{g}] = E_e[\bar{w}\Delta\bar{g}] + var_e[\bar{g}] - cov_e[\bar{g}, \bar{g}']$.

Because $var_e[\bar{g}]$ equals

$$\begin{aligned}
&4 \sum_{j=0}^M \sum_{k=0}^M E_e[\hat{\beta}_j \hat{\beta}_k p_j p_k] - 4 \sum_{j=0}^M \sum_{k=0}^M \beta_j \beta_k E_e[p_j] E_e[p_k] \\
&= 4 \sum_{j=0}^M \sum_{k=0}^M E_e[\hat{\beta}_j \hat{\beta}_k p_j p_k] - 4 \sum_{j=0}^M \sum_{k=0}^M \beta_j \beta_k \left(E_e[p_j p_k] - cov_e[p_j, p_k] \right) \\
&= 4 \sum_{j=0}^M \sum_{k=0}^M \beta_j \beta_k cov_e[p_j, p_k] + 4 \sum_{j=0}^M \sum_{k=0}^M cov_e[\hat{\beta}_j, \hat{\beta}_k p_j p_k]
\end{aligned} \quad (18)$$

(where the factors of 4 are due to each mean breeding value being twice the linear combinations of allele frequencies and average effects), adding $var_e[\bar{g}] - cov_e[\bar{g}, \bar{g}']$ to both sides of Equation 16 cancels the second two terms on the right hand side, and transforms $E_e[\bar{w}\Delta\bar{g}]$ to $E_e[\bar{g}]E_e[\Delta\bar{g}]$ on the left hand side, leaving

$$\boxed{E_e[\bar{g}]E_e[\Delta\bar{g}] = \underset{N \rightarrow \infty}{var_i[g_i]} - cov_e[\bar{g}, \bar{g}']}, \quad (19)$$

which resembles Fisher's Fundamental Theorem of Natural Selection (Fisher, 1930; Price, 1972; Edwards, 2014; Frank, 1997), with the inclusion of an additional term showing how positive autocovariance of the mean breeding value in consecutive generations reduces its expected rate of increase.

Interaction of stochastic and deterministic effects through linkage disequilibrium

We now use the result above to analyse Hill-Robertson type effects on the expected change in mean fitness. By making the simplifying assumption that

the average effects are close to their expected values,

$$\text{cov}_e[\bar{g}, \bar{g}'] \approx 4 \sum_{j=0}^M \sum_{k=0}^M \beta_j \beta_k \text{cov}_e[p_j, p'_k], \quad (20)$$

where the factor of 4 is due to each mean breeding value being twice the linear combination of average effects and allele frequencies. The covariance between p_j and p'_k can be understood by applying the law of total covariance. Let D_{jk}^* be the coefficients of LD and p_j^* and p_k^* be the allele frequencies in an initial state for loci j and k , which are random variables produced by mutation or genetic drift. Assume that the initial state, represented by the asterisk, is in the previous generation, which is two generations prior to the offspring generation (denoted by prime). The allele frequencies evolve in the order $p_j^* \rightarrow p_j \rightarrow p'_j$, and the coefficients of LD likewise. We assume that we are able to specify a probability model for how the allele frequencies and linkage disequilibria are generated in the initial state generation. Then, the law of total covariance implies:

$$\text{cov}_e[p_j, p'_k] = \text{cov}_e[E_e[p_j|D_{jk}^*], E_e[p'_k|D_{jk}^*]] + E_e[\text{cov}_e[p_j, p'_k|D_{jk}^*]]. \quad (21)$$

The second term describes how genetic drift is affected by randomly produced LD. The covariance between p_j and p'_k due to genetic drift under the Wright-Fisher model, for a given fluctuation in LD described by D_{jk}^* , is

$$\begin{aligned} \text{cov}_e[p_j, p'_k|D_{jk}^*] &= \text{cov}_e[p_j, p_k|D_{jk}^*] + \text{cov}_e[p_j, \Delta p_k|D_{jk}^*] \\ &= \text{cov}_e[p_j, p_k|D_{jk}^*] + \text{cov}_e\left[p_j, \sum_{l=0}^M \frac{\hat{\beta}_{kl} D_{kl}}{\bar{g}} \middle| D_{jk}^*\right] \\ &= \text{cov}_e[p_j, p_k|D_{jk}^*] + \text{cov}_e\left[p_j, \frac{\hat{\beta}_{jk} D_{jk}}{\bar{g}} \middle| D_{jk}^*\right] \quad (22) \\ &\approx \text{cov}_e[p_j, p_k|D_{jk}^*] \\ &= \frac{D_{jk}^*}{2N}, \end{aligned}$$

the expectation of which is zero when $j \neq k$, as positive LD is assumed to be equally likely to occur as negative LD, and $E_e[p_j^*(1 - p_j^*)]/2N$ otherwise. The

above assumes that D_{kl} is independent from p_j when $l \neq j$, and the approximation in the fourth line is justified when the average effects are small due to there being many loci influencing fitness.

The first term on the right hand side of Equation 21 describes how drift and selection interact. The expectation of the allele frequency p_j under selection is $E_e[p_j^* + \sum_{l=0}^M \beta_l D_{jl}^* / \bar{g}^*]$. Conditioned on a randomly induced state of LD, D_{jk}^* , and assuming that the pairwise LD coefficients are independent from each other within the generation and from \bar{g}^* and p_j^* (which involves some approximation), the expectation is

$$\begin{aligned} E_e[p_j | D_{jk}^*] &= I(j = k)p_j^* + \beta_k \frac{D_{jk}^*}{E_e[\bar{g}^*]} + C \\ &\approx I(j = k)p_j^* + I(j \neq k) \left(\beta_k \frac{D_{jk}^*}{E_e[\bar{g}^*]} \right) + C, \end{aligned} \quad (23)$$

where C is a constant and the indicator function $I(j = k)$ equals 1 if the condition in parentheses is fulfilled or 0 otherwise. Henceforth C denotes any constant and does not refer to a particular variable. Here, C is an expectation that does not depend on the value of D_{jk}^* (or any other variable), and it can thus be ignored, as we later require its covariances with other variables, which must be zero. The indicator function is used because p_j^* is determined by $D_{jj}^* = p_j^*(1 - p_j^*)$, so is necessary to retain as a random variable when $j = k$. It is assumed not to depend on the LD coefficient when $j \neq k$ and becomes incorporated into the constant, as LD is equally likely to be positive as it is negative. The approximation in the second line is justified because p_j^* is expected to be much larger than the other term. In a similar way, the conditional expectation of p'_k is

$$\begin{aligned}
E_e[p'_k|D_{jk}^*] &= I(j=k)p_k^* + \beta_j D_{jk}^*/E_e[\bar{g}^*] + \beta_j \frac{E_e[D_{jk}|D_{jk}^*]}{E_e[\bar{g}]} + C \\
&\approx I(j=k)p_j^* + I(j \neq k) \left(\beta_j D_{jk}^*/E_e[\bar{g}^*] + \beta_j \frac{E_e[D_{jk}|D_{jk}^*]}{E_e[\bar{g}]} \right) + C,
\end{aligned} \tag{24}$$

where;

$$\begin{aligned}
\beta_j \frac{E_e[D_{jk}|D_{jk}^*]}{E_e[\bar{g}]} &= \frac{\beta_j}{E_e[\bar{g}]} \left(D_{jk}^* + (\beta_j + \beta_k I(j \neq k)) \frac{E_e[p_{A_j A_k}^* (1 - p_{A_j A_k}^*) | D_{jk}^*]}{E_e[\bar{g}^*]} \right. \\
&\quad \left. - \beta_j E_e[p_j^*] \frac{D_{jk}^*}{E_e[\bar{g}^*]} - \beta_k E_e[p_k^*] \frac{D_{jk}^*}{E_e[\bar{g}^*]} - \beta_j \beta_k \frac{D_{jk}^{*2}}{E_e[\bar{g}^{*2}]} + C \right),
\end{aligned} \tag{25}$$

and $p_{A_j A_k}$ is the frequency of the haplotype with major alleles ('A') at both loci. If the number of loci under selection is large, average effects of alleles at any one locus will be small. Terms containing products of average effects are then very small relative to those featuring no such product, so this simplifies to

$$\beta_j \frac{E_e[D_{jk}|D_{jk}^*]}{E_e[\bar{g}]} \approx \frac{\beta_j}{E_e[\bar{g}]} D_{jk}^*. \tag{26}$$

This states that the change in LD structure due to selection has negligible consequences for the separate allele frequencies in the subsequent generation. When recombination is present at rate r_{jk} ,

$$\beta_j \frac{E_e[D_{jk}|D_{jk}^*]}{E_e[\bar{g}]} \approx \frac{\beta_j}{E_e[\bar{g}]} D_{jk}^* (1 - r_{jk}). \tag{27}$$

Taking the covariances of non-constants in Equations 23 and 24, the first term in the Equation 21 is

$$I(j=k) \text{var}_e[p_j^*] + I(j \neq k) \beta_j \beta_k \text{var}_e[D_{jk}^*] \frac{E_e[\bar{g}^*](1 - r_{jk}) + E_e[\bar{g}]}{E_e[\bar{g}^*]^2 E_e[\bar{g}]}. \tag{28}$$

The second term of this describes the phenomena that are relevant to the advantage of recombination. When the major alleles at loci j and k are both beneficial and potentially in LD, variance in the level of LD causes the covariance in their allele frequencies to be positive between the parental and offspring generations, as variances cannot be negative, and the signs of β_j and β_k are identical. The same is true when the major alleles are both deleterious. When the alleles have opposite effects, their frequencies are expected to fluctuate in opposite directions. In essence, this is a short-term version of the Hill-Robertson effect. When a deleterious allele hitchhikes alongside a beneficial allele, the former proceeds to a higher than expected frequency and the latter to lower than expected. The reverse is true when they appear disproportionately on separate chromosomes; the negative allele tends to be accompanied by another deleterious allele so goes to lower than expected frequency, while the beneficial allele frequency increases more than expected. As shown by the approximation in Equation 26, this LD structure persists across the parental and offspring generations; inducing the negative covariance between the deleterious frequency in one generation and the beneficial frequency in the other. An equivalent process occurs when LD forms between beneficial alleles, though their frequencies will positively rather than negatively covary. As alleles with opposite effects on fitness have negative covariance in frequency, positive covariance is induced between the mean fitness breeding values (due to differing signs on the effects on fitness), and alleles with the same direction of effect, being positively covarying, also create positive mean breeding value covariances. This is shown by substituting Equation 22 and Expression 28 into Equation 20 (via Equation 21):

$$\begin{aligned}
cov_e[\bar{g}, \bar{g}'] \approx & 4 \sum_{j=0}^M \beta_j^2 \left(var_e[p_j^*] + \frac{E_e[D_{jj}^*]}{2N} \right) \\
& + 4 \sum_{j=0}^M \sum_{k \neq j}^M \beta_j^2 \beta_k^2 var_e[D_{jk}^*] \frac{E_e[\bar{g}^*](1 - r_{jk}) + E_e[\bar{g}]}{E_e[\bar{g}^*]^2 E_e[\bar{g}]} .
\end{aligned} \tag{29}$$

As all of these terms are positive, the short-term consequence of Hill-Robertson

type interference is to increase the autocovariance of mean breeding values between parent and offspring generations. As shown in Equation 19, breeding value autocovariance induced through genetic drift lowers the efficacy of selection, relative to expectation in an equivalent population of infinite size, or one where there were no random environmental influences on fitness. Breeding value autocovariance due to LD between loci can be controlled by increasing the recombination rate, r_{jk} , which is evident in Equation 29. It should be noted that, in a realistic model, $var_e[D_{jk}^*]$ is also affected by recombination. Using the formula for the variance of a product of independent variables it can be represented as $var_e[D_{jk}^\bullet]((1 - r_{jk})^2 + var_e[1 - \hat{r}_{jk}])$ where D_{jk}^\bullet is the LD coefficient that is randomly generated by mutation, drift, or both together, but before recombination produces the gametes entering the initial state (*) generation, and \hat{r}_{jk} is the realised recombination rate (i.e. observed fraction of recombinations between loci j and k). The component of autocovariance due to the effects of genetic drift at individual loci (first summation in Equation 29) is unaffected by recombination. This appears large relative to the LD effects in the second term, but it is important to consider that there are only M items to sum over, and $M(M - 1)$ in the second summation, accounting for all pairwise relationships between loci. The relative importance of these effects depends on the population size, as $var_e[D_{jk}^*]$ is governed by fluctuations in LD caused by genetic drift or mutation which may be negligible in very large populations, and the number of loci, as when M becomes large, $M(M - 1)$ becomes extremely large.

4. Discussion

An influential class of theories for explaining the existence of sexual reproduction is based on the isolation of beneficial alleles on separate genomes due to populations being finite-sized (Muller, 1932; Hill and Robertson, 1966; Fisher, 1930; Felsenstein, 1974). Negative LD of this form reduces between-individual variance in fitness, as it implies a paucity of individuals carrying large numbers of beneficial (or deleterious) mutations, which reduces the pace of selective change in mean fitness via Fisher’s Fundamental Theorem of Natural Selection

(Fisher, 1930). This is also a feature of the species interaction (Hamilton, 1980) and negative epistasis (Barton, 1995) theories, though these propose how negative LD can be generated in populations of infinite size. Negative LD is likely to accumulate in finite populations due to random mutation and genetic drift interfering with selection. Hill-Robertson interference (Hill and Robertson, 1966) is the tendency for genetic drift, though equally likely to generate positive or negative LD, to yield a preponderance of negative LD in the long term, as negative LD persists for longer than positive LD under natural selection. This process applies equally to LD generated through mutation, which is especially likely to cause negative LD due to beneficial (Fisher, 1930) or deleterious (Muller, 1932) mutations rarely accumulating together in the same individuals.

The work presented here suggests that, starting from a position of uncertainty about whether it is positive or negative, LD generated through a stochastic process such as genetic drift reduces short-term expected changes in mean fitness. This addresses a separate question to the classic formulation of the Hill-Robertson effect, which describes how negative LD accumulates in the long term. We find that, in the short term, haplotype structures that are inherited by a parental generation will be passed on mostly unaltered to the offspring generation when there is no recombination, whether the LD is positive or negative (shown by the approximation in Equation 26). When LD is then randomly positive or negative between alleles with given effects on fitness (with the expectation that it is neutral), this induces covariance between the parental and offspring generations' allele frequencies at different loci. For example, when an initial generation randomly receives a set of haplotypes in which there is a large amount of positive LD between two beneficial alleles, their frequencies will tend in the same direction (upwards) across multiple generations (provided LD is not rapidly broken down by recombination), and the reverse is true when the LD is negative, creating positive intergenerational covariance in both cases. Equivalently, when LD is positive between alleles with opposite effects, the deleterious allele will do better than expected and the beneficial allele worse (and vice versa

with negative LD). Again, these effects will largely persist across generations, as the LD remains mostly unaltered unless recombination is present, so covariance between generations again results, though in this case it is negative. Both kinds of allele frequency covariance (and standard genetic drift that affects a single locus) result in positive covariance in fitness between generations, which would be referred to in the spatiotemporal statistics literature as autocovariance. Equation 19 shows that positive autocovariance reduces the expected change in mean fitness, i.e. it is detrimental to the short-term evolutionary prospects of the population. The intuition behind this is that genetic drift both counteracts the effects selection of selection on mean fitness, and causes positive autocovariance in traits between generations. In terms of its effects on autocovariance, randomly generated LD, when interfering with selection, is similar to genetic drift occurring at a single locus, as shown in Equation 29. It can therefore be seen as contributing to the total negative effect of genetic drift thereby harming the efficacy of selection, but with the added possibility of being cleared by recombination.

The assumptions required for the foregoing analysis are that there are no gene-environment correlations, no environmental correlations between parents and offspring, and that all alleles confer their effects on fitness additively. The analysis of short-term Hill-Robertson type effects also assumes that the average effects remain approximately constant. Further work could investigate the sensitivity of the results to these assumptions. The biological phenomena affecting autocovariance in mean fitness might also be explored, and these could incorporate more complex stochastic processes such as species interactions. We have not investigated the consequences of this work for the evolution of a sexually reproducing group of organisms alongside a competing asexual strain. This would be necessary to establish whether any benefits to recombination due the mechanism proposed can outweigh the two-fold cost of sex, and whether alleles for recombination were thus able to invade an asexual population, or to resist invasion themselves.

Conflict of interest

The author declares no conflict of interest

Acknowledgements

We thank Walter Bodmer for his advice and comments on the manuscript. We also thank three anonymous reviewers.

References

- Barton, N. H. (1995). A general model for the evolution of recombination. *Genet Res*, 65(2):123–45.
- Barton, N. H. (2010). Mutation and the evolution of recombination. *Philos Trans R Soc Lond B Biol Sci*, 365(1544):1281–94.
- Barton, N. H. and Otto, S. P. (2005). Evolution of recombination due to random drift. *Genetics*, 169(4):2353–70.
- Bodmer, W. F. (1970). The evolutionary significance of recombination in prokaryotes. *Symp Soc Gen Microbiol*, 20:279–294.
- Charlesworth, B. (1976). Recombination modification in a fluctuating environment. *Genetics*, 83(1):181–195.
- de Visser, J. A. G. M. and Elena, S. F. (2007). The evolution of sex: empirical insights into the roles of epistasis and drift. *Nat Rev Genet*, 8(2):139–49.
- Edwards, A. W. F. (2014). R.a. fisher’s gene-centred view of evolution and the fundamental theorem of natural selection. *Biol Rev Camb Philos Soc*, 89(1):135–47.

- Felsenstein, J. (1974). The evolutionary advantage of recombination. *Genetics*, 78(2):737–56.
- Fisher, R. A. (1930). *The genetical theory of natural selection*. The Clarendon press, Oxford.
- Frank, S. A. (1997). The price equation, fisher’s fundamental theorem, kin selection, and causal analysis. *Evolution*, 51(6):1712–1729.
- Frank, S. A. (1998). *Foundations of social evolution*. Princeton University Press, Princeton, New Jersey.
- Frank, S. A. (2012). Natural selection. iv. the price equation. *J Evol Biol*, 25(6):1002–19.
- Grafen, A. (2000). Developments of the price equation and natural selection under uncertainty. *Proc Biol Sci*, 267(1449):1223–7.
- Hamilton, W. D. (1980). Sex versus non-sex versus parasite. *Oikos*, 35(2):282–290.
- Hill, W. G. and Robertson, A. (1966). The effect of linkage on limits to artificial selection. *Genet Res*, 8:269–294.
- Iles, M. M., Walters, K., and Cannings, C. (2003). Recombination can evolve in large finite populations given selection on sufficient loci. *Genetics*, 165(4):2249–58.
- Keightley, P. D. and Otto, S. P. (2006). Interference among deleterious mutations favours sex and recombination in finite populations. *Nature*, 443(7107):89–92.
- Kondrashov, A. S. (1988). Deleterious mutations and the evolution of sexual reproduction. *Nature*, 336(6198):435–40.
- Maynard Smith, J. (1978). *The evolution of sex*. Cambridge University Press.

- Muller, H. J. (1932). Some genetic aspects of sex. *The American Naturalist*, 66(703):118–138.
- Otto, S. P. (2009). The evolutionary enigma of sex. *Am Nat*, 174 Suppl 1:S1–S14.
- Otto, S. P. and Gerstein, A. C. (2006). Why have sex? the population genetics of sex and recombination. *Biochem Soc Trans*, 34(Pt 4):519–22.
- Price, G. R. (1970). Selection and covariance. *Nature*, 227(5257):520–521.
- Price, G. R. (1972). Fisher’s ‘fundamental theorem’ made clear. *Ann Hum Genet*, 36(2):129–40.
- Robertson, A. (1966). A mathematical model of the culling process in dairy cattle. *Anim Prod*, 8:95–108.
- Sohail, M., Vakhrusheva, O. A., Sul, J. H., Pulit, S. L., Francioli, L. C., Genome of the Netherlands Consortium, Alzheimer’s Disease Neuroimaging Initiative, van den Berg, L. H., Veldink, J. H., de Bakker, P. I. W., Bazykin, G. A., Kondrashov, A. S., and Sunyaev, S. R. (2017). Negative selection in humans and fruit flies involves synergistic epistasis. *Science*, 356(6337):539–542.

Appendix

A: Average effects have zero covariance with genetic variables

We assume that the true underlying fitness function is of the form

$$w_i = \beta_0 + \sum_{j=1}^M \beta_j x_{ij} + \delta_i, \quad (\text{A1})$$

where the δ_i are environmental effects on fitness, which are independent and identically distributed among the different individuals and in different generations, with an expected value of zero, analogous to the effects of genetic drift. Using least squares regression theory, the vector of covariances between a scalar

genetic variable, referred to here as G , and the vector of average effects $\hat{\beta}$ (length M) is

$$\begin{aligned} E_e[G(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T w] - E_e[G] \beta \\ = E_e[G(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T (\mathbf{X} \beta + \delta)] - E_e[G] \beta \end{aligned} \quad (\text{A2})$$

where \mathbf{X} is the $N \times M$ matrix of minor-allele genotype scores (0, 1 or 2), w is the vector of fitnesses (length N), β the vector of expected average effects (length M) and δ the vector of environmental effects (length N). As $(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{X} \beta = \beta$, this is

$$E_e[G \beta] + E_e[G(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \delta] - E_e[G] \beta, \quad (\text{A3})$$

where β is a constant that can be moved outside of expectations, so the first and last terms cancel. When there are no correlations between environmental and genetic variables (including with complicated products and functions of genetic variables e.g. $G(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T$), δ can be moved into a separate expectation and $E_e[\delta] = 0$ by definition of the model, so the second term is zero, and hence genetic scalars such as D_{jk} , $E_i[x_{ij} x_{ik}]$ or $p_j p_k$ do not covary with $\hat{\beta}_l \forall \{j, k, l\}$.

B: Application to meiotic organisms

We here show that, under quite general conditions, $E_e[E_i[w_i \Delta \hat{g}_i]] = \sum_j^M E_e[(\hat{\beta}'_j - \hat{\beta}_j) E_i[w_i x_{ij}]]$ regardless of whether the population in question is asexually or meiotically reproducing, thus extending the results in the main text to a broad range of organisms. To do so we assume that meiosis is fair, i.e. that the probability of being transmitted to a gamete is equal for all alleles, and that there is no natural selection acting on gametes. These two assumptions are jointly referred to as *representativeness* (Grafen, 2000), as they imply that successful gametes will be statistically representative of the parental genome. First, note that

$$E_e[E_i[w_i\Delta\hat{g}_i]] = \sum_j^M E_e[E_i[w_i(\hat{\beta}'_j x'_{ij} - \hat{\beta}_j x_{ij})]], \quad (\text{B1})$$

where x'_{ij} is equal to $\sum_l^{w_i} x_{ijl}/w_i$, the mean genotype among the successful gametes of individual i at locus j , identified by l . In the asexual case, $x'_{ij} = x_{ij}$ (as we have chosen to ignore *de novo* mutations) but this is not necessarily true of meiotically reproducing organisms. In Equation B1, $E_e[E_i[w_i(\hat{\beta}'_j x'_{ij} - \hat{\beta}_j x_{ij})]]$ is equal to $E_e[E_i[(\sum_l^{w_i} \hat{\beta}'_j x_{ijl} - \hat{\beta}_j x_{ij})]]$, which is also a series of $\frac{1}{2}N\bar{w}$ terms of the form

$$2E_e[\hat{\beta}'_j x_{Aj}] - E_e[\hat{\beta}_j x_{Aj}] - E_e[\hat{\beta}_j x_{Bj}], \quad (\text{B2})$$

where A is the father and B the mother of the gamete l

To find the desired result then, it is sufficient to show for any i and j that

$$E_e[w_i\hat{\beta}'_j x'_{ij}] = E_e[w_i\hat{\beta}_j x_{ij}], \quad (\text{B3})$$

equivalent to $E_e[w_i\hat{\beta}'_j \Delta x_{ij}] = 0$. We now follow a similar procedure to Grafen (Grafen, 2000), partitioning the environmental effects, e , into two types depending on whether or not they affect random events at meiosis and fertilisation e.g. recombination, independent assortment of parental chromosomes and competition among spermatozoa. Representing these as τ and all other random events as ζ , $E_e[w_i\hat{\beta}'_j \Delta x_{ij}]$ can be rewritten $E_\zeta[w_i E_\tau[\hat{\beta}'_j \Delta x_{ij}|\zeta]]$. Note that w_i has been placed outside of the inner expectation as it is not affected by events at meiosis or fertilisation.

Representativeness ensures that $E_\tau[x'_{ij}|\zeta] = x_{ij}|\zeta$, so $E_\tau[\hat{\beta}'_j \Delta x_{ij}|\zeta]$ is equal to $\text{cov}_\tau[\hat{\beta}'_j, x'_{ij}|\zeta]$. Then, assuming that environmental effects in the parental and offspring generations are independent, the events unrelated to meiosis or fertilisation, ζ , are split into ζ_1 and ζ_2 representing the first and second generation (corresponding with non-primed and primed variables). Under reasonable

assumptions, the order of expectations can be interchanged, and the genotypes x'_{ij} must be independent from the environmental effects in that generation, so

$$E_{\zeta}[w_i \text{cov}_{\tau}[\hat{\beta}'_j, x'_{ij}|\zeta]] = E_{\zeta_1}[w_i \text{cov}_{\tau}[E_{\zeta_2}[\hat{\beta}'_j|\zeta_1, \tau], x'_{ij}|\zeta_1]]. \quad (\text{B4})$$

When there are no gene-environment correlations, $E_{\zeta_2}[\hat{\beta}'_j|\zeta_1, \tau] = \beta_j$ which is a constant independent from x'_{ij} , giving $\text{cov}_{\tau}[E_{\zeta_2}[\hat{\beta}'_j|\zeta_1, \tau], x'_{ij}|\zeta_1] = 0$ and $E_e[w_i \hat{\beta}'_j \Delta x_{ij}] = 0$, as desired.

C: Average effects in parental and offspring generations have zero covariance

Using a standard result from least-squares regression theory, the covariance matrix between the vectors of average effects in parental and offspring generations, $\hat{\beta}$ and $\hat{\beta}'$ is

$$\begin{aligned} & E_e[((\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T w)((\mathbf{X}'^T \mathbf{X}')^{-1} \mathbf{X}'^T w')^T] - \beta \beta^T \\ &= E_e[((\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T (\mathbf{X} \beta + \delta))((\mathbf{X}'^T \mathbf{X}')^{-1} \mathbf{X}'^T (\mathbf{X}' \beta' + \delta'))^T] - \beta \beta^T \\ &= E_e[(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \delta] \beta^T + E_e[(\mathbf{X}'^T \mathbf{X}')^{-1} \mathbf{X}'^T \delta'] \beta^T \\ &+ E_e[((\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \delta)((\mathbf{X}'^T \mathbf{X}')^{-1} \mathbf{X}'^T \delta')^T], \end{aligned} \quad (\text{C1})$$

using the same variable definitions as in Appendix A. Prime indicators denote the equivalent random variable in the offspring generation. It is assumed that the expected vectors of average effects, β and β' are the same in each generation. In reality, this is likely to be a simplification, though reasonable over the span of two generations. We now demonstrate that the three terms in the third line are each equal to zero. Assuming no gene-environment correlations in the parental generation, the first term is zero, as

$$E_e[(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \delta] = E_e[(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T] E_e[\delta], \quad (\text{C2})$$

and $E_e[\delta] = 0$ by definition of the model. If there are no gene-environment correlations in the offspring generation, the second term is also zero, leaving the

final term,

$$E_e[(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \delta \delta'^T \mathbf{X}' (\mathbf{X}'^T \mathbf{X}')^{-1}], \quad (\text{C3})$$

where $\delta \delta'^T$ is an $N \times N'$ matrix (where N' is the number of individuals in the offspring generation) describing the environmental sample covariances between individuals in the parental and offspring generations. Element $\{a, b\}$ of the $M \times M$ matrix C3 is equal to

$$\begin{aligned} & E_e \left[\sum_{l=0}^M \sum_{i=1}^N (\mathbf{X}^T \mathbf{X})_{al}^{-1} x_{il} \delta_i \sum_{m=0}^M \sum_{j=1}^{N'} (\mathbf{X}'^T \mathbf{X}')_{bm}^{-1} x'_{jm} \delta'_j \right] \\ &= \sum_{l=0}^M \sum_{i=1}^N \sum_{m=0}^M \sum_{j=1}^{N'} E_e \left[(\mathbf{X}^T \mathbf{X})_{al}^{-1} x_{il} (\mathbf{X}'^T \mathbf{X}')_{bm}^{-1} x'_{jm} \right] E_e[\delta_i \delta'_j]. \end{aligned} \quad (\text{C4})$$

When there are no parent-offspring correlations in environment (or any environmental correlation structures between the two generations) $E_e[\delta_i \delta'_j] = 0$, so each element $\{a, b\}$ is zero, hence the third term on the right hand side of Equation C1 is zero. This completes that proof that the covariance matrix for the vectors $\hat{\beta}$ and $\hat{\beta}'$ is zero. It can then be verified, in a similar way, that this is also true for $\hat{\beta} E_i[x_{is} x_{it}]$ and $\hat{\beta}'$ (where $E_i[x_{is} x_{it}]$ is a scalar expectation derived from genetic variables s and t and taken across individuals i): the result that is required to produce Equation 10. Using Appendix A, $E_e[\hat{\beta} E_i[x_{is} x_{it}]] \beta^T = E_e[E_i[x_{is} x_{it}]] \beta \beta^T$, so the covariance matrix is

$$\begin{aligned} & E_e[E_i[x_{is} x_{it}] ((\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T w) ((\mathbf{X}'^T \mathbf{X}')^{-1} \mathbf{X}'^T w')^T] - E_e[E_i[x_{ij} x_{ik}]] \beta \beta^T \\ &= \sum_{l=0}^M \sum_{i=1}^N \sum_{m=0}^M \sum_{j=1}^{N'} E_e \left[E_i[x_{is} x_{it}] (\mathbf{X}^T \mathbf{X})_{al}^{-1} x_{il} (\mathbf{X}'^T \mathbf{X}')_{bm}^{-1} x'_{jm} \right] E_e[\delta_i \delta'_j] \quad (\text{C5}) \\ &= 0. \end{aligned}$$